

Learning Parameters in Discrete Naive Bayes models by Computing Fibers of the Parametrization Map

Vincent Auvray and Louis Wehenkel

EE & CS Dept. and GIGA-R, University of Liège, Belgium

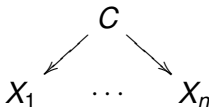
NIPS - AML 2008 - Whistler

Naive Bayesian networks

A discrete naive Bayesian network (or latent class model) with m classes is a distribution p over discrete variables X_1, \dots, X_n such that

$$p(X_1 = x_1, \dots, X_n = x_n) = \sum_{t=1}^m p(C = t) \prod_{i=1}^n p(X_i = x_i | C = t).$$

Graphically, the independencies between C, X_1, \dots, X_n are encoded by



Problem statement

Given a naive Bayesian network p , compute the parameters $p(C = t)$, $p(X_i = x_i | C = t)$ for $i = 1, \dots, n$, $t = 1, \dots, m$, and $x_i \in \mathcal{X}_i$ mapped to p .

Why?

- better understanding of the model
- estimation of parameters
- model selection
- study of parameter identifiability

Outline

Mathematical Results

Applications

Extension

Some notation

Given parameters of a naive Bayesian distribution, we define new parameters

$$w_t = p(C = t),$$
$$A_{x_i}^t = p(X_i = x_i | C = t) - p(X_i = x_i).$$

Given a distribution p , let

$$q(x_{i_1}, \dots, x_{i_k}) = p(x_{i_1}, \dots, x_{i_k})$$
$$- \sum_{\{X_{j_1}, \dots, X_{j_l}\} \subsetneq \{X_{i_1}, \dots, X_{i_k}\}} q(x_{j_1}, \dots, x_{j_l})$$
$$\prod_{X_k \in \{X_{i_1}, \dots, X_{i_k}\} \setminus \{X_{j_1}, \dots, X_{j_l}\}} p(x_k).$$

Some notation

For example, we have

$$\begin{aligned}q(x_i) &= 0, \\q(x_i, x_j) &= p(x_i, x_j) - p(x_i)p(x_j), \\q(x_i, x_j, x_k) &= p(x_i, x_j, x_k) - p(x_i)p(x_j, x_k) - p(x_j)p(x_i, x_k) \\&\quad - p(x_k)p(x_i, x_j) + 2p(x_i)p(x_j)p(x_k).\end{aligned}$$

With this notation, one can see that

$$\begin{aligned}q(x_{i_1}, \dots, x_{i_k}) &= \sum_{t=1}^m w_t \prod_{j=1}^k A_{x_{i_j}}^t, \\w^T A_{x_i} &= 0,\end{aligned}$$

where $w = (w_1 \ \dots \ w_m)^T$ and $A_{x_i} = (A_{x_i}^1 \ \dots \ A_{x_i}^m)^T$.

w is normal to the hyperplane spanned by the A_{x_i}

Consider the parameters of a naive Bayesian distribution.
Given vectors $A_{u_1}, \dots, A_{u_{m-1}}$, we have

$$(-1)^t \det (A_{u_1} \ \dots \ A_{u_{m-1}})^{\hat{t}} = w_t \det (1 \ A_{u_1} \ \dots \ A_{u_{m-1}}),$$

where the superscript \hat{t} denotes the removal of the t th row.

In other words, if

$$\det (1 \ A_{u_1} \ \dots \ A_{u_{m-1}}) \neq 0,$$

then w is the normal to the hyperplane spanned by $A_{u_1}, \dots, A_{u_{m-1}}$ and whose components sum to 1.

The components of A_{x_i} are the roots of a degree m polynomial

For $m = 2$, we have

$$\begin{aligned} s^2 q(u_1, v_1) + s q(x_i, u_1, v_1) - q(x_i, u_1) q(x_i, v_1) \\ = q(u_1, v_1)(s + A_{x_i}^1)(s + A_{x_i}^2). \end{aligned}$$

