



# 2110773 Data Mining Chapter2: Data Preprocessing

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

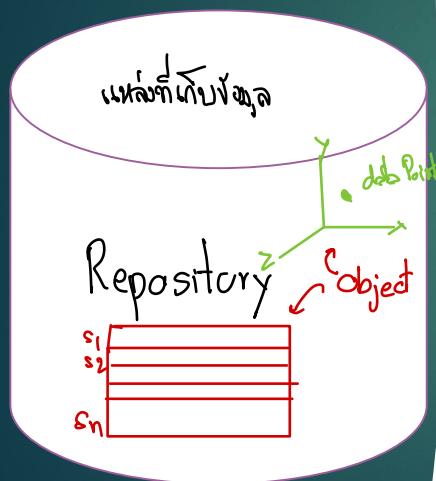
↓  
skill จากทฤษฎี

↓  
ค้นพบ/ความรู้ที่ซ่อน

▶ รศ. ดร. ญาใจ ลีมีปิยะกรณ์



## Types of Dataset



### Record

Relational records  
Document data: text documents  
Transaction



### Graph and network

World Wide Web → *Traversal path mining*  
Social or information networks  
Molecular Structures



### Others

Image  
Video data: sequence of images  
Temporal/ Time-series → *Trend Analysis*  
Spatial data: maps

2

2110773-2-2/66

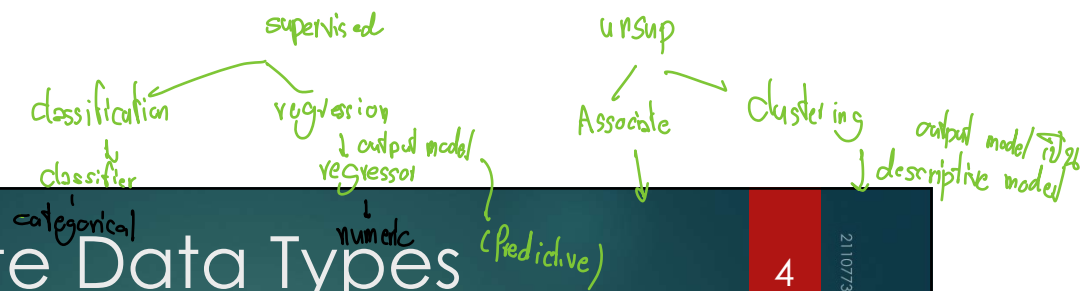
# Data Object

3

2110773-2 2/66

- ▶ Data sets are made up of data objects. *บรรจ data object ไปด้วย*
- ▶ A **data object** represents an entity. For examples:
  - ▶ medical database: patients, treatments *ชื่อ entity เปลี่ยนตามลักษณะ*
  - ▶ university database: students, professors, courses *conceptual design - ER diagram*
- ▶ Also called *samples, examples, instances, data points, objects, tuples.* *ชื่ออื่น ๆ ที่ใช้แทน object*
- ▶ Data objects are described by **attributes**. *ชื่อของตัว set of Attribute*
- ▶ 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*

*1 point = 1 object ที่คุณไปซื้อ 7 ชาติ*



# Attribute Data Types

4

2110773-2 2/66

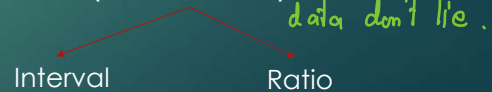
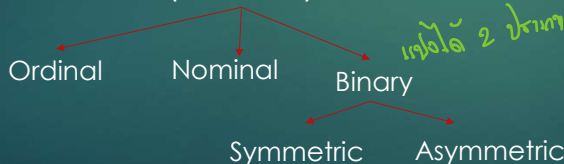
1. Qualitative/ Quantitative
2. Categorical/ Numeric
3. Discrete/ Continuous

*datum = หน่วยของ data*

- 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.

