# Data Mining: Introduction

Tom Heskes

## Instructors

Tom Heskes
tomh@cs.ru.nl
lectures/general stuff

Tom Claassen
tomc@cs.ru.nl
lectures

Roel Bouman
roel.bouman@ru.nl
coordination practical sessions/grading homework

student assistants (practical sessions, grading):

Lizzy Grootjen, David Leeftink, Steffen Ricklin, Ron Hommelsheim, Nienke Wessel, Janneke Verbeek, Linda Schmeitz, Frederik Stoel, Tamara Verbeek

# Course Outline

- 6 ects = 8 hours per week

- All, except for the exams, will be on-line

- No new lectures: video lectures from last year

- Q&A lectures on Tuesdays through Webex/Zoom

- Practical sessions on Wednesdays/Thursdays through Discord

- "Learning tasks": optional reading material + exercises for self study

- Lots of info on Brightspace

# Evaluation

- Two multiple-choice exams: one mid-term and one end-term*

- Project (more details later)

- Six homework assignments

- Mandatory: score ≥ 5.5 for at least 4 out of 6 homework assignments!

- Final grade: $0.35 \times$ Exam1 + $0.35 \times$ Exam2 + $0.3 \times$ Project**

- Re-exam: $0.7 \times$ Re-exam + $0.3 \times$ Project *

- If all else fails: full repeat next year

*On campus if corona permits. If not, we may consider dropping the midterm. Also: average of the exams needs to be at least 5.0 to pass the course.

**If final grade ≥ 5.5 (pass), average of homework assignments replaces half of the average exam grade or the project grade if it helps to give you a higher grade;.

# Lectures

- Video lectures from last year (and the year before) available in Brightspace

- Q&A sessions on Tuesday at the scheduled time

- Ask questions through the discussion board before or chat during the Q&A session

- NB: Q&A sessions will not be recorded

# Homework assignments

- Apply and understand data mining algorithms on (small) data sets

- Python:
  - open source, growing fast
  - many great data mining / machine learning packages

- Assistants at practical sessions are there to help you!

- Submit alone or as a pair; strictly follow guidelines on Brightspace!

# Practical sessions

- Starting this week on Discord on Wednesday and Thursday

- The place to ask for help, with homework or other exercises

- Feel free to pick the slot that suits you best.

- This week: get used to Python

- In preparation: download and install the 3.8 version at https://www.anaconda.com/download/

# Learning tasks

- Course Content → Learning Tasks

- To keep track and to study for exam

- Exercises are meant to practice

- Ask feedback when stuck on one of the exercises!

## Background

Data mining is the art and science of extracting knowledge out of databases. This is a rather vague and general definition that will be made more specific. What kind of data? What type of problems? What kind of techniques are available? How does data mining relate to other fields?

## Objectives

After completing this task you will be able to

- describe the objectives of data mining, its challenges, and its relationship with other fields of science;
- subdivide data mining tasks into different categories and give examples of problems for each of these.

## Instructions

1. Read and study chapter 1 of TSK.
2. Make exercises 1 through 3 of TSK, section 1.7.
3. What is the definition of data mining in TSK? Find two other definitions for data mining and compare them.
4. Find at least two examples of data mining applications that appeared in the press (the more recent and the closer to home, the better...). Describe these. What data mining tasks are involved?
5. *Data mining is very closely related to machine learning. Check out this note to learn about its aims, success, and challenges.

## Products

- Answers to the exercises.
- Three different definitions of data mining.
- Two "real-world" examples of data mining.

## Reflection

- Can you explain the difference between data mining and statistics, knowledge discovery, machine learning, and so on?
- Given a particular problem, can you tell what data mining task it belongs to?
- Can you describe some challenges in machine learning/data mining?

Radboud University Nijmegen

# Contents

- We will closely follow the **first edition** of the book "Introduction to Data Mining" by Tan, Steinbach and Kumar

- See

  https://www-users.cs.umn.edu/~kumar001/dmbook/firsted.php/
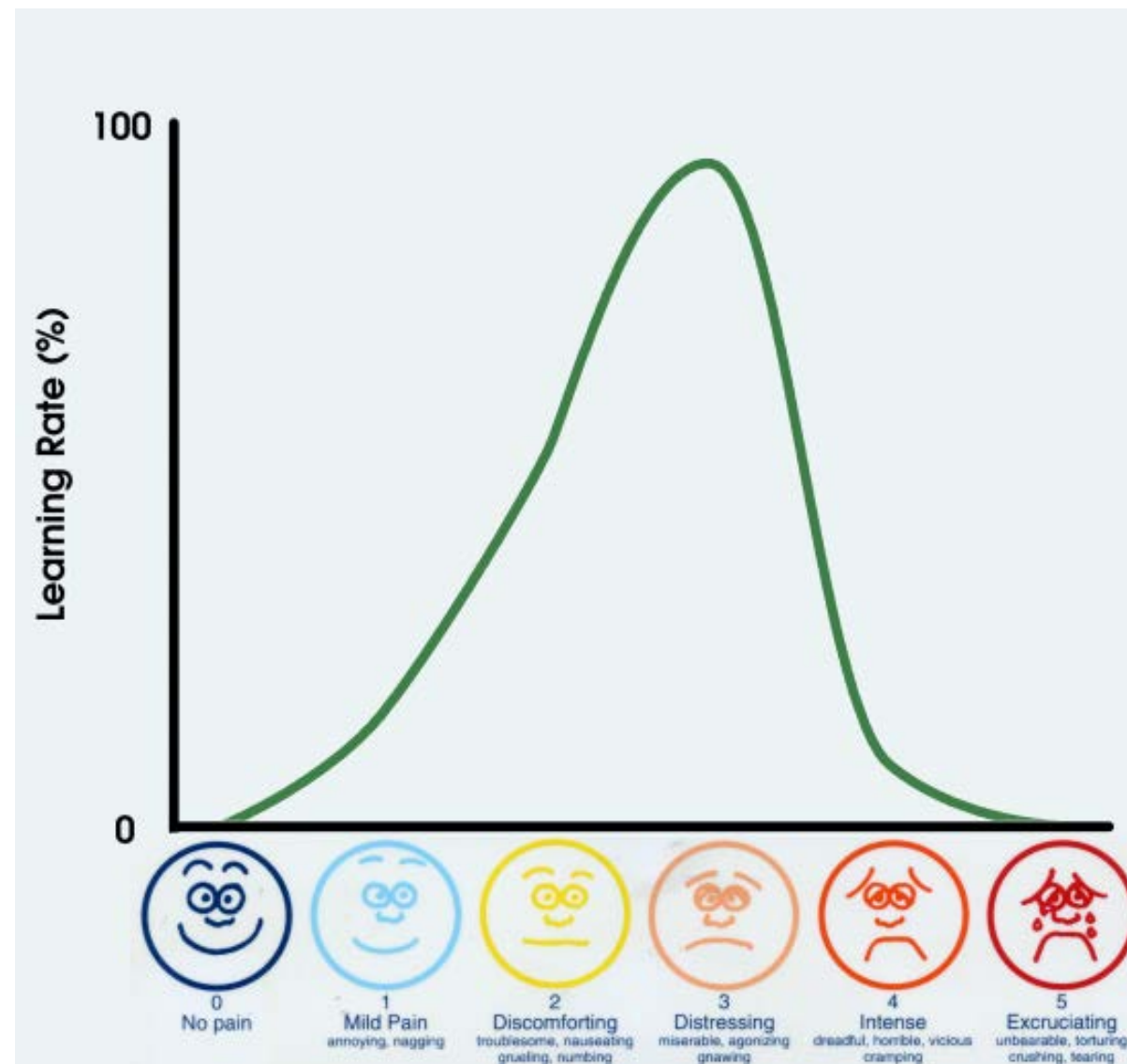
  where you can find slides, errata and some chapters

- Chapters 1, 2, 3, 4, 8, 10, 6, 5

# Advice

- Keep track

- Go through the video lecture + accompying slides before the Q&A session

- Start looking at the homework at least 2 weeks before the deadline

- Practice with some of the exercises mentioned in the learning tasks

- If you're stuck, formulate why and ask!

# Theory of Pain and Learning



*Struble, 2004*

# Data Mining: Introduction
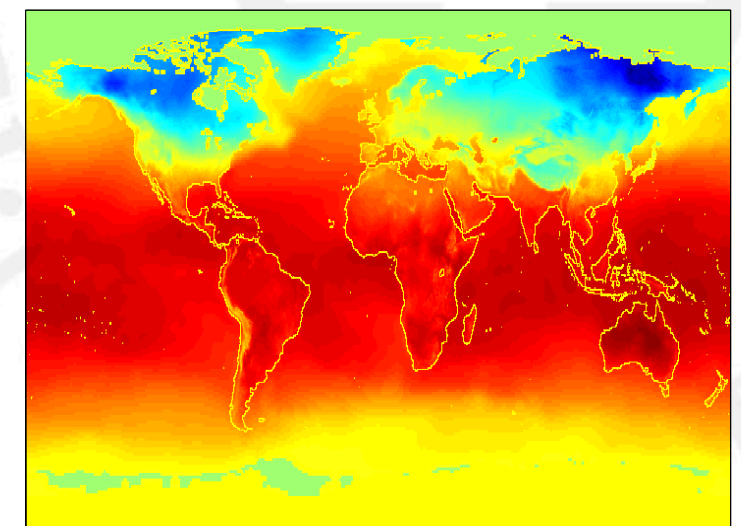
- Motivation

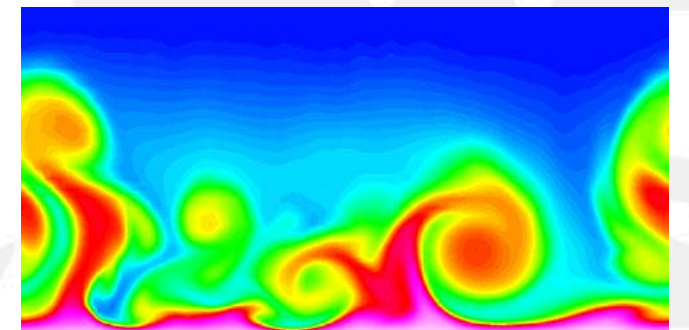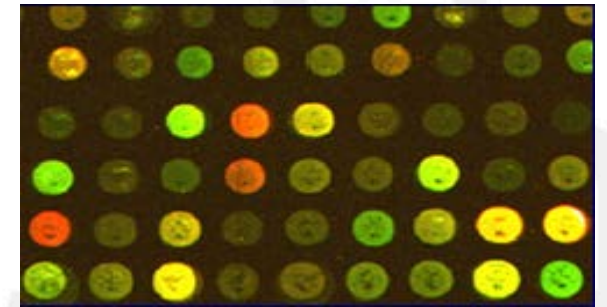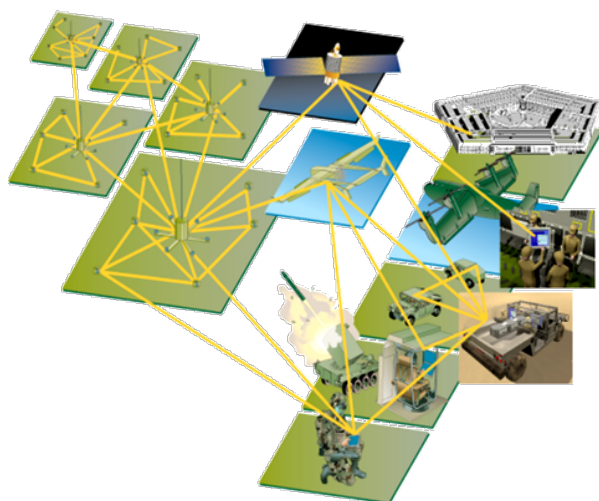- Examples

- Bit of history

- Challenges

# Why Mine Data? Commercial Viewpoint

- Lots of data is being collected and warehoused
  - web data, e-commerce
  - purchases at department/ grocery stores
  - bank/credit card transactions

- Computers have become cheaper and more powerful

- Competitive pressure is strong
  - provide better, customized services for an *edge* (e.g. in customer relationship management)

# Why Mine Data? Scientific Viewpoint
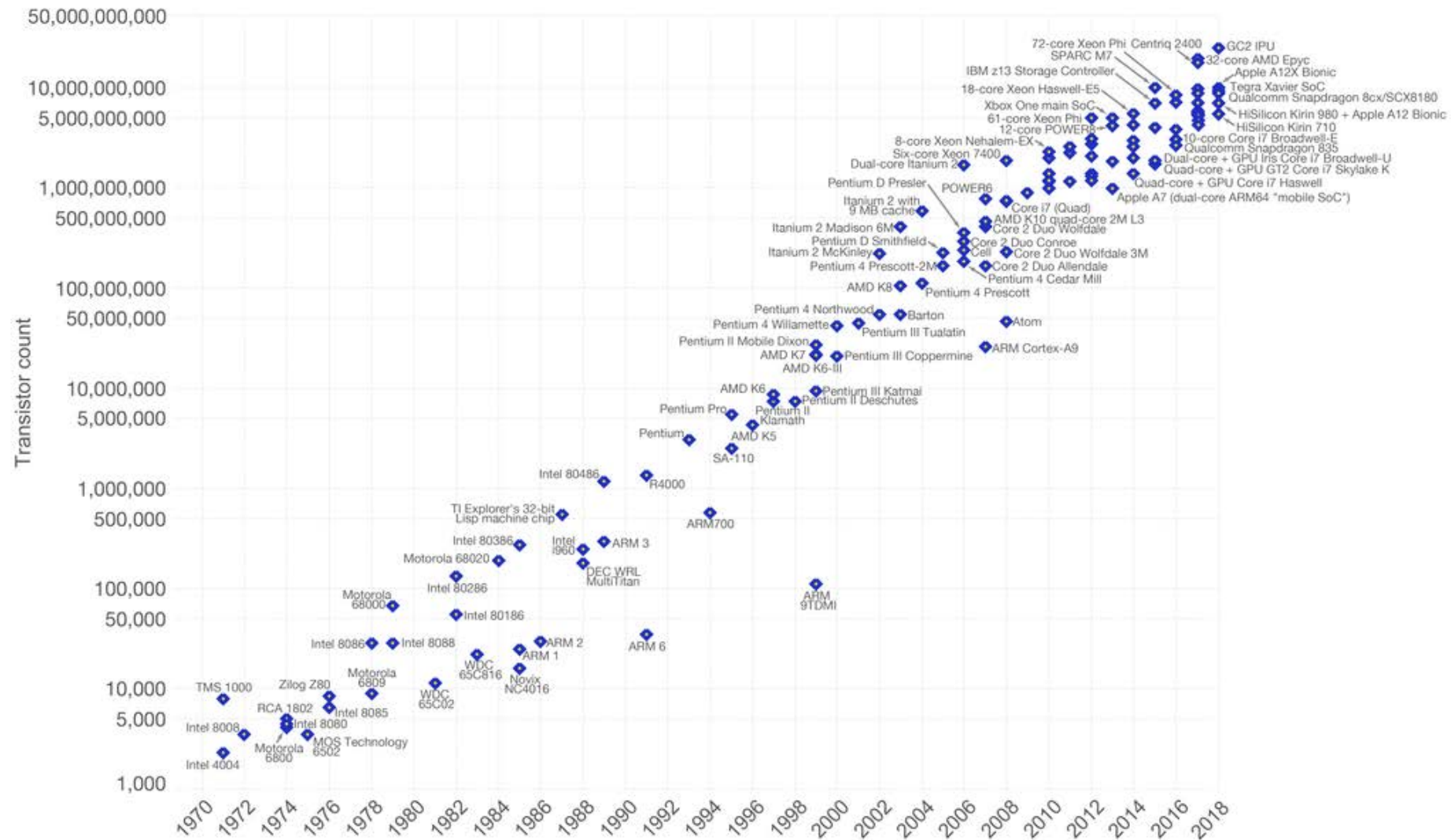
- Data collected and stored at enormous speeds (Gb/hour)
    - remote sensors on a satellite
    - telescopes scanning the skies
    - microarrays generating gene expression data
    - scientific simulations generating terabytes of data

- Traditional techniques infeasible for raw data

- Data mining may help scientists in classifying and segmenting data in hypothesis formation

# Moore's law



Moore's Law – The number of transistors on integrated circuit chips (1971-2018)

Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years. This advancement is important as other aspects of technological progress – such as processing speed or the price of electronic products – are linked to Moore's law.

