

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.

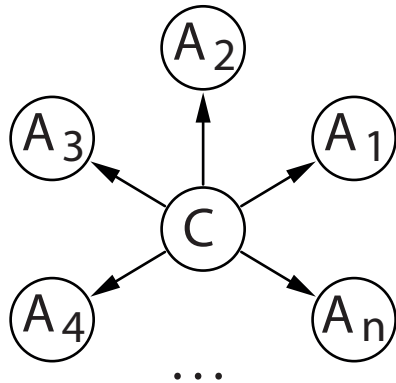
Desired: Good values for the numeric parameters of the model.

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.

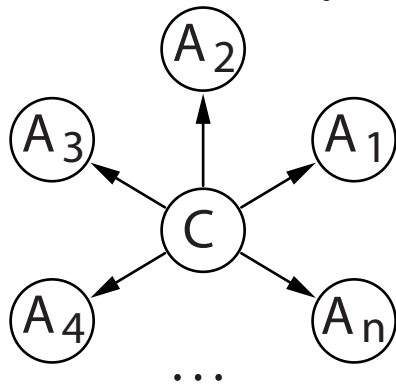
Store distributions $P(C = c_i)$ and $\forall 1 \leq k \leq m : P(A_k = a_k | C = c_i)$.

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:

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

Estimation of Probabilities:

Nominal/Categorical Attributes

$$\hat{P}(A_k = a_k \mid C = c_i) = \frac{\#(A_k = a_k, C = c_i) + \gamma}{\#(C = c_i) + n_{A_k} \gamma}$$

$\#(\varphi)$ is the number of example cases that satisfy the condition φ

n_{A_j} is the number of values of the attribute A_j .

γ is called **Laplace correction**

$\gamma = 0$: Maximum likelihood estimation.

Common choices: $\gamma = 1$ or $\gamma = \frac{1}{2}$.

Laplace correction help to avoid problems with attribute values that do not occur with some class in the given data.

It also introduces a bias towards a uniform distribution.

Naive Bayes Classifiers: Parameter Estimation

Estimation of Probabilities:

Metric/Numeric Attributes: Assume a normal distribution.

$$P(A_k = a_k \mid C = c_i) = \frac{1}{\sqrt{2\pi}\sigma_k(c_i)} \exp\left(-\frac{(a_k - \mu_k(c_i))^2}{2\sigma_k^2(c_i)}\right)$$

Estimate of mean value

$$\hat{\mu}_k(c_i) = \frac{1}{\#(C = c_i)} \sum_{j=1}^{\#(C=c_i)} a_k(j)$$

Estimate of variance

$$\hat{\sigma}_k^2(c_i) = \frac{1}{\xi} \sum_{j=1}^{\#(C=c_i)} (a_k(j) - \hat{\mu}_k(c_i))^2$$

$\xi = \#(C = c_i)$: Maximum likelihood estimation

$\xi = \#(C = c_i) - 1$: Unbiased estimation

Naive Bayes Classifiers: Simple Example 1

| No | Sex | Age | Blood pr. | Drug |
|----|--------|-----|-----------|------|
| 1 | male | 20 | normal | A |
| 2 | female | 73 | normal | B |
| 3 | female | 37 | high | A |
| 4 | male | 33 | low | B |
| 5 | female | 48 | high | A |
| 6 | male | 29 | normal | A |
| 7 | female | 52 | normal | B |
| 8 | male | 42 | low | B |
| 9 | male | 61 | normal | B |
| 10 | female | 30 | normal | A |
| 11 | female | 26 | low | B |
| 12 | male | 54 | high | A |

| $P(\text{Drug})$ | A | B |
|--|-------|-------|
| | 0.5 | 0.5 |
| $P(\text{Sex} \mid \text{Drug})$ | A | B |
| male | 0.5 | 0.5 |
| female | 0.5 | 0.5 |
| $P(\text{Age} \mid \text{Drug})$ | A | B |
| μ | 36.3 | 47.8 |
| σ^2 | 161.9 | 311.0 |
| $P(\text{Blood Pr.} \mid \text{Drug})$ | A | B |
| low | 0 | 0.5 |
| normal | 0.5 | 0.5 |
| high | 0.5 | 0 |

A simple database and estimated (conditional) probability distributions.