*Subjective เราบอกกันไม่ได้*  
Qualitative/ Categorical  
(Discrete)

*objective คำนวณได้*  
Quantitative/ Numeric  
(Continuous)



*2 เท่า  
6°C 41°F  
0°C 50°F  
ไม่ได้อยู่แค่ 2 เท่า*

*ไม่ได้อยู่แค่ 2 เท่า  
C, F → Interval แต่ 0 ไม่ใช่ 0 แท้จริง*

*ก. หน่วยรวม  
ที่นับไม่ได้* → **Quantify** (IQ, คะแนน) → **หน่วย** (norm: unit of measure)  
*C, F โดยไม่ต้องเน้น  
ในสิ่งตัวบุคคล*

2 นาที ( 80 kg 176 lb ) 2 นาทีนี้คือส่วนสูง ซึ่งเป็น Ratio  
40 kg 88 lb

## Attribute Types

5

21/07/3-2 2/66

- ▶ **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°).

Kevin คือของ no heat หรือ

## Data Type Examples

6

21/07/3-2 2/66

Data Type	Examples
Nominal	color, bloodType, zipCode, ID#, occupation, political party
Ordinal	the Gold, Silver, Bronze "รางวัลชนะเลิศ" (ทองคำ) รางวัลอันดับ 2, 3 medal, satisfaction, grade, frequency, academic ranking - อันดับมหาวิทยาลัย
Binary- symmetric	gender = { M, F }
Binary- asymmetric	labTest = { +, - } +, - ความหมายต่างกัน
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.

สิ่งที่มันเป็นเชิงสัญลักษณ์ เปรียบเทียบไม่ได้  
color = { red, green, ... }  
bloodType = { A, B, O, AB }

%	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 <i>มีชื่อ</i>	✓	✓	✓	✓
Order <i>ลำดับ</i>	✗	✓	✓	✓
Measurable Difference <i>มีค่าต่าง</i>	✗	✗	✓	✓
True Zero Starting Point	✗	✗	✗	✓

*เพราะมันมีชื่อ*

อันดับอาจมีชื่อ - ไม่มีชื่อ  
 อันดับที่มีความหมาย 1-5 ก็ไม่ถือว่าเป็น  
 number

# Survey

9

2110773-2 2/66

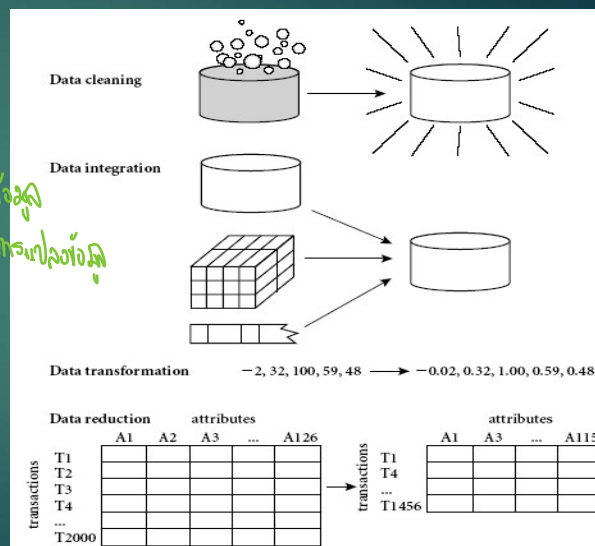
1. How old are you? \_\_\_\_\_ years *Ratio*
2. Are you: Male Female *Binary symmetric*
3. How much do you spend on groceries each week? \_\_\_\_\_ Baht *Ratio*
4. How many cups of coffee do you buy in a week? \_\_\_\_\_ *Ratio*
5. Which type of coffee do you like most?  
Latte Espresso Cappuccino Americano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

# Data Preprocessing

10

2110773-2 2/66

- ▶ Data Cleaning
- ▶ Data Integration *บูรณาการข้อมูล*
- ▶ Data Transformation *การแปลงข้อมูล*
- ▶ Data Reduction *ลดข้อมูล*



# Data Cleaning

11

21/07/3-2 21/66

- ▶ Fill in missing data *error data*
- ▶ Smooth noisy data- random error or variance in a measured variable
- ▶ Identify or remove outliers
- ▶ Resolve inconsistencies *ข้อมูลไม่สอดคล้อง*
  - ▶ 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

