# Data Mining

# Motivation

# Motivation

- **Since the nineties, more and more data was stored because hard drivers became cheaper and with higher capacity**

- **Today: "90% of world's data generated over last two years"**

  - **http://www.sciencedaily.com/releases/2013/05/130522085217.htm (article from 22/5/2013)**

- **Question: how to extract useful information from all this data?**

# Motivation

**Examples:**

1. How to select customers with similar purchase profiles in order to adjust the promotion campaigns to these natural groups of customers?
   - It is necessary to know the history of purchases of these customers

2. How to suggest a new book to a client?
   - It is necessary to know the history of this and other shopping customers

3. How to detect breast cancer?
   - It is necessary to have measurable characteristics of the cyst/lump

# Data Mining is …

# Data Mining is ...

**Sources of data are increasing everyday**:
- Sensors
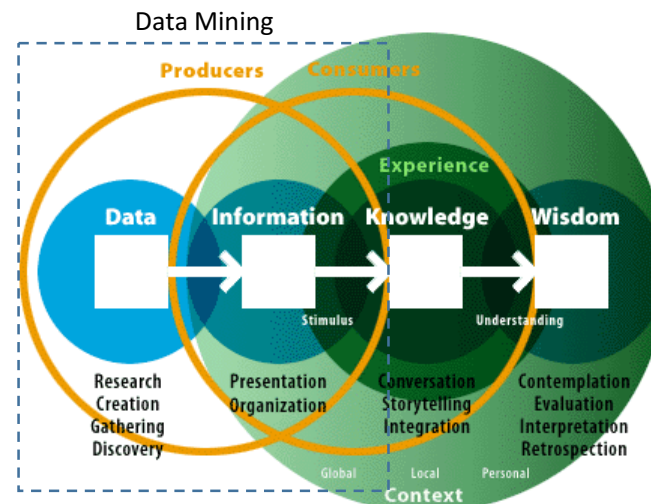- Internet, e.g. social networks
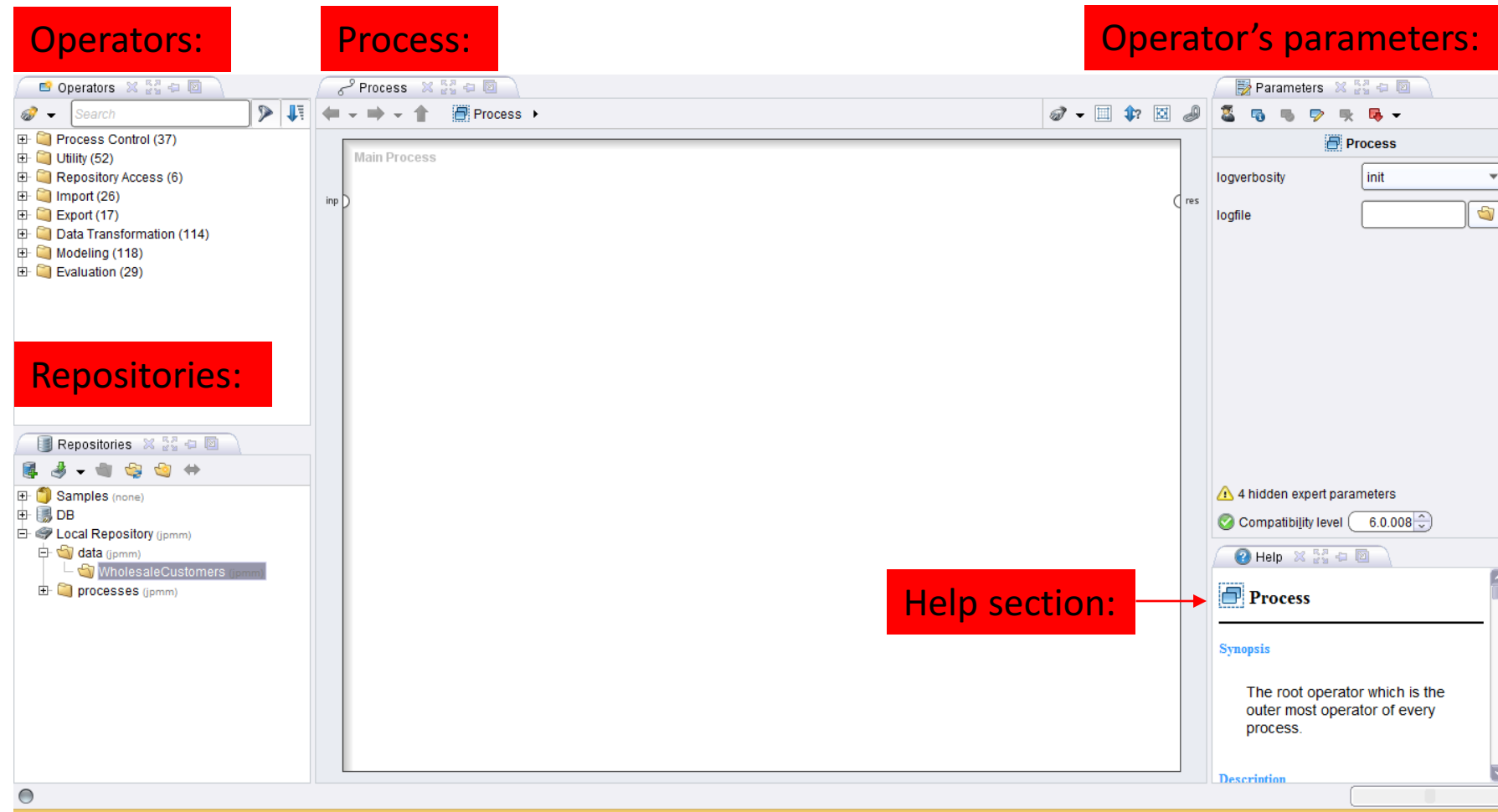- Data warehousing systems



