

# **Artificial Intelligence/ Inteligência Artificial**

## **Lecture 5c: Machine Learning – Data Preprocessing** (adapted from Tan et al, 2020)

**Luís Paulo Reis**

[lpreas@fe.up.pt](mailto:lpreas@fe.up.pt)

Director of LIACC – Artificial Intelligence and Computer Science Lab.  
Associate Professor at DEI/FEUP – Informatics Engineering Department,  
Faculty of Engineering of the University of Porto, Portugal  
President of APPIA – Portuguese Association for Artificial Intelligence

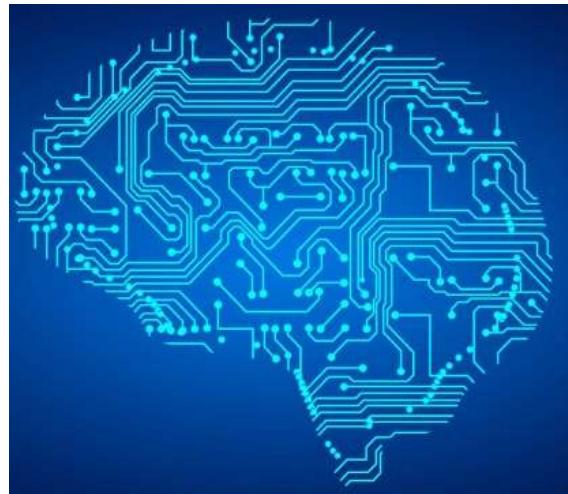


# Agenda

- Large scale data
- Commercial and Scientific interest and vast amount of opportunities
- Data Mining, Classification, Regression
- Knowledge Extraction Methodologies – CRISP-DM
- Data and Data Quality
- Data Pre-Processing (Aggregation, Sampling, Discretization and Binarization, Attribute Transformation, Dimensionality Reduction, Feature subset selection, Feature creation)
- Data Exploration

# Machine Learning

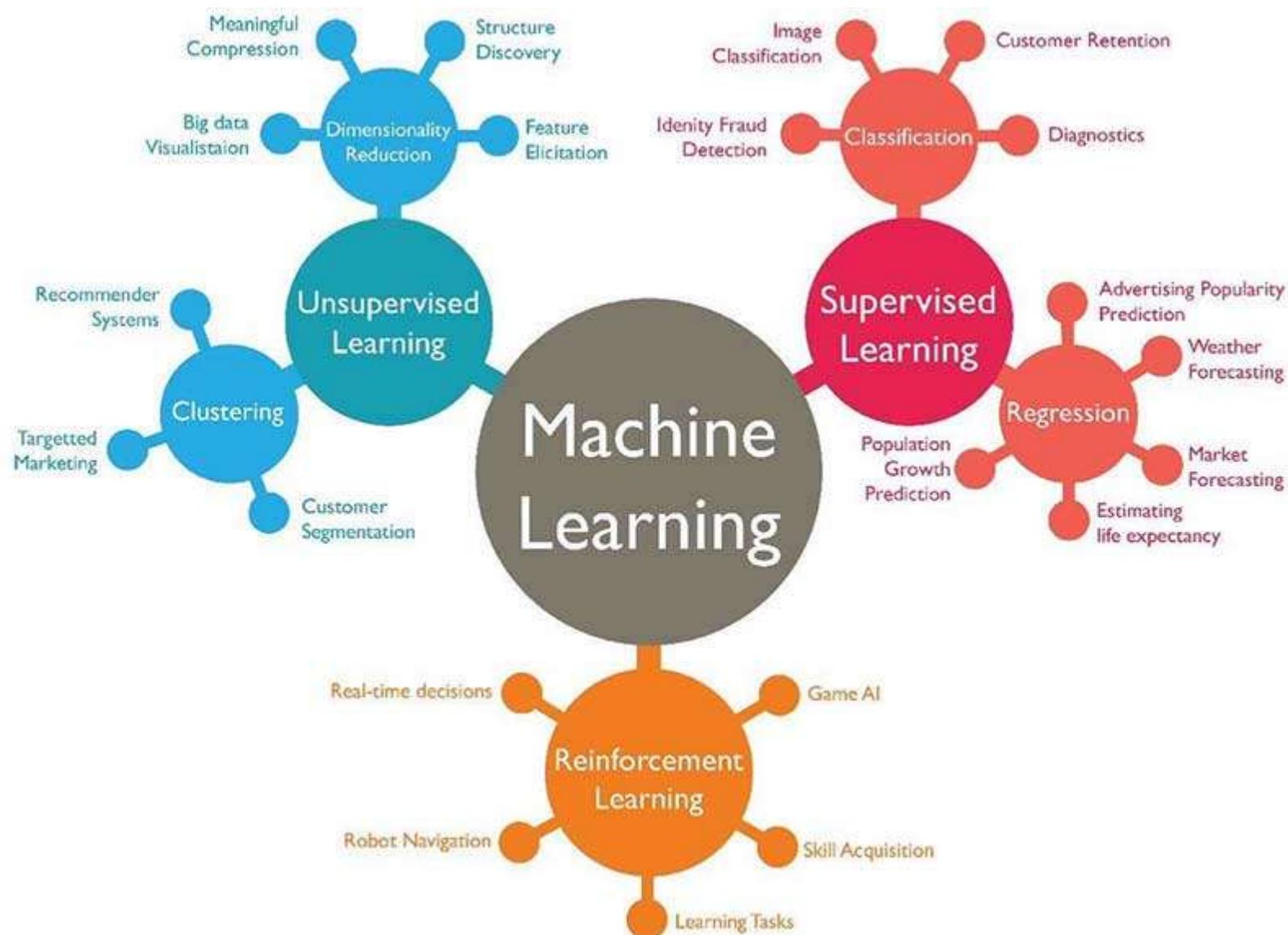
- Machine learning is a field of artificial intelligence that gives computer systems the **ability to "learn"** (e.g., progressively **improve performance** on a specific task) from data/results of their actions, without being explicitly programmed



# Machine Learning

- Machine Learning (ML) Tasks:
  - **Supervised learning:** Example inputs and desired outputs are available/given by a "teacher", and the goal is to learn how to map inputs to outputs (possibility semi-supervised)
  - **Reinforcement learning:** Data (in form of rewards and punishments) are given only as feedback to the computer/agent actions in a dynamic environment
  - **Unsupervised learning:** No labels/outputs are given to the learning algorithm, leaving it on its own to find structure in its input

# Machine Learning



# Large-scale Data is Everywhere!

- There has been enormous data growth in both commercial and scientific databases due to advances in data generation and collection technologies
- Gather whatever data you can whenever and wherever possible
- Expectations: Gathered data will have value either for the purpose collected or for a purpose not envisioned



Cyber Security

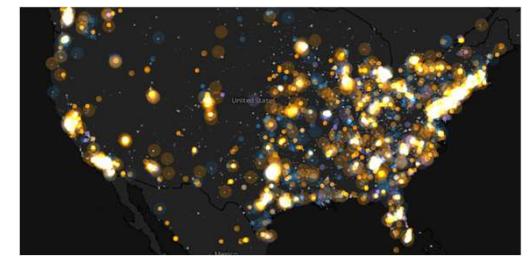


E-Commerce

[HEINSFACCTOR.NETWORK]



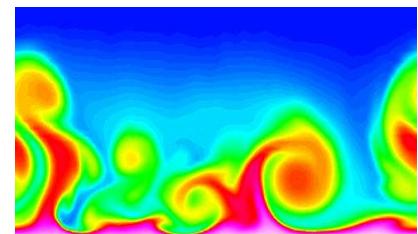
Traffic Patterns



Social Networking: Twitter



Sensor Networks



Computational Simulations

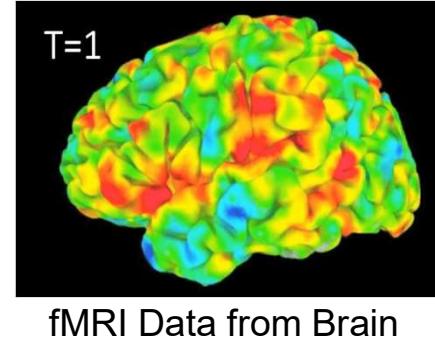
# Data Mining - Commercial Viewpoint

- Lots of data is being collected and warehoused
  - Web data
    - Google has Peta Bytes of web data
    - Facebook has billions of active users
  - Purchases at department/grocery stores, e-commerce
    - Amazon handles millions of visits/day
  - 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)



# Data Mining - Scientific Viewpoint

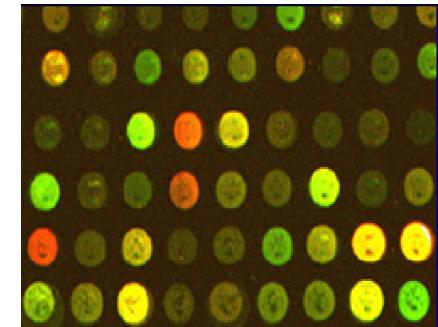
- Data collected and stored at enormous speeds
  - remote sensors on a satellite
    - NASA EOSDIS archives over petabytes of earth science data / year
  - telescopes scanning the skies
    - Sky survey data
  - High-throughput biological data
  - scientific simulations
    - terabytes of data generated in a few hours
- Data mining helps scientists
  - in automated analysis of massive datasets
  - In hypothesis formation



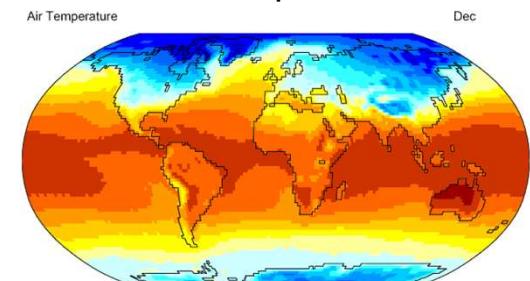
fMRI Data from Brain



Sky Survey Data



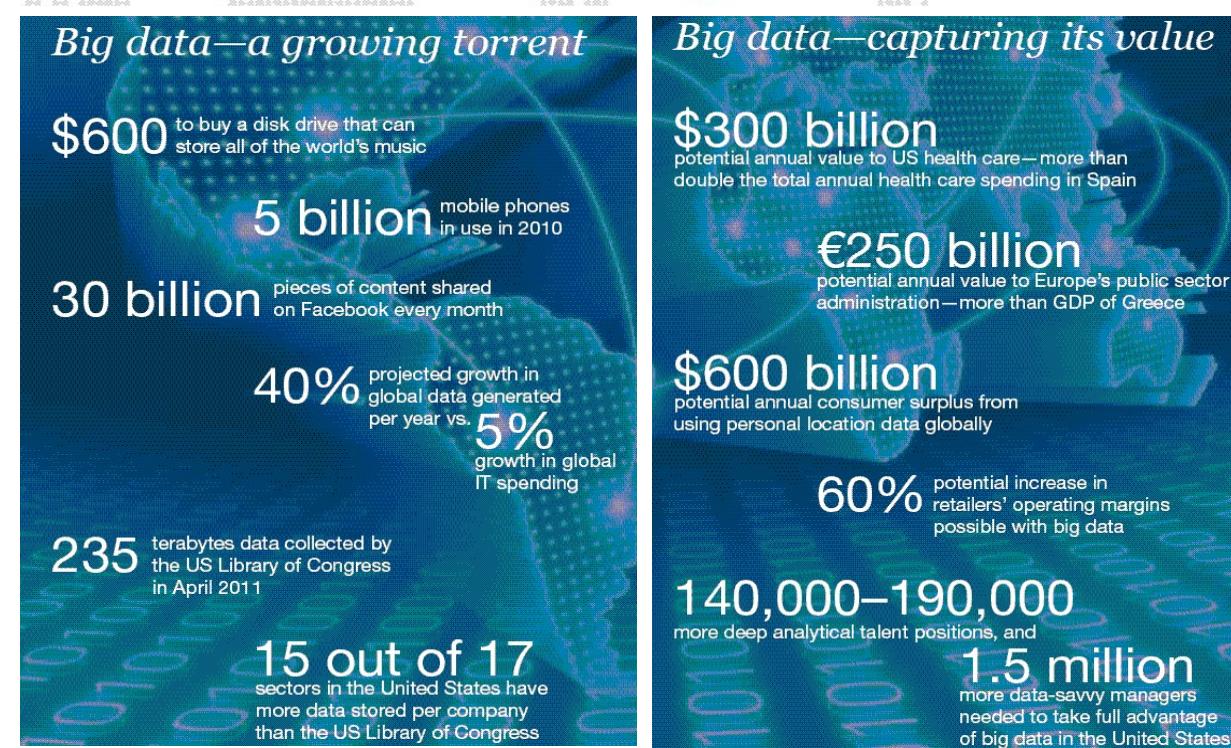
Gene Expression Data



# Opportunities to improve productivity

McKinsey Global Institute

## Big data: The next frontier for innovation, competition, and productivity



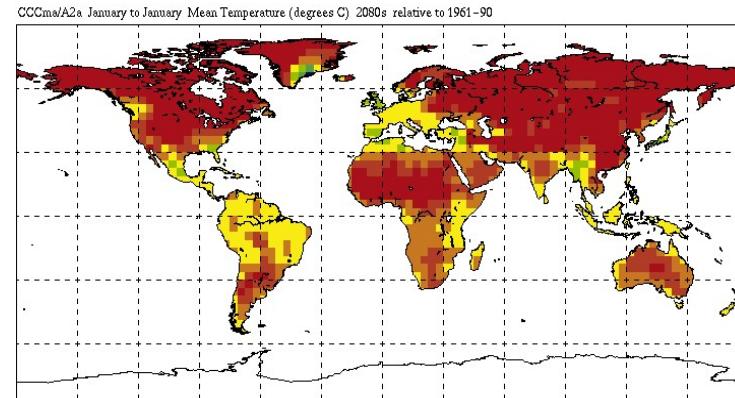
# Opportunities to Solve Society's Major Problems



