

Université

de Strasbourg

Introduction to data mining

Maja Temerinac-Ott

temerinacott@unistra.fr



Slides adopted from Germain Forestier

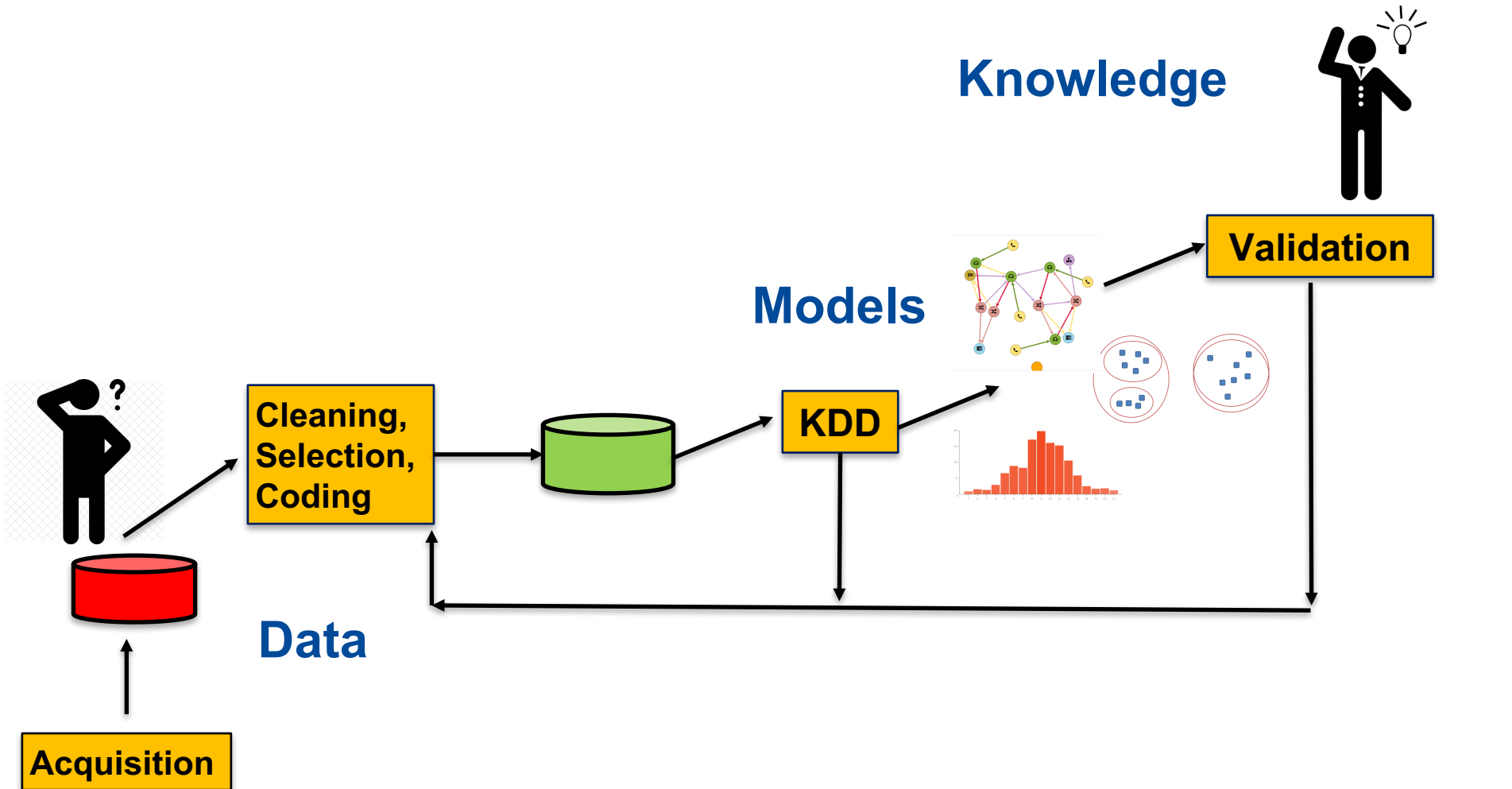
- Growing availability of huge amounts of data
 - Technology is available
 - Data collection: barcode, scanners, satellites, logs servers, etc.
 - To help store: database, data warehouses, digital libraries, www
- Why should we extract knowledge?
 - Economic necessity
 - E-commerce
 - High degree of competition
 - Customization, customer loyalty, market segmentation
- Customer data
- Digitalization of text, images, video, voice, etc.
- Internet and online catalogs

- The data quantity is too big to be treated manually or by classical algorithms:
 - The number of entries is millions to billions
 - Multi-dimensional data
 - Heterogeneous sources of data
- The user is full of data, but does not know how to understand it:
 - “The greatest problem of today is how to teach people to ignore the irrelevant, how to refuse to know things, before they are suffocated. For too many facts are as bad as none at all.”
(W.H. Auden)
- What do we need?
 - Extract interesting and useful knowledge from the data : rules, regularity, irregularities, patterns, constraints

- Extraction of implicit original (non-trivial) information, previously unknown and potentially useful from databases:
 - **Not trivial**: otherwise knowledge is not useful
 - **Implicit**: hidden knowledge is difficult to observe
 - **Unknown until now**: obvious!
 - **Potentially useful**: usable, understandable
- Whole process of discovery and interpretation of regularity in data
- Other names:
 - Knowledge Discovery in Databases (KDD)
 - Knowledge extraction
 - Data/pattern analysis
 - Data Analytics
 - Big Data

- Form groups of 4-5 people
- You will get a list of steps in the data mining process
- Please order the steps in chronological order!
- You have 5 min for the task.

1. Identify the problem
2. Find data
3. Clean the data
4. Coding of data, actions on variables
5. Search for a model, knowledge, etc.
6. Validation and interpretation of the result,
with possible return to the results in
previous steps
7. Integration of new knowledge

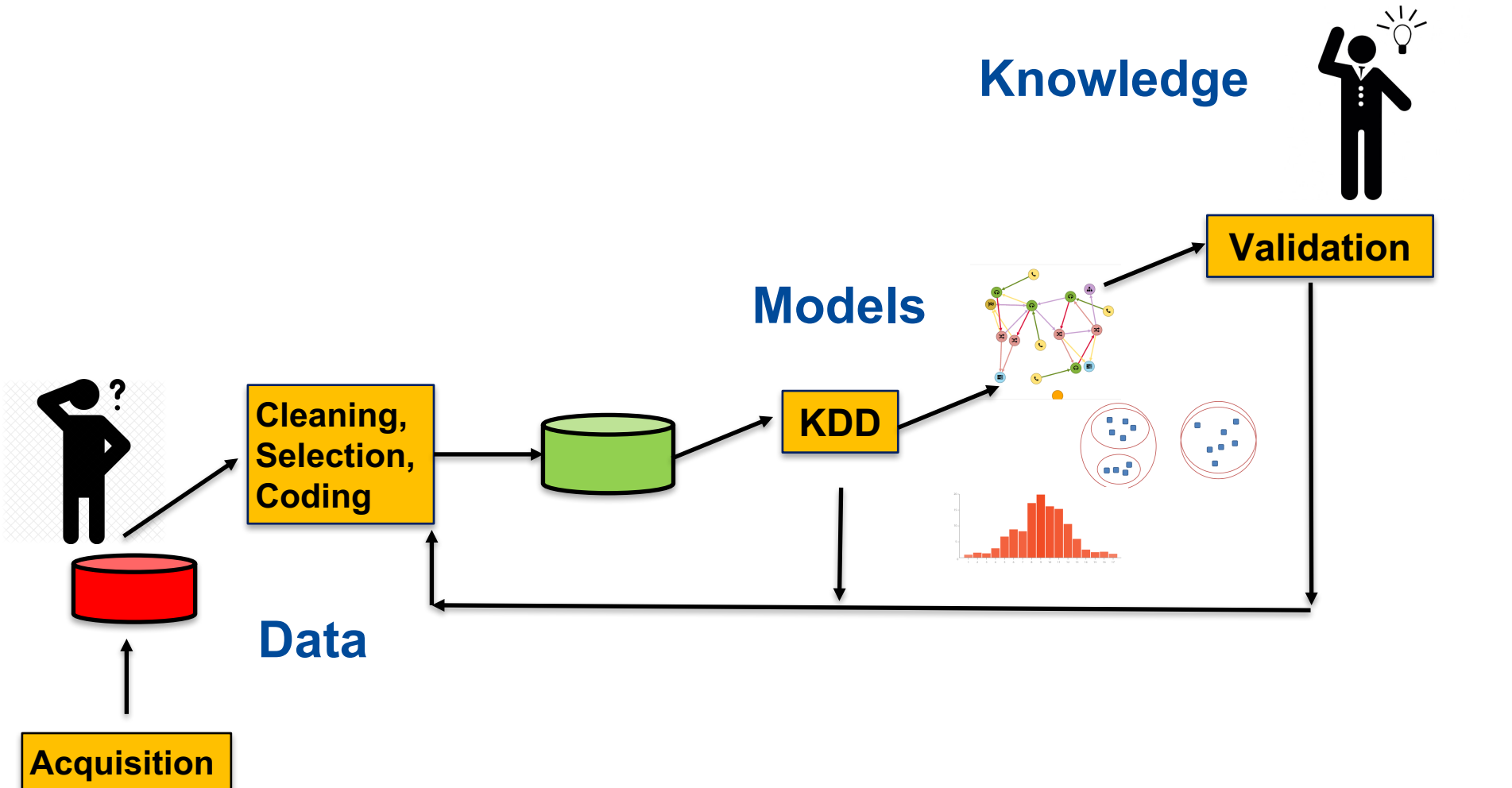


- Existing or needed data:
 - Files: information contained in one or more independant files
 - Relational Database: information contained in several files by a common `key`
 - Database transaction
- Data cleaning:
 - duplicates, input errors, outliers, missing information (ignored observations, average values, mean values on class, regression, etc.)

- **Data warehouses:** collection of data collected from multiple heterogeneous sources
- Data is recorded, cleaned, transformed and integrated
- Usually modeled by a multidimensional data structure (cube):
 - data is structured along several lines of analysis (dimensions of the cube) such as time, location, etc.
 - A cell is the intersection of different dimensions
 - The calculation of each cell is carried out at loading
 - The response time is thus stable whatever is requested

- The cubes are well suited for quick searches and analysis of data: On-Line Analytical Processing (OLAP)
 - What is the number of pairs of shoes sold by the store "OnVendDesChaussuresIci" in May 2016 AND Compare the sales with the same month of 2015 and 2014?
 - What are the components of production machines that have more unpredictable incidents during the period 2015-2016?
- Responses to OLAP requests may take a few seconds to several minutes.

