

## Supervised learning for text mining

Master Data Mining

Julien Velcin



Warning

*I won't mention deep learning  
in this course*

## Text classification

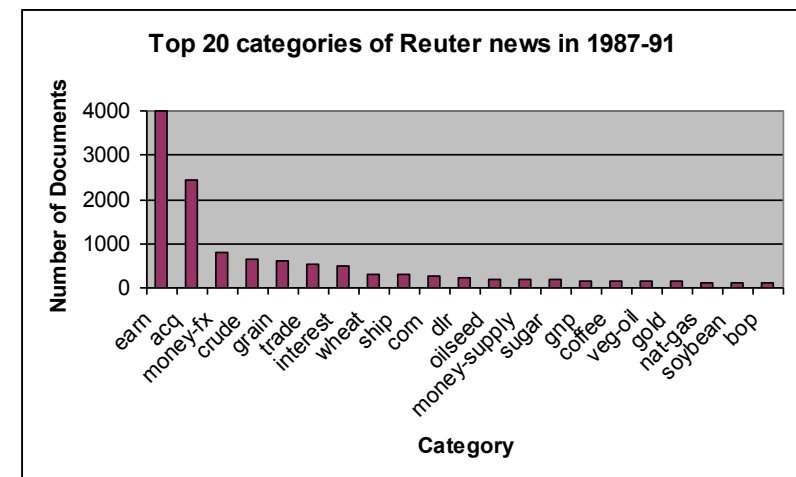
- Automatically classify documents into **predefined** classes  
e.g., digital libraries



Sport  
Politics  
Entertain.  
Economics  
etc.

- Application areas:
  - Email SPAM filtering
  - Internet directory construction (ex.: Yahoo!)
  - Automatic indexing ...

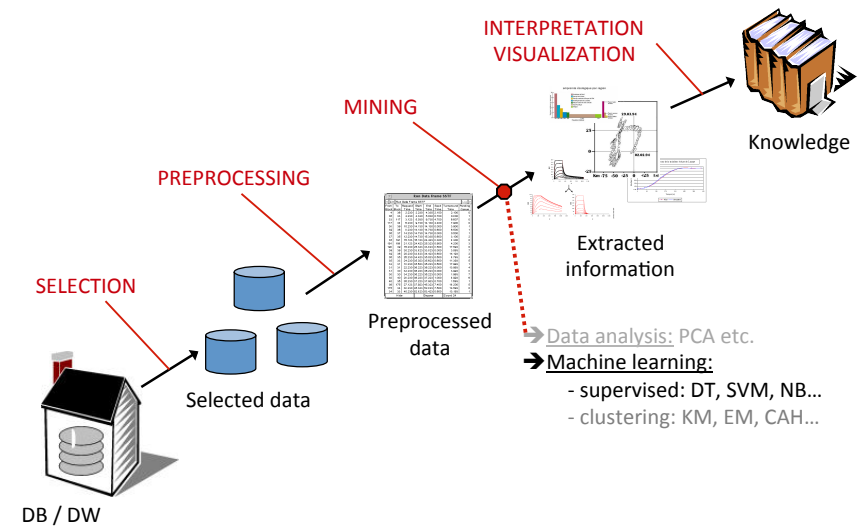
## A very classical dataset



## Basic ML approach

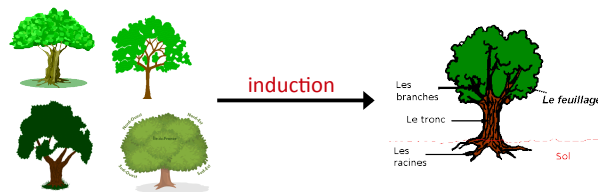
- Prepare a set of **training data**
  - Getting access to the data
  - Need of effective preprocessing
- Create a classifier
  - Apply a ML algorithm: NB, ANN, SVM etc.
- Classify new documents using the **classifier**

