

# ① Measure of Central Tendency

① Mean

② Median [EDA and Feature Engineering]

③ Mode



## ① Mean

Population (N)

Sample (n)

$X : \{1, 1, 2, 2, 3, 3, 4, 5, 5, 6\} \rightarrow$

$n=10$

$$\text{Population mean } (\mu) = \sum_{i=1}^N \frac{x_i}{N}$$

$$\text{Sample mean } (s) = \sum_{i=1}^n \frac{x_i}{n}$$

$$\begin{aligned} \mu &= \frac{1+1+2+2+3+3+4+5+5+6}{10} \\ &= \frac{32}{10} = 3.2 \end{aligned}$$

## ② Median

4, 5, 2, 3, 2, 1

Sort  $\rightarrow$  1, 2, 2, 3, 4, 5

Median

Even Count

1, 2, 2, 3, 4, 5

Odd Count

{1, 2, 2, 3, 4, 5, 7}



$$\downarrow$$

$$\frac{2+3}{2} = 2.5$$

$$\boxed{\text{Median} = 3}$$

$$\boxed{\text{Median} = 2.5}$$

Why Median?

$$\{1, 2, 3, 4, 5\}$$

$$S = \frac{1+2+3+4+5}{5} = \frac{15}{5} = \underline{\underline{3}}$$

$$\{1, 2, \boxed{3}, 4, 5\}$$

$\downarrow$

$$\text{Median} = \underline{\underline{3}}$$

$$\{1, 2, 3, 4, 5, \boxed{100}\}$$

$$S = \frac{1+2+3+4+5+100}{6} = \frac{115}{6} \approx 19.17$$

$$\{1, 2, \boxed{3}, \boxed{4}, 5, \boxed{100}\}$$

$$\text{Median} = \underline{\underline{3.5}}$$

③ Mode = Frequency Maximum

$$\{2, \underline{1}, \underline{1}, \underline{1}, 4, 5, 7, 8, 9, 10\}$$

$$\boxed{\text{Mode} = 1}$$

KDA And Feature Engineering

$\downarrow$

Type of Flower

Age

Lily

10

Rose

3

→

Rose

5

Sunflower

Mean or Median

Rose

8

Outliers

## ② Measure of Dispersion

① Variance

② Standard deviation

### ① Variance

#### Population Variance

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

$x_i$  = Data points

$\mu$  = Population Mean

$N$  = Population size

#### Sample Variance

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

Bessel's  
Correction  
↗

WHY DOES THE SAMPLE VARIANCE HAVE  $N-1$  IN THE DENOMINATOR? The reason we use  $n-1$  rather than  $n$  is so that the sample variance will be what is called an unbiased estimator of the population variance

$x_i$  = Data points

$\bar{x}$  = Sample mean

$n$  → Sample size

Eg:  $\{1, 2, 3, 4, 5\} \Rightarrow$  Sample

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

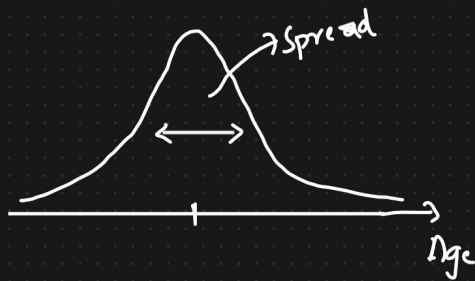
$$\bar{x} = \frac{1+2+3+4+5}{5} = 3$$

$x_i$	$\bar{x}$	$(x_i - \bar{x})^2$
1	3	4
2	3	1
3	3	0
4	3	1
5	3	4

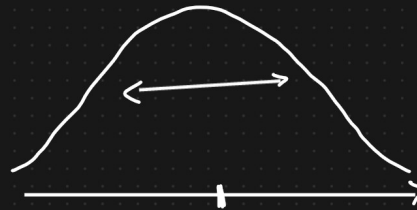
$$s^2 = \frac{10^5}{4^2} = \underline{\underline{2.5}}$$