- Selection of the data
 - Data sampling
 - Selection of sources
- Dimensionality reduction:
 - Selection or transformation of attributes
 - Weighting
- Coding:
 - Aggregation (sum, average), discretization, coding of attributes discrete, unification of benchmarks and standardization



- Goal: Learn something new!
 - **Concepts**: regrouping of data based on shared characteristics
 - **Associations**: correlations between attributes or data
- Principles:
 - Getting the highest level of abstraction possible
 - Rules or truths that are the basis for other truths

- Different approaches:

- 1. Estimation:**

- create a model that best describes a prediction variable for real data

- 2. Classification:**

- create a function that classifies an element among several pre-existing classes

- 3. Clustering** (regrouping):

- search to identify a finite set of categories or groups which can describe the data

- 4. Dependency modelling:**

- find a model that describes significant dependencies between variables

- **Supervised learning:**
 - Inductive model where the learner considers a set of grouped examples representative for the learning task (class of belonging, ownership, etc.): the examples are labeled beforehand
- **Predictive data mining:**
 - Divide / group instances into special classes for future predictions
 - Predicting unknown or missing values
- **Algorithms:**
 - Decision trees, classifications, genetic algorithms, regression (linear and non-linear)

■ Induction:

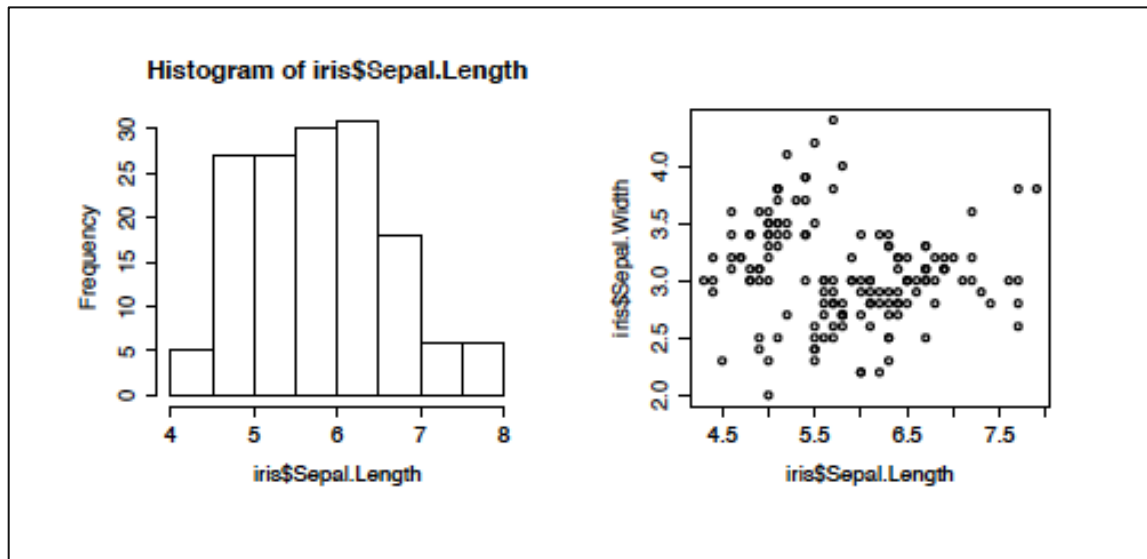
- It is a commonly used technique
- Generalization of an observation or reasoning established from singular cases.
- It consists in drawing conclusions from a series of facts

■ Example:

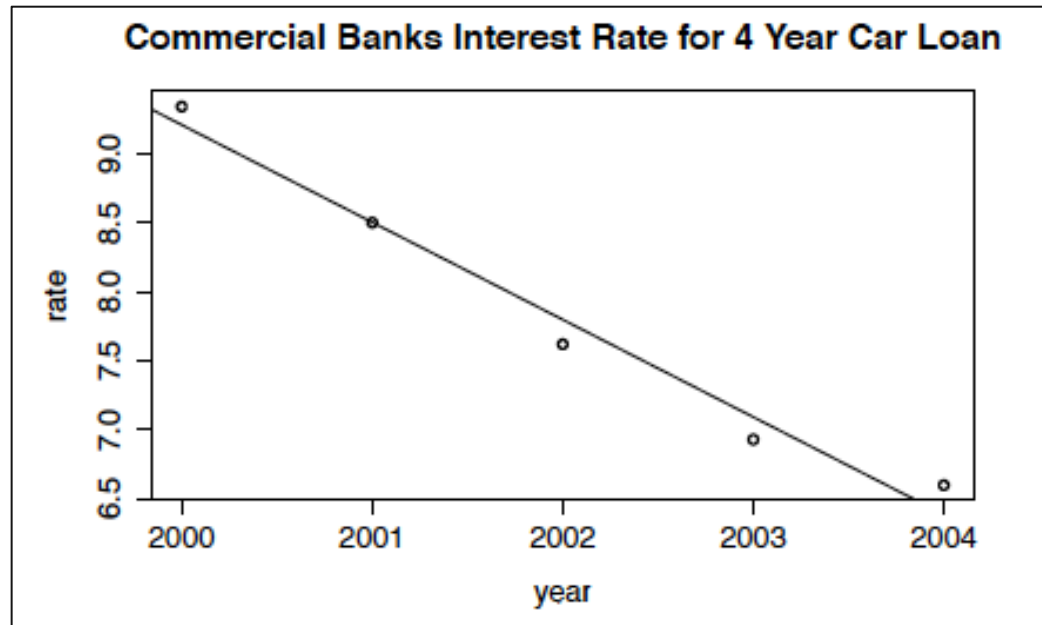
- Induction:
 - water, oil and milk freeze when they are cooled, we will infer that all liquids must freeze, provided the cold is rather intense
- Deduction:
 - all liquids are likely to freeze; so, if mercury is a liquid, it can freeze

- **Non-supervised learning:**
 - Construction of a model and discovery of relations is given without reference to other data
 - There is no prior information on the data
- **Explanation:**
 - Grouping instances into special classes based on their resemblance or the sharing of properties. The classes are unknown and are therefore created: they are used to explain "or summarize the data
- **Algorithms:**
 - Segmentation, grouping, discovery of associations and rules

- Obtain a visual representation of the data
- Not always possible depending on the data type
- Not always possible depending on the amount of data

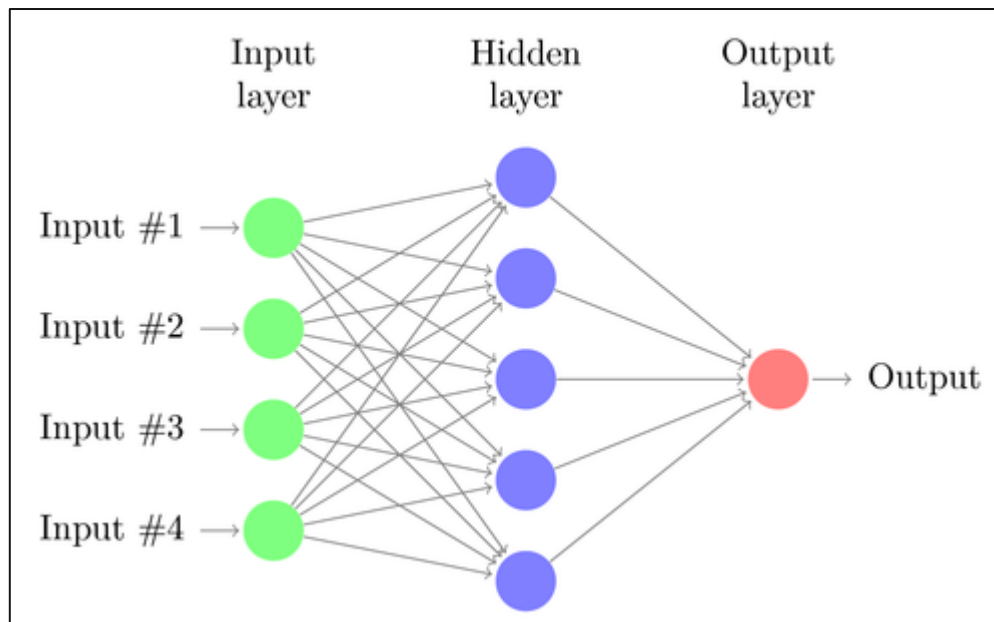


- **Regression:**
 - Analyze the relationship of one variable vs. one or more other variables
 - Least-squared method



■ Regression:

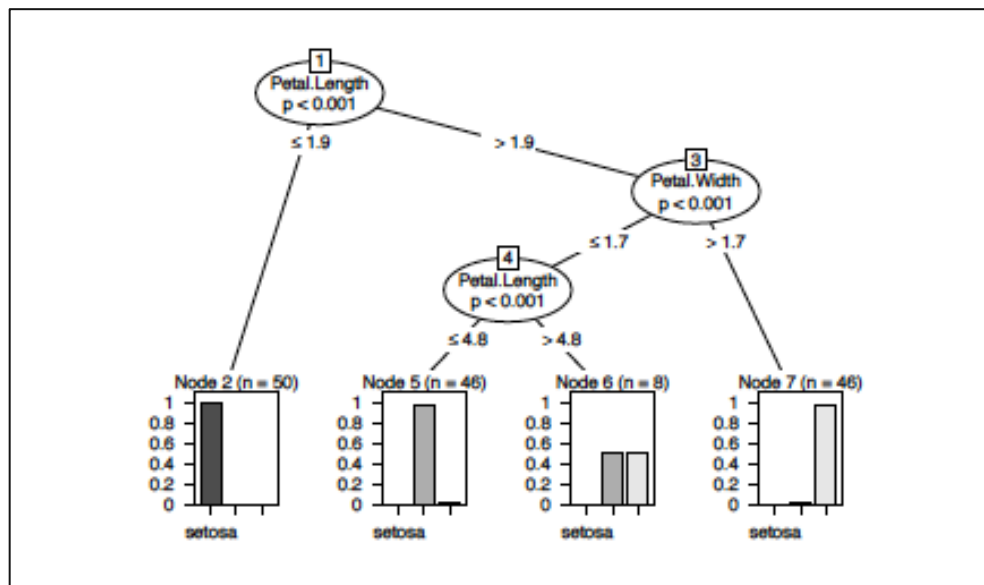
- Analyze the relationship of one variable vs. one or more other variables
- Neural network



- Division of the data set into disjoint classes
- **Goal:**
 - search for a set of predicates characterizing a class of objects and which can be applied to unknown objects in order to identify their class of belonging.
- **Principales techniques:**
 - Decision trees
 - Bayesian classifier
 - k-nearest neighbor
 - Neuronal network
 - SVM
 - Genetic algorithm

Decision trees:

- Classify objects into subclasses by hierarchical divisions
- Automatic construction from a sample
- There are several techniques to build the tree

