
Data Mining and Analysis

Analiza i eksploracja danych



Lecturer: JERZY STEFANOWSKI
Institute of Computing Science
Poznań University of Technology

Software Engineering – Master Course
Computer Science, PUT, revised 2015 / 2016

Background literature [Polish translations]

Translations:

- Larose D., Odkrywanie wiedzy z danych. Wprowadzanie do eksploracji danych, PWN, 2006.
- Larose D., Metody i modele eksploracji danych, PWN 2008.
- Hand D., Mannila H., Smyth P. Eksploracja danych, WNT, 2005 (Principles of Data Mining, MIT Press, 2001).



Polskie książki – Polish language books

- Tadeusz Morzy, Eksploracja danych. Metody i algorytmy. PWN 2013!!!!
- Koronacki J., Ćwik J., Statystyczne systemy uczące się, WNT 2005 (kolejne wydanie w drodze).
- Krawiec K, Stefanowski J., Uczenie maszynowe i sieci neuronowe, Wyd. PP, 2003.

Background literature (English)

- Han Jiawei and Kamber M. Data mining: Concepts and techniques, Morgan Kaufmann, 2001 (1 ed.), there is 2d
- Hand D., Mannila H., Smyth P. Principles of Data Mining, MIT Press, 2001.
- Kononenko I., Kukar M., Machine Learning and Data Mining: Introduction to Principles and Algorithms. Horwood Pub, 2007.
- Maimon O., Rokach L., The data mining and knowledge discovery Handbook, Springer 2005.
- Witten I., Eibe Frank, Data Mining, Practical Machine Learning Tools and Techniques with Java Implementations, Morgan Kaufmann, 1999.



Lecture 1 a.

Data Mining: Introduction

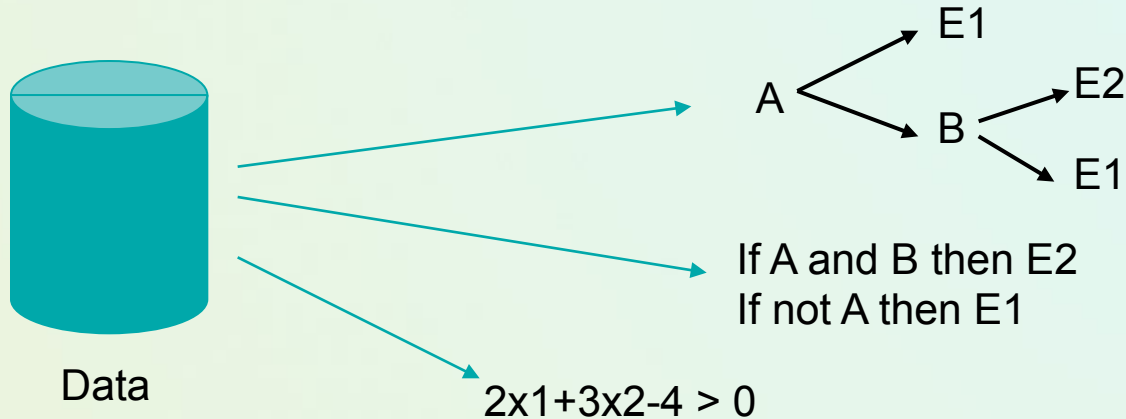
Motivations - data explosion problem

- Automated data collection tools and mature database technology lead to tremendous amounts of data stored in databases, data warehouses and other information repositories.
- More data is generated:
 - Bank, telecom, other business transactions ...
 - Scientific data: astronomy, biology, etc
 - Web, text, and e-commerce
- Very little data will ever be looked at by a human!
- We are drowning in data, but starving for knowledge!

Data Flood and Answers

- Data mining?
 - Extraction of useful information patterns from data
 - More than typical data analysis, machine learning or classical decision support!
- Knowledge Discovery is **NEEDED** to make sense and use of data.

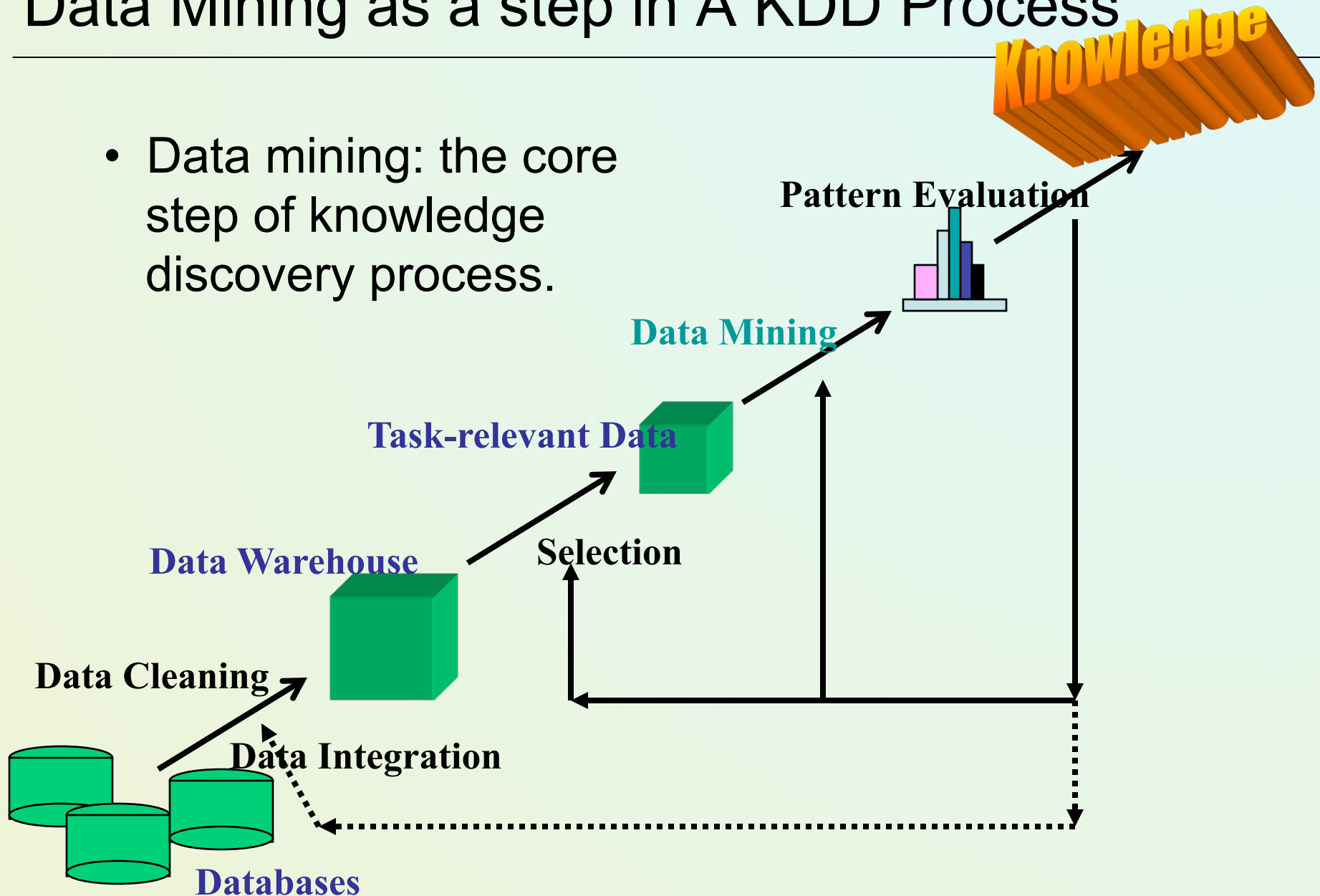
Data mining: what is it?



- **Data mining is**
 - Extraction of useful **patterns** from data sources, e.g., databases, texts, web, images.
- Patterns (**knowledge representation**) must be:
 - Valid, novel, potentially useful, understandable to the users.

Data Mining as a step in A KDD Process

- Data mining: the core step of knowledge discovery process.



Data Mining: On What Kind of Data?

- **Attribute-value tables (standard form / data table)**
- Multi-relational data / first order predicate calculus
- Structured data (graphs, workflows, ontologies, ...)
- Sequence data bases
- Other more complex data repositories
 - Object-oriented and object-relational databases
 - Spatial databases
 - Time-series data and temporal data
 - Data streams
 - Text databases and multimedia databases
 - WWW resources
 -



Flat files

- Actually the most common data source for data mining, especially at the research level.
- Simple data files in text or binary format with a structure known by the data mining algorithm to be applied.
- The data in these files can be transactions, time-series data, scientific measurements, etc.
- Big data – efficiency of access and management.

Instance	f_1	...	f_k	Y
x_1	$V_{1,1}$...	$V_{1,k}$	$V_{1,k+1}$
...
x_i	$V_{i,1}$...	$V_{i,k}$	$V_{i,k+1}$
...
x_n	$V_{n,1}$...	$V_{n,k}$	$V_{n,k+1}$

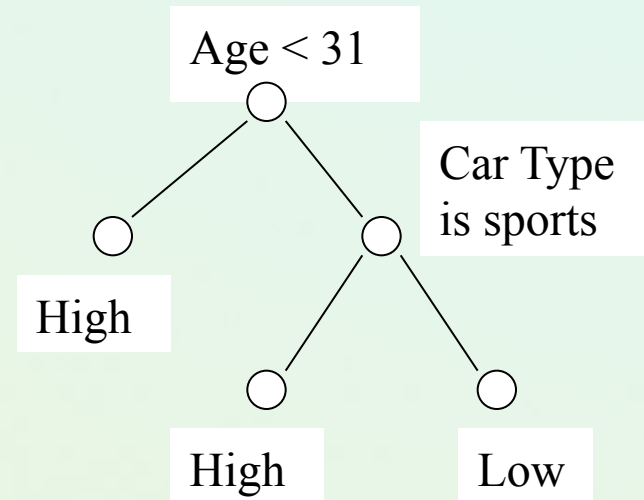
Types of attributes

- The most common distinction comes from measurement scale and statistics:
 - Nominal (also binary)
 - Ordinal
 - Interval-scaled
 - Ratio-scaled.
- Other names:
 - Categorical vs. numeric/continuous ones.
- Other types:
 - Criteria (preference-ordered), hierarchical, ...

Decision trees

- Typical approach to the classification task.

Age	Car Type	Risk
20	Combi	High
18	Sports	High
40	Sports	High
50	Family	Low
35	Minivan	Low
30	Combi	High
32	Family	Low
40	Combi	Low



Numeric prediction – regression function

- Example: 209 different computer configurations

	Cycle time (ns)	Main memory (Kb)		Cache (Kb)	Channels		Performance
	MYCT	MMIN	MMAX	CACH	CHMIN	CHMAX	PRP
1	125	256	6000	256	16	128	198
2	29	8000	32000	32	8	32	269
...							
208	480	512	8000	32	0	0	67
209	480	1000	4000	0	0	0	45

- Linear regression function

$$\text{PRP} = -55.9 + 0.0489 \text{ MYCT} + 0.0153 \text{ MMIN} + 0.0056 \text{ MMAX} \\ + 0.6410 \text{ CACH} - 0.2700 \text{ CHMIN} + 1.480 \text{ CHMAX}$$

Transforming text documents into a standard form

- Transformation into Vector Representation

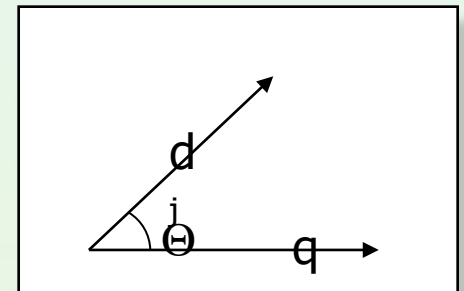
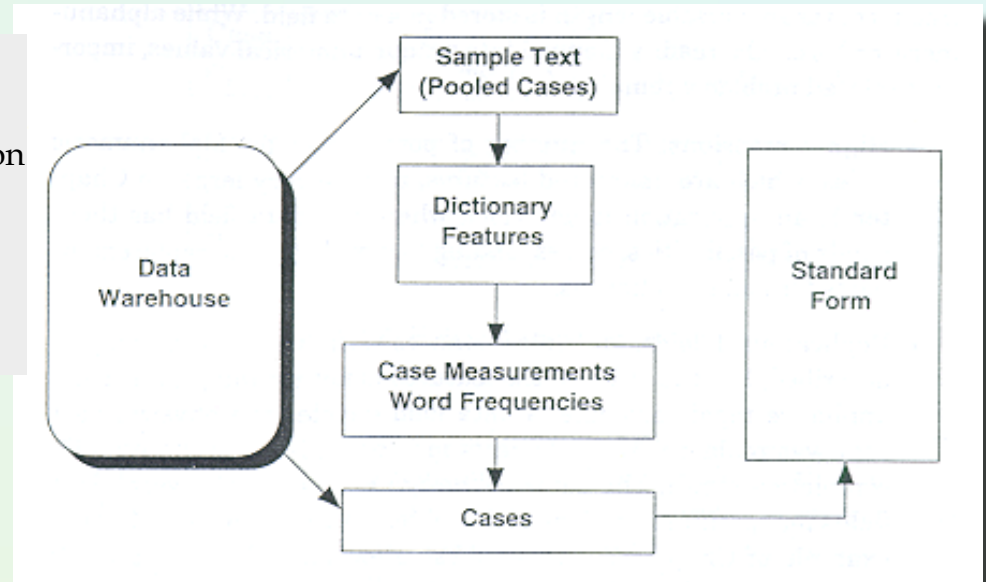
The $d=7$ documents:

