# Data Mining & Data Warehouse IS-372

**1**

**Dr. Mohammed Al-Sarem**

**Taibah University**

**IS Department**

Data Mining: Concepts and Techniques

Dr. Mohammed Al-Sarem

# CHAPTER 2

Data in Data Mining

# Outlines

- What is Data in Data Mining?

- Type of Attributes, Attributes Values

- Data Set in Data Mining

- Data Quality: Noisy Data, Outliers, Missing Values, Duplicated Data.

- Data pre-processing: Introduction
  - Aggregation
  - Sampling
  - Dimensionality Reduction

- Similarity & Dissimilarity

# Data in Data Mining

- Data: a ==collection of facts usually obtained as the result of experiences,== ==observations,== or experiments

- Data ==may consist of numbers,== ==words, images, …==

- Data: ==lowest level of abstraction== (from which information and knowledge are derived)

# 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

**Attributes**

**Objects**

| 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

- 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

# Types of Attributes

**Nominal**

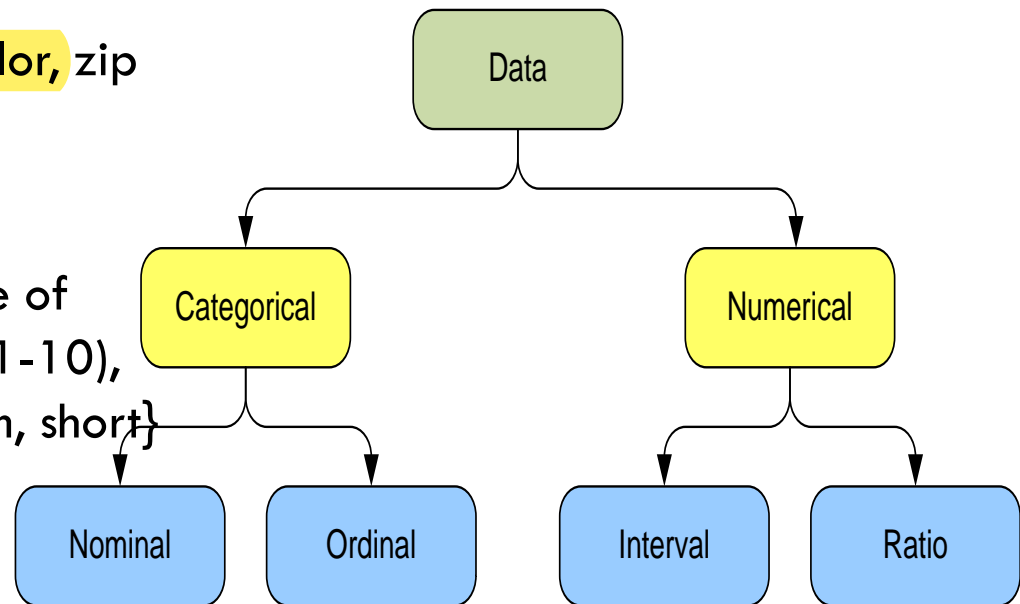- Examples: ID numbers, eye color, zip codes

**Ordinal**

- Examples: rankings (e.g., taste of potato chips on a scale from 1-10), grades, height in {tall, medium, short}

**Interval**

- Examples: calendar dates, temperatures in Celsius or Fahrenheit.

**Ratio**

- Examples: temperature in Kelvin, length, time, counts

```
                    Data
                     |
        -------------------------------
        |                             |
   Categorical                    Numerical
        |                             |
   ----------                   ----------
   |        |                   |        |
Nominal  Ordinal            Interval   Ratio
```

- DM with different data types?

- Other data types?

# Properties of Attribute Values

□ The type of an attribute depends on which of the following properties it possesses:

- Distinctness: $=$ $\neq$
- Order: $<$ $>$
- Addition: $+$ $-$
- Multiplication: $*$ $/$

- Nominal attribute: distinctness
- Ordinal attribute: distinctness & order
- Interval attribute: distinctness, order & addition
- Ratio attribute: all 4 properties

# 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.

# In class exercise #2

- Classify the following attributes as binary, discrete, or continuous. Also classify them as qualitative (nominal or ordinal) or quantitative (interval or ratio). Some cases may have more than one interpretation, so briefly indicate your reasoning if you think there may be some ambiguity.

