# Data Mining in Ship Construction & Operation

## A review of innovative methods and Open Software's
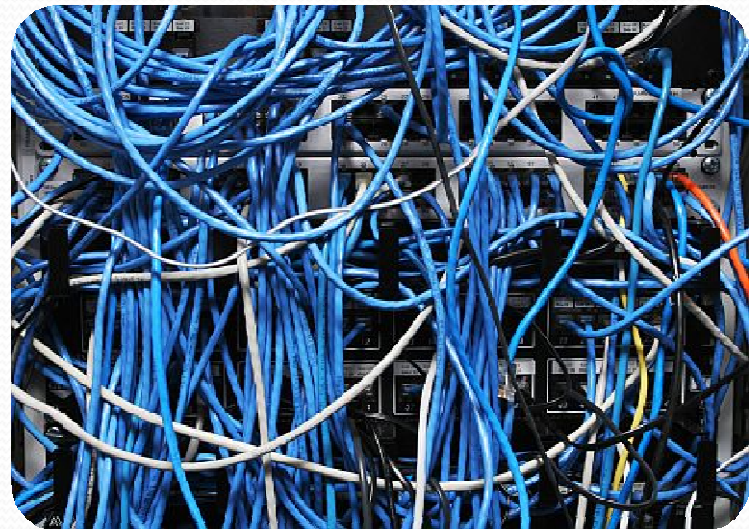


*Dr. Jean-David Caprace*
*Post PhD at UFRJ*

*January 2011*

# Summary
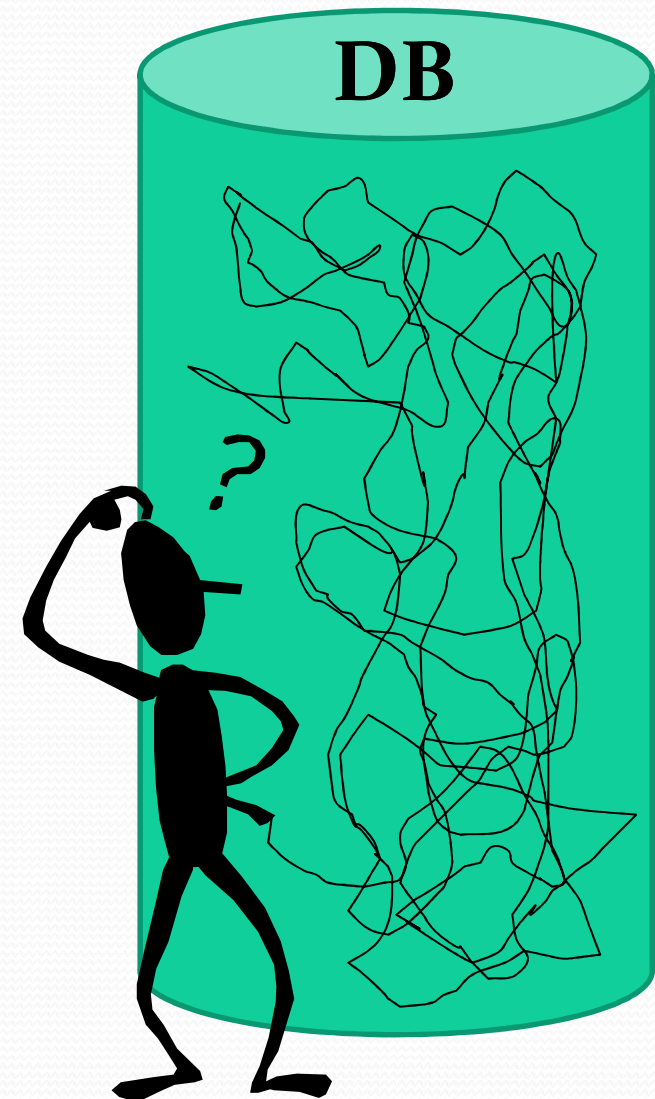## Outline of the presentation



- Data Mining
- Mathematical models
- Applications
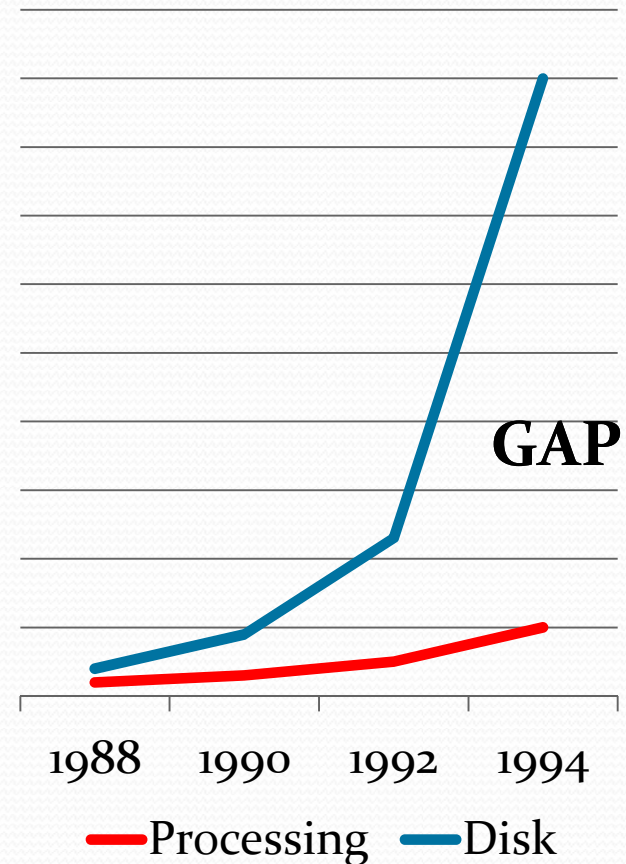- Open software's comparison

# Data Mining (DM)

## Why?

- Complexity of modern manufacturing processes
- Massive investment in automation & monitoring systems
  - ➔ Generation of large DBs
  - ➔ DBs underused
  - ➔ Human analysts take weeks to discover useful information

