



# **CS 412 Intro. to Data Mining**

15 w.u.

## **Chapter 2. Getting to Know Your Data**

**Jiawei Han, Computer Science, Univ. Illinois at Urbana-Champaign, 2017**





# Chapter 2. Getting to Know Your Data

---

ឧបករណ៍ → បង្ហាញ និង អនករចនា



- Data Objects and Attribute Types
- Basic Statistical Descriptions of Data
- Data Visualization
- Measuring Data Similarity and Dissimilarity
- Summary

# Types of Data Sets: (1) Record Data

- Relational records
  - Relational tables, highly structured
- Data matrix, e.g., numerical matrix, crosstabs

	China	England	France	Japan	USA	Total
Active Outdoors Crochet Glove		12.00	4.00	1.00	240.00	257.00
Active Outdoors Lycra Glove		10.00	6.00		323.00	339.00
InFlux Crochet Glove	3.00	6.00	8.00		132.00	149.00
InFlux Lycra Glove		2.00			143.00	145.00
Triumph Pro Helmet	3.00	1.00	7.00		333.00	344.00
Triumph Vertigo Helmet		3.00	22.00		474.00	499.00
Xtreme Adult Helmet	8.00	8.00	7.00	2.00	251.00	276.00
Xtreme Youth Helmet		1.00			76.00	77.00
<b>Total</b>	14.00	43.00	54.00	3.00	1,972.00	2,086.00

Data ធនាគារ កំណត់ព័ត៌មាន ឱ្យមាន

ឈរអ៊ី 2  
ម៉ាទ្រិស  
ការងារ  
រូបីន  
Row ក្នុង  
Column

Person:

Pers_ID	Surname	First_Name	City
0	Miller	Paul	London
1	Ortega	Alvaro	Valencia
2	Huber	Urs	Zurich
3	Blanc	Gaston	Paris
4	Bertolini	Fabrizio	Rom

no relation

Car:

Car_ID	Model	Year	Value	Pers_ID
101	Bentley	1973	100000	0
102	Rolls Royce	1965	330000	0
103	Peugeot	1993	500	3
104	Ferrari	2005	150000	4
105	Renault	1998	2000	3
106	Renault	2001	7000	3
107	Smart	1999	2000	2

នៅ: Row នឹងការបង្ហាញ 1 រួច  
Record

- Transaction data ឱ្យមានរាយការ

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

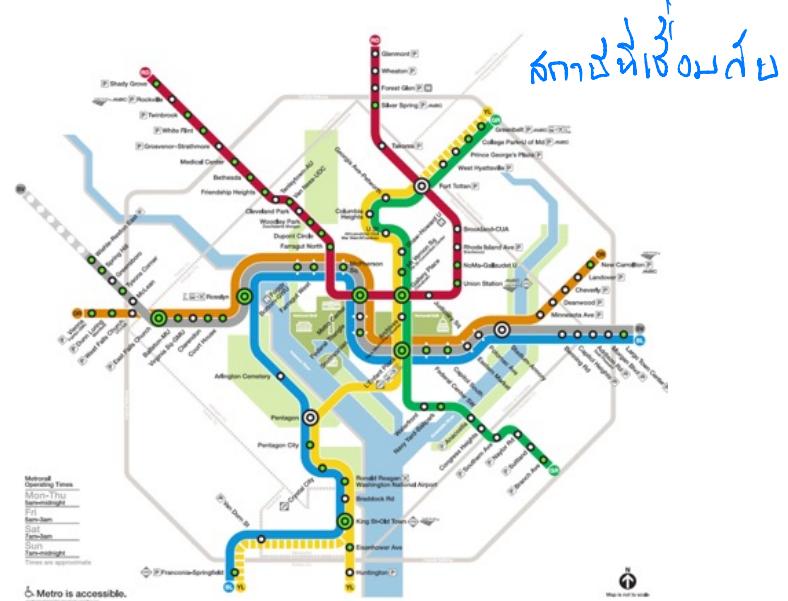
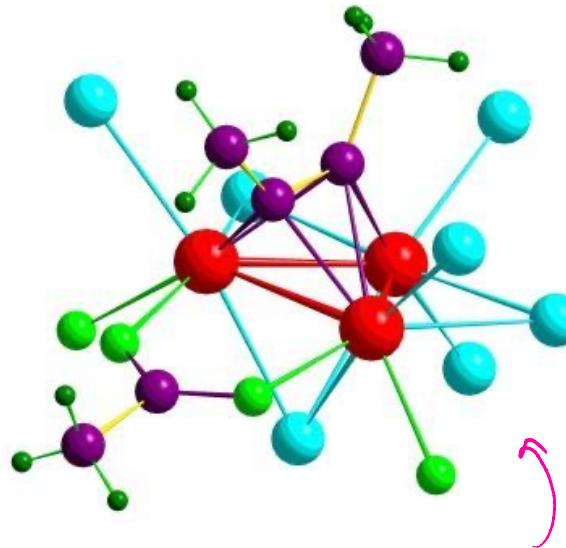
	team	coach	y	pla	ball	score	game	n	wi	lost	timeout	season
Document 1				3	0	5	0	2	6	0	2	0
Document 2				0	7	0	2	1	0	0	3	0
Document 3				0	1	0	0	1	2	2	0	3

នៅ

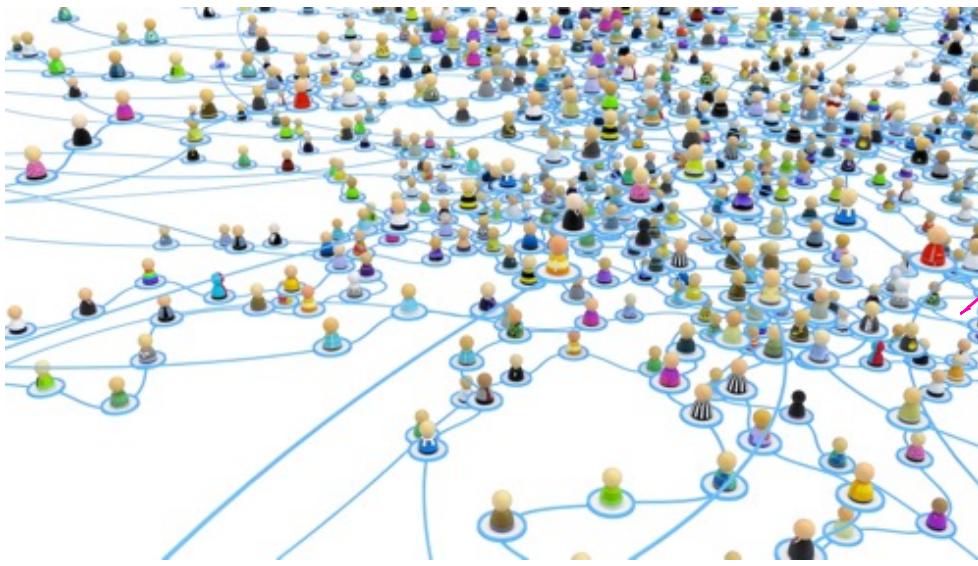
- Document data: Term-frequency vector (matrix) of text documents

# **Types of Data Sets: (2) Graphs and Networks**

- ❑ Transportation network
  - ❑ World Wide Web



- Molecular Structures  
分子構造
  - Social or information networks  
社会或信息网络



- Social networks
- Ansprechpartner

# Types of Data Sets: (3) Ordered Data

① ពេលវេលាក្នុងតម្លៃ  
attribute  
at time

សង្កែរតាមប័ណ្ណ  
សង្គមទូរគមនា  
តម្លៃសម្រាប់

ការអនុវត្តបញ្ជីប្រចាំថ្ងៃ  
តែក្នុងការិយាល័យប៉ុណ្ណោះណាមីនេះ

- Video data: sequence of images

ក្នុងរៀងរាល់

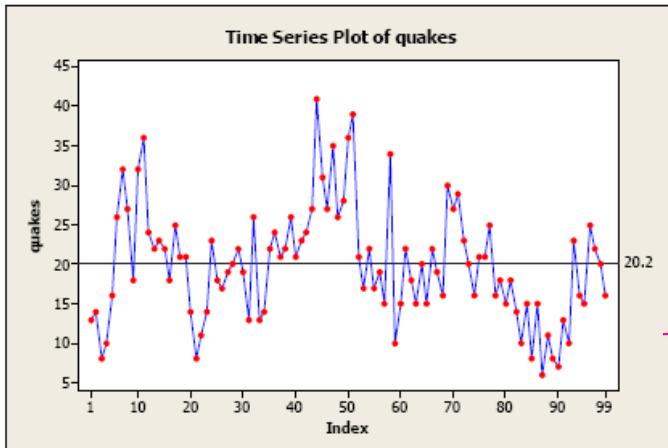
video

គរណនាយករាជការជាមួយក្រុងក្រោម

ការអនុវត្តបញ្ជីប្រចាំថ្ងៃ  
តែក្នុងការិយាល័យប៉ុណ្ណោះណាមីនេះ



- Temporal data: time-series



→ នៅ 2 រៀង ដុំខ្សោន

ការងារ  
នៅរាជរាជក្រឹត់

- Sequential Data: transaction sequences

សម្រាប់លើក  
តាមប័ណ្ណ

ការិយាល័យ  
តិចនាក់រាយជូន

- Genetic sequence data

Human	GTTTGAGG	- - - ATGTTCAACAAATGCTCCTTCATTCCCTATTTACAGACCTGCCGCA
Chimpanzee	GTTTGAGG	- - - ATGTTCAATAATGCTGCTTCACTCCCTATTTACAGACCTGCCGCA
Macaque	GTTTGAGG	- - - ATGCTCAATAATGCTCCTTCATTCCCTCATTACAAACTTGCCGCA
Human	GACAATTCTGCTAGCAGCCTTGTGCTATTATCTGTTTCTAAACCTTAGTAATTGAGTGT	
Chimpanzee	GACAATTCTGCTAGCAGCCTTGTGCTATTATCTGTTTCTAAACCTTAGTAATTGAGTGT	
Macaque	GACAATTCTGCTAGCAGCCTTGTGCTATTATCTGTTTCTAAACCTTAGTAATTGAGTGT	
Human	GATCTGGAGACTAA	- CTC TGA A A A A A A A A G C T G A T T T T T T T T C T C A A A C A A
Chimpanzee	GATCTGGAGACTAA	- CTC TGA A A A A A A A A G C T G A T T T T T T T C T C A A A C A A
Macaque	TATCTGGAGACTAA	- CTC TGA A A A A A A A A G C T G A T T T T T T T T C T C A A A C A A
Human	CAGAACACGATTAGCAAATTACTCTTAAGATATTATTTACATTTCATATTCTCTA	
Chimpanzee	CAGAACACGATTAGCAAATTACTCTTAAGATACTATTTCACATTTCATATTCTCTA	
Macaque	CAGAACATGATTAGCAAATTACCTCTTAAGATATTATTTGCACCTTCATATTCTCTA	
Human	CCCTGAGTTGATGTTGAGCAATATGTCACCTTCATAAAGCCAGGTATACAC	
Chimpanzee	CCCTGAGTTGATGTTGAGCCGATGTCACCTTCATAAAGCCAGGTATACAC	
Macaque	CCCTGAGTTGATGTTGAGCAATATGTCACCTCCACAAAGCCAGGTATATACATTACG	
Human	GACAGGTAAGTAAAAACATATTATTTACGTTTGTCCAAAGAATTAAATTTC	
Chimpanzee	GACAGGTAAGTAAAAACATATTATTTACGTTTGTCCAAAGAATTAAATTTC	
Macaque	GACAGGTAAGTAAAAACATATTATTTACGTTTGTCCAAAGAATTAAATTTC	
Human	AACTGTTGCGCGTGTGGTAA	- - - TGTAAAAACAAACTCAGTACA
Chimpanzee	AACTGTTGCGCGTGTGGTAA	- - - TGTAAAAACAAACTCAGTACA
Macaque	AACTGTTGCGCGTGTGGTAA	- - - CGTAAAAACAAACTCAGTACA

# Types of Data Sets: (4) Spatial, image and multimedia Data

Temporal → ເສີງວາດ  
ເປັນພනໍຫ

- Spatial data: maps



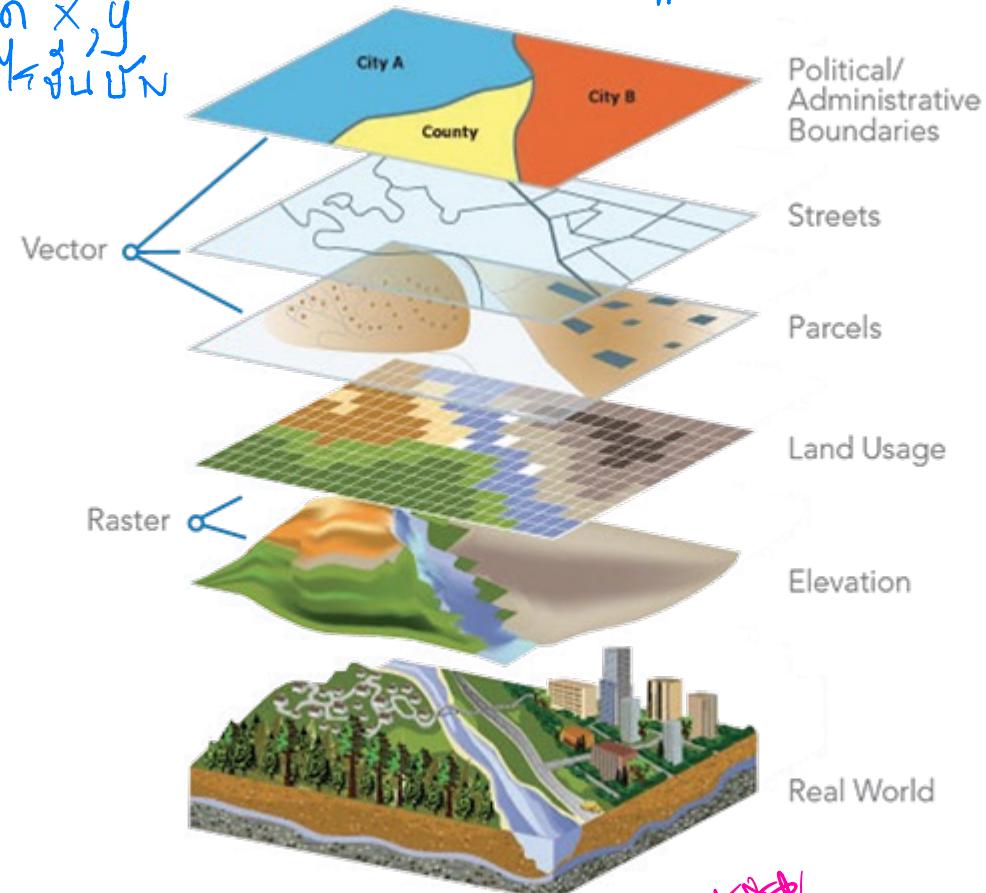
X<sub>1</sub>Y<sub>1</sub>, X<sub>2</sub>Y<sub>2</sub>,  
ເລກຕົວ  
ຈາກບົນ  
ດູບ

