

University of Liège
Faculty of Applied Sciences
Department of Electrical Engineering & Computer Science
Montefiore Institute

PhD Thesis in Engineering Sciences

**EXPLOITING RANDOM PROJECTIONS AND SPARSITY WITH
RANDOM FORESTS AND GRADIENT BOOSTING METHODS**

Application to multi-label and multi-output learning, random forest model
compression and leveraging input sparsity

ARNAUD JOLY



ADVISOR: LOUIS WEHENKEL
CO-ADVISOR: PIERRE GEURTS

December 2016

© 2017
ARNAUD JOLY
ALL RIGHTS RESERVED

JURY MEMBERS

DAMIEN ERNST, Professor at the Université de Liège (President);

LOUIS WEHENKEL, Professor at the Université de Liège (Advisor);

PIERRE GEURTS, Professor at the Université de Liège (Co-Advisor);

QUENTIN LOUVEAUX, Professor at Université de Liège;

ASHWIN ITTOO, Professor at the Université de Liège;

GRIGORIOS TSOUMAKAS, Professor at the Aristotle University of Thessaloniki;

CELINE VENS, Professor at the Katholieke Universiteit Leuven;

ABSTRACT

Within machine learning, the supervised learning field aims at modeling the input-output relationship of a system, from past observations of its behavior. Decision trees characterize the input-output relationship through a series of nested if – then – else questions, the testing nodes, leading to a set of predictions, the leaf nodes. Several of such trees are often combined together for state-of-the-art performance: random forest ensembles average the predictions of randomized decision trees trained independently in parallel, while tree boosting ensembles train decision trees sequentially to refine the predictions made by the previous ones.

The emergence of new applications requires scalable supervised learning algorithms in terms of computational power and memory space with respect to the number of inputs, outputs, and observations without sacrificing accuracy. In this thesis, we identify three main areas where decision tree methods could be improved for which we provide and evaluate original algorithmic solutions: (i) learning over high dimensional output spaces, (ii) learning with large sample datasets and stringent memory constraints at prediction time and (iii) learning over high dimensional sparse input spaces.

A first approach to solve *learning tasks with a high dimensional output space*, called binary relevance or single target, is to train one decision tree ensemble per output. However, it completely neglects the potential correlations existing between the outputs. An alternative approach called multi-output decision trees fits a single decision tree ensemble targeting simultaneously all the outputs, assuming that all outputs are correlated. Nevertheless, both approaches have (i) exactly the same computational complexity and (ii) target extreme output correlation structures. In our first contribution, we show how to combine random projection of the output space, a dimensionality reduction method, with the random forest algorithm decreasing the learning time complexity. The accuracy is preserved, and may even be improved by reaching a different bias-variance tradeoff. In our second contribution, we first formally adapt the gradient boosting ensemble method to multi-output supervised learning tasks such as multi-output regression and multi-label classification. We then propose to combine single random projections of the output space with gradient boosting on such tasks to adapt automatically to the output correlation structure.

The random forest algorithm often generates large ensembles of complex models thanks to the availability of a large number of observations. However, the space complexity of such models, proportional

to their total number of nodes, is often prohibitive, and therefore these models are not well suited under *stringent memory constraints at prediction time*. In our third contribution, we propose to compress these ensembles by solving a ℓ_1 -based regularization problem over the set of indicator functions defined by all their nodes.

Some supervised learning tasks have a *high dimensional but sparse input space*, where each observation has only a few of the input variables that have non zero values. Standard decision tree implementations are not well adapted to treat sparse input spaces, unlike other supervised learning techniques such as support vector machines or linear models. In our fourth contribution, we show how to exploit algorithmically the input space sparsity within decision tree methods. Our implementation yields a significant speed up both on synthetic and real datasets, while leading to exactly the same model. It also reduces the required memory to grow such models by exploiting sparse instead of dense memory storage for the input matrix.

RÉSUMÉ

Parmi les techniques d'apprentissage automatique, l'apprentissage supervisé vise à modéliser les relations entrée-sortie d'un système, à partir d'observations de son fonctionnement. Les arbres de décision caractérisent cette relation entrée-sortie à partir d'un ensemble hiérarchique de questions appelées les noeuds tests amenant à une prédiction, les noeuds feuilles. Plusieurs de ces arbres sont souvent combinés ensemble afin d'atteindre les performances de l'état de l'art: les ensembles de forêts aléatoires calculent la moyenne des prédictions d'arbres de décision randomisés, entraînés indépendamment et en parallèle alors que les ensembles d'arbres de boosting entraînent des arbres de décision séquentiellement, améliorant ainsi les prédictions faites par les précédents modèles de l'ensemble.

L'apparition de nouvelles applications requiert des algorithmes d'apprentissage supervisé efficaces en terme de puissance de calcul et d'espace mémoire par rapport au nombre d'entrées, de sorties, et d'observations sans sacrifier la précision du modèle. Dans cette thèse, nous avons identifié trois domaines principaux où les méthodes d'arbres de décision peuvent être améliorées pour lequel nous fournissons et évaluons des solutions algorithmiques originales: (i) apprentissage sur des espaces de sortie de haute dimension, (ii) apprentissage avec de grands ensembles d'échantillons et des contraintes mémoires strictes au moment de la prédiction et (iii) apprentissage sur des espaces d'entrée creux de haute dimension.

Une première approche pour résoudre des *tâches d'apprentissage avec un espace de sortie de haute dimension*, appelée «binary relevance» ou «single target», est l'apprentissage d'un ensemble d'arbres de décision par sortie. Toutefois, cette approche néglige complètement les corrélations potentiellement existantes entre les sorties. Une approche alternative, appelée «arbre de décision multi-sorties», est l'apprentissage d'un seul ensemble d'arbres de décision pour toutes les sorties, faisant l'hypothèse que toutes les sorties sont corrélées. Cependant, les deux approches ont (i) exactement la même complexité en temps de calcul et (ii) visent des structures de corrélation de sorties extrêmes. Dans notre première contribution, nous montrons comment combiner des projections aléatoires (une méthode de réduction de dimensionnalité) de l'espace de sortie avec l'algorithme des forêts aléatoires diminuant la complexité en temps de calcul de la phase d'apprentissage. La précision est préservée, et peut même être améliorée en atteignant un compromis biais-variance différent. Dans notre seconde contribution, nous adaptons d'abord formellement la méthode d'ensemble «gradient boosting» à la régression

multi-sorties et à la classification multi-labels. Nous proposons ensuite de combiner une seule projection aléatoire de l'espace de sortie avec l'algorithme de «gradient boosting» sur de telles tâches afin de s'adapter automatiquement à la structure des corrélations existant entre les sorties.

Les algorithmes de forêts aléatoires génèrent souvent de grands ensembles de modèles complexes grâce à la disponibilité d'un grand nombre d'observations. Toutefois, la complexité mémoire, proportionnelle au nombre total de noeuds, de tels modèles est souvent prohibitive, et donc ces modèles ne sont pas adaptés à des *contraintes mémoires fortes lors de la phase de prédiction*. Dans notre troisième contribution, nous proposons de compresser ces ensembles en résolvant un problème de régularisation basé sur la norme ℓ_1 sur l'ensemble des fonctions indicatrices défini par tous leurs noeuds.

Certaines tâches d'apprentissage supervisé ont un *espace d'entrée de haute dimension mais creux*, où chaque observation possède seulement quelques variables d'entrée avec une valeur non-nulle. Les implémentations standards des arbres de décision ne sont pas adaptées pour traiter des espaces d'entrée creux, contrairement à d'autres techniques d'apprentissage supervisé telles que les machines à vecteurs de support ou les modèles linéaires. Dans notre quatrième contribution, nous montrons comment exploiter algorithmiquement le creux de l'espace d'entrée avec les méthodes d'arbres de décision. Notre implémentation diminue significativement le temps de calcul sur des ensembles de données synthétiques et réelles, tout en fournissant exactement le même modèle. Cela permet aussi de réduire la mémoire nécessaire pour apprendre de tels modèles en exploitant des méthodes de stockage appropriées pour la matrice des entrées.

ACKNOWLEDGMENTS

This PhD thesis started with the trust granted by Prof. Louis Wehenkel, joined soon after by Prof. Pierre Geurts. I would like to express my sincere gratitude for their continuous encouragements, guidance and support. I have without doubt benefitted from their motivations, patience and knowledge. Our insightful discussions and interactions definitely moved the thesis forward.

I would like to thank the University of Liège, the FRS-FNRS, Belgium, the EU Network of Excellence PASCAL2, and the IUAP DYSCO, initiated by the Belgian State, Science Policy Office to have funded this research. Computational resources have been provided by the Consortium des Équipements de Calcul Intensif (CÉCI), funded by the Fonds de la Recherche Scientifique de Belgique (F.R.S.-FNRS) under Grant No. 2.5020.11.

The presented research would not have been the same without my co-authors (here in alphabetic order): Jean-Michel Begon, Mathieu Blondel, Lars Buitinck, Pierre Damas, Céline Delierneux, Damien Ernst, Hedayati Fares, Alexandre Gramfort, Pierre Geurts, André Gothot, Olivier Grisel, Jaques Grobler, Alexandre Hego, Bryan Holt, Justine Huart, Vincent François-Lavet, Nathalie Layios, Robert Layton, Christelle Lecut, Gilles Louppe, Andreas Mueller, Vlad Niculae, Cécile Oury, Panagiotis Papadimitriou, Fabian Pedregosa, Peter Prettenhofer, Zixiao Aaron Qiu, François Schnitzler, Antonio Sutera, Jake Vanderplas, Gael Varoquaux, and Louis Wehenkel.

I would like to thank the members of the jury, who take interests in my work, and took the time to read this dissertation.

Diane Zander and Sophie Cimino have been of an invaluable help with all the administrative procedures. I would like to thank them for their patience and availability. I would also like to thank David Colignon and Alain Empain for their helpfulness about anything related to super-computers.

I would like to thank my colleagues from the Montefiore Institute, Department of Electrical Engineering and Computer Science from the University of Liège, whom have created a pleasant, rich and stimulating environment (in alphabetic order): Samir Azrour, Tom Barbette, Julien Beckers, Jean-Michel Begon, Kyrylo Bessonov, Hamid Soleimani Bidgoli, Vincent Botta, Kridsakorn Chaichoompu, Célia Châtel, Julien Confetti, Mathilde De Becker, Renaud Detry, Damien Ernst, Ramouna Fouladi, Florence Fonteneau, Raphaël Fonteneau, Vincent François-Lavet, Damien Gérard, Quentin Gemine, Pierre Geurts, Samuel Hiard, Renaud Hoyoux, Fabien Heuze, Van Anh Huynh-Thu, Efthymios Karangelos, Philippe Latour, Gilles Louppe,

Francis Maes, Alejandro Marcos Alvarez, Benjamin Laugraud, Antoine Lejeune, Raphael Liégeois, Quentin Louveaux, Isabelle Mainz, Raphael Marée, Sébastien Mathieu, Axel Mathei, Romain Mormont, Frédéric Olivier, Julien Osmalsky, Sébastien Pierard, Zixiao Aaron Qiu, Loïc Rollus, Marie Schrynemackers, Oliver Stern, Benjamin Stévens, Antonio Sutera, David Taralla, François Van Lishout, Rémy Vandaele, Philippe Vanderbemden, and Marie Wehenkel.

I would like to thank the *scikit-learn* community who has shared with me their passion about computer science, machine learning and Python. By contributing to this open source project, I have learnt much since my first contribution.

I also offer my regards and blessing to all the people near and dear to my heart for their continuous support, and to all of those who supported in any respect during the completion of this project.

CONTENTS

1	INTRODUCTION	1
1.1	Publications	3
1.2	Outline	4
i	BACKGROUND	6
2	SUPERVISED LEARNING	7
2.1	Introduction	8
2.2	Classes of supervised learning algorithms	11
2.2.1	Linear models	12
2.2.2	(Deep) Artificial neural networks	15
2.2.3	Neighbors based methods	16
2.2.4	Decision tree models	17
2.2.5	From single to multiple output models	18
2.3	Evaluation of model prediction performance	23
2.4	Criteria to assess model performance	26
2.4.1	Metrics for binary classification	27
2.4.2	Metrics for multi-class classification	31
2.4.3	Metrics for multi-label classification and ranking	32
2.4.4	Regression metrics	36
2.5	Hyper-parameter optimization	38
2.6	Unsupervised projection methods	40
2.6.1	Principal components analysis	41
2.6.2	Random projection	43
2.6.3	Kernel functions	43
3	DECISION TREES	46
3.1	Decision tree model	48
3.2	Growing decision trees	50
3.2.1	Search among node splitting rules	51
3.2.2	Leaf labelling rules	55
3.2.3	Stop splitting criterion	55
3.3	Right decision tree size	57
3.4	Decision tree interpretation	58
3.5	Multi-output decision trees	60
4	BIAS-VARIANCE AND ENSEMBLE METHODS	63
4.1	Bias-variance error decomposition	64
4.2	Averaging ensembles	66
4.2.1	Variance reduction	67
4.2.2	Generic randomization induction methods	71
4.2.3	Randomized forest model	73
4.3	Boosting ensembles	74

4.3.1	Adaboost and variants	76
4.3.2	Functional gradient boosting	79
ii	LEARNING IN COMPRESSED SPACE THROUGH RANDOM PROJECTIONS	84
5	RANDOM PROJECTIONS OF THE OUTPUT SPACE	85
5.1	Methods	87
5.1.1	Multi-output regression trees in randomly projected output spaces	87
5.1.2	Exploitation in the context of tree ensembles	90
5.2	Bias/variance analysis	91
5.2.1	Single random trees.	92
5.2.2	Ensembles of t random trees.	94
5.3	Experiments	95
5.3.1	Effect of the size q of the Gaussian output space	95
5.3.2	Systematic analysis over 24 datasets	96
5.3.3	Input vs output space randomization	100
5.3.4	Alternative output dimension reduction techniques	101
5.3.5	Learning stage computing times	104
5.4	Conclusions	104
6	RANDOM OUTPUT SPACE PROJECTIONS FOR GRADIENT BOOSTING	106
6.1	Introduction	106
6.2	Gradient boosting with multiple outputs	108
6.2.1	Standard extension of gradient boosting to multi-output tasks	108
6.2.2	Adapting to the correlation structure in the output-space	111
6.2.3	Effect of random projections	114
6.2.4	Convergence when $M \rightarrow \infty$	120
6.3	Experiments	121
6.3.1	Experimental protocol	121
6.3.2	Experiments on synthetic datasets with known output correlation structures	123
6.3.3	Effect of random projection	130
6.3.4	Systematic analysis over real world datasets	136
6.4	Conclusions	143
iii	EXPLOITING SPARSITY FOR GROWING AND COMPRESSING DECISION TREES	145
7	ℓ_1-BASED COMPRESSION OF RANDOM FOREST MODELS	146

7.1	Compressing tree ensembles by ℓ_1 -norm regularization	147
7.2	Empirical analysis	149
7.2.1	Overall performances	150
7.2.2	Effect of the regularization parameter t .	151
7.2.3	Influence of the Extra-Tree meta parameters n_{\min} and M .	152
7.3	Conclusion	153
8	EXPLOITING INPUT SPARSITY WITH DECISION TREE	154
8.1	Tree growing	155
8.1.1	Standard node splitting algorithm	155
8.1.2	Splitting rules search on sparse data	157
8.1.3	Partitioning sparse data	163
8.2	Tree prediction	164
8.3	Experiments	165
8.3.1	Effect of the input space density on synthetic datasets	166
8.3.2	Effect of the input space density on real datasets	167
8.3.3	Algorithm comparison on 20 newsgroup	169
8.4	Conclusion	172
9	CONCLUSIONS	173
9.1	Conclusions	173
9.2	Perspectives and future works	175
9.2.1	Learning in compressed space through random projections	175
9.2.2	Growing and compressing decision trees	176
9.2.3	Learning in high dimensional and sparse input-output spaces	177
iv	APPENDIX	179
A	DESCRIPTION OF THE DATASETS	180
A.1	Synthetic datasets	180
A.2	Regression dataset	180
A.3	Multi-label dataset	180
A.4	Multi-output regression datasets	181
	BIBLIOGRAPHY	184

INTRODUCTION

Progress in information technology enables the acquisition and storage of growing amounts of rich data in many domains including science (biology, high-energy physics, astronomy, etc.), engineering (energy, transportation, production processes, etc.), and society (environment, commerce, etc.). Connected objects, such as smartphones, connected sensors or intelligent houses, are now able to record videos, images, audio signals, object localizations, temperatures, social interactions of the user through a social network, phone calls or user to computer program interactions such as voice help assistant or web search queries. The accumulating datasets come in various forms such as images, videos, time-series of measurements, recorded transactions, text etc. WEB technology often allows one to share locally acquired datasets, and numerical simulation often allows one to generate low cost datasets on demand. Opportunities exist thus for combining datasets from different sources to search for generic knowledge and enable robust decision.

All these rich datasets are of little use without the availability of automatic procedures able to extract relevant information from them in a principled way. In this context, the field of machine learning aims at developing theory and algorithmic solutions for the extraction of synthetic patterns of information from all kinds of datasets, so as to help us to better understand the underlying systems generating these data and hence to take better decisions for their control or exploitation.

Among the machine learning tasks, supervised learning aims at modeling a system by observing its behavior through samples of pairs of inputs and outputs. The objective of the generated model is to predict with high accuracy the outputs of the system given previously unseen inputs. A genomic application of supervised learning would be to model how a DNA sequence, a biological code, is linked to some genetic diseases. The samples used to fit the model are the input-output pairs obtained by sequencing the genome, the inputs, of patients with known medical records for the studied genetic diseases, the outputs. The objective is here twofold: (i) to understand how the DNA sequence influences the appearing of the studied genetic diseases and (ii) to use the predictive models to infer the probability of contracting the genetic disease.

The emergence of new applications, such as image annotation, personalized advertising or 3D image segmentation, leads to high dimensional data with a large number of inputs and outputs. It requires

scalable supervised learning algorithms in terms of computational power and memory space without sacrificing accuracy.

Decision trees (Breiman et al., 1984) are supervised learning models organized in the form of a hierarchical set of questions each one typically based on one input variable leading to a prediction. Used in isolation, trees are generally not competitive in terms of accuracy, but when combined into ensembles (Breiman, 2001; Friedman, 2001), they yield state-of-the-art performances on standard benchmarks (Caruana et al., 2008; Fernández-Delgado et al., 2014; Madjarov et al., 2012). They however suffer from several limitations that make them not always suited to address modern applications of machine learning techniques in particular involving high dimensional input and output spaces.

In this thesis, we identify three main areas where random forest methods could be improved and for which we provide and evaluate original algorithmic solutions: (i) learning over high dimensional output spaces, (ii) learning with large sample datasets and stringent memory constraints at prediction time and (iii) learning over high dimensional sparse input spaces. We discuss each one of these solutions in the following paragraphs.

HIGH DIMENSIONAL OUTPUT SPACES New applications of machine learning have multiple output variables, potentially in very high number (Agrawal et al., 2013; Dekel and Shamir, 2010), associated to the same set of input variables. A first approach to address such multi-output tasks is the so-called binary relevance / single target method (Spyromitros-Xioufis et al., 2016; Tsoumakas et al., 2009), which separately fits one decision tree ensemble for each output variable, assuming that the different output variables are independent. A second approach called multi-output decision trees (Blockeel et al., 2000; Geurts et al., 2006b; Kocev et al., 2013) fits a single decision tree ensemble targeting simultaneously all the outputs, assuming that all outputs are correlated. However in practice, (i) the computational complexity is the same for both approaches and (ii) we have often neither of these two extreme output correlation structures. As our first contribution, we show how to make random forest faster by exploiting random projections (a dimensionality reduction technique) of the output space. As a second contribution, we show how to combine gradient boosting of tree ensembles with single random projections of the output space to automatically adapt to a wide variety of correlation structures.

MEMORY CONSTRAINTS ON MODEL SIZE Even with a large number of training samples n , random forest ensembles have good computational complexity ($O(n \log n)$) and are easily parallelizable leading to the generation of very large ensembles. However, the resulting

models are big as the model complexity is proportional to the number of samples n and the ensemble size. As our third contribution, we propose to compress these tree ensembles by solving an appropriate optimization problem.

HIGH DIMENSIONAL SPARSE INPUT SPACES Some supervised learning tasks have very high dimensional input spaces, but only a few variables have non zero values for each sample. The input space is said to be “sparse”. Instances of such tasks can be found in text-based supervised learning, where each sample is often mapped to a vector of variables corresponding to the (frequency of) occurrence of all words (or multigrams) present in the dataset. The problem is sparse as the size of the text is small compared to the number of possible words (or multigrams). Standard decision tree implementations are not well adapted to treat sparse input spaces, unlike models such as support vector machines (Cortes and Vapnik, 1995; Scholkopf and Smola, 2001) or linear models (Bottou, 2012). Decision tree implementations are indeed treating these sparse variables as dense ones raising the memory needed. The computational complexity also does not depend upon the fraction of non zero values. As a fourth contribution, we propose an efficient decision tree implementation to treat supervised learning tasks with sparse input spaces.

1.1 PUBLICATIONS

This dissertation features several publications about random forest algorithms:

- (Joly et al., 2014) A. Joly, P. Geurts, and L. Wehenkel. *Random forests with random projections of the output space for high dimensional multi-label classification*. In Machine Learning and Knowledge Discovery in Databases, pages 607–622. Springer Berlin Heidelberg, 2014.
- (Joly et al., 2012) A. Joly, F. Schnitzler, P. Geurts, and L. Wehenkel. *L1-based compression of random forest models*. In European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning, 2012.
- (Buitinck et al., 2013) L. Buitinck, G. Louppe, M. Blondel, F. Pedregosa, A. Mueller, O. Grisel, V. Niculae, P. Prettenhofer, A. Gramfort, J. Grobler, R. Layton, J. Vanderplas, A. Joly, B. Holt, and G. Varoquaux. *Api design for machine learning software: experiences from the scikit-learn project*. arXiv preprint arXiv:1309.0238, 2013.

and also the following submitted article:

- H. Fares, A. Joly, and P. Papadimitriou. *Scalable Learning of Tree-Based Models on Sparsely Representable Data*.

Some collaborations were made during the thesis, but are not discussed within this manuscript:

- (Sutera et al., 2014) A. Sutera, A. Joly, V. François-Lavet, Z. A. Qiu, G. Louppe, D. Ernst, and P. Geurts. *Simple connectome inference from partial correlation statistics in calcium imaging*. In JMLR: Workshop and Conference Proceedings, pages 1–12, 2014.
- (Delierneux et al., 2015a) C. Delierneux, N. Layios, A. Hego, J. Huart, A. Joly, P. Geurts, P. Damas, C. Lecut, A. Gothot, and C. Oury. *Elevated basal levels of circulating activated platelets predict icu-acquired sepsis and mortality: a prospective study*. *Critical Care*, 19(Suppl 1):P29, 2015a.
- (Delierneux et al., 2015b) C. Delierneux, N. Layios, A. Hego, J. Huart, A. Joly, P. Geurts, P. Damas, C. Lecut, A. Gothot, and C. Oury. *Prospective analysis of platelet activation markers to predict severe infection and mortality in intensive care units*. In *journal of thrombosis and haemostasis*, volume 13, pages 651–651.
- (Begon et al., 2016) J.-M. Begon, A. Joly, and P. Geurts. *Joint learning and pruning of decision forests*. In *Belgian-Dutch Conference On Machine Learning*, 2016.

The following article has been submitted:

- C. Delierneux, N. Layios, A. Hego, J. Huart, C. Gosset, C. Lecut, N. Maes, P. Geurts, A. Joly, P. Lancellotti, P. Damas, A. Gothot, and C. Oury. *Incremental value of platelet markers to clinical variables for sepsis prediction in intensive care unit patients: a prospective pilot study*.

1.2 OUTLINE

In Part [i](#) of this thesis, we start by introducing in Chapter [2](#) the key concepts about supervised learning: (i) what are the most popular supervised learning models, (ii) how to assess the prediction performance of a supervised learning model and (iii) how to optimize the hyper-parameters of these models. We also present some unsupervised projection methods, such as random projections, which transform the original space to another one. We describe more in detail the decision tree model classes in Chapter [3](#). More specifically, we describe the methodology to grow and to prune such trees. We also show how to adapt decision tree growing and prediction algorithms to multi-output tasks. In Chapter [4](#), we show why and how to combine models into ensembles either by learning models independently with averaging methods or sequentially with boosting methods.

In Part [ii](#), we first show how to grow an ensemble of decision trees on very high dimensional output spaces by projecting the original output space onto a random sub-space of lower dimension. In [Chapter 5](#), it turns out that for random forest models, an averaging ensemble of decision trees, the learning time complexity can be reduced without affecting the prediction performance. Furthermore, it may lead to accuracy improvement ([Joly et al., 2014](#)). In [Chapter 6](#), we propose to combine random projections of the output space and the gradient tree boosting algorithm, while reducing learning time and automatically adapting to any output correlation structure.

In Part [iii](#), we leverage sparsity in the context of decision tree ensembles. In [Chapter 7](#), we exploit sparsifying optimization algorithms to compress random forest models while retaining their prediction performances ([Joly et al., 2012](#)). In [Chapter 8](#), we show how to leverage input sparsity to speed up decision tree induction.

During the thesis, I made significant contributions to the open source scikit-learn project ([Buitinck et al., 2013](#); [Pedregosa et al., 2011](#)) and developed my own open source libraries `random-output-trees`¹, containing the work presented in [Chapter 5](#) and [Chapter 6](#), and `clusterlib`², containing the tools to manage jobs on supercomputers.

¹ <https://github.com/arjoly/random-output-trees>

² <https://github.com/arjoly/clusterlib>

Part I

BACKGROUND

OUTLINE

In the field of machine learning, supervised learning aims at finding the best function which describes the input-output relation of a system only from observations of this relationship. Supervised learning problems can be broadly divided into classification tasks with discrete outputs and into regression tasks with continuous outputs. We first present major supervised learning methods for both classification and regression. Then, we show how to estimate their performance and how to optimize the hyper-parameters of these models. We also introduce unsupervised projection techniques used in conjunction with supervised learning methods.

Supervised learning aims at modeling an input-output system from observations of its behavior. The applications of such learning methods encompass a wide variety of tasks and domains ranging from image recognition to medical diagnosis tools. Supervised learning algorithms analyze the input-output pairs and learn how to predict the behavior of a system (see Figure 2.1) by observing its responses, described by output variables y_1, \dots, y_d , also called targets, to its environment described by input variables x_1, \dots, x_p , also called features. The outcome of the supervised learning is a function f modeling the behavior of the system.

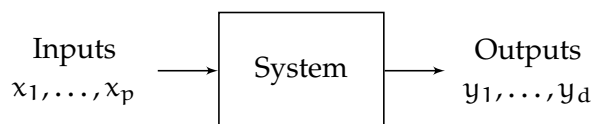


Figure 2.1: Input-output view of a system.

Supervised learning has numerous applications in the multimedia, in biology, in engineering or in the societal domain:

- Identification of digits from photos, such as house number from street photos or digit post code from letters.
- Automatic image annotation such as detecting tumorous cells or identifying people in photos.
- Detection of genetic diseases from DNA screening.
- Disease diagnostic based on clinical and biological data of a patient.

- Automatic text translation from a source language to a target language such as from French to English.
- Automatic voice to text transcription from audio records.
- Market price prediction on the basis of economical and performance indicators.

We introduce the supervised learning framework in Section 2.1. We describe in Section 2.2 the most common classes of supervised learning models used to map the outputs of the system to its inputs. We introduce how to assess their performances in Section 2.3, how to compare the model predictions to a ground truth in Section 2.4 and how to select the best hyper-parameters of such models in Section 2.5. We also show some input space projection methods in Section 2.6, often used in combination with supervised learning models improving the computational time and / or the accuracy of the model.

2.1 INTRODUCTION

The goal of supervised learning is to learn the function f mapping an input vector $x = (x_1, \dots, x_p)$ of a system to a vector of system outputs $y = (y_1, \dots, y_d)$, only from observations of input-output pairs. The set of possible input (resp. output) vectors form the input space \mathcal{X} (resp. output space \mathcal{Y}).

Once we have identified the input and output variables, we start to collect input-output pairs, also called samples. Table 2.1 displays 5 samples collected from a system with 4 inputs and 3 outputs. We distinguish three types of variables: binary variables taking only two different values, like the variables x_1 and y_1 ; categorical variables taking two or more possible values, like variables x_2 and y_2 , and numerical variables having numerical values, like x_2 , x_4 and y_3 . A binary variable is also a categorical variable. For simplicity, we will assume in the following without loss of generality that binary and categorical variables have been mapped from the set of their k original values to a set of integers of the same cardinality $\{0, \dots, k - 1\}$.

Table 2.1: A dataset formed of samples of pairs of four inputs and three outputs.

x_1	x_2	x_3	x_4	y_1	y_2	y_3
0	0.25	A	0.25	True	Small	1.8
?	-2	B	3.	True	Average	1.7
0	3	C	2.	False	?	1.65
1	10.7	?	-3.	False	Big	1.59
1	0.	A	2.	False	Big	?

When we collect data, some input and/or output values might be missing or unavailable. Tasks with missing input values are said to have missing data. Missing values are marked by a “?” in Table 2.1.

We classify supervised learning tasks into two main families based on their output domains. Classification tasks have either binary outputs as in disease prediction ($y \in \{\text{Healthy, sick}\}$) or categorical outputs as in digits recognition ($y \in \{0, \dots, 9\}$). Regression tasks have numerical outputs ($y \in \mathbb{R}$) such as in house price predictions. A classification task with only one binary output (resp. categorical output) is called a binary classification task (resp. multi-class classification task). A multi-class classification task is assumed to have more than two classes, otherwise it is a binary classification task. In the presence of multiple outputs, we further distinguish multi-label classification tasks which associate multiple binary output values to each input vector. In the multi-label context, the output variables are also called “labels” and the output vectors are called “label set”. From a modeling perspective, multi-class classification tasks are multi-label classification problems whose labels are mutually exclusive. Table 2.2 summarizes the different supervised learning tasks.

Table 2.2: The output domain determines the supervised learning task.

Supervised learning task	Output domain
Binary classification	$\mathcal{Y} = \{0, 1\}$
Multi-class classification	$\mathcal{Y} = \{0, 1, \dots, k-1\}$ with $k > 2$
Multi-label classification	$\mathcal{Y} = \{0, 1\}^d$ with $d > 1$
Multi-output multi-class classification	$\mathcal{Y} = \{0, 1, \dots, k-1\}^d$ with $k > 2, d > 1$
Regression	$\mathcal{Y} = \mathbb{R}$
Multi-output regression	$\mathcal{Y} = \mathbb{R}^d$ with $d > 1$

We will denote by \mathcal{X} an input space, and by \mathcal{Y} an output space. We denote by $P_{\mathcal{X}, \mathcal{Y}}$ the joint (unknown) sampling density over $\mathcal{X} \times \mathcal{Y}$. Superscript indices (x^i, y^i) denote (input, output) vectors of an observation $i \in \{1, \dots, n\}$. Subscript indices (e.g. x_j, y_k) denote components of vectors. With these notations supervised learning can be defined as follows:

SUPERVISED LEARNING

Given a learning sample $((x^i, y^i) \in (\mathcal{X} \times \mathcal{Y}))_{i=1}^n$ of n observations in the form of input-output pairs, a supervised learning task is defined as searching for a function $f^* : \mathcal{X} \rightarrow \mathcal{Y}$ in a hypothesis space $\mathcal{H} \subset \mathcal{Y}^{\mathcal{X}}$

Table 2.3: Common losses to measure the discrepancy between a ground truth y and either a prediction or a score y' . In classification, we assume here that the ground truth y is encoded with $\{-1, 1\}$ values.

Regression loss	
Square loss	$\ell(y, y') = \frac{1}{2}(y - y')^2$
Absolute loss	$\ell(y, y') = y - y' $
Binary classification loss	
0-1 loss	$\ell(y, y') = 1(y \neq y')$
Hinge loss	$\ell(y, y') = \max(0, 1 - yy')$
Logistic loss	$\ell(y, y') = \log(1 + \exp(-2yy'))$

that minimizes the expectation of some loss function $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$ over the joint distribution of input / output pairs:

$$f^* = \arg \min_{f \in \mathcal{H}} \mathbb{E}_{P_{x,y}} \{\ell(f(x), y)\}. \quad (2.1)$$

The choice of the loss function ℓ depends on the property of the supervised learning task (see Table 2.3 for their definitions):

- In regression ($\mathcal{Y} = \mathbb{R}$), we often use the squared loss, except when we want to be robust to the presence of outliers, samples with abnormal output values, where we prefer other losses such as the absolute loss.
- In classification tasks ($\mathcal{Y} = \{0, \dots, k-1\}$), the reported performance is commonly the average 0-1 loss, called the error rate. However, the model does not often directly minimize the 0-1 loss as it leads to non convex and non continuous optimization problems with often high computational cost. Instead, we can relax the multi-class or the binary constraint by optimizing a smoother loss such as the hinge loss or the logistic loss. To get a binary or multi-class prediction, we can threshold the predicted value $f(x)$.

Figure 2.2 plots several loss discrepancies $\ell(1, y')$ whenever the ground truth is $y = 1$ and as a function of the value y' predicted by the model. The 0-1 loss is a step function with a discontinuity at $y' = 1$. The hinge loss has a linear behavior whenever $y' \leq 1$ and is a constant with $y' \geq 1$. The logistic loss strongly penalizes any mistake and is zero only if the model is correct with an infinite score. The plot also highlights that we can use regression losses for classification

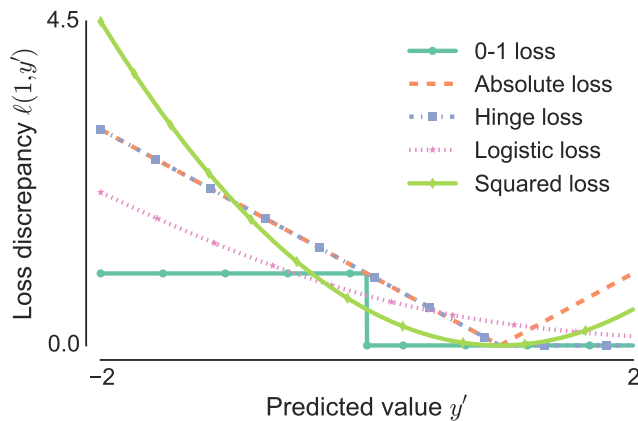


Figure 2.2: Loss discrepancies $\ell(1, y')$ with $y = 1$. (Adapted from (Hastie et al., 2009))

tasks. It shows that regression losses penalize any predicted value y' different from the ground truth y . However, this is not always the desired behavior. For instance whenever $y = 1$ (resp. $y = 0$), regression losses penalize any score greater than $y' > 1$ (resp. smaller than $y' < 0$), while the model truly believes that the output is positive (resp. negative). This is often the reason why regression losses are avoided for classification tasks.

2.2 CLASSES OF SUPERVISED LEARNING ALGORITHMS

Supervised learning aims at finding the best function f in a hypothesis space \mathcal{H} to model the input-output function of a system. If there is no restriction on the hypothesis space \mathcal{H} , the model f can be any function $f \in y^x$.

Consider a binary function f which has p binary inputs. The binary function is uniquely defined by knowing the output values of the 2^p possible input vectors. The hypothesis space of all binary functions contains 2^{2^p} binary functions. If we observe n different input-output pair assignments, there remain $2^{2^p - n}$ possible binary functions. For a binary function of $p = 5$ inputs, we have $2^{2^5} = 4294967296$ possible binary functions. If we observe $n = 16$ input-output pair assignments among the $2^5 = 32$ possible ones, we still have $2^{2^5 - 16} = 65536$ possible binary functions. The number of possible functions highly increases with the cardinality of each variable. The hypothesis space will be even larger with a stochastic function, where different output values are possible for each possible input assignment.

By making assumptions on the model class \mathcal{H} , we can largely reduce the size of the hypothesis space. For instance in the previous example, if we assume that 2 out of the 5 binary input variables are independent of the output, there remain $2^{2^3} = 256$ possible functions.

The correct function would be uniquely identified by observing the 8 possible assignments.

Given the data stochasticity, those model classes can directly model the input-output mapping $f : \mathcal{X} \rightarrow \mathcal{Y}$, but also the conditional probability $P(y|x)$ and predictions are made through

$$f(x) = \arg \min_{\hat{y} \in \mathcal{Y}} E_{y|x} [L(y, \hat{y})] = \arg \min_{\hat{y}} \int_{\mathcal{Y}} L(y, \hat{y}) dP(y|x). \quad (2.2)$$

We will present some of the most popular model classes: linear models in Section 2.2.1; artificial neural networks in Section 2.2.2 which are inspired from the neurons in the brain; neighbors-based models in Section 2.2.3 which find the nearest samples in the training set; decision tree based-models in Section 2.2.4 (and in more details in Chapter 3). Note that we introduce ensemble methods in Section 2.2.4 and discuss them more deeply in Chapter 4.

2.2.1 Linear models

Let us denote by $x \in \mathbb{R}^p = [x_1 \ \dots \ x_p]^T$ a vector of input variables. A linear model \hat{f} is a model having the following form

$$\hat{f}(x) = \beta_0 + \sum_{j=1}^p \beta_j x_j = \begin{bmatrix} 1 & x^T \end{bmatrix} \beta$$

where the vector $\beta \in \mathbb{R}^{1+p}$ is a concatenation of the intercept β_0 and the coefficients β_j .

Given a set of n input-output pairs $\{(x^i, y^i) \in (\mathcal{X}, \mathcal{Y})\}_{i=1}^n$, we retrieve the coefficient vector β of the linear model by minimizing a loss $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$:

$$\min_{\beta} \sum_{i=1}^n \ell(y^i, \hat{f}(x^i)) = \min_{\beta} \sum_{i=1}^n \ell(y^i, \begin{bmatrix} 1 & x^{iT} \end{bmatrix} \beta). \quad (2.3)$$

With the square loss $\ell(y, y') = \frac{1}{2}(y - y')^2$, there exists an analytical solution to Equation 2.3 called ordinary least squares. Let us denote by $\mathbf{X} \in \mathbb{R}^{n \times (1+p)}$ the concatenation of the input vectors with a first column of \mathbf{X} full of ones to model the intercept β_0 and by $\mathbf{y} \in \mathbb{R}^n$ the concatenation of the output values. We can now express the sum of squares in matrix notation:

$$\sum_{i=1}^n \ell(y, \hat{f}(x^i)) = \frac{1}{2} \sum_{i=1}^n (y^i - \hat{f}(x^i))^2 \quad (2.4)$$

$$= \frac{1}{2} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) \quad (2.5)$$

The first order differentiation of the sum of squares with respect to β yields to

$$\frac{\partial}{\partial \beta} \sum_{i=1}^n \ell(y, \hat{f}(x^i)) = \mathbf{X}^T(\mathbf{y} - \mathbf{X}\beta). \quad (2.6)$$

The vector minimizing the square loss is thus

$$\beta = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}. \quad (2.7)$$

The solution exists only if $\mathbf{X}^T \mathbf{X}$ is invertible.

Whenever the number of inputs plus one $p + 1$ is greater than the number of samples n , the analytical solution is ill posed as the matrix $\mathbf{X}^T \mathbf{X}$ is rank deficient ($\text{rank}(\mathbf{X}^T \mathbf{X}) < p + 1$). To ensure a unique solution, we can add a regularization penalty R with a multiplying constant $\lambda \in \mathbb{R}^+$ on the coefficients β of the linear model:

$$\min_{\beta} \sum_{i=1}^n L(y^i, \beta_0 + \sum_{j=1}^p \beta_j x_j^i) + \lambda R(\beta_1, \dots, \beta_p). \quad (2.8)$$

With a ℓ_2 -norm constraint on the coefficients, we transform the ordinary least square model into a ridge regression model (Hoerl and Kennard, 1970):

$$\min_{\beta} \sum_{i=1}^n \left(y^i - \beta_0 - \sum_{j=1}^p \beta_j x_j^i \right)^2 + \lambda \sum_{j=1}^p \beta_j^2. \quad (2.9)$$

One can show (see Section 3.4.1 of (Hastie et al., 2009)) that the constant λ controls the maximal value of all coefficients β_j in the ridge regression solution.

With a ℓ_1 -norm constraint ($R(\beta_1, \dots, \beta_p) = \sum_{j=1}^p |\beta_j|$) on the coefficient β_1, \dots, β_p , we have the Lasso model (Tibshirani, 1996b):

$$\min_{\beta} \sum_{i=1}^n \left(y^i - \beta_0 - \sum_{j=1}^p \beta_j x_j^i \right)^2 + \lambda \sum_{j=1}^p |\beta_j|. \quad (2.10)$$

Contrarily to the ridge regression, the Lasso has no closed formed analytical solution even though the resulting optimization problem remains convex. However, we gain that the ℓ_1 -norm penalty sparsifies the coefficients β_j of the linear model. If the constant λ tends towards infinity, all the coefficients will be zero $\beta_j = 0$. While with $\lambda = 0$, we have the ordinary least square formulation. With λ moving from $+\infty$ to 0, we progressively add variables to the linear model with a magnitude $\sum_{j=1}^p |\beta_j|$ depending on λ . The monotone Lasso (Hastie et al., 2007) further restricts the coefficient to monotonous variation with respect to λ and has been shown to perform better whenever the input variables are correlated.

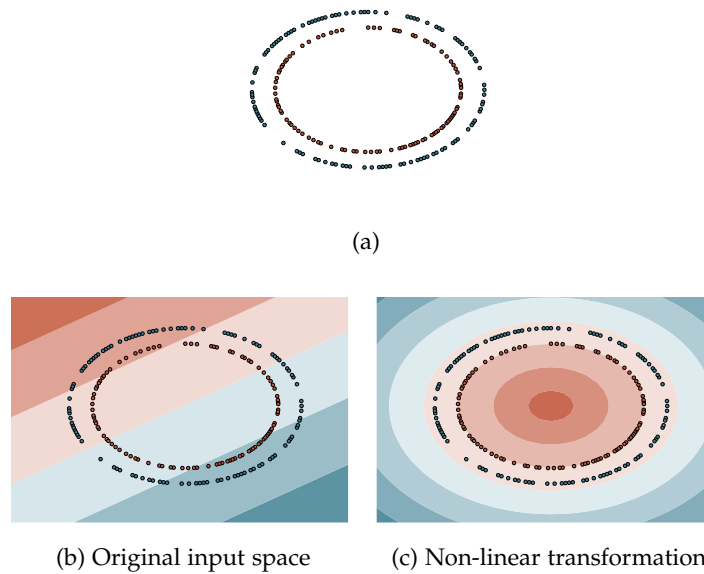


Figure 2.3: The logistic linear model on the bottom left is unable to find a separating hyperplane. If we fit the linear model on the distance from the center of circle, we separate perfectly both classes as shown in the bottom right.

A combination of the ℓ_1 -norm and the ℓ_2 -norm constraints on the coefficients is called an elastic net penalty (Zou and Hastie, 2005). It shares both the property of the Lasso and the ridge regression: sparsely selecting coefficients as in Lasso and considering groups of correlated variables together as in the ridge regression. With a careful design of the penalty term R , we can enforce further properties such as selecting variables in groups of pre-defined variables with the group Lasso (Meier et al., 2008; Yuan and Lin, 2006) or taking into account the variable locality in the coefficient vector β while adding a new variable to the linear model with the fused Lasso (Tibshirani et al., 2005).

By selecting an appropriate loss and penalty term, we have a wide variety of linear models at our disposal with different properties. In regression, an absolute loss leads to the least absolute deviation algorithm (Bloomfield and Steiger, 2012) which is robust to outliers. In classification, we can use a logistic loss to model the class probability distribution leading to the logistic regression model. With a hinge loss, we aim at finding a hyperplane which maximizes the separations between the classes leading to the support vector machine algorithm (Cortes and Vapnik, 1995).

A linear model can handle non linear problems by applying first a non linear transformation to the input space \mathcal{X} . For instance, consider the classification task of Figure 2.3a where each class is located on a concentric circle. Given the non linearity of the problem, we can not find a straight line separating both classes in the cartesian plane

as shown in Figure 2.3b. If instead we fit a linear model on the distance from the origin $\sqrt{x_1^2 + x_2^2}$ as illustrated in Figure 2.3c, we find a model separating perfectly both classes. We often use linear models in conjunction with kernel functions (presented in Section 2.6.3), which provide a range of ways to achieve non-linear transformations of the input space.

2.2.2 (Deep) Artificial neural networks

An artificial neural network is a statistical model mimicking the structure of the brain and composed of artificial neurons. A neuron, as shown in Figure 2.4a, is composed of three parts: the *soma*, the cell body, processes the information from its *dendrites* and transmits its results to other neurons through the *axon*, a nerve fiber. An artificial neuron follows the same structure (see Figure 2.4b) replacing biological processing by numerical computations. The basic neuron (Rosenblatt, 1958) used for supervised learning consists in a linear model of parameters $\beta \in \mathbb{R}^{p+1}$ followed by an activation function ϕ :

$$\hat{f}_{\text{neuron}}(x) = \phi \left(\beta_0 + \sum_{j=1}^p \beta_j x_j \right).$$

The activation function replicates artificially the non linear activation of real neurons. It is a scalar function such as a hyperbolic tangent $\phi(x) = \tanh(x)$, a sigmoid $\phi(x) = (1 + e^{-x})^{-1}$ or a rectified linear function $\phi(x) = \max(0, x)$.

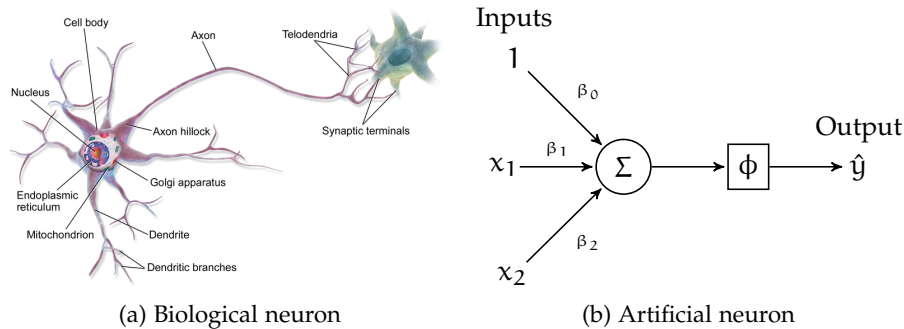


Figure 2.4: A biological neuron (on the left) and an artificial neuron (on the right).

More complex artificial neural networks are often structured into layers of artificial neurons. The inputs of a layer are the input variables or the outputs of the previous layer. Each neuron of the layer has one output. The neural network is divided into three parts as in Figure 2.5: the first and last layers are respectively the *input layer* and the *output layer*, while the layers in between are the *hidden layers*. The hidden layer of Figure 2.5 is called a fully connected layer as all the

neurons (here the input variables) from the previous layer are connected to each neuron of the layer. Other layer structures exist such as convolutional layers (Krizhevsky et al., 2012; LeCun et al., 2004) which mimic the visual cortex (Hubel and Wiesel, 1968). A network is not necessarily feed forward, but can have a more complex topology for example recurrent neural networks (Boulanger-Lewandowski et al., 2012; Graves et al., 2013) mimic the brain memory by forming internal cycles of neurons. Neural networks with many layers are also known (LeCun et al., 2015) as deep neural networks.

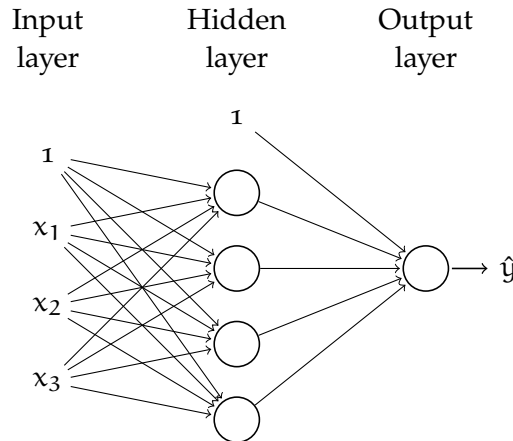


Figure 2.5: A neural network with an input layer, a fully connected hidden layer and an output layer.

Artificial neurons form a graph of variables. Through this representation, we can learn such models by applying gradient based optimization techniques (Bengio, 2012; Glorot and Bengio, 2010; LeCun et al., 2012) to find the coefficient vector associated to each neuron minimizing a given loss function.

2.2.3 Neighbors based methods

The k -nearest neighbors model is defined by a distance metric d and a set of samples. At learning time, those samples are stored in a database. We predict the output of an unseen sample by aggregating the outputs of the k -nearest samples in the input space according to the distance metric d , with k being a user-defined parameter.

More precisely, given a training set $((x^i, y^i) \in (\mathcal{X} \times \mathcal{Y}))_{i=1}^n$ and a distance measure $d : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}^+$, an unseen sample with value in the input space x is assigned a prediction through the following procedure:

1. Compute the distances $d(x^i, x)$ in the input space, $\forall i = 1, \dots, n$, between the training samples x^i and the input vector x .
2. Search for the k samples in the training set which have the smallest distance to the vector x .

- In classification, compute the proportion of samples of each class among these k -nearest neighbors: the final prediction is the class with the highest proportion. This corresponds to a majority vote over the k nearest neighbors. In regression, the prediction is the average output of the k -nearest neighbors.

The k -nearest neighbor method adapts to a wide variety of scenarios by selecting or by designing a proper distance metric such as the euclidean distance or the Hamming distance.

2.2.4 Decision tree models

A decision tree model is a hierarchical set of questions leading to a prediction. The internal nodes, also called test nodes, test the value of a feature. In Figure 2.6, the starting node, also called root node, tests whether the feature “Petal width” is bigger or smaller than 0.7cm. According to the answer, you follow either the right branch ($> 0.7\text{cm}$) leading to another test node or the left branch ($\leq 0.7\text{cm}$) leading to an external node, also called a leaf. To predict an unseen sample, you start at the root node and follow the tree structure until reaching a leaf labelled with a prediction. With the decision tree of Figure 2.6, an iris with petal width smaller than 0.7cm is an iris Setosa.

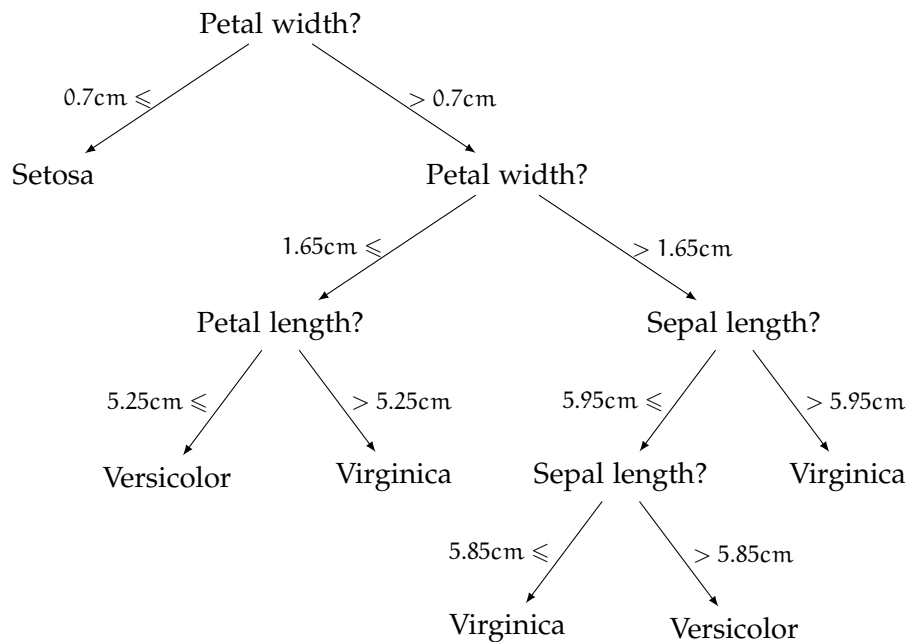


Figure 2.6: A decision tree classifying iris flowers into its Setosa, Versicolor or Virginica varieties according to the width and length of its petals and sepals.

A classification or a regression tree (Breiman et al., 1984) is built using all the input-output pairs $((x^i, y^i) \in (\mathcal{X} \times \mathcal{Y}))_{i=1}^n$ as follows: for each test node, the best split (S_r, S_l) of the local subsample S reaching

the node is chosen among the p input features combined with the selection of an optimal cut point. The best sample split (S_r, S_l) of S minimizes the average reduction of impurity

$$\begin{aligned} & \Delta I((y^i)_{i \in S}, (y^i)_{i \in S_l}, (y^i)_{i \in S_r}) \\ &= I((y^i)_{i \in S}) - \frac{|S_l|}{|S|} I((y^i)_{i \in S_l}) - \frac{|S_r|}{|S|} I((y^i)_{i \in S_r}), \end{aligned} \quad (2.11)$$

where I is the impurity of the output such as the entropy in classification or the variance in regression. The decision tree growth continues until we reach a stopping criterion such as no impurity $I((y^i)_{i \in S}) = 0$.

To avoid over-fitting, we can stop earlier the tree growth by adding further stopping criteria such as a maximal depth or a minimal number of samples to split a node.

Instead of a single decision tree, we often train an ensemble of such models:

- Averaging-based ensemble methods grow an ensemble by randomizing the tree growth. The random forest method (Breiman, 2001) trains decision trees on bootstrap copies of the training set, i.e. by sampling with replacement from the training dataset, and it randomizes the best split selection by searching this split among k out of the p features at each nodes ($k \leq p$).
- Boosting-based methods (Freund and Schapire, 1997; Friedman, 2001) build iteratively a sequence of weak models such as shallow trees which perform only slightly better than random guessing. Each new model refines the prediction of the ensemble by focusing on the wrongly predicted training input-output pairs.

We further discuss decision tree models in Chapter 3 and ensemble methods in Chapter 4.

2.2.5 From single to multiple output models

With multiple outputs supervised learning tasks, we have to infer the values of a set of d output variables y_1, \dots, y_d (instead of a single one) from a set of p input variables x_1, \dots, x_p . We hope to improve the accuracy and / or computational performance by exploiting the correlation structure between the outputs. There exist two main approaches to solve multiple output tasks: problem transformation presented in Section 2.2.5.1 and algorithm adaptation in Section 2.2.5.2. We present here a non exhaustive selection of both approaches. The interested reader will find a broader review of the multi-label literature in (Gibaja and Ventura, 2014; Madjarov et al., 2012; Tsoumakas et al., 2009; Zhang and Zhou, 2014) and of the multi-output regression literature in (Borchani et al., 2015; Spyromitros-Xioufis et al., 2016).

2.2.5.1 Problem transformation

The problem transformation approach transforms the original multi-output task into a set of single output tasks. Each of these single output tasks is then solved by classical classifiers or regressors. The possible output correlations are exploited through a careful reformulation of the original task.

INDEPENDENT ESTIMATORS The simplest way to handle multi-output learning is to treat all outputs in an independent way. We break the prediction of the d outputs into d independent single output prediction tasks. A model is fitted on each output. At prediction time, we concatenate the predictions of these d models. This is called the binary relevance method (Tsoumakas et al., 2009) in multi-label classification and the single target method (Spyromitros-Xioufis et al., 2016) in multi-output regression. Since we consider the outputs independently, we neglect the output correlation structure. Some methods may however benefit from sharing identical computations needed for the different outputs. For instance, the k -nearest neighbor method can share the search for the k -nearest neighbors in the input space, and the ordinary linear least squares method can share the computation of $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ in Equation 2.7.

ESTIMATOR CHAIN If the outputs are dependent, the model of a single output might benefit from the values of the correlated outputs. In the estimator chain method, we sequentially learn a model for each output by providing the predictions of the previously learnt models as auxiliary inputs. This is called a classifier chain (Read et al., 2011) in classification and a regressor chain (Spyromitros-Xioufis et al., 2016) in regression.

More precisely, the estimator chain method first generates an order o on the outputs for instance based on prior knowledge, the output density, the output variance or at random. Then with the training samples and the output order o , it sequentially learns d estimators: the l -th estimator f_{o_l} aims at predicting the o_l -th output using as inputs the concatenation of the input vectors with the predictions of the models learnt for the $l - 1$ previous outputs. To reduce the model variance, we can generate an ensemble of estimator chains by randomizing the chain order (and / or the underlying base estimator), and then we average their predictions.

In multi-label classification, Cheng et al. (2010) formulates a Bayes optimal classifier chain by modeling the conditional probability of $P_{y|x}(y|x)$. Under the chain rule, we have

$$P_{y|x}(y|x) = P_{y_1|x}(y_1|x) \prod_{j=2}^d P_{y_j|x, y_1, \dots, y_{j-1}}(y_j|x, y_1, \dots, y_{j-1}). \quad (2.12)$$

Each estimator of the chain approximates a probability factor of the chain rule decomposition. Using the estimation of $P_{y|x}$ made by the chain and a given loss function ℓ , we can perform Bayes optimal prediction:

$$h^*(x) = \arg \min_{y'} E_{y|x} \ell(y', y). \quad (2.13)$$

ERROR CORRECTING CODES Error correcting codes are techniques from information and coding theory used to properly deliver a message through a noisy channel. It first codes the original message, and then corrects the errors made during the transmission at decoding time. This idea have been applied to multi-class classification (Dietterich and Bakiri; Guruswami and Sahai, 1999), multi-label classification (Cisse et al., 2013; Ferng and Lin, 2011; Guo et al., 2008; Hsu et al., 2009; Kajdanowicz and Kazienko, 2012; Kapoor et al., 2012; Kouzani and Nasireding, 2009; Zhang and Schneider, 2011) and multi-output regression (Tsoumakas et al., 2014; Yu et al., 2006) tasks by viewing the predictions made by the supervised learning model(s) as a message transmitted through a noisy channel. It transforms the original task by encoding the output values with a binary error correcting code or output projections. One classifier is then fitted for each bit of the code or output projection. At prediction time, we concatenate the predictions made by each estimator and decode them by solving the inverse problem. Note that the output coding might also have for objective to reduce the dimensionality of the output space (Hsu et al., 2009; Kapoor et al., 2012).

PAIRWISE COMPARISONS In multi-label tasks, the ranking by pairwise comparison approach (Hüllermeier et al., 2008) aims to generate a ranking of the labels by making all the pairwise label comparisons. The original tasks is transformed into $d(d-1)/2$ binary classification tasks where we compare if a given label is more likely to appear than another label. The datasets comparing each label pair is obtained by collecting all the samples where only one of the outputs is true, but not both. This approach is similar to the one-versus-one approach (Park and Fürnkranz, 2007) in multi-class classification task, however we can not directly transform the ranking into a prediction, i.e. label set. To decrease the prediction time, alternative ranking construction schemes have been proposed (Mencia and Fürnkranz, 2008; Mencia and Fürnkranz, 2010) requiring less than $d(d-1)/2$ classifier predictions.

The Calibrated label ranking method (Brinker et al., 2006; Fürnkranz et al., 2008) extends the previous approach by adding a virtual label which will serve as a split point between the true and the false labels. For each label, we add a new tasks using all the samples comparing the label i to the virtual label whose value is the opposite of the label i . To the $d(d-1)/2$ tasks, we effectively add d tasks.

LABEL POWER SET For multi-label classification tasks, the label power set method (Tsoumakas et al., 2009) encodes each label set in the training set as a class. It transforms the original task into a multi-class classification task. At prediction time, the class predicted by the multi-class classifier is decoded thanks to the one-to-one mapping of the label power set encoding. The drawback of this approach is to generate a large number of classes due to the large number of possible label sets. For n samples and d labels, the maximal number of classes is $\max(2^d, n)$. This leads to accuracy issues if some label sets are not well represented in the training set. To alleviate the explosion of classes, raket (Tsoumakas and Vlahavas, 2007) generates an ensemble of multi-class classifiers by subsampling the output space and then applying the label power set transformation.

2.2.5.2 Algorithm adaptation

The algorithm adaptation approach modifies existing supervised learning algorithms to handle multiple output tasks. We show here how to extend the previously presented models classes to multi-output regression and to multi-label classification tasks.

LINEAR-BASED MODELS Linear-based models have been adapted to multi-output tasks by reformulated their mathematical formulation using multi-output losses and (possibly) regularization constraints enforcing assumptions on the input-output and the output-output correlation structures. The proposed methods are based for instance on extending least-square regression (Baldassarre et al., 2012; Breiman and Friedman, 1997; Dayal and MacGregor, 1997; Evgeniou et al., 2005; Similä and Tikka, 2007; Zhou and Tao, 2012) (with possibly regularization), canonical correlation analysis (Izenman, 1975; Van Der Merwe and Zidek, 1980), support vector machine (Elisseeff and Weston, 2001; Evgeniou and Pontil, 2004; Evgeniou et al., 2005; Jiang et al., 2008; Xu, 2012), support vector regression (Liu et al., 2009; Sánchez-Fernández et al., 2004; Vazquez and Walter, 2003; Xu et al., 2013), and conditional random fields (Ghamrawi and McCalum, 2005).

(DEEP) ARTIFICIAL NEURAL NETWORKS Neural networks handles multi-output tasks by having one node on the output layer per output variable. The network minimizes a global error function defined over all the outputs (Ciarelli et al., 2009; Nam et al., 2014; Specht, 1991; Zhang, 2009; Zhang and Zhou, 2006). The output correlation are taken into account by sharing the input and the hidden layers between all the outputs.