  ➔ **How to solve this problem?**
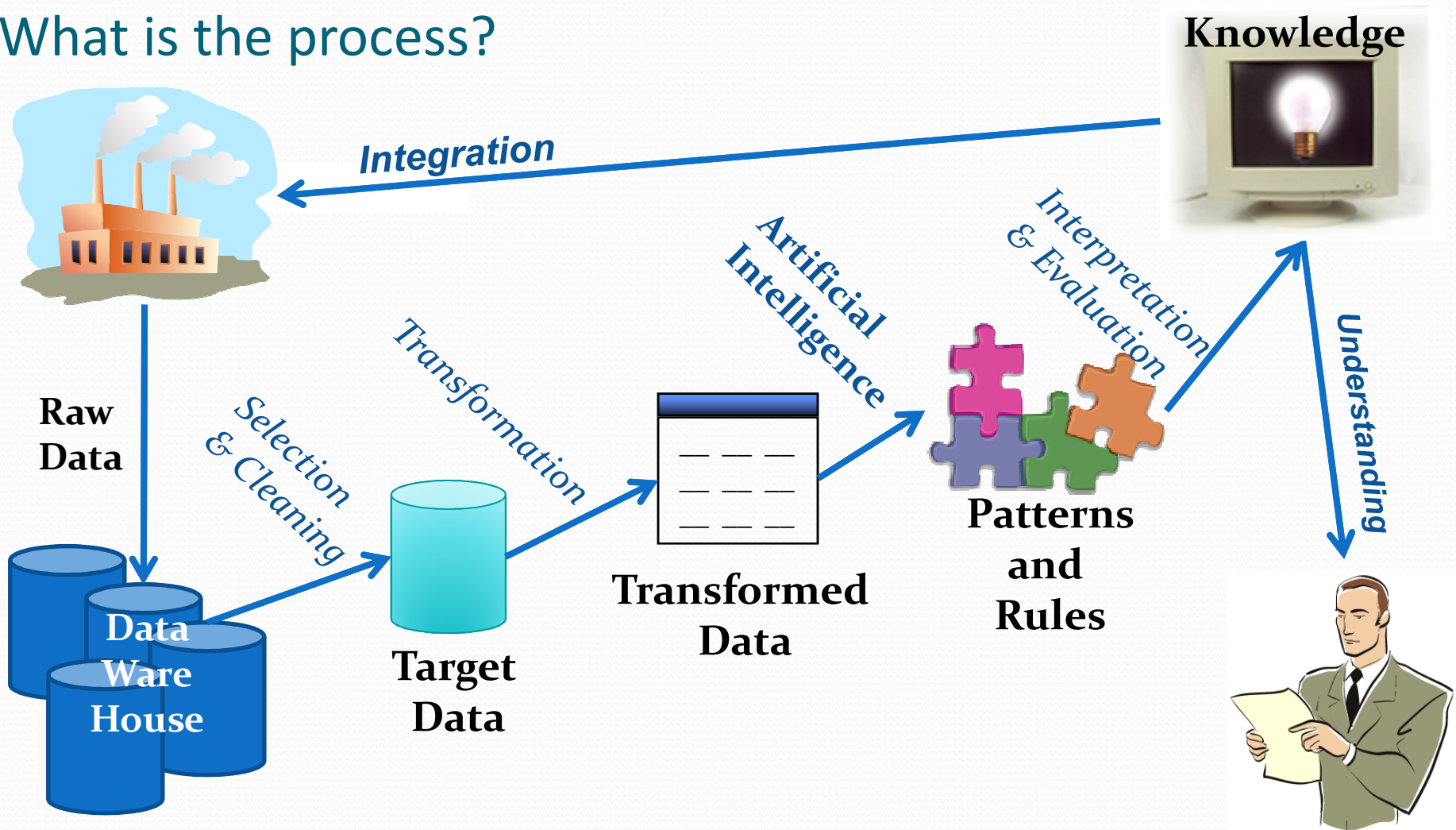
**DB**

# Data Mining (DM)

## Why?

- Moore's law
  - Computer speed doubles every 18 months
- Storage law
  - Total storage doubles every 9 months
- Consequence
  - very little data will ever be analyzed by humans
- Knowledge discovery is **NEEDED** to make sense and use of data ➜ Data Mining

**GAP**

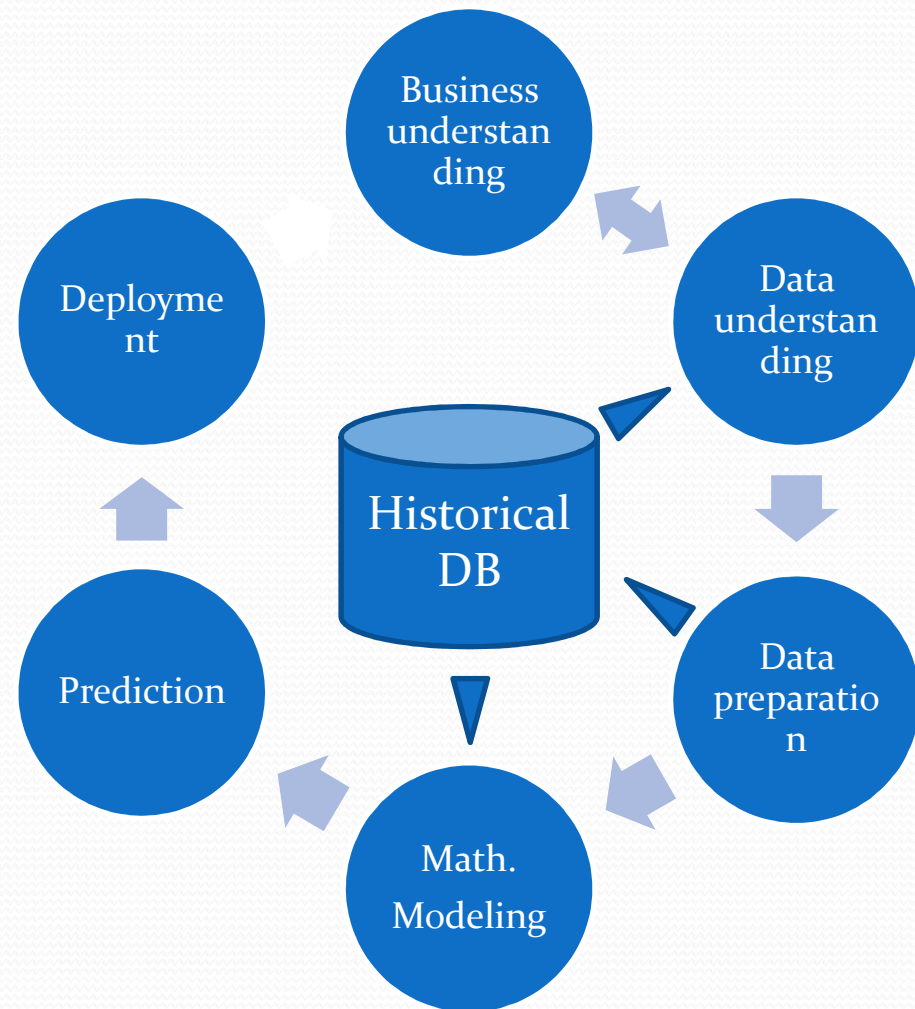1988    1990    1992    1994

—Processing   —Disk

# Data Mining (DM)

## What is the process?

**Knowledge**

*Integration*

**Raw Data**

*Selection & Cleaning*

*Transformation*

*Artificial Intelligence*

*Interpretation & Evaluation*

*Understanding*

**Data Ware House**

**Target Data**

**Transformed Data**

**Patterns and Rules**

# Data Mining (DM)

## What is the methodology?

- Easiness to retrieve the knowledge
- Detect **hidden** and **complex** relationships
- The crisp DM process ([www.crisp-dm.org](www.crisp-dm.org))
  - CRoss Industry Standard Process for DM = **World Standard**
  - Step-by-step data mining guide