(space)

ເຜັນພනໍຫ

ທີ່ນີ້ຄົດ X, Y  
ເທິດວ່າໃຊ້ມູນບັນ

ພາຫນີ



- Image data:

- Video data:

Spatio-temporal  
ມີມານີ້ວາດ  
ເປັນພනໍຫພລະອົກ

ໄຟຟ້າ  
ມີມານີ້ວາດ

# Important Characteristics of Structured Data

វិធី

វាយកដល់សម្រាប់ប័ណ្ណវិទ្យា

## □ Dimensionality

### □ Curse of dimensionality

មិនអាចពិនិត្យបានប្រចាំបីម៉ោង តែ

## □ Sparsity

តាមតាម

(ជាតិប៊ូល 0)

→ នឹង... នាយកន ជាតិប៊ូល 0 ត្រូវបានប្រើប្រាស់បាន

### □ Only presence counts

តាម មានឬមេន ឱ្យយូរ ក្នុងបណ្តុះបណ្តាល

## □ Resolution

គម្រោង

ខ្លួន តិច

ត្រូវ... អីបាន តាមឯណី លម្អិត

គម្រោងនៃការបង្កើត

គន្លឹននៃគោរោង សេរីភាព និងការបង្កើត

រាយការ - ទូរសព្ទ - រាយបោះ - ស្របតាមការបង្កើត

គឺ តិចនៅ ស្ថាន ម. ភូរោះ

គ.បោះពុម្ព —> ន.បោះពុម្ព

ន.បោះពុម្ព - គ.បោះពុម្ព

ការបង្កើត

អនុទាន - oppo - Oppo Flip - ស្ថាន

## □ Distribution

ការកំណត់

### □ Centrality and dispersion

ការបង្កើត និងការបង្កើត

ការបង្កើត

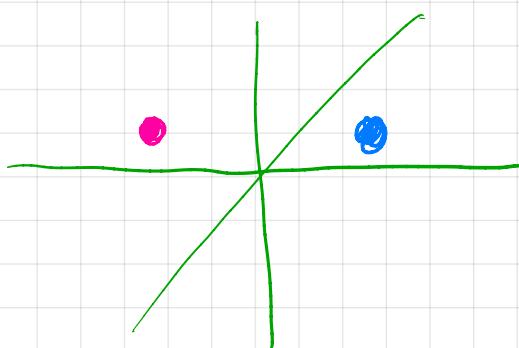
អនុទាន - oppo - Oppo Flip - ស្ថាន

## ຂະໜາດ Dimension

### Dimensionality

#### Curse of dimensionality

ສົ່ງເລີກຕົວທີ່ໄດ້ປັບປຸງໃນ



ຕະຫຼາດ = Dimension

ເລືດ = Dimension

3 Dimension  
9 Dimension

# Data Objects

- Data sets are made up of data objects

- A **data object** represents an entity

□ → data point → ព័ត៌មាន ឈ្មោះ

- Examples:

□ sales database: customers, store items, sales

អ្នកលក់ទូរតារ

Database ព័ត៌មាន  
ក្រសួង: សាខា

□ medical database: patients, treatments

បាន ពេន្យាល់

លេខ.

បានបង់

នាទី

អតិថជន  
អាជីវកម្ម

□ university database: students, professors, courses

- Also called *samples*, *examples*, *instances*, *data points*, *objects*, *tuples*

- Data objects are described by **attributes**

- Database *rows* → data objects; *columns* → attributes

data points

# Attributes

Data mining

ទីតាំងទទួលខ្លួន

អាជីវកម្ម

ប្រភេទ

សរុប

## □ Attribute (or dimensions, features, variables)

- A data field, representing a characteristic or feature of a data object.
- E.g., *customer\_ID, name, address*
- Types:
  - Nominal (e.g., red, blue) ជីវកម្ម និមិត្តុ ពេល
  - Binary (e.g., {true, false}) ជីវកម្ម និមិត្តុ ពេល (ទីផ្សារ)
  - Ordinal (e.g., {freshman, sophomore, junior, senior})  
ជីវកម្ម Nominal និង Ordinal គឺជាប្រភេទ សង្កែសារ (ទីផ្សារ)
  - Numeric: quantitative ជីវកម្ម 1 2 3 4
  - Interval-scaled:  $100^{\circ}\text{C}$  is interval scales  
ជីវកម្ម ក្នុង 0 ដល់ 100  $\rightarrow 0^{\circ}$  តែ វិវាទ
  - Ratio-scaled:  $100^{\circ}\text{K}$  is ratio scaled since it is twice as high as  $50^{\circ}\text{K}$
- Q1: Is student ID a nominal, ordinal, or interval-scaled data?
- Q2: What about eye color? Or color in the color spectrum of physics?

# Attribute Types

- **Nominal:** categories, states, or “names of things”

- *Hair\_color* = {*auburn, black, blond, brown, grey, red, white*}
- marital status, occupation, ID numbers, zip codes

- **Binary**

- Nominal attribute with only 2 states (0 and 1)

- Symmetric binary: both outcomes equally important

- e.g., gender

Tom ชาย 0  
Jerry 女性 1

เกต ไม่เกต  
ชื่อ นามสกุล

หมวดบดบัง 10%  
ผู้สูงวัย 99.99%  
↓ ความปั่นผื้นที่ทางเพศ  
↑ ทุกอาชญากรรม

- Asymmetric binary: outcomes not equally important.

- e.g., medical test (positive vs. negative)

↑ ดูเหมือนคนไข้เป็น  
↓ เจ้าเป็นผู้เชี่ยวชาญด้าน หายใจ  
Ex. 1 คันถัง 100 ลบ

- Convention: assign 1 to most important outcome (e.g., HIV positive)

- **Ordinal**

- Values have a meaningful order (ranking) but magnitude between successive values is not known
- *Size* = {*small, medium, large*}, grades, army rankings

# Numeric Attribute Types

---

- Quantity (integer or real-valued)
- Interval
  - Measured on a scale of **equal-sized units**
  - Values have order
    - E.g., *temperature in C° or F°, calendar dates*
  - No true zero-point  
*Handwritten note: 0 mm*
- Ratio
  - Inherent **zero-point**  
*Handwritten note: 0 nm*
  - We can speak of values as being an order of magnitude larger than the unit of measurement (10 K° is twice as high as 5 K°).
    - e.g., *temperature in Kelvin, length, counts, monetary quantities*

# Discrete vs. Continuous Attributes

## Discrete Attribute

ມີຊັບສິນນີ້ມາດູ - ມອົງດູ  
ພ່າຍຕູ - ຢຸມມົມວົງ

- Has only a finite or countably infinite set of values
  - E.g., zip codes, profession, or the set of words in a collection of documents
  - Sometimes, represented as integer variables
  - Note: Binary attributes are a special case of discrete attributes

## Continuous Attribute

- Has real numbers as attribute values
  - E.g., 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

# **Chapter 2. Getting to Know Your Data**

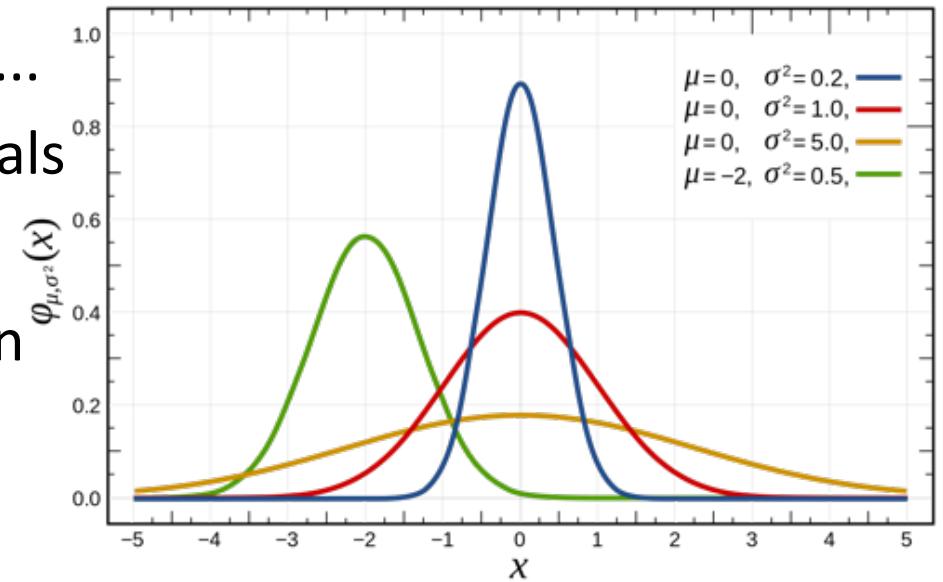
---

- Data Objects and Attribute Types
- Basic Statistical Descriptions of Data
- Data Visualization
- Measuring Data Similarity and Dissimilarity
- Summary



# Basic Statistical Descriptions of Data

- Motivation
  - To better understand the data: central tendency, variation and spread
- Data dispersion characteristics
  - Median, max, min, quantiles, outliers, variance, ...
- Numerical dimensions correspond to sorted intervals
  - Data dispersion:
    - Analyzed with multiple granularities of precision
    - Boxplot or quantile analysis on sorted intervals
- Dispersion analysis on computed measures
  - Folding measures into numerical dimensions
  - Boxplot or quantile analysis on the transformed cube



# Measuring the Central Tendency: (1) Mean

- Mean (algebraic measure) (sample vs. population):

Note:  $n$  is sample size and  $N$  is population size.

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad \mu = \frac{\sum x}{N}$$

- Weighted arithmetic mean:

$$\bar{x} = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$$

- Trimmed mean:

ជាមួយ<sup>ជាមួយ</sup>  
ស្រុក<sup>ស្រុក</sup>  
ត្រូវបាន<sup>ត្រូវបាន</sup>  
ស្វែងរក<sup>ស្វែងរក</sup>  
ដើម្បី<sup>ដើម្បី</sup>  
បង្ហាញ<sup>បង្ហាញ</sup>

- Chopping extreme values (e.g., Olympics gymnastics score computation)

ការបង្ហាញ<sup>ការបង្ហាញ</sup>  
ដើម្បី<sup>ដើម្បី</sup>  
បង្ហាញ<sup>បង្ហាញ</sup>

# Measuring the Central Tendency: (2) Median

## □ Median:

- Middle value if odd number of values, or average of the middle two values otherwise
- Estimated by interpolation (for *grouped data*):

age	frequency
1–5	200
6–15	450
16–20	300
21–50	1500
51–80	700
81–110	44

Median

សម្រាប់បង្ហាញ តាមរយៈបន្ទាន់ 1000

ក្រុងនេះ

Approximate  
median

Sum before the median interval

$$\text{median} = L_1 + \left( \frac{n/2 - (\sum \text{freq})_l}{\text{freq}_{\text{median}}} \right) \text{width}$$



Low interval limit

Interval width ( $L_2 - L_1$ )



# Measuring the Central Tendency: (3) Mode

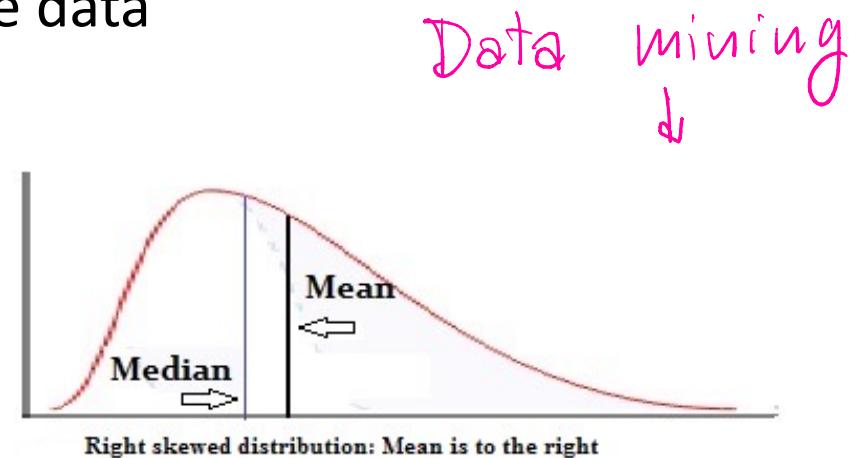
- Mode: Value that occurs most frequently in the data

ជាតុលេខាដែលបានបង្ហាញបានច្បាស់

- Unimodal → មិនមែនមែនស្ថាមួយទេ

- Empirical formula:

$$\text{mean} - \text{mode} = 3 \times (\text{mean} - \text{median})$$



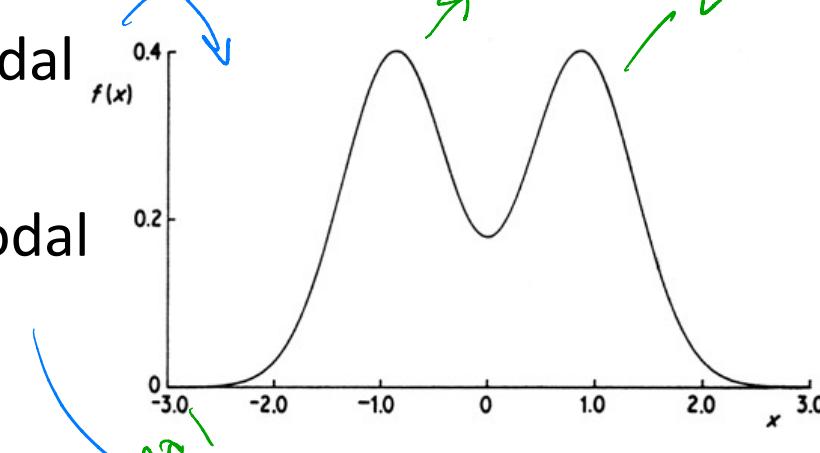
- Multi-modal

- Bimodal

- Trimodal

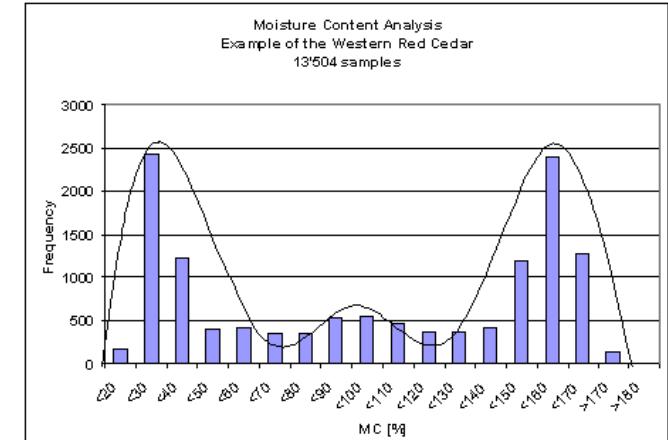
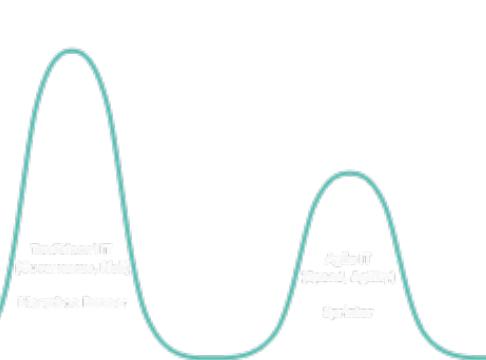
Assumption  
ការពេលរដ្ឋ