- D1: Large Scale Singular Value Computations
- D2: Software for the Sparse Singular Value Decomposition
- D3: Introduction to Modern Information Retrieval
- D4: Linear Algebra for Intelligent Information Retrieval
- D5: Matrix Computations
- D6: Singular Value Analysis of Cryptograms
- D7: Automatic Information Organization

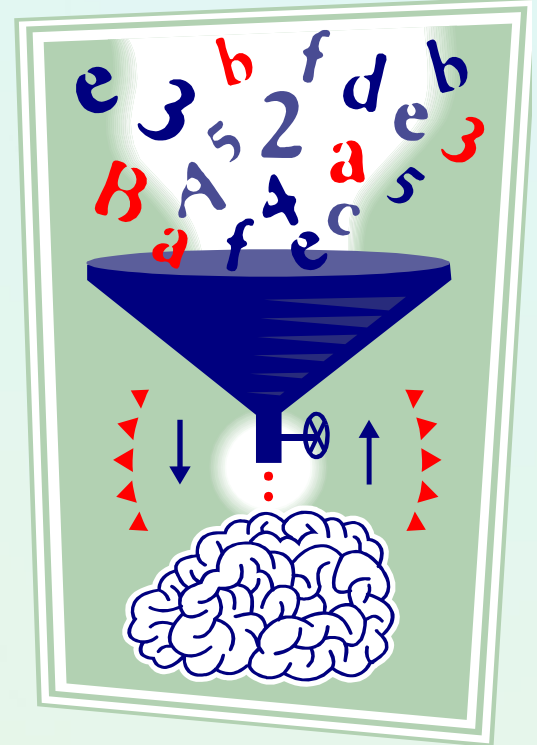
The $t=5$ terms:

- T1: Information
- T2: Singular
- T3: Value
- T4: Computations
- T5: Retrieval

$$A = \begin{pmatrix} 0.00 & 0.00 & 0.71 & 0.71 & 0.00 & 0.00 & 1.00 \\ 0.58 & 0.71 & 0.00 & 0.00 & 0.00 & 0.71 & 0.00 \\ 0.58 & 0.71 & 0.00 & 0.00 & 0.00 & 0.71 & 0.00 \\ 0.58 & 0.00 & 0.00 & 0.00 & 1.00 & 0.00 & 0.00 \\ 0.00 & 0.00 & 0.71 & 0.71 & 0.00 & 0.00 & 0.00 \end{pmatrix}$$



Data Preparation for Knowledge Discovery



A crucial issue: The majority of time / effort is put there.

Data Understanding: Quantity

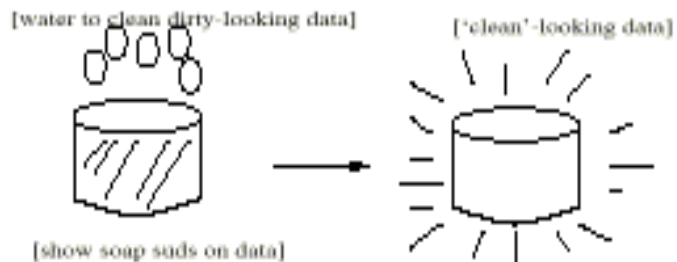
- Number of instances (records)
 - *Rule of thumb: 5,000 or more desired*
 - if less, results are less reliable; use special methods (bootstrap sampling, ...)
- Number of attributes (fields)
 - *Rule of thumb: for each field (attribute) find 10 or more instances*
 - If more fields, use feature reduction and selection
- Number of targets
 - *Rule of thumb: >100 for each class*
 - if very unbalanced, use stratified sampling or specific preprocessing (SMOTE, NCR, etc.)

Why Data Preprocessing?

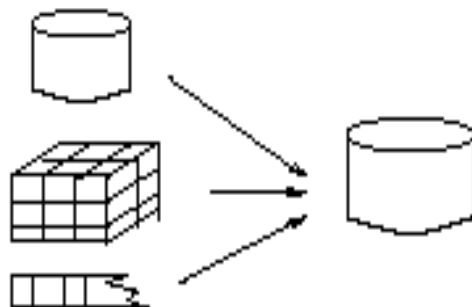
- Data in the real world is „dirty” ...
 - **incomplete**: lacking attribute values, lacking certain attributes of interest, or containing only aggregate data
 - e.g., occupation=""
 - **noisy**: containing errors or outliers
 - e.g., Salary="-10"
 - **inconsistent**: containing discrepancies (disagreements) in codes or names
 - e.g., Age="42" Birthday="03/07/1997"
 - e.g., Was rating "1,2,3", now rating "A, B, C"
 - e.g., discrepancy between duplicate records

Basic forms of data preprocessing

Data Cleaning



Data Integration



Data Transformation

-2, 32, 100, 59, 48 → -0.02, 0.32, 1.00, 0.59, 0.48

Data Reduction



From J.Han's book

Basis problems in „Data Cleaning”

- Data „acquisition” / integration and metadata
- Unified formats and other transformations
- Erroneous values
- Missing values
- Data validation and statistics

Erroneous / Incorrect values

- What suspicious can you see in this table?

Dane: TPDdatacleaning.STA 7v * 10c

	1 ID_CUST	2 CODEPOST	3 SEX	4 INCOME	5 AGE	6 MARTIALS	7 TRANS_SU
1	1001	10048	M	75 000		C	M 5000,00
2	1002	74002	F	40 000	40	W	4000,00
3	1003	90210		50 000	54	S	5400,00
4	1004	J2S7K7	F	-40 500	34	S	4500,00
5	1005	6269	M	54 000	37	M	6500,00
6	1006	45210	F	?	23	D	4500,00
7	1007	60210	M	99 450	0	M	3000,00
8	1008	65430	m	10000000	56	S	1000,00
9	1009	60211	M	3000	43	S	2400,00
10	1009	60211	M	3000	43	S	2400,00

Statystyki opisowe

Zmienne: AGE

Szczegółowe statystyki opisowe

Opcje

- Usuwanie BD przypadkami
- Wyświetl długie nazwy zmiennych
- Obliczenia zwiększonej precyzji

Rozkład

Tabele liczebności Histogramy

- Normalne częstości oczekiwane
- Testy normalności K-S i Lillieforsa
- Test W Shapiro-Wilka

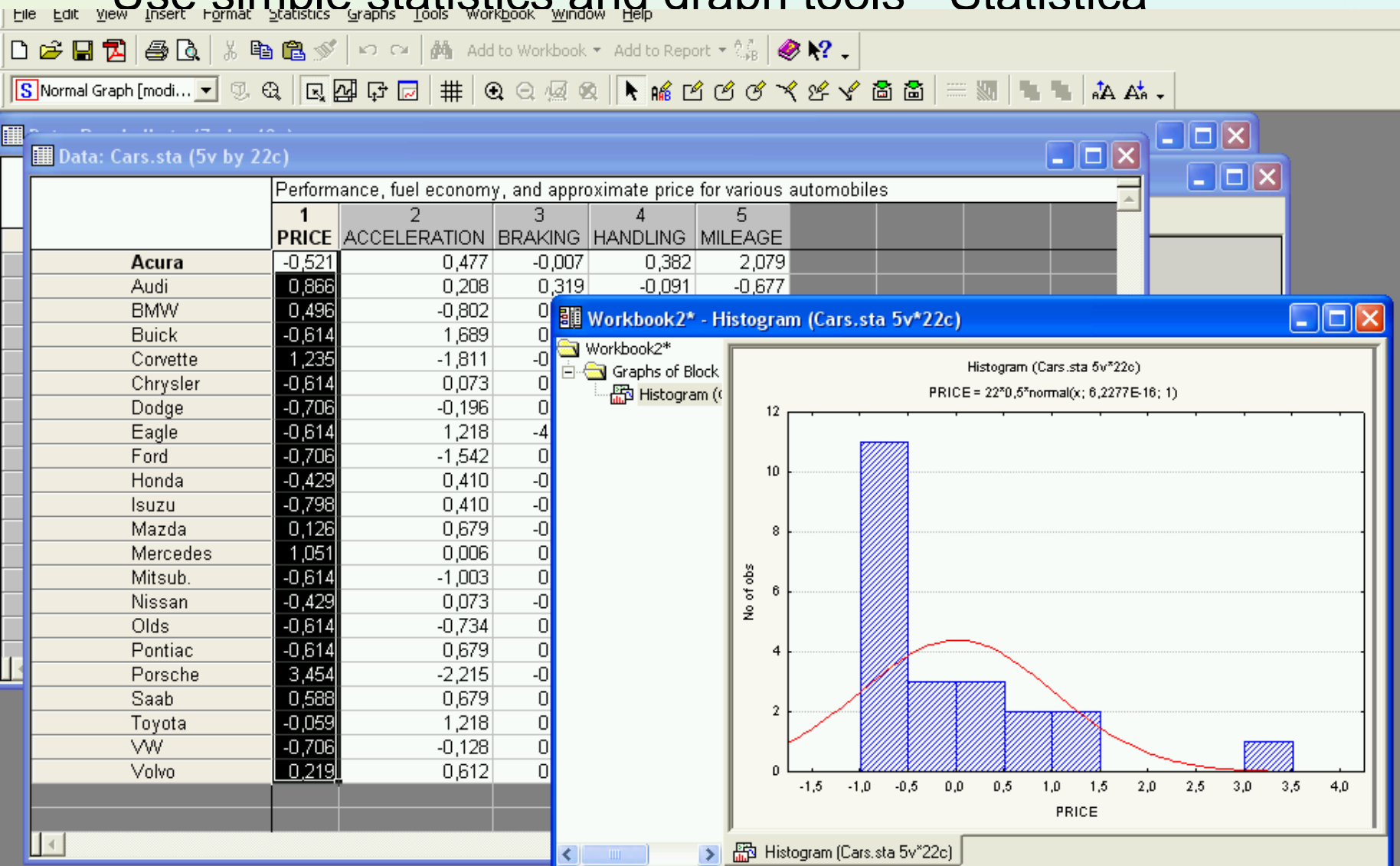
Incorrect values

- Reason: data has not been collected for mining it
- Result: errors and omissions that don't affect original purpose of data (e.g. age of customer)
- Typographical errors in nominal attributes \Rightarrow values need to be checked for consistency
- Typographical and measurement errors in numeric attributes \Rightarrow outliers need to be identified
- Errors may be deliberate (e.g. wrong zip codes)
- Other problems: duplicates, ...

Tools?

Outliers – graphical identification

- Use simple statistics and graph tools - Statistica



Redundant Data

- Redundant data occur often when integration of multiple databases
 - The same attribute may have different names in different databases
 - One attribute may be a “derived” attribute in another table, e.g., annual revenue
- Redundant data may be able to be detected by correlation analysis
- Large number of redundant data may slow-down or confuse knowledge discovery process.

