# Data Mining and Knowledge Discovery

## Part of
## Jožef Stefan IPS Programme - ICT3
## Part overlapping with ICT2
## Statistics Programme

## 2018 / 2019

# Nada Lavrač

Jožef Stefan Institute

Ljubljana, Slovenia

# Data Mining 2018/2019 Logistics: Course participants

Contacts:  http://kt.ijs.si/petra_kralj/dmkd.html

- Nada Lavrač: nada.lavrac@ijs.si

- Petra Kralj Novak: petra.kralj.novak@ijs.si

| IPS ICT3 students<br><br>Data and text mining<br>Knowledge Technologies Module | Živa Prelog<br>Blaž Škrlj<br>Junoš Lukan<br>Luka Žnidarič<br>Tadej Krivec<br>Tine Kolenik<br>Urban Škvorc |
|---|---|
| IPS ICT2 students<br><br>Data mining and knowledge discovery | Andrejaana Andova<br>Iztok Renčelj<br>Martin Molan<br>Patrik Zajec |
| Statistics<br><br>Podatkovno rudarjenje in odkrivanje zakonitosti v podatkih | Maja Buhin Pandur<br>Tina Grbac |

# Course Schedule – 2016/17

| | | | | | IKT2 | IKT3 | STAT |
|---|---|---|---|---|---|---|---|
| torek | 6.11. | 17-19 | MPS | Nada Lavrac | ✓ | ✓ | ✓ |
| sreda | 7.11. | 16-19 | MPS | Bojan Cestnik | ✓ | | |
| četrtek | 8.11. | 17-19 | Oranžna | Petra Kralj Novak | ✓ | ✓ | ✓ |
| torek | 13.11. | 15-17 | MPS | Nada Lavrac | ✓ | ✓ | ✓ |
| četrtek | 15.11. | 15-18 | Oranžna | Petra Kralj Novak | ✓ | ✓ | ✓ |
| sreda | 21.11. | 15-19 | MPŠ | Dunja Mladenić | ✓ | | |
| četrtek | 22.11. | 17-19 | Oranžna | Nada Lavrac | ✓ | ✓ | ✓ |
| četrtek | 29.11. | 15-18 | Oranžna | Petra Kralj Novak | ✓ | ✓ | ✓ |
| četrtek | 6.12. | 15-17 | Oranžna | Petra Kralj Novak | ✓ | ✓ | ✓ |
| ponedelje | 10.12. | 16-18 | Oranžna | Dunja Mladenić | ✓ | | |
| petek | 14.12. | 15-18 | Oranžna | Martin Žnidaršič | ✓ | ✓ | ✓ |
| sreda | 19.12. | 16-18 | Oranžna | Petra Kralj Novak | ✓ | ✓ | ✓ |
| četrtek | 10.1. | 15-17 | Oranžna | Petra Kralj Novak | ✓ | ✓ | ✓ |
| ponedelje | 14.1. | 17-19 | MPŠ | Dunja Mladenić | ✓ | | |

# Data Mining: PhD Credits and Coursework

- Attending lectures
- Attending practical exercises
  - Theory exercises and hands-on (intro to WEKA by dr. Petra Kralj Novak)
- **Written exam (40%)**
- **Seminar (60%):**
  - Data analysis of your own data (e.g., using WEKA for questionnaire data analysis)
  - Implementing a selected data mining workflow in the ClowdFlows data mining platform
  - …. own initiative is welcome …

# Data Mining: PhD Credits and coursework

**Exam:** Written exam (60 minutes) - Theory

**Seminar: topic selection + results presentation**

- One hour available for seminar topic discussion – one page written proposal defining the task and the selected dataset
- Deliver written report + electronic copy (4 pages in Information Society paper format, instructions on the web)
  - Report on data analysis of own data needs to follow the CRISP-DM methodology
  - Report on DM SW development needs to include SW compatible with the ClowdFlows I/O requirements
  - Presentation of your seminar results (15 minutes each: 10 minutes presentation + 5 minutes discussion)

# Data Mining: ICT2 Credits and Coursework

- 20 credits (8 Lavrač + 4 Cestnik + 8 Mladenić)

# Course Outline

**I. Introduction**
- – Data Mining and KDD process
- – Introduction to Data Mining
- – Data Mining platforms

**II. Predictive DM Techniques**
- – Decision Tree learning
- - Bayesian classifier
- – Classification rule learning
- – Classifier Evaluation

**III. Regression**

**IV. Descriptive DM**
- – Predictive vs. descriptive induction
- – Subgroup discovery
- – Association rule learning
  Hierarchical clustering

**V. Relational Data Mining**
- – RDM and Inductive Logic Programming
- – Propositionalization
- – Semantic data mining

**VI. Advanced Topics**

# Part I. Introduction

➡ Data Mining and the KDD process

- Introduction to Data Mining

- Data Mining platforms

# Machine Learning and Data Mining

- Machine Learning (ML) – computer algorithms/machines that learn predictive models from class-labeled data

- Data Mining (DM) – extraction of useful information from data: discovering relationships and patterns that have not previously been known, and use of ML techniques applied to solving real-life data analysis problems

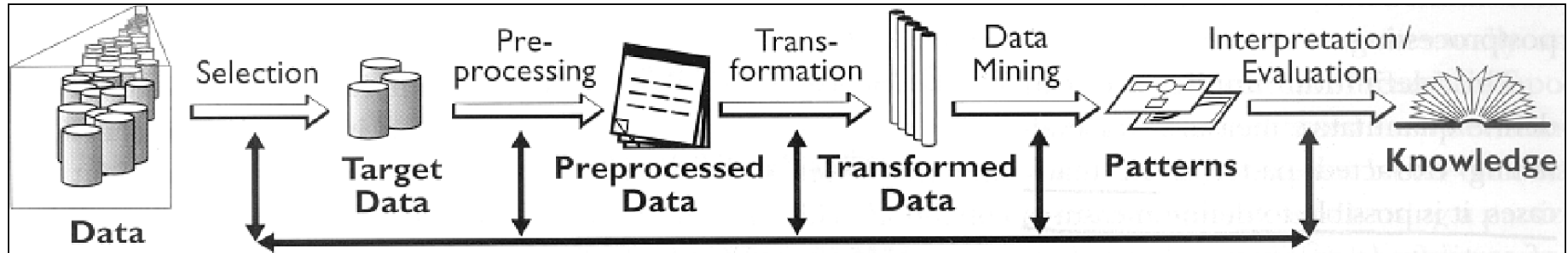- Knowledge discovery in databases (KDD) – the process of knowledge discovery

# Machine Learning and Data Mining

- Machine Learning (ML) – computer algorithms/machines that learn predictive models from class-labeled data

- Data Mining (DM) – extraction of useful information from data: discovering relationships and patterns that have not previously been known, and use of ML techniques applied to solving real-life data analysis problems

- Knowledge Discovery in Databases (KDD) – the process of knowledge discovery

# Data Mining and KDD

- Buzzword since 1996

- KDD is defined as "the process of identifying valid, novel, potentially useful and ultimately understandable models/patterns in data." *

- Data Mining (DM) is the key step in the KDD process, performed by using data mining techniques for extracting models or interesting patterns from the data.

*Usama M. Fayyad, Gregory Piatesky-Shapiro, Pedhraic Smyth: The KDD Process for Extracting Useful Knowledge form Volumes of Data. Comm ACM, Nov 96/Vol 39 No 11*

# KDD Process: CRISP-DM

KDD process of discovering useful knowledge from data
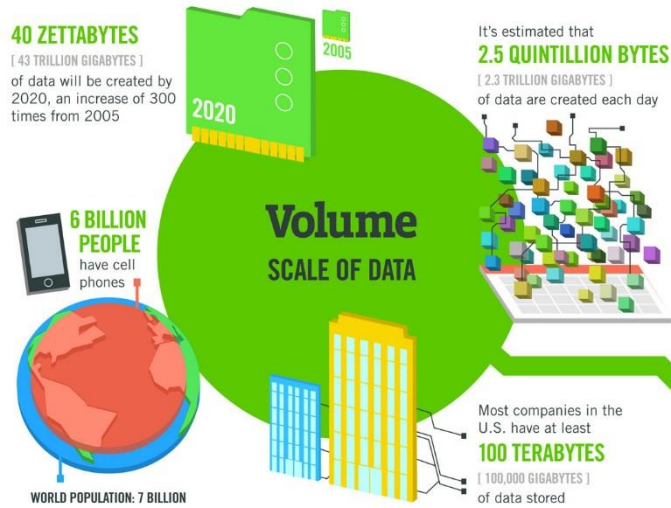


- KDD process involves several phases:
  - data preparation
  - data mining (machine learning, statistics)
  - evaluation and use of discovered patterns
- Data mining is the key step, but represents only 15%-25% of the entire KDD process

# **Big Data**

- Big Data – Buzzword since 2008 (special issue of Nature on Big Data)

  - data and techniques for dealing with very large volumes of data, possibly dynamic data streams

  - requiring large data storage resources, special algorithms for parallel computing architectures.

# The 4 Vs of Big Data



The FOUR V's of Big Data

**40 ZETTABYTES**
[ 43 TRILLION GIGABYTES ]
of data will be created by 2020, an increase of 300 times from 2005

2005
2020

It's estimated that
**2.5 QUINTILLION BYTES**
[ 2.3 TRILLION GIGABYTES ]
of data are created each day

**6 BILLION PEOPLE**
have cell phones

WORLD POPULATION: 7 BILLION

**Volume**
SCALE OF DATA

Most companies in the U.S. have at least
**100 TERABYTES**
[ 100,000 GIGABYTES ]
of data stored

From traffic patterns and music downloads to web history and medical records, data is recorded, stored, and analyzed to enable the technology and services that the world relies on every day. But what exactly is big data, and how can these massive amounts of data be used?

As a leader in the sector, IBM data scientists break big data into four dimensions: **Volume, Velocity, Variety and Veracity**

Depending on the industry and organization, big data encompasses information from multiple internal and external sources such as transactions, social media, enterprise content, sensors and mobile devices. Companies can leverage data to adapt their products and services to better meet customer needs, optimize operations and infrastructure, and find new sources of revenue.

By 2015
**4.4 MILLION IT JOBS**
will be created globally to support big data, with 1.9 million in the United States

As of 2011, the global size of data in healthcare was estimated to be
**150 EXABYTES**
[ 161 BILLION GIGABYTES ]

**30 BILLION PIECES OF CONTENT**
are shared on Facebook every month

**Variety**
DIFFERENT FORMS OF DATA

By 2014, it's anticipated there will be
**420 MILLION WEARABLE, WIRELESS HEALTH MONITORS**

**4 BILLION+ HOURS OF VIDEO**
are watched on YouTube each month

**400 MILLION TWEETS**
are sent per day by about 200 million monthly active users

The New York Stock Exchange captures
**1 TB OF TRADE INFORMATION**
during each trading session

Modern cars have close to
**100 SENSORS**
that monitor items such as fuel level and tire pressure

**Velocity**
ANALYSIS OF STREAMING DATA

By 2016, it is projected there will be
**18.9 BILLION NETWORK CONNECTIONS**
– almost 2.5 connections per person on earth

**1 IN 3 BUSINESS LEADERS**
don't trust the information they use to make decisions

**27% OF RESPONDENTS**
in one survey were unsure of how much of their data was inaccurate

**Veracity**
UNCERTAINTY OF DATA

Poor data quality costs the US economy around
**$3.1 TRILLION A YEAR**

**Sources:** McKinsey Global Institute, Twitter, Cisco, Gartner, EMC, SAS, IBM, MEPTEC, QAS

IBM

# Data Science

- Data Science – buzzword since 2012 when Harvard Business Review called it "The Sexiest Job of the 21st Century"

  - an interdisciplinary field that uses scientific methods, processes, algorithms and systems to extract knowledge and insights from data in various forms, both structured and unstructured, similar to data mining.

  - used interchangeably with earlier concepts like business analytics, business intelligence, predictive modeling, and statistics.
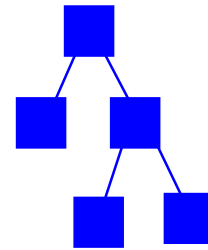
# Data Mining in a Nutshell

| Person | Age | Spect. presc. | Astigm. | Tear prod. | Lenses |
|---|---|---|---|---|---|
| O1 | 17 | myope | no | reduced | NONE |
| O2 | 23 | myope | no | normal | SOFT |
| O3 | 22 | myope | yes | reduced | NONE |
| O4 | 27 | myope | yes | normal | HARD |
| O5 | 19 | hypermetrope | no | reduced | NONE |
| O6-O13 | ... | ... | ... | ... | ... |
| O14 | 35 | hypermetrope | no | normal | SOFT |
| O15 | 43 | hypermetrope | yes | reduced | NONE |
| O16 | 39 | hypermetrope | yes | normal | NONE |
| O17 | 54 | myope | no | reduced | NONE |
| O18 | 62 | myope | no | normal | NONE |
| O19-O23 | ... | ... | ... | ... | ... |
| O24 | 56 | hypermetrope | yes | normal | NONE |

data
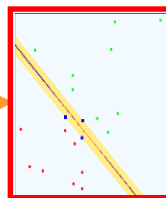
knowledge discovery
from data

Data Mining

model, patterns, …

**Given:** transaction data table, relational database, text
documents, Web pages
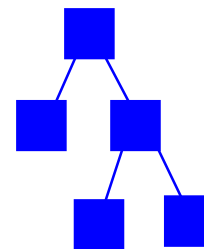**Find:** a classification model, a set of interesting patterns

# Data Mining in a Nutshell

| Person | Age | Spect. presc. | Astigm. | Tear prod. | Lenses |
|--------|-----|---------------|---------|------------|--------|
| O1 | 17 | myope | no | reduced | NONE |
| O2 | 23 | myope | no | normal | SOFT |
| O3 | 22 | myope | yes | reduced | NONE |
| O4 | 27 | myope | yes | normal | HARD |
| O5 | 19 | hypermetrope | no | reduced | NONE |
| O6-O13 | ... | ... | ... | ... | ... |
| O14 | 35 | hypermetrope | no | normal | SOFT |
| O15 | 43 | hypermetrope | yes | reduced | NONE |
| O16 | 39 | hypermetrope | yes | normal | NONE |
| O17 | 54 | myope | no | reduced | NONE |
| O18 | 62 | myope | no | normal | NONE |
| O19-O23 | ... | ... | ... | ... | ... |
| O24 | 56 | hypermetrope | yes | normal | NONE |

data

knowledge discovery
from data

Data Mining

model, patterns, …

**Given:** transaction data table, relational database, text documents, Web pages

**Find:** a classification model, a set of interesting patterns

new unclassified instance → classified instance

black box classifier
no explanation

symbolic model
symbolic patterns

explanation

# Simplified example: Learning a classification model from contact lens data

| Person | Age | Spect. presc. | Astigm. | Tear prod. | Lenses |
|--------|-----|---------------|---------|------------|--------|
| O1 | 17 | myope | no | reduced | NONE |
| O2 | 23 | myope | no | normal | SOFT |
| O3 | 22 | myope | yes | reduced | NONE |
| O4 | 27 | myope | yes | normal | HARD |
| O5 | 19 | hypermetrope | no | reduced | NONE |
| O6-O13 | … | … | … | … | … |
| O14 | 35 | hypermetrope | no | normal | SOFT |
| O15 | 43 | hypermetrope | yes | reduced | NONE |
| O16 | 39 | hypermetrope | yes | normal | NONE |
| O17 | 54 | myope | no | reduced | NONE |
| O18 | 62 | myope | no | normal | NONE |
| O19-O23 | … | … | … | … | … |
| O24 | 56 | hypermetrope | yes | normal | NONE |

# Simplified example: Learning a classification model from contact lens data

| Person | Age | Spect. presc. | Astigm. | Tear prod. | Lenses |
|---|---|---|---|---|---|
| O1 | young | myope | no | reduced | NONE |
| O2 | young | myope | no | normal | SOFT |
| O3 | young | myope | yes | reduced | NONE |
| O4 | young | myope | yes | normal | HARD |
| O5 | young | hypermetrope | no | reduced | NONE |
| O6-O13 | ... | ... | ... | ... | ... |
| O14 | pre-presbyo | hypermetrope | no | normal | SOFT |
| O15 | pre-presbyo | hypermetrope | yes | reduced | NONE |
| O16 | pre-presbyo | hypermetrope | yes | normal | NONE |
| O17 | presbyopic | myope | no | reduced | NONE |
| O18 | presbyopic | myope | no | normal | NONE |
| O19-O23 | ... | ... | ... | ... | ... |
| O24 | presbyopic | hypermetrope | yes | normal | NONE |

Data Mining

# Task reformulation: Binary Class Values

| Person | Age | Spect. presc. | Astigm. | Tear prod. | Lenses |
|--------|-----|---------------|---------|------------|--------|
| O1 | 17 | myope | no | reduced | NO |
| O2 | 23 | myope | no | normal | YES |
| O3 | 22 | myope | yes | reduced | NO |
| O4 | 27 | myope | yes | normal | YES |
| O5 | 19 | hypermetrope | no | reduced | NO |
| O6-O13 | … | … | … | … | … |
| O14 | 35 | hypermetrope | no | normal | YES |
| O15 | 43 | hypermetrope | yes | reduced | NO |
| O16 | 39 | hypermetrope | yes | normal | NO |
| O17 | 54 | myope | no | reduced | NO |
| O18 | 62 | myope | no | normal | NO |
| O19-O23 | … | … | … | … | … |
| O24 | 56 | hypermetrope | yes | normal | NO |

Binary classes (positive vs. negative examples of Target class)
 - for Concept learning – classification and class description
  - for Subgroup discovery – exploring patterns
    characterizing groups of instances of target class

# Learning from Numeric Class Data

| Person | Age | Spect. presc. | Astigm. | Tear prod. | LensPrice |
|--------|-----|---------------|---------|------------|-----------|
| O1 | 17 | myope | no | reduced | 0 |
| O2 | 23 | myope | no | normal | 8 |
| O3 | 22 | myope | yes | reduced | 0 |
| O4 | 27 | myope | yes | normal | 5 |
| O5 | 19 | hypermetrope | no | reduced | 0 |
| O6-O13 | … | … | … | … | … |
| O14 | 35 | hypermetrope | no | normal | 5 |
| O15 | 43 | hypermetrope | yes | reduced | 0 |
| O16 | 39 | hypermetrope | yes | normal | 0 |
| O17 | 54 | myope | no | reduced | 0 |
| O18 | 62 | myope | no | normal | 0 |
| O19-O23 | … | … | … | … | … |
| O24 | 56 | hypermetrope | yes | normal | 0 |

Numeric class values – regression analysis

# Learning from Unlabeled Data

| Person | Age | Spect. presc. | Astigm. | Tear prod. | Lenses |
|--------|-----|---------------|---------|------------|--------|
| O1 | 17 | myope | no | reduced | NONE |
| O2 | 23 | myope | no | normal | SOFT |
| O3 | 22 | myope | yes | reduced | NONE |
| O4 | 27 | myope | yes | normal | HARD |
| O5 | 19 | hypermetrope | no | reduced | NONE |
| O6-O13 | … | … | … | … | … |
| O14 | 35 | hypermetrope | no | normal | SOFT |
| O15 | 43 | hypermetrope | yes | reduced | NONE |
| O16 | 39 | hypermetrope | yes | normal | NONE |
| O17 | 54 | myope | no | reduced | NONE |
| O18 | 62 | myope | no | normal | NONE |
| O19-O23 | … | … | … | … | … |
| O24 | 56 | hypermetrope | yes | normal | NONE |

Unlabeled data - clustering: grouping of similar instances
- association rule learning

# Data Mining, ML and Statistics

- All three areas have a long tradition of developing inductive techniques for data analysis.
  - reasoning from properties of a data sample to properties of a population
- **DM vs. ML - Viewpoint in this course**:
  - Data Mining is the application of Machine Learning techniques to  hard real-life data analysis problems

# Data Mining, ML and Statistics

- All three areas have a long tradition of developing inductive techniques for data analysis.
  - reasoning from properties of a data sample to properties of a population
- **DM vs. Statistics:**
  - **Statistics**
    - Hypothesis testing when certain theoretical expectations about the data distribution, independence, random sampling, sample size, etc. are satisfied
    - Main approach: best fitting all the available data
  - **Data mining**
    - Automated construction of understandable patterns, and structured models
    - Main approach: structuring the data space, heuristic search for decision trees, rules, … covering (parts of) the data space

# Why learn and use symbolic models

**Given:** the learned classification model
(a decision tree or a set of rules)

**Find:**  the class label for a new unlabeled instance

# Why learn and use symbolic models

**Given:** the learned classification model
      (a decision tree or a set of rules)

**Find:** the class label for a new unlabeled instance

new unclassified instance              classified instance

# Why learn and use symbolic models

**Given:** the learned classification model
(a decision tree or a set of rules)

**Find:** - the class label for a new unlabeled instance

new unclassified instance → classified instance

- use the model for the explanation of classifications of new data instances
- use the discovered patterns for data exploration