តារាង  
ជាពេលរដ្ឋ



ជាតុលេខាណាមួយទេ

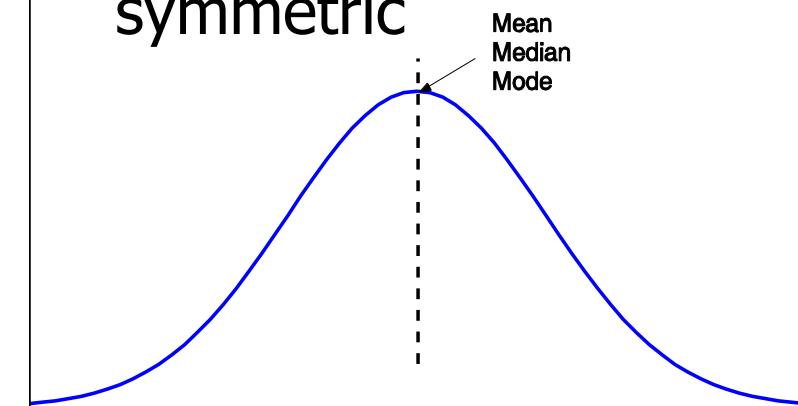
mode → គឺជាដាកណែស្ថាមួយទេ



# Symmetric vs. Skewed Data

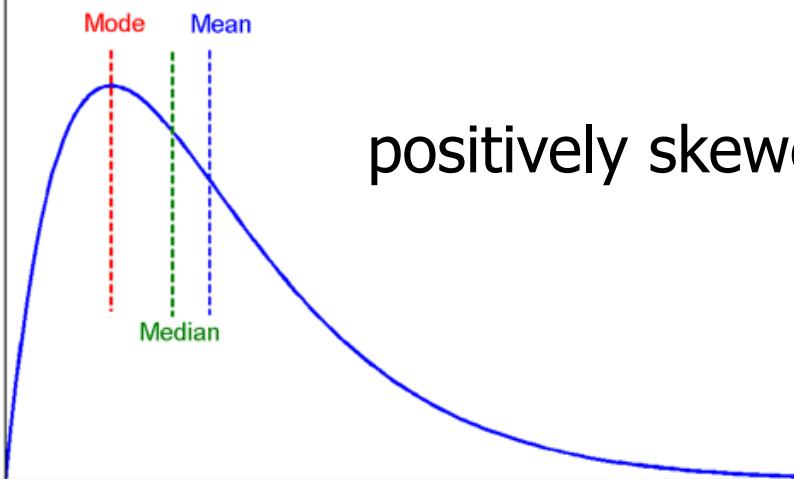
- Median, mean and mode of symmetric, positively and negatively skewed data

symmetric

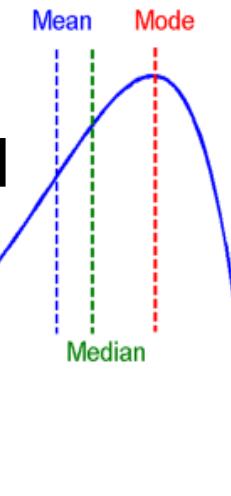


positively skewed

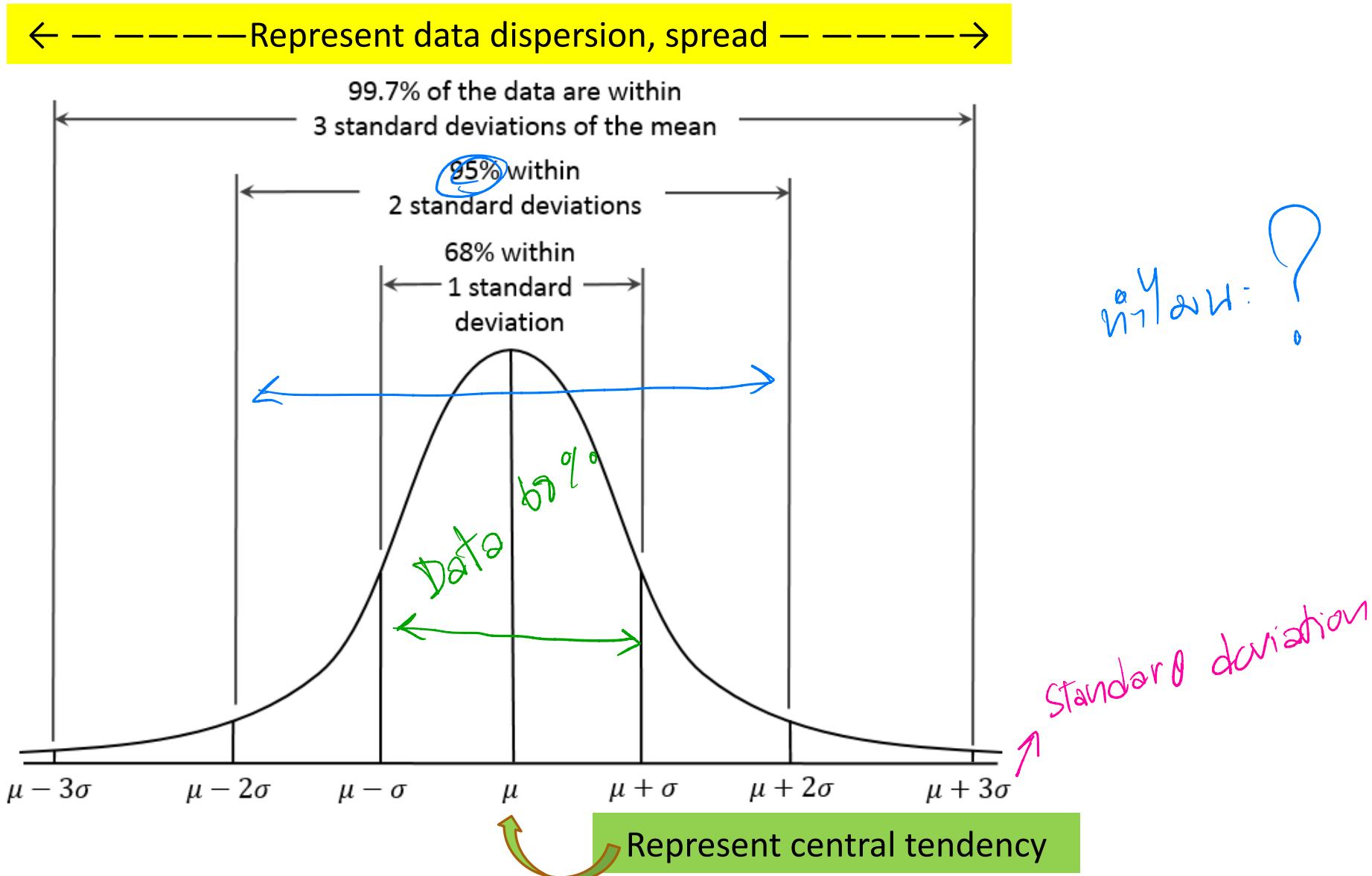
↑  
mean



negatively skewed



# Properties of Normal Distribution Curve



# Measures Data Distribution: Variance and Standard Deviation

- ❑ Variance and standard deviation (*sample: s, population: σ*)

- ❑ **Variance:** (algebraic, scalable computation)

- ❑ Q: Can you compute it incrementally and efficiently?

ມີ ມູນຄົງກໍາຕັ້ງ  
ກວ່າ ..

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n-1} \left[ \sum_{i=1}^n x_i^2 - \frac{1}{n} \left( \sum_{i=1}^n x_i \right)^2 \right]$$

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^n (x_i - \mu)^2 = \frac{1}{N} \sum_{i=1}^n x_i^2 - \mu^2$$

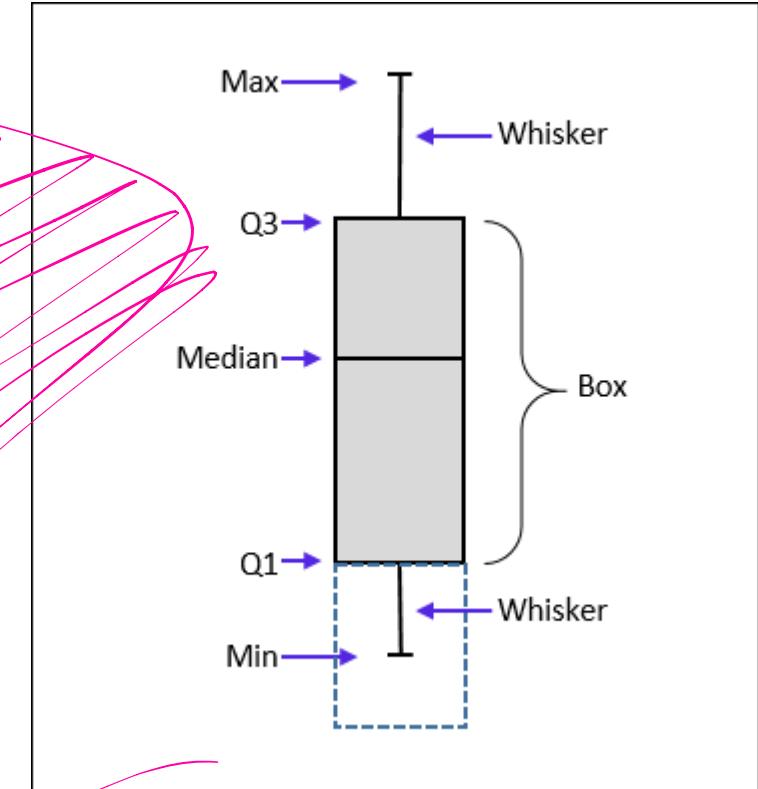
- ❑ **Standard deviation s (or σ)** is the square root of variance  $s^2$  (or  $\sigma^2$ )

# Graphic Displays of Basic Statistical Descriptions

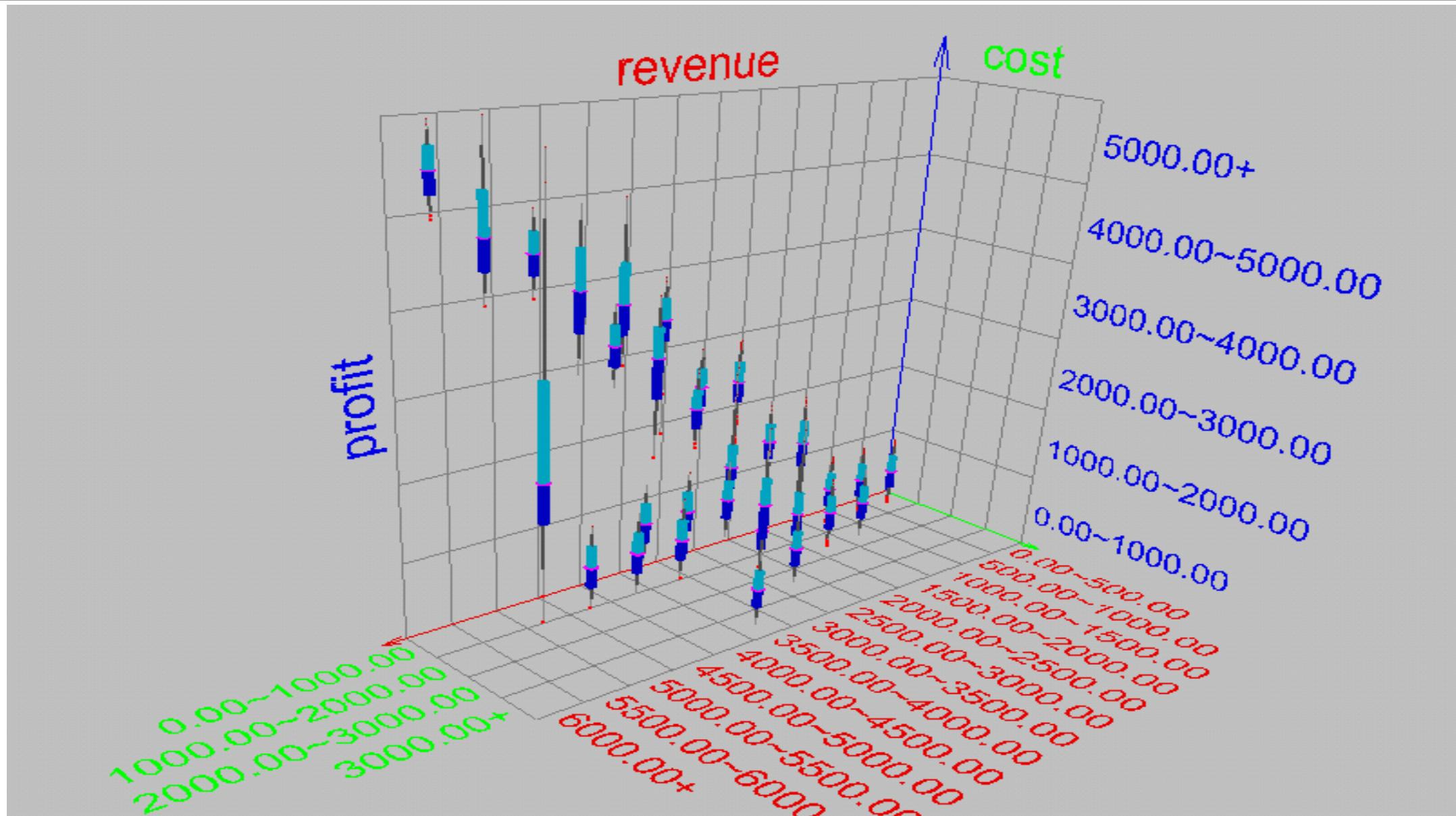
- plot ດ້ວຍສິນ, 5 ດ້ວຍ ມັດງານ  $Q_1 - Q_3$  ພັດ ມີມູນຄົງ ມີມູນກົງ
- Boxplot:** graphic display of five-number summary
- Histogram:** x-axis are values, y-axis represents frequencies
- Quantile plot:** each value  $x_i$  is paired with  $f_i$ , indicating that approximately  $100 f_i \%$  of data are  $\leq x_i$
- Quantile-quantile (q-q) plot:** graphs the quantiles of one univariate distribution against the corresponding quantiles of another
- Scatter plot:** each pair of values is a pair of coordinates and plotted as points in the plane

