

Basics To Intermediate Stats

- ① Measure of Central Tendency
- ② Measure of dispersion
- ③ Gaussian Distr
- ④ Z score
- ⑤ Standard Normal Distr

① Arithmetic Mean for Population & Sample

Mean (Average)

Population (N)

$$\downarrow$$

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

Sample (n)

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

$$= \underline{\underline{3.2}}$$

$$\mu = \sum_{i=1}^N \frac{x_i}{N}$$

$$= \frac{1+1+2+2+3+3+4+5+5+6}{10}$$

$$= \frac{32}{10} = 3.2$$

④ Central Tendency

- ① Mean ✓
- ② Median ✓
- ③ Mode ✓

Refers to the measure used to determine the centre of the distribution of data.

$\{1, 1, 2, 2, 3, 3, 4, 5, 5, 6, \boxed{100}\}$

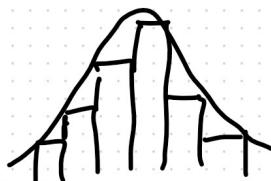
$$\text{Mean} = \frac{32 + 100}{11}$$

$$= \frac{132}{11} = 12$$

$$\begin{aligned} M &= 3.2 \\ &+ 100 \\ M &= 12 \end{aligned}$$

$\xrightarrow{\quad}$ outliers

Distribution



Median

$\{1, 1, 2, 2, 3, \boxed{4}, 5, 5, 6, \boxed{100}, 112\}$

1) Sort the numbers

Odd number = 11

$$\boxed{M=12}$$

$$\left\{ \begin{array}{l} \text{Median} = 3 \\ \hline \text{Median} = 3.5 \end{array} \right\} \text{ Avg} = \frac{3+4}{2} = 3.5$$

$$\begin{array}{l} M = 3.2 \\ \hline M = 12 \end{array} \rightarrow \boxed{\text{Median} = 3} \rightarrow \boxed{\text{Median} = 3.5}$$

$\{$ Median works well with outlier $\}$

$\xlongequal{\quad}$

② Mode = $\{1, 2, \underbrace{2, 2}_{2}, 3, 4, 5, \underbrace{6, 6, 6}_{3}, 7, 8, \underbrace{100, 100, 100, 100}_{2}\}$

Mode = {Most frequent Element}

$\xlongequal{\quad}$

$\boxed{\text{Mode} = 6} \rightarrow \text{Measure of Central Tendency}$

Mode

Type of flower

Rose

Lily

Sunflower

-

petal length

petal width

DATA SET
 \downarrow

10% Missing data

$\left\{ \begin{array}{l} \text{Missing} \\ \text{Value} \end{array} \right\} \rightarrow \text{Most frequent element}$

\Downarrow

Ages of Students }
 =
 Age
 25
 26
 =
 =
 32
 34
 38

Mean? ✓
 Median?
 Mode ??

(f) Measure of Dispersion

① Variance

② Standard Deviation

① Variance ✓

$$\mu \stackrel{=} \rightarrow \{1, 1, 2, 2, 4\} \quad \frac{80}{5} = 2$$

$$\mu \stackrel{=} \rightarrow \{2, 2, 2, 2, 2\} \quad \frac{10}{5} = 2$$

↓
Dispersion
Spread

Population Variance

$$\sigma^2 = \frac{N}{n} \sum_{i=1}^N (x_i - \mu)^2 = \frac{10.84}{6} = 1.81$$

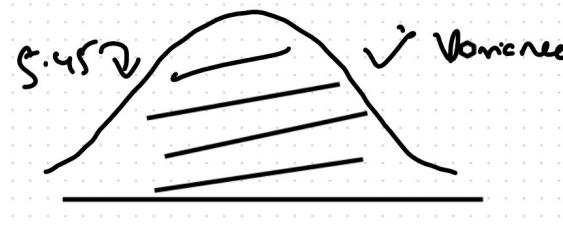
x	μ	$x - \mu$	$(x - \mu)^2$
1	2.83	-1.83	3.34
2	2.83	-0.83	0.6889
2	2.83	+0.83	0.6889
3	2.83	0.17	0.03
4	2.83	1.17	1.37
5	2.83	2.17	4.71
	$\mu = 2.83$		10.84

Sample Variance

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



Variance is more ??