Looking for correlated columns

STATISTICA - Workbook2* - [Correlations (EnginePerformance.sta)]

File Edit View Insert Format Statistics Graphs Tools Data Workbook Window Help

Add to Workbook Add to Report

Arial 10 B I U

Data: EnginePerformance.sta (79v by 128c)

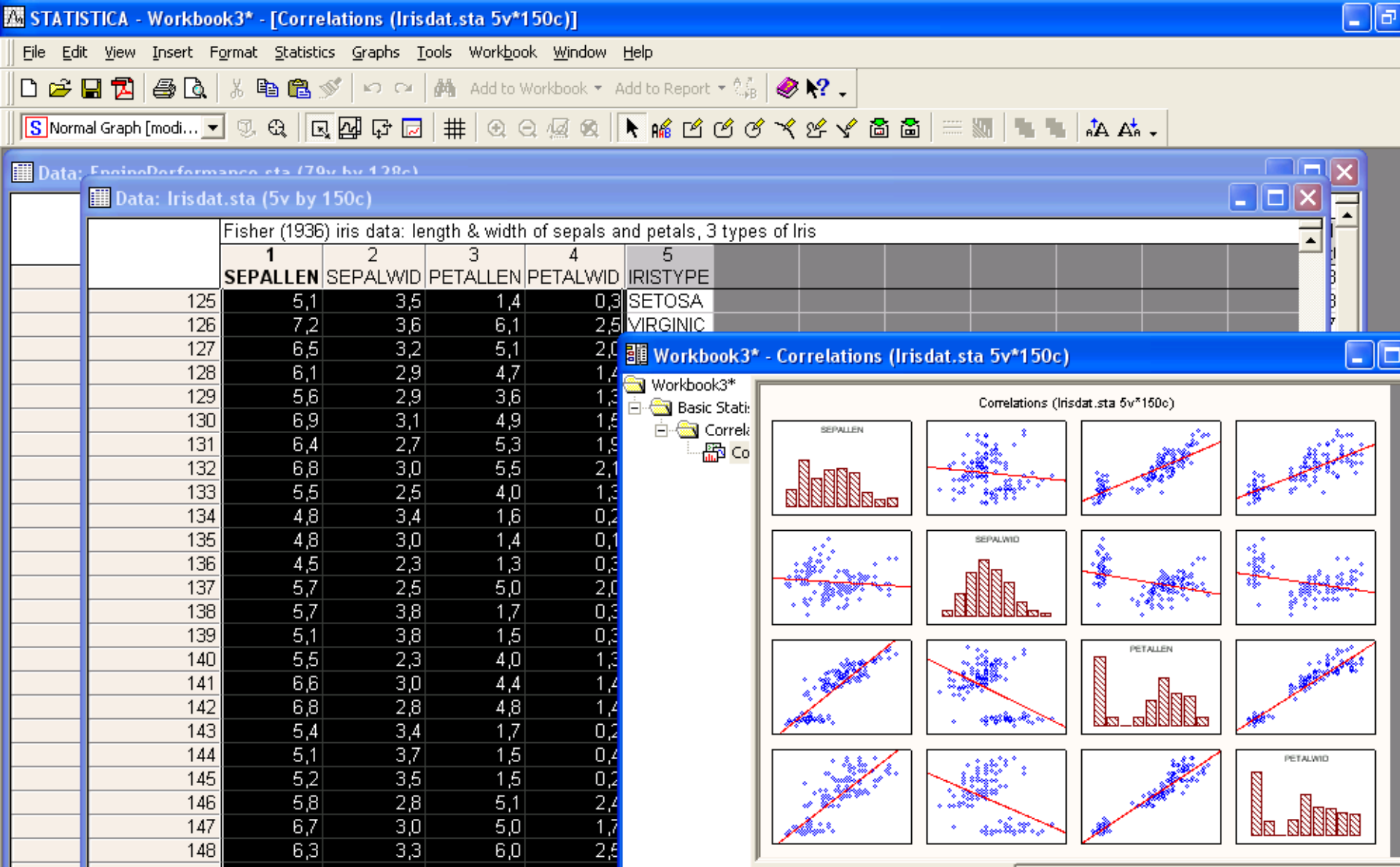
	1	2	3	4	5	6	7	8	9	10	11
	Serial Number	Efficiency	Fuel Economy(%)	Power(%)	Input01	Input02	Input03	Input04	Input05	Input06	Input07
1	#25457	102,384	100,066	99,814	100,186545	16,6255147	11,9297997	15,4501075	16,7199319	17,4754064	20,753
2	#25458	81,405	89,798	110,392	98,4136317	16,3445083	13,5326772	14,0013087	15,6347214	17,050197	20,303
3	#25459	94,070	92,072	87,917	98,7403916	16,5964348	12,0007502	15,5077475	15,7857113	18,6175749	20,527
4	#25460	108,855	89,369	90,945	99,5529412	16,7615965	12,0610633	14,2580726	13,8695801	17,8851961	19,81
5	#25461	107,903	89,453	95,912	98,8236109	16,6525248	12,2789147	14,6501313	20,634384	17,1218605	21,11
6	#25462	86,475	94,063								
7	#25463	105,583	94,868								
8	#25464	109,303	96,652								
9	#25465	103,633	91,181								
10	#25466	95,300	93,490								
11	#25467	102,334	90,320								
12	#25468	94,456	118,944								
13	#25469	109,349	107,956	1							
14	#25470	105,943	89,392								
15	#25471	101,390	102,309								
16	#25472	105,911	107,008	1							
17	#25473	78,027	91,527								
18	#25474	107,266	89,611								
19	#25475	99,571	101,998	1							
20	#25476	107,466	102,613	1							
21	#25477	109,327	95,364	1							
22	#25478	104,091	91,369								
23	#25479	95,655	90,542								
24	#25480	107,033	96,745								
25	#25481	108,802	107,768	1							
26	#25482	98,975	117,309	1							
27	#25483	104,152	100,064	1							
28	#25484	67,792	116,900								

Workbook2* - Correlations (EnginePerformance.sta)

Correlations (EnginePerformance.sta)
Marked correlations are significant at $p < ,05000$
N=128 (Casewise deletion of missing data)

Variable	Efficiency	Fuel Economy(%)	Power(%)	Input01	Input02	Input03
Efficiency	1,00	-0,09	0,12	0,12	0,19	0,
Fuel Economy(%)	-0,09	1,00	0,53	0,67	0,50	0,
Power(%)	0,12	0,53	1,00	0,26	0,14	0,
Input01	0,12	0,67	0,26	1,00	0,83	-0,
Input02	0,19	0,50	0,14	0,83	1,00	-0,
Input03	0,06	0,10	0,12	-0,01	-0,05	1,
Input04	-0,07	-0,08	0,00	-0,20	-0,23	-0,
Input05	-0,00	-0,00	0,06	-0,10	-0,04	0,
Input06	0,15	0,11	0,17	0,14	0,16	0,

Scatterplot matrix



Data Cleaning: Missing Values

- Missing data can appear in several forms:
 - <empty field> ? “0” “.” “999” “NA” ...
- Missing data may be due to
 - equipment malfunction
 - inconsistent with other recorded data and thus deleted
 - data not entered due to misunderstanding
 - certain data may not be considered important at the time of entry
 - not register history or changes of the data
- Missing data may need to be inferred – imputation!

Missing and other absent values of attributes

- Value may be missing because it is unrecorded or because it is inapplicable
- In medical data, value for **Pregnant?** attribute for **Jane** or **Anna** is missing, while for **Joe** should be considered **Not applicable**
- Don't care values

Hospital Check-in Database

Name	Age	Sex	Pregnant	..
Mary	25	F	N	
Jane	27	F	?	
Joe	30	M	-	
Anna	2	F	?	

Handle Missing Values

- Ignore / delete the instance: (not effective when the percentage of missing values per attribute varies considerably).
- Fill in the missing value manually: expert based + infeasible?
- Fill in a more advanced way :
 - a global constant : e.g., “unknown”, a new class? – **don't use it!**
 - the attribute mean or the most common value.
 - the attribute mean for all examples belonging to the same class.
 - the most probable value: inference-based such as Bayesian formula or decision tree // prediction - regression model
 - result of global closest fit (distance base approaches)
 - Use a prediction technique

Data Transformation

- Smoothing: remove noise from data
- Aggregation: summarization, data cube construction
- Generalization: concept hierarchy climbing
- Normalization: scaled to fall within a smaller, specified range
 - min-max normalization
 - z-score normalization
- Attribute/feature construction
 - New attributes constructed from the given ones

Data Transformation: Normalization

- min-max normalization

$$v' = \frac{v - \mathit{min}_A}{\mathit{max}_A - \mathit{min}_A} (\mathit{new_max}_A - \mathit{new_min}_A) + \mathit{new_min}_A$$

- z-score normalization

$$v' = \frac{v - \mathit{mean}_A}{\mathit{stand_dev}_A}$$

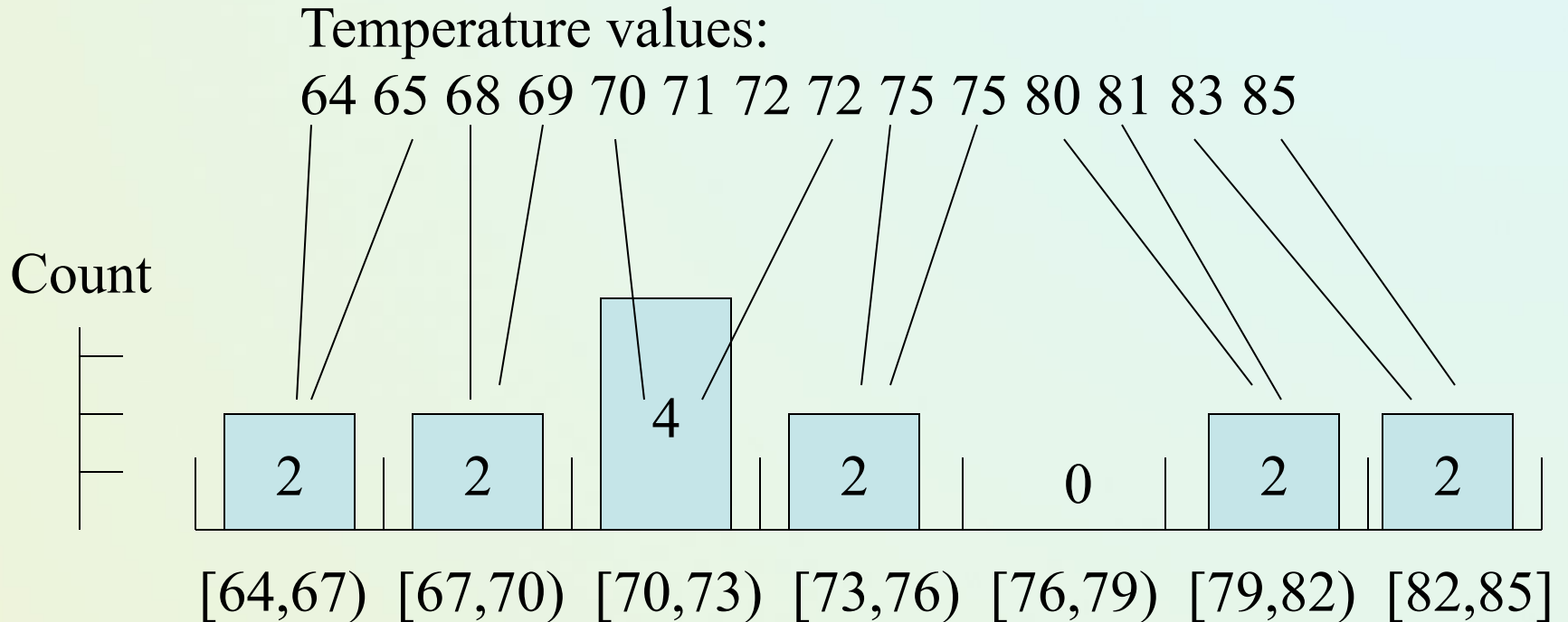
- normalization by decimal scaling

$$v' = \frac{v}{10^j} \quad \text{Where } j \text{ is the smallest integer such that } \text{Max}(|v'|) < 1$$