# Measuring the Dispersion of Data: Quartiles & Boxplots

- **Quartiles:**  $Q_1$  ( $25^{\text{th}}$  percentile),  $Q_3$  ( $75^{\text{th}}$  percentile)
- **Inter-quartile range:**  $IQR = Q_3 - Q_1$
- **Five number summary:** min,  $Q_1$ , median,  $Q_3$ , max
- **Boxplot:** Data is represented with a box
  - $Q_1$ ,  $Q_3$ , IQR: The ends of the box are at the first and third quartiles, i.e., the height of the box is  $IQR$
  - Median ( $Q_2$ ) is marked by a line within the box
  - Whiskers: two lines outside the box extended to Minimum and Maximum
  - Outliers: points beyond a specified outlier threshold, plotted individually
  - **Outlier:** usually, a value higher/lower than  $1.5 \times IQR$

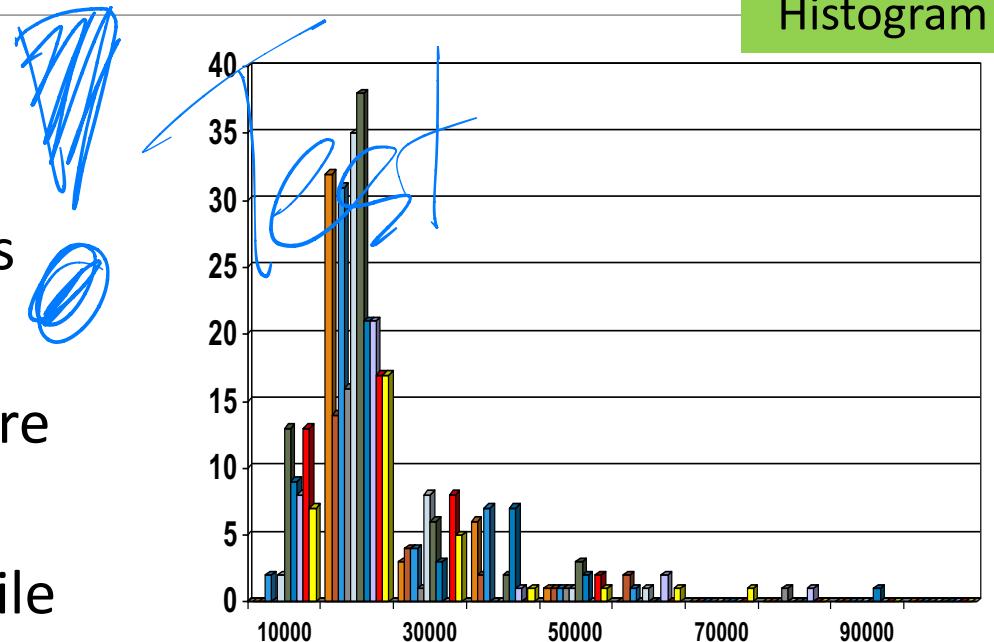


# Visualization of Data Dispersion: 3-D Boxplots

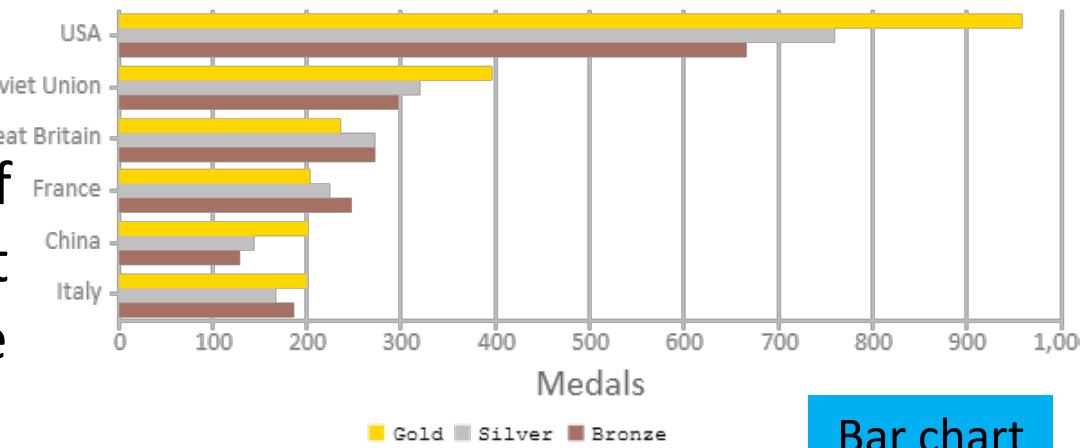


# Histogram Analysis

- ❑ Histogram: Graph display of tabulated frequencies, shown as bars
- ❑ Differences between histograms and bar charts
  - ❑ Histograms are used to show distributions of variables while bar charts are used to compare variables
  - ❑ Histograms plot binned quantitative data while bar charts plot categorical data
  - ❑ Bars can be reordered in bar charts but not in histograms
  - ❑ Differs from a bar chart in that it is the area of the bar that denotes the value, not the height as in bar charts, a crucial distinction when the categories are not of uniform width



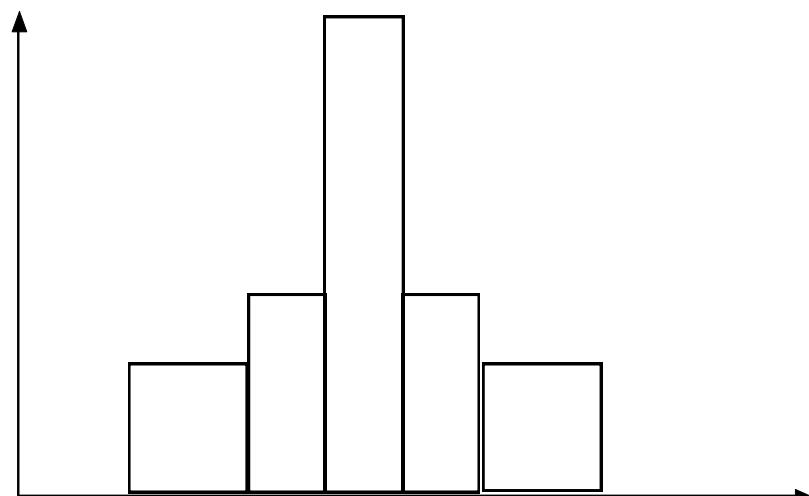
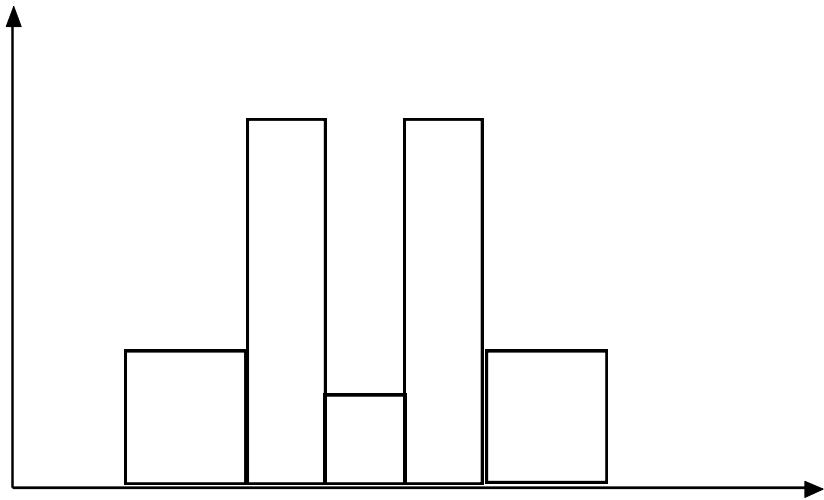
Olympic Medals of all Times (till 2012 Olympics)



Bar chart

# Histograms Often Tell More than Boxplots

---

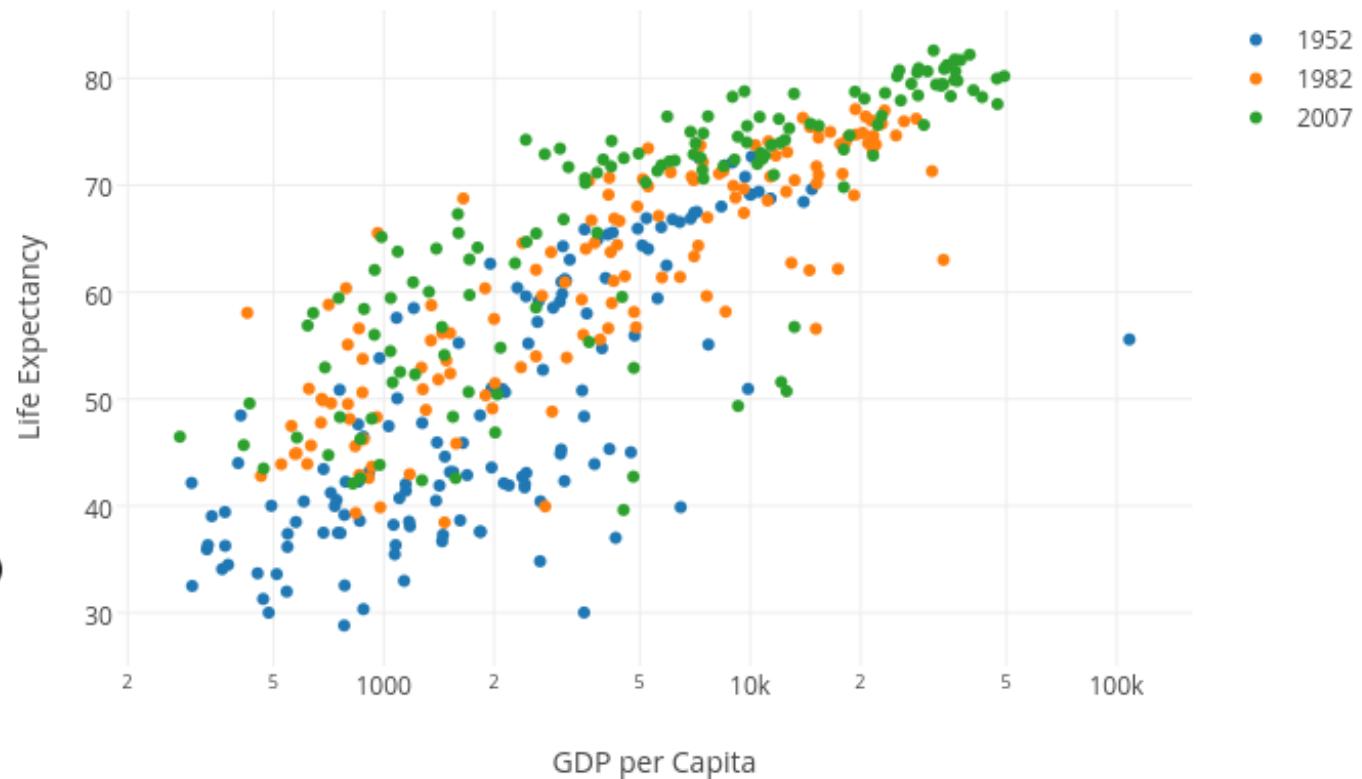
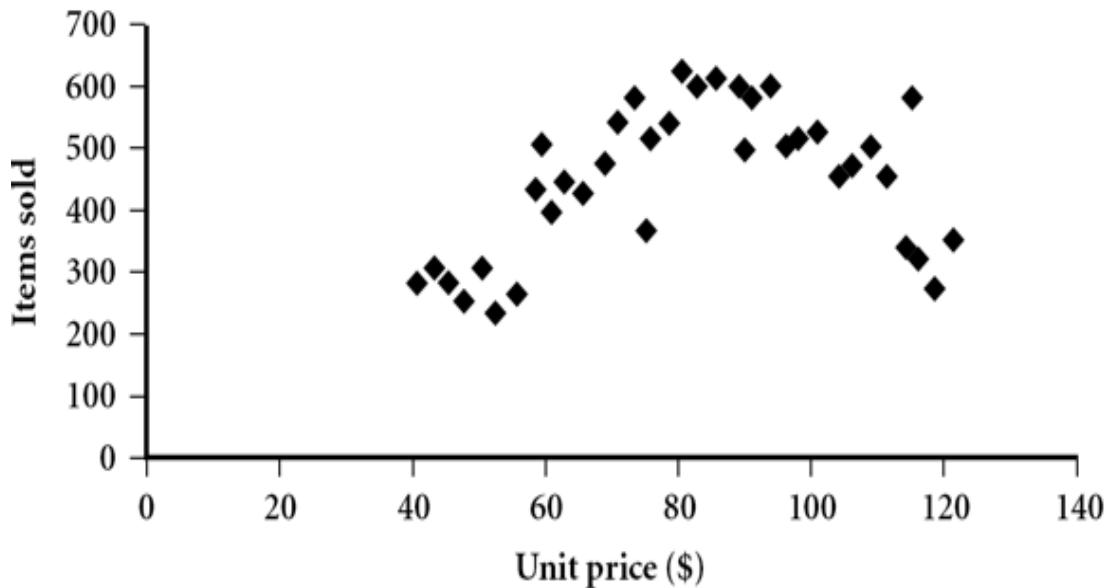


- The two histograms shown in the left may have the same boxplot representation
- The same values for: min, Q1, median, Q3, max
- But they have rather different data distributions



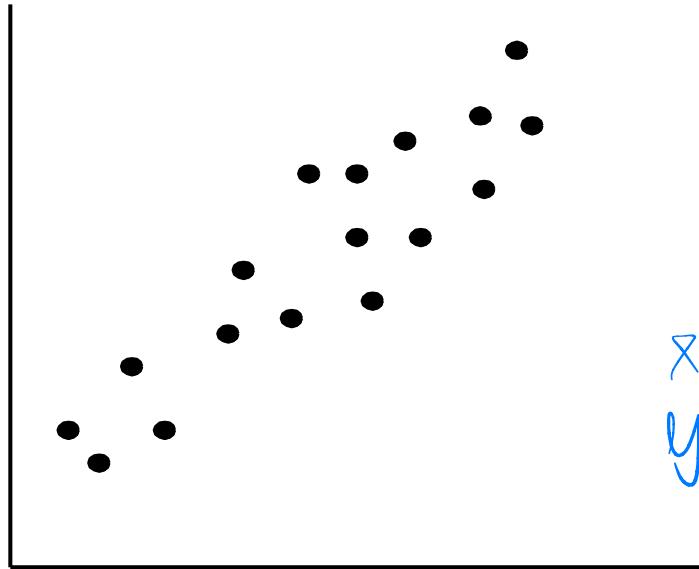
# Scatter plot

- ❑ Provides a first look at bivariate data to see clusters of points, outliers, etc.
- ❑ Each pair of values is treated as a pair of coordinates and plotted as points in the plane

