

## 2. Bayesian classifiers

- Probabilistic classification
  - Probabilistic classification problems
  - Naive Bayesian classifiers
  - Augmented Bayesian classifiers
  - Semi-naive Bayesian classifiers
  
- Multi-label classification
  - Multi-dimensional classification problems
  - Basic approaches
  - Chain classifiers

# Outline

- Probabilistic classification
  - Probabilistic classification problems
  - Naive Bayesian classifiers
  - Augmented Bayesian classifiers
  - Semi-naive Bayesian classifiers
  
- Multi-label classification
  - Multi-dimensional classification problems
  - Basic approaches
  - Chain classifiers

## Probabilistic classification problems

given a set of samples  $X$  and a set of class labels  $Y$  ( $\text{dom } Y \subseteq \mathbf{Z}_+$ )

- 'ordinary' classifier:  $f: X \rightarrow Y$
- probabilistic classifier:

$$f(x) = (\dots, \mathbf{P}(y_i | x), \dots), \quad i = 1, \dots, |Y|$$

$$- \mathbf{1}^T f(x) = 1$$

$$- \hat{y} = \operatorname{argmax}_y \mathbf{P}(y | x)$$

## Probabilistic classification problems

### Bayesian approach

$$\mathbf{P}(y | x) = \frac{\mathbf{P}(x | y) \mathbf{P}(y)}{\mathbf{P}(x)}$$

- $\mathbf{P}(x)$ : normalizing constant independent of labels
- $\mathbf{P}(y)$ : prior on class labels
- $\mathbf{P}(x | y)$ : likelihood of sample  $x$  under label  $y$

$$\begin{aligned}\mathbf{P}(x | y) &= \mathbf{P}(x_1, \dots, x_n | y) \\ &= \mathbf{P}(x_1 | y) \mathbf{P}(x_2 | x_1, y) \cdots \mathbf{P}(x_n | x_{n-1}, \dots, x_1, y)\end{aligned}$$

– can be difficult to calculate

## Naive Bayesian classifiers

**assumption:**  $x_1, \dots, x_n$  are independent given  $y$

$$\mathbf{P}(y | x) = \frac{\mathbf{P}(x | y) \mathbf{P}(y)}{\mathbf{P}(x)} \propto \mathbf{P}(y) \mathbf{P}(x_1, \dots, x_n | y) = \mathbf{P}(y) \prod_{i=1}^n \mathbf{P}(x_i | y)$$

### parameter learning

- prior  $\mathbf{P}(y)$ :

$$\mathbf{P}(y_i) = \frac{1}{|Y|} \quad \text{or} \quad \mathbf{P}(y_i) = \frac{\# \text{ samples in class } y_i}{\# \text{ samples in total}}$$

- likelihood  $\mathbf{P}(x | y)$ :

$$\mathbf{P}(x_k | y_i) = \frac{\# \text{ samples in class } y_i \text{ with feature } x_k}{\# \text{ samples in class } y_i}$$

for all  $k = 1, \dots, n$ ,  $y_i \in Y$

## Naive Bayesian classifiers

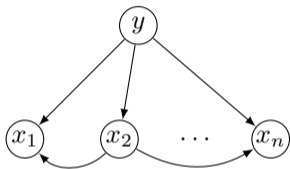
**handling continuous features:** Gaussian naive Bayes

$$p(x_k | y_i) = \frac{1}{\sqrt{2\pi}\sigma_{k|y_i}} \exp\left(-\frac{(x_k - \mu_{k|y_i})^2}{2\sigma_{k|y_i}^2}\right), \quad k = 1, \dots, n$$

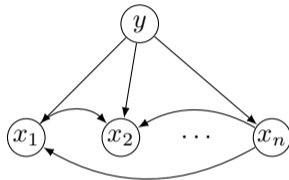
- $\mu_{k|y_i} = \mathbf{E}(X_k | y_i)$ ,  $k = 1, \dots, n$
- $\sigma_{k|y_i} = \sqrt{\mathbf{var}(X_k | y_i)} = \sqrt{\mathbf{E}\left(\left(X_k - \mathbf{E}(X_k | y_i)\right)^2 \mid y_i\right)}$ ,  $k = 1, \dots, n$
- $\mathbf{P}(x_k | y_i) \propto p(x_k | y_i)$

## Augmented Bayesian classifiers

**assumption:** some dependency structure (tree, DAG, ...) exists between  $x_1, \dots, x_n$  given  $y$



tree augmented Bayesian classifiers



Bayesian network augmented Bayesian classifiers

**parameter learning**

$$\mathbf{P}(x | y) = \mathbf{P}(x_1, \dots, x_n | y) = \prod_{i=1}^n \mathbf{P}(x_i | \mathbf{pa}(x_i), y)$$

## Semi-naive Bayesian classifiers

**basic idea:** naive Bayes + feature selection

- eliminate or join interdependent features given the class label

**feature selection metrics**

- local measure: e.g., mutual information
- global measure: e.g., performance of the classifier with and without the feature