  - Number of telephones in your house
  - Size of French Fries (Medium or Large or X-Large)
  - Ownership of a cell phone
  - Number of local phone calls you made in a month
  - Length of longest phone call
  - Length of your foot
  - Price of your textbook
  - Zip code
  - Temperature in degrees Fahrenheit
  - Temperature in degrees Celsius
  - Temperature in kelvins

# 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

# Important Characteristics of Structured Data

- Dimensionality
  - Curse of Dimensionality

- Sparsity
  - Only presence counts

- Resolution
  - Patterns depend on the scale

# Record Data

تمـيـن

□ Data that consists of a collection of records, each of which consists of a fixed set of 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 |

# 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 |

# Transigan 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.

| TID | Items |
|-----|-------|
| 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 and HTML Links



<a href="papers/papers.html#bbbb">
Data Mining </a>
<li>
<a href="papers/papers.html#aaaa">
Graph Partitioning </a>
<li>
<a href="papers/papers.html#aaaa">
Parallel Solution of Sparse Linear System of Equations </a>
<li>
<a href="papers/papers.html#ffff">
N-Body Computation and Dense Linear System Solvers
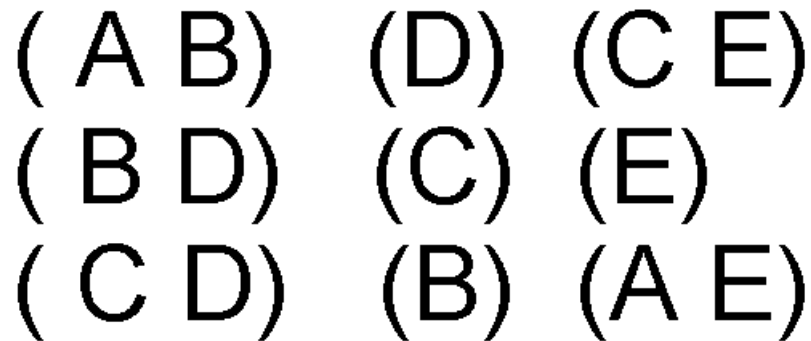
# Chemical Data

☐ **Benzene Molecule:** $C_6H_6$
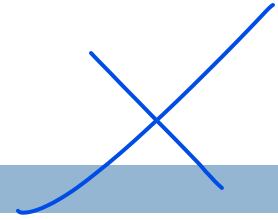
# Ordered Data

- **Sequences of transactions**

Items/Events

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

An element of
the sequence

# Ordered Data

☐ Genomic sequence data

GGTTCCGCCTTCAGCCCCGCGCC
CGCAGGGCCCGCCCCGCGCCGTC
GAGAAGGGCCCGCCTGGCGGGCG
GGGGGAGGCGGGGCCGCCCGAGC
CCAACCGAGTCCGACCAGGTGCC
CCCTCTGCTCGGCCTAGACCTGA
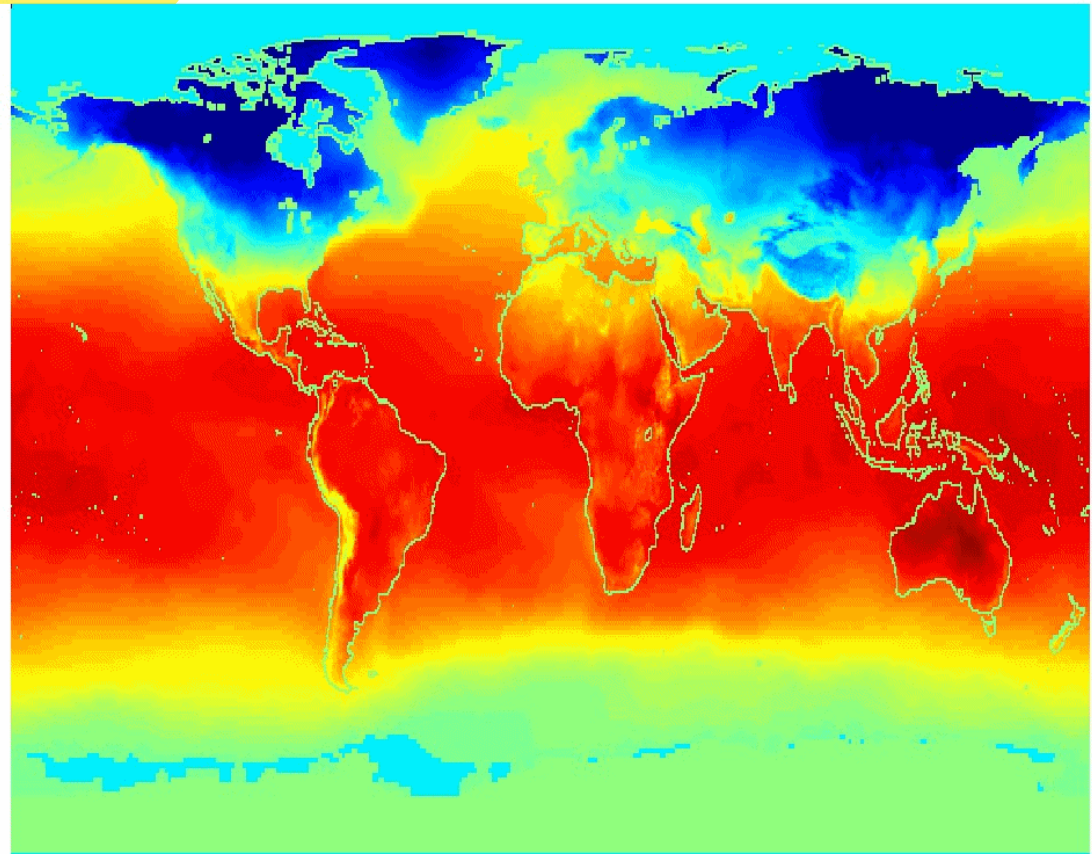GCTCATTAGGCGGCAGCGGACAG
GCCAAGTAGAACACGCGAAGCGC
TGGGCTGCCTGCTGCGACCAGGG

