



# 2110773

## Data Mining

### Chapter2:

## Data Preprocessing

Data Cleaning

noise

missing data

Data Integration -

inconsistency សារុក្រាស និងនិមួយ

- ▶ GARBAGE IN → GARBAGE OUT
- ▶ IMPORTANT & TIME-CONSUMING TASK IN KDD
- ▶ PRACTICE IS EVERYTHING

▶ រត. លោក លីម ពិយាយកន្តឺ



# Types of Dataset



## Record

Relational records  
Document data: text documents  
Transaction



## Graph and network

World Wide Web  
Social or information networks  
Molecular Structures



## Others

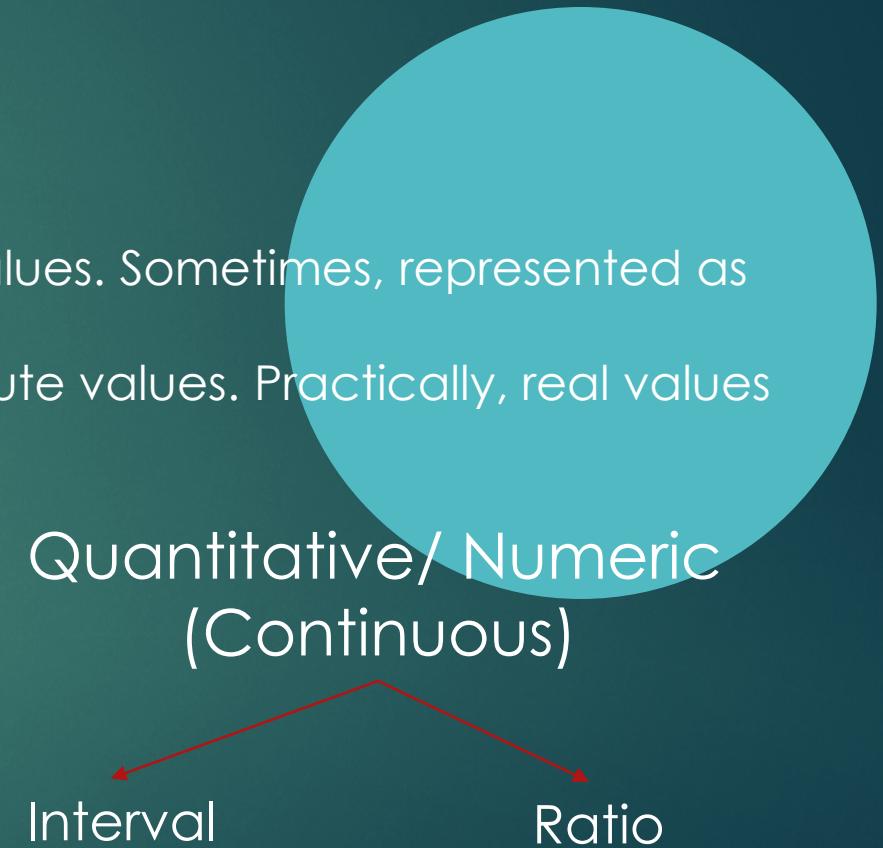
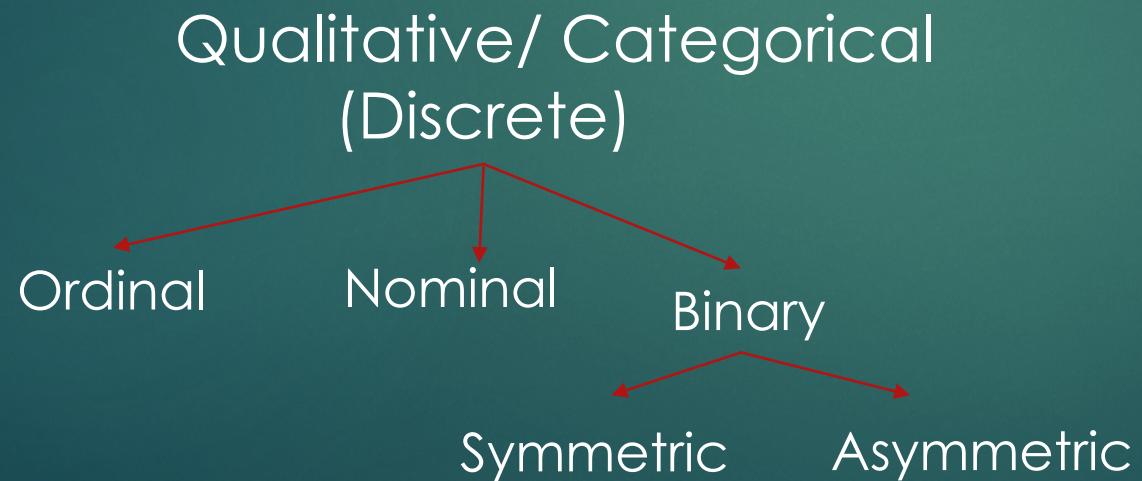
Image  
Video data: sequence of images  
Temporal/ Time-series  
Spatial data: maps

# Data Object

- ▶ Data sets are made up of data objects.
- ▶ A **data object** represents an entity. For examples:
  - ▶ medical database: patients, treatments
  - ▶ 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.
- ▶ Attribute (or **dimension, feature, variable**): a data field, representing a characteristic or feature of a data object, e.g., *customer \_ID, name, address, phone*

# Attribute Data Types

1. Qualitative/ Quantitative
2. Categorical/ Numeric
3. Discrete/ Continuous
  - Discrete: Has only a finite or countably infinite set of values. Sometimes, represented as integer variables
  - Continuous: Has real numbers (floating-point) as attribute values. Practically, real values can only be measured.



# Attribute Types

- ▶ **Nominal:** categories, states, or “names of things”. Categories cannot be compared
- ▶ **Binary:** Nominal attribute with only 2 states (0 and 1)
  - ▶ *Symmetric binary:* both outcomes equally important
  - ▶ *Asymmetric binary:* outcomes not equally important. Convention: assign 1 to most important outcome (e.g., covid19 positive)
- ▶ **Ordinal:** Values have a meaningful order (ranking) but magnitude between successive values is not known.  
Categories with an implied order
- ▶ **Quantity** (integer or real-valued)
- ▶ **Interval**
  - ▶ Measured on a scale of **equal-sized units**
  - ▶ Values have order
  - ▶ No true zero-point
- ▶ **Ratio**
  - ▶ Inherent **zero-point**
  - ▶ 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°).

# Data Type Examples

6

2110773-2 2/66

Data Type	Examples
Nominal	color, bloodType, zipCode, ID#, occupation, political party
Ordinal	medal, satisfaction, grade, frequency, academic ranking
Binary- symmetric	gender
Binary- asymmetric	labTest
Interval	celcius, farenheit, pH,
Ratio	kelvin, exam score, weight, height, pulse, monetary quantities

**Interval Data:** No true zero, differences (subtraction) are interpretable.

Data can be added/ subtracted at interval scale but nonsense be multiplied/ divided.

Ex. If a day's temperature in celcius/ farenheit is twice than the other day,  
we cannot say that one day is twice as hot as another day.

**Ratio Data:** True zero exists. Zero means none of that variable value, e.g. zero kelvin means no heat.  
The ratio of two measurements has a meaningful interpretation.

\*\* A scale is an ordered set of values, continuous or discrete, or a set of categories to which an attribute is mapped.

%	Adverb of Frequency	Example
100%	<b>Always</b>	I always study after class
90%	<b>Usually</b>	I usually walk to work
80%	<b>Normally / Generally</b>	I normally get good marks
70%	<b>Often / Frequently</b>	I often read in bed at night
50%	<b>Sometimes</b>	I sometimes sing in the shower
30%	<b>Occasionally</b>	I occasionally go to bed late
10%	<b>Seldom</b>	I seldom put salt on my food
5%	<b>Hardly ever / Rarely</b>	I hardly ever get angry
0%	<b>Never</b>	Vegetarians never eat meat

# Scales of Measurement

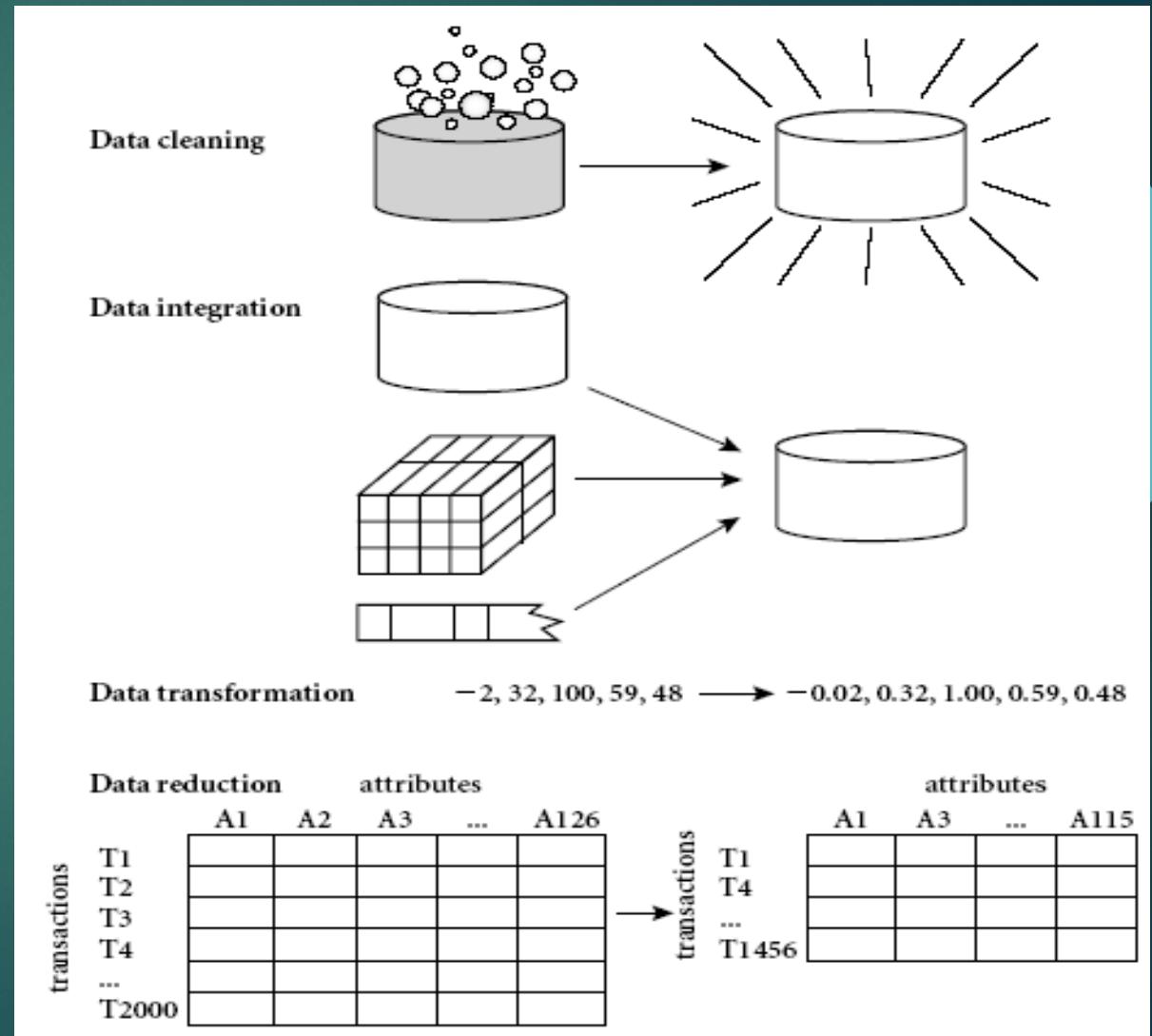
Data	Nominal	Ordinal	Interval	Ratio
Labeled				
		Order		
			Measurable Difference	
				True Zero Starting Point

# Survey

1. How old are you? \_\_\_\_\_ years
  
2. Are you:    Male    Female
  
3. How much do you spend on groceries each week? \_\_\_\_\_ Baht
  
4. How many cups of coffee do you buy in a week? \_\_\_\_\_
  
5. Which type of coffee do you like most?  
Latte              Espresso              Cappuccino              Americano
  
6. How likely are you to buy more than a cup of coffee per day?  
Very Likely        Likely        Not Likely        Very Unlikely

# Data Preprocessing

- ▶ Data Cleaning
- ▶ Data Integration
- ▶ Data Transformation
- ▶ Data Reduction



# Data Cleaning

- ▶ Fill in missing data
- ▶ Smooth noisy data- random error or variance in a measured variable
- ▶ Identify or remove outliers
- ▶ Resolve inconsistencies
  - ▶ Same name means differently (BL= blue/ black)
  - ▶ Different names appear the same (Bill vs. Williams)
  - ▶ Inappropriate values (Male-Pregnant; born Feb 29, 2562; age=41 birthday=28/08/2010)
  - ▶ Due to inconsistent Unit of Measure

# Missing Data

## ► Various reasons:

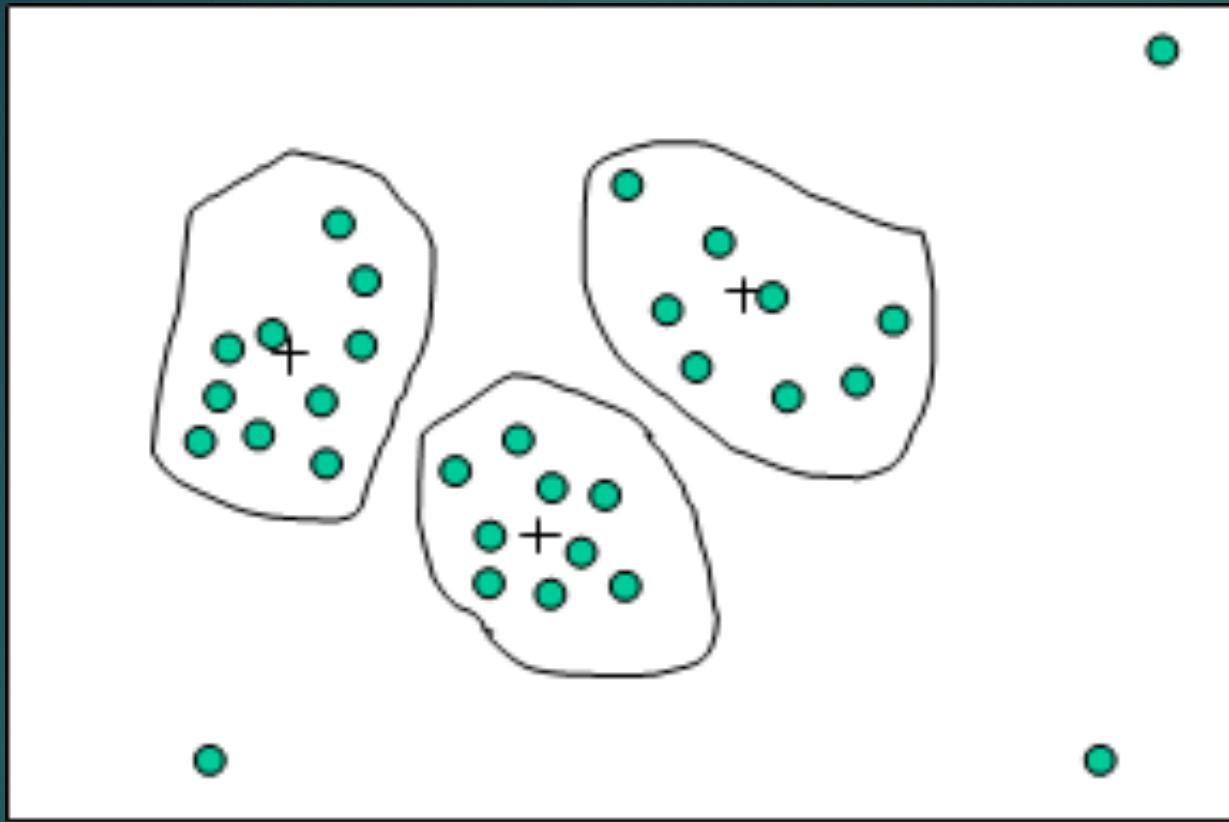
- ▶ truly missed/ impossible to always have a value
- ▶ Intentional (disguised missing data)
- ▶ not measured due to no equipment or not able to measure in the past
- ▶ Inconvenient, expensive

## ► Some methods

- ▶ Leave as is, however, some algo can't deal w/ missing values and the program may refuse to continue or lead to inaccurate results
- ▶ Remove the instance with missing value (e.g. in case of huge dataset or missing class label)
- ▶ A global constant, e.g. 999,999 (valid values are much smaller) or -1 (valid values are non-negative). Watch out for zeros as some features can use this as the boolean representation! or “unknown” can be treated as a new class ?!
- ▶ Imputing :
  - ❖ Attribute mean/median (Numerical variables); mode (Categorical variables)
  - ❖ Substitute w/ valid values of a certain feature e.g. fill in the seasonal averages of temperature for a certain location for missing temperature values given a date
  - ❖ Model-based/ inference-based: Regression, Decision Tree, k-nearest neighbor, Bayesian ...)