## ML in data mining



## Inductive principle

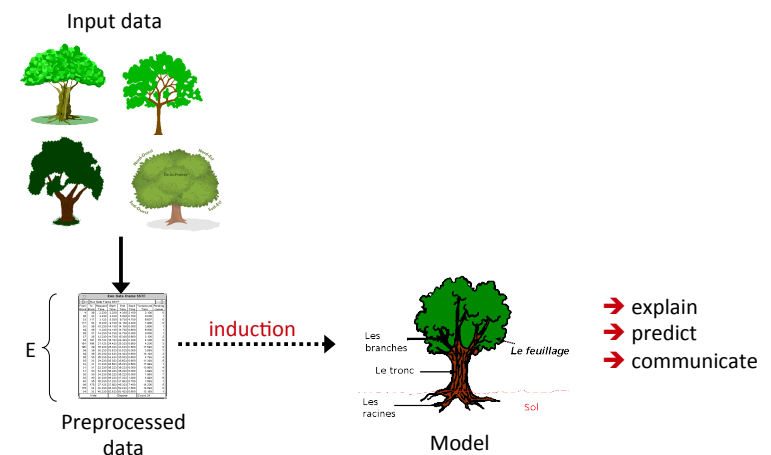
- Aristotle: knowledge comes from the (observable) world



A tree = roots + trunk + branches + leaves

+ **relations** between: roots and ground, roots and trunk etc.

## Induction using a machine



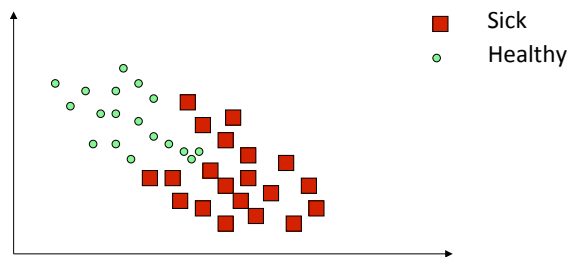
# Supervised learning

- Inductive learning based on **supervision**
- Given:
  - o a representation language  $L$ ,
  - o examples  $e_i \in E$  described using  $L$ ,
  - o for each  $e_i$ , a class  $\Phi(e_i)$  taken from  $\{c_1, c_2, \dots, c_p\}$
- The objective is to find:
  - a function (machine)  $h$  that relates each description built on  $L$  to a class in  $\{c_1, c_2, \dots, c_p\}$  → classifieur
  - a function  $h$  that relates each description built on  $L$  to a real value → régression

# A lot of existing algorithms

- K-Nearest Neighbours (KNN)
  - Roccio's algorithm
  - Decision Trees (DT)
  - Artificial Neuronal Networks (ANN)
  - Naive Bayes (NB)
  - Support Vector Machines (SVM)
  - Ensemble methods (bagging, boosting)
- etc.

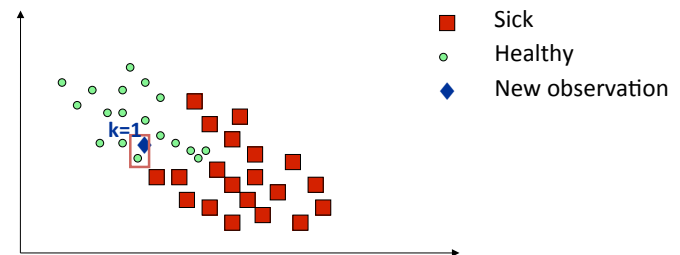
## K-Nearest Neighbours



- Using a metric  $d(e_i, e_j)$  that calculates the dissimilarity between two individuals  $e_i$  and  $e_j$ .

## K-Nearest Neighbours (cont'd)

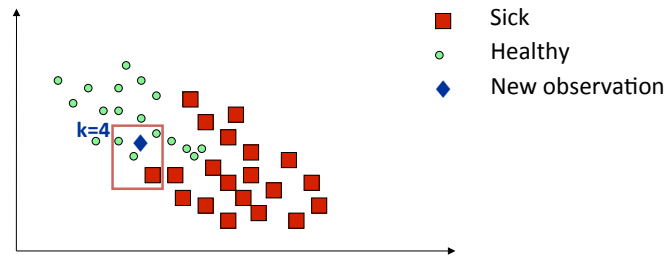
Given a new observation: ♦



**Conclusion:** we predict that this new person is healthy.

## K-Nearest Neighbours (cont'd)

- With  $k=4$ :



- **Same conclusion:** we predict that this new person is healthy.

