

Building Bayes Networks: Parameter Learning

Learning Naive Bayes Classifier

Given: A database of samples from domain of interest.

The graph underlying a graphical model for the domain.

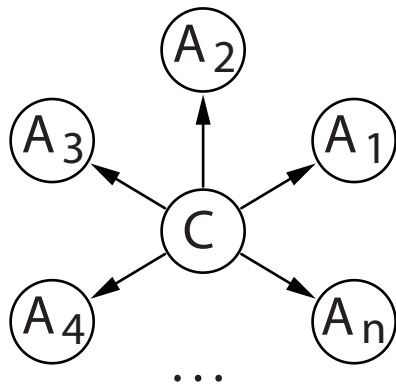
Desired: Good values for the numeric parameters of the model.

Example: Naive Bayes Classifiers

A naive Bayes classifier is a Bayesian network with star-like structure.

The class attribute is the only unconditional attribute.

All other attributes are conditioned on the class only



The structure of a naive Bayes classifier is fixed once the attributes have been selected. The only remaining task is to estimate the parameters of the needed probability distributions.

Probabilistic Classification

A classifier is an algorithm that assigns a class from a predefined set to a case or object, based on the values of descriptive attributes.

An optimal classifier maximizes the probability of a correct class assignment.

- Let C be a class attribute with $\text{dom}(C) = \{c_1, \dots, c_{n_C}\}$, which occur with probabilities p_i , $1 \leq i \leq n_C$.
- Let q_i be the probability with which a classifier assigns class c_i . ($q_i \in \{0, 1\}$ for a deterministic classifier)
- The probability of a correct assignment is

$$P(\text{correct assignment}) = \sum_{i=1}^{n_C} p_i q_i.$$

- Therefore the best choice for the q_i is

$$q_i = \begin{cases} 1, & \text{if } p_i = \max_{k=1}^{n_C} p_k, \\ 0, & \text{otherwise.} \end{cases}$$

Probabilistic Classification

Consequence: An optimal classifier should assign the **most probable class**.

This argument does not change if we take descriptive attributes into account.

- Let $U = \{A_1, \dots, A_m\}$ be a set of descriptive attributes with domains $\text{dom}(A_k)$, $1 \leq k \leq m$.
- Let $A_1 = a_1, \dots, A_m = a_m$ be an instantiation of the descriptive attributes.
- An optimal classifier should assign the class c_i for which

$$P(C = c_i \mid A_1 = a_1, \dots, A_m = a_m) = \max_{j=1}^{n_C} P(C = c_j \mid A_1 = a_1, \dots, A_m = a_m)$$

Problem: We cannot store a class (or the class probabilities) for every possible instantiation $A_1 = a_1, \dots, A_m = a_m$ of the descriptive attributes. (The table size grows exponentially with the number of attributes.)

Therefore: **Simplifying assumptions are necessary.**

Bayes' Rule and Bayes' Classifiers

Bayes' classifiers: Compute the class probabilities as

$$P(C = c_i | A_1 = a_1, \dots, A_m = a_m) = \frac{P(A_1 = a_1, \dots, A_m = a_m | C = c_i) \cdot P(C = c_i)}{P(A_1 = a_1, \dots, A_m = a_m)}.$$

Looks unreasonable at first sight: Even more probabilities to store.

Naive Bayes Classifiers

Naive Assumption:

The descriptive attributes are conditionally independent given the class.

Bayes' Rule:

$$P(C = c_i | \omega) = \frac{P(A_1 = a_1, \dots, A_m = a_m | C = c_i) \cdot P(C = c_i)}{P(A_1 = a_1, \dots, A_m = a_m)} \quad \leftarrow p_0$$

abbrev. for the
normalizing constant

Chain Rule of Probability:

$$P(C = c_i | \omega) = \frac{P(C = c_i)}{p_0} \cdot \prod_{k=1}^m P(A_k = a_k | A_1 = a_1, \dots, A_{k-1} = a_{k-1}, C = c_i)$$

Conditional Independence Assumption:

$$P(C = c_i | \omega) = \frac{P(C = c_i)}{p_0} \cdot \prod_{k=1}^m P(A_k = a_k | C = c_i)$$

Naive Bayes Classifiers (continued)

Consequence: Manageable amount of data to store.

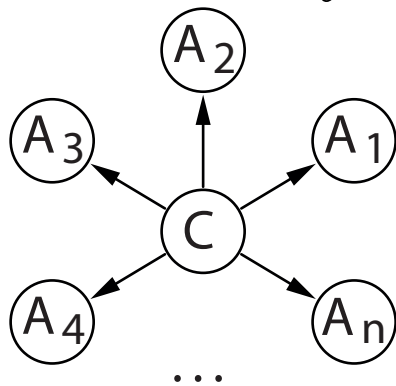
Store distributions $P(C = c_i)$ and $\forall 1 \leq k \leq m : P(A_k = a_k | C = c_i)$.