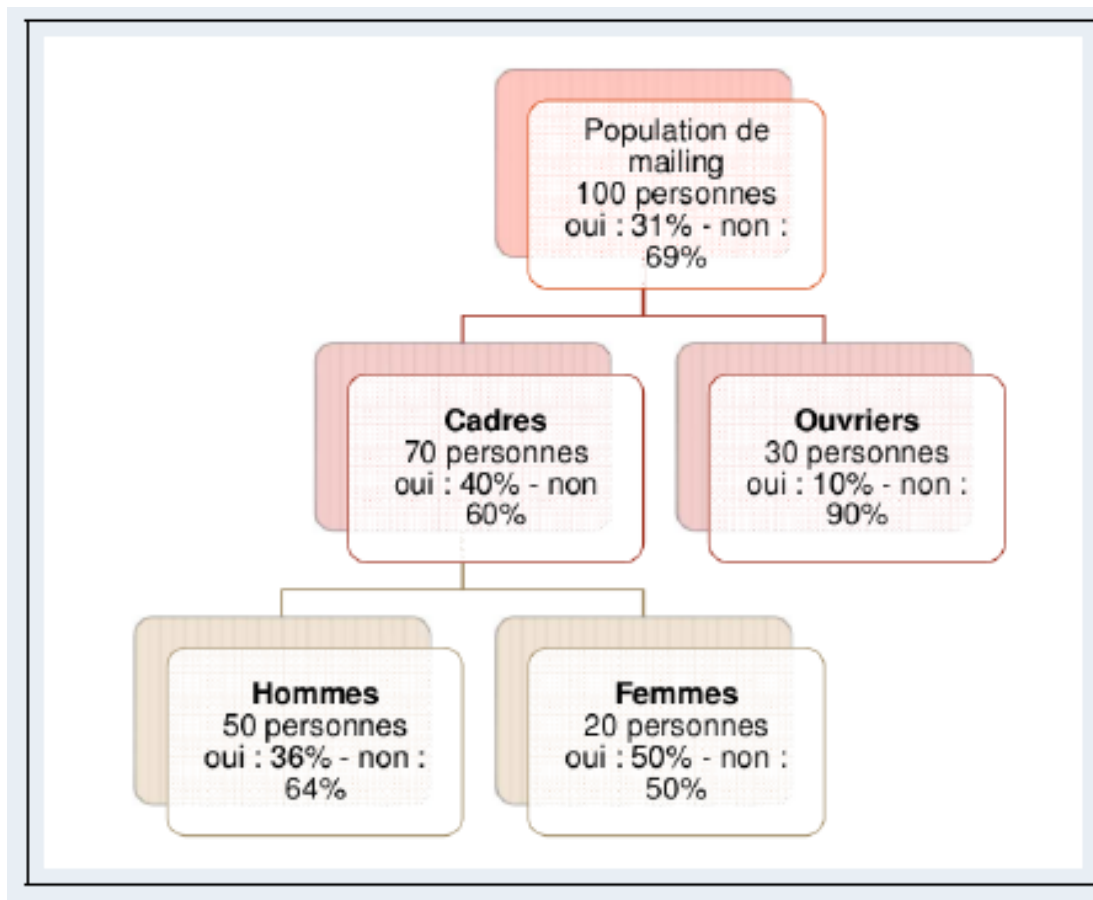


■ Decision Tree Example:

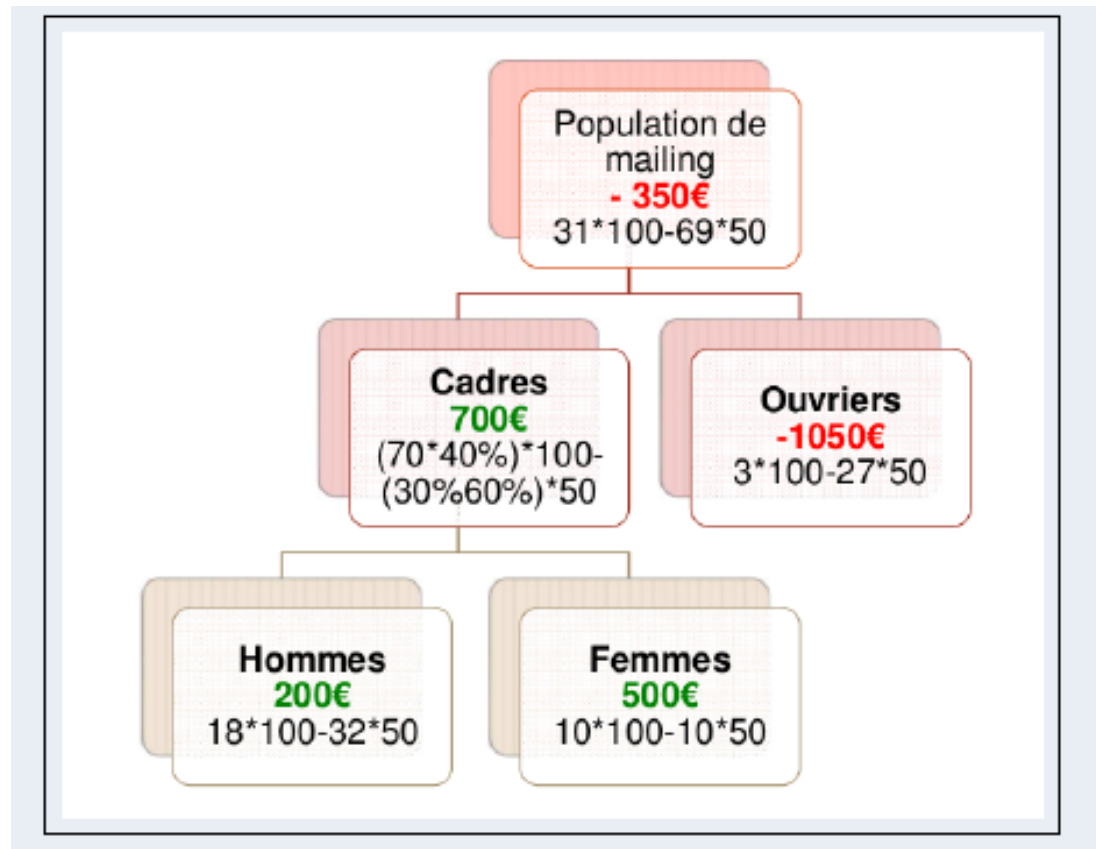
- You are working in a company who sends gifts to potential customers, who then can respond to place an order.
- If the customer does not respond to the gift the company pays 50 euro; otherwise it earns 100 euro.
- If you forget to send a gift to a responsive customer, you lose 100 euro
- Given data of past customers in a given table, decide which group of people you should target in the future

Nom	Prénom	Sexe	Profession	Réponse
Martin	Jeanne	F	Cadre	ok
Berluchette	Huguette	F	Ouvrière	ok
Sarkau	Sy	M	Ouvrier	non
Vil	Dominique	M	Cadre	non
Maitre	Kanter	M	Cadre	ok

Decision Tree Example:



- Decision Tree Example:



Mail only to executives or only female executives

■ Bayesian classifier:

$$P(c|x) = \frac{P(x|c)P(c)}{P(x)}$$

Diagram labels:

- Likelihood: $P(x|c)$
- Class Prior Probability: $P(c)$
- Posterior Probability: $P(c|x)$
- Predictor Prior Probability: $P(x)$

$$P(c|X) = P(x_1|c) \times P(x_2|c) \times \dots \times P(x_n|c) \times P(c)$$

- $P(c|x)$ is the posterior probability of *class* (c , *target*) given *predictor* (x , *attributes*).
- $P(c)$ is the prior probability of *class*.
- $P(x|c)$ is the likelihood which is the probability of *predictor* given *class*.
- $P(x)$ is the prior probability of *predictor*.

- **Bayesian classifier** example: Male or female?

Training set:

Sex	height (feet)	weight (lbs)	foot size(inches)
male	6	180	12
male	5.92 (5'11")	190	11
male	5.58 (5'7")	170	12
male	5.92 (5'11")	165	10
female	5	100	6
female	5.5 (5'6")	150	8
female	5.42 (5'5")	130	7
female	5.75 (5'9")	150	9

- **Bayesian classifier** example: Male or female?

Test case:

Sex	height (feet)	weight (lbs)	foot size(inches)
sample	6	130	8

- **Bayesian classifier** example: Male or female?

Sex	mean (height)	variance (height)	mean (weight)	variance (weight)	mean (foot size)	variance (foot size)
male	5.855	3.5033e-02	176.25	1.2292e+02	11.25	9.1667e-01
female	5.4175	9.7225e-02	132.5	5.5833e+02	7.5	1.6667e+00

Gaussian naive bayes:

$$p(x = v \mid c) = \frac{1}{\sqrt{2\pi\sigma_c^2}} e^{-\frac{(v-\mu_c)^2}{2\sigma_c^2}}$$

$$P(\text{male}) = 0.5$$

$$p(\text{height} \mid \text{male}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-(6 - \mu)^2}{2\sigma^2}\right) \approx 1.5789,$$

- **Bayesian classifier** example: Male or female?

$$P(\text{male}) = 0.5$$

$$p(\text{height} \mid \text{male}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-(6 - \mu)^2}{2\sigma^2}\right) \approx 1.5789,$$

$$p(\text{weight} \mid \text{male}) = 5.9881 \cdot 10^{-6}$$

$$p(\text{foot size} \mid \text{male}) = 1.3112 \cdot 10^{-3}$$

$$\text{posterior numerator (male)} = \text{their product} = 6.1984 \cdot 10^{-9}$$

$$P(\text{female}) = 0.5$$

$$p(\text{height} \mid \text{female}) = 2.2346 \cdot 10^{-1}$$

$$p(\text{weight} \mid \text{female}) = 1.6789 \cdot 10^{-2}$$

$$p(\text{foot size} \mid \text{female}) = 2.8669 \cdot 10^{-1}$$

$$\text{posterior numerator (female)} = \text{their product} = 5.3778 \cdot 10^{-4}$$

- **Bayesian classifier** example:

- Two classes:

▶ $c1 = 01100, 11001, 10110, 10101, 10010$
 ▶ $c2 = 01010, 11111, 11010, 11101, 10101$

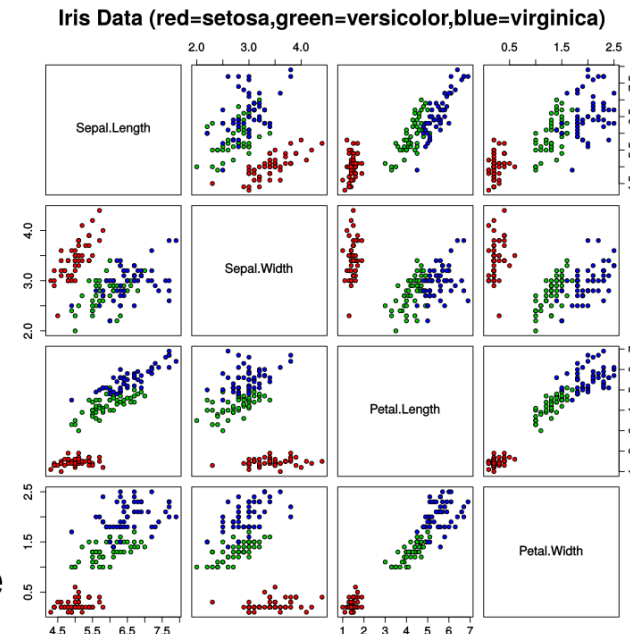
- Classify $x = 00111$

■ Bayesian classifier:

- Often used in text classification (i.e. spam)
- Works with little data, updates are possible