For $m = 3$, we have

$$\begin{aligned} s^3 \det \begin{pmatrix} q(u_1, v_1) & q(u_1, v_2) \\ q(u_2, v_1) & q(u_2, v_2) \end{pmatrix} + s^2 \left[\det \begin{pmatrix} q(u_1, v_1) & q(u_1, v_2) \\ q(x_i, u_2, v_1) & q(x_i, u_2, v_2) \end{pmatrix} \right. \\ \left. + \det \begin{pmatrix} q(x_i, u_1, v_1) & q(x_i, u_1, v_2) \\ q(u_2, v_1) & q(u_2, v_2) \end{pmatrix} \right] + s \left[- \det \begin{pmatrix} q(u_1, v_1) & q(u_1, v_2) \\ q(x_i, u_2) q(x_i, v_1) & q(x_i, u_2) q(x_i, v_2) \end{pmatrix} \right. \\ \left. - \det \begin{pmatrix} q(x_i, u_1) q(x_i, v_1) & q(x_i, u_1) q(x_i, v_2) \\ q(u_2, v_1) & q(u_2, v_2) \end{pmatrix} + \det \begin{pmatrix} q(x_i, u_1, v_1) & q(x_i, u_1, v_2) \\ q(x_i, u_2, v_1) & q(x_i, u_2, v_2) \end{pmatrix} \right] \\ - \det \begin{pmatrix} q(x_i, u_1, v_1) & q(x_i, u_1, v_2) \\ q(x_i, u_2) q(x_i, v_1) & q(x_i, u_2) q(x_i, v_2) \end{pmatrix} - \det \begin{pmatrix} q(x_i, u_1) q(x_i, v_1) & q(x_i, u_1) q(x_i, v_2) \\ q(x_i, u_2, v_1) & q(x_i, u_2, v_2) \end{pmatrix} \\ = \det \begin{pmatrix} q(u_1, v_1) & q(u_1, v_2) \\ q(u_2, v_1) & q(u_2, v_2) \end{pmatrix} (s + A_{x_i}^1)(s + A_{x_i}^2)(s + A_{x_i}^3). \end{aligned}$$

The components of A_{x_i} are the roots of a degree m polynomial

Given $\mathbf{u} = \{u_1, \dots, u_{m-1}\}$ and $\mathbf{v} = \{v_1, \dots, v_{m-1}\}$, consider the polynomial of degree m

$$\alpha_{X, \mathbf{u}, \mathbf{v}}(s) = s^m \det \begin{pmatrix} q(u_1, v_1) & \cdots & q(u_1, v_{m-1}) \\ \vdots & & \vdots \\ q(u_{m-1}, v_1) & \cdots & q(u_{m-1}, v_{m-1}) \end{pmatrix} + s^{m-1} \cdots + s \cdots + \cdots$$

whose coefficients are sums of determinants. We have

$$\alpha_{X, \mathbf{u}, \mathbf{v}}(s) = \det \begin{pmatrix} q(u_1, v_1) & \cdots & q(u_1, v_{m-1}) \\ \vdots & & \vdots \\ q(u_{m-1}, v_1) & \cdots & q(u_{m-1}, v_{m-1}) \end{pmatrix} \prod_{t=1}^m (s + A_{X_i}^t).$$

The parameters satisfy simple polynomial equations

Consider values $\{x_1, \dots, x_k\}$. The following equation holds

$$\det \begin{pmatrix} \prod_{j=1}^k A_{x_j}^t & q(x_1, \dots, x_k, v_1) & \dots & q(x_1, \dots, x_k, v_{m-1}) \\ A_{u_1}^t & q(u_1, v_1) & \dots & q(u_1, v_{m-1}) \\ \vdots & \vdots & & \vdots \\ A_{u_{m-1}}^t & q(u_{m-1}, v_1) & \dots & q(u_{m-1}, v_{m-1}) \end{pmatrix}$$
$$= q(x_1, \dots, x_k) \det \begin{pmatrix} q(u_1, v_1) & \dots & q(u_1, v_{m-1}) \\ \vdots & & \vdots \\ q(u_{m-1}, v_1) & \dots & q(u_{m-1}, v_{m-1}) \end{pmatrix}.$$

The parameters satisfy simple polynomial equations

For $\{x_1, \dots, x_k\} = \{u_0\}$, we have

$$\det \begin{pmatrix} A_{u_0}^t & q(u_0, v_1) & \dots & q(u_0, v_{m-1}) \\ A_{u_1}^t & q(u_1, v_1) & \dots & q(u_1, v_{m-1}) \\ \vdots & \vdots & & \vdots \\ A_{u_{m-1}}^t & q(u_{m-1}, v_1) & \dots & q(u_{m-1}, v_{m-1}) \end{pmatrix} = 0.$$