NEAREST NEIGHBORS The k-nearest neighbors algorithm predicts an unseen sample x by aggregating the output value of the k near-

est neighbors of x . This algorithm is adapted to multi-output tasks by sharing the nearest neighbors search among all outputs. If we just share the search, this is called binary relevance of k -nearest neighbors in classification and single target of k -nearest neighbors in regression. Multi-output extensions of the k -nearest neighbors modifies how the output values of the nearest neighbors are aggregated for the predictions for instance it can utilize the maximum a posteriori principle (Cheng and Hüllermeier, 2009; Younes et al., 2011; Zhang and Zhou, 2007) or it can re-interpret the output aggregation as a ranking problem (Brinker and Hüllermeier, 2007; Chiang et al., 2012),

DECISION TREES The decision tree model is a hierarchical structure partitioning the input space and associating a prediction to each partition. The growth of the tree structure is done by maximizing the reduction of an impurity measure computed in the output space. When the tree growth is stopped at a leaf, we associate a prediction to this final partition by aggregating the output values of the training samples. We adapt the decision tree algorithm to multi-output tasks in two steps (Blockeel et al., 2000; Clare and King, 2001; De'Ath, 2002; Noh et al., 2004; Segal, 1992; Vens et al., 2008; Zhang, 1998): (i) multi-output impurity measures are used to grow the structure as the sum over the output space of the entropy or the variance; (ii) the leaf predictions are obtained by computing a constant minimizing a multi-output loss function such as the ℓ_2 -norm loss in regression or the Hamming loss in classification. We discuss in more details how to adapt the decision tree algorithm to multi-output tasks in Section 3.5.

Instead of growing a single decision tree, they are often combined together to improve their generalization performance. Random forest models (Breiman, 2001; Geurts et al., 2006a) averages the predictions of several randomized decision trees and has been studied in the context of multi-output learning (Joly et al., 2014; Kocev et al., 2007, 2013; Madjarov et al., 2012; Segal and Xiao, 2011).

ENSEMBLES Ensemble methods aggregate the predictions of multiple models into a single one so to improve its generalization performance. We discuss how the averaging and boosting approaches have been adapted to multi-output supervised learning tasks.

Averaging ensemble methods have been straightforwardly adapted by averaging the prediction of multi-output models. Instead of averaging scalar predictions, it averages (Joly et al., 2014; Kocev et al., 2007, 2013; Madjarov et al., 2012; Segal and Xiao, 2011) the vector predictions of each model of the ensemble. If the learning algorithm is not inherently multi-output, we could use one the problem transformation techniques as in raket (Tsoumakas and Vlahavas, 2007), which uses the label power set transformation, or ensemble of estimator chain (Read et al., 2011).

Boosting ensembles are grown by sequentially adding weak models minimizing a selected loss, such as the Hamming loss (Schapire and Singer, 2000), the ranking loss (Schapire and Singer, 2000), the ℓ_2 -norm loss (Geurts et al., 2006b) or any differentiable loss function (see Chapter 6).

2.3 EVALUATION OF MODEL PREDICTION PERFORMANCE

For a given supervised learning model f trained on a set of samples $((x^i, y^i) \in (\mathcal{X} \times \mathcal{Y}))_{i=1}^n$, we want a model having good generalization able to predict unseen samples. Otherwise said, the model f should have minimal generalization error over the input-output pair distribution, where the generalization error is defined as:

$$\text{Generalization error} = \mathbb{E}_{P_{x,y}}\{\ell(f(x), y)\} \quad (2.14)$$

for a given loss function $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$.

Evaluating Equation 2.14 is generally unfeasible, except in the rare cases where (i) the input-output distribution $P_{x,y}$ is fully known and (ii) for restricted classes of models. In practice, neither of these conditions are met. However, we still need a principle way to approximate the generalization error.

A first approach to approximate Equation 2.14 is to evaluate the error of the model f on the training samples $\mathcal{L} = ((x^i, y^i) \in (\mathcal{X} \times \mathcal{Y}))_{i=1}^n$ leading to the resubstitution error:

$$\text{Resubstitution error} = \frac{1}{|\mathcal{L}|} \sum_{(x,y) \in \mathcal{L}} \ell(f(x), y) \quad (2.15)$$

A model with a high resubstitution error often has a high generalization error and indeed underfits the data. The linear model shown in Figure 2.7a underfits the data as it is not complex enough to fit the non linear data (here a second degree polynomial). Instead, we can fit a high order polynomial model to have a zero resubstitution error as illustrated in Figure 2.7b. This complex model has poor generalization error as it perfectly fits the noisy samples unable to retrieve the second order parabola. Such overly complex models with zero resubstitution error and non zero generalization error are said to overfit the data. Since a zero resubstitution error does not imply a low generalization error, it is a poor proxy of the generalization error.

Since we assess the quality of the model with the training samples, the resubstitution error is optimistic and biased. Furthermore, it favors overly complex models (as depicted in Figure 2.7). To improve the approximation of the generalization error, we need to use techniques which avoid to use the training samples for performance evaluation. They are either based on sample partitioning methods, such as hold out methods and cross validation techniques, or sample

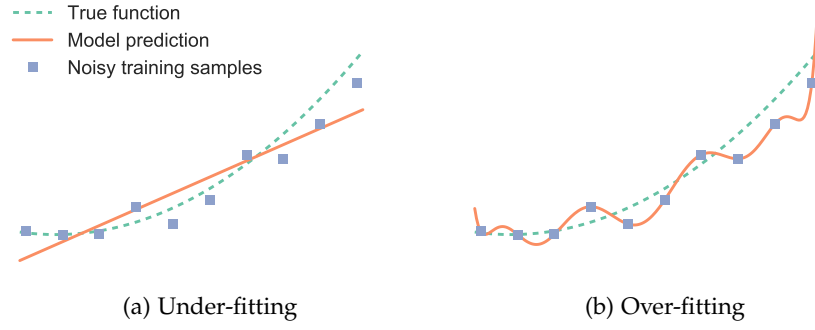


Figure 2.7: The linear model on the left figure underfits the training samples, while the high order polynomial model on the right overfits the training samples.

resampling methods, such as bootstrap estimation methods. Since the amount of available data and time are fixed for both the model training and the model assessment, there is a trade-off between (i) the quality of the error estimate, (ii) the number of samples available to learn a model and (iii) the amount of computing time available for the whole process.

The hold out evaluation method splits the samples into a training set LS , also called learning set, and a testing set TS commonly with a ratio of $2/3 - 1/3$. The hold out error is given by

$$\text{Hold out error} = \frac{1}{|TS|} \sum_{(x,y) \in TS} \ell(f(x), y). \quad (2.16)$$

This methods requires a high number of samples as a large part of the data is devoted to the model assessment impeding the model training. If too few samples are allocated to the testing set, the hold out estimate becomes unreliable as its confidence intervals widen (Kohavi et al., 1995). Since the hold out error is a random number depending on the sample partition, we can improve the error estimation by (i) generating B random partitions $(LS_b, TS_b)_{b=1}^B$ of the available samples, (ii) fitting a model f^b on each learning set LS_b and (iii) averaging the performance of the B models over their respective testing sets TS_b :

$$\text{Random subsampling error} = \frac{1}{B} \sum_{b=1}^B \frac{1}{|TS_b|} \sum_{(x,y) \in TS_b} \ell(f^b(x), y). \quad (2.17)$$

To improve the data usage efficiency, we can resort to cross-validation methods, also called rotation estimation, which split the samples into k folds $\{TS_1, \dots, TS_k\}$ approximately of the same size. Cross validation methods average the performance of k models

$(f^l)_{l=1}^k$ each tested on one of the k folds and trained using the $k - 1$ remaining folds:

$$\text{CV error} = \frac{1}{k} \sum_{l=1}^k \frac{1}{|TS_l|} \sum_{(x,y) \in TS_l} \ell(f^l(x), y). \quad (2.18)$$

The number of folds k is usually 5 or 10. If k is equal to the number of samples ($k = n$), it is called leave-one-out cross validation.

Given that the folds do not overlap for cross validation methods, we are tempted to assess the performance over the pooled cross validation estimates with a given metric obtained by concatenating the predictions made by each model f^l over each of the k -folds

$$\text{Pooled CV error} = \text{metric}((f^1(x), y)_{(x,y) \in TS_1} \frown \dots \frown (f^k(x), y)_{(x,y) \in TS_k}), \quad (2.19)$$

where \frown is the concatenation operator. There is no difference for sample-wise losses such as the square loss. However, this is not the case for metrics comparing a whole set of predictions to their ground truth. Depending on the metrics, it has been showed that pooling may or may not bias the error estimation (Airola et al., 2011; Forman and Scholz, 2010; Parker et al., 2007).

We can improve the quality of the estimate by repeating the cross validation procedures over B different k -fold partitions, averaging the performance of the models $f^{b,l}$ over each associated testing set $TS_{b,l}$:

$$\text{Repeated CV error} = \frac{1}{Bk} \sum_{b=1}^B \sum_{l=1}^k \frac{1}{|TS_{b,l}|} \sum_{(x,y) \in TS_{b,l}} \ell(f^{b,l}(x), y). \quad (2.20)$$

If all combinations are tested exhaustively as in the leave-one-out case, it is called complete cross validation. Since it is often too expensive (Kohavi et al., 1995), we can instead draw several sets of folds at random.

The bootstrap method (Efron, 1983) draws B bootstrap datasets $\{B_1, \dots, B_B\}$ by sampling with replacement n samples from the original dataset of size n . Each samples has a probability of $1 - (1 - \frac{1}{n})^n$ to be selected in a bootstrap which is approximately 0.632 for large n . A first approach to estimate the error is to train a model f^b on each bootstrap dataset and use the original dataset as a testing set:

$$\text{Bootstrap error} = \frac{1}{nB} \sum_{b=1}^B \sum_{(x,y) \in B_b} \ell(f^b(x), y). \quad (2.21)$$

This leads to over optimistic results, given the overlap between the training and the test data.

A better approach (discussed in Chapter 7.11 of (Hastie et al., 2009)) is to imitate cross validation methods by fitting on each bootstrap

dataset a model f^b and using the unused samples as a testing set. This approach is called bootstrap leave-one-out:

$$\text{LOO Bootstrap error} = \frac{1}{n} \sum_{i=1}^n \frac{1}{|C^{-i}|} \sum_{b \in C^{-i}} \ell(f^b(x^i), y^i), \quad (2.22)$$

where C^{-i} gives the bootstrap indices where the sample i was not drawn. It is similar to a 2-fold repeated cross validation or random subsampling error with a ratio of $2/3 - 1/3$ for the training and testing set. The estimation is thus biased as it uses approximately $0.632n$ training samples instead of n . We can alleviate this bias due to the sampling procedure through the “0.632” estimator which averages the training error and LOO Bootstrap error:

$$\begin{aligned} \text{0.632 estimator} &= 0.632 \times \text{LOO Bootstrap error} \\ &+ 0.368 \times \text{Resubstitution error.} \end{aligned} \quad (2.23)$$

Note that with very low sample size, it has been shown (Braga-Neto and Dougherty, 2004) that the bootstrap approach yields better error estimate than the cross validation approach.

Until now, we have assumed that the samples are independent and identically distributed. Whenever this is no longer true, such as with time series of measurements, we have to modify the assessment procedure to avoid biasing the error estimation. For instance, the hold out estimate would train the model on the oldest samples and test the model on the more recent samples. Similarly in the medical context if we have several samples for one patient, we should keep these samples together in the training set or in the testing set.

Partition-based methods (hold out, cross validation) break the assumption in classification that the samples from the training set are independent from the samples in the testing set as they are drawn without replacement from a pool of samples. The representation of each class in the testing set is thus not guaranteed to be the same as in the training set. It is advised (Kohavi et al., 1995) to perform stratified splits by keeping the same proportion of classes in each set.

2.4 CRITERIA TO ASSESS MODEL PERFORMANCE

Assessing the performance of a model requires evaluation metrics which will compare the ground truth to a prediction, a score or a probability estimate. The selection of an appropriate scoring or error measure is essential and is dependent of the supervised learning task and the goal behind the modeling.

A first approach to assess a model is to define a goal for the model and to quantify its realization. For instance, a company wants to maximize its benefits and consider that the revenue must exceed the data analysis cost of gathering samples, fitting a model and exploiting its

predictions. Unfortunately, this model optimization criterion is hardly expressible into economical terms. We could instead consider the effectiveness of the model such as the click-through-rate, used by on-line advertising companies, which counts the number of clicks on a link to the number of opportunities that users have to click on this link. However, it is hard to formulate a model optimizing directly this score and it requires to put the model into a production setting (or at least simulate its behavior). Other optimization criteria exist that are more amenable to mathematical analysis and numerical computation such as the square loss or the logistic loss. Knowing the properties of such criteria is necessary to make a proper choice.

We present binary classification metrics in Section 2.4.1. Then, we show how to extend these metrics to multi-class classification tasks in Section 2.4.2 and to multi-label classification tasks in Section 2.4.3. We introduce metrics for regression tasks and multi-output regression tasks in Section 2.4.4.

More details or alternative descriptions of these metrics can be found in the following references (Ferri et al., 2009; Hossin and Sulaiman, 2015; Sokolova and Lapalme, 2009). Note that I made significant contributions to the implementations and the documentations of these metrics in the scikit-learn library (Buitinck et al., 2013; Pedregosa et al., 2011).

2.4.1 Metrics for binary classification

Given a set of n ground truth values $(y^i \in \{0, 1\})_{i=1}^n$ and their associated model predictions $(\hat{y}^i \in \{0, 1\})_{i=1}^n$, we can distinguish in binary classification four categories of predictions (as shown in Table 2.4). We denote by true positives (TP) and true negatives (TN) the predictions where the model accurately predicts the target respectively as true or false:

$$\text{TP} = \sum_{i=1}^n 1(y^i = 1; \hat{y}^i = 1), \quad (2.24)$$

$$\text{TN} = \sum_{i=1}^n 1(y^i = 0; \hat{y}^i = 0). \quad (2.25)$$

Whenever the model wrongly predicts the samples, we call false positives (FP) samples predicted as true while their labels are false and false negatives (FN) samples predicted as false while their labels are true:

$$\text{FN} = \sum_{i=1}^n 1(y^i = 1; \hat{y}^i = 0), \quad (2.26)$$

$$\text{FP} = \sum_{i=1}^n 1(y^i = 0; \hat{y}^i = 1). \quad (2.27)$$

Table 2.4: For a binary classification task, the prediction of a model is divided into four categories leading to a confusion matrix.

	<i>Truly positive</i>	<i>Truly negative</i>
<i>Predicted positive</i>	True positive	False positive
<i>Predicted negative</i>	False negative	True negatives

Together, the true positive, true negatives, false negatives and false positives form the so called confusion or contingency matrix shown in Table 2.4.

Two common metrics to assess classification performance are the error rate, the average of the 0 – 1 loss, and its complement the accuracy:

$$\text{Error rate} = \frac{1}{n} \sum_{i=1}^n 1(y^i \neq \hat{y}^i), \quad (2.28)$$

$$\text{Accuracy} = 1 - \text{Error rate} = \frac{1}{n} \sum_{i=1}^n 1(y^i = \hat{y}^i). \quad (2.29)$$

Both metrics can be expressed in term of the confusion matrix:

$$\text{Error rate} = \frac{\text{FP} + \text{FN}}{n}, \quad (2.30)$$

$$\text{Accuracy} = \frac{\text{TN} + \text{TP}}{n}. \quad (2.31)$$

The error rate does not distinguish the false negatives from the false positives. Similarly, the accuracy does not differentiate true positives from true negatives. Thus, two classifiers may have exactly the same accuracy or error rate, while leading to a totally different outcome by increasing either the number of misses (false negatives) or the number of false alarms (false positives). Furthermore, the error rate and the accuracy can be overly optimistic whenever there is a high class imbalance. A classification task with 99.99% of samples in one of the classes would easily lead to an accuracy of 99.99% (and an error rate of 0.01%) by always predicting the most common class. The choice of an appropriate metric thus depends on the properties of the classification task, such as the class imbalance.

To differentiate false positives from false negative, we can assess separately the proportion of correctly classified positive and negative samples. This leads to the *true positive rate* (resp. *true negative rate*) which computes the proportion of correctly classified positive (resp. negative) samples:

$$\text{True positive rate} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (2.32)$$

$$\text{True negative rate} = \frac{\text{TN}}{\text{TN} + \text{FP}} \quad (2.33)$$

The complement of the true positive rate (resp. true negative rate) is the false negative rate (resp. false positive rate):

$$\text{False negative rate} = 1 - \text{True positive rate} = \frac{\text{FN}}{\text{TP} + \text{FN}}, \quad (2.34)$$

$$\text{False positive rate} = 1 - \text{True negative rate} = \frac{\text{FP}}{\text{TN} + \text{FP}}. \quad (2.35)$$

The true positive rate is also called *sensitivity* and tests the ability of the classifier to correctly classify all positive samples as true. A test with 100% sensitivity implies that all positive samples are correctly classified. However, this does not imply that all samples are correctly classified. A classifier predicting all samples as true leads to 100% sensitivity and totally neglects false positives. We have to look to the true negative rate, also called *specificity*, which tests the ability of the classifier to correctly classify all negative samples as negative. A perfect classifier should thus have a high sensitivity and a high specificity. In the medical domain, the sensitivity and the specificity are often used to characterize and to choose the behavior of diagnosis tests such as pregnancy tests.

The average of the specificity and sensitivity is called the *balanced accuracy*:

$$\text{Balanced accuracy} = \frac{\text{True positive rate} + \text{True negative rate}}{2} \quad (2.36)$$

$$= \frac{\text{specificity} + \text{sensitivity}}{2} \quad (2.37)$$

$$= \frac{1}{2} \frac{\text{TP}}{\text{TP} + \text{FN}} + \frac{1}{2} \frac{\text{TN}}{\text{TN} + \text{FP}}. \quad (2.38)$$

In the information retrieval context, a user sets a query to an information system, e.g. a web search engine, to detect which documents are relevant among a collection of such documents. In such systems, the collection of documents is often extremely large with only a few relevant documents to a given query. Due to the small proportion of relevant documents, we want to maximize the *precision*, the fraction of correctly predicted documents among the predicted documents. Binary classification tasks with a high class imbalance can be viewed as an information retrieval problems. In the context of binary classification tasks, the precision is expressed as

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}, \quad (2.39)$$

To have a perfect precision, one could predict all documents or samples as negative (as irrelevant documents). In parallel, we want also to maximize the recall, the proportion of correctly predicted true samples among the true samples. The recall is a synonym for true positive rate and sensitivity.

The precision and recall are often combined into a single number by computing the F_1 score, the harmonic mean of the precision and recall,

$$F_1 = \frac{2}{\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}}. \quad (2.40)$$

Some classifiers associate a score or a probability $\hat{f}(x)$ to a sample instead of a class label. We can threshold these continuous predictions by a constant τ to compute the number of true positives, false positives, false negatives and true negatives:

$$\text{TP}(\tau) = \sum_{i=1}^n 1(y^i = 1; f(x^i) \geq \tau), \quad (2.41)$$

$$\text{TN}(\tau) = \sum_{i=1}^n 1(y^i = 0; f(x^i) \leq \tau), \quad (2.42)$$

$$\text{FN}(\tau) = \sum_{i=1}^n 1(y^i = 0; f(x^i) \geq \tau), \quad (2.43)$$

$$\text{FP}(\tau) = \sum_{i=1}^n 1(y^i = 1; f(x^i) \leq \tau). \quad (2.44)$$

By varying τ , we can first derive performance curves to analyze the prediction performance of those more models and then select a classifier performance point with pre-determined classification performance.

The receiver operating characteristic (ROC) curve (Fawcett, 2006) plots the true positive rate as a function of the false positive rate by varying the threshold τ as shown in Figure 2.8a. The receiver, the model user, can indeed choose any point on the curve to operate at a given model specificity / sensitivity tradeoff. A random estimator has its performance on the line $((0, 0), (1, 1))$, while a perfect classifier has the points $(0, 1)$ with 0% of false positive rate and 100% of true positive rate on its curve. Any curve below the random line can be reversed symmetrically to the line $((0, 0), (1, 1))$ by flipping the classifier prediction. The ROC curve is often used in the clinical domain (Metz, 1978) and coupled to a cost analysis to determine the proper threshold τ . The area under the ROC curve can be interpreted as (Hanley and McNeil, 1982) the probability to rank with a higher score one true sample than one false sample chosen at random.

The precision-recall (PR) curve is the precision as a function of the recall as shown in Figure 2.8b. The ROC curve and the PR curves are linked as there is a one to one mapping between points in the ROC space and in the precision-recall space (Davis and Goadrich, 2006). However conversely to the ROC curve, the precision recall curve is sensitive to the class imbalance between the positive and negative classes. Since both the precision and recall do not take into account

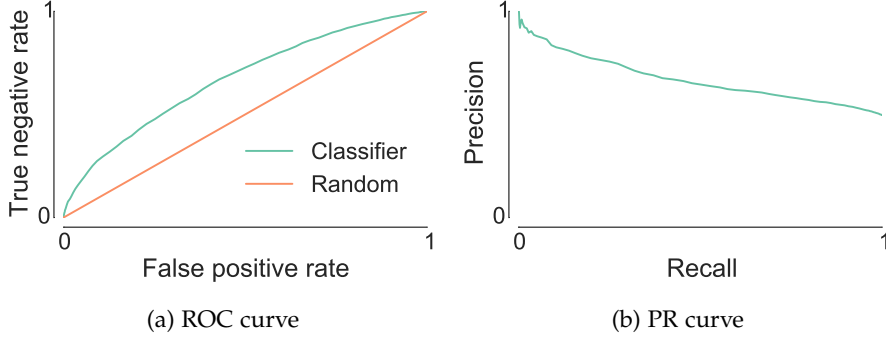


Figure 2.8: A receiver operating characteristic curve and a precision-recall curve of a classifier and a random model.

the amount of true negatives, the precision-recall curve (compared to the ROC curve) focuses on how well the estimator is able to classify correctly the positive class.

2.4.2 Metrics for multi-class classification

From binary classification to multi-class classification, the output value is no more restricted to two classes and can go up to k -classes. Given the ground truths $(y^i \in \{1, \dots, k\})_{i=1}^n$ and the associated model predictions $(\hat{y}^i \in \{1, \dots, k\})_{i=1}^n$, we can now divide the model predictions into k^2 categories leading to a $k \times k$ confusion matrix:

$$c_{l_1, l_2} = \sum_{i=1}^n 1(y^i = l_1; \hat{y}^i = l_2) \quad \forall l_1, l_2 \in \{1, \dots, k\}. \quad (2.45)$$

Metrics such as the accuracy, the error rate or the log loss (see Table 2.3) naturally extend to multi-class classification tasks. To extend other binary classification metrics (such as those developed in Section 2.4.1), we need to break the $k \times k$ confusion matrix into a set of 2×2 confusion matrices.

A first approach is to consider that each class l is in turns the positive class while the remaining labels form together the negative class. We thus have k confusion matrices whose true positives TP_l , true negatives TN_l , false negatives FN_l and false positives FP_l for the class $l \in \{1, \dots, k\}$ are

$$TP_l = c_{l, l}, \quad (2.46)$$

$$TN_l = \sum_{j=1, j \neq l}^k c_{j, j}, \quad (2.47)$$

$$FN_l = \sum_{j=1, j \neq l}^k c_{j, l}, \quad (2.48)$$

$$FP_l = \sum_{j=1, j \neq l}^k c_{k,j}. \quad (2.49)$$

By averaging a metric M computed on each derived confusion matrix, we have the so called macro-averaged (Sokolova and Lapalme, 2009) of the corresponding binary classification metric

$$\text{macro} - M = \frac{1}{k} \sum_{l=1}^k M(TP_l, TN_l, FN_l, FP_l).$$

Note that the balanced accuracy in binary classification is thus equal to the macro-specificity or macro-sensitivity in multi-class classification.

Another useful averaging is the micro-averaging (Sokolova and Lapalme, 2009). It uses as true positives TP_μ and true negatives TN_μ the sum of the diagonal elements of the confusion matrix and as false negatives FN_μ (resp. false positives FP_μ) the sum of the lower (resp. upper) triangular part of the confusion matrix:

$$TP_\mu = \sum_{l=1}^k c_{l,l}, \quad (2.50)$$

$$TN_\mu = \sum_{l=1}^k c_{l,l}, \quad (2.51)$$

$$FN_\mu = \sum_{l=1}^k \sum_{j=1; j < l}^k c_{l,j}, \quad (2.52)$$

$$FP_\mu = \sum_{l=1}^k \sum_{j=1; j > l}^k c_{l,j}. \quad (2.53)$$

Each averaging has its own properties: the macro-averaging considers that each class has the same importance and the micro-averaging reduces the importance given to the minority classes.

2.4.3 Metrics for multi-label classification and ranking

From binary to multi-label classification, the ground truths $(y^i \in \{0, 1\}^d)_{i=1}^n$ and the model predictions $(\hat{y}^i \in \{0, 1\}^d)_{i=1}^n$ are no longer scalars, but vectors of size d or label sets. Both representations are interchangeable. Usually, the number of labels associated to a sample is small compared to the total number of labels.

The accuracy (Ghamrawi and McCallum, 2005), also called subset accuracy, has a direct extension in multi-label classification

$$\text{Accuracy} = \frac{1}{n} \sum_{i=1}^n 1(y^i = \hat{y}^i), \quad (2.54)$$

and requires for each prediction that the predicted label set matches exactly the ground truth. This is an overly pessimistic metric, especially for high dimensional label space, as it penalizes any single mistake made for one sample. The complement of the subset accuracy is called the subset 0-1 loss ().

In information theory, the Hamming distance compares the number of differences between two coded messages. The Hamming error metric (Schapire and Singer, 1999) averages the Hamming distance between the ground truth and the model prediction over the samples

$$\text{Hamming error} = \frac{1}{n} \frac{1}{d} \sum_{i=1}^n \sum_{j=1}^d 1(y_j^i \neq \hat{y}_j^i). \quad (2.55)$$

By contrast to the subset accuracy, the Hamming error is an optimistic metric when the label space is sparse. For a sufficiently large number of samples and a label density¹ $\epsilon \rightarrow 0$, a (useless) model predicting always the presence of a label if its frequency of apparition is higher than 0.5 in the training set will roughly have a Hamming error of ϵ . In some situations, the label density ϵ is so small that (more useful) models have hardly an Hamming error lower than ϵ .

Both the Hamming error and the subset accuracy ignore the sparsity of the label space leading to either overly optimistic or pessimistic error. Multi-label metrics should be aware of the label space sparsity.

In statistics, the Jaccard index J or Jaccard similarity coefficient computes the similarity between two sets. Given two sets A and B , the Jaccard index is defined as

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} \text{ with } J(\emptyset, \emptyset) = 1. \quad (2.56)$$

With label sets encoded as boolean vectors $x, y \in \{0, 1\}^d$, the Jaccard index becomes

$$J(x, y) = \frac{x^T y}{1_d^T x + 1_d^T y - x^T y}, \quad (2.57)$$

where 1_d is a vector of ones of size d . The Jaccard similarity score (Godbole and Sarawagi, 2004), also sometimes called accuracy, averages over the samples the Jaccard index between the ground truths and the model predictions:

$$\text{Jaccard similarity score} = \frac{1}{n} \sum_{i=1}^n J(y^i, \hat{y}^i). \quad (2.58)$$

By contrast to the Hamming loss, the Jaccard similarity score puts more emphasis on the labels in the ground truth and the ones predicted by the models. Moreover, it totally ignores all the negative

¹ The label density is the average number of labels per samples on the ground truth divided by the size of the label space.

labels. The Jaccard similarity score can be viewed as “local” measure of similarity and the Hamming loss a “global” measure of distance.

A fitted model f applied to an input vector x can go beyond label prediction and associate to each label j a score or a probability estimate $f(x)_j$. When the density of the label space ϵ is small and the size of the label space d is very high, it is often hard to correctly predict all labels. Instead, the classifier can rank or score all the labels. We developed here metrics for such classifiers with different possible goals, e.g. to predict correctly the label with the highest score $f(x)_j$.

Note that in the following, we use indifferently the notation $|\cdot|$ to express the cardinality of a set or the ℓ_1 -norm of a vector.

If only the top scored label has to be correctly predicted, we are minimizing the one error (Schapire and Singer, 1999) which computes the fraction of labels with the highest score or probability that are incorrectly predicted:

$$\text{One error} = \frac{1}{n} \sum_{i=1}^n \mathbb{I} \left(y_j^i \neq 1 : j = \arg \max_{j \in \{1, \dots, d\}} f(x^i)_j \right). \quad (2.59)$$

If we want to discover all the true labels at the expense of some false labels, the coverage error (Schapire and Singer, 2000) is the metrics to minimize. It counts the average number of labels with the highest scores or probabilities to consider to cover all true labels:

$$\text{Coverage error} = \frac{1}{n} \sum_{i=1}^n \max_{j: y_j^i \neq 1} |\{k : f(x^i)_k \geq f(x^i)_j\}|. \quad (2.60)$$

For a label density of ϵ , the best coverage error is thus ϵd and the worst is d .

If we want to ensure that pairwise label comparisons are correctly made by the classifier, we will minimize the (pairwise) ranking loss metrics (Schapire and Singer, 1999). It counts for each sample the number of wrongly made pairwise comparisons divided by the number of true labels and false labels

$$\text{Ranking loss} = \frac{1}{n} \sum_{i=1}^n \frac{1}{|y^i|} \frac{1}{d - |y^i|} |\{(k, l) : f(x^i)_k < f(x^i)_l, y_k^i = 1, y_l^i = 0\}| \quad (2.61)$$

The ranking loss is between 0 and 1. A ranking loss of 0 (resp. 1) indicates that all pairwise comparisons are correct (resp. wrong).

If we want that the classifier gives on average a higher score to true labels, we will use the label ranking average precision metric (Schapire and Singer, 2000) to assess the accuracy of the models. For each samples y^i , it averages over each true labels j the ratio between (i) the number of true label (i.e. $y^i = 1$) with higher scores or probabilities than the label j to (ii) the number of labels (y^i) with higher score $f(x^i)$ than the label j . Mathematically, we average the

LRAP of all pairs of ground truth y^i and its associated prediction $f(x^i)$:

$$\text{LRAP}(\hat{f}) = \frac{1}{|\text{TS}|} \sum_{i=1}^n \frac{1}{|y^i|} \sum_{j \in \{k: y_k^i = 1\}} \frac{|\mathcal{L}_j^i(y^i)|}{|\mathcal{L}_j^i(1_d)|}, \quad (2.62)$$

where

$$\mathcal{L}_j^i(q) = \{k : q_k = 1 \text{ and } \hat{f}(x^i)_k \geq \hat{f}(x^i)_j\}.$$

The best possible average precision is thus 1. Note that the LRAP score is equal to fraction of positive labels if all labels are predicted with the same score or all negative labels have a score higher than the positive one.

Let us illustrate the computation of the previous metrics with a numerical example. We compare the ground truth \mathbf{y} of $n = 2$ samples in a label of size $d = 5$ to the probability score $\mathbf{f}(\mathbf{x})$ given by the classifier:

$$\mathbf{y} = \begin{bmatrix} 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \end{bmatrix},$$

$$\mathbf{f}(\mathbf{x}) = \begin{bmatrix} 0.75 & 0.6 & 0.1 & 0.8 & 0.15 \\ 0.25 & 0.8 & 0.1 & 0.15 & 0.3 \end{bmatrix}.$$

Thresholding $\mathbf{f}(\mathbf{x})$ at 0.5 yields the prediction $\hat{\mathbf{y}}$ of the classifier:

$$\hat{\mathbf{y}} = \mathbf{f}(\mathbf{x}) \leq 0.5 \begin{bmatrix} 1 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 \end{bmatrix}$$

Here, you will find the detailed computation of all previous metrics:

$$\begin{aligned} \text{Accuracy} &= 0 + 0 = 0, \\ \text{Hamming loss} &= \frac{1}{2} \frac{1}{5} = 0.5, \\ \text{Jaccard similarity score} &= \frac{1}{2} \left(\frac{1}{4} + \frac{0}{2} \right) = 0.125, \\ \text{Top error} &= \frac{1}{2} (1 + 1) = 2, \\ \text{Coverage error} &= \frac{5 + 3}{2} = 4 \\ \text{Ranking loss} &= \frac{1}{2} \left(\frac{1}{2} \frac{1}{3} (1 + 3) + \frac{1}{4} \frac{1}{1} 2 \right) \approx 0.583, \\ \text{LRAP} &= \frac{1}{2} \left(\frac{1}{2} \left(\frac{1}{2} + \frac{2}{5} \right) + \frac{1}{1} \frac{1}{3} \right) \approx 0.392. \end{aligned}$$

While the previous metrics are suited to assess multi-label classification models, we can complement these metrics with those developed for binary classification tasks, e.g. specificity, precision, ROC

AUC, ... (see Section 2.4.1). They are well understood in their respective domains and have attractive properties such as a good handling of class imbalance. We extend those metrics in three steps: (i) we break the ground truth and the model prediction vectors into its elements, (ii) we concatenate the elements into groups such as all predictions associated to a given sample or all samples associated to a given label and (iii) we average the binary classification metrics over each group. We will focus here on three averaging methods: macro-averaging, micro-averaging and sample-averaging. Each averaging method stems from a vision and different sets of assumptions.

If we view the multi-label classification task as a set of independent binary classification tasks, we compute the metrics M over each output separately and average the performance over all d labels leading the *macro-averaging* version (Yang, 1999) of the metrics M :

$$\text{macro-}M((y^i)_{i=1}^n, (\hat{y}^i)_{i=1}^n) = \frac{1}{d} \sum_{j=1}^d M((y_j^i)_{i=1}^n, (\hat{y}_j^i)_{i=1}^n). \quad (2.63)$$

If instead we view each sample as the result of a query (like in a search engine), we want to evaluate the quality of each query (or sample) separately. Under this perspective, the *sample-averaging* approach (Godbole and Sarawagi, 2004) computes and averages the metric M over each sample separately:

$$\text{sample-}M((y^i)_{i=1}^n, (\hat{y}^i)_{i=1}^n) = \frac{1}{n} \sum_{i=1}^n M(y^i, \hat{y}^i). \quad (2.64)$$

The *micro-averaging* approach (Yang, 1999) views all label-sample pairs as forming an unique binary classification task. It compute the metric M as if all label predictions were independent:

$$M\text{-micro}((y^i)_{i=1}^n, (\hat{y}^i)_{i=1}^n) = M((y_j^i)_{i,j=(1,\dots,n)}, (\hat{y}_k^i)_{i,j=(1,\dots,n)}). \quad (2.65)$$

2.4.4 Regression metrics

Given a set of n ground truths $(y^i \in \mathbb{R})_{i=1}^n$ and their associated model predictions $(\hat{y}^i \in \mathbb{R})_{i=1}^n$, regression tasks are often assessed using the mean square error (MSE), the average of the square loss, expressed by

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y^i - \hat{y}^i)^2. \quad (2.66)$$

From the mean square error, we can derive the r^2 score, also called the coefficient of determination. It is the fraction of variance explained by the model:

$$r^2 = 1 - \frac{\text{MSE}}{\text{Output variance}} \quad (2.67)$$

$$= 1 - \frac{\sum_{i=1}^n (y^i - \hat{y}^i)^2}{\sum_{i=1}^n (y^i - \frac{1}{n} \sum_{l=1}^n y^l)^2} \quad (2.68)$$

The r^2 score is normally between 0 and 1. A r^2 score of zero indicates that the model is no better than a constant, while a r^2 of one indicates that the model perfectly explains the output given the inputs. A negative r^2 score might occur and it indicates that the model is worse than a constant model.

Square-based metrics are highly sensitive to the presence of outliers with abnormally high prediction errors. The mean absolute error (MAE), the average of the absolute loss, is often suggested as a robust replacement of the MSE:

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y^i - \hat{y}^i|. \quad (2.69)$$

These single output metrics naturally extend to multi-output regression tasks. The multi-output mean squared error and mean absolute error for an output space size d is given by

$$\text{MSE} = \frac{1}{n} \frac{1}{d} \sum_{i=1}^n \|y^i - \hat{y}^i\|_{\ell_2}^2, \quad (2.70)$$

$$\text{MAE} = \frac{1}{n} \frac{1}{d} \sum_{i=1}^n \|y^i - \hat{y}^i\|_{\ell_1}. \quad (2.71)$$

These measures average the metrics over all outputs assuming they are independent.

Similarly, averaging the r^2 score over each output leads to the macro- r^2 score:

$$\text{macro-}r^2 = 1 - \frac{1}{d} \sum_{j=1}^d \frac{\sum_{i=1}^n (y_j^i - \hat{y}_j^i)^2}{\sum_{i=1}^n (y_j^i - \frac{1}{n} \sum_{l=1}^n y_j^l)^2}. \quad (2.72)$$

An alternative extension of the r^2 score is to consider the total fraction of the output variance, or more strictly the sum of the variance over each output, explained by the model

$$\text{variance-}r^2 = 1 - \frac{\text{MSE}}{\text{Total output variance}} \quad (2.73)$$

$$= 1 - \frac{\sum_{i=1}^n \|y^i - \hat{y}^i\|_{\ell_2}^2}{\sum_{i=1}^n \|y^i - \frac{1}{n} \sum_{l=1}^n y^l\|_{\ell_2}^2}, \quad (2.74)$$

which is equal to 1 minus the fraction of explained variance ([Bakker and Heskes, 2003](#)).

The variance- r^2 is a variance weighted average of the r^2 score. We can reformulate the variance- r^2 as:

$$\text{variance-}r^2 = 1 - \frac{\sum_{i=1}^n \sum_{j=1}^d (y_j^i - \hat{y}_j^i)^2}{\text{Total output variance}} \quad (2.75)$$

$$= 1 - \sum_{j=1}^d w_j \frac{\sum_{i=1}^n (y_j^i - \hat{y}_j^i)}{\sum_{i=1}^n (y_j^i - \frac{1}{n} \sum_{l=1}^n y_l^i)^2} \quad (2.76)$$

with $w_j = \frac{\sum_{i=1}^n (y_j^i - \frac{1}{n} \sum_{l=1}^n y_l^i)^2}{\text{Total output variance}}$. By contrast, the macro-r² score would have uniform weights $w_j = \frac{1}{d} \forall j$ in Equation 2.76.

2.5 HYPER-PARAMETER OPTIMIZATION

Supervised learning algorithms can be viewed as a function $A : (\mathcal{X} \times \mathcal{Y})_{i=1}^n \times \mathcal{A} \rightarrow \mathcal{H}$ taking as input a learning set $\mathcal{L} = ((x^i, y^i) \in (\mathcal{X} \times \mathcal{Y}))_{i=1}^n$ and a set of hyper-parameters $\alpha \in \mathcal{A}$ and outputting a function f in a hypothesis space $\mathcal{H} \subset \mathcal{Y}^{\mathcal{X}}$. The hypothesis space \mathcal{A} is defined through one or several hyper-parameter variables that can be either discrete, like the number of neighbors for a nearest neighbors model, or continuous, like the multiplying constant of a penalty loss in penalized linear models.

We need hyper-parameter tuning methods to find the best hyper-parameter set $\alpha^* \in \mathcal{A}$ that minimizes the expectation of some loss function $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$ over the joint distribution of input / output pairs $P_{\mathcal{X}, \mathcal{Y}}$:

$$\alpha^* = \arg \min_{\alpha \in \mathcal{A}} E_{P_{\mathcal{X}, \mathcal{Y}}} \{\ell(A(\mathcal{L}, \alpha)(x), y)\}. \quad (2.77)$$

Directly optimizing Equation 2.77 is in general not possible as it consists in minimizing the generalization error over unseen samples. Thus, we resort to validation techniques to split the samples into one (or more) validation set(s) S^{valid} to estimate the generalization error (see Section 2.3) and to select the best set of hyper-parameter α^* :

$$\alpha^* \approx \arg \min_{\alpha \in \mathcal{A}} \sum_{(x, y) \in S^{\text{valid}}} \{\ell(A(\mathcal{L}, \alpha)(x), y)\}. \quad (2.78)$$

Note that we can optimize a metric defined over a set of samples instead a loss as the area under the ROC curve.

In its simplest form, the hyper-parameter search assesses all possible hyper-parameter sets $\alpha \in \mathcal{A}$. While it is optimal on the validation set(s), this is impractical as the size of the hyper-parameter space \mathcal{A} is often unbounded. The hyper-parameter space often consists of continuous hyper-parameter variables leading to an infinite number of possible hyper-parameter sets. Whenever the number of hyper-parameter sets is finite ($|\mathcal{A}| < \infty$), we are limited by computational budget constraints. Instead, we resort to evaluate a subset of the hyper-parameter space $\mathcal{A}^- \subset \mathcal{A}$:

$$\alpha^* \approx \arg \min_{\alpha \in \mathcal{A}^-} \sum_{(x, y) \in S^{\text{valid}}} \{\ell(A(\mathcal{L}, \alpha)(x), y)\}. \quad (2.79)$$

The classical approach to design a finite and reduced subspace \mathcal{A}^- is to sample the hyper-parameter space \mathcal{A} through a manually defined grid. A too coarse grid will miss the optimum hyper-parameter set α^* , while a too fine grid will be very costly. In (Hsu et al., 2003), Hsu et al. suggests a two-stage approach: (i) a coarse parameter grid first identifies regions of interest in the hyper-parameter space, and then a finer grid locates the optimum. Nevertheless, we might still miss the optimal hyper-parameter set α^* since the objective function of Equation 2.79 is not necessarily convex nor concave

HOW TO WRONGLY OPTIMIZE AND / OR TO WRONGLY VALIDATE A MODEL?

Given a set of samples S and a supervised learning algorithm A , one wants simultaneously to find the hyper-parameter set α^* and estimate the generalization error of the associated model f .

A wrong approach would be to use directly one of the validation techniques presented in Section 2.3 dividing the sample set S into (multiple) pair(s) of a training set and a test set $(S^{\text{train}}, S^{\text{test}})$. If we select the best hyper-parameter set α^* based on the test set(s) S^{test} , then the approximation of the generalization error on S^{test} is biased: the hyper-parameter set α^* has been selected on the same test set(s). Another approach would be to repeat independently the described process using different partitions of the sample set S to first select the best model and then to estimate the generalization error. However, the generalization error is still biased: we might use the same samples to train, to select or to validate the model.

The correct approach is to use *nested validation techniques*. We first divide the sample set into (multiple) pair(s) of a test set S^{test} and training-validation set $S^{\text{train-valid}}$. Then we again apply a validation technique to split the training-validation set into (multiple) pair(s) of a training set S^{train} and a validation set S^{valid} . The models f with hyper-parameter α are first trained on S^{train} , then we select the best hyper-parameter set α^* on S^{valid} . We finally estimate the generalization error of the overall model training and selection procedure by re-training a model on $S^{\text{train-valid}}$ using the best hyper-parameter set α^* on the testing set S^{test} .

Proper validation is necessary and comes at the expense of the sample efficiency and computing time. Note that nested validation methods are not needed if we want solely either to select the best model or to estimate the generalization error of a given model.

In the grid search approach, we first sample each hyper-parameter variable and then build all possible combinations of hyper-parameter sets. However, some of these hyper-parameter variables have no or small influence on the performances of the models. In these conditions, large hyper-parameter grids are doomed to fail due to the

explosion of hyper-parameter sets. Random search techniques (Solis and Wets, 1981) tackles such optimization problems by (i) defining a distribution over the optimization variables, (ii) drawing random values from this distribution and (iii) selecting the best one out of these values. (Bergstra and Bengio, 2012) have shown that random hyper-parameter search scales better than grid search as the search is not affected by the hyper-parameter variables having few or no influence on the model performance. As an illustration, let us consider a model with one parameter and one without impact on its generalization error. Sampling 9 random hyper-parameter sets would yield more information than making a 3×3 grid as we evaluate 9 different values of the dependent variable in the random search instead of 3 in the grid.

For a continuous loss and a continuous hyper-parameter space, Bayesian hyper-parameter optimization (Bergstra et al., 2011; Hutter et al., 2011; Snoek et al., 2012) goes beyond random search and models the performance of a supervised learning algorithm A with hyper-parameters α . Starting from an initial Gaussian prior over the hyper-parameter space, it refines a posterior distribution of the model error with each new tested sets of hyper-parameters. New hyper-parameter sets are drawn to minimize the overall uncertainty and the model error.

2.6 UNSUPERVISED PROJECTION METHODS

Supervised learning aims at finding the best function f which maps the input space \mathcal{X} to the output space \mathcal{Y} given a set of n samples $((x^i, y^i) \in (\mathcal{X} \times \mathcal{Y}))_{i=1}^n$. However with very high dimensional input space, we need a very high number of samples n to find an accurate function f . This is the so-called curse of dimensionality. Another problem arises if the model f is unable to model the input-output relationship because the model classes \mathcal{H} is too restricted, for instance a linear model will fail to model quadratic data.

Unsupervised projection methods lift the original space \mathcal{X} of size p to another space \mathcal{Z} of size q . If the projection lowers the size of the original space ($q < p$), this is a dimensionality reduction technique. In the context of supervised learning, we hope to break the curse of dimensionality with such projection methods while speeding up the model training. If the projections perform non linear transformations of the input space, it might also improve the model performance. For instance, a linear estimator will be able to fit quadratic data if we enrich the input variables with their quadratic and bilinear forms. Note that projecting the input space to two or three dimensions ($q \in \{2, 3\}$) is an opportunity to get insights on the data through visualization.

We present three popular unsupervised projection methods and discuss their properties: (i) the principal component analysis ap-

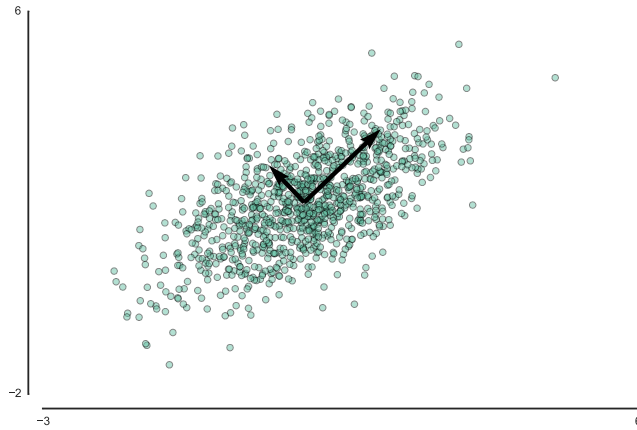


Figure 2.9: The principal components, shown as black arrows, are the orthogonal vectors maximizing the variance of the samples drawn here from a multivariate Gaussian distribution.

proach in Section 2.6.1, which aims to find a subspace maximizing the total variance of the data; (ii) random projection methods in Section 2.6.2, which project the original space onto a lower dimensional space while approximately preserving pairwise euclidean distances, and (iii) kernel functions in Section 2.6.3, which compute pairwise sample similarities lifting the original space to a non linear one.

2.6.1 Principal components analysis

The principal component analysis (PCA) method (Jolliffe, 2002) is a technique to find from a set of samples $(x^i \in \mathcal{X})_{i=1}^n$ an orthogonal linear transformation Z which maximizes the variance along each axis of the transformed space as shown in Figure 2.9.

Principal component analysis reduces the dimensionality of the dataset by keeping a fraction of the principal components vectors which totalize a large amount of the total variance. If we keep only two components, PCA allows to visualize high dimensional datasets as illustrated in Figure 2.10 with digits.

Mathematically, we want to find the first principal component vector u^1 which maximizes the variance along its direction:

$$\begin{aligned} u^1 &= \arg \max_u \sum_{i=1}^n \left\| u^T x^i - u^T \sum_{l=1}^n \frac{x^l}{n} \right\|_{\ell_2}^2 \\ &\text{s.t. } u^T u = 1. \end{aligned} \quad (2.80)$$

Given that the covariance matrix C is given by

$$C = \sum_{i=1}^n \left(x^i - \sum_{l=1}^n \frac{x^l}{n} \right) \left(x^i - \sum_{l=1}^n \frac{x^l}{n} \right)^T,$$

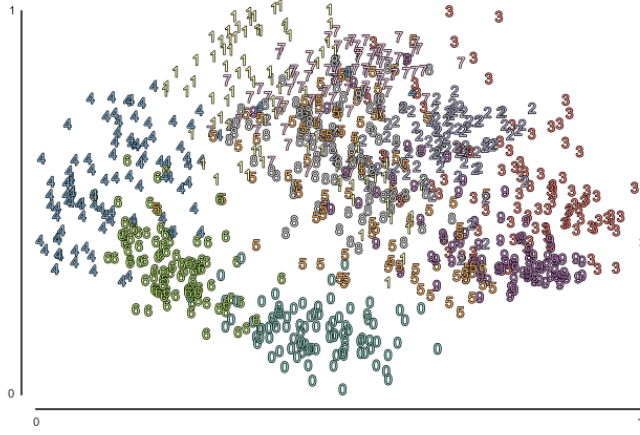


Figure 2.10: We project the digits dataset (Lichman, 2013) from its 8×8 pixel space on the two principal components with the largest variance (29% of the total variance). Digits such as 4 or 0 form well defined clusters on this two dimensional space.

we have

$$\mathbf{u}^1 = \arg \max_{\mathbf{u}} \mathbf{u}^T \mathbf{C} \mathbf{u} + \lambda_1 (1 - \mathbf{u}^T \mathbf{u}), \quad (2.81)$$

where λ_1 is the Lagrange multiplier of the normalization constraint.

By derivating with respect to \mathbf{u} and setting the first derivative to zero, we have that the maximum is indeed an eigen vector of the covariance matrix:

$$\mathbf{C} \mathbf{u}^1 = \lambda_1 \mathbf{u}^1 \quad (2.82)$$

We also note that the variance along \mathbf{u}^1 is given by $\mathbf{u}^{1T} \mathbf{C} \mathbf{u}^1 = \lambda_1$. Thus \mathbf{u}^1 is the eigenvector with the highest eigen value.

The following vector \mathbf{u}^{m+1} maximizing the variance are obtained by imposing that the $m + 1$ -th vector is orthogonal to the m previous one:

$$\begin{aligned} \mathbf{u}^{m+1} &= \arg \max_{\mathbf{u}^{m+1}} \mathbf{u}^{m+1T} \mathbf{C} \mathbf{u}^{m+1} \\ \text{s.t. } &\mathbf{u}^{m+1T} \mathbf{u}^{m+1} = 1, \\ &\mathbf{u}^{m+1T} \mathbf{u}^l = 0 \quad \forall l \in \{1, \dots, m\}, \end{aligned} \quad (2.83)$$

or alternatively in Lagrangian form

$$\arg \max_{\mathbf{u}^{m+1}} \mathbf{u}^{m+1T} \mathbf{C} \mathbf{u}^{m+1} + \lambda_{m+1} (1 - \mathbf{u}^{m+1T} \mathbf{u}^{m+1}) + \sum_{l=1}^m \mu_l \mathbf{u}^{m+1T} \mathbf{u}^l. \quad (2.84)$$

By differentiating with respect to \mathbf{u}^l and multiplying by \mathbf{u}^{m+1} , we have that the Lagrangian constants of the orthogonality constraints

are equal to zero $\mu_l = 0 \forall l \in \{1, \dots, m\}$. And it follows that the $m + 1$ -th principal component is the $m + 1$ -th eigen vector with the $m + 1$ -th largest eigen value λ_{m+1} since

$$Cu^{m+1} = \lambda_{m+1}u^{m+1}, u^{m+1T}Cu^{m+1} = \lambda_{m+1}. \quad (2.85)$$

2.6.2 Random projection

Random projection is a dimensionality reduction method which projects the space onto a smaller random space. For randomly projection, the Johnson-Lindenstrauss lemma gives the conditions of existence such that the distance between pairs of points is approximately preserved.

JOHNSON-LINDENSTRAUSS LEMMA (JOHNSON AND LINDENSTRAUSS, 1984)

Given $\epsilon > 0$ and an integer n , let q be a positive integer such that $q \geq 8\epsilon^{-2} \ln n$. For any sample $(y^i)_{i=1}^n$ of n points in \mathbb{R}^d there exists a matrix $\Phi \in \mathbb{R}^{q \times d}$ such that for all $i, j \in \{1, \dots, n\}$

$$(1 - \epsilon)\|y^i - y^j\|^2 \leq \|\Phi y^i - \Phi y^j\|^2 \leq (1 + \epsilon)\|y^i - y^j\|^2. \quad (2.86)$$

Moreover, when d is sufficiently large, several random matrices satisfy Equation 2.86 with high probability. In particular, we can consider Gaussian matrices whose elements are drawn *i.i.d.* in $\mathcal{N}(0, 1/q)$, as well as (sparse) Rademacher matrices whose elements are drawn in the finite set $\left\{-\sqrt{\frac{s}{q}}, 0, \sqrt{\frac{s}{q}}\right\}$ with probability $\left\{\frac{1}{2s}, 1 - \frac{1}{s}, \frac{1}{2s}\right\}$, where $1/s \in (0, 1]$ controls the sparsity of Φ . If $s = 3$, we will say that those projections are Achlioptas random projections (Achlioptas, 2003). When the size of the original space is p and $s = \sqrt{p}$, then we will say that we have sparse random projection as in (Li et al., 2006). Note a random sub-space (Ho, 1998) is also a random projection scheme (Candes and Plan, 2011): the projection matrix Φ is obtained by sub-sampling with or without replacement the identity matrix. Theoretical results proving 2.86 with high probability for each random projection matrix can be found in the relevant paper.

The choice of the number of random projections q is a trade-off between the quality of the approximation and the size of the resulting embedding as illustrated in Figure 2.11.

2.6.3 Kernel functions

A kernel function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ computes the similarity between pairs of samples (usually in the input space). Machine learning algorithms relying solely on dot products, such as support vector

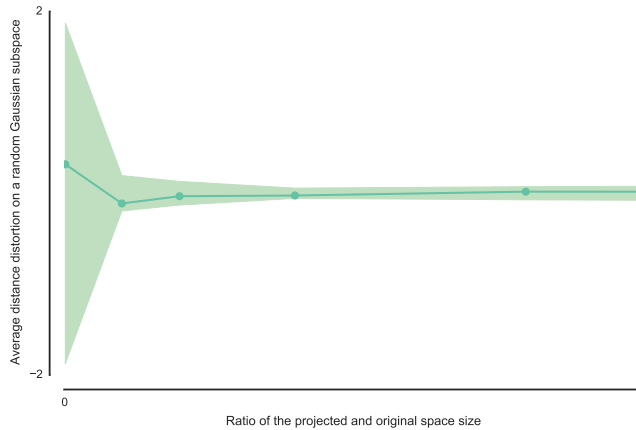


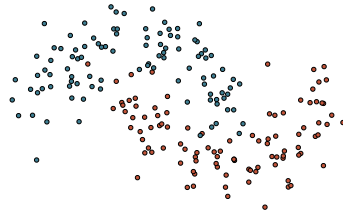
Figure 2.11: Given a set of samples $(x^i)_{i=0}^{2000} \in \mathbb{R}^{500}$ drawn from a Gaussian distribution, a few random projections preserve on average the distance between pairs of points up to a distortion ϵ .

Table 2.5: Some common kernel functions between two vectors x^i and x^j .

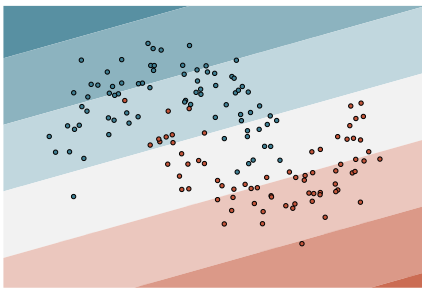
Kernels	
Linear kernel	$k(x^i, x^j) = x^{i\top} x^j$
Polynomial	$k(x^i, x^j) = (x^{i\top} x^j + r)^d$
Gaussian radial basis function	$k(x^i, x^j) = \exp(-\gamma \ x^i - x^j\ _2^2)$
Hyperbolic tangent	$k(x^i, x^j) = \tanh(\kappa x^{i\top} x^j + r)$

machine (Cortes and Vapnik, 1995) or principal components analysis (Jolliffe, 2002), are indeed using the linear kernel $k(x^i, x^j) = x^{i\top} x^j$. We can kernelize these algorithms by replacing their dot product with a kernel presented in Table 2.5. This is the so called kernel trick (Scholkopf and Smola, 2001). Kernel functions define non linear projection schemes lifting the original space to the one defined by the chosen kernel. It has been used in classification (Hsu et al., 2003), in regression (Jaakkola and Haussler, 1999) and in clustering (Schölkopf et al., 1997). Random kernels (Rahimi and Recht, 2007, 2009) can be used to compress the input space.

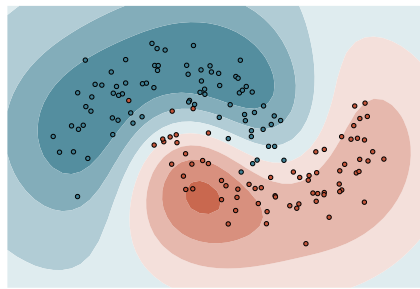
The task shown in Figure 2.12b requires to classify points belonging to one of two interleaved moons. Given the non linearities, we can not linearly separate both classes as illustrated in Figure 2.12b with a (linear) support vector machine. By lifting the linear kernel to the radial basis function kernel, the support vector machine algorithm now separates both classes as shown in Figure 2.12c. Effectively, kernel functions enable machine learning algorithms to handle a wide varieties of structured and unstructured data such as sequences, graphs, texts or vectors through an appropriate choice and design of kernel functions.



(a) Classification task



(b) Linear SVM



(c) Radial basis function SVM

Figure 2.12: The classification task consists in discriminating points belonging to one of the two interleaved moons (the blue or the red dots). Given the non linearities of the data, a linear support vector machine is not able to find a hyperplane separating both classes as shown in Figure 2.12b. By lifting the original input space through the radial basis function kernel, we are now able to separate both classes as illustrated in Figure 2.12c retrieving the interleaved moons.

3

DECISION TREES

OUTLINE

Decision trees are non parametric supervised learning models mapping the input space to the output space. The model is a hierarchical structure made of test and leaf nodes. Starting at the root node, the top of the tree, the test nodes lead the reasoning through the tree structure until reaching a leaf node outputting a prediction. In this chapter, we first describe the decision tree model and show how to predict unseen samples. Then, we present the methods and the techniques to grow and to prune these tree structures. We also introduce how to interpret a decision tree model to gain insights on the studied systems and phenomena.

A decision tree is comparable to a human reasoning organized through a hierarchical set of yes/no questions. As in medical diagnosis, an expert (here the doctor) diagnoses the patient state ("Is the patient healthy or sick?") by screening the patient body and by retrieving important past patient history through a directed set of questions. We can view each step of the reasoning as a branch of a decision tree structure. Each answer leads either to another question refining a set of hypotheses or finally to a conclusion, a prediction.

The binary questions at test nodes can target binary variables, like "Do you smoke?", categorical variables, like "Do you prefer pear or apple to orange or lemon?", or numerical variables, like "Is the outside temperature below 25 degree Celsius (77 degree Fahrenheit)?". Note that we can formulate multi-way questions as a set of binary questions. For instance, the multi-way question "Do you want to eat a pear, a peach or an apple?" is equivalent to ask sequentially "Do you want to eat a pear or one fruit among peach and apple?", then you would also ask "Do you want to eat a peach or an apple?" if you answered "a peach or an apple".

With only numerical input variables, questions that are typically asked are in the form "Is the value of variable x lower than a given value?". The decision tree is then a geometric structure partitioning the input space into a recursive set of p -dimensional (hyper)rectangles (also called p -orthotopes). The root node first divides the whole input space into two half-spaces. Each of those may further be divided into smaller (hyper)rectangles. The partition structure highlights the hierarchical nature of a decision tree. An artistic example of such partitioning in a two dimensional space is the

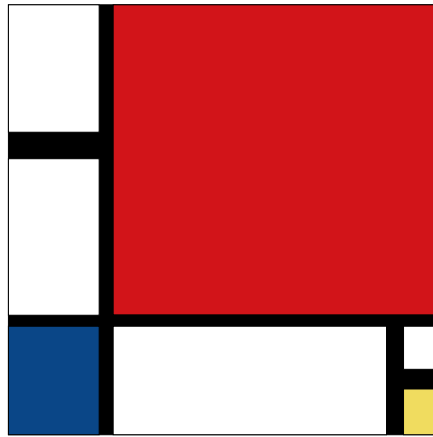


Figure 3.1: Reproduction of “Composition II in Red, Blue, and Yellow” of Piet Mondrian.

“Composition II in Red, Blue, and Yellow” by Piet Mondrian¹ shown in Figure 3.1. Here, Piet Mondrian divides hierarchically the whole painting into colored rectangles through heavy thick black lines. Each black line is conceptually a testing node of a decision tree, while the colored rectangles would be the predictions of leaves node.