Business understanding

Data understanding

Deployment

Historical DB

Data preparation

Prediction

Math. Modeling

# Data Mining (DM)

## What is the methodology?

| Business Understanding | Data Understanding | Data Preparation | Modeling | Evaluation | Deployment |
|---|---|---|---|---|---|
| **Determine Business Objectives**<br>*Background*<br>*Business Objectives*<br>*Business Success Criteria*<br><br>**Situation Assessment**<br>*Inventory of Resources*<br>*Requirements, Assumptions, and Constraints*<br>*Risks and Contingencies*<br>*Terminology*<br>*Costs and Benefits*<br><br>**Determine Data Mining Goal**<br>*Data Mining Goals*<br>*Data Mining Success Criteria*<br><br>**Produce Project Plan**<br>*Project Plan*<br>*Initial Asessment of Tools and Techniques* | **Collect Initial Data**<br>*Initial Data Collection Report*<br><br>**Describe Data**<br>*Data Description Report*<br><br>**Explore Data**<br>*Data Exploration Report*<br><br>**Verify Data Quality**<br>*Data Quality Report* | *Data Set*<br>*Data Set Description*<br><br>**Select Data**<br>*Rationale for Inclusion / Exclusion*<br><br>**Clean Data**<br>*Data Cleaning Report*<br><br>**Construct Data**<br>*Derived Attributes*<br>*Generated Records*<br><br>**Integrate Data**<br>*Merged Data*<br><br>**Format Data**<br>*Reformatted Data* | **Select Modeling Technique**<br>*Modeling Technique*<br>*Modeling Assumptions*<br><br>**Generate Test Design**<br>*Test Design*<br><br>**Build Model**<br>*Parameter Settings*<br>*Models*<br>*Model Description*<br><br>**Assess Model**<br>*Model Assessment*<br>*Revised Parameter Settings* | **Evaluate Results**<br>*Assessment of Data Mining Results w.r.t. Business Success Criteria*<br>*Approved Models*<br><br>**Review Process**<br>*Review of Process*<br><br>**Determine Next Steps**<br>*List of Possible Actions*<br>*Decision* | **Plan Deployment**<br>*Deployment Plan*<br><br>**Plan Monitoring and Maintenance**<br>*Monitoring and Maintenance Plan*<br><br>**Produce Final Report**<br>*Final Report*<br>*Final Presentation*<br><br>**Review Project**<br>*Experience Documentation* |

# Data Mining (DM)

Result of DM includes …

- **Forecasting** what may happen in the future
- **Classifying** objects into groups by recognizing patterns
- **Clustering** objects into groups based on their attributes
- **Associating** what events are likely to occur together
- **Sequencing** what events are likely to lead to later events

# Data Mining (DM)

Is **not** …

- Brute-force crunching of bulk data
- "Blind" application of algorithms
- Going to find relationships where none exist
- Presenting data in different ways
- A database intensive task
- A complex technology requiring an advanced degree in computer science

# Data Mining (DM)

## Versus statistical analysis

- **Statistics**
  - Core of Data Mining
  - Help to make the difference between noise and significant findings

- **Data Mining**
  - Cover the entire data analysis process
  - Knowledge extraction

What're the events occurring together?