# Noisy Data

- ▶ Random error or variance in a measured variable
  - ▶ Regression- smooth by fitting the data into regression functions
- ▶ Outliers are noisy data or data points inconsistent with the majority of data, e.g. one's age = 200 year, height=3 metre, widely deviated points
  - ▶ Detect and remove outliers- Clustering
  - ▶ Truncate outliers- Bell curve, Box plots



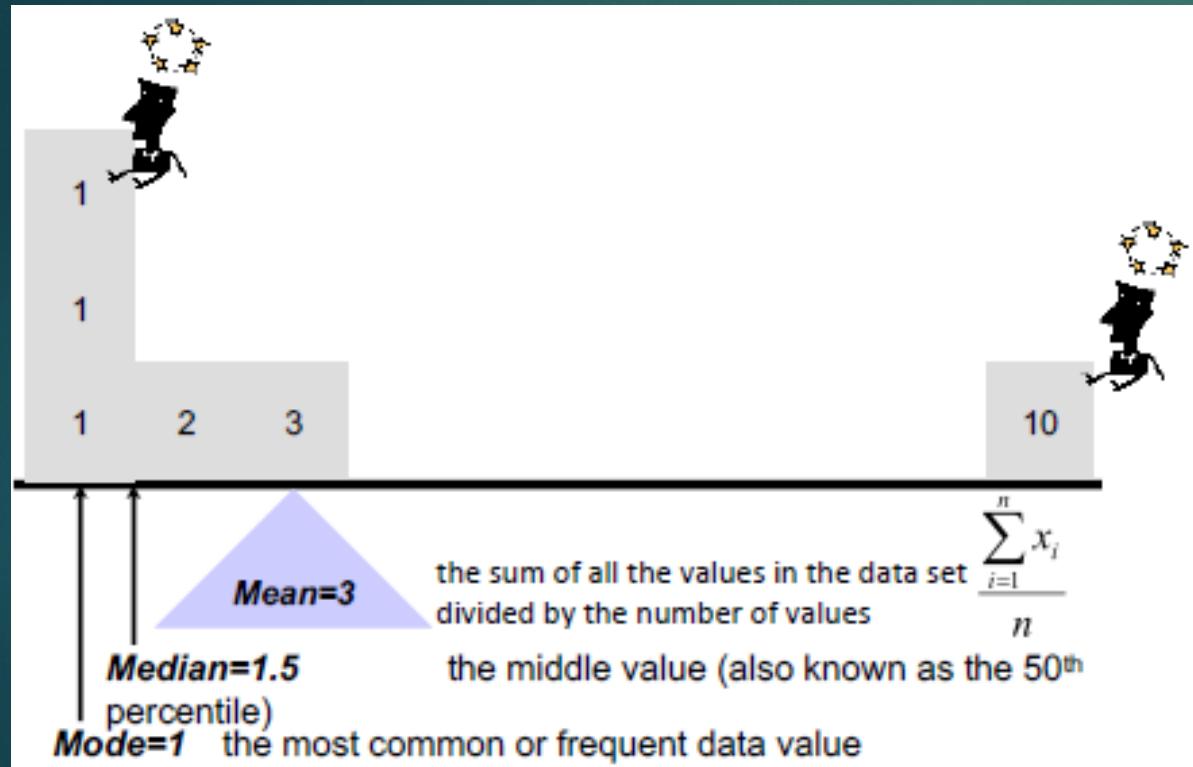
# Clustering

# Data Distribution

15

2110773-2 2/66

## 1. Central Tendency/ Center



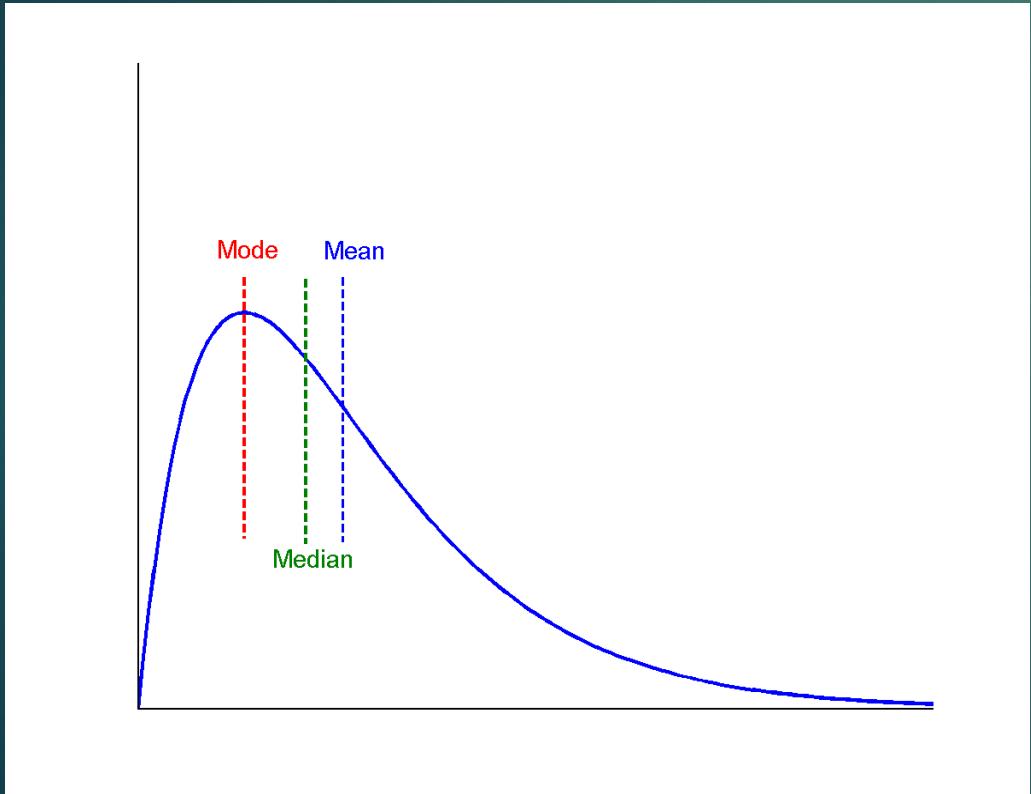
## 2. Spread/ Dispersion

Measure	Definition
<b>Range</b>	the difference between the maximum and minimum data values
<b>Interquartile Range</b>	the difference between the 25th and 75th percentiles
<b>Variance</b>	a measure of dispersion of the data around the mean
<b>Standard Deviation</b>	a measure of dispersion expressed in the same units of measurement as your data (the square root of the variance)

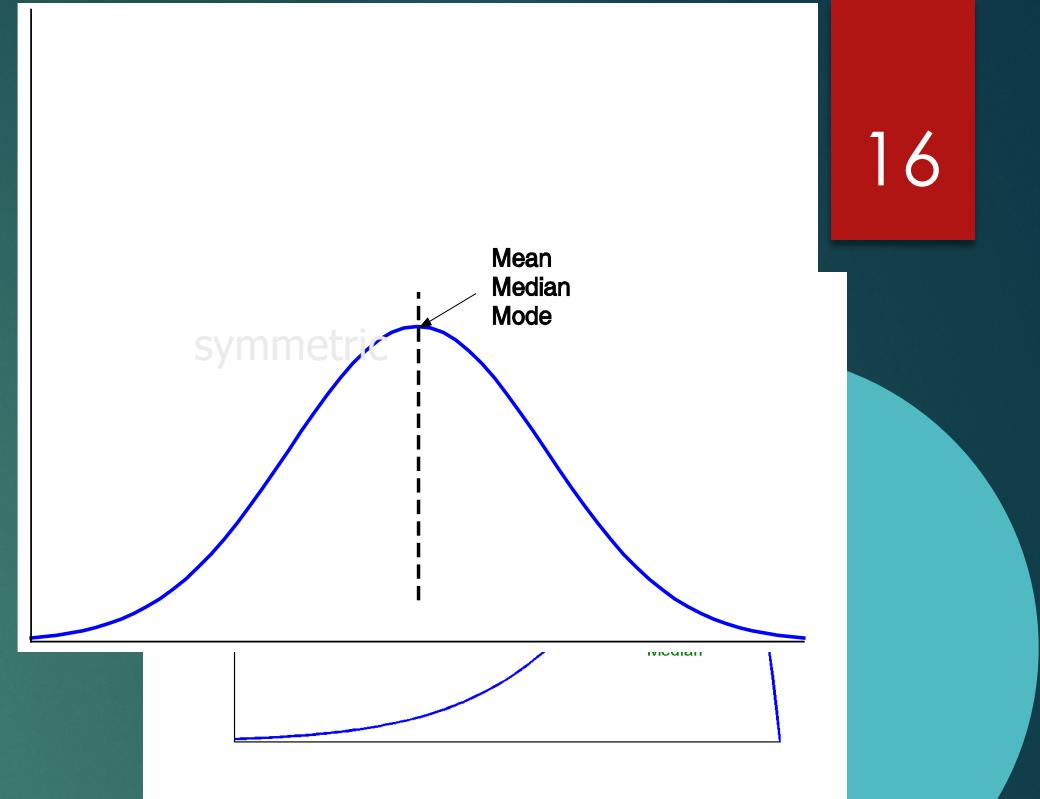
Measure of Central Tendency (Representative value):  
Mean, Median, Mode

# Symmetric vs. Skewed Data

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



positively skewed



negatively skewed

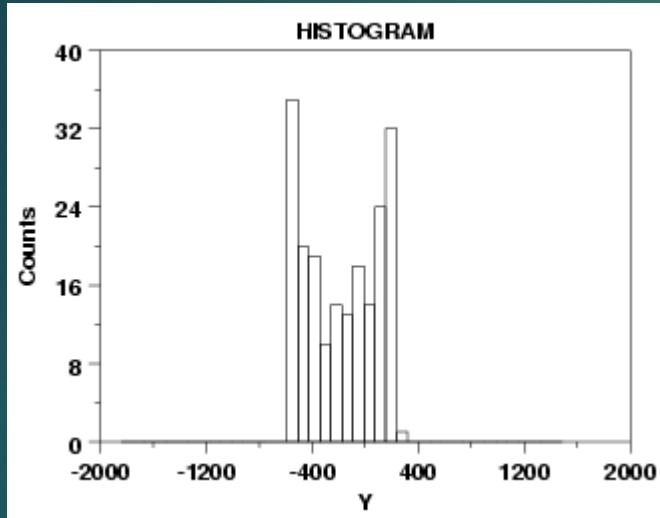
Type of Variable	Best measure of central tendency
Nominal	Mode
Ordinal	Median
Interval/Ratio (not skewed)	Mean
Interval/Ratio (skewed)	Median

When to use Mean, Median, Mode

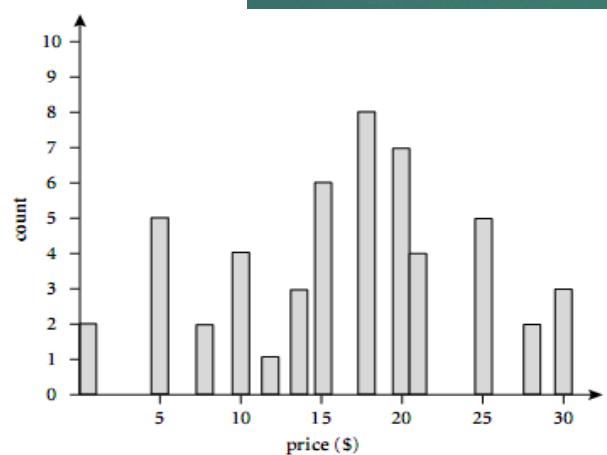
# Graphical Displays of Distribution

18

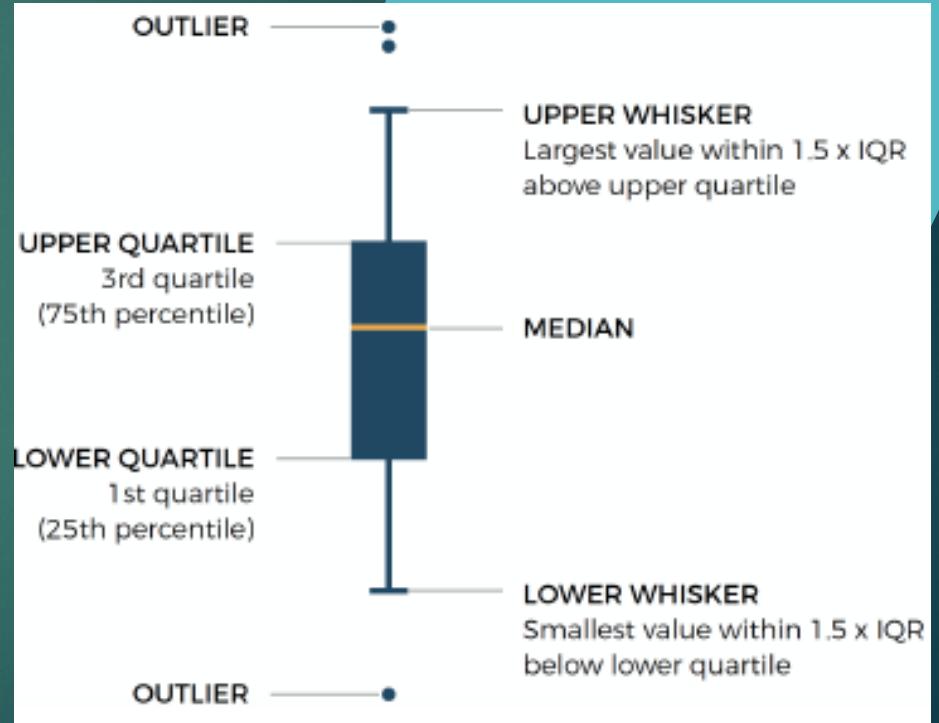
- ▶ Histogram Graph display of frequencies, shown as bars with numeric values on X axis