Decision trees are popular machine learning model, because of several nice properties:

- The hierarchical nature of the model takes into account non linear effects between the inputs and outputs, as well as conditional dependencies among inputs and outputs.
- Growing a decision tree is computationally fast.
- Decision tree techniques work with heterogeneous data combining binary, categorical or numerical input variables.
- Decision trees are interpretable models giving insights on the relationship between the inputs and the outputs.
- The tree training algorithm can be adapted to handle missing input values (Breiman et al., 1984; Friedman, 1977; Quinlan, 1989).

In Section 3.1, we present the structure of such models and how to exploit these to predict unseen samples. We show in Section 3.2 how to train a decision tree model. In Section 3.3, we describe the techniques used to prune a fully grown decision tree to the right size: too shallow trees tend to under-fit the data as they might lack predicting power, while large trees might overfit the data as they are too complex. In Section 3.4, we show how to interpret a decision tree to gain insights over the input-output relationships: (i) through input

¹ Piet Mondrian (1872-1944) is a painter famous for his grid-based paintings partitioning the tableau through black lines into colored rectangles usually blue, red, yellow and white

variable importance measures and (ii) through the conversion of the tree structure to a set of rules. In Section 3.5, we show how to extend decision trees to handle multi-output tasks.

3.1 DECISION TREE MODEL

The binary decision tree model is a tree structure built by recursively partitioning the input space. The root node is the node at the top of the tree. We distinguish two types of nodes: (i) the test nodes, also called internal nodes or branching nodes, and (ii) the leaves outputting predictions, also called external nodes or terminal nodes. A test node N_t has two children called the left child and the right child; it furthermore has a splitting rule s_t testing whether or not a sample belongs to its left or right child. For a continuous or categorical ordered input, the splitting rules are typically of the form $s_t(x) = x_{F_t} \leq \tau_t$ testing whether or not the input value x_{F_t} is smaller or equal to a constant τ_t . For a binary or categorical input, the splitting rule is of the form $s_t(x) = x_{F_t} \in B_t$ testing whether or not the input value x_{F_t} belongs to the subset of values B_t .

The decision tree of Figure 3.2 first splits the input space into two disjoint partitions $A_2 = \{x : x_{F_1} \leq \tau_1\}$ and $A_3 = \{x : x_{F_2} > \tau_1\}$ at the root node N_1 . The node N_1 has two children: N_2 the left child and N_3 the right child. The node N_2 is a leaf and thus a terminal partition. The test node N_3 further splits the input space based on a categorical set B_3 into partitions $A_4 = \{x \in A_3, x_{F_3} \in B_3\}$ and $A_5 = \{x \in A_3, x_{F_3} \in B_3\}$ with $A_3 = A_4 \cup A_5$. The input space is further split with 3 more testing nodes with two continuous splitting rules (N_4, N_5) or one categorical splitting rule N_8 . The remaining nodes N_6, N_7, N_9, N_{10} and N_{11} are leaf nodes. In total, the decision tree defines a partition of the input space into 11 (hyper)rectangles (A_1, \dots, A_{11}). A one to one relationship exists between the leaf nodes and the subsets of this input space partition. Note that all partitions of the input space are not expressible as a decision tree structure.

A decision tree predicts the output of an unseen sample by following the tree structure as described by Algorithm 3.1. The recursive process starts at the root node, then the splitting rules of the testing nodes send the sample further down the tree structure. We traverse the tree structure until reaching a leaf, a terminal node, outputting its associated prediction value. A decision tree model $\hat{f} : \mathcal{X} \rightarrow \mathcal{Y}$ is then expressible as a sum of indicator functions, denoted by 1 , over the $|\mathcal{T}|$ tree nodes:

$$\hat{f}(x) = \sum_{t=1}^{|\mathcal{T}|} \beta_t 1(x \in A_t) 1(t \text{ is a leaf}) \quad (3.1)$$

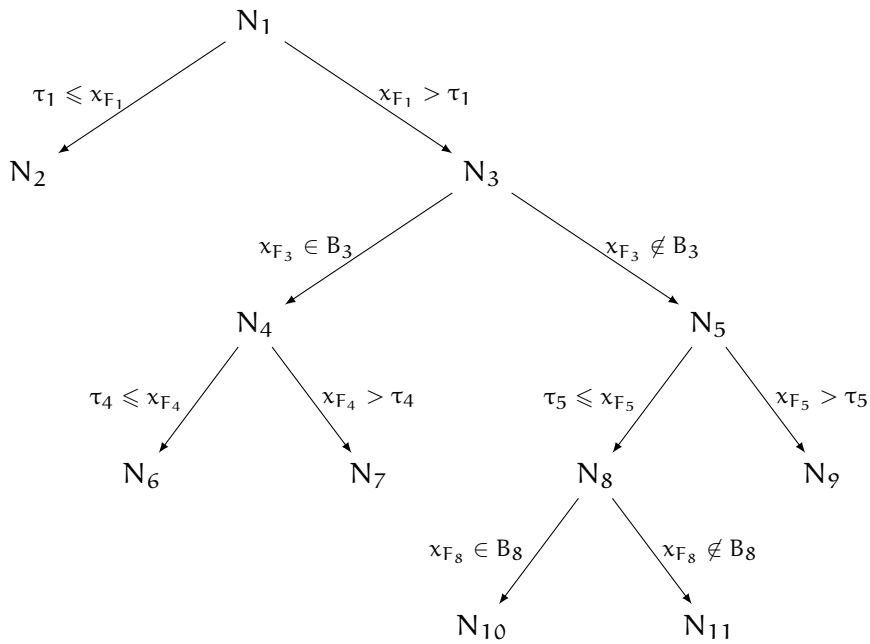


Figure 3.2: A binary decision tree structure containing 11 nodes: 5 test nodes and 6 leaves nodes.

where β_t is the prediction associated to a node N_t . The computational complexity of predicting an unseen sample is thus proportional to the depth of the followed branch.

Algorithm 3.1 Predict a sample x with a decision tree.

```

1: function TREE_PREDICT(tree,  $x$ )
2:    $t \leftarrow$  Index of the root node of the tree.
3:   while the node  $N_t$  is not a leaf. do
4:     if The splitting rule  $s_t(x)$  of node  $N_t$  is true then
5:        $t \leftarrow$  Index of the left child of node  $N_t$ .
6:     else
7:        $t \leftarrow$  Index of the right child node  $N_t$ .
8:     end if
9:   end while
10:  return  $\beta_t$ .
11: end function

```

BINARY VERSUS MULTI-WAY PARTITIONS

Decision trees do not have to respect a binary tree structure. Each one of their internal nodes could have more than two children with multi-way splitting rules. However such multi-way partitions are equivalent to a set of consecutive binary partitions. Furthermore, multi-way splits tends to quickly fragment the training data during the decision tree growth impeding its generalization performance. In prac-

tice, decision trees are therefore most of the time restricted to be binary. (Hastie et al., 2009)

3.2 GROWING DECISION TREES

We grow a decision tree using a set of samples of input-output pairs. Tree growth starts at the root node and divides recursively the input space through splitting rules until we reach a stopping criterion such as a maximal tree depth or minimum sample size in a node. For each new testing node, we search for the best splitting rule to divide the sample set at that node into two subsets. We hope to make partitions “purer” at each new division. The decision tree growing procedure has three main elements:

- a splitting rule search algorithm (see Section 3.2.1);
- stop splitting criteria (see Section 3.2.3) which dictate whenever we stop the development of a branch;
- a leaf labelling rule (see Section 3.2.2) determining the output value of a terminal partition.

Putting all those key elements together leads to the decision tree growing algorithm shown in Algorithm 3.2.

Algorithm 3.2 Grow a decision tree using the sample set $\mathcal{L} = \{(x^i, y^i) \in \mathcal{X} \times \mathcal{Y}\}_{i=1}^n$.

```

1: function GROW_TREE( $\mathcal{L}$ )
2:    $q = \text{EMPTYQUEUE}()$ 
3:   Initialize the tree structure with the root node ( $N_1$ )
4:    $q.\text{ENQUEUE}((1, \mathcal{L}))$ .
5:   while  $q$  is not empty do
6:      $(t, \mathcal{L}_t) \leftarrow q.\text{DEQUEUE}()$ 
7:     if Node  $N_t$  satisfies one stopping criterion then
8:       Label node  $t$  as a leaf using samples  $\mathcal{L}_t$ 
9:     else
10:      Search for the best splitting rule  $s_t$  using samples  $\mathcal{L}_t$ .
11:      Split  $\mathcal{L}_t$  into  $\mathcal{L}_{t,r}$  and  $\mathcal{L}_{t,l}$  using the splitting rule  $s_t$ .
12:      Label node  $t$  as a test node with the splitting rule  $s_t$ .
13:       $q.\text{ENQUEUE}((2t, \mathcal{L}_{t,l}))$ .
14:       $q.\text{ENQUEUE}((2t + 1, \mathcal{L}_{t,r}))$ .
15:     end if
16:   end while
17:   return The grown decision tree.
18: end function

```

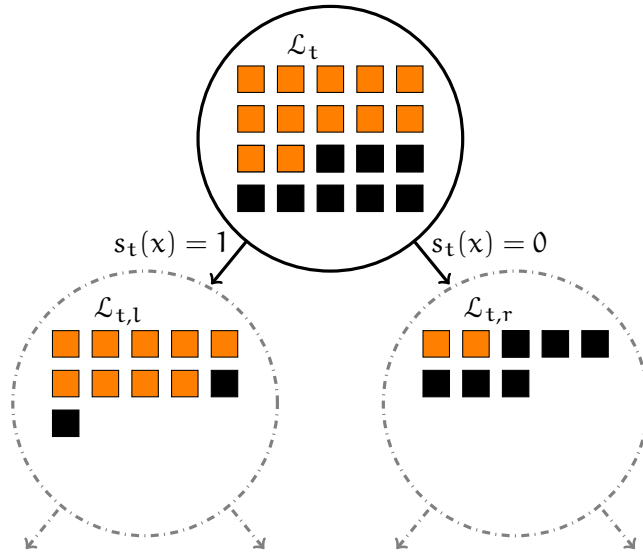


Figure 3.3: During tree growing (here for a binary classification task with orange and black classes), we search for the best splitting rule s_t to divide the sample set \mathcal{L}_t reaching node t into a left $\mathcal{L}_{t,l}$ and a right $\mathcal{L}_{t,r}$ subsets.

3.2.1 Search among node splitting rules

During tree growing, we recursively partition the input space \mathcal{X} and the sample set $\mathcal{L} = \{(x^i, y^i) \in \mathcal{X} \times \mathcal{Y}\}_{i=1}^n$. At each testing node t , we split the sample set \mathcal{L}_t reaching node t into two smaller subsets $\mathcal{L}_{t,l}$ and $\mathcal{L}_{t,r}$ using a binary splitting rule $s_t : \mathcal{X} \rightarrow \{0, 1\}$ as shown in Figure 3.3. This raises two questions (i) what is the set of available binary and axis-wise splitting rules $\Omega(\mathcal{L}_t)$ given a sample set \mathcal{L}_t and (ii) how to select the best one among all of them so as to make the descendants “purer” than the parent node.

For a variable $x_j \in \mathcal{X}_j$ of cardinality c_j , the associated family $Q(x_j)$ of splitting rules consists of all possible subsets of \mathcal{X}_j :

$$Q(x_j) = \{s_t(x) \equiv 1(x_j \in \mathcal{X}') : \mathcal{X}' \subset \mathcal{X}_j\}. \quad (3.2)$$

The size of the splitting rule family is increasing exponentially with the total number of possible values ($|Q(x_j)| = 2^{|\mathcal{X}_j| - 1}$).

If the possible values of the variable x_j are also ordered, we can reduce the size of the splitting rule family $Q(x_j)$ from an exponential number of candidates to a linear number of splitting rules ($|Q(x_j)| = |\mathcal{X}_j| - 1$):

$$Q(x_j) = \left\{ s_t(x) \equiv 1(x_j \leq \tau) : \tau \in \mathcal{X}_j \right\}. \quad (3.3)$$

With a numerical variable $x_j \in \mathbb{R}$, the number of possible splitting rules is infinite. However, the training set is of finite size. We consider

a family of splitting rules similar to Equation 3.3 with the possible values \tilde{x}_j available in the training set.

The selected splitting rule s_t should split the sample set \mathcal{L}_t such that the following conditions hold: the sample sets $\mathcal{L}_{t,r}$ and $\mathcal{L}_{t,l}$ are non empty ($\mathcal{L}_{t,r} \neq \emptyset, \mathcal{L}_{t,l} \neq \emptyset$) forming a disjoint ($\mathcal{L}_{t,l} \cap \mathcal{L}_{t,r} = \emptyset$) and non overlapping partition ($\mathcal{L}_t = \mathcal{L}_{t,l} \cup \mathcal{L}_{t,r}$) of the original sample set \mathcal{L}_t . During the expansion of a test node t into a left child and a right child, we thus select a splitting rule s_t among all possible splitting rules $\Omega(\mathcal{L}_t)$:

$$s_t \in \Omega(\mathcal{L}_t) = \left\{ s : s \in \bigcup_{j \in \{1, \dots, p\}} Q(x_j), \right. \\ \mathcal{L}_{t,l} = \{(x, y) \in \mathcal{L}_t : s(x) = 1\}, \\ \mathcal{L}_{t,r} = \{(x, y) \in \mathcal{L}_t : s(x) = 0\}, \\ \left. \mathcal{L}_{t,l} \neq \emptyset, \mathcal{L}_{t,r} \neq \emptyset \right\}. \quad (3.4)$$

We strive to select the “best” possible local splitting rule s_t for the split at node t leading ideally to good generalization performance. However, it is impossible to minimize directly the generalization error. Thus instead, we are going to minimize the resubstitution error, the error over the training set. However, obtaining such a tree is trivial and it has poor generalization performance. A more meaningful criterion is to search for the smallest tree minimizing the resubstitution error. However, this optimization problem is a NP-complete (Hyafil and Rivest, 1976). Instead, we greedily grow the tree by maximizing the reduction of an impurity measure function $I : (\mathcal{X} \times \mathcal{Y}) \times \dots \times (\mathcal{X} \times \mathcal{Y}) \rightarrow \mathbb{R}$. Mathematically, we define the impurity reduction ΔI obtained by dividing the sample set \mathcal{L}_t into two partitions ($\mathcal{L}_{t,r}, \mathcal{L}_{t,l}$) as

$$\Delta I(\mathcal{L}_t, \mathcal{L}_{t,l}, \mathcal{L}_{t,r}) = I(\mathcal{L}_t) - \frac{|\mathcal{L}_{t,l}|}{|\mathcal{L}_t|} I(\mathcal{L}_{t,l}) - \frac{|\mathcal{L}_{t,r}|}{|\mathcal{L}_t|} I(\mathcal{L}_{t,r}). \quad (3.5)$$

The splitting rule selection problem (line 10 of the tree growing Algorithm 3.2) is thus written as

$$s_t = \arg \max_{s \in \Omega(\mathcal{L}_t)} \Delta I(\mathcal{L}_t, \mathcal{L}_{t,l}, \mathcal{L}_{t,r}) \quad (3.6)$$

Intuitively, the impurity I should be minimal whenever all samples have the same output value. The node is then said to be “pure”

Given the additivity of the impurity reduction, the best split at node t according to the impurity reduction computed locally is also the best split at this node in terms of global impurity. The remaining tree impurity $\text{Imp}(T)$ of a tree T is the sum of the remaining impurities of all leaf nodes:

$$\text{Imp}(T) = \sum_{t \in T} p(t) I(\mathcal{L}_t) \mathbf{1}(t \text{ is a leaf}) \quad (3.7)$$

with $p(t) = |\mathcal{L}_t|/|\mathcal{L}|$ the proportion of samples reaching node t . If we develop the leaf node t into a test node, it leads to a new tree T' with a new splitting rule s_t having a left t_l and a right t_r children node. The overall impurity decreases from the original tree T to the bigger tree T' and the impurity decrease is given by

$$\text{Imp}(T) - \text{Imp}(T') = p(t)I(\mathcal{L}_t) - p(t_l)I(\mathcal{L}_{t,l}) - p(t_r)I(\mathcal{L}_{t,r}) \quad (3.8)$$

$$= p(t)\Delta I(\mathcal{L}_t, \mathcal{L}_{t,l}, \mathcal{L}_{t,r}) \quad (3.9)$$

Thus, the decision tree growing procedure is a repeated process aiming at decreasing the total impurity as quickly as possible by suitably choosing the local splitting rules.

In classification, a node is pure if all samples have the same class. Given a sample set \mathcal{L}_t reaching node t , we will denote by $p(y = l|t) = \frac{1}{|\mathcal{L}_t|} \sum_{(x,y) \in \mathcal{L}_t} \mathbf{1}(y = l)$ the proportion of samples reaching node t having the class l . A node will be pure if $p(y = l|t)$ is equal to $\mathbf{1}$ for a class l and zero for the others. The node impurity should increase whenever samples with different classes are mixed together. We require that the impurity measure I in classification satisfies the following properties:

1. I is minimal only whenever the node is pure $p(y = l|t) = 1$ and $p(y = m|t) = 0 \quad \forall m \in \{1, \dots, l-1, l+1, \dots, k\}$;
2. I is maximal only whenever $p(y = l|t) = 1/k \quad \forall l \in \{1, \dots, k\}$;
3. I is a symmetric function with respect to the class proportions $p(y = 1|t), \dots, p(y = k|t)$ so as not to favor any class.

A first function satisfying those three properties is the misclassification error rate:

$$\text{Error rate}(\mathcal{L}_t) = 1 - \max_{l \in \{1, \dots, k\}} p(y = l|t). \quad (3.10)$$

However, this is not an appropriate criterion. In practice, many candidate splitting rules have often the same error rate reduction, especially in the multi-class classification where only the number of samples of the majority class matters.

Consider the following split selection problem, we have a binary classification task with 500 negative and 500 positive samples. The first splitting rule leads to a left child with 125 positive and 375 negative samples, while the right child has 375 positive and 125 negative samples. The misclassification error reduction is thus given by:

$$\frac{500}{1000} - 2 \frac{500}{1000} \left(1 - \frac{375}{500}\right) = 0.25.$$

Now, let's consider a second splitting rule leading to a pure node with 250 positive samples and another node with 250 positive and 500

negative samples. This second split has the same impurity reduction score leading to a tie:

$$\frac{500}{1000} - \frac{750}{1000} \left(1 - \frac{500}{750}\right) - \frac{250}{1000} \left(1 - \frac{250}{250}\right) = 0.25.$$

The misclassification does not discriminate enough node purity as it varies linearly with the fraction of the majority class. To solve this issue, we add a fourth required properties to classification impurity functions (Breiman et al., 1984):

4. I must be a strictly concave function with respect to the class proportion $p(y = l|t)$.

This fourth property will increase the granularity of impurity reduction scores leading fewer ties in splitting rule scores. It will reduce the tree instability with respect to the training set.

Two more suitable classification impurity criteria satisfying all four properties are the Gini index, a statistical dispersion measure, and the entropy measure, an information theory measure.

$$\text{Gini}(\mathcal{L}_t) = \sum_{l=1}^k p(y = l|t)(1 - p(y = l|t)) \quad (3.11)$$

$$\text{Entropy}(\mathcal{L}_t) = - \sum_{l=1}^k p(y = l|t) \log p(y = l|t) \quad (3.12)$$

By minimizing the Gini index, we minimize the class dispersion. While selecting the splitting rule based on the entropy measure minimizes the unpredictability of the target, the remaining unexplained information of the output variable.

Given the strict concavity of the Gini index and entropy, we can now discriminate the two splits of the previous example. The first split would have an impurity reduction with the Gini index of 0.0625 and the entropy of ≈ 0.131 . The second split with a pure node would have an impurity reduction of ≈ 0.083 with the Gini index and of ≈ 0.216 with the entropy. Based on these criteria, both measures would choose the second split.

In regression tasks, we consider a node as pure if the dispersion of the output values is zero. We require the impurity regression criterion to be zero only if all output values have the same value. A common dispersion measure used to grow regression trees is the empirical variance

$$\text{Variance}(\mathcal{L}_t) = \frac{1}{|\mathcal{L}_t|} \sum_{(x,y) \in \mathcal{L}_t} (y - \bar{y})^2, \text{ with } \bar{y} = \frac{1}{|\mathcal{L}_t|} \sum_{(x,y) \in \mathcal{L}_t} y \quad (3.13)$$

By maximizing the variance reduction, we are searching for a splitting rule minimizing the square loss $\ell(y, y') = \frac{1}{2}(y - y')^2$. Note that

the Gini index and the empirical variance lead to the same impurity measure for binary classification tasks and multi-label classification tasks with output classes encoded with $\{0, 1\}$ numerical variables.

3.2.2 Leaf labelling rules

When tree growing is stopped by the activation of a stop splitting criterion, the newly created leaf needs to be labeled by an output value (see line 8 of Algorithm 3.2). It is a constant β_t chosen to minimize a given loss function ℓ over the samples $\mathcal{L}_t = \{(x, y) \in (\mathcal{X}, \mathcal{Y})\}$ reaching the node t :

$$\beta_t = \arg \min_{\beta} \sum_{(x, y) \in \mathcal{L}_t} \ell(y, \beta). \quad (3.14)$$

In regression tasks, we want to find the constant β_t minimizing the square loss in single output regression :

$$\beta_t = \arg \min_{\beta} \sum_{(x, y) \in \mathcal{L}_t} \frac{1}{2} (y - \beta)^2 \quad (3.15)$$

By setting the first derivative to zero, we have

$$\sum_{(x, y) \in \mathcal{L}_t} (y - \beta) = 0 \quad (3.16)$$

$$\beta_t = \frac{1}{|\mathcal{L}_t|} \sum_{(x, y) \in \mathcal{L}_t} y. \quad (3.17)$$

The constant leaf model minimizing the square loss is the average output value of the samples reaching node t .

In classification tasks, the constant β_t minimizing the 0 – 1 loss ($\ell_{0-1}(y, y') = 1(y \neq y')$) is the most frequent class. The decision tree can also be a class probability estimator by outputting the proportion $p(l|t)$ of samples of class l reaching node t from the sample set \mathcal{L}_t .

BEYOND CONSTANT LEAF MODELING

The leaf labelling rule can go beyond a constant model with supervised learning models such as linear models (Frank et al., 1998; Landwehr et al., 2005; Quinlan et al., 1992; Wang and Witten, 1996), kernel-based methods (Torgo, 1997), probabilistic models (Kohavi, 1996) or even tree-ensemble models (Matthew et al., 2015). It increases the modeling power of the decision tree at the expense of computing time and new hyper-parameters.

3.2.3 Stop splitting criterion

The tree growth at a node t naturally stops if all the samples \mathcal{L}_t reaching the node t (i) share the same output value (zero impurity) or

(ii) share the same input values (but not necessarily the same output) as in this case we can not find a valid split of the data. In both cases, we can not find a splitting rule to grow the tree further due to a lack of data.

A tree developed in such ways is then said to be fully developed. The question is “Should we stop sooner the tree growth?”. A testing node t splits the data \mathcal{L}_t into two partitions $(\mathcal{L}_{t,l}, \mathcal{L}_{t,r})$ leading to a left child node t_l and a right child node t_r . If we denote by \hat{f}_t , $\hat{f}_{t,r}$ and $\hat{f}_{t,l}$ the leaf models that would be assigned to the nodes t , t_r or t_l , we have that the resubstitution error reduction ΔErr associated to a loss function ℓ is given by:

$$\begin{aligned} \Delta\text{Err} &= \sum_{(x,y) \in \mathcal{L}_t} \ell(y, \hat{f}_t(x)) - \sum_{(x,y) \in \mathcal{L}_{t,r}} \ell(y, \hat{f}_{t,r}(x)) \\ &\quad - \sum_{(x,y) \in \mathcal{L}_{t,l}} \ell(y, \hat{f}_{t,l}(x)) \end{aligned} \quad (3.18)$$

$$\begin{aligned} &= \sum_{(x,y) \in \mathcal{L}_{t,r}} [\ell(y, \hat{f}_t(x)) - \ell(y, \hat{f}_{t,r}(x))] \\ &\quad + \sum_{(x,y) \in \mathcal{L}_{t,l}} [\ell(y, \hat{f}_t(x)) - \ell(y, \hat{f}_{t,l}(x))] \end{aligned} \quad (3.19)$$

Since we choose $\hat{f}_{t,r}$ and $\hat{f}_{t,l}$ so as to minimize the resubstitution error on their respective training data $(\mathcal{L}_{t,l}$ and $\mathcal{L}_{t,r})$, the resubstitution error never increases through node splitting. With only “natural” splitting rule, decision trees are optimally fitting the training data.

Stop splitting criteria avoid over-fitting by stopping earlier the tree growth. They are either based on (i) structural properties or on (ii) data statistics. Criteria computed on the left and the right children can also discard splitting rules, for example requiring a minimal number of samples in the left and right children to split a node.

Structural-based stop splitting criteria regularize the tree growth by explicitly limiting the tree complexity, by restricting for example:

- branch depths or
- the total number of nodes.

In the second case, the order in which the tree nodes are split starts to matter and can be chosen so as to maximize the total impurity reduction.

Data-based stop splitting criteria stop the tree growth if some statistical properties computed on the data used to split the node are below a threshold such as

- the number of samples reaching the node,
- the number of samples reaching the left and right children obtained after splitting,

- the impurity reduction or
- the p-value of a significance test, such as a Chi-square test, testing the independence of the split and the output variable.

3.3 RIGHT DECISION TREE SIZE

To find the right decision tree size, there are two main families of complexity reduction techniques, also called pruning techniques: (i) pre-pruning techniques stop the tree growth before the tree is fully developed (line 5 of Algorithm 3.2 and presented in Section 3.2.3) and (ii) post-pruning techniques remove tree nodes a posteriori setting a trade-off between the tree size and the resubstitution error. Both approaches lead to smaller decision trees aiming to improve generalization performance and to simplify decision tree interpretation.

Pre-pruning criteria are straightforward tools to control the decision tree size. However, it is unclear which pruning level (or hyperparameter values) leads to the best generalization performance. Too “strict” stop splitting criteria will grow shallow trees under-fitting the data. While too “loose” stop splitting criteria have the opposite effect, i.e., growing overly complex trees over-fitting the data.

While pre-pruning techniques select the tree complexity a priori, post-pruning techniques select the optimal complexity a posteriori. A naive approach to post-pruning would be to build independently a sequence of decision trees with different complexity by varying the stop splitting criteria, and then to select the one minimizing an approximation of the generalization error such as the error on a hold out sample set. However, this is not computationally efficient as it re-grows each time a (new) decision tree.

Post-pruning techniques first grow a single decision tree T with very loose or no stop splitting criterion. This decision tree clearly overfits the training data. Then, they select a posteriori a subtree $T^* \subseteq T$ among all possible subtrees of T . The original decision tree is thus pruned by collapsing nodes from the original tree into new leaf nodes. The post-pruning method minimizes a tradeoff between the decision tree error over a sample set S and a function measuring the decision tree complexity such as the number of nodes:

$$T^*(\lambda) = \arg \min_{\check{T} \subseteq T} \text{Error}(S|\check{T}) + \lambda \text{Complexity}(\check{T}). \quad (3.20)$$

The cost complexity pruning method (Breiman et al., 1984), also known as the weakest link pruning, implements Equation 3.20 through a complexity coefficient C_α measuring a tradeoff between the resubstitution error of a tree \check{T} and its complexity $|\check{T}|$ defined by the number of leaves:

$$C_\alpha(\check{T}) = \text{resubstitution error}(\mathcal{L}|\check{T}) + \alpha|\check{T}|. \quad (3.21)$$

For each α , there exists a unique tree \check{T}_α minimizing the cost complexity coefficient C_α . Large values of α lead to small trees, while conversely small values of α allow bigger sub-trees. For the extreme case $\alpha = 0$ (resp. $\alpha = \infty$), we have the original decision tree T_0 (resp. the subtree containing only the root node). One can show (Breiman et al., 1984) that we can sequentially obtain the \check{T}_α from the original tree T by removing the node with the smallest increase in resubstitution error. We select the optimal subtree T^* among all subtrees $\check{T}_\alpha \subseteq T$ by minimizing an approximation of the generalization error using for instance cross-validation methods.

The reduced error pruning method (Quinlan, 1987), another post pruning technique, splits the learning set into a training set and a pruning set. It grows on the training set a large decision tree. During the pruning phase, it first computes the error reduction of pruning each node and its descendants on the pruning set, then removes greedily the node reducing the most the error. It repeats this two steps procedure until the error on the pruning set starts increasing.

Other post pruning methods have been developed with pruning criteria based on statistical procedures (Quinlan, 1987, 1993) or on cost complexity criteria based on information theory (Mehta et al., 1995; Quinlan and Rivest, 1989; Wallace and Patrick, 1993). Instead of relying on greedy processes, authors (Almuallim, 1996; Bohanec and Bratko, 1994) have proposed dynamic programming algorithms to find an optimal sub-tree sequence minimizing the resubstitution error with increasing tree complexity at the expense of computational complexity.

3.4 DECISION TREE INTERPRETATION

A strength of the decision tree model is its interpretability. A closer inspection reveals that we can convert a decision tree model to a set of mutually exclusive classification or regression rules. We get these rules by following the path from each leaf to the root node. We have converted the decision tree shown in Figure 3.4 to three sets of predicting rules, one for each class of iris flower (Versicolor, Virginica and Setosa):

1. "If Petal width ≤ 0.7 cm, then Setosa"
2. "If Petal width > 0.7 cm and Petal width ≤ 1.65 cm and Petal length ≤ 5.25 cm, then Versicolor."
3. "If Petal width > 0.7 cm and Petal width > 1.65 cm and Sepal length ≤ 5.95 cm and Sepal length > 5.85 cm, then Versicolor."
4. "If Petal width > 0.7 cm and Petal width ≤ 1.65 cm and Petal length > 5.25 cm, then Virginica."

5. "If Petal width > 0.7cm and Petal width > 1.65cm and Sepal length ≤ 5.95cm and Sepal length ≤ 5.85, then Virginica."
6. "If Petal width > 0.7cm and Petal width > 1.65cm and Sepal length > 5.95cm, then Virginica."

Remark that given the binary hierarchical structure, some rules are redundant and can be further simplified. For instance, we can collapse the constraints "Petal width > 0.7cm and Petal width > 1.65cm" into "Petal width > 1.65cm" for the 6-th rule.

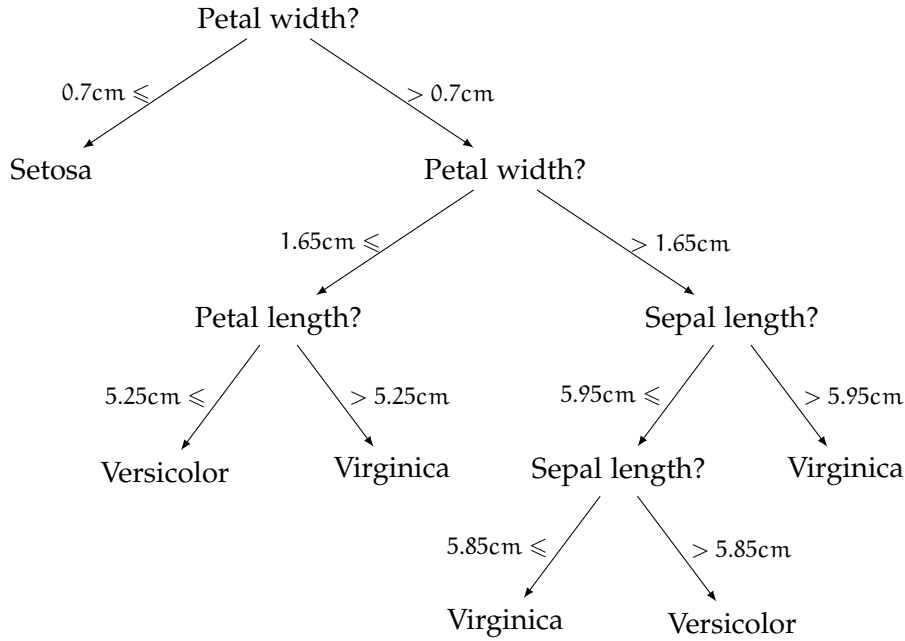


Figure 3.4: A decision tree is an interpretable set of rules organized hierarchically. For instance, we can recognize a Versicolor iris with two sets of rules. Among the four input variables (petal width, petal length, sepal width and sepal length), the decision tree shows that only one input variable is necessary to classify the Setosa iris variety.

The decision tree model also shows which input variables x_1, \dots, x_p are important to predict the output variable(s) y . During the decision tree growth, we select at each node t an axis-wise splitting rules s_t minimizing the reduction of an impurity measure ΔI dividing the samples \mathcal{L}_t reaching the node into two subsets $(\mathcal{L}_{t,l}, \mathcal{L}_{t,r})$. The mean decrease of impurity (MDI) of a variable x_j (Breiman et al., 1984) sums, over the nodes of a decision tree T where x_j is used to split, the total reduction of impurity associated to the split weighted by the probability of reaching that node over the training set \mathcal{L} :

$$MDI(x_j) = \sum_{\{t \in T: x_j \text{ is tested by } s_t\}} \frac{|\mathcal{L}_t|}{|\mathcal{L}|} \Delta I(\mathcal{L}_t, \mathcal{L}_{t,l}, \mathcal{L}_{t,r}). \quad (3.22)$$

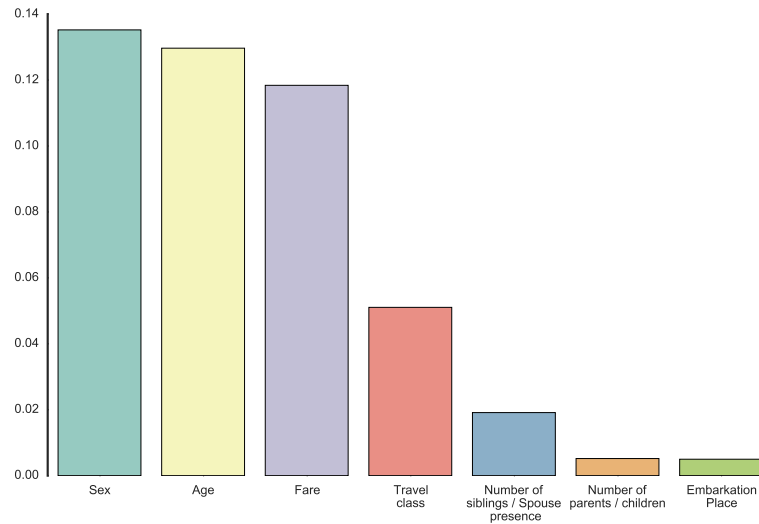


Figure 3.5: Mean impurity decrease for each input variable of the Titanic dataset for a decision tree whose complexity was chosen to maximize the accuracy score on a hold out pruning sample set.

The mean decrease of impurity scores and ranks the input variables according to their importances during the decision tree growth process as illustrated in Figure 3.5. It takes into account variable correlations, multivariate and non linear effects. As such, decision trees are often used as pre-processing tools to select a fraction of the top most important variables.

3.5 MULTI-OUTPUT DECISION TREES

Decision trees naturally extend from single output tasks to multiple output tasks (Blockeel et al., 2000; Clare and King, 2001; De’Ath, 2002; Noh et al., 2004; Segal, 1992; Siciliano and Mola, 2000; Vens et al., 2008; Zhang, 1998) such as multi-output regression, multi-label classification or multi-class classification. No core modification are needed. Instead, we need appropriate impurity measures and leaf labelling rules for the tree prediction Algorithm 3.1 and the tree growth Algorithm 3.2. Note that a multi-output decision tree can still be pruned (Struyf and Džeroski, 2005).

MULTI-OUTPUT IMPURITY MEASURES During the decision tree growth, we aim to select a splitting rule dividing the sample set \mathcal{L}_t reaching the node t into a left and a right sample sets $(\mathcal{L}_{t,l}, \mathcal{L}_{t,r})$. The best multi-output splitting rule is the one maximizing the reduction

of a multi-output impurity measure I . We can use native multi-output impurity measures such as the variance in regression (Segal, 1992):

$$\text{Variance}(\mathcal{L}_t) = \frac{1}{|\mathcal{L}_t|} \sum_{(x,y) \in \mathcal{L}_t} |y - \bar{y}|_2^2 \quad \text{with } \bar{y} = \frac{1}{|\mathcal{L}_t|} \sum_{(x,y) \in \mathcal{L}_t} y. \quad (3.23)$$

or any impurity criterion derived from an appropriate distance measure (Blockeel et al., 2000).

We can also extend known impurity measures, such as the Gini index or the entropy (see Section 3.2.1), by summing the impurity measures over each output (De'Ath, 2002):

$$I_{\text{mo}}(\mathcal{S}) = \sum_{j=1}^d I(\{(x, y_j) \in \mathcal{S}\}), \quad (3.24)$$

where $\mathcal{S} = \{(x^i, y^i) \in (\mathcal{X} \times \mathcal{Y})\}_{i=1}^n$ is a sample set.

Since we can define an impurity measure on any set of outputs, we can derive the mean decrease of impurity MDI (see Section 3.4) either on all or a subset of the outputs.

LEAF LABELLING AND PREDICTION RULE In the multi-output context, the leaf prediction β_t of a node t is a constant vector of output values chosen so as to minimize a multi-output loss function ℓ over the samples $\mathcal{L}_t = \{(x^i, y^i) \in (\mathcal{X}, \mathcal{Y})\}_{i=1}^n$ reaching the node:

$$\beta_t = \arg \min_{\beta} \sum_{(x,y) \in \mathcal{L}_t} \ell(y, \beta). \quad (3.25)$$

In multi-output regression, the loss ℓ is commonly the ℓ_2 -norm loss $\ell_2(y, y') = \frac{1}{2} \|y - y'\|_2^2$ the multi-output extension of the square loss. The constant β_t minimizing the ℓ_2 -norm loss is the average output vector

$$\beta_t = \frac{1}{|\mathcal{L}_t|} \sum_{(x,y) \in \mathcal{L}_t} y. \quad (3.26)$$

Whenever we extend the label rule assignment to multi-label and to multi-output classification tasks, there are two common possibilities either minimizing the subset 0–1 loss which is equal to zero if only if all outputs are correctly predicted $\ell_{\text{subset } 0-1}(y, y') = 1(y \neq y')$ and the Hamming loss which counts among the d outputs the number of wrongly predicted outputs $\ell_{\text{Hamming}}(y, y') = \sum_{j=1}^d 1(y_j \neq y'_j)$.

Minimizing the subset 0–1 loss takes into account output correlations. The constant vector β_t minimizing this loss is the most frequent output value combination. Note that the constant β_t might not be unique as we might have several output combinations with the same frequency of appearance in the sample set reaching the leaf.

The Hamming loss makes the assumption that all outputs are independent. The constant β_t minimizing this loss is the vector containing the most frequent class of each output.

When the trees are fully developed with only pure leaves, minimizing the Hamming loss or the subset 0 – 1 loss leads to identical leaf prediction rules.

OUTLINE

Ensemble methods combine supervised learning models together so as to improve generalization performance. We present two families of ensemble methods: averaging methods and boosting methods. Averaging ensembles grow independent unstable estimators and average their predictions. Boosting methods increase sequentially their total complexity by adding biased and stable estimators. In this chapter, we first show how to decompose the generalization error of supervised learning estimators into their bias, variance and irreducible error components. Then we show how to exploit averaging techniques to reduce variance and boosting techniques to sequentially decrease bias.

Ensemble methods fit several supervised learning models instead of a single one and combine their predictions. The goal is to reduce the generalization error by solving the same supervised learning task multiple times. We hope that the errors made by the different models will compensate and thereby improve the overall accuracy whenever we consider them together.

Real life examples of “ensemble methods” in the human society are democratic elections. Each eligible person is asked to cast its vote for instance to choose between political candidates. This approach considers each person of the committee as an independent expert and averages simultaneously their opinions. In supervised learning, these kinds of voting mechanism are called “averaging methods”.

Instead of querying all experts independently, we can instead collect their opinions sequentially. We ask to each new expert to refine the predictions made by the previous ones. The expert sequence is chosen so that each element of the sequence improves the accuracy focusing on the unexplained phenomena. For instance in medical diagnosis, a person itself is the first one to assess its health status. The next expert in the line is the general practitioner followed by a series of specialists. We call these ensemble methods “boosting methods”.

The “averaging” approach aims to reduce the variability in the expert pool by averaging their predictions. At the other end, the “boosting” approach carefully refines its predictions by cumulating the predictions of each expert.

In Section 4.1, we show how to decompose the error made by supervised learning models into three terms: a variance term due to the

variability of the model with respect the learning sample set, a bias term due to a lack of modeling power of the estimator and an irreducible error term due to the nature of the supervised learning task. Averaging methods presented in Section 4.2 are variance reducing techniques growing independently supervised learning estimators. In Section 4.3, we show how to learn sequentially a series of estimators through boosting methods increasing the overall ensemble complexity and reducing the ensemble model bias.

4.1 BIAS-VARIANCE ERROR DECOMPOSITION

The expected error or generalization error Err associated to a loss $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$ of a supervised learning algorithm is a random variable depending on the learning samples $\mathcal{L} = \{(x^i, y^i) \in \mathcal{X} \times \mathcal{Y}\}_{i=1}^n$ drawn independently and identically from a distribution $P_{\mathcal{X}, \mathcal{Y}}$ and used to fit a model $f_{\mathcal{L}} : \mathcal{X} \rightarrow \mathcal{Y}$ in a hypothesis space $\mathcal{H} \subset \mathcal{Y}^{\mathcal{X}}$. We want here to analyze the expectation of the generalization error Err over the distribution of learning samples defined as follows:

$$E_{\mathcal{L}}(\text{Err}) = E_{\mathcal{L}} E_{P_{\mathcal{X}, \mathcal{Y}}} \{\ell(f_{\mathcal{L}}(x), y)\} \quad (4.1)$$

$$= E_{\mathcal{L}} E_{P_{\mathcal{X}}} E_{P_{\mathcal{Y}|\mathcal{X}}} \{\ell(f_{\mathcal{L}}(x), y)\} \quad (4.2)$$

Let us denote the Bayes model by $f_{\text{Bayes}}(x) \in \mathcal{Y}^{\mathcal{X}}$, the best model possible minimizing the generalization error:

$$f_{\text{Bayes}}(x) = \arg \min_{f \in \mathcal{Y}^{\mathcal{X}}} E_{P_{\mathcal{X}, \mathcal{Y}}} \{\ell(f(x), y)\}. \quad (4.3)$$

For the squared loss $\ell(y, y') = \frac{1}{2}(y - y')^2$, we can decompose the expected error over the learning set $E_{\mathcal{L}}(\text{Err})$ into three terms (see (German et al., 1992) for the proof):

$$E_{\mathcal{L}}(\text{Err}) \triangleq E_{\mathcal{L}} E_{P_{\mathcal{X}, \mathcal{Y}}} \{(f_{\mathcal{L}}(x) - y)^2\} \quad (4.4)$$

$$= E_{P_{\mathcal{X}}} \left\{ \text{Var}_{\mathcal{L}} \{f_{\mathcal{L}}(x)\} + \text{Bias}^2(f_{\mathcal{L}}(x)) + \text{Var}_{P_{\mathcal{Y}|\mathcal{X}}} \{y\} \right\}, \quad (4.5)$$

where

$$\text{Var}_{\mathcal{L}} \{f_{\mathcal{L}}(x)\} = E_{\mathcal{L}} \left\{ (f_{\mathcal{L}}(x) - f_{\text{avg}}(x))^2 \right\}, \quad (4.6)$$

$$\text{Bias}^2(f_{\mathcal{L}}(x)) = (E_{\mathcal{L}} \{f_{\mathcal{L}}(x)\}(x) - f_{\text{Bayes}}(x))^2. \quad (4.7)$$

We interpret each term of Equation 4.5 as follows:

- The variance of a supervised learning algorithm $E_{P_{\mathcal{X}}} \text{Var}_{\mathcal{L}} \{f_{\mathcal{L}}(x)\}$ describes the variability of the model with a randomly drawn learning sample set \mathcal{L} . Supervised learning algorithms with a high variance have often a high complexity which makes them overfit the data.

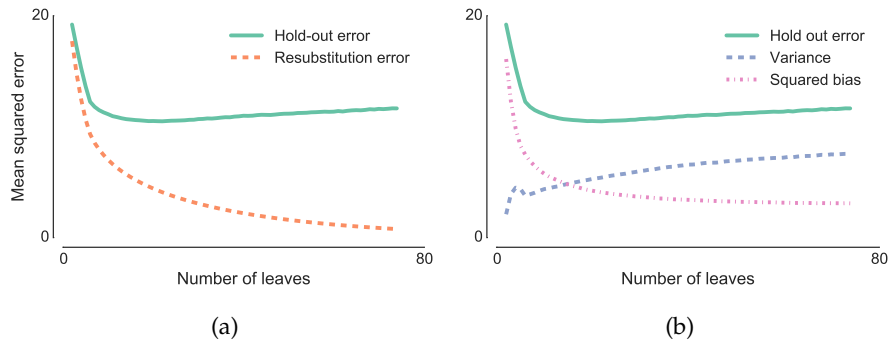


Figure 4.1: Bias-variance decomposition of decision tree models of increasing complexity on the Friedman1 dataset.

- The square bias $E_{P_{\mathcal{X}}} \{ \text{Bias}^2(f_{\mathcal{L}}(x)) \}$ is the distance between the average model and the Bayes model. Biased models are not complex enough to model the input-output function. They are indeed underfitting the data.
- The irreducible error $E_{P_{\mathcal{X}}} \text{Var}_{P_{y|x}} \{y\}$ is the variance of the target around its true mean. It is the minimal attainable error on a supervised learning problem.

The bias-variance decomposition allows to analyze and to interpret the effect of hyper-parameters. It highlights and gives insights on their effects on the bias and the variance. In Figure 4.1, we have fitted decision tree models with an increasing number of leaves on the Friedman1 dataset (Friedman, 1991), a simulated regression dataset. We first assess the resubstitution error over 300 samples and an approximation of the generalization error, the hold out error, computed on an independent testing set of 20000 samples. We compute these errors (see Figure 4.1a) by averaging the performance of decision tree models over 100 learning sets \mathcal{L} drawn from the same distribution $P_{\mathcal{X},y}$. By increasing the number of leaves, the resubstitution error decreases up to zero with fully developed trees. On the other hand, the hold out error starts increasing beyond 20 leaves indicating that the model is under-fitting with less than 20 leaves and over-fitting with more than 20 leaves. By increasing the number of leaves, we decrease the bias as we grow more complex models as shown in Figure 4.1b. It also increases the variance as the tree structures become more unstable with the learning set \mathcal{L} .

In general, we have the following trends for a single decision tree model. Large decision trees overfit and are unstable with respect to the learning set \mathcal{L} which corresponds to a high variance and a small bias. Shallow decision trees, on the other hand, underfit the learning set \mathcal{L} and have stable structures, which corresponds to a small variance and a high bias. The pruning technique presented in Section 3.3

allows to select a tradeoff between the variance and the bias of the algorithm by adjusting the tree complexity.

The bias-variance decomposition of the square loss is by far the most studied decomposition, but there nevertheless exist similar decompositions for other losses, e.g. see (Domingos, 2000) for a decomposition of the polynomial loss $\ell(y, y') = |y - y'|^p$, see (Domingos, 2000; Friedman, 1997; Kohavi et al., 1996; Tibshirani, 1996a) for the 0 – 1 loss, or see (James, 2003) for losses in general.

4.2 AVERAGING ENSEMBLES

An averaging ensemble model $f_{\theta_1, \dots, \theta_M}$ builds a set of M supervised learning models $\{f_{\theta_m}\}_{\theta_m=1}^M$, instead of a single one. Each model f_{θ_m} of the ensemble is different as we randomize and perturb the original supervised learning algorithm at fitting time. We describe entirely the induced randomization of one model f_{θ_m} by drawing i.i.d. a random vector of parameters θ_m from a distribution of model parameters P_θ .

In regression, the averaging ensemble predicts an unseen sample by averaging the predictions of each model of the ensemble:

$$f_{\theta_1, \dots, \theta_M}(x) = \sum_{m=1}^M f_{\theta_m}(x). \quad (4.8)$$

It minimizes the square loss (or its extension the ℓ_2 -norm loss) between the ensemble model and its members:

$$f_{\theta_1, \dots, \theta_M}(x) = \arg \min_{y \in \mathcal{Y}} \sum_{m=1}^M (y - f_{\theta_m}(x))^2. \quad (4.9)$$

In classification, the averaging ensemble combines the predictions of its members to minimize the 0-1 loss by a majority vote of all its members:

$$f_{\theta_1, \dots, \theta_M}(x) = \arg \min_{c \in \mathcal{Y}} \sum_{m=1}^M 1(f_{\theta_m}(x) \neq c). \quad (4.10)$$

An alternative approach, called soft voting, is to classify according to the average of the probability estimates $\hat{P}_{f_{\theta_m}(x)}$ provided by the ensemble members:

$$f_{\theta_1, \dots, \theta_M}(x) = \arg \max_{c \in \mathcal{Y}} \sum_{m=1}^M \hat{P}_{f_{\theta_m}(x)}(Y = c). \quad (4.11)$$

Both approaches have been studied and yield almost exactly the same result, but soft voting provides smoother probability class estimates than majority vote (Breiman, 1996a; Zhou, 2012). The multi-output extension to ensemble predictions often minimizes the Hamming loss applying either soft-voting or majority voting to each output independently. Minimizing the subset 0 – 1 loss for multi-label tasks would lead to predict the most frequent label set.

AMBIGUITY DECOMPOSITION

The ambiguity decomposition (Krogh et al., 1995) of the square loss shows that the generalization error of an ensemble $f_{\theta_1, \dots, \theta_M}$ of M models $\{f_{\theta_m}\}_{m=1}^M$ is always lower or equal than the average generalization error \bar{E} of its members:

$$E_{P_x} E_{P_{y|x}} \{(y - f_{\theta_1, \dots, \theta_M}(x))^2\} = \bar{E} - \bar{A} \leq \bar{E} \quad (4.12)$$

with

$$\bar{E} = \frac{1}{M} \sum_{m=1}^M E_{P_x} E_{P_{y|x}} \{(y - f_{\theta_m}(x))^2\} \quad (4.13)$$

$$\bar{A} = \frac{1}{M} \sum_{m=1}^M E_{P_x} \{(f_{\theta_m}(x) - f_{\theta_1, \dots, \theta_M}(x))^2\} \quad (4.14)$$

The ambiguity term \bar{A} is the variance of the ensemble around its average model $f_{\theta_1, \dots, \theta_M}$. The equality occurs only if all average models are identical $f_{\theta_1, \dots, \theta_M} = f_{\theta_m} \quad \forall m \in \{1, \dots, M\}$.

Averaging ensemble models are obtained by first perturbing supervised learning models and then combining them. They aim to reduce the generalization error of the ensemble compared to the original single model by reducing the variance of the learning algorithm. Let us illustrate the effects of an averaging method called bagging on the bias and variance of fully grown decision trees. The bagging method fits each estimator of the ensemble on a bootstrap copy of the learning set. In Figure 4.2, we show the variance and the bias as function of the number of fully grown decision trees in the bagging ensemble. With a single decision tree, the variance is the dominating error component. By increasing the size of the Bagging ensemble, the variance and the hold out error are reduced, while leaving the bias mostly unchanged.

In Section 4.2.1, we show how the bias-variance decomposition of a randomized supervised learning algorithm is affected by the ensemble averaging method. In Section 4.2.2, we present how to induce randomization without modifying the original supervised learning algorithm. In Section 4.2.3, we describe specific randomization schemes for decision tree methods leading to random forest models.

4.2.1 Variance reduction

Let us first study the bias-variance decomposition for a model $f_{\mathcal{L}, \theta}$ trained on a learning set \mathcal{L} whose randomness is entirely captured by a random vector of parameters θ . The model $f_{\mathcal{L}, \theta}$ is thus a function of two random variables θ and the learning set \mathcal{L} . It admits the

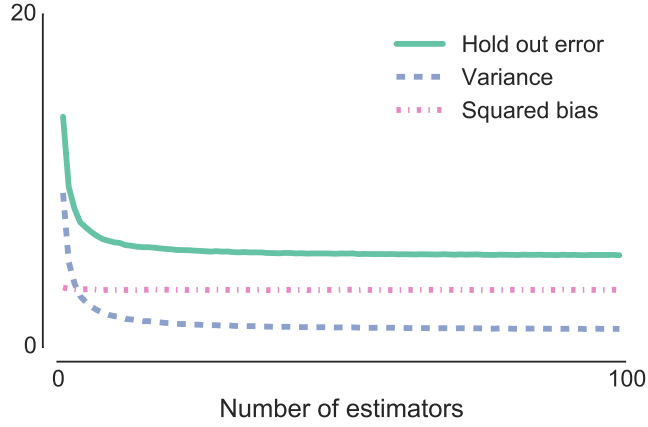


Figure 4.2: Bias-variance decomposition of an ensemble of randomized and fully developed decision tree models fitted on bootstrap copies of the learning set (bagging method) for increasing ensemble size on the Friedman1 dataset. The hold out error, variance and squared bias are averaged over 100 independent learning and testing sets of respective size 300 and 20000 samples.

following bias-variance decomposition of the generalization error for the average square loss (Geurts, 2002):

$$E_{\mathcal{L},\theta}(\text{Err}) = E_{P_x} \left\{ \text{Var}_{\mathcal{L},\theta} \{f_{\mathcal{L},\theta}(x)\} + \text{Bias}^2(f_{\mathcal{L},\theta}(x)) + \text{Var}_{P_{y|x}} \{y\} \right\}, \quad (4.15)$$

where

$$\text{Var}_{\mathcal{L},\theta} \{f_{\mathcal{L},\theta}(x)\} = E_{\mathcal{L},\theta} \left\{ (f_{\mathcal{L},\theta}(x) - E_{\mathcal{L}} E_{\theta} \{f_{\theta}(x)\})^2 \right\}, \quad (4.16)$$

$$\text{Bias}^2(f_{\mathcal{L},\theta}(x)) = (E_{\mathcal{L}} E_{\theta} \{f_{\mathcal{L},\theta}(x)\} - f_{\text{Bayes}}(x))^2. \quad (4.17)$$

By comparison to the bias-variance decomposition of an unperturbed model (see Equation 4.15), we have two main differences:

1. The squared bias is now the distance between the Bayes model f_{Bayes} and the average model $E_{\mathcal{L}} E_{\theta} \{f_{\mathcal{L},\theta}(x)\}$ over both the learning set \mathcal{L} and the randomization parameter θ . Note that the average model of the randomized algorithm is different from the non-randomized model $E_{\mathcal{L}} \{f_{\mathcal{L}}(x)\}$. The randomization of the original algorithm might increase the squared bias.
2. The variance $\text{Var}_{\mathcal{L},\theta} \{f_{\mathcal{L},\theta}(x)\}$ of the algorithm now depends on the two random variables \mathcal{L} and θ . With the law of total variance, we can further decompose the variance term into two terms:

$$\text{Var}_{\mathcal{L},\theta} \{f_{\mathcal{L},\theta}(x)\} = \text{Var}_{\mathcal{L}} \left\{ E_{\theta|\mathcal{L}} \{f_{\mathcal{L},\theta}(x)\} \right\} + E_{\mathcal{L}} \text{Var}_{\theta|\mathcal{L}} \{f_{\mathcal{L},\theta}(x)\} \quad (4.18)$$

The first term is the variance brought by the learning sets of the average model over all parameter vectors θ . The second term describes the variance brought by the parameter vector θ averaged over all learning sets \mathcal{L} .

Now, we can study the bias-variance decomposition of the generalization error for an ensemble model $f_{\mathcal{L},\theta_1,\dots,\theta_m}$ whose constituents $\{f_{\mathcal{L},\theta_m}\}_{m=1}^M$ depend each on the learning set \mathcal{L} and a random parameter vector $\{\theta_m\}_{m=1}^M$ capturing the randomness of the models. The bias-variance decomposition of the ensemble model $f_{\theta_1,\dots,\theta_m}$ is given by

$$\begin{aligned} E_{\mathcal{L},\theta_1,\dots,\theta_M}(\text{Err}) &= E_{P_x} \text{Var}_{\mathcal{L}} \{E_{\theta_1,\dots,\theta_M|\mathcal{L}} \{f_{\mathcal{L},\theta_1,\dots,\theta_M}(x)\}\} \\ &\quad + E_{P_x} E_{\mathcal{L}} \text{Var}_{\theta_1,\dots,\theta_M|\mathcal{L}} \{f_{\mathcal{L},\theta_1,\dots,\theta_M}(x)\} \\ &\quad + E_{P_x} \text{Bias}^2(f_{\mathcal{L},\theta_1,\dots,\theta_M}(x)) \\ &\quad + E_{P_x} \text{Var}_{P_{y|x}} \{y\}. \end{aligned} \quad (4.19)$$

Let us compare the decomposition for a single random model (Equations 4.15-4.18) to the decomposition for an ensemble of random models (Equation 4.19). We are going to expand the bias-variance decomposition using the ensemble prediction formula $f_{\theta_1,\dots,\theta_m}(x) = \frac{1}{M} \sum_{m=1}^M f_{\theta_m}(x)$ (the demonstration follows (Geurts, 2002)). As previously, the variance $\text{Var}_{P_{y|x}}(y)$ is irreducible as this term does not depend on the supervised learning model.

The average ensemble model of an ensemble of randomized models is equal to the average model of a single model $f_{\mathcal{L},\theta}$ of random parameter vector θ :

$$E_{\mathcal{L},\theta_1,\dots,\theta_M} \{f_{\mathcal{L},\theta_1,\dots,\theta_M}(x)\} = \frac{1}{M} \sum_{m=1}^M E_{\mathcal{L}} E_{\theta_m} \{f_{\mathcal{L},\theta_m}(x)\}, \quad (4.20)$$

$$= E_{\mathcal{L}} E_{\theta} \{f_{\mathcal{L},\theta}(x)\}. \quad (4.21)$$

The squared bias of the ensemble is thus unchanged compared to a single randomized model.

Now, let us consider the two variance terms. The first one depends on the variability of the learning set \mathcal{L} . With an ensemble of randomized models, it becomes:

$$\begin{aligned} &E_{P_x} \text{Var}_{\mathcal{L}} \{E_{\theta_1,\dots,\theta_M|\mathcal{L}} \{f_{\mathcal{L},\theta_1,\dots,\theta_M}(x)\}\} \\ &= E_{P_x} \text{Var}_{\mathcal{L}} \left\{ \frac{1}{M} \sum_{m=1}^M E_{\theta_m|\mathcal{L}} \{f_{\mathcal{L},\theta_m}(x)\} \right\}, \end{aligned} \quad (4.22)$$

$$= E_{P_x} \text{Var}_{\mathcal{L}} \{E_{\theta|\mathcal{L}} \{f_{\mathcal{L},\theta}(x)\}\}. \quad (4.23)$$

The variance of the ensemble of randomized model with respect to the learning set \mathcal{L} drawn from the input-output pair distribution $P_{x,y}$ is not affected by the averaging and is equal to the variance of a single randomized model.

Let us developed the second variance term of the decomposition describing the variance with respect to the set of random parameter vectors $\theta_1, \dots, \theta_M$:

$$\begin{aligned} & E_{P_x} E_{\mathcal{L}} \text{Var}_{\theta_1, \dots, \theta_m | \mathcal{L}} \{f_{\mathcal{L}, \theta_1, \dots, \theta_M}(x)\} \\ &= E_{P_x} E_{\mathcal{L}} \text{Var}_{\theta_1, \dots, \theta_M | \mathcal{L}} \left\{ \frac{1}{M} \sum_{m=1}^M f_{\mathcal{L}, \theta_m}(x) \right\}, \end{aligned} \quad (4.24)$$

$$= \frac{1}{M^2} E_{P_x} E_{\mathcal{L}} \left\{ \sum_{m=1}^M \text{Var}_{\theta_1, \dots, \theta_M | \mathcal{L}} \{f_{\mathcal{L}, \theta_m}(x)\} \right\}, \quad (4.25)$$

$$= \frac{1}{M^2} E_{P_x} E_{\mathcal{L}} \left\{ \sum_{m=1}^M \text{Var}_{\theta_m | \mathcal{L}} \{f_{\mathcal{L}, \theta_m}(x)\} \right\}, \quad (4.26)$$

$$= \frac{1}{M} E_{P_x} E_{\mathcal{L}} \text{Var}_{\theta | \mathcal{L}} \{f_{\mathcal{L}, \theta}(x)\}, \quad (4.27)$$

where we use the following properties: (i) $\text{Var}\{\alpha x\} = \alpha^2 \text{Var}\{x\}$ where α is a constant, (ii) at a fixed learning set \mathcal{L} , the models $f_{\mathcal{L}, \theta_m} \forall m$ are independent, (iii) the variance of a sum of independent random variables is equal to the sum of the variance of each independent random variables ($\text{Var}\{\sum_{i=1}^n x_i\} = \sum_{i=1}^n \text{Var}\{x_i\}$).