**Extracting information from data**:
- Without doing assumptions about data distribution
- Discovering unknown information
- Using computational resources

Image obtained in 31-5-2014 from
http://www.nathan.com/thoughts/course.html

# Rapid Miner

Two simple examples

# Rapid Miner

# Rapid Miner

The explanation of the operator's inputs and outputs is given in the help section for the highlighted operator.

In the example:

- tra is the training set
- mod is the model
- exa is the example set

# Rapid Miner

1. How to select customers with similar purchase profiles in order to adjust the promotion campaigns to these natural groups of customers?

   - It is a clustering task
   - We are using the dataset *WholesaleCustomers*
     - Use only the quantitative variables

# Rapid Miner

# Rapid Miner

1. Loading data to the workspace

2. Insertion of the k-means operator

3. Definition of the parameters

4. Execution

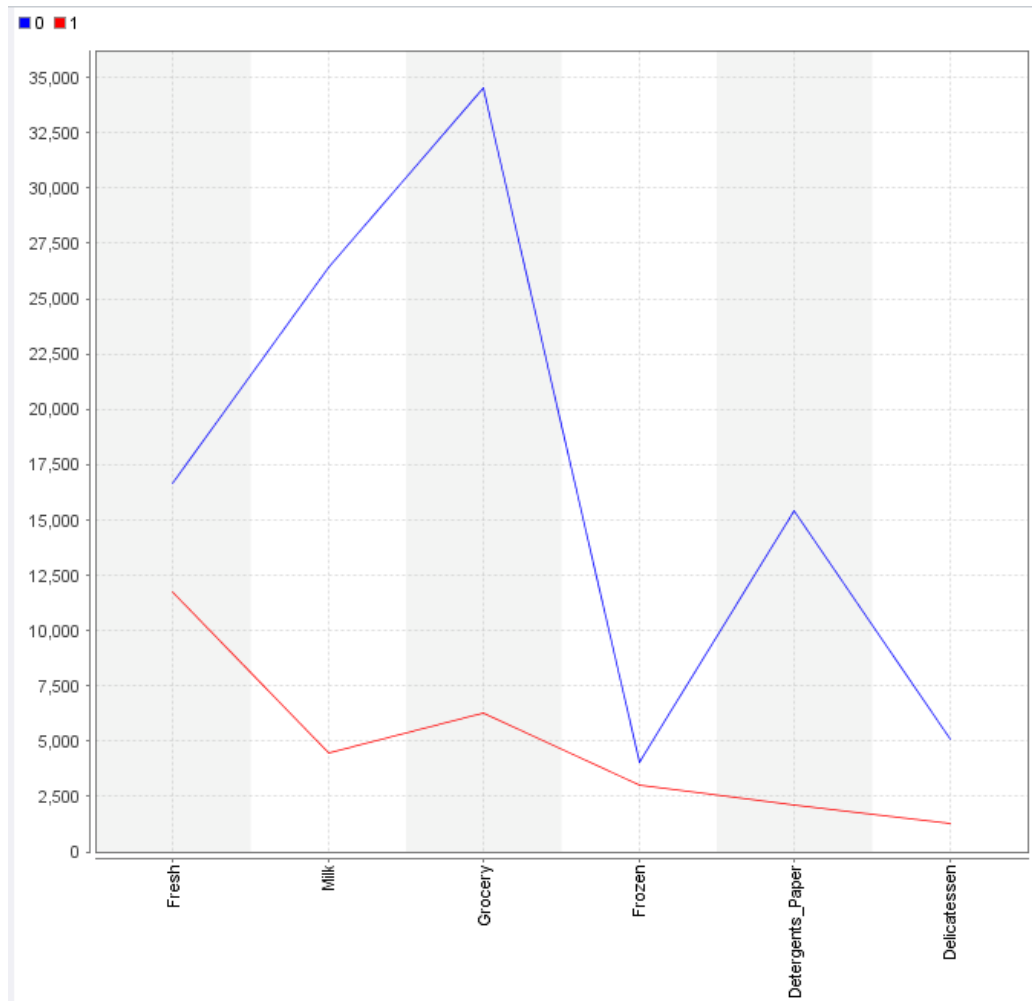5. Interpretation of the results

# Rapid Miner

# Rapid Miner

5



Analising and interpreting the results (see the ones in the image).
Save the process.

# Rapid Miner

2. How to determine the benignity/malignancy of a breast cyst?
   - Is a prediction task, more precisely, a classification task
   - Let's use the dataset BreastCancerWisconsin

# Rapid Miner

1. Loading data to the workspace
   o You should define the target variable
2. Insertion of the decision tree operator
3. Definition of the parameters
4. Execution
5. Interpretation of the results

# Rapid Miner

# Rapid Miner

# Rapid Miner

# Process model

# Process models

# The KDD process model

1. Selection
2. Preprocessing
3. Transformation
4. *Data Mining*
5. Interpretation / Evaluation



Fayyad, U. M. et al. 1996. From data mining to knowledge discovery: an overview. In Fayyad, U. M. et al (Eds.),  Advances in knowledge discovery and data mining. AAAI Press / The MIT Press.

# The KDD process model

- **Selection**: identification and selection of all external and internal sources of information and selection of the subset of data or variables needed for the KDD process.

- **Preprocessing**: includes the removal of data with extreme values (outliers), filling in missing values, etc.

- **Transformation**: converting data into a format suitable for Data Mining algorithms.