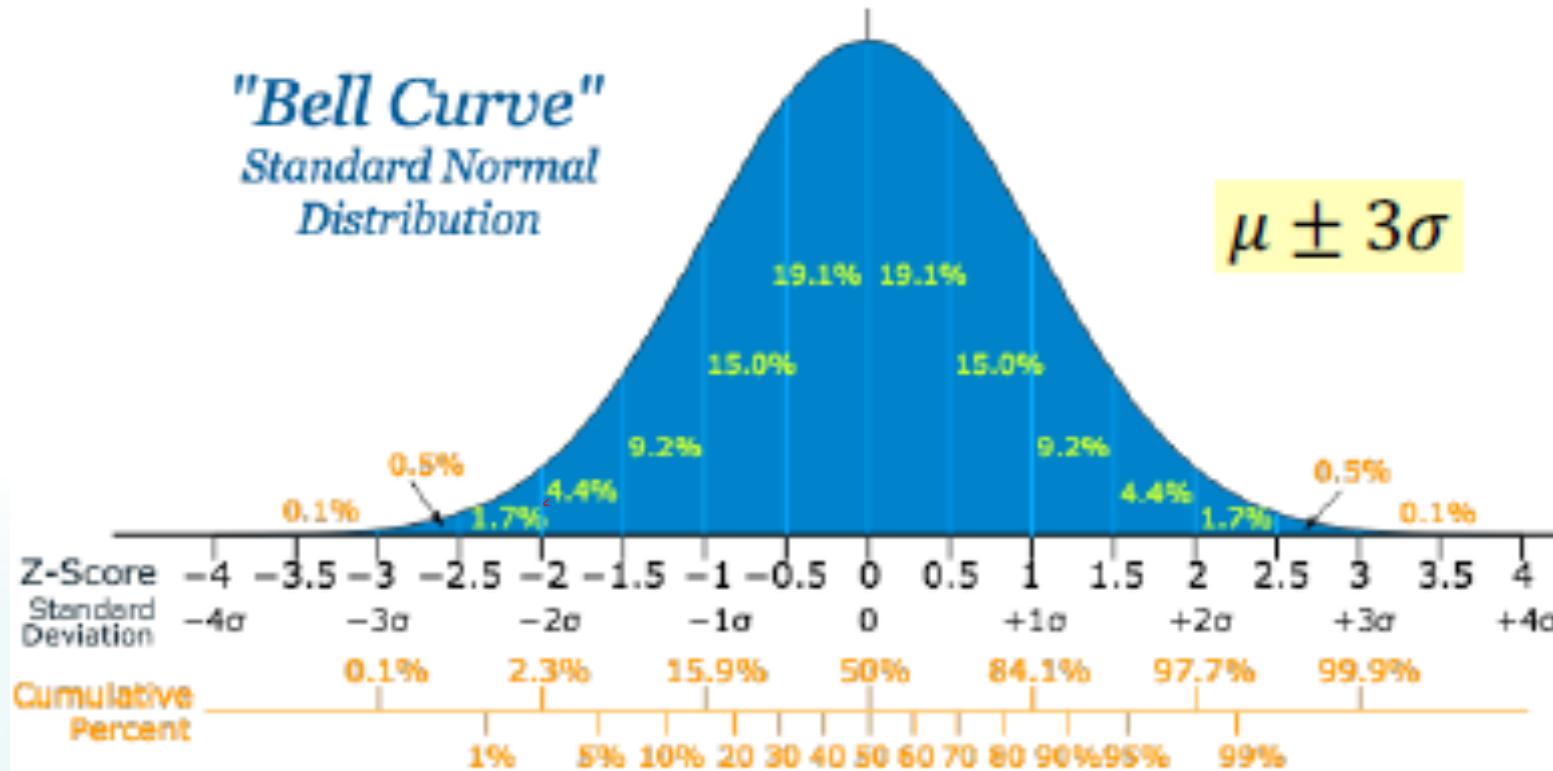


## Singleton Histogram



- ## ► Box plots





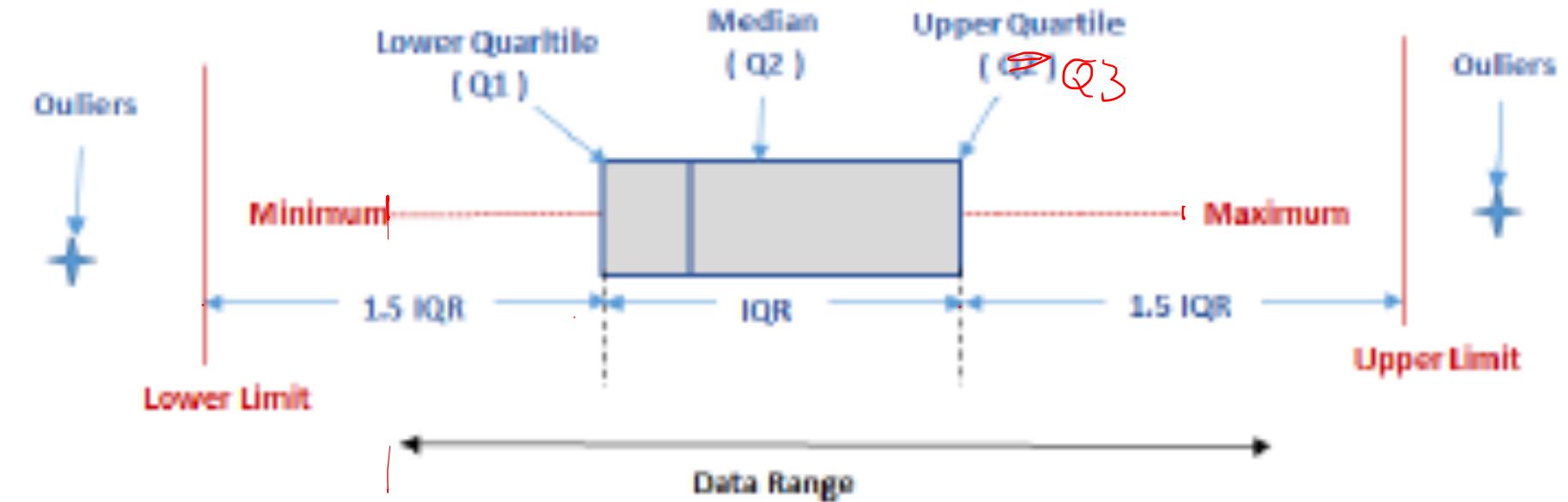
Truncate Outliers:  
Bell Curve

Variance and standard deviation  
(sample:  $s$ , population:  $\sigma$ )

Standard deviation is the square root of variance

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

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



Truncate Outliers:  
Box Plots

# Interquartile Range

- ▶ **IQR is a measure of spread indicating where the bulk of the values lie.**
  - ❖ **Quartiles:**  $Q_1$  (25<sup>th</sup> percentile),  $Q_3$  (75<sup>th</sup> percentile)
  - ❖ **Inter-quartile range:**  $IQR = Q_3 - Q_1$
  - ❖ **Five number summary:** min,  $Q_1$ , median,  $Q_3$ , max
  - ❖ **Boxplot:** ends of the box are the quartiles; median is marked; add whiskers, and plot outliers individually
  - ❖ **Outlier:** usually, a value higher/lower than  $1.5 \times IQR$

# IQR Calculation

## Odd set of numbers

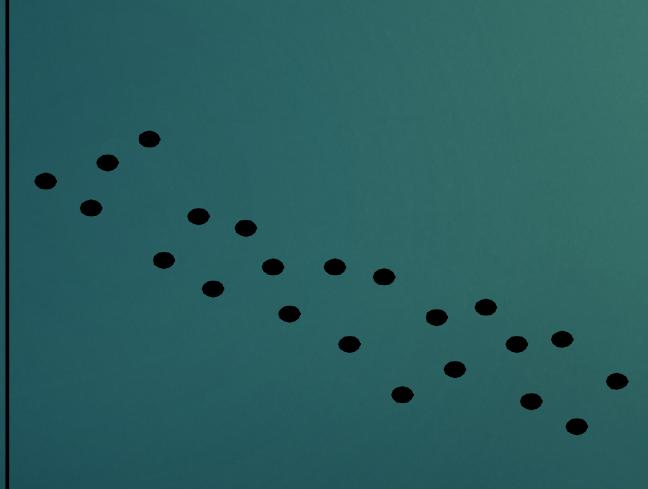
- ▶ Step 1: **Put the numbers in order.**  
1, 2, 5, 6, 7, 9, 12, 15, 18, 19, 27.
- ▶ Step 2: **Find the median.**  
1, 2, 5, 6, 7, **9**, 12, 15, 18, 19, 27.
- ▶ Step 3: **Place parentheses around the numbers above and below the median.**  
Not necessary **statistically**, but it makes Q1 and Q3 easier to spot.  
 $(1, 2, 5, 6, 7), 9, (12, 15, 18, 19, 27)$ .
- ▶ Step 4: **Find Q1 and Q3**  
Think of Q1 as a median in the lower half of the data and think of Q3 as a median for the upper half of data.  
 $(1, 2, \textbf{5}, 6, 7), \textbf{9}, (12, 15, \textbf{18}, 19, 27)$ . Q1 = 5 and Q3 = 18.
- ▶ Step 5: **Subtract Q1 from Q3 to find the interquartile range.**  
 $18 - 5 = 13$ .

## Even set of numbers

- ▶ Step 1: **Put the numbers in order.**  
3, 5, 7, 8, 9, 11, 15, 16, **20**, 21.
- ▶ Step 2: **Make a mark in the center of the data:**  
3, 5, 7, 8, 9, | 11, 15, 16, 20, 21.
- ▶ Step 3: **Place parentheses around the numbers above and below the mark you made in Step 2—it makes Q1 and Q3 easier to spot.**  
 $(3, 5, 7, 8, 9), | (11, 15, 16, 20, 21)$ .
- ▶ Step 4: **Find Q1 and Q3**  
Q1 is the median (the middle) of the lower half of the data, and Q3 is the median (the middle) of the upper half of the data.  
 $(3, 5, \textbf{7}, 8, 9), | (11, 15, \textbf{16}, 20, 21)$ . Q1 = 7 and Q3 = 16.
- ▶ Step 5: **Subtract Q1 from Q3.**  
 $16 - 7 = 9$ .

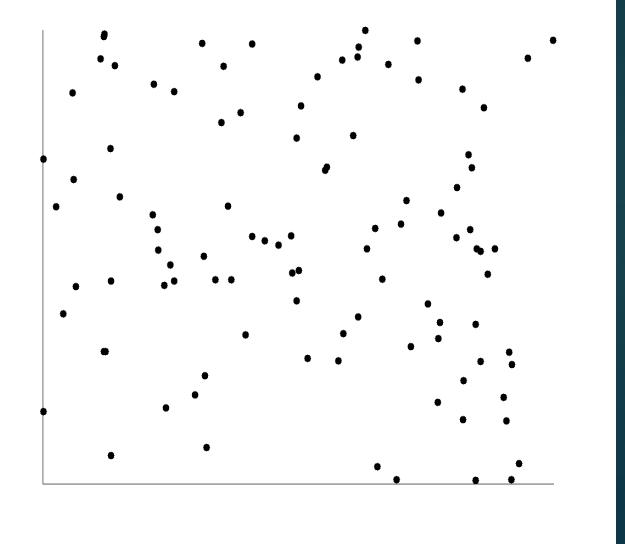
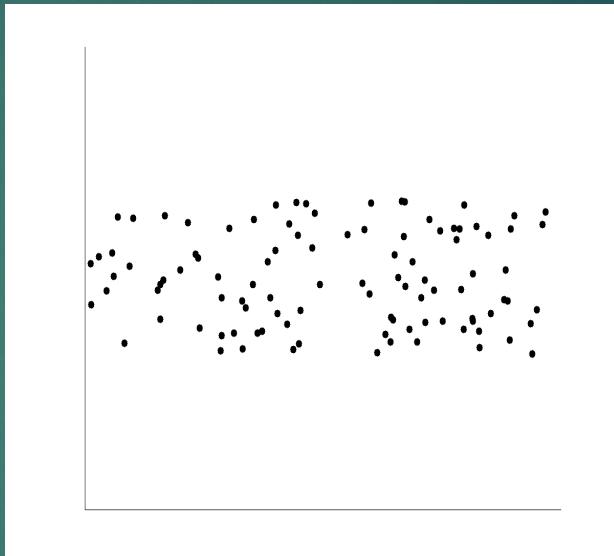
# Correlated Data

Positively



Negatively

Uncorrelated Data



# Regression

- ▶ Linear Regression

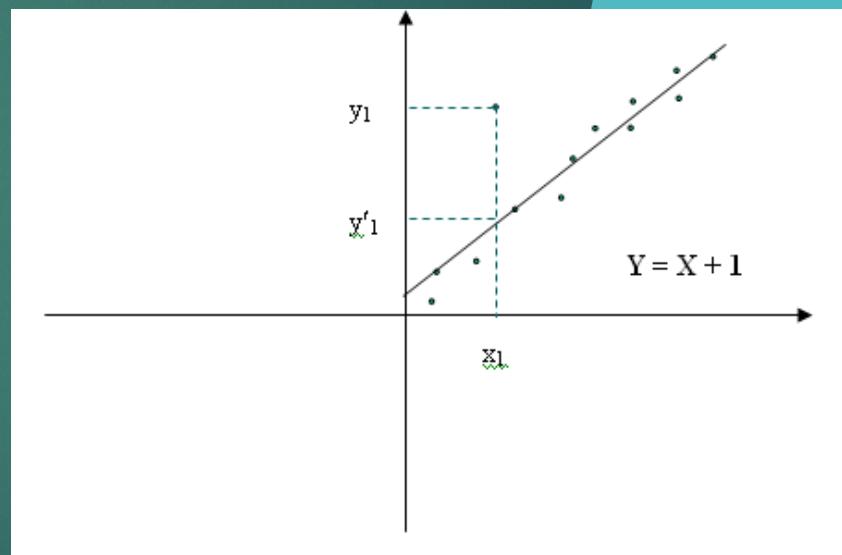
$$Y = \alpha + \beta X$$

- ▶ Multiple Linear Regression

$$Y = b_0 + b_1 X_1 + b_2 X_2 + \dots + b_m X_m$$

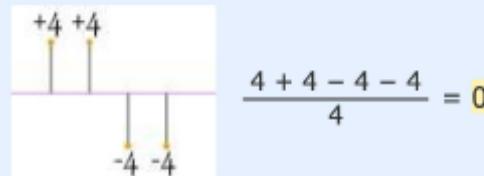
- ▶ Smooth out noise

- ▶ Fill in missing value

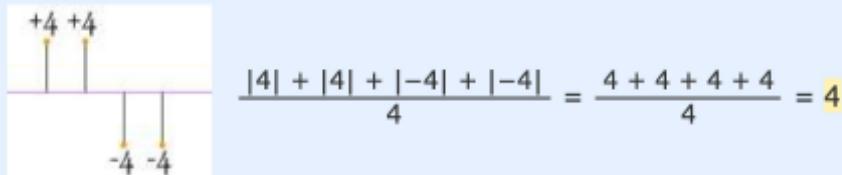


### \*Footnote: Why square the differences?

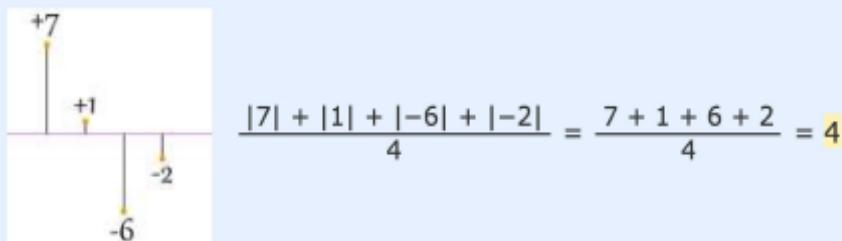
If we just add up the differences from the mean ... the negatives cancel the positives:



So that won't work. How about we use [absolute values](#)?

