

# Introduction to Machine Learning

## NPFL 054

<http://ufal.mff.cuni.cz/course/npfl054>

Barbora Hladká  
hladka@ufal.mff.cuni.cz

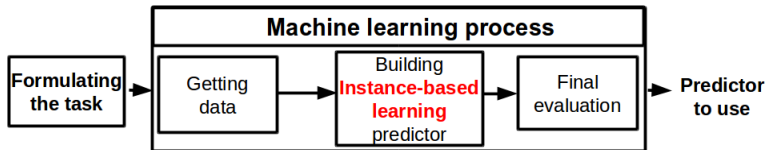
Martin Holub  
holub@ufal.mff.cuni.cz

Charles University,  
Faculty of Mathematics and Physics,  
Institute of Formal and Applied Linguistics

# Outline

- Instance-based learning
- Naïve Bayes algorithm
- Bayesian networks

# Instance-based learning



# Instance-based learning

## Key idea

- IBL methods initially store the training data (that is why IBL methods are often referred to as "**lazy**" methods).
- For a new instance, prediction is based on local similarity, i.e. a set of similar instances are retrieved and used for prediction
- IBL methods can construct a different approximation of a target function for each distinct test instance.
- Both classification and regression.

# Instance-based learning

## Key points

- ① A distance metric
- ② How many nearby neighbours look at?
- ③ A weighting function
- ④ How to fit with local points?

# Instance-based learning

## Distance metric

The most common ones

- **Euclidean distance**

$$E(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{\sum_{r=1}^m (x_{i_r} - x_{j_r})^2} \quad (1)$$

- **Manhattan distance**

$$M(\mathbf{x}_i, \mathbf{x}_j) = \sum_{r=1}^m |x_{i_r} - x_{j_r}| \quad (2)$$

# Instance-based learning

## Learning algorithms

- k-Nearest Neighbour
- Distance weighted k-NN
- Locally weighted linear regression
- ...

# Instance-based learning

## *k*-Nearest Neighbour algorithm

- 1 **A distance metric:** Euclidian (most widely used)
  - 2 **How many nearby neighbours look at?**  $k$
  - 3 **A weighting function:** unused
  - 4 **How to fit with local points?**
- **k-NN classification**

$$h(\mathbf{x}) = \operatorname{argmax}_{v \in Y} \sum_{i=1}^k \delta(v, y_i), \quad (3)$$

where  $\delta(a, b) = 1$  if  $a = b$ , otherwise 0

- **k-NN regression**

$$h(\mathbf{x}) = \frac{\sum_{i=1}^k y_i}{k} \quad (4)$$



# Instance-based learning

## Distance-weighted $k$ -NN algorithm

- 1 **A distance metric:** Euclidian (most widely used)
- 2 **How many nearby neighbours look at?**  $k$
- 3 **A weighting function:** greater weight closer neighbours

$$w_i(\mathbf{x}) \equiv \frac{1}{d(\mathbf{x}, \mathbf{x}_i)^2}$$

### 4 How to fit with local points?

- **Classification**

$$h(\mathbf{x}) = \operatorname{argmax}_{v \in Y} \sum_{i=1}^k w_i(\mathbf{x}) \delta(v, y_i) \quad (5)$$

- **Regression**

$$h(\mathbf{x}) = \frac{\sum_{i=1}^k w_i(\mathbf{x}) y_i}{\sum_{i=1}^k w_i(\mathbf{x})} \quad (6)$$

# Instance-based learning

## Distance-weighted $k$ -NN algorithm

### Shepard's method

- Classification

$$h(\mathbf{x}) = \operatorname{argmax}_{v \in Y} \sum_{i=1}^n w_i(\mathbf{x}) \delta(v, y_i) \quad (7)$$

- Regression

$$h(\mathbf{x}) = \frac{\sum_{i=1}^n w_i(\mathbf{x}) y_i}{\sum_{i=1}^n w_i(\mathbf{x})} \quad (8)$$