Radboud University Nijmegen

# Mining Large Data Sets - Motivation

- Current data volume estimated at ~40 Zettabyte ($10^{13}$ GB),  and doubling every 1.5 years to reach 44 ZB in 2020

- There is often information "hidden" in the data that is not readily evident

- Human analysts may take weeks to discover useful information

- Much of the data is never analyzed at all

# Examples of Massive Datasets



SKA Science Archive
searches on Google 98PB
uploads to facebook 180PB
LOFAR Long Term Archive 23PB
YouTube 15PB
CERN 15PB
SKA Phase1 Science Archive 300PB
PER YEAR
• 1 Petabyte

- **Pubmed text database**
  - Records for >30 million published articles

- **Web search engines**
  - 60 billion Web pages indexed
  - 100's of millions of site visitors per day

- **CALTRANS loop sensor data (traffic)**
  - Every 30 seconds, thousands of sensors, 2 Gbytes per day

- **NASA MODIS satellite**
  - Coverage at 250m resolution, 37 bands, whole earth, every day

- **Retail transaction data**
  - Ebay, Amazon, Walmart: >100 million transactions per day
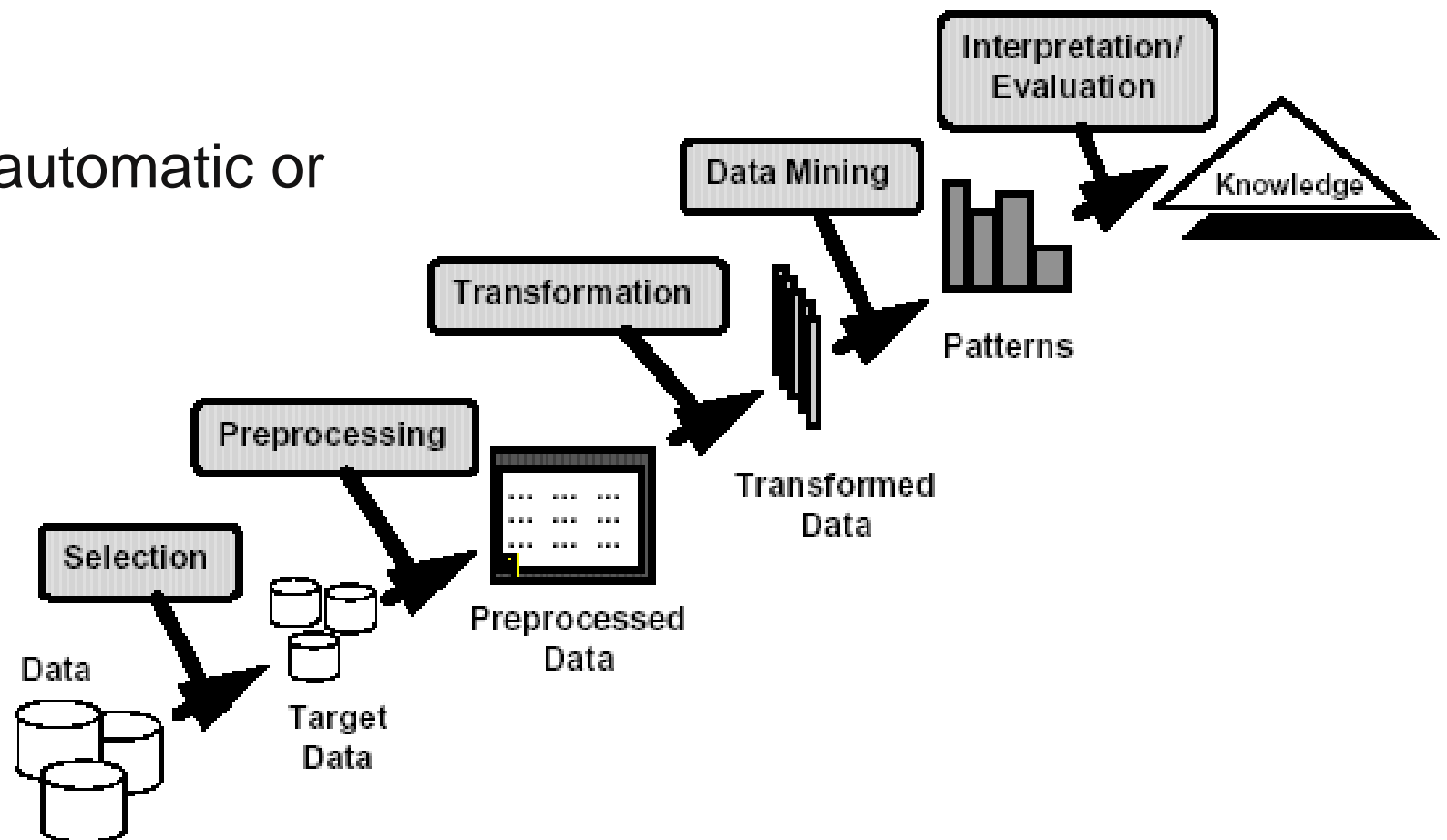  - Visa, Mastercard: similar or larger numbers
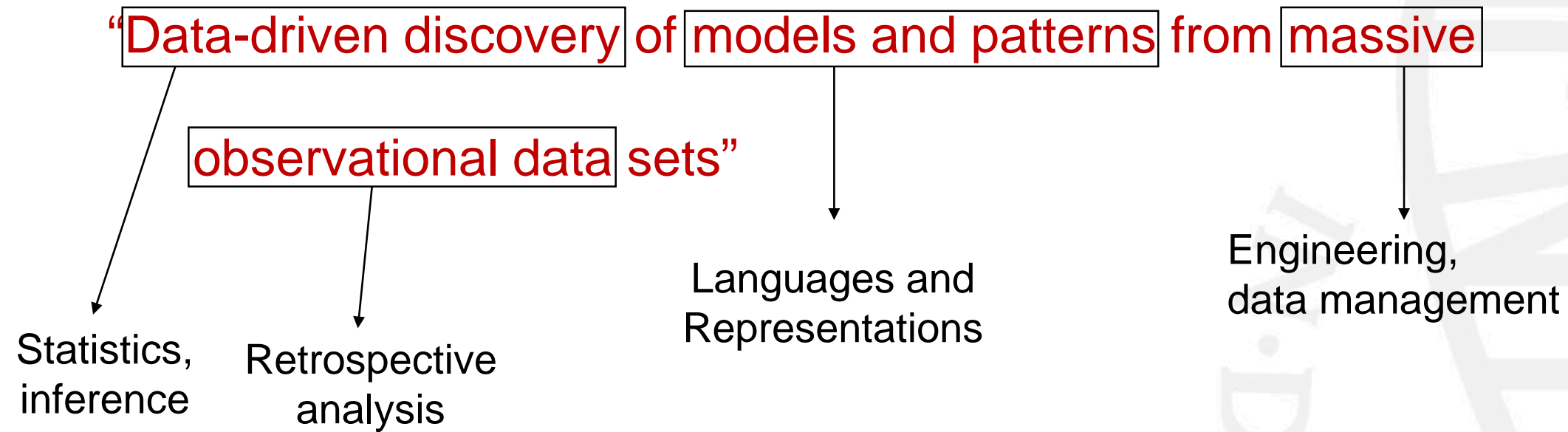
# Harvard Business Review

# What is Data Mining?

Many definitions:

- Non-trivial extraction of implicit, previously unknown and potentially useful information from data

- Exploration & analysis, by automatic or semi-automatic means, of large quantities of data in order to discover meaningful patterns

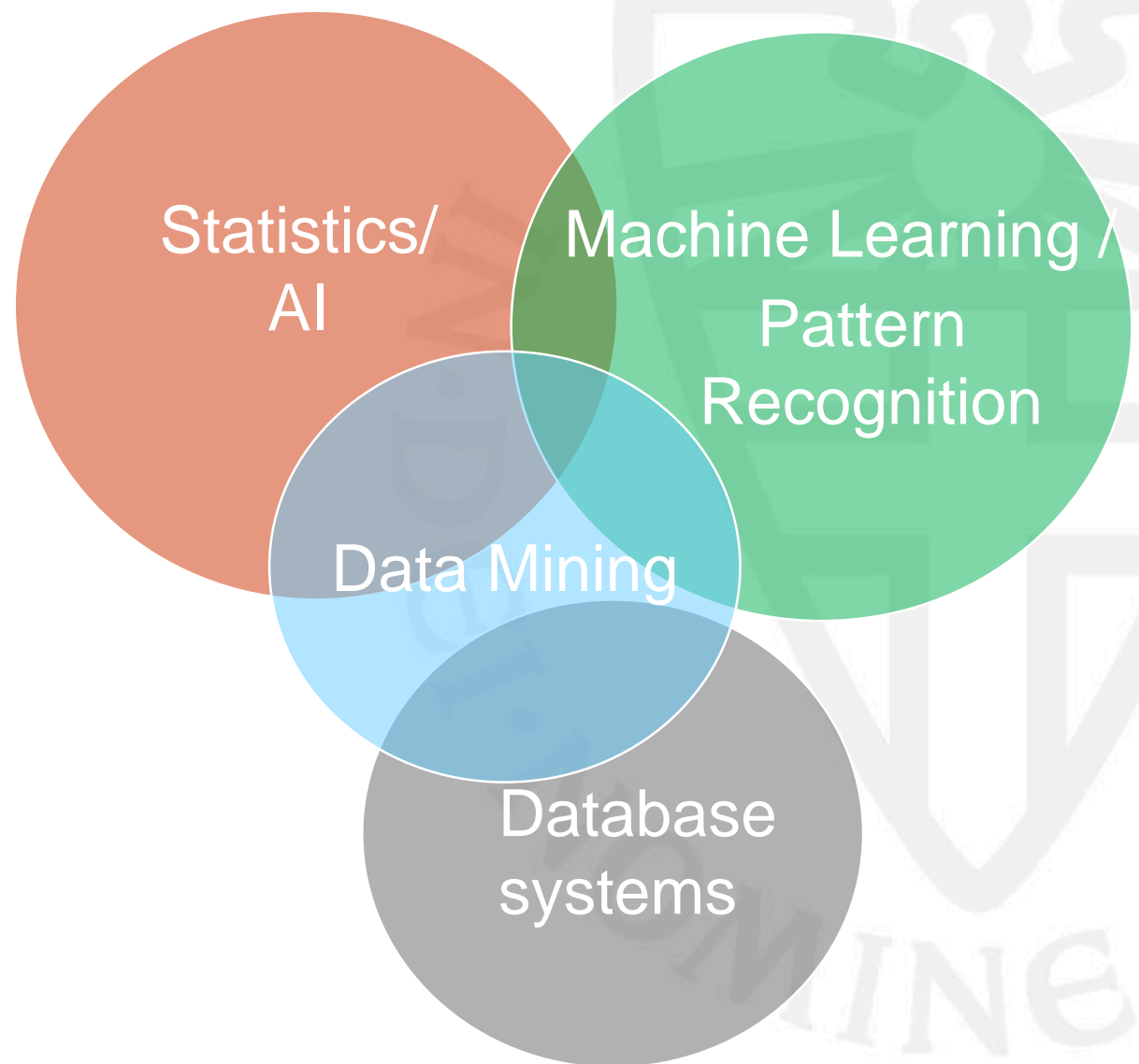# Another Definition

"Data-driven discovery of models and patterns from massive observational data sets"

Statistics, inference

Retrospective analysis

Languages and Representations

Engineering, data management

*Smyth, 2003*

Radboud University Nijmegen

# What is (not) Data Mining?

- What is not Data Mining?

  - Look up phone number in phone directory
  - Query a Web search engine for information about "Amazon"

- What is Data Mining?