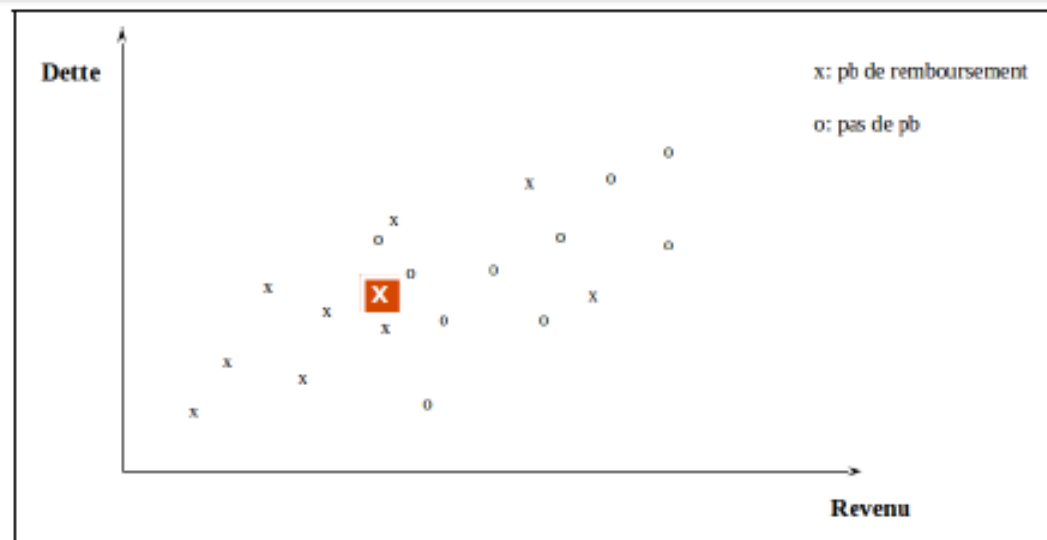


Use the four properties to predict the type
Of iris



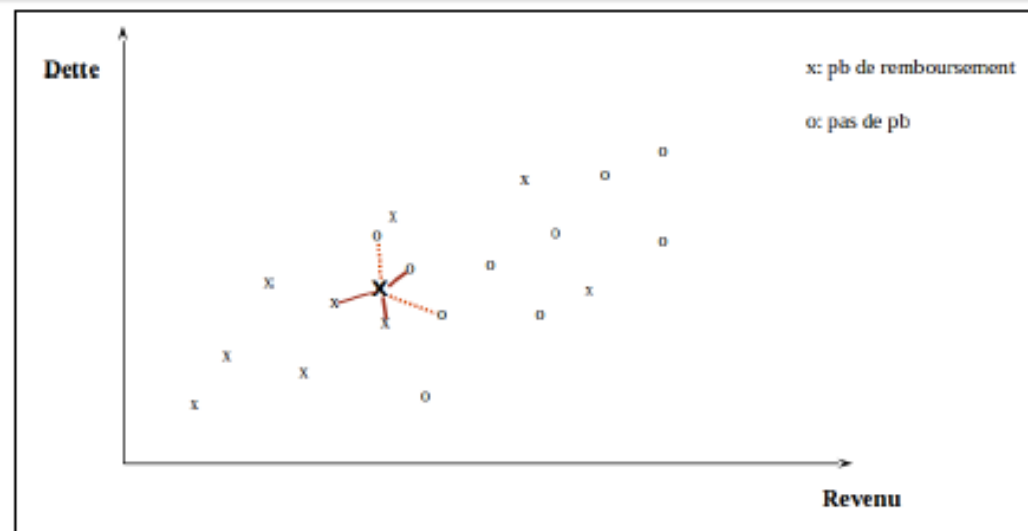
■ K-nearest neighbor:

- All distances between the point X to be classified and all labeled points is computed
- We keep the K closest labeled points. The Majority class in this set is attributed to X.



■ K-nearest neighbor:

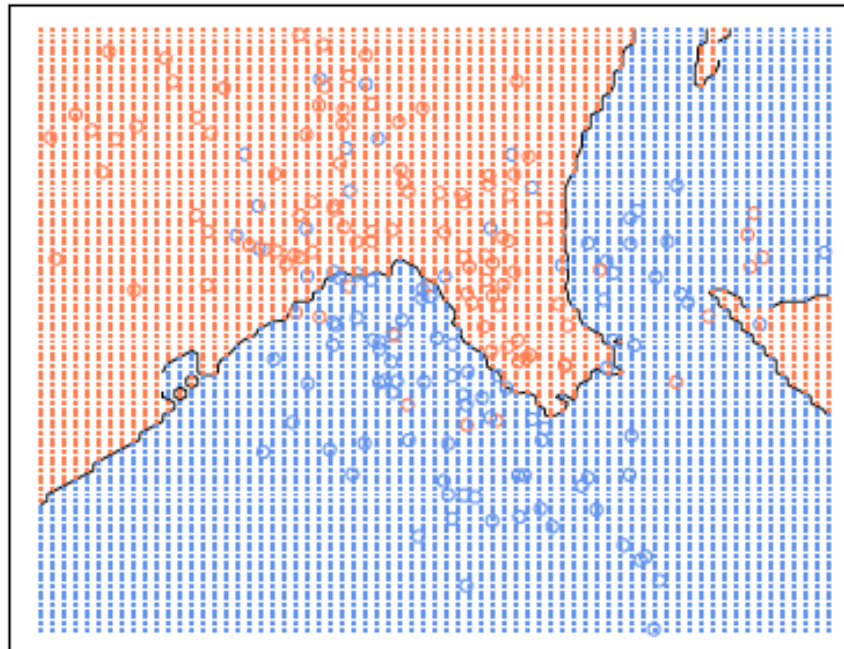
- All distances between the point X to be classified and all labeled points is computed
- We keep the K closest labeled points. The Majority class in this set is attributed to X.



$K = 3 \rightarrow x, K = 5 \rightarrow o$

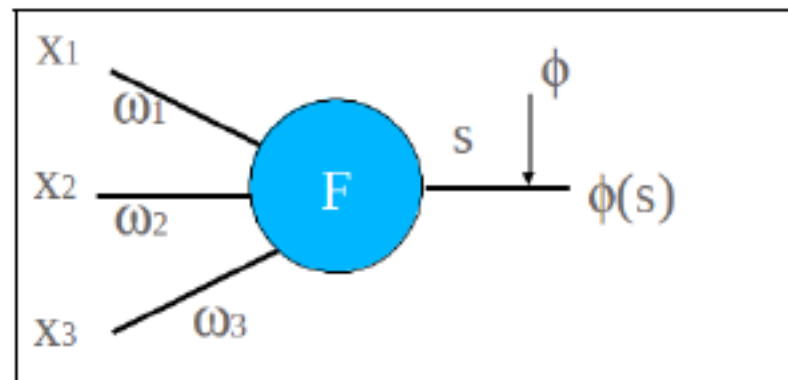
- **K-nearest neighbor:**
 - No hypothesis on the distribution of classes
 - Complexity increasing with the size of the training base

15-nearest neighbour



■ Neural networks:

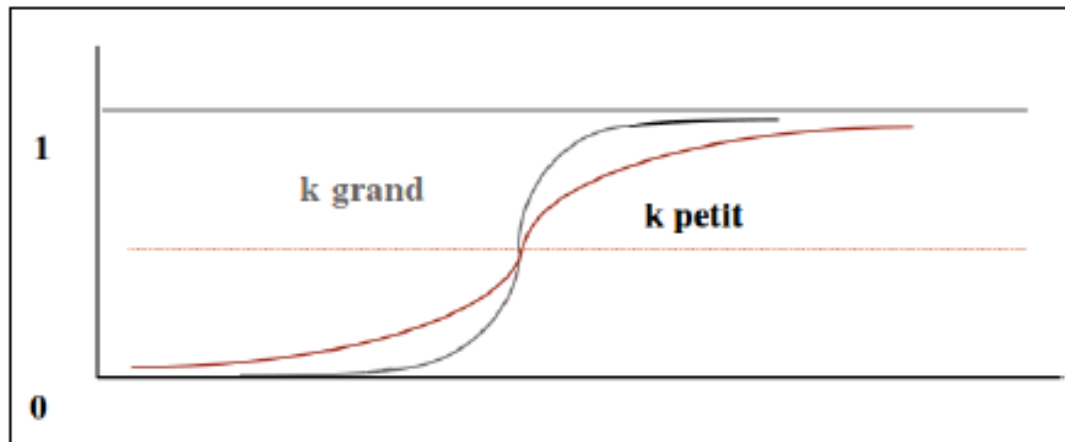
- Inspired by the structure of the nervous system
- A large number of connected neurons that process information
- The response of the neuron depends on its state and the weights of the connections
- The weights (or forces) are developed by experiment



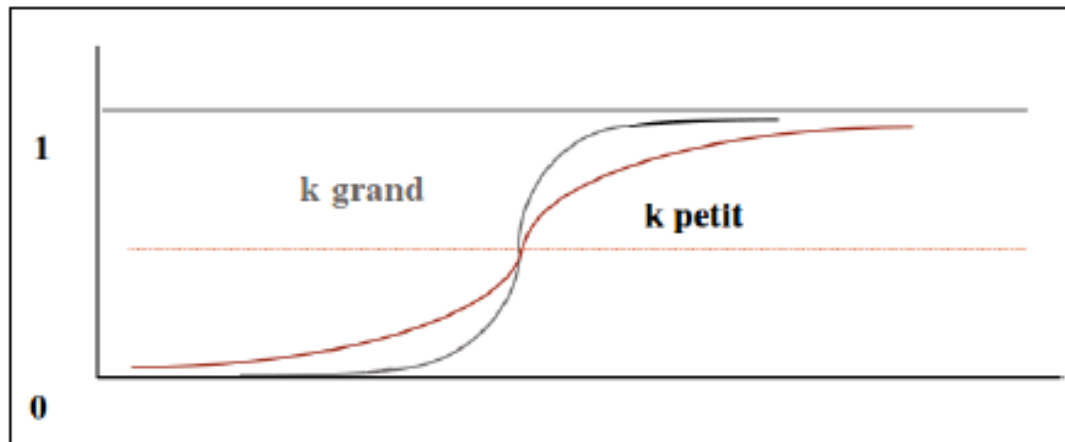
- **Neural networks** principle:
 - Construction of a network of simple computational units (neurons) linked by connections
 - Learning network parameters (weight of connections) using a set of examples
- Components of a neuron unite:
 - **inputs** (incoming connections or input variables)
 - **weights** on incoming connections
 - a **function F** which computes an output as a function of the inputs and the weight on the inlets
 - **activation function** which modifies the amplitude of the output of the node

■ Activation function:

- ▶ $\Phi(s) = \text{linéaire}$
- ▶ $\Phi(s) = \text{seuil}$
 - ▶ $\Phi(s) = 0$ si $s \leq a$
 - ▶ $\Phi(s) = 1$ si $s \geq a$
- ▶ $\Phi(s) = 1/(1 + e^{-ks})$