Putting all together the bias-variance decomposition of Equation 4.19 becomes (Geurts, 2002):

$$\begin{aligned} E_{\mathcal{L}, \theta_1, \dots, \theta_m}(\text{Err}) &= E_{P_x} \text{Var}_{\mathcal{L}} \{E_{\theta | \mathcal{L}} \{f_{\mathcal{L}, \theta}(x)\}\} \\ &\quad + \frac{1}{M} E_{P_x} E_{\mathcal{L}} \text{Var}_{\theta | \mathcal{L}} \{f_{\mathcal{L}, \theta}(x)\} \\ &\quad + E_{P_x} (E_{\mathcal{L}} E_{\theta} \{f_{\mathcal{L}, \theta}(x)\} - f_{\text{Bayes}}(x))^2 \\ &\quad + E_{P_x} \text{Var}_{P_{y|x}} \{y\}. \end{aligned} \quad (4.28)$$

The bias variance decomposition of an ensemble of randomized models $f_{\mathcal{L}, \theta_1, \dots, \theta_M}$ (see Equation 4.28) shows that averaging M models reduces the variance related to the randomization θ by a factor $1/M$ over a single randomized model $f_{\mathcal{L}, \theta}$ (see Equation 4.15) without modifying the other terms. Note that we can not compare the bias variance decomposition of an ensemble of randomized models $f_{\mathcal{L}, \theta_1, \dots, \theta_M}$ to its non randomized counterparts $f_{\mathcal{L}}$ (see Equation 4.5). The bias and variance terms are indeed not comparable.

In practice, we first perturb the learning algorithm which increases the variance of the models and then we combine them through averaging. The variance reduction effect is expected to be higher than the added variance at training time. The bias is either unaffected or increased through the randomization induction. Perturbing the algorithm is thus a tradeoff between the reduction in variance and the increase in bias. An ensemble of randomized models $f_{\mathcal{L}, \theta_1, \dots, \theta_M}$ will have better performance than its non perturbed counterparts $f_{\mathcal{L}}$ if the bias increase is compensated by the variance reduction. In (Louppe,

2014), the authors have shown that the generalization error is reduced if the randomization induction decorrelates the models of the ensemble.

The previous decomposition does not apply to the 0 – 1 loss in classification. However, the main conclusions remains valid (Breiman, 1996a; Domingos, 2000; Geurts, 2002).

4.2.2 Generic randomization induction methods

In this section, we first discuss generic randomization methods to perturb a supervised learning algorithm without modifying the original algorithm through (i) the learning sample set $\mathcal{L} = \{(x^i, y^i) \in \mathcal{X} \times \mathcal{Y}\}_{i=1}^n$ available at fitting time, (ii) the input space \mathcal{X} or (iii) the output space \mathcal{Y} . Those three perturbation principles can be either applied separately or together.

We present in succession these three model agnostic randomization principles (perturbing \mathcal{L} , \mathcal{X} or \mathcal{Y}).

4.2.2.1 Sampling-based randomization

One of the earliest randomization method, called bagging (Breiman, 1996a), fits independent models on bootstrap copies of the learning set. A bootstrap copy (Efron, 1979) is obtained by sampling with replacement $|\mathcal{L}|$ samples from the learning set \mathcal{L} . The original motivation was a first theoretical development and empirical experiments showing that bagging reduces the error of an unstable estimator such as a decision tree. Bootstrap sampling totally ignores the class distribution in the original sample set \mathcal{L} and might lead to highly unbalanced bootstraps. A partial solution is to use stratified bootstraps or to bootstrap (Chen et al., 2004) separately the minority and majority classes. In the bagging approach only a fraction of the dataset is provided as training set to each estimator, the wagging approach (Bauer and Kohavi, 1999) fits instead each estimator on the entire training set with random weights. Instead of using bootstrap copies of training set, Büchlmann and Yu (2002) proposes to subsample the training set, i.e. to sample without replacement the training set.

To take into account the input space structure, Kuncheva et al. (2007) proposes to generate random input space partition with a random hyperplane. An estimator is then build for each partition. To get an ensemble, Kuncheva et al. (2007) repeats this process multiple times.

4.2.2.2 Input-based randomization

Input-based randomization techniques are often based on dimensionality reduction techniques. The random subspace method (Ho, 1998) builds each model on a random subset of the input space \mathcal{X} obtained

by sub-sampling inputs without replacement. It was later combined with the bagging method in (Panov and Džeroski, 2007), by bootstrapping the learning set \mathcal{L} before learning an estimator, and with sub-sampling techniques with/without replacement in (Louppe and Geurts, 2012) generating random (sample-input) patches of the data. Note that while we reduce the input space size, we can also over-sample the learning set. For instance, Maree et al. (2005) apply a supervised learning algorithm on random sub-windows extracted from an image, which effectively (i) increases the sample size available to train each model; (ii) reduces the input space size, (iii) and takes into account spatial (pixel) correlation in the images.

Since decision tree made their split orthogonally to the input space, authors have proposed to randomize such ensembles by randomly projecting the input space. Rotation forest (Rodríguez et al., 2006) is an ensemble method combining bagging with principal component analysis. For each bootstrap copy, it first slices the p input variables into q subsets, then projects each subset of inputs of size $\frac{p}{q}$ on its principal components and finally grows q models (one on each subset). Kuncheva and Rodríguez (2007) further compares three input dimensionality reduction techniques (described in Section 2.6): (i) the PCA approach of Rodríguez et al. (2006), (ii) Gaussian random projections and (iii) sparse Gaussian random projections. On their benchmark, they find that the PCA-based rotation matrices yield the best results and also that sparse random projections are strictly better than dense random projections. The idea of using dense Rademacher or Gaussian random projections was again re-discovered by Schclar and Rokach (2009). Similarly, Blaser and Fryzlewicz (2015) proposed to make ensembles through random rotation of the input space.

4.2.2.3 Output-based randomization

Output-based randomization methods directly perturb the output space \mathcal{Y} of each member of the ensemble.

In regression, we can induce randomization to an output variable through the addition of an independent Gaussian noise (Breiman, 2000). We fit each model of the ensemble on the perturbed output $y' = y + \epsilon_m$ with $\epsilon_m \sim \mathcal{N}(0; \sigma)$.

In classification, we perturb the output of each model of the ensemble by having a non zero probability to randomly switch the class associated to each sample (Breiman, 2000; Martínez-Muñoz and Suárez, 2005; Martínez-Muñoz et al., 2008).

For multi-label tasks and multi-output tasks, supervised learning algorithms, such as random k -label subset (see also Section 2.2.5) (Tsoumakas and Vlahavas, 2007), randomizes the ensemble by building each model of the ensemble on a subset of the output space or the label sets present in the learning set.

4.2.3 Randomized forest model

The decision tree algorithm has a high variance, due to the instability of its structure. Large decision trees, such as fully developed trees, are often very unstable, especially at the bottom of the tree. The selected splitting rules depend on the samples reaching those nodes. Small changes in the learning set might lead to very different tree structures. Authors have proposed randomization schemes to perturb the search and selection of the best splitting rule improving the generalization error through averaging methods.

One of the first propositions to perturb the splitting rule search (Dietterich and Kong, 1995) was to select randomly at each node one splitting rule among the top k splitting rules with the highest impurity reduction. The variance of the algorithm increases with the number k of splitting rule candidates, leaving the bias unchanged. Later in the context of digit recognition, Amit et al. (1997) randomized the tree growth by restricting the splitting rule search at each node to a random subset of k input variables out of the p available. The original motivation was to drastically reduce the splitting rule search space as the number of input variables is very high in digit recognition tasks. This randomization scheme increases more the variance of the algorithm than the one of Dietterich and Kong, at the expense of increasing the bias. Note the similarity with the random subspace approach (Ho, 1998) which subsamples the input space prior fitting a new estimator in an averaging ensemble.

Breiman got inspired by the work of Amit et al. and combined its bagging method (Breiman, 1996a) with the input variable sub-sampling leading to the well known¹ “random forest” algorithm (Breiman, 2001). The combination of both randomization schemes has led to one of the best of the shelf estimator for many supervised learning tasks (Caruana et al., 2008; Fernández-Delgado et al., 2014).

Later on, Geurts et al. (2006a) randomized the cut point and input variable selection of the splitting rules. At each test node, it draws one random splitting rule for k randomly selected input variables (without replacement) and then selects the best one. This randomized tree ensemble is called extremely randomized trees or extra trees. For a splitting rule $s(x) = \mathbb{1}_{x_j \leq \tau}$ associated to an ordered variable, the algorithm draws uniformly at random the threshold τ between the minimum and maximum of the possible cut point values. Similarly for an unordered variable, the algorithm draws a non empty subset B among the possible values to generate a splitting rule of the form $s(x) = \mathbb{1}_{x_j \in B}$. Empirically, it has been shown (Geurts, 2002) that the variance of the decision tree algorithm is due to the variability of

¹ The random forest method usually refers to the the algorithm of Breiman, however any averaging ensemble of randomized trees is also a random forest.

the cut point selection with respect to the learning set. We can view the perturbation of the cut point selection as a way to transfer the variance due to the learning set to the variance due to the randomization of the cut point selection. The hyper-parameter k controls the trade-off between the bias and variance of the algorithm.

Besides perturbing the binary and axis wise splitting rules, they have been some research to make splitting rule through random hyper-planes. [Breiman \(2001\)](#) proposed to select the best splitting rule obtained from random sparse linear input combinations with non zero values drawn uniformly in $[-1, 1]$. [Tomita et al. \(2015\)](#) proposed to use sparse random projection where non zero elements are drawn uniformly in $[-1, 1]$. Those approaches increase the variance, while also trying to reduce the bias by allowing random oblique splits. However, [Menze et al. \(2011\)](#) have shown that those random sparse hyper-planes are inferior to deterministic linear models such as ridge regressors or a linear discriminant analysis (LDA) models.

For supervised learning tasks with many outputs, we can also perturb the output space of each decision tree by randomly projecting the output space ([Joly et al., 2014](#)) onto a lower dimensional subspace or through random output sub-sampling. The leaves are later re-labelled on the original output space. This approach is developed in Chapter 5 of this thesis.

4.3 BOOSTING ENSEMBLES

Boosting methods originate from the following question: “How can we combine a set of weak models together, each one doing slightly better than random guessing, so as to get one stronger model having good generalization performance?”. A boosting model f answers this question through a weighted combination of M weak models $\{f_m\}_{m=1}^M$ leading to

$$f(x) = \sum_{m=1}^M \alpha_m f_m(x) \quad (4.29)$$

where the coefficients $\{\alpha_m \in \mathbb{R}\}_{m=1}^M$ highlight the contribution of each model $f_m(x)$ to the ensemble.

For a boosting ensemble, we usually want to minimize a loss function $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$ over a learning set $\mathcal{L} = \{(x^i, y^i) \in \mathcal{X} \times \mathcal{Y}\}_{i=1}^n$:

$$\min_{\{(\alpha_m, f_m) \in \mathbb{R} \times \mathcal{H}\}_{m=1}^M} \sum_{(x, y) \in \mathcal{L}} \ell \left(y, \sum_{m=1}^M \alpha_m f_m(x) \right) \quad (4.30)$$

where we select each model f_m over a hypothesis space \mathcal{H} . Solving this equation for many loss functions and models is either intractable

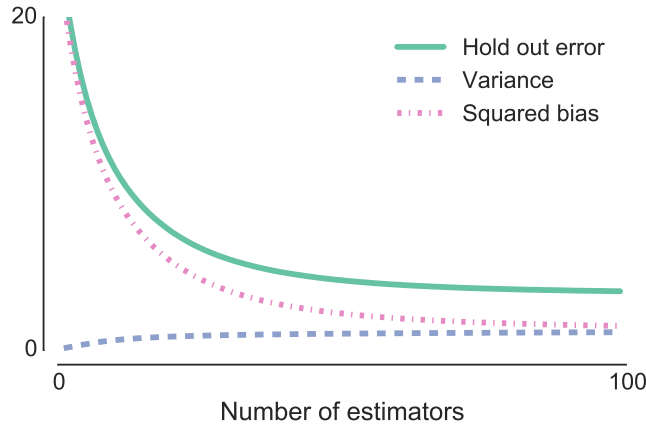


Figure 4.3: Bias-variance decomposition of a boosting ensemble minimizing the square loss with an increasing number of weak models on the Friedman1 dataset. The hold out error, variance and squared bias are averaged over 100 independent learning and testing sets of respective size 300 and 20000 samples.

or numerically too intensive for practical purpose. However, we can solve easily Equation 4.30 for a single model ($M = 1$).

So boosting methods develop iterative and tractable schemes to solve Equation 4.30 by adding sequentially models to the ensemble. A new model $f_m(x)$ builds over the work done by the previous $m - 1$ models to yield better predictions. It further minimizes the loss ℓ averaged over the training data:

$$\min_{\alpha_m, f_m \in \mathbb{R} \times \mathcal{H}} \sum_{(x, y) \in \mathcal{L}} \ell \left(y, \sum_{l=1}^{m-1} \alpha_l f_l(x) + \alpha_m f_m(x) \right). \quad (4.31)$$

To improve the predictions made by the $m - 1$ models, the new model f_m with coefficient α_m concentrates its efforts on the wrongly predicted samples.

From a bias-variance perspective, each newly added model aims to reduce the bias while leaving the variance term unmodified if possible. We choose the base model so that it has a high bias and a small variance such as a stump, i.e. a decision tree with only one testing node, or such as a linear model with only one non-zero coefficient. In Figure 4.3, we sequentially fit stumps to decrease the least square loss $\ell(y, y') = \frac{1}{2}(y - y')^2$. With a few stumps, the squared bias component of the generalization error dominates with a low variance. By adding more stumps to the ensemble, we drastically decrease the generalization error by diminishing the squared bias. The best performance is a trade-off between the bias reduction and the increase in variance.

We present the adaptive boosting and its variants in Section 4.3.1, which directly solve Equation 4.31, and the functional gradient boosting approach in Section 4.3.2, which approximately solves Equation 4.31 through the loss gradient.

4.3.1 Adaboost and variants

One of the most popular and influential boosting algorithms is the “AdaBoost” algorithm (Freund and Schapire, 1997). This supervised learning algorithm aims to solve binary classification tasks with $\mathcal{Y} = \{-1, 1\}$. The algorithm generates iteratively an ensemble of estimators $\{f_m\}_{m=1}^M$ by minimizing the exponential loss function:

$$\ell_{\text{exp}}(y, y') = \exp(-yy'), \quad (4.32)$$

assuming a binary response of the weak models $f_m(x) \in \{-1, 1\} \forall m$.

The prediction of an unseen sample $f(x)$ by an AdaBoost ensemble is a majority vote from its members:

$$f(x) = \text{sign} \left(\sum_{m=1}^M \alpha_m f_m(x) \right), \quad (4.33)$$

where the $\{\alpha_m\}_{m=1}^M$ are constant weights indicating the contribution of a model f_m to solve the binary classification task. The sign operator transforms the sum into an appropriate output value ($\mathcal{Y} = \{-1, 1\}$).

Given a learning set $\mathcal{L} = \{(x^i, y^i) \in \mathcal{X} \times \{-1, 1\}\}_{i=1}^n$, we iteratively fit a weak model f_m over the learning set \mathcal{L} by making the weak learner focuses on each sample with a weight $(w^i)_{i=1}^n$. The higher the value of w^i , the more the algorithm will concentrate to predict correctly the i -th sample. To design this algorithm, we need to answer to the following questions: (i) how to assess the contribution α_m of the m -th model f_m to the ensemble and (ii) how to update the weight w^i to reduce iteratively the exponential loss.

We can write the resubstitution error of the exponential loss as

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \ell_{\text{exp}} \left(y^i, \sum_{l=1}^{m-1} \alpha_l f_l(x^i) + \alpha_m f_m(x^i) \right) \\ &= \frac{1}{n} \sum_{i=1}^n \exp \left(-y^i \left(\sum_{l=1}^{m-1} \alpha_l f_l(x^i) + \alpha_m f_m(x^i) \right) \right) \end{aligned} \quad (4.34)$$

$$= \frac{1}{n} \sum_{i=1}^n \exp \left(-y^i \left(\sum_{l=1}^{m-1} \alpha_l f_l(x^i) \right) \right) \exp(-y^i \alpha_m f_m(x^i)) \quad (4.35)$$

$$= \frac{1}{n} \sum_{i=1}^n \ell_{\text{exp}} \left(y, \sum_{l=1}^{m-1} \alpha_l f_l(x^i) \right) \exp(-y^i \alpha_m f_m(x^i)) \quad (4.36)$$

$$= \sum_{i=1}^n w^i \exp(-y^i \alpha_m f_m(x^i)) \quad (4.37)$$

with

$$w^i = \frac{1}{n} \ell_{\text{exp}} \left(y, \sum_{l=1}^{m-1} \alpha_l f_l(x^i) \right) \forall i. \quad (4.38)$$

Note that the weight computation is expressible as a recursive equation starting with $w^i = 1/n$:

$$w^i \leftarrow w^i \exp(\alpha_m 1(y^i \neq f_m(x^i))). \quad (4.39)$$

The sample weight w^i highlights how well the i -th sample is predicted by the $m - 1$ first estimators of the boosting ensemble. A zero weight w^i means that the i -th sample is perfectly predicted. The m -th estimators should thus focus on the sample with high weight w^i to reduce the resubstitution error. Otherwise, it should minimize the weighted resubstitution error.

Let us now separate the resubstitution error of the correctly classified points from the misclassified ones:

$$\begin{aligned} & \sum_{i=1}^n w^i \exp(-y^i \alpha_m f_m(x^i)) \\ &= \sum_{i=1}^n w^i \exp(-\alpha_m) 1(y^i = f_m(x^i)) + \exp(\alpha_m) 1(y^i \neq f_m(x^i)) \end{aligned} \quad (4.40)$$

By derivating the last equation with respect to α_m and setting the derivative to zero, the α_m minimizing the resubstitution error of the exponential loss is

$$\alpha_m = \log\left(\frac{1 - \text{err}_m}{\text{err}_m}\right) \quad (4.41)$$

with

$$\text{err}_m = \frac{1}{2} \frac{\sum_{i=1}^n w^i 1(y^i \neq f_m(x^i))}{\sum_{i=1}^n w^i}. \quad (4.42)$$

The optimization of the constant α_m means that the resubstitution error is upper bounded and can not increase with the size of the ensemble on the learning set.

Putting everything together, we obtain the AdaBoost algorithm (see Algorithm 4.1). Many extensions and enhancements of this fundamental idea have been proposed. If the weak model is able to predict a probability estimate, [Friedman et al. \(2000\)](#) have proposed an appropriate extension called “Real Adaboost” by contrast to Algorithm 4.1 which they call “Discrete Adaboost”.

A direct multi-class extension, called AdaBoost.M1, of the Adaboost algorithm is to use a multi-class weak learner instead of a binary one. The AdaBoost.M1 ensemble predicts a new sample through:

$$f(x) = \arg \max_{k \in \mathcal{Y}} \sum_{m=1}^M \alpha_m 1(f_m(x) = k). \quad (4.43)$$

Algorithm 4.1 AdaBoost.M1 for binary classification $\mathcal{Y} = \{-1, 1\}$.

- 1: **function** ADABOOST($\mathcal{L} = \{x^i, y^i \in \mathcal{X} \times \mathcal{Y}\}_{i=1}^n$)
- 2: Initialize the sample weights $w^i \leftarrow 1/n \forall i \in \{1, \dots, n\}$.
- 3: **for** $m = 1$ to M **do**
- 4: Fit a model $f_m(x)$ to the learning set \mathcal{L} and $(w^i)_{i=1}^n$.
- 5: Compute the weighted error rate

$$\text{err}_m \leftarrow \frac{\sum_{i=1}^n w^i \mathbf{1}(y^i \neq f_m(x^i))}{\sum_{i=1}^n w^i}.$$

- 6: Compute $\alpha_m \leftarrow \frac{1}{2} \log \left(\frac{1 - \text{err}_m}{\text{err}_m} \right)$.
- 7: Update the weights

$$w^i \leftarrow w^i \exp(\alpha_m \mathbf{1}(y^i \neq f_m(x^i))).$$

- 8: **end for**
 - 9: **return** $\hat{f}(x) = \text{sign} \left(\sum_{m=1}^M \alpha_m f_m(x) \right)$
 - 10: **end function**
-

An improvement over this approach is to directly minimize the multi-class exponential loss as in the SAMME algorithm (Zhu et al., 2009). It replaces the line 6 of Algorithm 4.1 by

$$\alpha_m \leftarrow \frac{1}{2} \log \left(\frac{1 - \text{err}_m}{\text{err}_m} \right) + \log(|\mathcal{Y}| - 1) \quad (4.44)$$

We can minimize other losses than the exponential loss during the ensemble growth, such as the logistic loss $\ell_{\text{logistic}}(y, y') = \log(1 + \exp(-2yy'))$ with the LogitBoost algorithm (Collins et al., 2002) for binary classification tasks; the Hamming loss $\ell_{\text{Hamming}}(y, y') = \sum_{j=1}^n \mathbf{1}(y_j \neq y'_j)$ leading to the AdaBoost.MH algorithm (Schapire and Singer, 2000) and the pairwise ranking loss ℓ_{ranking}

$$\ell_{\text{ranking}}(y, y') = \frac{1}{|y^i|} \frac{1}{d - |y^i|} \left| \{(k, l) : y_k^i < y_l^i, y_k^i = 1, y_l^i = 0\} \right|$$

leading to the AdaBoost.MR algorithm (Schapire and Singer, 2000) for multi-label classification tasks and also a wide range of regression losses as proposed in (Drucker, 1997) for regression tasks.

HOW TO TAKE INTO ACCOUNT SAMPLE WEIGHTS IN SUPERVISED LEARNING ALGORITHMS?

For a set of learning samples $((x^i, y^i) \in (\mathcal{X} \times \mathcal{Y}))_{i=1}^n$ and a set of weights $(w^i \in \mathbb{R}^+)_{i=1}^n$, the weighted resubstitution error is given by

$$\text{Resubstitution error} = \frac{\sum_{i=1}^n w^i \ell(f(x^i), y^i)}{\sum_{i=1}^n w^i}. \quad (4.45)$$

Extending supervised learning algorithms to support sample weights means that we have to modify the learning algorithm so as to minimize the weighted resubstitution error:

- For linear models, we will minimize the weighted average of a given loss $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$ over the learning set $\mathcal{L} = \{(x^i, y^i) \in \mathcal{X} \times \mathcal{Y}\}_{i=1}^n$

$$\min_{\beta_0, \beta} \sum_{i=1}^n w^i L(y^i, \beta_0 + \beta^\top x), \quad (4.46)$$

where the $w^i \in \mathbb{R}^+$ are the weight associated to each sample. There is an analytical solution in the case of a ℓ_2 norm regularization penalty and algorithms for a ℓ_1 regularization penalty can be easily extended to accommodate for the weights.

- For a decision tree, we will use a weighted impurity criterion and a weight-aware leaf labelling rule assignment procedure. It also allows new stopping rule based on sample-weight, such as a minimal total weight to split a node.
- For a k-nearest neighbors, we will store the sample weight during fit and we will predict an unseen sample through a weighted aggregation of the nearest neighbors.

Conversely, to support unweighted supervised learning task with a weight-aware implementation, we can set the sample weights to a constant such as $w^i = 1/n \forall i$ prior the model training.

4.3.2 Functional gradient boosting

The AdaBoost algorithm has an analytical and closed-form solution to Equation 4.31 with the exponential loss. However, we would like to build boosting ensembles when such closed-form solutions are not available. Functional gradient boosting, a forward stagewise additive approach, approximately solves Equation 4.31 for a given loss $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$ by sequentially adding new basis function f_m , a regression model, with a weight α_m without modifying the previous models.

If we want to add a new model f_m to a boosting ensemble with $m - 1$ models while minimizing the square loss, the loss for a sample (x, y) is given by

$$\ell(y, f(x)) = \frac{1}{2}(y - f(x))^2 \quad (4.47)$$

$$= \frac{1}{2}\left(y - \sum_{l=1}^{m-1} \alpha_l f_l(x) - \alpha_m f_m(x)\right)^2 \quad (4.48)$$

$$= \frac{1}{2} (r_m(x) - \alpha_m f_m(x))^2 \quad (4.49)$$

where $r_m(x) = y - \sum_{l=1}^{m-1} \alpha_l f_l(x)$ is the remaining residual of the $m-1$ models to predict a sample x . Thus for the square loss, we can add a new models f_m by fitting the new model on the residuals left by the $m-1$ previous models. This approach is called least square regression boosting.

Solving Equation 4.31 is difficult for general loss functions. It requires to be able to expand a new basis function f_m while minimizing the chosen loss function. For instance in the context of decision trees, it would require a specific splitting criterion and a leaf labelling rule minimizing the chosen loss.

Instead of solving Equation 4.31, Friedman (2001) proposed a fast approximate solution for arbitrary differentiable losses inspired from numerical optimization. We can re-write the loss function minimization as

$$\hat{f} = \arg \min_f \ell(f) = \min_f \sum_{(x,y) \in \mathcal{L}} \ell(y, f(x)). \quad (4.50)$$

with the constraint that f is a sum of supervised learning models. Ignoring this constraint, the Equation 4.50 is an unconstrained minimization problem with $f \in \mathbb{R}^n$ being a n -dimensional vector. Iterative solvers solve such minimization problems by correcting an initial estimate through a recursive equation. The final solution is a sum of vectors

$$\hat{f} = \sum_{m=0}^M h_m, \quad h_m \in \mathbb{R}^n, \quad (4.51)$$

where h_0 is the initial estimate. The construction of the sequence of h_0, \dots, h_m depends on the chosen optimization algorithm.

The gradient boosting algorithm (Friedman, 2001) uses the same approach as the gradient descent method. The update rule of the gradient descent algorithm h_m is of the form

$$h_m = -\rho_m g_m \quad (4.52)$$

where ρ_m is a scalar and $g_m \in \mathbb{R}^n$ is the gradient of $L(f)$ with respect to f evaluated at the current approximate solution $\hat{f} = \sum_{l=1}^{m-1} \rho_l h_l$:

$$g_m^i = \left[\frac{\partial}{\partial y'} \ell(y^i, y') \right]_{y'=\hat{f}}. \quad (4.53)$$

The scalar ρ_m is the step length in the negative loss gradient direction $-g_m$ chosen so as to minimize the objective function $\ell(f)$:

$$\rho_m = \arg \min_{\rho \in \mathbb{R}} \ell(\hat{f} - \rho g_m). \quad (4.54)$$

Back to supervised learning, we can only compute the loss gradient for the training samples. To generalize to unseen data, the idea is to approximate the direction of the negative gradient using a regression model g_m selected within a hypothesis space \mathcal{H} of weak base-learners minimizing the square loss on the training data:

$$g_m = \arg \min_{g \in \mathcal{H}} \sum_{i=1}^n (-g_m^i - g(x^i))^2. \quad (4.55)$$

The gradient boosting approach can be summarized as follows: start at an initial constant estimate $\rho_0 \in \mathbb{R}$, then iteratively follows the negative gradient of the loss ℓ as estimated by a regression model g_m fitted over the training samples and make an optimal step length ρ_m minimizing the loss ℓ . The gradient boosting ensemble predicts an unseen sample through

$$f(x) = \rho_0 + \sum_{m=1}^M \rho_m g_m(x). \quad (4.56)$$

The whole procedure is given in Algorithm 4.2. The algorithm is completely defined once we have (i) a starting model, usually the constant minimizing the chosen loss (line 4) and (ii) the gradient of the loss (line 2).

We compute the optimal step length (line 6 of Algorithm 4.2) either analytically as for the square loss or numerically using, e.g., the Brent's method (Brent, 2013), a robust root-finding method allowing to minimize single unconstrained optimization problem, as for the logistic loss. Friedman (2001) advises to use one step of the Newton–Raphson method. However, the Newton–Raphson algorithm might not converge if the first and second derivative of the loss are small. These conditions occurs frequently in highly imbalanced supervised learning tasks.

A learning rate $\mu \in (0, 1]$ is often added to shrink the size of the gradient step ρ_m in the residual space in order to avoid overfitting the the training samples. Another possible modification is to induce randomization, e.g. by subsampling without replacement the samples available (from all learning samples) at each iteration (Friedman, 2002).

Table 4.1 gives an overview of regression and classification losses with their gradients, while Table 4.2 gives the starting constant models minimizing losses. The square loss in regression and the exponential loss in classification leads to nice gradient boosting algorithm (respectively the least square regression boosting algorithm and the exponential classification boosting algorithm (Zhu et al., 2009)). However, these losses are not robust to outlier. More robust losses can be used such as the absolute loss in regression and the logistic loss or the hinge loss in classification.

Algorithm 4.2 Gradient boosting algorithm

```

1: function GRADIENTBOOSTING( $\mathcal{L} = \{(x^i, y^i) \in \mathcal{X} \times \mathcal{Y}\}_{i=1}^n; \ell; \mathcal{H}; M$ )
2:    $f_0(x) = \rho_0 = \arg \min_{\rho \in \mathbb{R}} \sum_{i=1}^n \ell(y^i, \rho)$ .
3:   for  $m = 1$  to  $M$  do
4:     Compute the loss gradient for the training set points

           
$$g_m^i = \left[ \frac{\partial}{\partial y'} \ell(y^i, y') \right]_{y' = f_{m-1}(x)} \quad \forall i \in \{1, \dots, n\}.$$


5:     Find a correlated direction to the loss gradient

           
$$g_m = \arg \min_{g \in \mathcal{H}} \sum_{i=1}^n (-g_m^i - g(x^i))^2.$$


6:     Find an optimal step length in the direction  $g_m$ 

           
$$\rho_m = \arg \min_{\rho \in \mathbb{R}} \sum_{i=1}^n \ell(y^i, f_{m-1}(x^i) + \rho g_m(x^i)).$$


7:      $f_m(x) = f_{m-1}(x) + \mu \rho_m g_m(x)$ .
8:   end for
9:   return  $f_M(x)$ 
10: end function

```

Table 4.1: Regression loss ($\mathcal{Y} = \mathbb{R}$) and binary classification loss ($\mathcal{Y} = \{-1, 1\}$) their derivative with respect to a basis function $f(x)$.

Regression	$\ell(y, y')$	$-\partial \ell(y, y') / \partial y'$
Square	$\frac{1}{2}(y - y')^2$	$y - y'$
Absolute	$ y - y' $	$\text{sign}(y - y')$
Classification	$\ell(y, y')$	$-\partial \ell(y, y') / \partial y'$
Exponential	$\exp(-yy')$	$-y \exp(-yy')$
Logistic	$\log(1 + \exp(-2yy'))$	$\frac{2y}{1 + \exp(2yy')}$
Hinge	$\max(0, 1 - yy')$	$-y \mathbb{1}(yy' < 1)$

Table 4.2: Constant minimizers of regression losses ($\mathcal{Y} = \mathbb{R}$) and binary classification losses ($\mathcal{Y} = \{-1, 1\}$) given a set of samples $\mathcal{L} = \{(x^i, y^i) \in \mathcal{X} \times \mathcal{Y}\}_{i=1}^n$.

Regression	
Square	$f_0(x) = \frac{1}{n} \sum_{i=1}^n y^i$
Absolute	$f_0(x) = \text{median}(\{y^i\}_{i=1}^n)$
Classification	
Exponential	$f_0(x) = \log \left(\frac{\sum_{i=1}^n 1(y^i=1)}{\sum_{i=1}^n 1(y^i=-1)} \right)$
Logistic	$f_0(x) = \log \left(\frac{\sum_{i=1}^n 1(y^i=1)}{\sum_{i=1}^n 1(y^i=-1)} \right)$
Hinge	$f_0(x) = \text{sign} \left(\frac{1}{n} \sum_{i=1}^n 1(y^i = 1) - \frac{1}{2} \right)$

Part II

LEARNING IN COMPRESSED SPACE THROUGH RANDOM PROJECTIONS

RANDOM FORESTS WITH RANDOM PROJECTIONS OF THE OUTPUT SPACE FOR HIGH DIMENSIONAL MULTI-LABEL CLASSIFICATION

OUTLINE

We adapt the idea of random projections applied to the output space, so as to enhance tree-based ensemble methods in the context of multi-label classification. We show how learning time complexity can be reduced without affecting computational complexity and accuracy of predictions. We also show that random output space projections may be used in order to reach different bias-variance tradeoffs, over a broad panel of benchmark problems, and that this may lead to improved accuracy while reducing significantly the computational burden of the learning stage.

This chapter is based on previous work published in

Arnaud Joly, Pierre Geurts, and Louis Wehenkel. Random forests with random projections of the output space for high dimensional multi-label classification. In *Machine Learning and Knowledge Discovery in Databases*, pages 607–622. Springer Berlin Heidelberg, 2014.

Within supervised learning, the goal of multi-label classification is to train models to annotate objects with a subset of labels taken from a set of candidate labels. Typical applications include the determination of topics addressed in a text document, the identification of object categories present within an image, or the prediction of biological properties of a gene. In many applications, the number of candidate labels may be very large, ranging from hundreds to hundreds of thousands (Agrawal et al., 2013) and often even exceeding the sample size (Dekel and Shamir, 2010). The very large scale nature of the output space in such problems poses both statistical and computational challenges that need to be specifically addressed.

A simple approach to multi-label classification problems, called binary relevance, is to train independently a binary classifier for each label. Several more complex schemes have however been proposed to take into account the dependencies between the labels (see Section 2.2.5). In the context of tree-based methods, one way is to train multi-output trees (see Section 3.5), i.e. trees that can predict multiple outputs at once. With respect to binary relevance, the multi-output tree approach has the advantage of building a single model for all

labels. It can thus potentially take into account label dependencies and reduce memory requirements for the storage of the models. An extensive experimental comparison (Madjarov et al., 2012) shows that this approach compares favorably with other approaches, including non tree-based methods, both in terms of accuracy and computing times. In addition, multi-output trees inherit all intrinsic advantages of tree-based methods, such as robustness to irrelevant features, interpretability through feature importance scores, or fast computations of predictions, that make them very attractive to address multi-label problems. The computational complexity of learning multi-output trees is however similar to that of the binary relevance method. Both approaches are indeed $O(pdn \log n)$, where p is the number of input features, d the number of candidate output labels, and n the sample size; this is a limiting factor when dealing with large sets of candidate labels.

One generic approach to reduce computational complexity is to apply some compression technique prior to the training stage to reduce the number of outputs to a number q much smaller than the total number d of labels. A model can then be trained to make predictions in the compressed output space and a prediction in the original label space can be obtained by decoding the compressed prediction. As multi-label vectors are typically very sparse, one can expect a drastic dimensionality reduction by using appropriate compression techniques. This idea has been explored for example in (Hsu et al., 2009) using compressed sensing, and in (Cisse et al., 2013) using bloom filters, in both cases using regularized linear models as base learners. The approach obviously reduces computing times for training the model. Random projections are also exploited in (Tsoumakas et al., 2014) for multi-target regression. In this latter work however, they are not used to improve computing times by compression but instead to improve predictive performance. Indeed, more (sparse) random projections are computed than there are outputs and they are used each as an output to train some single target regressor. As in (Cisse et al., 2013; Hsu et al., 2009), the predictions of the regressors need to be decoded at prediction time to obtain a prediction in the original output space. This is achieved in (Tsoumakas et al., 2014) by solving an overdetermined linear system.

In this chapter, we explore the use of random output space projections for large-scale multi-label classification in the context of tree-based ensemble methods. We first explore the idea proposed for linear models in (Hsu et al., 2009) with random forests: a (single) random projection of the multi-label vector to a q -dimensional random subspace is computed and then a multi-output random forest is grown based on score computations using the projected outputs. We exploit however the fact that the approximation provided by a tree ensemble is a weighted average of output vectors from the training

sample to avoid the decoding stage: at training time all leaf labels are directly computed in the original multi-label space. We show theoretically and empirically that when q is large enough, ensembles grown on such random output spaces are equivalent to ensembles grown on the original output space. When d is large enough compared to n , this idea hence may reduce computing times at the learning stage without affecting accuracy and computational complexity of predictions.

Next, we propose to exploit the randomization inherent to the projection of the output space as a way to obtain randomized trees in the context of ensemble methods: each tree in the ensemble is thus grown from a different randomly projected subspace of dimension q . As previously, labels at leaf nodes are directly computed in the original output space to avoid the decoding step. We show, theoretically, that this idea can lead to better accuracy than the first idea and, empirically, that best results are obtained on many problems with very low values of q , which leads to significant computing time reductions at the learning stage. In addition, we study the interaction between input randomization (à la Random Forests) and output randomization (through random projections), showing that there is an interest, both in terms of predictive performance and in terms of computing times, to optimally combine these two ways of randomization. All in all, the proposed approach constitutes a very attractive way to address large-scale multi-label problems with tree-based ensemble methods.

The rest of the chapter is structured as follows: Section 5.1 presents the proposed algorithms and their theoretical properties; Section 5.2 analyses the proposed algorithm from a bias-variance perspective; Section 5.3 provides the empirical validations, whereas Section 5.4 discusses our work and provides further research directions.

5.1 METHODS

We first present how we propose to exploit random projections to reduce the computational burden of learning single multi-output trees in very high-dimensional output spaces. Then we present and compare two ways to exploit this idea with ensembles of trees.

5.1.1 *Multi-output regression trees in randomly projected output spaces*

The multi-output single tree algorithm described in Chapter 3 requires the computation of the sum of impurity criterion, such as the variance (or Gini), at each tree node and for each candidate split. When \mathcal{Y} is very high-dimensional, this computation constitutes the main computational bottleneck of the algorithm. We thus propose to approximate variance computations by using random projections of the output space. The multi-output regression tree algo-

rithm is modified as follows (denoting by \mathcal{L} the learning sample $\mathcal{L} = ((x^i, y^i) \in \mathcal{X} \times \mathcal{Y})_{i=1}^n$):

- First, a projection matrix Φ of dimension $q \times d$ is randomly generated.
- A new dataset $\mathcal{L}_m = ((x^i, \Phi y^i))_{i=1}^n$ is constructed by projecting each learning sample output using the projection matrix Φ .
- A tree (structure) \mathcal{T}_m is grown using the projected learning sample \mathcal{L}_m .
- Predictions \hat{y} at each leaf of \mathcal{T} are computed using the corresponding outputs in the original output space.

The resulting tree is exploited in the standard way to make predictions: an input vector x is propagated through the tree until it reaches a leaf from which a prediction \hat{y} in the original output space is directly retrieved.

If Φ satisfies the Jonhson-Lindenstrauss lemma (Equation 2.86), the following theorem shows that variance computed in the projected subspace is an ϵ -approximation of the variance computed over the original space.

Theorem 1. *Given $\epsilon > 0$, a sample $(y^i)_{i=1}^n$ of n points $y \in \mathbb{R}^d$, and a projection matrix $\Phi \in \mathbb{R}^{q \times d}$ such that for all $i, j \in \{1, \dots, n\}$ the condition given by Equation 2.86 holds, we have also:*

$$(1 - \epsilon) \text{Var}((y^i)_{i=1}^n) \leq \text{Var}((\Phi y^i)_{i=1}^n) \leq (1 + \epsilon) \text{Var}((y^i)_{i=1}^n). \quad (5.1)$$

Proof. The sum of the variances of n observations drawn from a random vector $y \in \mathbb{R}^d$ can be interpreted as a sum of squared euclidean distances between the pairs of observations

$$\text{Var}((y^i)_{i=1}^n) = \frac{1}{2n^2} \sum_{i=1}^n \sum_{j=1}^n \|y^i - y^j\|^2. \quad (5.2)$$

Starting from the defition of the variance, we have

$$\begin{aligned} & \text{Var}((y^i)_{i=1}^n) \\ & \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \|y^i - \frac{1}{n} \sum_{j=1}^n y^j\|^2 \end{aligned} \quad (5.3)$$

$$= \frac{1}{n} \sum_{i=1}^n (y^i - \frac{1}{n} \sum_{j=1}^n y^j)^T (y^i - \frac{1}{n} \sum_{k=1}^n y^k) \quad (5.4)$$

$$= \frac{1}{n} \sum_{i=1}^n \left(y^{iT} y^i - \frac{2}{n} \sum_{j=1}^n y^{iT} y^j + \frac{1}{n^2} \sum_{j=1}^n \sum_{k=1}^n y^{jT} y^k \right) \quad (5.5)$$

$$= \frac{1}{n} \sum_{i=1}^n y^{i\top} y^i - \frac{2}{n^2} \sum_{i=1}^n \sum_{j=1}^n y^{i\top} y^j + \frac{1}{n^2} \sum_{j=1}^n \sum_{k=1}^n y^{j\top} y^k \quad (5.6)$$

$$= \frac{1}{n} \sum_{i=1}^n y^{i\top} y^i - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n y^{i\top} y^j \quad (5.7)$$

$$= \frac{1}{2n} \sum_{i=1}^n y^{i\top} y^i + \frac{1}{2n} \sum_{j=1}^n y^{j\top} y^j - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n y^{i\top} y^j \quad (5.8)$$

$$= \frac{1}{2n^2} \sum_{i=1}^n \sum_{j=1}^n y^{i\top} y^i + \frac{1}{2n^2} \sum_{i=1}^n \sum_{j=1}^n y^{j\top} y^j - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n y^{i\top} y^j \quad (5.9)$$

$$= \frac{1}{2n^2} \sum_{i=1}^n \sum_{j=1}^n (y^{i\top} y^i + y^{j\top} y^j - 2y^{i\top} y^j) \quad (5.10)$$

$$= \frac{1}{2n^2} \sum_{i=1}^n \sum_{j=1}^n \|y^i - y^j\|^2. \quad (5.11)$$

From the Johnson-Lindenstrauss Lemma we have for any i, j

$$(1 - \epsilon) \|y^i - y^j\|^2 \leq \|\Phi y^i - \Phi y^j\|^2 \leq (1 + \epsilon) \|y^i - y^j\|^2. \quad (5.12)$$

By summing the three terms of Equation 5.12 over all pairs i, j and dividing by $1/(2n^2)$ and by then using Equation 5.2, we get Equation 5.1.

□

As a consequence, any split score approximated from the randomly projected output space will be ϵ -close to the unprojected scores in any subsample of the complete learning sample. Thus, if the condition given by Equation 2.86) is satisfied for a sufficiently small ϵ then the tree grown from the projected data will be identical to the tree grown from the original data¹.

For a given size q of the projection subspace, the complexity is reduced from $O(dn)$ to $O(qn)$ for the computation of one split score and thus from $O(dpn \log n)$ to $O(qpn \log n)$ for the construction of one full (balanced) tree, where one can expect q to be much smaller than d and at worst of $O(\epsilon^{-2} \log n)$. The whole procedure requires to generate the projection matrix and to project the training data. These two steps are respectively $O(dq)$ and $O(ndq)$ but they can often be significantly accelerated by exploiting the sparsity of the projection matrix and/or of the original output data, and they are called only once before growing the tree.

All in all, this means that when d is sufficiently large, the random projection approach may allow us to significantly reduce tree building complexity from $O(dpn \log n)$ to $O(qpn \log n + ndq)$, without impact on predictive accuracy (see Section 5.3, for empirical results).

¹ Strictly speaking, this is only the case when the optimum scores of test splits as computed over the original output space are isolated, i.e. when there is only one single best split, no tie.

5.1.2 *Exploitation in the context of tree ensembles*

The idea developed in the previous section can be directly exploited in the context of ensembles of randomized multi-output regression trees. Instead of building a single tree from the projected learning sample, one can grow a randomized ensemble of them. This “shared subspace” algorithm is described in pseudo-code in Algorithm 5.1.

Algorithm 5.1 Grow t decision trees on a single shared subspace Φ using learning samples $\mathcal{L} = ((x^i, y^i) \in (\mathbb{R}^p \times \mathbb{R}^d))_{i=1}^n$

```

1: function GROWFORESTSHAREDOUTPUTSUBSPACE( $\mathcal{L}, t$ )
2:   Generate a sub-space  $\Phi \in \mathbb{R}^{q \times d}$ ;
3:   for  $j = 1$  to  $t$  do
4:     Build a tree structure  $\mathcal{T}_j$  using  $((x^i, \Phi y^i))_{i=1}^n$ ;
5:     Label the leaves of  $\mathcal{T}_j$  using  $((x^i, y^i))_{i=1}^n$ ;
6:     Add the labelled tree  $\mathcal{T}_j$  to the ensemble;
7:   end for
8: end function

```

Another idea is to exploit the random projections used so as to introduce a novel kind of diversity among the different trees of an ensemble. Instead of building all the trees of the ensemble from a same shared output-space projection, one could instead grow each tree in the ensemble from a different output-space projection. Algorithm 5.2 implements this idea in pseudo-code. The randomization introduced by the output space projection can of course be combined with any existing randomization scheme to grow ensembles of trees. In this chapter, we will consider the combination of random projections with the randomizations already introduced in Random Forests and Extra Trees. The interplay between these different randomizations will be discussed theoretically in the next subsection by a bias/variance analysis and empirically in Section 5.3. Note that while when looking at single trees or shared ensembles, the size q of the projected subspace should not be too small so that condition (Equation 2.86) is satisfied, the optimal value of q when projections are randomized at each tree is likely to be smaller, as suggested by the bias/variance analysis in the next section.

From the computational point of view, the main difference between these two ways of transposing random-output projections to ensembles of trees is that in the case of Algorithm 5.2, the generation of the projection matrix Φ and the computation of projected outputs is carried out t times, while it is done only once for the case of Algorithm 5.1. These aspects will be empirically evaluated in Section 5.3.

Algorithm 5.2 Grow t decision trees on individual random subspaces $(\Phi_j)_{j=1}^t$ using learning samples $\mathcal{L} = ((x^i, y^i) \in (\mathbb{R}^p \times \mathbb{R}^d))_{i=1}^n$

```

1: function GROWFORESTOUTPUTSUBSPACE( $\mathcal{L}, t$ )
2:   for  $j = 1$  to  $t$  do
3:     Generate a sub-space  $\Phi_j \in \mathbb{R}^{q \times d}$ ;
4:     Build a tree structure  $\mathcal{T}_j$  using  $((x^i, \Phi_j y^i))_{i=1}^n$ ;
5:     Label the leaves of  $\mathcal{T}_j$  using  $((x^i, y^i))_{i=1}^n$ ;
6:     Add the labelled tree  $\mathcal{T}_j$  to the ensemble;
7:   end for
8: end function

```

5.2 BIAS/VARIANCE ANALYSIS

In this section, we adapt the bias/variance analysis carried out in Section 4.2.1 to take into account random output projections. The details of the derivations are reported in Section 5.2.1 for a single tree and in Section 5.2.2 for an ensemble of t randomized trees.

Let us denote by $f_{\mathcal{L}, \phi, \sigma}(\cdot) : \mathcal{X} \rightarrow \mathbb{R}^d$ a single multi-output tree obtained from a projection matrix ϕ (below we use Φ to denote the corresponding random variable), where σ is the value of a random variable capturing the random perturbation scheme used to build this tree (e.g., bootstrapping and/or random input space selection). The square error of this model at some point $x \in \mathcal{X}$ is defined by:

$$\text{Err}(f_{\mathcal{L}, \phi, \sigma}(x)) \stackrel{\text{def}}{=} \mathbb{E}_{Y|x} \{\|Y - f_{\mathcal{L}, \phi, \sigma}(x)\|^2\},$$

and its average can be decomposed in its residual error, (squared) bias, and variance terms denoted:

$$\mathbb{E}_{\mathcal{L}, \Phi, \sigma} \{\text{Err}(f_{\mathcal{L}, \phi, \sigma}(x))\} = \sigma_{\mathbb{R}}^2(x) + B^2(x) + V(x)$$

where the variance term $V(x)$ can be further decomposed as the sum of the following three terms:

$$\begin{aligned} V_{\mathcal{L}}(x) &= \text{Var}_{\mathcal{L}} \{\mathbb{E}_{\Phi, \sigma | \mathcal{L}} \{f_{\mathcal{L}, \phi, \sigma}(x)\}\} \\ V_{\text{Algo}}(x) &= \mathbb{E}_{\mathcal{L}} \{\mathbb{E}_{\Phi | \mathcal{L}} \{\text{Var}_{\sigma | \mathcal{L}, \Phi} \{f_{\mathcal{L}, \phi, \sigma}(x)\}\}\}, \\ V_{\text{Proj}}(x) &= \mathbb{E}_{\mathcal{L}} \{\text{Var}_{\Phi | \mathcal{L}} \{\mathbb{E}_{\sigma | \mathcal{L}, \Phi} \{f_{\mathcal{L}, \phi, \sigma}(x)\}\}\}, \end{aligned}$$

that measure errors due to the randomness of, respectively, the learning sample, the tree algorithm, and the output space projection (see Section 5.2.1).

Approximations computed respectively by Algorithm 5.1 and Algorithm 5.2 take the following forms:

- $f_{1; \mathcal{L}, \sigma^t, \phi}(x) = \frac{1}{t} \sum_{i=1}^t f_{\mathcal{L}, \phi, \sigma_i}(x)$
- $f_{2; \mathcal{L}, \sigma^t, \phi^t}(x) = \frac{1}{t} \sum_{i=1}^t f_{\mathcal{L}, \phi_i, \sigma_i}(x),$

where $\sigma^t = (\epsilon_1, \dots, \epsilon_t)$ and $\phi^t = (\phi_1, \dots, \phi_t)$ are vectors of i.i.d. values of the random variables σ and Φ respectively.

We are interested in comparing the average errors of these two algorithms, where the average is taken over all random parameters (including the learning sample). We show that these can be decomposed as follows (see Section 5.2.2):

$$\begin{aligned} & E_{\mathcal{L}, \Phi, \sigma^t} \{ \text{Err}(f_{1; \mathcal{L}, \Phi, \sigma^t}(x)) \} \\ &= \sigma_{\mathbb{R}}^2(x) + B^2(x) + V_{\mathcal{L}}(x) + \frac{V_{\text{Algo}}(x)}{t} + V_{\text{Proj}}(x), \\ & E_{\mathcal{L}, \Phi^t, \sigma^t} \{ \text{Err}(f_{2; \mathcal{L}, \Phi^t, \sigma^t}(x)) \} \\ &= \sigma_{\mathbb{R}}^2(x) + B^2(x) + V_{\mathcal{L}}(x) + \frac{V_{\text{Algo}}(x) + V_{\text{Proj}}(x)}{t}. \end{aligned}$$

From this result, it is hence clear that Algorithm 5.2 can not be worse, on the average, than Algorithm 5.1. If the additional computational burden needed to generate a different random projection for each tree is not problematic, then Algorithm 5.2 should always be preferred to Algorithm 5.1.

For a fixed level of tree randomization (σ), whether the additional randomization brought by random projections could be beneficial in terms of predictive performance remains an open question that will be addressed empirically in the next section. Nevertheless, with respect to an ensemble grown from the original output space, one can expect that the output-projections will always increase the bias term, since they disturb the algorithm in its objective of reducing the errors on the learning sample. For small values of q , the average error will therefore decrease (with a sufficiently large number t of trees) only if the increase in bias is compensated by a decrease of variance.

The value of q , the dimension of the projected subspace, that will lead to the best tradeoff between bias and variance will hence depend both on the level of tree randomization and on the learning problem. The more (resp. less) tree randomization, the higher (resp. the lower) could be the optimal value of q , since both randomizations affect bias and variance in the same direction.

5.2.1 Single random trees.

Let us denote by $f_{\mathcal{L}, \phi, \sigma} : \mathcal{X} \rightarrow \mathbb{R}^d$ a single multi-output (random) tree obtained from a projection matrix ϕ (below we use Φ to denote the corresponding random variable), where σ is the value of a random variable σ capturing the random perturbation scheme used to build this tree (e.g., bootstrapping and/or random input space selection). Denoting by $\text{Err}(f_{\mathcal{L}, \phi, \sigma}(x))$ the square error of this model at some point $x \in \mathcal{X}$ defined by:

$$E_{P_{y|x}} \{ \|y - f_{\mathcal{L}, \phi, \sigma}(x)\|^2 \}. \quad (5.13)$$

The average of this square error can be decomposed as follows:

$$\begin{aligned} & \mathbb{E}_{\mathcal{L}, \Phi, \sigma} \{ \text{Err}(f_{\mathcal{L}, \Phi, \sigma}(x)) \} \\ = & \sigma_{\mathbb{R}}^2(x) + \|f_{\text{Bayes}}(x) - \bar{f}(x)\|^2 + \text{Var}_{\mathcal{L}, \Phi, \sigma} \{ f_{\mathcal{L}, \Phi, \sigma}(x) \}, \end{aligned}$$

where

$$\begin{aligned} \bar{f}(x) & \stackrel{\text{def}}{=} \mathbb{E}_{\mathcal{L}, \Phi, \sigma} \{ f_{\mathcal{L}, \Phi, \sigma}(x) \} \\ f_{\text{Bayes}}(x) & = \mathbb{E}_{Y|x} \{ Y \} \\ \text{Var}_{\mathcal{L}, \Phi, \sigma} \{ f_{\mathcal{L}, \Phi, \sigma}(x) \} & \stackrel{\text{def}}{=} \mathbb{E}_{\mathcal{L}, \Phi, \sigma} \{ \|f_{\mathcal{L}, \Phi, \sigma}(x) - \bar{f}(x)\|^2 \}. \end{aligned}$$

The three terms of this decomposition are respectively the residual error, the bias, and the variance of this estimator (at x).

The variance term can be further decomposed as follows using the law of total variance:

$$\begin{aligned} & \text{Var}_{\mathcal{L}, \Phi, \sigma} \{ f_{\mathcal{L}, \Phi, \sigma}(x) \} \\ = & \text{Var}_{\mathcal{L}} \{ \mathbb{E}_{\Phi, \sigma | \mathcal{L}} \{ f_{\mathcal{L}, \Phi, \sigma}(x) \} \} \\ & + \mathbb{E}_{\mathcal{L}} \{ \text{Var}_{\Phi, \sigma | \mathcal{L}} \{ f_{\mathcal{L}, \Phi, \sigma}(x) \} \}. \end{aligned} \quad (5.14)$$

The first term is the variance due to the learning sample randomization and the second term is the average variance (over \mathcal{L}) due to both the random forest randomization and the random output projection. By using the law of total variance a second time, the second term of Equation 5.14) can be further decomposed as follows:

$$\begin{aligned} & \mathbb{E}_{\mathcal{L}} \{ \text{Var}_{\Phi, \sigma | \mathcal{L}} \{ f_{\mathcal{L}, \Phi, \sigma}(x) \} \} \\ = & \mathbb{E}_{\mathcal{L}} \{ \text{Var}_{\Phi | \mathcal{L}} \{ \mathbb{E}_{\sigma | \mathcal{L}, \Phi} \{ f_{\mathcal{L}, \Phi, \sigma}(x) \} \} \} \\ & + \mathbb{E}_{\mathcal{L}} \{ \mathbb{E}_{\Phi | \mathcal{L}} \{ \text{Var}_{\sigma | \mathcal{L}, \Phi} \{ f_{\mathcal{L}, \Phi, \sigma}(x) \} \} \}. \end{aligned} \quad (5.15)$$

The first term of this decomposition is the variance due to the random choice of a projection and the second term is the average variance due to the random forest randomization. Note that all these terms are non negative. In what follows, we will denote these three terms respectively $V_{\mathcal{L}}(x)$, $V_{\text{Algo}}(x)$, and $V_{\text{Proj}}(x)$. We thus have:

$$\text{Var}_{\mathcal{L}, \Phi, \sigma} \{ f_{\mathcal{L}, \Phi, \sigma}(x) \} = V_{\mathcal{L}}(x) + V_{\text{Algo}}(x) + V_{\text{Proj}}(x),$$

with

$$\begin{aligned} V_{\mathcal{L}}(x) & = \text{Var}_{\mathcal{L}} \{ \mathbb{E}_{\Phi, \sigma | \mathcal{L}} \{ f_{\mathcal{L}, \Phi, \sigma}(x) \} \} \\ V_{\text{Algo}}(x) & = \mathbb{E}_{\mathcal{L}} \{ \text{Var}_{\Phi | \mathcal{L}} \{ \text{Var}_{\sigma | \mathcal{L}, \Phi} \{ f_{\mathcal{L}, \Phi, \sigma}(x) \} \} \}, \\ V_{\text{Proj}}(x) & = \mathbb{E}_{\mathcal{L}} \{ \mathbb{E}_{\Phi | \mathcal{L}} \{ \text{Var}_{\sigma | \mathcal{L}, \Phi} \{ f_{\mathcal{L}, \Phi, \sigma}(x) \} \} \}, \end{aligned}$$

5.2.2 *Ensembles of t random trees.*

When the random projection is fixed for all t trees in the ensemble (Algorithm 5.1), the algorithm computes an approximation, denoted $f_1(x; \mathcal{L}, \phi, \sigma^t)$, that takes the following form:

$$f_{1; \mathcal{L}, \phi, \sigma^t}(x) = \frac{1}{t} \sum_{i=1}^t f_{\mathcal{L}, \phi, \sigma_i}(x),$$

where $\sigma^t = (\sigma_1, \dots, \sigma_t)$ is a vector of i.i.d. values of the random variable σ . When a different random projection is chosen for each tree (Algorithm 5.2), the algorithm computes an approximation, denoted by $f_2(x; \mathcal{L}, \phi^t, \sigma^t)$, of the following form:

$$f_{2; \mathcal{L}, \phi^t, \sigma^t}(x) = \frac{1}{t} \sum_{i=1}^t f_{\mathcal{L}, \phi_i, \sigma_i}(x),$$

where $\phi^t = (\phi_1, \dots, \phi_t)$ is also a vector of i.i.d. random projection matrices.

We would like to compare the average errors of these two algorithms with the average errors of the original single tree method, where the average is taken for all algorithms over their random parameters (that include the learning sample).

Given that all trees are grown independently of each other, one can show that the average models corresponding to each algorithm are equal:

$$\begin{aligned} \bar{f}(x) &= E_{\mathcal{L}, \Phi, \sigma^t} \{f_{1; \mathcal{L}, \Phi, \sigma^t}(x)\} \\ &= E_{\mathcal{L}, \Phi^t, \sigma^t} \{f_{2; \mathcal{L}, \Phi^t, \sigma^t}(x)\}. \end{aligned}$$

They thus all have the exact same bias (and residual error) and differ only in their variance.

Using the same argument, the first term of the variance decomposition in (5.14), ie. $V_{\mathcal{L}}(x)$, is the same for all three algorithms since:

$$\begin{aligned} &E_{\Phi, \sigma | \mathcal{L}} \{f_{\mathcal{L}, \Phi, \sigma}(x)\} \\ &= E_{\Phi, \sigma^t | \mathcal{L}} \{f_{1; \mathcal{L}, \Phi, \sigma^t}(x)\} \\ &= E_{\Phi^t, \sigma^t | \mathcal{L}} \{f_{2; \mathcal{L}, \Phi^t, \sigma^t}(x)\}. \end{aligned}$$

Their variance thus only differ in the second term of Equation 5.14.

Again, because of the conditional independence of the ensemble terms given the learning set \mathcal{L} and the projection matrix ϕ , Algorithm 5.1, which keeps the output projection fixed for all trees, is such that

$$E_{\sigma^t | \mathcal{L}, \Phi} \{f_{1; \mathcal{L}, \Phi, \sigma^t}(x)\} = E_{\sigma | \mathcal{L}, \Phi} \{f_{\mathcal{L}, \Phi, \sigma}(x)\}$$

and

$$\text{Var}_{\sigma^t | \mathcal{L}, \Phi} \{f_{1; \mathcal{L}, \Phi, \sigma^t}(x)\} = \frac{1}{t} \text{Var}_{\sigma | \mathcal{L}, \Phi} \{f_{\mathcal{L}, \Phi, \sigma}(x)\}.$$

It thus divides the second term of Equation 5.15 by the number t of ensemble terms. Algorithm 5.2, on the other hand, is such that:

$$\text{Var}_{\Phi^t, \sigma^t | \mathcal{L}} \{f_{2; \mathcal{L}, \Phi, \sigma^t}(x)\} = \frac{1}{t} \text{Var}_{\Phi, \sigma | \mathcal{L}} \{f_{\mathcal{L}, \Phi, \sigma}(x)\},$$

and thus divides the second term of Equation 5.14 by t .

Putting all these results together one gets that:

$$\begin{aligned} & E_{\mathcal{L}, \Phi, \sigma} \{\text{Err}(f_{\mathcal{L}, \Phi, \sigma}(x))\} \\ &= \sigma_{\mathbb{R}}^2(x) + B^2(x) + V_{\mathcal{L}}(x) + V_{\text{Algo}}(x) + V_{\text{Proj}}(x), \\ & E_{\mathcal{L}, \Phi, \sigma^t} \{\text{Err}(f_{1; \mathcal{L}, \Phi, \sigma^t}(x))\} \\ &= \sigma_{\mathbb{R}}^2(x) + B^2(x) + V_{\mathcal{L}}(x) + \frac{V_{\text{Algo}}(x)}{t} + V_{\text{Proj}}(x), \\ & E_{\mathcal{L}, \Phi^t, \sigma^t} \{\text{Err}(f_{2; \mathcal{L}, \Phi^t, \sigma^t}(x))\} \\ &= \sigma_{\mathbb{R}}^2(x) + B^2(x) + V_{\mathcal{L}}(x) + \frac{V_{\text{Algo}}(x) + V_{\text{Proj}}(x)}{t}. \end{aligned}$$

Given that all terms are positive, this result clearly shows that Algorithm 5.2 can not be worse than Algorithm 5.1.