  - Certain names are more prevalent in certain US locations (O'Brien, O'Rurke, O'Reilly… in Boston area)

  - Group together similar documents returned by search engine according to their context (e.g. Amazon rainforest vs. Amazon.com)
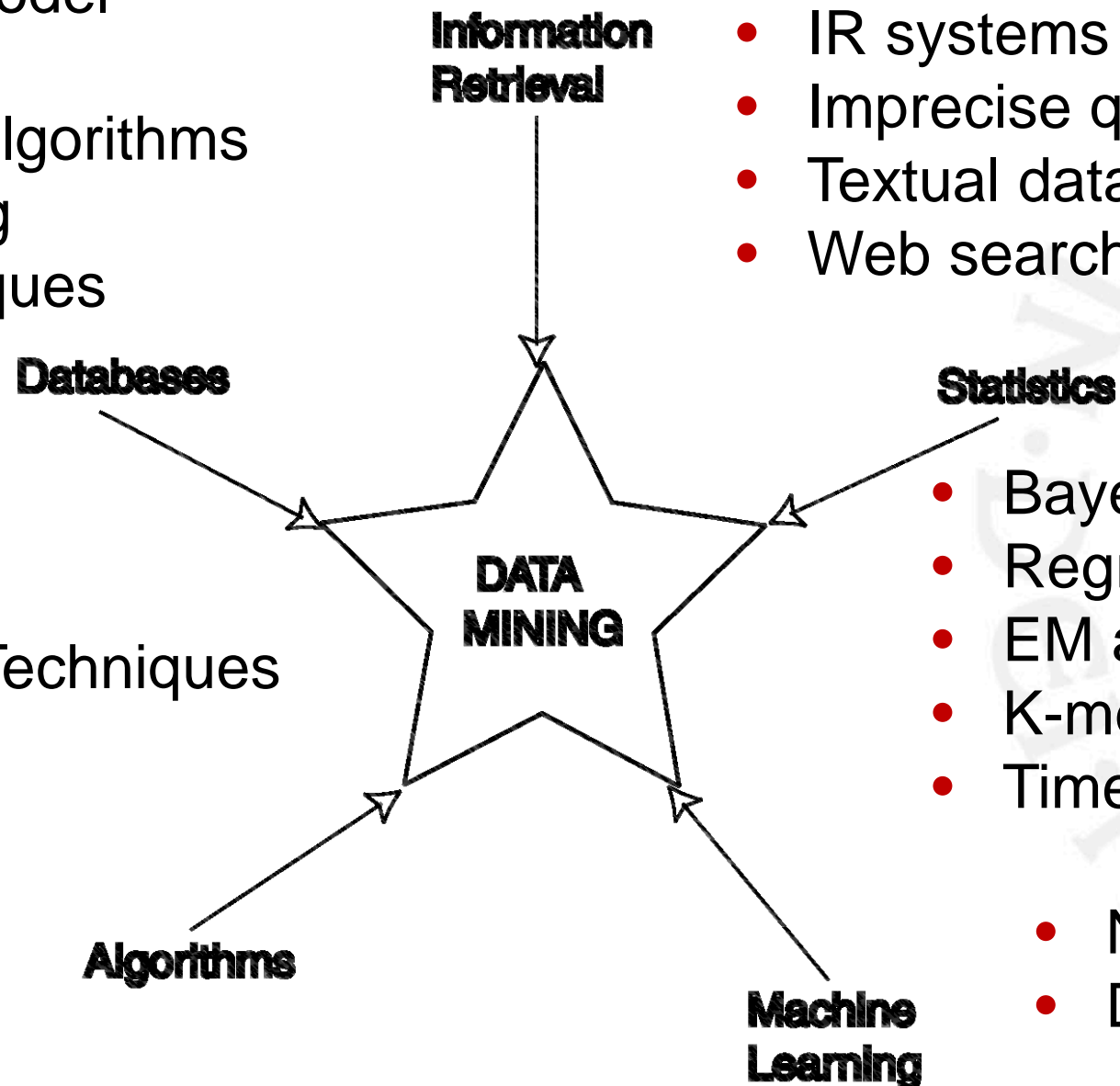
# Origins of Data Mining

- Draws ideas from machine learning/AI, pattern recognition, statistics, and database systems

- Traditional techniques may be unsuitable due to
  - Enormity of data
  - High dimensionality of data
  - Heterogeneous, distributed nature of data



Statistics/AI

Machine Learning / Pattern Recognition

Data Mining

Database systems

# Data Mining Development

- Relational data model
- SQL
- Association rule algorithms
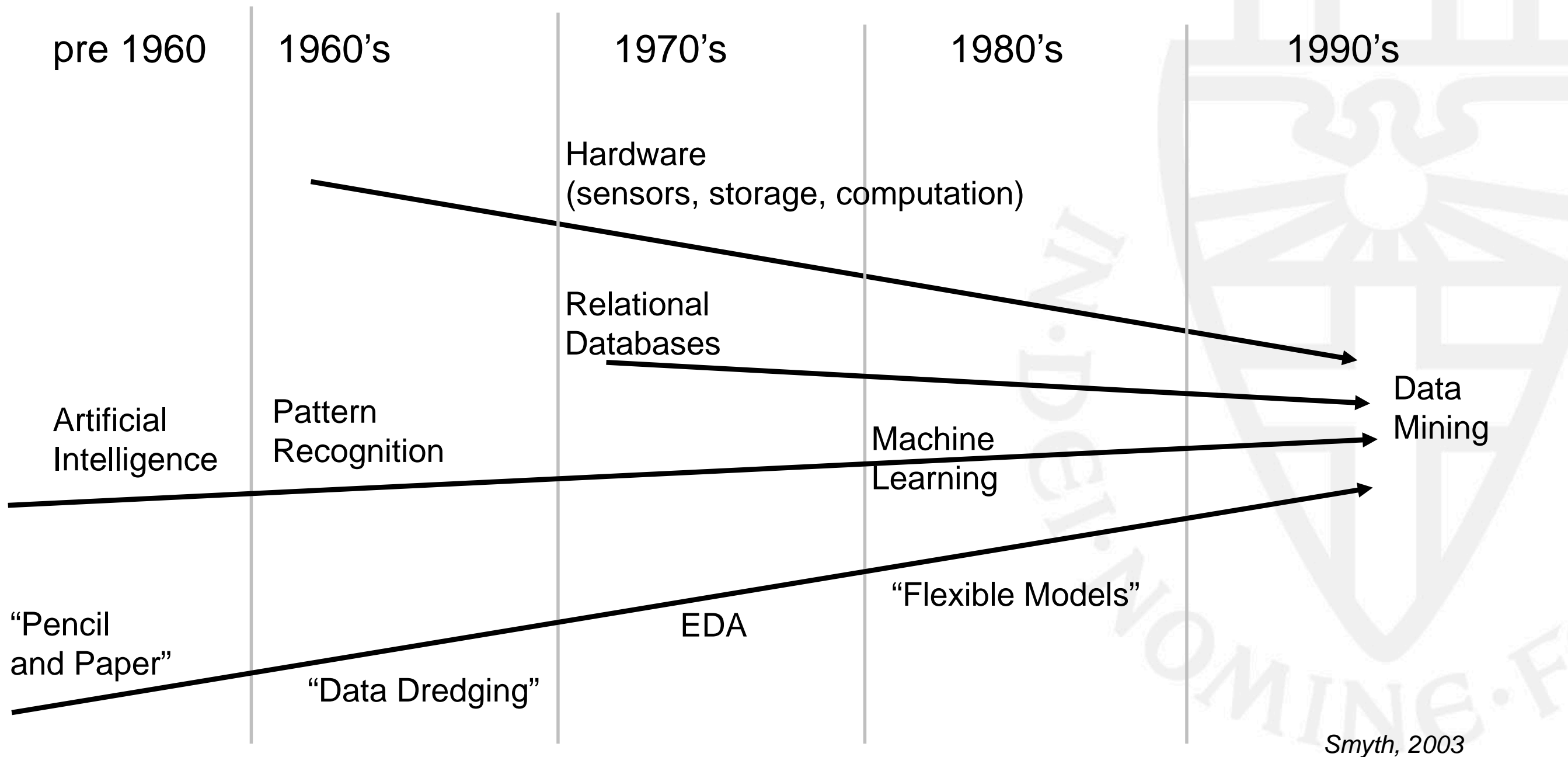- Data warehousing
- Scalability techniques

**Information Retrieval**

- Similarity measures
- Hierarchical clustering
- IR systems
- Imprecise queries
- Textual data
- Web search engines

**Databases**

**Statistics**

**DATA MINING**

- Bayes theorem
- Regression analysis
- EM algorithm
- K-means clustering
- Time series analysis

- Algorithm Design Techniques
- Algorithm Analysis
- Data Structures

**Algorithms**

**Machine Learning**

- Neural networks
- Decision tree algorithms

*Dunham, 2003*

# Origins of Data Mining

| pre 1960 | 1960's | 1970's | 1980's | 1990's |
|----------|--------|--------|--------|--------|

Hardware
(sensors, storage, computation)

Relational
Databases

Data
Mining

Artificial
Intelligence

Pattern
Recognition

Machine
Learning

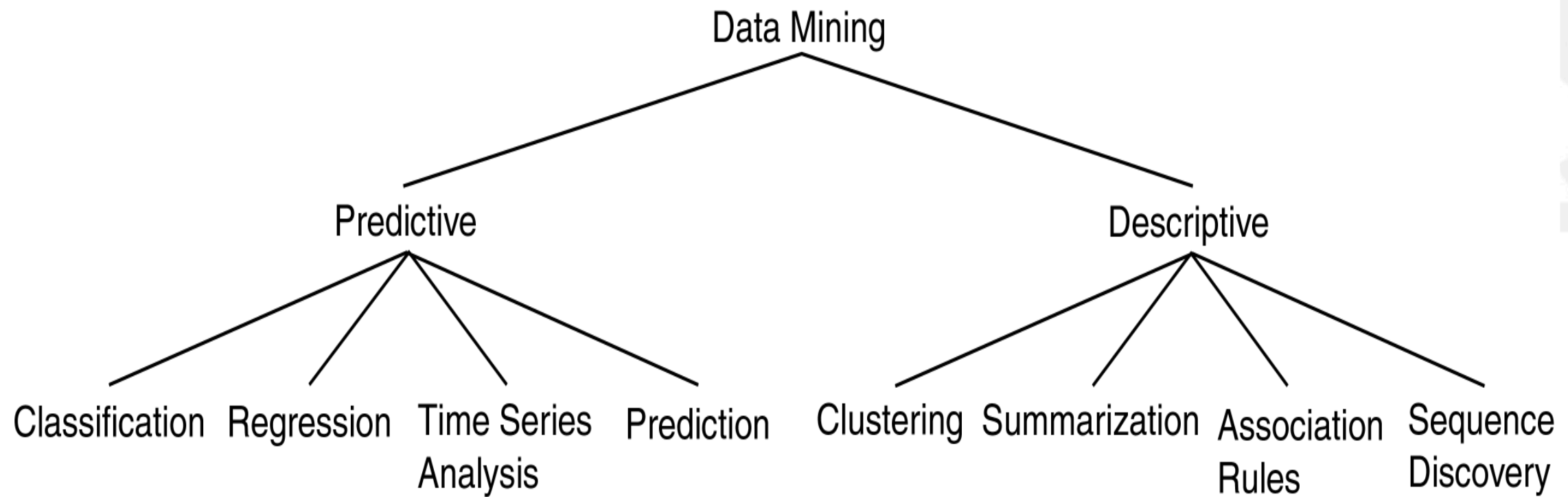"Flexible Models"

"Pencil
and Paper"

EDA

"Data Dredging"

*Smyth, 2003*

# Data Mining Tasks

- Prediction Methods

  - Use some variables to predict unknown or future values of other variables.

- Description Methods

  - Find human-interpretable patterns that describe the data.

# Data Mining Tasks



Dunham, 2003

# Data Mining Tasks...

- Classification [Predictive]

- Clustering [Descriptive]

- Association rule discovery [Descriptive]

- Regression [Predictive]
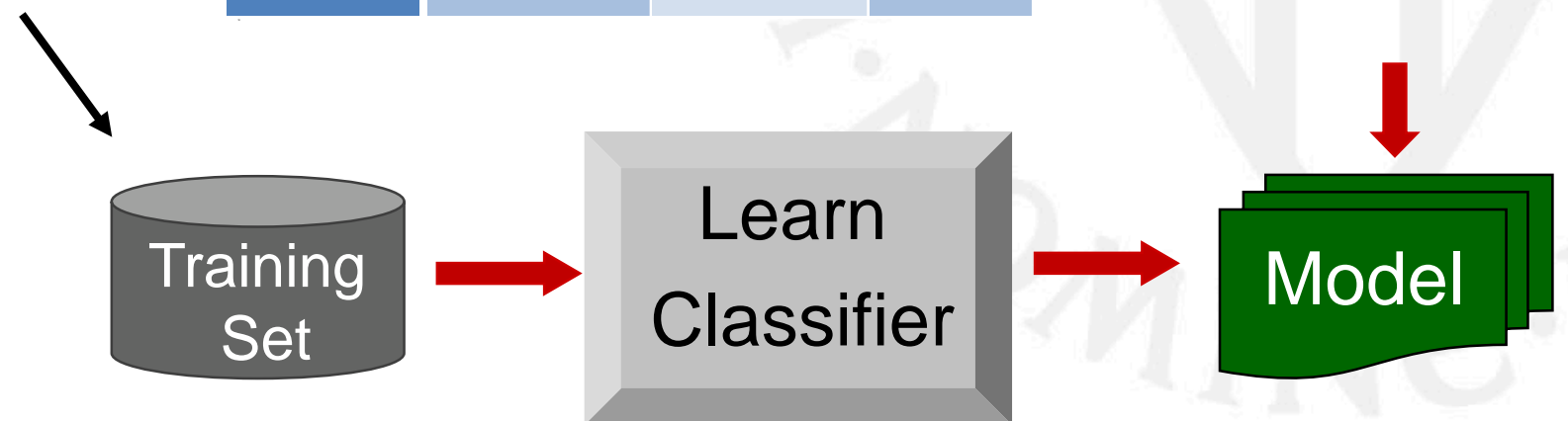
- Deviation detection [Predictive]

# Classification: Definition

- Given a collection of records (training set).

  - Each record contains a set of attributes, one of the attributes is the class.

- Find a model for class attribute as a function of the values of other attributes.

- Goal: previously unseen records should be assigned a class as accurately as possible.

  - A test set is used to determine the accuracy of the model. Usually, the given data set is divided into training and test sets, with training set used to build the model and test set used to validate it.

# Classification Example

| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|---------------|----------------|-------|
| 1 | Yes | Single | 125K | **No** |
| 2 | No | Married | 100K | **No** |
| 3 | No | Single | 70K | **No** |
| 4 | Yes | Married | 120K | **No** |
| 5 | No | Divorced | 95K | **Yes** |
| 6 | No | Married | 60K | **No** |
| 7 | Yes | Divorced | 220K | **No** |
| 8 | No | Single | 85K | **Yes** |
| 9 | No | Married | 75K | **No** |
| 10 | No | Single | 90K | **Yes** |

*categorical*   *categorical*   *continuous*   *class*

| Refund | Marital Status | Taxable Income | Cheat |
|--------|---------------|----------------|-------|
| No | Single | 75K | **?** |
| Yes | Married | 50K | **?** |
| No | Married | 150K | **?** |
| Yes | Divorced | 90K | **?** |
| No | Single | 40K | **?** |
| No | Married | 80K | **?** |

Test Set

Training Set → Learn Classifier → Model

# Classification: Application 1

Direct Marketing

- **Goal**: Reduce cost of mailing by targeting a set of consumers likely to buy a new cell-phone product.