# Data Mining

| Person | Age | Spect. presc. | Astigm. | Tear prod. | Lenses |
|--------|-----|---------------|---------|------------|--------|
| O1 | 17 | myope | no | reduced | NONE |
| O2 | 23 | myope | no | normal | SOFT |
| O3 | 22 | myope | yes | reduced | NONE |
| O4 | 27 | myope | yes | normal | HARD |
| O5 | 19 | hypermetrope | no | reduced | NONE |
| O6-O13 | ... | ... | ... | ... | ... |
| O14 | 35 | hypermetrope | no | normal | SOFT |
| O15 | 43 | hypermetrope | yes | reduced | NONE |
| O16 | 39 | hypermetrope | yes | normal | NONE |
| O17 | 54 | myope | no | reduced | NONE |
| O18 | 62 | myope | no | normal | NONE |
| O19-O23 | ... | ... | ... | ... | ... |
| O24 | 56 | hypermetrope | yes | normal | NONE |

data

knowledge discovery from data

Data Mining

model, patterns, …

**Given:** transaction data table, relational database, text documents, Web pages

**Find:** a classification model, a set of interesting patterns

new unclassified instance → classified instance

black box classifier
no explanation

symbolic model
symbolic patterns

explanation

# Contact lens data

| Person | Age | Spect. presc. | Astigm. | Tear prod. | Lenses |
|--------|-----|---------------|---------|------------|--------|
| O1 | 17 | myope | no | reduced | NONE |
| O2 | 23 | myope | no | normal | SOFT |
| O3 | 22 | myope | yes | reduced | NONE |
| O4 | 27 | myope | yes | normal | HARD |
| O5 | 19 | hypermetrope | no | reduced | NONE |
| O6-O13 | … | … | … | … | … |
| O14 | 35 | hypermetrope | no | normal | SOFT |
| O15 | 43 | hypermetrope | yes | reduced | NONE |
| O16 | 39 | hypermetrope | yes | normal | NONE |
| O17 | 54 | myope | no | reduced | NONE |
| O18 | 62 | myope | no | normal | NONE |
| O19-O23 | … | … | … | … | … |
| O24 | 56 | hypermetrope | yes | normal | NONE |

# Pattern discovery in Contact lens data

| Person | Age | Spect. presc. | Astigm. | Tear prod. | Lenses |
|--------|-----|---------------|---------|------------|--------|
| O1 | 17 | myope | no | reduced | NONE |
| O2 | 23 | myope | no | normal | SOFT |
| O3 | 22 | myope | yes | reduced | NONE |
| O4 | 27 | myope | yes | normal | HARD |
| O5 | 19 | hypermetrope | no | reduced | NONE |
| O6-O13 | ... | ... | ... | ... | ... |
| O14 | 35 | hypermetrope | no | normal | SOFT |
| O15 | 43 | hypermetrope | yes | reduced | NONE |
| O16 | 39 | hypermetrope | yes | normal | NONE |
| O17 | 54 | myope | no | reduced | NONE |
| O18 | 62 | myope | no | normal | NONE |
| O19-O23 | ... | ... | ... | ... | ... |
| O24 | 56 | hypermetrope | yes | normal | NONE |

**PATTERN**

**Rule:**

IF
Tear prod. = reduced

THEN
Lenses = NONE

# Learning a classification model from contact lens data

| Person | Age | Spect. presc. | Astigm. | Tear prod. | Lenses |
|--------|-----|---------------|---------|------------|--------|
| O1 | young | myope | no | reduced | NONE |
| O2 | young | myope | no | normal | SOFT |
| O3 | young | myope | yes | reduced | NONE |
| O4 | young | myope | yes | normal | HARD |
| O5 | young | hypermetrope | no | reduced | NONE |
| O6-O13 | … | … | … | … | … |
| O14 | pre-presbyc | hypermetrope | no | normal | SOFT |
| O15 | pre-presbyc | hypermetrope | yes | reduced | NONE |
| O16 | pre-presbyc | hypermetrope | yes | normal | NONE |
| O17 | presbyopic | myope | no | reduced | NONE |
| O18 | presbyopic | myope | no | normal | NONE |
| O19-O23 | … | … | … | … | … |
| O24 | presbyopic | hypermetrope | yes | normal | NONE |

Data Mining

# Decision tree classification model learned from contact lens data



nodes: attributes
arcs: values of attributes
leaves: classes

# Learning a classification model from contact lens data

| Person | Age | Spect. presc. | Astigm. | Tear prod. | Lenses |
|--------|-----|---------------|---------|------------|--------|
| O1 | 17 | myope | no | reduced | NONE |
| O2 | 23 | myope | no | normal | SOFT |
| O3 | 22 | myope | yes | reduced | NONE |
| O4 | 27 | myope | yes | normal | HARD |
| O5 | 19 | hypermetrope | no | reduced | NONE |
| O6-O13 | … | … | … | … | … |
| O14 | 35 | hypermetrope | no | normal | SOFT |
| O15 | 43 | hypermetrope | yes | reduced | NONE |
| O16 | 39 | hypermetrope | yes | normal | NONE |
| O17 | 54 | myope | no | reduced | NONE |
| O18 | 62 | myope | no | normal | NONE |
| O19-O23 | … | … | … | … | … |
| O24 | 56 | hypermetrope | yes | normal | NONE |

Data Mining



lenses=NONE ← tear production=red
lenses=NONE ← tear production=normal AND astigmatism=yes AND
    spect. pre.=hypermetrope
lenses=SOFT ← tear production=normal AND astigmatism=no
lenses=HARD ← tear production=normal AND astigmatism=yes AND
    spect. pre.=myope
lenses=NONE ←

# Classification rules model learned from contact lens data

lenses=NONE ← tear production=reduced

lenses=NONE ← tear production=normal AND
      astigmatism=yes AND
      spect. pre.=hypermetrope

lenses=SOFT ← tear production=normal AND
      astigmatism=no

lenses=HARD ← tear production=normal AND
      astigmatism=yes AND
      spect. pre.=myope

lenses=NONE ←

# Task reformulation: Binary Class Values

| Person | Age | Spect. presc. | Astigm. | Tear prod. | Lenses |
|--------|-----|---------------|---------|------------|--------|
| O1 | 17 | myope | no | reduced | NO |
| O2 | 23 | myope | no | normal | YES |
| O3 | 22 | myope | yes | reduced | NO |
| O4 | 27 | myope | yes | normal | YES |
| O5 | 19 | hypermetrope | no | reduced | NO |
| O6-O13 | … | … | … | … | … |
| O14 | 35 | hypermetrope | no | normal | YES |
| O15 | 43 | hypermetrope | yes | reduced | NO |
| O16 | 39 | hypermetrope | yes | normal | NO |
| O17 | 54 | myope | no | reduced | NO |
| O18 | 62 | myope | no | normal | NO |
| O19-O23 | … | … | … | … | … |
| O24 | 56 | hypermetrope | yes | normal | NO |

Binary classes (positive vs. negative examples of Target class)
- for Concept learning tasks
    - classification and class description
    - "one vs. all" multi-class learning
- for Subgroup discovery tasks

# Learning from Numeric Class Data

| Person | Age | Spect. presc. | Astigm. | Tear prod. | LensPrice |
|--------|-----|---------------|---------|------------|-----------|
| O1 | 17 | myope | no | reduced | 0 |
| O2 | 23 | myope | no | normal | 8 |
| O3 | 22 | myope | yes | reduced | 0 |
| O4 | 27 | myope | yes | normal | 5 |
| O5 | 19 | hypermetrope | no | reduced | 0 |
| O6-O13 | ... | ... | ... | ... | ... |
| O14 | 35 | hypermetrope | no | normal | 5 |
| O15 | 43 | hypermetrope | yes | reduced | 0 |
| O16 | 39 | hypermetrope | yes | normal | 0 |
| O17 | 54 | myope | no | reduced | 0 |
| O18 | 62 | myope | no | normal | 0 |
| O19-O23 | ... | ... | ... | ... | ... |
| O24 | 56 | hypermetrope | yes | normal | 0 |

Numeric class values – regression analysis

# Learning from Unlabeled Data

| Person | Age | Spect. presc. | Astigm. | Tear prod. | Lenses |
|--------|-----|---------------|---------|------------|--------|
| O1 | 17 | myope | no | reduced | NONE |
| O2 | 23 | myope | no | normal | SOFT |
| O3 | 22 | myope | yes | reduced | NONE |
| O4 | 27 | myope | yes | normal | HARD |
| O5 | 19 | hypermetrope | no | reduced | NONE |
| O6-O13 | … | … | … | … | … |
| O14 | 35 | hypermetrope | no | normal | SOFT |
| O15 | 43 | hypermetrope | yes | reduced | NONE |
| O16 | 39 | hypermetrope | yes | normal | NONE |
| O17 | 54 | myope | no | reduced | NONE |
| O18 | 62 | myope | no | normal | NONE |
| O19-O23 | … | … | … | … | … |
| O24 | 56 | hypermetrope | yes | normal | NONE |

Unlabeled data - clustering: grouping of similar instances
- association rule learning

# Why learn and use symbolic models

**Given:** the learned classification model
(a decision tree or a set of rules)

**Find:** - the class label for a new unlabeled instance

new unclassified instance → classified instance

- use the model for the explanation of classifications of new data instances
- use the discovered patterns for data exploration

# First Generation Data Mining

- **First machine learning algorithms for**
  - Decision tree and rule learning in 1970s and early 1980s by Quinlan, Michalski et al., Breiman et al., …
- **Characterized by**
  - Learning from data stored in a single data table
  - Relatively small set of instances and attributes
- **Lots of ML research followed in 1980s**
  - Numerous conferences ICML, ECML, … and ML sessions at AI conferences IJCAI, ECAI, AAAI, …
  - Extended set of learning tasks and algorithms addressed

# Second Generation Data Mining

- **Developed since 1990s:**
  - Focused on data mining tasks characterized by large datasets described by large numbers of attributes
  - Industrial standard: CRISP-DM methodology (1997)

# Second Generation Data Mining

- **Developed since 1990s:**
  - Focused on data mining tasks characterized by large datasets described by large numbers of attributes
  - Industrial standard: CRISP-DM methodology (1997)



  - New conferences on practical aspects of data mining and knowledge discovery: KDD, PKDD, …
  - New learning tasks and efficient learning algorithms:
    - Learning predictive models: Bayesian network learning,, relational data mining, statistical relational learning, SVMs, …
    - Learning descriptive patterns: association rule learning, subgroup discovery, …

# Second Generation Data Mining Platforms

## Orange, WEKA, KNIME, RapidMiner, …

# Second Generation Data Mining Platforms

Orange, WEKA, KNIME, RapidMiner, …



- include numerous data mining algorithms
- enable data and model visualization
- like Orange, Taverna, WEKA, KNIME, RapidMiner, also enable complex **workflow** construction

# Third Generation Data Mining

- **Orange4WS** (Podpečan et al. 2009), **ClowdFlows** (Kranjc et al. 2012) **and TextFlows** (Perovšek et al. 2016)
  - are service oriented (DM algorithms as web services)
  - user-friendly HCI: canvas for workflow construction
  - include functionality of standard data mining platforms
    - WEKA algorithms, implemented as Web services
  - Include new functionality
    - relational data mining
    - semantic data mining
    - NLP processing and text mining
  - enable simplified construction of Web services from available algorithms
  - ClowdFlows and TextFlows run in a browser – enables data mining, workflow construction and sharing on the web

# ClowdFlows platform

- **Large algorithm repository**
  - Relational data mining
  - All Orange algorithms
  - WEKA algorithms as web services
  - Data and results visualization
  - Text analysis
  - Social network analysis
  - Analysis of big data streams
- **Large workflow repository**
  - Enables access to our technology heritage

# ClowdFlows platform

- Large repository of algorithms
- Large repository of workflows



**Example workflow**:
Propositionalization with RSD
available in ClowdFlows at
http://clowdflows.org/workflow/611/

# TextFlows

- Motivation:
  - Develop an online text mining platform for composition, execution and sharing of text mining workflows

- TextFlows platform – fork of ClowdFlows.org:
  - Specialized on text mining
  - Web-based user interface
  - Visual programming
  - Big roster of existing workflow (mostly text mining) components
  - Cloud-based service-oriented architecture

# "Big Data" Use Case

- Real-time analysis of big data streams

- Example: semantic graph construction from news streams. http://clowdflows.org/workflow/1729/.



RSS Reader — Summarize news article — Triplet Extraction — WordNet lemmatizer on triplets — Sliding Window — Streaming triplet graph

- Example: news monitoring by graph visualization (graph of CNN RSS feeds)

  http://clowdflows.org/streams/data/31/1

# Part I: Summary

- KDD is the overall process of discovering useful knowledge in data
  - many steps including data preparation, cleaning, transformation, pre-processing
- Data Mining is the data analysis phase in KDD
  - DM takes only 15%-25% of the effort of the overall KDD process
  - employing techniques from machine learning and statistics
- Predictive and descriptive induction have different goals: classifier vs. pattern discovery
- Many application areas, many powerful tools available

# Course Outline

**I. Introduction**

- – Data Mining and KDD process
- – Introduction to Data Mining
- – Data Mining platforms

**II. Predictive DM Techniques**

- – Decision Tree learning
- - Bayesian classifier
- – Classification rule learning
- – Classifier Evaluation

**III. Regression**

**IV. Descriptive DM**

- – Predictive vs. descriptive induction
- – Subgroup discovery
- – Association rule learning Hierarchical clustering

**V. Relational Data Mining**

- – RDM and Inductive Logic Programming
- – Propositionalization
- – Semantic data mining

**VI. Advanced Topics**

# Part II. Predictive DM techniques

→ Decision tree learning

- Bayesian Classifier

- Rule learning

- Evaluation

# **Predictive DM - Classification**

- data are objects, characterized with attributes - they belong to different classes (discrete labels)

- given objects described with attribute values, induce a model to predict different classes

- decision trees, if-then rules, discriminant analysis, ...

# **Predictive DM - classification formulated as a machine learning task**

- Given a set of labeled **training examples** (n-tuples of attribute values, labeled by class name)

|  | A1 | A2 | A3 | Class |
|---|---|---|---|---|
| example1 | $v_{1,1}$ | $v_{1,2}$ | $v_{1,3}$ | $C_1$ |
| example2 | $v_{2,1}$ | $v_{2,2}$ | $v_{2,3}$ | $C_2$ |

. .

- Performing generalization from examples (induction)
- Find a **hypothesis** (a decision tree or classification rules) which explains the training examples, e.g. decision trees or classification rules of the form:

IF $(A_i = v_{i,k})$ & $(A_j = v_{j,l})$ & ... THEN Class = $C_n$

# Decision Tree Learning

| Person | Age | Spect. presc. | Astigm. | Tear prod. | Lenses |
|--------|-----|---------------|---------|------------|--------|
| O1 | young | myope | no | reduced | NONE |
| O2 | young | myope | no | normal | SOFT |
| O3 | young | myope | yes | reduced | NONE |
| O4 | young | myope | yes | normal | HARD |
| O5 | young | hypermetrope | no | reduced | NONE |
| O6-O13 | ... | ... | ... | ... | ... |
| O14 | pre-presbyc | hypermetrope | no | normal | SOFT |
| O15 | pre-presbyc | hypermetrope | yes | reduced | NONE |
| O16 | pre-presbyc | hypermetrope | yes | normal | NONE |
| O17 | presbyopic | myope | no | reduced | NONE |
| O18 | presbyopic | myope | no | normal | NONE |
| O19-O23 | ... | ... | ... | ... | ... |
| O24 | presbyopic | hypermetrope | yes | normal | NONE |

Data Mining

# Decision Tree classifier

# Decision tree learning algorithm

- ID3 (Quinlan 1979), CART (Breiman et al. 1984), C4.5, J48 in WEKA, ...
  - create the root node of the tree
  - if all examples from S belong to the same class Cj
    - then label the root with Cj
  - else
    - select the 'most informative' attribute **A** with values **v1, v2, … vn**
    - divide training set **S** into **S1,… , Sn** according to values **v1,…,vn**
    - recursively build sub-trees **T1,…,Tn** for **S1,…,Sn**

# Decision tree search heuristics

- Central choice in decision tree algorithms: Which attribute to test at each node in the tree ? The attribute that is most useful for classifying examples.

- Define a statistical property, called **information gain**, measuring how well a given attribute separates the training examples w.r.t their target classification.

- First define a measure commonly used in information theory, called **entropy**, to characterize the (im)purity of an arbitrary collection of examples.

# **Entropy**

- **S** - training set, **C$_1$,...,C$_N$** - classes
- **Entropy E(S)** – measure of the impurity of training set S

$$E(S) = -\sum_{c=1}^{N} p_c . \log_2 p_c$$

**p$_c$** - prior probability of class **C$_c$** (relative frequency of **C$_c$** in **S**)

- Entropy in binary classification problems

$$\mathbf{E(S) = - p_+ \log_2 p_+ - p_- \log_2 p_-}$$

# Entropy

- $E(S) = - p_+ \log_2 p_+ - p_- \log_2 p_-$

- The entropy function relative to a Boolean classification, as the proportion **p$_+$** of positive examples varies between 0 and 1

# Entropy – why ?

- **Entropy E(S) =** expected amount of information (in bits) needed to assign a class to a randomly drawn object in S (under the optimal, shortest-length code)

- Why ?

- Information theory: optimal length code assigns $-\log_2 p$ bits to a message having probability p

- So, in binary classification problems, the expected number of bits to encode + or − of a random member of S is:

$$p_+ \left( -\log_2 p_+ \right) + p_- \left( -\log_2 p_- \right) = - p_+ \log_2 p_+ - p_- \log_2 p_-$$

# **Entropy – example calculation**

- Training set S: 14 examples (9 pos., 5 neg.)
- Notation: S = [9+, 5-]
- E(S) = - $p_+$ log$_2$p$_+$ - p$_-$ log$_2$p$_-$
- Computing entropy, if probability is estimated by relative frequency

$$E(S) = -\left( \frac{|S_+|}{|S|} \cdot \log \frac{|S_+|}{|S|} \right) - \left( \frac{|S_-|}{|S|} \cdot \log \frac{|S_-|}{|S|} \right)$$

- E([9+,5-]) = - (9/14) log$_2$(9/14) - (5/14) log$_2$(5/14)

  = 0.940

# Information gain
# search heuristic

- **Information gain** measure is aimed to minimize the number of tests needed for the classification of a new object

- **Gain(S,A)** – expected reduction in entropy of S due to sorting on A

$$Gain(S, A) = E(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} \cdot E(S_v)$$

- **Most informative** attribute: **max Gain(S,A)**

# Information gain search heuristic

- Which attribute is more informative, A1 or A2 ?

[9+,5−],  E = 0.94

(A1)

[6+, 2−]
E=0.811

[3+, 3−]
E=1.00

[9+,5−],  E = 0.94

(A2)

[9+, 0−]
E=0.0

[0+, 5−]
E=0.0

- Gain(S,A1) = 0.94 − (8/14 x 0.811 + 6/14 x 1.00) = 0.048

- Gain(S,A2) = 0.94 − 0 = 0.94          A2 has max Gain

# Heuristic search in ID3

- **Search bias:** Search the space of decision trees from simplest to increasingly complex (greedy search, no backtracking, prefer small trees)
- **Search heuristics:** At a node, select the attribute that is most useful for classifying examples, split the node accordingly
- **Stopping criteria:** A node becomes a leaf
  - if all examples belong to same class $C_j$, label the leaf with $C_j$
  - if all attributes were used, label the leaf with the most common value $C_k$ of examples in the node
- **Extension to ID3:** handling noise - tree pruning

# **Pruning of decision trees**

- Avoid overfitting the data by tree pruning
- Pruned trees are
  - less accurate on training data
  - more accurate when classifying unseen data

# **Handling noise – Tree pruning**

Sources of imperfection

1. Random errors (noise) in training examples

   - erroneous attribute values

   - erroneous classification

2. Too sparse training examples (incompleteness)

3. Inappropriate/insufficient set of attributes (inexactness)

4. Missing attribute values in training examples

# **Handling noise – Tree pruning**

- Handling imperfect data

    - handling imperfections of type 1-3

        - pre-pruning (stopping criteria)

        - post-pruning / rule truncation

    - handling missing values

- Pruning avoids perfectly fitting noisy data: relaxing the completeness (fitting all +) and consistency (fitting all -) criteria in ID3

# Prediction of breast cancer recurrence: Tree pruning

Degree_of_malig

&lt; 3       ≥ 3

Tumor_size           Involved_nodes

&lt; 15     ≥ 15           &lt; 3        ≥ 3

Age       no_recur 125       no_recur 30       no_recur 27

recurrence 39       recurrence 18       recurrence 10

&lt; 40         ≥40

no_recur 4
recurrence 1     no_recur 4

no_rec 4     rec1

# Pruned decision tree for contact lenses recommendation

# **Accuracy and error**

- Accuracy: percentage of correct classifications
  - on the training set
  - on unseen instances

- How accurate is a decision tree when classifying unseen instances
  - An estimate of accuracy on unseen instances can be computed, e.g., by averaging over 4 runs:
    - split the example set into training set (e.g. 70%) and test set (e.g. 30%)
    - induce a decision tree from training set, compute its accuracy on test set

- Error = 1 - Accuracy

- High error may indicate data overfitting

# **Overfitting and accuracy**

- Typical relation between tree size and accuracy



- Question: how to prune optimally?