Improving health care and reducing costs



Finding alternative/ green energy sources



Predicting the impact of climate change

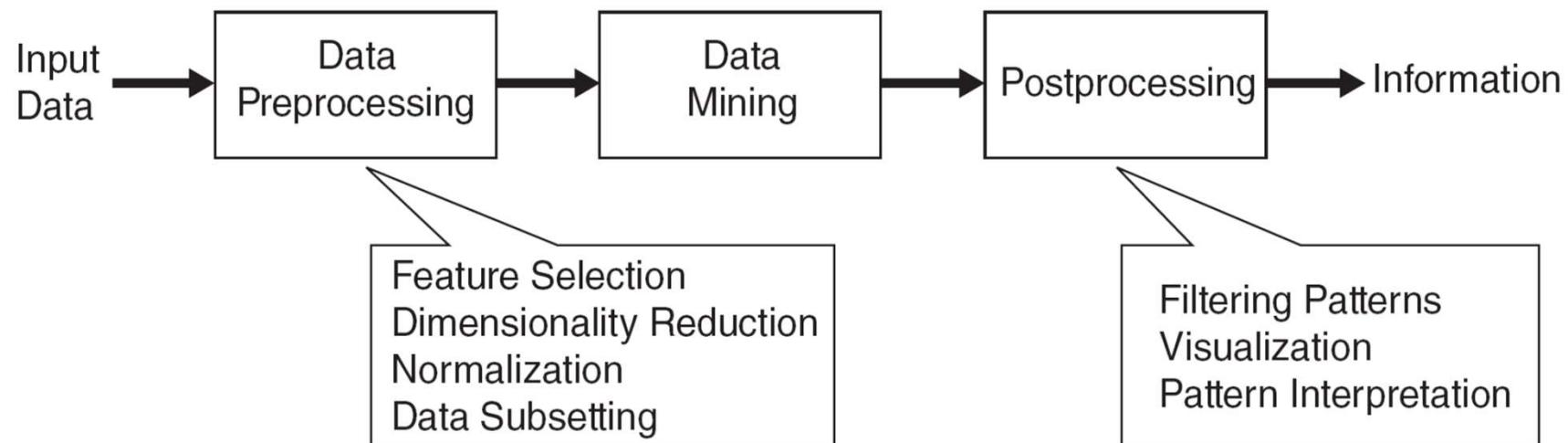


Reducing hunger and poverty by increasing agriculture production

# Data Mining

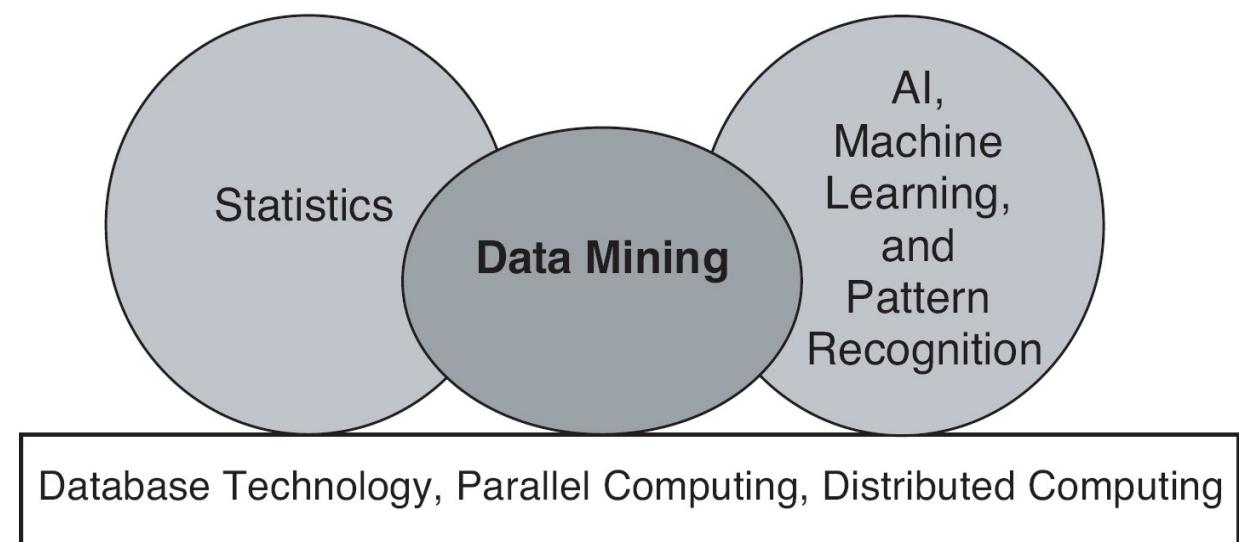
- Many Definitions

- 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



# Origins of Data Mining

- Draws ideas from machine learning/AI, pattern recognition, statistics, and database systems
- Traditional techniques may be unsuitable due to data that is
  - Large-scale
  - High dimensional
  - Heterogeneous
  - Complex
  - Distributed



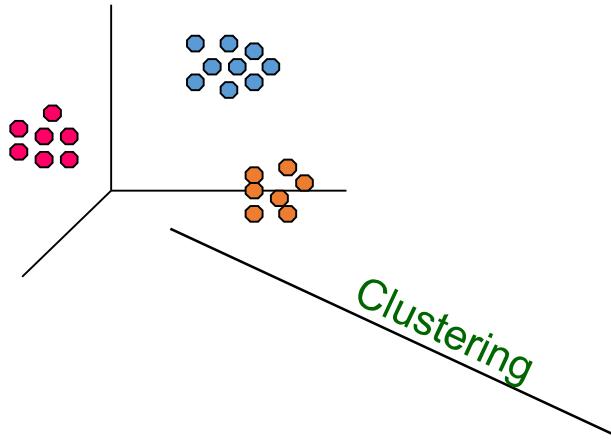
- A key component of the emerging field of data science and data-driven discovery

# Data Mining Tasks

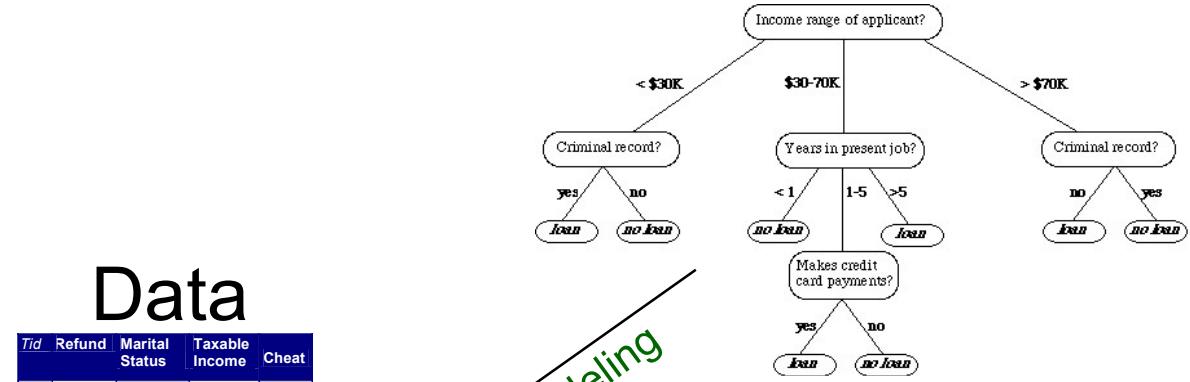
- Prediction Methods
  - Use some variables to predict unknown or future values of other variables
- Description Methods
  - Find human-interpretable patterns that describe the data

[Fayyad, et.al.] Advances in Knowledge Discovery and Data Mining, 1996

# Data Mining Tasks



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
11	No	Married	60K	No
12	Yes	Divorced	220K	No
13	No	Single	85K	Yes
14	No	Married	75K	No
15	No	Single	90K	Yes



Association  
Rules

Anomaly  
Detection

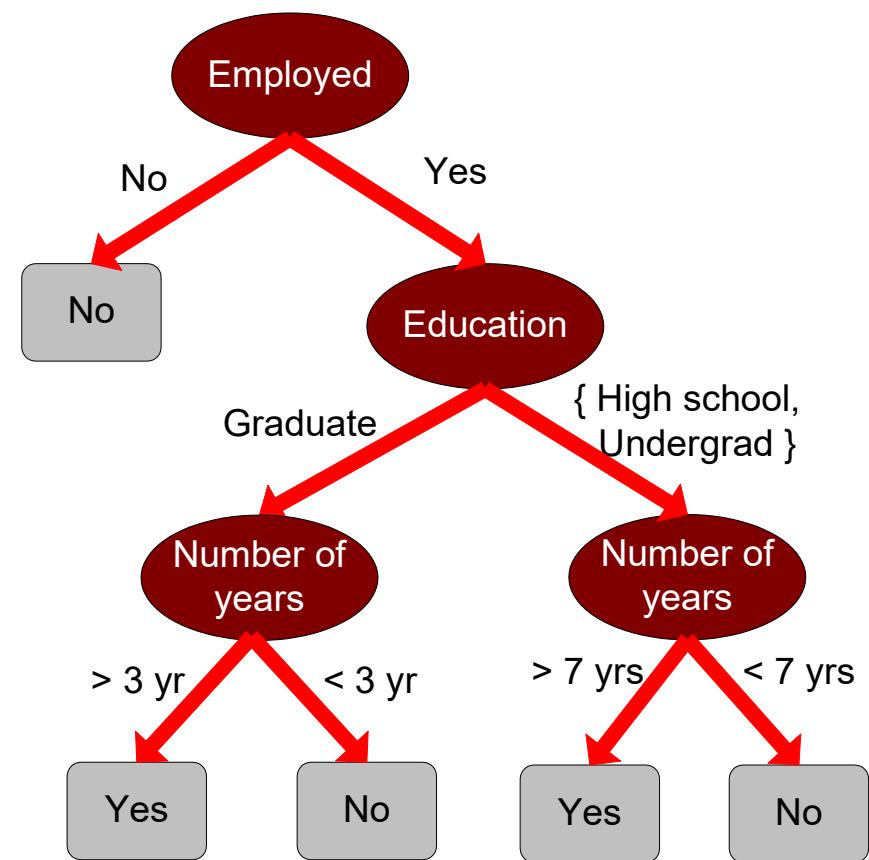
Predictive Modeling

# Predictive Modeling: Classification

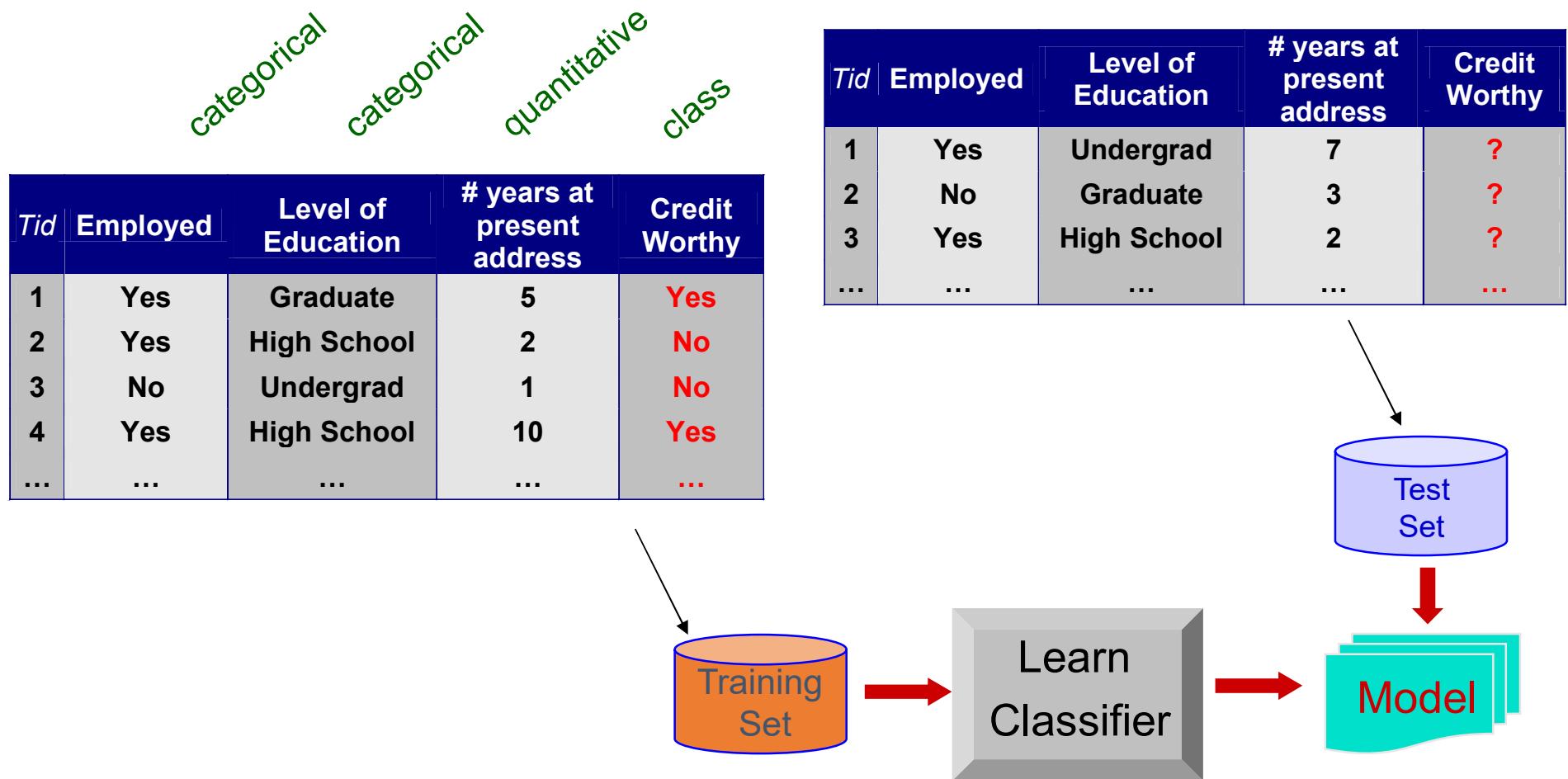
- Find a model for class attribute as a function of the values of other attributes

Tid	Employed	Level of Education	# years at present address	Credit Worthy
1	Yes	Graduate	5	Yes
2	Yes	High School	2	No
3	No	Undergrad	1	No
4	Yes	High School	10	Yes
...	...	...	...	...

Model for predicting credit worthiness

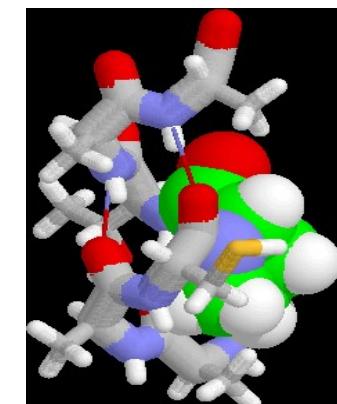


# Classification Example



# Examples of Classification Task

- Classifying credit card transactions as legitimate or fraudulent
- Classifying land covers (water bodies, urban areas, forests, etc.) using satellite data
- Categorizing news stories as finance, weather, entertainment, sports, etc
- Identifying intruders in the cyberspace
- Predicting tumor cells as benign or malignant
- Classifying secondary structures of protein as alpha-helix, beta-sheet, or random coil



# Classification: Application 1

- 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 2

- Churn prediction for telephone customers
  - **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

[Berry & Linoff] Data Mining Techniques, 1997

# Classification: Application 3

- Sky Survey Cataloging
  - **Goal:** To predict class (star or galaxy) of sky objects, especially visually faint ones, based on the telescopic survey images
    - 3000 images with  $23,040 \times 23,040$  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!