**Data mining**

**Statistics**

Is the relationship significant?
e.g. Use t-test

# Mathematical models

Different types

**Predictive models**
Predicting & Classifying

- Decision trees
- Regressions
- Neural networks
- Others

**Descriptive models**
Grouping and Associate

- Factorial analysis
- Clustering
- Associations

# Mathematical models

## Predictive models

| Decision trees | Regressions | Neural networks | Other supervised learning |
|---|---|---|---|
| <ul><li>Classification tree</li><li>ID3</li><li>**C4.5**</li><li>CHAID</li><li>ECHAID</li><li>CART</li><li>C5</li><li>J48</li><li>QUEST</li><li>M5P</li></ul> | <ul><li>Multi linear regression</li><li>Polynomial regression</li><li>Logistic regression</li><li>Proportional hazards models – Cox</li><li>Partial Least Squares regression - PLS</li></ul> | <ul><li>**Multi Layer Perceptron - MLP**</li><li>Probabilistic Neural Network – PNN</li><li>**Plenty** ➔ AHC, TDNN, ARP, AMF, ALN, GRNN, BSB, FCM, BM, MFT, RCC, BPTT, RTRL, EKF, AG, BAM, TAM, etc.</li></ul> | <ul><li>Bayesian networks</li><li>Support Vector Machine – SVM</li><li>SVM for Regression - SVR</li></ul> |

# Mathematical models

## Descriptive models

| Factorial analysis | Clustering | Association |
|---|---|---|
| • Principal Component analysis - PCA<br>• Independent Component Analysis - ICA<br>• Correspondence analysis<br>• Multiple correspondence analysis<br>• Multiple discriminant analysis | • Hierarchical clustering – dendrograms<br>• **K-Means**<br>• X-Means<br>• K-Medoids<br>• Fuzzy c-Means<br>• Self Organizing Maps – SOM<br>• Nearest Neighbor Search – NNS<br>• Expectation-Maximization<br>• Optics | • **Apriori**<br>• Generalized Rule Induction – GRI<br>• Carma<br>• Tertius<br>• Generalized Sequential Patern - GSP |

# Artificial Neural Networks

## What is it?

- Advantages
  - Learn from training experience
  - Extract non-linear relationships
  - Works with all data types
  - Accurate
- Drawbacks
  - Risk of 'overfit' the data
  - Extensive amount of training time
  - Black box → difficult interpretation

# Artificial Neural Networks

Possible applications ...

- Straightening cost/time assessment during ship design



Plate thickness, Stiffener scantling, Stiffener spacing, etc.

Ship/section name, Workload, Worker, etc.

**ANN OUTPUT** ( Correlation factor (**) : 0,827 )

STRAIGHTENING TIME

**9 inputs parameters**
**1 output parameter**
**R2 = 0.827**

# Artificial Neural Networks

## Possible applications …

- Blocks cost/time assessment during ship design



Plate thickness,
Stiffener scantling,
Stiffener spacing, etc.

Ship/section name,
Workload,
Worker, etc.

( Correlation factor (**) : 0,984 )

LOG_TIME_TP_CALC

12 inputs parameters
1 output parameter
R2 = 0.984

# Artificial Neural Networks

Possible applications …

- Place of corrosion prediction
- Part failure prediction

(Conditioned Based Maintenance)



Fast Fourier Transformation (Hanning Window)

Amplitude

Frequency

Original Data   y = sin(x)+sin(x*2)*2+sin(x*0.8)*0.5+sin(x*4)*0.2+sin(x*5)*0.7

Time

# Artificial Neural Networks

Possible applications ...

- Face detection
  - Count number of passenger in a cruise ship (evacuation)
  - Recognize passenger in a cruise ship
  - Detection of intrusions (terrorism)

# Decision trees

## What is it?

- Advantages
  - Intuitive outputs
  - Handle all types of attributes (numeric and symbolic)
  - Have value even with little hard data
  - Can be combined with other DM
- Drawbacks
  - Target must be symbolic



if X > 5 then blue
else if Y > 3 then blue
else if X > 2 then green
else blue

# Decision trees
# Possible applications…

- Identify the rules that generate costs in the design or the operation of a ship

- Prediction of symbolic attributes
  - What is the risk of (high, …, low)?
  - What is the cost of (high, …, low)?
  - Characterization of the complexity of a system (for maintenance)



**Goal Classification: TPMCARRE-BIN**

- HIGH : 504
- LOW : 513

# Clustering models

## What is it?

- Advantage
  - Data segmentation
    - Reduce the quantity of data for future analyze
  - Classification of the data in different groups
  - Extraction of knowledge about not known groups
- Drawbacks
  - Sometimes give different results for each run
  - Concept of mean is often required

# Clustering models

## Possible applications ...

- Classification of ship parts/section/blocks
- Identify different groups of ships in a fleet
- Gather different identical event sequences (maintenance/repair)

- Possibility to combine with another model for the prediction

# Clustering models

Possible applications ...