## K-Nearest Neighbours (cont'd)

- Pros:
  - Effective
  - Non-parametric
  - Pairwise (local) comparison of documents
- Cons:
  - Classification time is long (usually in  $O(n^2)$ )
  - Difficult to find an optimal value of  $k$

## Rocchio's algorithm

- Build prototype vector for each class
- Prototype vector: average vector over all training document vectors that belong to class  $c_i$

$$\vec{\mu}(c) = \frac{1}{|D_c|} \sum_{\vec{d} \in D_c} \vec{d}$$

- Calculate similarity between test document and each of prototype vectors
- Assign test document to the class with maximum similarity

## Rocchio's algorithm (cont'd)

- Pros:
  - Easy to implement
  - Very fast learner
- Cons:
  - Low classification accuracy
  - Only for linear classification boundaries

## Need of inductive bias in ML

- Bias = set of assumptions that the learner uses to predict outputs given inputs that it has not encountered [Mitchell, 80]
- Choices for reducing the solution (hypotheses) space; guide of the learning process
- In concept learning: way to favour a generalization rather than another one
- **Absolut need** of bias for learning [Mitchell, 80].
- Examples :
  - Occam's razor,
  - Conditional independance (cf. NB)
  - Maximum margin (cf. SVM)
  - etc.

## Learning as search

- Choosing a family of concepts (hypotheses)
- Searching for the best hypothesis possible
- To do this, you need fixing biases
- Occam's razor addresses the trade-off:

simplicity / efficiency

## Illustration of the Occam's razor

- Given the sequence: 1, 2, 3, 5 ... , a?

## Statistical machine learning

- Main concepts:
  - Data (examples, instances, observations)
  - Hypotheses (probabilistic theories on the process of data generation)
- Objective:
  - Explain the observed data
  - Predict new observations

## Example [Russel & Norvig, 2003]

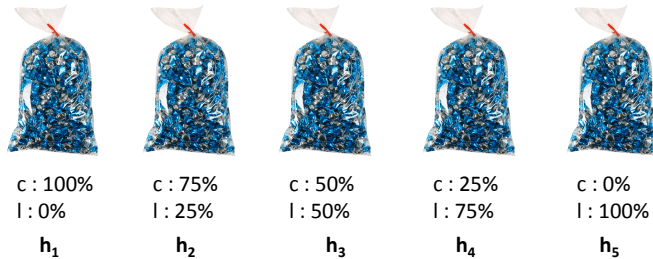
- Two flavours for canddies:

– Cherry

– Lime



- Canddies are sold in 5 different bags:



## Example (cont'd)

- Given a new “unknown” bag:



- H can take 5 values:

$$H \in \{h_1, h_2, h_3, h_4, h_5\}$$

- Observations: opening the wrappers!

Ex. : d1, d2, d3, d4, d5, d6, d7, d8, d9, d10, etc.



## Bayesian learning

- Hypothesis probability given the data:

$$P(h_i/D) = \frac{P(D/h_i)P(h_i)}{P(D)} \propto P(D/h_i)P(h_i)$$

- Making predictions:

$$\begin{aligned} P(X/D) &= \sum_i P(X/D, h_i)P(h_i/D) \\ &= \sum_i P(X/h_i)P(h_i/D) \end{aligned}$$

## Bayesian learning (cont'd)

$$P(X/D) \propto \sum_i P(X/h_i)P(D/h_i)P(h_i)$$

probability distribution (pointing to  $P(X/h_i)$ )  
likelihood (pointing to  $P(D/h_i)$ )  
prior (pointing to  $P(h_i)$ )

- How to calculate the likelihood?

