

# Data Science - Lecture 4

## Introduction To Data Science

Dr. Faisal Kamiran

# What is today's agenda?

Today we are going to learn following things :

- Data Understanding
- Data Preprocessing

# What is 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, or feature
- A collection of attributes describe an object
  - Object is also known as record, point, case, sample, entity, or instance

**Objects**

**Attributes**

<i>Tid</i>	<i>Refund</i>	<i>Marital Status</i>	<i>Taxable Income</i>	<i>Cheat</i>
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
- 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 values can be different
      - ID has no limit but age has a maximum and minimum value

# 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
  - 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 with limited precision.
  - Continuous attributes are typically represented as floating-point variables.

# Type of Attributes

- Categorical (Qualitative)
  - 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 in {tall, medium, short}
- Numeric (Quantitative)
  - Interval
    - ◆ Examples: temperatures in Celsius or Fahrenheit.
  - Ratio
    - ◆ Examples: temperature in Kelvin, length, time, counts

# Type of Attributes

OK to compute....	Nominal	Ordinal	Interval	Ratio
Mode, Entropy	Yes	Yes	Yes	Yes
median and percentiles.	No	Yes	Yes	Yes
add or subtract.	No	No	Yes	Yes
mean, standard deviation, standard error of the mean.	No	No	Yes	Yes
ratio, or coefficient of variation.	No	No	No	Yes

# 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 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.

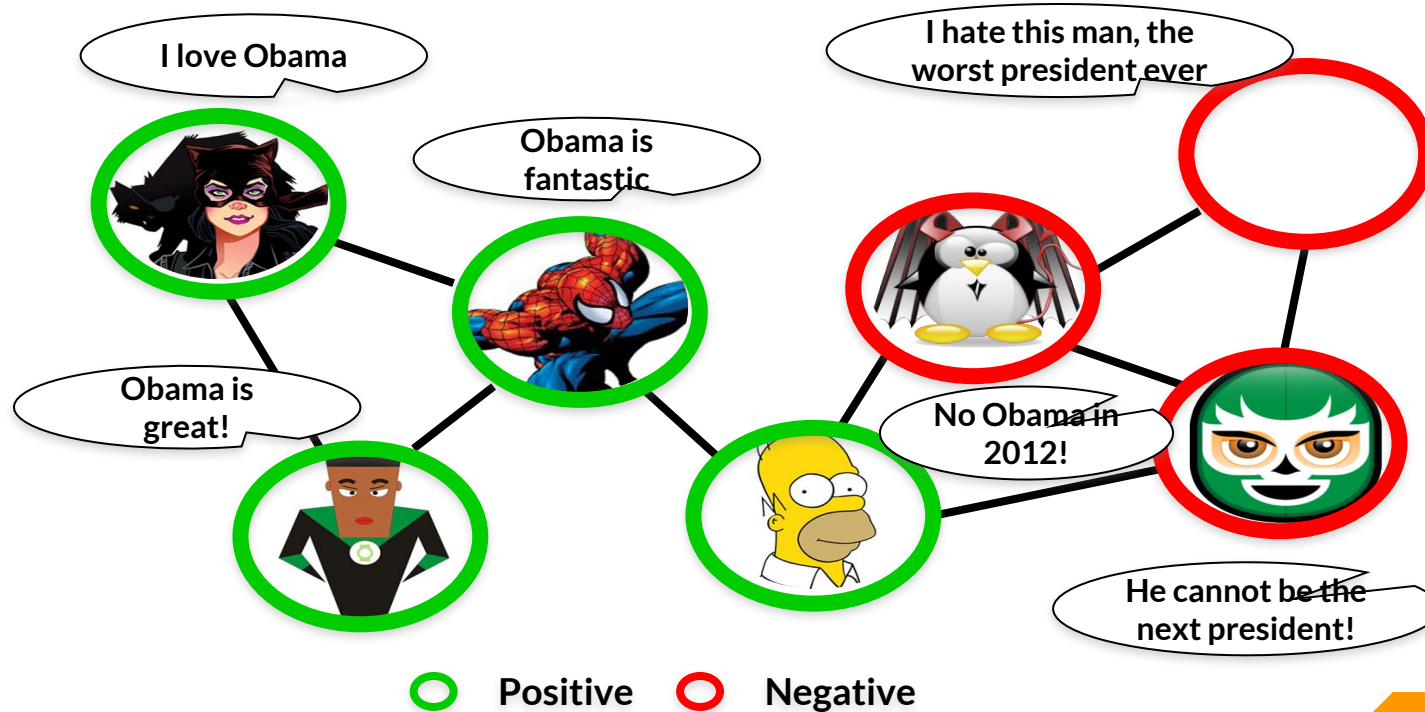
	team	coach	play	ball	score	game	win	lost	timeout	season
Document 1	3	0	5	0	2	6	0	2	0	2
Document 2	0	7	0	2	1	0	0	3	0	0
Document 3	0	1	0	0	1	2	2	0	3	0