# **Avoiding overfitting**

- How can we avoid overfitting?
  - Pre-pruning (forward pruning): stop growing the tree e.g., when data split not statistically significant or too few examples are in a split
  - Post-pruning: grow full tree, then post-prune

Pre-pruning

Post-pruning

- forward pruning considered inferior (myopic)
- post pruning makes use of sub trees

# Selected decision/regression tree learners

- Decision tree learners

  - ID3 (Quinlan 1979)

  - CART (Breiman et al. 1984)
  - Assistant (Cestnik et al. 1987)
  - C4.5 (Quinlan 1993), C5 (See5, Quinlan)
  - J48 (available in WEKA)

- Regression tree learners, model tree learners

  - M5, M5P (implemented in WEKA)

# Features of C4.5 and J48

- Implemented as part of the WEKA data mining workbench

- Handling noisy data: post-pruning

- Handling incompletely specified training instances: 'unknown' values (**?**)

  - in learning assign conditional probability of value v: $p(v|C) = p(vC) / p(C)$

  - in classification: follow all branches, weighted by prior prob. of missing attribute values

# Other features of C4.5

- Binarization of attribute values
  - for continuous values select a boundary value maximally increasing the informativity of the attribute: sort the values and try every possible split (done automaticaly)
  - for discrete values try grouping the values until two groups remain *

- 'Majority' classification in NULL leaf (with no corresponding training example)
  - if an example 'falls' into a NULL leaf during classification, the class assigned to this example is the majority class of the parent of the NULL leaf

* the basic C4.5 doesn't support binarisation of discrete attributes, it supports grouping

# Appropriate problems for decision tree learning

- Classification problems: classify an instance into one of a discrete set of possible categories (medical diagnosis, classifying loan applicants, …)
- Characteristics:
  - instances described by attribute-value pairs

    (discrete or real-valued attributes)
  - target function has discrete output values
    (boolean or multi-valued, if real-valued then regression trees)
  - disjunctive hypothesis may be required
  - training data may be noisy
    (classification errors and/or errors in attribute values)
  - training data may contain missing attribute values

# Classifier evaluation

- **Use of induced models**
  - discovery of new patterns, new knowledge
  - classification of new objects
- **Evaluating the quality of induced models**
  - Accuracy, Error = 1 - Accuracy
  - classification accuracy on testing examples = percentage of correctly classified instances
    - split the example set into training set (e.g. 70%) to induce a concept, and test set (e.g. 30%) to test its accuracy
    - more elaborate strategies: 10-fold cross validation, leave-one-out, ...
  - comprehensibility (compactness)
  - information contents (information score), significance

# n-fold cross validation

- A method for accuracy estimation of classifiers
- Partition set D into n disjoint, almost equally-sized folds $T_i$ where $U_i T_i = D$
- **for** i = 1, ..., n **do**
    - form a training set out of n-1 folds: $Di = D \backslash T_i$
    - induce classifier $H_i$ from examples in Di
    - use fold $T_i$ for testing the accuracy of $H_i$
- Estimate the accuracy of the classifier by averaging accuracies over 10 folds $T_i$

# Part II. Predictive DM techniques

- Decision tree learning
- Bayesian Classifier
- Rule learning
- Evaluation

# Bayesian methods

- Bayesian methods – simple but powerful classification methods
  - Based on Bayesian formula

$$p(H \mid D) = \frac{p(D \mid H)}{p(D)} \, p(H)$$

- Main methods:
  - Naive Bayesian classifier
  - Semi-naïve Bayesian classifier
  - Bayesian networks *

* Out of scope of this course

# Naïve Bayesian classifier

- Probability of class, for given attribute values

$$p(c_j \mid v_1...v_n) = p(c_j) \cdot \frac{p(v_1...v_n \mid c_j)}{p(v_1...v_n)}$$

- For all $C_j$ compute probability $p(C_j)$, given values $v_i$ of all attributes describing the example which we want to classify (assumption: conditional independence of attributes, when estimating $p(C_j)$ and $p(C_j \mid v_i)$)

$$p(c_j \mid v_1...v_n) \approx p(c_j) \cdot \prod_i \frac{p(c_j \mid v_i)}{p(c_j)}$$

- Output $C_{MAX}$ with maximal posterior probability of class:

$$C_{MAX} = \arg\max_{Cj} p(c_j \mid v_1...v_n)$$

# **Semi-naïve Bayesian classifier**

- Naive Bayesian estimation of probabilities (reliable)

$$\frac{p(c_j \mid v_i)}{p(c_j)} \cdot \frac{p(c_j \mid v_k)}{p(c_j)}$$

- Semi-naïve Bayesian estimation of probabilities (less reliable)

$$\frac{p(c_j \mid v_i, v_k)}{p(c_j)}$$

# **Probability estimation**

- ## Relative frequency:

$$p(c_j) = \frac{n(c_j)}{N} , \; p(c_j \mid v_i) = \frac{n(c_j, v_i)}{n(v_i)}$$

j = 1 . . k, for k classes

[6+,1-] (7) = 6/7                     problems with small samples

[2+,0-] (2) = 2/2 = 1

- ## Laplace estimate (prior probability):

$$p(c_j) = \frac{n(c_j) + 1}{N + k}$$

assumes uniform prior

distribution of k classes

[6+,1-] (7) = 6+1 / 7+2 = 7/9

[2+,0-] (2) = 2+1 / 2+2 = 3/4

# **Probability estimation**

- Relative frequency:

$$p(c_j) = \frac{n(c_j)}{N}, \ p(c_j \mid v_i) = \frac{n(c_j, v_i)}{n(v_i)} \qquad \text{j = 1 . . k, for k classes}$$

- Prior probability: Laplace law

$$p(c_j) = \frac{n(c_j) + 1}{N + k}$$

- m-estimate:

$$p(c_j) = \frac{n(c_j) + m \cdot p_a(c_j)}{N + m}$$

# Probability estimation: intuition

- Experiment with N trials, n successful
- Estimate probability of success of next trial
- **Relative frequency: n/N**
  - reliable estimate when number of trials is large
  - Unreliable when number of trials is small, e.g., 1/1=1
- **Laplace: (n+1)/(N+2), (n+1)/(N+k),** k classes
  - Assumes uniform distribution of classes
- **m-estimate: $(n+m.p_a)/(N+m)$**
  - Prior probability of success $p_a$, parameter m (weight of prior probability, i.e., number of 'virtual' examples )

# **Explanation of Bayesian classifier**

- Based on information theory
  - Expected number of bits needed to encode a message = optimal code length -log p for a message, whose probability is p (*)

- Explanation based of the sum of information gains of individual attribute values $v_i$ (Kononenko and Bratko 1991, Kononenko 1993)

$$-\log(p(c_j \mid v_1...v_n)) =$$

$$= -\log(p(c_j)) - \sum_{i=1}^{n}(-\log p(c_j) + \log(p(c_j \mid v_i))$$

\* log p denotes binary logarithm

# Example of explanation of semi-naïve Bayesian classifier

Hip surgery prognosis

Class = no ("no complications", most probable class, 2 class problem)

| Attribute value | For decision (bit) | Against (bit) |
|---|---|---|
| Age = 70-80 | 0.07 | |
| Sex = Female | | -0.19 |
| Mobility before injury = Fully mobile | 0.04 | |
| State of health before injury = Other | 0.52 | |
| Mechanism of injury = Simple fall | | -0.08 |
| Additional injuries = None | 0 | |
| Time between injury and operation > 10 days | 0.42 | |
| Fracture classification acc. To Garden = Garden III | | -0.3 |
| Fracture classification acc. To Pauwels = Pauwels III | | -0.14 |
| Transfusion = Yes | 0.07 | |
| Antibiotic profilaxies = Yes | | -0.32 |
| Hospital rehabilitation = Yes | 0.05 | |
| General complications = None | | 0 |
| **Combination:** | 0.21 | |
|    Time between injury and examination < 6 hours | | |
|    AND Hospitalization time between 4 and 5 weeks | | |
| **Combination:** | 0.63 | |
|   Therapy = Artroplastic AND anticoagulant therapy = Yes | | |

# Visualization of information gains for/against C_i

# Naïve Bayesian classifier

- Naïve Bayesian classifier can be used
  - when we have sufficient number of training examples for reliable probability estimation
- It achieves good classification accuracy
  - can be used as 'gold standard' for comparison with other classifiers
- Resistant to noise (errors)
  - Reliable probability estimation
  - Uses all available information
- Successful in many application domains
  - Web page and document classification
  - Medical diagnosis and prognosis, …

# Improved classification accuracy due to using m-estimate

|  | Primary tumor | Breast cancer | thyroid | Rheumatology |
|---|---|---|---|---|
| #instan | 339 | 288 | 884 | 355 |
| #class | 22 | 2 | 4 | 6 |
| #attrib | 17 | 10 | 15 | 32 |
| #values | 2 | 2.7 | 9.1 | 9.1 |
| majority | 25% | 80% | 56% | 66% |
| entropy | 3.64 | 0.72 | 1.59 | 1.7 |

|  | Relative freq. | m-estimate |
|---|---|---|
| Primary tumor | 48.20% | 52.50% |
| Breast cancer | 77.40% | 79.70% |
| hepatitis | 58.40% | 90.00% |
| lymphography | 79.70% | 87.70% |

# Part II. Predictive DM techniques

- Decision tree learning
- Bayesian Classifier

  Rule learning
- Evaluation

# Rule Learning

| Person | Age | Spect. presc. | Astigm. | Tear prod. | Lenses |
|--------|-----|---------------|---------|------------|--------|
| O1 | young | myope | no | reduced | NONE |
| O2 | young | myope | no | normal | SOFT |
| O3 | young | myope | yes | reduced | NONE |
| O4 | young | myope | yes | normal | HARD |
| O5 | young | hypermetrope | no | reduced | NONE |
| O6-O13 | ... | ... | ... | ... | ... |
| O14 | pre-presbyo | hypermetrope | no | normal | SOFT |
| O15 | pre-presbyo | hypermetrope | yes | reduced | NONE |
| O16 | pre-presbyo | hypermetrope | yes | normal | NONE |
| O17 | presbyopic | myope | no | reduced | NONE |
| O18 | presbyopic | myope | no | normal | NONE |
| O19-O23 | ... | ... | ... | ... | ... |
| O24 | presbyopic | hypermetrope | yes | normal | NONE |

data

knowledge discovery
from data

Rule learning

Model: a set of rules

Patterns: individual rules

**Given:** transaction data table, relational database (a set of objects, described by attribute values)
**Find:** a classification model in the form of a set of rules;
or a set of interesting patterns in the form of individual rules

# Rule set representation

- Rule base is a disjunctive set of conjunctive rules
- Standard form of rules:

  IF Condition THEN Class

  Class IF Conditions

  Class $\leftarrow$ Conditions

- Form of CN2 rules:

  IF Conditions THEN MajClass [ClassDistr]

- Rule base:   {R1, R2, R3, …, DefaultRule}

# Contact lens data: Classification rules

**Type of task:** prediction and classification
**Hypothesis language:** rules X ➔ C,  if X then C
        X conjunction of attribute values, C class

tear production=reduced → lenses=NONE
tear production=normal & astigmatism=yes &
        spect. pre.=hypermetrope → lenses=NONE
tear production=normal & astigmatism=no → lenses=SOFT
tear production=normal & astigmatism=yes &
        spect. pre.=myope → lenses=HARD
DEFAULT lenses=NONE

# **Rule learning**

- Two rule learning approaches:
  - Learn decision tree, convert to rules
  - Learn set/list of rules
    - Learning an unordered set of rules
    - Learning an ordered list of rules
- Heuristics, overfitting, pruning

# Contact lenses: convert decision tree to an unordered rule set



tear production=reduced **=>** lenses=NONE [S=0,H=0,N=12]
tear production=normal & astigmatism=yes & spect. pre.=hypermetrope **=>**
lenses=NONE  [S=0,H=1,N=2]
tear production=normal & astigmatism=no **=>** lenses=SOFT     [S=5,H=0,N=1]
tear production=normal & astigmatism=yes & spect. pre.=myope **=>** lenses=HARD
[S=0,H=3,N=2]
DEFAULT lenses=NONE                Order independent rule set (may overlap)

# Contact lenses: convert decision tree to decision list



IF tear production=reduced THEN lenses=NONE
ELSE /*tear production=normal*/
  IF astigmatism=no THEN lenses=SOFT
  ELSE /*astigmatism=yes*/
    IF spect. pre.=myope THEN lenses=HARD
    ELSE /* spect.pre.=hypermetrope*/
      lenses=NONE                    Ordered (order dependent) rule list

# Converting decision tree to rules, and rule post-pruning (Quinlan 1993)

- Very frequently used method, e.g., in C4.5 and J48

- Procedure:
  - grow a full tree (allowing overfitting)
  - convert the tree to an equivalent set of rules
  - prune each rule independently of others
  - sort final rules into a desired sequence for use

# Concept learning: Task reformulation for rule learning: (pos. vs. neg. examples of Target class)

| Person | Age | Spect. presc. | Astigm. | Tear prod. | Lenses |
|--------|-----|---------------|---------|------------|--------|
| O1 | 17 | myope | no | reduced | NO |
| O2 | 23 | myope | no | normal | YES |
| O3 | 22 | myope | yes | reduced | NO |
| O4 | 27 | myope | yes | normal | YES |
| O5 | 19 | hypermetrope | no | reduced | NO |
| O6-O13 | … | … | … | … | … |
| O14 | 35 | hypermetrope | no | normal | YES |
| O15 | 43 | hypermetrope | yes | reduced | NO |
| O16 | 39 | hypermetrope | yes | normal | NO |
| O17 | 54 | myope | no | reduced | NO |
| O18 | 62 | myope | no | normal | NO |
| O19-O23 | … | … | … | … | … |
| O24 | 56 | hypermetrope | yes | normal | NO |

# Original covering algorithm (AQ, Michalski 1969,86)

**Given** examples of N classes $C_1, \ldots, C_N$

**for** each class Ci **do**

– Ei := Pi U Ni (Pi pos., Ni neg.)

– RuleBase(Ci) := empty

– **repeat {learn-set-of-rules}**

  • **learn-one-rule** R covering some positive examples and no negatives

  • add R to RuleBase(Ci)

  • delete from Pi all pos. ex. covered by R

– **until** Pi = empty

# Covering algorithm

Positive examples

Negative examples

# Covering algorithm



Rule1: Cl=+ ← Cond2 AND Cond3

Positive examples

Negative examples

# Covering algorithm

Rule1: Cl=+ ← Cond2 AND Cond3

Positive examples

Negative examples

# Covering algorithm

Rule1: Cl=+ ← Cond2 AND Cond3

Positive examples

Negative examples



Rule2: Cl=+ ← Cond8 AND Cond6

# Probability estimates

- **Relative frequency** :
  - problems with small samples

$$p(Class \mid Cond) =$$

$$= \frac{n(Class.Cond)}{n(Cond)}$$

[6+,1-] (7) = 6/7
[2+,0-] (2) = 2/2 = 1

- **Laplace estimate** :
  - assumes uniform prior distribution of k classes

$$= \frac{n(Class.Cond)+1}{n(Cond)+k} \quad k = 2$$

[6+,1-] (7) = 6+1 / 7+2 = 7/9
[2+,0-] (2) = 2+1 / 2+2 = 3/4

# Learn-one-rule: search heuristics

- Assume a two-class problem

- Two classes (+,-),  learn rules for + class (Cl).

- Search for specializations R' of a rule R = Cl ← Cond from the RuleBase.

- Specializarion R' of rule R = Cl ← Cond

    has the form    R' = Cl ← Cond & Cond'

- Heuristic search for rules: find the 'best' Cond' to be added to the current rule R, such that rule accuracy is improved, e.g., such that $Acc(R') > Acc(R)$

  - where the expected **classification accuracy** can be estimated as $A(R) = p(Cl|Cond)$

# Learn-one-rule:
# Greedy vs. beam search

- learn-one-rule by greedy general-to-specific search, at each step selecting the `best' descendant, no backtracking
  - e.g., the best descendant of the initial rule

    lenses=NONE ←

  - is rule  lenses=NONE ← tear production=reduced


- beam search: maintain a list of k best candidates at each step; descendants (specializations) of each of these k candidates are generated, and the resulting set is again reduced to k best candidates

# What is "high" rule accuracy (rule precision) ?

- Rule evaluation measures:
  - aimed at maximizing classification accuracy
  - minimizing Error = 1 - Accuracy
  - avoiding overfitting
- BUT: Rule accuracy/precision should be traded off against the "default" accuracy/precision of the rule **Cl ←true**
  - 68% accuracy is OK if there are 20% examples of that class in the training set, but bad if there are 80%
- **Relative accuracy** *(relative precision)*
  - RAcc(Cl ←Cond) = p(Cl | Cond) – p(Cl)

# Learn-one-rule: search heuristics

- Assume two classes (+,-), learn rules for + class (Cl). Search for specializations of one rule $R = Cl \leftarrow Cond$ from RuleBase.
- Expected **classification accuracy**: $A(R) = p(Cl|Cond)$
- **Informativity** (info needed to specify that example covered by Cond belongs to Cl): $I(R) = -\log_2 p(Cl|Cond)$
- **Accuracy gain** (increase in expected accuracy):
    $AG(R',R) = p(Cl|Cond') - p(Cl|Cond)$
- **Information gain** (decrease in the information needed):
    $IG(R',R) = \log_2 p(Cl|Cond') - \log_2 p(Cl|Cond)$
- **Weighted** measures favoring more general rules: WAG, WIG
    $WAG(R',R) =$
        $p(Cond')/p(Cond) \cdot (p(Cl|Cond') - p(Cl|Cond))$
- **Weighted relative accuracy** trades off coverage and relative accuracy $WRAcc(R) = p(Cond).(p(Cl|Cond) - p(Cl))$

# Ordered set of rules: if-then-else rules

- rule  Class IF Conditions is learned by first determining Conditions and then Class
- **Notice:** mixed sequence of classes C1, …, Cn in RuleBase
- **But: ordered** execution when classifying a new instance: rules are sequentially tried and the first rule that `fires' (covers the example) is used for classification
- **Decision list {R1, R2, R3, …, D}:** rules Ri are interpreted as **if-then-else** rules
- If no rule fires, then DefaultClass (majority class in $E_{cur}$)

# Sequential covering algorithm

- RuleBase := empty
- $E_{cur}$ := E
- **repeat**
  - learn-one-rule R
  - RuleBase := RuleBase U R
  - $E_{cur}$ := $E_{cur}$ - {examples covered and correctly classified by R}     **(DELETE ONLY POS. EX.!)**
  - **until** performance(R, $E_{cur}$) < ThresholdR
- RuleBase := sort RuleBase by performance(R,E)
- return RuleBase

# Learn ordered set of rules (CN2, Clark and Niblett 1989)

- RuleBase := empty
- $E_{cur}$:= E
- **repeat**
  - learn-one-rule R
  - RuleBase := RuleBase U R
  - $E_{cur}$ := $E_{cur}$ - {all examples covered by R}
  **(NOT ONLY POS. EX.!)**
- **until** performance(R, $E_{cur}$) < ThresholdR
- RuleBase := sort RuleBase by performance(R,E)
- RuleBase := RuleBase U DefaultRule($E_{cur}$)

# Learn-one-rule: Beam search in CN2

- Beam search in CN2 learn-one-rule algo.:
  - construct BeamSize of best rule bodies (conjunctive conditions) that are statistically significant
  - BestBody - min. entropy of examples covered by Body
  - construct best rule R := Head $\leftarrow$ BestBody by adding majority class of examples covered by BestBody in rule Head

- performance (R, $E_{cur}$) : - Entropy($E_{cur}$)
  - performance(R, $E_{cur}$) < ThresholdR (neg. num.)
  - Why? Ent. > t is bad, Perf. = -Ent < -t is bad

# **Variations**

- Sequential vs. simultaneous covering of data (as in TDIDT): choosing between attribute-values vs. choosing attributes
- Learning rules vs. learning decision trees and converting them to rules
- Pre-pruning vs. post-pruning of rules
- What statistical evaluation functions to use
- Probabilistic classification

- Best performing rule learning algorithm: Ripper
- JRip implementation of Ripper in WEKA, available in ClowdFlows

# Probabilistic classification

- In the ordered case of standard CN2 rules are interpreted in an `IF-THEN-ELSE` fashion, and the first fired rule assigns the class.
- In the unordered case all rules are tried and all rules which fire are collected. If a clash occurs, a probabilistic method is used to resolve the clash.
- A simplified example:
  1. tear production=reduced **=>** lenses=NONE [S=0,H=0,N=12]
  2. tear production=normal & astigmatism=yes & spect. pre.=hypermetrope **=>** lenses=NONE  [S=0,H=1,N=2]
  3. tear production=normal & astigmatism=no **=>** lenses=SOFT [S=5,H=0,N=1]
  4. tear production=normal & astigmatism=yes & spect. pre.=myope **=>** lenses=HARD [S=0,H=3,N=2]
  5. DEFAULT lenses=NONE

Suppose we want to classify a person with normal tear production and astigmatism. Two rules fire: rule 2 with coverage [S=0,H=1,N=2] and rule 4 with coverage [S=0,H=3,N=2]. The classifier computes total coverage as [S=0,H=4,N=4], resulting in probabilistic classification into class H with probability 0.5 and N with probability 0.5. In this case, the clash can not be resolved, as both probabilities are equal.