☐ **Spatio-Temporal Data**

Jan

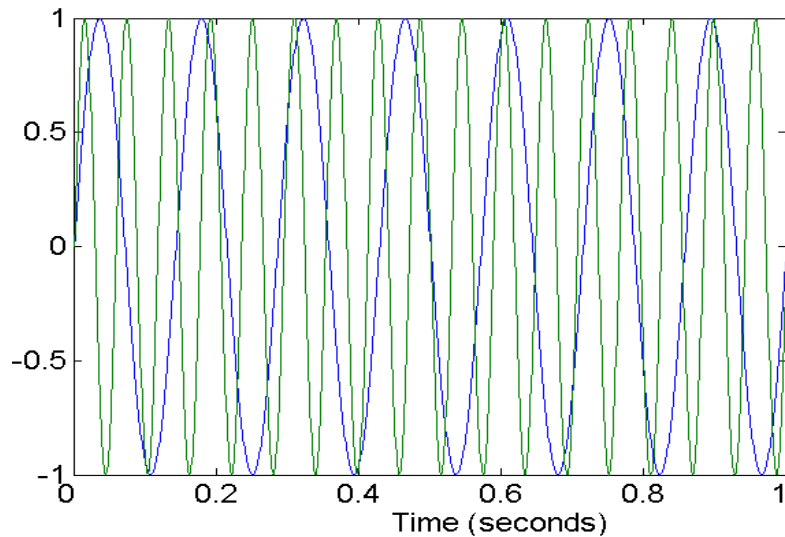Average Monthly
Temperature of
land and ocean

# 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:
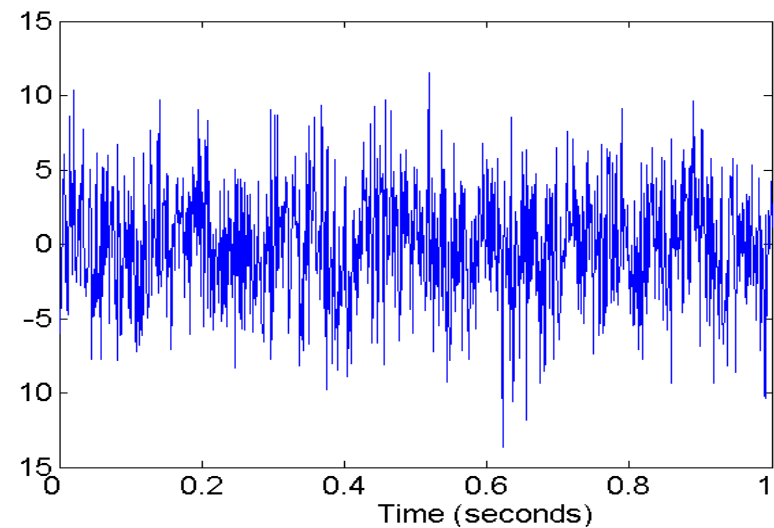  ▫ Noise and outliers
  ▫ missing values
  ▫ duplicate data