12

21/07/3-2 21/66

- ▶ 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 *As is = ตามสภาพ*
  - ▶ Leave as is, however, some algo can't deal w/ missing values and the program may refuse to continue or lead to inaccurate results *ไม่สนใจ / บาง Algo มีความสามารถในการจัดการ*
  - ▶ 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?! *ใส่ค่าคงที่ทั่วๆ ไป*
  - ▶ Imputing : *ประมาณค่าที่หายไป*
    - ❖ Attribute mean/median (Numerical variables); mode (Categorical variables) *ค่าที่ถนัดที่สุด symbolic*
    - ❖ 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 *แทนค่าที่หายจาก attribute ขึ้นอยู่กับ เช่น อุณหภูมิเฉลี่ย ในฤดูกาลนั้นๆ*
    - ❖ Model-based/ inference-based: Regression, Decision Tree, k-nearest neighbor, Bayesian ... *ใช้ model ทำนายค่า*



# Noisy Data

13

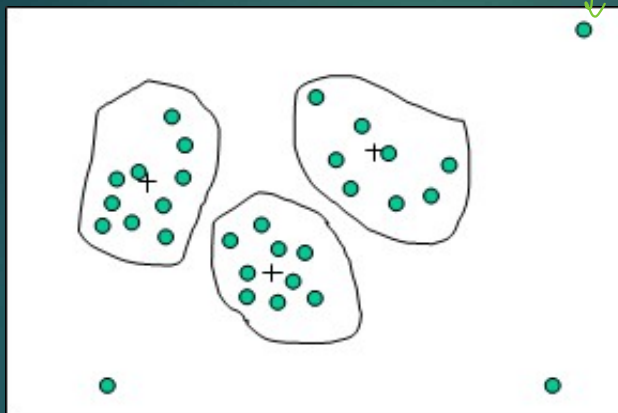
2110773-2 2/66

- ▶ 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

Normal Curve

14

2110773-2 2/66



ข้อมูลที่มีลักษณะต่างจากข้อมูลส่วนใหญ่  
Clustering

# Data Distribution

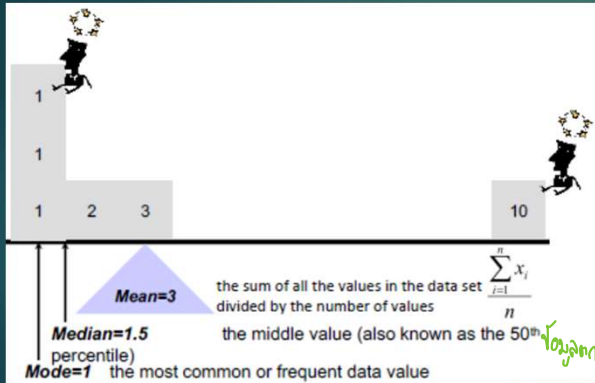
การกระจายตัวข้อมูล

15

2110773-2 21/66

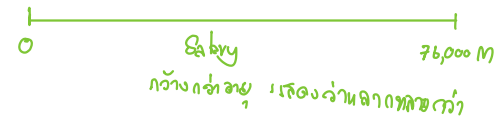
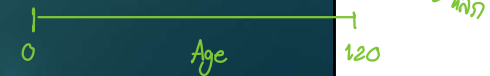
1. Central Tendency/ Center   
 ข้อมูลอยู่ใกล้ๆ ค่ากลาง เป็น bell curve   
 คู่มือค่า mean

2. Spread/ Dispersion การวัดการกระจาย, การวัด



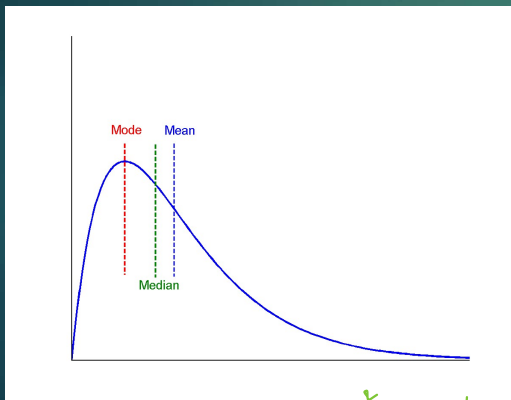
Measure	Definition
Range	the difference between the maximum and minimum data values
Interquartile Range	the difference between the 25th and 75th percentiles
Variance	a measure of dispersion of the data around the mean
Standard Deviation	a measure of dispersion expressed in the same units of measurement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

