

What is Statistics?

Science of collecting, organising and analyzing data → facts or pieces of information.

Types of Statistics

Descriptive Stats → organising and summarising data

Inferential Stats → consist of using data you have "measured" to form conclusion

What comes under what type of statistics

Descriptive

Inferential

- | | |
|---|-----------------------|
| ① Measure of central tendency (mean, median, mode) | ① Z-test |
| ② Measure of dispersion (Variance, Std Deviation) | ② t-test |
| ③ Different type of distribution of data. (PMF, CDF, Histogram) | ③ Chi Square Test |
| | ④ Hypothesis testing |
| | P-value, significance |

Population and Sampling

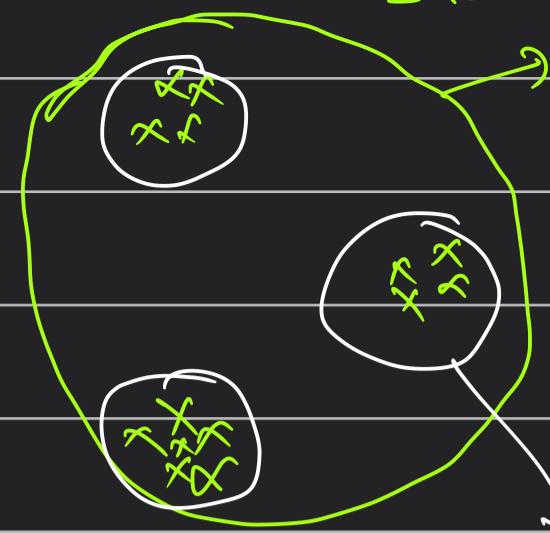
Population :- the group we are interested in studying

Sample :- a subset of population.

for Ex - Exit poll \sum Media \Rightarrow Will try to predict who will win?

Suppose a state 'A' is having election.

... ^ = Population



The media went to these region

and asked who did you vote

based upon response from these regions media

will predict who will win? These regions can be called Sample's.

Remember one thing Exit polls are not always free....

• Types of Sampling Techniques

Technique used to create a sample that is representative of entire population these techniques are called sampling techniques.

[Population is represented by " N " and sample is represented by " n "]

① Simple Random Sampling

We randomly select a sample " n " from population " N ".

When performing this, every member of population (N) has an equal chance of being selected for your sample (n)

② Stratified Sampling

\Downarrow keeping \Rightarrow Non overlapping groups

Population

Let's say I wanna do a survey where I need info from male group. (separate grouping, where groups don't have common thing to overlap) then we can use Stratified sampling.

③ Systematic Sampling

Involves selecting k^{th} element from population to be included in sample, where k is the sampling interval or skip interval.

Ex Select every 4th guy in a mall to fill a survey.

④ Convenience Sampling

↳ voluntary response sampling

sample is selected based on convenience of accessibility. Involves things which are readily available & easy to reach.

[We use different sampling technique based upon our different use cases]

■ What are variables and its types?

is a property that can take on many values.

Variable \Rightarrow Singular Mode.

Ages = [2, 10, 12, 14] \Rightarrow not a variable.

• Types of Variable

Quantitative Variable

↓
Discrete

use discrete
number
eg
No of bank
accounts,
No of
children
in family

Continuous

is continuous
in nature.
can be upto
infinity
Eg Height
153.7, 165.5
or weight

Qualitative / Categorical

↓

Based upon some property
we classify the variable

Eg Gender [male
↓ Female

Types of flower [Rose
↓ Lily
↓ Gheetes

We divide these things
based upon some characteristics

■ Measure of Central Tendency - Mean, Median & Mode

Central tendency refers to measure used to determine the "center" of distribution of data.

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

Mean = Average

↓

Population Data (N)

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

population
mean

Sample (n)

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

sample
mean

Median = Central value

(most called
measure of central
tendency)

① Sort all numbers

odd length = middle element

② Find central element

even length = Average of
two central
elements.

Mode :- most frequent elements

$\{1, 2, 2, 3, 3, 3, 4, 4, 5, 5, 5, 5, 5, 17, 9\}$ - 5 is repeated most so 5.

Use Cases

- (1) Outliers \Rightarrow median
- (2) Categorical features \Rightarrow mode, based upon if some category is repeating more.

Measure of Dispersion - Variance & Standard Deviation

Measure of Dispersion talks about spread.

↳ Determined two things Variance

& Standard deviation

We measure these things around Population & Sample.