# Instance-based learning

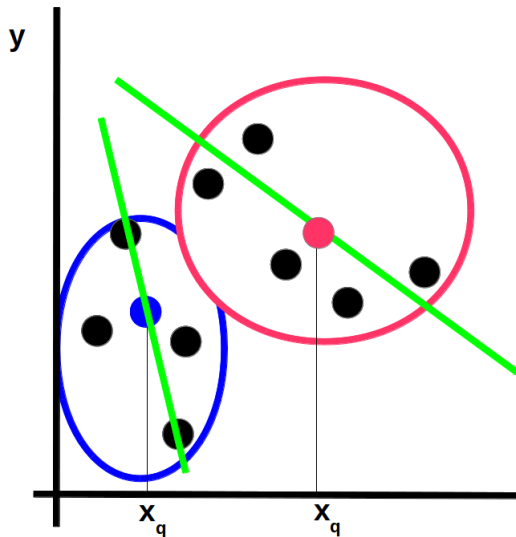
## Locally weighted linear regression

- 1 **A distance metric:** Euclidian (most widely used)
- 2 **How many nearby neighbours look at?**  $k$
- 3 **A weighting function:**  $w_i(\mathbf{x})$
- 4 **How to fit with local points?**

$$\Theta^* = \operatorname{argmin}_{\Theta} \sum_{i=1}^k w_i(\mathbf{x}) (\Theta^T \mathbf{x}_i - y_i)^2 \quad (9)$$

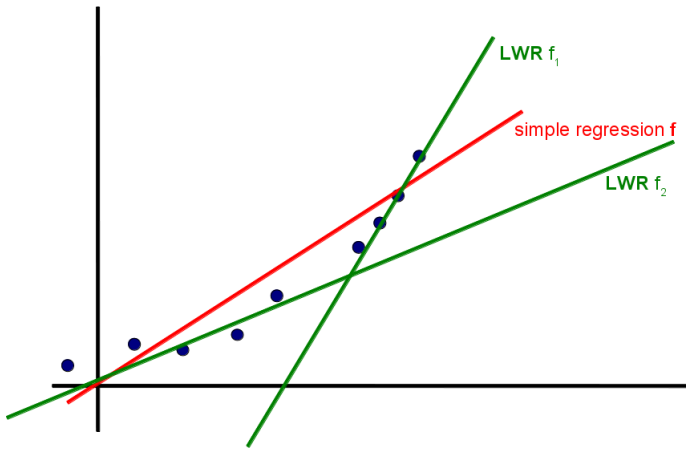
# Instance-based learning

## Locally weighted linear regression



# Instance-based learning

## LW linear regression vs. simple regression



# Naive Bayes classifier

## Bayes theorem

$$P(A|B) = \frac{P(A, B)}{P(B)}, P(B|A) = \frac{P(A, B)}{P(A)} \quad (10)$$

$$P(A|B) = \frac{P(A) * P(B|A)}{P(B)} \quad (11)$$

- $P(A)$  is the prior probability (marginal) probability of  $A$ . It does not take into account any information about  $B$ .
- $P(A|B)$  is the conditional probability of  $A$ , given  $B$ . So called the posterior probability because it depends upon the specified value of  $B$ ;  $P(B|A)$  is the conditional probability of  $B$  given  $A$ .
- $P(B)$  is the prior (marginal) probability of  $B$ , and acts as a normalizing constant.

# Naive Bayes classifier

## Bayes theorem

$$\Pr(Y | A_1, \dots, A_m) = \frac{\Pr(Y) \times \Pr(A_1, \dots, A_m | Y)}{\Pr(A_1, \dots, A_m)} \quad (12)$$

$$\text{posterior} = \frac{\text{prior} \times \text{likelihood}}{\text{evidence}} \quad (13)$$

# Naive Bayes classifier

## Conditional independence

Let  $X, Y$  and  $Z$  be three discrete random variables. We say that  $X$  is *conditionally independent* of  $Y$  given  $Z$  if

$\forall x_i, y_j, z_k, x_i \in \text{Values}(X), y_j \in \text{Values}(Y), z_k \in \text{Values}(Z) :$

$$\Pr(X = x_i | Y = y_j, Z = z_k) = \Pr(X = x_i | Z = z_k) \quad (14)$$

i.e.,  $P(X|Y, Z) = P(X|Z)$ .