Discretization

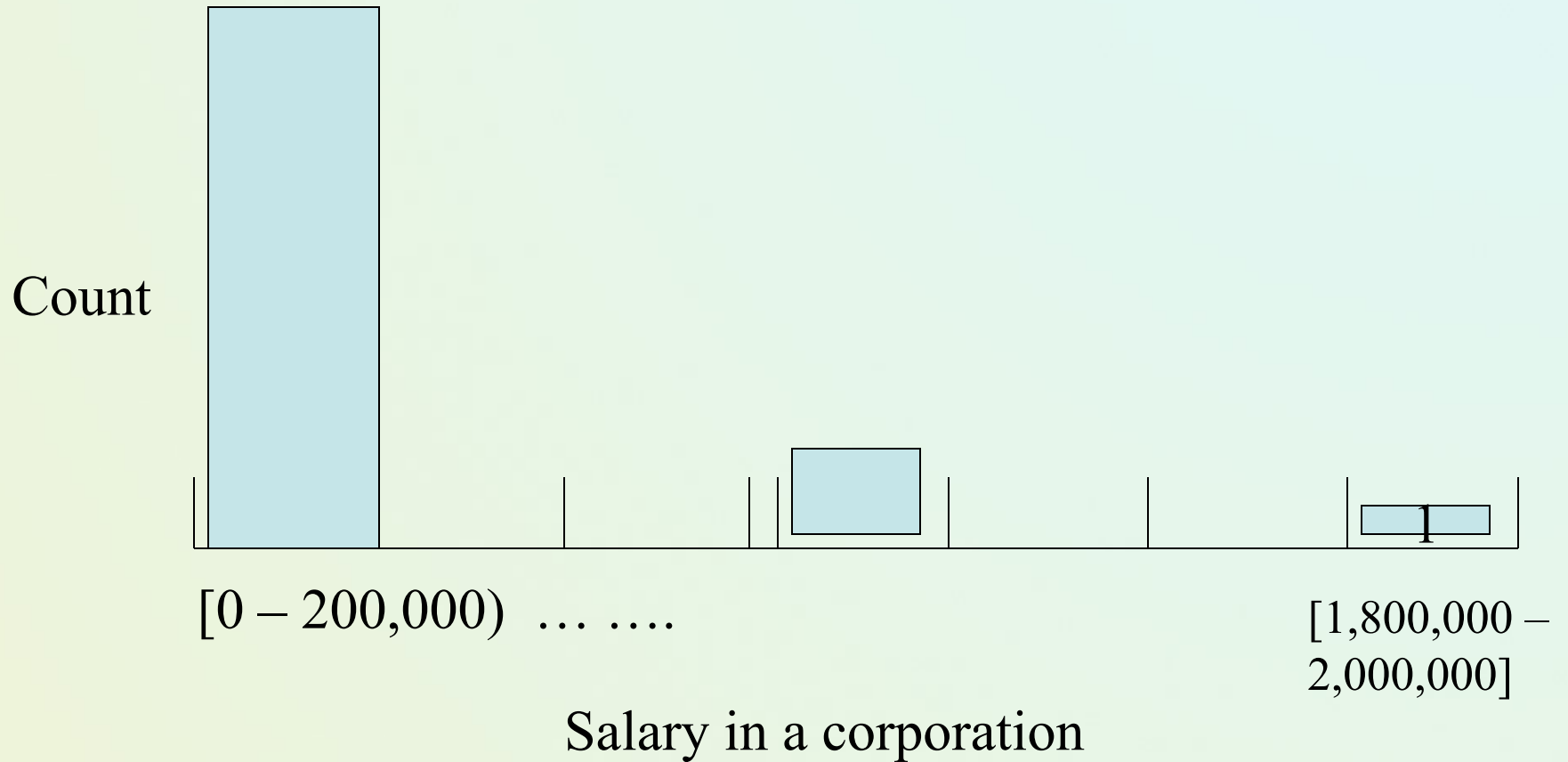
- Some methods require discrete values, e.g. most versions of Naïve Bayes, CHAID, Associations,
- Discretization → transformation of numerical values into codes / values of ordered subintervals defined over the domain of an attribute.
- Discretization is very useful for generating a summary of data
- Many approaches have been proposed:
 - Supervised vs. unsupervised,
 - Global vs. local (attribute point of view),
 - Dynamic vs. static choice of parameters

Discretization: Equal-Width (Length)



Equal Width, bins $\text{Low} \leq \text{value} < \text{High}$

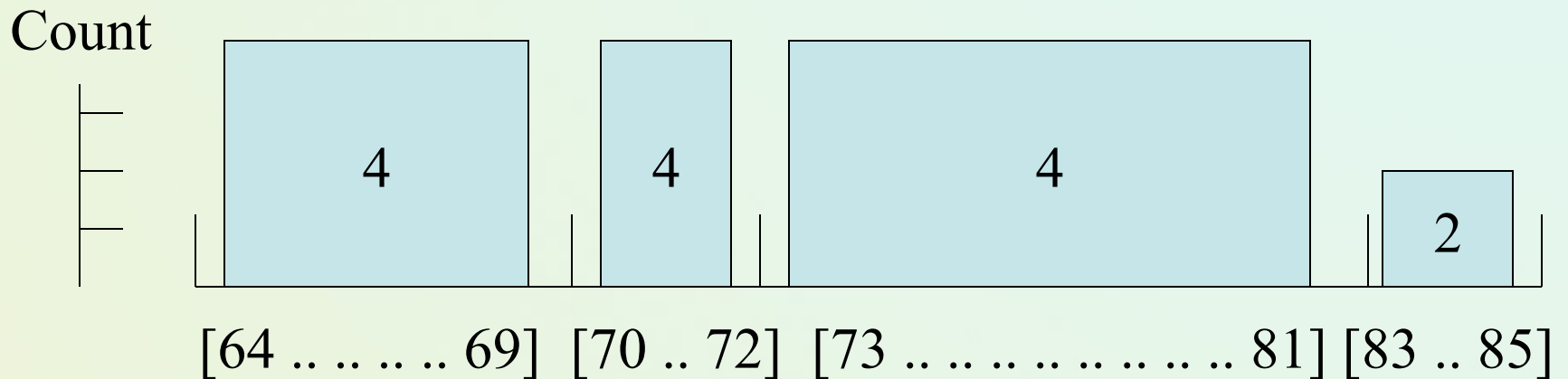
Discretization: Equal-Width may produce clumping



Discretization: Equal-Frequency

Temperature values:

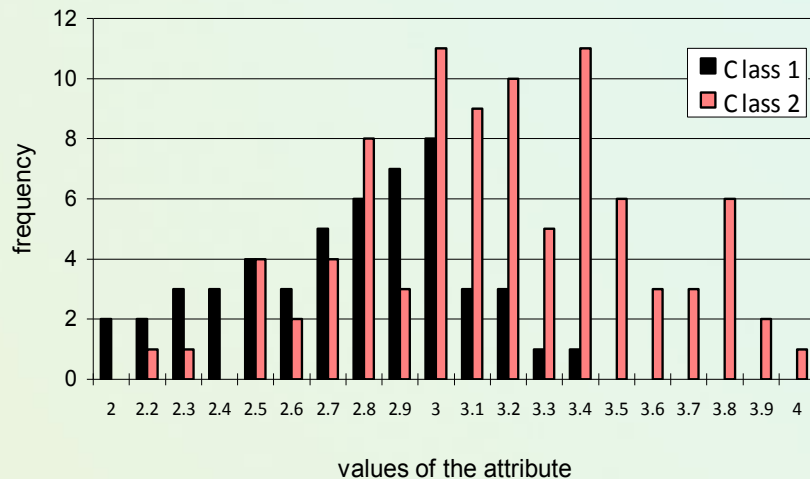
64 65 68 69 70 71 72 72 75 75 80 81 83 85



Equal Height = 4, except for the last bin

Supervised (class) discretization

- Use information about attribute value distribution + class assignment.

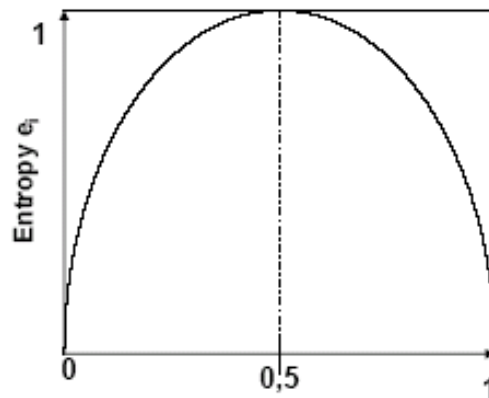


- Minimal entropy based approaches; Chi-Merge, others

Class entropy discretization

Evaluate purity of information about learning examples with **Entropy** (similar to decision tree split)

$$e_i = - \sum_{j=1}^k p_{ij} \log_2 p_{ij}$$



However, you also need a conditional entropy (with and attribute splitting)

Entropy-Based Discretization

- For learning examples S ; If S discretized into two subintervals S_1 i S_2 using (cut point) T , conditional entropy is defined as:

$$E(S, T) = \frac{|S_1|}{|S|} Ent(S_1) + \frac{|S_2|}{|S|} Ent(S_2)$$

- Scan all possible cut points
- Choose the one minimizing the entropy.
- Continue until a stopping conditions such as

$$Ent(S) - E(T, S) > \delta$$

- MDL principle could be also exploited

A Toy example

- Starting entropy $Ent(S) = -\frac{3}{6} \cdot \lg \frac{3}{6} - \frac{3}{6} \cdot \lg \frac{3}{6} = 1$
- Attribute *IQ* and a cut point $T=107$ (inter. $Left < T$)

105	107	107	109	113	115
yes	no	no	no	yes	yes

$$Ent(S | T) = \frac{1}{6}(-1 \cdot \lg 1) + \frac{5}{6}\left(-\frac{3}{5} \cdot \lg \frac{3}{5} - \frac{2}{5} \lg \frac{2}{5}\right) = 0.811$$

- Yet another cut point $T=113$ $Ent(S|T) = 0.541$ - the better choice.
- Fayyad and Irani theoretical advice – limit tested cut points

Outliers and Errors

- Outliers are values thought to be out of range.
- Approaches:
 - do nothing
 - enforce upper and lower bounds
 - let binning handle the problem

Examine Data Statistics

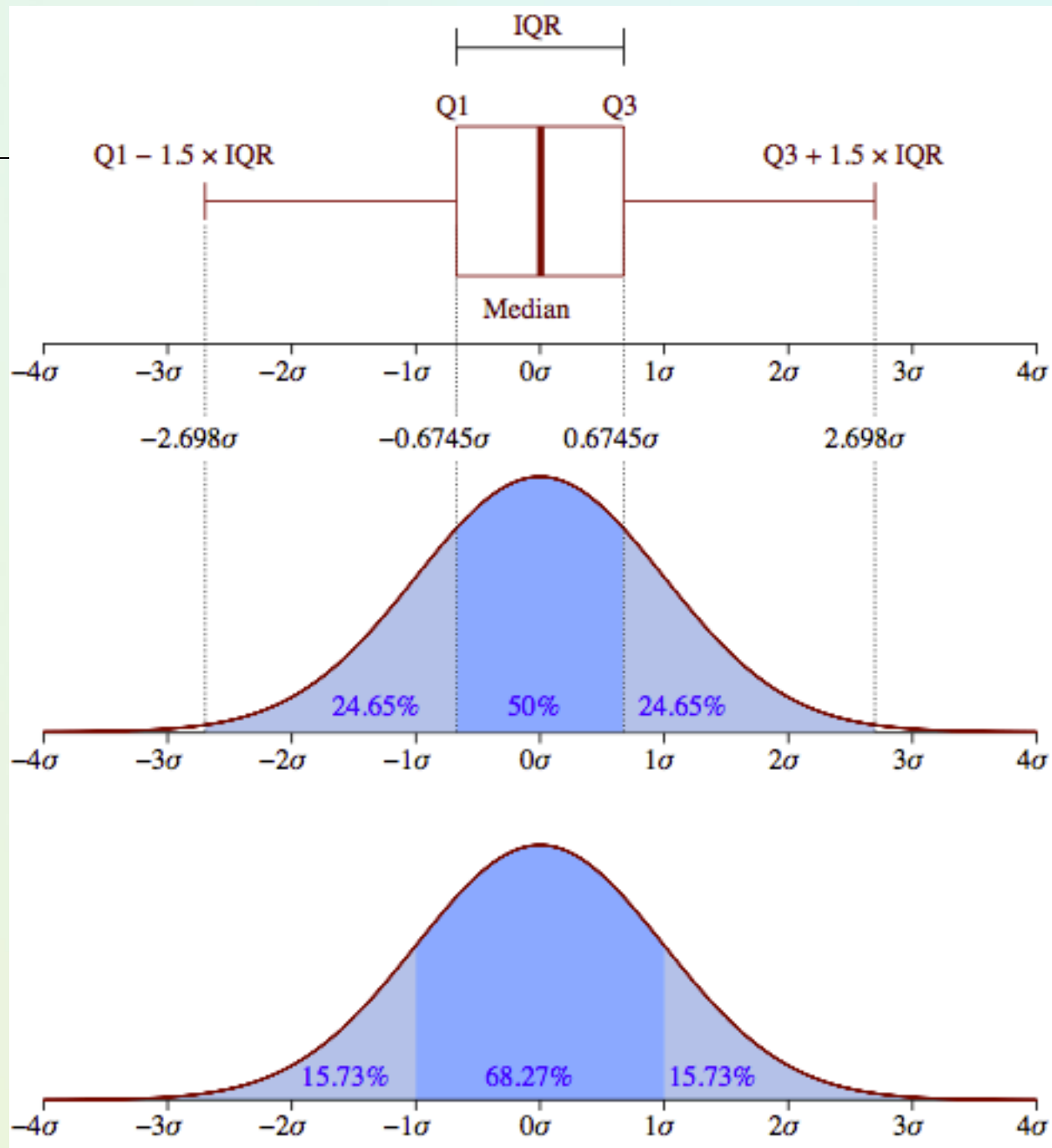
***** Field 9: MILES_ACCUMULATED

Total entries = 865636 (23809 different values). Contains non-numeric values. Missing data indicated by "" (and possibly others).