- **Data Mining**: in this step the specialized tools seek, through specialized algorithms, existing patterns in the data. This search can be performed automatically or interactively systems through the aid of the analyst responsible for the generation of the hypotheses. At the end of the process, the DM system should generate a report of the analysis carried out in order to enable analysts to verify the results obtained.
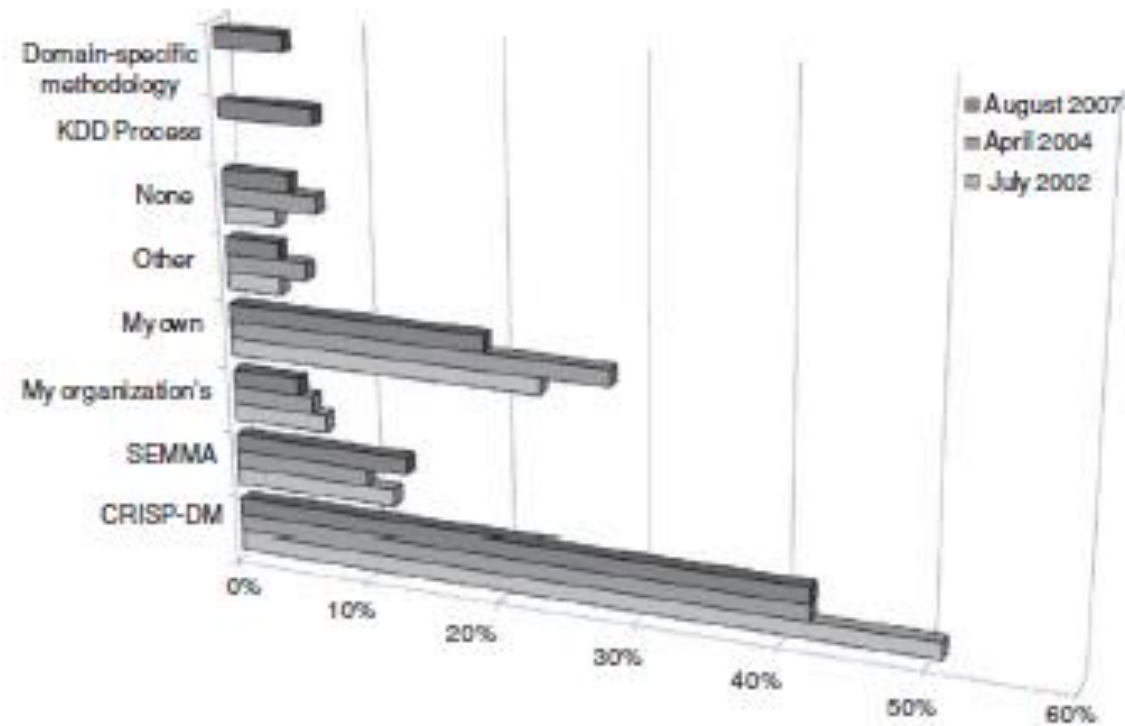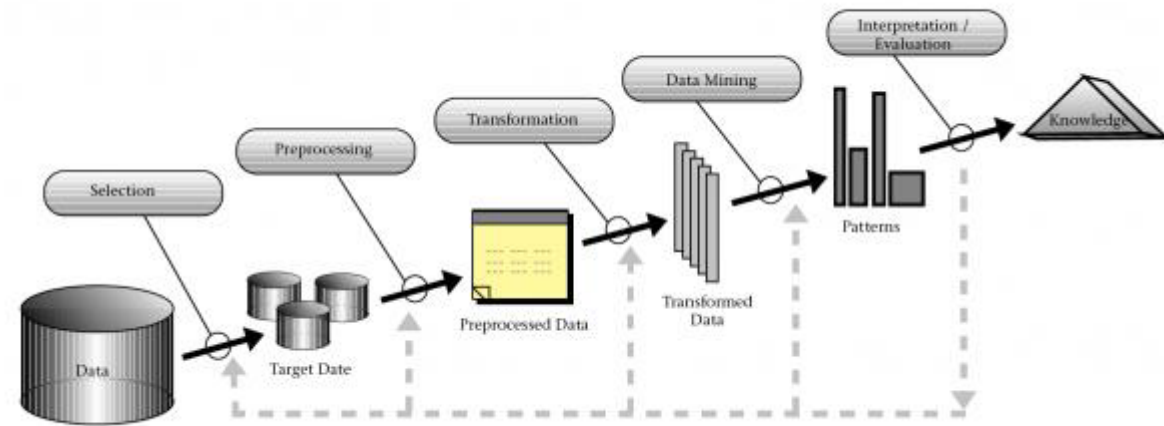
- **Interpretation/evaluation**: this step should be performed in conjunction with business analysts. If the knowledge generated is not satisfactory, analysts can form a new set of experiments giving rise to a new iteration of the process.

# The CRISP-DM process model

CRISP-DM (CRoss Industry Standard Process for Data Mining)