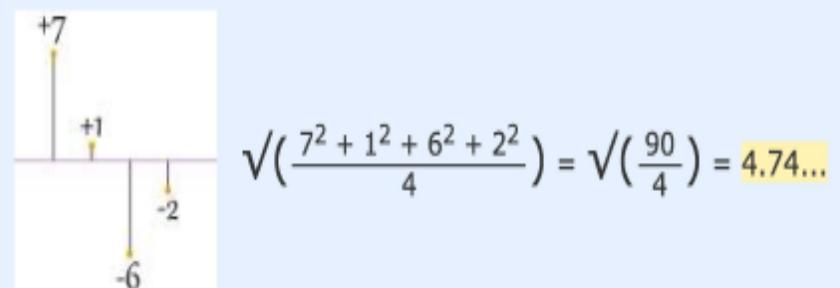
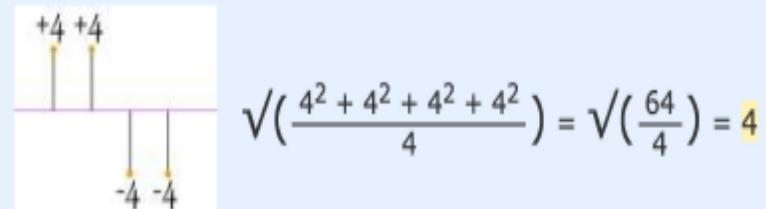


That looks good (and is the [Mean Deviation](#)), but what about this case:



Oh No! It also gives a value of 4, Even though the differences are more spread out.

So let us try squaring each difference (and taking the square root at the end):



That is nice! The Standard Deviation is bigger when the differences are more spread out ... just what we want.

In fact this method is a similar idea to [distance between points](#), just applied in a different way.

And it is easier to use algebra on squares and square roots than absolute values, which makes the standard deviation easy to use in other areas of mathematics.

# Data Integration

- ▶ Integration of multiple databases
- ▶ Handle data inconsistencies, majorly due to
  - ▶ Unit of Measure differences
  - ▶ Value differences
- ▶ Manage data redundancies
  - ▶ Correlation analysis

# Data Transformation<sub>1</sub>

ການມຳກັງຂໍ້ມູນ

27

2110773-2 2/66

- \* Many models implemented in Sklearn might perform poorly if the numeric features do not more or less follow a standard Gaussian (normal) distribution. Except for tree-based models, the objective function of Sklearn algorithms assumes the features follow a **normal distribution**. sklearn ສົມຜະລິດຕະໂປຣດິກ, ລາຍງືນຈະສົມຜະລິດຕະໂປຣດິກ ທີ່ມີຄວາມຖຸກຕະຫຼອດກົດໝາຍ
- \* **Standardization** or **Scaling** numeric features is required for distance-based algorithms e.g. SVM, kNN to achieve better results
- \* Scaling and Normalization are very similar and confusing, sometimes used interchangeably
- \* what's the difference?

# Data Transformation<sub>2</sub>

## ► Standardization (Scaling) / Normalization

- **numeric variables** are transformed in both cases
  - ❖ Min-max Scaling using min and max values of distribution → `MinMaxScaler()`
  - ❖ Z-score using variance and mean → `StandardScaler()`
  - ❖ Sigmoidal
  - ❖ Log Transforms → `PowerTransformer()`

## ► Data Type Conversion

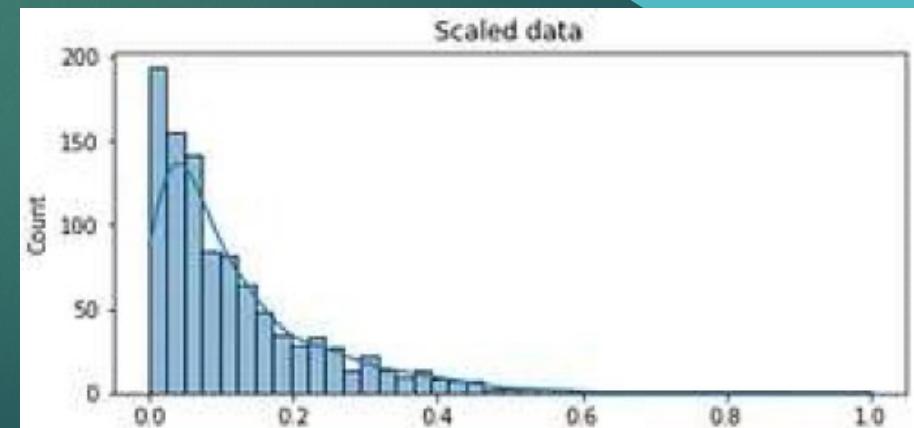
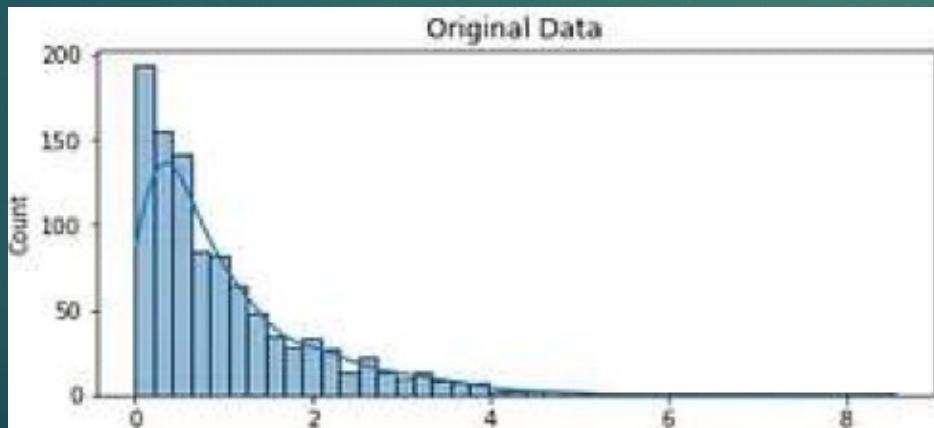
# Scaling

ສຳຄັນຈຳກົດປະໂຫວາດຮູ່ກາຕີ້ວ

- ▶ Scaling is a method to standardize the range of independent variables or features of data.
- ▶ Change **range** of data to same scale, e.g. 0-1, 0-100
- ▶ Applied in distance-based algorithms, e.g. SVM, kNN  
        ເລືອດທີ່ວ່າຈະສຳເນົາດ້ວຍກົດປະໂຫວາດຮູ່ກາຕີ້ວ (k-mean)
- ▶ Same importance for a change of “1” in any numeric features
- ▶ By scaling, variables are compared on equal footing

ເພື່ອສຳເນົາດ້ວຍກົດປະໂຫວາດຮູ່ກາຕີ້ວ

ກົດປະໂຫວາດຮູ່ກາຕີ້ວ ສຳເນົາດ້ວຍກົດປະໂຫວາດຮູ່ກາຕີ້ວ



# Scaling: case study

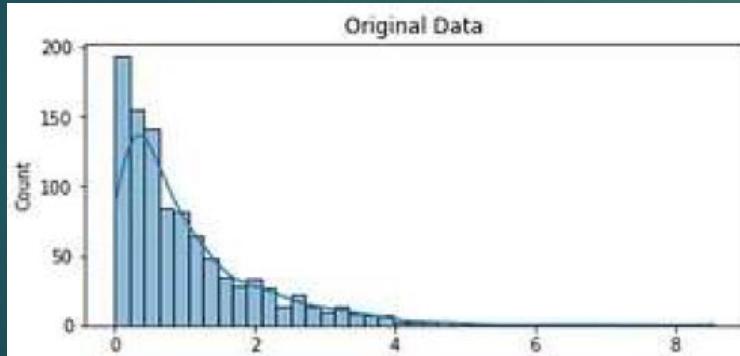
- ▶ Purpose: Change the values of numeric columns to a common scale
- ▶ Example: *age(x1)* ranges 0-100; *income(x2)* ranges 0-1,000,000
- ▶ Observing *income* will influence the result more due to its larger value
- ▶ Example of two deep neural network models w/ and w/o data scaling, accuracy = 88.93%, 48.80% respectively

without scaling

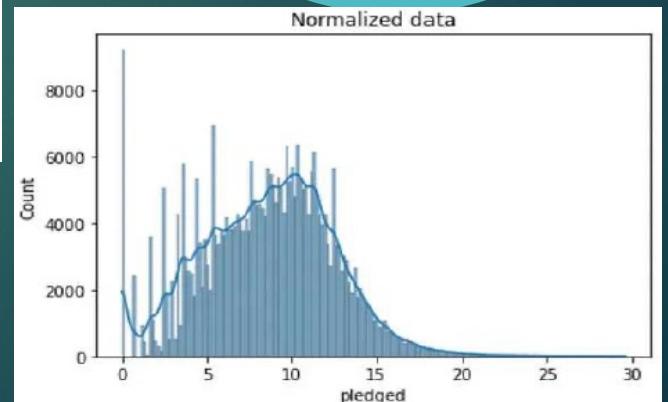
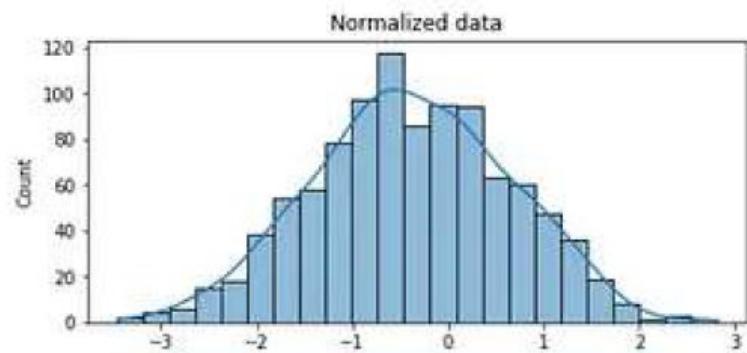
Elevation	Aspect	Slope	Horizontal_D	Vertical_Dist	Horizontal_D	Hillshade_9a	Hillshade_Nc	Hillshade_3p	Horizontal_Distance_To_Fire_Points
2596	51	3	258	0	510	221	232	148	6279
2590	56	2	212	-6	390	220	235	151	6225
2804	139	9	268	65	3180	234	238	135	6121
2785	155	18	242	118	3090	238	238	122	6211
2595	45	2	153	-1	391	220	234	150	6172
2579	132	6	300	-15	67	230	237	140	6031
2606	45	7	270	5	633	222	225	138	6256
2605	49	4	234	7	573	222	230	144	6228
2617	45	9	240	56	666	223	221	133	6244
2612	59	10	247	11	636	228	219	124	6230
2612	201	4	180	51	735	218	243	161	6222
2886	151	11	371	26	5253	234	240	136	4051
2742	134	22	150	69	3215	248	224	92	6091
2609	214	7	150	46	771	213	247	170	6211
2503	157	4	67	4	674	224	240	151	5600
2495	51	7	42	2	752	224	225	137	5576
2610	259	1	120	-1	607	216	239	161	6096
2517	72	7	85	6	595	228	227	133	5607
2504	0	4	95	5	691	214	232	156	5572

# Normalization

- ▶ Change **shape of distribution** การเปลี่ยนรูปทรงการกรุ่นกระจายตัว
- ▶ Change the observations so that they can be described as **Normal** distribution, also known as **Gaussian** distribution
- ▶ Histogram and Boxplot can be used for identifying the underlying distribution of features



From Kaggle source



Box-Cox Transformation

# Normalization- case study

	count	mean	std	min	25%	50%	75%	max
carat	53940.0	0.80	0.47	0.2	0.40	0.70	1.04	5.01
depth	53940.0	61.75	1.43	43.0	61.00	61.80	62.50	79.00
table	53940.0	57.46	2.23	43.0	56.00	57.00	59.00	95.00
price	53940.0	3932.80	3989.44	326.0	950.00	2401.00	5324.25	18823.00
x	53940.0	5.73	1.12	0.0	4.71	5.70	6.54	10.74
y	53940.0	5.73	1.14	0.0	4.72	5.71	6.54	58.90
z	53940.0	3.54	0.71	0.0	2.91	3.53	4.04	31.80

độ tuổi kim cương

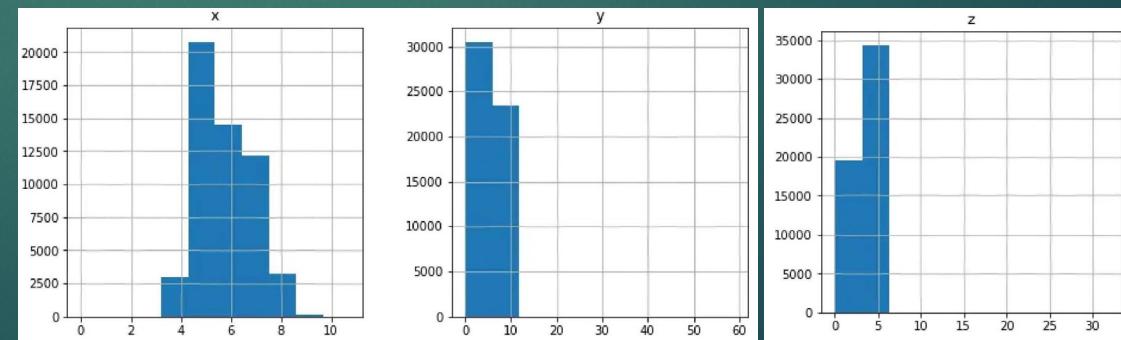
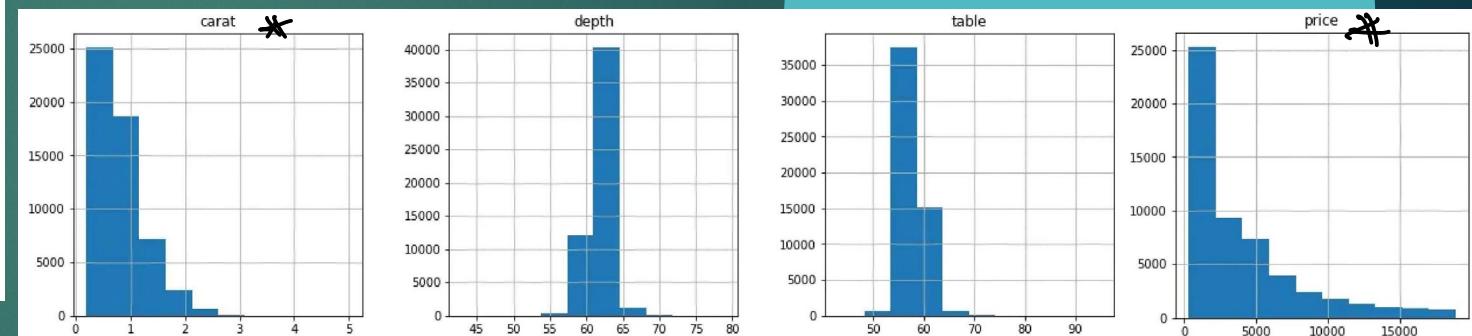
median

StandardScaler()

scaling method  $\rightarrow$  Z-score

$$\text{Z-score} = \frac{x - \bar{x}}{\sigma}$$

S.D.