Classification: Compute for all classes c_i

$$P(C = c_i | A_1 = a_1, \dots, A_m = a_m) \cdot p_0 = P(C = c_i) \cdot \prod_{j=1}^n P(A_j = a_j | C = c_i)$$

and predict the class c_i for which this value is largest.

Relation to Bayesian Networks:



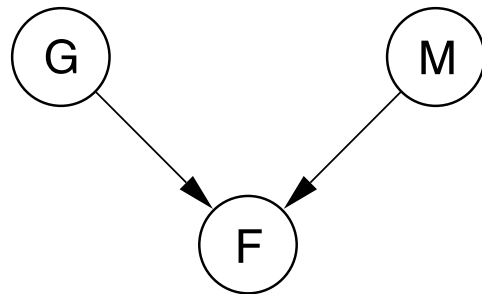
Decomposition formula:

$$\begin{aligned} &P(C = c_i, A_1 = a_1, \dots, A_n = a_n) \\ &= P(C = c_i) \cdot \prod_{j=1}^n P(A_j = a_j | C = c_i) \end{aligned}$$

Learning the parameters of a Graphical Model

$A_1 = G$	$Q_{11} = \phi$
$a_{11} = g$	
$a_{12} = \bar{g}$	

$A_2 = M$	$Q_{21} = \phi$
$a_{12} = m$	
$a_{22} = \bar{m}$	



$A_3 = F$	$Q_{31} = (g, m)$	$Q_{32} = (g, \bar{m})$	$Q_{33} = (\bar{g}, m)$	$Q_{34} = (\bar{g}, \bar{m})$
$a_{31} = f$				
$a_{32} = \bar{f}$				

$$V = \{G, M, F\}$$

$$\text{dom}(G) = \{g, \bar{g}\}$$

$$\text{dom}(M) = \{m, \bar{m}\}$$

$$\text{dom}(F) = \{f, \bar{f}\}$$

The potential tables' layout is determined by the graph structure.

The parameters (i. e. the table entries) can be easily estimated from the database, e. g.:

$$\hat{P}(f \mid g, m) = \frac{\hat{P}(f, g, m)}{\hat{P}(g, m)} = \frac{\frac{\#(g, m, f)}{|D|}}{\frac{\#(g, m)}{|D|}} = \frac{\#(g, m, f)}{\#(g, m)}$$

Likelihood of a Database

Flu G	\bar{g}	\bar{g}	\bar{g}	\bar{g}	g	g	g	g
Malaria M	\bar{m}	\bar{m}	m	m	\bar{m}	\bar{m}	m	m
Fever F	\bar{f}	f	\bar{f}	f	\bar{f}	f	\bar{f}	f
#	34	6	2	8	16	24	0	10

Database D with 100 entries for 3 attributes.

As the structure given by the graph of the previous slide suggests, the probability of $P(g, m, f)$ can be computed by:

$$P(\mathbf{g}, \mathbf{m}, \mathbf{f}) = P(\mathbf{g})P(\mathbf{m})P(\mathbf{f} \mid \mathbf{g}, \mathbf{m})$$

Necessary conditional probabilities can be calculated using the Bayes-Theorem:

$$\hat{P}(\mathbf{f} \mid \mathbf{g}, \mathbf{m}) = \frac{\hat{P}(\mathbf{f}, \mathbf{g}, \mathbf{m})}{\hat{P}(\mathbf{g}, \mathbf{m})} = \frac{\frac{\#(\mathbf{g}, \mathbf{m}, \mathbf{f})}{|D|}}{\frac{\#(\mathbf{g}, \mathbf{m})}{|D|}} = \frac{\#(\mathbf{g}, \mathbf{m}, \mathbf{f})}{\#(\mathbf{g}, \mathbf{m})} = \frac{10}{10} = 1.00$$

$$\hat{P}(\mathbf{f} \mid \bar{\mathbf{g}}, \bar{\mathbf{m}}) = \frac{\hat{P}(\mathbf{f}, \bar{\mathbf{g}}, \bar{\mathbf{m}})}{\hat{P}(\bar{\mathbf{g}}, \bar{\mathbf{m}})} = \frac{\frac{\#(\bar{\mathbf{g}}, \bar{\mathbf{m}}, \mathbf{f})}{|D|}}{\frac{\#(\bar{\mathbf{g}}, \bar{\mathbf{m}})}{|D|}} = \frac{\#(\bar{\mathbf{g}}, \bar{\mathbf{m}}, \mathbf{f})}{\#(\bar{\mathbf{g}}, \bar{\mathbf{m}})} = \frac{6}{40} = 0.15$$

Likelihood of a Database (2)

The likelihood of the calculated probabilities $P(D \mid B_S, B_P)$ can be computed under presence of three assumptions:

1. The data generation process can be described exactly by a bayesian network (B_S, B_P)
2. The single tuples of the dataset are independent of each other.
3. All tuples are complete, therefore no missing values hinder the probability inference

The first assumption legitimates the search of an appropriate bayesian network.