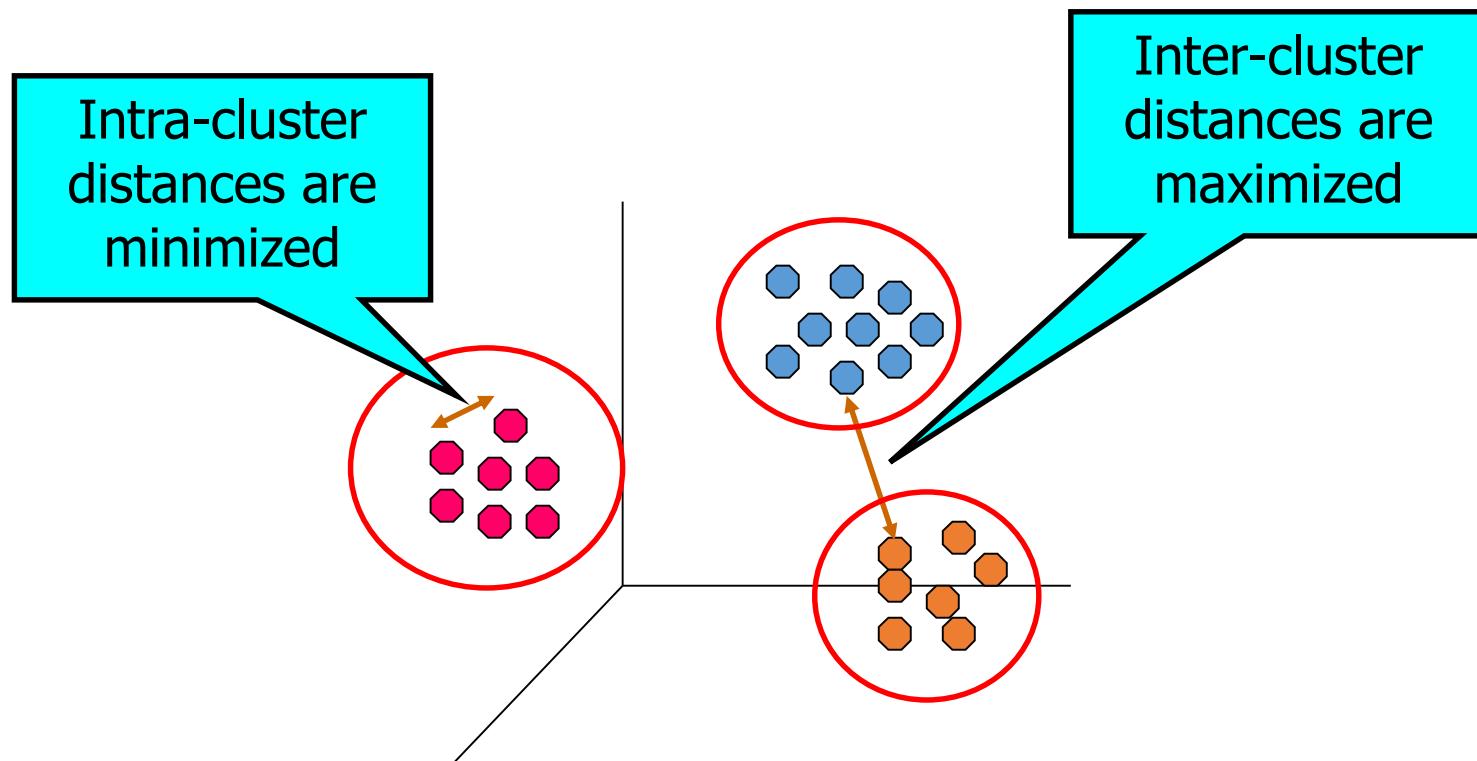
[Fayyad, et.al.] Advances in Knowledge Discovery and Data Mining, 1996

# 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
- Extensively studied in statistics, neural network fields
- 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

# Unsupervised Learning - Clustering

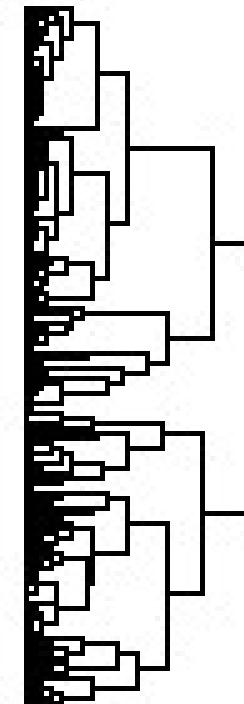
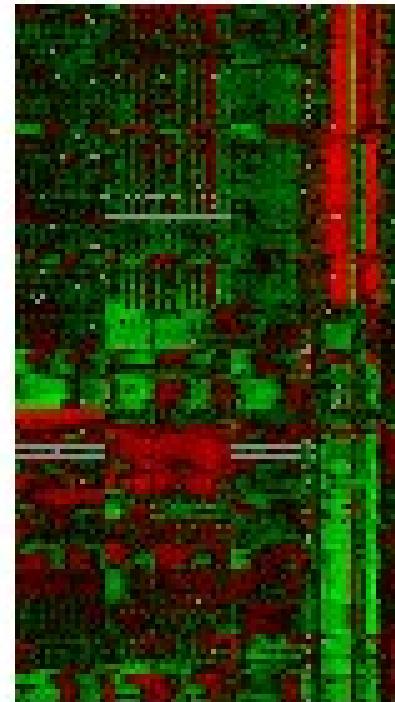
Finding groups of objects such that the objects in a group will be similar (or related) to one another and different from (or unrelated to) the objects in other groups



# Applications of Cluster Analysis

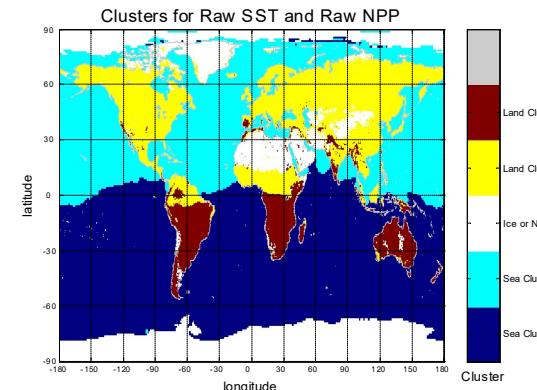
- **Understanding**

- Custom profiling for targeted marketing
- Group related documents for browsing
- Group genes and proteins that have similar functionality
- Group stocks with similar price fluctuations



- **Summarization**

- Reduce the size of large data sets



# 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



# Unsupervised Learning - Association Rule Discovery

- 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.

<i>TID</i>	<i>Items</i>
1	Bread, Coke, Milk
2	Beer, Bread
3	Beer, Coke, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Coke, Diaper, Milk

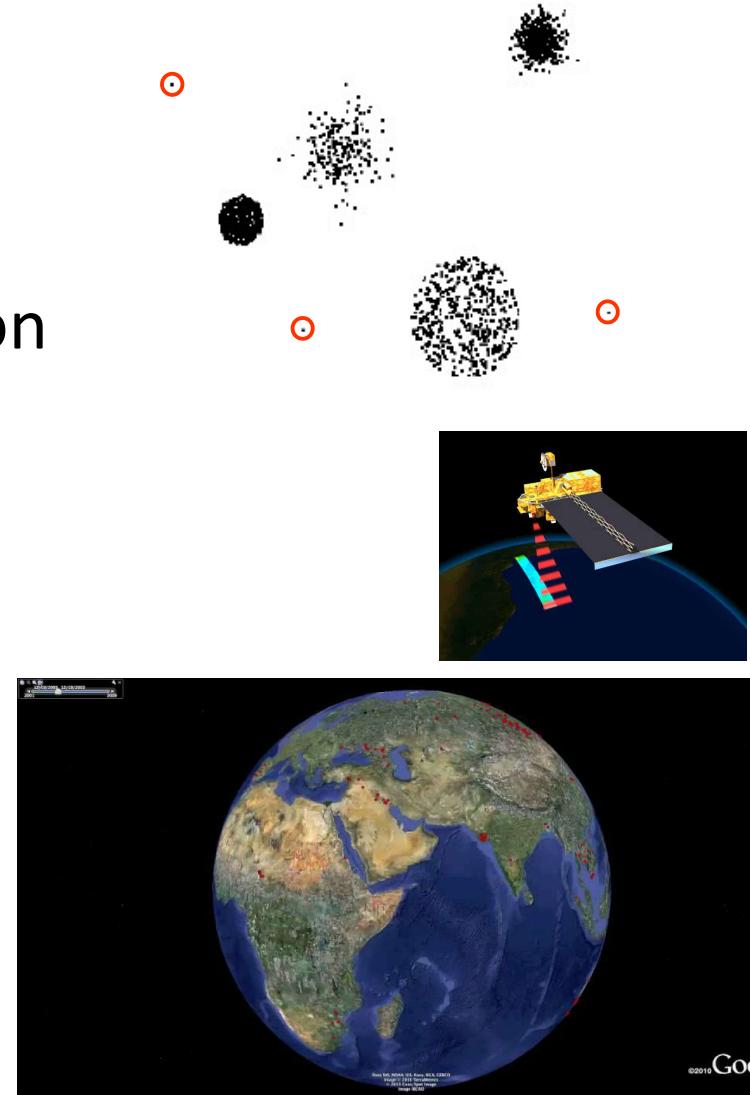
Rules Discovered:  
 $\{\text{Milk}\} \rightarrow \{\text{Coke}\}$   
 $\{\text{Diaper}, \text{Milk}\} \rightarrow \{\text{Beer}\}$

# Association Analysis: Applications

- Market-basket analysis
  - Rules are used for sales promotion, shelf management, and inventory management
- Telecommunication alarm diagnosis
  - Rules are used to find combination of alarms that occur together frequently in the same time period
- Medical Informatics
  - Rules are used to find combination of patient symptoms and test results associated with certain diseases

# Deviation/Anomaly/Change Detection

- Detect significant deviations from normal behavior
- Applications:
  - Credit Card Fraud Detection
  - Network Intrusion Detection
  - Identify anomalous behavior from sensor networks for monitoring and surveillance
  - Detecting changes in the global forest cover



# Motivating Challenges

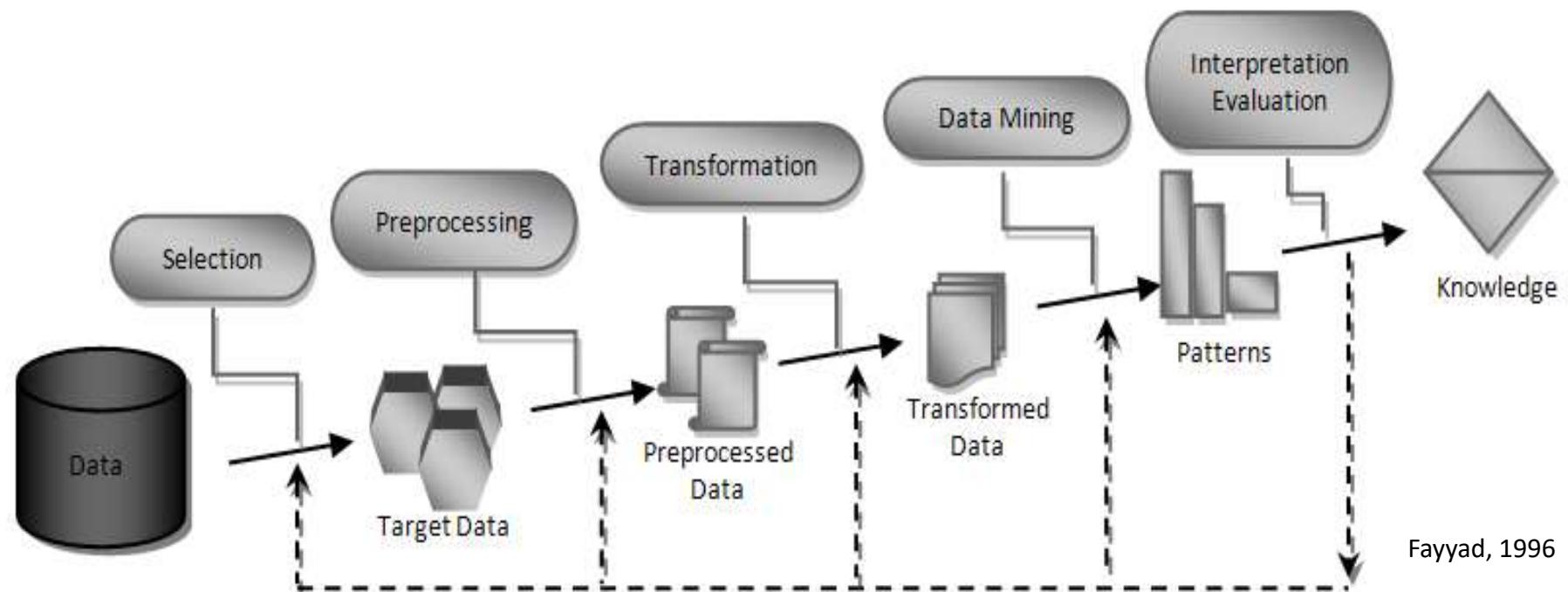
- Scalability
- High Dimensionality
- Heterogeneous and Complex Data
- Data Ownership and Distribution
- Non-traditional Analysis

# Knowledge Extraction Methodologies

- Extração de Conhecimento (Knowledge Discovery) and Mineração de Dados (Data Mining)
- Different approaches to extract knowledge from data:
  - **KDD** (Knowledge Discovery in Databases)
  - **SEMMA** (Sample, Explore, Modify, Model and Access)
  - **CRISP-DM** (Cross-Industry Standard for Data Mining)