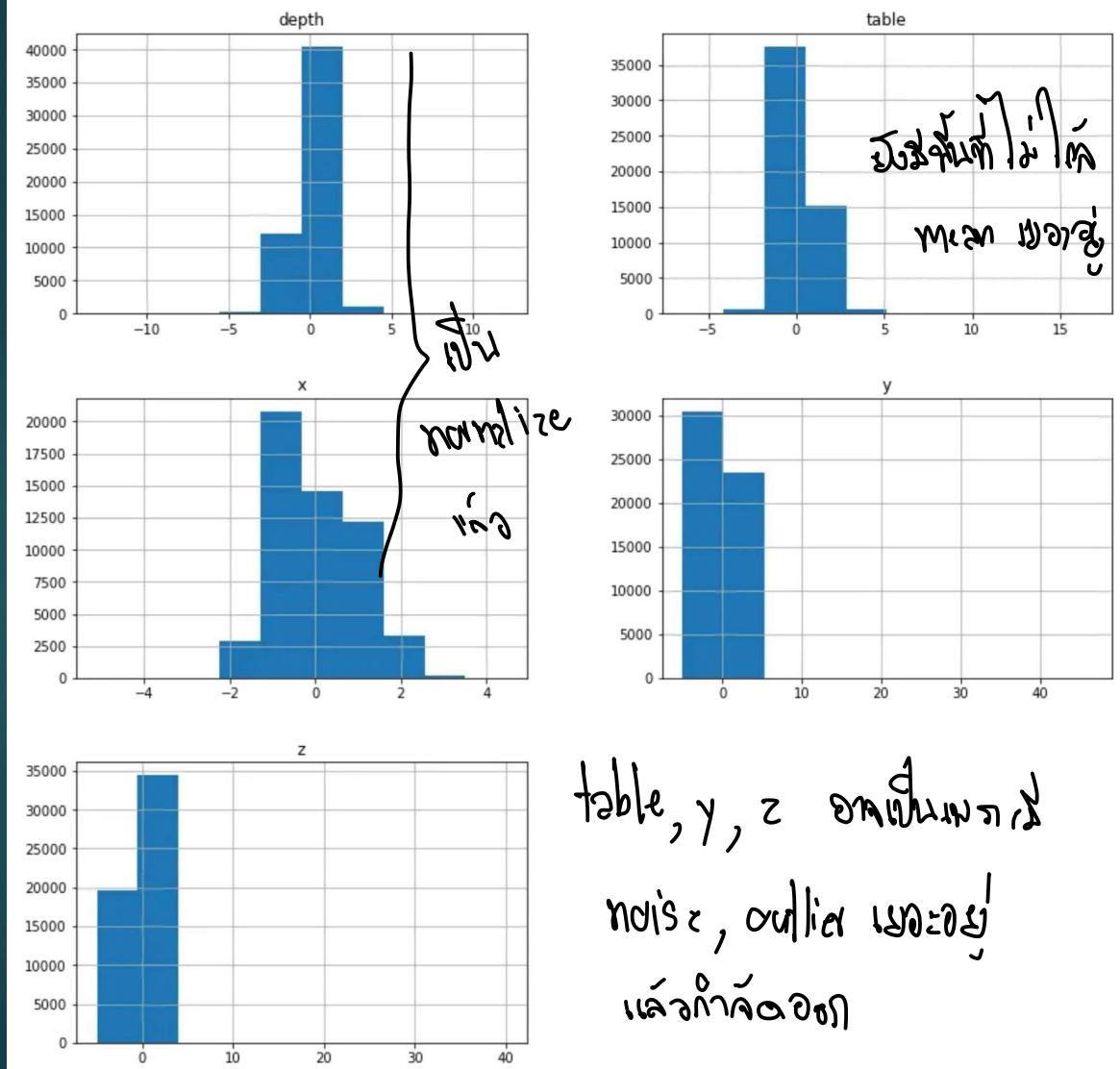


$$\begin{aligned} \text{carat} & \\ \sqrt{1} &\rightarrow \frac{v_1 - \bar{x}_{\text{carat}}}{\sigma_{\text{carat}}} \\ \sqrt{2} & \\ \vdots & \\ \sqrt{n} & \end{aligned}$$

# Normalization w/ StandardScaler()

33

mean = 0 , var<sup>2</sup> = 1



normalize ທີ່ໃຫຍ່ bell curve

```

1 >>> diamonds[to_scale].var()
2 depth      1.000019
3 table     1.000019
4 x         1.000019
5 y         1.000019
6 z         1.000019
7 dtype: float64
8
9 >>> diamonds[to_scale].mean().round(3)
10 depth    -0.0
11 table     0.0
12 x         0.0
13 y         -0.0
14 z         -0.0
15 dtype: float64

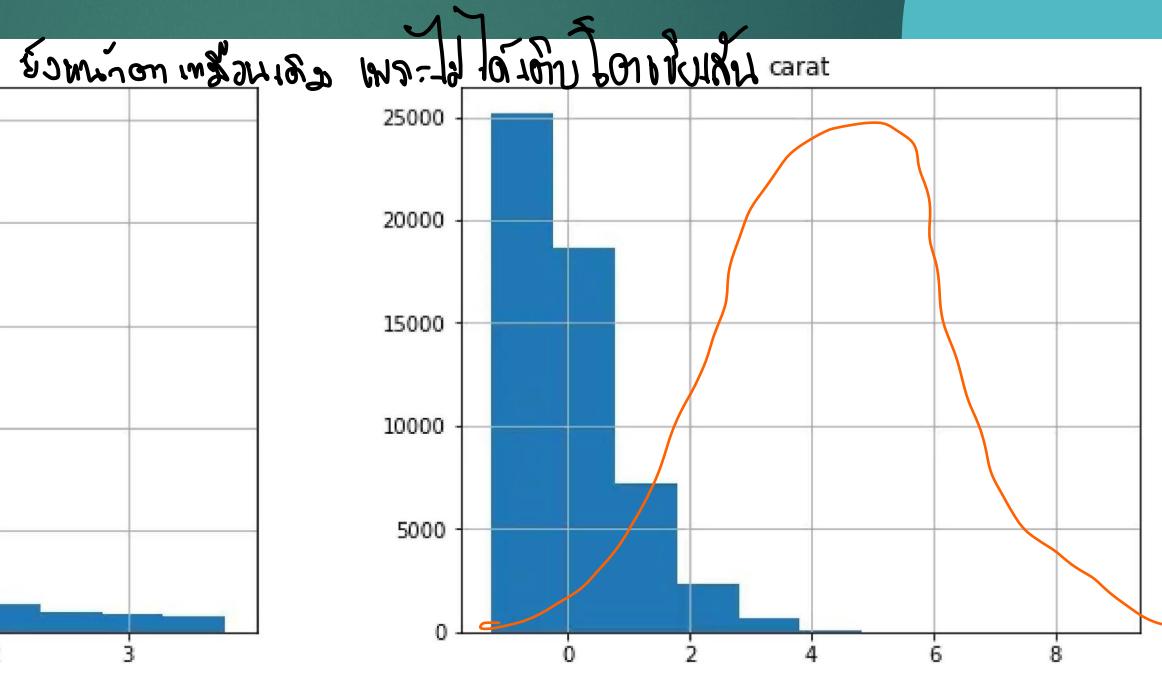
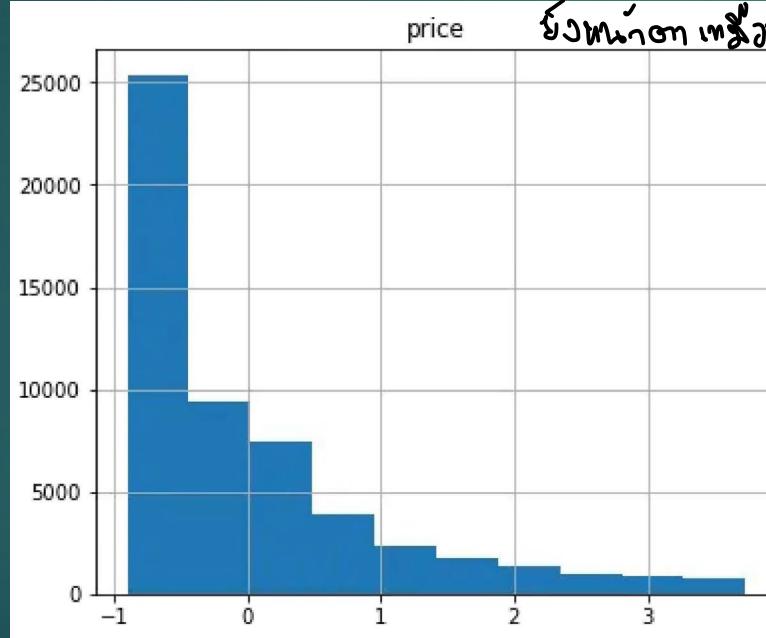
```

Depth and x now genuinely look like a Gaussian distribution. However, the features table, y, and z are still squished into the corner of their plots, suggesting the presence of outliers

# Skewed Distribution w/ StandardScaler()

ក្នុងពេលវិទ្យាល័យក៏នេះ ការប្រើប្រាស់លាក់ស្ថី មួយចំណាំ 1 2 3 4 នៅក្នុង exp box - Cox

- When a feature **does not follow a linear distribution**, it would be unwise to use the mean and the standard deviation to scale it.
- To implement non-linear transformations, **Sklearn offers a PowerTransformer()** using logarithmic functions to support Box-Cox and Yeo-Johnson transform.



Standard vs Power Law in Z-score

# Log Transform

వ్యక్తిగత పరిస్థితిలో నియమిత వ్యవహారం

log exp

- Base 2 — the base 2 logarithm of 8 is 3, because  $2^3 = 8$   $\log_2 8 = 3$
- Base 10 — the base 10 logarithm of 100 is 2, because  $10^2 = 100$   $\log_{10} 100 = 2$
- Natural Log — the base of the natural log is the mathematical constant “e” or Euler’s number which is equal to 2.718282. So, the natural log of 7.389 is 2, because  $e^2 = 7.389$ .

Natural log transformation function of NumPy

```
import numpy as np

x = [1, 2, 3, 4, 5]
y = np.log(x)
y
```

	before			After
	Income	Age	Department	log_income
0	15000	25	HR	9.615805
1	1800	18	Legal	7.495542
2	120000	42	Marketing	11.695247
3	10000	51	Management	9.210340

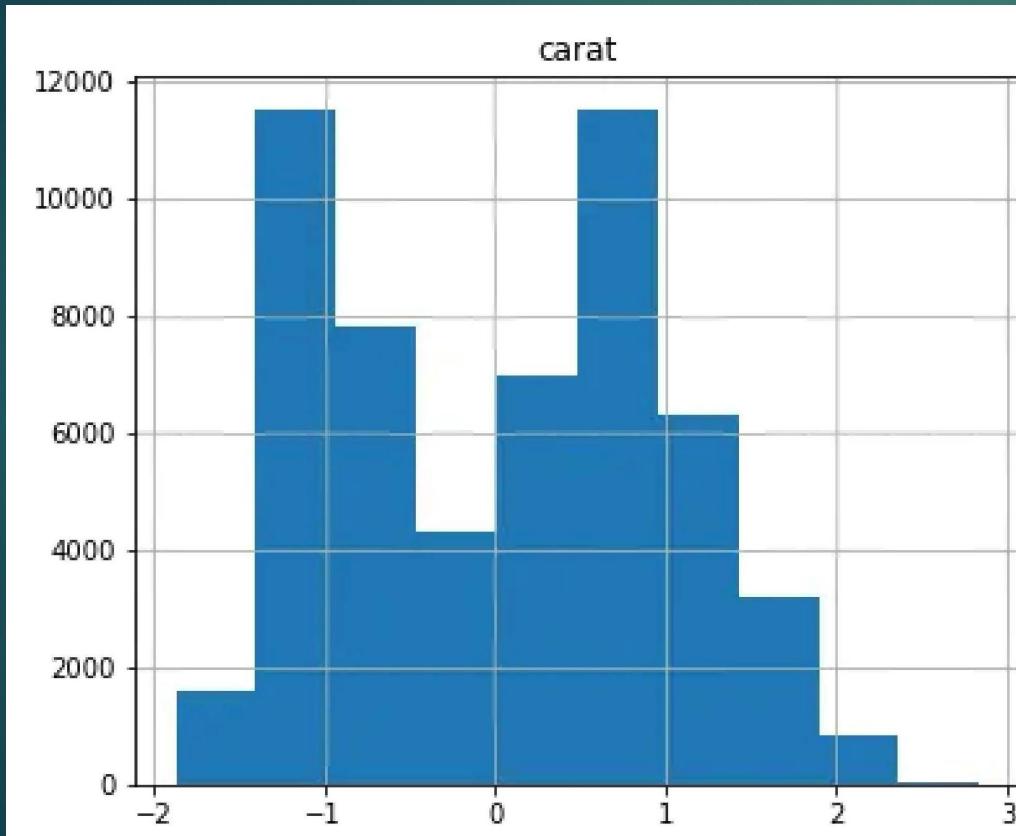
Normalization

# Skewed Distribution w/ PowerTransformer()

36

2110773-2 2/66

- The new features look much better than the old skewed ones.