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

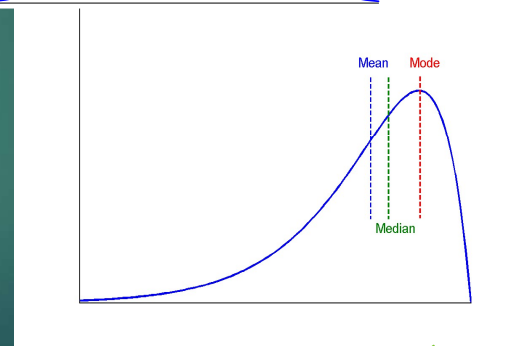
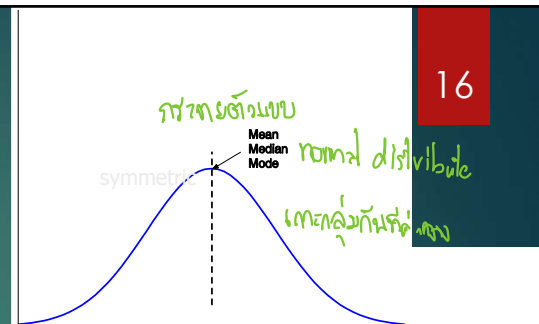


## Symmetric vs. Skewed Data

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



positively skewed เบ้ขวา (ดูที่หาง)



negatively skewed เบ้ซ้าย

เพราะข้อมูล การทางซ้ายมาก

จึงต้องตอบไปทางขวาเพื่อหาค่า normal distribution



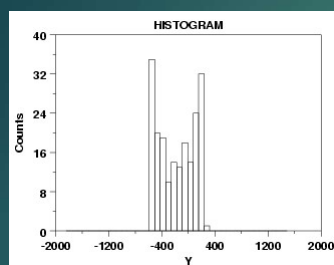
ตัวที่นิยมมากที่สุดเท่าที่ทราบ

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

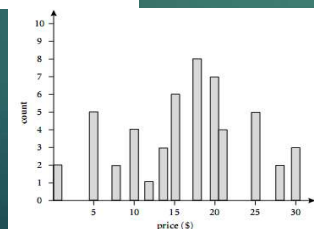
## When to use Mean, Median, Mode

## Graphical Displays of Distribution

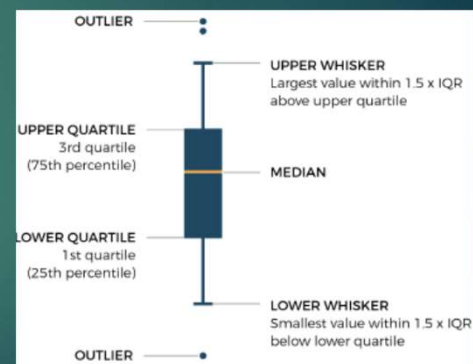
- ▶ Histogram Graph display of frequencies, shown as bars with numeric values on X axis