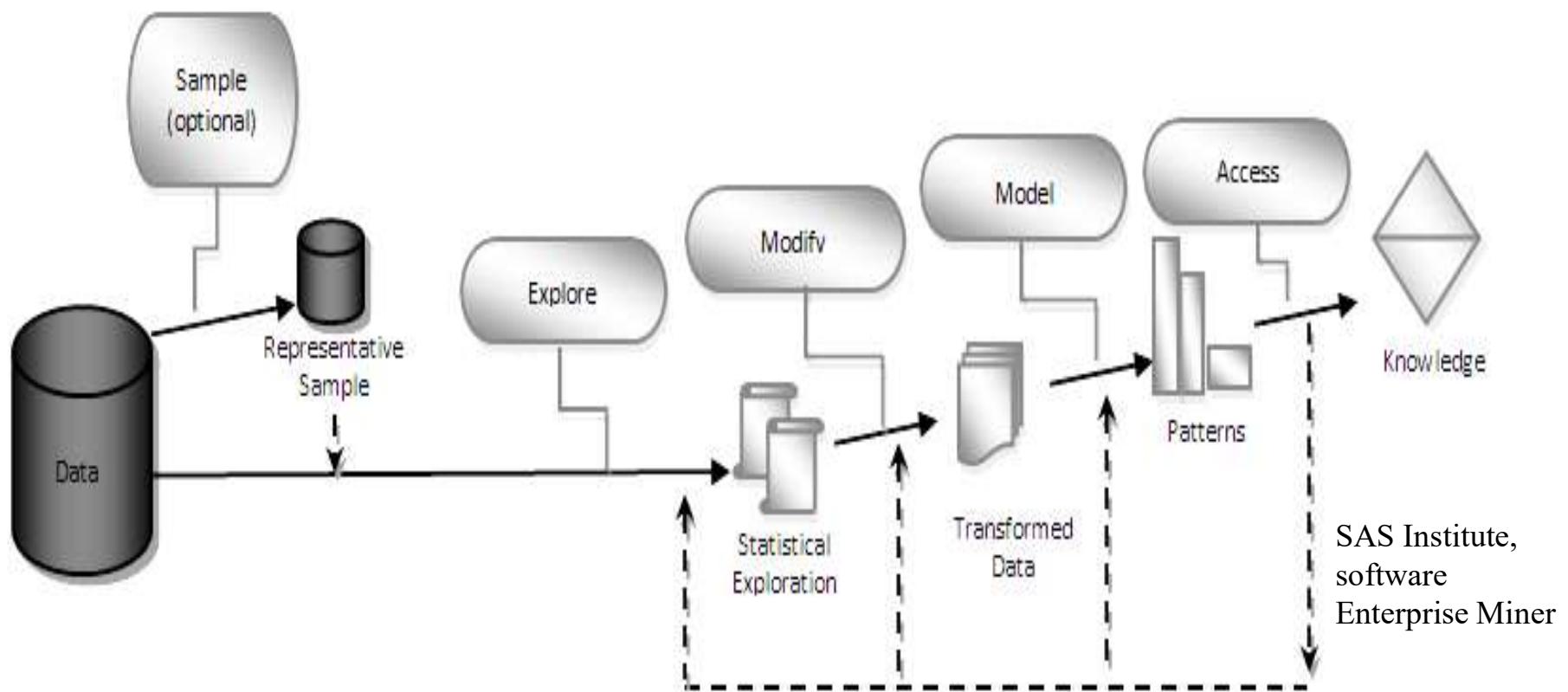
# KDD - Knowledge Discovery in Databases

- KDD (Knowledge Discovery in Databases)



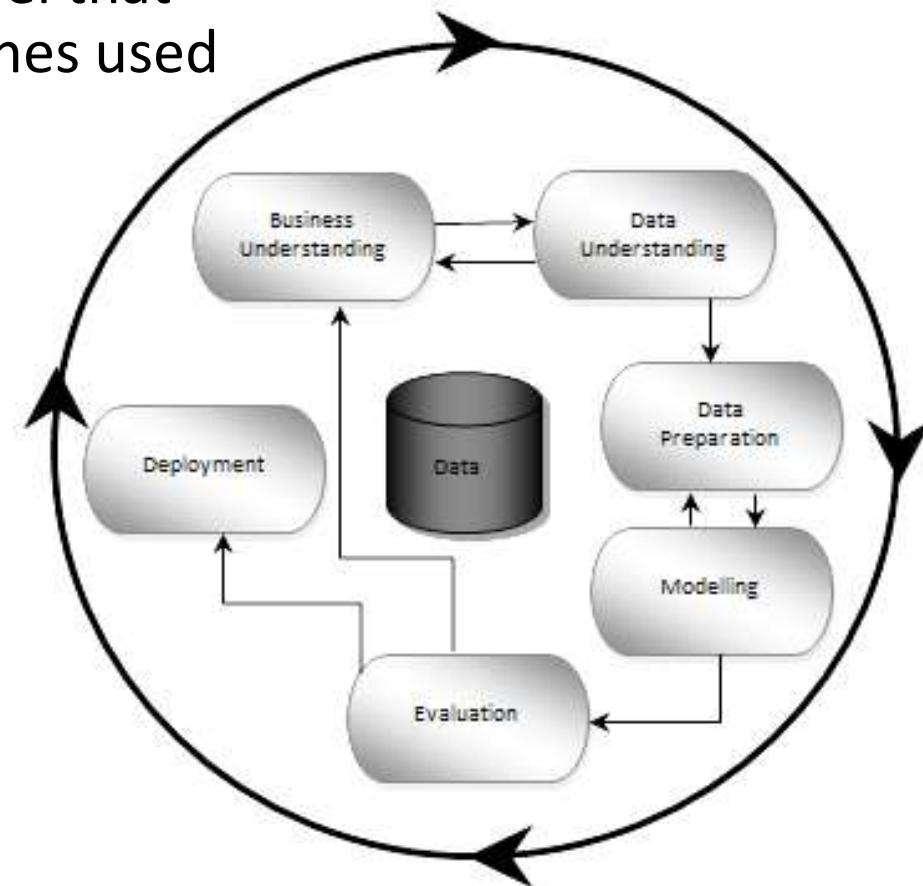
# SEMMA Phases

- SEMMA (Sample, Explore, Modify, Model and Access)



# CRISP-DM Phases

- CRISP-DM (Cross-Industry Standard for Data Mining)
- Open standard process model that describes common approaches used by data mining experts.
- It is the most widely-used analytics model
- European Commission and 4 Companies: Daimler-Benz (now DaimlerChrysler), Integral Solutions Ltd. (now IBM-SPSS), NCR and OHRA, 2000

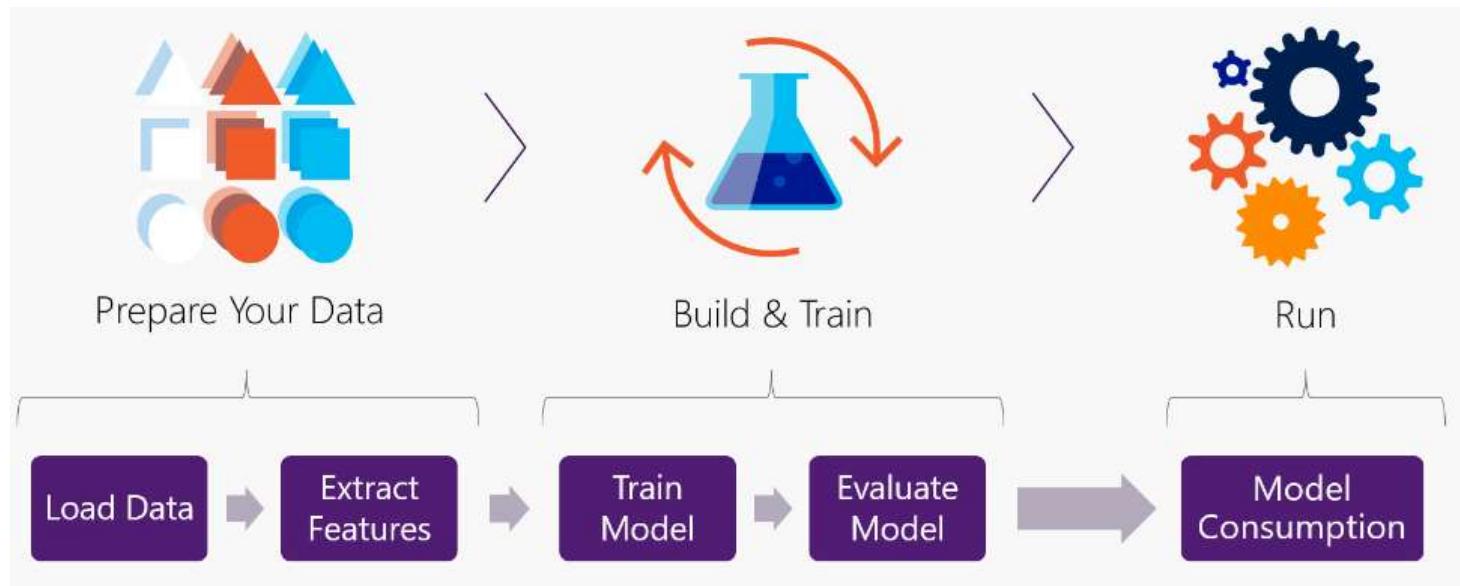


# Knowledge Extraction Methodologies

Data Mining Process Models	KDD	CRISP-DM	SEMMA
No. of Steps	9	6	5
Name of Steps	Developing and Understanding of the Application Creating a Target Data Set Data Cleaning and Pre-processing Data Transformation Choosing the suitable Data Mining Task Choosing the suitable Data Mining Algorithm Employing Data Mining Algorithm Interpreting Mined Patterns Using Discovered Knowledge	Business Understanding Data Understanding Data Preparation Modeling Evaluation Deployment	----- Sample Explore Modify Model Assessment -----

# Data and Data Pre-processing

- Attributes and Objects
- Types of Data
- Data Quality
- Similarity and Distance
- Data Preprocessing



# Data

- Collection of ***data objects*** and their ***attributes***

- An ***attribute*** is a property or characteristic of an object
  - Examples: eye color of a person, temperature, etc.
  - Attribute is also known as variable, field, characteristic, dimension, or feature
- A collection of attributes describe an ***object***
  - Object is also known as record, point, case, sample, entity, or instance

**Attributes**

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

# Attribute Values

- **Attribute values** are numbers or symbols assigned to an attribute for a particular object
- Distinction between attributes and attribute values
  - Same attribute can be mapped to different attribute values
    - Example: height can be measured in feet or meters
  - Different attributes can be mapped to the same set of values
    - Example: Attribute values for ID and age are integers
  - But properties of attribute can be different than the properties of the values used to represent the attribute

# Types of Attributes

- There are different types of attributes
  - Nominal
    - Examples: ID numbers, eye color, zip codes
  - Ordinal
    - Examples: rankings (e.g., taste of potato chips on a scale from 1-10), grades, height {tall, medium, short}
  - Interval
    - Examples: calendar dates, temperatures in Celsius or Fahrenheit.
  - Ratio
    - Examples: temperature in Kelvin, length, counts, elapsed time (e.g., time to run a race)

# Properties of Attribute Values

- The type of an attribute depends on which of the following properties/operations it possesses:
  - Distinctness:                            $= \neq$
  - Order:                                    $< >$
  - Differences are                        $+$   $-$   
meaningful :
  - Ratios are                               $*$   $/$   
meaningful
  - Nominal attribute: distinctness
  - Ordinal attribute: distinctness & order
  - Interval attribute: distinctness, order & meaningful differences
  - Ratio attribute: all 4 properties/operations

<b>Attribute Type</b>	<b>Description</b>	<b>Examples</b>	<b>Operations</b>
Categorical Qualitative	Nominal Nominal attribute values only distinguish. ( $=$ , $\neq$ )	zip codes, employee ID numbers, eye color, sex: { <i>male</i> , <i>female</i> }	mode, entropy, contingency correlation, $\chi^2$ test
	Ordinal Ordinal attribute values also order objects. ( $<$ , $>$ )	hardness of minerals, { <i>good</i> , <i>better</i> , <i>best</i> }, grades, street numbers	median, percentiles, rank correlation, run tests, sign tests
Numeric Quantitative	Interval For interval attributes, differences between values are meaningful. (+, -)	calendar dates, temperature in Celsius or Fahrenheit	mean, standard deviation, Pearson's correlation, <i>t</i> and <i>F</i> tests
	Ratio For ratio variables, both differences and ratios are meaningful. (*, /)	temperature in Kelvin, monetary quantities, counts, age, mass, length, current	geometric mean, harmonic mean, percent variation

<b>Attribute Type</b>	<b>Transformation</b>	<b>Comments</b>
Categorical Qualitative	Nominal	Any permutation of values If all employee ID numbers were reassigned, would it make any difference?
	Ordinal	An order preserving change of values, i.e., $new\_value = f(old\_value)$ where $f$ is a monotonic function An attribute encompassing the notion of good, better best can be represented equally well by the values {1, 2, 3} or by { 0.5, 1, 10}.
Numeric Quantitative Q	Interval	$new\_value = a * old\_value + b$ where a and b are constants Thus, the Fahrenheit and Celsius temperature scales differ in terms of where their zero value is and the size of a unit (degree).
	Ratio	$new\_value = a * old\_value$ Length can be measured in meters or feet.

# Discrete and Continuous Attributes

- Discrete Attribute
  - Has only a finite or countably infinite set of values
  - Examples: zip codes, counts, or the set of words in a collection of documents
  - Often represented as integer variables.
  - Note: **binary attributes** are a special case of discrete attributes
- Continuous Attribute
  - Has real numbers as attribute values
  - Examples: temperature, height, or weight.
  - Practically, real values can only be measured and represented using a finite number of digits.
  - Continuous attributes are typically represented as floating-point variables.