For $m = 3$ and $\{x_1, \dots, x_k\} = \{u_1, u_2\}$, we have

$$\begin{aligned} \det \begin{pmatrix} A_{u_1}^t & A_{u_2}^t & q(u_1, u_2, v_1) & q(u_1, u_2, v_2) \\ A_{u_1}^t & & q(u_1, v_1) & q(u_1, v_2) \\ A_{u_2}^t & & q(u_2, v_1) & q(u_2, v_2) \end{pmatrix} \\ = q(u_1, u_2) \det \begin{pmatrix} q(u_1, v_1) & q(u_1, v_2) \\ q(u_2, v_1) & q(u_2, v_2) \end{pmatrix}. \end{aligned}$$

Some determinants have an interpretable decomposition

Consider sets of values $\mathbf{s}_1, \dots, \mathbf{s}_{m-1}$. We have

$$\det \begin{pmatrix} q(\mathbf{s}_1, v_1) & \dots & q(\mathbf{s}_1, v_{m-1}) \\ \vdots & & \vdots \\ q(\mathbf{s}_{m-1}, v_1) & \dots & q(\mathbf{s}_{m-1}, v_{m-1}) \end{pmatrix} \\ = \left(\prod_{t=1}^m w_t \right) \det \begin{pmatrix} 1 & A_{v_1} & \dots & A_{v_{m-1}} \end{pmatrix} \det M,$$

where

$$M = \begin{pmatrix} 1 & \prod_{x \in \mathbf{s}_1} A_x^1 & \dots & \prod_{x \in \mathbf{s}_{m-1}} A_x^1 \\ \vdots & \vdots & & \vdots \\ 1 & \prod_{x \in \mathbf{s}_1} A_x^m & \dots & \prod_{x \in \mathbf{s}_{m-1}} A_x^m \end{pmatrix}$$

Simple implicit equations follow

Consider a naive Bayesian distribution with m classes and consider sets of values $\mathbf{s}_1, \dots, \mathbf{s}_{m'-1}$. If $m' > m$, we have

$$\det \begin{pmatrix} q(\mathbf{s}_1, v_1) & \dots & q(\mathbf{s}_1, v_{m'-1}) \\ \vdots & & \vdots \\ q(\mathbf{s}_{m'-1}, v_1) & \dots & q(\mathbf{s}_{m'-1}, v_{m'-1}) \end{pmatrix} = 0.$$

Consider sets of values $\mathbf{s}_1, \dots, \mathbf{s}_{m-1}$ and $\mathbf{r}_1, \dots, \mathbf{r}_{m-1}$. We have

$$\begin{aligned} & \det \begin{pmatrix} q(\mathbf{s}_1, v_1) & \dots & q(\mathbf{s}_1, v_{m-1}) \\ \vdots & & \vdots \\ q(\mathbf{s}_{m-1}, v_1) & \dots & q(\mathbf{s}_{m-1}, v_{m-1}) \end{pmatrix} \det \begin{pmatrix} q(\mathbf{r}_1, u_1) & \dots & q(\mathbf{r}_1, u_{m-1}) \\ \vdots & & \vdots \\ q(\mathbf{r}_{m-1}, u_1) & \dots & q(\mathbf{r}_{m-1}, u_{m-1}) \end{pmatrix} \\ = & \det \begin{pmatrix} q(\mathbf{s}_1, u_1) & \dots & q(\mathbf{s}_1, u_{m-1}) \\ \vdots & & \vdots \\ q(\mathbf{s}_{m-1}, u_1) & \dots & q(\mathbf{s}_{m-1}, u_{m-1}) \end{pmatrix} \det \begin{pmatrix} q(\mathbf{r}_1, v_1) & \dots & q(\mathbf{r}_1, v_{m-1}) \\ \vdots & & \vdots \\ q(\mathbf{r}_{m-1}, v_1) & \dots & q(\mathbf{r}_{m-1}, v_{m-1}) \end{pmatrix} \end{aligned}$$

Outline

Mathematical Results

Applications

Extension

Potential applications of our results

- Compute the set of parameters mapped to a given naive Bayesian distribution
- Estimate parameters from data by applying the previous computation to the distribution of observed frequencies
- Derive sufficient conditions for parameter identifiability and obtain results on the dimensionality of the model
- Building block in the computation of analytic asymptotic approximations to the marginal likelihood of the model
- Building block in model selection and learning of hidden causes