- Approach:
  - Use the data for a similar product introduced before.
  - We know which customers decided to buy and which decided otherwise. This {buy, don't buy} decision forms the class attribute.
  - Collect various demographic, lifestyle, and company-interaction related information about all such customers.
  - Type of business, where they stay, how much they earn, etc.
  - Use this information as input attributes to learn a classifier model.

# Classification: Application 2

Fraud Detection

- **Goal**: Predict fraudulent cases in credit card transactions.

- **Approach**:
  - Use credit card transactions and the information on its account-holder as attributes.
  - When does a customer buy, what does he buy, how often he pays on time, etc
  - Label past transactions as fraud or fair transactions. This forms the class attribute.
  - Learn a model for the class of the transactions.
  - Use this model to detect fraud by observing credit card transactions on an account.

# Classification: Application 3

Customer Attrition/Churn

- **Goal**: To predict whether a customer is likely to be lost to a competitor.

- **Approach**:
  - Use detailed record of transactions with each of the past and present customers, to find attributes.
  - How often the customer calls, where he calls, what time-of-the day he calls most, his financial status, marital status, etc.
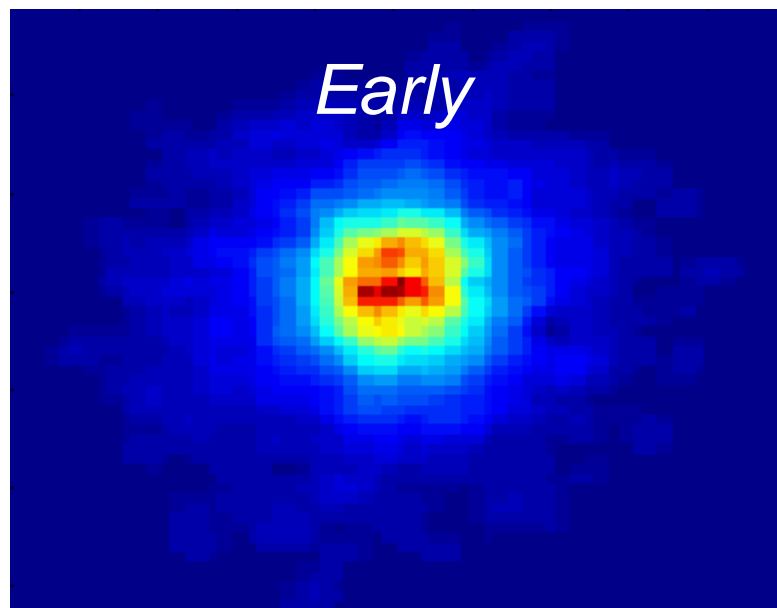  - Label the customers as loyal or disloyal.
  - Find a model for loyalty.

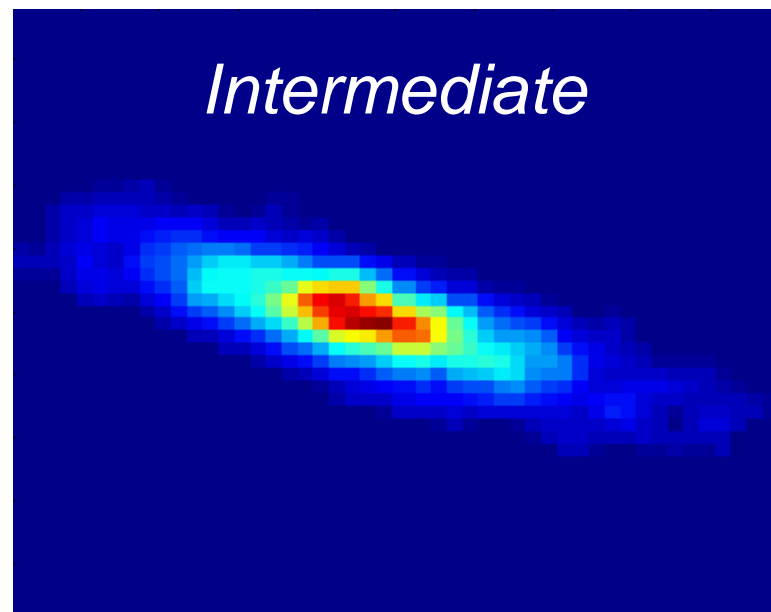# Classification: Application 4

Sky Survey Cataloging

- Goal: To predict class (star or galaxy) of sky objects, especially visually faint ones, based on the telescopic survey images (from Palomar Observatory).
    - 3000 images with 23,040 x 23,040 pixels per image.

- Approach:
    - Segment the image.
    - Measure image attributes (features) - 40 of them per object.
    - Model the class based on these features.
    - Success Story: Could find 16 new high red-shift quasars, some of the farthest objects that are difficult to find!
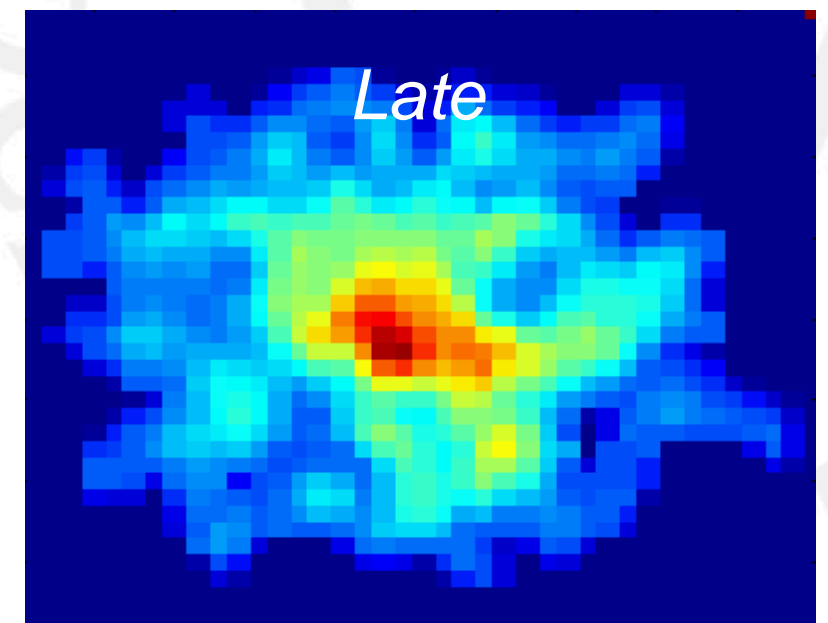
Radboud University Nijmegen

# Classifying Galaxies

*Early*

Class:
- Stages of formation

Attributes:
- Image features,
- Characteristics of light waves received, etc.


*Intermediate*


*Late*

Data size:
- 72 million stars, 20 million galaxies
- Object catalog: 9 GB
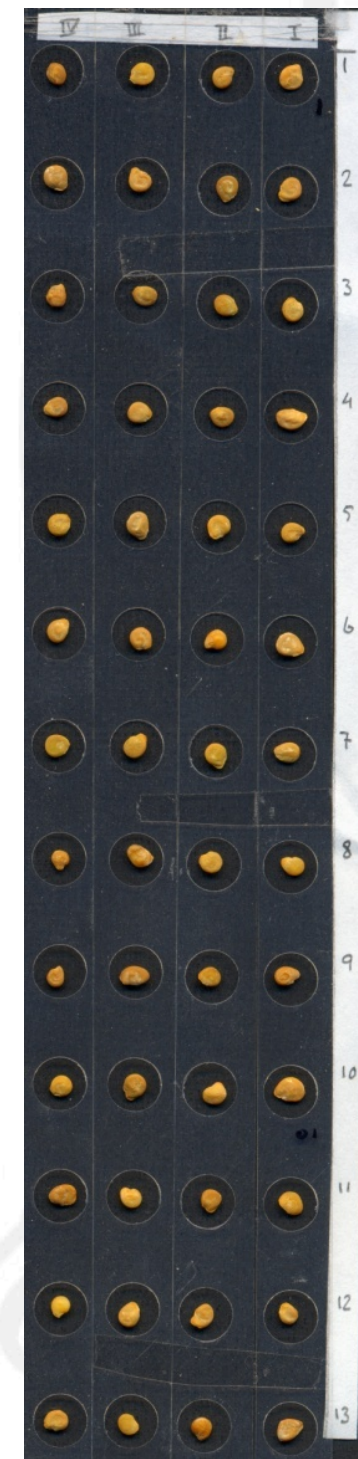- Image database: 150 GB
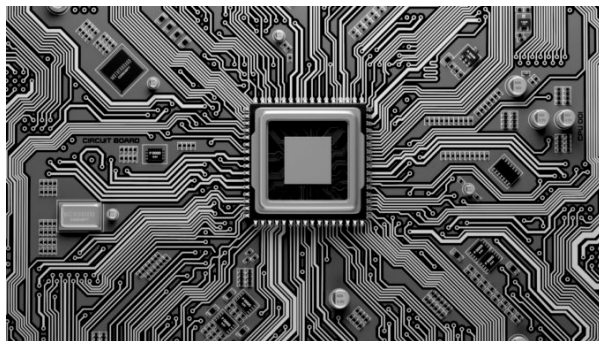
Radboud University Nijmegen

# Classification: Application 5

Classification of tomato seeds

- **Goal**: to predict whether tomato seeds germinate

- **Approach**:

    - "scan" the seeds
    - extract features
    - build a classifier
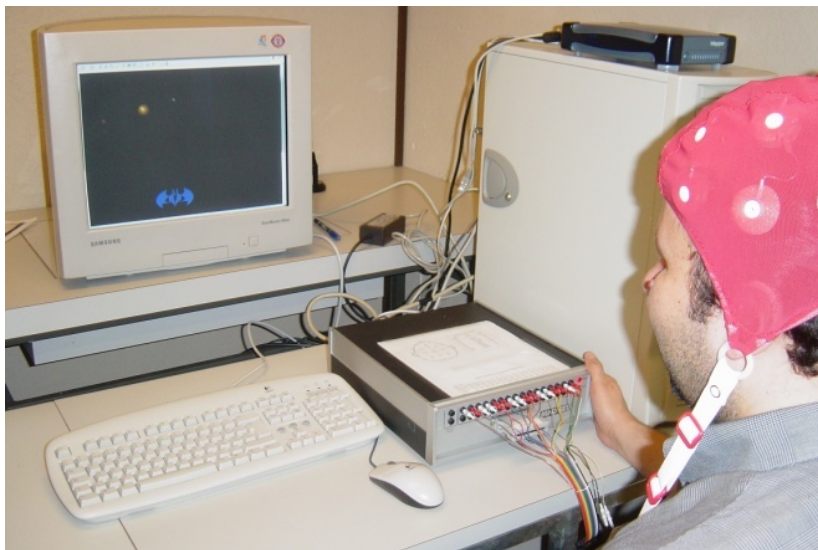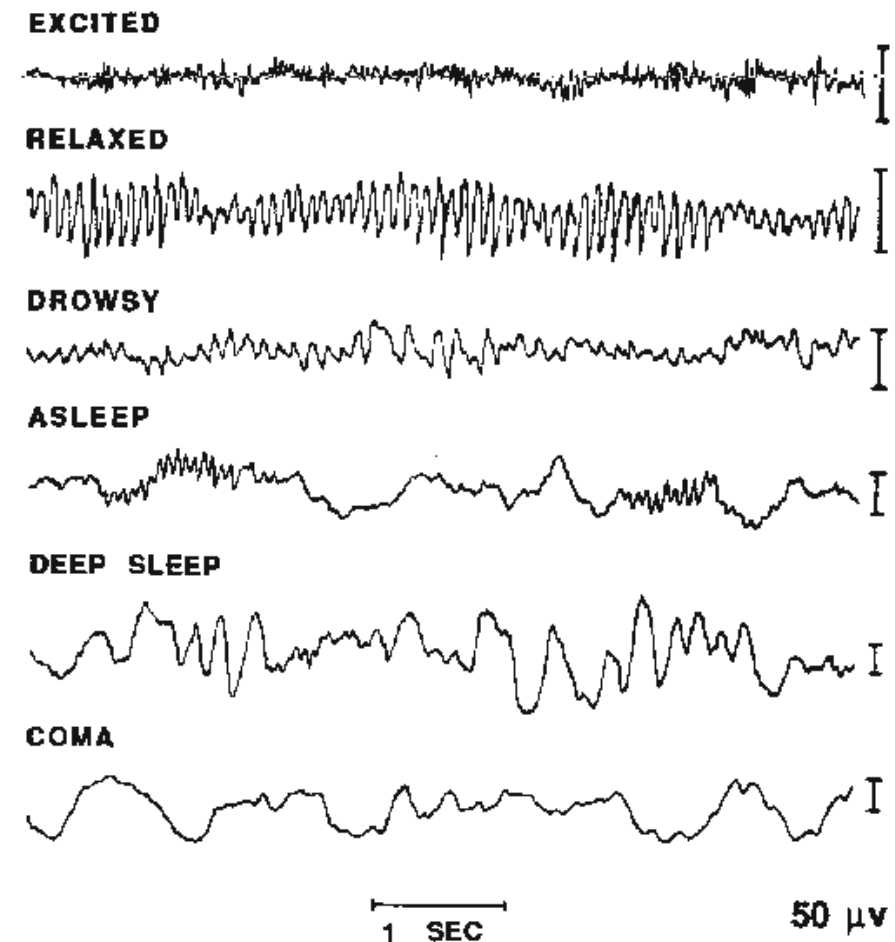    - use the classifier to blow away infertile seeds

# Classification: Example 6

Brain-computer interfacing

- Goal: read a person's mind

- Approach:
  - measure EEG signals
  - classify them



## EEG
### ElectroEncephaloGram

EXCITED

RELAXED

DROWSY

ASLEEP

DEEP SLEEP

COMA

1 SEC          50 µv

Volkskrant



Welke letters las u zonet?
De MRI-scanner weet het

Een team in Nijmegen is er voor het eerst in geslaagd om bij iemand die een woord ziet, te achterhalen welke letters hij heeft gelezen, gegeven welke stukjes hersenschors er oplichten. De crux zit hem in een wiskundig model.
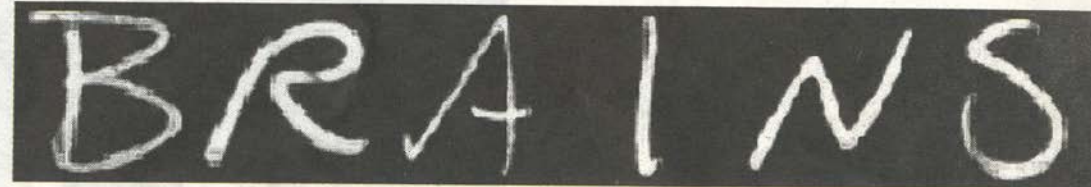
Van onze verslaggever
Bard van de Weijer

AMSTERDAM Derek Ogilvie zal zijn vingers erbij aflikken: onderzoekers van de Radboud Universiteit hebben een methode ontwikkeld waarmee uit iemands hersenactiviteit afgeleid kan worden welke letters hij ziet.
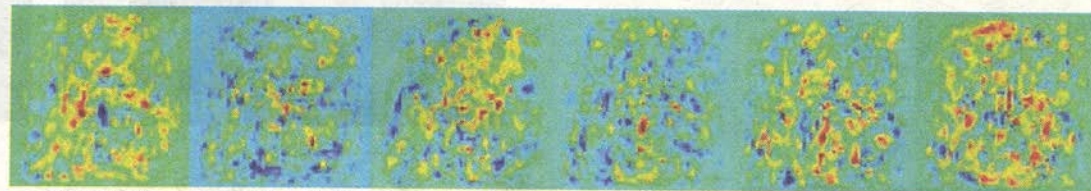
Het aflezen gebeurt met behulp van een MRI-scanner die kijkt naar de visuele cortex, het hersengebied waar beeldinformatie wordt verwerkt. Daartoe worden kubusjes brein van 2×2×2 millimeter in de visuele cortex geanalyseerd. Deze kubusjes, zogenoemde voxels, lichten op als ze worden geactiveerd door visuele informatie.