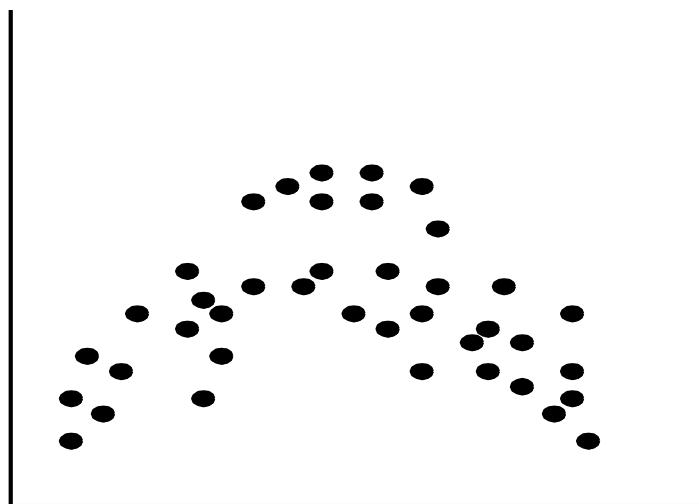


# Positively and Negatively Correlated Data

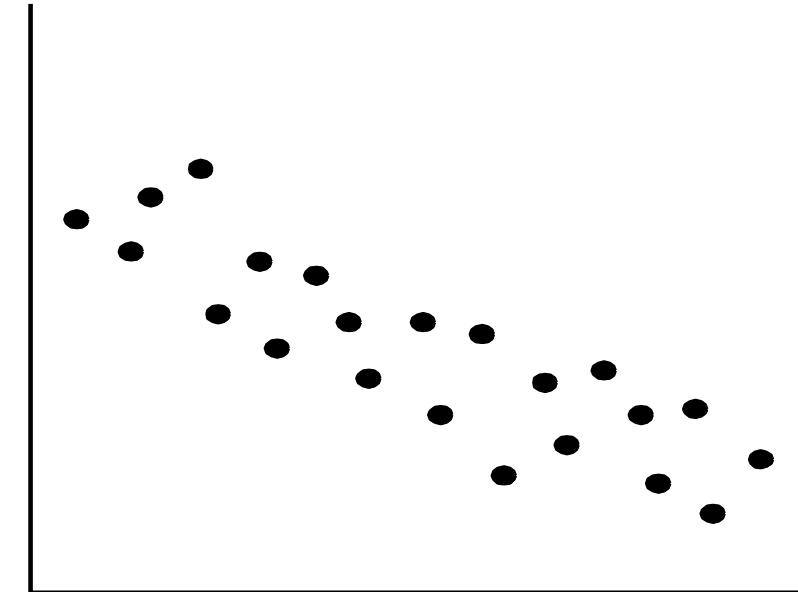
What  
Pain  
values?  
?



X  
Y  
What  
is  
this?



X  
Y  
What  
is  
this?



X  
Y  
?

- The left half fragment is positively correlated
- The right half is negative correlated

# Uncorrelated Data

---



မျှတော်လုပ်နည်



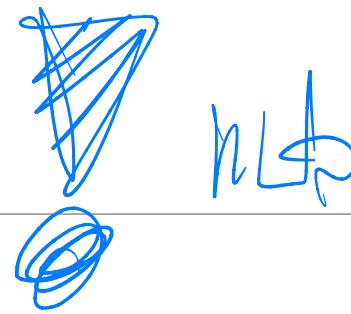
နဲ့ outlier ဖြစ်တဲ့ ပဲ

# **Chapter 2. Getting to Know Your Data**

---

- Data Objects and Attribute Types
- Basic Statistical Descriptions of Data
- Data Visualization 
- Measuring Data Similarity and Dissimilarity
- Summary

# Data Visualization



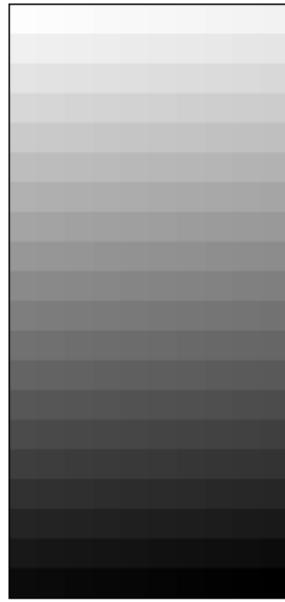
- Why data visualization?
  - Gain insight into an information space by mapping data onto graphical primitives
  - Provide qualitative overview of large data sets
  - Search for patterns, trends, structure, irregularities, relationships among data
  - Help find interesting regions and suitable parameters for further quantitative analysis      දැන්ත්වා යොමු කළ ඇති
  - Provide a visual proof of computer representations derived
- Categorization of visualization methods:
  - Pixel-oriented visualization techniques
  - Geometric projection visualization techniques
  - Icon-based visualization techniques
  - Hierarchical visualization techniques
  - Visualizing complex data and relations

# Pixel-Oriented Visualization Techniques

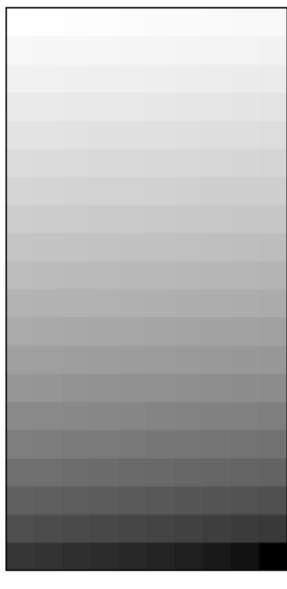
- For a data set of  $m$  dimensions, create  $m$  windows on the screen, one for each dimension
- The  $m$  dimension values of a record are mapped to  $m$  pixels at the corresponding positions in the windows
- The colors of the pixels reflect the corresponding values

维度值映射到像素

多个属性  
在并列的  
图上显示



(a) Income



(b) Credit Limit

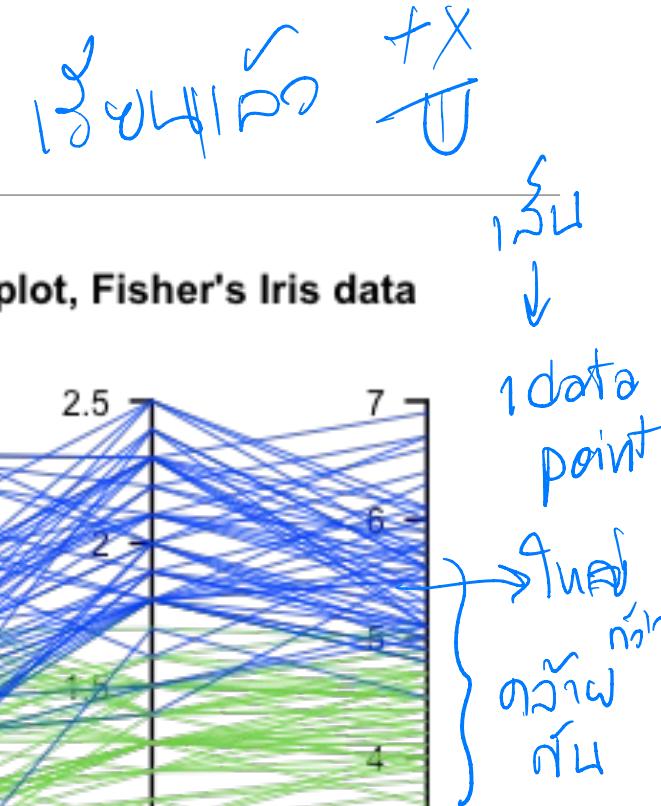


(c) transaction volume

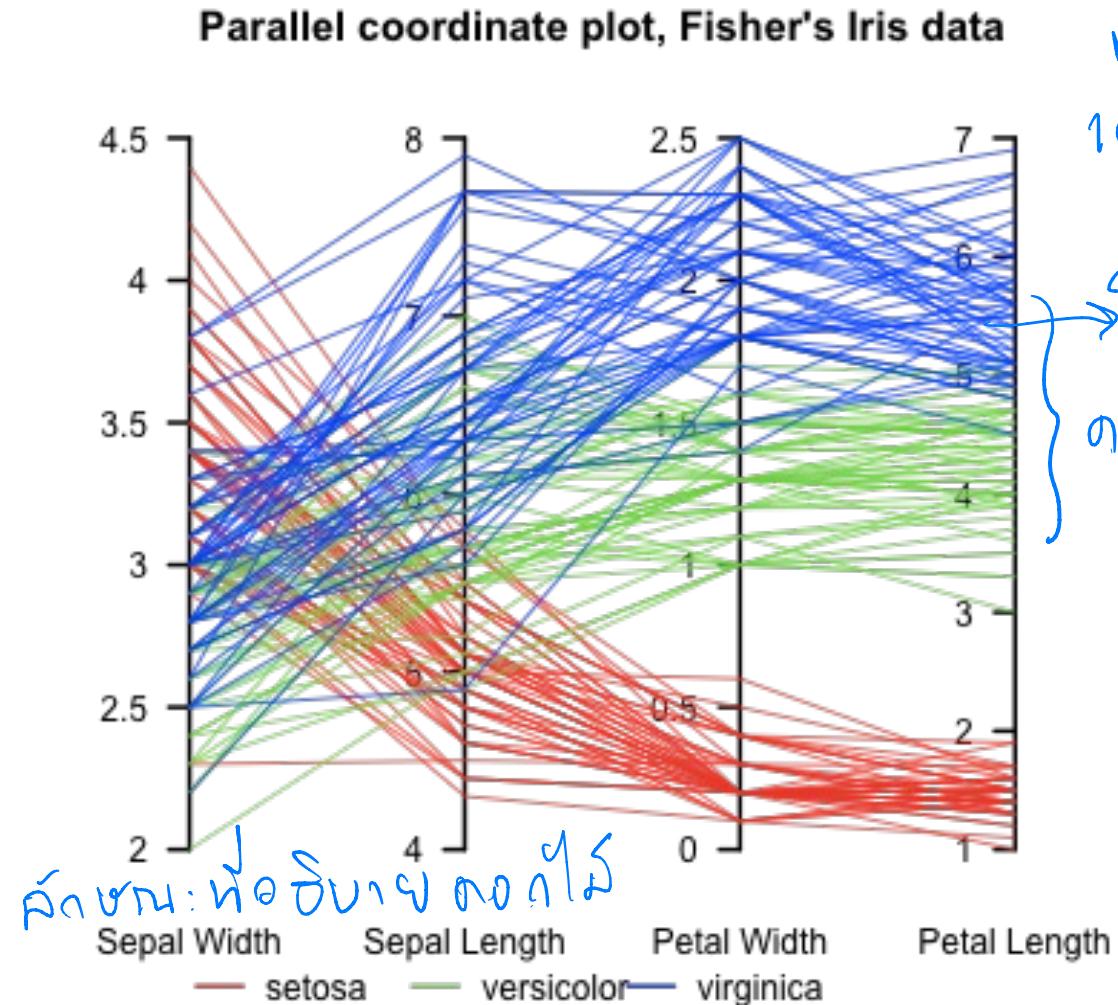


(d) age

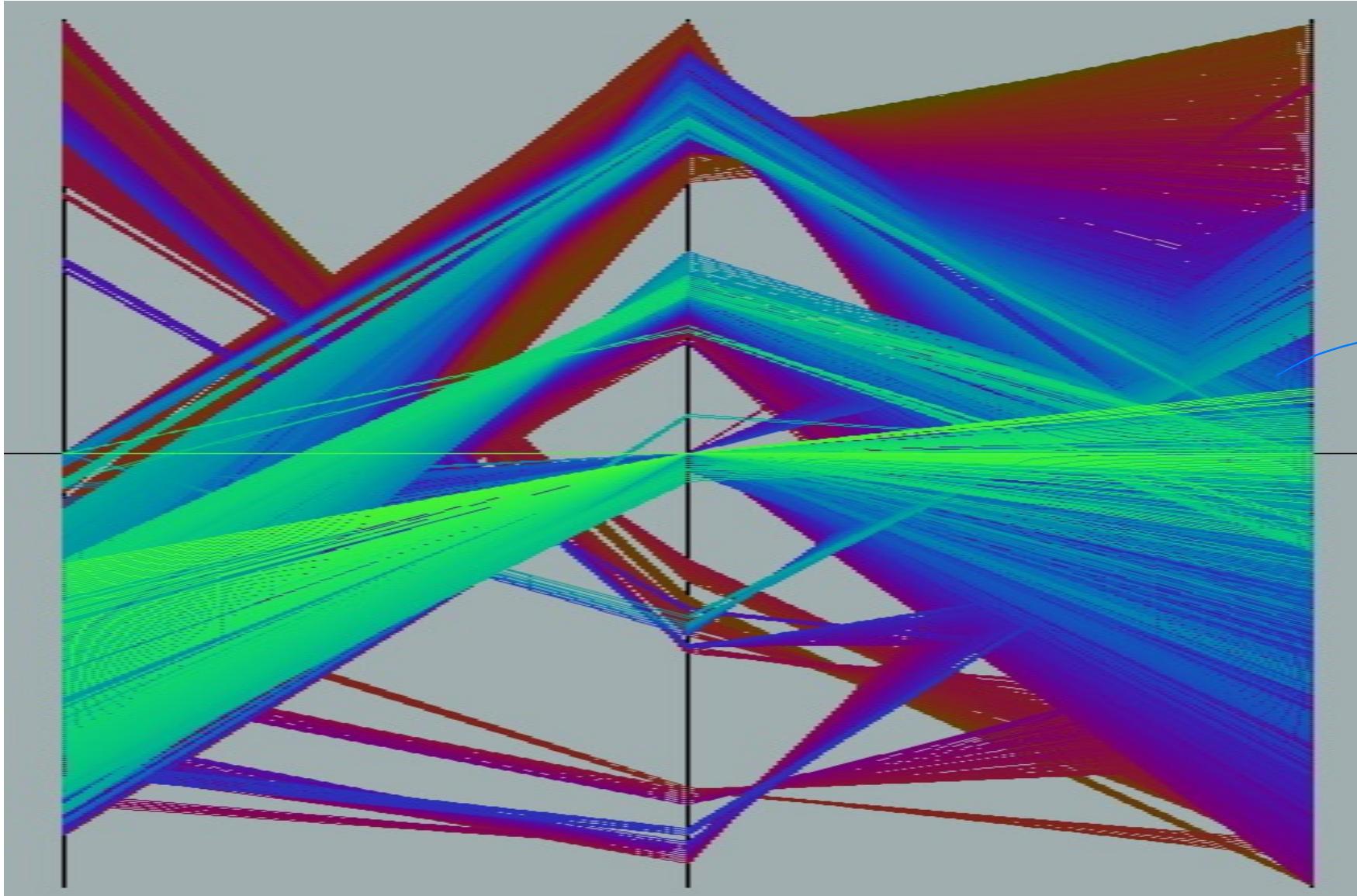
# Parallel Coordinates



- n equidistant axes which are parallel to one of the screen axes and correspond to the attributes
- The axes are scaled to the [minimum, maximum]: range of the corresponding attribute
- Every data item corresponds to a polygonal line which intersects each of the axes at the point which corresponds to the value for the attribute