An important hypothesis to compute the parameters

Suppose that we have a distribution p and sets of values

$$\mathbf{t} = \{t_1, \dots, t_{m-1}\},$$

$$\mathbf{u} = \{u_1, \dots, u_{m-1}\},$$

$$\mathbf{v} = \{v_1, \dots, v_{m-1}\}$$

such that

$$\det \begin{pmatrix} q(t_1, u_1) & \dots & q(t_1, u_{m-1}) \\ \vdots & & \vdots \\ q(t_{m-1}, u_1) & \dots & q(t_{m-1}, u_{m-1}) \end{pmatrix} \neq 0,$$

$$\det \begin{pmatrix} q(t_1, v_1) & \dots & q(t_1, v_{m-1}) \\ \vdots & & \vdots \\ q(t_{m-1}, v_1) & \dots & q(t_{m-1}, v_{m-1}) \end{pmatrix} \neq 0,$$

$$\det \begin{pmatrix} q(u_1, v_1) & \dots & q(u_1, v_{m-1}) \\ \vdots & & \vdots \\ q(u_{m-1}, v_1) & \dots & q(u_{m-1}, v_{m-1}) \end{pmatrix} \neq 0.$$

Computation of w from $A_{u_1}, \dots, A_{u_{m-1}}$

Our hypothesis amounts to

$$\left(\prod_{i=1}^m w_i \right) \det \begin{pmatrix} 1 & A_{t_1} & \dots & A_{t_{m-1}} \end{pmatrix} \\ \det \begin{pmatrix} 1 & A_{u_1} & \dots & A_{u_{m-1}} \end{pmatrix} \\ \det \begin{pmatrix} 1 & A_{v_1} & \dots & A_{v_{m-1}} \end{pmatrix} \neq 0.$$

Hence, we have

$$w_i = \frac{(-1)^i \det \begin{pmatrix} A_{u_1} & \dots & A_{u_{m-1}} \end{pmatrix}^{\hat{i}}}{\det \begin{pmatrix} 1 & A_{u_1} & \dots & A_{u_{m-1}} \end{pmatrix}}.$$

Computation of A_x from $A_{u_1}, \dots, A_{u_{m-1}}$

Since

$$\det \begin{pmatrix} A_{u_1}^t & q(u_1, v_1) & \dots & q(u_1, v_{m-1}) \\ \vdots & \vdots & & \vdots \\ A_{u_{m-1}}^t & q(u_{m-1}, v_1) & \dots & q(u_{m-1}, v_{m-1}) \\ A_x^t & q(x, v_1) & \dots & q(x, v_{m-1}) \end{pmatrix} = 0,$$

we have, for all values x distinct of v_1, \dots, v_{m-1} ,

$$A_x^T = (q(x, v_1) \quad \dots \quad q(x, v_{m-1})) \\ \begin{pmatrix} q(u_1, v_1) & \dots & q(u_1, v_{m-1}) \\ \vdots & & \vdots \\ q(u_{m-1}, v_1) & \dots & q(u_{m-1}, v_{m-1}) \end{pmatrix}^{-1} \\ (A_{u_1} \quad \dots \quad A_{u_{m-1}})^T.$$

Computation of $A_{u_1}, \dots, A_{u_{m-1}}$

Find the roots of the polynomials $\nu_{u_j, \mathbf{t}, \mathbf{v}}$ to obtain

$$\begin{aligned} & \{A_{u_1}^1, \dots, A_{u_1}^m\}, \\ & \quad \vdots \\ & \{A_{u_{m-1}}^1, \dots, A_{u_{m-1}}^m\}. \end{aligned}$$

Note that these sets are not ordered: we are not able to assign each element to its hidden class.

There is some trivial non-identifiability due to the fact that classes can be permuted freely. To remove this degree of freedom from the analysis, we order the set $\{A_{u_1}^1, \dots, A_{u_1}^m\}$ arbitrarily.

Computation of $A_{u_1}, \dots, A_{u_{m-1}}$: a brute force approach

For each ordering of each set $\{A_{u_i}^1, \dots, A_{u_i}^m\}$ with $i = 2, \dots, m - 1$

1. compute a candidate parameter with the previous procedure
2. test if the candidate satisfies the constraints to be a parameter and if it is mapped to the distribution

However, there are $(m!)^{m-2}$ candidate parameters to test.