Als een proefpersoon de letter G ziet, lichten andere voxels op dan bij de letter T. De MRI-scanner meet dus voor elke letter een ander activatiepatroon. Een algoritme kan uit deze patronen de letters reconstrueren die de proefpersoon in de scanner ziet. Het gaat om handgeschreven letters, in allerlei variaties, die alle door het systeem worden herkend.
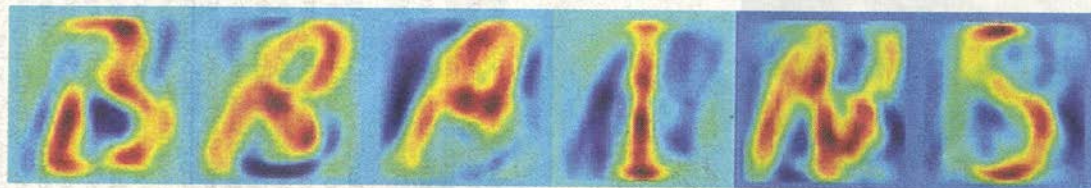
'Het is geen gedachten lezen', zegt cognitief neurowetenschapper Marcel van Gerven van het Donders Instituut van de Radboud Universiteit. 'We reconstrueren perceptie, dus wat iemand ziet, niet wat hij denkt.' Een belangrijk verschil, omdat gedachten

We vermoeden dat het brein ook op deze manier werkt

over 'alles' kunnen gaan en het analyseren van visuele informatie – letters in dit geval – het aantal mogelijkheden beperkt. Het algoritme is getraind op het herkennen van letters. Als een proefpersoon een afbeelding van een vliegtuig wordt voorgehouden, zal dat niet herkend worden.

Tot zover is er volgens Van Gerven nog niet veel nieuws onder de zon. 'We zijn niet de eersten die met MRI-scans beeldpatronen in de visuele cortex kunnen herkennen. Het is wel voor het eerst gelukt om met een wiskundig model het oorspronkelijke beeld met hoge kwaliteit te reconstrueren.'

Dit gebeurt door twee bronnen te combineren: de onderzoekers kijken in een gebiedje van duizend voxels hoe deze reageren op externe stimuli. Deze gegevens – de wat gruizige afbeeldingen hierboven – worden gecombineerd met voorkennis over de eigenschappen van letters. Door de data van de MRI-scan te vergelijken met deze 'kennis' kan worden herleid welke letters de proefpersoon waarneemt.
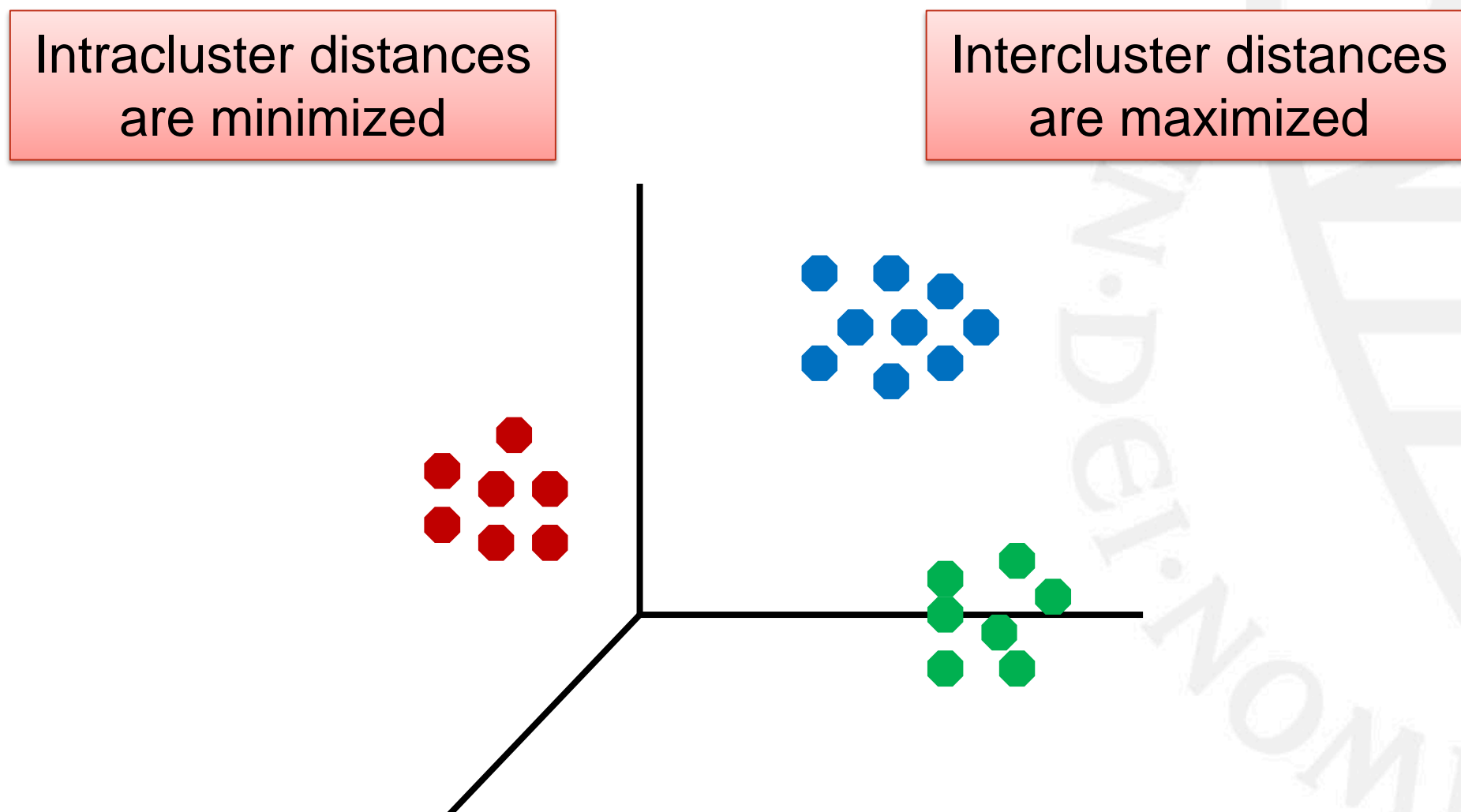
'We vermoeden dat het brein ook op deze manier werkt', zegt Van Gerven. 'Je kunt al die lijntjes en bochtjes niet begrijpen voor je hebt leren lezen. Pas als sprake is van een zekere context kun je letters onderscheiden.' De onderzoekers hopen met hun onderzoek meer te weten te komen over de werking van het brein. Hoewel het bedenken van praktische toepassingen niet het eerste doel is, ziet de onderzoeker wel mogelijkheden. 'Er is een relatie tussen perceptie en verbeelding. Je zou wellicht een reconstructie kunnen maken van een beeld dat iemand zich in gedachten voorstelt. Denk aan een getuige die zich de verdachte inbeeldt en dat je dat beeld dan kunt visualiseren. Maar dat is echt de verre toekomst.'

Een selectie van de oorspronkelijke handgeschreven letters ...

... wat de MRI-scanner ziet oplichten in de visuele cortex ...

... en de reconstructie van de letters door het algoritme.

Illustraties Radboud Universiteit

# Clustering Definition

- Given a set of data points, each having a set of attributes, and a similarity measure among them, find clusters such that

  - Data points in one cluster are more similar to one another.
  - Data points in separate clusters are less similar to one another.

- Similarity measures:

  - Euclidean distance if attributes are continuous.
  - Other problem-specific measures.

# Illustrating Clustering

- Euclidean distance based clustering in 3-D space.

Intracluster distances are minimized

Intercluster distances are maximized

# Clustering: Application 1

Market segmentation

- **Goal**: subdivide a market into distinct subsets of customers where any subset may conceivably be selected as a market target to be reached with a distinct marketing mix.

- **Approach**:
    - Collect different attributes of customers based on their geographical and lifestyle related information.
    - Find clusters of similar customers.
    - Measure the clustering quality by observing buying patterns of customers in same cluster vs. those from different clusters.

# Clustering: Application 2

Document clustering

- **Goal**: To find groups of documents that are similar to each other based on the important terms appearing in them.

- **Approach**: To identify frequently occurring terms in each document. Form a similarity measure based on the frequencies of different terms. Use it to cluster.

- **Gain**: Information retrieval can utilize the clusters to relate a new document or search term to clustered documents.

# Illustrating Document Clustering

- Data points: 3204 articles of Los Angeles Times.

- Similarity measure: How many words are common in these documents (after some word filtering).

| Category | Total Articles | Correctly Placed |
|---|---|---|
| Financial | 555 | 364 |
| Foreign | 341 | 260 |
| National | 273 | 36 |
| Metro | 943 | 746 |
| Sports | 738 | 573 |
| Entertainment | 354 | 278 |

# Clustering of S&P 500 Stock Data

- Observe stock movements on a daily basis.
- Data points: time series of stock-{up/down}
- Similarity measure: Two points are more similar if the events described by them frequently happen together on the same day.

| | Discovered Clusters | Industry Group |
|---|---|---|
| 1 | Applied-Matl-DOWN,Bay-Network-Down,3-COM-DOWN,Cabletron-Sys-DOWN,CISCO-DOWN,HP-DOWN,DSC-Comm-DOWN,INTEL-DOWN,LSI-Logic-DOWN,Micron-Tech-DOWN,Texas-Inst-Down,Tellabs-Inc-Down,Natl-Semiconduct-DOWN,Oracl-DOWN,SGI-DOWN,Sun-DOWN | Technology1-DOWN |
| 2 | Apple-Comp-DOWN,Autodesk-DOWN,DEC-DOWN,ADV-Micro-Device-DOWN,Andrew-Corp-DOWN,Computer-Assoc-DOWN,Circuit-City-DOWN,Compaq-DOWN, EMC-Corp-DOWN, Gen-Inst-DOWN,Motorola-DOWN,Microsoft-DOWN,Scientific-Atl-DOWN | Technology2-DOWN |
| 3 | Fannie-Mae-DOWN,Fed-Home-Loan-DOWN,MBNA-Corp-DOWN,Morgan-Stanley-DOWN | Financial-DOWN |
| 4 | Baker-Hughes-UP,Dresser-Inds-UP,Halliburton-HLD-UP,Louisiana-Land-UP,Phillips-Petro-UP,Unocal-UP,Schlumberger-UP | Oil-UP |

# Association Rule Discovery: Definition

- Given a set of records each of which contain some number of items from a given collection: Produce dependency rules which will predict occurrence of an item based on occurrences of other items.

| TID | Items |
|-----|-------|
| 1 | Bread, Coke, Milk |
| 2 | Beer, Bread |
| 3 | Beer, Coke, Diaper, Milk |
| 4 | Beer, Bread, Diaper, Milk |
| 5 | Coke, Diaper, Milk |

Rules Discovered:
    {Milk} --> {Coke}
    {Diaper, Milk} --> {Beer}

# Association Rule Discovery: Application 1

Marketing and sales promotion

- Suppose the discovered rule is
                {Bagels, … } → {Potato Chips}

- Potato Chips as consequent: Can be used to determine what should be done to boost its sales.

- Bagels in the antecedent: Can be used to see which products would be affected if the store discontinues selling bagels.

- Bagels in antecedent *and* Potato chips in consequent: Can be used to see what products should be sold with Bagels to promote sale of Potato chips!

# Association Rule Discovery: Application 2

Supermarket shelf management

- Goal: To identify items that are bought together by sufficiently many customers.

- Approach: Process the point-of-sale data collected with barcode scanners to find dependencies among items.