# Parallel Coordinates of a Data Set

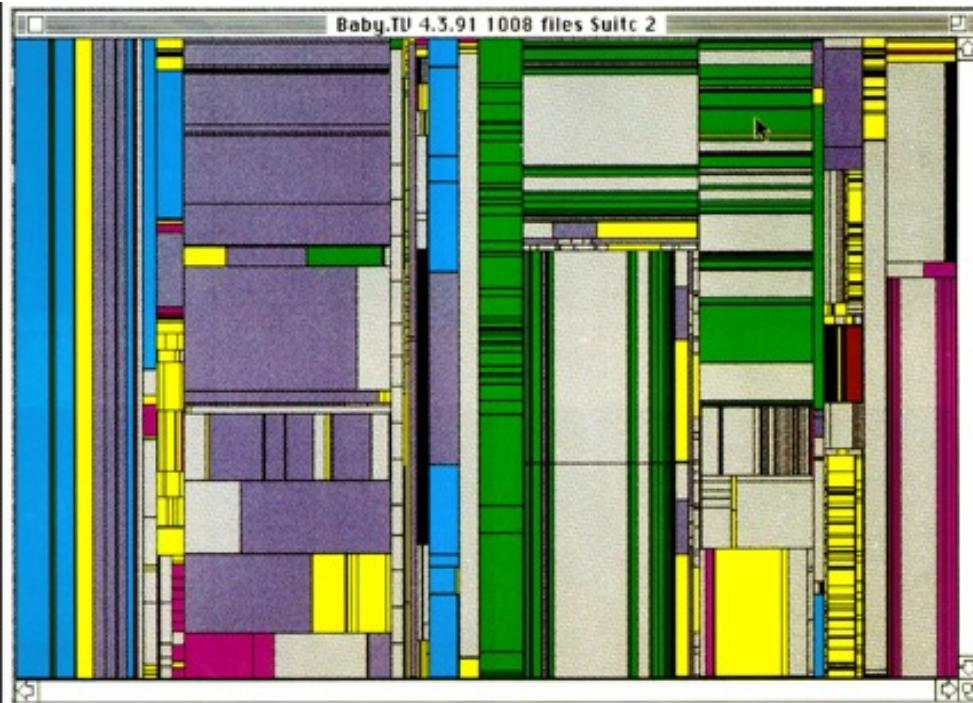


କ୍ରମିକ  
ପରିବହନ

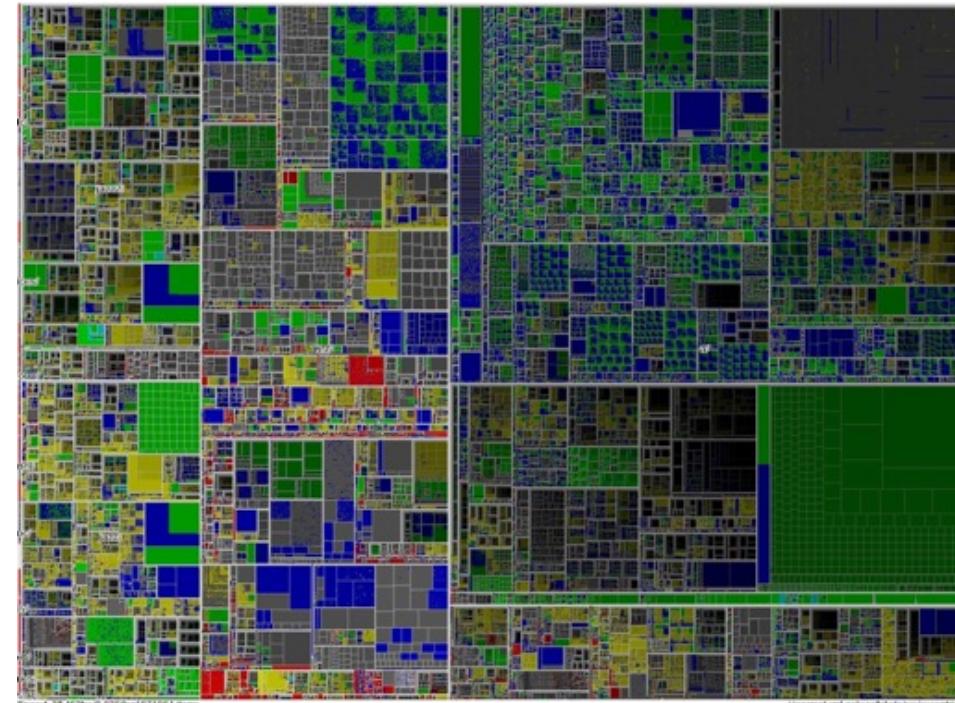
# Tree-Map

W W W 100%

- Screen-filling method which uses a hierarchical partitioning of the screen into regions depending on the attribute values
- The x- and y-dimension of the screen are partitioned alternately according to the attribute values (classes)



Schneiderman@UMD: Tree-Map of a File System



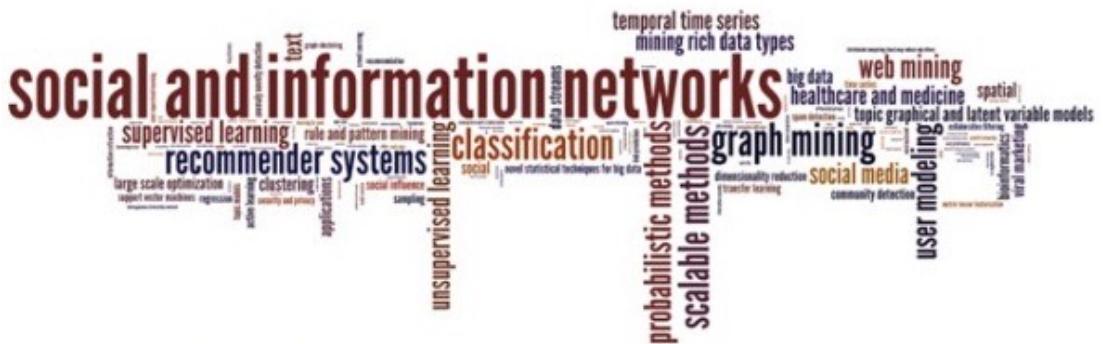
Schneiderman@UMD: Tree-Map to support large data sets of a million items

# Visualizing Complex Data and Relations: Tag Cloud

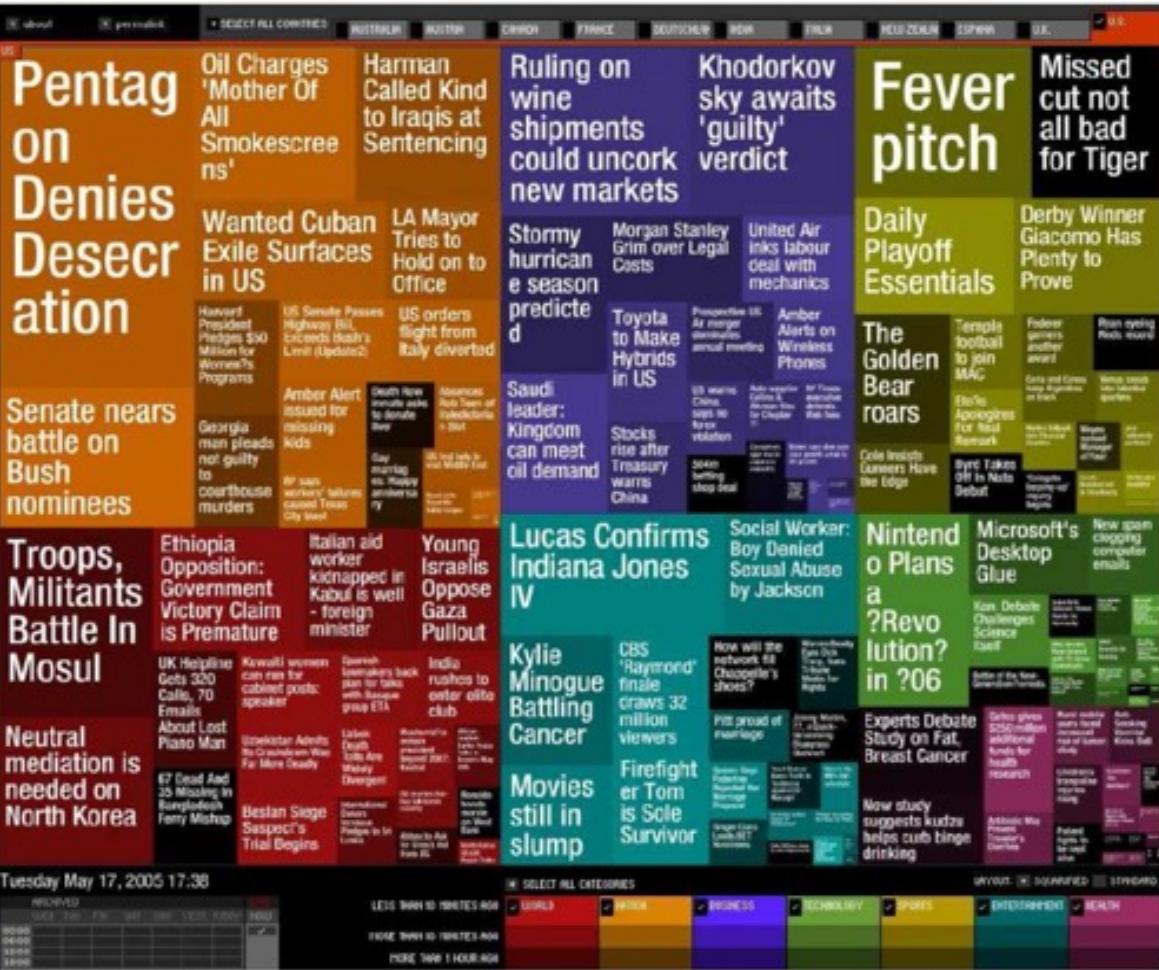
- ❑ Tag cloud: Visualizing user-generated tags
  - ❑ The importance of tag is represented by font size/color
  - ❑ Popularly used to visualize word/phrase distributions

↗ ମହିଳା / ମହାରାଜୀ

עכטניא  
דכלה



KDD 2013 Research Paper Title Tag Cloud



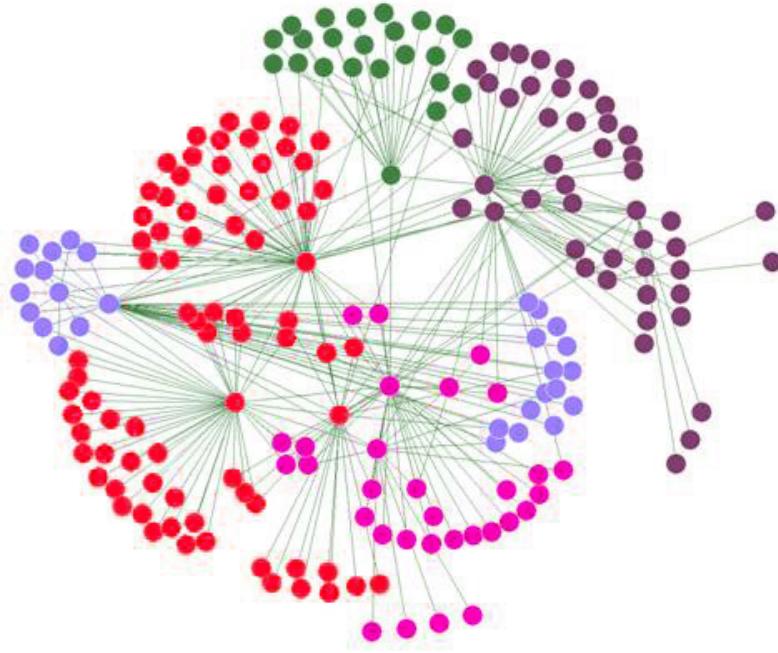
Newsmap: Google News Stories in 2005



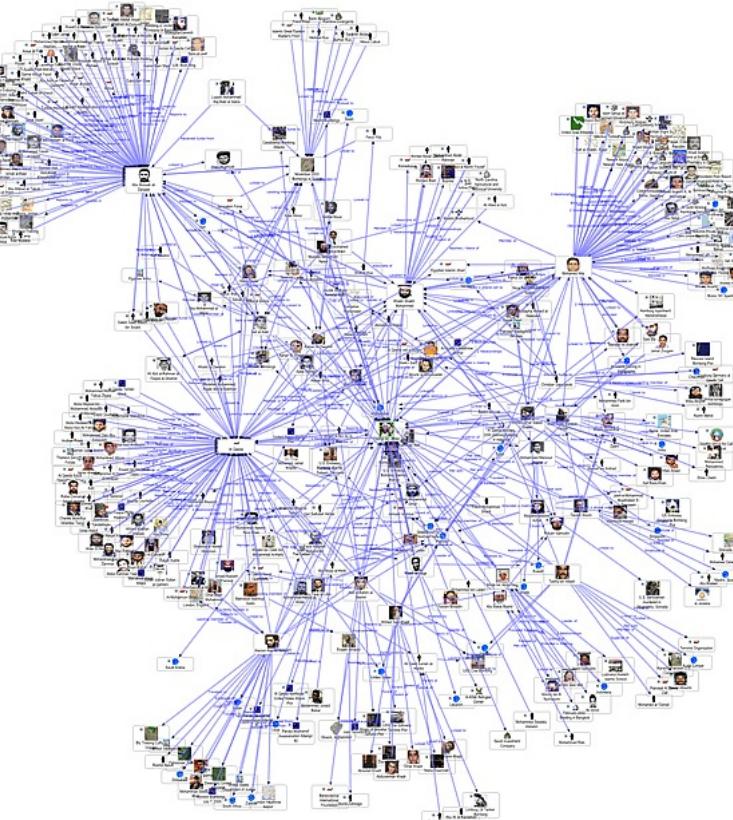
Ex. Treemap bins for  
information

# Visualizing Complex Data and Relations: Social Networks

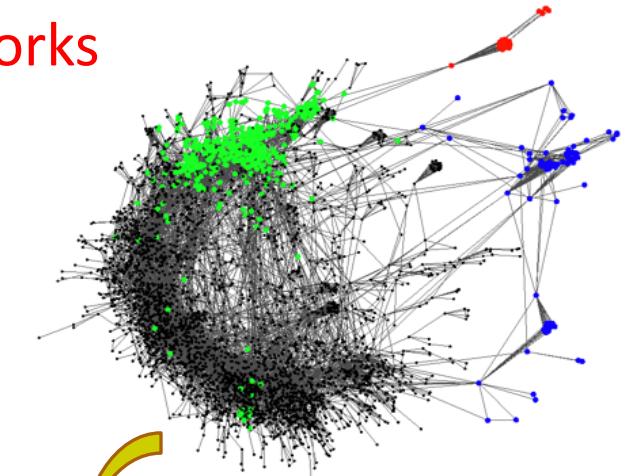
- Visualizing non-numerical data: social and information networks



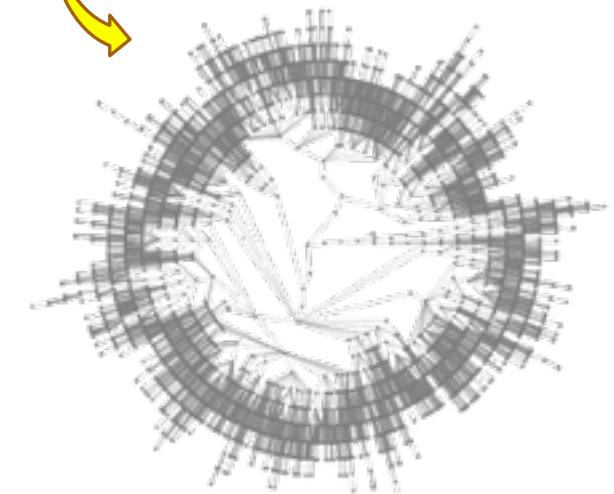
A typical network structure



A social network



organizing  
information networks



# **Chapter 2. Getting to Know Your Data**

---

- Data Objects and Attribute Types
- Basic Statistical Descriptions of Data
- Data Visualization
- Measuring Data Similarity and Dissimilarity
- Summary



# Similarity, Dissimilarity, and Proximity

- **Similarity measure or similarity function** ?
  - A real-valued function that quantifies the similarity between two objects
  - Measure how two data objects are alike: The higher value, the more alike
  - Often falls in the range [0,1]: 0: no similarity; 1: completely similar
- **Dissimilarity (or distance) measure** ?
  - Numerical measure of how different two data objects are
  - In some sense, the inverse of similarity: The lower, the more alike
  - Minimum dissimilarity is often 0 (i.e., completely similar)
  - Range [0, 1] or [0,  $\infty$ ), depending on the definition
- **Proximity** usually refers to either similarity or dissimilarity

நிலைப்பு

விரைவுமன்றங்கள்

?