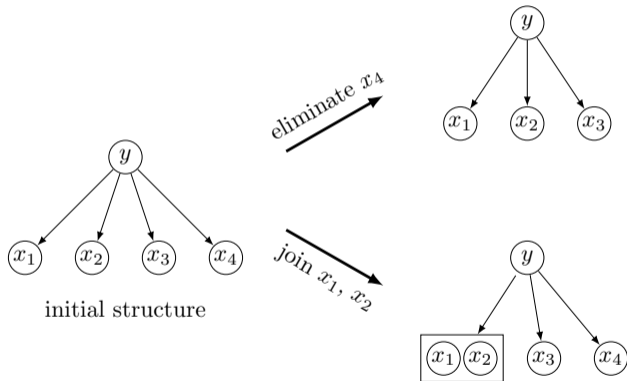
**model structure learning process**

- bottom-up: start from an empty structure and add features
- top-down: from a full structure with all the features and eliminate (or combine) features

**parameter learning:** the same as naive Bayesian classifiers

## Semi-naive Bayesian classifiers

**example:** top-down structure learning



# Outline

- Probabilistic classification
  - Probabilistic classification problems
  - Naive Bayesian classifiers
  - Augmented Bayesian classifiers
  - Semi-naive Bayesian classifiers
- Multi-label classification
  - Multi-dimensional classification problems
  - Basic approaches
  - Chain classifiers

## Multi-dimensional classification problems

given a set of samples  $X$  and a set of class labels  $Y$  ( $\text{dom } Y \subseteq \mathbf{Z}_+^m$ )

probabilistic classifier  $f$ :

$$f(x) = (\dots, \mathbf{P}(y \mid x), \dots) = (\dots, \mathbf{P}(y_1, \dots, y_m \mid x), \dots)$$

- $y \in Y$  is a  $m$ -dimensional vector
- $\mathbf{1}^T f(x) = 1$
- $\hat{y} = \operatorname{argmax}_y \mathbf{P}(y \mid x)$

**multi-label classification:**  $\text{dom } Y_i = \{0, 1\}, i = 1, \dots, m$

## Basic approaches

### binary relevance

- assumption: no dependencies between all pairs of classes
- solve  $m$  independent binary classification problems
- a classifier is independently learnt for each class  $Y_1, \dots, Y_m$
- final prediction is a simple concatenation of results from all classifier,  $\hat{y} = (\hat{y}_1, \dots, \hat{y}_m)$

### label power-set

- basic idea: transform multi-label classification to single-class scenario
- define a mapping  $g: Y \rightarrow Y'$  from  $\mathbf{dom} Y \subseteq \mathbf{Z}_+^m$  to  $\mathbf{dom} Y' \subseteq \mathbf{Z}_+$
- learn a single-class classifier on  $Y'$  given  $X$
- interactions between classes are implicitly considered
- $|Y'|$  increases exponentially w.r.t.  $m$

## Chain classifiers

**basic idea:** generalize the binary relevance approach to considering some dependencies between classes

- $m$  binary classifiers ( $f_1, \dots, f_m$ ) linked in a chain, each corresponding to one class
- the predictions  $\hat{y}_1, \dots, \hat{y}_{i-1}$  from  $f_1, \dots, f_{i-1}$  is incorporated into the features of  $f_i$

$$\hat{y}_1 = \operatorname{argmax}_{y_1} \mathbf{P}(y_1 | x)$$
$$\hat{y}_i = \operatorname{argmax}_{y_i} \mathbf{P}(y_i | x, \hat{y}_1, \dots, \hat{y}_{i-1}), \quad i = 2, \dots, m$$

- model performance depends on the order of classes in the chain

## Chain classifiers

### circular chain classifier

- $(f_1, \dots, f_m)$  are connected in a circular way
- the first cycle:

$$\hat{y}_1 = \operatorname{argmax}_{y_1} \mathbf{P}(y_1 | x)$$
$$\hat{y}_i = \operatorname{argmax}_{y_i} \mathbf{P}(y_i | x, \hat{y}_1, \dots, \hat{y}_{i-1}), \quad i = 2, \dots, m$$

- from the second cycle:

$$\hat{y}_i = \operatorname{argmax}_{y_i} \mathbf{P}(y_i | x, \hat{y}_{-i}), \quad i = 1, \dots, m$$

each binary classifier in the chain receives the predictions of all other classifiers as additional feature

- repeated for a prefixed number of cycles or until convergence

## Chain classifiers

### Bayesian chain classifier

- connection between  $(f_1, \dots, f_m)$  represented as a DAG

$$\mathbf{P}(y \mid x) = \mathbf{P}(y_1, \dots, y_m \mid x) = \prod_{i=1}^m \mathbf{P}(y_i \mid \mathbf{pa}(y_i), x)$$

- to get final prediction  $\hat{y}$ , approximate the hard combinatorial optimization problem

$$\text{maximize (over } y) \quad \prod_{i=1}^m \mathbf{P}(y_i \mid \mathbf{pa}(y_i), x)$$

with a sequence of independent optimization problems

$$\text{maximize (over } y_i) \quad \mathbf{P}(y_i \mid \mathbf{pa}(y_i), x)$$

for all  $i = 1, \dots, m$

## Chain classifiers

example