5.3 EXPERIMENTS

5.3.1 Effect of the size q of the Gaussian output space

To illustrate the behaviour of our algorithms, we first focus on the “Delicious” dataset (Tsoumakas et al., 2008a), which has a large number of labels ($d = 983$), of input features ($p = 500$), and of training ($n_{\text{LS}} = 12920$) and testing ($n_{\text{TS}} = 3185$) samples.

The top part of figure 5.1 shows, when Gaussian output-space projections are combined with the standard CART algorithm building a single tree, how the precision converges (cf Theorem 1) when q increases towards d . We observe that in this case, convergence is reached around $q = 200$ at the expense of a slight decrease of accuracy, so that a compression factor of about 5 is possible with respect to the original output dimension $d = 983$.

The bottom part of figure 5.1 shows, on the same dataset, how the method behaves when combined with Random Forests. Let us first notice that the Random Forests grown on the original output space (green line) are significantly more accurate than the single trees, their accuracy being almost twice as high. We also observe that Algorithm 5.2 (orange curve) converges much more rapidly than Algorithm 5.1 (blue curve) and slightly outperforms the Random Forest grown on the original output space. It needs only about $q = 25$ components to converge, while Algorithm 5.1 needs about $q = 75$ of them. These results are in accordance with the analysis of Section 5.2, showing that Algorithm 5.2 can’t be inferior to Algorithm 5.1. In the rest of this chapter we will therefore focus on Algorithm 5.2.

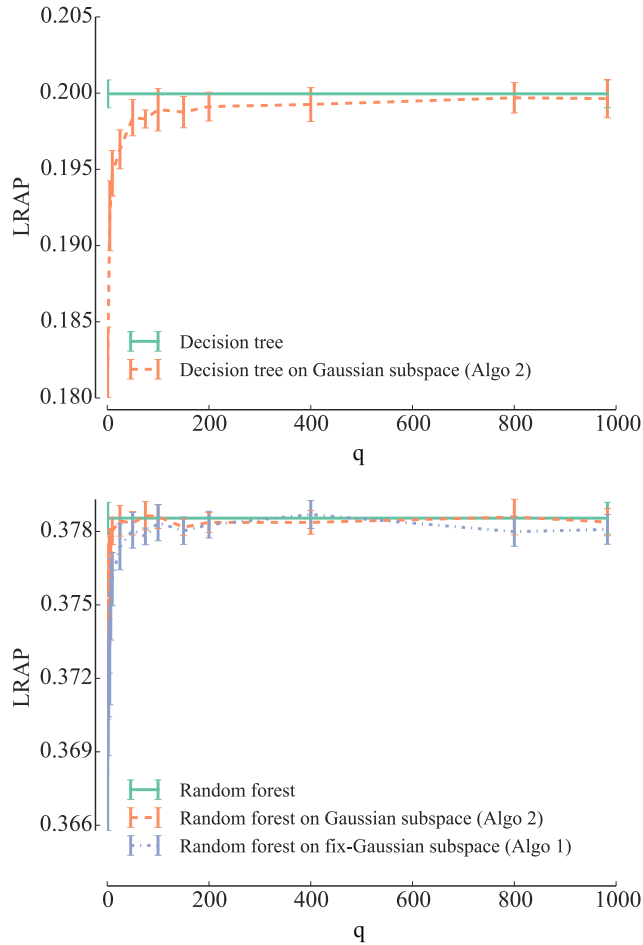


Figure 5.1: Models built for the “Delicious” dataset ($d = 983$) for growing numbers q of Gaussian projections. Top: single unpruned CART trees ($n_{\min} = 1$); Bottom: Random Forests ($k = \sqrt{p}$, $t = 100$, $n_{\min} = 1$). The curves represent average values (and standard deviations) obtained from 10 applications of the randomized algorithms over a same single LS/TS split.

5.3.2 Systematic analysis over 24 datasets

To assess our methods, we have collected 24 different multi-label classification datasets from the literature (see Section D of the supplementary material, for more information and bibliographic references to these datasets) covering a broad spectrum of application domains and ranges of the output dimension ($d \in [6; 3993]$, see Table 5.1). For 21 of the datasets, we made experiments where the dataset is split randomly into a learning set of size n_{LS} , and a test set of size n_{TS} , and are repeated 10 times (to get average precisions and standard deviations), and for 3 of them we used a ten-fold cross-validation scheme (see Table 5.1).

Table 5.1 shows our results on the 24 multi-label datasets, by comparing Random Forests learnt on the original output space with those

learnt by Algorithm 5.2 combined with Gaussian subspaces of size $q \in \{1, d, \ln d\}^2$. In these experiments, the three parameters of Random Forests are set respectively to $k = \sqrt{p}$, $n_{\min} = 1$ (default values, see (Geurts et al., 2006a)) and $t = 100$ (reasonable computing budget). Each model is learnt ten times on a different shuffled train/testing split, except for the 3 EUR-lex datasets where we kept the original 10 folds of cross-validation.

We observe that for all datasets (except maybe SCOP-GO), taking $q = d$ leads to a similar average precision to the standard Random Forests, i.e. no difference superior to one standard deviation of the error. On 11 datasets, we see that $q = 1$ already yields a similar average precision (values not underlined in column $q = 1$). For the 13 remaining datasets, increasing q to $\ln d$ significantly decreases the gap with the Random Forest baseline and 3 more datasets reach this baseline. We also observe that on several datasets such as “Drug-interaction” and “SCOP-GO”, better performance on the Gaussian subspace is attained with high output randomization ($q = \{1, \ln d\}$) than with $q = d$. We thus conclude that the optimal level of output randomization (i.e. the optimal value of the ratio q/d) which maximizes accuracy performances, is dataset dependent.

While our method is intended for tasks with very high dimensional output spaces, we however notice that even with relatively small numbers of labels, its accuracy remains comparable to the baseline, with suitable q .

To complete the analysis, let’s carry out the same experiments with a different base-learner combining Gaussian random projections (with $q \in \{1, \ln d, d\}$) with the Extra Trees method of (Geurts et al., 2006a). Results on 23 datasets are compiled in Table 5.2.

Like for Random Forests, we observe that for all 23 datasets taking $q = d$ leads to a similar average precision to the standard Random Forests, i.e. no difference superior to one standard deviation of the error. This is already the case with $q = 1$ for 12 datasets and with $q = \ln d$ for 4 more datasets. Interestingly, on 3 datasets with $q = 1$ and 3 datasets with $q = \ln d$, the increased randomization brought by the projections actually improves average precision with respect to standard Random Forests (bold values in Table 5.2).

² $\ln d$ is rounded to the nearest integer value; in Table 5.1 the values of $\ln d$ vary between 2 for $d = 6$ and 8 for $d = 3993$.

Table 5.1: High output space compression ratio is possible, with no or negligible average precision reduction ($t = 100$, $n_{\min} = 1$, $k = \sqrt{p}$). Each dataset has n_{LS} training samples, n_{TS} testing samples, p input features and d labels. Label ranking average precisions are displayed in terms of their mean values and standard deviations over 10 random LS/TS splits, or over the 10 folds of cross-validation. Mean scores in the last three columns are underlined if they show a difference with respect to the standard Random Forests of more than one standard deviation.

Datasets Name	Random Forests	Random Forests on Gaussian sub-space		
		$q = 1$	$q = \lfloor 0.5 + \ln d \rfloor$	$q = d$
emotions	0.800 ±0.014	0.800 ±0.010	0.810 ±0.014	0.810 ±0.016
scene	0.870 ±0.003	0.875 ±0.007	0.872 ±0.004	0.872 ±0.004
yeast	0.759 ±0.008	<u>0.748</u> ±0.006	0.755 ±0.004	0.758 ±0.005
tmc2017	0.756 ±0.003	<u>0.741</u> ±0.003	<u>0.748</u> ±0.003	0.757 ±0.003
genbase	0.992 ±0.004	0.994 ±0.002	0.994 ±0.004	0.993 ±0.004
reuters	0.865 ±0.004	0.864 ±0.003	0.863 ±0.004	0.862 ±0.004
medical	0.848 ±0.009	<u>0.836</u> ±0.011	0.842 ±0.014	0.841 ±0.009
enron	0.683 ±0.009	0.680 ±0.006	0.685 ±0.009	0.686 ±0.008
mediamill	0.779 ±0.001	<u>0.772</u> ±0.001	0.777 ±0.002	0.779 ±0.002
Yeast-GO	0.420 ±0.010	<u>0.353</u> ±0.008	<u>0.381</u> ±0.005	0.420 ±0.010
bibtex	0.566 ±0.004	<u>0.513</u> ±0.006	<u>0.548</u> ±0.007	0.564 ±0.008
CAL500	0.504 ±0.011	0.504 ±0.004	0.506 ±0.007	0.502 ±0.010
WIPO	0.490 ±0.010	<u>0.430</u> ±0.010	<u>0.460</u> ±0.010	0.480 ±0.010
EUR-Lex (subj.)	0.840 ±0.005	<u>0.814</u> ±0.004	<u>0.828</u> ±0.005	0.840 ±0.004
bookmarks	0.453 ±0.001	<u>0.436</u> ±0.002	<u>0.445</u> ±0.002	0.453 ±0.002
diatoms	0.700 ±0.010	<u>0.650</u> ±0.010	<u>0.670</u> ±0.010	0.710 ±0.020
corel5k	0.303 ±0.012	0.309 ±0.011	0.307 ±0.011	0.299 ±0.013
EUR-Lex (dir.)	0.814 ±0.006	<u>0.782</u> ±0.008	<u>0.796</u> ±0.009	0.813 ±0.007
SCOP-GO	0.811 ±0.004	0.808 ±0.005	0.811 ±0.004	<u>0.806</u> ±0.004
delicious	0.384 ±0.004	0.381 ±0.003	0.382 ±0.002	0.383 ±0.004
drug-interaction	0.379 ±0.014	0.384 ±0.009	0.378 ±0.013	0.367 ±0.016
protein-interaction	0.330 ±0.015	0.337 ±0.016	0.337 ±0.017	0.335 ±0.014
Expression-GO	0.235 ±0.005	<u>0.211</u> ±0.005	<u>0.219</u> ±0.005	0.232 ±0.005
EUR-Lex (desc.)	0.523 ±0.008	<u>0.485</u> ±0.008	<u>0.497</u> ±0.009	0.523 ±0.007

Table 5.2: Experiments with Gaussian projections and Extra Trees ($t = 100$, $n_{\min} = 1$, $k = \sqrt{p}$). Mean scores in the last three columns are underlined if they show a negative difference with respect to the standard Random Forests of more than one standard deviation. Bold values highlight improvement over standard RF of more than one standard deviation.

Datasets	Extra trees	Extra trees on Gaussian sub-space		
		$q = 1$	$q = \lfloor 0.5 + \ln d \rfloor$	$q = d$
emotions	0.81 ± 0.01	0.81 ± 0.014	0.80 ± 0.013	0.81 ± 0.014
scene	0.873 ± 0.004	0.876 ± 0.003	0.877 ± 0.007	0.878 ± 0.006
yeast	0.757 ± 0.008	<u>0.746</u> ± 0.004	0.752 ± 0.009	0.757 ± 0.01
tmc2017	0.782 ± 0.003	<u>0.759</u> ± 0.004	<u>0.77</u> ± 0.002	0.779 ± 0.002
genbase	0.987 ± 0.005	0.991 ± 0.004	0.992 ± 0.001	0.992 ± 0.005
reuters	0.88 ± 0.003	0.88 ± 0.003	0.878 ± 0.004	0.88 ± 0.004
medical	0.855 ± 0.008	0.867 ± 0.009	0.872 ± 0.006	0.862 ± 0.008
enron	0.66 ± 0.01	0.65 ± 0.01	0.663 ± 0.008	0.66 ± 0.01
mediamill	0.786 ± 0.002	<u>0.778</u> ± 0.002	<u>0.781</u> ± 0.002	0.784 ± 0.001
Yeast-GO	0.49 ± 0.009	<u>0.47</u> ± 0.01	0.482 ± 0.008	0.48 ± 0.01
bibtex	0.584 ± 0.005	<u>0.538</u> ± 0.005	<u>0.564</u> ± 0.004	0.583 ± 0.004
CAL500	0.5 ± 0.007	0.502 ± 0.008	0.499 ± 0.007	0.503 ± 0.009
WIPO	0.52 ± 0.01	<u>0.474</u> ± 0.007	<u>0.49</u> ± 0.01	0.515 ± 0.006
EUR-Lex (subj.)	0.845 ± 0.006	<u>0.834</u> ± 0.004	<u>0.838</u> ± 0.003	0.845 ± 0.005
bookmarks	0.453 ± 0.002	<u>0.436</u> ± 0.002	<u>0.444</u> ± 0.003	0.452 ± 0.002
diatoms	0.73 ± 0.01	<u>0.69</u> ± 0.01	<u>0.71</u> ± 0.01	0.73 ± 0.01
corel5k	0.285 ± 0.009	0.313 ± 0.011	0.309 ± 0.009	0.285 ± 0.011
EUR-Lex (dir.)	0.815 ± 0.007	<u>0.805</u> ± 0.006	<u>0.807</u> ± 0.009	0.815 ± 0.007
SCOP-GO	0.778 ± 0.005	0.782 ± 0.004	0.782 ± 0.006	0.778 ± 0.005
delicious	0.354 ± 0.003	0.36 ± 0.004	0.358 ± 0.004	0.355 ± 0.003
drug-interaction	0.353 ± 0.011	0.375 ± 0.017	0.364 ± 0.014	0.355 ± 0.016
protein-interaction	0.299 ± 0.013	0.307 ± 0.009	0.305 ± 0.012	0.306 ± 0.017
Expression-GO	0.231 ± 0.007	<u>0.218</u> ± 0.005	0.228 ± 0.005	0.235 ± 0.005

5.3.3 *Input vs output space randomization*

We study in this section the interaction of the additional randomization of the output space with that concerning the input space already built in the Random Forest method.

To this end, we consider the “Drug-interaction” dataset ($p = 660$ input features and $d = 1554$ output labels (Yamanishi et al., 2011)), and we study the effect of parameter k controlling the input space randomization of the Random Forest method with the randomization of the output space by Gaussian projections controlled by the parameter q . To this end, Figure 5.2 shows the evolution of the accuracy for growing values of k (i.e. decreasing strength of the input space randomization), for three different quite low values of q (in this case $q \in \{1, \ln d, 2 \ln d\}$). We observe that Random Forests learned on a very low-dimensional Gaussian subspace (red, blue and pink curves) yield essentially better performances than Random Forests on the original output space, and also that their behaviour with respect to the parameter k is quite different. On this dataset, the output-space randomisation makes the method completely immune to the ‘over-fitting’ phenomenon observed for high values of k with the baseline method (green curve).

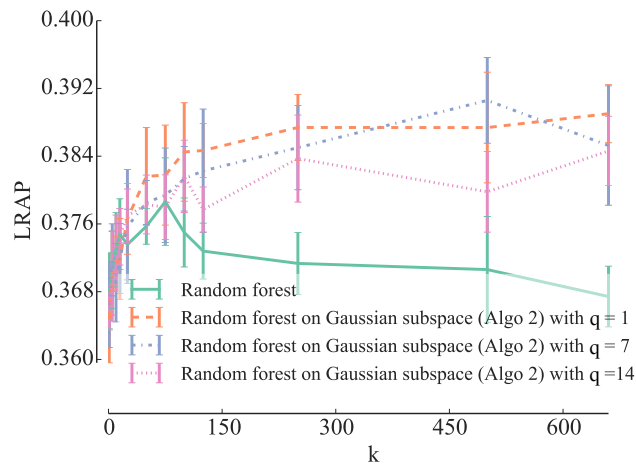


Figure 5.2: Output randomization with Gaussian projections yield better average precision than the original output space on the “Drug-Interaction” dataset ($n_{\min} = 1$, $t = 100$).

We carry out the same experiment, but on the “Delicious” dataset. Figure 5.3 shows the evolution of the accuracy for growing values of k (i.e. decreasing strength of the input space randomization), for three different values of q (in this case $q \in \{1, \ln d, 2 \ln d\}$) on a Gaussian output space.

Like on “Drug-interaction” (see Figure 5.2), using low-dimensional output spaces makes the method more robust with respect to over-fitting as k increases. However, unlike on “Drug-interaction”, it is not really possible to improve over baseline Random Forests by tuning

jointly input and output randomization. This shows that the interaction between q and k may be different from one dataset to another.

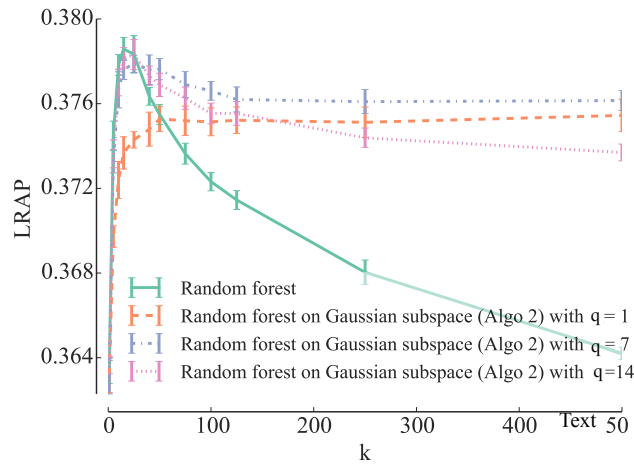


Figure 5.3: “Delicious” dataset: $n_{\min} = 1$; $t = 100$.

It is thus advisable to jointly optimize the value of q and k , so as to maximise the tradeoff between accuracy and computing times in a problem and algorithm specific way.

5.3.4 Alternative output dimension reduction techniques

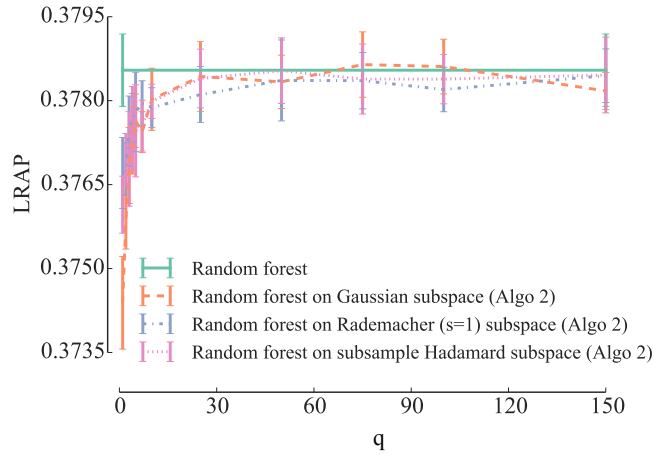
In this section, we study Algorithm 5.2 when it is combined with alternative output-space dimensionality reduction techniques. We focus again on the “Delicious” dataset, but similar trends could be observed on other datasets.

Figure 5.4a first compares Gaussian random projections with two other dense projections: Rademacher matrices with $s = 1$ (cf. Section 2.2) and compression matrices obtained by sub-sampling (without replacement) Hadamard matrices (Candes and Plan, 2011). We observe that Rademacher and subsample-Hadamard sub-spaces behave very similarly to Gaussian random projections.

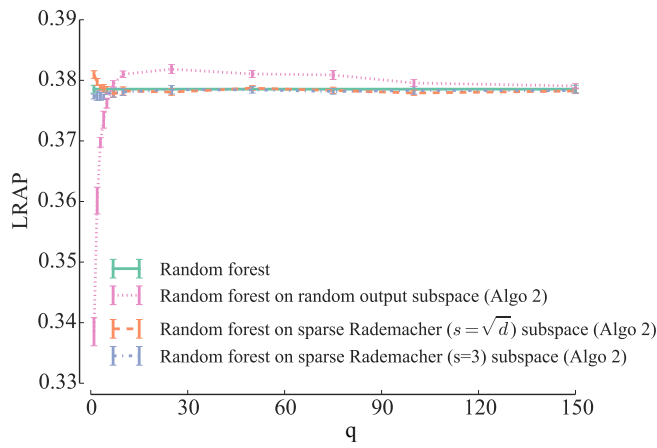
In a second step, we compare Gaussian random projections with two (very) sparse projections: first, sparse Rademacher sub-spaces obtained by setting the sparsity parameter s to 3 and \sqrt{d} , selecting respectively about 33% and 2% of the original outputs to compute each component, and second, sub-sampled identity subspaces, similar to (Tsoumakos and Vlahavas, 2007), where each of the q selected components corresponds to a randomly chosen original label and also preserve sparsity. Sparse projections are very interesting from a computational point of view as they require much less operations to compute the projections but the number of components required for condition (2.86) to be satisfied is typically higher than for dense projections (Candes and Plan, 2011; Li et al., 2006). Figure 5.4b compares these three projection methods with standard Random Forests

on the “delicious” dataset. All three projection methods converge to plain Random Forests as the number of components q increases but their behaviour at low q values are very different. Rademacher projections converge faster with $s = 3$ than with $s = 1$ and interestingly, the sparsest variant ($s = \sqrt{d}$) has its optimum at $q = 1$ and improves in this case over the Random Forests baseline. Random output subspaces converge slower but they lead to a notable improvement of the score over baseline Random Forests. This suggests that although their theoretical guarantees are less good, sparse projections actually provide on this problem a better bias/variance tradeoff than dense ones when used in the context of Algorithm 5.2.

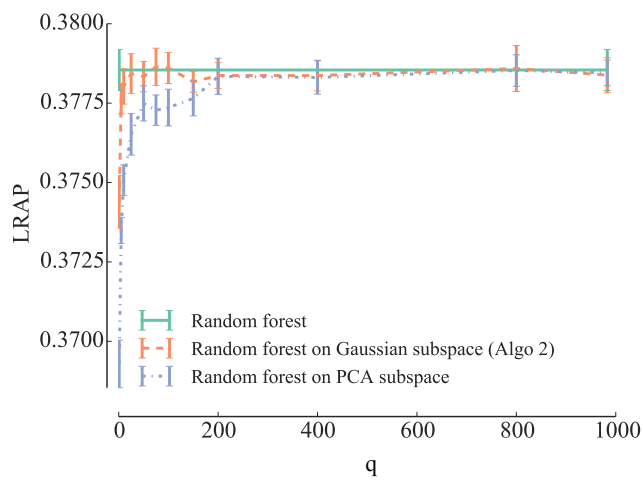
Another popular dimension reduction technique is the principal component analysis (PCA). In Figure 5.4c, we repeat the same experiment to compare PCA with Gaussian random projections. Concerning PCA, the curve is generated in decreasing order of eigenvalues, according to their contribution to the explanation of the output-space variance. We observe that this way of doing is far less effective than the random projection techniques studied previously.



(a) Computing the impurity criterion on a dense Rademacher or on a subsample-Hadamard output sub-space is another efficient way to learn tree ensembles.



(b) Sparse random projections output sub-space yield better average precision than on the original output space.



(c) PCA compared with Gaussian subspaces.

Figure 5.4: “Delicious” dataset, $t = 100$, $k = \sqrt{p}$, $n_{\min} = 1$.

5.3.5 Learning stage computing times

Our implementation of the learning algorithms is based on the *scikit-learn* Python package version 0.14-dev (Buitinck et al., 2013; Pedregosa et al., 2011). To fix ideas about computing times, we report these obtained on a Mac Pro 4.1 with a dual Quad-Core Intel Xeon processor at 2.26 GHz, on the “Delicious” dataset. Matrix operation, such as random projections, are performed with the BLAS and the LAPACK from the Mac OS X *Accelerate* framework. Reported times are obtained by summing the user and sys time of the UNIX *time* utility.

The reported timings correspond to the following operation: (i) load the dataset in memory, (ii) execute the algorithm. All methods use the same code to build trees. In these conditions, learning a random forest on the original output space ($t = 100$, $n_{\min} = 1$, $k = \sqrt{d}$) takes 3348 s; learning the same model on a Gaussian output space of size $q = 25$ requires 311 s, while $q = 1$ and $q = 250$ take respectively 236 s and 1088 s. Generating a Gaussian sub-space of size $q = 25$ and projecting the output data of the training samples is done in less than 0.25 s, while $q = 1$ and $q = 250$ takes around 0.07 s and 1 s respectively. The time needed to compute the projections is thus negligible with respect to the time needed for the tree construction.

We see that a speed-up of an order of magnitude could be obtained, while at the same time preserving accuracy with respect to the baseline Random Forests method. Equivalently, for a fixed computing time budget, randomly projecting the output space allows to build more trees and thus to improve predictive performances with respect to standard Random Forests.

5.4 CONCLUSIONS

This chapter explores the use of random output space projections combined with tree-based ensemble methods to address large-scale multi-label classification problems. We study two algorithmic variants that either build a tree-based ensemble model on a single shared random subspace or build each tree in the ensemble on a newly drawn random subspace. The second approach is shown theoretically and empirically to always outperform the first in terms of accuracy. Experiments on 24 datasets show that on most problems, using gaussian projections allows to reduce very drastically the size of the output space, and therefore computing times, without affecting accuracy. Remarkably, we also show that by adjusting jointly the level of input and output randomizations and choosing appropriately the projection method, one could also improve predictive performance over the standard Random Forests, while still improving very significantly computing times. As future work, it would be very interesting to pro-

pose efficient techniques to automatically adjust these parameters, so as to reach the best tradeoff between accuracy and computing times on a given problem.

To the best of our knowledge, our work is the first to study random output projections in the context of multi-output tree-based ensemble methods. The possibility with these methods to relabel tree leaves with predictions in the original output space makes this combination very attractive. Indeed, unlike similar works with linear models (Cisse et al., 2013; Hsu et al., 2009), our approach only relies on Johnson-Lindenstrauss lemma, and not on any output sparsity assumption, and also does not require to use any output reconstruction method. Besides multi-label classification, we would like to test our method on other, not necessarily sparse, multi-output prediction problems.

6

GRADIENT BOOSTING WITH RANDOM OUTPUT PROJECTIONS FOR MULTI-LABEL AND MULTI-OUTPUTS REGRESSION TASKS

OUTLINE

We first formally adapt the gradient boosting ensemble method for multi-output supervised learning tasks such as multi-output regression and multi-label classification. We then propose to combine single random projections of the output space with gradient boosting on such tasks to adapt automatically to the output correlation structure. The idea of this method is to train each weak model on a single random projection of the output space and then to exploit the predictions of the resulting model to approximate the gradients of all other outputs. Through weak model sharing and random projection of the output space, we implicitly take into account the output correlations. We perform extensive experiments with these methods both on artificial and real problems using tree-based weak learners. Randomly projecting the output space shows to provide a better adaptation to different output correlation patterns and is therefore competitive with the best of the other methods in most settings. Thanks to the model sharing, the convergence speed is also faster, reducing the computing times to reach a specific accuracy.

This contribution is a joint work with Pierre Geurts and Louis Wehenkel from the University of Liège.

6.1 INTRODUCTION

Multi-output supervised learning aims to model the input-output relationship of a system from observations of input-output pairs whenever the output space is a vector of random variables. Multi-output classification and regression tasks have numerous applications in domains ranging from biology to multimedia.

The most straightforward way to address multi-output tasks is to apply standard single output methods separately and independently on each output. Although simple, this method, called binary relevance (Tsoumakas et al., 2009) in multi-label classification or single target (Spyromitros-Xioufifis et al., 2016) in multi-output regression, is often suboptimal as it does not exploit potential correlations that might exist between the outputs. For this reason, several approaches have been proposed in the literature that improve over binary relevance by exploiting output correlations. These approaches

include for example the explicit construction of the output dependency graph (Dembczynski et al., 2010; Gasse et al., 2015; Zhang and Zhang, 2010) or the sharing of models learnt for one output to the other outputs (Huang et al., 2012; Read et al., 2011; Yan et al., 2007). Our contribution falls into the latter category.

Classification and regression trees (Breiman et al., 1984) are popular supervised learning methods that provide state-of-the-art accuracy when exploited in the context of ensemble methods, namely Random forests (Breiman, 2001) and gradient boosting (Friedman, 2001). Classification and regression trees have been extended by several authors to the joint prediction of multiple outputs (see, e.g., Blockeel et al., 2000; Segal, 1992)). These extensions build a single tree to predict all outputs at once. They adapt the score measure used to assess splits during the tree growth to take into account all outputs and label each tree leaf with a vector of values, one for each output (see Section 3.5 for more information). Like standard classification or regression trees, multiple output trees can be exploited in the context of random forests (Barutcuoglu et al., 2006; Joly et al., 2014; Kocev et al., 2007, 2013; Segal and Xiao, 2011) or boosting (Geurts et al., 2007) ensembles, which often offer very significant accuracy improvements with respect to single trees. Multiple output trees have been shown to be competitive with other multiple output methods (Madjarov et al., 2012), but, to the best of our knowledge, it has not been studied as extensively in the context of gradient boosting.

Binary relevance / single target of single output tree models and multiple output tree models represent two extremes in terms of tree structure learning: the former builds a separate tree ensemble structure for each output, while the latter builds a single tree ensemble structure for all outputs. Building separate ensembles for each output may be rather inefficient when the outputs are strongly correlated. Correlations between the outputs could indeed be exploited either to reduce model complexity (by sharing the tree structures between several outputs) or to improve accuracy by regularization. Trying to fit a single tree structure for all outputs seems however counterproductive when the outputs are independent. Indeed, in the case of independent outputs, simultaneously fitting all outputs with a single tree structure may require a much more complex tree structure than the sum of the individual tree complexities required to fit the individual outputs. Since training a more complex tree requires a larger learning sample, multiple output trees are expected to be outperformed by binary relevance / single target in this situation.

In this chapter, we first formally adapt gradient boosting to multiple output tasks. We then propose a new method that aims at circumventing the limitations of both binary relevance / single target and multiple output methods, in the specific context of tree-based base-learners. Our method is an extension of gradient tree boosting

that can adapt itself to the presence or absence of correlations between the outputs. At each boosting iteration, a single regression tree structure is grown to fit a single random projection of the outputs, or more precisely, of their residuals with respect to the previously built models. Then, the predictions of this tree are fitted linearly to the current residuals of all the outputs (independently). New residuals are then computed taking into account the resulting predictions and the process repeats itself to fit these new residuals. Because of the linear fit, only the outputs that are correlated with the random projection at each iteration will benefit from a reduction of their residuals, while outputs that are independent of the random projection will remain mostly unaffected. As a consequence, tree structures will only be shared between correlated outputs as one would expect. Another variant that we explore consists in replacing the linear global fit by a relabelling of all tree leaves for each output in turn.

The chapter is structured as follows. We show how to extend the gradient boosting algorithms to multi-output tasks in Section 6.2. We provide for these algorithms a convergence proof on the training data and discuss the effect of the random projection of the output space. We study empirically the proposed approach in Section 6.3. Our first experiments compare the proposed approaches to binary relevance / single target on artificial datasets where the output correlation is known. We also highlight the effect of the choice and size of the random projection space. We finally carry out an empirical evaluation of these methods on 21 real-world multi-label and 8 multi-output regression tasks. We draw our conclusions in Section 6.4.

6.2 GRADIENT BOOSTING WITH MULTIPLE OUTPUTS

Starting from a multi-output loss, we show in Section 6.2.1 how to extend the standard gradient boosting algorithm to solve multi-output tasks, such as multi-output regression and multi-label classification, by exploiting existing weak model learners suited for multi-output prediction. In Section 6.2.2, we then propose to combine single random projections of the output space with gradient boosting to automatically adapt to the output correlation structure on these tasks. We discuss and compare the effect of the random projection of the output space in Section 6.2.3. We give a convergence proof on the training data for the proposed algorithms in Section 6.2.4.

6.2.1 *Standard extension of gradient boosting to multi-output tasks*

A loss function $\ell(y, y') \in \mathbb{R}^+$ computes the difference between a ground truth y and a model prediction y' . It compares scalars with single output tasks and vectors with multi-output tasks. The two most common regression losses are the square loss $\ell_{\text{square}}(y, y') =$

$\frac{1}{2}(y - y')^2$ and the absolute loss $\ell_{\text{absolute}}(y, y') = |y - y'|$. Their multi-output extensions are the ℓ_2 -norm and ℓ_1 -norm losses:

$$\ell_2(y, y') = \frac{1}{2} \|y - y'\|_{\ell_2}^2, \quad (6.1)$$

$$\ell_1(y, y') = \|y - y'\|_{\ell_1}. \quad (6.2)$$

In classification, the most commonly used loss to compare a ground truth y to the model prediction $f(x)$ is the 0–1 loss $\ell_{0-1}(y, y') = 1(y \neq y')$, where 1 is the indicator function. It has two standard multiple output extensions (i) the Hamming loss ℓ_{Hamming} and (ii) the subset 0–1 loss $\ell_{\text{subset } 0-1}$:

$$\ell_{\text{Hamming}}(y, y') = \sum_{j=1}^d 1(y_j \neq y'_j), \quad (6.3)$$

$$\ell_{\text{subset } 0-1}(y, y') = 1(y \neq y'). \quad (6.4)$$

Since these losses are discrete, they are non-differentiable and difficult to optimize. Instead, we propose to extend the logistic loss $\ell_{\text{logistic}}(y, y') = \log(1 + \exp(-2yy'))$ used for binary classification tasks to the multi-label case, as follows:

$$\ell_{\text{logistic}}(y, y') = \sum_{j=1}^d \log(1 + \exp(-2y_j y'_j)), \quad (6.5)$$

where we suppose that the d components y_j of the target output vector belong to $\{-1, 1\}$, while the d components y'_j of the predictions may belong to \mathbb{R} .

Given a training set $\mathcal{L} = ((x^i, y^i) \in \mathcal{X} \times \mathcal{Y})_{i=1}^n$ and one of these multi-output losses ℓ , we want to learn a model f_M expressed in the following form

$$f_M(x) = \rho_0 + \sum_{m=1}^M \rho_m \odot g_m(x), \quad (6.6)$$

where the terms g_m are selected within a hypothesis space \mathcal{H} of weak multi-output base-learners, the coefficients $\{\rho_m \in \mathbb{R}^d\}_{m=0}^M$ are d -dimensional vectors highlighting the contributions of each term g_m to the ensemble, and where the symbol \odot denotes the Hadamard product. Note that the prediction $f_M(x) \in \mathbb{R}^d$ targets the minimization of the chosen loss ℓ , but a transformation might be needed to have a prediction in \mathcal{Y} , e.g. we would apply the logit function to each output for the multi-output logistic loss to get a probability estimate of the positive classes.

The gradient boosting method builds such a model in an iterative fashion, as described in Algorithm 6.1, and discussed below.

Algorithm 6.1 Gradient boosting with multi-output regressor weak models.

```

1: function GB-MO( $\mathcal{L} = ((x^i, y^i))_{i=1}^n; \ell; \mathcal{H}; M$ )
2:    $f_0(x) = \rho_0 = \arg \min_{\rho \in \mathbb{R}^d} \sum_{i=1}^n \ell(y^i, \rho)$ 
3:   for  $m = 1$  to  $M$  do
4:     Compute the loss gradient for the learning set samples
       
$$g_m^i \in \mathbb{R}^d = [\nabla_{y'} \ell(y^i, y')]_{y'=f_{m-1}(x^i)} \quad \forall i \in \{1, \dots, n\}.$$

5:     Fit the negative loss gradient
       
$$g_m = \arg \min_{g \in \mathcal{H}} \sum_{i=1}^n \|-g_m^i - g(x^i)\|_{\ell_2}^2.$$

6:     Find an optimal step length in the direction of  $g_m$ 
       
$$\rho_m = \arg \min_{\rho \in \mathbb{R}^d} \sum_{i=1}^n \ell(y^i, f_{m-1}(x^i) + \rho \odot g_m(x^i)).$$

7:      $f_m(x) = f_{m-1}(x) + \rho_m \odot g_m(x)$ 
8:   end for
9:   return  $f_M(x)$ 
10: end function

```

To build the ensemble model, we first initialize it with the constant model defined by the vector $\rho_0 \in \mathbb{R}^d$ minimizing the multi-output loss ℓ (line 2):

$$\rho_0 = \arg \min_{\rho \in \mathbb{R}^d} \sum_{i=1}^n \ell(y^i, \rho). \quad (6.7)$$

At each subsequent iteration m , the multi-output gradient boosting approach adds a new multi-output weak model $g_m(x)$ with a weight ρ_m to the current ensemble model by approximating the minimization of the multi-output loss ℓ :

$$(\rho_m, g_m) = \arg \min_{(\rho, g) \in \mathbb{R}^d \times \mathcal{H}} \sum_{i=1}^n \ell(y^i, f_{m-1}(x^i) + \rho \odot g(x^i)). \quad (6.8)$$

To approximate Equation 6.8, it first fits a multi-output weak model g_m to model the negative gradient g_m^i of the multi-output loss ℓ

$$g_m^i \in \mathbb{R}^d = [\nabla_{y'} \ell(y^i, y')]_{y'=f_{m-1}(x^i)} \quad (6.9)$$

associated to each sample $i \in \mathcal{L}$ of the training set, by minimizing the ℓ_2 -loss:

$$g_m = \arg \min_{g \in \mathcal{H}} \sum_{i=1}^n \|-g_m^i - g(x^i)\|_{\ell_2}^2. \quad (6.10)$$

It then computes an optimal step length vector $\rho_m \in \mathbb{R}^d$ in the direction of the weak model g_m to minimize the multi-output loss ℓ :

$$\rho_m = \arg \min_{\rho \in \mathbb{R}^d} \sum_{i=1}^n \ell(y^i, f_{m-1}(x^i) + \rho \odot g_m(x^i)). \quad (6.11)$$

6.2.2 Adapting to the correlation structure in the output-space

Binary relevance / single target of gradient boosting models and gradient boosting of multi-output models (Algorithm 6.1) implicitly target two extreme correlation structures. On the one hand, binary relevance / single target predicts all outputs independently, thus assuming that outputs are not correlated. On the other hand, gradient boosting of multi-output models handles them all together, thus assuming that they are all correlated. Both approaches thus exploit the available dataset in a rather biased way. To remove this bias, we propose a more flexible approach that can adapt itself automatically to the correlation structure among output variables.

Our idea is that a weak learner used at some step of the gradient boosting algorithm could be fitted on a single random projection of the output space, rather than always targeting simultaneously all outputs or always targeting a single a priori fixed output.

We thus propose to first generate at each iteration of the boosting algorithm one random projection vector of size $\phi_m \in \mathbb{R}^{1 \times d}$. The weak learner is then fitted on the projection of the current residuals according to ϕ_m reducing dimensionality from d outputs to a single output. A weight vector $\rho_m \in \mathbb{R}^d$ is then selected to minimize the multi-output loss ℓ . The whole approach is described in Algorithm 6.2. If the loss is decomposable, non zero components of the weight vector ρ_m highlight the contribution of the current m -th model to the overall loss decrease. Note that sign flips due to the projection are taken into account by the additive weights ρ_m . A single output regressor can now handle multi-output tasks through a sequence of single random projections.

The prediction of an unseen sample x by the model produced by Algorithm 6.2 is now given by

$$f(x) = \rho_0 + \sum_{m=1}^M \rho_m g_m(x), \quad (6.12)$$

where $\rho_0 \in \mathbb{R}^d$ is a constant prediction, and the coefficients $\{\rho_m \in \mathbb{R}^d\}_{m=1}^M$ highlight the contribution of each model g_m to the ensemble. Note that it is different from Equation 6.6 (no Hadamard product), since here the weak models g_m produce single output predictions.

Whenever we use decision trees as models, we can grow the tree structure on any output space and then (re)label it in another one as in Chapter 5 Section 5.1.1 by (re)propagating the training samples

Algorithm 6.2 Gradient boosting on randomly projected residual spaces.

- 1: **function** GB-RPO($\mathcal{L} = ((x^i, y^i))_{i=1}^n; \ell; \mathcal{H}; M$)
- 2: $f_0(x) = \rho_0 = \arg \min_{\rho \in \mathbb{R}^d} \sum_{i=1}^n \ell(y^i, \rho)$
- 3: **for** $m = 1$ to M **do**
- 4: Compute the loss gradient for the learning set samples

$$g_m^i \in \mathbb{R}^d = [\nabla_{y'} \ell(y^i, y')]_{y' = f_{m-1}(x^i)} \quad \forall i \in \{1, \dots, n\}.$$

- 5: Generate a random projection $\phi_m \in \mathbb{R}^{1 \times d}$.
- 6: Fit the projected loss gradient

$$g_m = \arg \min_{g \in \mathcal{H}} \sum_{i=1}^n (-\phi_m g_m^i - g(x^i))^2.$$

- 7: Find an optimal step length in the direction of g_m .

$$\rho_m = \arg \min_{\rho \in \mathbb{R}^d} \sum_{i=1}^n \ell(y^i, f_{m-1}(x^i) + \rho g_m(x^i)),$$

- 8: $f_m(x) = f_{m-1}(x) + \rho_m g_m(x)$
 - 9: **end for**
 - 10: **return** $f_M(x)$
 - 11: **end function**
-

in the tree structure. This idea of leaf relabelling could be readily applied to Algorithm 6.2 leading to Algorithm 6.3. After fitting the decision tree on the random projection(s) and before optimizing the additive weights ρ_m , we relabel the tree structure in the original residual space (line 7). More precisely, each leaf is labelled by the average unprojected residual vector of all training examples falling into that leaf. The prediction of an unseen sample is then obtained with Equation 6.6 as for Algorithm 6.1. We will investigate whether it is better or not to relabel the decision tree structure in the experimental section. Note that Algorithm 6.3 can be straightforwardly used in a multiple random projection context ($q \geq 1$) using a random projection matrix $\phi_m \in \mathbb{R}^{q \times d}$. The resulting algorithm with arbitrary q corresponds to the application to gradient boosting of the idea explored in Chapter 5 in the context of random forests. We will study in Section 6.2.3.3 and Section 6.3.3.2 the effect of the size of the projected space q .

Algorithm 6.3 Gradient boosting on randomly projected residual spaces with relabelled decision trees as weak models.

- 1: **function** GB-RELABEL-RPO($\mathcal{L} = ((x^i, y^i))_{i=1}^n; \ell; \mathcal{H}; M; q$)
- 2: $f_0(x) = \rho_0 = \arg \min_{\rho \in \mathbb{R}^d} \sum_{i=1}^n \ell(y^i, \rho)$
- 3: **for** $m = 1$ to M **do**
- 4: Compute the loss gradient for the learning set samples

$$g_m^i \in \mathbb{R}^d = [\nabla_{y'} \ell(y^i, y')]_{y' = f_{m-1}(x^i)} \quad \forall i \in \{1, \dots, n\}.$$

- 5: Generate a random projection $\phi_m \in \mathbb{R}^{q \times d}$.
- 6: Fit a single-output tree g_m on the projected negative loss gradients

$$g_m = \arg \min_{g \in \mathcal{H}} \sum_{i=1}^n \|\phi_m g_m^i - g(x^i)\|_{\ell_2}^2.$$

- 7: Relabel each leaf of the tree g_m in the original (unprojected) residual space, by averaging at each leaf the g_m^i vectors of all examples falling into that leaf.
- 8: Find an optimal step length in the direction of g'_m .

$$\rho_m = \arg \min_{\rho \in \mathbb{R}^d} \sum_{i=1}^n \ell(y^i, f_{m-1}(x^i) + \rho \odot g'_m(x^i)).$$

- 9: $f_m(x) = f_{m-1}(x) + \rho_m \odot g'_m(x)$
 - 10: **end for**
 - 11: **return** $f_M(x)$
 - 12: **end function**
-

To the three presented algorithms, we also add a constant learning rate $\mu \in (0, 1]$ to shrink the size of the gradient step ρ_m in the residual space. Indeed, for a given weak model space \mathcal{H} and a loss ℓ , optimizing both the learning rate μ and the number of steps M typically improves generalization performance.

6.2.3 Effect of random projections

Randomly projecting the output space in the context of the gradient boosting approach has two direct consequences: (i) it strongly reduces the size of the output space, and (ii) it randomly combines several outputs. We will consider here the following random projection matrices $\phi \in \mathbb{R}^{q \times d}$ ordered from the sparsest to the densest ones:

- **Random output subsampling matrices** is obtained by sampling random lines from the identity matrix.
- **(Sparse) Rademacher matrices** is obtained by drawing its elements in $\left\{-\sqrt{\frac{s}{q}}, 0, \sqrt{\frac{s}{q}}\right\}$ with probability $\left\{\frac{1}{2s}, 1 - \frac{1}{s}, \frac{1}{2s}\right\}$, where $1/s \in (0, 1]$ controls the sparsity of ϕ . With $s = 1$, we have (dense) **Rademacher random projections**. If $s = 3$, we will call them **Achlioptas random projections** (Achlioptas, 2003). When $s = \sqrt{d}$, we will say that we have **sparse random projections** as in (Li et al., 2006).
- **Gaussian matrices** are obtained by drawing their elements *i.i.d.* in $\mathcal{N}(0, 1/q)$.

We discuss in more details the random sub-sampling projection in Section 6.2.3.1 and the impact of the density of random projection matrices in Section 6.2.3.2. We study the benefit to use more than a single random projection of the output space ($q > 1$) in Section 6.2.3.3. We highlight the difference in model representations between tree ensemble techniques, *i.e.* the gradient tree boosting approaches and the random forest approaches, in Section 6.2.3.4.

6.2.3.1 ℓ_2 -norm loss and random output sub-sampling

The gradient boosting method has an analytical solution when the loss is the square loss or its extension the ℓ_2 -norm loss $\ell_2(y, y') = \frac{1}{2}\|y - y'\|^2$:

- The constant model f_0 minimizing this loss is the average output value of the training set given by

$$f_0(x) = \rho_0 = \arg \min_{\rho \in \mathbb{R}^d} \sum_{i=1}^n \frac{1}{2} \|y^i - \rho\|_{\ell_2}^2 = \frac{1}{n} \sum_{i=1}^n y^i. \quad (6.13)$$

- The gradient of the ℓ_2 -norm loss for the i -th sample is the difference between the ground truth y^i and the prediction of the ensemble f at the current step m ($\forall i \in \{1, \dots, n\}$):

$$g_m^i = [\nabla_{y'} \ell(y^i, y')]_{y'=f_{m-1}(x^i)} = y^i - f_{m-1}(x^i). \quad (6.14)$$

- Once a new weak estimator g_m has been fitted on the loss gradient g_m^i or the projected gradient $\phi_m g_m^i$ with or without relabelling, we have to optimize the multiplicative weight vector ρ_m of the new weak model in the ensemble. For Algorithm 6.1 and Algorithm 6.3 that exploit multi-output weak learners, this amounts to

$$\rho_m = \arg \min_{\rho \in \mathbb{R}^d} \sum_{i=1}^n \frac{1}{2} \|y^i - f_m(x^i) - \rho \odot g_m(x^i)\|^2 \quad (6.15)$$

$$= \arg \min_{\rho \in \mathbb{R}^d} \sum_{i=1}^n \frac{1}{2} \|g_m^i - \rho \odot g_m(x^i)\|^2 \quad (6.16)$$

which has the following solution:

$$\rho_{m,j} = \frac{\sum_{i=1}^n g_{m,j}^i g_m(x^i)_j}{\sum_{i=1}^n g_m(x^i)_j} \quad \forall j \in \{1, \dots, d\}. \quad (6.17)$$

For Algorithm 6.2, we have to solve

$$\rho_m = \arg \min_{\rho \in \mathbb{R}^d} \sum_{i=1}^n \frac{1}{2} \|y^i - f_m(x^i) - \rho g_m(x^i)\|^2 \quad (6.18)$$

$$= \arg \min_{\rho \in \mathbb{R}^d} \sum_{i=1}^n \frac{1}{2} \|g_m^i - \rho g_m(x^i)\|^2 \quad (6.19)$$

which has the following solution

$$\rho_{m,j} = \frac{\sum_{i=1}^n g_{m,j}^i g_m(x^i)}{\sum_{i=1}^n g_m(x^i)} \quad \forall j \in \{1, \dots, d\}. \quad (6.20)$$

From Equation 6.17 and Equation 6.20, we have that the weight $\rho_{m,j}$ is proportional to the correlation between the loss gradient of the output j and the weak estimator g_m . If the model g_m is independent of the output j , the weight $\rho_{m,j}$ will be close to zero and g_m will thus not contribute to the prediction of this output. On the opposite, a high magnitude of $|\rho_{m,j}|$ means that the model g_m is useful to predict the output j .

If we subsample the output space at each boosting iteration (Algorithm 6.2 with random output sub-sampling), the weight $\rho_{m,j}$ is then proportional to the correlation between the model fitted on the sub-sampled output and the output j . If correlations exist between the outputs, the optimization of the constant ρ_m allows to share the

trained model at the m -th iteration on the sub-sampled output to all the other outputs. In the extreme case where all outputs are independent given the inputs, the weight ρ_m is expected to be nearly zero for all outputs except for the sub-sampled output, and Algorithm 6.2 would be equivalent to the binary relevance / single target approach. If all outputs are strictly identical, the elements of the constant vector ρ_m would have the same value, and Algorithm 6.2 would be equivalent to the multi-output gradient boosting approach (Algorithm 6.1). Algorithm 6.2 would also produce in this case the exact same model as binary relevance / single target approach asymptotically but it would require d times less trees to reach similar performance, as each tree would be shared by all d outputs.

Algorithm 6.3 with random output sub-sampling is a gradient boosting approach fitting one decision tree at each iteration on a random output space and relabelling the tree in the original output space. The leaf relabelling procedure minimizes the ℓ_2 -norm loss over the training samples by averaging the output values of the samples reaching the corresponding leaves. In this case, the optimization of the weight ρ_m is unnecessary, as it would lead to an all ones vector. For similar reasons if the multi-output gradient boosting method (Algorithm 6.1) uses decision trees as weak estimators, the weight ρ_m is also an all ones vector as the leaf predictions already minimize the ℓ_2 -norm loss. The difference between these two algorithms is that Algorithm 6.3 grows trees using a random output at each iteration instead of all of them with Algorithm 6.1.

6.2.3.2 Density of the random projections

In Chapter 5, we have combined the random forest method with a wide variety of random projection schemes. While the algorithms presented in this chapter were originally devised with random output sub-sampling in mind (see Section 6.2.3.1), it seems natural to also combine the proposed approaches with random projection schemes such as Gaussian random projections or (sparse) Rademacher random projections.

With random output sub-sampling, the projection matrix $\phi_m \in \mathbb{R}^{1 \times d}$ is extremely sparse as only one element is non zero. With denser random projections, the weak estimators of Algorithm 6.2 and Algorithm 6.3 are fitted on the projected gradient loss $\{(x^i, \phi_m g_m^i)\}_{i=1}^n$. It means that a weak estimator g_m is trying to model the direction of a weighted combination of the gradient loss.

Otherwise said, the weak model fitted at the m -th step approximates a projection ϕ_m of the gradient losses given the input vector. We can interpret the weight $\rho_{m,j}$ when minimizing the ℓ_2 -norm loss as the correlation between the output j and a weighted approximation of the output variables ϕ_m . With an extremely sparse projection having only one non zero element, we have the situation described in

the previous section. If we have two non zero elements, we have the following extreme cases: (i) both combined outputs are identical and (ii) both combined outputs are independent given the inputs. In the first situation, the effect is identical to the case where we sub-sample only one output. In the second situation, the weak model makes a compromise between the independent outputs given by ϕ_m . Between those two extremes, the loss gradient direction ϕ_m approximated by the weak model is useful to predict both outputs. The random projection of the output space will indeed prevent over-fitting by inducing some variance in the learning process. The previous reasoning can be extended to more than two output variables.

Dense random projection schemes, such as Gaussian random projection, consider a higher number of outputs together and is hoped to speed up convergence by increasing the correlation between the fitted tree in the projected space and the residual space. Conversely, sparse random projections, such as random output sub-sampling, make the weak model focus on few outputs.

6.2.3.3 *Gradient tree boosting and multiple random projections*

The gradient boosting multi-output strategy combining random projections and tree relabelling (Algorithm 6.3) can use random projection matrices $\phi_m \in \mathbb{R}^{q \times d}$ with more than one line ($q \geq 1$).

The weak estimators are multi-output regression trees using the variance as impurity criterion to grow their tree structures. With an increasing number of projections q , we have the theoretical guarantee (see Chapter 5) that the variance computed in the projected space is an approximation of the variance in the original output space.

When the projected space is of infinite size $q \rightarrow \infty$, the decision trees grown on the original space or on the projected space are identical as the approximation of the variance is exact. We thus have that Algorithm 6.3 is equivalent to the gradient boosting with multi-output regression tree method (Algorithm 6.1).

Whenever the projected space is of finite size ($q < \infty$), Algorithm 6.3 is thus an approximation of Algorithm 6.1. We study empirically the effect of the number of projections q in Algorithm 6.3 in Section 6.3.3.2.

6.2.3.4 *Representation bias of decision tree ensembles*

Random forests and gradient tree boosting build an ensemble of trees either independently or sequentially, and thus offer different bias/variance tradeoffs. The predictions of all these ensembles can be expressed as a weighted combination of the ground truth outputs of the training set samples. In the present section, we discuss the differences between single tree models, random forest models and gradient tree boosting models in terms of the representation biases of the obtained

models. We also highlight the differences between single target models and multi-output tree models.

SINGLE TREE MODELS. The prediction of a regression tree learner can be written as a weighted linear combination of the training samples $\mathcal{L} = \{(x^i, y^i) \in \mathcal{X} \times \mathcal{Y}\}_{i=1}^n$. We associate to each training sample (x^i, y^i) a weight function $w^i: \mathcal{X} \rightarrow \mathbb{R}$ which gives the contribution of a ground truth y^i to predict an unseen sample x . The prediction of a *single output tree* f is given by

$$f(x) = \sum_{i=1}^n w^i(x) y^i. \quad (6.21)$$

The weight function $w^i(x)$ is non zero if both the samples (x^i, y^i) and the unseen sample x reach the same leaf of the tree. If both (x^i, y^i) and x end up in the same leaf of the tree, $w^i(x)$ is equal to the inverse of the number of training samples reaching that leaf. The weight $w^i(x)$ can thus be rewritten as $k(x^i, x)$ and the function $k(\cdot, \cdot)$ is actually a positive semi-definite kernel (Geurts et al., 2006a).

We can also express multi-output models as a weighted sum of the training samples. With a *single target regression tree*, we have an independent weight function w_j^i for each sample of the training set and each output as we fit one model per output. The prediction of this model for output j is given by:

$$f(x)_j = \sum_{i=1}^n w_j^i(x) y_j^i. \quad (6.22)$$

With a *multi-output regression tree*, the decision tree structure is shared between all outputs so we have a single weight function w^i for each training sample:

$$f(x)_j = \sum_{i=1}^n w^i(x) y_j^i. \quad (6.23)$$

RANDOM FOREST MODELS. If we have a *single target random forest model*, the prediction of the j -th output combines the predictions of the M models of the ensemble in the following way:

$$f(x)_j = \frac{1}{M} \sum_{m=1}^M \sum_{i=1}^n w_{m,j}^i(x) y_j^i, \quad (6.24)$$

with one weight function $w_{m,j}^i$ per tree, sample and output. We note that we can combine the weights of the individual trees into a single one per sample and per output

$$w_j^i(x) = \frac{1}{M} \sum_{m=1}^M w_{m,j}^i(x). \quad (6.25)$$

The prediction of the j -th output for an ensemble of independent models has the same form as a single target regression tree model:

$$f(x)_j = \frac{1}{M} \sum_{m=1}^M \sum_{i=1}^n w_{m,j}^i(x) y_j^i = \sum_{i=1}^n w_j^i(x) y_j^i. \quad (6.26)$$

We can repeat the previous development with a *multi-output random forest model*. The prediction for the j -th output of an unseen sample x combines the predictions of the M trees:

$$f(x)_j = \frac{1}{M} \sum_{m=1}^M \sum_{i=1}^n w_m^i(x) y_j^i = \sum_{i=1}^n w^i(x) y_j^i \quad (6.27)$$

with

$$w^i(x) = \frac{1}{M} \sum_{m=1}^M w_m^i(x). \quad (6.28)$$

With this framework, the prediction of an ensemble model has the same form as the prediction of a single constituting tree.

GRADIENT TREE BOOSTING MODELS. The prediction of a *single output gradient boosting tree ensemble* is given by

$$f(x) = \rho_0 + \sum_{m=1}^M \mu \rho_m g_m(x), \quad (6.29)$$

but also as

$$f(x) = \sum_{m=1}^M \sum_{i=1}^n w_m^i(x) y^i = \sum_{i=1}^n w^i(x) y^i, \quad (6.30)$$

where the weight $w^i(x)$ takes into account the learning rate μ , the prediction of all tree models g_m and the associated ρ_m . Given the similarity between gradient boosting prediction and random forest model, we deduce that the *single target gradient boosting tree ensemble* has the form of Equation 6.22 and that *multi-output gradient tree boosting* (Algorithm 6.1) and *gradient boosting tree with projection of the output space and relabelling* (Algorithm 6.3) has the form of Equation 6.23.

However, we note that the prediction model of the *gradient tree boosting with random projection of the output space* (Algorithm 6.2) is not given by Equation 6.22 and Equation 6.23 as the prediction of a single output j can combine the prediction of all d outputs. More formally, the prediction of the j -th output is given by:

$$f(x)_j = \sum_{m=1}^M \sum_{i=1}^n \sum_{k=1}^d w_{m,j,k}^i(x) y_k^i \quad (6.31)$$

where the weight function $w_{m,j,k}^i$ takes into account the contribution of the m -th model fitted on a random projection ϕ_m of the output space to predict the j -th output using the k -th outputs and the i -th sample. The triple summation can be simplified by using a single weight to summarize the contribution of all M models:

$$f(x)_j = \sum_{m=1}^M \sum_{i=1}^n \sum_{k=1}^d w_{m,j,k}^i(x) y_k^i = \sum_{i=1}^n \sum_{k=1}^d w_{j,k}^i(x) y_k^i. \quad (6.32)$$

Between the studied methods, we can distinguish three groups of multi-output tree models. The first one considers that all outputs are independent as with binary relevance / single target trees, random forests or gradient tree boosting models. The second group with multi-output random forests, gradient boosting of multi-output tree and gradient boosting with random projection of the output space and relabelling share the tree structures between all outputs, but the leaf predictions are different for each output. The last and most flexible group is the gradient tree boosting with random projection of the output space sharing both the tree structures and the leaf predictions. We will highlight in the experiments the impact of these differences in representation biases.

6.2.4 Convergence when $M \rightarrow \infty$

Similarly to (Geurts et al., 2007), we can prove the convergence of the training-set loss of the gradient boosting with multi-output models (Algorithm 6.1), and gradient boosting on randomly projected spaces with (Algorithm 6.2) or without relabelling (Algorithm 6.3).

Since the loss function is lower-bounded by 0, we merely need to show that the loss ℓ is non-increasing on the training set at each step m of the gradient boosting algorithm.

For Algorithm 6.1 and Algorithm 6.3, we note that

$$\begin{aligned} \sum_{i=1}^n \ell(y^i, f_m(x^i)) &= \min_{\rho \in \mathbb{R}^d} \sum_{i=1}^n \ell(y^i, f_{m-1}(x^i) + \rho \odot g_m(x^i)) \\ &\leq \sum_{i=1}^n \ell(y^i, f_{m-1}(x^i)). \end{aligned} \quad (6.33)$$

and the learning-set loss is hence non increasing with M if we use a learning rate $\mu = 1$. If the loss $\ell(y, y')$ is convex in its second argument y' (which is the case for those loss-functions that we use in practice), then this convergence property actually holds for any value $\mu \in (0; 1]$ of the learning rate. Indeed, we have

$$\sum_{i=1}^n \ell(y^i, f_{m-1}(x^i))$$

$$\begin{aligned}
&\geq (1 - \mu) \sum_{i=1}^n \ell(y^i, f_{m-1}(x^i)) + \mu \sum_{i=1}^n \ell(y^i, f_{m-1}(x^i) + \rho_m \odot g_m(x^i)) \\
&\geq \sum_{i=1}^n \ell(y^i, f_{m-1}(x^i) + \mu \rho_m \odot g_m(x^i)).
\end{aligned}$$

given Equation 6.33 and the convexity property.

For Algorithm 6.2, we have a weak estimator g_m fitted on a single random projection of the output space ϕ_m with a multiplying constant vector $\rho_m \in \mathbb{R}^d$, and we have:

$$\begin{aligned}
\sum_{i=1}^n \ell(y^i, f_m(x^i)) &= \min_{\rho \in \mathbb{R}^d} \sum_{i=1}^n \ell(y^i, f_{m-1}(x^i) + \rho g_m(x^i)) \\
&\leq \sum_{i=1}^n \ell(y^i, f_{m-1}(x^i)). \tag{6.34}
\end{aligned}$$

and the error is also non increasing for Algorithm 6.2, under the same conditions as above.

The previous development shows that Algorithm 6.1, Algorithm 6.2 and Algorithm 6.3 are converging on the training set for a given loss ℓ . The binary relevance / single target of gradient boosting regression trees admits a similar convergence proof. We expect however the convergence speed of the binary relevance / single target to be lower assuming that it fits one weak estimator for each output in a round robin fashion.

6.3 EXPERIMENTS

We describe the experimental protocol in Section 6.3.1. Our first experiments in Section 6.3.2 illustrate the multi-output gradient boosting methods on synthetic datasets where the output correlation structure is known. The effect of the choice and / or the number of random projections of the output space is later studied for Algorithm 6.2 and Algorithm 6.3 in Section 6.3.3. We compare multi-output gradient boosting approaches and multi-output random forest approaches in Section 6.3.4 over 29 real multi-label and multi-output datasets.