# Key Messages for Attribute Types

- The types of operations you choose should be “meaningful” for the type of data you have
  - Distinctness, order, meaningful intervals, and meaningful ratios are only four (among many possible) properties of data
  - The data type you see – often numbers or strings – may not capture all the properties or may suggest properties that are not present
  - Analysis may depend on these other properties of the data
    - Many statistical analyses depend only on the distribution
  - In the end, what is meaningful can be specific to domain

# Important Characteristics of Data

- Dimensionality (number of attributes)
  - High dimensional data brings a large number of challenges
- Sparsity
  - Only presence counts
- Resolution
  - Patterns depend on the scale
- Size
  - Type of analysis may depend on size of data. Data of large size poses several challenges

# Types of data sets

- Record
  - Data Matrix
  - Document Data
  - Transaction Data
- Graph
  - World Wide Web
  - Molecular Structures
- Ordered
  - Spatial Data
  - Temporal Data
  - Sequential Data
  - Genetic Sequence Data

# Record Data

- Data that consists of a collection of records, each of which consists of a fixed set of attributes

<i>Tid</i>	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

# Data Matrix

- If data objects have the same fixed set of numeric attributes, then the data objects can be thought of as points in a multi-dimensional space, where each dimension represents a distinct attribute
- Such a data set can be represented by an  $m$  by  $n$  matrix, where there are  $m$  rows, one for each object, and  $n$  columns, one for each attribute

Projection of x Load	Projection of y load	Distance	Load	Thickness
10.23	5.27	15.22	2.7	1.2
12.65	6.25	16.22	2.2	1.1

# Document Data

- Each document becomes a ‘term’ vector
  - Each term is a component (attribute) of the vector
  - The value of each component is the number of times the corresponding term occurs in the document.

	team	coach	play	ball	score	game	win	lost	timeout	season
Document 1	3	0	5	0	2	6	0	2	0	2
Document 2	0	7	0	2	1	0	0	3	0	0
Document 3	0	1	0	0	1	2	2	0	3	0

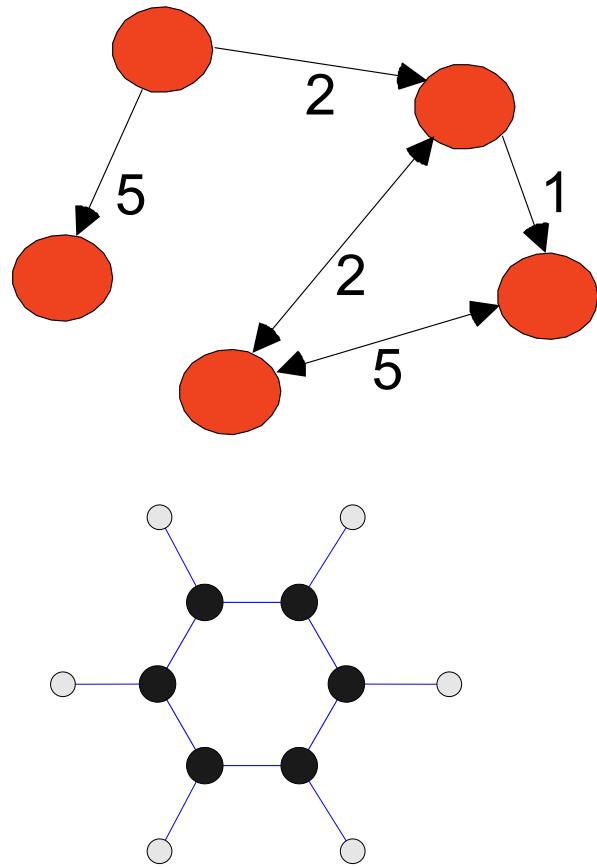
# Transaction Data

- A special type of data, where
  - Each transaction involves a set of items.
  - For example, consider a grocery store. The set of products purchased by a customer during one shopping trip constitute a transaction, while the individual products that were purchased are the items.
  - Can represent transaction data as record data

<i>TID</i>	<i>Items</i>
1	Bread, Coke, Milk
2	Beer, Bread
3	Beer, Coke, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Coke, Diaper, Milk

# Graph Data

- Examples: Generic graph, a molecule, and webpages



**Useful Links:**

- [Bibliography](#)
- Other Useful Web sites
  - [ACM SIGKDD](#)
  - [KDnuggets](#)
  - [The Data Mine](#)

**Knowledge Discovery and Data Mining Bibliography**  
(Gets updated frequently, so visit often!)

- [Books](#)
- [General Data Mining](#)

**Book References in Data Mining and Knowledge Discovery**

Usama Fayyad, Gregory Piatetsky-Shapiro, Padhraic Smyth, and Ramasamy Iyer, "Advances in Knowledge Discovery and Data Mining", AAAI Press/the MIT Press, 1996.

J. Ross Quinlan, "C4.5: Programs for Machine Learning", Morgan Kaufmann Publishers, 1993.

Michael Berry and Gordon Linoff, "Data Mining Techniques (For Marketing, Sales, and Customer Support)", John Wiley & Sons, 1997.

**General Data Mining**

Usama Fayyad, "Mining Databases: Towards Algorithms for Knowledge Discovery", Bulletin of the IEEE Computer Society Technical Committee on Data Engineering, vol. 21, no. 1, March 1998.

Christopher Matheus, Philip Chan, and Gregory Piatetsky-Shapiro, "Systems for knowledge Discovery in databases", IEEE Transactions on Knowledge and Data Engineering, 5(6):903-913, December 1993.

# Ordered Data

- Sequences of transactions

## Items/Events

( A B )    ( D )    ( C E )  
( B D )    ( C )    ( E )  
( C D )    ( B )    ( A E )



An element of  
the sequence

- Genomic sequence data

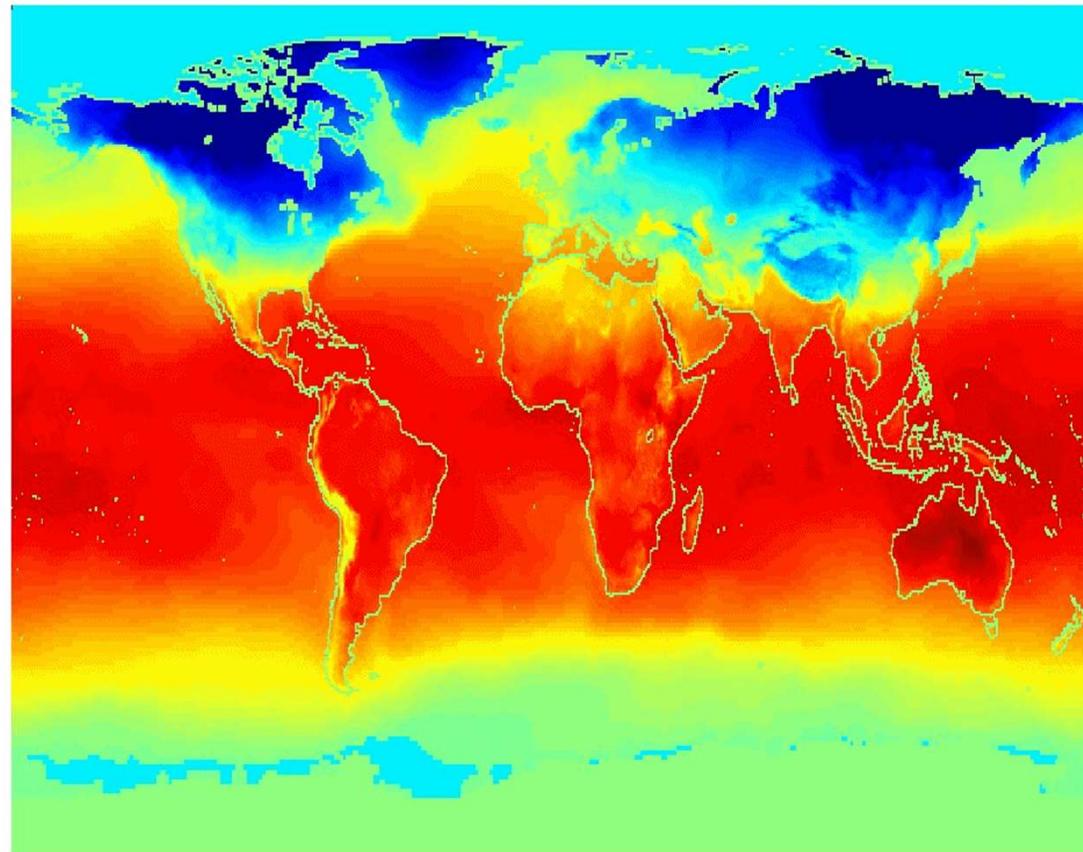
GGTTCCGCCTTCAGCCCCGCGCC  
CGCAGGGCCC GCCCGCGCCGTC  
GAGAAGGGCCC GCCTGGCGGGCG  
GGGGGAGGC GGGGCCGCCGAGC  
CCAACCGAGTCCGACCAGGTGCC  
CCCTCTGCTCGGCCTAGACCTGA  
GCTCATTAGGC GGCA GCGGACAG  
GCCAAGTAGAACACACGCGAAGCGC  
TGGGCTGCCTGCTGCGACCAGGG

# Ordered Data

- Spatio-Temporal Data

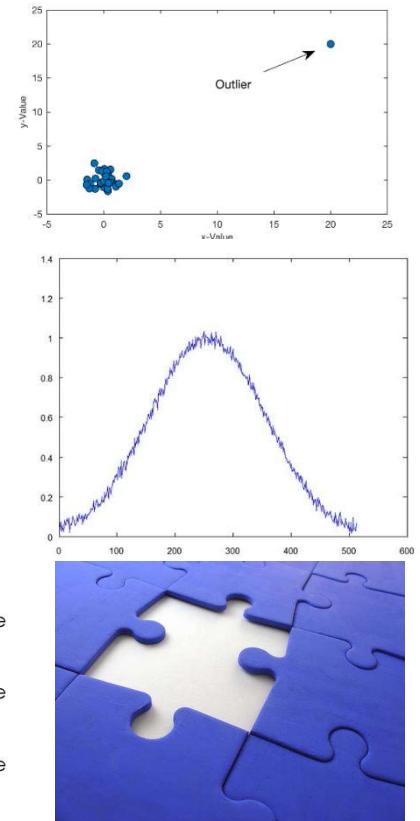
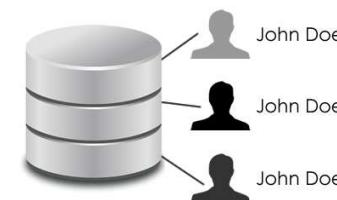
**Average Monthly Temperature of land and ocean**

Jan



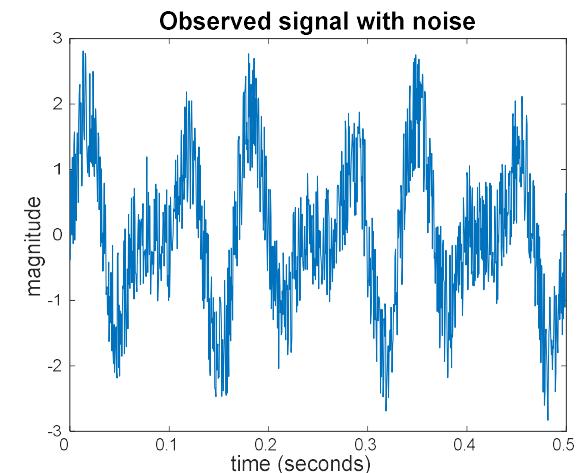
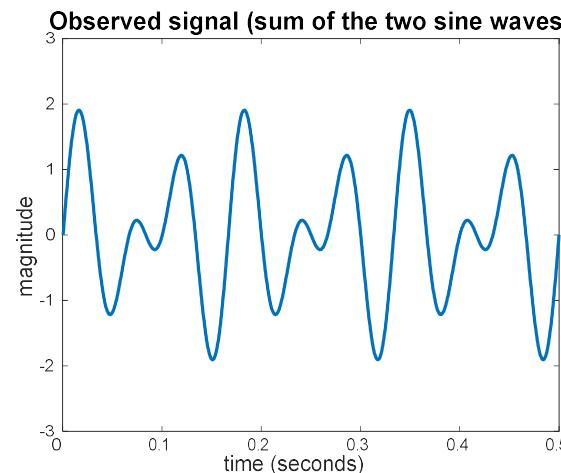
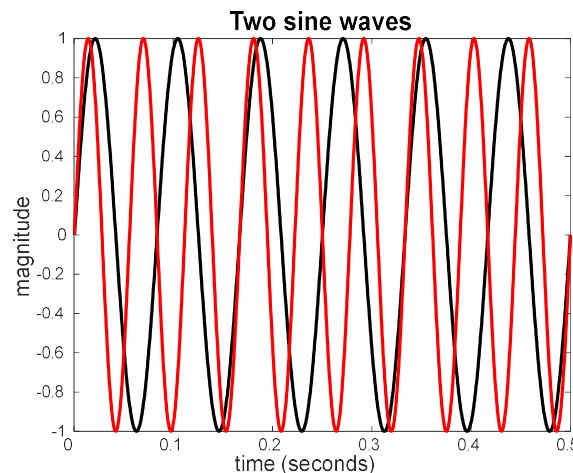
# Data Quality

- Poor data quality negatively affects many data processing efforts
- Data mining example: a classification model for detecting people who are loan risks is built using poor data
  - Some credit-worthy candidates are denied loans
  - More loans are given to individuals that default
- What kinds of data quality problems?
- How can we detect problems with the data?
- What can we do about these problems?
- Examples of data quality problems:
  - Noise and outliers
  - Wrong data
  - Fake data
  - Missing values
  - Duplicate data



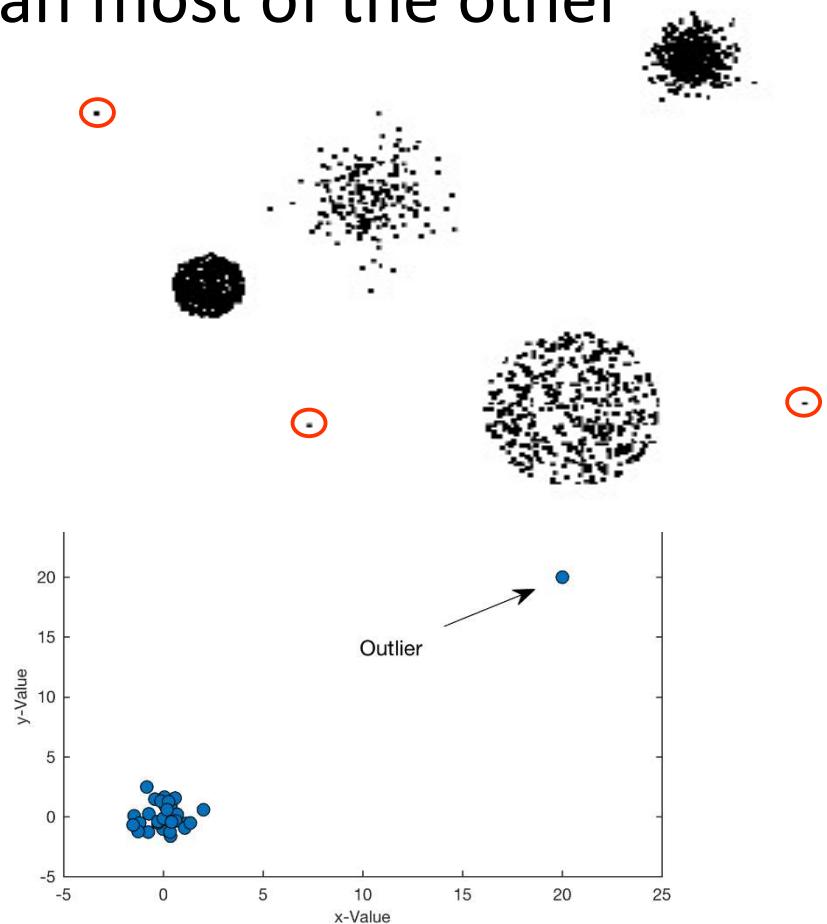
# Noise

- For objects, noise is an extraneous object
- For attributes, noise refers to modification of original values
  - Examples: distortion of a person's voice when talking on a poor phone and "snow" on television screen
  - The figures below show two sine waves of the same magnitude and different frequencies, the waves combined, and the two sine waves with random noise
    - The magnitude and shape of the original signal is distorted



# Outliers

- **Outliers** are data objects with characteristics that are considerably different than most of the other data objects in the data set
  - **Case 1:** Outliers are noise that interferes with data analysis
  - **Case 2:** Outliers are the goal of our analysis
    - Credit card fraud
    - Intrusion detection
- Causes?