# 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
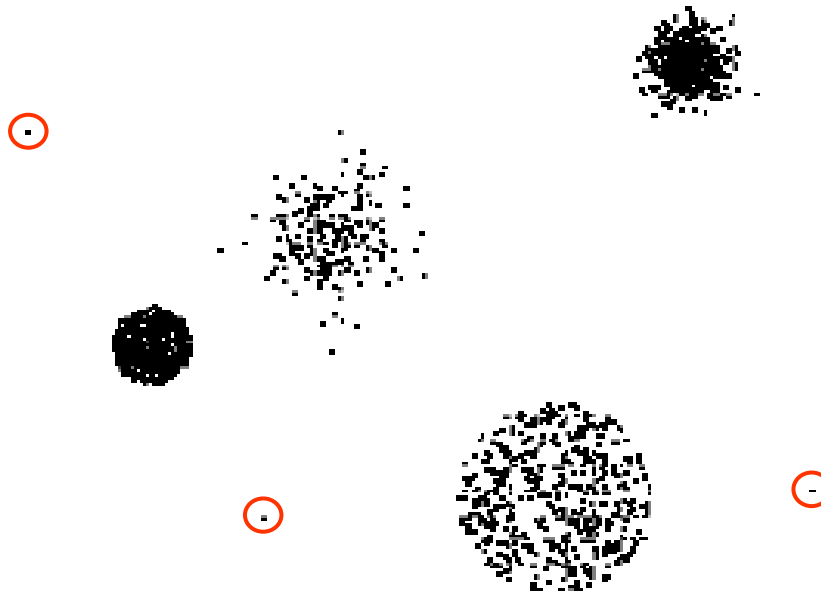
Two Sine Waves

Two Sine Waves + Noise

# Outliers

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

# 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 heterogeous sources

- Examples:
  - Same person with multiple email addresses

- Data cleaning
  - Process of dealing with duplicate data issues

# Similarity and Dissimilarity

- Similarity
  - Numerical measure of how alike two data objects are.
  - Is higher when objects are more alike.
  - Often falls in the range [0,1]
- Dissimilarity
  - Numerical measure of how different are two data objects
  - 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

$p$ and $q$ are the attribute values for two data objects.

| Attribute Type | Dissimilarity | Similarity |
|---|---|---|
| Nominal | $d = \begin{cases} 0 & \text{if } p = q \\ 1 & \text{if } p \neq q \end{cases}$ | $s = \begin{cases} 1 & \text{if } p = q \\ 0 & \text{if } p \neq q \end{cases}$ |
| Ordinal | $d = \frac{\lvert p-q \rvert}{n-1}$ (values mapped to integers $0$ to $n-1$, where $n$ is the number of values) | $s = 1 - \frac{\lvert p-q \rvert}{n-1}$ |
| Interval or Ratio | $d = \lvert p - q \rvert$ | $s = -d,\ s = \frac{1}{1+d}$ or $s = 1 - \frac{d - min\_d}{max\_d - min\_d}$ |

**Table 5.1.** Similarity and dissimilarity for simple attributes

# Euclidean Distance

□ Euclidean Distance

$$dist = \sqrt{\sum_{k=1}^{n}(p_k - q_k)^2}$$

Where $n$ is the number of dimensions (attributes) and $p_k$ and $q_k$ are, respectively, the k[th] attributes (components) or data objects $p$ and $q$.
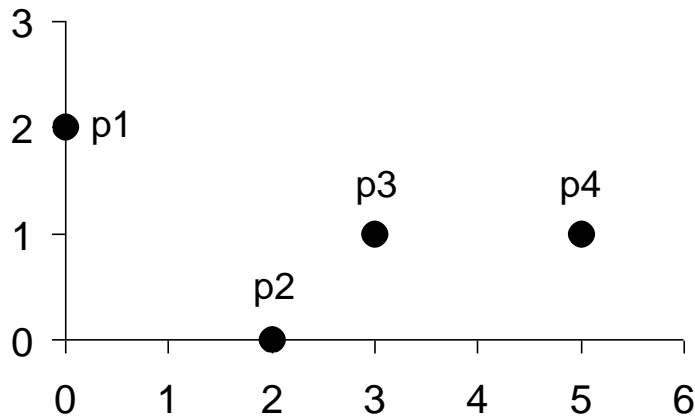
□ Standardization is necessary, if scales differ.