6.3.1 Experimental protocol

We describe the metrics used to assess the performance of the supervised learning algorithms in Section 6.3.1.1. The protocol used to optimize hyper-parameters is given in Section 6.3.1.2.

Note that the datasets used in the following experiments are described in Appendix A. Whenever the number of testing samples is not given, we use half of the data as training set and half of the data as testing set.

6.3.1.1 Accuracy assessment protocol

We assess the accuracy of the predictors on a test set using the “Label Ranking Average Precision (LRAP)” (defined in Section 2.4.3) for multi-label classification tasks and the “macro- r^2 score” (defined in Section 2.4.4) for multi-output regression tasks.

6.3.1.2 Hyper-parameter optimization protocol

The hyper-parameters of the supervised learning algorithms are optimized as follows: we define an hyper-parameter grid and the best hyper-parameter set is selected using 20% of the training samples as a validation set. The results shown are averaged over five random split of the dataset while preserving the training-testing set size ratio.

For the boosting ensembles, we optimize the learning rate μ among $\{1., 0.5, 0.2, 0.1, 0.05, 0.02, 0.01\}$ and use decision trees as weak models whose hyper-parameters are also optimized: the number of features drawn at each node k during the tree growth is selected among $k \in \{\sqrt{p}, 0.1p, 0.2p, 0.5p, p\}$, the maximum number of tree leaves n_{\max_leaves} grown in best-first fashion is chosen among $n_{\max_leaves} \in \{2, \dots, 8\}$. Note that a decision tree with $n_{\max_leaves} = 2$ and $k = p$ is called a stump. We add new weak models to the ensemble by minimizing either the square loss or the absolute loss (or their multi-output extensions) in regression and either the square loss or the logistic loss (or their multi-output extensions) in classification, the choice of the loss being an additional hyper-parameter tuned on the validation set.

We also optimize the number of boosting steps n_{iter} of each gradient boosting algorithm over the validation set. However note that the number of steps has a different meaning depending on the algorithm. For binary relevance / single target gradient boosting, the number of boosting steps n_{iter} gives the number of weak models fitted per output. The implemented algorithm here fits weak models in a round robin fashion over all outputs. For all other (multi-output) methods, the number of boosting steps n_{iter} is the total number of weak models for all outputs as only one model is needed to fit all outputs. The computing time of one boosting iteration is thus different between the approaches. We will set the budget, the maximal number of boosting steps n_{iter} , for each algorithm to $n_{iter} = 10000$ on synthetic experiments (see Section 6.3.2) so that the performance of the estimator is not limited by the computational power. On the real datasets however, this setting would have been too costly. We decided instead to limit the computing time allocated to each gradient boosting algorithm on each classification (resp. regression) problem to $100 \times T$ (resp. $500 \times T$), where T is the time needed on this specific problem for one iteration of multi-output gradient boosting (Algorithm 6.1) with stumps and the ℓ_2 -norm loss. The maximum number of iterations, n_{iter} , is thus

set independently for each problem and each hyper-parameter setting such that this time constraint is satisfied. As a consequence, all approaches thus receive approximately the same global time budget for model training and hyper-parameter optimization.

For the random forest algorithms, we use the default hyper-parameter setting suggested in (Hastie et al., 2009), which corresponds in classification to 100 totally developed trees with $k = \sqrt{p}$ and in regression to 100 trees with $k = p/3$ and a minimum of 5 samples to split a node ($n_{\min} = 5$).

The base learner implementations are based on the random-output-trees¹ (Joly et al., 2014) version 0.1 and on the scikit-learn (Buitinck et al., 2013; Pedregosa et al., 2011) of version 0.16 Python package. The algorithms presented in this chapter will be provided in random-output-trees version 0.2.

6.3.2 Experiments on synthetic datasets with known output correlation structures

We study here the proposed boosting approaches on synthetic datasets whose output correlation structures are known. The datasets are first presented in Section 6.3.2.1. We then compare on these datasets multi-output gradient boosting approaches in terms of their convergence speed in Section 6.3.2.2 and in terms of their best performance whenever hyper-parameters are optimized in Section 6.3.2.3.

6.3.2.1 Synthetic datasets

To illustrate multi-output boosting strategies, we use three synthetic datasets with a specific output structure: (i) chained outputs, (ii) totally correlated outputs and (iii) fully independent outputs. Those tasks are derived from the **friedman1** regression dataset which consists in solving the following single target regression task (Friedman, 1991)

$$f(x) = 10 \sin(\pi x_1 x_2) + 20(x_3 - 0.5)^2 + 10x_4 + 5x_5 y = f(x) + \epsilon \quad (6.35)$$

with $x \in \mathbb{R}^5 \sim \mathcal{N}(0; I_5)$ and $\epsilon \sim \mathcal{N}(0; 1)$ where I_5 is an identity matrix of size 5×5 .

The **friedman1-chain** problem consists in d regression tasks forming a chain obtained by cumulatively adding independent standard Normal noise. We draw samples from the following distribution

$$y_1 = f(x) + \epsilon_1, \quad (6.36)$$

$$y_j = y_{j-1} + \epsilon_j \quad \forall j \in \{2, \dots, d\} \quad (6.37)$$

¹ <https://github.com/arjoly/random-output-trees>

with $x \sim \mathcal{N}(0; I_5)$ and $\epsilon \sim \mathcal{N}(0; I_d)$. Given the chain structure, the output with the least amount of noise is the first one of the chain and averaging a subset of the outputs would not lead to any reduction of the output noise with respect to the first output, since total noise variance accumulates more than linearly with the number of outputs. The optimal multi-output strategy is thus to build a model using only the first output and then to replicate the prediction of this model for all other outputs.

The **friedman1-group** problem consists in solving d regression tasks simultaneously obtained from one **friedman1** problem without noise where an independent normal noise is added. Given $x \sim \mathcal{N}(0; I_5)$ and $\epsilon \sim \mathcal{N}(0; I_d)$, we have to solve the following task:

$$y_j = f(x) + \epsilon_j \quad \forall j \in \{1, \dots, d\}. \quad (6.38)$$

If the output-output structure is known, the additive noises $\epsilon_j, \forall j \in \{1, \dots, d\}$, can be filtered out by averaging all outputs. The optimal strategy to address this problem is thus to train a single output regression model to fit the average output. Predictions on unseen data would be later done by replicating the output of this model for all outputs.

The **friedman1-ind** problem consists in d independent **friedman1** tasks. Drawing samples from $x \sim \mathcal{N}(0; I_{5d})$ and $\epsilon \sim \mathcal{N}(0; I_d)$, we have

$$y_j = f(x_{5j+1:5j+5}) + \epsilon_j \quad \forall j \in \{1, \dots, d\}. \quad (6.39)$$

where $x_{5j+1:5j+5}$ is a slice of feature vector from feature $5j + 1$ to $5j + 5$. Since all outputs are independent, the best multi-output strategy is single target: one independent model fits each output

For each multi-output **friedman** problem, we consider 300 training samples, 4000 testing samples and $d = 16$ outputs.

6.3.2.2 Convergence with known output correlation structure

We first study the macro- r^2 score convergence as a function of time (see Figure 6.1) for three multi-output gradient boosting strategies: (i) single target of gradient tree boosting (st-gbrt), (ii) gradient boosting with multi-output regression tree (gbmort, Algorithm 6.1) and (iii) gradient boosting with output subsampling of the output space (gbrt-rpo-subsample, Algorithm 6.2). We train each boosting algorithm on the three **friedman1** artificial datasets with the same set of hyperparameters: a learning rate of $\mu = 0.1$ and stumps as weak estimators (a decision tree with $k = p$, $n_{\max_leaves} = 2$) while minimizing the square loss.

On the **friedman1-chain** (see Figure 6.1a) and **friedman1-group** (see Figure 6.1b), gbmort and gbrt-rpo-subsampled converge more than 100 times faster (note the logarithmic scale of the abscissa) than single target. Furthermore, the optimal macro- r^2 is slightly better for

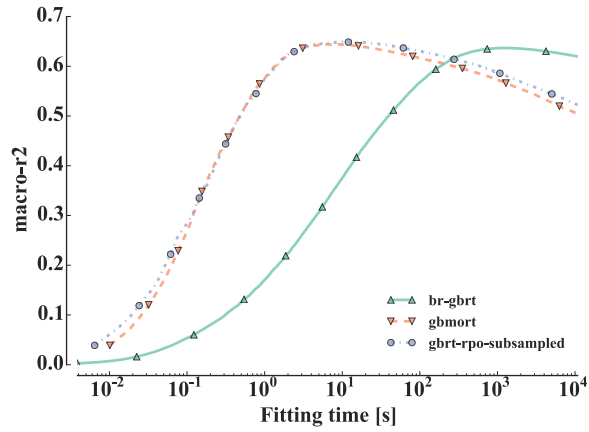
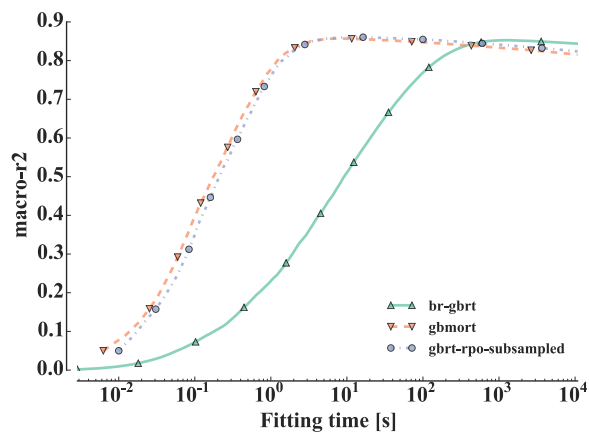
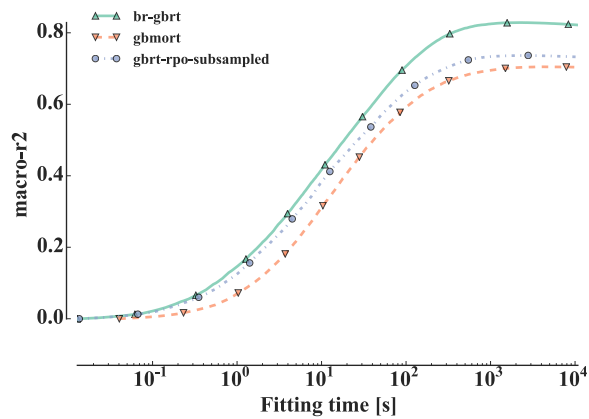
(a) Friedman₁-chain(b) Friedman₁-group(c) Friedman₁-ind

Figure 6.1: The convergence speed and the optimum reached are affected by the output correlation structure. The gbmort and gbrt-rpo-subsampled algorithms both exploit output correlations, which yields faster convergence and slightly better performance than st-gbrt on friedman₁-chain and friedman₁-group. However, st-gbrt converges to a better optimum than gbmort and gbrt-rpo-subsample when there is no output correlation as in friedman₁-ind. (Model parameters: $k = p$, $n_{\max_leaves} = 2$, $\mu = 0.1$)

gbmort and gbrt-rpo-subsampled than st-gbrt. Since all outputs are correlated on both datasets, gbmort and gbrt-rpo-subsampled are exploiting the output structure to have faster convergence. The gbmort method exploits the output structure by filtering the output noise as the stump fitted at each iteration is the one that maximizes the reduction of the average output variance. By contrast, gbrt-rpo-subsampled detects output correlations by optimizing the ρ_m constant and then shares the information obtained by the current weak model with all other outputs.

On the `friedman1-ind` dataset (see Figure 6.1c), all three methods converge at the same rate. However, the single target strategy converges to a better optimum than gbmort and gbrt-rpo-subsampled. Since all outputs are independent, single target enforces the proper correlation structure (see Figure 6.1c). The gbmort method has the worst performance as it assumes the wrong set of hypotheses. The gbrt-rpo-subsampled method pays the price of its flexibility by overfitting the additive weight associated to each output, **but less than gbmort**.

This experiment confirms that enforcing the right correlation structure yields faster convergence and the best accuracy. Nevertheless, the output structure is unknown in practice. We need flexible approaches such as gbrt-rpo-subsampled that automatically detects and exploits the correlation structure.

6.3.2.3 Performance and output modeling assumption

The presence or absence of structures among the outputs have shown to affect the convergence speed of multi-output gradient boosting methods. As discussed in (Cheng et al., 2010), we talk about conditionally independent outputs when:

$$P(y_1, \dots, y_q | x) = P(y_1 | x) \cdots P(y_q | x)$$

and about unconditionally independent outputs when:

$$P(y_1, \dots, y_q) = P(y_1) \cdots P(y_q).$$

When the outputs are not conditionally independent and the loss function can not be decomposed over the outputs (eg., the subset $0 - 1$ loss), one might need to model the joint output distribution $P(y_1, \dots, y_q | x)$ to obtain a Bayes optimal prediction. If the outputs are conditionally independent however or if the loss function can be decomposed over the outputs, then a Bayes optimal prediction can be obtained by modeling separately the marginal conditional output distributions $P(y_j | x)$ for all j . This suggests that in this case, binary relevance / single target is not really penalized asymptotically with respect to multiple output methods for not considering the outputs jointly. In the case of an infinite sample size, it is thus expected to

Table 6.1: All methods compared on the 3 artificial datasets. Exploiting the output correlation structure (if it exists) allows beating single target in a finite sample size, decomposable metric and conditionally independent output.

Dataset	friedman ₁ -chain	friedman ₁ -group	friedman ₁ -ind
artificial-gbrt	0.654±0.015(1)	0.889±0.009(1)	0.831±0.004(1)
st-gbrt	0.626±0.016(5)	0.873±0.008(5)	0.830±0.003(2)
gbmort	0.640±0.008(4)	0.874±0.012(4)	0.644±0.010(5)
gbrt-relabel-rpo-subsampled	0.648±0.015(2)	0.880±0.009(2)	0.706±0.009(4)
gbrt-rpo-subsampled	0.645±0.013(3)	0.876±0.007(3)	0.789±0.003(3)

provide as good models as all the multiple output methods. *Since in practice we have to deal with finite sample sizes, multiple output methods may provide better results by better controlling the bias/variance trade-off.*

Let us study this question on the three synthetic datasets: friedman₁-chain, friedman₁-group or friedman₁-ind. We optimize the hyper-parameters with a computational budget of 10000 weak models per hyper-parameter set. Five strategies are compared (i) the artificial-gbrt method, which assumes that the output structure is known and implements the optimal strategy on each problem as explained in Section 6.3.2.1, (ii) single target of gradient boosting regression trees (st-gbrt), (iii) gradient boosting with multi-output regression tree (gbmort, Algorithm 6.1) and gradient boosting with randomly sub-sampled outputs (iv) without relabelling (gbrt-rpo-subsampled, Algorithm 6.2) and (v) with relabelling (gbrt-relabel-rpo-subsampled, Algorithm 6.3). All boosting algorithms minimize the square loss, the absolute loss or their multi-outputs extension the ℓ_2 -norm loss.

We give the performance on the three tasks for each estimator in Table 6.1 and the p-value of Student’s paired t-test comparing the performance of two estimators on the same dataset in Table 6.2.

As expected, we obtain the best performance if the output correlation structure is known with the custom strategies implemented with artificial-gbrt. Excluding this artificial method, the best boosting methods on the two problems with output correlations, friedman₁-chain and friedman₁-group, are the two gradient boosting approaches with output subsampling (gbrt-relabel-rpo-subsampled and gbrt-rpo-subsampled).

In friedman₁-chain, the output correlation structure forms a chain as each new output is the previous one in the chain with a noisy output. Predicting outputs at the end of the chain, without using the previous ones, is a difficult task. The single target approach is thus expected to be sub-optimal. And indeed, on this problem, artificial-gbrt, gbrt-relabel-rpo-subsampled and gbrt-rpo-subsampled are sig-

Table 6.2: P-values given by Student's paired t-test on the synthetic datasets. We highlight p-values inferior to $\alpha = 0.05$ in bold. Note that the sign $<$ (resp. $>$) indicates that the estimator in the row has better (resp. lower) score than the column estimator.

	artificial-gbrt	st-gbrt	gbmort	gbrt-relabel-rpo-subsampled	gbrt-rpo-subsampled
Dataset friedman1-chain					
artificial-gbrt		0.003 ($>$)	0.16	0.34	0.24
st-gbrt	0.003 ($<$)		0.11	0.04 ($<$)	0.03 ($<$)
gbmort	0.16	0.11		0.38	0.46
gbrt-relabel-rpo-subsampled	0.34	0.04 ($>$)	0.38		0.57
gbrt-rpo-subsampled	0.24	0.03 ($>$)	0.46	0.57	
Dataset friedman1-group					
artificial-gbrt		0.005 ($>$)	0.009 ($>$)	0.047 ($>$)	0.006 ($>$)
st-gbrt	0.005 ($<$)		0.56	0.046 ($<$)	0.17
gbmort	0.009 ($<$)	0.56		0.15	0.63
gbrt-relabel-rpo-subsampled	0.047 ($<$)	0.046 ($>$)	0.15		0.04 ($>$)
gbrt-rpo-subsampled	0.006 ($<$)	0.17	0.63	0.04 ($<$)	
Dataset friedman1-ind					
artificial-gbrt		0.17	2e-06 ($>$)	2e-06 ($>$)	1e-05 ($>$)
st-gbrt	0.17		2e-06 ($>$)	4e-06 ($>$)	4e-06 ($>$)
gbmort	2e-06 ($<$)	2e-06 ($<$)		9e-05 ($<$)	6e-06 ($<$)
gbrt-relabel-rpo-subsampled	2e-06 ($<$)	4e-06 ($<$)	9e-05 ($>$)		3e-05 ($<$)
gbrt-rpo-subsampled	1e-05 ($<$)	4e-06 ($<$)	6e-06 ($>$)	3e-05 ($>$)	

nificantly better than st-gbrt (with $\alpha = 0.05$). All the multi-output methods, including gbmort, are indistinguishable from a statistical point of view, but we note that gbmort is however not significantly better than st-gbrt.

In `friedman1-group`, among the ten pairs of algorithms, four are not significantly different, showing a p-value greater than 0.05 (see Table 6.2). We first note that gbmort is not better than st-gbrt while exploiting the correlation. Secondly, the boosting methods with random output sub-sampling are the best methods. They are however not significantly better than gbmort and significantly worse than artificial-gbrt, which assumes the output structure is known. Note that gbrt-relabel-rpo-subsampled is significantly better than gbrt-rpo-subsampled.

In `friedman1-ind`, where there is no correlation between the outputs, the best strategy is single target which makes independent models for each output. From a conceptual and statistical point of view, there is no difference between artificial-gbrt and st-gbrt. The gbmort algorithm, which is optimal when all outputs are correlated, is here significantly worse than all other methods. The two boosting methods with output subsampling (gbrt-rpo-subsampled and gbrt-relabel-rpo-subsampled method), which can adapt themselves to the absence of correlation between the outputs, perform better than gbmort, but they are significantly worse than st-gbrt. For these two algorithms, we note that not relabelling the leaves (gbrt-rpo-subsampled) leads to superior performance than relabelling them (gbrt-relabel-rpo-subsampled). Since in `friedman1-ind` the outputs have disjoint feature support, the test nodes of a decision tree fitted on one output will partition the samples using these features. Thus, it is not surprising that relabeling the trees leaves actually deteriorates performance.

In the previous experiment, all the outputs were dependent of the inputs. However in multi-output tasks with very high number of outputs, it is likely that some of them have few or no links with the inputs, i.e., are pure noise. Let us repeat the previous experiments with the main difference that we add to the original 16 outputs 16 purely noisy outputs obtained through random permutations of the original outputs. We show the results of optimizing each algorithm in Table 6.3 and the associated p-values in Table 6.4. We report the macro- r^2 score computed either on all outputs (macro- r^2) including the noisy outputs or only on the 16 original outputs (half-macro- r^2). P-value were computed between each pair of algorithms using Student's t-test on the macro r^2 score computed on all outputs.

We observe that the gbrt-rpo-subsampled algorithm has the best performance on `friedman1-chain` and `friedman1-group` and is the second best on the `friedman1-ind`, below st-gbrt. Interestingly on `friedman1-chain` and `friedman1-group`, this algorithm is significantly better than all the others, including gbmort. Since this latter method

Table 6.3: Friedman datasets with noisy outputs.

friedman ₁ -chain	half-macro-r ²	macro-r ²
st-gbrt	0.611 (4)	0.265 ± 0.006 (4)
gbmort	0.617 (3)	0.291 ± 0.012 (3)
gbrt-relabel-rpo-subsampled	0.628 (2)	0.292 ± 0.006 (2)
gbrt-rpo-subsampled	0.629 (1)	0.303 ± 0.007 (1)
Friedman ₁ -group	half-macro-r ²	macro-r ²
st-gbrt	0.840 (3)	0.364 ± 0.007 (4)
gbmort	0.833 (4)	0.394 ± 0.004 (3)
gbrt-relabel-rpo-subsampled	0.855 (2)	0.395 ± 0.005 (2)
gbrt-rpo-subsampled	0.862 (1)	0.414 ± 0.006 (1)
Friedman ₁ -ind	half-macro-r ²	macro-r ²
st-gbrt	0.806 (1)	0.3536 ± 0.0015 (1)
gbmort	0.486 (4)	0.1850 ± 0.0081 (4)
gbrt-relabel-rpo-subsampled	0.570 (3)	0.2049 ± 0.0033 (3)
gbrt-rpo-subsampled	0.739 (2)	0.3033 ± 0.0021 (2)

tries to fit all outputs simultaneously, it is the most disadvantaged by the introduction of the noisy outputs.

6.3.3 Effect of random projection

With the gradient boosting and random projection of the output space approaches (Algorithms 6.2 and 6.3), we have considered until now only sub-sampling a single output at each iteration as random projection scheme. In Section 6.3.3.1, we show empirically the effect of other random projection schemes such as Gaussian random projection. In Section 6.3.3.2, we study the effect of increasing the number of projections in the gradient boosting algorithm with random projection of the output space and relabelling (parameter q of Algorithm 6.3). We also show empirically the link between Algorithm 6.3 and gradient boosting with multi-output regression tree (Algorithm 6.1).

6.3.3.1 Choice of the random projection scheme

Beside random output sub-sampling, we can combine the multi-output gradient boosting strategies (Algorithms 6.2 and 6.3) with other random projection schemes. A key difference between random output sub-sampling and random projections such as Gaussian and (sparse) Rademacher projections is that the latter combines together several outputs.

Table 6.4: P-values given by Student's paired t-test on the synthetic datasets. We highlight p-values inferior to $\alpha = 0.05$ in bold. Note that the sign $<$ (resp. $>$) indicates that the estimator in the row has better (resp. lower) score than the column estimator.

	st-gbrt	gbmort	gbrt-relabel-rpo-subsampled	gbrt-rpo-subsampled
Dataset friedman1-chain				
st-gbrt		0.0009 ($<$)	0.0002 ($<$)	0.0002 ($<$)
gbmort	0.0009 ($>$)		0.86	0.04 ($<$)
gbrt-relabel-rpo-subsampled	0.0002 ($>$)	0.86		0.03 ($<$)
gbrt-rpo-subsampled	0.0002 ($>$)	0.04 ($>$)	0.03 ($>$)	
Dataset friedman1-group				
st-gbrt		0.0002 ($<$)	0.0006 ($<$)	0.0003 ($<$)
gbmort	0.0002 ($>$)		0.74	0.008 ($<$)
gbrt-relabel-rpo-subsampled	0.0006 ($>$)	0.74		0.002 ($<$)
gbrt-rpo-subsampled	0.0003 ($>$)	0.008 ($>$)	0.002 ($>$)	
Dataset friedman1-ind				
st-gbrt		1e-06 ($>$)	1e-07 ($>$)	1e-06 ($>$)
gbmort	1e-06 ($<$)		0.02 ($<$)	6e-06 ($<$)
gbrt-relabel-rpo-subsampled	1e-07 ($<$)	0.02 ($>$)		2e-06 ($<$)
gbrt-rpo-subsampled	1e-06 ($<$)	6e-06 ($>$)	2e-06 ($>$)	

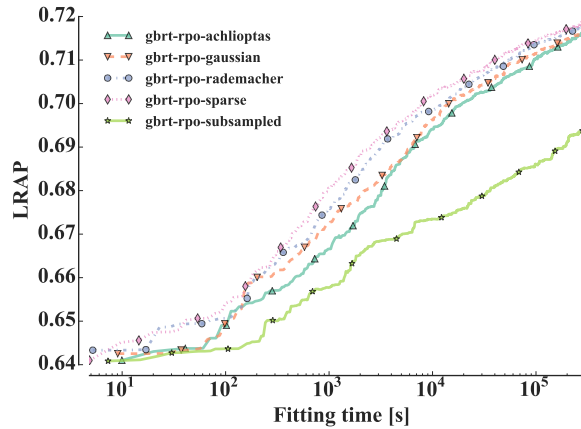


Figure 6.2: On the mediamill dataset, Gaussian, Achlioptas and sparse random projections with gbrt-rpo (Algorithm 6.2) show 10 times faster convergence in terms of LRAP score, than sub-sampling one output variable at each iteration. ($k = p$, stumps, $\mu = 0.1$, logistic loss)

We show in Figures 6.2, 6.3 and 6.4 the LRAP or macro- r^2 score convergence of gradient boosting with randomly projected outputs (gbrt-rpo, Algorithm 6.2) respectively on the mediamill, delicious, and Friedman₁-ind datasets with different random projection schemes.

The impact of the random projection scheme on convergence speed of gbrt-rpo (Algorithm 6.2) is very problem dependent. On the mediamill dataset, Gaussian, Achlioptas, or sparse random projections all improve convergence speed by a factor of 10 (see Figure 6.2) compared to subsampling randomly only one output. On the delicious (Figure 6.3) and friedman₁-ind (Figure 6.4), this is the opposite: subsampling leads to faster convergence than all other projections schemes. Note that we have the same behavior if one relabels the tree structure grown at each iteration as in Algorithm 6.3 (results not shown).

Dense random projections, such as Gaussian random projections, force the weak model to consider several outputs jointly and it should thus only improve when outputs are somewhat correlated (which seems to be the case on mediamill). When all of the outputs are independent or the correlation is less strong, as in friedman₁-ind or delicious, this has a detrimental effect. In this situation, sub-sampling only one output at each iteration leads to the best performance.

6.3.3.2 Effect of the size of the projected space

The multi-output gradient boosting strategy combining random projections and tree relabelling (Algorithm 6.3) can use more than one random projection ($q \geq 1$) by using multi-output trees as base learners. In this section, we study the effect of the size of the projected

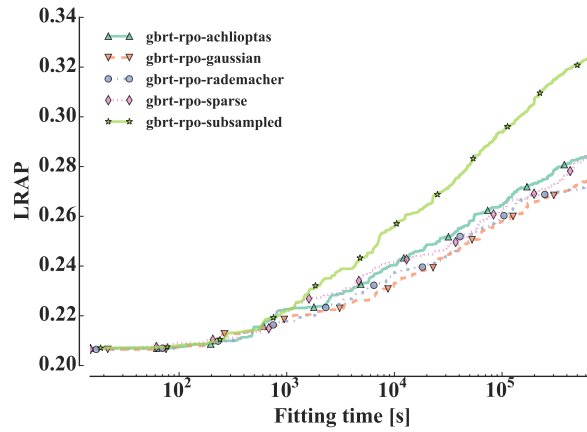


Figure 6.3: On the delicious dataset, Gaussian, Achlioptas and sparse random projections with gbrt-rpo (Algorithm 6.2) show 10 times faster convergence in terms of LRAP score, than sub-sampling one output variable at each iteration. ($k = p$, stumps, $\mu = 0.1$, logistic loss)

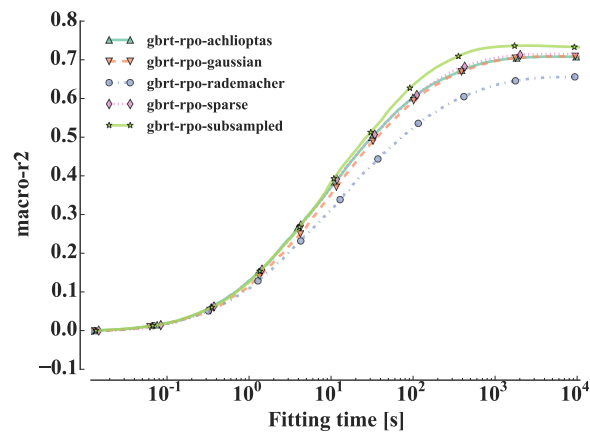
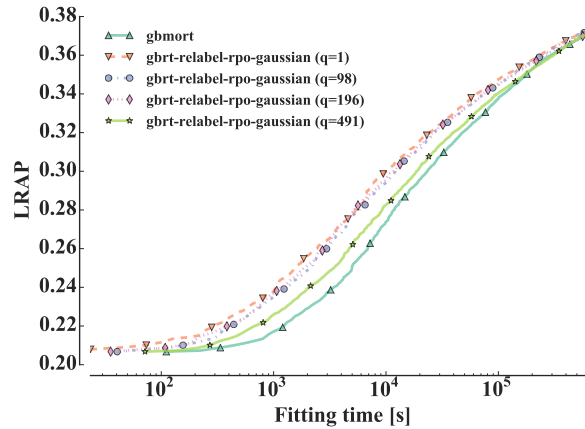
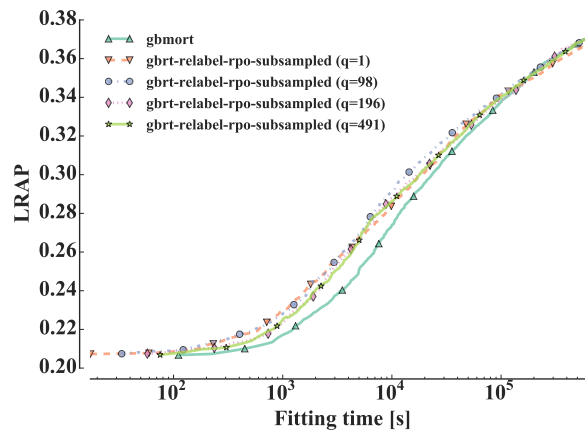


Figure 6.4: On the friedman1-ind dataset where there is no output correlation, gbrt-rpo (Algorithm 6.2) with one random subsampled output leads to a higher macro- r^2 score than using Gaussian, Achlioptas or sparse random projections. ($k = p$, stumps, $\mu = 0.1$, square loss)



(a)

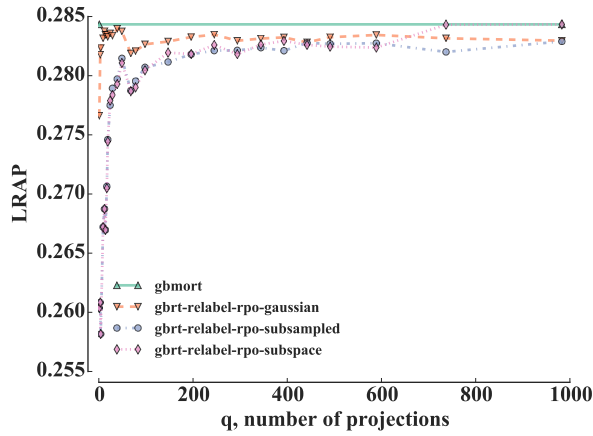


(b)

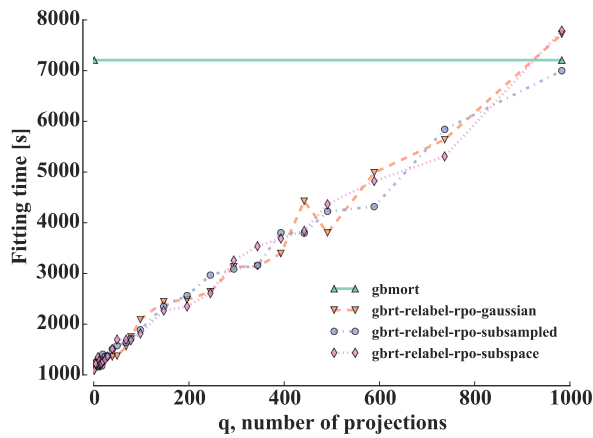
Figure 6.5: On the delicious dataset, LRAP score as a function of the boosting ensemble fitting time for gbrt-rpo-gaussian-relabel and gbrt-rpo-subsampled-relabel with different number of projections q . ($k = p$, stumps, $\mu = 0.1$, logistic loss)

space q in Algorithm 6.3. This approach corresponds to the one developed in Chapter 5 for random forest.

Figure 6.5 shows the LRAP score as a function of the fitting time for gbmort (Algorithm 6.1) and gbrt-relabel-rpo (Algorithm 6.3) with either Gaussian random projection (see Figure 6.5a) or output subsampling (see Figure 6.5b) for a number of projections $q \in \{1, 98, 196, 491\}$ on the delicious dataset. In Figure 6.5a and Figure 6.5b, one Gaussian random projection or one sub-sampled output has faster convergence than their counterparts with a higher number of projections q and gbmort at fixed computational budget. Note that when the number of projections q increases, gradient boosting with random projection of the output space and relabeling becomes similar to gbmort.



(a)



(b)

Figure 6.6: On delicious, increasing the number of random projections q allows to reach the same LRAP score as gbmort at a significantly reduced computational cost. ($k = p$, stumps, $\mu = 0.1$, $M = 100$, logistic loss)

Instead of fixing the computational budget as a function of the training time, we now set the computational budget to 100 boosting steps. On the delicious dataset, gbrt-relabel-rpo (Algorithm 6.3) with Gaussian random projection yields approximately the same performance as gbmort with $q \geq 20$ random projections as shown in Figure 6.6a and reduces computing times by a factor 7 at $q = 20$ projections (see Figure 6.6b).

These experiments show that gradient boosting with random projection and relabelling (gbrt-relabel-rpo, Algorithm 6.3) is indeed an approximation of gradient boosting with multi-output trees (gbmort, Algorithm 6.1). The number of random projections q influences simultaneously the bias-variance tradeoff and the convergence speed of Algorithm 6.3.

6.3.4 Systematic analysis over real world datasets

We perform a systematic analysis over real world multi-label classification and multi-output regression datasets. For this study, we evaluate the proposed algorithms: gradient boosting of multi-output regression trees (gbmort, Algorithm 6.1), gradient boosting with random projection of the output space (gbrt-rpo, Algorithm 6.2), and gradient boosting with random projection of the output space and relabelling (gbrt-relabel-rpo, Algorithm 6.3). For the two latter algorithms, we consider two random projection schemes: (i) Gaussian random projection, a dense random projection, and (ii) random output sub-sampling, a sparse random projection. They will be compared to three common and well established tree-based multi-output algorithms: (i) binary relevance / single target of gradient boosting regression tree (br-gbrt / st-gbrt), (ii) multi-output random forest (mo-rf) and (iii) binary relevance / single target of random forest models (br-rf / st-rf).

We will compare all methods on multi-label tasks in Section 6.3.4.1 and on multi-output regression tasks in Section 6.3.4.2. Following the recommendations of (Demšar, 2006), we use the Friedman test and its associated Nemenyi post-hoc test. Pairwise comparisons are also carried out using the Wilcoxon signed ranked test.

6.3.4.1 Multi-label datasets

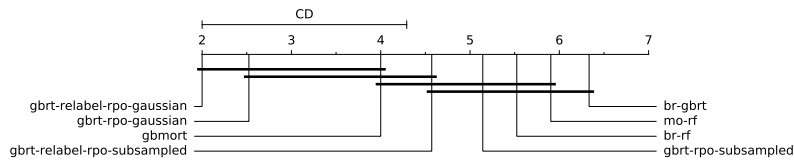
Table 6.5 and Table 6.6 show the performance of the random forest models and the boosting algorithms over the 21 multi-label datasets. The critical distance diagram of Figure 6.7 gives the ranks of the algorithms and has an associated Friedman test p-value of 1.36×10^{-10} with a critical distance of 2.29 given by the Nemenyi post-hoc test ($\alpha = 0.05$). Thus, we can reject the null hypothesis that all methods are equivalent. Table 6.7 gives the outcome of the pairwise Wilcoxon signed ranked tests.

The best average performer is gbrt-relabel-rpo-gaussian which is significantly better according to the Nemenyi post-hoc test than all methods except gbrt-rpo-gaussian and gbmort.

Gradient boosting with the Gaussian random projection has a significantly better average rank than the random output sub-sampling projection. Relabelling tree leaves allows to have better performance on the 21 multi-label dataset. Indeed, both gbrt-relabel-rpo-gaussian and gbrt-relabel-rpo-subsampled are better ranked and significantly better than their counterparts without relabelling (gbrt-rpo-gaussian and gbrt-rpo-subsampled). These results somewhat contrast with the results obtained on the artificial datasets, where relabelling was always counterproductive.

Among all compared methods, br-gbrt has the worst rank and it is significantly less good than all gbrt variants according to the

Figure 6.7: Critical difference diagram between algorithms on the multi-label datasets.



Wilcoxon signed rank test. This might be actually a consequence of the constant budget in time that was allocated to all methods (see Section 6.3.1). All methods were given the same budget in time but, given the very slow convergence rate of br-gbrt, this budget may not allow to grow enough trees per output with this method to reach competitive performance.

We notice also that both random forests based methods (mo-rf and br-rf) are less good than all gbrt variants, most of the time significantly, except for br-gbrt. It has to be noted however that no hyperparameter was tuned for the random forests. Such tuning could slightly change our conclusions, although random forests often work well with default setting.

6.3.4.2 Multi-output regression datasets

Table 6.8 shows the performance of the random forest models and the boosting algorithms over the 8 multi-output regression datasets. The critical distance diagram of Figure 6.8 gives the rank of each estimator. The associated Friedman test has a p-value of 0.3. Given the outcome of the test, we can therefore not reject the null hypothesis that the estimator performances can not be distinguished. Table 6.9 gives the outcomes of the pairwise Wilcoxon signed ranked tests. They confirm the fact that all methods are very close to each other as only two comparisons show a p-value lower than 0.05 (st-rf is better than st-gbrt and gbrt-rpo-subsampled). This lack of statistical power is probably partly due here to the smaller number of datasets included in the comparison (8 problems versus 21 problems in classification).

If we ignore statistical tests, as with multi-label tasks, gbrt-relabel-rpo-gaussian has the best average rank and st-gbrt the worst average rank. This time however, gbrt-relabel-rpo-gaussian is followed by the random forest based algorithms (st-rf and mo-rf) and gbmort. Given the lack of statistical significance, this ranking should however be interpreted cautiously.

Table 6.5: LRAP scores over 21 multi-label datasets (part 1).

	CAL500	bibtex	birds
br-gbrt	0.505 ± 0.002 (3.5)	0.587 ± 0.007 (6)	0.787 ± 0.009 (6)
br-rf	0.484 ± 0.002 (8)	0.542 ± 0.005 (8)	0.802 ± 0.013 (1)
gbmort	0.501 ± 0.005 (6)	0.595 ± 0.005 (4.5)	0.772 ± 0.007 (8)
gbrt-relabel-rpo-gaussian	0.507 ± 0.009 (1)	0.607 ± 0.005 (1)	0.800 ± 0.017 (2)
gbrt-relabel-rpo-subsampled	0.499 ± 0.008 (7)	0.596 ± 0.005 (3)	0.790 ± 0.016 (4)
gbrt-rpo-gaussian	0.505 ± 0.006 (3.5)	0.600 ± 0.003 (2)	0.793 ± 0.017 (3)
gbrt-rpo-subsampled	0.506 ± 0.006 (2)	0.595 ± 0.007 (4.5)	0.779 ± 0.018 (7)
mo-rf	0.502 ± 0.003 (5)	0.553 ± 0.005 (7)	0.789 ± 0.012 (5)
	bookmarks	corel5k	delicious
br-gbrt	0.4463 ± 0.0038 (7)	0.291 ± 0.006 (7)	0.347 ± 0.002 (8)
br-rf	0.4472 ± 0.0019 (6)	0.273 ± 0.012 (8)	0.373 ± 0.004 (6.5)
gbmort	0.4855 ± 0.0016 (2)	0.312 ± 0.009 (3.5)	0.384 ± 0.003 (3.5)
gbrt-relabel-rpo-gaussian	0.4893 ± 0.0003 (1)	0.315 ± 0.007 (1.5)	0.389 ± 0.003 (1)
gbrt-relabel-rpo-subsampled	0.4718 ± 0.0034 (4)	0.310 ± 0.007 (5)	0.384 ± 0.003 (3.5)
gbrt-rpo-gaussian	0.4753 ± 0.0022 (3)	0.315 ± 0.010 (1.5)	0.386 ± 0.004 (2)
gbrt-rpo-subsampled	0.4621 ± 0.0026 (5)	0.312 ± 0.006 (3.5)	0.377 ± 0.003 (5)
mo-rf	0.4312 ± 0.0023 (8)	0.294 ± 0.010 (6)	0.373 ± 0.004 (6.5)
	diatoms	drug-interaction	emotions
br-gbrt	0.623 ± 0.007 (7.5)	0.271 ± 0.018 (8)	0.800 ± 0.022 (7)
br-rf	0.623 ± 0.011 (7.5)	0.310 ± 0.009 (5)	0.816 ± 0.009 (1)
gbmort	0.656 ± 0.012 (4)	0.304 ± 0.005 (7)	0.794 ± 0.014 (8)
gbrt-relabel-rpo-gaussian	0.725 ± 0.010 (1)	0.326 ± 0.008 (1)	0.802 ± 0.017 (5.5)
gbrt-relabel-rpo-subsampled	0.685 ± 0.012 (3)	0.322 ± 0.009 (3)	0.808 ± 0.021 (3)
gbrt-rpo-gaussian	0.702 ± 0.014 (2)	0.323 ± 0.011 (2)	0.804 ± 0.009 (4)
gbrt-rpo-subsampled	0.653 ± 0.013 (5.5)	0.312 ± 0.013 (4)	0.802 ± 0.007 (5.5)
mo-rf	0.653 ± 0.010 (5.5)	0.308 ± 0.007 (6)	0.810 ± 0.010 (2)
	enron	genbase	mediamill
br-gbrt	0.685 ± 0.006 (6)	0.989 ± 0.009 (8)	0.7449 ± 0.0020 (8)
br-rf	0.683 ± 0.005 (7)	0.994 ± 0.005 (2)	0.7819 ± 0.0009 (1)
gbmort	0.705 ± 0.004 (2.5)	0.990 ± 0.004 (6)	0.7504 ± 0.0013 (7)
gbrt-relabel-rpo-gaussian	0.705 ± 0.003 (2.5)	0.993 ± 0.006 (3)	0.7660 ± 0.0021 (3)
gbrt-relabel-rpo-subsampled	0.697 ± 0.004 (5)	0.990 ± 0.010 (6)	0.7588 ± 0.0013 (5)
gbrt-rpo-gaussian	0.706 ± 0.004 (1)	0.992 ± 0.007 (4)	0.7608 ± 0.0008 (4)
gbrt-rpo-subsampled	0.699 ± 0.005 (4)	0.990 ± 0.005 (6)	0.7519 ± 0.0006 (6)
mo-rf	0.676 ± 0.004 (8)	0.995 ± 0.004 (1)	0.7793 ± 0.0015 (2)

Table 6.6: LRAP scores over 21 multi-label datasets (part 2).

	medical	protein-interaction reuters	
br-gbrt	0.864 ± 0.006 (3)	0.294 ± 0.007 (6)	0.939 ± 0.0033 (7)
br-rf	0.821 ± 0.007 (8)	0.293 ± 0.006 (7)	0.9406 ± 0.0016 (6)
gbmort	0.867 ± 0.011 (1.5)	0.310 ± 0.007 (2.5)	0.9483 ± 0.0014 (3)
gbrt-relabel-rpo-gaussian	0.867 ± 0.019 (1.5)	0.310 ± 0.009 (2.5)	0.9508 ± 0.0009 (1)
gbrt-relabel-rpo-subsampled	0.856 ± 0.012 (5)	0.303 ± 0.003 (4.5)	0.9441 ± 0.0016 (4)
gbrt-rpo-gaussian	0.859 ± 0.017 (4)	0.311 ± 0.007 (1)	0.9486 ± 0.0021 (2)
gbrt-rpo-subsampled	0.851 ± 0.009 (6)	0.303 ± 0.003 (4.5)	0.9430 ± 0.0031 (5)
mo-rf	0.827 ± 0.006 (7)	0.288 ± 0.009 (8)	0.9337 ± 0.0021 (8)
	scene	scop-go	sequence-funcat
br-gbrt	0.880 ± 0.003 (4)	0.716 ± 0.047 (8)	0.678 ± 0.008 (6)
br-rf	0.876 ± 0.003 (6)	0.798 ± 0.004 (2)	0.658 ± 0.008 (7)
gbmort	0.886 ± 0.004 (1)	0.796 ± 0.007 (3)	0.699 ± 0.005 (3)
gbrt-relabel-rpo-gaussian	0.884 ± 0.006 (2.5)	0.788 ± 0.006 (4)	0.703 ± 0.007 (2)
gbrt-relabel-rpo-subsampled	0.879 ± 0.008 (5)	0.770 ± 0.010 (6)	0.685 ± 0.008 (5)
gbrt-rpo-gaussian	0.884 ± 0.005 (2.5)	0.775 ± 0.018 (5)	0.706 ± 0.007 (1)
gbrt-rpo-subsampled	0.875 ± 0.006 (7)	0.723 ± 0.016 (7)	0.691 ± 0.006 (4)
mo-rf	0.865 ± 0.003 (8)	0.800 ± 0.006 (1)	0.643 ± 0.003 (8)
	wipo	yeast	yeast-go
br-gbrt	0.706 ± 0.009 (6)	0.756 ± 0.009 (8)	0.499 ± 0.009 (4.5)
br-rf	0.633 ± 0.013 (7)	0.760 ± 0.008 (3.5)	0.463 ± 0.010 (7)
gbmort	0.762 ± 0.011 (3)	0.760 ± 0.007 (3.5)	0.504 ± 0.015 (3)
gbrt-relabel-rpo-gaussian	0.776 ± 0.012 (1)	0.762 ± 0.007 (2)	0.524 ± 0.012 (1)
gbrt-relabel-rpo-subsampled	0.751 ± 0.017 (4)	0.758 ± 0.005 (5.5)	0.496 ± 0.013 (6)
gbrt-rpo-gaussian	0.763 ± 0.010 (2)	0.763 ± 0.005 (1)	0.522 ± 0.012 (2)
gbrt-rpo-subsampled	0.724 ± 0.011 (5)	0.758 ± 0.008 (5.5)	0.499 ± 0.011 (4.5)
mo-rf	0.624 ± 0.018 (8)	0.757 ± 0.008 (7)	0.415 ± 0.014 (8)

Figure 6.8: Critical difference diagram between algorithm on the multi-output regression datasets.

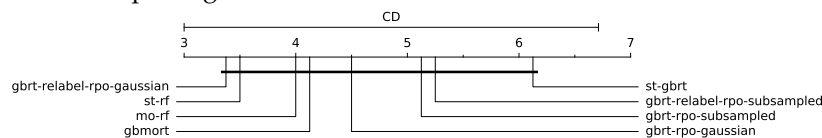


Table 6.7: P-values given by the Wilcoxon signed rank test on the multi-label datasets. We bold p-values below $\alpha = 0.05$. Note that the sign $>$ (resp. $<$) indicates that the row estimator has superior (resp. inferior) LRAP score than the column estimator.

	br-gbrt	br-rf	gbmort	gbrt-relabel-rpo-gaussian	gbrt-relabel-rpo-subsampled	gbrt-rpo-gaussian	gbrt-rpo-subsampled	mo-rf
br-gbrt		0.79	0.001 ($<$)	6e-05 ($<$)	0.001 ($<$)	0.0001 ($<$)	0.003 ($<$)	0.54
br-rf	0.79		0.01 ($<$)	0.002 ($<$)	0.02 ($<$)	0.005 ($<$)	0.07	0.29
gbmort	0.001 ($>$)	0.01 ($>$)		0.001 ($<$)	0.36	0.04 ($<$)	0.03 ($>$)	0.02 ($>$)
gbrt-relabel-rpo-gaussian	6e-05 ($>$)	0.002 ($>$)	0.001 ($>$)		0.0002 ($>$)	0.005 ($>$)	7e-05 ($>$)	0.0007 ($>$)
gbrt-relabel-rpo-subsampled	0.001 ($>$)	0.02 ($>$)	0.36	0.0002 ($<$)		0.0002 ($<$)	0.02 ($>$)	0.008 ($>$)
gbrt-rpo-gaussian	0.0001 ($>$)	0.005 ($>$)	0.04 ($>$)	0.005 ($<$)	0.0002 ($>$)		7e-05 ($>$)	0.002 ($>$)
gbrt-rpo-subsampled	0.003 ($>$)	0.07	0.03 ($<$)	7e-05 ($<$)	0.02 ($<$)	7e-05 ($<$)		0.04 ($>$)
mo-rf	0.54	0.29	0.02 ($<$)	0.0007 ($<$)	0.008 ($<$)	0.002 ($<$)	0.04 ($<$)	

Table 6.8: Performance over 8 multi-output regression dataset

	atp1d	atp7d	edm
gbmort	$0.80 \pm 0.03(5.5)$	$0.63 \pm 0.03(2)$	$0.39 \pm 0.16(3)$
gbrt-relabel-rpo-gaussian	$0.81 \pm 0.03(3.5)$	$0.66 \pm 0.04(1)$	$0.25 \pm 0.28(8)$
gbrt-relabel-rpo-subsampled	$0.79 \pm 0.04(7)$	$0.54 \pm 0.13(7)$	$0.35 \pm 0.10(5)$
gbrt-rpo-gaussian	$0.80 \pm 0.04(5.5)$	$0.54 \pm 0.20(7)$	$0.36 \pm 0.04(4)$
gbrt-rpo-subsampled	$0.81 \pm 0.04(3.5)$	$0.54 \pm 0.16(7)$	$0.31 \pm 0.27(7)$
mo-rf	$0.82 \pm 0.03(2)$	$0.6 \pm 0.06(4)$	$0.51 \pm 0.02(1)$
st-gbrt	$0.78 \pm 0.05(8)$	$0.59 \pm 0.08(5)$	$0.34 \pm 0.14(6)$
st-rf	$0.83 \pm 0.02(1)$	$0.61 \pm 0.07(3)$	$0.47 \pm 0.04(2)$
	oes10	oes97	scm1d
gbmort	$0.77 \pm 0.05(3.5)$	$0.67 \pm 0.07(8)$	$0.908 \pm 0.003(4.5)$
gbrt-relabel-rpo-gaussian	$0.75 \pm 0.04(7.5)$	$0.71 \pm 0.07(2.5)$	$0.910 \pm 0.004(2.5)$
gbrt-relabel-rpo-subsampled	$0.75 \pm 0.06(7.5)$	$0.68 \pm 0.07(6)$	$0.912 \pm 0.003(1)$
gbrt-rpo-gaussian	$0.77 \pm 0.03(3.5)$	$0.68 \pm 0.08(6)$	$0.910 \pm 0.004(2.5)$
gbrt-rpo-subsampled	$0.76 \pm 0.02(5.5)$	$0.71 \pm 0.08(2.5)$	$0.908 \pm 0.004(4.5)$
mo-rf	$0.76 \pm 0.04(5.5)$	$0.69 \pm 0.05(4)$	$0.898 \pm 0.004(8)$
st-gbrt	$0.79 \pm 0.03(1.5)$	$0.68 \pm 0.07(6)$	$0.905 \pm 0.003(7)$
st-rf	$0.79 \pm 0.03(1.5)$	$0.72 \pm 0.05(1)$	$0.907 \pm 0.004(6)$
	scm2od	water-quality	
gbmort	$0.856 \pm 0.006(2)$	$0.14 \pm 0.01(4.5)$	
gbrt-relabel-rpo-gaussian	$0.862 \pm 0.006(1)$	$0.15 \pm 0.01(2)$	
gbrt-relabel-rpo-subsampled	$0.854 \pm 0.007(3)$	$0.14 \pm 0.02(4.5)$	
gbrt-rpo-gaussian	$0.852 \pm 0.006(4)$	$0.14 \pm 0.01(4.5)$	
gbrt-rpo-subsampled	$0.850 \pm 0.007(5)$	$0.13 \pm 0.02(7.5)$	
mo-rf	$0.849 \pm 0.007(6.5)$	$0.16 \pm 0.01(1)$	
st-gbrt	$0.836 \pm 0.006(8)$	$0.13 \pm 0.02(7.5)$	
st-rf	$0.849 \pm 0.006(6.5)$	$0.14 \pm 0.01(4.5)$	

Table 6.9: P-value given by the Wilcoxon signed rank test on multi-output regression datasets. We bold p-values below $\alpha = 0.05$. Note that the sign $>$ (resp. $<$) indicates that the row estimator has superior (resp. inferior) macro- r^2 score than the column estimator.

	st-gbrt	st-rf	gbmort	gbrt-relabel-rpo-gaussian	gbrt-relabel-rpo-subsampled	gbrt-rpo-gaussian	gbrt-rpo-subsampled	mo-rf
st-gbrt		0.02(<)	0.26	0.58	0.78	0.67	0.89	0.16
st-rf	0.02(>)		0.33	0.67	0.09	0.09	0.04(>)	0.58
gbmort	0.26	0.33		0.4	0.16	0.67	0.4	0.48
gbrt-relabel-rpo-gaussian	0.58	0.67	0.4		0.16	0.4	0.48	0.58
gbrt-relabel-rpo-subsampled	0.78	0.09	0.16	0.16		0.16	1	0.07
gbrt-rpo-gaussian	0.67	0.09	0.67	0.4	0.16		0.48	0.26
gbrt-rpo-subsampled	0.89	0.04(<)	0.4	0.48	1	0.48		0.4
mo-rf	0.16	0.58	0.48	0.58	0.07	0.26	0.4	

6.4 CONCLUSIONS

In this chapter, we have first formally extended the gradient boosting algorithm to multi-output tasks leading to the “multi-output gradient boosting algorithm” (gbmort). It sequentially minimizes a multi-output loss using multi-output weak models considering that all outputs are correlated. By contrast, binary relevance / single target of gradient boosting models fit one gradient boosting model per output considering that all outputs are independent. However in practice, we do not expect to have either all outputs independent or all outputs dependent. So, we propose a more flexible approach which adapts automatically to the output correlation structure called “gradient boosting with random projection of the output space” (gbrt-rpo). At each boosting step, it fits a single weak model on a random projection of the output space and optimize a multiplicative weight separately for each output. We have also proposed a variant of this algorithm (gbrt-relabel-rpo) only valid with decision trees as weak models: it fits a decision tree on the randomly projected space and then it relabels tree leaves with predictions in the original (residual) output space. The combination of the gradient boosting algorithm and the random projection of the output space yields faster convergence by exploiting existing correlations between the outputs and by reducing the dimensionality of the output space. It also provides new bias-variance-convergence trade-off potentially allowing to improve performance.

We have evaluated in depth these new algorithms on several artificial and real datasets. Experiments on artificial problems highlighted that gb-rpo with output subsampling offers an interesting tradeoff between single target and multi-output gradient boosting. Because of its capacity to automatically adapt to the output space structure, it outperforms both methods in terms of convergence speed and accuracy when outputs are dependent and it is superior to gbmort (but not st-rt) when outputs are fully independent. On the 29 real datasets, gbrt-relabel-rpo with the denser Gaussian projections turns out to be the best overall approach on both multi-label classification and multi-output regression problems, although all methods are statistically undistinguishable on the regression tasks. Our experiments also show that gradient boosting based methods are competitive with random forests based methods. Given that multi-output random forests were shown to be competitive with several other multi-label approaches in (Madjarov et al., 2012), we are confident that our solutions will be globally competitive as well, although a broader empirical comparison should be conducted as future work. One drawback of gradient boosting with respect to random forests however is that its performance is more sensitive to its hyper-parameters that thus require careful tuning. Although not discussed in this chapter, besides pre-

dictive performance, `gbrt-rpo` (without relabeling) has also the advantage of reducing model size with respect to `mo-rf` (multi-output random forests) and `gbmort`, in particular in the presence of many outputs. Indeed, in `mo-rf` and `gbmort`, one needs to store a vector of the size of the number of outputs per leaf node. In `gbrt-rpo`, one needs to store only one real number (a prediction for the projection) per leaf node and a vector of the size of the number of outputs per tree (ρ_m). At fixed number of trees and fixed tree complexity, this could lead to a strong reduction of the model memory requirement when the number of labels is large. Note that the approach proposed in Chapter 5 does not solve this issue because of leaf node relabeling. This could be addressed by deactivating leaf relabeling and inverting the projection at prediction time to obtain a prediction in the original output space, as done for example in (Hsu et al., 2009; Kapoor et al., 2012; Tsoumakas et al., 2014). However, this would be at the expense of computing times at prediction time and of accuracy because of the potential introduction of errors at the decoding stage. Finally, while we restricted our experiments here to tree-based weak learners, Algorithms 6.1 and 6.2 are generic and could exploit respectively any multiple output and any single output regression method. As future work, we believe that it would be interesting to evaluate them with other weak learners.

Part III

EXPLOITING SPARSITY FOR GROWING AND COMPRESSING DECISION TREES

7

ℓ_1 -BASED COMPRESSION OF RANDOM FOREST MODELS

OUTLINE

Random forests are effective supervised learning methods applicable to large-scale datasets. However, the space complexity of tree ensembles, in terms of their total number of nodes, is often prohibitive, specially in the context of problems with large sample sizes and very high-dimensional input spaces. We propose to study their compressibility by applying a ℓ_1 -based regularization to the set of indicator functions defined by all their nodes. We show experimentally that preserving or even improving the model accuracy while significantly reducing its space complexity is indeed possible.

This chapter extends on previous work published in

Arnaud Joly, François Schnitzler, Pierre Geurts, and Louis Wehenkel. L1-based compression of random forest models. In European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning, 2012.

High-dimensional supervised learning problems, *e.g.* in image exploitation and bioinformatics, are more frequent than ever. Tree-based ensemble methods, such as random forests (Breiman, 2001) and extremely randomized trees (Geurts et al., 2006a), are effective variance reduction techniques offering in this context a good trade-off between accuracy, computational complexity, and interpretability. The number of nodes of a tree ensemble grows as nM (n being the size of the learning sample and M the number of trees in the ensemble). Empirical observations show that the variance of individual trees increases with the dimension p of the original feature space used to represent the inputs of the learning problem. Hence, the number $M(p)$ of ensemble terms yielding near-optimal accuracy, which is proportional to this variance, also increases with p . The net result is that the space complexity of these tree-based ensemble methods will grow as $nM(p)$, which may jeopardize their practicality in large scale problems, or when memory is limited.

While pruning of single tree models is a standard approach, less work has been devoted to pruning ensembles of trees. On the one hand, Geurts (2000) proposes to transpose the classical cost-complexity pruning of individual trees to ensembles. On the other

hand, [Friedman and Popescu \(2008\)](#); [Meinshausen \(2010\)](#); [Meinshausen et al. \(2009\)](#) propose to improve model interpretability by selecting optimal rule subsets from tree-ensembles. Another approach to reduce complexity and/or improve accuracy of ensembles of trees is to merely select an optimal subset of trees from a very large ensemble generated in a random fashion at the first hand (see, e.g. ([Bernard et al., 2009](#); [Martinez-Muoz et al., 2009](#))).

To further investigate the feasibility of reducing the space complexity of tree-based ensemble models, we consider in this chapter the following method ([Joly et al., 2012](#)): (i) build an ensemble of trees; (ii) apply to this ensemble a ‘compression step’ by reformulating the tree-ensemble based model as a linear model in terms of node indicator functions and by using an ℓ_1 -norm regularization approach - à la Lasso ([Tibshirani, 1996b](#)) - to select a minimal subset of these indicator functions while maintaining predictive accuracy. We propose an algorithmic framework and an empirical investigation of this idea, based on three complementary datasets, and we show that indeed it is possible to so compress significantly tree-based ensemble models, both in regression and in classification problems. We also observe that the compression rate and the accuracy of the compressed models further increase with the ensemble size M , even beyond the number $M(p)$ of terms required to ensure convergence of the variance reduction effect.

The rest of this chapter is organized as follows: Section 7.1 introduces the ℓ_1 -norm based compression algorithm of random forests; Section 7.2 provides our empirical study and Section 7.3 concludes and describes further perspectives.

7.1 COMPRESSING TREE ENSEMBLES BY ℓ_1 -NORM REGULARIZATION

From an ensemble of M decision trees, one can extract a set of node indicator functions as follows: each indicator function $1_{m,l}(x)$ is a binary variable equal to 1 if the input vector x reaches the l th node in the m th tree, 0 otherwise. Using these indicator functions, the output predicted by the model may be rewritten as ([Geurts et al., 2006a](#); [Vens and Costa, 2011](#)):

$$\hat{f}(x) = \frac{1}{M} \sum_{m=1}^M \sum_{l=1}^{N_m} w_{m,l} 1_{m,l}(x), \quad (7.1)$$

where N_m is the number of nodes in the m th tree and $w_{m,l}$ is equal to the leaf-label if node (m, l) is a leaf and to zero if it is an internal node.