# Naive Bayes classifier

## Conditional independence

### Thunder & Rain & Lighting

Assume three variables: Thunder, Rain, and Lighting.

Thunder is conditionally independent of Rain given Lighting:

$$\Pr(\text{Thunder}|\text{Rain}, \text{Lighting}) = \Pr(\text{Thunder}|\text{Lighting})$$

# Naive Bayes classifier

## Discriminative vs. generative classifiers

- **discriminative classifier** does not care about how the data was generated. It simply classifies a given example.
- **generative classifier** models how the data was generated in order to classify an example. It asks the question *Based on the generation assumptions, which class is most likely to generate this example?*

# Naïve Bayes classifier

## Discriminative vs. generative classifiers

$\Pr(y|\mathbf{x}) = ?$

- Logistic regression classifier is a **discriminative classifier**

$$h_{\Theta}(\mathbf{x}) = p(y = 1|\mathbf{x}, \Theta)$$

- Naïve Bayes classifier is a **generative classifier**

① Learn  $\Pr(\mathbf{x}|y)$  and  $\Pr(y)$

② Apply Bayes rule to get

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

③

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

# Naive Bayes classifier

If we work with two features  $A_1, A_2$  and we assume that they are conditionally independent given the target class  $Y$ , then

$$\begin{aligned} \Pr(A_1, A_2 | Y) &\stackrel{\text{product rule}}{=} \Pr(A_1 | A_2, Y) * \Pr(A_2 | Y) \stackrel{\text{conditional independence assumption}}{=} \\ &= \Pr(A_1 | Y) * \Pr(A_2 | Y) \end{aligned}$$

# Naïve Bayes classifier

Assume conditional independence of features  $A_1, \dots, A_m$  given  $Y$ . Thus

$$\begin{aligned}\Pr(\mathbf{x}|y) &= \Pr(x_1, x_2, \dots, x_m|y) \stackrel{\text{chain rule}}{=} \prod_{j=1}^m \Pr(x_j|x_1, x_2, \dots, x_{j-1}, y) \stackrel{\text{c. i. a.}}{=} \\ &= \prod_{j=1}^m \Pr(x_j|y)\end{aligned}$$

**Naïve Bayes classifier**

$$\hat{y} = \underset{y_k \in Y}{\operatorname{argmax}} \Pr(y_k) \prod_{j=1}^m \Pr(x_j|y_k) \quad (15)$$

# Naive Bayes classifier

Naive assumption of feature conditional independence given a target class is rarely true in real world applications. Nevertheless, Naive Bayes classifier surprisingly often shows good performance in classification.

# Naive Bayes Classifier is a linear classifier

NB classifier gives a method for predicting rather than an explicit classifier.

## Prediction function

$$\hat{y} = \operatorname{argmax}_{y_k \in Y} \Pr(y_k) \prod_{j=1}^m \Pr(x_j | y_k)$$

We focus on **binary classification**  $Y = \{0, 1\}$  with binary features  $A_1, \dots, A_m$ .

We predict  $\hat{y} = 1$  iff

$$\frac{\Pr(y = 1) \prod_{j=1}^m \Pr(x_j | y = 1)}{\Pr(y = 0) \prod_{j=1}^m \Pr(x_j | y = 0)} > 1$$

# Naive Bayes Classifier is a linear classifier

**Denote**  $p_j = \Pr(x_j = 1|y = 1)$ ,  $q_j = \Pr(x_j = 1|y = 0)$

Then

$$\frac{\Pr(y = 1) \prod_{j=1}^m p_j^{x_j} (1 - p_j)^{1-x_j}}{\Pr(y = 0) \prod_{j=1}^m q_j^{x_j} (1 - q_j)^{1-x_j}} > 1$$

$$\frac{\Pr(y = 1) \prod_{j=1}^m (1 - p_j) \left(\frac{p_j}{1-p_j}\right)^{x_j}}{\Pr(y = 0) \prod_{j=1}^m (1 - q_j) \left(\frac{q_j}{1-q_j}\right)^{x_j}} > 1$$

**Take logarithm**



# Naive Bayes Classifier is a linear classifier

$$\log \frac{\Pr(y = 1)}{\Pr(y = 0)} + \sum_{j=1}^m \log \frac{1 - p_j}{1 - q_j} + \sum_{j=1}^m \left( \log \frac{p_j}{1 - p_j} - \log \frac{q_j}{1 - q_j} \right) x_j > 0$$

NB classifier as a linear classifier where

$$\Theta_j = \log \frac{p_j}{1 - p_j} - \log \frac{q_j}{1 - q_j}$$

# Bayesian belief networks

## Motivation

Task: Will students fall asleep during the lecture?