# Euclidean Distance

$$\sqrt{|0 - 2|^2 + |2 - 0|^2}$$
$$\sqrt{4 + 4} = 8$$
$$\sqrt{8}$$
$$\approx 2.82$$

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

|     | 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

# Minkowski Distance

☐ Minkowski Distance is a generalization of Euclidean Distance

$$dist = (\sum_{k=1}^{n} | p_k - q_k |^r)^{\frac{1}{r}}$$

Where $r$ is a parameter, $n$ is the number of dimensions (attributes) and $p_k$ and $q_k$ are, respectively, the kth attributes (components) or data objects $p$ and $q$.

# Minkowski Distance: Examples

- $r = 1$. City block (Manhattan, taxicab, $L_1$ norm) distance.

  - A common example of this 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

- Do not confuse $r$ with $n$, i.e., all these distances are defined for all numbers of dimensions.

# Minkowski Distance

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

| 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 |

| L2 | 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 |

| $L_\infty$ | p1 | p2 | p3 | p4 |
|----|----|----|----|----|
| p1 | 0 | 2 | 3 | 5 |
| p2 | 2 | 0 | 1 | 3 |
| p3 | 3 | 1 | 0 | 2 |
| p4 | 5 | 3 | 2 | 0 |

Distance Matrix

# Common Properties of a Distance

□ Distances, such as the Euclidean distance, have some well known properties.

1. $d(p, q) \geq 0$ for all $p$ and $q$ and $d(p, q) = 0$ only if $p = q$. (Positive definiteness)

2. $d(p, q) = d(q, p)$ for all $p$ and $q$. (Symmetry)

3. $d(p, r) \leq d(p, q) + d(q, r)$ for all points $p$, $q$, and $r$. (Triangle Inequality)

where $d(p, q)$ is the distance (dissimilarity) between points (data objects), $p$ and $q$.

□ A distance that satisfies these properties is a <span style="color:red">metric</span>

# Common Properties of a Similarity

□ Similarities, also have some well known properties.

1. $s(p, q) = 1$ (or maximum similarity) only if $p = q$.

2. $s(p, q) = s(q, p)$ for all $p$ and $q$. (Symmetry)

where $s(p, q)$ is the similarity between points (data objects), $p$ and $q$.

# Similarity Between Binary Vectors

- Common situation is that objects, *p* and *q*, have only binary attributes

- Compute similarities using the following quantities

  $M_{01}$ = the number of attributes where p was 0 and q was 1

  $M_{10}$ = the number of attributes where p was 1 and q was 0

  $M_{00}$ = the number of attributes where p was 0 and q was 0

  $M_{11}$ = the number of attributes where p was 1 and q was 1

- Simple Matching and Jaccard Coefficients

  SMC = number of matches / number of attributes

  $= (M_{11} + M_{00}) / (M_{01} + M_{10} + M_{11} + M_{00})$

  J = number of 11 matches / number of not-both-zero attributes values

  $= (M_{11}) / (M_{01} + M_{10} + M_{11})$

$p = $ 1 0 0 0 0 0 0 0 0 0

$q = $ 0 0 0 0 0 0 1 0 0 1

$M_{01}$ = 2  (the number of attributes where p was 0 and q was 1)

$M_{10}$ = 1  (the number of attributes where p was 1 and q was 0)

$M_{00}$ = 7  (the number of attributes where p was 0 and q was 0)

$M_{11}$ = 0  (the number of attributes where p was 1 and q was 1)

SMC = $(M_{11} + M_{00})/(M_{01} + M_{10} + M_{11} + M_{00})$ = (0+7) / (2+1+0+7) = 0.7

J = $(M_{11}) / (M_{01} + M_{10} + M_{11})$ = 0 / (2 + 1 + 0) = 0

# Cosine Similarity

□ If $d_1$ and $d_2$ are two document vectors, then

$$\cos(d_1, d_2) = (d_1 \bullet d_2) / ||d_1|| \, ||d_2|| ,$$

where $\bullet$ indicates vector dot product and $|| d ||$ is the length of vector $d$.