Variance : Spread of the data.

$$s^2 = \underline{\underline{2.5}}$$



$$s^2 = \underline{\underline{6.5}} \uparrow \uparrow$$



## (2) Standard deviation

Population std

$$\sigma = \sqrt{\text{Variance}}$$

Sample std

$$s = \sqrt{\text{Sample variance}}$$

$$s^2 = 2.5$$

$$\sqrt{s^2} = \text{Sample Std}$$

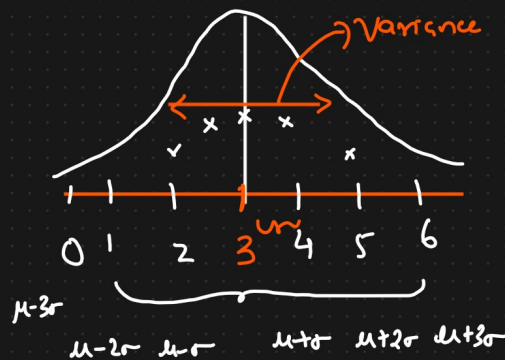
Consider

$\{1, 2, 3, 4, 5\}$

$$\rightarrow \mu = 3 \checkmark$$

$$\rightarrow \sigma = 1 \checkmark$$

5





### ③ Random Variables

$$\begin{cases} x + 5 = 7 \\ y + x = 10 \end{cases} \quad \begin{array}{l} x = 2 \\ y + 2 = 10 \\ y = 8 // \end{array}$$

Random Variable is a process of mapping the output of a random process or experiment to a number

Eg: Tossing a Coin

Rolling a dice

Measure the Temperature for the next day

$$X = \begin{cases} 0 & \text{if H} \\ 1 & \text{if T} \end{cases} \quad \begin{array}{l} \text{Quantifying a Random} \\ \text{Process} \end{array}$$

$$Y = \begin{cases} \text{Sum of the rolling of dice 7 times} \\ \{4, 5, 6, 1, 2, 2\} = 20 \end{cases}$$

$$\underline{P(Y \geq 15)}$$

$$\underline{P(H)}$$

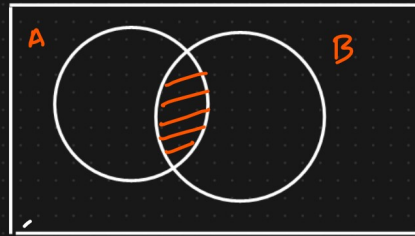
#### ④ Sets

$$A = \{1, 2, 3, 4, 5, 6, 7, 8\}$$

$$B = \{3, 4, 5, 6, 7\}$$

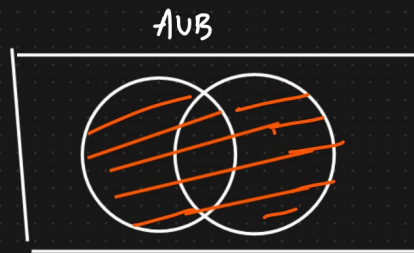
##### ① Intersection

$$A \cap B = \{3, 4, 5, 6, 7\}$$



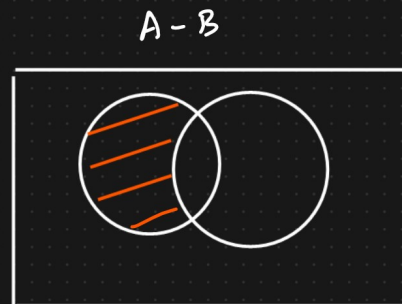
##### ② Union

$$A \cup B = \{1, 2, 3, 4, 5, 6, 7, 8\}$$



##### ③ Difference

$$A - B = \{1, 2, 8\}$$



##### ④ Subset

$$A \rightarrow B \Rightarrow \text{FALSE}$$

$$B \rightarrow A \Rightarrow \text{TRUE}$$

##### ⑤ Superset

$$A \rightarrow B \Rightarrow \text{True}$$

$$B \rightarrow A \Rightarrow \text{False.}$$