$Attr = \{\text{It's raining, They are tired, They were at the party last night}\}$

$Values(\text{It's raining}) = \{\text{Yes, No}\}$

$Values(\text{They are tired}) = \{\text{Yes, No}\}$

$Values(\text{They were at the party last night}) = \{\text{Yes, No}\}$

$Y = \text{FallAsleep}, Values(\text{FallAsleep}) = \{\text{Yes, No}\}$

# Bayesian belief networks

## Motivation

Naïve Bayes assumption

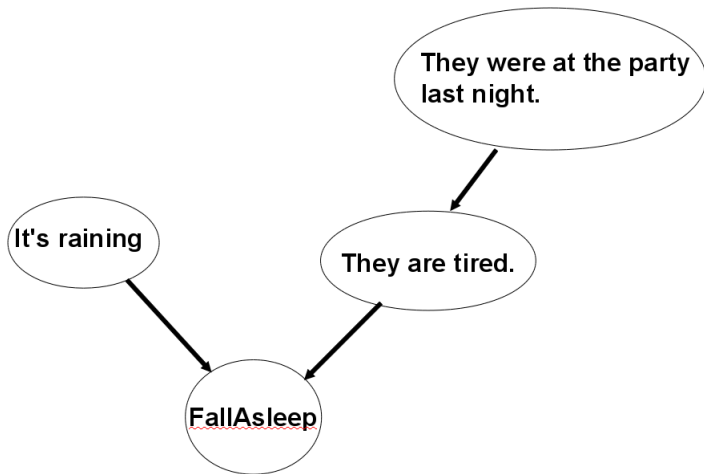
– features are conditionally independent given the value of the target class.



# Bayesian belief networks

## Motivation

... but



# Bayesian belief networks

## Motivation

- Naïve Bayes classifier assumes that ALL features are conditionally independent given the value of the target class.
- A Bayesian network is a graphical model that encodes probabilistic relationships among attributes of interest.
- BBNs allow stating conditional independence assumptions that apply to SUBSETS of the attributes.
- Dependencies are modeled as graph where nodes correspond to attributes and edges go from cause to effect.
- BBNs combine prior knowledge with observed data.
- BBNs are less constraining than the global assumption by NB.

# Bayesian belief networks

## Settings

Consider an arbitrary set of random variables  $X_1, X_2, \dots, X_m$ . Each variable  $X_i$  can take on the set of possible values  $Values(X_i)$ .

We define the **joint space** of the variables  $X_1, X_2, \dots, X_m$  to be the cross product  $Values(X_1) \times Values(X_2) \times Values(X_3) \times \dots \times Values(X_m)$ .

The probability distribution over the joint space is called the **joint probability distribution**  $\Pr(x_1, x_2, \dots, x_m)$  where  $x_1 \in Values(X_1), x_2 \in Values(X_2), \dots, x_n \in Values(X_m)$ .

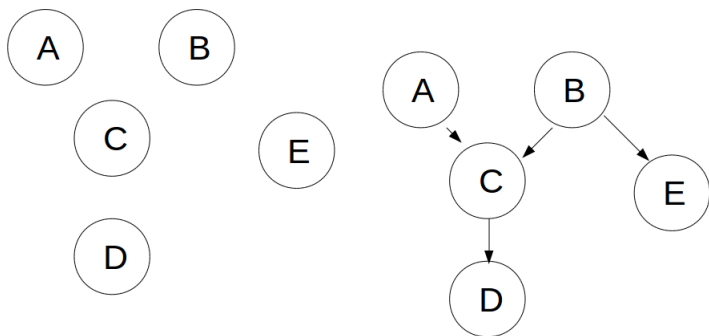
BBN describes the joint probability distribution for a set of variables by specifying a set of conditional independence assumptions together with sets of local conditional probabilities.