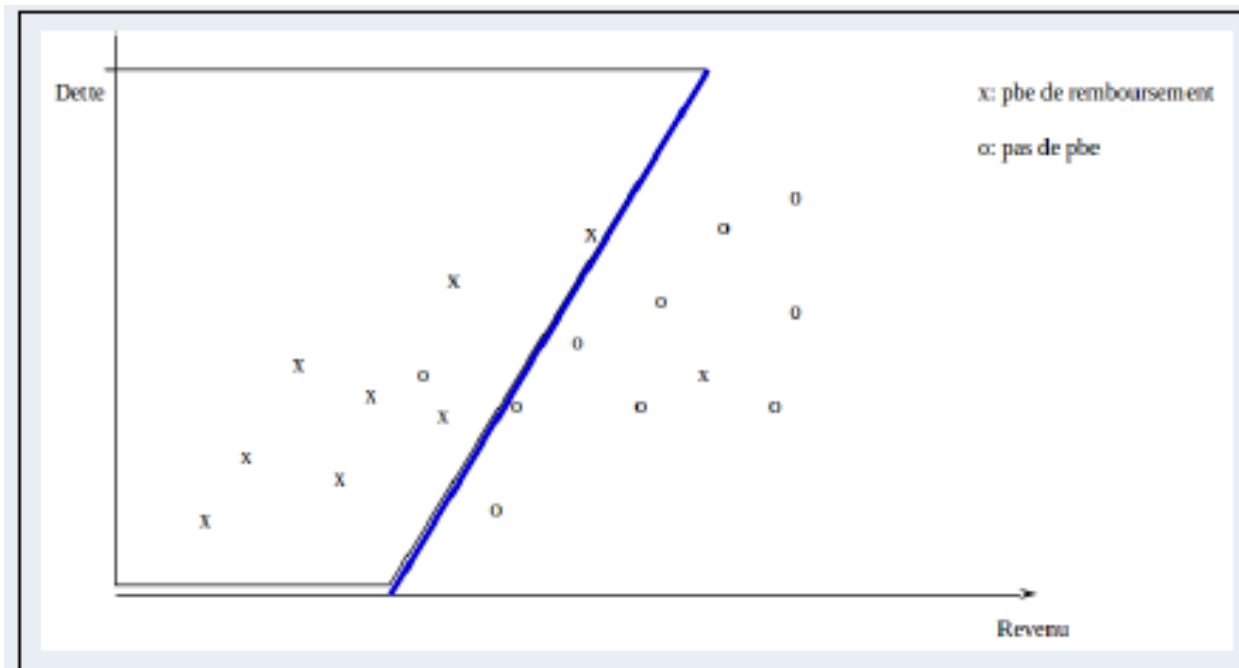


- **Activation function:** $\Phi(s) = 1/(1 + e^{-ks})$
 - If the coefficient k is large, then the output is almost always close to 0 or 1: relatively symbolic neural network
 - If the coefficient k of $1/(1 + e^{-ks})$ is small, then the strength of each cell is well distributed between 0 and 1
 - Another implicit parameter is the center of the sigmoid function

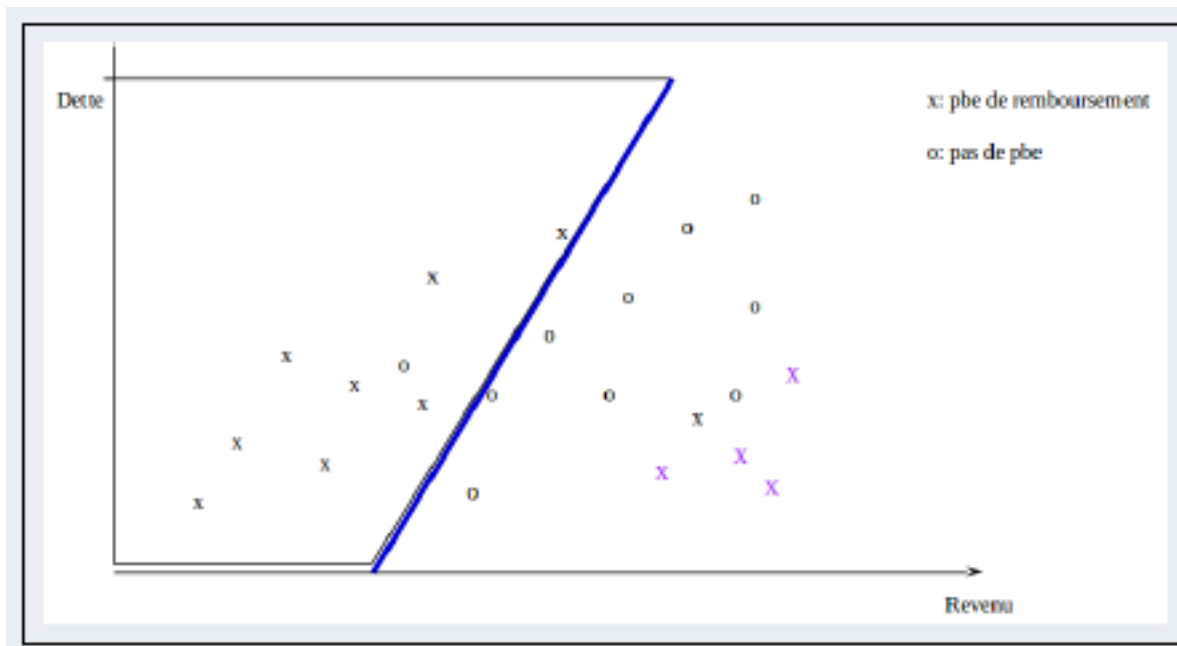


- **Simplest neural network:**
 - Single neuron
 - F = weighted sum of inputs
 - Activation function = thresholding
 - $s = 1$ if $w_1x_1 + w_2x_2 + \dots > a$
 - $s = 0$ if $w_1x_1 + w_2x_2 + \dots \leq a$
 - Equation of a hyperplane!

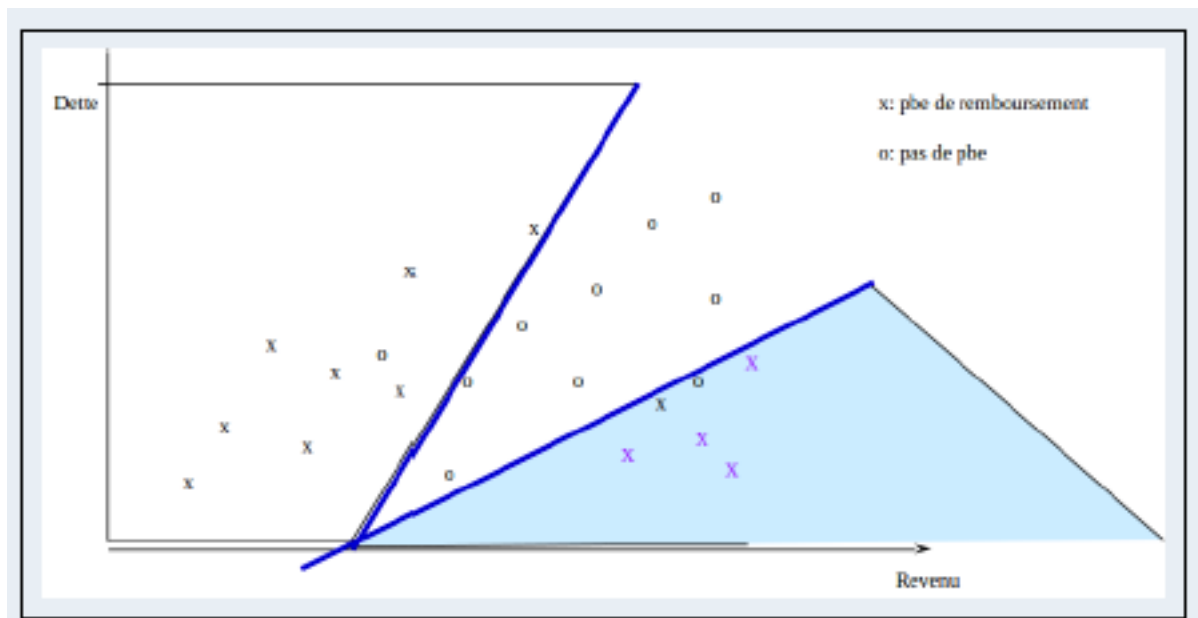
- **Simplest neural network:**
 - Linear separation



- **Simplest neural network:**
 - Additional examples
 - How to find a neuronal network discriminating between the two classes?



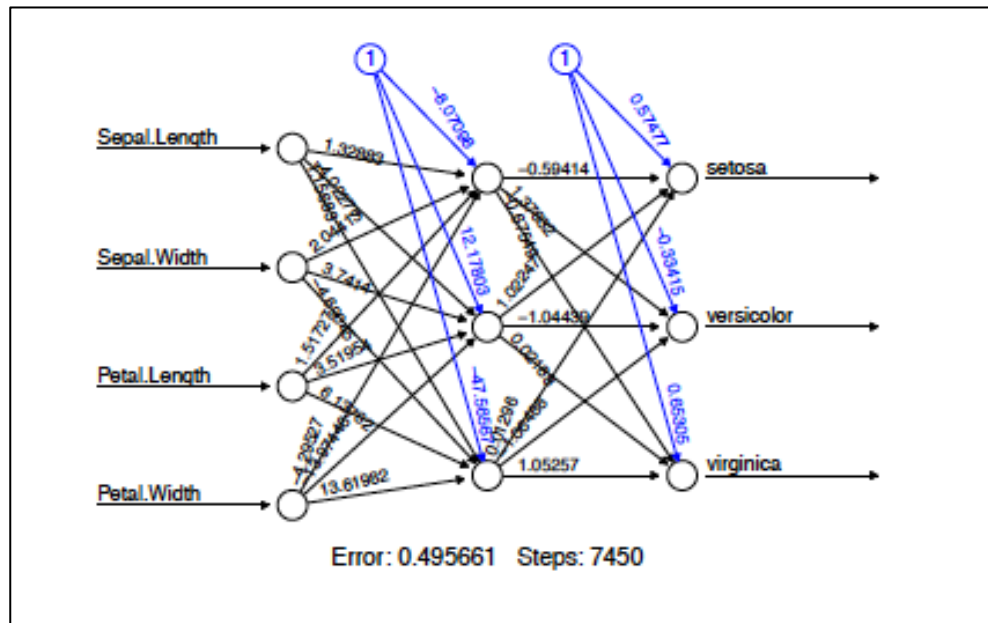
- **Simplest neural network:**
 - Additional examples



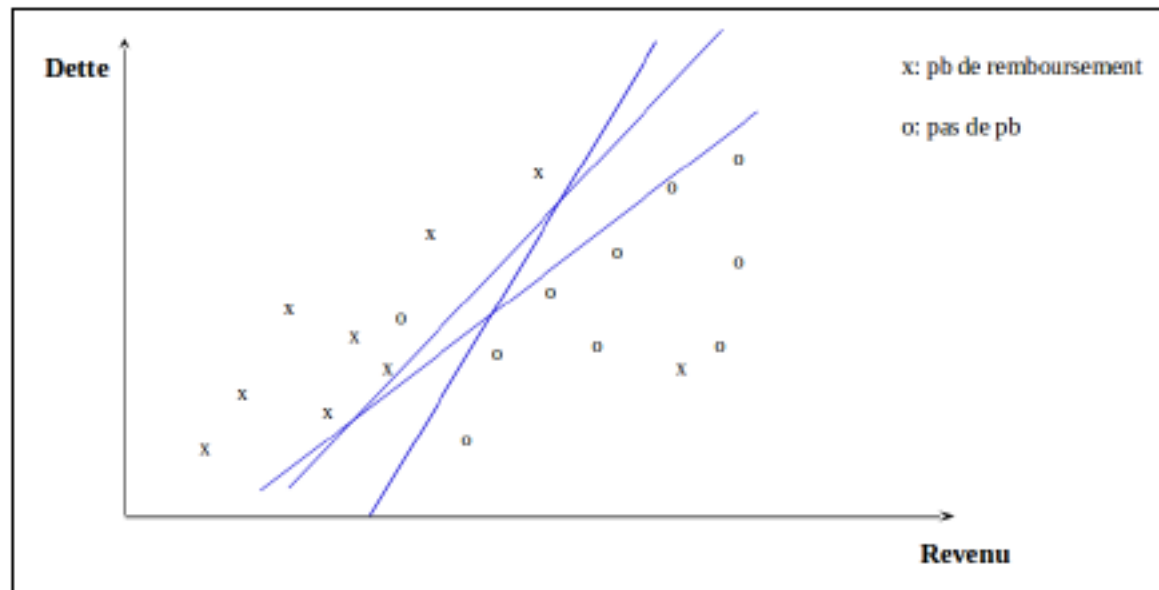
- NN In practice:
 - Choose a **calculation function** and an **activation function**
 - Choose an **architecture**:
 - Number of inputs
 - Number of outputs
 - Number of internal layers
 - Number of neurons of each of the internal layers
 - Select a **cost function**
 - Decide **when to stop** training