# Part II. Predictive DM techniques

- Decision tree learning
- Bayesian Classifier
- Rule learning

→ Evaluation

# Classifier evaluation

- Accuracy and Error
- n-fold cross-validation
- Confusion matrix
- ROC

# **Evaluating hypotheses**

- **Use of induced hypotheses**
  - discovery of new patterns, new knowledge
  - classification of new objects
- **Evaluating the quality of induced hypotheses**
  - Accuracy, Error = 1 - Accuracy
  - classification accuracy on testing examples = percentage of correctly classified instances
    - split the example set into training set (e.g. 70%) to induce a concept, and test set (e.g. 30%) to test its accuracy
    - more elaborate strategies: 10-fold cross validation, leave-one-out, ...
  - comprehensibility (compactness)
  - information contents (information score), significance

# **n-fold cross validation**

- A method for accuracy estimation of classifiers
- Partition set D into n disjoint, almost equally-sized folds $T_i$ where $U_i T_i = D$
- **for** i = 1, ..., n **do**
  - form a training set out of n-1 folds: $Di = D \backslash T_i$
  - induce classifier $H_i$ from examples in Di
  - use fold $T_i$ for testing the accuracy of $H_i$
- Estimate the accuracy of the classifier by averaging accuracies over 10 folds $T_i$

- Partition

- Partition

- Train

$D$

$T_1$  $T_2$  $T_3$

$D \backslash T_1 = D_1$  $D \backslash T_2 = D_2$  $D \backslash T_3 = D_3$

- Partition

- Train

$D \setminus T_1 = D_1$  $D \setminus T_2 = D_2$  $D \setminus T_3 = D_3$

• Partition

• Train

• Test

D

$T_1$ $T_2$ $T_3$

$D \backslash T_1 = D_1$  $D \backslash T_2 = D_2$  $D \backslash T_3 = D_3$

$T_1$ $T_2$ $T_3$

# Confusion matrix and rule (in)accuracy

- Accuracy of a classifier is measured as TP+TN / N.
- Suppose two rules are both 80% accurate on an evaluation dataset, are they always equally good?
  - e.g., Rule 1 correctly classifies 40 out of 50 positives and 40 out of 50 negatives; Rule 2 correctly classifies 30 out of 50 positives and 50 out of 50 negatives
  - on a test set which has more negatives than positives, Rule 2 is preferable;
  - on a test set which has more positives than negatives, Rule 1 is preferable; unless…
  - …the proportion of positives becomes so high that the 'always positive' predictor becomes superior!
- Conclusion: classification accuracy is not always an appropriate rule quality measure

# Confusion matrix

|  | Predicted positive | Predicted negative |  |
|---|---|---|---|
| Positive examples | **True positives** | **False negatives** |  |
| Negative examples | **False positives** | **True negatives** |  |
|  |  |  |  |

- also called *contingency table*

## Classifier 1

|  | Predicted positive | Predicted negative |  |
|---|---|---|---|
| Positive examples | **40** | **10** | 50 |
| Negative examples | **10** | **40** | 50 |
|  | 50 | 50 | 100 |

## Classifier 2

|  | Predicted positive | Predicted negative |  |
|---|---|---|---|
| Positive examples | **30** | **20** | 50 |
| Negative examples | **0** | **50** | 50 |
|  | 30 | 70 | 100 |

# ROC space

- ***True positive rate*** =
  #true pos. / #pos.
  - $TPr_1 = 40/50 = 80\%$
  - $TPr_2 = 30/50 = 60\%$
- ***False positive rate***
  = #false pos. / #neg.
  - $FPr_1 = 10/50 = 20\%$
  - $FPr_2 = 0/50 = 0\%$
- ***ROC space*** has
  - FPr on X axis
  - TPr on Y axis

## Classifier 1

|  | Predicted positive | Predicted negative |  |
|---|---|---|---|
| Positive examples | **40** | **10** | 50 |
| Negative examples | **10** | **40** | 50 |
|  | 50 | 50 | 100 |

## Classifier 2

|  | Predicted positive | Predicted negative |  |
|---|---|---|---|
| Positive examples | **30** | **20** | 50 |
| Negative examples | **0** | **50** | 50 |
|  | 30 | 70 | 100 |

# The ROC space

# The ROC convex hull

# Course Outline

## I. Introduction
- – Data Mining and KDD process
- – Introduction to Data Mining
- – Data Mining platforms

## II. Predictive DM Techniques
- – Decision Tree learning
- - Bayesian classifier
- – Classification rule learning
- – Classifier Evaluation
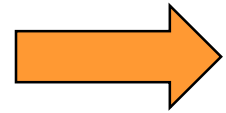
## III. Regression

## IV. Descriptive DM
- – Predictive vs. descriptive induction
- – Subgroup discovery
- – Association rule learning Hierarchical clustering

## V. Relational Data Mining
- – RDM and Inductive Logic Programming
- – Propositionalization
- – Semantic data mining

## VI. Advanced Topics

# III. Predictive DM – Regression

- often referred to as estimation or regression
- data are objects, characterized with attributes (discrete or continuous), classes of objects are continuous (numeric)
- given objects described with attribute values, induce a model to predict the numeric class value
- regression trees, linear and logistic regression, ANN, kNN, ...

# Estimation/regression example: Customer data

| Customer | Gender | Age | Income | Spent | |
|----------|--------|-----|--------|-------|---|
| c1 | male | 30 | 214000 | 18800 | |
| c2 | female | 19 | 139000 | 15100 | |
| c3 | male | 55 | 50000 | 12400 | |
| c4 | female | 48 | 26000 | 8600 | |
| c5 | male | 63 | 191000 | 28100 | |
| O6-O13 | … | … | … | … | |
| c14 | female | 61 | 95000 | 18100 | |
| c15 | male | 56 | 44000 | 12000 | |
| c16 | male | 36 | 102000 | 13800 | |
| c17 | female | 57 | 215000 | 29300 | |
| c18 | male | 33 | 67000 | 9700 | |
| c19 | female | 26 | 95000 | 11000 | |
| c20 | female | 55 | 214000 | 28800 | |

# Customer data: regression tree



**In the nodes one usually has
Predicted value +- st. deviation**

# Predicting algal biomass: regression tree

# Predicting algal biomass: regression tree

| Regression | Classification |
|---|---|
| **Data**: attribute-value description | |
| **Target variable**: <br> Continuous | **Target variable**: <br> Categorical (nominal) |
| **Evaluation**: cross validation, separate test set, … | |
| **Error**: <br> MSE, MAE, RMSE, … | **Error**: <br> 1-accuracy |
| **Algorithms**: <br> Linear regression, regression trees,… | **Algorithms**: <br> Decision trees, Naïve Bayes, … |
| **Baseline predictor**: <br> Mean of the target variable | **Baseline predictor**: <br> Majority class |

# Example regression problem

- data about 80 people: Age and Height



| Age | Height |
|-----|--------|
| 3   | 1.03   |
| 5   | 1.19   |
| 6   | 1.26   |
| 9   | 1.39   |
| 15  | 1.69   |
| 19  | 1.67   |
| 22  | 1.86   |
| 25  | 1.85   |
| 41  | 1.59   |
| 48  | 1.60   |
| 54  | 1.90   |
| 71  | 1.82   |
| …   | …      |

# Test set

| Age | Height |
|-----|--------|
| 2   | 0.85   |
| 10  | 1.4    |
| 35  | 1.7    |
| 70  | 1.6    |

# Baseline numeric model

- Average of the target variable

# Baseline numeric predictor

- Average of the target variable is 1.63



| Age | Height | Baseline |
|-----|--------|----------|
| 2 | 0.85 | |
| 10 | 1.4 | |
| 35 | 1.7 | |
| 70 | 1.6 | |

# Linear Regression Model

Height =    0.0056 * Age + 1.4181

# Regression tree

# kNN – K nearest neighbors

- Looks at K closest examples (by age) and predicts the average of their target variable
- K=3

# Which predictor is the best?

| Age | Height | Baseline | Linear regression | Regression tree | Model tree | kNN |
|---|---|---|---|---|---|---|
| 2 | 0.85 | 1.63 | 1.43 | 1.39 | 1.20 | 1.01 |
| 10 | 1.4 | 1.63 | 1.47 | 1.46 | 1.47 | 1.51 |
| 35 | 1.7 | 1.63 | 1.61 | 1.71 | 1.71 | 1.67 |
| 70 | 1.6 | 1.63 | 1.81 | 1.71 | 1.75 | 1.81 |

# Course Outline

**I. Introduction**
- – Data Mining and KDD process
- – Introduction to Data Mining
- – Data Mining platforms

**II. Predictive DM Techniques**
- – Decision Tree learning
- - Bayesian classifier
- – Classification rule learning
- – Classifier Evaluation

**III. Regression**

**IV. Descriptive DM**
- – Predictive vs. descriptive induction
- – Subgroup discovery
- – Association rule learning Hierarchical clustering

**V. Relational Data Mining**
- – RDM and Inductive Logic Programming
- – Propositionalization
- – Semantic data mining

**VI. Advanced Topics**

# Part IV. Descriptive DM techniques

→ • Predictive vs. descriptive induction

• Subgroup discovery

• Association rule learning

• Hierarchical clustering

# Descriptive DM:
# Subgroup discovery example - Customer data

| Customer | Gender | Age | Income | Spent | BigSpender |
|----------|--------|-----|--------|-------|------------|
| c1 | male | 30 | 214000 | 18800 | yes |
| c2 | female | 19 | 139000 | 15100 | yes |
| c3 | male | 55 | 50000 | 12400 | no |
| c4 | female | 48 | 26000 | 8600 | no |
| c5 | male | 63 | 191000 | 28100 | yes |
| O6-O13 | … | … | … | … | … |
| c14 | female | 61 | 95000 | 18100 | yes |
| c15 | male | 56 | 44000 | 12000 | no |
| c16 | male | 36 | 102000 | 13800 | no |
| c17 | female | 57 | 215000 | 29300 | yes |
| c18 | male | 33 | 67000 | 9700 | no |
| c19 | female | 26 | 95000 | 11000 | no |
| c20 | female | 55 | 214000 | 28800 | yes |

# Customer data: Subgroup discovery

**Type of task:** description (pattern discovery)
**Hypothesis language:** rules **X ➔ Y,** if X then Y
X is conjunctions of items, Y is target class

Age $>$ 52 & Sex = male ➔ BigSpender = no

Age $>$ 52 & Sex = male & Income $\leq$ 73250
➔ BigSpender = no

# Descriptive DM: Association rule learning example - Customer data

| Customer | Gender | Age | Income | Spent | BigSpender |
|----------|--------|-----|--------|-------|------------|
| c1 | male | 30 | 214000 | 18800 | yes |
| c2 | female | 19 | 139000 | 15100 | yes |
| c3 | male | 55 | 50000 | 12400 | no |
| c4 | female | 48 | 26000 | 8600 | no |
| c5 | male | 63 | 191000 | 28100 | yes |
| O6-O13 | … | … | … | … | … |
| c14 | female | 61 | 95000 | 18100 | yes |
| c15 | male | 56 | 44000 | 12000 | no |
| c16 | male | 36 | 102000 | 13800 | no |
| c17 | female | 57 | 215000 | 29300 | yes |
| c18 | male | 33 | 67000 | 9700 | no |
| c19 | female | 26 | 95000 | 11000 | no |
| c20 | female | 55 | 214000 | 28800 | yes |

# Customer data: Association rules

**Type of task:** description (pattern discovery)

**Hypothesis language:** rules **X ➜ Y,** if X then Y

   X, Y conjunctions of items

1. Age $>$ 52 & BigSpender = no ➜ Sex = male
2. Age $>$ 52 & BigSpender = no ➜

      Sex = male & Income $\leq$ 73250

3. Sex = male & Age $>$ 52 & Income $\leq$ 73250 ➜

      BigSpender = no

# Descriptive DM:
# Clustering and association rule learning example - Customer data

| Customer | Gender | Age | Income | Spent | BigSpender |
|----------|--------|-----|--------|-------|------------|
| c1 | male | 30 | 214000 | 18800 | yes |
| c2 | female | 19 | 139000 | 15100 | yes |
| c3 | male | 55 | 50000 | 12400 | no |
| c4 | female | 48 | 26000 | 8600 | no |
| c5 | male | 63 | 191000 | 28100 | yes |
| O6-O13 | … | … | … | … | … |
| c14 | female | 61 | 95000 | 18100 | yes |
| c15 | male | 56 | 44000 | 12000 | no |
| c16 | male | 36 | 102000 | 13800 | no |
| c17 | female | 57 | 215000 | 29300 | yes |
| c18 | male | 33 | 67000 | 9700 | no |
| c19 | female | 26 | 95000 | 11000 | no |
| c20 | female | 55 | 214000 | 28800 | yes |

# Predictive vs. descriptive induction

- **Predictive induction:** Inducing classifiers for solving classification and prediction tasks,
  - Classification rule learning, Decision tree learning, ...
  - Bayesian classifier, ANN, SVM, ...
  - Data analysis through hypothesis generation and testing
- **Descriptive induction:** Discovering interesting regularities in the data, uncovering patterns, ... for solving KDD tasks
  - Symbolic clustering, Association rule learning, Subgroup discovery, ...
  - Exploratory data analysis

# **Descriptive DM**

- Often used for preliminary explanatory data analysis

- User gets feel for the data and its structure

- Aims at deriving descriptions of characteristics of the data

- Visualization and descriptive statistical techniques can be used

# Predictive vs. descriptive DM: Summary from a rule learning perspective

- **Predictive DM:** Induces **rulesets** acting as classifiers for solving classification and prediction tasks
- **Descriptive DM:** Discovers **individual rules** describing interesting regularities in the data

- **Therefore:** Different goals, different heuristics, different evaluation criteria

# Descriptive DM

- **Description**
  - Data description and summarization: describe elementary and aggregated data characteristics (statistics, …)
  - Dependency analysis:
    - describe associations, dependencies, …
    - discovery of properties and constraints

- **Segmentation**
  - Clustering: separate objects into subsets according to distance and/or similarity (clustering, SOM, visualization, ...)
  - Subgroup discovery: find unusual subgroups that are significantly different from the majority (deviation detection w.r.t. overall class distribution)

# Part IV. Descriptive DM techniques

- Predictive vs. descriptive induction
- Subgroup discovery
- Association rule learning
- Hierarchical clustering

# Subgroup Discovery

| Person | Age | Spect. presc. | Astigm. | Tear prod. | Lenses |
|--------|-----|---------------|---------|-----------|--------|
| O1 | 17 | myope | no | reduced | NO |
| O2 | 23 | myope | no | normal | YES |
| O3 | 22 | myope | yes | reduced | NO |
| O4 | 27 | myope | yes | normal | YES |
| O5 | 19 | hypermetrope | no | reduced | NO |
| O6-O13 | ... | ... | ... | ... | ... |
| O14 | 35 | hypermetrope | no | normal | YES |
| O15 | 43 | hypermetrope | yes | reduced | NO |
| O16 | 39 | hypermetrope | yes | normal | NO |
| O17 | 54 | myope | no | reduced | NO |
| O18 | 62 | myope | no | normal | NO |
| O19-O23 | ... | ... | ... | ... | ... |
| O24 | 56 | hypermetrope | yes | normal | NO |

Subgroup Discovery

Class YES          Class NO

2

1      3

- A task in which individual interpretable patterns in the form of rules are induced from data, labeled by a predefined property of interest.

- SD algorithms learn several independent rules that describe groups of target class examples
    - subgroups must be large and significant

# Classification versus Subgroup Discovery

- **Classification (predictive induction) - constructing sets of classification rules**
  - aimed at learning a model for classification or prediction
  - rules are dependent
- **Subgroup discovery (descriptive induction) – constructing individual subgroup describing rules**
  - aimed at finding interesting patterns in target class examples
    - large subgroups (high target class coverage)
    - with significantly different distribution of target class examples (high TP/FP ratio, high significance, high WRAcc
  - each rule (pattern) is an independent chunk of knowledge

# Classification versus Subgroup discovery



Class YES          Class NO

# Subgroup discovery in
# High CHD Risk Group Detection

**Input:** Patient records described by anamnestic, laboratory and ECG attributes

**Task**: Find and characterize population subgroups with high CHD risk (large enough, distributionaly unusual)

From **best induced descriptions**, five were selected by the expert as **most actionable** for CHD risk screening (by GPs):

    high-CHD-risk $\leftarrow$ male & pos. fam. history & age > 46

    high-CHD-risk $\leftarrow$ female & bodymassIndex > 25 & age > 63

    high-CHD-risk $\leftarrow$ ...

    high-CHD-risk $\leftarrow$ ...

    high-CHD-risk $\leftarrow$ ...

(Gamberger & Lavrač, JAIR 2002)

# Subgroup Discovery: Medical Use Case

- **Find and characterize population subgroups with high risk for coronary heart disease (CHD)** (Gamberger, Lavrač, Krstačić)

- **A1** for males: **principal risk factors**

    CHD ← pos. fam. history & age > 46

- **A2** for females: **principal risk factors**

    CHD ← bodyMassIndex > 25 & age >63

- **A1, A2** (anamnestic info only), **B1, B2** (an. and physical examination), **C1** (an., phy. and ECG)

- **A1: supporting factors** (found by statistical analysis): psychosocial stress, as well as cigarette smoking, hypertension and overweight

# Subgroup discovery in functional genomics

- Functional genomics is a typical scientific discovery domain, studying genes and their functions

- Very large number of attributes (genes)

- Interesting subgroup describing patterns discovered by SD algorithm

CancerType = Leukemia

IF  KIAA0128 = DIFF. EXPRESSED

AND  prostoglandin d2 synthase = NOT_ DIFF.  EXPRESSED

- Interpretable by biologists

    D. Gamberger, N. Lavrač, F. Železný, J. Tolar

    Journal of Biomedical Informatics 37(5):269-284, 2004

# Subgroups vs. classifiers

- Classifiers:
  - Classification rules aim at pure subgroups
  - A set of rules forms a domain model
- Subgroups:
  - Rules describing subgroups aim at significantly higher proportion of positives
  - Each rule is an independent chunk of knowledge
- Link
  - SD can be viewed as cost-sensitive classification
  - Instead of *FNcost* we aim at increased *TPprofit*

# Classification Rule Learning for Subgroup Discovery: Deficiencies

- Only first few rules induced by the covering algorithm have sufficient support (coverage)

- Subsequent rules are induced from smaller and strongly biased example subsets (pos. examples not covered by previously induced rules), which hinders their ability to detect population subgroups

- 'Ordered' rules are induced and interpreted sequentially as a **if-then-else** decision list

# CN2-SD: Adapting CN2 Rule Learning to Subgroup Discovery

- Weighted covering algorithm

- Weighted relative accuracy (WRAcc) search heuristics, with added example weights

- Probabilistic classification

- Evaluation with different interestingness measures

# CN2-SD: CN2 Adaptations

- General-to-specific search  (beam search) for best rules
- Rule quality measure:
  - CN2: Laplace: $\text{Acc}(\text{Class} \leftarrow \text{Cond}) =$
    $$= p(\text{Class}|\text{Cond}) = (n_c + 1) / (n_{rule} + k)$$
  - CN2-SD: Weighted Relative Accuracy

    $\text{WRAcc}(\text{Class} \leftarrow \text{Cond}) =$
    $$p(\text{Cond}) \, (p(\text{Class}|\text{Cond}) - p(\text{Class}))$$
- Weighted covering approach (example weights)
- Significance testing (likelihood ratio statistics)
- Output: Unordered rule sets (probabilistic classification)

# CN2-SD: Weighted Covering

- Standard covering approach:
  covered examples are deleted from current training set

- Weighted covering approach:
  - weights assigned to examples
  - covered pos. examples are re-weighted:
    in all covering loop iterations, store
    count i how many times (with how many
    rules induced so far) a pos. example has
    been covered: w(e,i), w(e,0)=1
    - **Additive weights: `w(e,i) = 1/(i+1`**
      `w(e,i)` **– pos. example e being covered `i` times**

# Subgroup Discovery

Positive examples

Negative examples

# Subgroup Discovery

Rule1: Cl=+ ← Cond6 AND Cond2

Positive examples

Negative examples

# Subgroup Discovery

Positive examples

Negative examples



Rule2: Cl=+ ← Cond3 AND Cond4

# Subgroup Discovery



Positive examples

Negative examples

# CN2-SD: Weighted WRAcc Search Heuristic

- **Weighted relative accuracy (WRAcc) search heuristics, with added example weights**

  $WRAcc(Cl \leftarrow Cond) = p(Cond) (p(Cl|Cond) - p(Cl))$

  increased coverage, decreased # of rules, approx. equal accuracy (PKDD-2000)

- In WRAcc computation, probabilities are estimated with relative frequencies, adapt:

  $WRAcc(Cl \leftarrow Cond) = p(Cond) (p(Cl|Cond) - p(Cl)) =$

  $n'(Cond)/N' ( n'(Cl.Cond)/n'(Cond) - n'(Cl)/N' )$

  - N' : sum of weights of examples
  - n'(Cond) : sum of weights of all covered examples
  - n'(Cl.Cond) : sum of weights of all correctly covered examples

# SD algorithms in the Orange DM Platform

- **Orange** data mining toolkit
  - classification and subgroup discovery algorithms
  - data mining workflows
  - visualization



- **SD Algorithms in Orange**
  - SD (Gamberger & Lavrač, JAIR 2002)
  - Apriori-SD (Kavšek & Lavrač, AAI 2006)
  - CN2-SD (Lavrač et al., JMLR 2004): Adapting CN2 classification rule learner to Subgroup Discovery

# Part IV. Descriptive DM techniques

- Predictive vs. descriptive induction
- Subgroup discovery
→ • Association rule learning
- Hierarchical clustering

# Association Rule Learning

**Rules: X =>Y,  if X then Y**

X and Y are itemsets (records, conjunction of items), where items/features are binary-valued attributes)

**Given:** Transactions

itemsets (records)

| | i1 | i2 | .................... | i50 |
|---|---|---|---|---|
| t1 | 1 | 1 | | 0 |
| t2 | 0 | 1 | | 0 |
| ... | ... | ... | .................. | ... |

**Find:** A set of association rules in the form X =>Y

**Example:** Market basket analysis

beer & coke **=>** peanuts & chips (0.05, 0.65)

- Support:  $Sup(X,Y) = \#XY/\#D = p(XY)$

- Confidence: $Conf(X,Y) = \#XY/\#X = Sup(X,Y)/Sup(X) =$
$$= p(XY)/p(X) = p(Y|X)$$

# Association Rule Learning: Examples

- Market basket analysis
  - beer & coke $\Rightarrow$ peanuts & chips  (5%, 65%)

    (IF beer AND coke THEN peanuts AND chips)
  - Support 5%: 5% of all customers buy all four items
  - Confidence 65%: 65% of customers that buy beer and coke also buy peanuts and chips
- Insurance
  - mortgage & loans & savings $\Rightarrow$ insurance (2%, 62%)
  - Support 2%: 2% of all customers have all four
  - Confidence 62%: 62% of all customers that have mortgage, loan and savings also have insurance