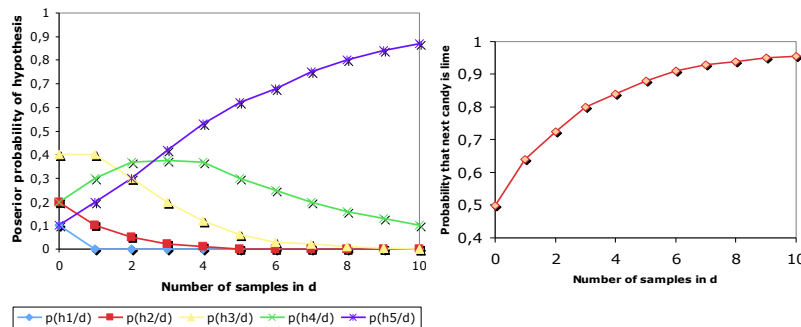
$$P(D/h) = \prod_i P(d_i/h)$$

**i.i.d. hypothesis** (independently identically distributed)

- Example of observed data:

$P(D/h_3) = 0,5^{10}$

## Evolution of learning



Here are the priors:  $\langle 0,1 ; 0,2 ; 0,4 ; 0,2 ; 0,1 \rangle$

## Role of priors

- $P(h_i)$  very important in Bayesian learning
- Bounding the complexity of models
- Complex models (**hypotheses**) are less likely *a priori*
- **Trade-off** between complexity and expressivity
- If you have no idea on priors: uniform

## Decision in bayesian learning

- Given  $p(X/D)$  what decision make?
  - the simplest way: taking the most likely solution
  - another solution: drawing a random number
- MAP = Maximum A Posteriori
  - here:  $h_5$ .

$$\arg \max_i P(h_i / d)$$

$$P(X / D) \approx P(X / h_{MAP})$$

## Link with information theory

- Choosing  $h_{MAP}$  for maximizing  $P(D/h_i)P(h_i)$  means minimizing:
 
$$-\log_2 P(D/h_i) - \log_2 P(h_i)$$
- If you split this formula into 2 parts:
  - $-\log_2 P(h_i)$  coding hypothesis  $h_i$
  - $-\log_2 P(D/h_i)$  additional coding for  $D$

MAP → maximum compression of the data!

## Bayesian learning of parametric models

- To begin with: we assume that the data are **complete** (no missing values)
- The objective is to estimate the values of the parameters of a **given model**
- Classical approach: **Maximum-Likelihood Estimation** (MLE) for estimating the model parameters

## Learning of discrete models

**Example:** Selling bags containing candies with two different flavours

$\theta$  with the cherry flavour

$1-\theta$  with the lime flavour

- The problem hypothesis =  $h_\theta$ .
- Assume uniform priors
- Just one random variable: *Flavour*.



## Example (cont'd)

- N observations:

o c candies with cherry flavour



o l candies with lime flavour



- Estimating the likelihood:

$$P(D/h_\theta) = \prod_{j=1}^N P(d_j/h_\theta) = \theta^c \cdot (1-\theta)^l$$

- $h_{\text{MAP}} \rightarrow$  taking  $\theta$  for maximizing this formula  
= maximizing the log-likelihood

## Calculating $h_\theta$

- Likelihood:  $P(D/h_\theta) = \prod_{j=1}^N P(d_j/h_\theta) = \theta^c \cdot (1-\theta)^l$
- Log-likelihood:

$$L(D/h_\theta) = \log P(D/h_\theta) = \sum_{j=1}^N \log P(d_j/h_\theta) = c \log \theta + l \log(1-\theta)$$

Which  $\theta$  maximizes this formula?

$$\frac{\partial L(D/h_\theta)}{\partial \theta} = \frac{c}{\theta} - \frac{l}{1-\theta} = 0 \quad \rightarrow \quad \theta = \frac{c}{c+l} = \frac{c}{N}$$



## First conclusion

$$\theta = \frac{c}{c+l} = \frac{c}{N}$$

- The **real rate**  $\theta$  of cherry candies is estimated to the rate of **observed** cherry candies...
- Quite an obvious result! But we have:
  - written a (log)likelihood function,
  - derived this function wrt  $\theta$ ,
  - found the best value for the parameter  $\theta$ .