# Missing Values

- Reasons for missing values
  - Information is not collected  
(e.g., people decline to give their age and weight)
  - Attributes may not be applicable to all cases  
(e.g., annual income is not applicable to children)
- Handling missing values
  - Eliminate data objects or variables
  - Estimate missing values
    - Example: time series of temperature
    - Example: census results
  - Ignore the missing value during analysis

# Duplicate Data

- Data set may include data objects that are duplicates, or almost duplicates of one another
  - Major issue when merging data from heterogeneous sources
- Examples:
  - Same person with multiple email addresses
- Data cleaning
  - Process of dealing with duplicate data issues
- When should duplicate data not be removed?

# Similarity and Dissimilarity Measures

- Similarity measure
  - Numerical measure of how alike two data objects are.
  - Is higher when objects are more alike.
  - Often falls in the range [0,1]
- Dissimilarity measure
  - Numerical measure of how different two data objects are
  - Lower when objects are more alike
  - Minimum dissimilarity is often 0
  - Upper limit varies
- **Proximity** refers to a similarity or dissimilarity

# Similarity/Dissimilarity for Simple Attributes

Similarity and dissimilarity between two objects,  $x$  and  $y$ , with respect to a single, simple attribute.

Attribute Type	Dissimilarity	Similarity
Nominal	$d = \begin{cases} 0 & \text{if } x = y \\ 1 & \text{if } x \neq y \end{cases}$	$s = \begin{cases} 1 & \text{if } x = y \\ 0 & \text{if } x \neq y \end{cases}$
Ordinal	$d =  x - y /(n - 1)$ (values mapped to integers 0 to $n-1$ , where $n$ is the number of values)	$s = 1 - d$
Interval or Ratio	$d =  x - y $	$s = -d, s = \frac{1}{1+d}, s = e^{-d}, s = 1 - \frac{d - \min_d}{\max_d - \min_d}$

# Euclidean Distance

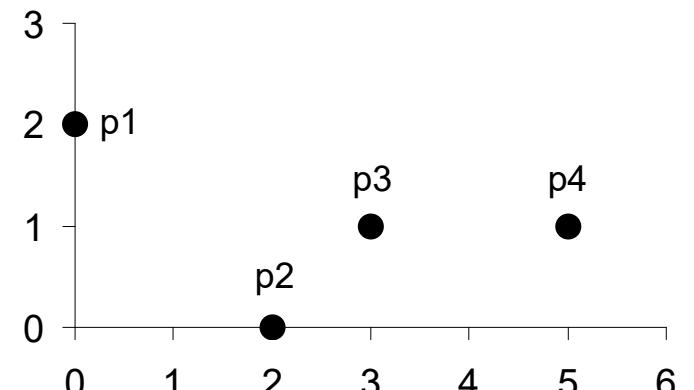
- Euclidean Distance  $d(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{k=1}^n (x_k - y_k)^2}$

where  $n$  is the number of dimensions (attributes) and  $x_k$  and  $y_k$  are, respectively, the  $k^{th}$  attributes (components) or data objects  $\mathbf{x}$  and  $\mathbf{y}$ .

Standardization is necessary,  
if scales differ

	p1	p2	p3	p4
p1	0	2.828	3.162	5.099
p2	2.828	0	1.414	3.162
p3	3.162	1.414	0	2
p4	5.099	3.162	2	0

Distance Matrix



point	x	y
p1	0	2
p2	2	0
p3	3	1
p4	5	1

# Minkowski Distance, Manhattan Distance

Minkowski Distance is a generalization of Euclidean Distance

Where  $r$  is a parameter,  $n$  is the number of dimensions (attributes) and  $x_k$  and  $y_k$  are, respectively, the  $k^{\text{th}}$  attributes (components) or data objects  $x$  and  $y$ .

$$d(\mathbf{x}, \mathbf{y}) = \left( \sum_{k=1}^n |x_k - y_k|^r \right)^{1/r}$$

L1	p1	p2	p3	p4
p1	0	4	4	6
p2	4	0	2	4
p3	4	2	0	2
p4	6	4	2	0

- $r = 1$ . City block (Manhattan, taxicab,  $L_1$  norm) distance.
  - A common example of this for binary vectors is the Hamming distance, which is just the number of bits that are different between two binary vectors
- $r = 2$ . Euclidean distance
- $r \rightarrow \infty$ . “supremum” ( $L_{\max}$  norm,  $L_\infty$  norm) distance.
  - This is the maximum difference between any component of the vectors

# Common Properties of a Distance

- Distances, such as the Euclidean distance, have some well known properties.
  1.  $d(\mathbf{x}, \mathbf{y}) \geq 0$  for all  $\mathbf{x}$  and  $\mathbf{y}$  and  $d(\mathbf{x}, \mathbf{y}) = 0$  if and only if  $\mathbf{x} = \mathbf{y}$ .
  2.  $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})$  for all  $\mathbf{x}$  and  $\mathbf{y}$ . (Symmetry)
  3.  $d(\mathbf{x}, \mathbf{z}) \leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z})$  for all points  $\mathbf{x}$ ,  $\mathbf{y}$ , and  $\mathbf{z}$ . (Triangle Inequality)
- A distance that satisfies these properties is a **metric**

# Common Properties of a Similarity

- Similarities, also have some well known properties.
  1.  $s(\mathbf{x}, \mathbf{y}) = 1$  (or maximum similarity) only if  $\mathbf{x} = \mathbf{y}$ .  
(does not always hold, e.g., cosine)
  2.  $s(\mathbf{x}, \mathbf{y}) = s(\mathbf{y}, \mathbf{x})$  for all  $\mathbf{x}$  and  $\mathbf{y}$ . (Symmetry)

where  $s(\mathbf{x}, \mathbf{y})$  is the similarity between points (data objects),  $\mathbf{x}$  and  $\mathbf{y}$ .

# Correlation measures the linear relationship between objects

$$\text{corr}(\mathbf{x}, \mathbf{y}) = \frac{\text{covariance}(\mathbf{x}, \mathbf{y})}{\text{standard\_deviation}(\mathbf{x}) * \text{standard\_deviation}(\mathbf{y})} = \frac{s_{xy}}{s_x s_y}, \quad (2.11)$$

where we are using the following standard statistical notation and definitions

$$\text{covariance}(\mathbf{x}, \mathbf{y}) = s_{xy} = \frac{1}{n - 1} \sum_{k=1}^n (x_k - \bar{x})(y_k - \bar{y}) \quad (2.12)$$

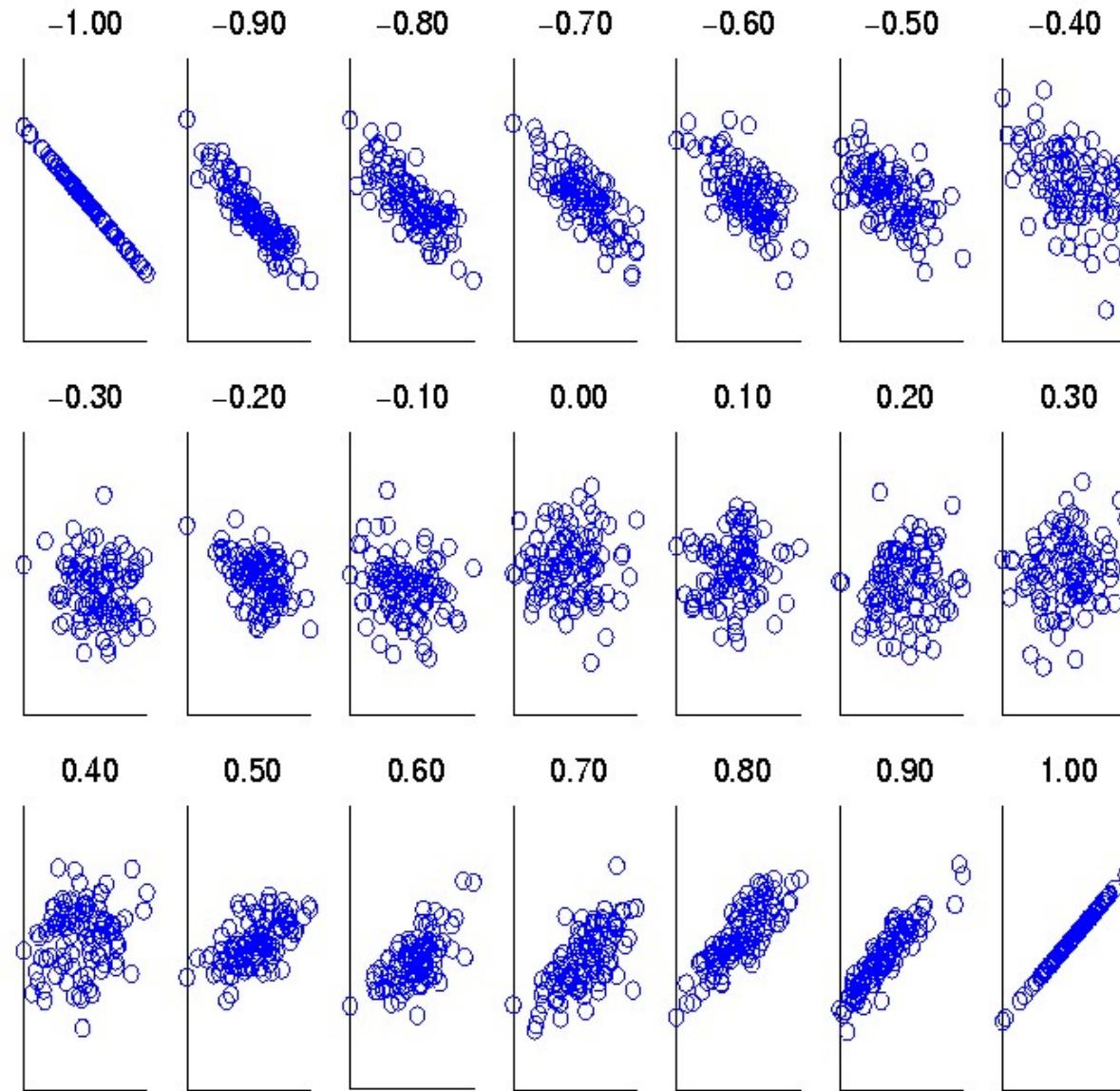
$$\text{standard\_deviation}(\mathbf{x}) = s_x = \sqrt{\frac{1}{n - 1} \sum_{k=1}^n (x_k - \bar{x})^2}$$

$$\text{standard\_deviation}(\mathbf{y}) = s_y = \sqrt{\frac{1}{n - 1} \sum_{k=1}^n (y_k - \bar{y})^2}$$

$$\bar{x} = \frac{1}{n} \sum_{k=1}^n x_k \text{ is the mean of } \mathbf{x}$$

$$\bar{y} = \frac{1}{n} \sum_{k=1}^n y_k \text{ is the mean of } \mathbf{y}$$

# Visually Evaluating Correlation



**Scatter plots  
showing the  
similarity from  
-1 to 1.**

# Comparison of Proximity Measures

- Domain of application
  - Similarity measures tend to be specific to the type of attribute and data
  - Record data, images, graphs, sequences, 3D-protein structure, etc. tend to have different measures
- However, one can talk about various properties that you would like a proximity measure to have
  - Symmetry is a common one
  - Tolerance to noise and outliers is another
  - Ability to find more types of patterns?
  - Many others possible
- The measure must be applicable to the data and produce results that agree with domain knowledge

# Data Preprocessing

- Aggregation
- Sampling
- Discretization and Binarization
- Attribute Transformation
- Dimensionality Reduction
- Feature subset selection
- Feature creation

# Aggregation

- Combining two or more attributes (or objects) into a single attribute (or object)
- Purpose
  - Data reduction - reduce the number of attributes or objects
  - Change of scale
    - Cities aggregated into regions, states, countries, etc.
    - Days aggregated into weeks, months, or years
  - More “stable” data - aggregated data tends to have less variability

**Table 2.4.** Data set containing information about customer purchases.

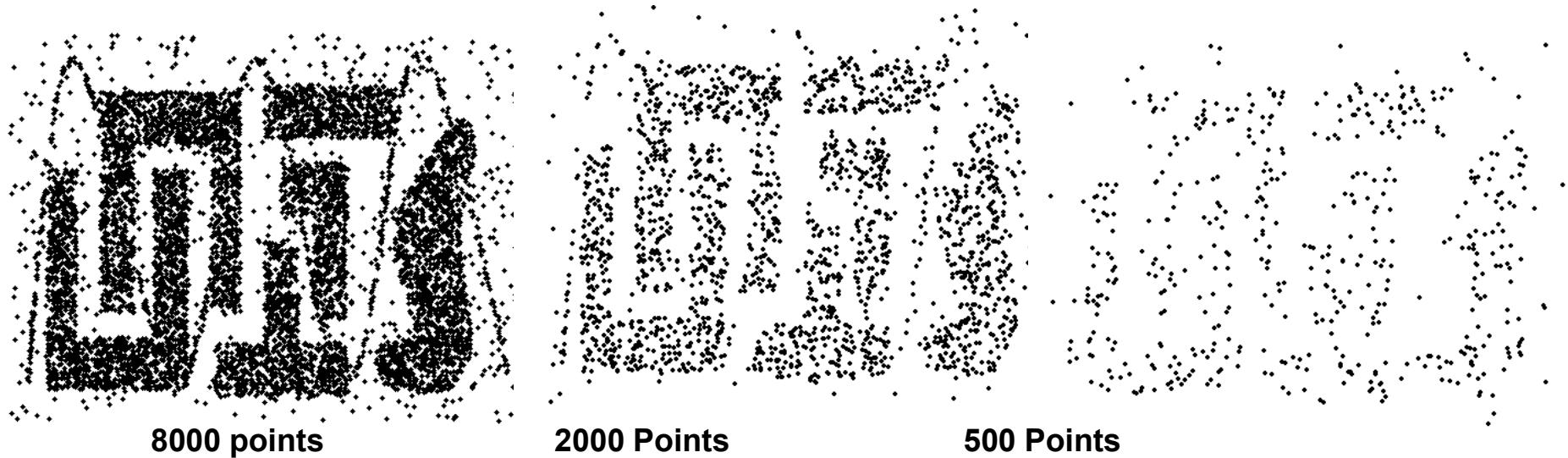
Transaction ID	Item	Store Location	Date	Price	...
:	:	:	:	:	
101123	Watch	Chicago	09/06/04	\$25.99	...
101123	Battery	Chicago	09/06/04	\$5.99	...
101124	Shoes	Minneapolis	09/06/04	\$75.00	...
:	:	:	:	:	