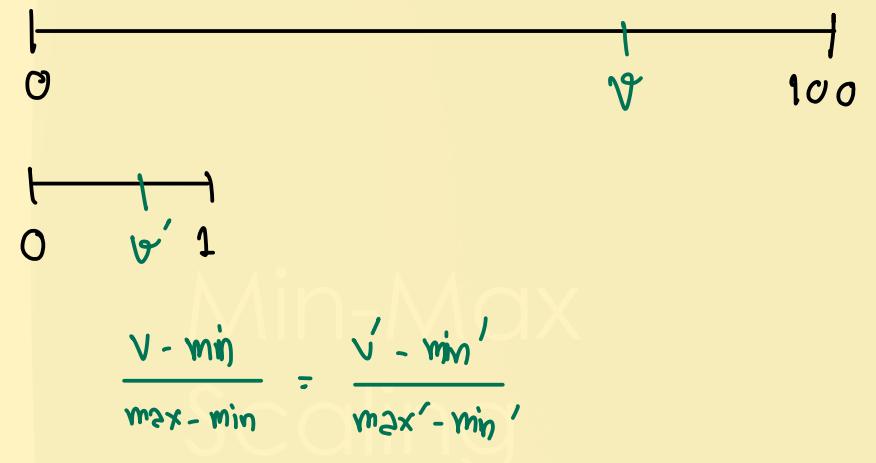
เป็นการแปลงข้อมูลเชิงเส้นจากช่วงที่เป็นไปได้เดิมของค่าอินพุต ให้เป็นช่วงข้อมูลใหม่ที่กำหนด ปกติคือช่วง [0-1]

กำหนดให้  $v$  คือค่าคุณลักษณะเดิม;  $v'$  คือค่าคุณลักษณะใหม่

$\min_A, \max_A$  คือ ค่าต่ำสุดและสูงสุดเดิมของคุณลักษณะ  $A$

$\text{new\_min}_A, \text{new\_max}_A$  คือ ค่าต่ำสุดและสูงสุดใหม่ของคุณลักษณะ  $A$  จะได้ว่า

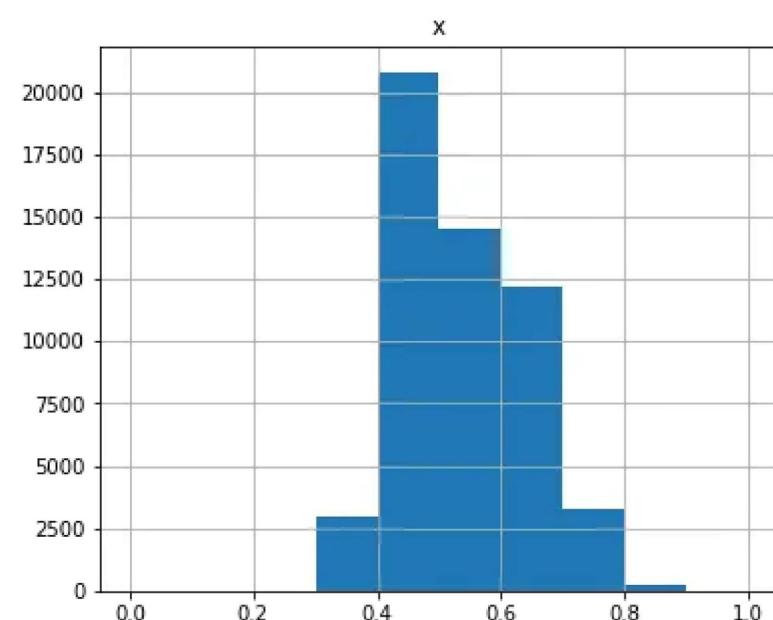
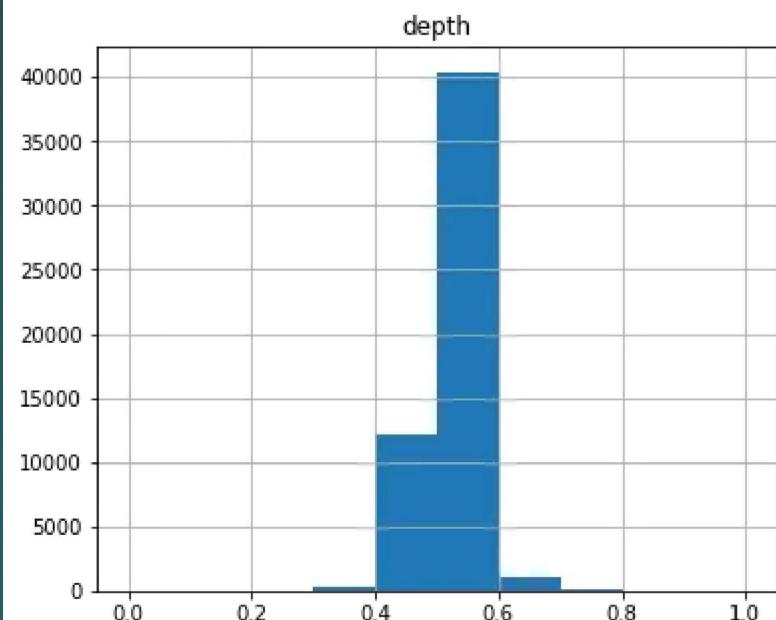
$$v' = \frac{v - \min_A}{\max_A - \min_A} (\text{new\_max}_A - \text{new\_min}_A) + \text{new\_min}_A \quad (\text{สูตรที่ 1})$$



min-max scaling

# Normalization w/ MinMaxScaler()

ມີຄົນ outlier ດັວງໄໝ່ອຸ່ນ  
MinMaxScaler does not work well with features with outliers.



Even though it forces features to follow a normal distribution, the features won't have unit variance and a mean of 0:

# StandardScaler/ PowerTransformer/ MinMaxScaler

- ▶ Scale data using StandardScaler, a transformer used when we want a feature to follow a normal distribution with mean 0 and unit variance. Used most often with distributions without too many outliers.
- ▶ Log transform data using PowerTransformer, a transformer used when we want a heavily skewed feature to be transformed into a normal distribution as close as possible.
- ▶ Normalize data using MinMaxScaler, a transformer used when we want the feature values to lie within specific min and max values. It doesn't work well with many outliers and is prone to unexpected behaviors if values go out of the given range in the test set. It is a less popular alternative to scaling.

# Data Scaling: Sigmoidal

- ▶ แปลงค่าอินพุตให้อยู่ในช่วง -1 ถึง 1 โดยใช้ฟังก์ชันซิกมอยด์ ซึ่งไม่ใช่ฟังก์ชันเชิงเส้น ข้อดีของวิธีนี้ คือ จะยังคงมีการเก็บรักษาค่าแปลกแยกไว้ การคำนวณหาค่าข้อมูลใหม่ ( $y'$ ) ที่มาลงบน

$$y' = \frac{1 - e^{-\alpha}}{1 + e^{-\alpha}}$$

โดยที่

$$\alpha = \frac{y - \text{mean}}{\text{stddev}}$$

## Sigmoidal Normalization Example

Alpha	Y	Sig Y'	Average	Std-dev
0.24	45	0.12		
0.10	35	0.05		
0.58	70	0.28		
0.57	69	0.28		
-0.08	22	-0.04		
-0.25	10	-0.12		
-0.32	5	-0.16		
-0.18	15	-0.09		
0.03	30	0.01		
3.78	300	0.96		
-0.10	21	-0.05		
-0.23	11	-0.12		
-0.19	14	-0.10		
-0.01	27	-0.01		
-1.08	-50	-0.49		
-1.43	-75	-0.61		
-0.36	2	-0.18		
-0.36	2	-0.18		
-0.36	2	-0.18		

ถ้า 300 คือ outlier  
ก็จะ sigmoid

เป็น bell curve ทาง sklearn  
ก็จะสามารถ train model ได้

ອາກົ້າໃຫ້ສໍາເລັດຕາວ່າຈະ ranking ສົງ

# Data Type Conversion: Label encoding

- ▶ CATEGORICAL → NUMERIC
- ▶ IN CASE THE ALGORITHM NEEDS NUMERICAL VALUES
- ▶ THE METHOD CAN BE PROBLEMATIC AS THE LEARNER MAY CONCLUDE THAT THERE IS AN ORDER. FOR EXAMPLE, AFRICA AND NORTH AMERICA DIFFER BY 4.

Label	Encoded Label
Africa	1
Asia	2
Europe	3
South America	4
North America	5
Other	6

The diagram shows three vectors  $v_1$ ,  $v_2$ , and  $v_3$  originating from the same point on a horizontal axis. A vertical line represents a plane. The projections of these vectors onto the plane are shown as shorter blue segments. Specifically, the projection of  $v_1$  is labeled  $v_1'$ ,  $v_2$ 's projection is  $v_2'$ , and  $v_3$ 's projection is  $v_3'$ . The projections are parallel to each other, indicating they lie in the same direction.

- The encoding produces a sparse matrix (grid of numbers) w/ lots of zeroes (false values) and occasional ones (true values).

# How Categorical w/ binary sparse Matrix

	is_africa	is_asia	is_europe	is_sam	is_nam
Africa	1	0	0	0	0
Asia	0	1	0	0	0
Europe	0	0	1	0	0
South America	0	0	0	1	0
North America	0	0	0	0	1
Other	0	0	0	0	0

# Binning 1

## Data Type

- conversion from numeric → categorical

- First **sort** data and **partition** into bins សរស់ ចំណាំ ការពារ ហើយ

- Label each bin w/ a symbol or value

- Given attribute values (for one attribute e.g., age):

  - 0, 4, 12, 16, 16, 18, 24, 26, 28

ແປងពេញ

- Equi-width binning – for bin width of e.g., 10:

  - Bin 1: 0, 4

[-,10) bin ចាយទីផលិត

  - Bin 2: 12, 16, 16, 18

[10,20) bin ចិត្តរុន្ត

  - Bin 3: 24, 26, 28

[20,+) bin សម្រួល្យូវ

\*\* – to denote negative infinity, + for positive infinity

- Alternative Equi-width:  $\text{Width} = (\text{Max} - \text{Min}) / \# \text{intervals}$  ស្ថិតិថាគារបង្កើត ចំណាំ ការពារ

ស្ថិតិថាគារបង្កើត ចំណាំ ការពារ

$$(28 - 0) / 3 = 7$$

# Binning<sub>2</sub>

- ① Data cleaning
- ② Data Reduction
- ③ Data conversion

▶ Equi-depth/  
Equi-frequency:  
ใช้ความถี่ในการแบ่งข้อมูลออกเป็น N ช่วง

▶ Equi-frequency binning –  
for bin density of e.g., 3:

- ▶ Bin 1: 0, 4, 12 [-,14) bin
- ▶ Bin 2: 16, 16, 18 [14,21) bin
- ▶ Bin 3: 24, 26, 28 [21,+] bin

smooth out noise

Histogram

“nearest neighbor” คือเดียวกัน = class neighbor

เจาะปัญหา noise จัดการให้มีส่วนลดค่าความหลาภูมิ

- ▶ Given a list of product prices  
4, 8, 9, 15, 21, 21, 24, 25, 26, 28, 29, 34
- ▶ แบ่งข้อมูลโดยวิธีแบ่งเป็นความลึกที่เท่ากัน **默記 4 ตัว/ห้าม**  
Bin1: 4, 8, 9, 15; Bin2: 21, 21, 24, 25; Bin3: 26, 28, 29, 34
- ▶ ปรับเรียบโดยใช้ค่า bin means (ค่าเฉลี่ยของแต่ละบิน)  
Bin 1: 9, 9, 9, 9; Bin 2: 23, 23, 23, 23; Bin 3: 29, 29, 29, 29
- ▶ ปรับเรียบโดยใช้ค่า bin boundaries (ค่าขอบของแต่ละบินที่ใกล้มากกว่า)  
Bin1: 4, 4, 4, 15; Bin2: 21, 21, 25, 25; Bin3: 26, 26, 26, 34
- Note: 
- ▶ Binning inevitably leads to loss of information, however, it reduces the chance of overfitting.
- ▶ Certainly, there will be improvements in speed and reduction of memory or storage requirements and redundancy.

# Data Reduction

ការបញ្ចូល

សំរាប់ model តើមួយ performance ពេលវេលា

- ▶ Obtain a reduced representation of the data set that is much smaller in volume but yet produces the same (or almost the same) analytical results
- ▶ Complex data analysis may take a very long time to run on the complete/ huge data set.

## ▶ Data Reduction Strategies

- ❖ Data Aggregation
- ❖ Dimensionality Reduction/ Feature selection
- ❖ Numerosity Reduction
- ❖ Discretization and Concept Hierarchy Generation

# ການ Data Aggregation

ໄຟສົ່ງເຮືອມກະຫຼວງ ຂອງຫ້ອງລ

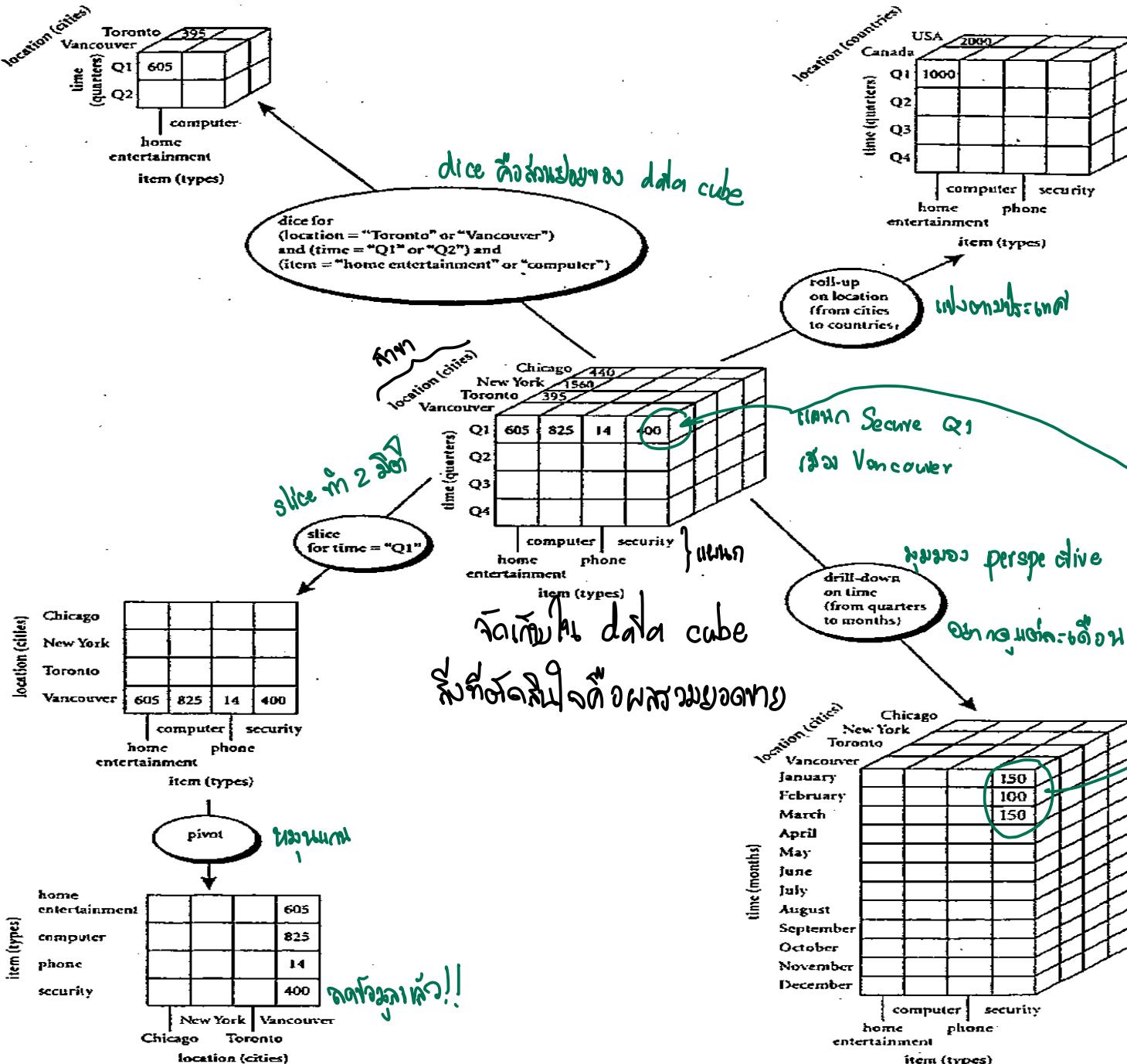
46

2110773-2 2/66

- ▶ ການລດຂໍ້ມູນໂດຍໃຊ້ຄ່າຜລວມ
- ▶ Data Cube in Data Warehouse
- ▶ ມິຕີຂໍ້ມູນ ຄື່ອ ມຸນມອງ (perspective) ທີ່ຈຶ່ງອົງຄໍກຣສນໃຈ ຕ້ອງການເກີບບັນທຶກຂໍ້ມູນໄວ້ ເຊັ່ນ ເວລາ ສານທີ່ ປະເທດ
- ▶ ແທນທີ່ຈະເກີບຂໍ້ມູນດີບຂອງຮາຍກາຮ່າຍທີ່ກຳນົດທີ່ເກີດຂຶ້ນ ວົງຄໍກຈະລົດປົກມານຂໍ້ມູນໂດຍຈັດເກີບຂໍ້ມູນຮວມຂອງຍອດຂາຍສໍາຮັບແຕ່ລະມິຕີທີ່ນໍາສົນໃຈໃນໂຄຮງສ້າງກາຮ່າຍຈັດເກີບແບບລູກບາສກຂໍ້ມູນ (data cube)

# OLAP

## Online Analytical Processing



# ការតទួល column = Dimensionality Reduction/ Feature Selection

- ❖ feature ត្រូវបានកែសម្រាប់ការ train (irrelevant)
- ❖ select  $m$  from  $n$  features,  $m \leq n$ , saving in search space
- ❖ suggestion:
  - ✓ remove key/ ID attribute
  - ✓ remove attributes with (too many) unique values
  - ✓ remove attributes with missing values  $> 50\%$
  - ✓ remove irrelevant, redundant features. Be cautious of removing relevant features as it is harmful.