## Estimating the likelihood

- Unwrapping  $N$  candies:
  - $c$  with cherry flavour,  $l$  with lime flavour,
  - cherry candies:  $r_c$  red wrappers,  $g_c$  green wrappers,
  - lime candies:  $r_l$  red wrappers,  $g_l$  green wrappers.

$$P(D/h_{\theta,\theta_1,\theta_2}) = \theta^c \cdot (1-\theta)^l \cdot \theta_1^{r_c} \cdot (1-\theta_1)^{g_c} \cdot \theta_2^{r_l} \cdot (1-\theta_2)^{g_l}$$

$$L = [c \cdot \log \theta + l \cdot \log(1-\theta)] + [r_c \log \theta_1 + g_c \log(1-\theta_1)] + [r_l \log \theta_2 + g_l \log(1-\theta_2)]$$

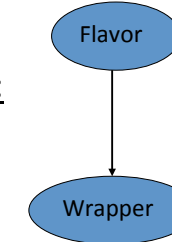
## New more complicated example

Adding wrappers:  
red and green

Model with 3 parameters:

$\theta, \theta_1, \theta_2$

$P(F=\text{cherry}) : \theta$



$F$	$P(W=\text{red}/F)$
cherry	$\theta_1$
lime	$\theta_2$

What is the likelihood of a cherry candy in a green wrapper?

$$P(F = c, W = g / h_{\theta,\theta_1,\theta_2})$$

$$= P(F = c / h_{\theta,\theta_1,\theta_2}) P(W = g / F = c, h_{\theta,\theta_1,\theta_2}) = \theta \cdot (1 - \theta_1)$$

## Calculating $h_{\theta,\theta_1,\theta_2}$

- Derivations:

$$\frac{\partial L}{\partial \theta} = \frac{c}{\theta} - \frac{l}{1-\theta} = 0 \quad \rightarrow \quad \theta = \frac{c}{c+l}$$

$$\frac{\partial L}{\partial \theta_1} = \frac{r_c}{\theta_1} - \frac{g_c}{1-\theta_1} = 0 \quad \rightarrow \quad \theta_1 = \frac{r_c}{r_c + g_c}$$

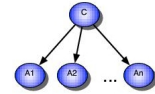
$$\frac{\partial L}{\partial \theta_2} = \frac{r_l}{\theta_2} - \frac{g_l}{1-\theta_2} = 0 \quad \rightarrow \quad \theta_2 = \frac{r_l}{r_l + g_l}$$

- Same conclusion...

## Preliminary conclusions

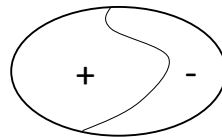
- Results fit our expectations
- Using complete data and  $h_{\text{MAP}}$ , Bayesian learning split the problem into **subproblems** (one for each model parameter)
- **Mind** the “zero probability” situation (especially for small datasets)

## Naive Bayes



- Naive Bayes (NB) is a classical classifier in ML
- Based on the **Bayes' Theorem**
- Strong independence assumptions
- Easy and quick to implement
- Good results in practice

## Basics on NB



- Example:  $e = \text{tuple}(x_1, x_2, \dots, x_m) + \text{class}$
- Case of binary classification:  $c \in \{+, -\}$
- Learning set:  
$$E = \{(e_1, +), (e_2, +), (e_3, -), \dots, (e_n, +)\}$$
- We are looking for a function  $f_B$  so that:  
$$f_B(e) = \frac{p(C = +/e)}{p(C = -/e)}$$
- If  $f_B \geq 1$  then  $h(e)=+$  else  $h(e)=-$

## Calculating $f_B$

- How to calculate  $p(c/e)$ ?
- Bayes' rule:

$$p(c/e) = \frac{p(e/c) \times p(c)}{p(e)}$$

- And now...  $p(e/c)$  ?

$$p(e/c) = p(x_1, x_2, x_3, \dots, x_m / c)$$

- What else?