# Sampling

- Sampling is the main technique employed for data reduction
  - It is often used for both the preliminary investigation of the data and the final data analysis.
- Statisticians often sample because **obtaining** the entire set of data of interest is too expensive or time consuming
- Sampling is typically used in data mining because **processing** the entire set of data of interest is too expensive or time consuming

# Sampling

- Key principle for effective sampling :
  - Using a sample will work almost as well as using the entire data set, if the sample is **representative**
  - A sample is **representative** if it has approximately the same properties (of interest) as the original set of data



- Simple Random Sampling vs Stratified Sampling

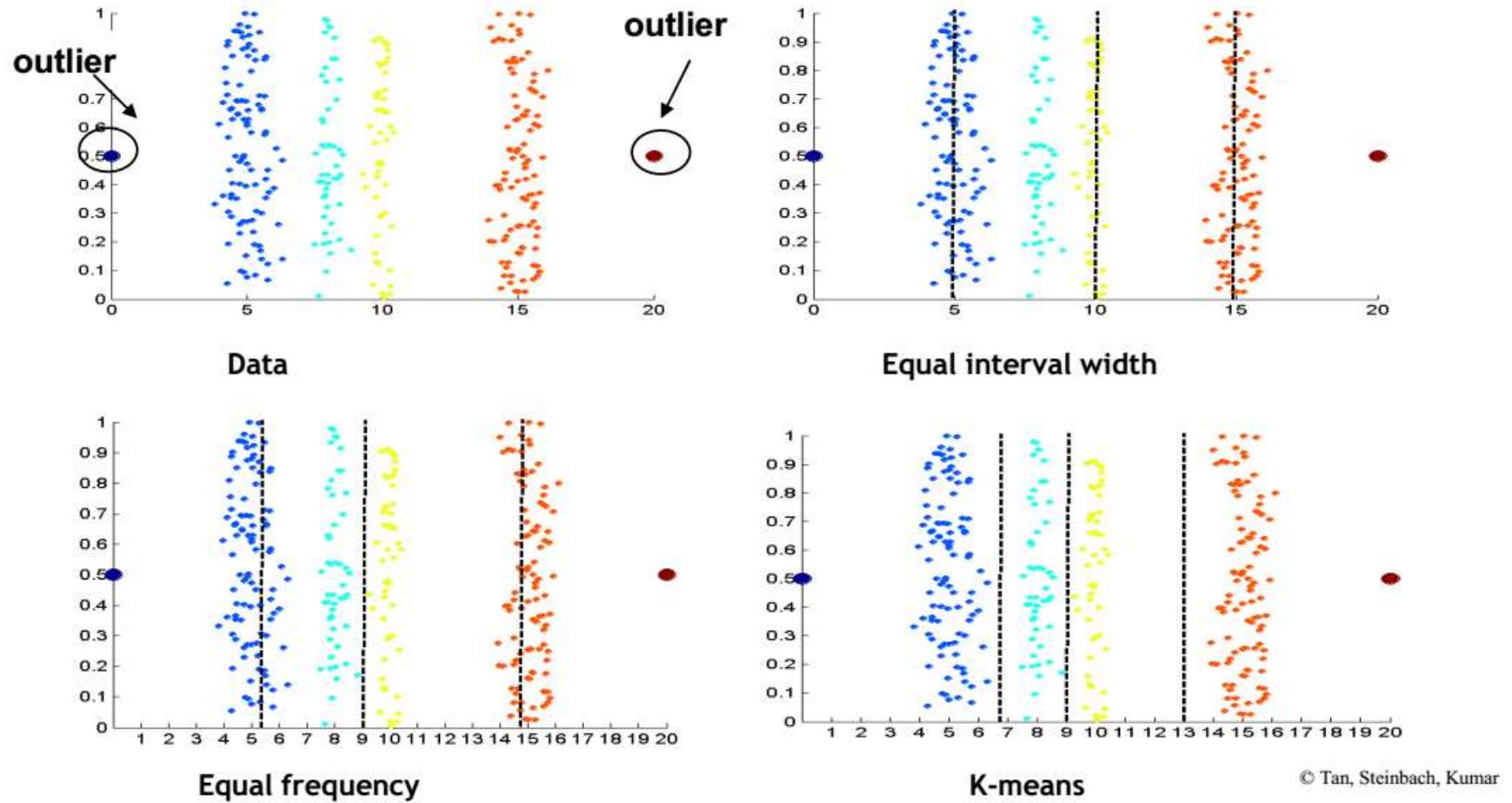
# Types of Sampling

- Simple Random Sampling
  - There is an equal probability of selecting any particular item
  - Sampling without replacement
    - As each item is selected, it is removed from the population
  - Sampling with replacement
    - Objects are not removed from the population as they are selected for the sample.
    - In sampling with replacement, the same object can be picked up more than once
- Stratified sampling
  - Split the data into several partitions; then draw random samples from each partition

# Discretization

- Discretization is the process of converting a continuous attribute into an ordinal attribute
  - A potentially infinite number of values are mapped into a small number of categories
  - Discretization is used in both unsupervised and supervised settings

# Unsupervised Discretization

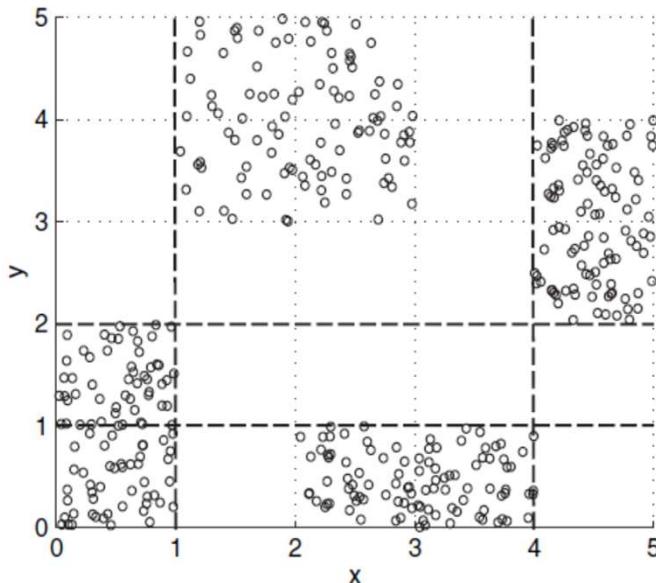


**Discretization to obtain 4 values**

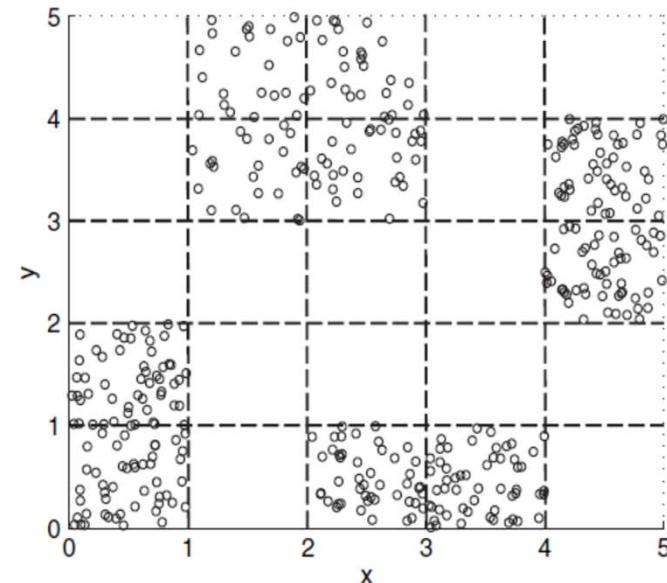
**Data consists of four groups of points and two outliers. Data is one-dimensional, but a random y component is added to reduce overlap.**

# Discretization in Supervised Settings

- Many classification algorithms work best if both the independent and dependent variables have only a few values
- We give an illustration of the usefulness of discretization using the following example.



(a) Three intervals



(b) Five intervals

**Figure 2.14.** Discretizing  $x$  and  $y$  attributes for four groups (classes) of points.

# Binarization

- Binarization maps a continuous or categorical attribute into one or more binary variables

**Table 2.6.** Conversion of a categorical attribute to five asymmetric binary attributes.

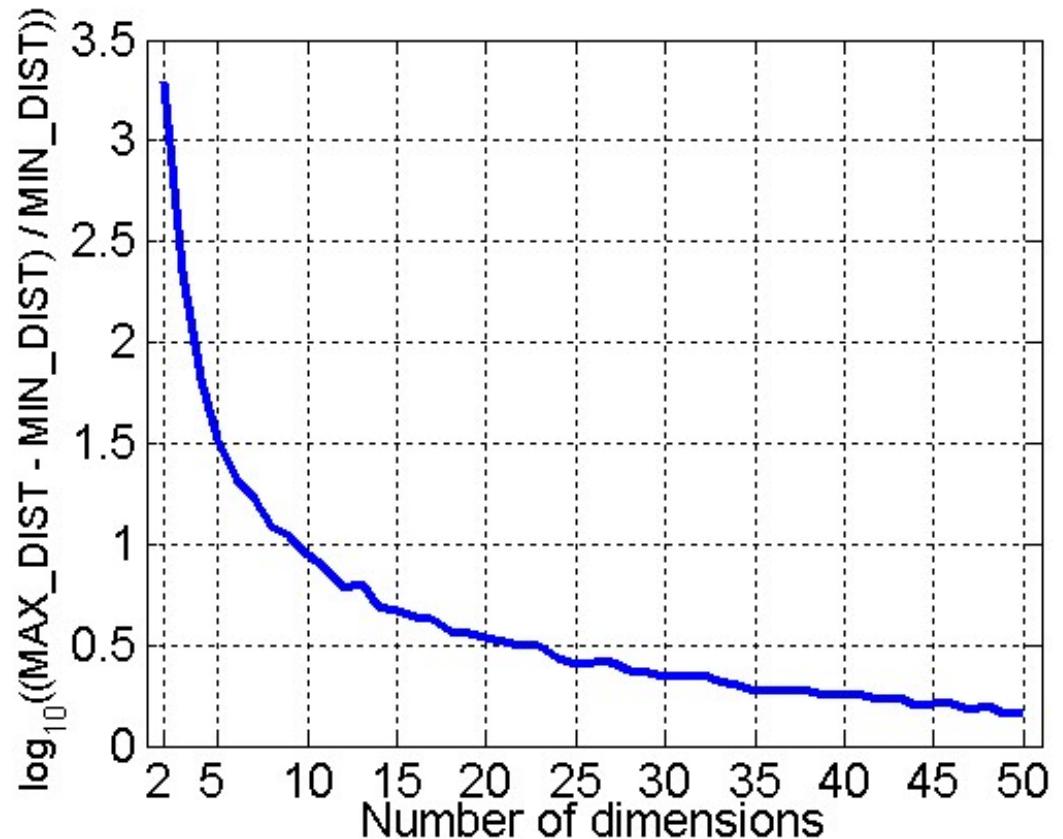
Categorical Value	Integer Value	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$
<i>awful</i>	0	1	0	0	0	0
<i>poor</i>	1	0	1	0	0	0
<i>OK</i>	2	0	0	1	0	0
<i>good</i>	3	0	0	0	1	0
<i>great</i>	4	0	0	0	0	1

# Attribute Transformation

- An **attribute transform** is a function that maps the entire set of values of a given attribute to a new set of replacement values such that each old value can be identified with one of the new values
  - Simple functions:  $x^k$ ,  $\log(x)$ ,  $e^x$ ,  $|x|$
  - **Normalization**
    - Refers to various techniques to adjust to differences among attributes in terms of frequency of occurrence, mean, variance, range
    - Take out unwanted, common signal, e.g., seasonality
  - In statistics, **standardization** refers to subtracting off the means and dividing by the standard deviation

# Curse of Dimensionality

- When dimensionality increases, data becomes increasingly sparse in the space that it occupies
- Definitions of density and distance between points, which are critical for clustering and outlier detection, become less meaningful



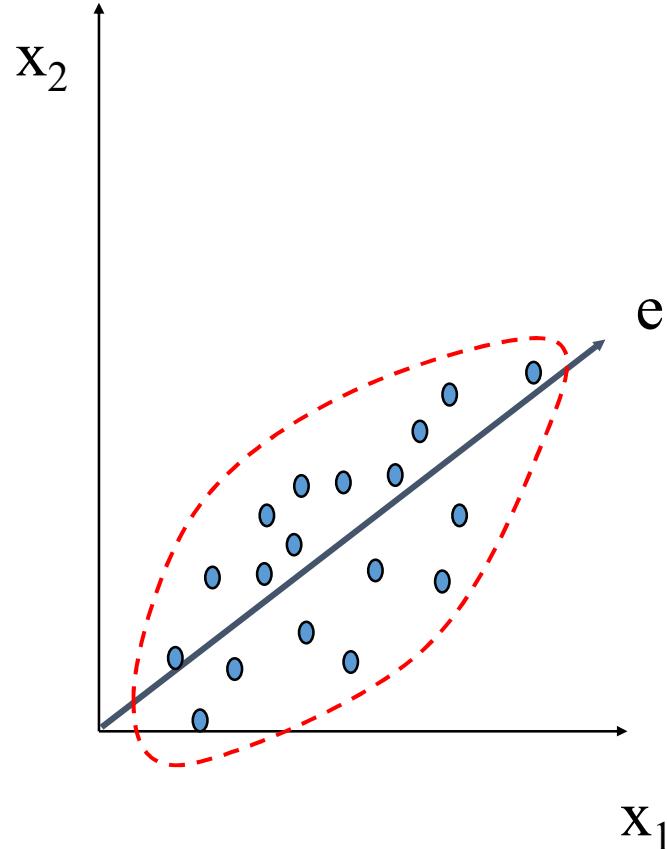
- Randomly generate 500 points
- Compute difference between max and min distance between any pair of points