# Association Rule Learning

**Given:** a set of transactions D

**Find:** all association rules that hold on the set of transactions that have

- user defined minimum support, i.e., support > MinSup, and
- user defined minimum confidence, i.e., confidence > MinConf

It is a form of exploratory data analysis, rather than hypothesis verification

# Searching for the associations

- Find all large itemsets

- Use the large itemsets to generate association rules

- If XY is a large itemset, compute

  r = support(XY) / support(X)

- If r > MinConf, then X $\Rightarrow$ Y holds

  (support > MinSup, as XY is large)

# **Large itemsets**

- Large itemsets are itemsets that appear in at least MinSup transaction

- All subsets of a large itemset are large itemsets (e.g., if A,B appears in at least MinSup transactions, so do A and B)

- This observation is the basis for very efficient algorithms for association rules discovery (linear in the number of transactions)

# Association vs. Classification rules

- Exploration of dependencies
- Different combinations of dependent and independent attributes
- Complete search (all rules found)

- Focused prediction
- Predict one attribute (class) from the others
- Heuristic search (subset of rules found)

# Part IV. Descriptive DM techniques

- Predictive vs. descriptive induction
- Subgroup discovery
- Association rule learning
- Hierarchical clustering

# Hierarchical clustering

- Algorithm (agglomerative hierarchical clustering):

| |
|---|
| **Each instance is a cluster;** |
| |
| **repeat** |
|     **find** *nearest* **pair** $C_i$ **in** $C_j$; |
|     *fuse* $C_i$ **in** $C_j$ **in a new cluster** |
|         $C_r = C_i \cup C_j$; |
|     **determine** *dissimilarities* **between** |
|         $C_r$ **and other clusters;** |
| |
| **until** **one cluster left;** |

- Dendogram:

# **Hierarchical clustering**

- Fusing the nearest pair of clusters



Where the diagram shows: $C_i$, $C_j$, $C_k$ with $d(C_i, C_k)$, $d(C_i, C_j)$, $d(C_j, C_k)$

- Minimizing intra-cluster similarity
- Maximizing inter-cluster similarity

- Computing the dissimilarities from the "new" cluster

# Hierarchical clustering: example



a) sample problem

b) dissimilarity matrix

c) dissimilarity matrix after 'fusing' elements **x** and **y**

d) dissimilarity matrix after 'fusing' elements **w** and **v**

e) dissimilarity matrix after 'fusing' cluster **(x,y)** and element **z**

f) dendrogram

# Results of clustering



A dendogram of resistance vectors

[Bohanec et al., "PTAH: A system for supporting nosocomial infection therapy", IDAMAP book, 1997]

# Course Outline

## I. Introduction

- Data Mining and KDD process
- Introduction to Data Mining
- Data Mining platforms

## II. Predictive DM Techniques

- Decision Tree learning
- Bayesian classifier
- Classification rule learning
- Classifier Evaluation

## III. Regression

## IV. Descriptive DM

- Predictive vs. descriptive induction
- Subgroup discovery
- Association rule learning Hierarchical clustering

## V. Relational Data Mining

- RDM and Inductive Logic Programming
- Propositionalization
- Semantic data mining

## VI. Advanced Topics

# Part V:
# Relational Data Mining

What is RDM

- Propositionalization techniques
- Semantic Data Mining

# Relational Data Mining
# (Inductive Logic Programming) task



knowledge discovery
from data

Relational Data Mining

model, patterns, …

Relational representation of customers, orders and stores.

**Given:** a relational database, a set of tables. sets of logical facts, a graph, …
**Find:** a classification model, a set of interesting patterns

# Relational data mining

- **ILP, relational learning, relational data mining**
  - Learning from complex multi-relational data



Relational representation of customers, orders and stores.

# Relational data mining

- **ILP, relational learning, relational data mining**
  - Learning from complex multi-relational data
  - Learning from complex structured data: e.g., molecules and their biochemical properties



Relational representation of customers, orders and stores.

# Sample problem: East-West trains

# RDM knowledge representation (database)

**LOAD_TABLE**

| LOAD | CAR | OBJECT | NUMBER |
|------|-----|--------|--------|
| l1 | c1 | circle | 1 |
| l2 | c2 | hexagon | 1 |
| l3 | c3 | triangle | 1 |
| l4 | c4 | rectangle | 3 |
| … | … | … | |

**TRAIN_TABLE**

| TRAIN | EASTBOUND |
|-------|-----------|
| t1 | **TRUE** |
| t2 | **TRUE** |
| … | … |
| t6 | **FALSE** |
| … | … |

**CAR_TABLE**

| CAR | TRAIN | SHAPE | LENGTH | ROOF | WHEELS |
|-----|-------|-------|--------|------|--------|
| c1 | t1 | rectangle | short | none | 2 |
| c2 | t1 | rectangle | long | none | 3 |
| c3 | t1 | rectangle | short | peaked | 2 |
| c4 | t1 | rectangle | long | none | 2 |
| … | … | … | | | … |

# ER diagram for East-West trains

# Relational data mining

- Relational data mining is characterized by using background knowledge (domain knowledge) in the data mining process

- Selected approaches:
  - Inductive logic programming - ILP (Muggleton, 1991; Lavrač & Džeroski 1994), …
  - Relational learning (Quinlan,1993)
  - Learning in DL (Lisi 2004), …
  - Relational Data Mining (Džeroski & Lavrač, 2001),
  - Statistical relational learning (Domingos, De Raedt…)
  - Propositionalization approach to RDM (Lavrač et al.)

# Our early work:
# Semantic subgroup discovery

- Propositionalization approach: Using relational subgroup discovery in the SDM context
  - General purpose system **RSD** for **Relational Subgroup Discovery**, using a propositionalization approach to relational data mining
  - Applied to semantic data mining in a biomedical application by using the Gene Ontology as background knowledge in analyzing microarray data

(Železny and Lavrač, MLJ 2006)

# Part V:
# Relational Data Mining

- What is RDM

  ➡ Propositionalization techniques

- Semantic Data Mining

# Relational Data Mining through Propositionalization

| customer | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| ID | Zip | Sex | SoSt | Income | Age | Club | Resp |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 3478 | 34677 | m | si | 60-70 | 32 | me | nr |
| 3479 | 43666 | f | ma | 80-90 | 45 | nm | re |
| ... | ... | ... | ... | ... | ... | ... | ... |

| order | | | | |
|---|---|---|---|---|
| Customer ID | Order ID | Store ID | Delivery Mode | Paymt Mode |
| ... | ... | ... | ... | ... |
| 3478 | 2140267 | 12 | regular | cash |
| 3478 | 3446778 | 12 | express | check |
| 3478 | 4728386 | 17 | regular | check |
| 3479 | 3233444 | 17 | express | credit |
| 3479 | 3475886 | 12 | regular | credit |
| ... | ... | ... | ... | ... |

| store | | | |
|---|---|---|---|
| Store ID | Size | Type | Location |
| ... | ... | ... | ... |
| 12 | small | franchise | city |
| 17 | large | indep | rural |
| ... | ... | ... | ... |

Relational representation of customers, orders and stores.

**Step 1**

Propositionalization

| | f1 | f2 | f3 | f4 | f5 | f6 | ... | | | ... | fn |
|---|---|---|---|---|---|---|---|---|---|---|---|
| g1 | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 0 | 1 | 1 |
| g2 | 0 | 1 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | 0 |
| g3 | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 1 |
| g4 | 1 | 1 | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 1 | 0 |
| g5 | 1 | 1 | 1 | 0 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 |
| g1 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 |
| g2 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 1 | 1 |
| g3 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 0 |
| g4 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 1 |

# Relational Data Mining through Propositionalization



customer

| ID | Zip | Sex | SoSt | Income | Age | Club | Resp |
|----|-----|-----|------|--------|-----|------|------|
| ... | ... | ... | ... | ... | ... | ... | ... |
| 3478 | 34677 | m | si | 60-70 | 32 | me | nr |
| 3479 | 43666 | f | ma | 80-90 | 45 | nm | re |
| ... | ... | ... | ... | ... | ... | ... | ... |

order

| Customer ID | Order ID | Store ID | Delivery Mode | Paymt Mode |
|-------------|----------|----------|---------------|------------|
| ... | ... | ... | ... | ... |
| 3478 | 2140267 | 12 | regular | cash |
| 3478 | 3446778 | 12 | express | check |
| 3478 | 4728386 | 17 | regular | check |
| 3479 | 3233444 | 17 | express | credit |
| 3479 | 3475886 | 12 | regular | credit |
| ... | ... | ... | ... | ... |

store

| Store ID | Size | Type | Location |
|----------|------|------|----------|
| ... | ... | ... | ... |
| 12 | small | franchise | city |
| 17 | large | indep | rural |
| ... | ... | ... | ... |

Relational representation of customers, orders and stores.

**Step 1**

Propositionalization

1. constructing relational features
2. constructing a propositional table

| | f1 | f2 | f3 | f4 | f5 | f6 | ... | | | ... | fn |
|-----|----|----|----|----|----|----|-----|---|---|-----|----|
| g1 | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 0 | 1 | 1 |
| g2 | 0 | 1 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | 0 |
| g3 | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 1 |
| g4 | 1 | 1 | 1 | 0 | 1 | 1 | 0 | 1 | 1 | 1 | 0 | 0 |
| g5 | 1 | 1 | 1 | 0 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 |
| g1 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 |
| g2 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 1 | 1 |
| g3 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 0 |
| g4 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 1 |

# Relational Data Mining through Propositionalization



Relational representation of customers, orders and stores.

**Step 1**

Propositionalization

**Step 2**

Data Mining

model, patterns, …

# Relational Data Mining through Propositionalization

**Step 1**

Propositionalization

**Step 2**

Data Mining

Relational representation of customers, orders and stores.

```
target(A) :-
    'Doctor'(A), 'Italy'(A).

target(A) :-
    'Public'(A), 'Gold'(A).

target(A) :-
    'Poland'(A), 'Deposit'(A), 'Gold'(A).

target(A) :-
    'Germany'(A), 'Insurance'(A).

target(A) :-
    'Service'(A), 'Germany'(A).
```

patterns (set of rules)

# Sample ILP problem: East-West trains



1. TRAINS GOING EAST

2. TRAINS GOING WEST

# Relational data representation



| LOAD | CAR | OBJECT | NUMBER |
|---|---|---|---|
| l1 | c1 | circle | 1 |
| l2 | c2 | hexagon | 1 |
| l3 | c3 | triangle | 1 |
| l4 | c4 | rectangle | 3 |
| ... | ... | ... | |

**TRAIN_TABLE**

| TRAIN | EASTBOUND |
|---|---|
| t 1 | TRUE |
| t 2 | TRUE |
| … | … |
| t 6 | FALSE |
| … | … |

| CAR | TRAIN | SHAPE | LENGTH | ROOF | WHEELS |
|---|---|---|---|---|---|
| c1 | t 1 | rectangle | short | none | 2 |
| c2 | t 1 | rectangle | long | none | 3 |
| c3 | t 1 | rectangle | short | peaked | 2 |
| c4 | t 1 | rectangle | long | none | 2 |
| … | … | … | | | … |

# Propositionalization in a nutshell

**Propositionalization task**

**Transform** a multi-relational
(**multiple-table**)
representation to a
propositional representation
(**single table**)

Proposed in ILP systems
LINUS (Lavrac et al. 1991, 1994),
1BC (Flach and Lachiche 1999), ...

| LOAD | CAR | OBJECT | NUMBER |
|------|-----|--------|--------|
| l1 | c1 | circle | 1 |
| l2 | c2 | hexagon | 1 |
| l3 | c3 | triangle | 1 |
| l4 | c4 | rectangle | 3 |
| ... | ... | ... | |

**TRAIN_TABLE**

| TRAIN | EASTBOUND |
|-------|-----------|
| t1 | TRUE |
| t2 | TRUE |
| ... | ... |
| t6 | FALSE |
| ... | ... |

| CAR | TRAIN | SHAPE | LENGTH | ROOF | WHEELS |
|-----|-------|-------|--------|------|--------|
| c1 | t1 | rectangle | short | none | 2 |
| c2 | t1 | rectangle | long | none | 3 |
| c3 | t1 | rectangle | short | peaked | 2 |
| c4 | t1 | rectangle | long | none | 2 |
| ... | ... | ... | | | ... |

# Propositionalization in a nutshell

**Main propositionalization step: first-order feature construction**

f1(T):-hasCar(T,C),clength(C,short).

f2(T):-hasCar(T,C), hasLoad(C,L),
     loadShape(L,circle)

f3(T) :- ....

**Propositional learning:**

t(T) ← f1(T), f4(T)

**Relational interpretation:**

eastbound(T) ←
hasShortCar(T),hasClosedCar(T).

| LOAD | CAR | OBJECT | NUMBER |
|---|---|---|---|
| l1 | c1 | circle | 1 |
| l2 | c2 | hexagon | 1 |
| l3 | c3 | triangle | 1 |
| l4 | c4 | rectangle | 3 |
| ... | ... | ... | |

**TRAIN_TABLE**

| TRAIN | EASTBOUND |
|---|---|
| t1 | TRUE |
| t2 | TRUE |
| ... | ... |
| t6 | FALSE |
| ... | ... |

| CAR | TRAIN | SHAPE | LENGTH | ROOF | WHEELS |
|---|---|---|---|---|---|
| c1 | t1 | rectangle | short | none | 2 |
| c2 | t1 | rectangle | long | none | 3 |
| c3 | t1 | rectangle | short | peaked | 2 |
| c4 | t1 | rectangle | long | none | 2 |
| ... | ... | ... | | | ... |

**PROPOSITIONAL TRAIN_TABLE**

| train(T) | f1(T) | f2(T) | f3(T) | f4(T) | f5(T) |
|---|---|---|---|---|---|
| t1 | t | t | f | t | t |
| t2 | t | t | t | t | t |
| t3 | f | f | t | f | f |
| t4 | t | f | t | f | f |
| ... | ... | ... | | | ... |

# Part V:
# Relational Data Mining

- What is RDM
- Propositionalization techniques

Semantic Data Mining

# Semantic data mining

- **ILP, relational learning, relational data mining**
  - Learning from complex multi-relational data
  - Learning from complex structured data: e.g., molecules and their biochemical properties
  - Learning by using domain knowledge in the form of ontologies = **semantic data mining**



Relational representation of customers, orders and stores.

# Using domain ontologies in Semantic Data Mining

Using domain ontologies as background knowledge, e.g., using the Gene Ontology (GO)

- GO is a database of terms, describing gene sets in terms of their
  - functions (12,093)
  - processes (1,812)
  - components (7,459)
- Genes are annotated to GO terms
- Terms are connected (is_a, part_of)
- Levels represent terms generality

# What is Semantic Data Mining

- Ontology-driven (semantic) data mining is an emerging research topic

- Semantic Data Mining (SDM) - a new term denoting:

  - the new challenge of mining semantically annotated resources, with ontologies used as background knowledge to data mining

  - approaches with which semantic data are mined

# What is Semantic Data Mining

SDM task definition



ontologies

annotations, mappings

Semantic data mining

data

```
target(A) :-
    'Doctor'(A), 'Italy'(A).

target(A) :-
    'P

target

target        'Gold'(A).

target
    'Germany'(A), 'Insurance'(A).

target(A) :-
    'Service'(A), 'Germany'(A).
```

model, patterns

**Given:**

- transaction data table, relational database, text documents, Web pages, …

- one or more domain ontologies

**Find:** a classification model, a set of patterns

# Using domain ontologies (e.g. Gene Ontology) as background knowledge for Data Mining

## Gene Ontology

**12093 biological process**
**1812 cellular components**
**7459 molecular functions**

**Joint work with**
**Igor Trajkovski**
**and Filip Zelezny**

# Using domain ontologies (e.g. Gene Ontology) as background knowledge for Data Mining

**First-order features, describing**

**gene properties and relations between genes, can be viewed as generalisations of individual genes**

# Semantic subgroup discovery with RSD

1. Take ontology terms represented as logical facts in Prolog, e.g.

```
component(gene2532,'GO:0016020').
function(gene2534,'GO:0030554').
process(gene2534,'GO:0007243').
interaction(gene2534,gene4803).
```

2. Automatically generate generalized relational features:

```
f(2,A):-component(A,'GO:0016020').
f(7,A):-function(A,'GO:0030554').
f(11,A):-process(A,'GO:0007243').
f(224,A):- interaction(A,B), function(B,'GO:0016787'),
           component(B,'GO:0043231').
```

3. Propositionalization: Determine truth values of features

4. Learn rules by a subgroup discovery algorithm CN2-SD

# Step 2: RSD feature construction

Construction of first order features, with support > *min_support*

f(7,A):-function(A,'GO:0046872').
f(8,A):-function(A,'GO:0004871').
f(11,A):-process(A,'GO:0007165').
f(14,A):-process(A,'GO:0044267').
f(15,A):-process(A,'GO:0050874').
f(20,A):-function(A,'GO:0004871'), process(A,'GO:0050874').
f(26,A):-component(A,'GO:0016021').
f(29,A):- function(A,'GO:0046872'), component(A,'GO:0016020').
f(122,A):-interaction(A,B),function(B,'GO:0004872').
f(223,A):-interaction(A,B),function(B,'GO:0004871'),
    process(B,'GO:0009613').
f(224,A):-interaction(A,B),function(B,'GO:0016787'),
    component(B,'GO:0043231').

existential

# Step 3: RSD Propositionalization

diffexp g1 (gene64499)          random g1 (gene7443)
diffexp g2 (gene2534)           random g2 (gene9221)
diffexp g3 (gene5199)           random g3 (gene2339)
diffexp g4 (gene1052)           random g4 (gene9657)
diffexp g5 (gene6036)           random g5 (gene19679)

….                              ….

|    | f1 | f2 | f3 | f4 | f5 | f6 | … |   |   |   | … | fn |
|----|----|----|----|----|----|----|---|---|---|---|---|----|
| g1 | 1  | 0  | 0  | 1  | 1  | 1  | 0 | 0 | 1 | 0 | 1 | 1  |
| g2 | 0  | 1  | 1  | 0  | 1  | 1  | 0 | 0 | 0 | 1 | 1 | 0  |
| g3 | 0  | 1  | 1  | 1  | 0  | 0  | 1 | 1 | 0 | 0 | 0 | 1  |
| g4 | 1  | 1  | 1  | 0  | 1  | 1  | 0 | 0 | 1 | 1 | 1 | 0  |
| g5 | 1  | 1  | 1  | 0  | 0  | 1  | 0 | 1 | 1 | 0 | 1 | 0  |
| g1 | 0  | 0  | 1  | 1  | 0  | 0  | 0 | 1 | 0 | 0 | 0 | 1  |
| g2 | 1  | 1  | 0  | 0  | 1  | 1  | 0 | 1 | 0 | 1 | 1 | 1  |
| g3 | 0  | 0  | 0  | 0  | 1  | 0  | 0 | 1 | 1 | 1 | 0 | 0  |
| g4 | 1  | 0  | 1  | 1  | 1  | 0  | 1 | 0 | 0 | 1 | 0 | 1  |

# Step 4: RSD rule construction with CN2-SD

|     | f1 | f2 | f3 | f4 | f5 | f6 | ... |   |   |   | ... | fn |
|-----|----|----|----|----|----|----|-----|---|---|---|-----|----|
| g1  | 1  | 0  | 0  | 1  | 1  | 1  | 0   | 0 | 1 | 0 | 1   | 1  |
| g2  | 0  | 1  | 1  | 0  | 1  | 1  | 0   | 0 | 0 | 1 | 1   | 0  |
| g3  | 0  | 1  | 1  | 1  | 0  | 0  | 1   | 1 | 0 | 0 | 0   | 1  |
| g4  | 1  | 1  | 1  | 0  | 1  | 1  | 0   | 0 | 1 | 1 | 1   | 0  |
| g5  | 1  | 1  | 1  | 0  | 0  | 1  | 0   | 1 | 1 | 0 | 1   | 0  |
| g1  | 0  | 0  | 1  | 1  | 0  | 0  | 0   | 1 | 0 | 0 | 0   | 1  |
| g2  | 1  | 1  | 0  | 0  | 1  | 1  | 0   | 1 | 0 | 1 | 1   | 1  |
| g3  | 0  | 0  | 0  | 0  | 1  | 0  | 0   | 1 | 1 | 1 | 0   | 0  |
| g4  | 1  | 0  | 1  | 1  | 1  | 0  | 1   | 0 | 0 | 1 | 0   | 1  |

Over-
expressed

IF

f2 and f3

[4,0]

diffexp(A) :- interaction(A,B) & function(B,'GO:0004871')

# Subgroup Discovery

diff. exp. genes

Not diff. exp. genes

# Subgroup Discovery

Cl=YES ← f2 and f3

diff. exp. genes

Not diff. exp. genes



In RSD (using propositional learner CN2-SD):

Quality of the rules = Coverage  x  Precision

*Coverage = sum of the covered weights

*Precision = purity of the covered genes

# Subgroup Discovery



diff. exp. genes

Not diff. exp. genes

RSD naturally uses gene weights in its procedure for repetitive subgroup generation, via its heuristic rule evaluation: weighted relative accuracy

# RSD Lessons learned

Efficient propositionalization can be applied to individual-centered, multi-instance learning problems:

– one free global variable (denoting an individual, e.g. molecule M)

– one or more structural predicates: (e.g. has_atom(M,A)), each introducing a new existential local variable (e.g. atom A), using either the global variable (M) or a local variable introduced by other structural predicates (A)

– one or more utility predicates defining properties of individuals or their parts, assigning values to variables

feature121(M):- hasAtom(M,A), atomType(A,21)

feature235(M):- lumo(M,Lu), lessThr(Lu,-1.21)

mutagenic(M):- feature121(M), feature235(M)