# Transaction Data

- A special type of record data, where
  - each record (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.

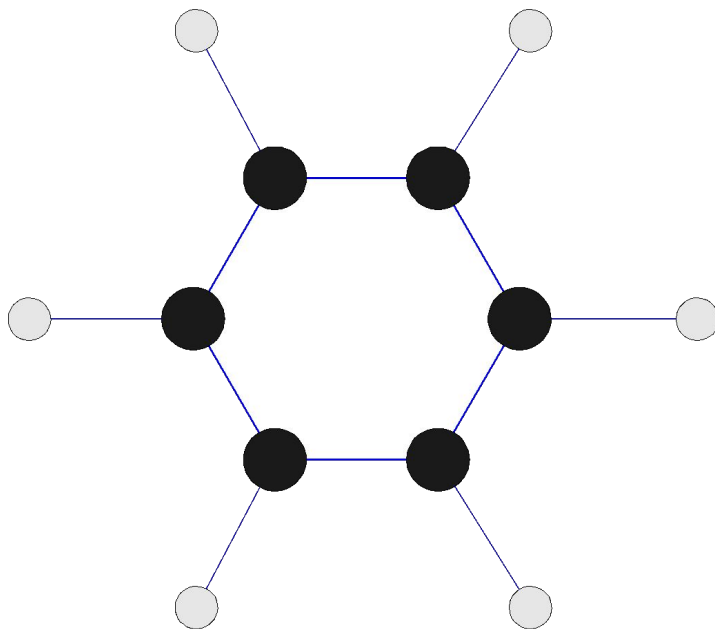
<i><b>TID</b></i>	<i><b>Items</b></i>
<b>1</b>	<b>Bread, Coke, Milk</b>
<b>2</b>	<b>Beer, Bread</b>
<b>3</b>	<b>Beer, Coke, Diaper, Milk</b>
<b>4</b>	<b>Beer, Bread, Diaper, Milk</b>
<b>5</b>	<b>Coke, Diaper, Milk</b>