- A classic rule --
  - If a customer buys diaper and milk, then he is very likely to buy beer (on Thursday)

# Association Rule Discovery: Application 3

Market basket analysis at Schuitema (now Jumbo)

- Goal: find and visualize clusters of products that are "similar", i.e., are typically bought together with the same products

- Approach: self-organizing map using specific similarity measure based on co-occurrence

# Self-organizing Map
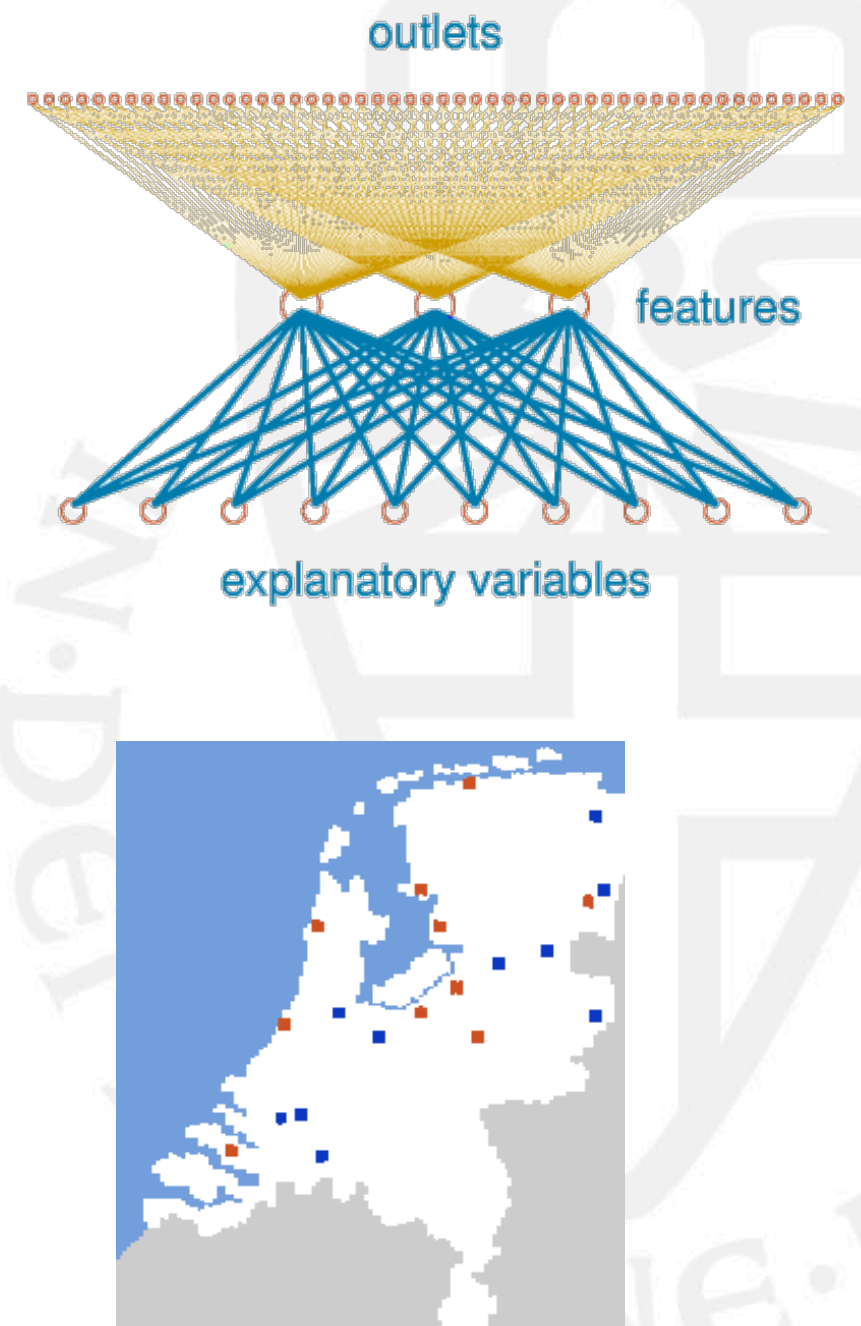
# Regression

- Predict a value of a given continuous valued variable based on the values of other variables, assuming a linear or nonlinear model of dependency.

- Greatly studied in statistics, neural networks.

- Examples:

    - Predicting sales amounts of new product based on advertising expenditure.
    - Predicting wind velocities as a function of temperature, humidity, air pressure, etc.
    - Time series prediction of stock market indices.
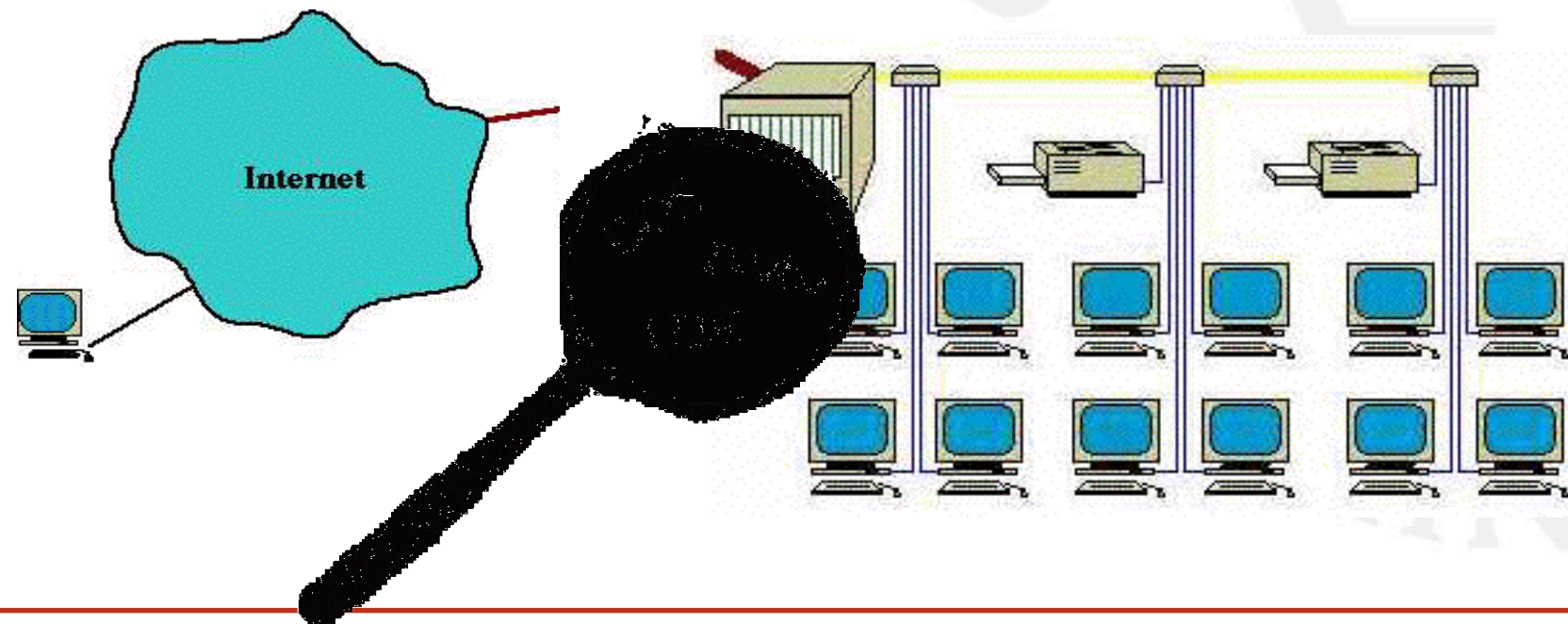
# Regression: Application 1

Predicting newspaper sales

- **Goal**: optimize single-copy sales of De Telegraaf

- **Approach**:
  - learn from past sales
  - let outlets learn from each other



outlets

features

explanatory variables

- better weather, more sales
- worse weather, more sales
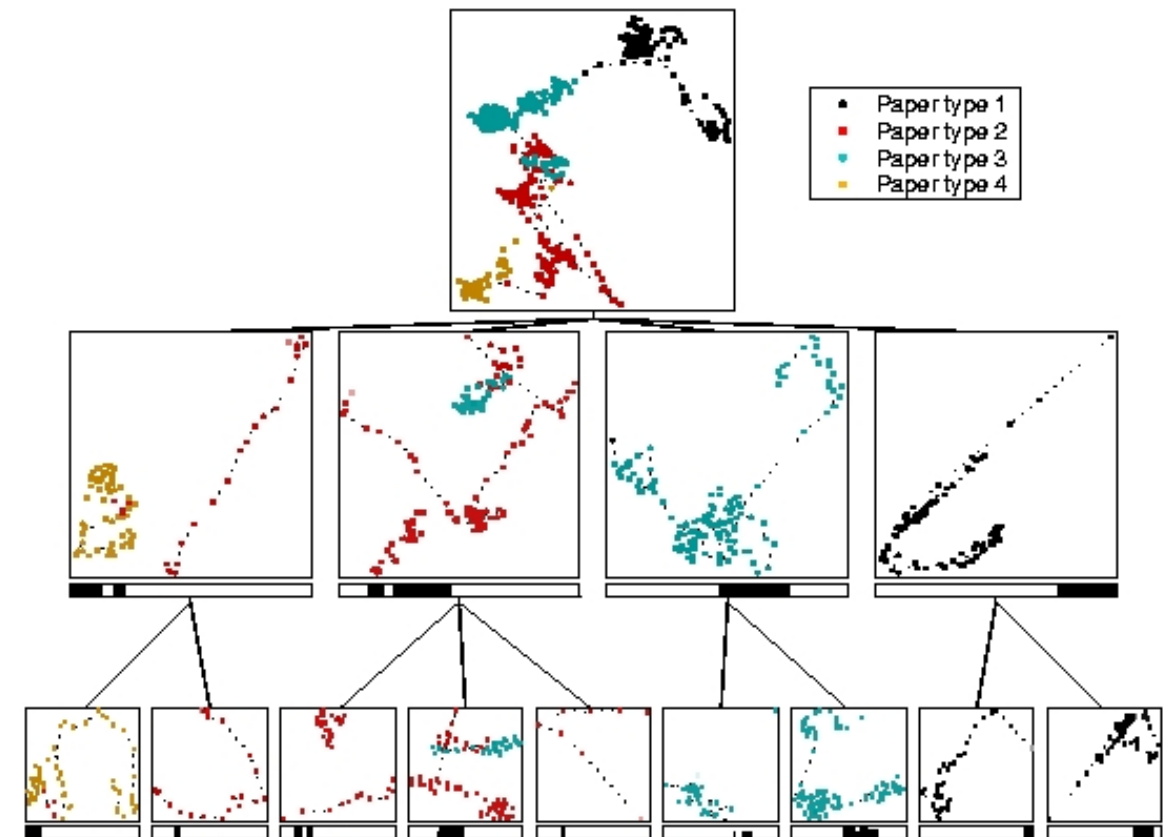
# Deviation / Anomaly Detection

- Detect significant deviations from normal behavior

- Applications:
  - Credit card fraud detection
  - Network intrusion detection

# Deviation / Anomaly Detection: Application 1

Monitoring paper mills

- Goal: alert operators when the paper mill starts behaving "weirdly"

- Approach: visualize the dynamics by cleverly projecting the measurements of hundreds of sensors



Paper type 1
Paper type 2
Paper type 3
Paper type 4

# Data mining?

- Dividing the customers of a company according to their gender.

- Predicting the profitability of customers.

- Computing the total sales of a company.

- Sorting a student database based on student identification numbers.

- Predicting the outcomes of tossing a fair pair of dice.

- Predicting the outcomes of tossing a possibly unfair pair of dice after having seen some amount of tosses.
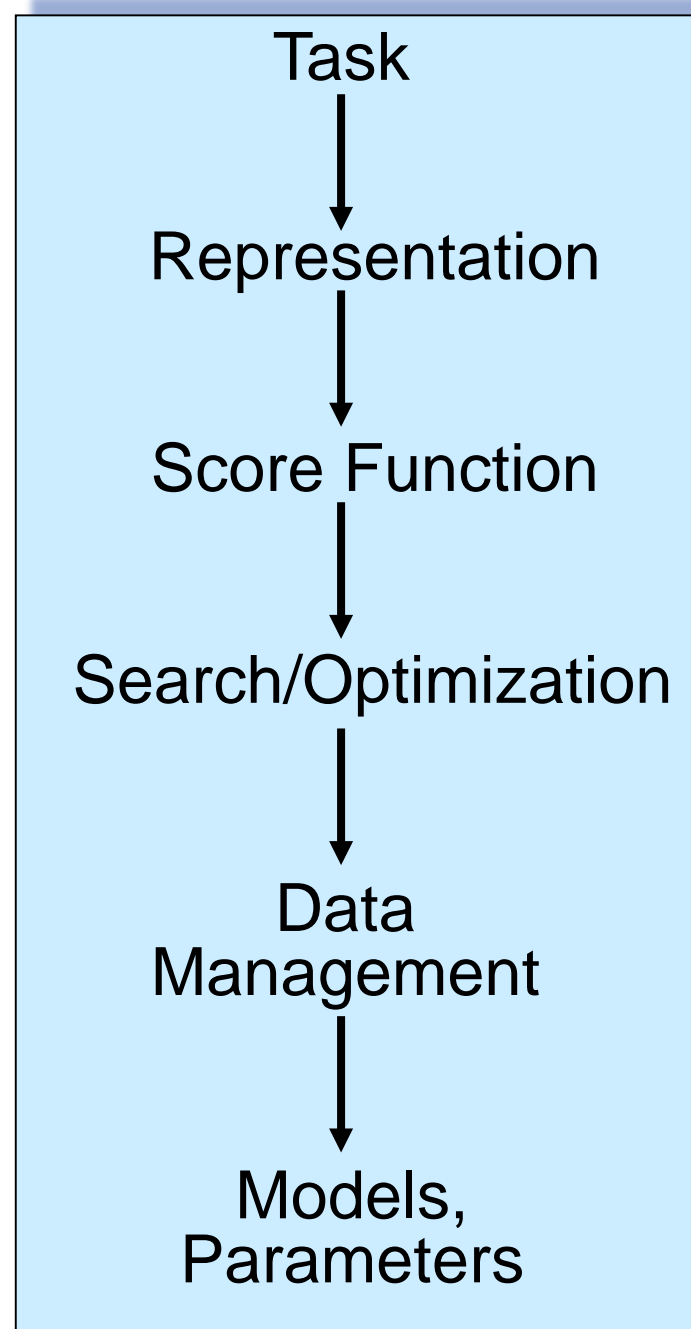
# Data mining?

- Predicting the future stock price of a company using historical records.

- Monitoring the heart rate of a patient for abnormalities given observations of both abnormal and normal behavior.

- Monitoring the heart rate of a patient for abnormalities given observations of only normal behavior.

- Monitoring seismic waves for earthquake activities.