# Dimensionality Reduction

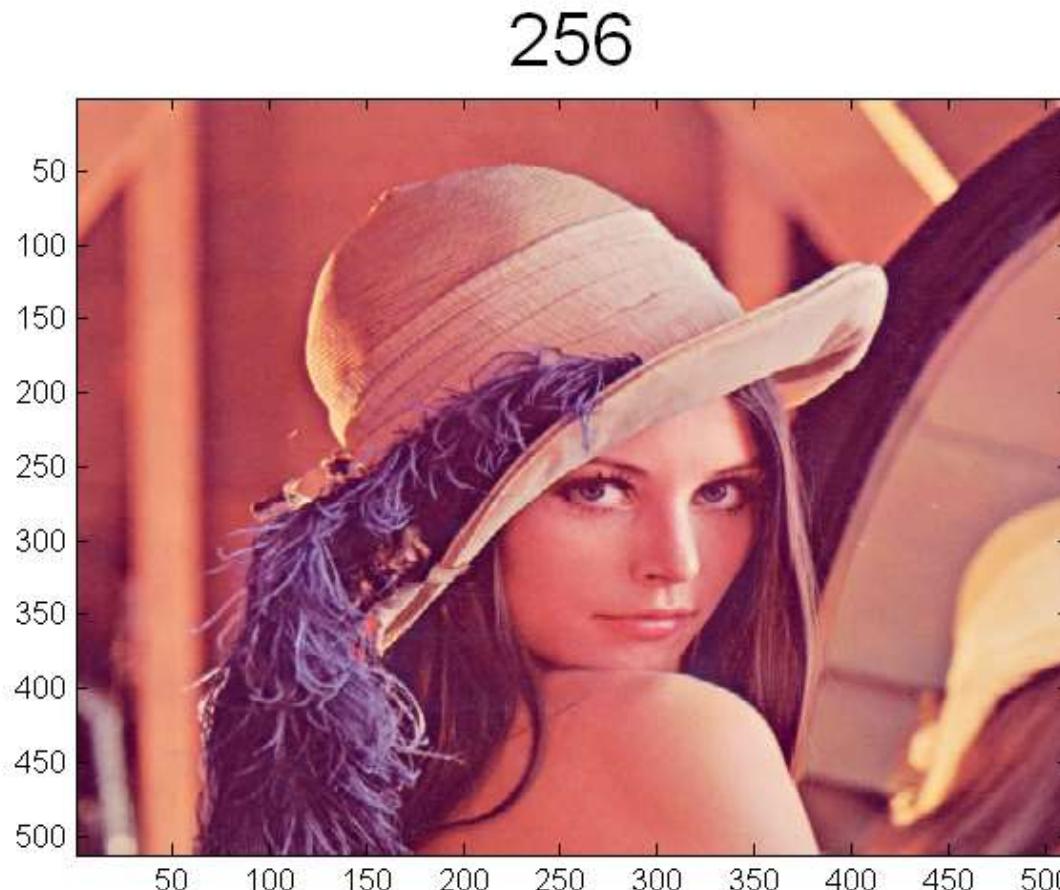
- Purpose:
  - Avoid curse of dimensionality
  - Reduce amount of time and memory required by data mining algorithms
  - Allow data to be more easily visualized
  - May help to eliminate irrelevant features or reduce noise
- Techniques
  - Principal Components Analysis (PCA)
  - Singular Value Decomposition
  - Others: supervised and non-linear techniques

# Dimensionality Reduction: PCA

- Goal is to find a projection that captures the largest amount of variation in data



# Dimensionality Reduction: PCA



# Feature Subset Selection

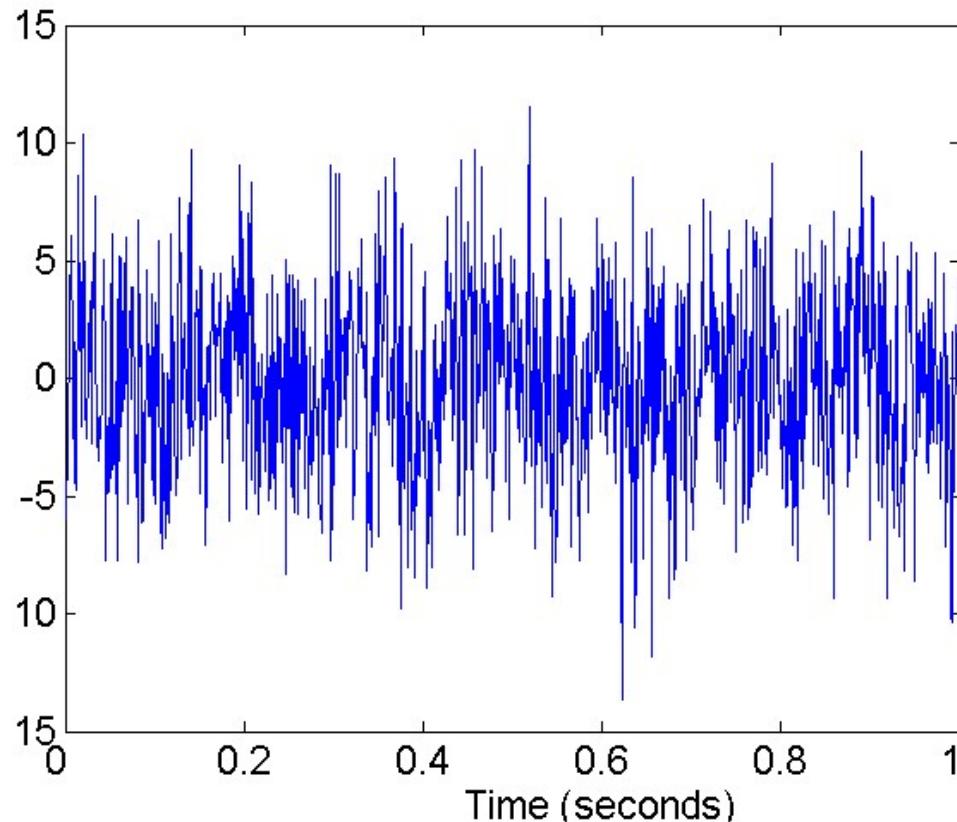
- Another way to reduce dimensionality of data
- Redundant features
  - Duplicate much or all of the information contained in one or more other attributes
  - Example: purchase price of a product and the amount of sales tax paid
- Irrelevant features
  - Contain no information that is useful for the data mining task at hand
  - Example: students' ID is often irrelevant to the task of predicting students' GPA
- Many techniques developed, especially for classification

# Feature Creation

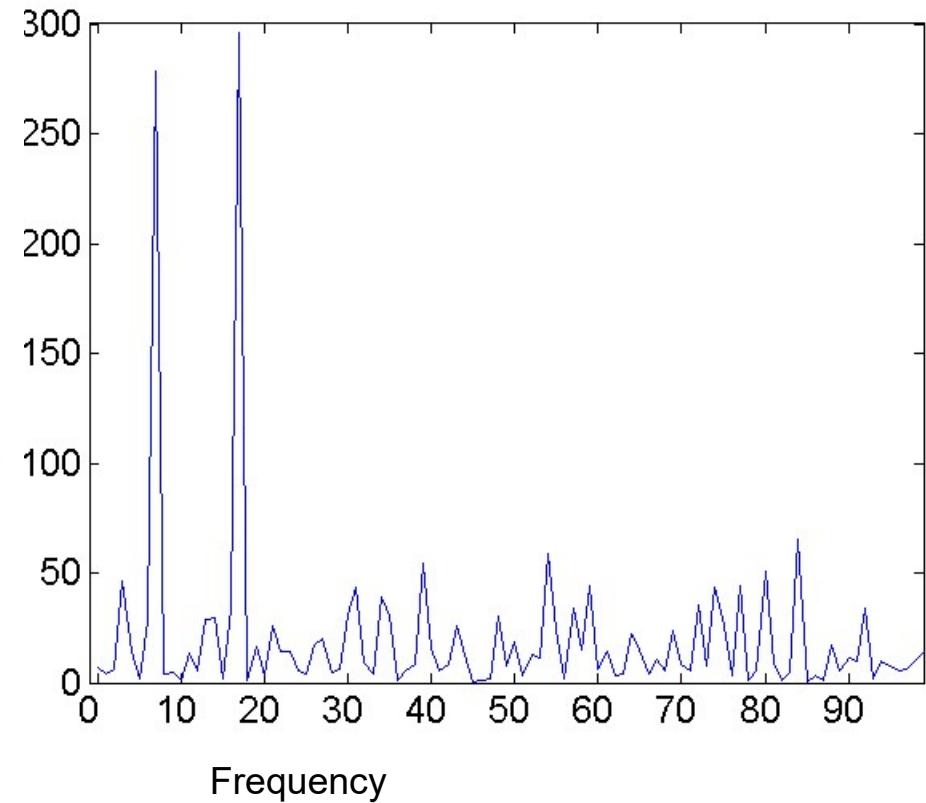
- Create new attributes that can capture the important information in a data set much more efficiently than the original attributes
- Three general methodologies:
  - Feature extraction
    - Example: extracting edges from images
  - Feature construction
    - Example: dividing mass by volume to get density
  - Mapping data to new space
    - Example: Fourier and wavelet analysis

# Mapping Data to a New Space

## □ Fourier and wavelet transform



**Two Sine Waves + Noise**



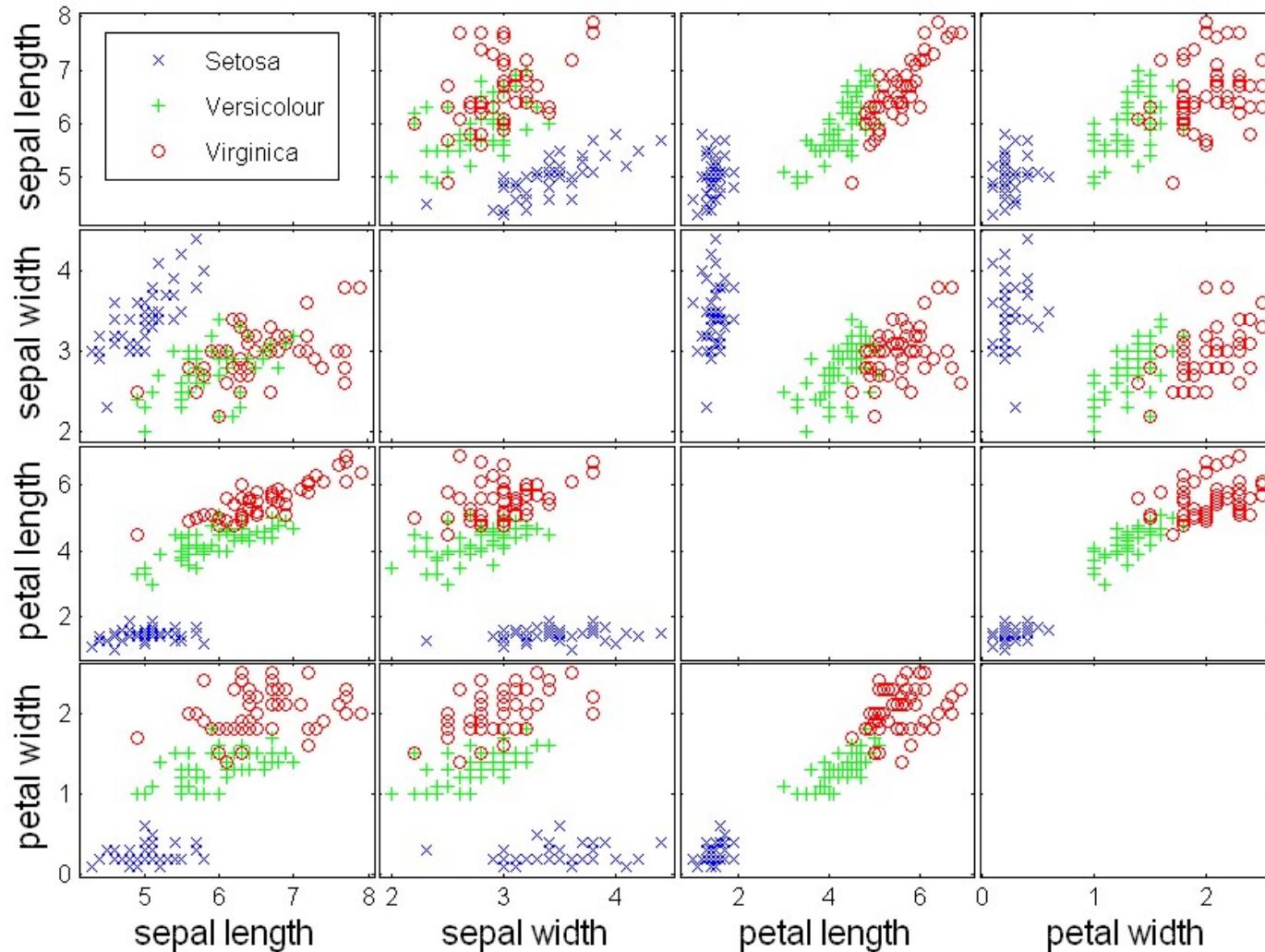
**Frequency**

# Data Exploration

**A preliminary exploration of the data to better understand its characteristics**

- Key motivations of data exploration include
  - Helping to select the right tool for preprocessing or analysis
  - Making use of humans' abilities to recognize patterns
    - People can recognize patterns not captured by data analysis tools
- Related to the area of Exploratory Data Analysis (EDA)
  - Created by statistician John Tukey - book "Exploratory Data Analysis"
- In EDA, as originally defined by Tukey
  - Focus on visualization
  - Clustering and anomaly detection were viewed as exploratory techniques
- Modern Data exploration:
  - Summary statistics
  - Visualization
  - Online Analytical Processing (OLAP)

# Example of Scatter Plot of Iris Attributes



# Conclusions/Summary

- Large scale data
- Commercial and Scientific interest and vast amount of opportunities
- Classification, Regression
- Knowledge Extraction Methodologies – CRISP-DM
- Data and Data Quality
- Data Pre-Processing (Aggregation, Sampling, Discretization and Binarization, Attribute Transformation, Dimensionality Reduction, Feature subset selection, Feature creation)
- Data Exploration
- Model Creation is the Next Step

# Bibliografia

- Tan, P., Steinbach, M., Karpatne, A. & Kumar, V. (2019). Introduction to Data Mining, 2<sup>nd</sup> Ed Pearson Addison-Wesley.
- Tan, P., Steinbach, M., Karpatne, A. & Kumar, V. (2020) Introduction to Data Mining Slides, <http://www-users.cs.umn.edu/~kumar/dmbook/index.php>
- Brígida Mónica Faria, (2018), Slides de Extração de Conhecimento, Instituto Politécnico do Porto, 2018
- Gladys Castillo (2008), Slides de Aprendizagem Computacional (Machine Learning), Universidade de Aveiro, 2008
- RapidMiner: Data Science Platform, (2017), <https://rapidminer.com/>
- Towards Data Science, (2020), <https://towardsdatascience.com/>
- DataCamp, <https://campus.datacamp.com/courses/data-science-for-everyone/>
- Pandas Website, <https://pandas.pydata.org/>
- Scikit-Learn Website, <https://scikit-learn.org/>
- Matplotlib Website, <https://matplotlib.org>
- Seaborn Website, <https://seaborn.pydata.org/>

# **Artificial Intelligence/ Inteligência Artificial**

## **Lecture 5c: Machine Learning – Data Preprocessing** (adapted from Tan et al, 2020)

**Luís Paulo Reis**

[lpreas@fe.up.pt](mailto:lpreas@fe.up.pt)

Director of LIACC – Artificial Intelligence and Computer Science Lab.  
Associate Professor at DEI/FEUP – Informatics Engineering Department,  
Faculty of Engineering of the University of Porto, Portugal  
President of APPIA – Portuguese Association for Artificial Intelligence