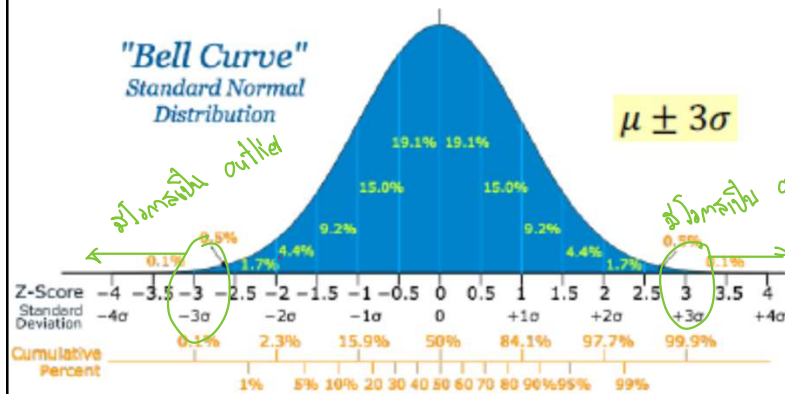


Singleton Histogram



- ▶ Box plots





## Truncate Outliers: Bell Curve

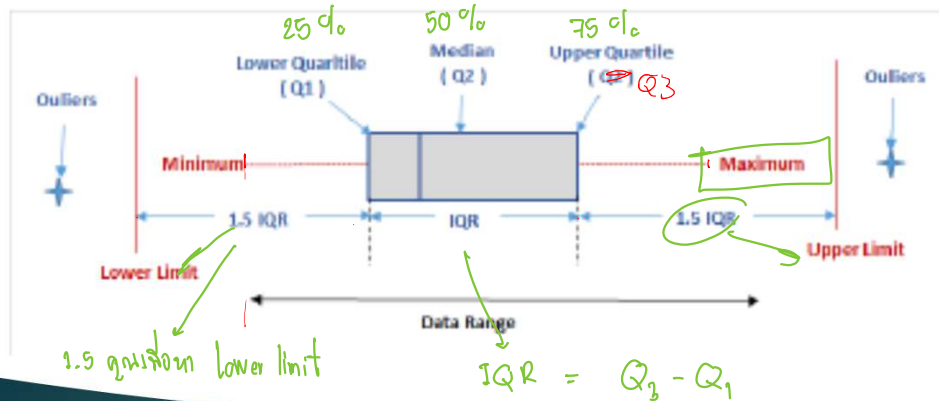
Variance and standard deviation  
(sample:  $s$ , population:  $\sigma$ )

Standard deviation is the square  
root of variance

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

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

2110773-2 2/66



## Truncate Outliers: Box Plots

9th Quartile

2110773-2 2/66

# Interquartile Range : IQR

21

2110773-2-2/66

- ▶ 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

22

2110773-2-2/66

### 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  $Q_1$  and  $Q_3$  easier to spot.  
(1, 2, 5, 6, 7), 9, (12, 15, 18, 19, 27).
- ▶ Step 4: **Find  $Q_1$  and  $Q_3$**  median of the lower half of the data and think of  $Q_3$  as a median for the upper half of data.  
(1, 2, 5, 6, 7), 9, (12, 15, 18, 19, 27).  $Q_1 = 5$  and  $Q_3 = 18$ .
- ▶ Step 5: **Subtract  $Q_1$  from  $Q_3$  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  $Q_1$  and  $Q_3$  easier to spot.**  
(3, 5, 7, 8, 9), | (11, 15, 16, 20, 21).
- ▶ Step 4: **Find  $Q_1$  and  $Q_3$**   
 $Q_1$  is the median (the middle) of the lower half of the data, and  $Q_3$  is the median (the middle) of the upper half of the data.  
(3, 5, 7, 8, 9), | (11, 15, 16, 20, 21).  $Q_1 = 7$  and  $Q_3 = 16$ .
- ▶ Step 5: **Subtract  $Q_1$  from  $Q_3$ .**  
 $16 - 7 = 9$ .

## Correlated Data

Positively

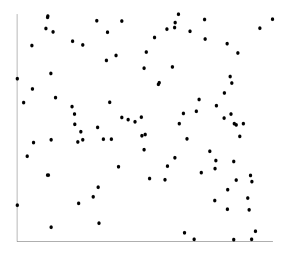


Negatively



## Uncorrelated Data

Smear ৗৗৗ



23

2110773-2 2/66

## Regression

### ► Linear Regression

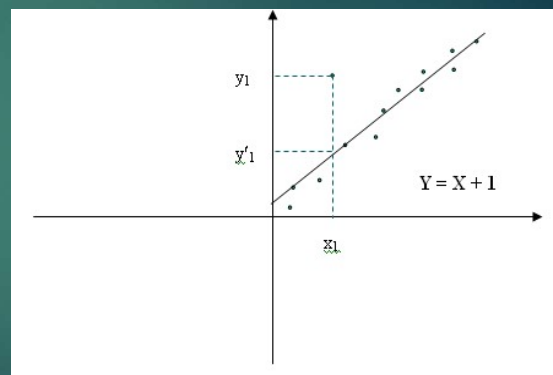
$$Y = \alpha + \beta X$$

### ► 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

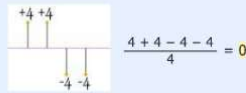


24

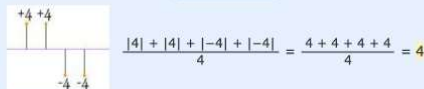
2110773-2 2/66

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

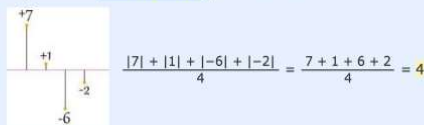
If we just add up the differences from the mean ... the negatives cancel the positives:



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

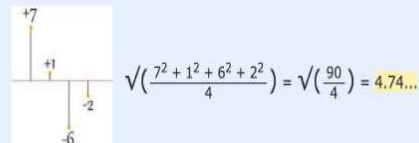
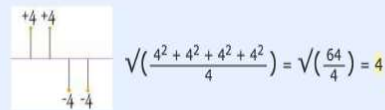


That looks good (and is the [Mean Deviation](#)), but what about this case:



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

27

2110773-2 2/66

- ▶ make the learning algorithms understanding easier, and achieving better results
  - Standardization (Scaling) / Normalization
    - ❖ Min-max Scaling
    - ❖ Log Transform
    - ❖ Z-score/ Zero-mean
    - ❖ Sigmoidal
  - Data Type Conversion

<https://towardsdatascience.com/how-to-differentiate-between-scaling-normalization-and-log-transformations-69873d365a94>

## Scaling vs. Normalization

28

2110773-2 2/66

- \* very similar, confusing
- \* sometimes used interchangeably
- \* **numeric variables** are transformed in both cases
- \* what's the difference?

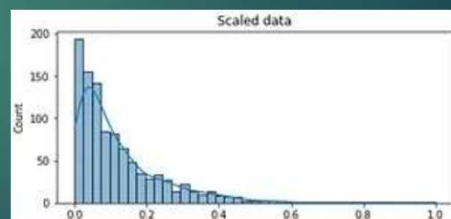
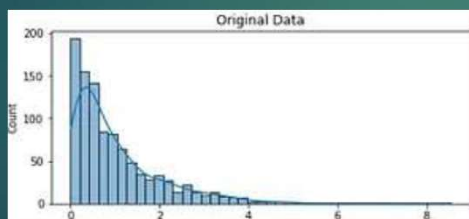


# Scaling

29

2110773-2 2/66

- ▶ Numeric variables
- ▶ Change **range** of data e.g. 0-1, 0-100
- ▶ Applied in distance-based algorithms, e.g. SVM, kNN
- ▶ Same importance for a change of "1" in any numeric feature
- ▶ By scaling, variables are compared on equal footing



From Kaggle source

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

30

2110773-2 2/66

# Scaling: case study

31

2110773-2-2/66

- 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

Elevation	Aspect	Slope	Horizontal_Dist	Vertical_Dist	Horizontal_Dist	Hillshade_3a	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	616	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	155	5572

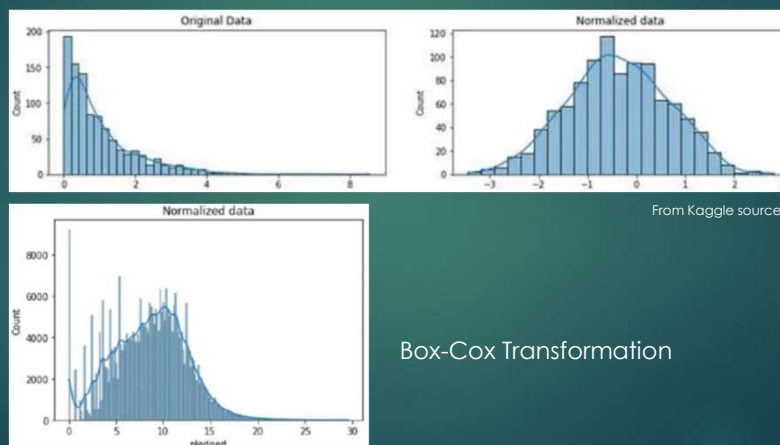
<https://medium.com/@urvashiluniya/why-data-normalization-is-necessary-for-machine-learning-models-681b65a05029>

# Normalization

32

2110773-2-2/66

- Change **shape of distribution**
- Change the observations so that they can be described as **Normal** distribution, also known as **Gaussian** distribution



Box-Cox Transformation