Numeric items = 165161, high = 418187.000, low = -95050.000
mean = 4194.557, std = 10505.109, skew = 7.000

Most frequent entries:

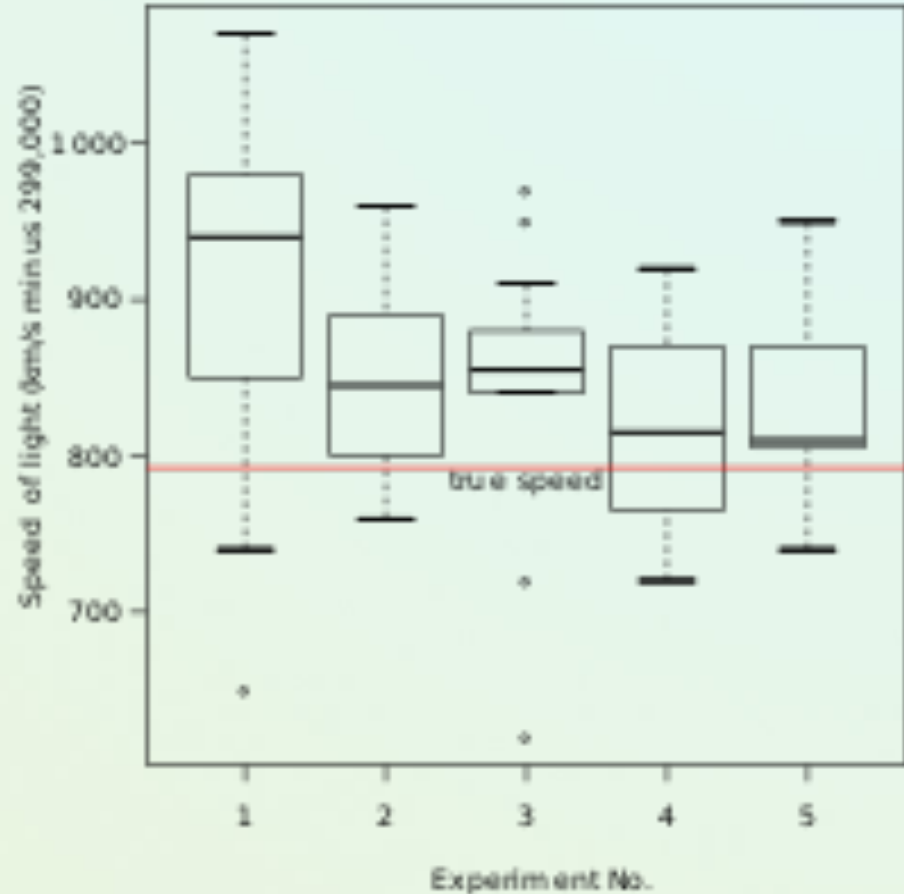
Value	Total
:	700474 (80.9%)
0:	32748 (3.8%)
1:	416 (0.0%)
2:	337 (0.0%)
10:	321 (0.0%)
8:	284 (0.0%)
5:	269 (0.0%)
6:	267 (0.0%)
12:	262 (0.0%)
7:	246 (0.0%)
4:	237 (0.0%)



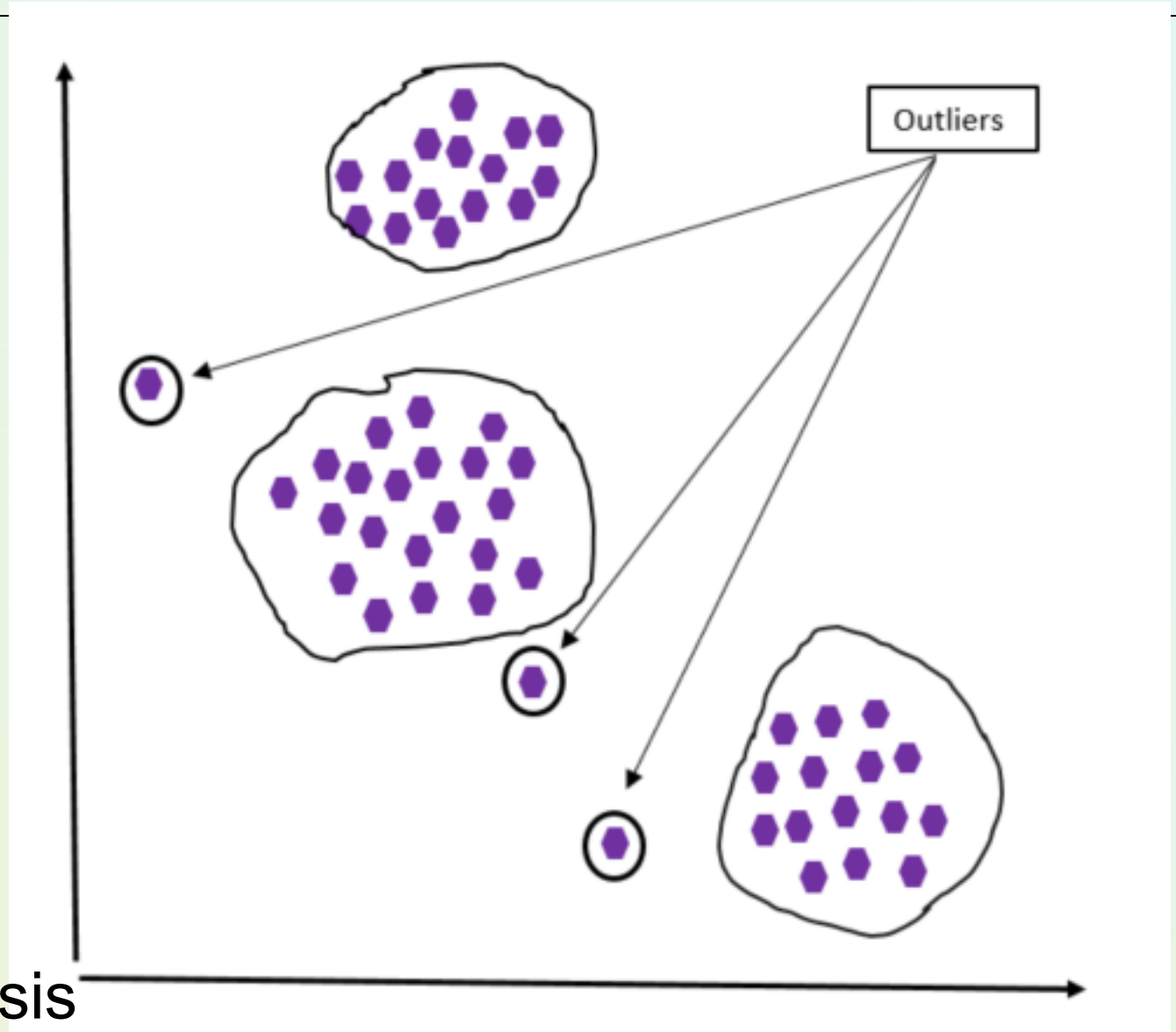
Outliers – quite far from typical distributions

Box-plots

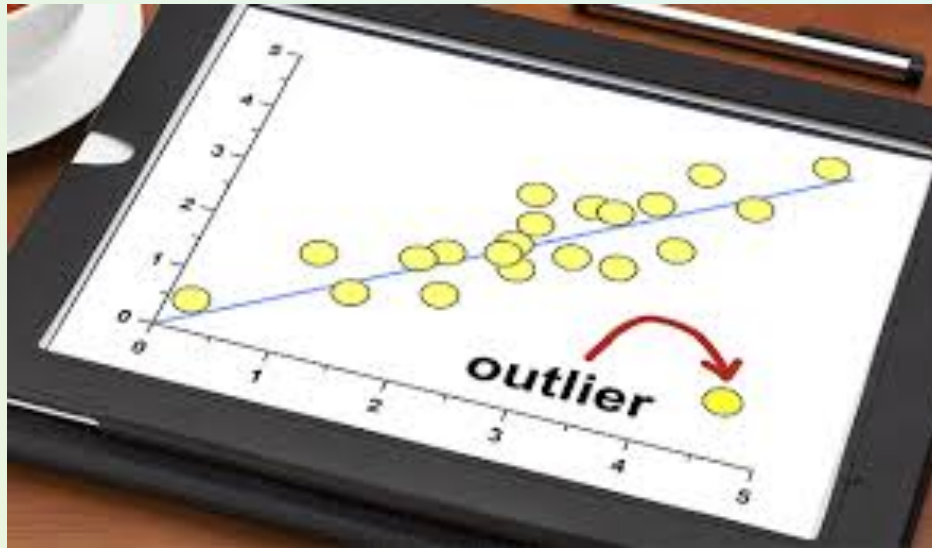
- Box plot of data from the Michelson–Morley experiment displaying four outliers in the middle column, as well as one outlier in the first column.



Multi-dimensional case



Another perspective - regression



- Numeric prediction – regression model
- Linear model $y = a_1x_1 + a_2x_2 + \dots + a_mx_m$

Data Preprocessing: Attribute Selection

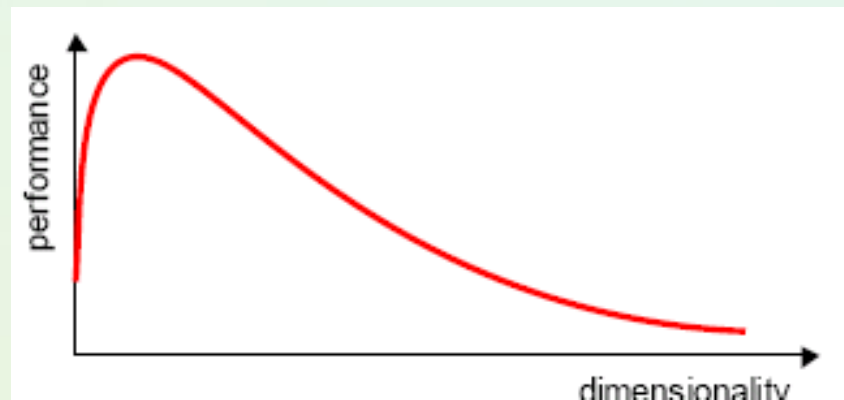
First: Remove fields with no or little variability

- Examine the number of distinct field values
 - *Rule of thumb: remove a field where almost all values are the same (e.g. null), except possibly in $minp$ % or less of all records.*
 - *$minp$ could be 0.5% or more generally less than 5% of the number of targets of the smallest class*
- More sophisticated (statistical or ML) techniques specific for data mining tasks
 - In WEKA see attribute selection

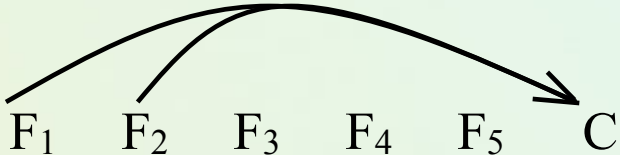
Too many attributes!