# SEGS: using RSD approach

- The SEGS approach enables to discover new medical knowledge from the combination of gene expression data with public gene annotation databases

- The SEGS approach proved effective in several biomedical applications (JBI 2008, …)

  - The work on semantic data mining - using ontologies as background knowledge for subgroup discovery with SEGS - was done in collaboration with I.Trajkovski, F. Železny and J. Tolar

- Recent work: Semantic subgroup discovery implemented in Orange4WS

# Semantic subgroup discovery with SEGS

- SEGS workflow is implemented in the Orange4WS data mining environment



- SEGS is also implemented also as a Web applications

  (Trajkovski et al., IEEE TSMC 2008, Trajkovski et al., JBI 2008)

# From SEGS to SDM-SEGS: Generalizing SEGS

- SDM-SEGS: a general semantic data mining



- Discovers subgroups both for ranked and labeled data

- Exploits input ontologies in OWL format

- Is also implemented in Orange4WS

# Relational Data Mining in Orange4WS

- service for propositionalization through efficient first-order feature construction (Železny and Lavrač, MLJ 2006)

    f121(M):- hasAtom(M,A), atomType(A,21)

    f235(M):- lumo(M,Lu), lessThr(Lu,1.21)

- subgroup discovery using CN2-SD

    mutagenic(M) ← feature121(M), feature235(M)

# Semantic Data Mining in Orange4WS

- A special purpose Semantic Data Mining algorithm SEGS
  - discovers interesting gene group descriptions as conjunctions of ontology concepts from GO, KEGG and Entrez
  - integrates public gene annotation data through relational features
  - SEGS algorithm (Trajkovski, Železny, Lavrač and Tolar, JBI 2008) is available in Orange4WS
- Recent developments:
  - Special purpose SDM algorithms: RSD, SDM-SEGS, SDM-Aleph, Hedwig
  - Implemented in web based DM platform ClowdFlows

# Third Generation Data Mining Platform: ClowdFlows

- **ClowdFlows** - browsed-based DM platform for data mining in the cloud and workflow sharing on the web (Kranjc et al. 2012)

- RSD, SDM-SEGS, SDM-Aleph, Hedwig are available as ingredients of elaborate data mining workflows in ClowdFlows

- **Example workflow**: Propositionalization with RSD available in ClowdFlows at http://clowdflows.org/workflow/611/

# Sample biomedical application of Hedwig

- Semantic subgroup discovery and semantic explanation of subgroups on breast cancer data (Vavpetič et al., JIIS 2014)



- The workflow, implemented in ClowdFlows, is available at http://clowdflows.org/workflow/1283/

# Semantic Data Mining

- Semantic subgroup discovery (Vavpetič et al., 2012)

# Course Outline

**I. Introduction**
- Data Mining and KDD process
- Introduction to Data Mining
- Data Mining platforms

**II. Predictive DM Techniques**
- Decision Tree learning
- Bayesian classifier
- Classification rule learning
- Classifier Evaluation

**III. Regression**

**IV. Descriptive DM**
- Predictive vs. descriptive induction
- Subgroup discovery
- Association rule learning Hierarchical clustering

**V. Relational Data Mining**
- RDM and Inductive Logic Programming
- Propositionalization
- Semantic data mining

**VI. Advanced Topics**

# Advanced Topics I.

ClowdFlows Data Mining Platform
(PhD of Janez Kranjc, demo Martin Žnidaršič)

- Outlier detection with NoiseRank
(PhD of Borut Sluban)

# Open data science platform ClowdFlows

- Third generation platform for the creation and execution of complex data mining workflows
  - Algorithms as web services (in the cloud)
  - No need for platform installation
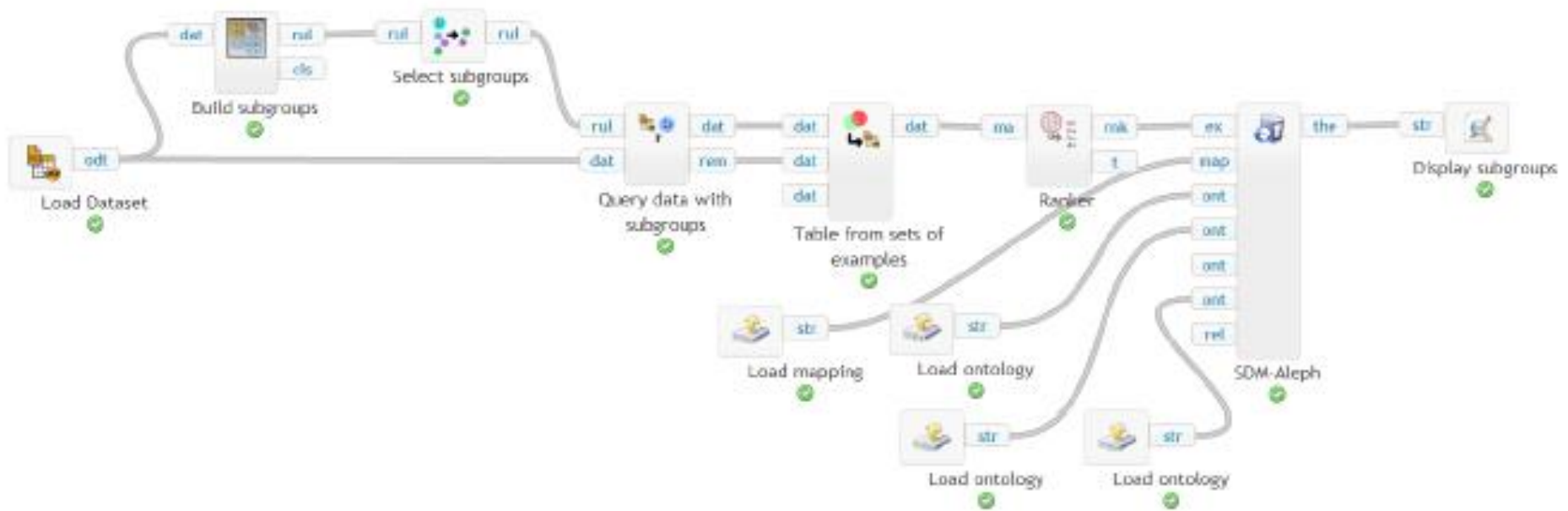  - Workflows are openly accessible and executable from any modern web browser by a web site klick

http://clowdflows.org/workflow/1283/

# ClowdFlows platform

- is service oriented (DM algorithms as web services)

- includes functionality of other DM platforms, e.g. WEKA algorithms, implemented as Web services

- includes new functionality, e.g. relational data mining, semantic data mining, big data analytics, text mining, …

- enables simplified construction of Web services from available algorithms

- runs in any browser, enabling workflow construction and sharing on the web
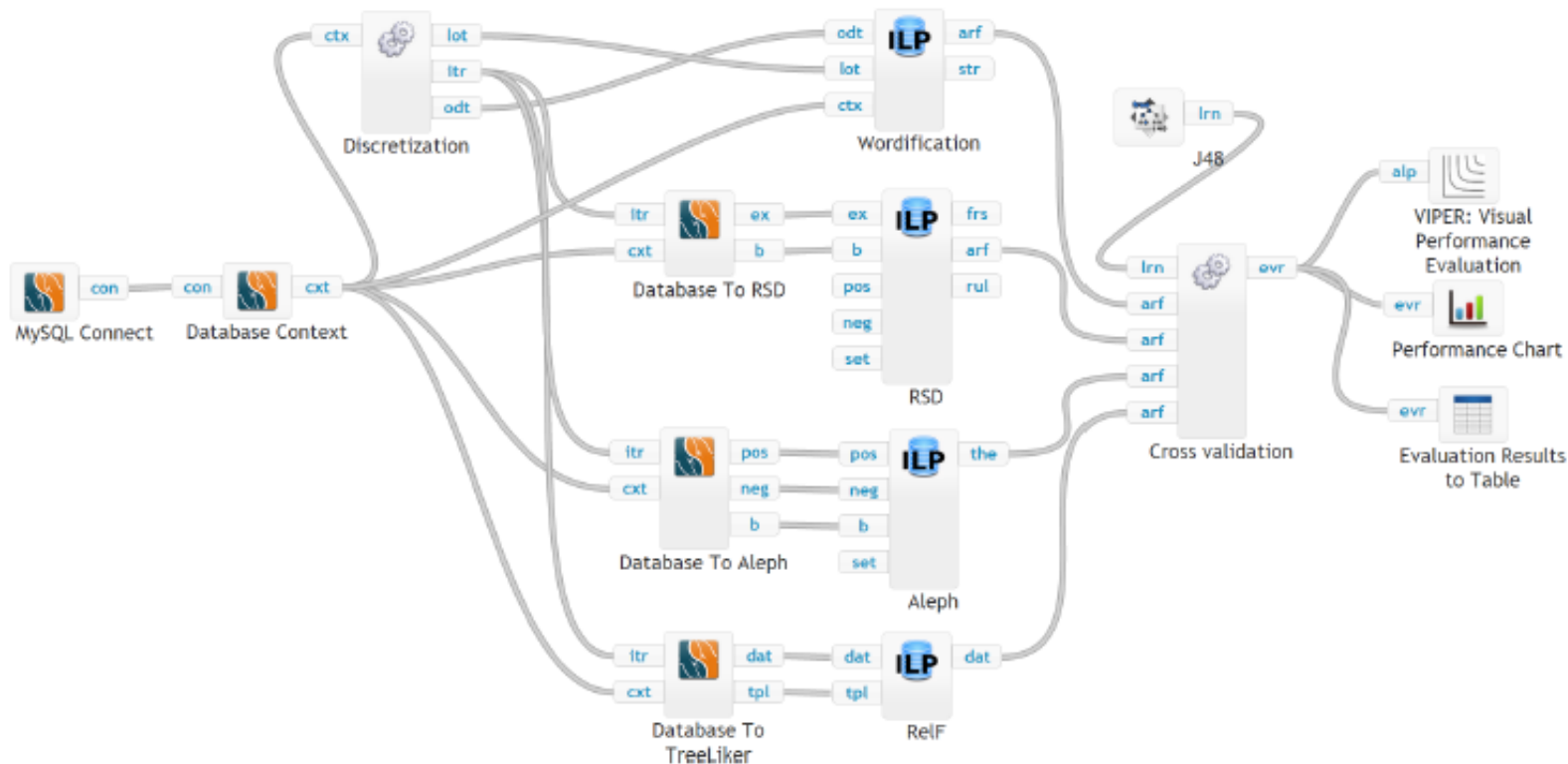
- user-friendly HCI: canvas for workflow construction

# SDM in ClowdFlows

- Semantic subgroup discovery and semantic explanation of subgroups on breast cancer data (Vavpetič et al., JIIS 2014)



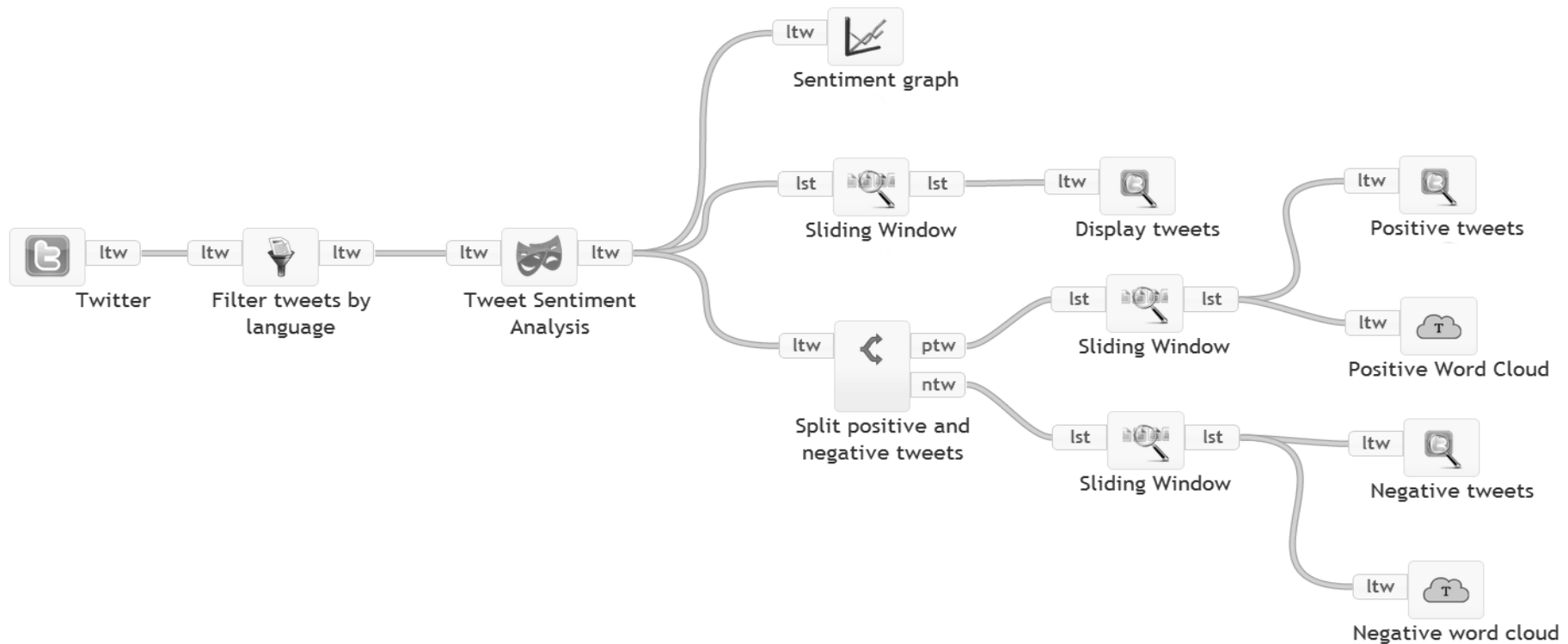- The workflow, implemented in ClowdFlows, is available for sharing at http://clowdflows.org/workflow/1283/

# Propositionalization and Wordification in ClowdFlows



Wordification and propositionalization algorithms comparison, available at http://clowdflows.org/workflow/1456/

# Analysis of Big data in ClowdFlows

- Big data analysis in real time

- Example: Semantic graph construction from a stream of web news http://clowdflows.org/workflow/1729/.



RSS Reader · Summarize news article · Triplet Extraction · WordNet lemmatizer on triplets · Sliding Window · Streaming triplet graph

- Example: news monitoring by visualization of graph constructed from sports news (CNN RSS feeds)

http://clowdflows.org/streams/data/31/15524/

# Analysis of Big data in ClowdFlows

- Analysis of positive/negative sentiment in tweets in real time  http://clowdflows.org/workflow/1041/.

# Advanced Topics I.

- ClowdFlows Data Mining Platform
  (PhD of Janez Kranjc, demo Martin Žnidaršič)
- Outlier detection with NoiseRank
  (PhD of Borut Sluban)

# Noise and outliers

- Errors in the data – noise

  – Animals of white color



- Exceptions or Outliers

  – Herd of sheep

# Noise and outliers



- ## Data in nature
  - follows certain patters
  - adheres to the laws of physics
  - is not random

- ## Build models to Identify the "laws" of the data

  Patterns and rules =

  = "laws" of the data



- ## Errors and outliers
  - Do NOT obey the laws (models)

# Noise and outlier detection

- **Noise** in data negatively affect data mining results. (Zhu et al., 2004)

- False medical diagnosis (**classification noise**) can have serious consequences (Gamberger et al. 2003)

- **Outlier** detection proved to be effective in detection of network intrusion and bank fraud. (Aggarwal and Yu, 2001)

# Detecting noise and outliers

- Errors and exceptions are:

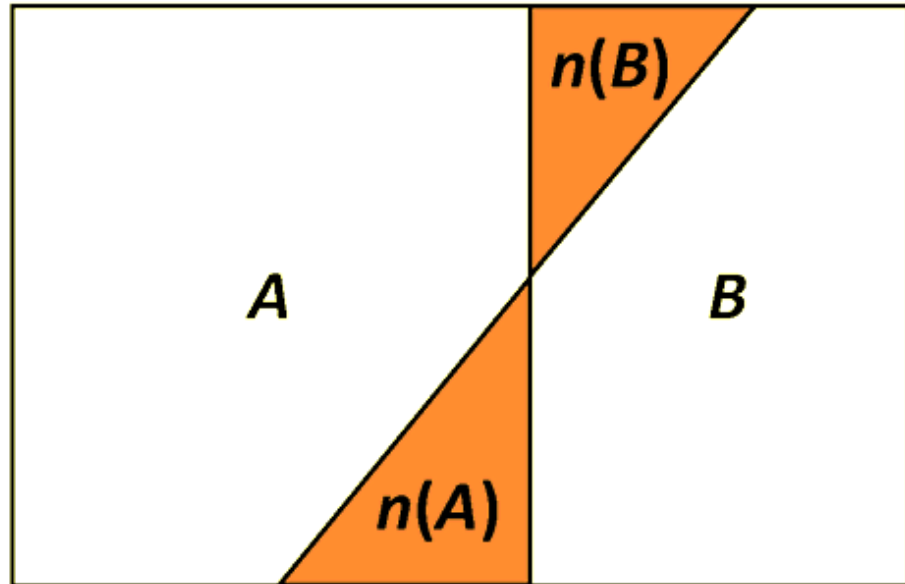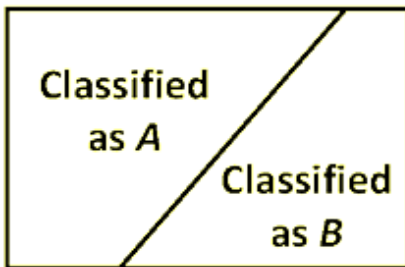  - Inconsistencies with common patterns

  - Great deviations from expected values
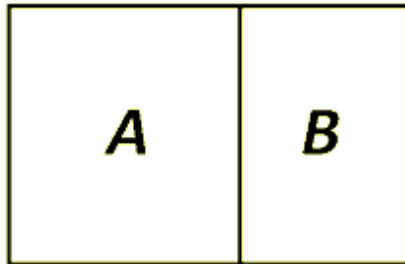
  - Hard to describe
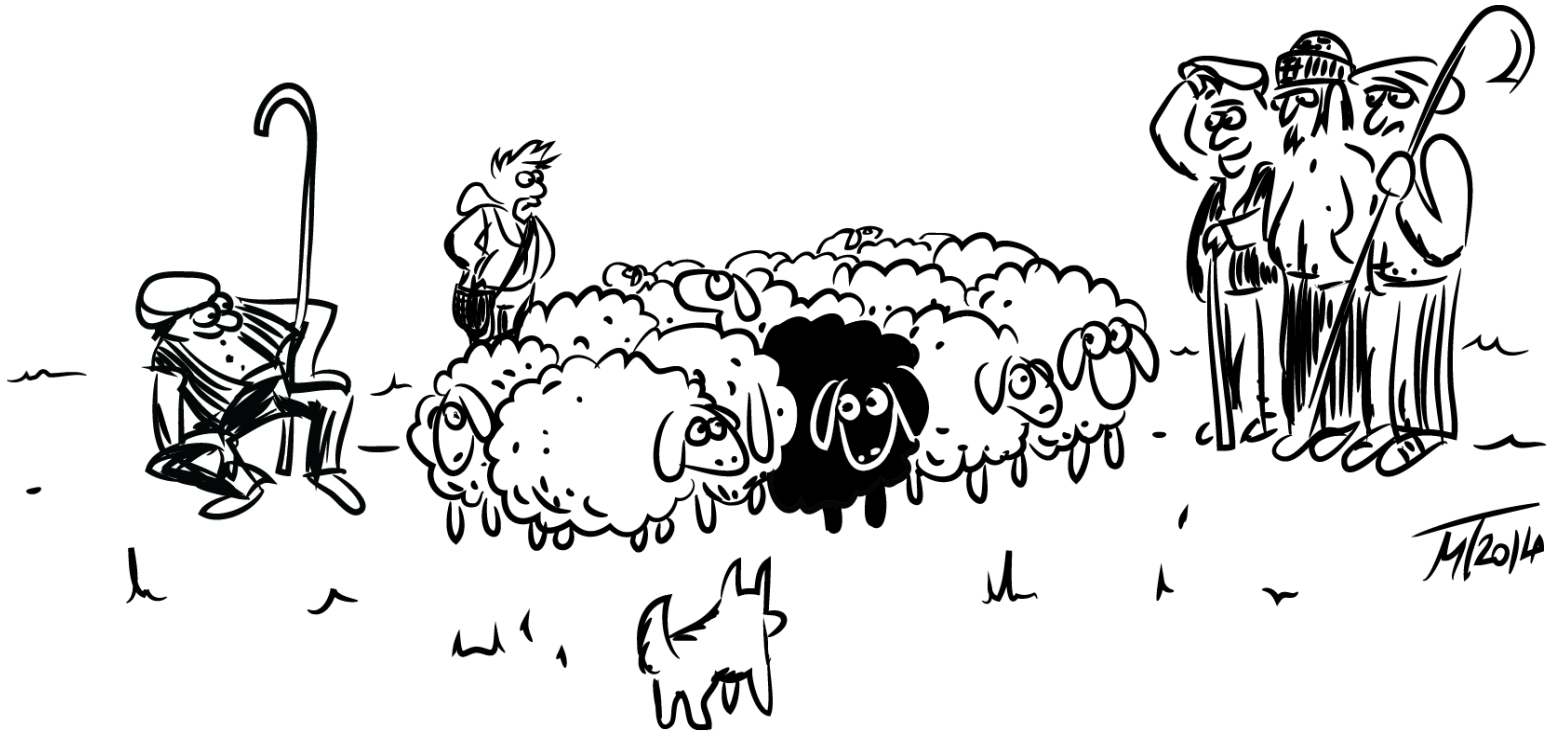
# Classification noise filtering