# conditional independence

We state the following fundamental hypothesis:

**Hypothesis:** We assume that all the attributes are **independent** given the associated class.

So:

$$p(x_1 / x_2, x_3, \dots, x_m, c) = p(x_1 / c)$$

$$p(e / c) = p(x_1, x_2, x_3, \dots, x_m / c) = \prod_{i=1}^m p(x_i / c)$$

## Illustration (Quinlan,83)



Attributs	Pif	Temp	Humid	Vent	Golf
Valeurs	soleil, couvert, pluie	chaud, bon, frais	normale, haute	vrai, faux	jouer, ne_pas_jouer

1	soleil	chaud	haute	faux	ne_pas...
2	soleil	chaud	haute	vrai	ne_pas...
3	couvert	chaud	haute	faux	jouer
4	pluie	bon	haute	faux	jouer
5	pluie	frais	normale	faux	jouer
6	pluie	frais	normale	vrai	ne_pas...
7	couvert	frais	normale	vrai	jouer
8	soleil	bon	haute	faux	ne_pas...
9	soleil	frais	normale	faux	jouer
10	pluie	bon	normale	faux	jouer
11	soleil	bon	normale	vrai	jouer
12	couvert	bon	haute	vrai	jouer
13	couvert	chaud	normale	faux	jouer
14	pluie	bon	haute	vrai	ne_pas...

class

# The Naive Bayes classifier (NB)

- With this hypothesis,  $f_B \rightarrow f_{NB}$ :

$$f_{NB}(e) = \frac{p(C = +)}{p(C = -)} \prod_{i=1}^m \frac{p(x_i / C = +)}{p(x_i / C = -)}$$

with:

- o  $p(C=+)$ : proportion of positive instances,
- o  $p(C=-)$ : proportion of negative instances,
- o  $p(x_i/C=+)$ : proportion of positive instances with  $X=x_i$
- o  $p(x_i/C=-)$ : proportion of negative instances with  $X=x_i$

jouer = classe + (9 exemples sur 14)

1	soleil	chaud	haute	faux	ne_pas...
2	soleil	chaud	haute	vrai	ne_pas...
3	couvert	chaud	haute	faux	jouer
4	pluie	bon	haute	faux	jouer
5	pluie	frais	normale	faux	jouer
6	pluie	frais	normale	vrai	ne_pas...
7	couvert	frais	normale	vrai	jouer
8	soleil	bon	haute	faux	ne_pas...
9	soleil	frais	normale	faux	jouer
10	pluie	bon	normale	faux	jouer
11	soleil	bon	normale	vrai	jouer
12	couvert	bon	haute	vrai	jouer
13	couvert	chaud	normale	faux	jouer
14	pluie	bon	haute	vrai	ne_pas...

Attributs	Pif	Temp	Humid	Vent	Golf
Valeurs	Soleil (2), couvert (4), pluie (3)	chaud (2), bon (4), frais (3)	normale (6), haute (3)	vrai (3), faux (6)	jouer, ne_pas_jouer

ne\_pas\_jouer = classe - (5 exemples sur 14)

1	soleil	chaud	haute	faux	ne_pas...
2	soleil	chaud	haute	vrai	ne_pas...
3	couvert	chaud	haute	faux	jouer
4	pluie	bon	haute	faux	jouer
5	pluie	frais	normale	faux	jouer
6	pluie	frais	normale	vrai	ne_pas...
7	couvert	frais	normale	vrai	jouer
8	soleil	bon	haute	faux	ne_pas...
9	soleil	frais	normale	faux	jouer
10	pluie	bon	normale	faux	jouer
11	soleil	bon	normale	vrai	jouer
12	couvert	bon	haute	vrai	jouer
13	couvert	chaud	normale	faux	jouer
14	pluie	bon	haute	vrai	ne_pas...