# Graph Data : “Love Obama”



# Chemical Data

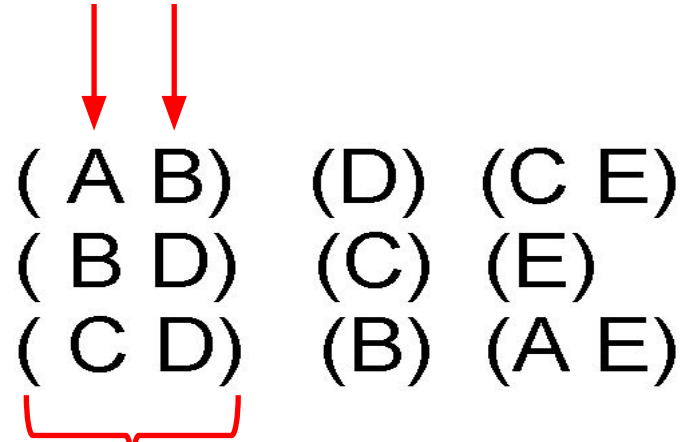
- Benzene Molecule:  $C_6H_6$



# Ordered Data

- Sequences of transactions

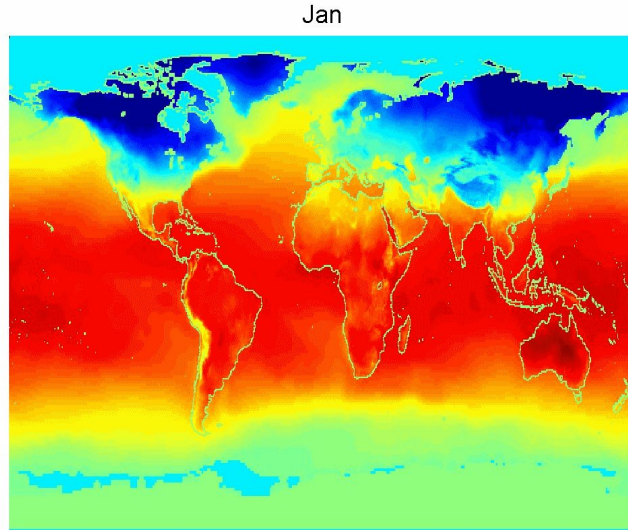
Items/Events



An element of the  
sequence

# Ordered Data

- Spatio-Temporal Data



**Average Monthly Temperature of Land and Ocean**

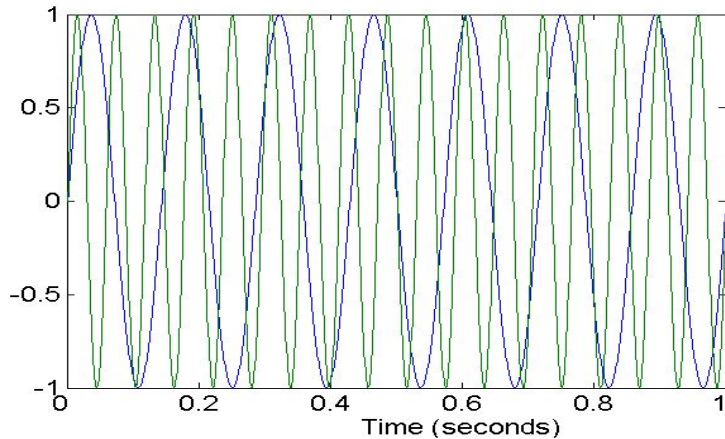


# Data Quality

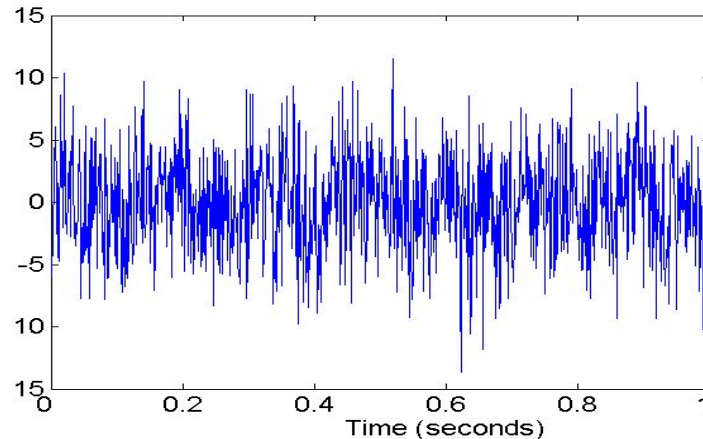
- What kinds of data quality problems?
  - How can we detect problems with the data?
  - What can we do about these problems?
- 
- Examples of data quality problems:
    - Noise and outliers
    - missing values
    - duplicate data

# Noise

- 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



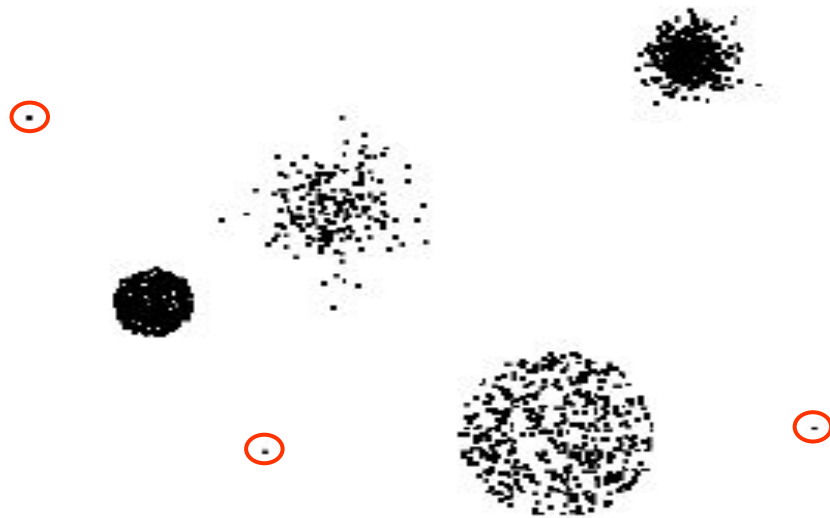
**Two Sine Waves**



**Two Sine Waves + Noise**

# Outliers

- Outliers are data objects with characteristics that are considerably different than most of the other data objects in the data set