## Representation

- 1 A directed acyclic graph  $G = (V, E)$ 
  - nodes are random variables
  - arcs between nodes represent probabilistic dependencies
  - arcs are drawn from cause to effect
  - $Y$  is a *descendant* of  $X$  if there is a directed path from  $X$  to  $Y$ .
- 2 The network arcs represent the assertion that the variable  $X$  is conditionally independent of its nondescendants given its immediate predecessors  $Parents(X)$ ;  $\Pr(X|X_i)_{X_i \in Parents(X)}$
- 3 A set of tables for each node in the graph - a conditional probability table is given for each variable; it describes the probability distribution for that variable given the values of its immediate predecessors.

# Building a Bayes net

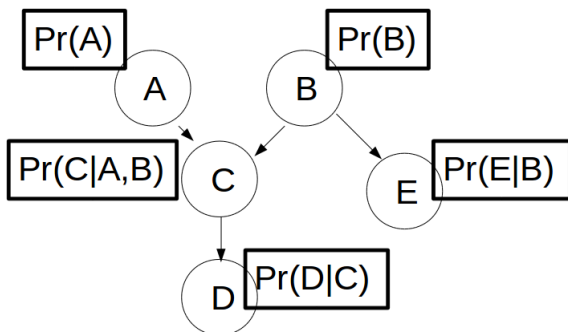
1. Choose the variables to be included in the net:  $A, B, C, D, E$
2. Add the links





# Building a Bayes net

3. Add a probability table for each root node  $\Pr(X)$  and nonroot node  $\Pr(X|X_i)_{X_i \in Parents(X)}$



## Once the net is built ...

The joint probability of any assignment of values  $x_1, x_2, \dots, x_m$  to the tuple of network variables  $X_1, X_2, \dots, X_m$  can be computed by the formula

$$\Pr(x_1, x_2, \dots, x_m) = \Pr(X_1 = x_1 \wedge X_2 = x_2 \wedge \dots \wedge X_m = x_m) = \prod_{i=1}^m \Pr(x_i | \text{Parents}(X_i)) \quad (16)$$

# Bayesian belief networks

Two components

- ① A function for evaluating a given network based on the data.
- ② A method for searching through the space of possible networks.

Learning the network structure

- searching through the space of possible sets of edges
- estimating the conditional probability tables for each set
- computing the quality of the network

# K2 algorithm

This 'search and score' algorithm heuristically searches for the most probable belief-network structure given a training data.

It starts by assuming that a node has no parents, after which, in every step it adds incrementally the parent whose addition mostly increase the probability of the resulting structure. K2 stops adding parents to the nodes when the addition of a single parent cannot increase the probability of the network given the data.

In general, the BBNs deal with probability propagation that consists of updating the probability values of the variables in a dependence graph, given some variables that have been observed.

# K2 algorithm

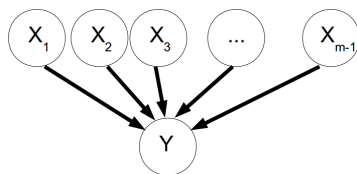
## INPUT

- a set of  $m$  nodes (i.e., attributes)
- an ordering on the nodes  $X_1, X_2, \dots, X_m$  (parents are before their kids)
- an upper bound  $u$  on the number of parents a node may have,
- training data  $D$ ,  $|D| = n$

## OUTPUT

- for each node, a printout of the parent nodes

Note on the initial nodes ordering: the Naïve Bayes Classifier is a network with an edge leading from the target feature to each other features. This network can be used as a starting point for the search.



# Summary of Examination Requirements

- Key points of instance-based learning – distance metric, number of neighbours, weighting function, fitting local points
- $k$ -NN (weighted) algorithm
- Locally weighted linear regression
- Discriminative and generative classifiers
- Naïve Bayes Classifier – conditional independence, linear decision boundary
- Bayesian belief networks – structure, conditional probabilities