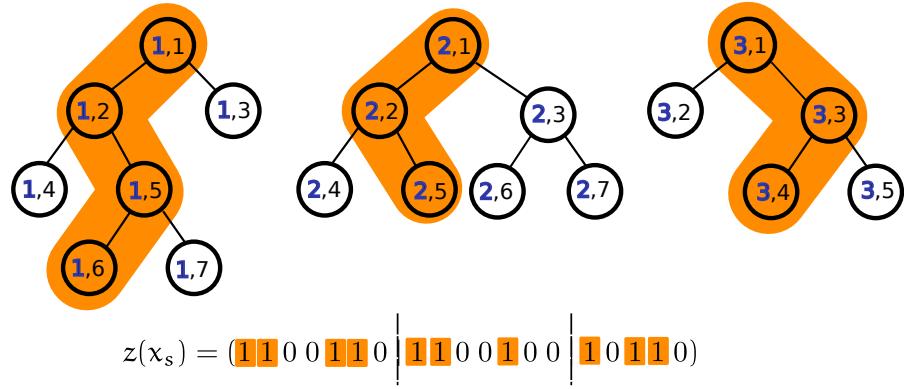


Figure 7.1: From a random forest model, one can lift the original input space representation of a sample x_s toward the node indicator space \mathcal{Z} .

We can therefore interpret a tree building algorithm as the (random) inference of a new representation which lifts the original input space \mathcal{X} towards the space \mathcal{Z} of dimension $q = \sum_{m=1}^M N_m$ by

$$z(x) = (1_{1,1}(x), \dots, 1_{1,N_1}(x), \dots, 1_{M,1}(x), \dots, 1_{M,N_M}(x)).$$

As an illustration, Figure 7.1 shows a set of three decision trees with respective sizes 7, 7 and 5 nodes. The propagation of a sample x_s through the forest makes it pass through nodes 1.1, 1.2, 1.5 in the left tree, nodes 2.1, 2.2, 2.5 in the middle tree and nodes 3.1, 3.3 and 3.4 in the left tree (highlighted in orange).

We propose to compress the tree ensemble by applying a variable selection method to its induced feature space \mathcal{Z} . Namely, by ℓ_1 -regularization we can search for a linear model by solving the following optimization problem:

$$\begin{aligned}
 (\beta_j^*(t))_{j=0}^q &= \arg \min_{\beta} \sum_{i=1}^n \left(y^i - \beta_0 - \sum_{j=1}^q \beta_j z_j(x^i) \right)^2 \\
 \text{s.t. } \sum_{j=1}^q |\beta_j| &\leq t.
 \end{aligned} \tag{7.2}$$

This optimization problem, also called Lasso (Tibshirani, 1996b) (see Section 2.2.1), has received much attention in the past decade and is particularly successful in high dimension. The ℓ_1 -norm constraint leads to a sparse solution: only a few weights β_j will be non zero, and their number tends to zero with $t \rightarrow 0$; the optimal value t^* of t is problem specific and is typically adjusted by cross-validation.

In order to solve Equation 7.2 for growing values of t , we use the ‘incremental forward stagewise regression’ algorithm (Hastie et al., 2007) solving the monotone Lasso which imposes that each $\beta_j^*(t)$ increases monotonically with t . This version deals indeed better with many correlated variables, which is relevant in our setting, since each

node indicator function is highly correlated with those of its neighbor nodes in the tree from which it originates. The final weights $\beta_j^*(t^*)$ may be exploited to prune the randomized tree ensemble: a test node can be deleted if all its descendants correspond to $\beta_j^*(t^*) = 0$.

Starting from a forest model \hat{f} , a value of parameter t , and a sample \mathcal{S} , the tree ensemble compression procedure is described in Algorithm 7.1.

Algorithm 7.1 ℓ_1 -based compression of tree ensemble model \hat{f} using a sample $\mathcal{S} = \{x^i, y^i \in \mathcal{X} \times \mathcal{Y}\}_{i=1}^n$

- 1: **function** FORESTCOMPRESSION(\mathcal{S}, \hat{f}, t)
- 2: Lift the sample \mathcal{S} to the random forest space \mathcal{Z}

$$\mathcal{S}_z = \{(z(x^i), y^i) \in \mathcal{Z} \times \mathcal{Y}\}_{(x,y) \in \mathcal{S}}$$

with the induced feature space by the forest model \hat{f}

$$z(x) = (1_{1,1}(x), \dots, 1_{1,N_1}(x), \dots, 1_{M,1}(x), \dots, 1_{M,N_M}(x)).$$

- 3: Select weight vector $\beta^*(t)$ over \mathcal{Z} through ℓ_1 minimization

$$\begin{aligned} (\beta_j^*(t))_{j=0}^q &= \arg \min_{\beta} \sum_{i=1}^n \left(y^i - \beta_0 - \sum_{j=1}^q \beta_j z_j(x^i) \right)^2 \\ \text{s.t. } \sum_{j=1}^q |\beta_j| &\leq t. \end{aligned}$$

- 4: Compress the random forest model \hat{f} using vector $\beta^*(t)$
 - 5: **return** The compressed model.
 - 6: **end function**
-

Note that in practice both the forest construction and the generation of its sequence of compressed versions for growing values of t may use the same sample (the learning set). A separate validation set is however required to select the optimal value of parameter t . This is similar to what is done with the pruning of a single decision tree (see Section 3.3).

7.2 EMPIRICAL ANALYSIS

In the following experiments, datasets are pre-whitened: input/output data are translated to zero mean and rescaled to unit variance. All results shown are averaged over 50 runs in order to avoid randomization artifacts.

Each one of these runs consisted of first generating a training set, and a testing set, and then working as follows. When using the monotone Lasso, we apply the incremental forward stagewise algorithm

with a step size $\epsilon = 0.01$. The optimal number of steps n_{step}^* or the optimal point $t^* = n_{\text{step}}^* \epsilon$ was chosen by ten-fold cross-validation t_{cv}^* over the training set (to this end, we used a quadratic loss in regression and a 0–1 loss in classification). More precisely, the training set is first divided ten times through cross-validation into a learning set, used both to fit a forest model and to run the incremental forward stagewise algorithm on it, and into a validation set, to estimate the losses of the resulting sequence of compressed forests. For each fold, we assess the model fitted over the training set using the validation set with increasing values of t by steps of ϵ . For each value of t , the ten model losses are averaged. The optimal value of t^* and the corresponding model compression level are those leading to the best average loss over the ten folds. The model is then refitted using the entire training set with $t = t^*$.

Below, we will apply our approach while using the extremely randomized trees method (Geurts et al., 2006a) to grow the forests (abbreviated by “ET”) and we denote their ℓ_1 -regularization-based compressed version “rET”.

We present an overall performance analysis in Section 7.2.1. Later on, we enhance our comprehension of the pruning algorithm by studying the effect of the regularization parameter t in Section 7.2.2 and of the complexity of the initial forest model by varying the pre-pruning rule values n_{min} , the minimum number of samples to split, and M , the number of trees, in Section 7.2.3. While in these last two sections, we focus our analysis on models obtained on the Friedman1 problem, we notice that similar conclusions can also be drawn for Two-norm and SEFTi datasets.

7.2.1 Overall performances

We have evaluated our approach on two regression datasets Friedman1 and SEFTi and one classification dataset Two-norm (see Appendix A for their description).

We have used a set of representative meta-parameter values (K , n_{min} and M) of the Extra-Trees algorithm (see Table 7.1). Accuracies are measured on the test sample and complexity is measured by the number of test nodes of the ET and rET models (the compression factor being the ratio of the former to the latter). We observe a compression factor between 9 and 34, a slightly lower error for the rET model than for the ET model on the two regression problems (Friedman1 and SEFTi) and the opposite on Two-norm. To compare, we show the results obtained with the linear Lasso based on the original features (its complexity is measured by the number of kept features): it is much less accurate than both ET and rET on the (non-linear) regression problems (Friedman1 and SEFTi), but superior on the (linear) classification problem (Two-norm).

Table 7.1: Overall assessment (parameters of the Extra-Tree method: $M = 100$; $K = p$; $n_{\min} = 1$ on Friedman₁ and Two-norm, $n_{\min} = 10$ on SEFTi).

Datasets	Error			Complexity			
	ET	rET	Lasso	ET	rET	ET/rET	Lasso
Friedman ₁	0.19587	0.18593	0.282441	29900	885	34	4
Two-norm	0.04177	0.06707	0.033500	4878	540	9	20
SEFTi	0.86159	0.84131	0.988031	39436	2055	19	14

Side experiments (results not provided) show that changing the value of parameter K does not influence significantly the final accuracy and complexity on the Two-norm and Friedman₁ datasets, while for SEFTi, accuracy increases strongly with K (presumably due to a large number of noisy and/or irrelevant features) with however little impact on the final complexity.

7.2.2 Effect of the regularization parameter t .

The complexity of the regularized ET model is shrunk with the ℓ_1 -norm constraint of Equation 7.2 in a way depending on the value of t . As shown in Figure 7.2(a), an increase of t decreases the error of rET until $t = 3$, leading to a complexity (Figure 7.2(b)) of about 900 test nodes. Notice that in general the rET model eventually overfits when t becomes large, although this is not visible on the range of values displayed in Figure 7.2(a) as the algorithm stops before.

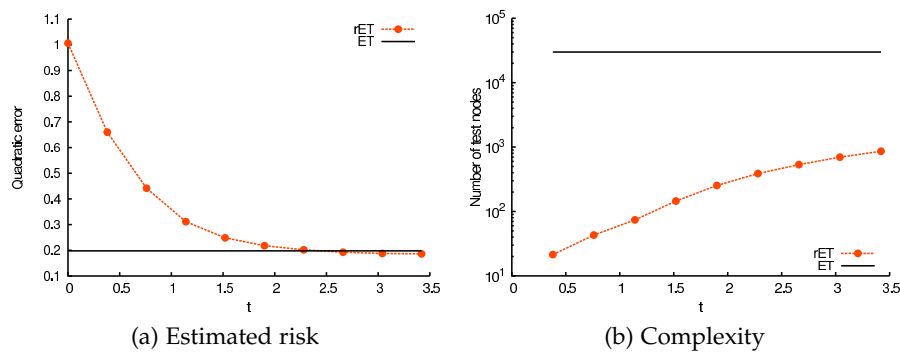


Figure 7.2: An increase of t decreases the error of rET until $t = 3$ with drastic pruning (Friedman₁, $M = 100$, $K = p = 10$ and $n_{\min} = 1$).

7.2.3 Influence of the Extra-Tree meta parameters n_{\min} and M .

The complexity of an ET model grows (linearly) with the size of the ensemble M and is inversely proportional to its pre-pruning parameter n_{\min} .

Figure 7.3 shows the effect of n_{\min} on both ET and rET. Interestingly, the accuracy and the complexity of the rET model are both more robust with respect to the choice of the precise value of n_{\min} than those of the ET model, specially for the smaller values of n_{\min} ($n_{\min} \leq 10$, in Figures 7.3).

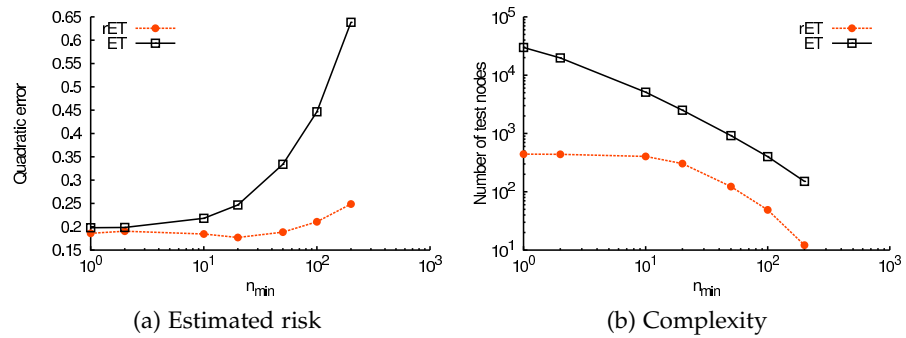


Figure 7.3: The accuracy and complexity of an rET model does not depend on n_{\min} , for n_{\min} small enough (Friedman1, $M = 100$, $K = p = 10$ and $t = t_{cv}^*$).

Figure 7.4 shows the effect of M on both ET and rET models. We observe that increasing the value of M beyond the value $M(p)$ where variance reduction has stabilized ($M(p) \simeq 100$ in Figure 7.4) allows to further improve the accuracy of the rET model without increasing its complexity.

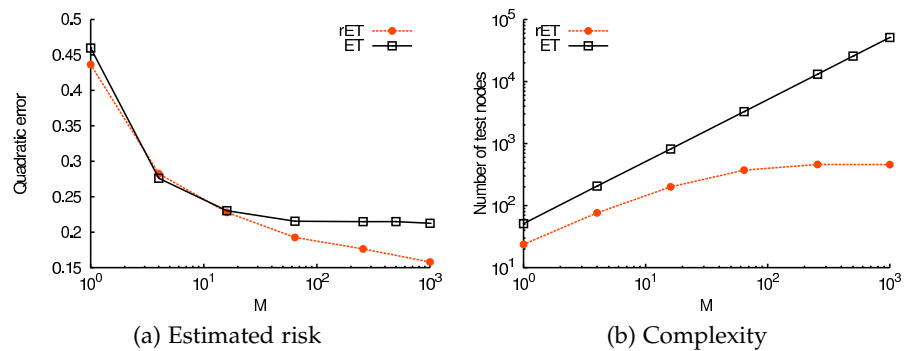


Figure 7.4: After variance reduction has stabilized ($M \simeq 100$), further increasing M keeps enhancing the accuracy of the rET model without increasing complexity (Friedman1, $n_{\min} = 10$, $K = p = 10$ and $t = t_{cv}^*$).

7.3 CONCLUSION

Compression of randomized tree ensembles with ℓ_1 -norm regularization leads to a drastic reduction of space complexity while preserving accuracy. The complexity of the pruned model does not seem to be directly related to the complexity of the original forest, *i.e.* the number and complexity of each randomized tree, as long as this forest has explored a large enough space of variable interactions.

The strong compressibility of large randomized tree ensemble models suggests that it could be possible to design novel algorithms based on tree-based randomization which would scale in a better way to very high-dimensional input spaces than the existing methods. To achieve this, one open question is how to get the compressed tree ensemble directly, *i.e.* without generating a huge randomized tree ensemble and then pruning it.

Tree-based ensemble models may be interpreted as lifting the original input space towards a (randomly generated) high-dimensional discrete and sparse representation, where each induced feature corresponds to the indicator function of a particular tree node, and takes the value 1 for a given observation if this observation reaches this node, and 0 otherwise. The dimension of this representation is on the order of $nM(p)$, but the number s of non-zero components for a given observation is only on the order of $M(p) \log n$. Compressed sensing theory (Candès and Wakin, 2008) tells us that high-dimensional sparsely representable observations may be compressed by projecting them on a random subspace of dimension proportional to $s \log p$, where p is the original dimension of the observations and $s \ll p$ is the number of non-zero terms in their sparse representation basis. This suggests that one could reduce the space complexity of tree-based method by applying compressed sensing to their original input feature space if its dimension is high, and/or to their induced feature space if $nM(p)$ is too large.

Since the publication of our work on this subject, several authors have proposed similar ideas to post-prune a fully grown random forest model: Ren et al. (2015) propose to iteratively remove or re-weight the leaves of the random forest model, while Duroux and Scornet (2016) study the impact of pre-pruning on random forest models.

8

EXPLOITING INPUT SPARSITY WITH DECISION TREE

OUTLINE

Many supervised learning tasks, such as text annotation, are characterized by high dimensional and sparse input spaces where the input vectors of each sample has only a few non zero values. We show how to exploit algorithmically the input space sparsity within decision tree methods. It leads to significant speed up both on synthetics and real datasets, while leading to exactly the same model. We also reduce the required memory to grow such models by exploiting sparse memory storage instead of dense memory storage for the input matrix.

This contribution is a joint work with Fares Hedayati and Panagiotis Papadimitriou, working at www.upwork.com. The outcome of this research has been proposed and merged in the scikit-learn ([Buitinck et al., 2013](#); [Pedregosa et al., 2011](#)) open source package.

Many machine learning tasks such as text annotation usually require training over very big datasets with millions of web documents. Such tasks require defining a mapping between the raw input space and the output space. For example in text classification, a text document (raw input space) is usually mapped to a vector whose dimensions correspond to all of the possible words in a dictionary and the values of the vector elements are determined by the frequency of the words in the document. Although such vectors have many dimensions, they are often sparsely representable. For instance, the number of unique words associated to a text document is actually small compared to the number of words of a given language. We describe those samples with sparse input vectors as having a few non zero values.

Exploiting the low density, i.e. the fraction of non zero elements, and the high sparsity, i.e. the fraction of zero elements, is key to address such high dimensional supervised learning tasks. Many models directly formulate their entire algorithm to exploit the input sparsity. Linear models such as logistic regression or support vector machine harness the sparsity by expressing most of their operations as a set of linear algebra operations such as dot products who directly exploit the sparsity to speed up computations.

Unfortunately, decision tree methods are not expressible only through linear algebra operations. Decision tree methods are recursively partitioning the input space by searching for the best possible splitting rules. As a consequence, most machine learning packages ei-

ther do not support sparse input vectors for tree-based methods, only support stumps (decision tree with only one internal node) or have a sub-optimal implementation through the simulation of a random access memory as in the dense case. The only solution is often to densify the input space which leads first to severe memory constraints and then to slow training time.

In Section 8.1, we present efficient algorithms to grow vanilla decision trees, boosting and random forest methods on sparse input data. In Section 8.2, we describe how to adapt the prediction algorithm of these models to sparse input data. In Section 8.3, we show empirically the speed up obtained by fitting decision trees with this input sparsity aware implementation.

8.1 TREE GROWING

During the decision tree fitting, the tree growing algorithm (see Algorithm 3.2) interacts with the input space at two key points:

1. during the search of a splitting rule s_t using a sample set \mathcal{L}_t at the expansion of node t (see line 10 of Algorithm 3.2);
2. during the data partitioning of the sample \mathcal{L}_t into a left and a right partition following the splitting rule s_t at node t (see line 11 of Algorithm 3.2).

In this section, we show how to adapt decision tree at these three key points to handle sparsely expressed data. While at the same time, we will show how to harness the sparsity to speed up the original algorithm. We first explain how node splitting is implemented in standard decision trees in Section 8.1.1 and then explain our efficient implementation for sparse input data in Section 8.1.2. In Section 8.1.3, we further describe how to propagate samples with sparse input during the decision tree growing.

8.1.1 Standard node splitting algorithm

During the decision tree growth, the crux of the tree growing algorithm in high dimensional input space is the search of the best possible local splitting rule s_t (as described in Section 3.2.1). Given a learning set \mathcal{L}_t reaching a node t , we search for the splitting rule s_t among all the possible binary and axis-wise splitting rules $\Omega(\mathcal{L}_t)$. We strive to maximize the impurity reduction ΔI obtained by dividing the sample set \mathcal{L}_t into two partitions $(\mathcal{L}_{t,r}, \mathcal{L}_{t,l})$. The splitting rule selection problem (line 10 of the tree growing Algorithm 3.2) is written as

$$s_t = \arg \max_{s \in \Omega(\mathcal{L}_t)} \Delta I(\mathcal{L}_t, \mathcal{L}_{t,l}, \mathcal{L}_{t,r}) \quad (8.1)$$

with

$$\Omega(\mathcal{L}_t) = \left\{ s : s \in \bigcup_{j \in \{1, \dots, p\}} Q(x_j), \right. \\ \mathcal{L}_{t,l} = \{(x, y) \in \mathcal{L}_t : s(x) = 1\}, \\ \mathcal{L}_{t,r} = \{(x, y) \in \mathcal{L}_t : s(x) = 0\}, \\ \left. \mathcal{L}_{t,l} \neq \emptyset, \mathcal{L}_{t,r} \neq \emptyset \right\} \quad (8.2)$$

where $Q(x_j)$ is the set of all splitting rules associated to an input variable x_j .

Decision tree libraries carefully optimize this part of the algorithm to have the proper computational complexity with a low constant. A careful design of the algorithm for instance does not move around samples according to the partitions, but instead move an identification number linked to each sample. The learning set \mathcal{L} is implemented as an array of row indices L linked to the rows of the input matrix X and the output matrix Y . The sample set \mathcal{L}_t reaching a node t is implemented as a slice of the array $L[\text{start}_t : \text{end}_t[$ where the elements from the start_t to the ' $\text{end}_t - 1$ ' indices gives the indices of the samples reaching node t .

Let us take a small example with a set of 10 training samples $L = [0, 1, \dots, 9]$ illustrating the management of the array L . During the tree growth (see Algorithm 3.2) when we partition the sample set $\mathcal{L}_t = \{1, \dots, 10\}$ into two sample sets $\mathcal{L}_{t,l} = \{9, 1, 5, 3\}$ and $\mathcal{L}_{t,r} = \{2, 7, 6, 4, 8, 0\}$. In practice, we modify the array L such that from 0 to $|\mathcal{L}_{1,l}| = 4$ (resp. from $|\mathcal{L}_{1,l}| = 4$ to $|\mathcal{L}_1| = 10$) are located the samples of the left child (resp. right child). It leads to

$$L = [9, 1, 5, 3, 2, 7, 6, 4, 8, 0].$$

We represent each sample set \mathcal{L}_t as a slice $[\text{start} : \text{end}[$, a chunk, of the array L . The sample set $\mathcal{L}_{1,l}$ is the slice $[0 : 4[$ of L , while the sample set $\mathcal{L}_{1,r}$ is the slice $[4 : 10[$ of L . Now if the right node $\mathcal{L}_{1,r}$ is further split into a left node with samples $\{6, 0\}$ and a right node with samples $\{2, 7, 4, 8\}$ (in orange), then $L[4 : 10[$ is further modified to reflect the split:

$$L = [9, 1, 5, 3, 2, 7, 4, 8, 6, 0].$$

To further speed up the best splitting rule search with ordered variables, we sort the possible thresholds associated to an ordered input variable x_j (programmatically $\text{sort}(X_j[L[\text{start}_t : \text{end}_t[[]])$). By sorting the possible thresholds sets, we can evaluate the impurity measure I and the impurity reduction ΔI in an online fashion. For instance, the Gini index and entropy criteria can be computed by updating the class frequency in the left and right split when moving from one splitting threshold to the next.

8.1.2 Splitting rules search on sparse data

To handle sparse input data with decision trees, we need efficient procedures to select the best splitting rules s_t among all the possible splitting rules knowing that we have a high proportion of zeros in the input matrix X . In this section, we propose an efficient method to exploit the sparsity of the input space with decision tree models. Our method takes advantage of the input sparsity by avoiding sorting sample sets of a node along a feature unless there are non zero elements at this feature. This approach speeds up training substantially as extracting the possible threshold values and sorting them is a costly but essential and ubiquitous component of tree-based models.

The splitting rule search algorithm for an ordered variable x_j at a node t is divided in two parts (see Algorithm 8.1): (i) to extract efficiently the non zero values associated to x_j in the sample partition \mathcal{L}_t (line 3) and (ii) to search separately among the splitting rules with the positive, negative or zero threshold (line 4).

Algorithm 8.1 Search for the best splitting rule s_t^* given a sparse input matrix X and a set of samples \mathcal{L}_t

```

1: function FINDBESTSPARSESPLIT( $X, \mathcal{L}_t$ )
2:   for  $j = 1$  to  $p$  do
3:     Extract strictly positive  $X_{j, pos}$  and strictly negatives  $X_{j, neg}$ 
       values from  $X_j$  given  $\mathcal{L}_t$ .
4:     Search for the best splitting rule of the form  $s_j^*(x) = x_j \leq$ 
        $\tau$  with  $\tau \in X_{j, pos} \cup \{0\} \cup X_{j, neg}$  maximizing the impurity
       reduction  $\Delta I$  over the sample set  $\mathcal{L}_t$ .
5:     Update  $s^*$  if the splitting rule  $s_j^*$  leads to higher impurity
       reduction.
6:   end for
7:   return  $s^*$ 
8: end function

```

To extract the non zero values of a sparse input matrix $X \in \mathbb{R}^{n \times p}$ with sparsity s in the context of the decision tree growth, we need to perform efficiently two operations on matrices: (i) the column indexing for a given input variable j and (ii) the extraction of the non zero row values associated to the set of samples \mathcal{L}_t reaching the node t in this column j . The overall cost of extracting $|\mathcal{L}_t|$ samples at a column j from the input matrix X should be proportional to the number of non zero elements and not to $|\mathcal{L}_t|$.¹

Among the different sparse matrix representations (Barrett et al., 1994; Hwu and Kirk, 2009; Pissanetzky, 1984), the sparse csc matrix format is the most appropriate for tree growing as it allows efficient

¹ We assume here that the matrix X is uniformly sparse.

column indexing, as required during node splitting. Let us show how to perform an efficient extraction of the non zero values given the sample set \mathcal{L}_t using this matrix format. Note that using a compressed row storage sparse format² would not be appropriate during the tree growth as we need to be able to efficiently subsample input variables at each expansion of a new testing node.

COMPRESSED SPARSE COLUMN (CSC) MATRIX FORMAT

The sparse csc matrix with n_{nz} non zero elements is a data structure composed of three arrays:

`indices` $\in \mathbb{Z}^{n_{nz}}$ containing the row indices of the non zero elements.

`data` $\in \mathbb{R}^{n_{nz}}$ containing the values of the non zero elements.

`indptr` $\in \mathbb{Z}^p$ containing the slice of the non zero elements. For a column $j \in \{1, \dots, p\}$, the row index and the values of the non zero elements of columns j are stored from `indptr[j]` to `indptr[j + 1]` in the `indices` and `data` arrays.

The non zero values associated to an input variable j are located from `indptr[j]` to `indptr[j + 1] - 1` in the `indices` and the `data` arrays (when `indptr[j] = indptr[j + 1]`, the column thus contains only zeros). Extracting them requires to perform a set intersection between the sample set \mathcal{L}_t reaching node t and the non zero values `indices[indptr[j] : indptr[j + 1]]` of the column j .

For instance, the following matrix A has only 3 non zero elements, but we would use an array of 20 elements:

$$A \in \mathbb{R}^{4 \times 5} = \begin{bmatrix} a & 0 & 0 & b & 0 \\ 0 & 0 & 0 & c & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}. \quad (8.3)$$

The csc representation of this matrix A is given by

$$\begin{aligned} \text{data} &= [a \quad b \quad c], \\ \text{indices} &= [0 \quad 0 \quad 1], \\ \text{inptr} &= [0 \quad 1 \quad 1 \quad 1 \quad 3 \quad 3]. \end{aligned}$$

² The compressed row storage (csr) sparse array format is made of three arrays `indptr`, `indices` and `value`. The non zero elements of the i -th row of the sparse csc matrix are stored from `indptr[i]` to `indptr[i + 1]` in the `indices` arrays, giving the column indices, and `value` arrays, giving the stored values. It is the transposed version of the csc sparse format.

Let $n_{nz,j} = (\text{indptr}[j + 1] - \text{indptr}[j]) \forall j$ be the number of samples with non zero values for input variable j and let us assume that the indices of the input csc matrix array are sorted column-wise, i.e. for the j -th row, the elements of indices from $\text{indptr}[j]$ to $\text{indptr}[j + 1] - 1$ are sorted. Standard intersection algorithms have the following time complexity:

1. in $O(|\mathcal{L}_t| \log n_{nz,j})$ by performing $|\mathcal{L}_t|$ binary search on the sorted $n_{nz,j}$ nonzero elements;
2. in $O(|\mathcal{L}_t| \log |\mathcal{L}_t| + n_{nz,j})$ by sorting the sample set \mathcal{L}_t and retrieving the intersection by iterating over both arrays;
3. in $O(n_{nz,j})$ by maintaining a data structure such as a hash table of \mathcal{L}_t allowing to efficiently check if the elements of $\text{indices}[\text{indptr}[j]:\text{indptr}[j + 1][$ are contained in the sample partition \mathcal{L}_t .

The optimal intersection algorithm depends on the number of non zero elements for input variable j and the number of samples $|\mathcal{L}_t|$ reaching node t . During the decision tree growth, we have two opposite situations: either the size of the sample partition $|\mathcal{L}_t|$ is high with respect to the number of non zero elements ($|\mathcal{L}_t| \approx O(n) \gg n_{nz,j}$) or, typically at the bottom of the tree, the partition size is small with respect to the number of non zero elements ($n_{nz,j} \gg |\mathcal{L}_t|$). In the first case (i.e., at the top of the tree), the most efficient approach is thus approach (3), while in the second case (i.e., at the bottom of the tree), approach (1) should be faster. We first describe how to implement approach (3), then approach (1), and finally how to combine both approaches.

A straightforward implementation of approach (3) is, at each node, to allocate a hash table containing all training examples in that node (in $O(|\mathcal{L}_t|)$) and then to compute the intersection by checking if the non zero elements of the csc matrix belong to the hash table (in $O(n_{nz,j})$). We can however avoid the overhead required for the allocation, creation, and deallocation of the hash table by maintaining and exploiting a mapping between the csc matrix and the sample set \mathcal{L}_t . Since the array L is constantly modified during the tree growth, we propose to use an array, denoted `mapping`, to keep track of the position of a sample i in the array L as illustrated in Figure 8.1. During the tree growing, we keep the following invariant:

$$\text{mapping}[L[i]] = i. \tag{8.4}$$

In the above example, the array L was $[0, 1, \dots, 9]$ with `mapping` = $[0, 1, \dots, 9]$. After a few splits, the array L has become

$$L = [9, 1, 5, 3, 2, 7, 4, 8, 6, 0]$$

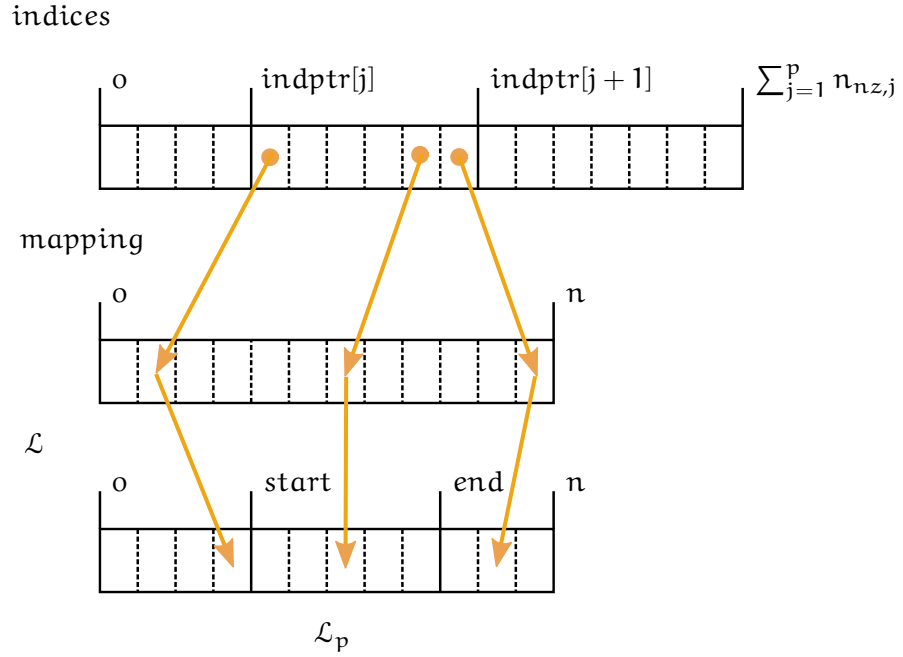


Figure 8.1: The array mapping allows to efficiently compute the intersection between the indices array of the csc matrix and a sample set \mathcal{L}_t .

and the associated mapping array is

$$\text{mapping} = [9, 1, 4, 3, 6, 2, 8, 5, 7, 0].$$

Thanks to the mapping array, we can now check in constant time $O(1)$ whether a sample i belongs to the sample set \mathcal{L}_t . Indeed, given that \mathcal{L}_t is represented by a slice from an index start to an index $\text{end} - 1$ in L , sample i belongs to \mathcal{L}_t if the following condition holds:

$$\text{start} \leq \text{mapping}[i] < \text{end}. \quad (8.5)$$

To extract the non zero values of the csc matrix associated to the j -th input variable in the sample set \mathcal{L}_t , we check the previous condition for all samples from $\text{indptr}[j]$ to $\text{indptr}[j+1] - 1$ in the indices array. Thus, we perform the intersection between \mathcal{L}_t and the $n_{nz,j}$ non zero values in $O(n_{nz,j})$. The whole method is described in Algorithm 8.2. Note that to swap samples in the array L , we use a modified swap function (see Algorithm 8.3) which preserves the mapping invariant.

In practice, the number of non zero elements $n_{nz,j}$ of feature j could be much greater than the size of a sample set \mathcal{L}_t . This is likely to happen near the leaf nodes. Whenever the tree is fully developed, there are only a few samples reaching these nodes. The approach (1) shown in Algorithm 8.4 exploits the relatively small size of the sample set and performs repeated binary search on the $n_{nz,j}$ non zero elements associated to the feature j .

Algorithm 8.2 Return the n_{neg} strictly negative and n_{pos} positive values $(X_{j,neg}, X_{j,pos})$ associated to the j -th variable from the sample set $L[start : end[$ through a *given* mapping *satisfying* $mapping[L[i]] = i$. The array L is modified so that $L[start : start + n_{neg}[$ contains the samples with negatives values, $L[start + n_{neg} : end - n_{pos}[$ contains the zero values and $L[end - n_{pos} : end[$ the samples with positives values.

```

1: function EXTRACT_NNZ_MAPPING( $X, j, L, start, end, mapping$ )
2:    $X_{j,pos} = []$ 
3:    $X_{j,neg} = []$ 
4:    $start_p = end$ 
5:    $end_n = start$ 
6:   for  $k \in [X.indptr[j]:X.indptr[j + 1][$  do
7:      $index = indices[k]$ 
8:      $value = data[k]$ 
9:     if  $start \leq mapping[index] < end$  then
10:       $i = mapping[index]$ 
11:      if  $value > 0$  then
12:         $X_{j,pos}.APPEND(value)$ 
13:         $start_p - = 1$ 
14:         $SWAP(\mathcal{L}, i, start_p, mapping)$ 
15:      else
16:         $X_{j,neg}.APPEND(value)$ 
17:         $SWAP(\mathcal{L}, i, end_n, mapping)$ 
18:         $end_n + = 1$ 
19:      end if
20:    end if
21:  end for
22:  return  $X_{j,pos}, X_{j,neg}, end_n - start, end - start_p$ 
23: end function

```

Algorithm 8.3 Swap two elements at positions p_1 and p_2 in the array L in place while maintaining the invariant of the mapping array.

```

1: function SWAP( $L, p_1, p_2, mapping$ )
2:    $L[p_1], L[p_2] = L[p_2], L[p_1]$ 
3:    $mapping[L[p_1]] = p_1$ 
4:    $mapping[L[p_2]] = p_2$ 
5: end function

```

Algorithm 8.4 Return the n_{neg} strictly negative and n_{pos} positive values ($X_{j,neg}, X_{j,pos}$) associated to the j -th variable from the sample set $L[start : end[$ through *repeated binary search*. The array L is modified so that $L[start : start + n_{neg}[$ contains the samples with negatives values, $L[start + n_{neg} : end - n_{pos}[$ contains the zero values and $L[end - n_{pos} : end[$ the samples with positives values.

```

1: function EXTRACT_NNZ_BSEARCH( $X, j, L, start, end, mapping$ )
2:    $X_{j,pos} = []$ 
3:    $X_{j,neg} = []$ 
4:    $start_p = end$ 
5:    $end_n = start$ 
6:    $indices_j = X.indices[X.indptr[j] : X.indptr[j + 1][$ 
7:    $data_j = X.data[X.indptr[j] : X.indptr[j + 1][$ 
8:    $\mathcal{L} = \text{SORT}(\mathcal{L}, start, end)$ 
9:   for  $i \in [start : end[$  do
10:     // Get the position of  $L[i]$  in  $indices_j$ , and -1 if it is not
    found:  $p = \text{BINARYSEARCH}(L[i], indices_j)$ 
11:     if  $p \neq -1$  then
12:       if  $data_j[p] > 0$  then
13:          $start_{p-} = 1$ 
14:          $X_{j,pos}.APPEND(data_j[p])$ 
15:          $\text{SWAP}(\mathcal{L}, i, start_p, mapping)$ 
16:       else
17:          $X_{j,neg}.APPEND(data_j[p])$ 
18:          $\text{SWAP}(\mathcal{L}, i, end_n, mapping)$ 
19:          $end_{n+} = 1$ 
20:       end if
21:     end if
22:   end for
23:   return  $X_{j,pos}, X_{j,neg}, end_n - start, end - start_p$ 
24: end function

```

The optimal extraction of non zero values is a hybrid approach combining the mapping-based algorithm (Algorithm 8.2) and the binary search algorithm (Algorithm 8.4). Empirical experiments have shown that it is advantageous to use the mapping-based algorithm whenever

$$|\mathcal{L}_t| \times \log(n_{nz,j}) < 0.1 \times n_{nz,j}. \quad (8.6)$$

and the binary search otherwise (see Algorithm 8.5). The formula is based on the computational complexity of both algorithms. We have determined the constant of 0.1 empirically.

Algorithm 8.5 Return the n_{neg} strictly negative and n_{pos} positive values $(X_{j,neg}, X_{j,pos})$ associated to the j -th variable from the sample set $L[start : end[$. The array L is modified so that $L[start : start + n_{neg}[$ contains the samples with negatives values, $L[start + n_{neg} : end - n_{pos}[$ contains the zero values and $L[end - n_{pos} : end[$ the samples with positives values.

```

function EXTRACT_NNZ( $X, j, L, start, end, mapping$ )
  Let  $n_{nz,j}$  be the number of non zero values in column  $j$  of  $X$ .
  if  $(end - start) \times \log(n_{nz,j}) < 0.1 \times n_{nz,j}$  then
    return EXTRACT_NNZ_MAPPING( $X, j, L, start, end, mapping$ )
  else
    return EXTRACT_NNZ_BSEARCH( $X, j, L, start, end, mapping$ )
  end if
end function

```

Note that after extracting the non zero values, we need to sort the thresholds of the splitting rules to search efficiently for the best one. Thanks to Algorithm 8.5, Algorithm 8.2 and Algorithm 8.4, we have already made a three way partition pivot as in the quicksort on the value 0. This speeds up the overall splitting algorithm (the line 4 of Algorithm 8.1). Instead of sorting the thresholds in the sample set \mathcal{L}_t in $O(|\mathcal{L}_t \log(\mathcal{L}_t)|)$, we can perform the sort in $O(d\mathcal{L}_t \log(d\mathcal{L}_t))$ given an input space density d .

As a further refinement, let us note that we can sometimes significantly speed up the decision tree growth by avoiding to search for splitting rules on constant input variables. To do so, we can cache the input variables that were found constant during the expansion of the parents of the node of t . If an input variable is found constant, caching this information avoids the overhead of searching for a splitting rule when no valid one exists.

8.1.3 Partitioning sparse data

During the tree growth, we need to partition a sample set $\mathcal{L}_t = \{(x, y) \in \mathcal{X} \times \mathcal{Y}\}_{i=1}^n$ at a testing node t according to a splitting rule

$s_t(x)$. The splitting rule s_t associated to an ordered input variable is of the form $s_t(x) = 1(x_{F_t} \leq \tau_t)$, where τ_t is a threshold constant on the F_t -th input variable.

During the tree growth, we have the constraint that the input data matrix is in the csc sparse format. We can not convert the current sparse format to another one as it would require to store both the new and old representations into memory. An efficient way to split the sample set \mathcal{L}_t into its left $\mathcal{L}_{t,l}$ and right $\mathcal{L}_{t,r}$ subset is to use the Algorithm 8.5. It will extract the non zero values of a given input variable, but also partition the array L representing the sample set \mathcal{L}_t into three parts: (i) $L[\text{start} : \text{start} + n_{\text{neg}}[$ contains the n_{neg} samples with negatives values, (ii) $L[\text{start} + n_{\text{neg}} : \text{end} - n_{\text{pos}}[$ contains the elements with zero values and (iii) $L[\text{end} - n_{\text{pos}} : \text{end}]$ the n_{pos} samples with positives values. Once the non zero values have been extracted, we have to partition the samples either with negative values ($L[\text{start} : \text{start} + n_{\text{neg}}]$) or with positive values ($L[\text{end} - n_{\text{pos}} : \text{end}]$) according to the sign of the threshold τ_t .

The complexity to partition once the data is $O(\min(n_{nz,j}, |\mathcal{L}_t| \times \log(n_{nz,j})))$ for a batch of \mathcal{L}_t samples instead of the usual $O(|\mathcal{L}_t|)$ with dense input data.

8.2 TREE PREDICTION

The prediction of an unseen sample x by a decision tree (see Algorithm 3.1) is done by traversing the tree from the top to the bottom. At each test node, a splitting rule of the form tests whether or not the sample x should go in the left or the right branch. The splitting rule s_t associated to an ordered input variable is of the form $s_t(x) = 1(x_{F_t} \leq \tau_t)$, where τ_t is a threshold constant on the F_t -th input variable.

We need to have an efficient row and column indexing of the input data matrix. We discuss here two options: (i) using the dictionary of key (dok) sparse matrix format and (ii) using a csr sparse matrix format³.

The dictionary of key (dok) sparse matrix format store the non zero values in a hash table whose keys are the pairs formed from the row and the column index. It is straightforward to apply Algorithm 3.1 with the dok format. The computational complexity is thus unchanged.

To predict one or several unseen samples using a csr array, we need a procedure to efficiently access to both the row and the column index without densifying the csr matrix. We allocate two ar-

³ The compressed row storage (csr) sparse array format (Barrett et al., 1994; Hwu and Kirk, 2009; Pissanetzky, 1984) is made of three arrays `indptr`, `indices` and `value`. The non zero elements of the i -th row of the sparse csc matrix are stored from `indptr[i]` to `indptr[i + 1]` in the `indices` arrays, giving the column indices, and `value` arrays, giving the stored values. It is the transposed version of the csc sparse format.

rays $\text{nz_mask} \in \mathbb{Z}^p$ with all elements having the value “-1” and $\text{nz_value} \in \mathbb{R}^p$ of size p . To predict the i -th sample from a test set, we set in the array nz_mask the value i to the non zero values $\text{indices}[\text{indptr}[i] : \text{indptr}[i + 1]]$ associated to this sample. The array nz_value is modified to contain the non zero values of the i -th sample, i.e. $\text{values}[\text{indptr}[i] : \text{indptr}[i + 1]]$. During the tree traversal (see Algorithm 3.1), we get the j -th input value by first checking if it is zero with $\text{nz_mask}[j] \neq i$, otherwise the value is stored at $\text{nz_value}[j]$. Assuming a proportion of zero elements s for a batch of n test samples with p input variables, the extra cost of using this approach is $O(p + (1 - s)np)$. Note that using the nz_mask and the nz_value arrays is more efficient than densifying each sample as it would add an extra cost of $O(np)$.

The csr sparse format leads to a worse computational complexity than the dok format, which has no extra computing cost. However, the csr approach was chosen in the scikit-learn machine learning python library. The standard implementation of the dok format in scipy is indeed currently implemented using the python dict object. While with the csr format, we can work only with low level c arrays and we can also easily release the global interpreter lock (GIL). Note that here the csr format is also better suited than the csc format as the complexity is independent of the number of samples to predict at once.

8.3 EXPERIMENTS

In this section, we compare the training time and prediction time of decision tree growing and prediction algorithms using either dense data representation or sparse data representation. More specifically, we compare three input matrix layouts using scikit-learn version 0.17: (i) the dense c array layout, a row major order layout whose consecutive and contiguous elements are row values, (ii) the dense fortran array layout, a column major order layout whose consecutive and contiguous elements are column values, and (iii) the sparse array layout, the csc sparse format during tree growing and the csr sparse format during tree prediction (as proposed respectively in Sections 8.1.2 and 8.1.3). The comparison will be made with stumps, a decision tree with a single test node, and fully grown decision trees. These results will be indicative of what could be gained in the context of boosting methods using decision tree of low complexity such as stumps and random forest methods using deep decision trees. Note that all splitting algorithms lead exactly to the same decision tree structure and have the same generalization performance.

We assess the effect of the input space density on synthetic datasets in Section 8.3.1. Then, we compare training times using each of the three input matrix layouts on real datasets in Section 8.3.2.

8.3.1 Effect of the input space density on synthetic datasets

As a synthetic experiment, we compare the decision tree growing algorithm and tree prediction algorithm on synthetic regression tasks with $n = 10^5$ samples and $p = 10^3$ features. The input matrices are sparse random matrices whose non zero elements are drawn uniformly in $\mathcal{N}(0;1)$. The output vector is drawn uniformly at random in $[0, 1]$. The input space density ranges from 10^{-3} to 1. Each point is an average over 20 experiments.

Figure 8.2a shows in logarithmic scale the computing times to grow a single stump. We first note that column-based layouts (the fortran and csc format) are more appropriate to grow decision tree on sparse data. While the density is ranging from 10^{-3} to 10^{-1} , the most expensive part of the tree growing for a single stump is to retrieve the input values from a sample set and to sort them. For a set of n samples and p features with a sparsity s , the computational cost to grow a stump on dense data is $O(pn \log n)$. By contrast, growing a stump using a csc sparse input matrix has a computational complexity of

$$\begin{aligned} &O(p(1-s)n \log((1-s)n) + \min\{pn \log(n(1-s)), pn(1-s)\}) \\ &= O(p(1-s)n \log((1-s)n)). \end{aligned} \quad (8.7)$$

The first term of Equation 8.7 corresponds to the sorting of non zero elements. The second term highlights the contribution of the retrieval of the non zero values, which is always less costly than the sorting operation. If the density is 1 ($s = 0$), both the dense and sparse formats have the same complexity as shown in the right point of Figure 8.2a. Overall with a sparse dataset, the csc format is significantly faster as it leverages the sparsity. The bad performance of the dense c layout compared to the fortran layout or csc layout can be explained by the higher number of cache misses.

The time required to predict the training set using the fitted stump is shown in Figure 8.2b. The only difference between the three input matrix layouts (c, fortran and csr) is the access to the non zero elements. The differences between the dense and the sparse input matrix format can be explained by a better exploitation of the cache for the sparse format, especially when the density is below 0.02. When the density is over 0.02, the cost of copying the non zero values to the arrays `nz_mask` and `nz_values` becomes dominant.

Figure 8.3 shows the time required to learn a fully grown decision tree as a function of the dataset density. Note that the maximal depth decreases as a function of the dataset density (see Figure 8.3c) as the decision tree becomes more balanced. As with the stump, the fortran layout is more appropriate than the c layout to grow a fully grown decision tree on sparse data. The sparse csc algorithm is significantly faster than the dense splitting algorithm if the density is sufficiently low (here below 0.02). The sparse splitting algorithm becomes slower

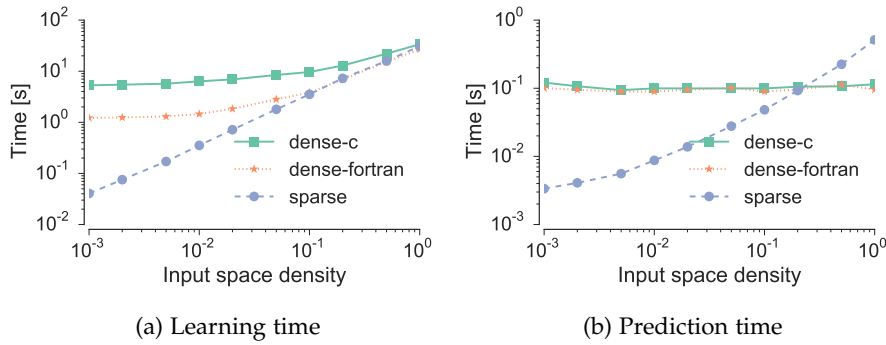


Figure 8.2: Learning and prediction time of stumps as a function of the input space density.

Table 8.1: Dataset properties ordered by input space density.

Dataset	n	p	Density
news20.binary (Keerthi and DeCoste, 2005)	19996	1355191	0.0003
20newsgroup (Lang, 1995; Rennie and Rifkin, 2001)	11314	130107	0.0012
rcv1 (Bekkerman and Scholz, 2008)	23149	47236	0.0016
sector-scale (Keerthi et al., 2008; Rennie and Rifkin, 2001)	9619	55197	0.0029
farm-ads-vect (Mesterharm and Pazzani, 2011)	4143	54877	0.0036
E2006-train (Kogan et al., 2009)	16087	150360	0.0083
mushrooms (Lichman, 2013)	8124	112	0.1875
mnist (LeCun et al., 1998)	70000	784	0.1914
covtype (Lichman, 2013)	581012	54	0.2110

as the density increases (here beyond 0.02) since the extraction of the non zero values becomes more costly than finding the right split.

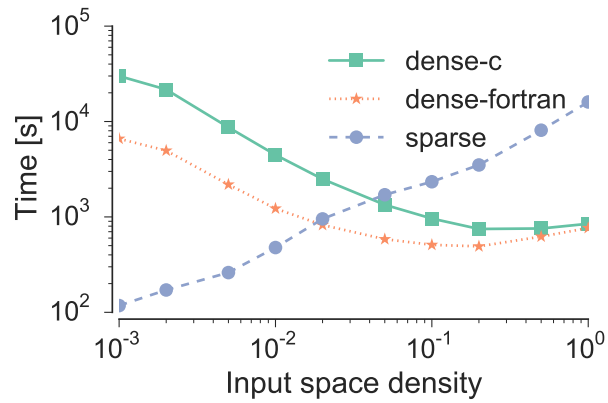
With fully developed decision trees, the prediction time (see Figure 8.3b) is similar between its dense and sparse version. The prediction time is lower when the input space density is high as the trees are more balanced. We note that the prediction time is correlated with the maximal depth of the tree.

Whenever the complexity of the decision tree lies between a stump and a fully developed tree, the behavior moves from one extreme to the other.

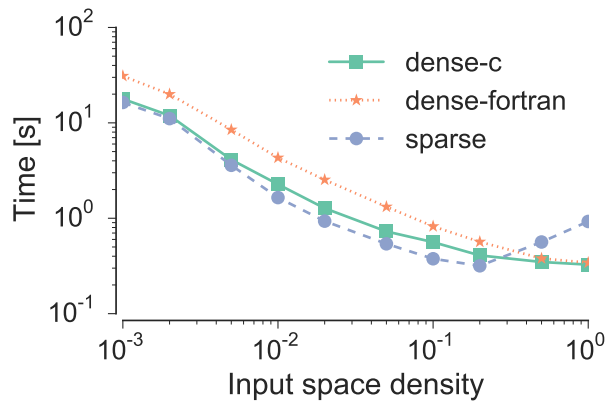
8.3.2 Effect of the input space density on real datasets

To further study the impact of the input space density, we have selected 9 datasets whose input space density ranges from 0.00034 to 0.22. These datasets are presented in Table 8.1 ordered by input space density.

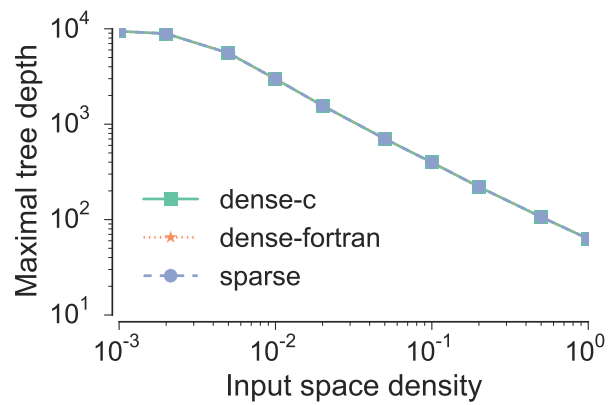
Table 8.2 shows the time to train a single stump. The fastest algorithm here is the tree growing algorithm with the input sparse csc



(a) Learning time



(b) Prediction time



(c) Maximal depth

Figure 8.3: Learning time, prediction time and maximal depth as a function of the input space as a function of the input space density for fully grown decision trees.

Table 8.2: The time (in second) required to train a stump using a sparse csc layout is always faster than with the fortran or c dense memory layout on all sparse selected datasets.

Dataset	c	fortran	sparse	fortran / sparse
news20.binary	N/A	N/A	3.31	N/A
20newsgroup	213.18	15.77	0.79	20.1
rcv1	197.61	28.90	15.97	1.8
sector-scale	58.56	7.10	1.08	6.6
farm-ads-vect	26.03	3.22	0.14	23.0
E2006-train	413.35	35.44	8.62	4.1
mushrooms	0.03	0.02	0.02	1.2
mnist	5.23	2.65	2.00	1.3
covtype	3.89	2.57	2.35	1.1

matrix. The speed up factor between the sparse and fortran memory layout ranges from 1.1 to 23 times. Note that the fortran layout is always faster than the c layout. The column major order layout is here better suited for sparse dataset. The difference could be explained by fewer cache misses with the fortran memory layout than with the c memory layout.

Table 8.3 shows the time required to grow a fully developed decision tree with c, fortran or sparse csc memory layout. The dense fortran layout is here always faster than the dense c layout. The sparse memory layout is faster by a factor between 1.3 and 7.5 than the fortran layout when the input space density is below 0.3%.

Note that we were unable to grow a decision tree with a dense memory layout on the news20.binary dataset as it would require 108.4 Gigabyte to only store the input matrix instead of 78 Megabytes.

8.3.3 Algorithm comparison on 20 newsgroup

Decision trees are rarely used in the context of sparse input datasets. One reason is the lack of implementations exploiting sparsity during the decision tree growth. With the previous experiments, we have shown that it increases significantly the computing time, but also the amount of memory needed. With the proposed tree growing and prediction algorithms, it is interesting to compare the training time, prediction time and accuracy of some tree based models, such as random forest or adaboost, with methods more commonly used in the presence of sparse data.

We compare tree based methods to methods more commonly used on sparse datasets on the 20 newsgroup dataset, which have

Table 8.3: The time required to train a fully grown decision tree using a c, a fortran or a sparse csc memory layout.

Dataset	c	fortran	sparse	fortran / sparse
news20.binary	N/A	N/A	428.89	N/A
20newsgroup	4281.02	518.34	69.06	7.5
rcv1	1458.45	562.19	442.15	1.3
sector-scale	5337.65	878.67	206.69	4.3
farm-ads-vect	227.24	45.71	7.95	5.7
E2006-train	2467.39	752.15	1083.62	0.7
mushrooms	0.07	0.05	0.04	1.3
mnist	80.09	56.11	178.59	0.3
covtype	53.51	33.82	240.26	0.1

$p = 130107$ input variables, 11314 training sample and 7532 testing samples.

The compared algorithms are optimized on their respective hyperparameters (see Table 8.4 for the details) using 5-fold cross validation strategy on the training samples.

The results obtained on the 20 newsgroup dataset are shown in Table 8.5 using scikit-learn version 0.17.1 and input sparsity-aware implementations. The algorithm with the highest accuracy is the linear estimator trained with ridge regression. It is closely followed by the random forest ($m = 1000$) model, the multinomial naives Bayes and extra-trees ($m = 1000$). More generally, tree-based ensemble methods (random forest, extra trees, and adaboost) show similar performance as linear methods (ridge, naive bayes, linear SVC and SGD), with all methods from these two families reaching at least 0.75 of accuracy. On the other hand, the k-nearest neighbors and the single decision tree perform very poorly (with an accuracy below 0.6).

We also note that increasing the number of trees from 100 to 1000 significantly improves the performance of both random forests and extra trees. Their accuracy increases respectively by 0.0573 and 0.0409 in absolute value. Building tree ensembles also very significantly improves the accuracy with respect to single trees (by at least 0.20). This further suggests that the variance of single trees is very high on this problem.

From a modeling perspective, growing decision tree ensemble on datasets with sparse inputs is possible. From a training time perspective, the time needed to grow and to optimize an ensemble of 100 trees is comparable to the time needed to train linear models, e.g., with SGD or ridge regression. Note that naive Bayes models are particularly fast to train compared to the other estimators. However note

Table 8.4: Hyper-parameters grids.

<i>k</i> -nearest neighbors	
<i>k</i> , number of neighbors	$\{1, \dots, 10\}$
<i>Decision tree</i>	
n_{\min} , min. number of samples to split a node	$\{2, 5, 10, 15\}$
<i>Extra trees, random forest</i>	
n_{\min} , min. number of samples to split a node	$\{2, 5, 10, 15\}$
<i>k</i> , number of features drawn at each nodes	$\{1, \log p, \sqrt{p}, 0.001p\}$
<i>M</i> , ensemble size	100 or 1000
<i>Adaboost with decision trees as weak estimators</i>	
n_{\min} , min. number of samples to split a node	$\{2, 5, 10, 15\}$
<i>k</i> , number of features drawn at each nodes	$\{1, \log p, \sqrt{p}, 0.001p\}$
<i>M</i> , ensemble size	100
μ , learning rate	$\{1, 0.1, 0.01\}$
<i>Bernoulli and multinomial naive Bayes</i>	
λ , additive smoothing parameter	$\{0, 0.001, 0.01, 0.1, 1\}$
<i>Ridge classifier</i>	
λ , penalty parameter	$\{0.0001, 0.001, 0.01, 0.1, 1, 10\}$
<i>Linear support vector machine classifier (SVC)</i>	
λ , penalty parameter	$\{10^i\}_{i=-5}^4$
<i>Stochastic gradient descent (SGD)</i>	
Loss	{hinge, logistic}
Penalty constraint	$\{\ell_1, \ell_2, \text{elastic net}\}$
λ , penalty parameter	$\{0.0001, 0.001, 0.01, 0.1, 1, 10\}$
Number of iterations	100

Table 8.5: Accuracy, training time (in second) and prediction time (in second) of algorithms on 20 newsgroup sorted by accuracy score.

Estimator	Training time	Prediction time	Accuracy
K-nearest neighbors	295	22.35	0.457
Decision tree	859	0.02	0.557
Adaboost	37867	8.84	0.752
Bernoulli naives Bayes	25	0.52	0.765
Random forest (m = 100)	11766	12.71	0.778
Extra trees (m = 100)	22186	28.07	0.794
SGD	42776	0.25	0.814
Linear SVC	2481	0.20	0.822
Extra trees (m = 1000)	213009	253.39	0.833
Multinomial naives Bayes	14	0.11	0.833
Random forest (m = 1000)	114472	125.92	0.835
Ridge	5737	0.13	0.844

that the comparison is dependent upon the chosen hyper-parameters, the implementation and the grid size. From the point of view of prediction time, tree ensemble methods are particularly slow compared to the other estimators.

8.4 CONCLUSION

We propose an algorithm to grow decision tree models whenever the input space is sparse. Our approach takes advantage of input sparsity to speed up the search and selection of the best splitting rule during the tree growing. It first speeds up the expansion of a tree node by extracting efficiently the non zero threshold of the splitting rules used to partition data. Secondly, the selection of best splitting rule is also faster as we avoid to sort data with zero values. We reduce the memory needed as we do not need to densify the input space. We also show how to predict samples with sparse inputs.

CONCLUSIONS

9.1 CONCLUSIONS

As we now gather or generate data at every moment, machine learning techniques are emerging as ubiquitous tools in sciences, engineering, or society. Within machine learning, this thesis focuses on supervised learning, which aims at modelling input-output relationships only from observations of input-output pairs, using tree-based ensemble methods, a popular method family exploiting tree structured input-output models. Modern applications of supervised learning raise new computational, memory, and modeling challenges to existing supervised learning algorithms. In this work, we identified and addressed the following questions in the context of tree-based supervised learning methods: (i) how to efficiently learn in high dimensional, and possibly sparse, output spaces? (ii) how to reduce the memory requirement of tree-based models at prediction time? (iii) how to efficiently learn in high dimensional and sparse input spaces? We summarize below our main contributions and conclusions around these three questions.

LEARNING IN HIGH DIMENSIONAL AND POSSIBLY SPARSE OUTPUT SPACES. Decision trees are grown by recursively partitioning the input space while selecting at each node the split maximizing the reduction of some impurity measure. Impurity measures have been extended to address with such models multi-dimensional output spaces, so as to solve multi-label classification or multi-target regression problems. However, when the output space is of high dimension, the computation of the impurity becomes a computational bottleneck of the tree growing algorithm. To speed up this algorithm, we propose to approximate impurity computations during tree growing through random projections of the output space. More precisely, before growing a tree, a few random projections of the output space are computed and the tree is grown to fit these projections instead of the original outputs. Tree leaves are then relabelled in the original output space to provide proper predictions at test time. We show theoretically that when the number of projections is large enough, impurity scores, and thereby the learned tree structures and their predictions, are not affected by this trick. We then exploit the randomization introduced by the projections in the context of random forests, by building each tree of the forest from a different randomly projected subspace. Through experiments on several multi-label clas-

sification problems, we show that randomly projecting the outputs can significantly reduce computing times at training without affecting predictive performance. On some problems, the randomization induced by the projections even allows to reach a better bias-variance tradeoff within random forests, which leads to improved overall performance. In contrast with existing works on random projections of the output, our proposed leaf relabelling strategy also allows to avoid any decoding step and thus preserves computational efficiency at prediction time with respect to standard unprojected multi-output random forests.