# 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
  - Estimate Missing Values
  - Ignore the Missing Value During Analysis
  - Replace with all possible values (weighted by their probabilities)

# 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

# Data Preprocessing

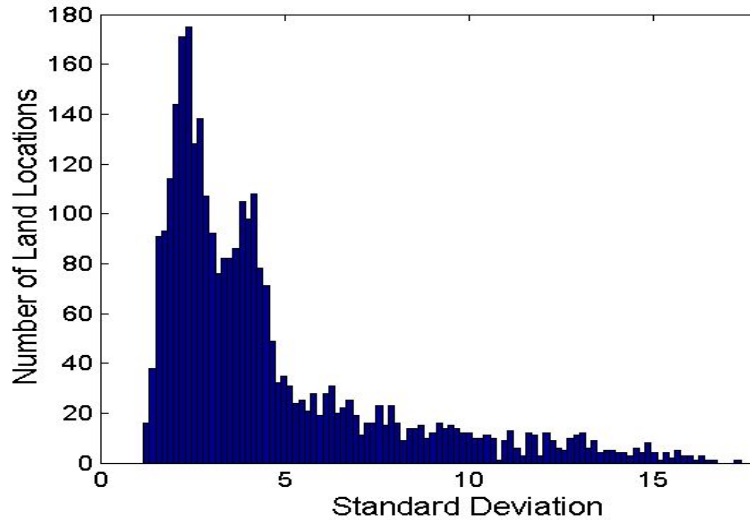
- Aggregation
- Sampling
- Dimensionality Reduction
- Feature subset selection
- Discretization and Binarization
- Attribute Transformation

# Aggregation

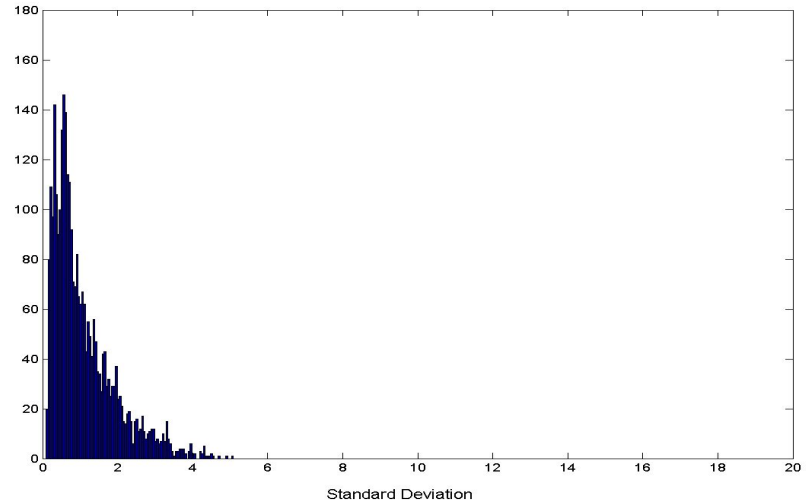
- Combining two or more attributes (or objects) into a single attribute (or object)
- Purpose
  - Data reduction
    - ◆ Reduce the number of attributes or objects
  - Change of scale
    - ◆ Cities aggregated into regions, states, countries, etc
  - More “stable” data
    - ◆ Aggregated data tends to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

# Aggregation

## Variation of Precipitation in Australia



Standard Deviation of Average  
Monthly Precipitation



Standard Deviation of Average  
Yearly Precipitation



# Sampling

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

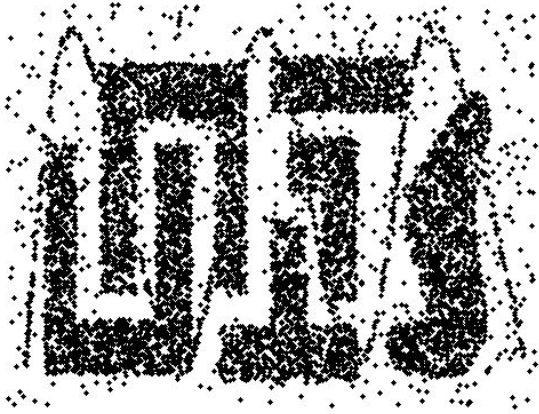
# Sampling

- The key principle for effective sampling is following:
  - using a sample will work almost as well as using the entire data sets, if the sample is representative.
  - A sample is representative if it has approximately the same property (of interest) as the original set of data .