- **Advantages** of neural networks:
 - Robust to noise
 - Classification or estimation is quick after the training is completed
 - Available in all data mining software
- **Disadvantages:**
 - Black box: difficult to interpret the obtained model
 - Significant learning time
 - Selection of parameters is difficult

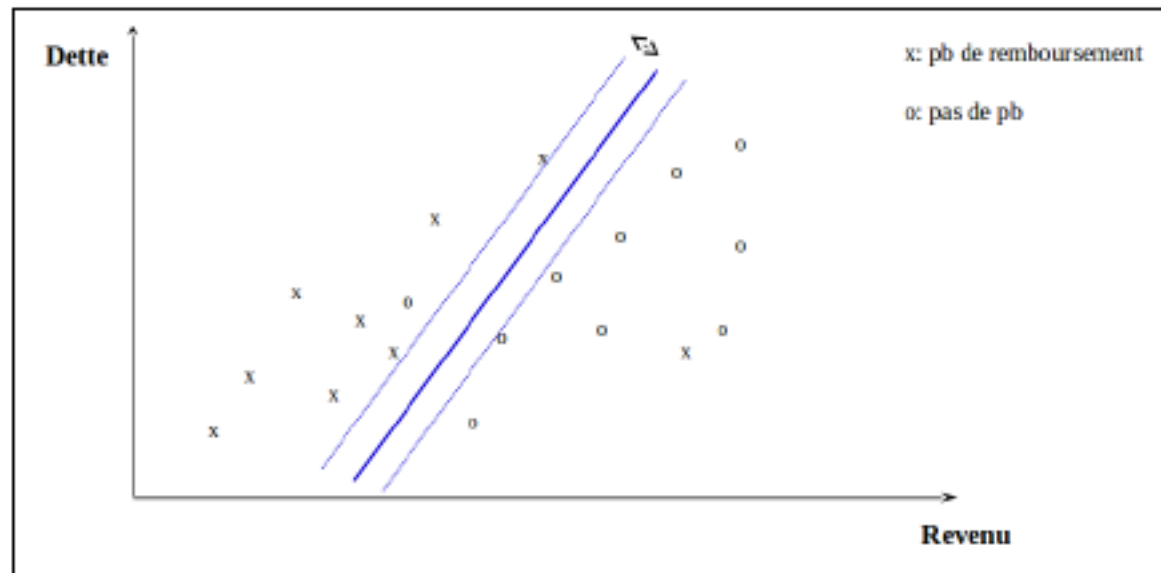
- Neural networks:
 - Renewed popularity with the arrival of Deep Learning
 - Frequently used in image analysis (Convolution Neural Network)



- **Support Vector Machine (SVM):**
 - Divide the data into two classes using a hyperplane
 - Maximize the gap between this hyperplane and the data
 - Find the line that maximizes the gap among several possible



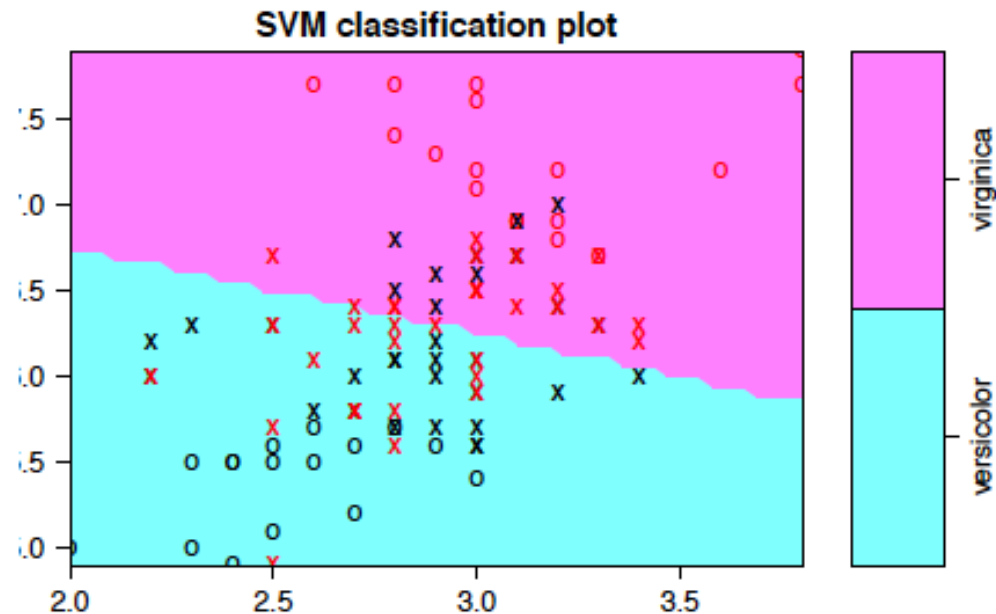
- **Support Vector Machine (SVM):**
 - Divide the data into two classes using a hyperplane
 - Maximize the gap between this hyperplane and the data
 - Find the line that maximizes the gap among several possible



■ Support Vector Machine (SVM):

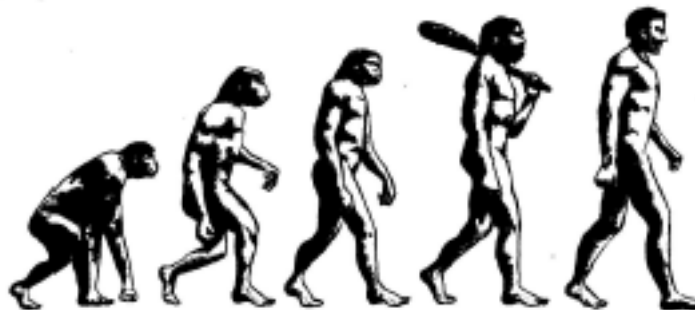
- Need to find a kernel suited to transform the data

- Example of the Gaussian kernel: $K(x, y) = \exp\left(-\frac{\|x-y\|^2}{2\sigma^2}\right)$



■ Genetic algorithms:

- Inspired by the theories of the evolution of Darwin, Lamarck or Baldwin
- General optimization method
- Can be used in classification or estimation



■ Genetic algorithms:

- General scheme:
 - We define the “parameters” to be optimized: range of values, thresholds, etc. The corresponding genotype (chromosomes) is defined.
 - We define the function of computation of the phenotype and the function of an individual
 - The mechanisms and rates of crossing and mutation are defined
 - The function of electing survivors is defined

- **Genetic algorithms:**
 - General scheme:
 1. Initialize the population
 2. Calculate the degree of adaptation $f(x)$ of each individual
 3. As long as not finished or no convergence:
 - a) Reproduction of parents
 - select 2 individuals at a time
 - apply genetic operators
 - b) calculate the degree of adaptation $f(x)$ of each child
 - c) select survivors from parents and children

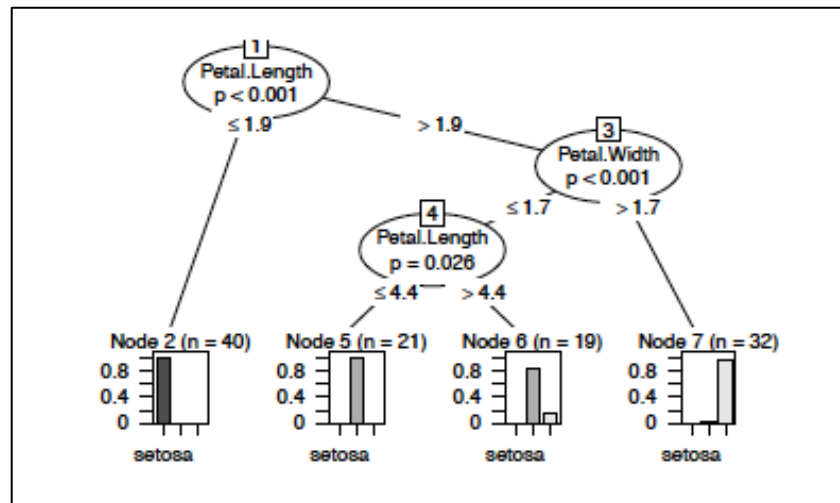
■ Validation:

- Validation by test:
 - Data = learning set + test set
 - Construction of a model on the learning set and test model on the test set for which the results are known
 - Cross-validation



■ Validation:

- Validation by test:
 - **Split test / train data**, in general more data for learning
 - The number of cross-validation folds depends on the volume of the data



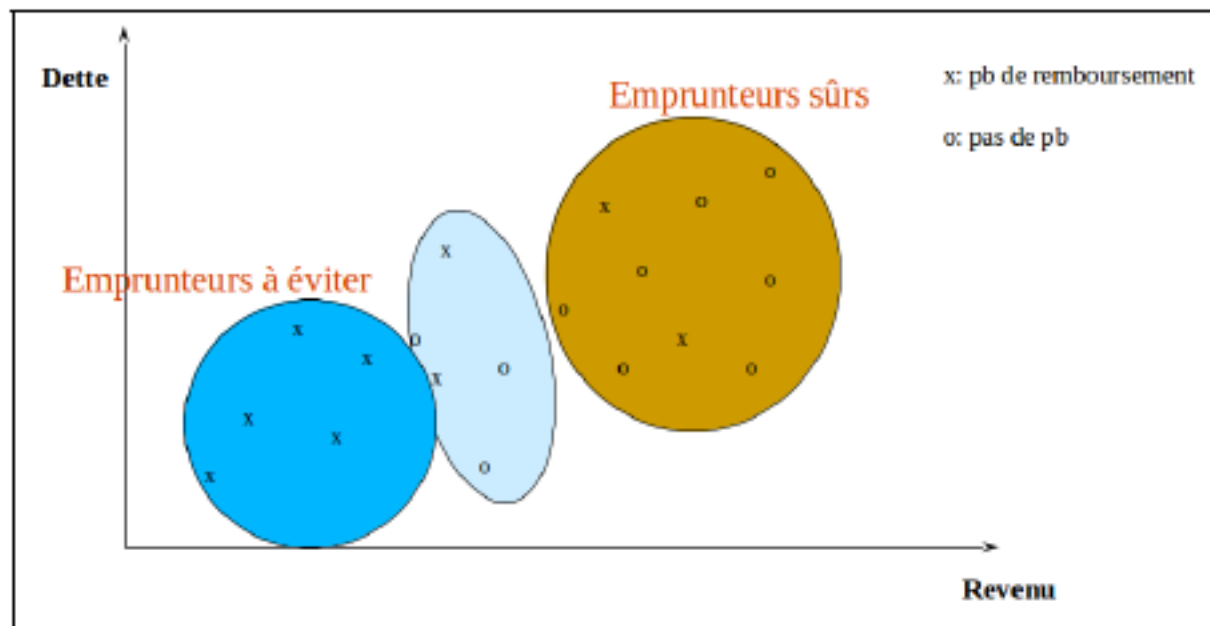
■ Clustering:

- Aim of the clustering: obtain a simplified representation (structuring) of the initial data
- Organization of a set of objects into a set of homogeneous and / or natural groupings

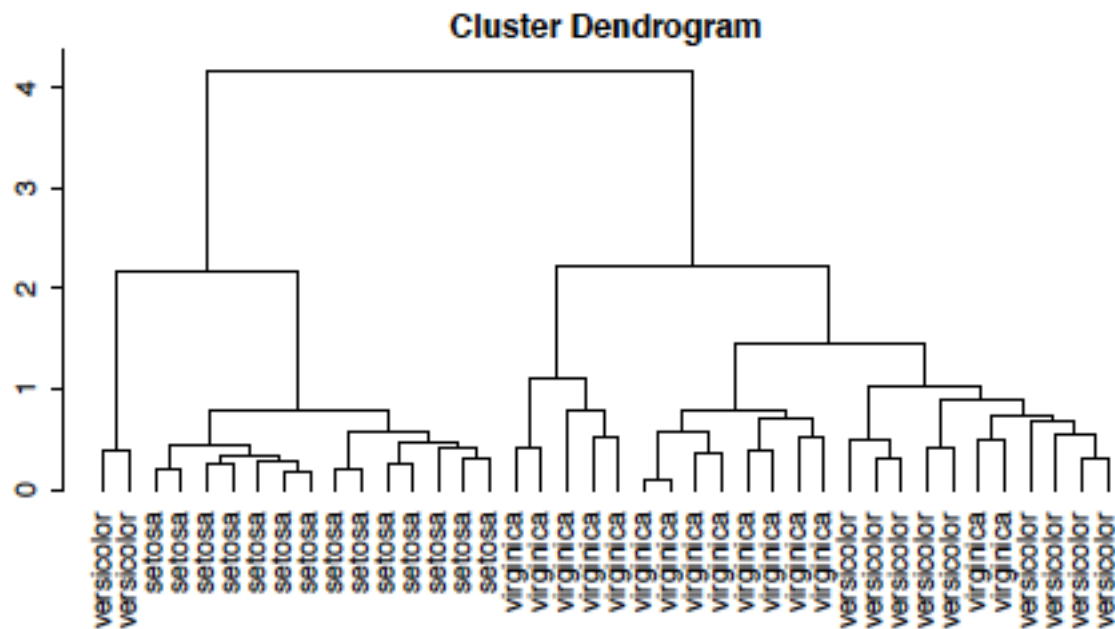


■ Clustering:

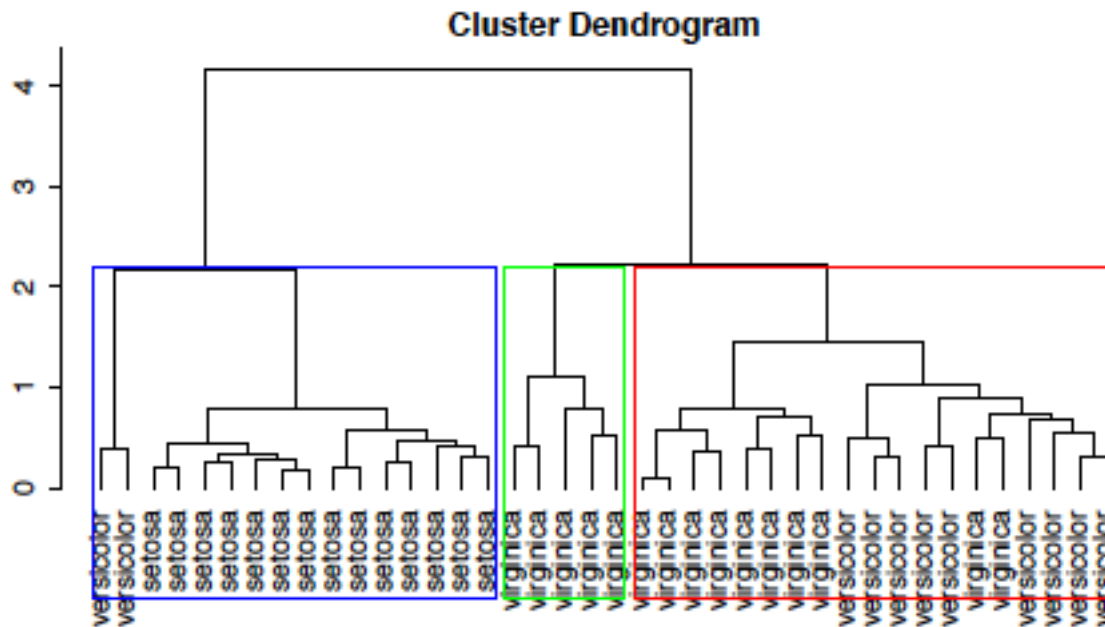
- Automatic partitioning from data
- No a priori semantic interpretation



- **Clustering:**
 - Hierarchical clustering



- **Clustering:**
 - Hierarchical clustering



■ Associations:

- Association rules: analysis of the shopping basket
- “On fridays, customers often buy beer packs and at the same time diapers”
- Are there any causal links between the purchase of a product P0 and a other product P1?

[illegible]

Carrefour market 

DIXIE FOURNE LES MARCHÉS
1^{er} - CHARENTAIS
LUNDI AU SAMEDI DE 10 h à 19h30
Charente
8 Avenue du Carrefour Market

DESCRIPTION	QST	MONTANT
SAISON		1,50
DES FICHES MARCHÉS		2,50
TOTAL SAISONNIER		3,98
LES MARCHÉS CROUS		3,75
TOTAL SAISONNIER		3,75
Donc TOTAL		7,73
3 MARCHÉS	TOTAL A PAYER	7,73

EXPRESS **EUR** **10,30K**
Nette Remise **- 2,58K**

000 000 00000 00000000 10000000

Carrefour MARKET ET FOURNE MARCHÉS
Pour l'achat de vos produits
à 0,90€

Géant
Création
CONCEPT, MAINTIEN, REPARATION
TOUTES LES SÉRIES
115 125 135 145 155 165 175 185 195 205 215 225 235 245 255 265 275 285 295 305 315 325 335 345 355 365 375 385 395 405 415 425 435 445 455 465 475 485 495 505 515 525 535 545 555 565 575 585 595 605 615 625 635 645 655 665 675 685 695 705 715 725 735 745 755 765 775 785 795 805 815 825 835 845 855 865 875 885 895 905 915 925 935 945 955 965 975 985 995 1005 1015 1025 1035 1045 1055 1065 1075 1085 1095 1105 1115 1125 1135 1145 1155 1165 1175 1185 1195 1205 1215 1225 1235 1245 1255 1265 1275 1285 1295 1305 1315 1325 1335 1345 1355 1365 1375 1385 1395 1405 1415 1425 1435 1445 1455 1465 1475 1485 1495 1505 1515 1525 1535 1545 1555 1565 1575 1585 1595 1605 1615 1625 1635 1645 1655 1665 1675 1685 1695 1705 1715 1725 1735 1745 1755 1765 1775 1785 1795 1805 1815 1825 1835 1845 1855 1865 1875 1885 1895 1905 1915 1925 1935 1945 1955 1965 1975 1985 1995 2005 2015 2025 2035 2045 2055 2065 2075 2085 2095 2105 2115 2125 2135 2145 2155 2165 2175 2185 2195 2205 2215 2225 2235 2245 2255 2265 2275 2285 2295 2305 2315 2325 2335 2345 2355 2365 2375 2385 2395 2405 2415 2425 2435 2445 2455 2465 2475 2485 2495 2505 2515 2525 2535 2545 2555 2565 2575 2585 2595 2605 2615 2625 2635 2645 2655 2665 2675 2685 2695 2705 2715 2725 2735 2745 2755 2765 2775 2785 2795 2805 2815 2825 2835 2845 2855 2865 2875 2885 2895 2905 2915 2925 2935 2945 2955 2965 2975 2985 2995 3005 3015 3025 3035 3045 3055 3065 3075 3085 3095 3105 3115 3125 3135 3145 3155 3165 3175 3185 3195 3205 3215 3225 3235 3245 3255 3265 3275 3285 3295 3305 3315 3325 3335 3345 3355 3365 3375 3385 3395 3405 3415 3425 3435 3445 3455 3465 3475 3485 3495 3505 3515 3525 3535 3545 3555 3565 3575 3585 3595 3605 3615 3625 3635 3645 3655 3665 3675 3685 3695 3705 3715 3725 3735 3745 3755 3765 3775 3785 3795 3805 3815 3825 3835 3845 3855 3865 3875 3885 3895 3905 3915 3925 3935 3945 3955 3965 3975 3985 3995 4005 4015 4025 4035 4045 4055 4065 4075 4085 4095 4105 4115 4125 4135 4145 4155 4165 4175 4185 4195 4205 4215 4225 4235 4245 4255 4265 4275 4285 4295 4305 4315 4325 4335 4345 4355 4365 4375 4385 4395 4405 4415 4425 4435 4445 4455 4465 4475 4485 4495 4505 4515 4525 4535 4545 4555 4565 4575 4585 4595 4605 4615 4625 4635 4645 4655 4665 4675 4685 4695 4705 4715 4725 4735 4745 4755 4765 4775 4785 4795 4805 4815 4825 4835 4845 4855 4865 4875 4885 4895 4905 4915 4925 4935 4945 4955 4965 4975 4985 4995 5005 5015 5025 5035 5045 5055 5065 5075 5085 5095 5105 5115 5125 5135 5145 5155 5165 5175 5185 5195 5205 5215 5225 5235 5245 5255 5265 5275 5285 5295 5305 5315 5325 5335 5345 5355 5365 5375 5385 5395 5405 5415 5425 5435 5445 5455 5465 5475 5485 5495 5505 5515 5525 5535 5545 5555 5565 5575 5585 5595 5605 5615 5625 5635 5645 5655 5665 5675 5685 5695 5705 5715 5725 5735 5745 5755 5765 5775 5785 5795 5805 5815 5825 5835 5845 5855 5865 5875 5885 5895 5905 5915 5925 5935 5945 5955 5965 5975 5985 5995 6005 6015 6025 6035 6045 6055 6065 6075 6085 6095 6105 6115 6125 6135 6145 6155 6165 6175 6185 6195 6205 6215 6225 6235 6245 6255 6265 6275 6285 6295 6305 6315 6325 6335 6345 6355 6365 6375 6385 6395 6405 6415 6425 6435 6445 6455 6465 6475 6485 6495 6505 6515 6525 6535 6545 6555 6565 6575 6585 6595 6605 6615 6625 6635 6645 6655 6665 6675 6685 6695 6705 6715 6725 6735 6745 6755 6765 6775 6785 6795 6805 6815 6825 6835 6845 6855 6865 6875 6885 6895 6905 6915 6925 6935 6945 6955 6965 6975 6985 6995 7005 7015 7025 7035 7045 7055 7065 7075 7085 7095 7105 7115 7125 7135 7145 7155 7165 7175 7185 7195 7205 7215 7225 7235 7245 7255 7265 7275 7285 7295 7305 7315 7325 7335 7345 7355 7365 7375 7385 7395 7405 7415 7425 7435 7445 7455 7465 7475 7485 7495 7505 7515 7525 7535 7545 7555 7565 7575 7585 7595 7605 7615 7625 7635 7645 7655 7665 7675 7685 7695 7705 7715 7725 7735 7745 7755 7765 7775 7785 7795 7805 7815 7825 7835 7845 7855 7865 7875 7885 7895 7905 7915 7925 7935 7945 7955 7965 7975 7985 7995 8005 8015 8025 8035 8045 8055 8065 8075 8085 8095 8105 8115 8125 8135 8145 8155 8165 8175 8185 8195 8205 8215 8225 8235 8245 8255 8265 8275 8285 8295 8305 8315 8325 8335 8345 8355 8365 8375 8385 83

■ Associations:

- Rules of association: premise -> conclusion
- Questions:
 - beurre -> pain ?
 - poisson, viande -> lait ?
 - fromage, pâtes -> vin ?