Naive Bayes Classifiers: Simple Example 1

$$P(\text{Drug A} \mid \text{male}, 61, \text{normal})$$

$$\begin{aligned} &= c_1 \cdot P(\text{Drug A}) \cdot P(\text{male} \mid \text{Drug A}) \cdot P(61 \mid \text{Drug A}) \cdot P(\text{normal} \mid \text{Drug A}) \\ &\approx c_1 \cdot 0.5 \cdot 0.5 \cdot 0.004787 \cdot 0.5 = c_1 \cdot 5.984 \cdot 10^{-4} = 0.219 \end{aligned}$$

$$P(\text{Drug B} \mid \text{male}, 61, \text{normal})$$

$$\begin{aligned} &= c_1 \cdot P(\text{Drug B}) \cdot P(\text{male} \mid \text{Drug B}) \cdot P(61 \mid \text{Drug B}) \cdot P(\text{normal} \mid \text{Drug B}) \\ &\approx c_1 \cdot 0.5 \cdot 0.5 \cdot 0.017120 \cdot 0.5 = c_1 \cdot 2.140 \cdot 10^{-3} = 0.781 \end{aligned}$$

$$P(\text{Drug A} \mid \text{female}, 30, \text{normal})$$

$$\begin{aligned} &= c_2 \cdot P(\text{Drug A}) \cdot P(\text{female} \mid \text{Drug A}) \cdot P(30 \mid \text{Drug A}) \cdot P(\text{normal} \mid \text{Drug A}) \\ &\approx c_2 \cdot 0.5 \cdot 0.5 \cdot 0.027703 \cdot 0.5 = c_2 \cdot 3.471 \cdot 10^{-3} = 0.671 \end{aligned}$$

$$P(\text{Drug B} \mid \text{female}, 30, \text{normal})$$

$$\begin{aligned} &= c_2 \cdot P(\text{Drug B}) \cdot P(\text{female} \mid \text{Drug B}) \cdot P(30 \mid \text{Drug B}) \cdot P(\text{normal} \mid \text{Drug B}) \\ &\approx c_2 \cdot 0.5 \cdot 0.5 \cdot 0.013567 \cdot 0.5 = c_2 \cdot 1.696 \cdot 10^{-3} = 0.329 \end{aligned}$$