தொழில்துறை

தொழில்!

முக்கேண்டுமின்று வருகிறது

# Data Matrix and Dissimilarity Matrix

- Data matrix

- A data matrix of  $n$  data points with  $l$  dimensions

- Dissimilarity (distance) matrix

- $n$  data points, but registers only the distance  $d(i, j)$   
(typically metric)

- Usually symmetric, thus a triangular matrix

- **Distance functions** are usually different for real, boolean, categorical, ordinal, ratio, and vector variables

- Weights can be associated with different variables based on applications and data semantics

$$D = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1l} \\ x_{21} & x_{22} & \dots & x_{2l} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nl} \end{pmatrix}$$

$$\begin{pmatrix} 0 & & & \\ d(2,1) & 0 & & \\ \vdots & \vdots & \ddots & \\ d(n,1) & d(n,2) & \dots & 0 \end{pmatrix}$$

# Standardizing Numeric Data

---

- Z-score:

$$z = \frac{x - \mu}{\sigma}$$

- X: raw score to be standardized,  $\mu$ : mean of the population,  $\sigma$ : standard deviation
- the distance between the raw score and the population mean in units of the standard deviation
- negative when the raw score is below the mean, “+” when above
- An alternative way: Calculate the mean absolute deviation

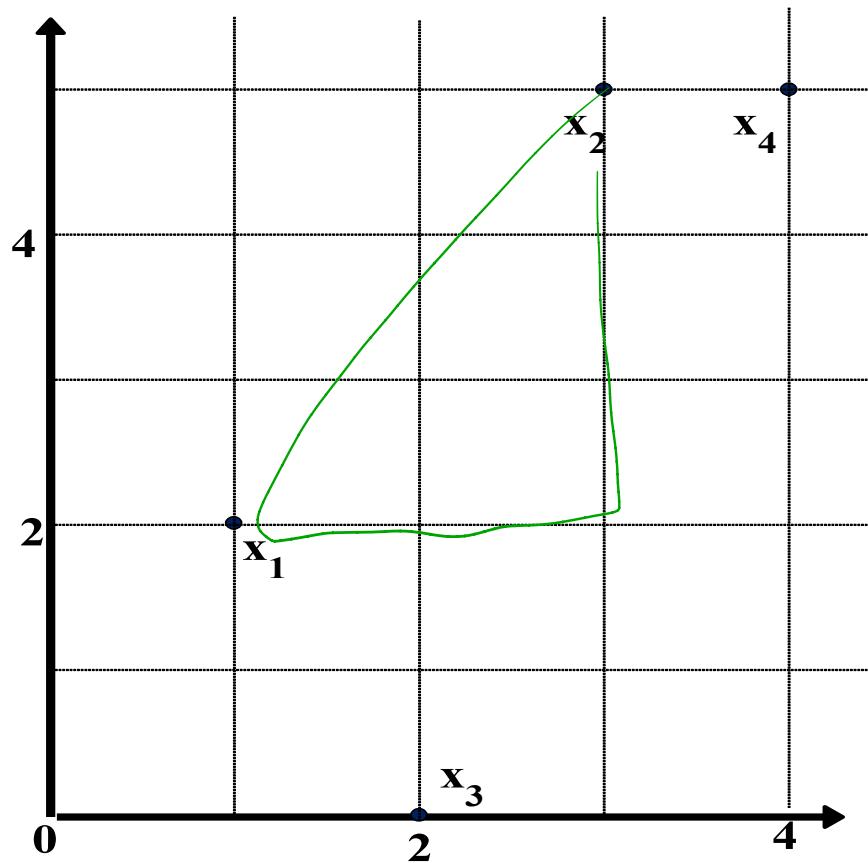
$$s_f = \frac{1}{n}(|x_{1f} - m_f| + |x_{2f} - m_f| + \dots + |x_{nf} - m_f|)$$

where

$$m_f = \frac{1}{n}(x_{1f} + x_{2f} + \dots + x_{nf})$$

- standardized measure (z-score): 
$$z_{if} = \frac{x_{if} - m_f}{s_f}$$
- Using mean absolute deviation is more robust than using standard deviation

# Example: Data Matrix and Dissimilarity Matrix



Data Matrix

point	attribute1	attribute2
$x1$	1	2
$x2$	3	5
$x3$	2	0
$x4$	4	5

ນາມລາຍການດັບພື້ນຖານ

Dissimilarity Matrix (by Euclidean Distance)

	$x1$	$x2$	$x3$	$x4$
$x1$	0			
$x2$	3.61	0		
$x3$	2.24	5.1	0	
$x4$	4.24	1	5.39	0

ເພີ່ມຂຽນຕະຫຼານນີ້ໃນນິຕ່ານ

# Distance on Numeric Data: Minkowski Distance

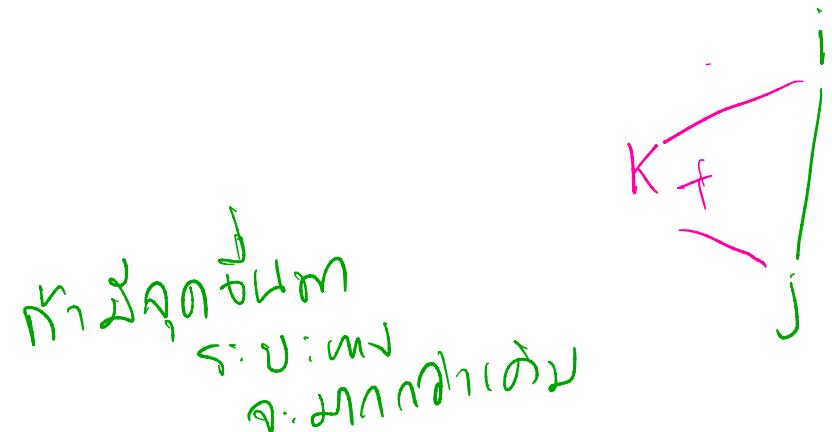
- Minkowski distance: A popular distance measure

$$d(i, j) = \sqrt[p]{|x_{i1} - x_{j1}|^p + |x_{i2} - x_{j2}|^p + \dots + |x_{il} - x_{jl}|^p}$$

where  $i = (x_{i1}, x_{i2}, \dots, x_{il})$  and  $j = (x_{j1}, x_{j2}, \dots, x_{jl})$  are two  $l$ -dimensional data objects, and  $p$  is the order (the distance so defined is also called L- $p$  norm)

- Properties

- $d(i, j) > 0$  if  $i \neq j$ , and  $d(i, i) = 0$  (Positivity)
- $d(i, j) = d(j, i)$  (Symmetry)
- $d(i, j) \leq d(i, k) + d(k, j)$  (Triangle Inequality)
- A distance that satisfies these properties is a **metric**
- Note: There are nonmetric dissimilarities, e.g., set differences



# Special Cases of Minkowski Distance

---

- $p = 1$ : ( $L_1$  norm) Manhattan (or city block) distance
  - E.g., the Hamming distance: the number of bits that are different between two binary vectors

$$d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{il} - x_{jl}|$$

- $p = 2$ : ( $L_2$  norm) Euclidean distance

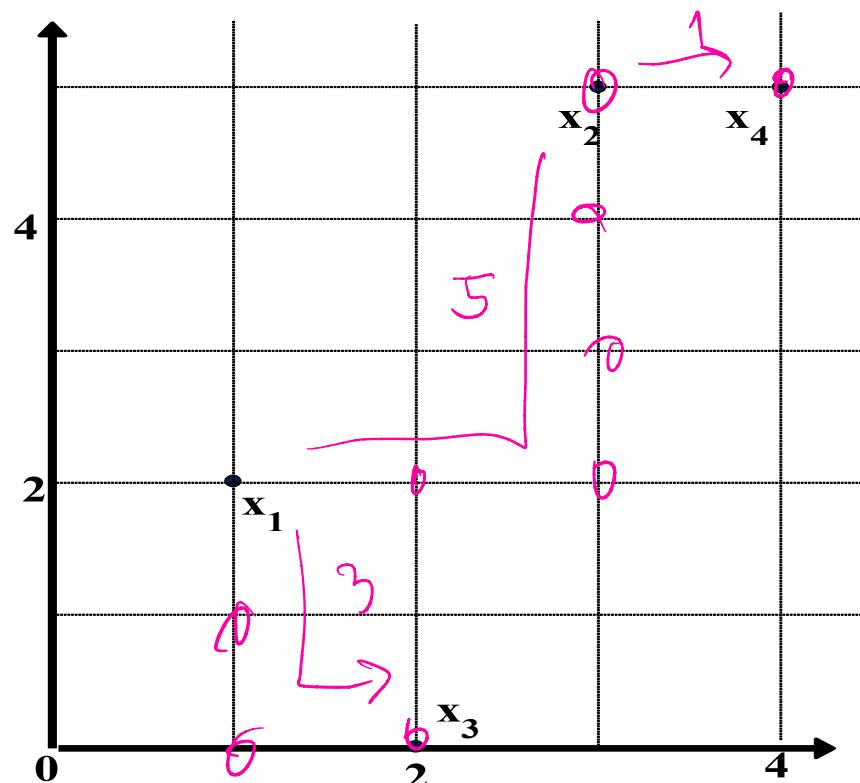
$$d(i, j) = \sqrt{|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + \dots + |x_{il} - x_{jl}|^2}$$

- $p \rightarrow \infty$ : ( $L_{\max}$  norm,  $L_\infty$  norm) “supremum” distance
  - The maximum difference between any component (attribute) of the vectors

$$d(i, j) = \lim_{p \rightarrow \infty} \sqrt[p]{|x_{i1} - x_{j1}|^p + |x_{i2} - x_{j2}|^p + \dots + |x_{il} - x_{jl}|^p} = \max_{f=1}^l |x_{if} - x_{jf}|$$

# Example: Minkowski Distance at Special Cases

point	attribute 1	attribute 2
x1	1	2
x2	3	5
x3	2	0
x4	4	5



## Manhattan ( $L_1$ )

L	x1	x2	x3	x4
x1	0			
x2	5	0		
x3	3	6	0	
x4	6	1	7	0

ເຕັມຕະຫຼານ

city block

## Euclidean ( $L_2$ )

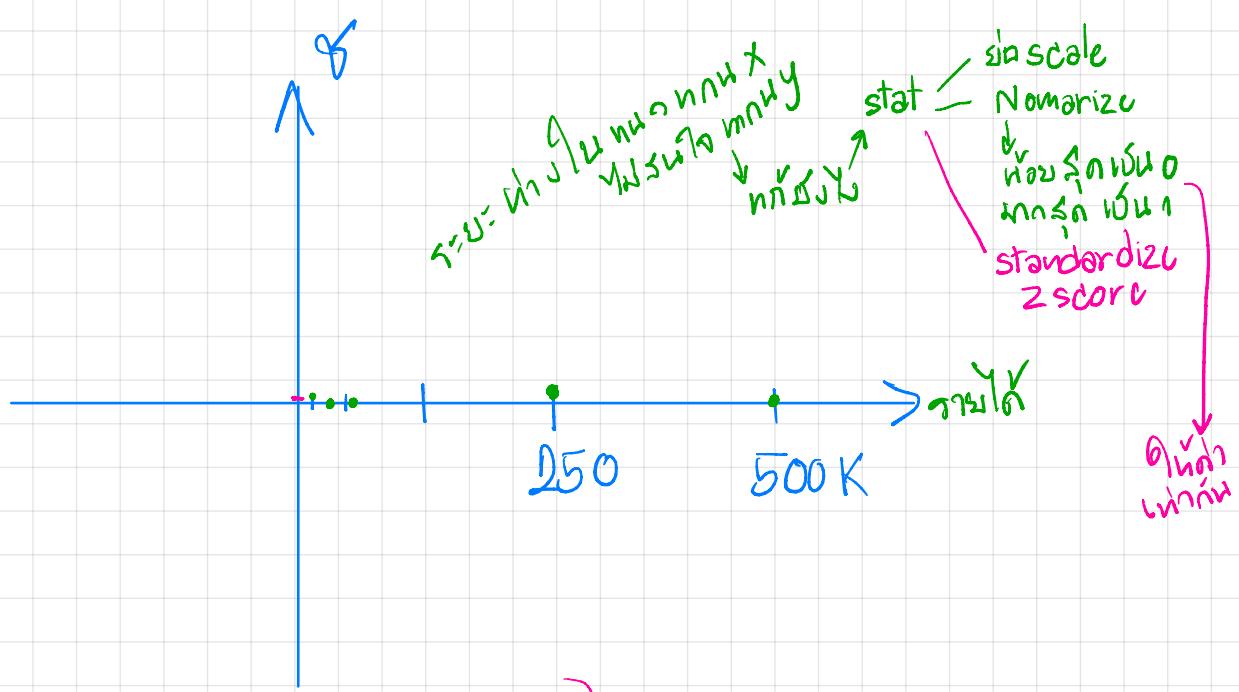
L2	x1	x2	x3	x4
x1	0			
x2	3.61	0		
x3	2.24	5.1	0	
x4	4.24	1	5.39	0

## Supremum ( $L_\infty$ )

$L_\infty$	x1	x2	x3	x4
x1	0			
x2	3	0		
x3	2	5	0	
x4	3	1	5	0

ກົດລາຍລະອຽດ

ລ	ລາຍລະອຽດ
0	40K
1	8K
1	500K
0	15K
1	250K



- Normalize  
Max ນີ້ Min ມີນ  
ກົດສູງກິດໂປ່ງ ທີ່ມີຄວາມ ພັນຍາ
- Standardize  
-1 ກົ່າ 1 Z score
- Standard deviation