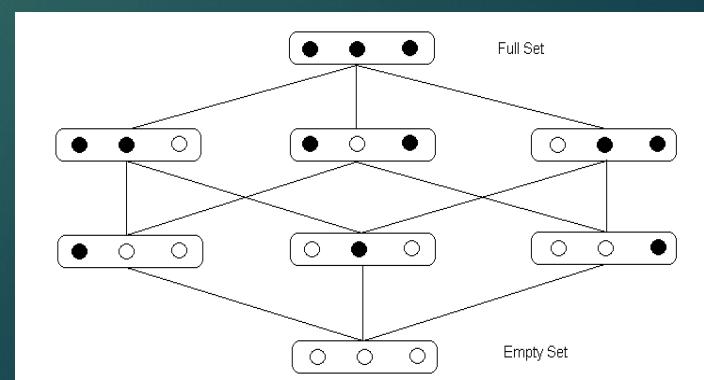
ការតទួល column = Dimensionality Reduction/ Feature Selection

	ID	new	GPA	Math	Science	English	class
A							
B+							
B							
C							
F							

មិនអាចរាយការណ៍បានការតទួល column = Dimensionality Reduction/ Feature Selection

$$\therefore 2^3 = 8$$

2 ពាណិជ្ជកម្ម  
3 features



# Principal Component Analysis (PCA)

អ្នក់ eigenvalue និង eigenvector

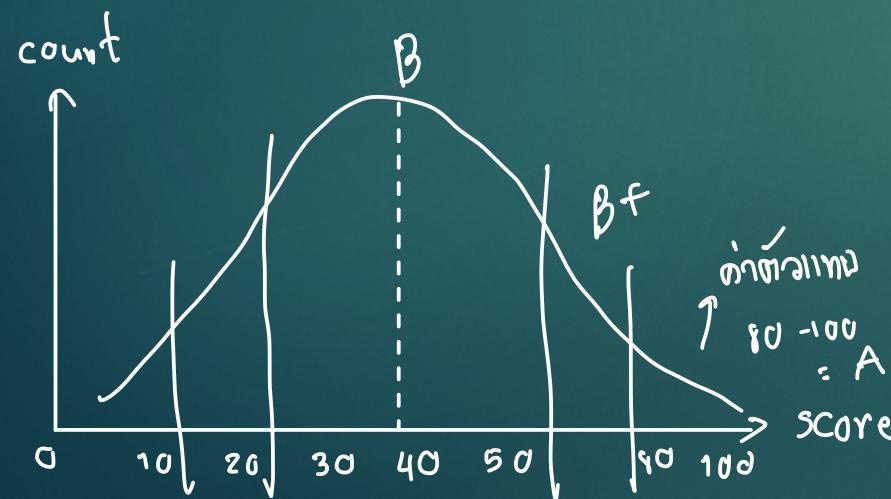
- ▶ Datasets typically contain a large number of features, but such high-dimensional feature spaces are not always helpful.
- ▶ In general, all the features are *not* equally important.
- ▶ Dimensionality reduction algorithms aim to reduce the dimension of the feature space to a fraction of the original number of dimensions.
- ▶ Principal Component Analysis (PCA) is linear dimensionality reduction technique.
- ▶ PCA is one of the most popular dimensionality reduction algorithms that takes advantage of existing correlations between the input variables in the dataset and combines those correlated variables into a new smaller set of uncorrelated variables called **principal components**. នេះ BMI ត្រូវបានអនុវត្តន៍. សម្រាប់  
▶ PCA requires feature scaling if there is a significant difference in the scale between the features of the dataset.  
ឧប. សង្គមពីរគឺ correlate  
ដែលបានការពារ

# Numerosity Reduction

ແພນດີ່າວົງ ນຳມະວາງທຸກໆທີ່ຈິງຈານກວດຂົງ

- ▶ Replace original data by smaller form of data representation

- ▶ ໃຊ້ເຄື່ອງນີ້ອ ເຊັ່ນແຜນກາພີສໂຕແກຣມ (Histogram) ອີ່ວິວິກາຮຈັດກລຸ່ມ (Clustering) ປ່າຍ  
ແສດງກາຮຈາຍຂອງຂໍ້ມູນ ແລະ ເກີບຄ່າຕ້ວແໜກລຸ່ມແທນຄ່າຂໍ້ມູນຈົງ ອີ່ອາລໃຫ້ວິທາງ  
ສົດໃຫ້ ເຊັ່ນ ກາຮສຸ່ມຕ້ວອຢ່າງ (Sampling/ Instance selection)



ຫຼັງຈາກນີ້ຈະໄຟ້

ມີມະວາງ ມີສົມຈຳຕໍດະເນັນເກົ່ານົກ

$$\text{Class GPA} = 4.00$$

$$\text{Data mining} = 2.75$$

# Instance Selection

Use portion rather than the whole huge dataset

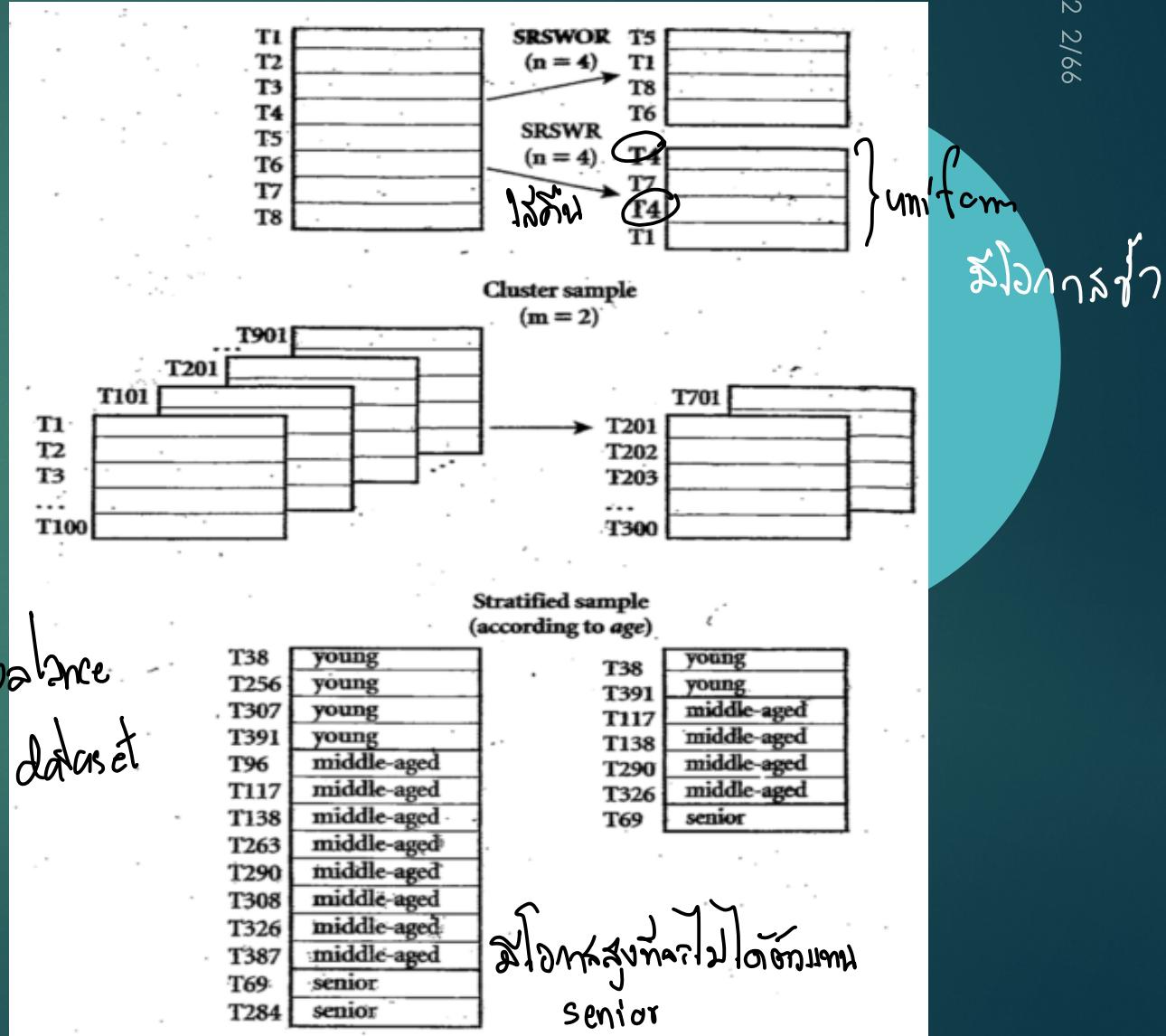
51

2110773-2 2/66

## Sampling methods

- Simple Random Sample Without Replacement (SRSWOR) ពិនិត្យស្ថិតិ
- Simple Random Sample With Replacement (SRSWR)
- Cluster Sample
- Stratified sampling treating each stratum as a population
  - ❖ Proportional allocation ការណែនាំចាត់រាង
  - ❖ Equal sample sizes  $2:4:1$   
 $y:m:s$