„Curse of dimensionality” [Bellman 1961]

- For a given sample size, there is a maximum number of features above which the performance of our classifier will degrade rather than improve!
- „the number of samples required per variable increases exponentially with the number of variables”



Toy classification example [D.Mladenec 2005]

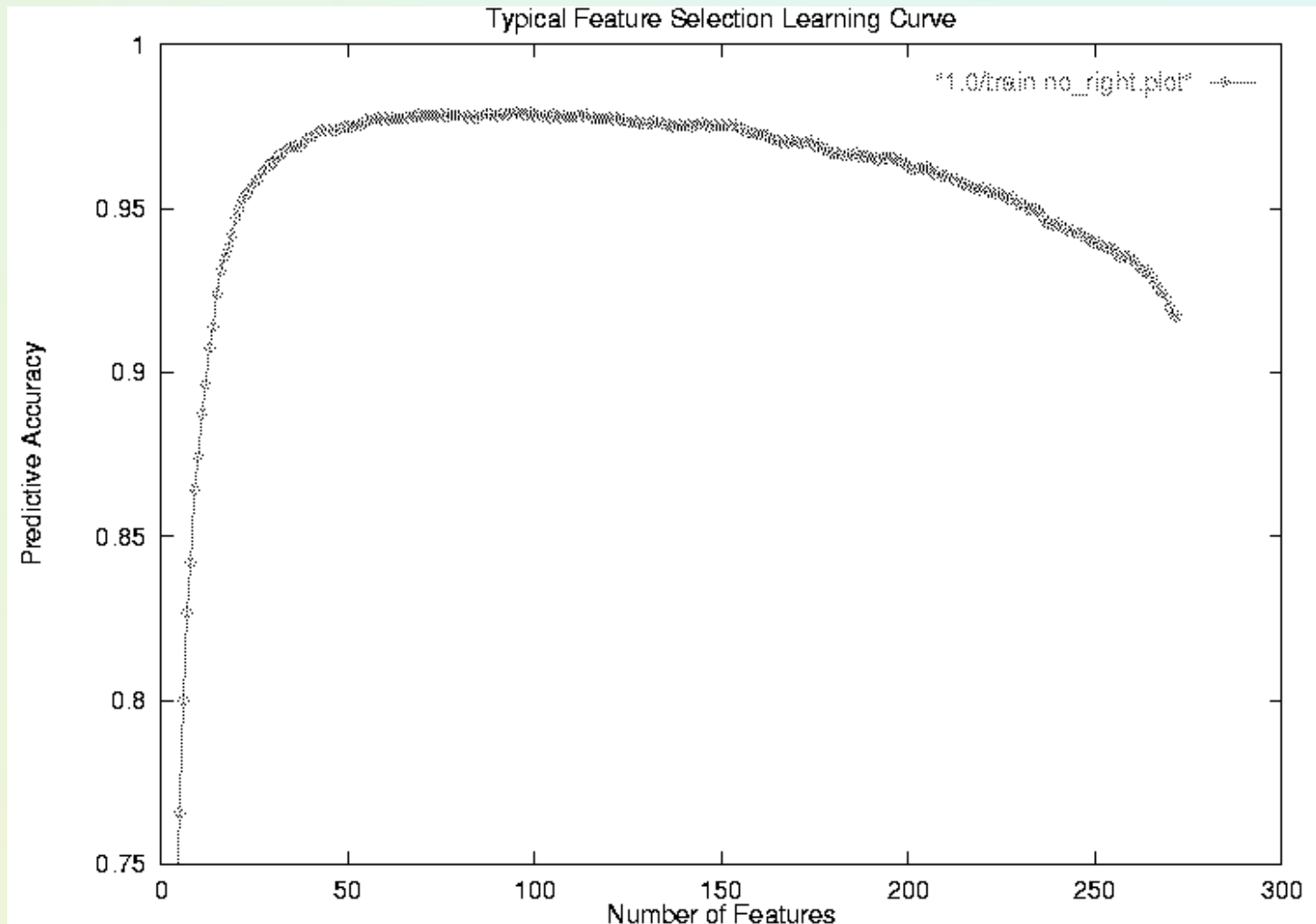


F ₁	F ₂	F ₃	F ₄	F ₅	C
0	0	1	0	1	0
0	1	0	0	1	1
1	0	1	0	1	1
1	1	0	0	1	1
0	0	1	1	0	0
0	1	0	1	0	1
1	0	1	1	0	1
1	1	0	1	0	1

- Data set
 - Five Boolean features
 - $C = F_1 \vee F_2$
 - $F_3 = \neg F_2, F_5 = \neg F_4$
 - Optimal subset:
 $\{F_1, F_2\}$ or $\{F_1, F_3\}$
- optimization in space of all feature subsets 2^F (possibilities)

(tutorial on genomics [Yu 2004])

Real working K-NN with many attributes



Different attribute selection methods

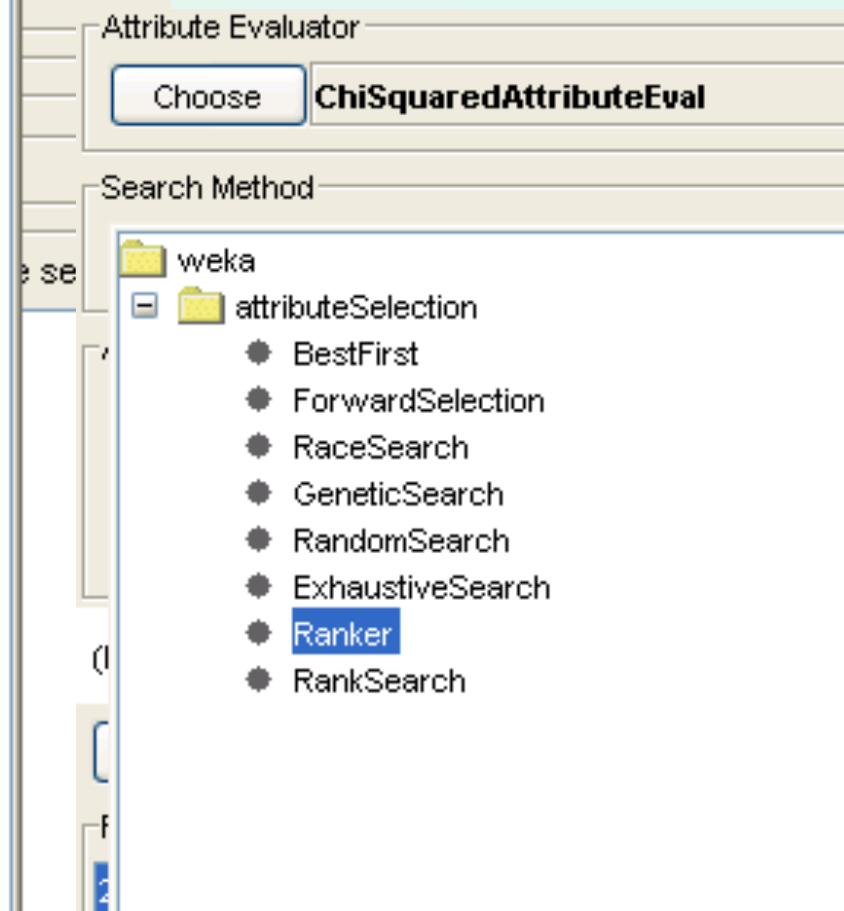
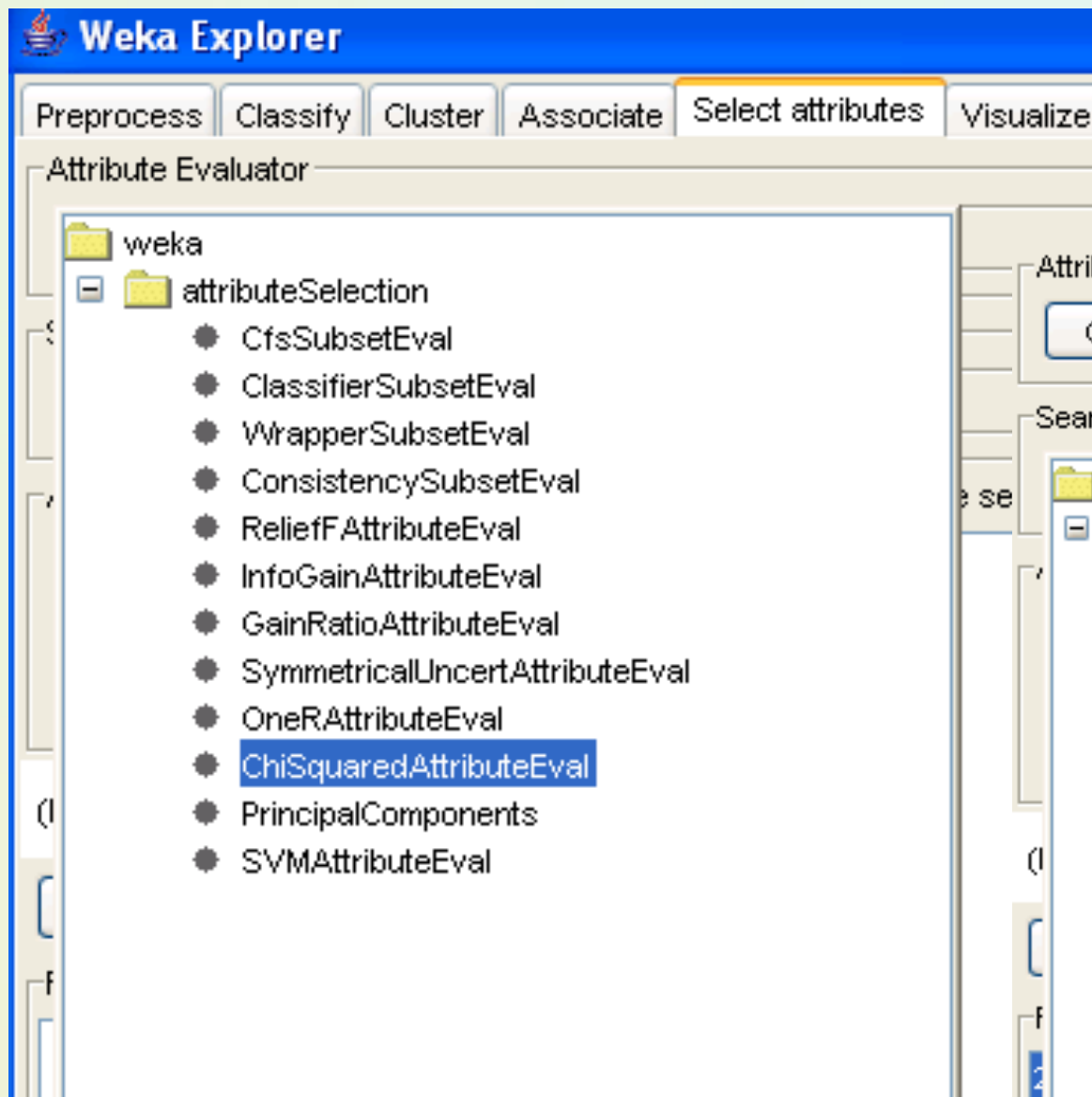
Filtering (single attributes)

- Correlation-based measure.
- Contextual-merit.
- Info-Gain.
 - Gain ratio
 - Chi-squared statistic
 - Liu Consistency measure