Identifiant	Transaction
1	beurre fruits lait pain
2	fruits lait pain
3	beurre fromage pain pâtes viande vin
4	fromage fruits lait légumes pain pâtes poisson
5	beurre fruits lait légumes pain pâtes poisson viande
6	beurre fromage légumes pain pâtes viande vin
7	beurre fromage lait légumes pain pâtes viande vin
8	fruits légumes poisson
9	beurre fromage lait pain pâtes viande vin
10	beurre fromage fruits lait légumes pain poisson viande

■ Associations:

- Formally:
 - Given a set of transactions D , find all the association rules $X \rightarrow Y$ having support and confidence above the minimum thresholds predicted by the user
 - A transaction is a set of attributes (butter, fruit, milk, bread)
 - **Support**: % of transactions in D that contain X and Y
 - **Confidence**: % of transactions that contain X which also contain Y

■ **Associations:**

- Interpretation:
 - $R : X \rightarrow Y (A\%, B\%)$
 - A% of all transactions show that X and Y have been purchased at the same time (support of the rule) and B% of clients who purchased X have also purchased Y (confidence in the rule).

■ Associations:

- Two sub problems:
 - FIS: Find all frequent ensembles (**item sets**) that have support greater than or equal to a minimum value “**minsup**”
 - Generate all the association rules having confidence greater or equal to “**minconf**”

$$\text{support}_{A \Rightarrow B} = \frac{|\{t: A \cup B \subseteq t\}|}{|T|} \quad \text{confidence}_{A \Rightarrow B} = \frac{|\{t: A \cup B \subseteq t\}|}{|\{t: A \subseteq t\}|}$$

■ Associations:

Tickets	Produits achetés
1	beurre fruits lait pain
2	fruits lait pain
3	beurre fromage pain pâtes viande vin
4	fromage fruits lait légumes pain pâtes poisson
5	beurre fruits lait légumes pain pâtes poisson viande
6	beurre fromage légumes pain pâtes viande vin
7	beurre fromage lait légumes pain pâtes viande vin
8	fruits légumes poisson
9	beurre fromage lait pain pâtes viande vin
10	beurre fromage fruits lait légumes pain poisson viande

Support =	$\frac{\text{green box et blue box}}{\text{Tous tickets}}$
Confiance =	$\frac{\text{green box et blue box}}{\text{green box} + \text{red box}}$

beurre → pain	Support	Confiance
poisson viande → lait	70%	100%
fromage pâtes → vin	20%	100%
	40%	80%

■ Associations:

Tickets	Produits achetés
1	beurre fruits lait pain
2	fruits lait pain
3	beurre fromage pain pâtes viande vin
4	fromage fruits lait légumes pain pâtes poisson vin
5	beurre fruits lait légumes pain pâtes poisson viande
6	beurre fromage légumes pain pâtes viande vin
7	beurre fromage lait légumes pain pâtes viande vin
8	fruits légumes poisson
9	beurre fromage lait pain pâtes viande vin
10	beurre fromage fruits lait légumes pain poisson viande

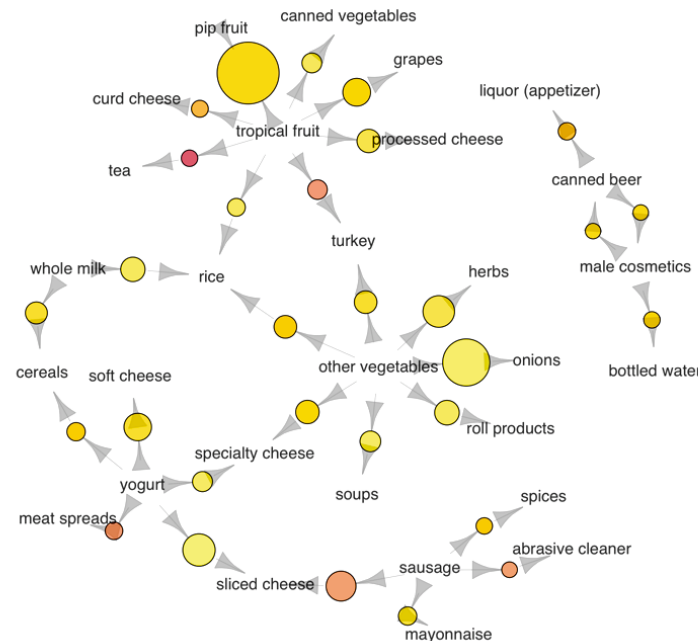
Support = $\frac{\text{fromage et pain}}{\text{Tous tickets}}$

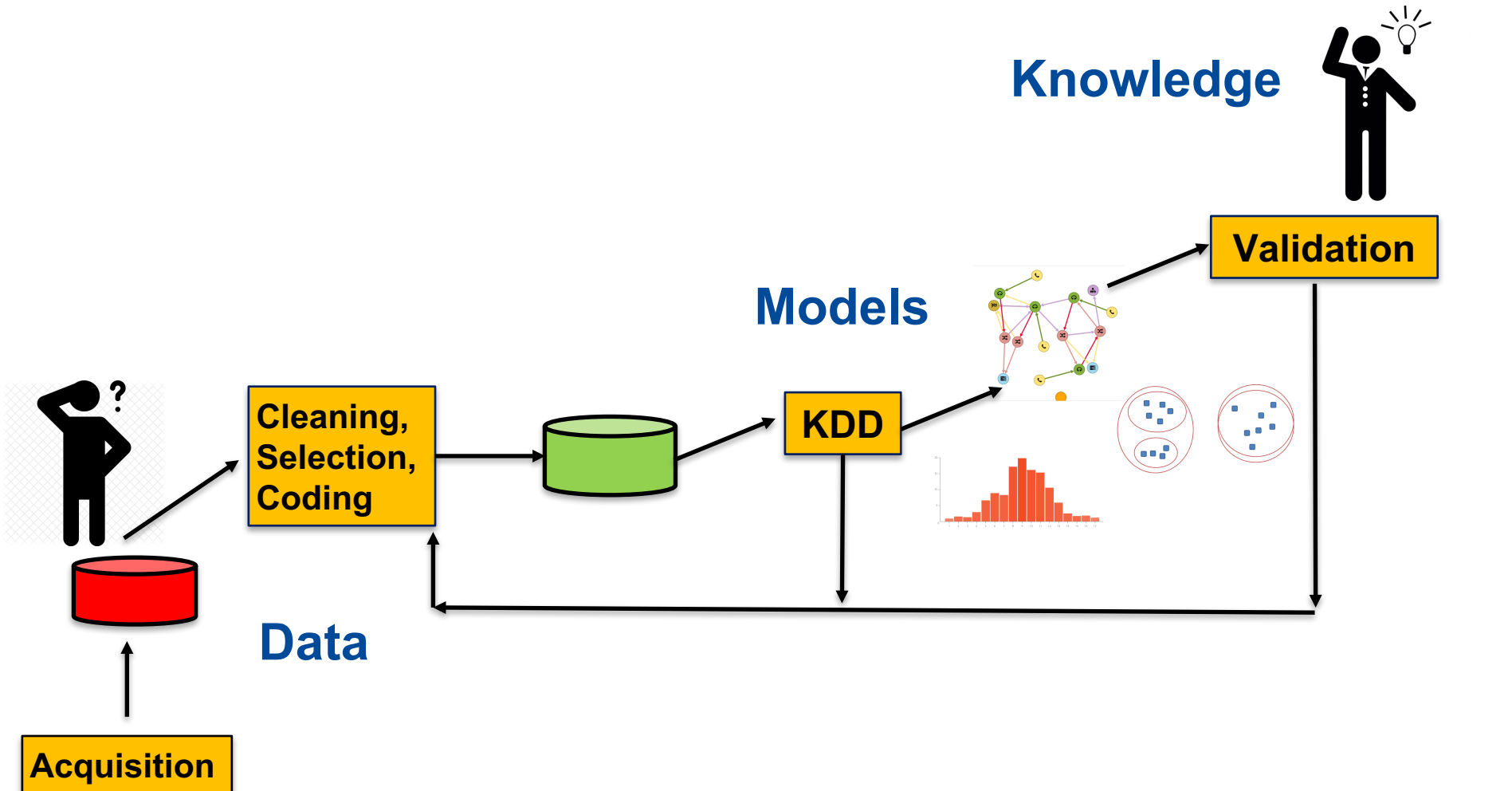
Confiance = $\frac{\text{fromage et pain}}{\text{fromage} + \text{pain}}$

		Support	Confiance
beurre	→ pain	70%	100%
poisson viande	→ lait	20%	100%
fromage pâtes	→ vin	40%	80%

■ Associations:

- Association rules:
 - Numerous criteria for evaluating the interest of a rule
 - Difficult to scale to large volumes of data





- Generation of a large number of models
- Is the generated model interesting?
- How to measure the interest of a model:
 - Novelty
 - Easy to understand
 - Can be validated by new data (with a certain confidence)
 - Usefulness
 - Confirms (or rejects) the hypothesis of an expert

- Evaluation of a model
 - subjective: expert
 - objective: statistics and structure of the model
- Can we find all models? (completeness)
- Can we only generate the interesting models? (optimization):
 - Generating all the models and iterating according to certain measures and characteristics: nonrealistic
 - Generate only the models that satisfy a particular condition

- Form groups of 4-5 people
- You will get a statement about data mining.
- What do you think? Is this statement correct or wrong? Why?
- You have 5 min for the task.

- “Data mining methods are **more inductive then methods based on hypothesis** because there is no a priori knowledge on the data”
- **False**: condition of application of the methods, choice of data, coding data, choice of explanatory variables, order of input of variables in the algorithm,...

- “**All data available** must be used to generate the model”
- **False**: coding of the data, order of input of the variables in the algorithm, irregular effects, outliers, influence of redundancies, correlations, the computer data model, saturation, instability ...

- “With all these methods we will **always find out something great!**”
- **False:** common sense solutions need to be found (specialists, experts). In fact, we need to find the best solution (among n) for difficult data

- “Data mining is **revolutionary!**”
- **False:** traditional data analysis methods + more specific methods (neural networks). Optimization of existing techniques because of the large amount of data.

- **Question:**
 - “Why so many algorithms?”
- **Answer:**
 - Because none is optimal in all cases
- As they are in practice complementary to one another, combining them intelligently (by constructing models) it is possible to achieve very significant performance gains