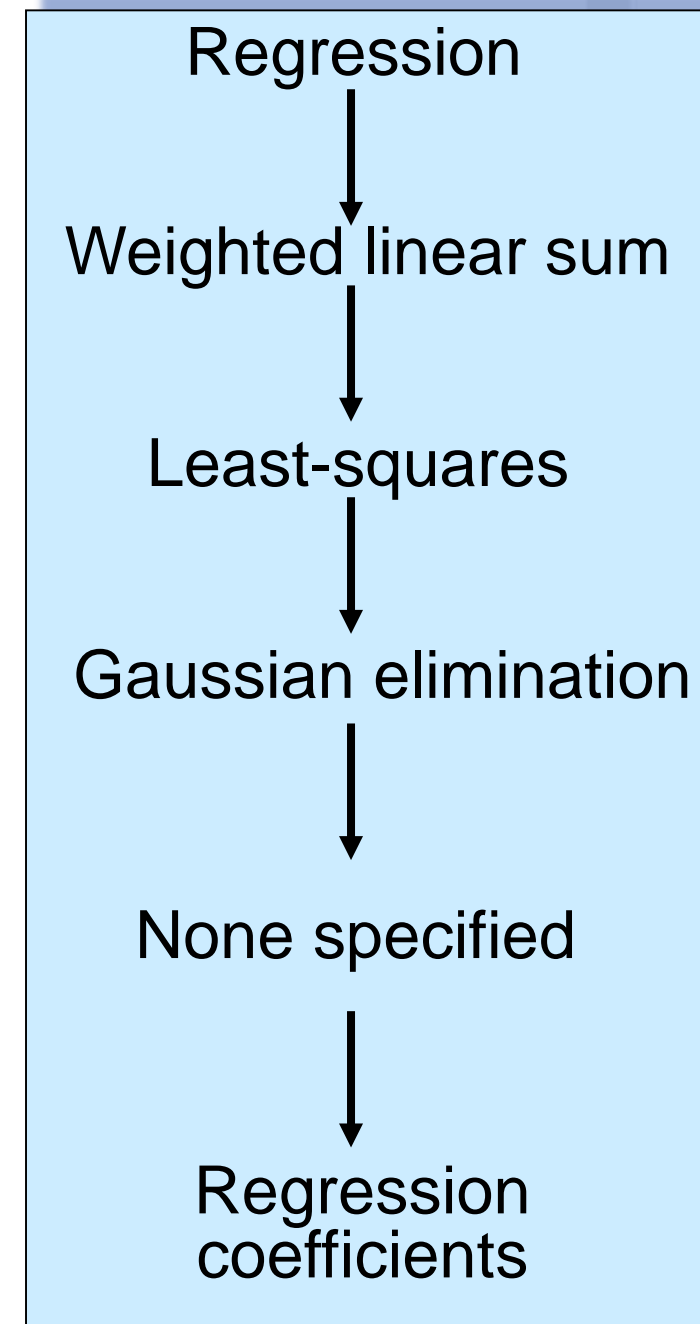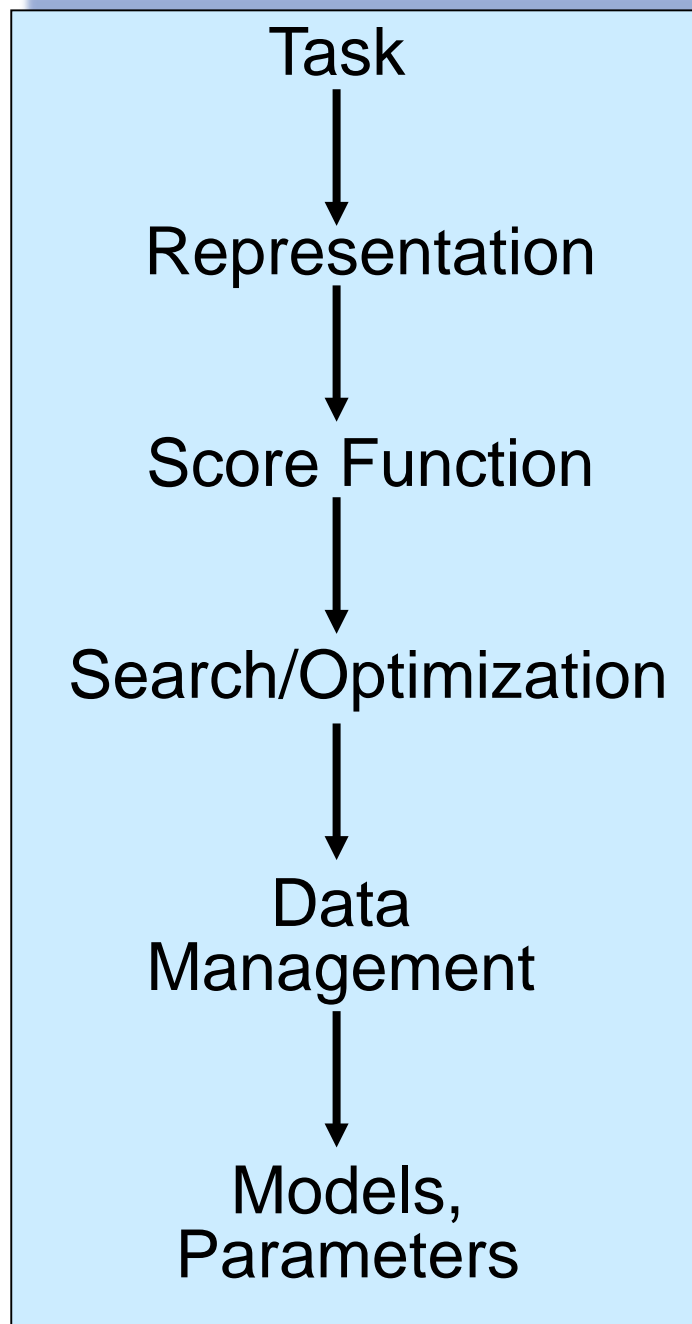- Extracting the frequencies of a sound wave.

# Components of Data Mining Algorithms

- Representation:
  - Determining the nature and structure of the representation to be used

- Score function:
  - quantifying and comparing how well different representations fit the data

- Search/Optimization method:
  - Choosing an algorithmic process to optimize the score function

- Data Management:
  - Deciding what principles of data management are required to implement the algorithms efficiently
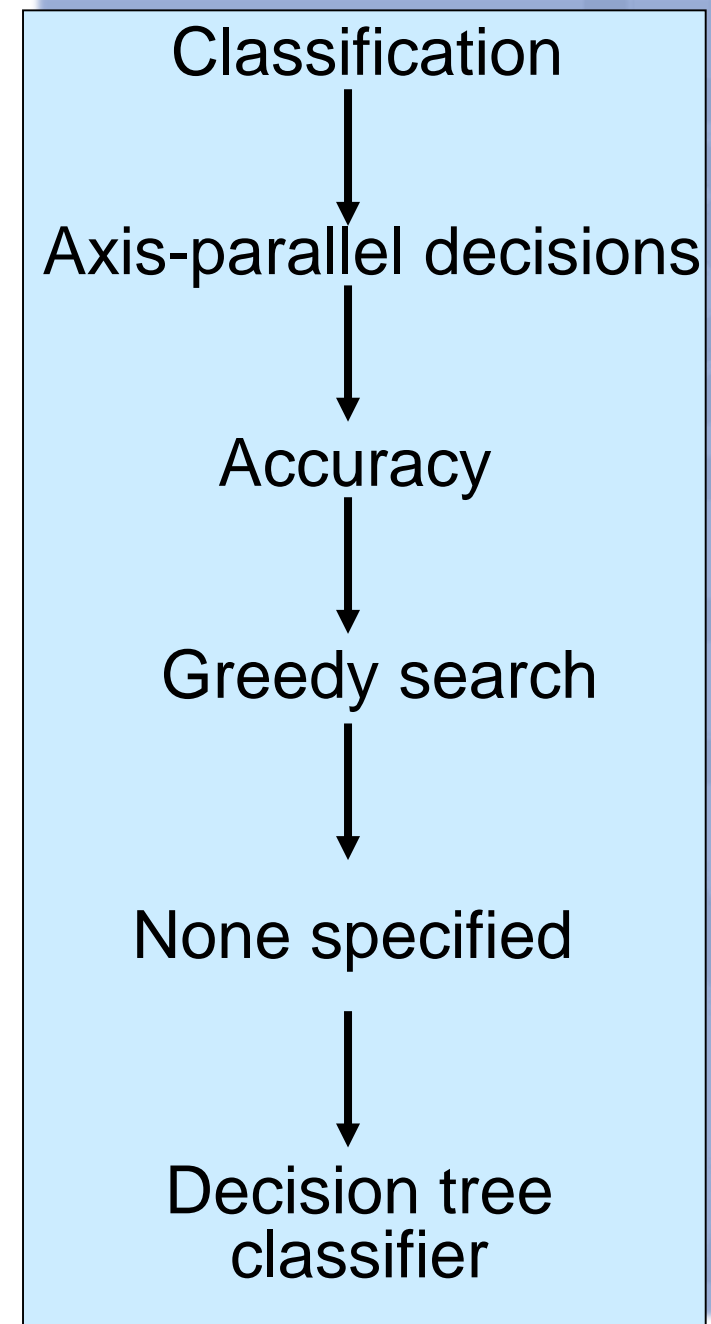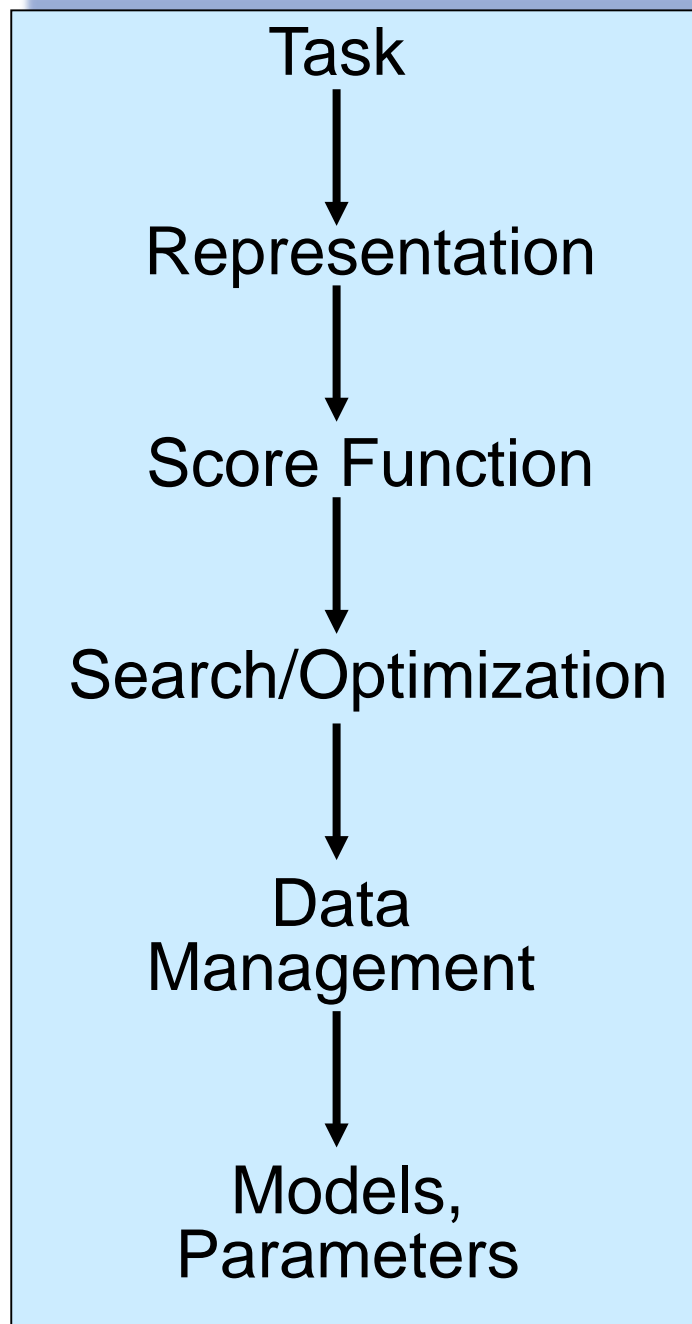
# What's in a Data Mining Algorithm?

Task

↓

Representation

↓

Score Function

↓

Search/Optimization

↓

Data
Management

↓

Models,
Parameters

# Multivariate Linear Regression

Task

↓

Representation

↓

Score Function

↓

Search/Optimization

↓

Data
Management

↓

Models,
Parameters

Regression

↓

Weighted linear sum

↓

Least-squares

↓

Gaussian elimination

↓

None specified

↓

Regression
coefficients

# Decision trees (CART, ID3, …)

| | |
|---|---|
| Task | Classification |
| ↓ | ↓ |
| Representation | Axis-parallel decisions |
| ↓ | ↓ |
| Score Function | Accuracy |
| ↓ | ↓ |
| Search/Optimization | Greedy search |
| ↓ | ↓ |
| Data Management | None specified |
| ↓ | ↓ |
| Models, Parameters | Decision tree classifier |

# Hierarchical Clustering

| | |
|---|---|
| Task | Clustering |
| ↓ | ↓ |
| Representation | Tree of clusters |
| ↓ | ↓ |
| Score Function | Various |
| ↓ | ↓ |
| Search/Optimization | Greedy search |
| ↓ | ↓ |
| Data Management | None specified |
| ↓ | ↓ |
| Models, Parameters | Dendogram |