Subsets and more advanced search

- Relief method
- Wrapper model

WEKA – attribute selection tools



Ranking with ...? WEKA

The screenshot shows the Weka Explorer interface with the 'Select attributes' tab selected. The 'Attribute Evaluator' is set to 'ChiSquaredAttributeEval' and the 'Search Method' is 'Ranker -T -1.7976931348623157E308 -N -1'. The 'Attribute Selection Mode' is set to 'Use full training set' with 'Folds' set to 10 and 'Seed' set to 1. The 'Attribute selection output' window displays the following text:

```
AS:
D1:
Evaluation mode:  evaluate on all training data

=== Attribute Selection on all input data ===

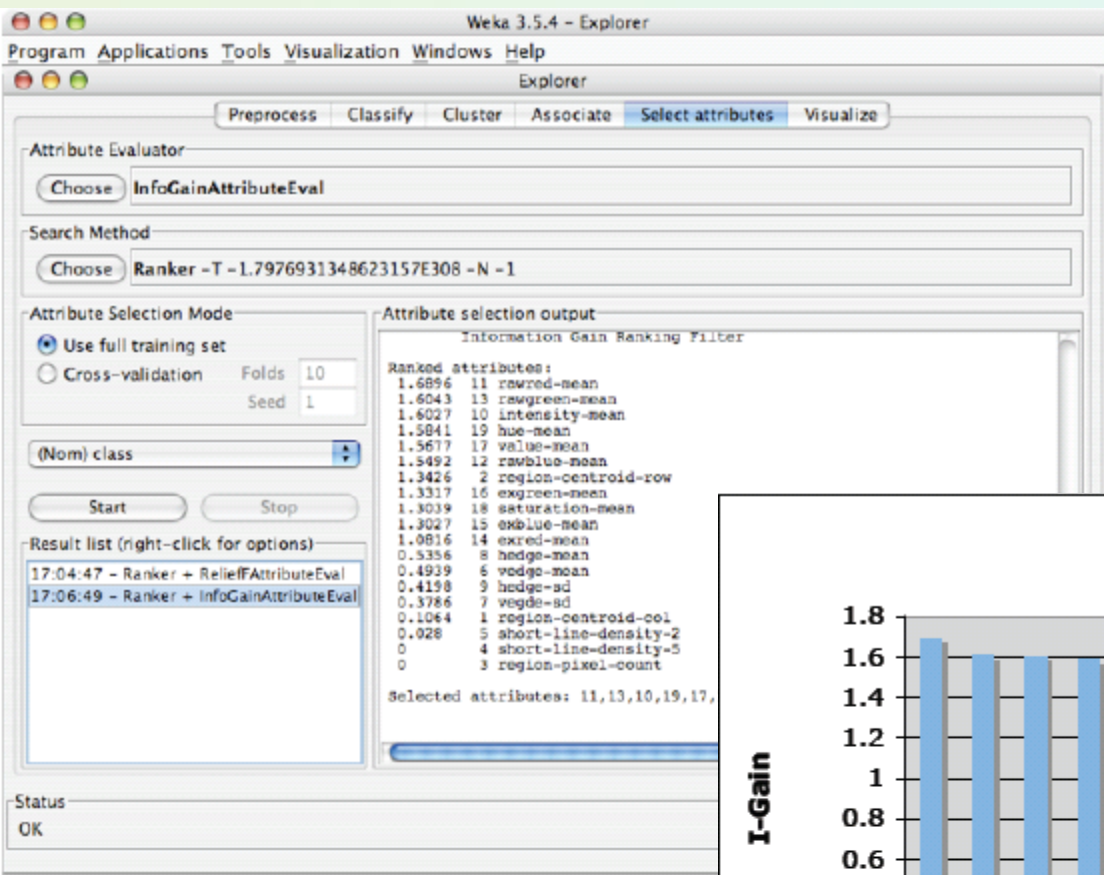
Search Method:
  Attribute ranking.

Attribute Evaluator (supervised, Class (nominal): 9 D1:):
  Chi-squared Ranking Filter

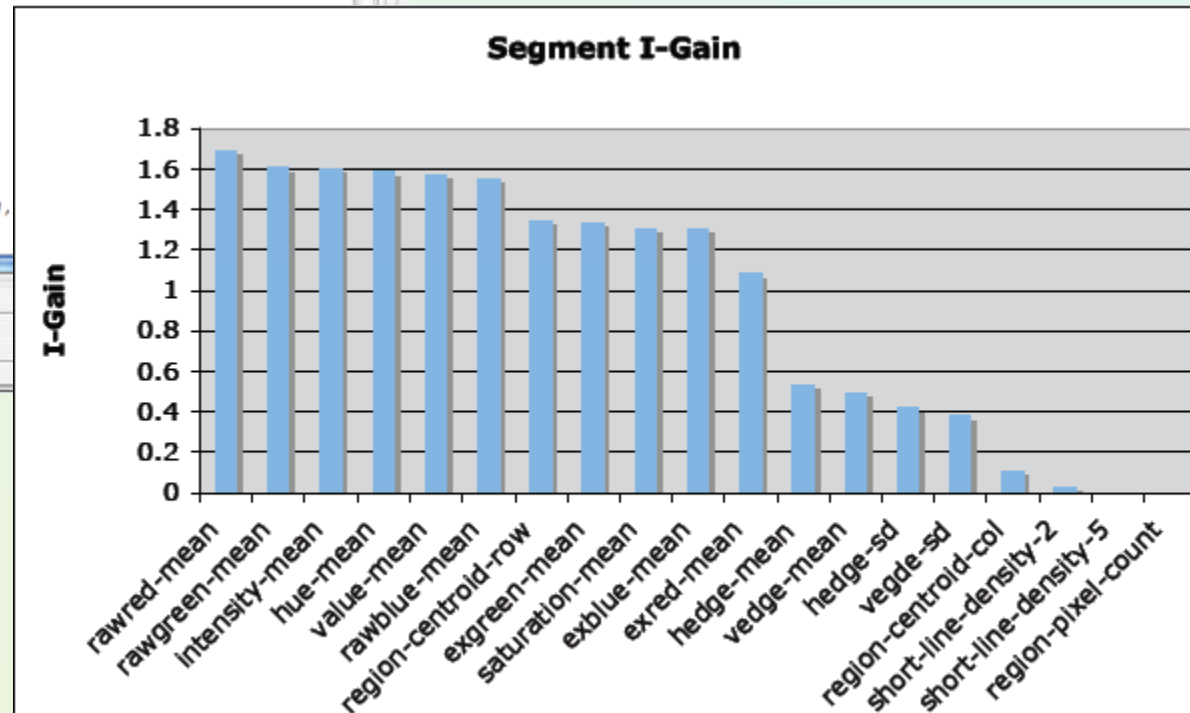
Ranked attributes:
71.9035  2  A3:
68.5634  1  A2:
67.8595  4  A5:
67.629   8  A9:
64.2122  7  A8:
64.0766  3  A4:
18.9905  5  A6:
14.0986  6  A7:

Selected attributes: 2,1,4,8,7,3,5,6 : 8
```

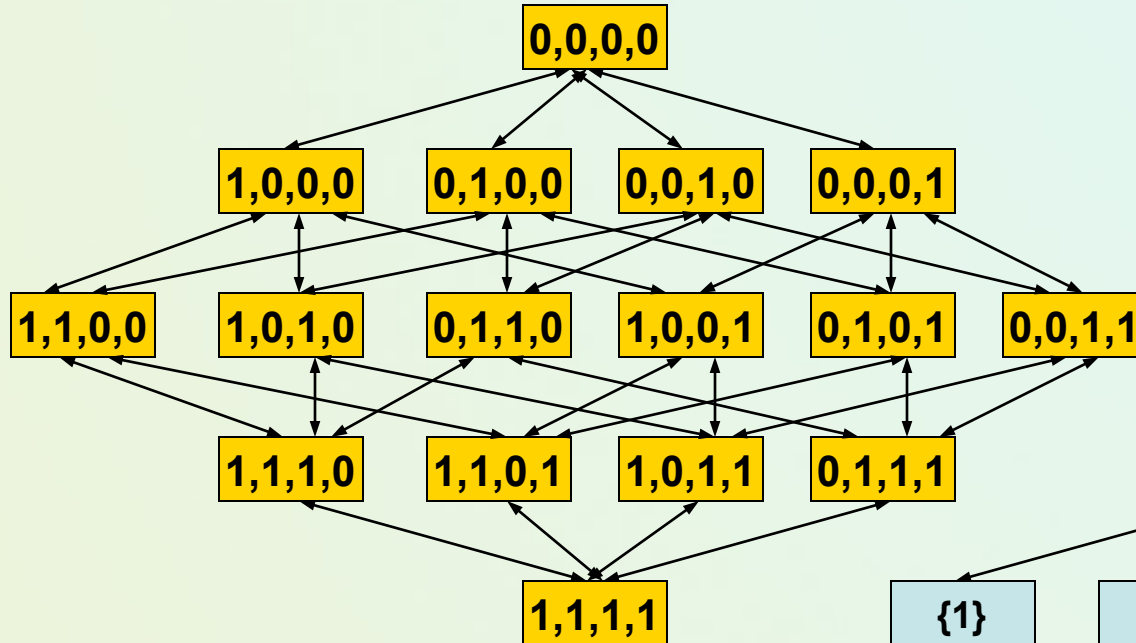
How could you exploit ranking



- Find threshold τ
- Median or smth else?

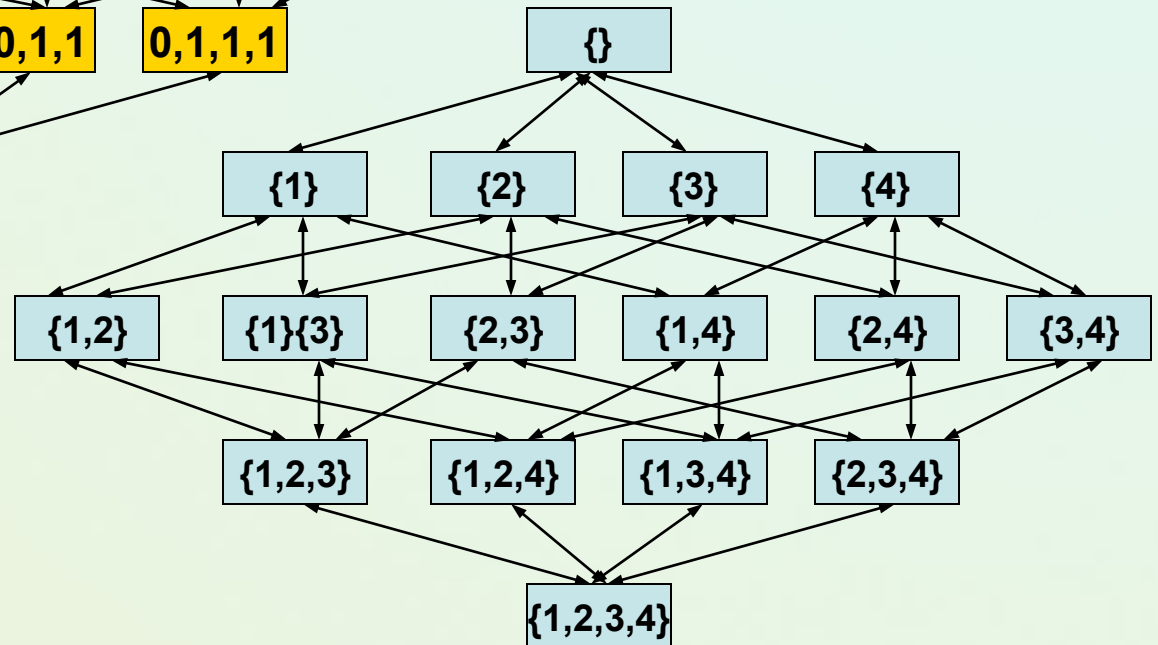


Search in Subset Space



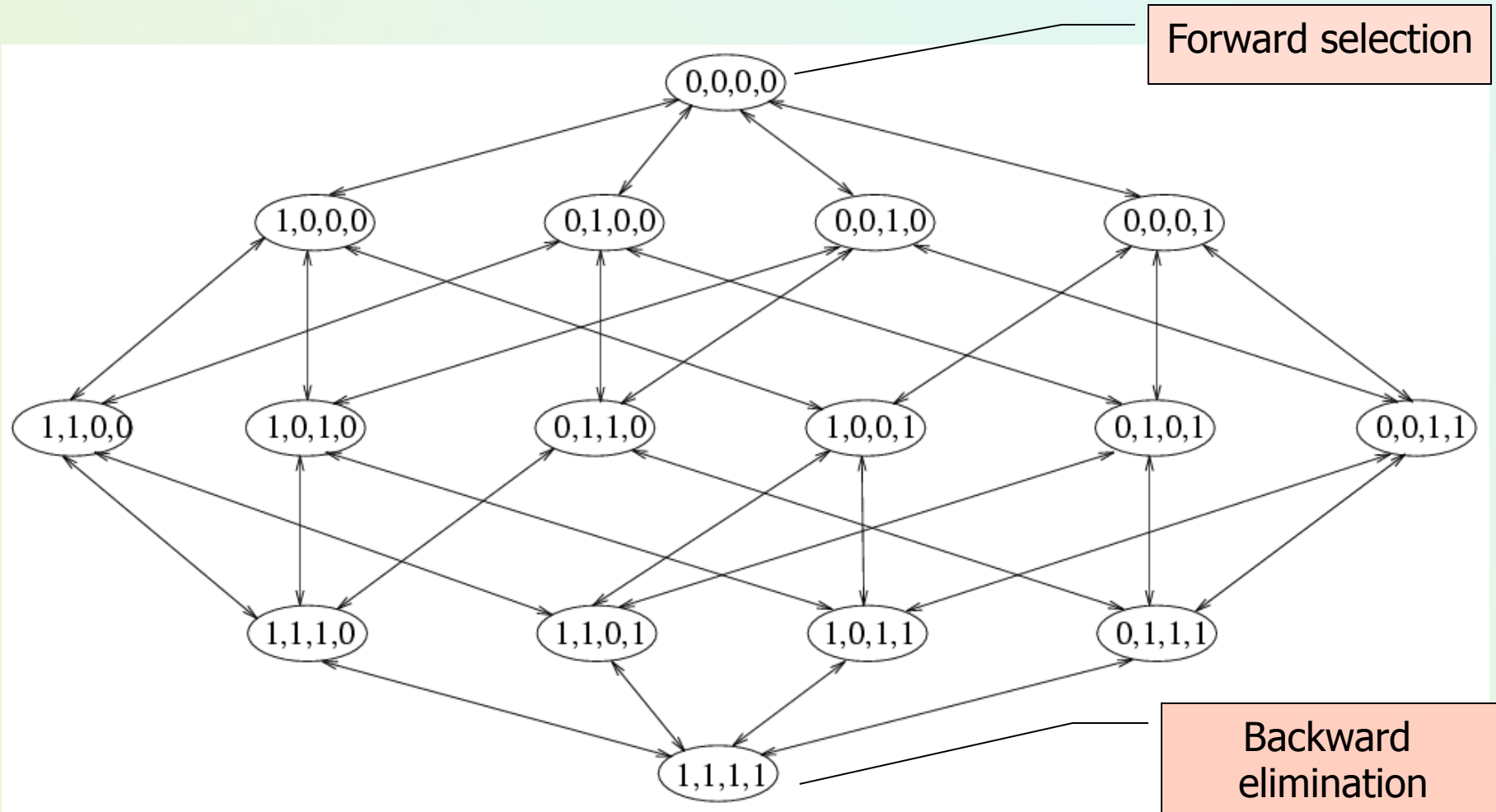
Subset Inclusion State Space
Poset Relation: Set Inclusion
 $A \leq B =$ "B is a subset of A"