ດີນທີ່  
ກົດສູງກິດໂປ່ງ ທີ່ມີຄວາມ ພັນຍາ

# Proximity Measure for Binary Attributes

- A contingency table for binary data

		Object <i>j</i>		sum
		0	1	
Object <i>i</i>	0	<i>q</i>	<i>r</i>	<i>q+r</i>
	1	<i>s</i>	<i>t</i>	<i>s+t</i>
sum	<i>q+s</i>	<i>r+t</i>	<i>p</i>	

ສຳເນົາການແບບ

- Distance measure for symmetric binary variables

$$d(i, j) = \frac{r + s}{q + r + s + t}$$

- Distance measure for asymmetric binary variables:

$$d(i, j) = \frac{r + s}{q + r + s}$$

- Jaccard coefficient (*similarity* measure for asymmetric binary variables):

$$\text{sim}_{\text{Jaccard}}(i, j) = \frac{q}{q + r + s}$$

- Note: Jaccard coefficient is the same as

(a concept discussed in Pattern Discovery)

$$\text{coherence}(i, j) = \frac{\text{sup}(i, j)}{\text{sup}(i) + \text{sup}(j) - \text{sup}(i, j)} = \frac{q}{(q + r) + (q + s) - q}$$

# Example: Dissimilarity between Asymmetric Binary Variables

Name	Gender	Fever	Cough	Test-1	Test-2	Test-3	Test-4
Jack	M	Y	N	P	N	N	N
Mary	F	Y	N	P	N	P	N
Jim	M	Y	P	N	N	N	N

- Gender is a symmetric attribute (not counted in)
- The remaining attributes are asymmetric binary
- Let the values Y and P be 1, and the value N be 0
- Distance:  $d(i, j) = \frac{r + s}{q + r + s}$

$$d(jack, mary) = \frac{0+1}{2+0+1} = 0.33$$

$$d(jack, jim) = \frac{1+1}{1+1+1} = 0.67$$

$$d(jim, mary) = \frac{1+2}{1+1+2} = 0.75$$

		Mary		
		1	0	$\Sigma_{\text{row}}$
Jack		1	2	0
		0	1	3
$\Sigma_{\text{col}}$		3	3	6

		Jim		
		1	0	$\Sigma_{\text{row}}$
Jack		1	1	2
		0	1	3
$\Sigma_{\text{col}}$		2	4	6

		Mary		
		1	0	$\Sigma_{\text{row}}$
Jim		1	1	2
		0	2	4
$\Sigma_{\text{col}}$		3	3	6

Bramper.

# Proximity Measure for Categorical Attributes

ສູ່ມະນຸຍາກົດຕະຫຼາດ

ວິທີ X

Group ເລືດ  
ອັນດາ

- Categorical data, also called nominal attributes

- Example: Color (red, yellow, blue, green), profession, etc.

- Method 1: Simple matching

- $m$ : # of matches,  $p$ : total # of variables

$$d(i, j) = \frac{p - m}{p}$$

	color	Nationality	Blood
C	Black	Thai	B
N	Blue	Thai	O
R	Blue	B	A

$$\begin{aligned} d(i, j) &= \frac{P - m}{P} \cdot \frac{m=2}{P=4} \\ &= \frac{4-2}{4} = \frac{2}{4} = \frac{1}{2} = 0.5 \end{aligned}$$

- Method 2: Use a large number of binary attributes

- Creating a new binary attribute for each of the  $M$  nominal states

dummy  $\rightarrow$  0, 1  
variable

C	N	R
Black	Thai	B
Blue	Thai	B

$c = [\text{black}, \text{blue}, \text{Green}]$

	$c$	$c\text{-block}$	$c\text{-blue}$	$c\text{-green}$
- black	1	0	1	0
- blue	0	1	0	0

Stat - Dummy Variable  
Mining - One-hot - Encoder

Now  
new  
question  
remain

# Ordinal Variables

- An ordinal variable can be discrete or continuous
  - Order is important, e.g., rank (e.g., freshman, sophomore, junior, senior)
  - Can be treated like interval-scaled
    - Replace *an ordinal variable value* by its rank:  $r_{if} \in \{1, \dots, M_f\}$
    - Map the range of each variable onto [0, 1] by replacing *i*-th object in the *f*-th variable by
- သုတေသနကြောင်း၊ ပုဂ္ဂနိုင်ရုံး၊ အမျိုးအစား၊ လူမှာ
- $$z_{if} = \frac{r_{if} - 1}{M_f - 1}$$
- ပုဂ္ဂနိုင်၊ ပုဂ္ဂနိုင်ရုံး၊ အမျိုးအစား၊ လူမှာ
- Example: freshman: 0; sophomore: 1/3; junior: 2/3; senior 1
  - Then distance:  $d(\text{freshman}, \text{senior}) = 1$ ,  $d(\text{junior}, \text{senior}) = 1/3$
  - Compute the dissimilarity using methods for interval-scaled variables

# Attributes of Mixed Type

- A dataset may contain all attribute types
  - Nominal, symmetric binary, asymmetric binary, numeric, and ordinal  
*ជំនាញបែកចុះសាស្ត្រ*
- One may use a weighted formula to combine their effects:

$$d(i, j) = \frac{\sum_{f=1}^p w_{ij}^{(f)} d_{ij}^{(f)}}{\sum_{f=1}^p w_{ij}^{(f)}}$$

- If  $f$  is numeric: Use the normalized distance
- If  $f$  is binary or nominal:  $d_{ij}^{(f)} = 0$  if  $x_{if} = x_{jf}$ ; or  $d_{ij}^{(f)} = 1$  otherwise
- If  $f$  is ordinal
  - Compute ranks  $z_{if}$  (where  $z_{if} = \frac{r_{if} - 1}{M_f - 1}$ )  
*ផ្លូវការកំណត់លទ្ធផល, Dummy ទៅការ  
ក្នុងរាយការ*
  - Treat  $z_{if}$  as interval-scaled

ຄົນໃຫຍ່  
ຂັ້ນຂາຍ

ກົດຕູມໃຫຍ່ແລກວິທະຍາ: ມອງກົມຂອງຈຸດ

# Cosine Similarity of Two Vectors

ສົກມູນລາກາຖີ່ຢູ່ໂທລາຍນໍາ  
ເວັບໄຊ

- A **document** can be represented by a bag of terms or a long vector, with each attribute recording the *frequency* of a particular term (such as word, keyword, or phrase) in the document

ອີງຕາມດີເຊື່ອ: ລັບກ່າຍທົກຕາມເຊີ້ນຫຼັບເຈັບປັດ

Document	team	coach	hockey	baseball	soccer	penalty	score	win	loss	season
Document1	5	0	3	0	2	0	0	2	0	0
Document2	3	0	2	0	1	1	0	1	0	1
Document3	0	7	0	2	1	0	0	3	0	0
Document4	0	1	0	0	1	2	2	0	3	0

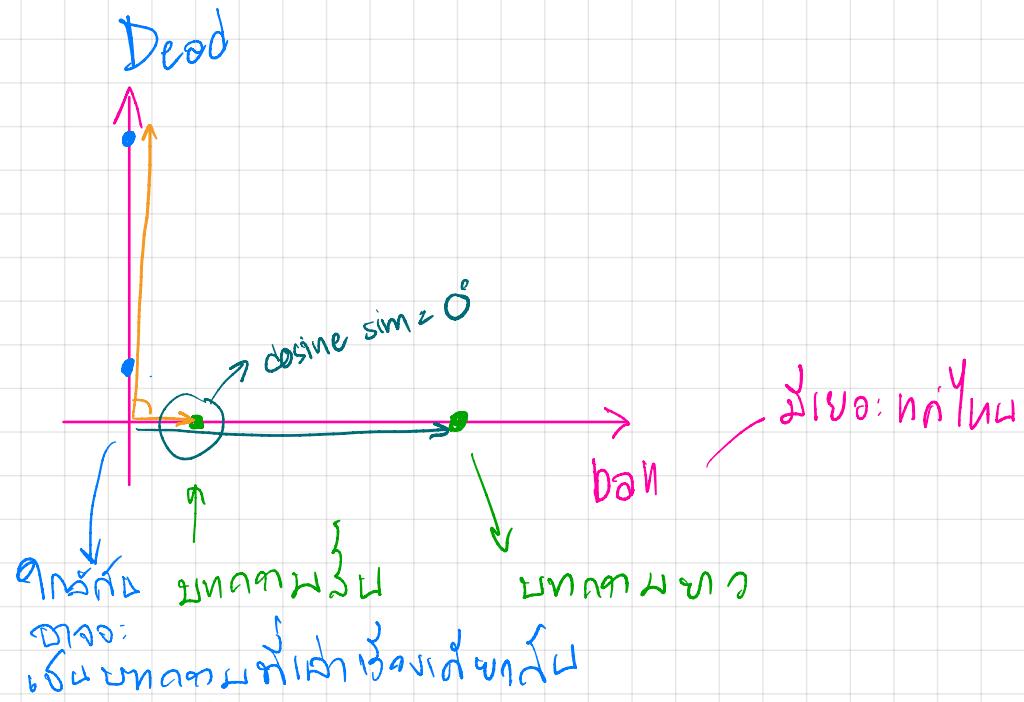
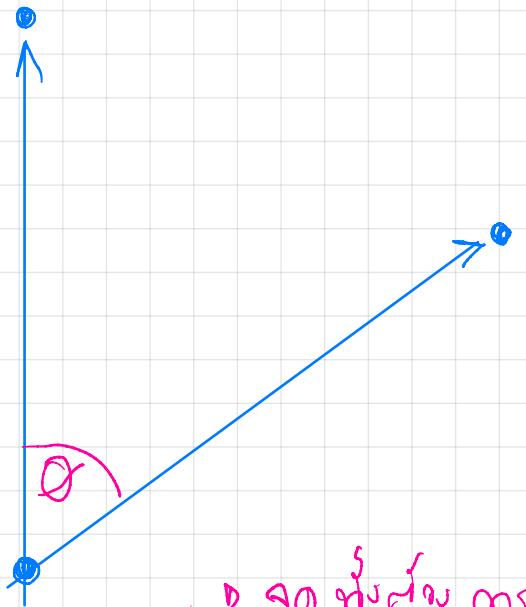
- Other vector objects: Gene features in micro-arrays
- Applications: Information retrieval, biologic taxonomy, gene feature mapping, etc.
- Cosine measure: If  $d_1$  and  $d_2$  are two vectors (e.g., term-frequency vectors), then

$$\cos(d_1, d_2) = \frac{d_1 \bullet d_2}{\|d_1\| \times \|d_2\|}$$

ວິນດອນ vec  $d_1$   
= ຢາມວິນດອນ

where  $\bullet$  indicates vector dot product,  $\|d\|$ : the length of vector  $d$

Dosine Sim



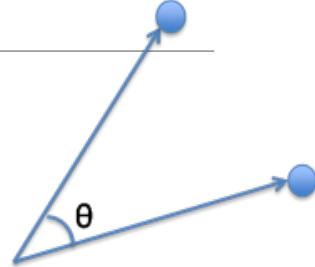
ระยะห่าง cosine distance ↴  
m

# Example: Calculating Cosine Similarity

- Calculating Cosine Similarity:

$$\cos(d_1, d_2) = \frac{d_1 \bullet d_2}{\|d_1\| \times \|d_2\|}$$

$$sim(A, B) = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|}$$



where • indicates vector dot product, ||d||: the length of vector d

- Ex: Find the **similarity** between documents 1 and 2. / ឧទេរស ការណ៍

$$d_1 = (5, 0, 3, 0, 2, 0, 0, 2, 0, 0)$$

$$d_2 = (3, 0, 2, 0, 1, 1, 1, 0, 1, 0, 1)$$

feature vector

ពីរប្រព័ន្ធដែល  
ត្រូវការគុណភាព  
ខ្លា...

- First, calculate vector dot product

$$d_1 \bullet d_2 = 5 \times 3 + 0 \times 0 + 3 \times 2 + 0 \times 0 + 2 \times 1 + 0 \times 1 + 0 \times 1 + 2 \times 1 + 0 \times 0 + 0 \times 1 = 25$$

- Then, calculate  $\|d_1\|$  and  $\|d_2\|$

$$\|d_1\| = \sqrt{5 \times 5 + 0 \times 0 + 3 \times 3 + 0 \times 0 + 2 \times 2 + 0 \times 0 + 0 \times 0 + 2 \times 2 + 0 \times 0 + 0 \times 0} = 6.481$$

$$\|d_2\| = \sqrt{3 \times 3 + 0 \times 0 + 2 \times 2 + 0 \times 0 + 1 \times 1 + 1 \times 1 + 0 \times 0 + 1 \times 1 + 0 \times 0 + 1 \times 1} = 4.12$$

- Calculate cosine similarity:  $\cos(d_1, d_2) = 25 / (6.481 \times 4.12) = 0.94$

# **Chapter 2. Getting to Know Your Data**

---

- Data Objects and Attribute Types
- Basic Statistical Descriptions of Data
- Data Visualization
- Measuring Data Similarity and Dissimilarity
- Summary



*Preprocessing*

# Summary

*Two main parts*



*Introducing Data Mining features*

- Data attribute types: nominal, binary, ordinal, interval-scaled, ratio-scaled
- Many types of data sets, e.g., numerical, text, graph, Web, image.
- Gain insight into the data by:
  - Basic statistical data description: central tendency, dispersion, graphical displays
  - Data visualization: map data onto graphical primitives
  - Measure data similarity
- Above steps are the beginning of data preprocessing
- Many methods have been developed but still an active area of research

# References

---

- W. Cleveland, Visualizing Data, Hobart Press, 1993
- T. Dasu and T. Johnson. Exploratory Data Mining and Data Cleaning. John Wiley, 2003
- U. Fayyad, G. Grinstein, and A. Wierse. Information Visualization in Data Mining and Knowledge Discovery, Morgan Kaufmann, 2001
- L. Kaufman and P. J. Rousseeuw. Finding Groups in Data: an Introduction to Cluster Analysis. John Wiley & Sons, 1990.
- H. V. Jagadish et al., Special Issue on Data Reduction Techniques. Bulletin of the Tech. Committee on Data Eng., 20(4), Dec. 1997
- D. A. Keim. Information visualization and visual data mining, IEEE trans. on Visualization and Computer Graphics, 8(1), 2002
- D. Pyle. Data Preparation for Data Mining. Morgan Kaufmann, 1999
- S. Santini and R. Jain," Similarity measures", IEEE Trans. on Pattern Analysis and Machine Intelligence, 21(9), 1999
- E. R. Tufte. The Visual Display of Quantitative Information, 2<sup>nd</sup> ed., Graphics Press, 2001
- C. Yu, et al., Visual data mining of multimedia data for social and behavioral studies, Information Visualization, 8(1), 2009