Naive Bayes Classifiers: Simple Example 2

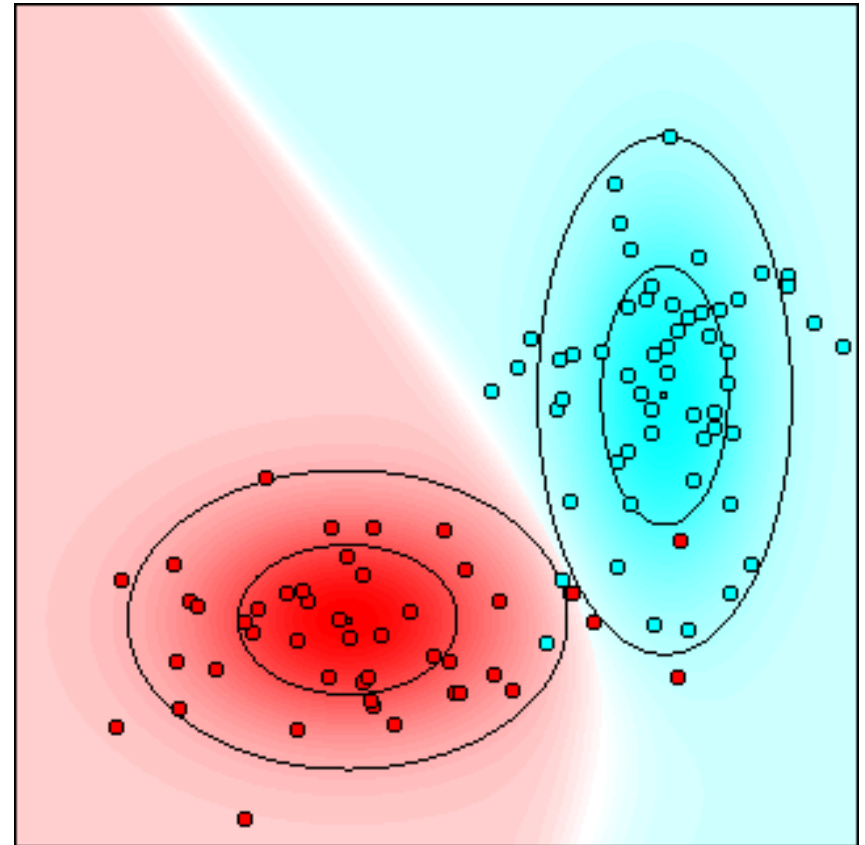
100 data points, 2 classes

Small squares: mean values

Inner ellipses:
one standard deviation

Outer ellipses:
two standard deviations

Classes overlap:
classification is not perfect



Naive Bayes Classifier

Naive Bayes Classifiers: Simple Example 3

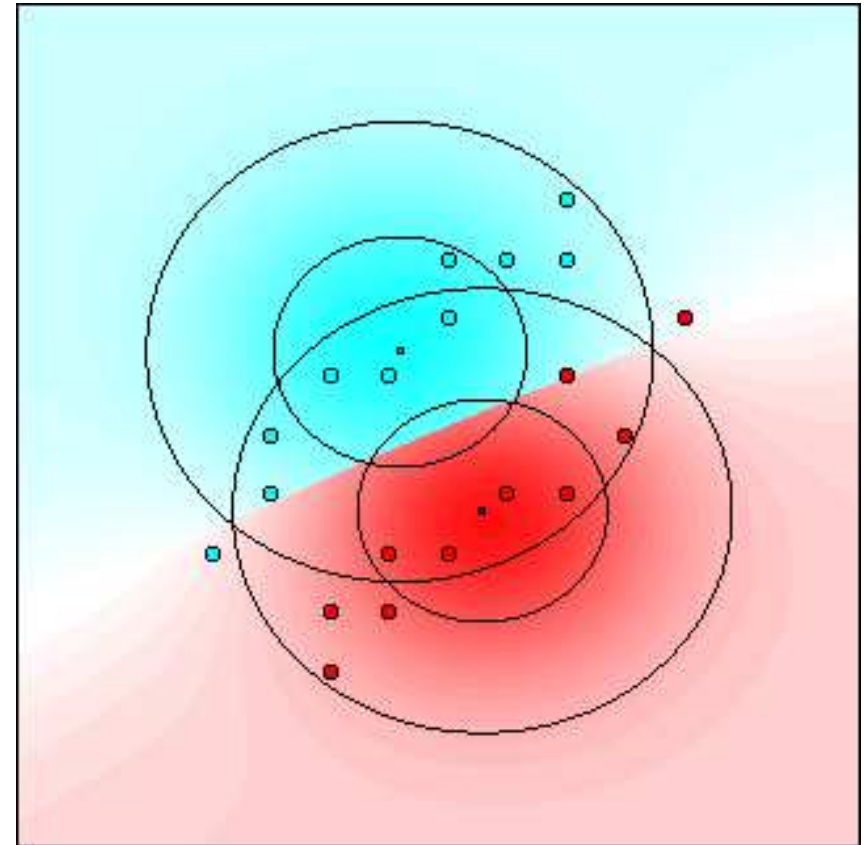
20 data points, 2 classes

Small squares: mean values

Inner ellipses:
one standard deviation

Outer ellipses:
two standard deviations

Attributes are not conditionally
independent given the class



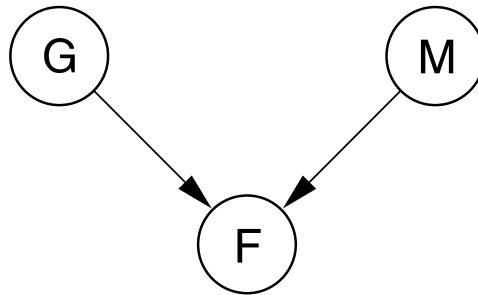
Naive Bayes Classifier

Learning the parameters of a Graphical Model

Probability values can be estimated using methods of inductive statistics.

| | P(G) |
|--------------------|------|
| $a_{11} = g$ | |
| $a_{12} = \bar{g}$ | |

| | P(M) |
|--------------------|------|
| $a_{12} = m$ | |
| $a_{22} = \bar{m}$ | |



| P(F G,M) | g, m | g, \bar{m} | \bar{g} , m | \bar{g} , \bar{m} |
|--------------------|------|--------------|---------------|-----------------------|
| $a_{31} = f$ | | | | |
| $a_{32} = \bar{f}$ | | | | |

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

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

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

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

Learning the parameters of a Graphical Model

| | | | | | | | | |
|------------------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----|
| 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:

$$P(g, m, f) = P(g)P(m)P(f | g, m)$$

Estimates for these probabilities can be calculated, e.g. using the database

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

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

Likelihood of a Database

Let B_P be the description of the parameters, B_S be the given structure and D the data.

The likelihood of the calculated probabilities $P(D | B_S, B_P)$ 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

| | | | | | | | | |
|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----|
| 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 c_h 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

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

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 $\bar{\mathbf{f}}$).

Bayes Theorem gives the likelihood $P(B_P \mid D, B_S)$.

Maximum likelihood approach gives a good estimate for \hat{B}_P .

Likelihood of a Database

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 | $\theta_{iq_ir_i}$ |

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

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