"Up" operator: DELETE
"Down" operator: ADD



How to move in the space

- An example of search space (*John & Kohavi 1997*)

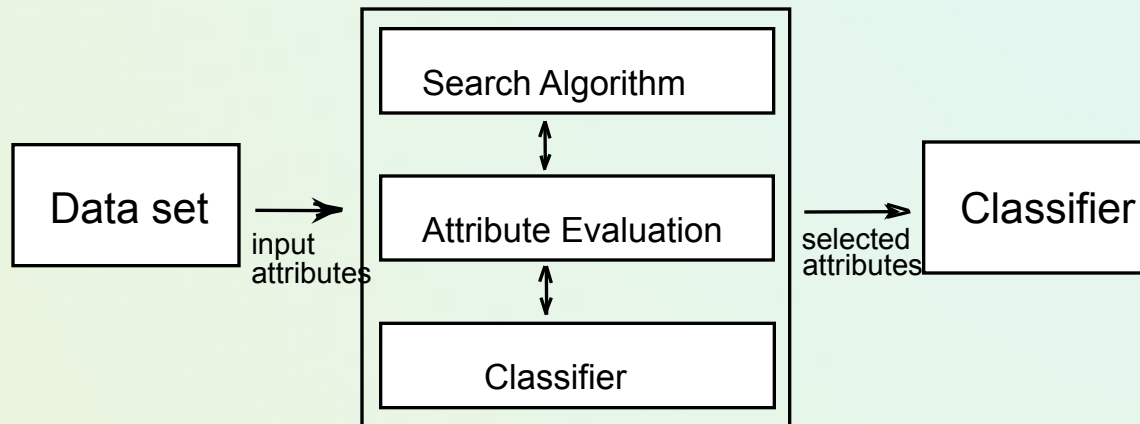


Heuristic Subset Search Techniques

- There are 2^d possible sub-features of d features
- Several heuristic feature selection methods:
 - Best step-wise feature selection (forward):
 - The best single-feature is picked first
 - Then next best feature condition to the first, ...
 - Step-wise feature elimination (backward):
 - Repeatedly eliminate the worst feature
 - Best combined feature selection and elimination
 - Partly non-deterministic search (genetic and other techniques)

Wrapper approach

- Filter vs. Wrapper approach (Kohavi et al. 94, and ...)



- The classifier is used by the evaluation function
- Search algorithms:
 - Forward selection
 - Backward elimination
 - ...

Constructing new attribute

- Following A.Berge – find new attributes

- In general - two approaches for dimensionality reduction

- Feature selection: choose a subset of the features

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ \vdots \\ x_n \end{bmatrix} \longrightarrow \begin{bmatrix} x_{i_1} \\ x_{i_2} \\ \vdots \\ x_{i_m} \end{bmatrix}$$

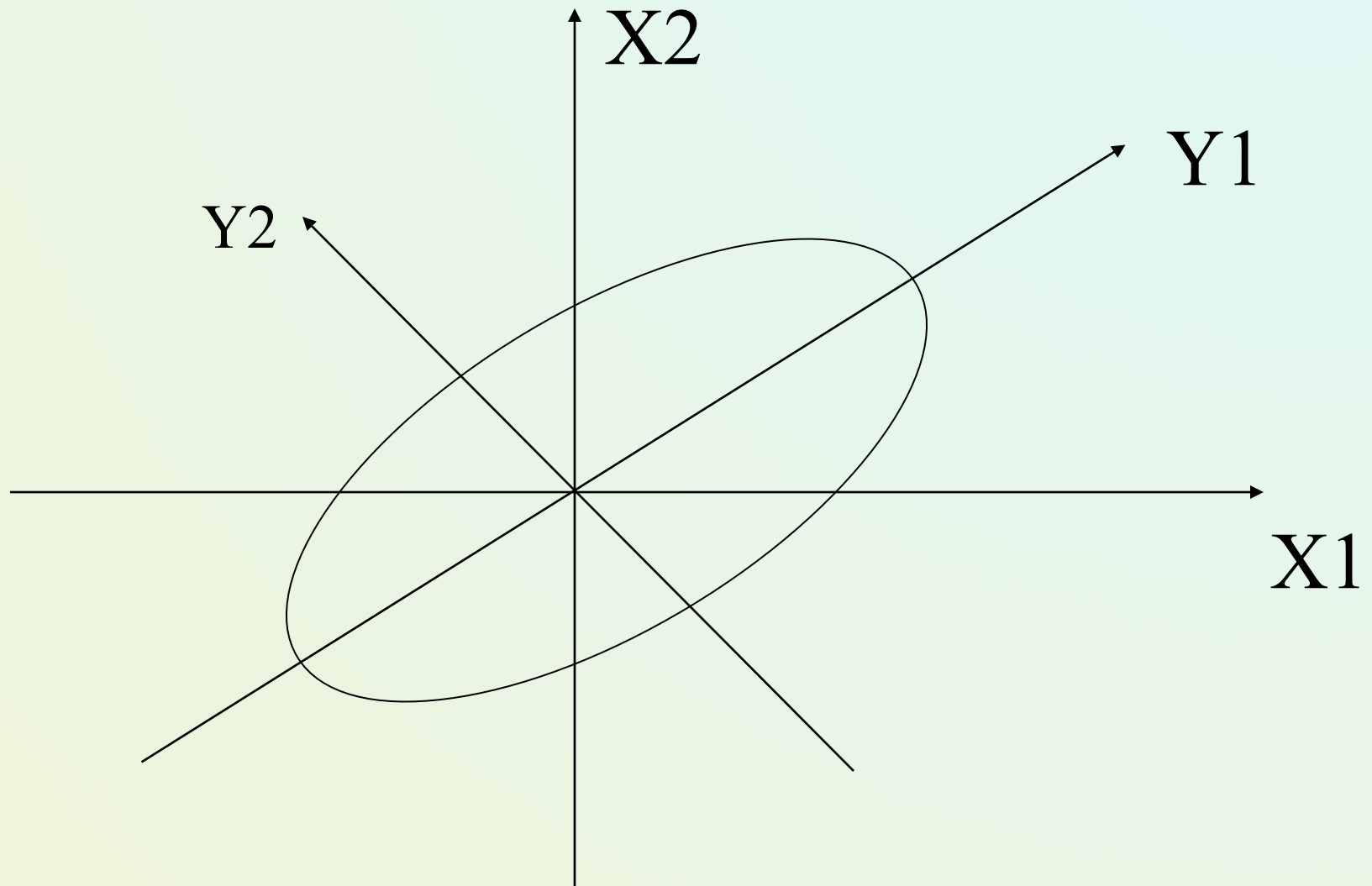
- Feature extraction: create a subset of new features by combining existing features

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ \vdots \\ x_n \end{bmatrix} \longrightarrow \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} = f \left(\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ \vdots \\ x_n \end{bmatrix} \right)$$

Principal Component Analysis (PCA)

- Given N data vectors from n -dimensions, find $k \leq n$ orthogonal vectors (*principal components*) that can be best used to represent data
- Steps
 - Normalize input data: Each attribute falls within the same range
 - Compute k orthonormal (unit) vectors, i.e., *principal components*
 - Each input data (vector) is a linear combination of the k principal component vectors
 - The principal components are sorted in order of decreasing “significance” or strength
 - Since the components are sorted, the size of the data can be reduced by eliminating the weak components, i.e., those with low variance. (i.e., using the strongest principal components, it is possible to reconstruct a good approximation of the original data)
- Works for numeric data only
- Used when the number of dimensions is large

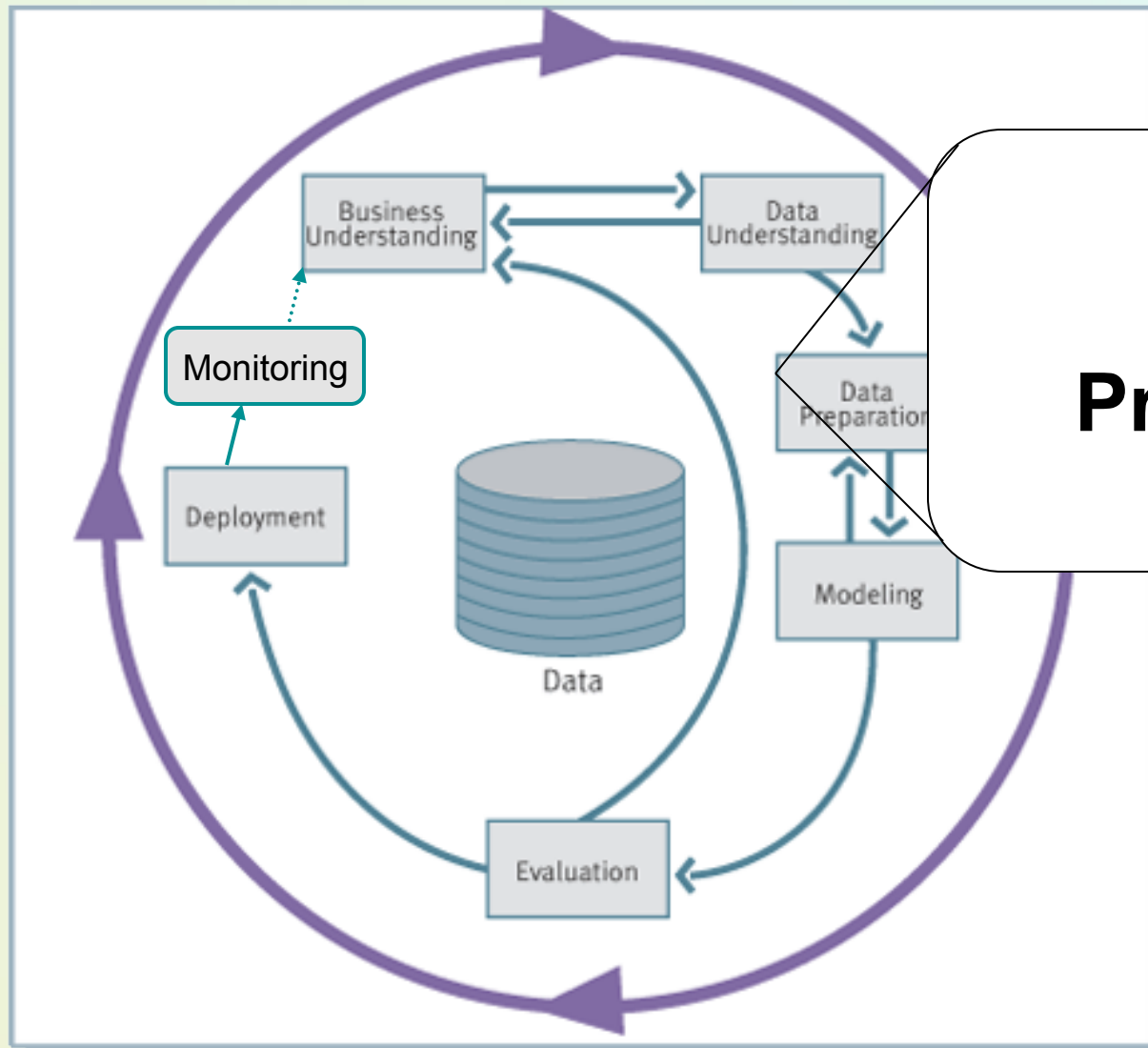
Principal Component Analysis



Summary

- Data preparation or preprocessing is a big issue for both data warehousing and data mining
- Descriptive data summarization is need for quality data preprocessing
- Data preparation includes
 - Data cleaning and data integration
 - Data reduction and feature selection
 - Discretization
- A lot a methods have been developed but data preprocessing still an active area of research

Knowledge Discovery Process, in practice



Data Preparation

Data Preparation estimated to take 70-80% of the time and effort

Any questions, remarks?