Attributs	Pif	Temp	Humid	Vent	Golf
Valeurs	Soleil (3), couvert (0), pluie (2)	chaud (2), bon (2), frais (1)	normale (1), haute (4)	vrai (3), faux (2)	jouer, nepas_jouer

## Predicting the class of a new observation

- Bob would like to know whether he can go to play golf. He observes that the weather is **sunny**, the temperature is **high**, the humidity is **normal** and there is **no wind**.

$o = \{(Pif=soleil), (Temp=bon), (Humid=normal), (Vent=faux)\}$

$$f_{NB}(o) = \frac{p(C=+)}{p(C=-)} \prod_{i=1}^m \frac{p(x_i / C=+)}{p(x_i / C=-)}$$

## Calculating the probabilities

$$p(C=+) = p(\text{jouer}) = 9/14 \approx 0,64$$

$$p(C=-) = p(\text{ne\_pas\_jouer}) = 5/14 \approx 0,36$$

Attributs	Pif	Temp	Humid	Vent	Golf
Valeurs	Soleil (2), couvert (4), pluie (3)	chaud (2), bon (4), frais (3)	normale (6), haute (3)	vrai (3), faux (6)	jouer
$p(x_i / C=+)$	0,22 ; 0,45 ; 0,33	0,22 ; 0,45 ; 0,33	0,67 ; 0,33	0,33 ; 0,67	

Attributs	Pif	Temp	Humid	Vent	Golf
Valeurs	Soleil (3), couvert (0), pluie (2)	chaud (2), bon (2), frais (1)	normale (1), haute (4)	vrai (3), faux (2)	nepas_jouer
$p(x_i / C=-)$	0,6 ; 0 ; 0,4	0,4 ; 0,4 ; 0,2	0,2 ; 0,8	0,6 ; 0,4	

$o = \{(P=s), (T=b), (H=n), (V=f)\}$

## Predicting the class of a new observation (cont'd)

$$f_{NB}(o) = \frac{p(C=+)}{p(C=-)} \prod_{i=1}^m \frac{p(x_i / C=+)}{p(x_i / C=-)} = \frac{9}{14} \frac{14}{5} \prod_{i=1}^m \frac{p(x_i / C=+)}{p(x_i / C=-)}$$

$$f_{NB}(o) = \frac{9}{5} \left[ \frac{p(P=s/+)}{p(P=s/-)} \cdot \frac{p(T=b/+)}{p(T=b/-)} \cdot \frac{p(H=n/+)}{p(H=n/-)} \cdot \frac{p(V=f/+)}{p(V=f/-)} \right]$$

$$f_{NB}(o) \approx \frac{9}{5} \left[ \frac{0,22}{0,6} \cdot \frac{0,45}{0,4} \cdot \frac{0,67}{0,2} \cdot \frac{0,67}{0,4} \right]$$

$$f_{NB}(o) \approx 1,8 \times [0,37 \times 1,13 \times 3,35 \times 1,68] \approx 4,24 \geq 1$$

**Conclusion:** Bob can go to play golf!

## Pros of NB

- Simplicity of the theoretical background
- Memory complexity = probability tables
- Calculus are really fast: mainly x
- Often a good baseline for experiments
- Good results for real datasets

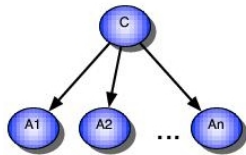
## Cons of NB

- Strong hypotheses (e.g., the conditional independence)
- “Zero-probability” issue  
E.g., test with  $\phi = \{(P=c), (T=ch), (H=h), (V=v)\}$
- Theoretical results not really well known

## Link with Bayesian networks

- Very simple form of a Bayesian network

$$\forall X_i, p(X_i = x_{ij} / X_1, X_2, \dots, X_m, C) = p(X_i = x_{ij} / C)$$



## Managing continuous variables

- Basic NB deal only with discrete variables
- Possible solutions for continuous variables:
  - discretizing the variables,
  - adding an **hypothesis** on the data distribution (e.g., using a Gaussian distribution)

$$P(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$