- Model the data
- What can't be modeled is considered noise
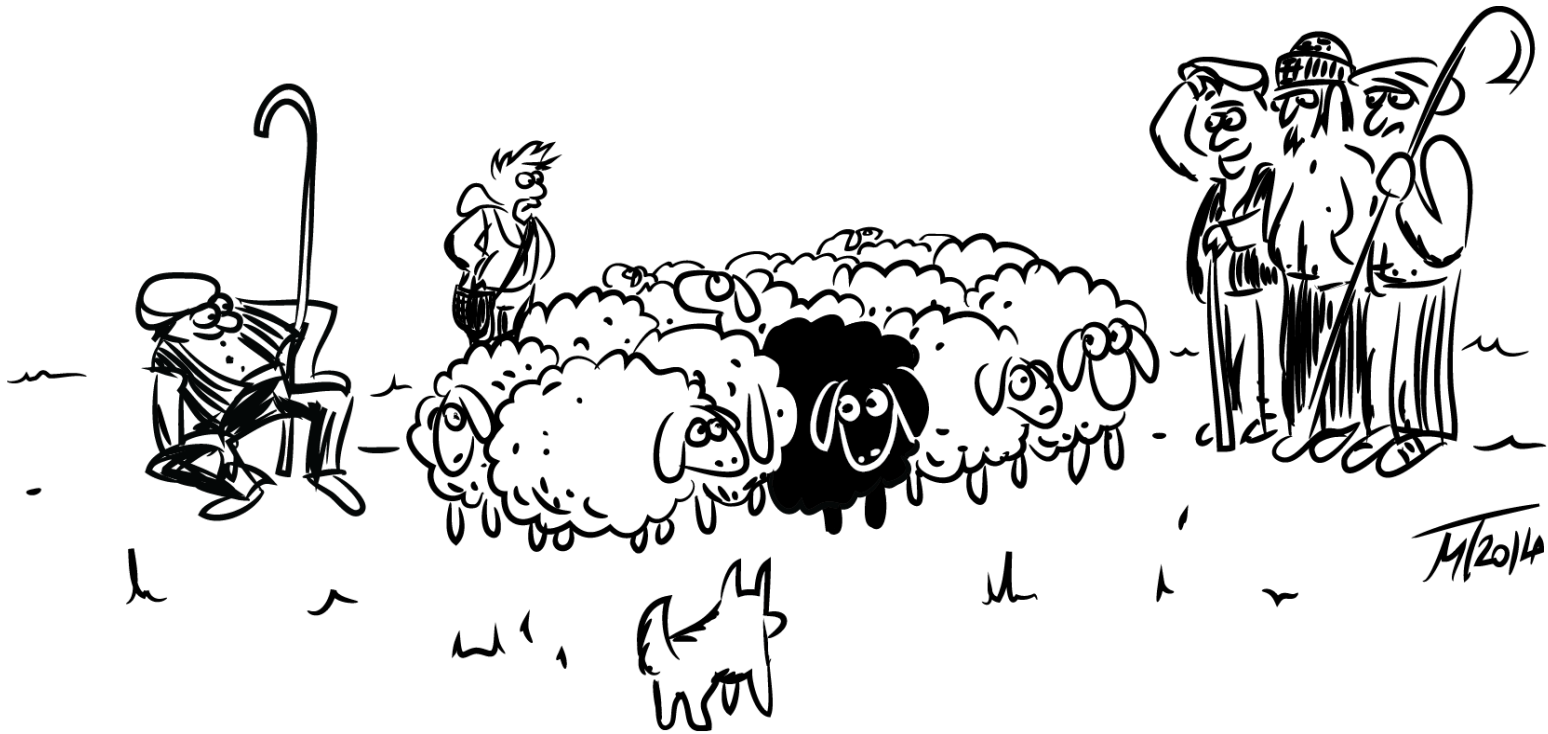
# Classification noise filtering

- Model the data, using any learning algorithm
- What can't be modeled is considered noise
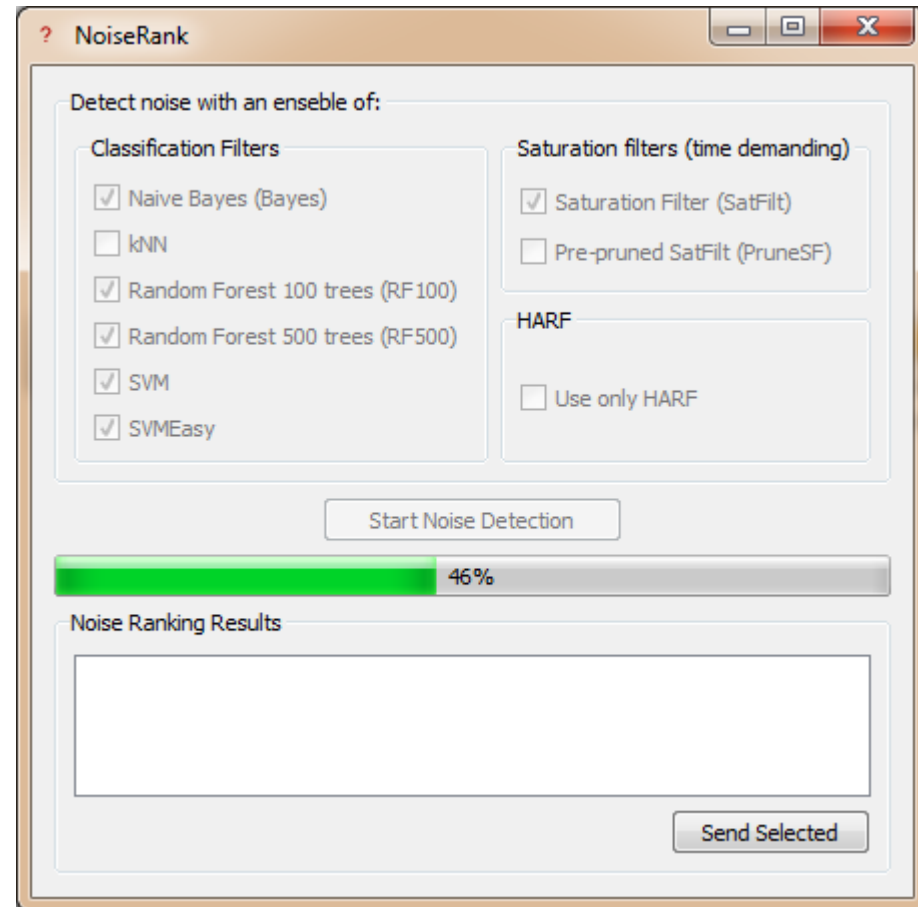
# Ensembles of classifiers

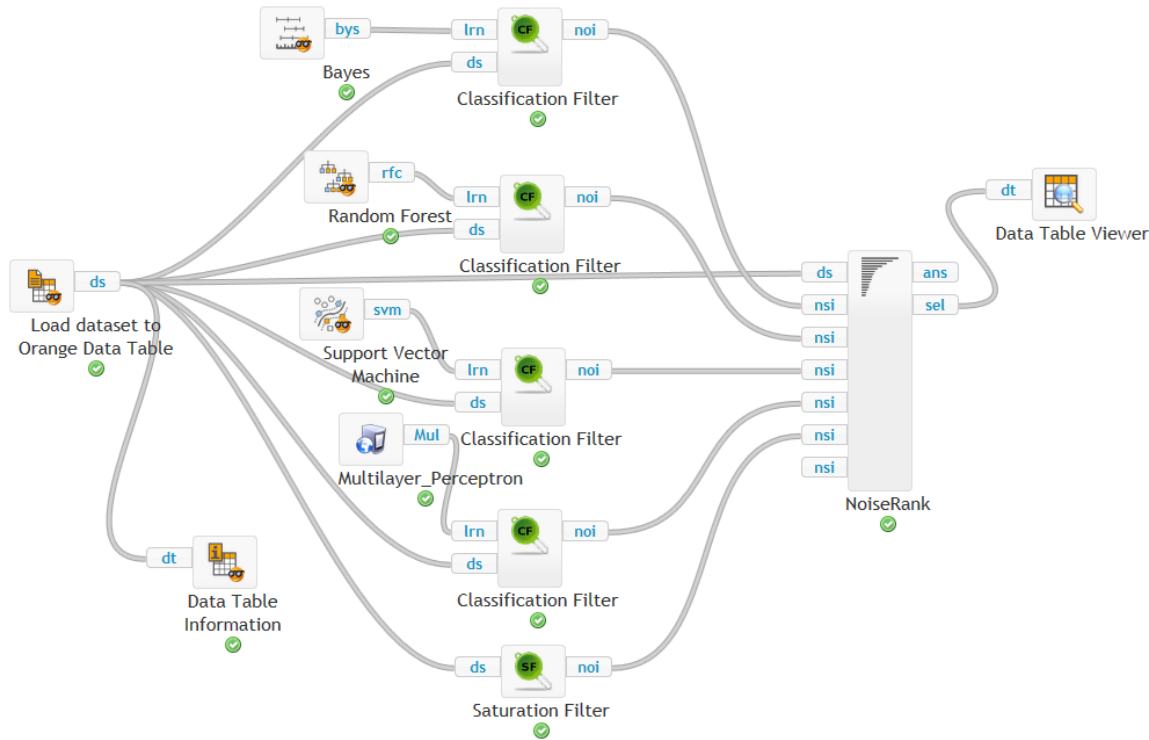# Ensembles of classifiers



- Combine predictions of various models

- To overcome weaknesses or bias of individual models

- Averaging, Majority voting, Consensus voting, Ranking, etc.
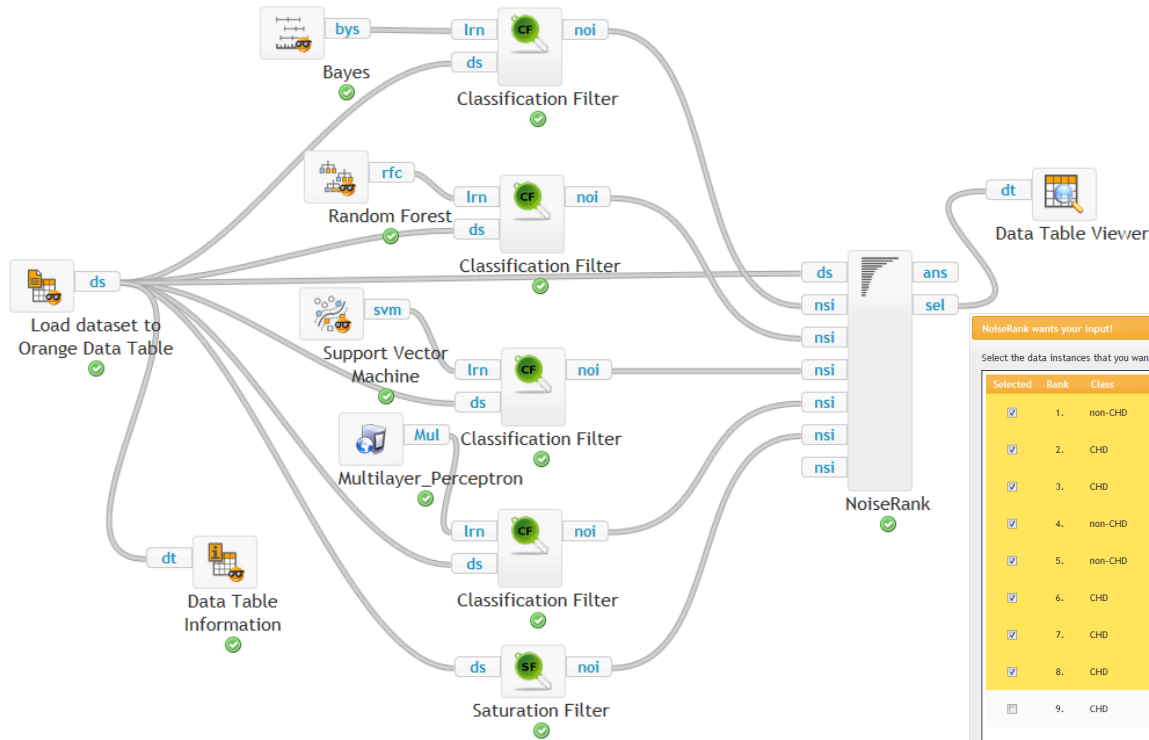
# NoiseRank: Ensemble-based noise and outlier detection

- Misclassified document detection by an ensemble of diverse classifiers (e.g., Naive Bayes, Random Forest, SVM, … classifiers)

- Ranking of misclassified documents by "voting" of classifiers

# NoiseRank Workflows

# NoiseRank Workflows

# NoiseRank: Ranked List of Noisy instances/Outliers

# Try it out

- NoiseRank
  - http://clowdflows.org/workflow/115/

- Clowdflows:
  - Noise Handling
  - Orange, Weka classification
  - Performance evaluation

- Noise filtering using ensembles (with performance evaluation)
  - http://clowdflows.org/workflow/245/

# Noise filtering using ensembles (with performance evaluation)

# **Advanced Topics II.**

→ Text mining: An introduction

- Document clustering and outlier detection

- Wordification approach to relational data mining

# Background: Data mining

| Person | Age | Spect. presc. | Astigm. | Tear prod. | Lenses |
|--------|-----|---------------|---------|------------|--------|
| O1 | 17 | myope | no | reduced | NONE |
| O2 | 23 | myope | no | normal | SOFT |
| O3 | 22 | myope | yes | reduced | NONE |
| O4 | 27 | myope | yes | normal | HARD |
| O5 | 19 | hypermetrope | no | reduced | NONE |
| O6-O13 | ... | ... | ... | ... | ... |
| O14 | 35 | hypermetrope | no | normal | SOFT |
| O15 | 43 | hypermetrope | yes | reduced | NONE |
| O16 | 39 | hypermetrope | yes | normal | NONE |
| O17 | 54 | myope | no | reduced | NONE |
| O18 | 62 | myope | no | normal | NONE |
| O19-O23 | ... | ... | ... | ... | ... |
| O24 | 56 | hypermetrope | yes | normal | NONE |

data

knowledge discovery from data

Data Mining

model, patterns, clusters, …

**Given:** transaction data table, a set of text documents, …

**Find:** a classification model, a set of interesting patterns

# Data mining: Task reformulation

| Person | Young | Myope | Astigm. | Reuced tea | Lenses |
|--------|-------|-------|---------|------------|--------|
| O1 | 1 | 1 | 0 | 1 | NO |
| O2 | 1 | 1 | 0 | 0 | YES |
| O3 | 1 | 1 | 1 | 1 | NO |
| O4 | 1 | 1 | 1 | 0 | YES |
| O5 | 1 | 0 | 0 | 1 | NO |
| O6-O13 | ... | ... | ... | ... | ... |
| O14 | 0 | 0 | 0 | 0 | YES |
| O15 | 0 | 0 | 1 | 1 | NO |
| O16 | 0 | 0 | 1 | 0 | NO |
| O17 | 0 | 1 | 0 | 1 | NO |
| O18 | 0 | 1 | 0 | 0 | NO |
| O19-O23 | ... | ... | ... | ... | ... |
| O24 | 0 | 0 | 1 | 0 | NO |

Binary features and class values

# Text mining:
# Words/terms as binary features

| Document | Word1 | Word2 | … | WordN | Class |
|---|---|---|---|---|---|
| d1 | 1 | 1 | 0 | 1 | NO |
| d2 | 1 | 1 | 0 | 0 | YES |
| d3 | 1 | 1 | 1 | 1 | NO |
| d4 | 1 | 1 | 1 | 0 | YES |
| d5 | 1 | 0 | 0 | 1 | NO |
| d6-d13 | … | … | … | … | … |
| d14 | 0 | 0 | 0 | 0 | YES |
| d15 | 0 | 0 | 1 | 1 | NO |
| d16 | 0 | 0 | 1 | 0 | NO |
| d17 | 0 | 1 | 0 | 1 | NO |
| d18 | 0 | 1 | 0 | 0 | NO |
| d19-d23 | … | … | … | … | … |
| d24 | 0 | 0 | 1 | 0 | NO |

Instances = documents
Words and terms = Binary features

# Text Mining from unlabeled data

| Document | Word1 | Word2 | … | WordN | Class |
|----------|-------|-------|-----|-------|-------|
| d1 | 1 | 1 | 0 | 1 | NO |
| d2 | 1 | 1 | 0 | 0 | YES |
| d3 | 1 | 1 | 1 | 1 | NO |
| d4 | 1 | 1 | 1 | 0 | YES |
| d5 | 1 | 0 | 0 | 1 | NO |
| d6-d13 | … | … | … | … | … |
| d14 | 0 | 0 | 0 | 0 | YES |
| d15 | 0 | 0 | 1 | 1 | NO |
| d16 | 0 | 0 | 1 | 0 | NO |
| d17 | 0 | 1 | 0 | 1 | NO |
| d18 | 0 | 1 | 0 | 0 | NO |
| d19-d23 | … | … | … | … | … |
| d24 | 0 | 0 | 1 | 0 | NO |

Unlabeled data - clustering: grouping of similar instances
- association rule learning

# Text mining

**Step 1**

BoW vector construction

1. BoW features construction
2. Table of BoW vectors construction

| Document | Word1 | Word2 | … | WordN | Class |
|----------|-------|-------|-----|-------|-------|
| d1 | 1 | 1 | 0 | 1 | NO |
| d2 | 1 | 1 | 0 | 0 | YES |
| d3 | 1 | 1 | 1 | 1 | NO |
| d4 | 1 | 1 | 1 | 0 | YES |
| d5 | 1 | 0 | 0 | 1 | NO |
| d6-d13 | … | … | … | … | … |
| d14 | 0 | 0 | 0 | 0 | YES |
| d15 | 0 | 0 | 1 | 1 | NO |
| d16 | 0 | 0 | 1 | 0 | NO |
| d17 | 0 | 1 | 0 | 1 | NO |
| d18 | 0 | 1 | 0 | 0 | NO |
| d19-d23 | … | … | … | … | … |
| d24 | 0 | 0 | 1 | 0 | NO |

| Document | Word1 | Word2 | … | WordN | Class |
|----------|-------|-------|-----|-------|-------|
| d1 | 1 | 1 | 0 | 1 | NO |
| d2 | 1 | 1 | 0 | 0 | YES |
| d3 | 1 | 1 | 1 | 1 | NO |
| d4 | 1 | 1 | 1 | 0 | YES |
| d5 | 1 | 0 | 0 | 1 | NO |
| d6-d13 | … | … | … | … | … |
| d14 | 0 | 0 | 0 | 0 | YES |
| d15 | 0 | 0 | 1 | 1 | NO |
| d16 | 0 | 0 | 1 | 0 | NO |
| d17 | 0 | 1 | 0 | 1 | NO |
| d18 | 0 | 1 | 0 | 0 | NO |
| d19-d23 | … | … | … | … | … |
| d24 | 0 | 0 | 1 | 0 | NO |

**Step 2**

Data Mining

model, patterns, clusters,

…

# Text Mining

- Feature construction
  - StopWords elimination
  - Stemming or lemmatization
  - Term construction by frequent N-Grams construction
  - Terms obtained from thesaurus (e.g., WordNet)

- BoW vector construction

- Mining of BoW vector table
  - Feature selection, Document similarity computation
  - Text mining: Categorization, Clustering, Summarization, …

# Stemming and Lemmatization

- Different forms of the same word usually problematic for text data analysis
  - because they have different spelling and similar meaning (e.g. learns, learned, learning,…)
  - usually treated as completely unrelated words
- Stemming is a process of transforming a word into its stem
  - cutting off a suffix (eg., smejala -> smej)
- Lemmatization is a process of transforming a word into its normalized form
  - replacing the word, most often replacing a suffix (eg., smejala -> smejati)

# Bag-of-Words document representation

# Word weighting

- In bag-of-words representation each word is represented as a separate variable having numeric weight.
- The most popular weighting schema is normalized word frequency TFIDF:

$$tfidf(w) = tf \cdot \log\left(\frac{N}{df(w)}\right)$$

  - Tf(w) – term frequency (number of word occurrences in a document)
  - Df(w) – document frequency (number of documents containing the word)
  - N – number of all documents
  - Tfidf(w) – relative importance of the word in the document

The word is more important if it appears several times in a target document

The word is more important if it appears in less documents

# Cosine similarity between document vectors

- Each document D is represented as a vector of TF-IDF weights

- Similarity between two vectors is estimated by the similarity between their vector representations (cosine of the angle between the two vectors):

$$Similarity\ (D_1, D_2) = \frac{\sum\limits_i x_{1i} x_{2i}}{\sqrt{\sum\limits_j x_j^2}\sqrt{\sum\limits_k x_k^2}}$$

# Advanced Topics II.

- Text mining: An introduction

  Document clustering and outlier detection

- Wordification approach to relational data mining

# Document clustering

- Clustering is a process of finding natural groups in data in a unsupervised way (no class labels pre-assigned to documents)

- Document similarity is used

- Most popular clustering methods:
  - K-Means clustering
  - Agglomerative hierarchical clustering
  - EM (Gaussian Mixture)
  - …

# Document clustering with OntoGen
# ontogen.ijs.si



**Slide adapted from D. Mladenić, JSI**

# Using OntoGen for clustering PubMed articles on autism

Work by
Petrič et al. 2009



[www.ontogen.si](http://www.ontogen.si)
Fortuna, Mladenić,
Grobelnik 2006

# K-Means clustering in OntoGen

OntoGen uses k-Means clustering for semi-automated topic ontology construction

- Given:
  - set of documents (eg., word-vectors with TFIDF),
  - distance measure (eg., cosine similarity)
  - K - number of groups
- For each group initialize its centroid with a random document
- While not converging
  - each document is assigned to the nearest group (represented by its centroid)
  - for each group calculate new centroid (group mass point, average document in the group)

# Detecting outlier documents

- By classification noise detection on a domain pair dataset, assuming two separate document corpora A and C

# Outlier detection for cross-domain knowledge discovery



*2-dimensional projection of documents (about autism (red) and calcineurin (blue). Outlier documents are bolded for the user to easily spot them.*

***Our research has shown that most domain bridging terms appear in outlier documents.***
(Lavrač, Sluban, Grčar, Juršič 2010)

# Using OntoGen for outlier document identification



**Slide adapted from D. Mladenić, JSI**

# NoiseRank: Ensemble-based noise and outlier detection

- Misclassified document detection by an ensemble of diverse classifiers (e.g., Naive Bayes, Random Forest, SVM, … classifiers)

- Ranking of misclassified documents by "voting" of classifiers

# NoiseRank on news articles

Articles on Kenyan elections: local vs. Western media

```
Rank | Class | ID | Detected by:
-----------------------------------------------------------------|
 1.      WE     352    __Bayes____RF100____RF500_____SVM____SVMEasy__SatFilt_
-----------------------------------------------------------------|
 2.      LO      25    __Bayes____RF100____RF500_____SVM____SVMEasy_
 3.      LO     101    __Bayes____RF100____RF500_____SVM____SVMEasy_
 4.      LO     173    __Bayes____RF100____RF500_____SVM____SVMEasy_
 5.      WE     348    __Bayes____RF100____RF500_____SVM____SVMEasy_
 6.      WE     326    __Bayes____RF100____RF500_____SVM____SVMEasy_
 7.      WE     357    __Bayes____RF100____RF500_____SVM____SatFilt_
 8.      WE     410    __Bayes____RF100____RF500_____SVM____SVMEasy_
-----------------------------------------------------------------|
 9.      LO      21    __RF100____RF500_____SVM____SVMEasy_
10.      LO       4    __Bayes____RF500_____SVM____SVMEasy_
11.      LO      68    __RF100____RF500_____SVM____SVMEasy_
12.      LO     162    __Bayes____RF500_____SVM____SVMEasy_
13.      WE     358    __Bayes____RF100____RF500_____SVM___
14.      WE     464    __RF100____RF500_____SVM____SVMEasy_
-----------------------------------------------------------------|
15.      LO     153    __Bayes_____SVM____SVMEasy_
16.      LO     201    __RF100____RF500___SatFilt_
17.      WE     238    __RF100____RF500_____SVM___
18.      WE     364    __Bayes____RF500_____SVM___
19.      WE     370    __Bayes____RF100_____SVM___
20.      WE     379    __RF100____RF500___SVMEasy_
```

# NoiseRank on news articles

- **Article 352: Out of topic**
  The article was later indeed removed from the corpus used for further linguistic analysis, since it is not about Kenya(ns) or the socio-political climate but about British tourists or expatriates' misfortune.