- Image segmentation
  - Contour detection
- Edge distortion measurement?
- Corrosion surface detection

# Association rules

## What is it?

- Advantages
  - Rules are intuitive
  - Works with huge DB
  - Can detect event sequences

- Drawbacks
  - Huge number of rules
    - Need to be filtered

# Association rules

Possible applications …

- Ship maintenance and operation
  - Do certain faults/incidents lead to specific repairs?
  - Do certain repairs produce subsequent faults/incidents?
  - Are there repairs that lead to other repairs?

# Data mining software's

## Open source or commercial?

- Open source software's
  - Considerably improved
  - Integrates huge number of different algorithm's
  - Can manage huge DB
- Commercial software's
  - Better access to different DB
  - Better exploitation of models
  - Better reporting

**Knime**
**Tanagra**
**Weka**
**R**
**Orange**
**RapidMiner**
**...**

**SPAD**
**SAS**
**SPSS**
**STATISTICA**
**S-PLUS**
**...**

# Data mining software's

## Open source

- R - http://www.r-project.org/

# Data mining software's
## Open source

**oran e** (logo)

- Orange – http://www.ailab.si/orange/
  - (+) ... se, and stor...
  - (-) I... pro...
  - (-) I...
  - (-) I... iles)
  - (+) ...
  - (+) ...
  - (-) V...

# Data mining software's

## Open source

- Weka - http://www.cs.waikato.ac.nz/ml/weka/
  - (+) T... e and stor...
  - (+)
  - (-) ... prol...
  - (+) ... hm
  - (+-) ... program (fre...
  - (-) ...

# Data mining software's

## Open source

- Knime - http://www.knime.org/

# Data mining software's

## Open source

- Tanagra - http://eric.univ-
lyon2...
  - (+) ...he
    logi...
  - (-) ...
    prol...
  - (+) ...
    com...
  - (+) ...m
  - (-) ...
  - (-) N...

# Data mining software's

## Open source



- Rapid Miner - http://rapid-i.com/c...
  - (+) Tr... ...e and store...
  - (-) No... exam...
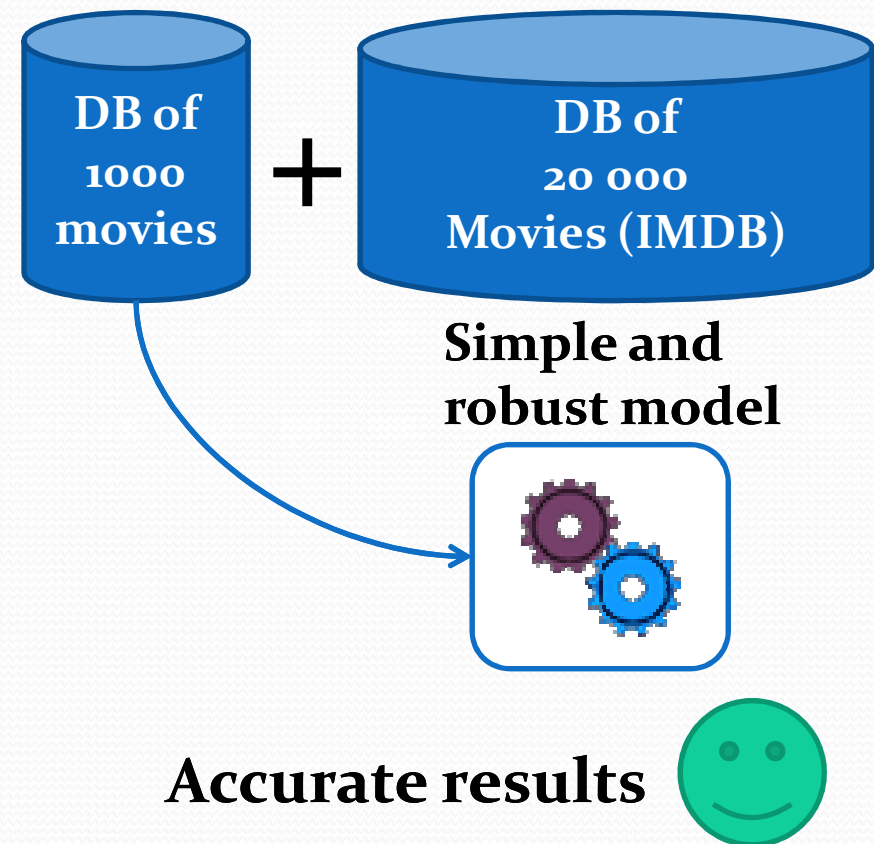  - (+) Bi...

# An interesting story …

Data mining to predict movie rating