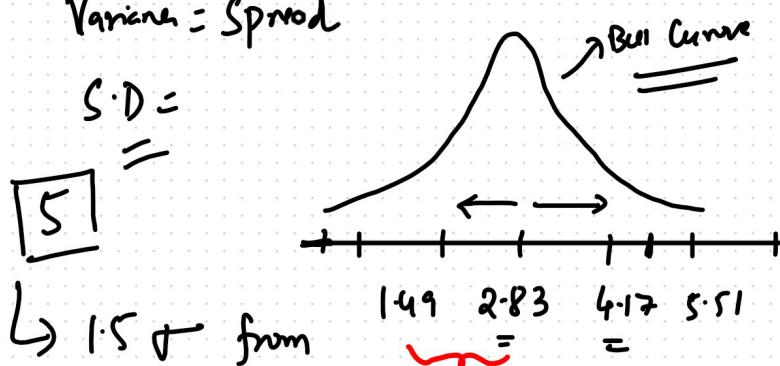
$$\sigma = \sqrt{\text{Variance}} = \sqrt{1.81} = 1.345$$

Variance = Spread

S.D. =

5

$$\begin{array}{r} 2.83 \\ 1.34 \\ \hline 4.17 \end{array} \quad \begin{array}{r} 2.83 \\ 1.34 \\ \hline 1.49 \end{array}$$



$$\begin{array}{r} 4.17 \\ 1.34 \\ \hline 5.51 \end{array}$$

the μ

④ Percentiles And Quartiles {Find outliers?}

Percentage : 1, 2, 3, 4, 5

% of the numbers that are odd?

% = # of numbers that are odd

$$\begin{aligned} & \text{Total Numbers} \\ &= \frac{3}{5} = 0.6 = 60\% \end{aligned}$$

Percentiles (GATE, CAT, GMAT, SAT)

↳ Dfn : A percentile is a value below which a certain percentage of observations lie.

Data set : 2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12

11?

What is the percentile ranking of 10? $n=20$

$x=10$
Percentile Rank of $x = \frac{\# \text{ of values below } x}{n} \times 100$

11?

$$= \frac{4}{20} \times 100 = 20\%$$

$$= \frac{17}{20} \times 100 = 85\%$$

② What value exists at percentile ranking
of 25%? $\boxed{75\%}$??

$$\text{Value} = \frac{\text{Percentile}}{100} \times (n+1)$$
$$= \frac{25}{100} \times (21) = 5.25 \rightarrow \text{Index Position} = \boxed{5} \rightarrow 25\%$$
$$= \frac{75}{100} \times (21) = 15.75 \Rightarrow \boxed{9} \text{ Answer} \rightarrow \text{index} =$$

Five Number Summary

- ① Minimum
- ② First Quartile (Q1)
- ③ Median
- ④ Third Quartile (Q3)
- ⑤ Maximum

Removing the outliers

Outlier?? [27 > 13]

{ 1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 27 }

$$25\% \quad \frac{25}{100} \times (19+1)$$

~50 ✓

[Lower fence \longleftrightarrow Higher fence]

$$\frac{25}{100} \times 20 = \underline{\underline{5}} \Rightarrow \text{index} = 3 \Rightarrow$$

$$\text{Lower fence} = Q_1 - 1.5(\text{IQR})$$

$$Q_1 = \underline{\underline{3}} \checkmark$$

$$\text{Upper fence} = Q_3 + 1.5(\text{IQR}) \quad Q_3 = (75\%)$$

$$Q_3 = \underline{\underline{7}} \checkmark$$

$$\text{Interquartile Range (IQR)} = Q_3 - Q_1$$

$$= 7 - 3 = \underline{\underline{4}}$$

$$\text{Lower Fence} = Q_1 - 1.5(\text{IQR})$$

[Lower Fence \longleftrightarrow Higher Fence]

$$= 3 - 1.5(4)$$

$$[-3 \longleftrightarrow 13]$$

$$= 3 - 6 = \boxed{-3} \checkmark$$

$$\text{Higher Fence} = Q_3 + 1.5(\text{IQR})$$

$$= 7 + 1.5(4)$$

$$= 7 + 6 = \underline{\underline{13}}$$

Remaining data

1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 27

$$\text{Minimum} = 1$$

$$Q_1 = 3$$

$$\text{Median} = 5$$

$$Q_3 = 7$$

$$\text{Max} = 9$$

DATA

Visualisation

→ 5 Number Summary

Box plot

27

Box plot



5

27

$$\sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1} \rightarrow \left. \begin{array}{l} \text{Bessel's correction} \\ \text{Degree of freedom} \end{array} \right\}$$

STATS

Why sample variance is $n-1$?