# Association Rules

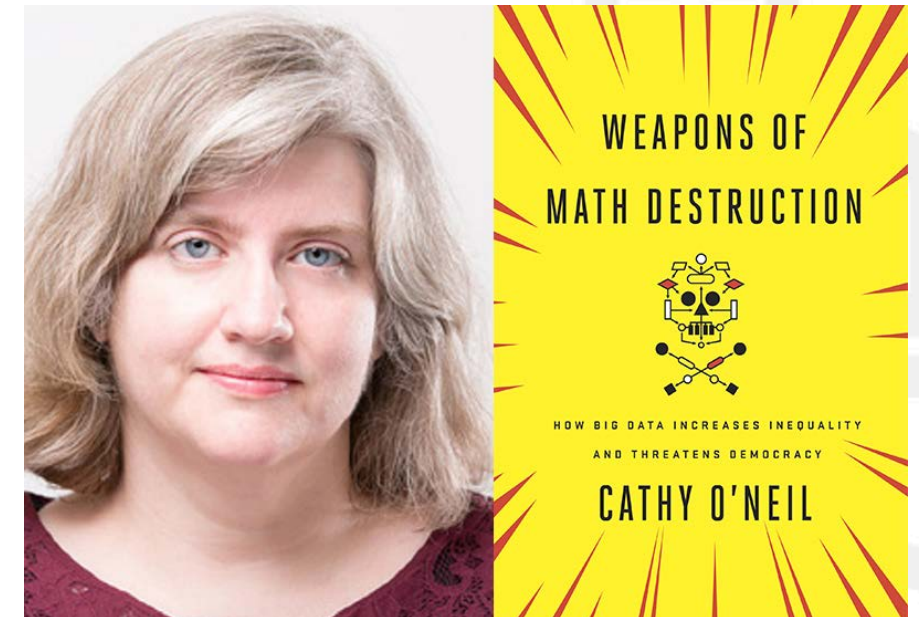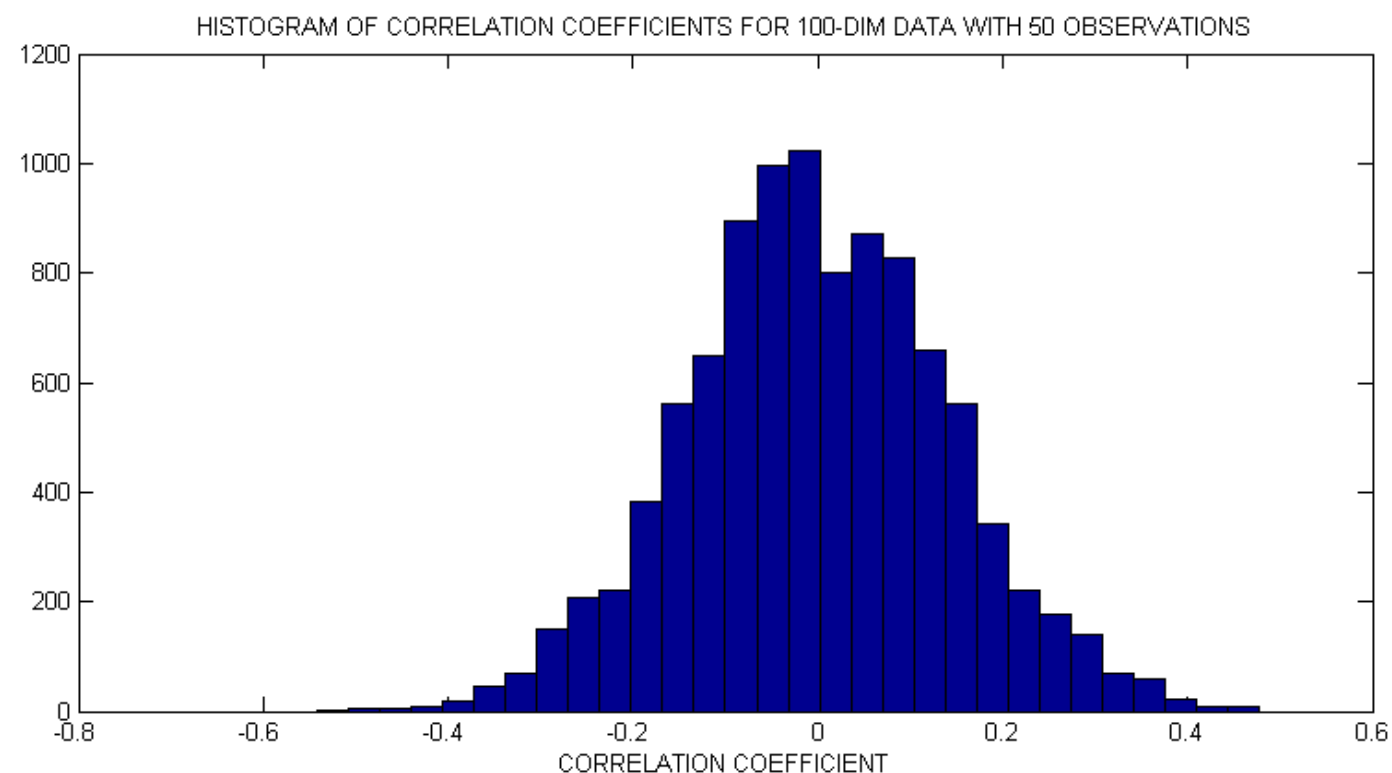| Task | Pattern Discovery |
|:---:|:---:|
| ↓ | ↓ |
| Representation | IF-THEN rules |
| ↓ | ↓ |
| Score Function | No explicit score |
| ↓ | ↓ |
| Search/Optimization | Systematic search |
| ↓ | ↓ |
| Data Management | Multiple linear scans |
| ↓ | ↓ |
| Models, Parameters | Set of rules |

# Data Mining: the Downside

- Hype

- One of the "weapons of math destruction"

- Data dredging, snooping and fishing
  - Finding spurious structure in data that is not real

- Historically, 'data mining' was a derogatory term in the statistics community
  - The Super Bowl fallacy
  - Bangladesh butter prices and the US stock market

- The challenges of being interdisciplinary
  - computer science, statistics, domain discipline

# Example of "Data Fishing"

- Example: data set with
  - 50 data vectors
  - 100 variables
  - Even if data are entirely random (no dependence) there is a very high probability some variables will appear dependent just by chance.



HISTOGRAM OF CORRELATION COEFFICIENTS FOR 100-DIM DATA WITH 50 OBSERVATIONS

# PYTHON Code for Correlations

```python
import numpy as np
import matplotlib.pyplot as plt

nObjects = 50
nVariables = 100
# Generate matrix with standard normal random variables
mu = 0.0
sigma = 1.0
x = np.random.normal(mu, sigma, (nVariables, nObjects))
# Compute correlations between variables
corrvector = np.array([])
for i in range(1, nVariables):
    for j in range(i+1, nVariables):
        # Numpy return the CC's in matrix format
        corMtrx = np.corrcoef(x[:,i], x[:,j])
        corrvector = np.append(corrvector, corMtrx[0,1])

# Plot the histogram
plt.hist(corrvector, 20)
plt.show()
```

# PYTHON Code for Correlations

```python
import numpy as np
import matplotlib.pyplot as plt

nObjects = 50
nVariables = 100

# Generate matrix with standard normal random variables
mu = 0.0
sigma = 1.0
x = np.random.normal(mu, sigma, (nVariables, nObjects))

# Compute correlations between variables
correlations = np.corrcoef(x)
dummy = np.triu(correlations, 1);
corrvector = dummy[dummy != 0];

# Plot the histogram
plt.hist(corrvector, 20)
plt.show()
```
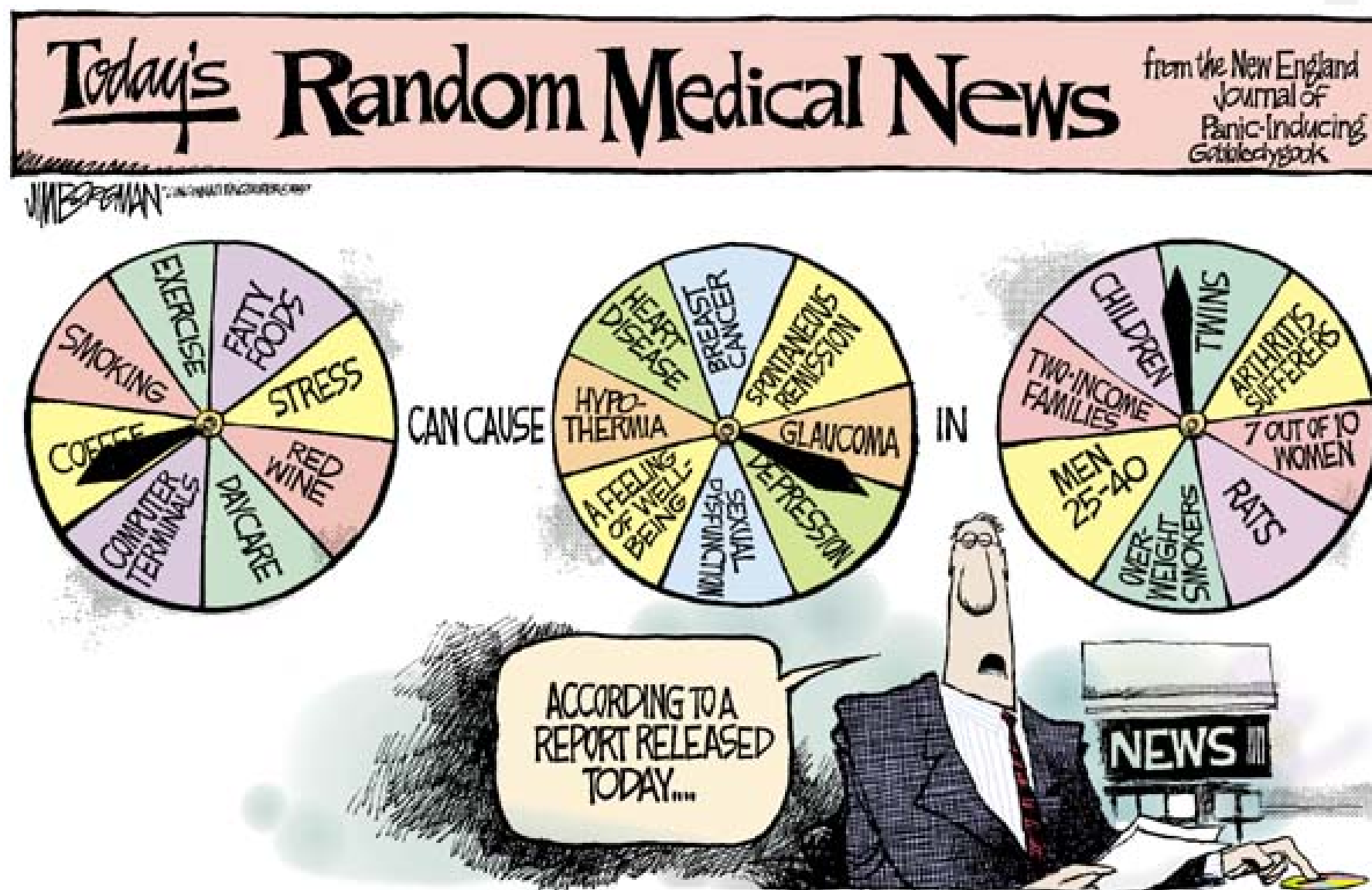
# Possible Pitfalls



Let the data speak…

The data may have quite a lot to say…..
but it may just be noise!

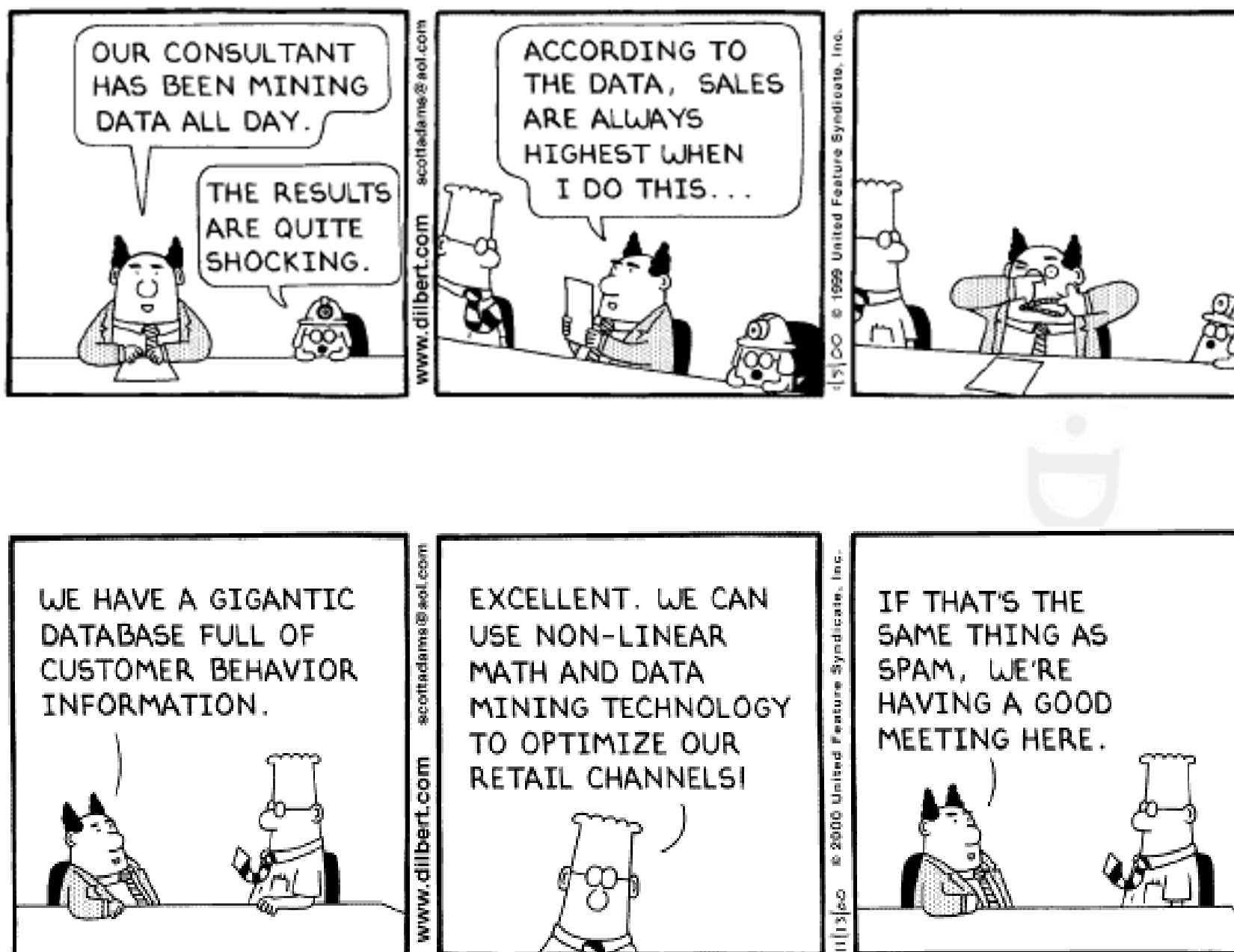*Smyth, 2003*

# Data Issues in Health Science

# Dilbert (1)

# Dilbert (2)

# Challenges of Data Mining

- Scalability

- Dimensionality

- Complex and heterogeneous data

- Data quality

- Data ownership and distribution

- Privacy preservation

- Fairness

- Streaming data (e.g., intrusion detection)



BRIEF HISTORY OF FAIRNESS IN ML

PAPERS

LOL FAIRNESS!!

OH, CRAP.

2011 2012 2013 2014 2015 2016 2017