□ Example:

$$d_1 = 3\ 2\ 0\ 5\ 0\ 0\ 0\ 2\ 0\ 0$$
$$d_2 = 1\ 0\ 0\ 0\ 0\ 0\ 0\ 1\ 0\ 2$$

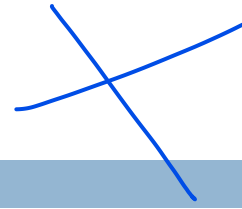$d_1 \bullet d_2 =$ 3*1 + 2*0 + 0*0 + 5*0 + 0*0 + 0*0 + 0*0 + 2*1 + 0*0 + 0*2 = 5

$||d_1|| = (3*3+2*2+0*0+5*5+0*0+0*0+0*0+2*2+0*0+0*0)^{0.5} = (42)^{0.5} = 6.481$

$||d_2|| = (1*1+0*0+0*0+0*0+0*0+0*0+0*0+1*1+0*0+2*2)^{0.5} = (6)^{0.5} = 2.245$

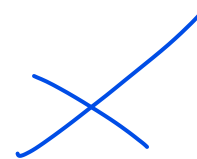$$\cos(d_1, d_2) = 5 / (6.481 * 2.245)$$

.3150

# Correlation

- Correlation measures the linear relationship between objects

- To compute correlation, we standardize data objects, p and q, and then take their dot product

$$p'_k = (p_k - mean(p)) / std(p)$$

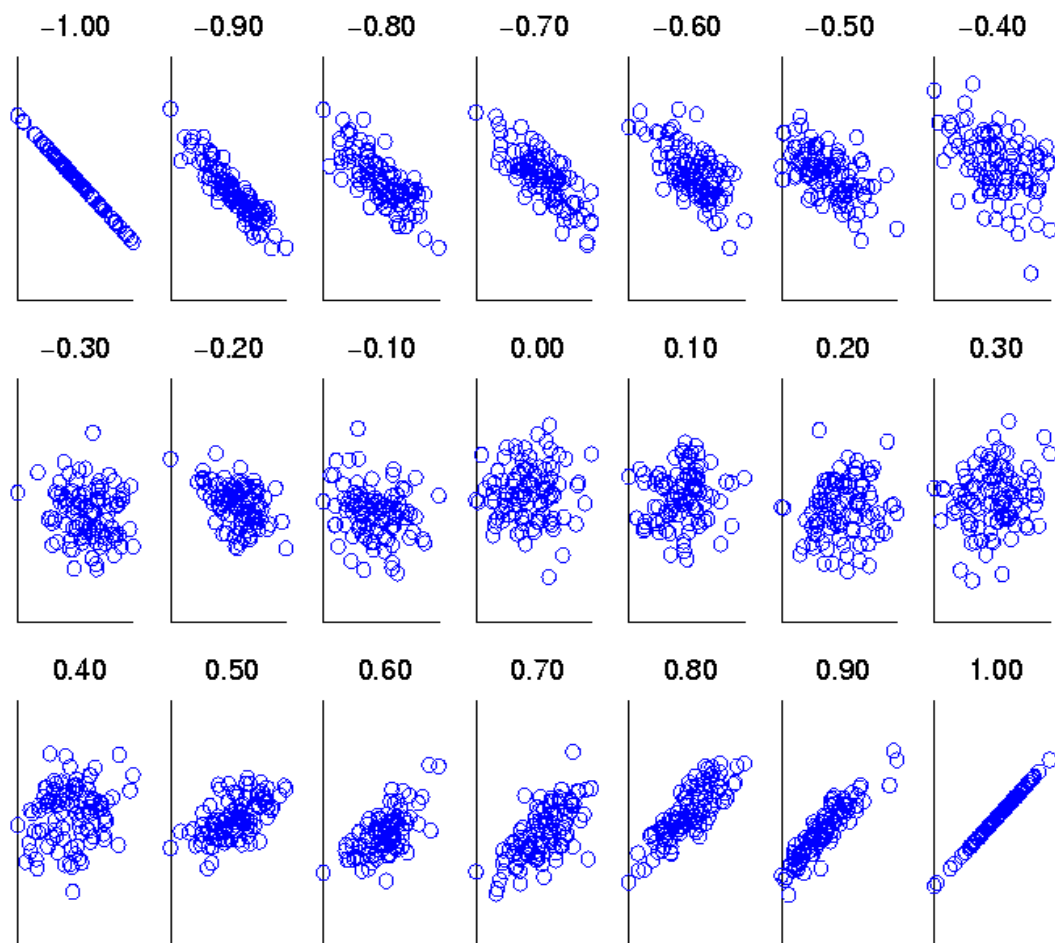$$q'_k = (q_k - mean(q)) / std(q)$$

$$correlation(p,q) = p' \bullet q'$$

# Visually Evaluating Correlation



Scatter plots showing the similarity from –1 to 1.