The second assumption is required for an unbiased observation of dataset tuples.

Assumption three ensures the inference of B_P using D and B_S as shown on the previous slides.

Likelihood of a Database (3)

Flu G	\bar{g}	\bar{g}	\bar{g}	\bar{g}	g	g	g	g
Malaria M	\bar{m}	\bar{m}	m	m	\bar{m}	\bar{m}	m	m
Fever F	\bar{f}	f	\bar{f}	f	\bar{f}	f	\bar{f}	f
#	34	6	2	8	16	24	0	10

Database D with 100 entries for 3 attributes.

$$P(D \mid B_S, B_P) = \prod_{h=1}^{100} P(c_h \mid B_S, B_P)$$

$$\begin{aligned}
 &= \underbrace{P(g, m, f) \cdots P(g, m, f)}_{\substack{\text{Case 1} \\ \text{Case 10} \\ \text{10 times}}} \cdots \underbrace{P(\bar{g}, m, f) \cdots P(\bar{g}, m, f)}_{\substack{\text{Case 51} \\ \text{Case 58} \\ \text{8 times}}} \cdots \underbrace{P(\bar{g}, \bar{m}, \bar{f}) \cdots P(\bar{g}, \bar{m}, \bar{f})}_{\substack{\text{Case 67} \\ \text{Case 100} \\ \text{34 times}}} \\
 &= \underbrace{P(g, m, f)^{10}} \cdots \underbrace{P(\bar{g}, m, f)^8} \cdots \underbrace{P(\bar{g}, \bar{m}, \bar{f})^{34}} \\
 &= \underbrace{P(f \mid g, m)^{10} P(g)^{10} P(m)^{10}} \cdots \underbrace{P(f \mid \bar{g}, m)^8 P(\bar{g})^8 P(m)^8} \cdots \underbrace{P(\bar{f} \mid \bar{g}, \bar{m})^{34} P(\bar{g})^{34} P(\bar{m})^{34}}
 \end{aligned}$$

Likelihood of a Database (4)

$$\begin{aligned} P(D \mid B_S, B_P) &= \prod_{h=1}^{100} P(c_h \mid B_S, B_P) \\ &= P(\mathbf{f} \mid \mathbf{g}, \mathbf{m})^{10} P(\bar{\mathbf{f}} \mid \mathbf{g}, \mathbf{m})^0 P(\mathbf{f} \mid \mathbf{g}, \bar{\mathbf{m}})^{24} P(\bar{\mathbf{f}} \mid \mathbf{g}, \bar{\mathbf{m}})^{16} \\ &\quad \cdot P(\mathbf{f} \mid \bar{\mathbf{g}}, \mathbf{m})^8 P(\bar{\mathbf{f}} \mid \bar{\mathbf{g}}, \mathbf{m})^2 P(\mathbf{f} \mid \bar{\mathbf{g}}, \bar{\mathbf{m}})^6 P(\bar{\mathbf{f}} \mid \bar{\mathbf{g}}, \bar{\mathbf{m}})^{34} \\ &\quad \cdot P(\mathbf{g})^{50} P(\bar{\mathbf{g}})^{50} P(\mathbf{m})^{20} P(\bar{\mathbf{m}})^{80} \end{aligned}$$

The last equation shows the principle of reordering the factors:

First, we sort by attributes (here: **F**, **G** then **M**).

Within the same attributes, factors are grouped by the parent attributes' values combinations (here: for **F**: (\mathbf{g}, \mathbf{m}) , $(\mathbf{g}, \bar{\mathbf{m}})$, $(\bar{\mathbf{g}}, \mathbf{m})$ and $(\bar{\mathbf{g}}, \bar{\mathbf{m}})$).

Finally, it is sorted by attribute values (here: for **F**: first **f**, then $\bar{\mathbf{f}}$).

Likelihood of a Database (5)

General likelihood of a database D given a known bayesian network structure B_S and the parameters B_P :

$$P(D \mid B_S, B_P) = \prod_{i=1}^n \prod_{j=1}^{q_i} \prod_{k=1}^{r_i} \theta_{ijk}^{\alpha_{ijk}}$$

General potential table:

A_i	Q_{i1}	\cdots	Q_{ij}	\cdots	Q_{iq_i}
a_{i1}	θ_{i11}	\cdots	θ_{ij1}	\cdots	θ_{iq_i1}
\vdots	\vdots	\ddots	\vdots	\ddots	\vdots
a_{ik}	θ_{i1k}	\cdots	θ_{ijk}	\cdots	θ_{iq_ik}
\vdots	\vdots	\ddots	\vdots	\ddots	\vdots
a_{ir_i}	θ_{i1r_i}	\cdots	θ_{ijr_i}	\cdots	$\theta_{iq_ir_i}$

$$P(A_i = a_{ik} \mid \text{parents}(A_i) = Q_{ij}) = \theta_{ijk}$$

$$\sum_{k=1}^{r_i} \theta_{ijk} = 1$$