Corollary: under our hypothesis, there are at most $(m!)^{m-1}$ parameters mapped to the distribution.

Computation of $A_{u_1}, \dots, A_{u_{m-1}}$: a second approach

We have

$$\det M \left(\sum_{p=1}^m \prod_{j=1}^k A_{x_{ij}}^p \right) = q(x_{i_1}, \dots, x_{i_k}) \det M \\ + \sum_{a=1}^{m-1} \sum_{b=1}^{m-1} (-1)^{a+b} q(x_{i_1}, \dots, x_{i_k}, t_a, v_b) \det M_{\hat{a}, \hat{b}},$$

where

$$M = \begin{pmatrix} q(t_1, v_1) & \dots & q(t_1, v_{m-1}) \\ \vdots & & \vdots \\ q(t_{m-1}, v_1) & \dots & q(t_{m-1}, v_{m-1}) \end{pmatrix} \quad (1)$$

We can constraint the orderings to those satisfying the above equation with $\{x_{i_1}, \dots, x_{i_k}\} = \{u_1, u_j\}$.

Computation of $A_{u_1}, \dots, A_{u_{m-1}}$

The previous algorithm do not make use of all our theoretical results. For $m = 3$, recall that we have

$$\det \begin{pmatrix} A_{u_1}^t A_{u_2}^t & q(u_1, u_2, v_1) & q(u_1, u_2, v_2) \\ A_{u_1}^t & q(u_1, v_1) & q(u_1, v_2) \\ A_{u_2}^t & q(u_2, v_1) & q(u_2, v_2) \end{pmatrix} \\ = q(u_1, u_2) \det \begin{pmatrix} q(u_1, v_1) & q(u_1, v_2) \\ q(u_2, v_1) & q(u_2, v_2) \end{pmatrix}.$$

We can derive $A_{u_2}^t$ from $A_{u_1}^t$ by solving the above equation.

We are currently investigating how to make use of all our results in the general case.

The inversion algorithms can be adapted to estimate parameters

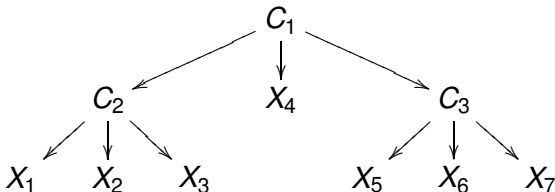
- Basic idea: apply the inversion algorithm to the observed distribution \hat{p} .
- Instead of testing whether a candidate parameter is mapped to p , we find the parameter minimizing the relative entropy to \hat{p} .
- Suppose that the unknown p is a naive Bayesian distribution with m classes satisfying our inversion assumption. As the sample size increases, \hat{p} converges to p and, by continuity, our estimate converges to a true parameter mapped to p .

Practical issues

The estimation procedure has several issues:

- The computational complexity grows extremely fast with m , but linearly with n .
- The estimates are numerically unstable and require large sample sizes. For smaller sample sizes, there may not even be a single candidate parameter satisfying the parameter constraints.
- There are many degrees of freedom in the choice of \mathbf{t} , \mathbf{u} and \mathbf{v} . Asymptotically, any choice is suitable. For small sample size, it is probably important.
- The results are not competitive with the E.M. algorithm.

Extension to hierarchical latent class models



The parameters mapped to a HLC distribution with the above structure can be derived from the parameters mapped to the naive Bayesian distributions over

$$\{X_1, X_2, X_3\}$$

$$\{X_1, X_4, X_5\}$$

$$\{X_5, X_6, X_7\}$$

obtained by marginalization.

Conclusion

We presented some simple and interesting polynomial equations constraining a naive Bayesian distribution and its parameters. These results may be applied to

- compute the parameters mapped to a naive Bayesian distribution,
- estimate parameters from data.

The implicit equation

$$\det \begin{pmatrix} q(u_1, v_1) & \dots & q(u_1, v_{m'-1}) \\ \vdots & & \vdots \\ q(u_{m'-1}, v_1) & \dots & q(u_{m'-1}, v_{m'-1}) \end{pmatrix} = 0$$

holding for $m' > m$ is similar to a tetrad constraint. A future research direction would investigate whether the constraint can indeed be used to learn hidden causes from data.