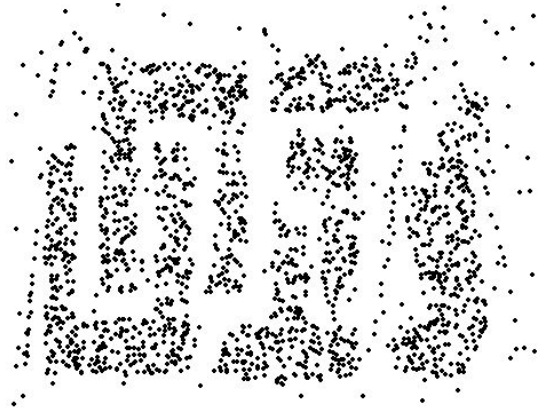
# Types of Sampling

- 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

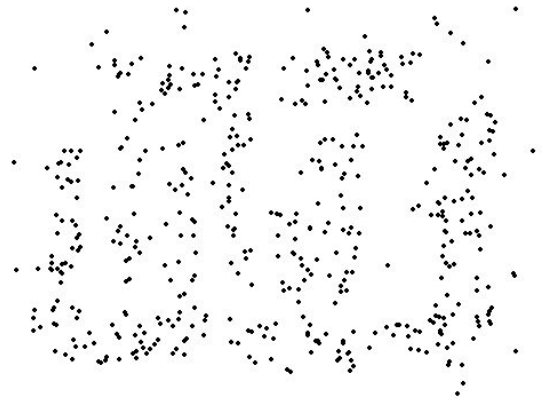
# Sample Size



**8000 points**



**2000 Points**



**500 Points**

# Curse of Dimensionality

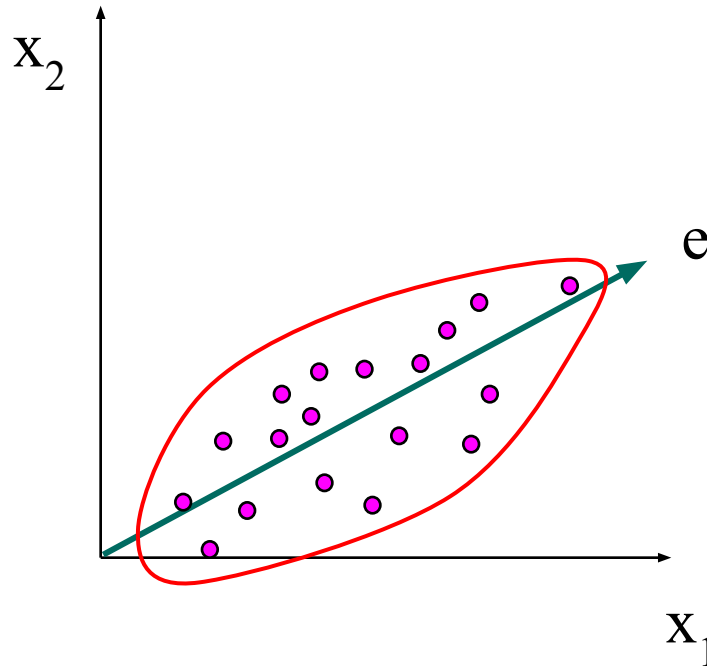
- When dimensionality increases, data becomes increasingly sparse in the space that it occupies
- Definitions of density and distance between points, which is critical for clustering and outlier detection, become less meaningful

# 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 Component Analysis
  - 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



# 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



# Feature Subset Selection

- Techniques:
  - Brute-force approach:
    - ◆ Try all possible feature subsets as input to data mining algorithm
  - Embedded approaches:
    - ◆ Feature selection occurs naturally as part of the data mining algorithm
  - Filter approaches:
    - ◆ Features are selected before data mining algorithm is run

“

*Questions ?*