- **Article 173: Guest journalist**
  Wrongly classified because it could be regarded as a "Western article" among the local Kenyan press. The author does not have the cultural sensitivity or does not follow the editorial guidelines requiring to be careful when mentioning words like tribe in negative contexts. One could even say that he has a kind of "Western" writing style.

# Advanced Topics III.

- Text mining: An introduction
- Document clustering and outlier

Wordification approach to relational data mining

# Propositionaization through Wordification: Motivation

- Develop a RDM technique inspired by **text mining**
- Using a large number of simple, easy to understand features (**words**)
- **I**mproved **scalability**, handling large datasets
- Used as a preprocessing step to propositional learners

# Wordification Methodology

- Transform a relational database to a document corpus
  - For each individual (row) in the main table, concatenate words generated for the main table with words generated for the other tables, linked through external keys

# Text mining: Words/terms as binary features

| Document | Word1 | Word2 | … | WordN | Class |
|---|---|---|---|---|---|
| d1 | 1 | 1 | 0 | 1 | NO |
| d2 | 1 | 1 | 0 | 0 | YES |
| d3 | 1 | 1 | 1 | 1 | NO |
| d4 | 1 | 1 | 1 | 0 | YES |
| d5 | 1 | 0 | 0 | 1 | NO |
| d6-d13 | … | … | … | … | … |
| d14 | 0 | 0 | 0 | 0 | YES |
| d15 | 0 | 0 | 1 | 1 | NO |
| d16 | 0 | 0 | 1 | 0 | NO |
| d17 | 0 | 1 | 0 | 1 | NO |
| d18 | 0 | 1 | 0 | 0 | NO |
| d19-d23 | … | … | … | … | … |
| d24 | 0 | 0 | 1 | 0 | NO |

Instances = documents
Words and terms = Binary features

# Text mining

## Step 1

BoW vector construction

1. BoW features construction
2. Table of BoW vectors construction

| Document | Word1 | Word2 | … | WordN | Class |
|----------|-------|-------|---|-------|-------|
| d1 | 1 | 1 | 0 | 1 | NO |
| d2 | 1 | 1 | 0 | 0 | YES |
| d3 | 1 | 1 | 1 | 1 | NO |
| d4 | 1 | 1 | 1 | 0 | YES |
| d5 | 1 | 0 | 0 | 1 | NO |
| d6-d13 | … | … | … | … | … |
| d14 | 0 | 0 | 0 | 0 | YES |
| d15 | 0 | 0 | 1 | 1 | NO |
| d16 | 0 | 0 | 1 | 0 | NO |
| d17 | 0 | 1 | 0 | 1 | NO |
| d18 | 0 | 1 | 0 | 0 | NO |
| d19-d23 | … | … | … | … | … |
| d24 | 0 | 0 | 1 | 0 | NO |

| Document | Word1 | Word2 | … | WordN | Class |
|----------|-------|-------|---|-------|-------|
| d1 | 1 | 1 | 0 | 1 | NO |
| d2 | 1 | 1 | 0 | 0 | YES |
| d3 | 1 | 1 | 1 | 1 | NO |
| d4 | 1 | 1 | 1 | 0 | YES |
| d5 | 1 | 0 | 0 | 1 | NO |
| d6-d13 | … | … | … | … | … |
| d14 | 0 | 0 | 0 | 0 | YES |
| d15 | 0 | 0 | 1 | 1 | NO |
| d16 | 0 | 0 | 1 | 0 | NO |
| d17 | 0 | 1 | 0 | 1 | NO |
| d18 | 0 | 1 | 0 | 0 | NO |
| d19-d23 | … | … | … | … | … |
| d24 | 0 | 0 | 1 | 0 | NO |

## Step 2

Data Mining

model, patterns, clusters,

…

# Wordification Methodology

- One individual of the main data table in the relational database ~ one text document

- Features (attribute values)  ~ the words of this document

- Individual words (called **word-items** or **witems**) are constructed as combinations of:

$$[table\ name]\_[attribute\ name]\_[value]$$

-  **n-grams** are constructed to model feature dependencies:

$$[witem_1]\_[witem_2]\_ ... \_[witem_n]$$

# Wordification Methodology

- Transform a relational database to a document corpus

- Construct BoW vectors with TF-IDF weights on words

   (optional: Perform feature selection)

- Apply text mining or propositional learning on BoW table

# Wordification

**TRAIN**

| trainID | eastbound |
|---------|-----------|
| t1 | east |
| … | … |
| t5 | west |
| … | … |

**CAR**

| carID | shape | roof | wheels | train |
|-------|-------|------|--------|-------|
| c11 | rectangle | none | 2 | t1 |
| c12 | rectangle | peaked | 3 | t1 |
| … | … | … | … | … |
| c51 | rectangle | none | 2 | t5 |
| c52 | hexagon | flat | 2 | t5 |
| … | … | … | … | … |

**t1:** [car_roof_none, car_shape_rectangle, car_wheels_2, car_roof_none__car_shape_rectangle, car_roof_none__car_wheels_2, car_shape_rectangle__car_wheels_2, car_roof_peaked, car_shape_rectangle, car_wheels_3, car_roof_peaked__car_shape_rectangle, car_roof_peaked__car_wheels_3, car_shape_rectangle__car_wheels_3], **east**

# Wordification

**t1:** [car_roof_none, car_shape_rectangle, car_wheels_2, car_roof_none__car_shape_rectangle, car_roof_none__car_wheels_2, car_shape_rectangle__car_wheels_2, car_roof_peaked, car_shape_rectangle, car_wheels_3, car_roof_peaked__car_shape_rectangle, car_roof_peaked__car_wheels_3, car_shape_rectangle__car_wheels_3], **east**

**t5:** [car_roof_none, car_shape_rectangle, car_wheels_2, car_roof_none__car_shape_rectangle, car_roof_none__car_wheels_2, car_shape_rectangle__car_wheels_2, car_roof_flat, car_shape_hexagon, car_wheels_2, car_roof_flat__car_shape_hexagon, car_roof_flat__car_wheels_2, car_shape_hexagon__car_wheels_2], **west**

## TF-IDF calculation for BoW vector construction:

|    | car_shape _rectangle | car_roof _peaked | car_wheels_3 | car_roof_peaked__ car_shape_rectangle | car_shape_rectangle __car_wheels_3 | ... | class |
|----|------|------|------|------|------|-----|-----|
| t1 | 0.000 | 0.693 | 0.693 | 0.693 | 0.693 | ... | east |
| ... | ... | ... | ... | ... | ... | ... | ... |
| t5 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | ... | west |
| ... | ... | ... | ... | ... | ... | ... | ... |

# TF-IDF weights

- No explicit use of existential variables in features, TF-IDF instead

- The weight of a word indicates how relevant is the feature for the given individual

- The TF-IDF weights can then be used either for filtering words with low importance or for using them directly by a propositional learner (e.g. J48)

# Experiments

- Cross-validation experiments on 8 relational datasets: Trains (in two variants), Carcinogenesis, Mutagenensis with 42 and 188 examples, IMDB, and Financial.

- Results (using J48 for propositional learning)

# Experiments

- Cross-validation experiments on 8 relational datasets: Trains (in two variants), Carcinogenesis, Mutagenensis with 42 and 188 examples, IMDB, and Financial.

- Results (using J48 for propositional learning)
  - first applying Friedman test to rank the algorithms,
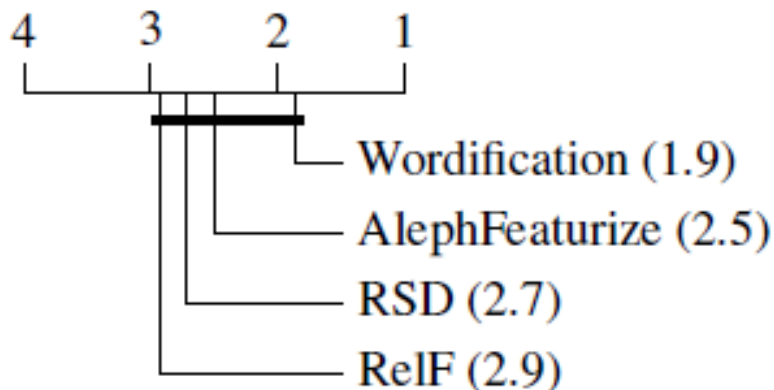  - then post-hoc test Nemenyi test to compare multiple algorithms to each other

# Experiments

- Cross-validation experiments on 8 relational datasets: Trains (in two variants), Carcinogenesis, Mutagenensis with 42 and 188 examples, IMDB, and Financial.

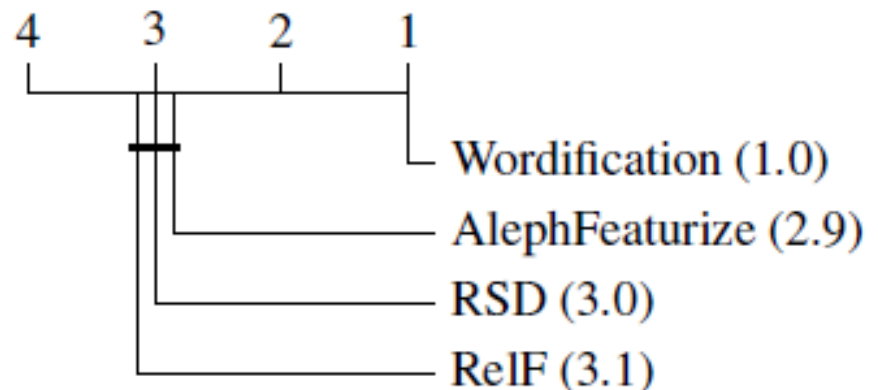- Results (using J48 for propositional learning)



MEASURE = CA

CD = 1.77

4    3    2    1

Wordification (1.9)
AlephFeaturize (2.5)
RSD (2.7)
RelF (2.9)

MEASURE = RUN-TIME

CD = 1.77

4    3    2    1

Wordification (1.0)
AlephFeaturize (2.9)
RSD (3.0)
RelF (3.1)

# Experiments

| Domain | Algorithm | J48-Accuracy[%] | J48-AUC | Run-time[s] |
|---|---|---|---|---|
| Trains | Wordification | 55.00 | 0.51 | **0.11** |
| without position | RelF | 65.00 | 0.65 | 1.04 |
| | RSD | 65.00 | 0.68 | 0.53 |
| | AlephFeaturize | **75.00** | **0.82** | 0.40 |
| | | | | |
| Trains | Wordification | **95.00** | **0.91** | **0.12** |
| | RelF | 65.00 | 0.62 | 1.06 |
| | RSD | 50.00 | 0.53 | 0.47 |
| | AlephFeaturize | 85.00 | 0.74 | 0.38 |
| | | | | |
| Mutagenesis42 | Wordification | **97.62** | **0.93** | **0.39** |
| | RelF | 80.95 | 0.59 | 2.11 |
| | RSD | **97.62** | **0.93** | 2.63 |
| | AlephFeaturize | **97.62** | **0.93** | 2.07 |
| | | | | |
| Mutagenesis188 | Wordification | **95.74** | 0.90 | **1.65** |
| | RelF | 75.53 | 0.79 | 7.76 |
| | RSD | 94.15 | **0.91** | 10.10 |
| | AlephFeaturize | 87.23 | 0.88 | 19.27 |
| | | | | |
| IMDB | Wordification | **84.34** | **0.79** | **1.23** |
| | RelF | 79.52 | 0.73 | 32.49 |
| | RSD | 73.49 | 0.47 | 4.33 |
| | AlephFeaturize | 73.49 | 0.47 | 4.96 |
| | | | | |
| Carcinogenesis | Wordification | **61.09** | **0.62** | **1.79** |
| | RelF | 54.71 | 0.53 | 16.44 |
| | RSD | 58.05 | 0.56 | 9.29 |
| | AlephFeaturize | 55.32 | 0.49 | 104.70 |
| | | | | |
| Financial | Wordification | 86.75 | 0.48 | **4.65** |
| | RelF | **97.00** | **0.91** | 260.93 |
| | RSD | 86.75 | 0.48 | 533.68 |
| | AlephFeaturize | 86.75 | 0.48 | 525.86 |

# Use Case: IMDB

- **IMDB subset:** Top 250 and bottom 100 movies
- Movies, actors, movie genres, directors, director genres
- Wordification methodology applied
- Association rules learned on BoW vector table

# Use Case: IMDB

goodMovie ← director_genre_drama, movie_genre_thriller,
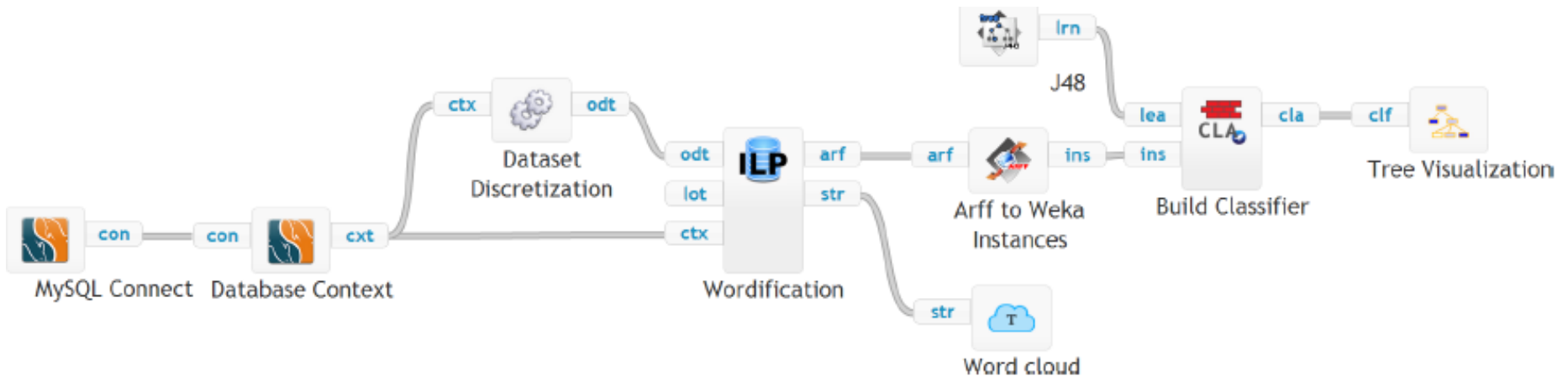            director_name_AlfredHitchcock. (Support: 5.38% Confidence: 100.00%)

movie_genre_drama ← goodMovie, actor_name_RobertDeNiro.
(Support: 3.59% Confidence: 100.00%)

director_name_AlfredHitchcock ← actor_name_AlfredHitchcock.
(Support: 4.79% Confidence: 100.00%)

director_name_StevenSpielberg ← goodMovie, movie_genre_adventure,
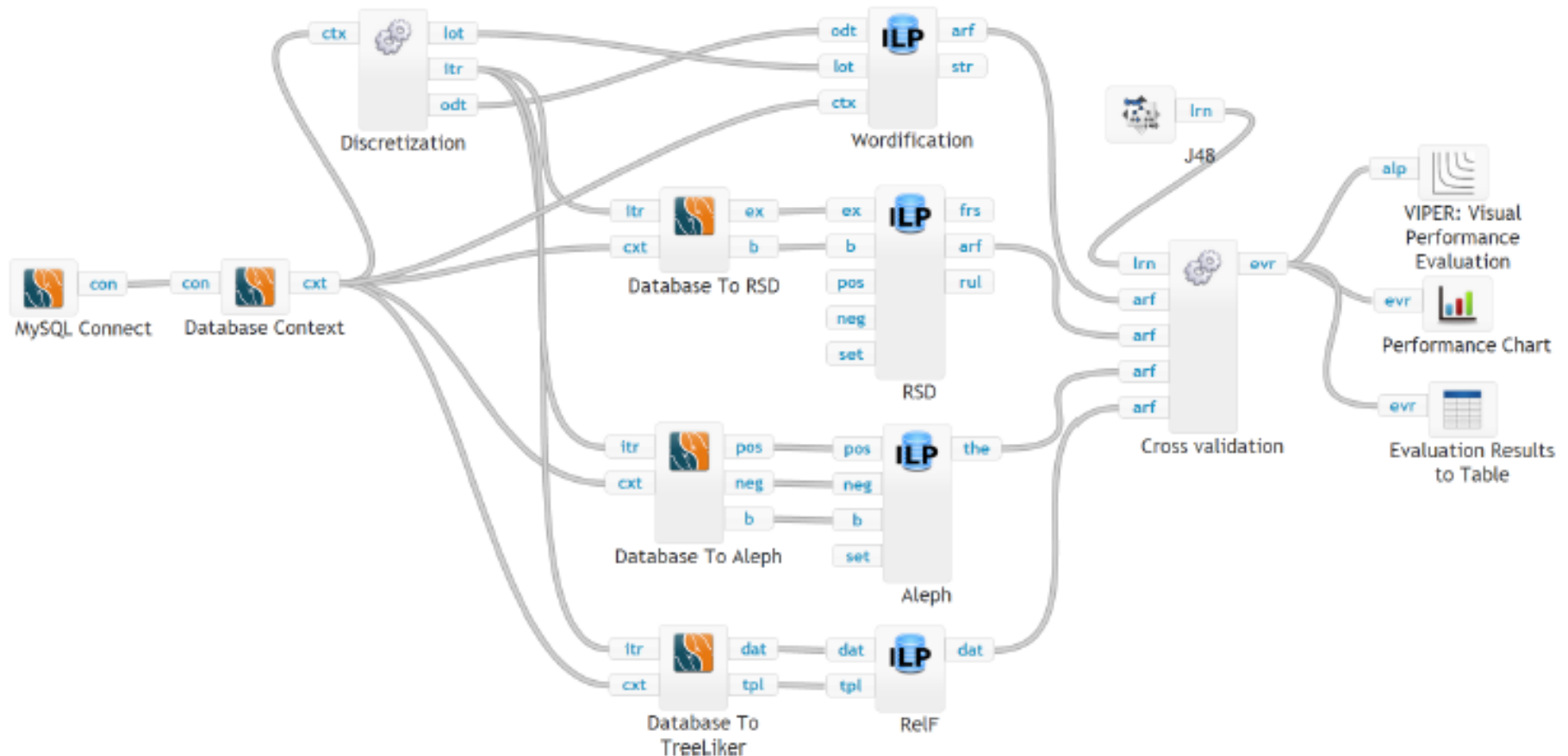(Support: 1.79% Confidence: 100.00%)                    actor_name_TedGrossman.

# Wordification implemented in ClowdFlows

- Propositionalization through wordification, available at http://clowdflows.org/workflow/1455/

# Evaluation implemented in ClowdFlows

- Wordification and propositionalization algorithms comparison, available at http://clowdflows.org/workflow/1456/

# Summary

– Wordification methodology

– Implemented in ClowdFlows

– Allows for solving non-standard RDM tasks, including RDM clustering, **word cloud visualization**, **association rule learning**, topic ontology construction, outlier detection, …