1. Business understanding
2. Data understanding
3. Data preparation
4. Modeling
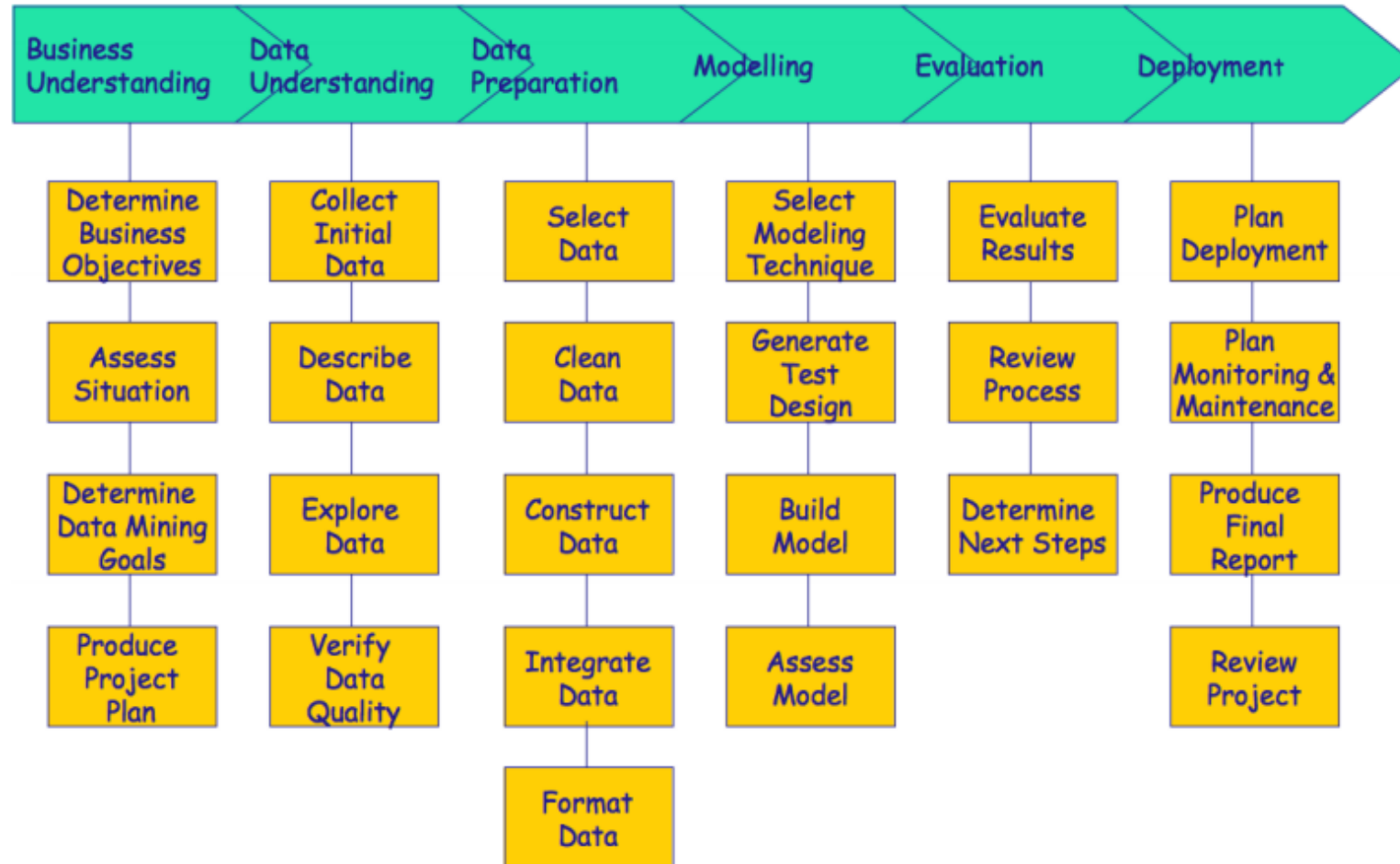5. Evaluation
6. Deployment

# The CRISP-DM process model

1. **Understand the Business**: focuses on understanding the objective of the project from a business perspective, defining a preliminary plan to achieve the goals.

2. **Understand the data**: data gathering and early activities to better understand the data, identifying problems or interesting sets.

3. **Data preparation**: construction of the final data set from the initial one. Normally occurs several times in the process.

4. **Modeling**: several modeling techniques are applied, and its parameters calibrated for optimization. Thus, it is common to return to Data Preparation during this phase.

5. **Evaluation**: a model that seems to have great quality in a data analysis perspective was constructed. However, it is necessary to verify if the model reaches the goals of the business.

6. **Deployment**: the knowledge acquired by the model is organized and presented in a way that the customer can use.

# The CRISP-DM process model

# Process models

| KDD | CRISP-DM |
|---|---|
| Pre-process KDD | Business understanding |
| Selection | Data understanding |
| Preprocessing | |
| Transformation | Data preparation |
| Data Mining | Modeling |
| Interpretation/Evaluation | Evaluation |
| Post-process KDD | Deployment |

# Process models