Multi-output random forests build a single tree ensemble to predict all outputs simultaneously. While often effective, this approach is justified only when the individual outputs are strongly dependent (conditionally to the inputs). On the other hand, building a separate ensemble for each output, as done in the binary relevance / single target approach, is justified only when the outputs are (conditionally) independent. In our second contribution, we build on gradient boosting and random projections to propose a new approach that tries to bridge the gap between these two extreme assumptions and can hopefully adapt automatically to any intermediate output dependency structure. The idea of this approach is to grow each tree of a gradient boosting ensemble to fit a random projection of the original (residual) output space and then to weight this model in the prediction of each output according to its “correlation” with this output. Through extensive experiments on several artificial and real problems, we show that the resulting method has a faster convergence than binary relevance and that it can adapt better than both binary relevance and multi-output gradient boosting to any output dependency structure. The resulting method is also competitive with multi-output random forests. Although we only carried out experiments with tree-based weak models, the resulting gradient boosting methods are generic and can thus be used with any base regressor.

REDUCING MEMORY REQUIREMENTS OF TREE-BASED MODELS AT PREDICTION TIME. One drawback of random forest methods is that they need to grow very large ensembles of unpruned trees to achieve optimal performance. The resulting models are thus potentially very complex, especially with large datasets, as the complexity of unpruned trees typically depends linearly on the dataset size. On the other hand, only very few nodes are required to make a prediction for a given test example. Our investigation of the question of ensemble compression started with the observation that the random forest model can be viewed as linear models in the node indicator space. Each of these binary variables defining this space indicates whether or not a sample reaches a given node of the forest. In the original linear representation of a forest in the indicator space, non

zero “coefficients” are given only to the leaf nodes. We propose to post-prune the random forest model by selecting and re-weighting the nodes of the linear model through the exploitation of sparse linear estimators. More precisely, from the tree ensemble, we first extract node indicator variables. Then, we project a sample set on this new representation and select a subset of these variables through a Lasso model. The non zero coefficients of the obtained Lasso model are later used to prune and to re-weight the decision tree structure. The resulting post-pruning approach is shown experimentally to reduce very significantly the size of random forests, while preserving, and even sometimes improving, their predictive performance.

LEARNING IN HIGH DIMENSIONAL AND SPARSE INPUT SPACES. Some supervised learning tasks (e.g., text classification) need to deal with high dimensional and sparse input spaces, where input variables have each only a few non zero values. Dealing with input sparsity in the context of decision tree ensembles is challenging computationally for two reasons: (i) it is more difficult algorithmic-wise to exploit sparsity during the tree growth than for example with linear models, leading to slow tree training algorithms requiring a high amount of memory, (ii) the decision tree structures are very unbalanced, which further affects computational complexity. For these two reasons, linear methods are often preferred to decision tree algorithms to learn with sparse datasets. In our last contribution, we specifically developed an efficient implementation of the tree growing algorithm to handle sparse input variables. While previous implementations required to densify the input data matrix, our implementation allows to directly fit decision trees on appropriate sparse input matrices. It speeds up decision tree training on sparse input data and saves memory by avoiding input data “densification”. We also show how to predict unseen sparse input samples with a decision tree model. Note that in this contribution we only focus on improving computing times without modifying the original algorithm.

9.2 PERSPECTIVES AND FUTURE WORKS

We collect in this section some future research directions based on the presented ideas of this thesis.

9.2.1 *Learning in compressed space through random projections*

- The combination of random forest models with random projections adds two new hyper-parameters to tree based methods: the choice and the size of the random output projection subspace. It is not clear yet what would be good default hyper-parameter choices. Extensive empirical studies and the Johnson-

Lindenstrauss lemma might help us to define good default values. These two hyper-parameters also introduce randomization in the output space. It would be interesting to further investigate how the input and output space randomizations jointly modify the bias-variance tradeoff of the ensemble.

- Single random projection of the output space with the gradient boosting algorithm is a generic multi-output method usable with any single output regressor. In this thesis, we specifically focused on tree-structured weak models. We suggest to investigate other kinds of weak models.
- We have combined a dimensionality reduction method with an ensemble method, random forest methods in Chapter 5 and with gradient boosting methods in Chapter 6, while keeping the generation of the random projection matrix independent from the supervised learning task. It would be interesting to investigate more adaptive projection schemes. A simple instance of this approach would be to draw a new random projection matrix according to the residuals, e.g. by sub-sampling an output variable with a probability proportional to the fraction of unexplained variance. An optimal instance of this approach, but computationally more expensive, would compute a projection maximizing the variance along each axis with the principal component analysis algorithm.
- Kernelizing the output of tree-based methods (Geurts et al., 2006b, 2007) allows one to treat complex outputs such as images, texts and graphs. It would be interesting to investigate how to combine output kernel tree-based methods with random projection of the output space to improve computational complexity and accuracy. This idea has been studied (Lopez-Paz et al., 2014) in the context of the kernel principal component algorithm and the kernel canonical correlation analysis algorithm.

9.2.2 *Growing and compressing decision trees*

- In the context of the ℓ_1 -based compression of random forests, we first grow a whole forest, and then prune it. The post-pruning step is costly in terms of computational time and memory as we start from a complex random forest model. A first study in collaboration with Jean-Michel Begon (Begon et al., 2016) shows that we actually do not need to start from the whole forest, and can grow a compressed random forest model greedily. Starting from a set of root nodes, the idea is to sequentially develop the nodes that reduce the most the error of the chosen loss function. The process continues until reaching a complexity constraint, saving time and memory.

- The space complexity of a decision tree model is linear in the number of (test and leaf) nodes, which is typically proportional to the training set size n . In the context of d outputs multi-output classification or regression tasks, or d classes multi-class classification tasks, a leaf node is a constant model stored as a vector of size d . The space complexity of a decision tree model is thus $O(nd)$. With high dimensional output spaces or many classes, it thus may become prohibitive to store decision tree or random forest models. We would like to investigate two further approaches to compress such models: (i) by adapting the (post)-pruning method developed in (Joly et al., 2012) and in Chapter 7 to multi-output tasks and multi-class classification tasks and (ii) by compressing exactly or approximately each constant leaf model. Both approaches can be used together. For approach (ii), an exact solution would compute the constant leaf models on-the-fly at prediction time by storing once the output matrix and the indices of the samples reaching the leaf at learning time. With totally developed trees, it should not modify the computational complexity of the prediction algorithm. If we agree to depart from the vanilla decision tree model, it is also possible to approximate the leaf model, for instance by keeping at the leaf nodes only the subset of the k output-values reducing the most the error either at the leaf level as in (Prabhu and Varma, 2014) or at the tree level. Also, if the output space is sparse, appropriate sparse data structures could help to further reduce the memory footprint of the models.

9.2.3 *Learning in high dimensional and sparse input-output spaces*

- In Chapter 8, we have shown how to improve tree growing efficiency in the case of sparse high-dimensional input spaces. However, the small fraction of non zero input values exploited for each split typically leads to highly unbalanced tree structure. On such tasks, it would be interesting to grow decision trees with multivariate splitting rules to make the tree balanced. For instance in text-based supervised learning, the input text is often converted to a vector through a bag-of-words, where each variable indicates the presence of a single word. In this context, each test node assesses the presence or the absence of a single word. The tree growing algorithm lead to unbalanced trees as each training sample has only a small fraction of all possible words. We propose to investigate node splitting methods that would combine several sparse input variables into a dense one. In the text example, we would generate new dense variables by collapsing several words together. In a more general context,

we could use random “or” or random “addition” functions of several sparse input variables.

Furthermore, while we have shown empirically that the implementation proposed in Chapter 8 indeed translates into a speed up and reduction of memory consumption, it would be interesting to study formally its computational complexity as a function of the input-space sparsity.

- In the multi-label classification task, the output space is often very large and *sparse* (as in Chapter 5 and Chapter 6), having few “non zero values”¹. It would be interesting to exploit the output space sparsity to speed up the decision tree algorithm. The algorithm interacts with the output space during the search of the best split and the training of the leaf models. The search for the best split for an ordered variable is done by first sorting the possible splitting rules according to their threshold values, and then the best split is selected by computing incrementally the reduction of impurity by moving samples from the right partition to the left partition. The leaf models are constant estimators obtained by aggregating output values. Both procedures require efficient sample-wise indexing or row-wise indexing as provided by the compressed row storage (csr) sparse matrix format. We propose to implement impurity functions and leaf model training procedures to work with csr sparse matrices.

¹ We assume that the majority class of each output is coded as “0” and the minority classes is coded with the value “1”.

Part IV

APPENDIX

A

DESCRIPTION OF THE DATASETS

A.1 SYNTHETIC DATASETS

- Friedman₁ (Friedman, 1991) is a regression problem with $p = 10$ independent input variables of uniform distribution $\mathcal{U}(0, 1)$. We try to estimate the output $y = 10 \sin(\pi x_1 x_2) + 20(x_3 - \frac{1}{2})^2 + 10x_4 + 6x_5 + \epsilon$, where ϵ is a Gaussian noise $\mathcal{N}(0, 1)$. There are 300 learning samples and 2000 testing samples.
- Two-norm (Breiman, 1996b) is a binary classification problem with $p = 20$ normally distributed (and class-conditionally independent) input variables: either from $\mathcal{N}(-\alpha, 1)$ if the class is 0 or from $\mathcal{N}(\alpha, 1)$ if the class is 1 (with $\alpha = \frac{2}{\sqrt{20}}$). There are 300 learning and 2000 testing samples.

A.2 REGRESSION DATASET

- SEFTi (AA&YA, 2008) is a (simulated) regression problem which concerns the tool level fault isolation in a semiconductor manufacturing. One quarter of the values are missing at random and were replaced by the median. There are $p = 600$ input variables, 2000 learning samples and 2000 testing samples.

A.3 MULTI-LABEL DATASET

Experiments are performed on several multi-label datasets: the yeast (Elisseeff and Weston, 2001) and the bird (Briggs et al., 2013) datasets in the biology domain; the corel5k (Duygulu et al., 2002) and the scene (Boutell et al., 2004) datasets in the image domain; the emotions (Tsoumakas et al., 2008b) and the CAL500 (Turnbull et al., 2008) datasets in the music domain; the bibtex (Katakis et al., 2008), the bookmarks (Katakis et al., 2008), the delicious (Tsoumakas et al., 2008a), the enron (Klimt and Yang, 2004), the EUR-Lex (subject matters, directory codes and eurovoc descriptors) (Mencía and Fürnkranz, 2010) the genbase (Diplaris et al., 2005), the medical¹, the tmc2007 (Srivastava and Zane-Ulman, 2005) datasets in the text domain and the mediamill (Snoek et al., 2006) dataset in the video domain.

¹ The medical dataset comes from the computational medicine center's 2007 medical natural language processing challenge <http://computationalmedicine.org/challenge/previous>.

Several hierarchical classification tasks are also studied to increase the diversity in the number of label and treated as multi-label classification task. Each node of the hierarchy is treated as one label. Nodes of the hierarchy which never occurred in the training or testing set were removed. The Reuters (Rousu et al., 2005), WIPO (Rousu et al., 2005) datasets are from the text domain. The Diatoms (Dimitrovski et al., 2012) dataset is from the image domain. SCOP-GO (Clare, 2003), Yeast-GO (Barutcuoglu et al., 2006) and Expression-GO (Vens et al., 2008) are from the biological domain. Missing values in the Expression-GO dataset were inferred using the median for continuous features and the most frequent value for categorical features using the entire dataset. The inference of a drug-protein interaction network (Yamanishi et al., 2011) is also considered either using the drugs to infer the interactions with the protein (drug-interaction), either using the proteins to infer the interactions with the drugs (protein-interaction).

Those datasets were selected to have a wide range of number of outputs d . Their basic characteristics are summarized at Table A.1. For more information on a particular dataset, please see the relevant paper.

A.4 MULTI-OUTPUT REGRESSION DATASETS

Multi-output regression is evaluated on several real world datasets: the edm (Karalič and Bratko, 1997) dataset in the industrial domain; the water-quality (Džeroski et al., 2000) dataset in the environmental domain; the atp1d (Spyromitros-Xioufis et al., 2016), the atp7d (Spyromitros-Xioufis et al., 2016), the scm1d (Spyromitros-Xioufis et al., 2016) and the scm2od (Spyromitros-Xioufis et al., 2016) datasets in the price prediction domain; the oes97 (Spyromitros-Xioufis et al., 2016) and the oes10 (Spyromitros-Xioufis et al., 2016) datasets in the human resource domain. The output of those datasets were normalized to have zero mean and unit variance.

If the number of testing samples is unspecified, we use a 50% of the samples as training and validation set and 50% of the samples as testing set.

Table A.1: Selected datasets have a number of labels d ranging from 6 up to 3993 in the biology, the text, the image, the video or the music domain. Each dataset has n_{LS} training samples, n_{TS} testing samples and p input features.

Datasets	n_{LS}	n_{TS}	p	d
emotions	391	202	72	6
scene	1211	1196	2407	6
yeast	1500	917	103	14
birds	322	323	260	19
tmc2007	21519	7077	49060	22
genbase	463	199	1186	27
reuters	2500	5000	19769	34
medical	333	645	1449	45
enron	1123	579	1001	53
mediamill	30993	12914	120	101
Yeast-GO	2310	1155	5930	132
bibtex	4880	2515	1836	159
CAL500	376	126	68	174
WIPO	1352	358	74435	188
EUR-Lex (subject matters)	19348	10-cv	5000	201
bookmarks	65892	21964	2150	208
diatoms	2065	1054	371	359
corel5k	4500	500	499	374
EUR-Lex (directory codes)	19348	10-cv	5000	412
SCOP-GO	6507	3336	2003	465
delicious	12920	3185	500	983
drug-interaction	1396	466	660	1554
protein-interaction	1165	389	876	1862
Expression-GO	2485	551	1288	2717
EUR-Lex (eurovoc descriptors)	19348	10-cv	5000	3993

Table A.2: Selected multi-output regression ranging from $d = 2$ to $d = 16$ outputs.

Datasets	n_{LS}	n_{TS}	p	d
atp1d	337		411	6
atp7d	296		411	6
edm	154		16	2
oes10	403		298	16
oes97	334		263	16
scm1d	8145	1658	280	16
scm2od	7463	1503	61	16
water-quality	1060		16	14

BIBLIOGRAPHY

- I. AA&YA. Manufacturing data: Semiconductor tool fault isolation, 11 2008. URL <mailto:causality@clopin.net>. (Cited on page 180.)
- D. Achlioptas. Database-friendly random projections: Johnson-lindenstrauss with binary coins. *Journal of computer and System Sciences*, 66(4):671–687, 2003. (Cited on pages 43 and 114.)
- R. Agrawal, A. Gupta, Y. Prabhu, and M. Varma. Multi-label learning with millions of labels: Recommending advertiser bid phrases for web pages. In *Proceedings of the 22nd international conference on World Wide Web*, pages 13–24. International World Wide Web Conferences Steering Committee, 2013. (Cited on pages 2 and 85.)
- A. Airola, T. Pahikkala, W. Waegeman, B. De Baets, and T. Salakoski. An experimental comparison of cross-validation techniques for estimating the area under the roc curve. *Computational Statistics & Data Analysis*, 55(4):1828–1844, 2011. (Cited on page 25.)
- H. Almuallim. An efficient algorithm for optimal pruning of decision trees. *Artificial Intelligence*, 83(2):347–362, 1996. (Cited on page 58.)
- Y. Amit, D. Geman, and K. Wilder. Joint induction of shape features and tree classifiers. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 19(11):1300–1305, 1997. (Cited on page 73.)
- B. Bakker and T. Heskes. Task clustering and gating for bayesian multitask learning. *Journal of Machine Learning Research*, 4(May):83–99, 2003. (Cited on page 37.)
- L. Baldassarre, L. Rosasco, A. Barla, and A. Verri. Multi-output learning via spectral filtering. *Machine Learning*, 87(3):259–301, 2012. (Cited on page 21.)
- R. Barrett, M. Berry, T. F. Chan, J. Demmel, J. Donato, J. Dongarra, V. Eijkhout, R. Pozo, C. Romine, and H. V. der Vorst. *Templates for the Solution of Linear Systems: Building Blocks for Iterative Methods, 2nd Edition*. SIAM, Philadelphia, PA, 1994. (Cited on pages 157 and 164.)
- Z. Barutcuoglu, R. E. Schapire, and O. G. Troyanskaya. Hierarchical multi-label prediction of gene function. *Bioinformatics*, 22(7):830–836, 2006. (Cited on pages 107 and 181.)
- E. Bauer and R. Kohavi. An empirical comparison of voting classification algorithms: Bagging, boosting, and variants. *Machine learning*, 36(1-2):105–139, 1999. (Cited on page 71.)

- J.-M. Begon, A. Joly, and P. Geurts. Joint learning and pruning of decision forests. In *Belgian-Dutch Conference On Machine Learning*, 2016. (Cited on pages 4 and 176.)
- R. Bekkerman and M. Scholz. Data weaving: Scaling up the state-of-the-art in data clustering. In *Proceedings of the 17th ACM conference on Information and knowledge management*, pages 1083–1092. ACM, 2008. (Cited on page 167.)
- Y. Bengio. Practical recommendations for gradient-based training of deep architectures. In *Neural Networks: Tricks of the Trade*, pages 437–478. Springer, 2012. (Cited on page 16.)
- J. Bergstra and Y. Bengio. Random search for hyper-parameter optimization. *The Journal of Machine Learning Research*, 13(1):281–305, 2012. (Cited on page 40.)
- J. S. Bergstra, R. Bardenet, Y. Bengio, and B. Kégl. Algorithms for hyper-parameter optimization. In *Advances in Neural Information Processing Systems*, pages 2546–2554, 2011. (Cited on page 40.)
- S. Bernard, L. Heutte, and S. Adam. On the selection of decision trees in random forests. In *Neural Networks, 2009. IJCNN 2009. International Joint Conference on*, pages 302–307. IEEE, 2009. (Cited on page 147.)
- R. Blaser and P. Fryzlewicz. Random rotation ensembles. *Journal of Machine Learning Research*, 2:1–15, 2015. (Cited on page 72.)
- H. Blockeel, L. De Raedt, and J. Ramon. Top-down induction of clustering trees. *arXiv preprint cs/0011032*, 2000. (Cited on pages 2, 22, 60, 61, and 107.)
- P. Bloomfield and W. Steiger. *Least absolute deviations: Theory, applications and algorithms*, volume 6. Springer Science & Business Media, 2012. (Cited on page 14.)
- M. Bohanec and I. Bratko. Trading accuracy for simplicity in decision trees. *Machine Learning*, 15(3):223–250, 1994. (Cited on page 58.)
- H. Borchani, G. Varando, C. Bielza, and P. Larrañaga. A survey on multi-output regression. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 5(5):216–233, 2015. (Cited on page 18.)
- L. Bottou. Stochastic gradient descent tricks. In *Neural Networks: Tricks of the Trade*, pages 421–436. Springer, 2012. (Cited on page 3.)
- N. Boulanger-Lewandowski, Y. Bengio, and P. Vincent. Modeling temporal dependencies in high-dimensional sequences: Application to polyphonic music generation and transcription. *arXiv preprint arXiv:1206.6392*, 2012. (Cited on page 16.)

- M. R. Boutell, J. Luo, X. Shen, and C. M. Brown. Learning multi-label scene classification. *Pattern recognition*, 37(9):1757–1771, 2004. (Cited on page 180.)
- U. M. Braga-Neto and E. R. Dougherty. Is cross-validation valid for small-sample microarray classification? *Bioinformatics*, 20(3):374–380, 2004. (Cited on page 26.)
- L. Breiman. Bagging predictors. *Machine learning*, 24(2):123–140, 1996a. (Cited on pages 66, 71, and 73.)
- L. Breiman. Bias, variance, and arcing classifiers. *Statistics*, 1996b. (Cited on page 180.)
- L. Breiman. Randomizing outputs to increase prediction accuracy. *Machine Learning*, 40(3):229–242, 2000. (Cited on page 72.)
- L. Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001. (Cited on pages 2, 18, 22, 73, 74, 107, and 146.)
- L. Breiman and J. H. Friedman. Predicting multivariate responses in multiple linear regression. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 59(1):3–54, 1997. (Cited on page 21.)
- L. Breiman, J. Friedman, C. J. Stone, and R. A. Olshen. *Classification and regression trees*. CRC press, 1984. (Cited on pages 2, 17, 47, 54, 57, 58, 59, and 107.)
- R. P. Brent. *Algorithms for minimization without derivatives*. Courier Corporation, 2013. (Cited on page 81.)
- F. Briggs, Y. Huang, R. Raich, K. Eftaxias, Z. Lei, W. Cukierski, S. F. Hadley, A. Hadley, M. Betts, X. Z. Fern, et al. The 9th annual mlsp competition: New methods for acoustic classification of multiple simultaneous bird species in a noisy environment. In *Machine Learning for Signal Processing (MLSP), 2013 IEEE International Workshop on*, pages 1–8. IEEE, 2013. (Cited on page 180.)
- K. Brinker and E. Hüllermeier. Case-based multilabel ranking. In *IJCAI*, pages 702–707, 2007. (Cited on page 22.)
- K. Brinker, J. Fürnkranz, and E. Hüllermeier. A unified model for multilabel classification and ranking. In *Proceedings of the 2006 conference on ECAI 2006: 17th European Conference on Artificial Intelligence August 29–September 1, 2006, Riva del Garda, Italy*, pages 489–493. IOS Press, 2006. (Cited on page 20.)
- P. Büchlmann and B. Yu. Analyzing bagging. *Annals of Statistics*, pages 927–961, 2002. (Cited on page 71.)

- L. Buitinck, G. Louppe, M. Blondel, F. Pedregosa, A. Mueller, O. Grisel, V. Niculae, P. Prettenhofer, A. Gramfort, J. Grobler, et al. Api design for machine learning software: experiences from the scikit-learn project. *arXiv preprint arXiv:1309.0238*, 2013. (Cited on pages 3, 5, 27, 104, 123, and 154.)
- E. Candès and M. Wakin. An introduction to compressive sampling. *Signal Processing Magazine, IEEE*, 25(2):21–30, 2008. (Cited on page 153.)
- E. J. Candes and Y. Plan. A probabilistic and ripless theory of compressed sensing. *Information Theory, IEEE Transactions on*, 57(11):7235–7254, 2011. (Cited on pages 43 and 101.)
- R. Caruana, N. Karampatziakis, and A. Yessenalina. An empirical evaluation of supervised learning in high dimensions. In *Proceedings of the 25th international conference on Machine learning*, pages 96–103. ACM, 2008. (Cited on pages 2 and 73.)
- C. Chen, A. Liaw, and L. Breiman. Using random forest to learn imbalanced data. *University of California, Berkeley*, 2004. (Cited on page 71.)
- W. Cheng and E. Hüllermeier. A simple instance-based approach to multilabel classification using the mallows model. In *Working Notes of the First International Workshop on Learning from Multi-Label Data*, pages 28–38, 2009. (Cited on page 22.)
- W. Cheng, E. Hüllermeier, and K. J. Dembczynski. Bayes optimal multilabel classification via probabilistic classifier chains. In *Proceedings of the 27th international conference on machine learning (ICML-10)*, pages 279–286, 2010. (Cited on pages 19 and 126.)
- T.-H. Chiang, H.-Y. Lo, and S.-D. Lin. A ranking-based knn approach for multi-label classification. *ACML*, 25:81–96, 2012. (Cited on page 22.)
- P. M. Ciarelli, E. Oliveira, C. Badue, and A. F. De Souza. Multi-label text categorization using a probabilistic neural network. *International Journal of Computer Information Systems and Industrial Management Applications*, 1(133-144):40, 2009. (Cited on page 21.)
- M. M. Cisse, N. Usunier, T. Artieres, and P. Gallinari. Robust bloom filters for large multilabel classification tasks. In *Advances in Neural Information Processing Systems*, pages 1851–1859, 2013. (Cited on pages 20, 86, and 105.)
- A. Clare. *Machine learning and data mining for yeast functional genomics*. PhD thesis, The University of Wales, 2003. (Cited on page 181.)

- A. Clare and R. D. King. Knowledge discovery in multi-label phenotype data. In *European Conference on Principles of Data Mining and Knowledge Discovery*, pages 42–53. Springer, 2001. (Cited on pages 22 and 60.)
- M. Collins, R. E. Schapire, and Y. Singer. Logistic regression, adaboost and bregman distances. *Machine Learning*, 48(1-3):253–285, 2002. (Cited on page 78.)
- C. Cortes and V. Vapnik. Support-vector networks. *Machine learning*, 20(3):273–297, 1995. (Cited on pages 3, 14, and 44.)
- J. Davis and M. Goadrich. The relationship between precision-recall and roc curves. In *Proceedings of the 23rd international conference on Machine learning*, pages 233–240. ACM, 2006. (Cited on page 30.)
- B. S. Dayal and J. F. MacGregor. Multi-output process identification. *Journal of Process Control*, 7(4):269–282, 1997. (Cited on page 21.)
- G. De’Ath. Multivariate regression trees: a new technique for modeling species–environment relationships. *Ecology*, 83(4):1105–1117, 2002. (Cited on pages 22, 60, and 61.)
- O. Dekel and O. Shamir. Multiclass-multilabel classification with more classes than examples. In *International Conference on Artificial Intelligence and Statistics*, pages 137–144, 2010. (Cited on pages 2 and 85.)
- C. Delierneux, N. Layios, A. Hego, J. Huart, A. Joly, P. Geurts, P. Damas, C. Lecut, A. Gothot, and C. Oury. Elevated basal levels of circulating activated platelets predict icu-acquired sepsis and mortality: a prospective study. *Critical Care*, 19(Suppl 1):P29, 2015a. ISSN 1364-8535. doi: 10.1186/cc14109. URL <http://ccforum.com/content/19/S1/P29>. (Cited on page 4.)
- C. Delierneux, N. Layios, A. Hego, J. Huart, A. Joly, P. Geurts, P. Damas, C. Lecut, A. Gothot, and C. Oury. Prospective analysis of platelet activation markers to predict severe infection and mortality in intensive care units. In *JOURNAL OF THROMBOSIS AND HAEMOSTASIS*, volume 13, pages 651–651. WILEY-BLACKWELL 111 RIVER ST, HOBOKEN 07030-5774, NJ USA, 2015b. (Cited on page 4.)
- K. Dembczynski, W. Waegeman, W. Cheng, and E. Hüllermeier. On label dependence in multi-label classification. In *Workshop proceedings of learning from multi-label data*, pages 5–12. Citeseer, 2010. (Cited on page 107.)
- J. Demšar. Statistical comparisons of classifiers over multiple data sets. *Journal of Machine learning research*, 7(Jan):1–30, 2006. (Cited on page 136.)

- T. G. Dietterich and G. Bakiri. Error-correcting output codes: A general method for improving multiclass inductive learning programs. Citeseer. (Cited on page 20.)
- T. G. Dietterich and E. B. Kong. Machine learning bias, statistical bias, and statistical variance of decision tree algorithms. Technical report, Technical report, Department of Computer Science, Oregon State University, 1995. (Cited on page 73.)
- I. Dimitrovski, D. Kocev, S. Loskovska, and S. Džeroski. Hierarchical classification of diatom images using ensembles of predictive clustering trees. *Ecological Informatics*, 7(1):19–29, 2012. (Cited on page 181.)
- S. Diplaris, G. Tsoumakas, P. A. Mitkas, and I. Vlahavas. Protein classification with multiple algorithms. In *Advances in Informatics*, pages 448–456. Springer, 2005. (Cited on page 180.)
- P. Domingos. A unified bias-variance decomposition. In *Proceedings of 17th International Conference on Machine Learning*, pages 231–238, 2000. (Cited on pages 66 and 71.)
- H. Drucker. Improving regressors using boosting techniques. In *ICML*, volume 97, pages 107–115, 1997. (Cited on page 78.)
- R. Duroux and E. Scornet. Impact of subsampling and pruning on random forests. *arXiv preprint arXiv:1603.04261*, 2016. (Cited on page 153.)
- P. Duygulu, K. Barnard, J. F. de Freitas, and D. A. Forsyth. Object recognition as machine translation: Learning a lexicon for a fixed image vocabulary. In *Computer Vision—ECCV 2002*, pages 97–112. Springer, 2002. (Cited on page 180.)
- S. Džeroski, D. Demšar, and J. Grbović. Predicting chemical parameters of river water quality from bioindicator data. *Applied Intelligence*, 13(1):7–17, 2000. (Cited on page 181.)
- B. Efron. Bootstrap methods: another look at the jackknife. *The annals of Statistics*, pages 1–26, 1979. (Cited on page 71.)
- B. Efron. Estimating the error rate of a prediction rule: improvement on cross-validation. *Journal of the American Statistical Association*, 78(382):316–331, 1983. (Cited on page 25.)
- A. Elisseeff and J. Weston. A kernel method for multi-labelled classification. In *Advances in neural information processing systems*, pages 681–687, 2001. (Cited on pages 21 and 180.)

- T. Evgeniou and M. Pontil. Regularized multi-task learning. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 109–117. ACM, 2004. (Cited on page 21.)
- T. Evgeniou, C. A. Micchelli, and M. Pontil. Learning multiple tasks with kernel methods. *Journal of Machine Learning Research*, 6(Apr): 615–637, 2005. (Cited on page 21.)
- T. Fawcett. An introduction to roc analysis. *Pattern recognition letters*, 27(8):861–874, 2006. (Cited on page 30.)
- M. Fernández-Delgado, E. Cernadas, S. Barro, and D. Amorim. Do we need hundreds of classifiers to solve real world classification problems? *The Journal of Machine Learning Research*, 15(1):3133–3181, 2014. (Cited on pages 2 and 73.)
- C.-S. Ferng and H.-T. Lin. Multi-label classification with error-correcting codes. In *ACML*, pages 281–295, 2011. (Cited on page 20.)
- C. Ferri, J. Hernández-Orallo, and R. Modroiu. An experimental comparison of performance measures for classification. *Pattern Recognition Letters*, 30(1):27–38, 2009. (Cited on page 27.)
- G. Forman and M. Scholz. Apples-to-apples in cross-validation studies: pitfalls in classifier performance measurement. *ACM SIGKDD Explorations Newsletter*, 12(1):49–57, 2010. (Cited on page 25.)
- E. Frank, Y. Wang, S. Inglis, G. Holmes, and I. H. Witten. Using model trees for classification. *Machine Learning*, 32(1):63–76, 1998. (Cited on page 55.)
- Y. Freund and R. E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of computer and system sciences*, 55(1):119–139, 1997. (Cited on pages 18 and 76.)
- J. Friedman. Multivariate adaptive regression splines. *The Annals of Statistics*, pages 1–67, 1991. (Cited on pages 65, 123, and 180.)
- J. Friedman, T. Hastie, R. Tibshirani, et al. Additive logistic regression: a statistical view of boosting (with discussion and a rejoinder by the authors). *The annals of statistics*, 28(2):337–407, 2000. (Cited on page 77.)
- J. H. Friedman. A recursive partitioning decision rule for nonparametric classification. *IEEE Transactions on Computers*, (4):404–408, 1977. (Cited on page 47.)
- J. H. Friedman. On bias, variance, 0/1—loss, and the curse-of-dimensionality. *Data mining and knowledge discovery*, 1(1):55–77, 1997. (Cited on page 66.)

- J. H. Friedman. Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pages 1189–1232, 2001. (Cited on pages 2, 18, 80, 81, and 107.)
- J. H. Friedman. Stochastic gradient boosting. *Computational Statistics & Data Analysis*, 38(4):367–378, 2002. (Cited on page 81.)
- J. H. Friedman and B. E. Popescu. Predictive learning via rule ensembles. *The Annals of Applied Statistics*, pages 916–954, 2008. (Cited on page 147.)
- J. Fürnkranz, E. Hüllermeier, E. L. Mencía, and K. Brinker. Multilabel classification via calibrated label ranking. *Machine learning*, 73(2):133–153, 2008. (Cited on page 20.)
- M. Gasse, A. Aussem, and H. Elghazel. On the optimality of multi-label classification under subset zero-one loss for distributions satisfying the composition property. In *Proceedings of the 32nd International Conference on Machine Learning (ICML-15)*, pages 2531–2539, 2015. (Cited on page 107.)
- S. Geman, E. Bienenstock, and R. Doursat. Neural networks and the bias/variance dilemma. *Neural computation*, 4(1):1–58, 1992. (Cited on page 64.)
- P. Geurts. Some enhancements of decision tree bagging. *Principles of Data Mining and Knowledge Discovery*, pages 141–148, 2000. (Cited on page 146.)
- P. Geurts. *Contributions to decision tree induction: bias/variance tradeoff and time series classification*. PhD thesis, University of Liège Belgium, 2002. (Cited on pages 68, 69, 70, 71, and 73.)
- P. Geurts, D. Ernst, and L. Wehenkel. Extremely randomized trees. *Machine learning*, 63(1):3–42, 2006a. (Cited on pages 22, 73, 97, 118, 146, 147, and 150.)
- P. Geurts, L. Wehenkel, and F. d’Alché Buc. Kernelizing the output of tree-based methods. In *Proceedings of the 23rd international conference on Machine learning*, pages 345–352. Acm, 2006b. (Cited on pages 2, 23, and 176.)
- P. Geurts, L. Wehenkel, and F. d’Alché Buc. Gradient boosting for kernelized output spaces. In *Proceedings of the 24th international conference on Machine learning*, pages 289–296. ACM, 2007. (Cited on pages 107, 120, and 176.)
- N. Ghamrawi and A. McCallum. Collective multi-label classification. In *Proceedings of the 14th ACM international conference on Information and knowledge management*, pages 195–200. ACM, 2005. (Cited on pages 21 and 32.)

- E. Gibaja and S. Ventura. Multi-label learning: a review of the state of the art and ongoing research. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 4(6):411–444, 2014. (Cited on page 18.)
- X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *International conference on artificial intelligence and statistics*, pages 249–256, 2010. (Cited on page 16.)
- S. Godbole and S. Sarawagi. Discriminative methods for multi-labeled classification. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, pages 22–30. Springer, 2004. (Cited on pages 33 and 36.)
- A. Graves, A.-r. Mohamed, and G. Hinton. Speech recognition with deep recurrent neural networks. In *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*, pages 6645–6649. IEEE, 2013. (Cited on page 16.)
- Y. Z. Guo, K. Ramamohanarao, and L. A. Park. Error correcting output coding-based conditional random fields for web page prediction. In *Web Intelligence and Intelligent Agent Technology, 2008. WI-IAT'08. IEEE/WIC/ACM International Conference on*, volume 1, pages 743–746. IEEE, 2008. (Cited on page 20.)
- V. Guruswami and A. Sahai. Multiclass learning, boosting, and error-correcting codes. In *Proceedings of the twelfth annual conference on Computational learning theory*, pages 145–155. ACM, 1999. (Cited on page 20.)
- J. A. Hanley and B. J. McNeil. The meaning and use of the area under a receiver operating characteristic (roc) curve. *Radiology*, 143(1):29–36, 1982. (Cited on page 30.)
- T. Hastie, J. Taylor, R. Tibshirani, G. Walther, et al. Forward stagewise regression and the monotone lasso. *Electronic Journal of Statistics*, 1: 1–29, 2007. (Cited on pages 13 and 148.)
- T. Hastie, R. Tibshirani, and J. Friedman. *The elements of statistical learning*. Springer series in statistics Springer, Berlin, 2009. (Cited on pages 11, 13, 25, 50, and 123.)
- T. K. Ho. The random subspace method for constructing decision forests. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 20(8):832–844, 1998. (Cited on pages 43, 71, and 73.)
- A. E. Hoerl and R. W. Kennard. Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12(1):55–67, 1970. (Cited on page 13.)

- M. Hossin and M. Sulaiman. A review on evaluation metrics for data classification evaluations. *International Journal of Data Mining & Knowledge Management Process*, 5(2):1, 2015. (Cited on page 27.)
- C.-W. Hsu, C.-C. Chang, C.-J. Lin, et al. A practical guide to support vector classification, 2003. (Cited on pages 39 and 44.)
- D. Hsu, S. Kakade, J. Langford, and T. Zhang. Multi-label prediction via compressed sensing. In *NIPS*, volume 22, pages 772–780, 2009. (Cited on pages 20, 86, 105, and 144.)
- S.-J. Huang, Y. Yu, and Z.-H. Zhou. Multi-label hypothesis reuse. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 525–533. ACM, 2012. (Cited on page 107.)
- D. H. Hubel and T. N. Wiesel. Receptive fields and functional architecture of monkey striate cortex. *The Journal of physiology*, 195(1): 215–243, 1968. (Cited on page 16.)
- E. Hüllermeier, J. Fürnkranz, W. Cheng, and K. Brinker. Label ranking by learning pairwise preferences. *Artificial Intelligence*, 172(16):1897–1916, 2008. (Cited on page 20.)
- F. Hutter, H. H. Hoos, and K. Leyton-Brown. Sequential model-based optimization for general algorithm configuration. In *Learning and Intelligent Optimization*, pages 507–523. Springer, 2011. (Cited on page 40.)
- W.-m. Hwu and D. Kirk. Programming massively parallel processors. *Special Edition*, 92, 2009. (Cited on pages 157 and 164.)
- L. Hyafil and R. L. Rivest. Constructing optimal binary decision trees is np-complete. *Information Processing Letters*, 5(1):15–17, 1976. (Cited on page 52.)
- A. J. Izenman. Reduced-rank regression for the multivariate linear model. *Journal of multivariate analysis*, 5(2):248–264, 1975. (Cited on page 21.)
- T. S. Jaakkola and D. Haussler. Probabilistic kernel regression models. In *AISTATS*, 1999. (Cited on page 44.)
- G. M. James. Variance and bias for general loss functions. *Machine Learning*, 51(2):115–135, 2003. (Cited on page 66.)
- A. Jiang, C. Wang, and Y. Zhu. Calibrated rank-svm for multi-label image categorization. In *2008 IEEE International Joint Conference on Neural Networks (IEEE World Congress on Computational Intelligence)*, pages 1450–1455. IEEE, 2008. (Cited on page 21.)

- W. B. Johnson and J. Lindenstrauss. Extensions of lipschitz mappings into a hilbert space. *Contemporary mathematics*, 26(189-206):1, 1984. (Cited on page 43.)
- I. Jolliffe. *Principal component analysis*. Wiley Online Library, 2002. (Cited on pages 41 and 44.)
- A. Joly, F. Schnitzler, P. Geurts, and L. Wehenkel. L1-based compression of random forest models. In *European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning*, 2012. (Cited on pages 3, 5, 147, and 177.)
- A. Joly, P. Geurts, and L. Wehenkel. Random forests with random projections of the output space for high dimensional multi-label classification. In *Machine Learning and Knowledge Discovery in Databases*, pages 607–622. Springer Berlin Heidelberg, 2014. (Cited on pages 3, 5, 22, 74, 107, and 123.)
- T. Kajdanowicz and P. Kazienko. Multi-label classification using error correcting output codes. *International Journal of Applied Mathematics and Computer Science*, 22(4):829–840, 2012. (Cited on page 20.)
- A. Kapoor, R. Viswanathan, and P. Jain. Multilabel classification using bayesian compressed sensing. In *Advances in Neural Information Processing Systems*, pages 2645–2653, 2012. (Cited on pages 20 and 144.)
- A. Karalič and I. Bratko. First order regression. *Machine Learning*, 26(2-3):147–176, 1997. (Cited on page 181.)
- I. Katakis, G. Tsoumakas, and I. Vlahavas. Multilabel text classification for automated tag suggestion. In *Proceedings of the ECML/PKDD*, 2008. (Cited on page 180.)
- S. S. Keerthi and D. DeCoste. A modified finite newton method for fast solution of large scale linear svms. *Journal of Machine Learning Research*, 6(Mar):341–361, 2005. (Cited on page 167.)
- S. S. Keerthi, S. Sundararajan, K.-W. Chang, C.-J. Hsieh, and C.-J. Lin. A sequential dual method for large scale multi-class linear svms. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 408–416. ACM, 2008. (Cited on page 167.)
- B. Klimt and Y. Yang. The enron corpus: A new dataset for email classification research. In *Machine learning: ECML 2004*, pages 217–226. Springer, 2004. (Cited on page 180.)
- D. Kocev, C. Vens, J. Struyf, and S. Džeroski. Ensembles of multi-objective decision trees. In *European Conference on Machine Learning*, pages 624–631. Springer, 2007. (Cited on pages 22 and 107.)

- D. Kocev, C. Vens, J. Struyf, and S. Džeroski. Tree ensembles for predicting structured outputs. *Pattern Recognition*, 46(3):817–833, 2013. (Cited on pages 2, 22, and 107.)
- S. Kogan, D. Levin, B. R. Routledge, J. S. Sagi, and N. A. Smith. Predicting risk from financial reports with regression. In *Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 272–280. Association for Computational Linguistics, 2009. (Cited on page 167.)
- R. Kohavi. Scaling up the accuracy of naive-bayes classifiers: A decision-tree hybrid. In *KDD*, volume 96, pages 202–207. Citeseer, 1996. (Cited on page 55.)
- R. Kohavi, D. H. Wolpert, et al. Bias plus variance decomposition for zero-one loss functions. In *ICML*, volume 96, pages 275–83, 1996. (Cited on page 66.)
- R. Kohavi et al. A study of cross-validation and bootstrap for accuracy estimation and model selection. In *Ijcai*, volume 14, pages 1137–1145, 1995. (Cited on pages 24, 25, and 26.)
- A. Z. Kouzani and G. Nasireding. Multilabel classification by bch code and random forests. *International journal of recent trends in engineering*, 2(1):113–116, 2009. (Cited on page 20.)
- A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C. Burges, L. Bottou, and K. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran Associates, Inc., 2012. URL <http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>. (Cited on page 16.)
- A. Krogh, J. Vedelsby, et al. Neural network ensembles, cross validation, and active learning. *Advances in neural information processing systems*, 7:231–238, 1995. (Cited on page 67.)
- L. Kuncheva, J. J. Rodriguez, et al. Classifier ensembles with a random linear oracle. *Knowledge and Data Engineering, IEEE Transactions on*, 19(4):500–508, 2007. (Cited on page 71.)
- L. I. Kuncheva and J. J. Rodríguez. An experimental study on rotation forest ensembles. In *Multiple Classifier Systems*, pages 459–468. Springer, 2007. (Cited on page 72.)
- N. Landwehr, M. Hall, and E. Frank. Logistic model trees. *Machine Learning*, 59(1-2):161–205, 2005. (Cited on page 55.)

- K. Lang. Newsweeder: Learning to filter netnews. In *Proceedings of the Twelfth International Conference on Machine Learning*, pages 331–339, 1995. (Cited on page 167.)
- Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11): 2278–2324, 1998. (Cited on page 167.)
- Y. LeCun, F. J. Huang, and L. Bottou. Learning methods for generic object recognition with invariance to pose and lighting. In *Computer Vision and Pattern Recognition, 2004. CVPR 2004. Proceedings of the 2004 IEEE Computer Society Conference on*, volume 2, pages II–97. IEEE, 2004. (Cited on page 16.)
- Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *Nature*, 521(7553): 436–444, 2015. (Cited on page 16.)
- Y. A. LeCun, L. Bottou, G. B. Orr, and K.-R. Müller. Efficient backprop. In *Neural networks: Tricks of the trade*, pages 9–48. Springer, 2012. (Cited on page 16.)
- P. Li, T. J. Hastie, and K. W. Church. Very sparse random projections. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 287–296. ACM, 2006. (Cited on pages 43, 101, and 114.)
- M. Lichman. UCI machine learning repository, 2013. URL <http://archive.ics.uci.edu/ml>. (Cited on pages 42 and 167.)
- G. Liu, Z. Lin, and Y. Yu. Multi-output regression on the output manifold. *Pattern Recognition*, 42(11):2737–2743, 2009. (Cited on page 21.)
- D. Lopez-Paz, S. Sra, A. J. Smola, Z. Ghahramani, and B. Schölkopf. Randomized nonlinear component analysis. In *ICML*, pages 1359–1367, 2014. (Cited on page 176.)
- G. Louppe. *Understanding Random Forests: From Theory to Practice*. PhD thesis, University of Liège, 2014. (Cited on page 70.)
- G. Louppe and P. Geurts. Ensembles on random patches. In *Machine Learning and Knowledge Discovery in Databases*, pages 346–361. Springer, 2012. (Cited on page 72.)
- G. Madjarov, D. Kocev, D. Gjorgjevikj, and S. Džeroski. An extensive experimental comparison of methods for multi-label learning. *Pattern Recognition*, 45(9):3084–3104, 2012. (Cited on pages 2, 18, 22, 86, 107, and 143.)
- R. Maree, P. Geurts, J. Piater, and L. Wehenkel. Random subwindows for robust image classification. In *Computer Vision and Pattern*

- Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*, volume 1, pages 34–40. IEEE, 2005. (Cited on page 72.)
- G. Martínez-Muñoz and A. Suárez. Switching class labels to generate classification ensembles. *Pattern Recognition*, 38(10):1483–1494, 2005. (Cited on page 72.)
- G. Martínez-Muñoz, A. Sánchez-Martínez, D. Hernández-Lobato, and A. Suárez. Class-switching neural network ensembles. *Neurocomputing*, 71(13):2521–2528, 2008. (Cited on page 72.)
- G. Martínez-Muñoz, D. Hernández-Lobato, and A. Suárez. An analysis of ensemble pruning techniques based on ordered aggregation. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 31(2): 245–259, 2009. (Cited on page 147.)
- T. Matthew, C.-S. Chen, J. Yu, and M. Wyle. Bayesian and empirical bayesian forests. In *Proceedings of the 32nd International Conference on Machine Learning (ICML-15)*, pages 967–976, 2015. (Cited on page 55.)
- M. Mehta, J. Rissanen, R. Agrawal, et al. Mdl-based decision tree pruning. In *KDD*, volume 21, pages 216–221, 1995. (Cited on page 58.)
- L. Meier, S. Van De Geer, and P. Bühlmann. The group lasso for logistic regression. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 70(1):53–71, 2008. (Cited on page 14.)
- N. Meinshausen. Node harvest. *The Annals of Applied Statistics*, pages 2049–2072, 2010. (Cited on page 147.)
- N. Meinshausen et al. Forest garrote. *Electronic Journal of Statistics*, 3: 1288–1304, 2009. (Cited on page 147.)
- E. L. Mencía and J. Fürnkranz. Efficient pairwise multilabel classification for large-scale problems in the legal domain. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 50–65. Springer, 2008. (Cited on page 20.)
- E. L. Mencía and J. Fürnkranz. *Efficient multilabel classification algorithms for large-scale problems in the legal domain*. Springer, 2010. (Cited on pages 20 and 180.)
- B. H. Menze, B. M. Kelm, D. N. Splitthoff, U. Koethe, and F. A. Hamprecht. On oblique random forests. In *Machine Learning and Knowledge Discovery in Databases*, pages 453–469. Springer, 2011. (Cited on page 74.)

- C. Mesterharm and M. J. Pazzani. Active learning using on-line algorithms. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 850–858. ACM, 2011. (Cited on page 167.)
- C. E. Metz. Basic principles of roc analysis. In *Seminars in nuclear medicine*, volume 8, pages 283–298. Elsevier, 1978. (Cited on page 30.)
- J. Nam, J. Kim, E. L. Mencía, I. Gurevych, and J. Fürnkranz. Large-scale multi-label text classification—revisiting neural networks. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 437–452. Springer, 2014. (Cited on page 21.)
- H. G. Noh, M. S. Song, and S. H. Park. An unbiased method for constructing multilabel classification trees. *Computational Statistics & Data Analysis*, 47(1):149–164, 2004. (Cited on pages 22 and 60.)
- P. Panov and S. Džeroski. *Combining bagging and random subspaces to create better ensembles*. Springer, 2007. (Cited on page 72.)
- S.-H. Park and J. Fürnkranz. Efficient pairwise classification. In *European Conference on Machine Learning*, pages 658–665. Springer, 2007. (Cited on page 20.)
- B. J. Parker, S. Günter, and J. Bedo. Stratification bias in low signal microarray studies. *BMC bioinformatics*, 8(1):326, 2007. (Cited on page 25.)
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, et al. Scikit-learn: Machine learning in python. *The Journal of Machine Learning Research*, 12:2825–2830, 2011. (Cited on pages 5, 27, 104, 123, and 154.)
- S. Pissanetzky. *Sparse Matrix Technology—electronic edition*. Academic Press, 1984. (Cited on pages 157 and 164.)
- Y. Prabhu and M. Varma. Fastxml: A fast, accurate and stable tree-classifier for extreme multi-label learning. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 263–272. ACM, 2014. (Cited on page 177.)
- J. R. Quinlan. Simplifying decision trees. *International journal of man-machine studies*, 27(3):221–234, 1987. (Cited on page 58.)
- J. R. Quinlan. Unknown attribute values in induction. In *Proc. of the Sixth Int. Workshop on Machine Learning*, pages 164–168, 1989. (Cited on page 47.)

- J. R. Quinlan. *C4.5: Programs for Machine Learning*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1993. ISBN 1-55860-238-0. (Cited on page 58.)
- J. R. Quinlan and R. L. Rivest. Inferring decision trees using the minimum description length principle. *Information and computation*, 80(3):227–248, 1989. (Cited on page 58.)
- J. R. Quinlan et al. Learning with continuous classes. In *5th Australian joint conference on artificial intelligence*, volume 92, pages 343–348. Singapore, 1992. (Cited on page 55.)
- A. Rahimi and B. Recht. Random features for large-scale kernel machines. In *Advances in neural information processing systems*, pages 1177–1184, 2007. (Cited on page 44.)
- A. Rahimi and B. Recht. Weighted sums of random kitchen sinks: Replacing minimization with randomization in learning. In *Advances in neural information processing systems*, pages 1313–1320, 2009. (Cited on page 44.)
- J. Read, B. Pfahringer, G. Holmes, and E. Frank. Classifier chains for multi-label classification. *Machine learning*, 85(3):333–359, 2011. (Cited on pages 19, 22, and 107.)
- S. Ren, X. Cao, Y. Wei, and J. Sun. Global refinement of random forest. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 723–730, 2015. (Cited on page 153.)
- J. D. Rennie and R. Rifkin. Improving multiclass text classification with the support vector machine. 2001. (Cited on page 167.)
- J. J. Rodriguez, L. I. Kuncheva, and C. J. Alonso. Rotation forest: A new classifier ensemble method. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 28(10):1619–1630, 2006. (Cited on page 72.)
- F. Rosenblatt. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological review*, 65(6): 386, 1958. (Cited on page 15.)
- J. Rousu, C. Saunders, S. Szedmak, and J. Shawe-Taylor. Learning hierarchical multi-category text classification models. In *Proceedings of the 22nd international conference on Machine learning*, pages 744–751. ACM, 2005. (Cited on page 181.)
- M. Sánchez-Fernández, M. de Prado-Cumplido, J. Arenas-García, and F. Pérez-Cruz. Svm multiregression for nonlinear channel estimation in multiple-input multiple-output systems. *IEEE transactions on signal processing*, 52(8):2298–2307, 2004. (Cited on page 21.)

- R. E. Schapire and Y. Singer. Improved boosting algorithms using confidence-rated predictions. *Machine learning*, 37(3):297–336, 1999. (Cited on pages 33 and 34.)
- R. E. Schapire and Y. Singer. Boostexter: A boosting-based system for text categorization. *Machine learning*, 39(2):135–168, 2000. (Cited on pages 23, 34, and 78.)
- A. Schclar and L. Rokach. Random projection ensemble classifiers. In *Enterprise information systems*, pages 309–316. Springer, 2009. (Cited on page 72.)
- B. Scholkopf and A. J. Smola. *Learning with kernels: support vector machines, regularization, optimization, and beyond*. MIT press, 2001. (Cited on pages 3 and 44.)
- B. Schölkopf, A. Smola, and K.-R. Müller. Kernel principal component analysis. In *International Conference on Artificial Neural Networks*, pages 583–588. Springer, 1997. (Cited on page 44.)
- M. Segal and Y. Xiao. Multivariate random forests. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 1(1):80–87, 2011. (Cited on pages 22 and 107.)
- M. R. Segal. Tree-structured methods for longitudinal data. *Journal of the American Statistical Association*, 87(418):407–418, 1992. (Cited on pages 22, 60, 61, and 107.)
- R. Siciliano and F. Mola. Multivariate data analysis and modeling through classification and regression trees. *Computational Statistics & Data Analysis*, 32(3):285–301, 2000. (Cited on page 60.)
- T. Similä and J. Tikka. Input selection and shrinkage in multiresponse linear regression. *Computational Statistics & Data Analysis*, 52(1):406–422, 2007. (Cited on page 21.)
- C. G. M. Snoek, M. Worring, J. C. Van Gemert, J.-M. Geusebroek, and A. W. M. Smeulders. The challenge problem for automated detection of 101 semantic concepts in multimedia. In *Proceedings of the 14th annual ACM international conference on Multimedia*, pages 421–430. ACM, 2006. (Cited on page 180.)
- J. Snoek, H. Larochelle, and R. P. Adams. Practical bayesian optimization of machine learning algorithms. In *Advances in neural information processing systems*, pages 2951–2959, 2012. (Cited on page 40.)
- M. Sokolova and G. Lapalme. A systematic analysis of performance measures for classification tasks. *Information Processing & Management*, 45(4):427–437, 2009. (Cited on pages 27 and 32.)

- F. J. Solis and R. J.-B. Wets. Minimization by random search techniques. *Mathematics of operations research*, 6(1):19–30, 1981. (Cited on page 40.)
- D. F. Specht. A general regression neural network. *IEEE transactions on neural networks*, 2(6):568–576, 1991. (Cited on page 21.)
- E. Spyromitros-Xioufis, G. Tsoumakas, W. Groves, and I. Vlahavas. Multi-label classification methods for multi-target regression. *Machine learning journal*, 2016. (Cited on pages 2, 18, 19, 106, and 181.)
- A. N. Srivastava and B. Zane-Ulman. Discovering recurring anomalies in text reports regarding complex space systems. In *Aerospace Conference, 2005 IEEE*, pages 3853–3862. IEEE, 2005. (Cited on page 180.)
- J. Struyf and S. Džeroski. Constraint based induction of multi-objective regression trees. In *International Workshop on Knowledge Discovery in Inductive Databases*, pages 222–233. Springer, 2005. (Cited on page 60.)
- A. Sutera, A. Joly, V. François-Lavet, Z. A. Qiu, G. Louppe, D. Ernst, and P. Geurts. Simple connectome inference from partial correlation statistics in calcium imaging. In *JMLR: Workshop and Conference Proceedings*, pages 1–12, 2014. (Cited on page 4.)
- R. Tibshirani. *Bias, variance and prediction error for classification rules*. Citeseer, 1996a. (Cited on page 66.)
- R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 267–288, 1996b. (Cited on pages 13, 147, and 148.)
- R. Tibshirani, M. Saunders, S. Rosset, J. Zhu, and K. Knight. Sparsity and smoothness via the fused lasso. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(1):91–108, 2005. (Cited on page 14.)
- T. M. Tomita, M. Maggioni, and J. T. Vogelstein. Randomer forests. *arXiv preprint arXiv:1506.03410*, 2015. (Cited on page 74.)
- L. Torgo. Functional models for regression tree leaves. In *ICML*, volume 97, pages 385–393. Citeseer, 1997. (Cited on page 55.)
- G. Tsoumakas and I. Vlahavas. Random k-labelsets: An ensemble method for multilabel classification. In *Machine learning: ECML 2007*, pages 406–417. Springer, 2007. (Cited on pages 21, 22, 72, and 101.)
- G. Tsoumakas, I. Katakis, and I. Vlahavas. Effective and efficient multilabel classification in domains with large number of labels. In

- Proc. ECML/PKDD 2008 Workshop on Mining Multidimensional Data (MMD'08)*, pages 30–44, 2008a. (Cited on pages 95 and 180.)
- G. Tsoumakas, I. Katakis, and I. Vlahavas. Mining multi-label data. In *Data mining and knowledge discovery handbook*, pages 667–685. Springer, 2009. (Cited on pages 2, 18, 19, 21, and 106.)
- G. Tsoumakas, E. Spyromitros-Xioufis, A. Vrekou, and I. Vlahavas. Multi-target regression via random linear target combinations. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 225–240. Springer, 2014. (Cited on pages 20, 86, and 144.)
- K. T. G. Tsoumakas, G. Kalliris, and I. Vlahavas. Multi-label classification of music into emotions. In *ISMIR 2008: Proceedings of the 9th International Conference of Music Information Retrieval*, page 325. Lulu. com, 2008b. (Cited on page 180.)
- D. Turnbull, L. Barrington, D. Torres, and G. Lanckriet. Semantic annotation and retrieval of music and sound effects. *Audio, Speech, and Language Processing, IEEE Transactions on*, 16(2):467–476, 2008. (Cited on page 180.)
- A. Van Der Merwe and J. Zidek. Multivariate regression analysis and canonical variates. *Canadian Journal of Statistics*, 8(1):27–39, 1980. (Cited on page 21.)
- E. Vazquez and E. Walter. Multi-output support vector regression. In *13th IFAC Symposium on System Identification*, pages 1820–1825. Citeseer, 2003. (Cited on page 21.)
- C. Vens and F. Costa. Random forest based feature induction. In *Data Mining (ICDM), 2011 IEEE 11th International Conference on*, pages 744–753. IEEE, 2011. (Cited on page 147.)
- C. Vens, J. Struyf, L. Schietgat, S. Džeroski, and H. Blockeel. Decision trees for hierarchical multi-label classification. *Machine Learning*, 73(2):185–214, 2008. (Cited on pages 22, 60, and 181.)
- C. S. Wallace and J. Patrick. Coding decision trees. *Machine Learning*, 11(1):7–22, 1993. (Cited on page 58.)
- Y. Wang and I. H. Witten. Induction of model trees for predicting continuous classes. 1996. (Cited on page 55.)
- J. Xu. An efficient multi-label support vector machine with a zero label. *Expert Systems with Applications*, 39(5):4796–4804, 2012. (Cited on page 21.)
- S. Xu, X. An, X. Qiao, L. Zhu, and L. Li. Multi-output least-squares support vector regression machines. *Pattern Recognition Letters*, 34(9):1078–1084, 2013. (Cited on page 21.)

- Y. Yamanishi, E. Pauwels, H. Saigo, and V. Stoven. Extracting sets of chemical substructures and protein domains governing drug-target interactions. *Journal of chemical information and modeling*, 51(5):1183–1194, 2011. (Cited on pages 100 and 181.)
- R. Yan, J. Tesic, and J. R. Smith. Model-shared subspace boosting for multi-label classification. In *Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 834–843. ACM, 2007. (Cited on page 107.)
- Y. Yang. An evaluation of statistical approaches to text categorization. *Information retrieval*, 1(1-2):69–90, 1999. (Cited on page 36.)
- Z. Younes, F. Abdallah, T. Denoeux, and H. Snoussi. A dependent multilabel classification method derived from the k-nearest neighbor rule. *EURASIP Journal on Advances in Signal Processing*, 2011(1): 1–14, 2011. (Cited on page 22.)
- S. Yu, K. Yu, V. Tresp, and H.-P. Kriegel. Multi-output regularized feature projection. *IEEE Transactions on Knowledge and Data Engineering*, 18(12):1600–1613, 2006. (Cited on page 20.)
- M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68(1):49–67, 2006. (Cited on page 14.)
- H. Zhang. Classification trees for multiple binary responses. *Journal of the American Statistical Association*, 93(441):180–193, 1998. (Cited on pages 22 and 60.)
- M.-L. Zhang. Ml-rbf: Rbf neural networks for multi-label learning. *Neural Processing Letters*, 29(2):61–74, 2009. (Cited on page 21.)
- M.-L. Zhang and K. Zhang. Multi-label learning by exploiting label dependency. In *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 999–1008. ACM, 2010. (Cited on page 107.)
- M.-L. Zhang and Z.-H. Zhou. Multilabel neural networks with applications to functional genomics and text categorization. *IEEE transactions on Knowledge and Data Engineering*, 18(10):1338–1351, 2006. (Cited on page 21.)
- M.-L. Zhang and Z.-H. Zhou. Ml-knn: A lazy learning approach to multi-label learning. *Pattern recognition*, 40(7):2038–2048, 2007. (Cited on page 22.)
- M.-L. Zhang and Z.-H. Zhou. A review on multi-label learning algorithms. *IEEE transactions on knowledge and data engineering*, 26(8): 1819–1837, 2014. (Cited on page 18.)

- Y. Zhang and J. G. Schneider. Multi-label output codes using canonical correlation analysis. In *AISTATS*, pages 873–882, 2011. (Cited on page 20.)
- T. Zhou and D. Tao. Multi-label subspace ensemble. In *International Conference on Artificial Intelligence and Statistics*, pages 1444–1452, 2012. (Cited on page 21.)
- Z.-H. Zhou. *Ensemble methods: foundations and algorithms*. CRC Press, 2012. (Cited on page 66.)
- J. Zhu, H. Zou, S. Rosset, and T. Hastie. Multi-class adaboost. *Statistics and its Interface*, 2(3):349–360, 2009. (Cited on pages 78 and 81.)
- H. Zou and T. Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(2):301–320, 2005. (Cited on page 14.)