$2:2:2$  រូបភាពខ្លួន



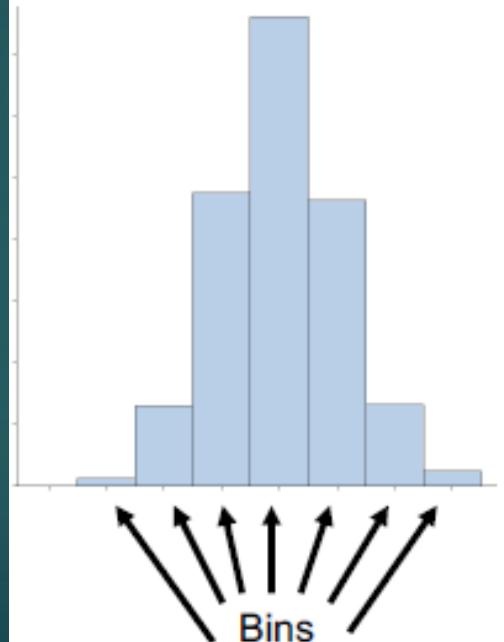
# Discretization & Con

## Discretization

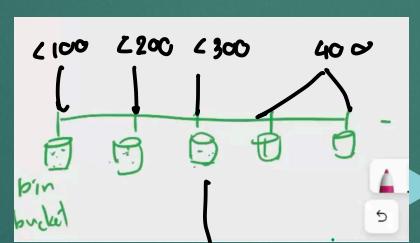
ทําให้พื้นที่  
บูรณาการน้อย

### ► Binning methods

- Equi-width/ Equal-width
- Equi-depth/ Equal-frequency

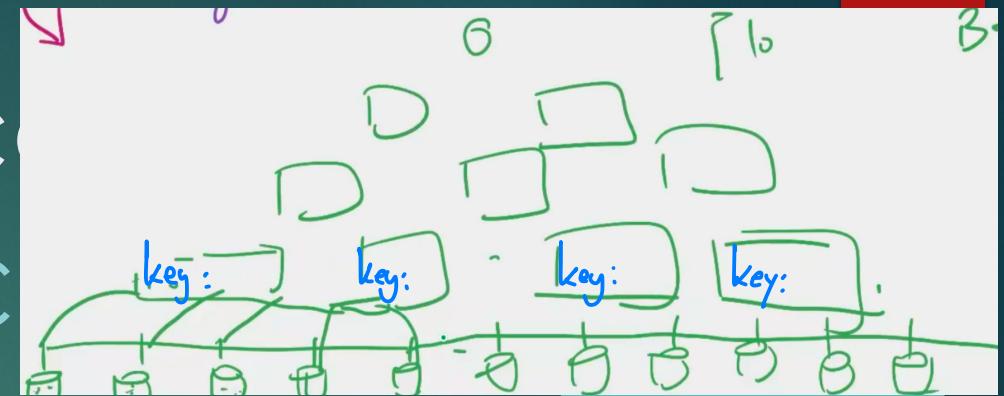


numerical → categorical

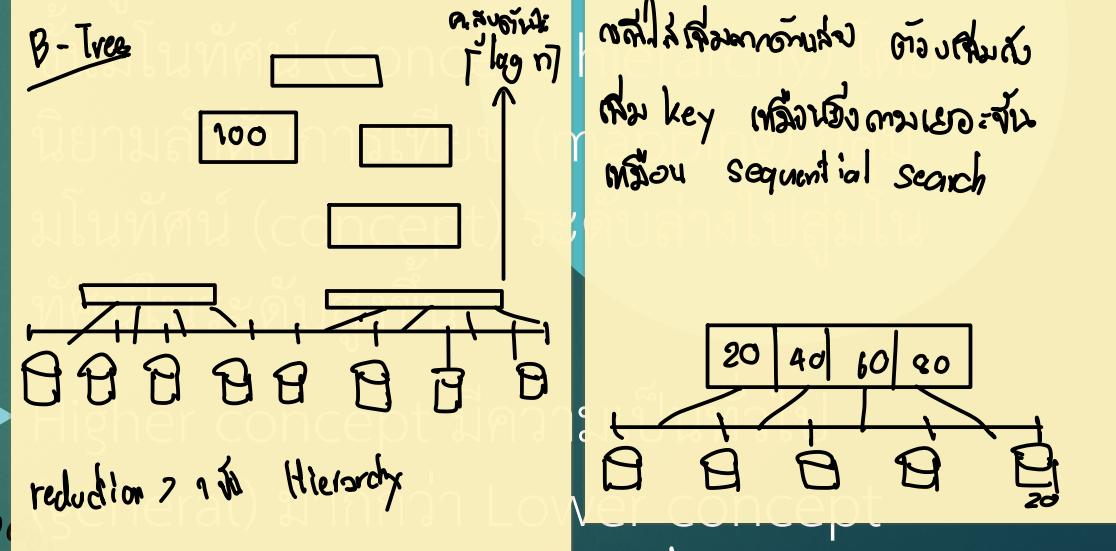


ใน histogram  
ใน distribution เนรื่องว่ามีตัวบ่ง

เพื่อให้ distribution เนรื่องว่ามีตัวบ่ง



การลดข้อมูลประเภท Categorical หรือ  
ข้อมูลที่ไม่ต้องเนื่องด้วยการสร้างเป็นลำดับ



B-Tree , B<sup>+</sup>-Tree

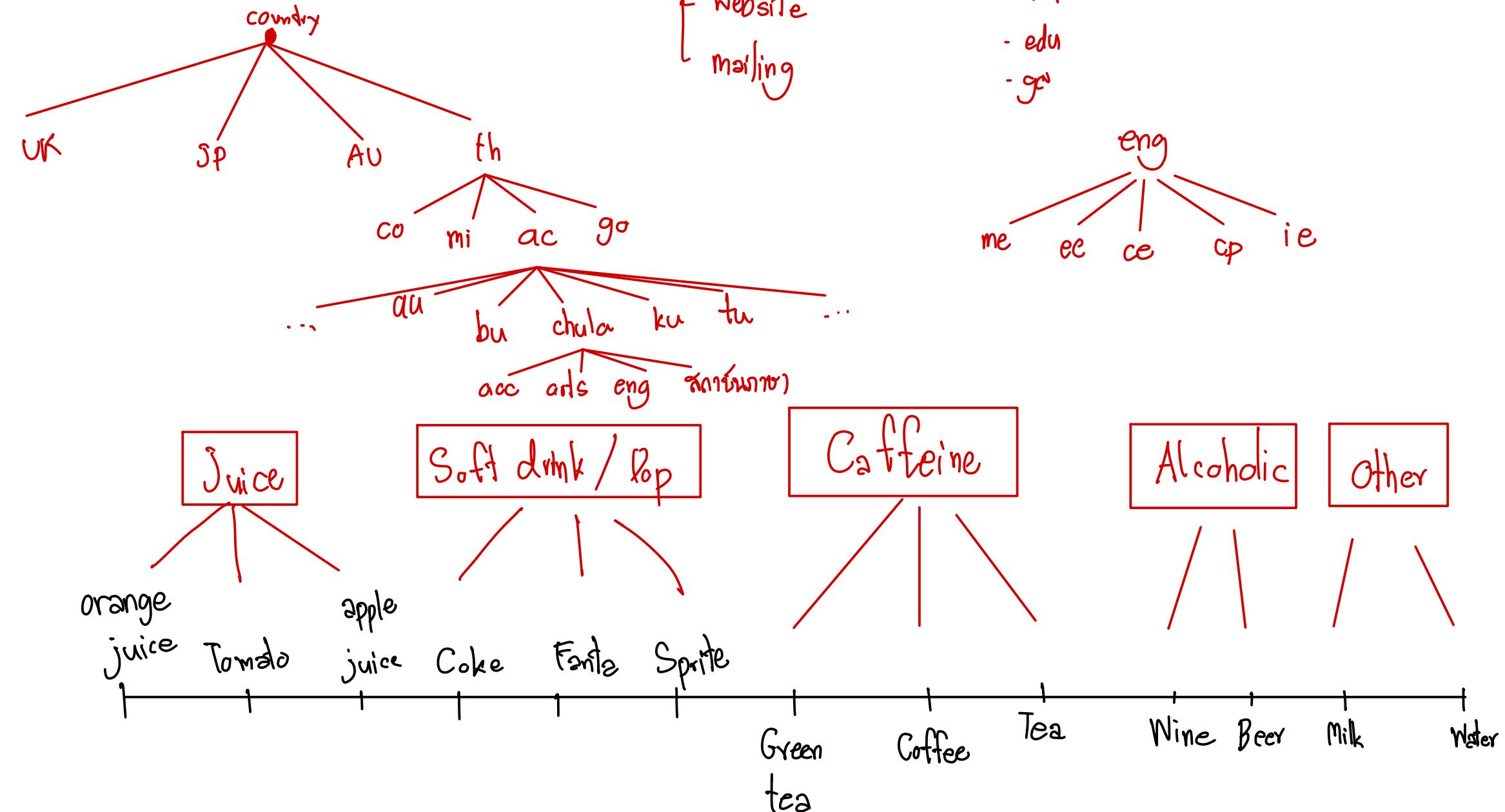
จึง indexing direct  
access search

URL: www.csail.mit.edu.ac.th

Addr.

[ email  
website  
mailing ]

.com  
.mil  
.edu  
.gov

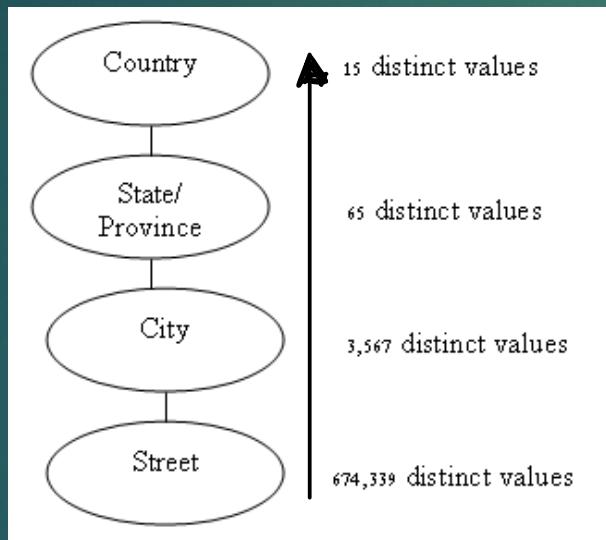


# Schema hierarchy ส่วนมากเป็นความสัมพันธ์

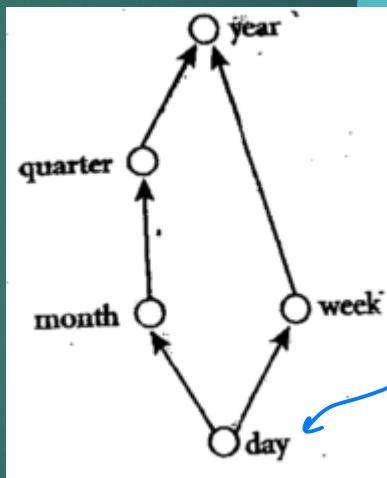
53

ระหว่างคุณลักษณะในฐานข้อมูล ซึ่งอาจเป็นความสัมพันธ์แบบ

Total Order ตามลำดับ distinct value



Partial Order

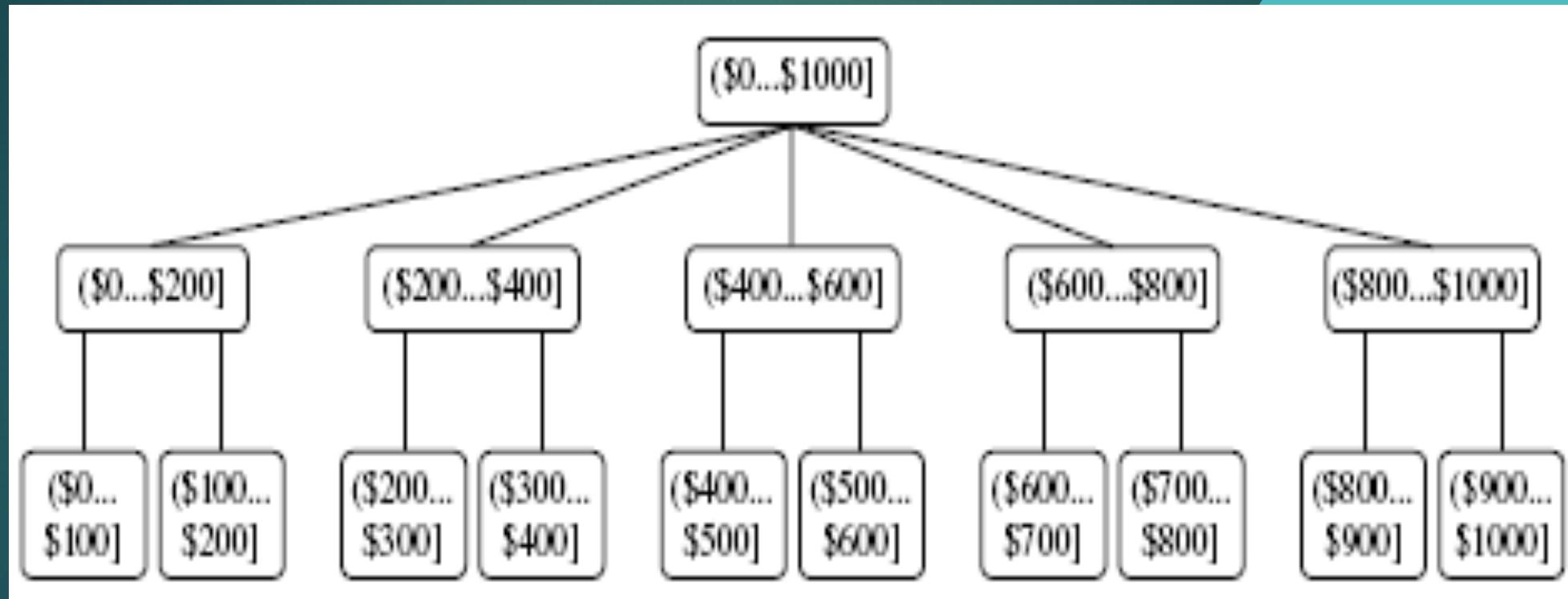


lattice โครงสร้างพื้นฐาน

เช่นเดียวกับ ที่นี่  
เดือนที่ 30 ก็ต้อง 31

# Set-grouping hierarchy การแบ่งค่าคุณลักษณะ

ออกเป็นช่วงๆ และการแบ่งช่วงค่าสามารถทำต่อเป็นลำดับชั้น



# Operation-derived Hierarchy

- ▶ การกำหนดลำดับชั้นมโนทัศน์ (concept hierarchy) จะขึ้นอยู่กับการใช้งานหรือการปฏิบัติงานของผู้ใช้/ ผู้เชี่ยวชาญ ตัวอย่างเช่น email address หรือ URL ของหน้าเว็บต่างๆ

# Rule-based Hierarchy [ការងារ ផ្តោះ និង ចំណែក]

- ▶ การกำหนดតម្លៃជូនទៅស្ថាន វាទេអីពីរក្នុងបញ្ហានេះ តាមរយៈរាយចំណែក  
กำหนดให้  $P_1$  = retail price of X;  $P_2$  = actual cost of X
- ▶ lowProfitMargin(X)  $\leftarrow$  price(X, P1) and cost(X, P2) and  
 $(P_1 - P_2) < \$50$       ពាណិជ្ជកម្ម - លក់ < \$50 ក្នុងមីនី
- ▶ mediumProfitMargin(X)  $\leftarrow$  price(X, P1) and cost(X, P2)  
and  $((P_1 - P_2) \geq \$50 \text{ and } (P_1 - P_2) \leq \$250)$
- ▶ highProfitMargin(X)  $\leftarrow$  price(X, P1) and cost(X, P2) and  
 $(P_1 - P_2) > \$250$

# Data Preparation

- 1) Examining the Data Set
- 2) Feature Selection/ Narrowing down columns manually

- ▶ Remove Id or key និង PK
- ▶ Remove irrelevant variables  
Feature understanding is extremely important! Require Domain Expert សមត្ថភាពអ្នកដំណឹង ឬជាប្រធានបទពិសេស
- ▶ Remove Calculated fields
- ▶ Remove flat values  
សមត្ថភាពដំណឹង

ឬស្ថាដែល, គើរចាតិ

derive ជាក field នៃ

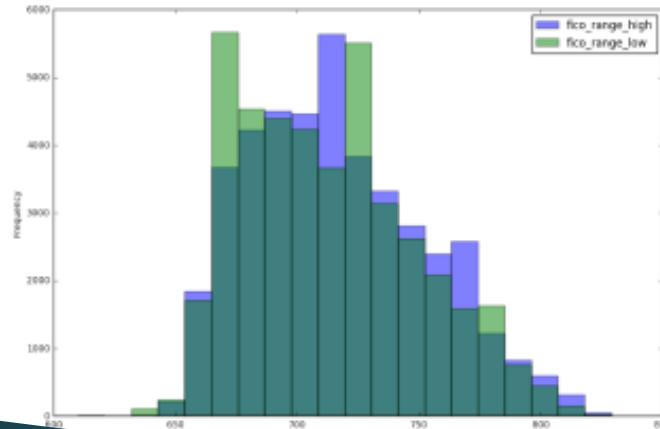
មនុស្សនៅ ពេកខាងក្រោមគេង

## 3) Preparing features

- ▶ Drop unqualified features
  - Variables with missing values > 50%
  - Too many unique values តុលាល័យត្រូវបានលើកឡើង 90-95%
- ▶ Handling missing values
- ▶ Investigate categorical features
  - Recode, consolidation (grouping)
  - Convert ordinal to numeric
  - Convert categorical to numeric → one hot, label encoding
- ▶ Check all numeric variables
  - ▶ Truncate outliers → bell curve, box-plot ងារបញ្ជី
  - ▶ Feature Transformation → scaling, normalization

- Numerical variables

- Out of ranges
- Distribution: histogram

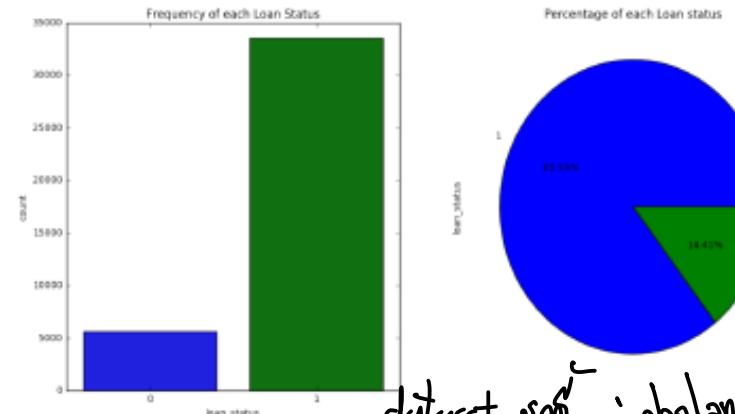


- Categorical variables

- Miscodes *911*
- Distribution: frequency table, bar chart

- Target variable *class, label*

- Understand proportion of each class:  
bar chart, pie chart



dataset *van* imbalance *en*voegen

# Examining the Dataset