Variance of Population

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

population
mean

Sample Variance

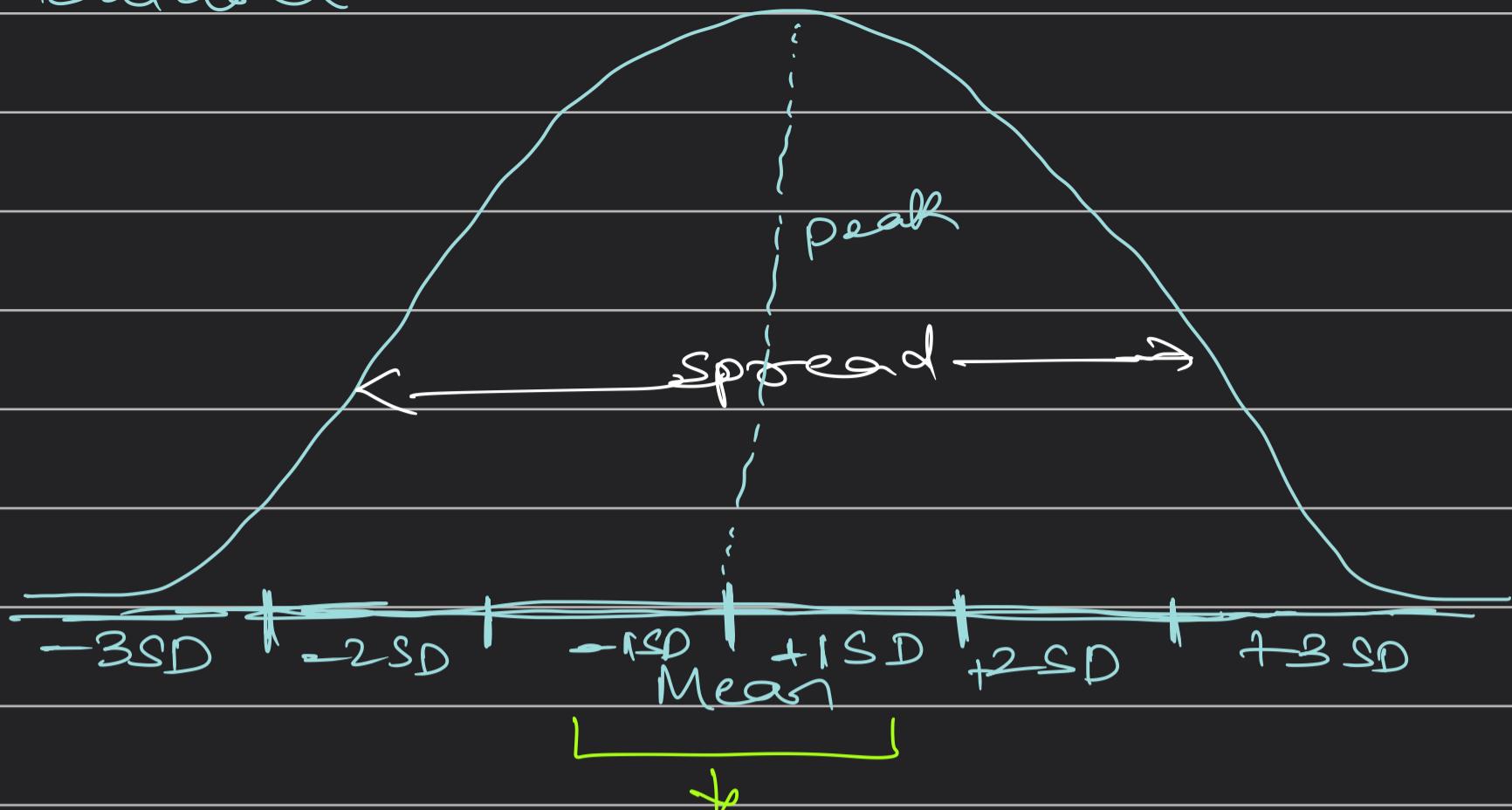
$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

Bessel's Correction
Degree of Freedom

Standard Deviation $= \sqrt{\text{Variance}}$

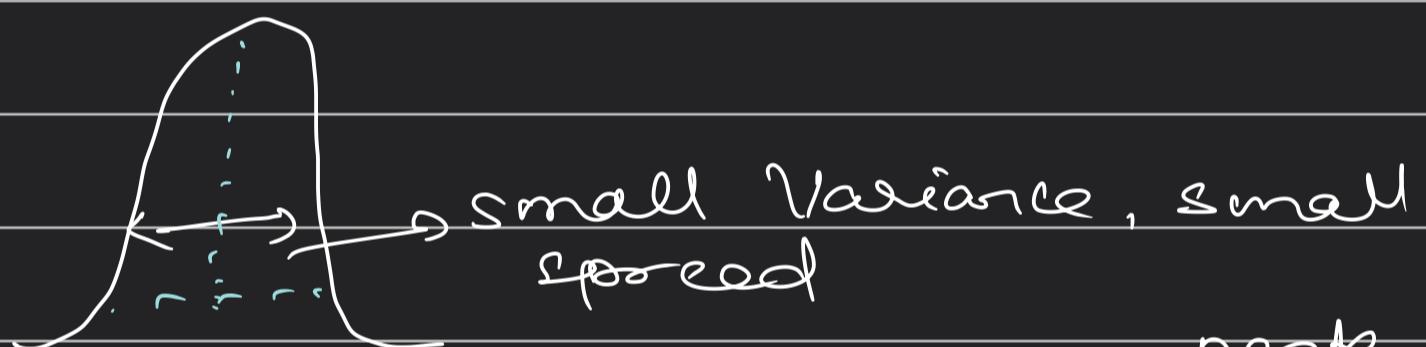
Let understand what these mean

Dataset



Any value around $+1SD$ $-1SD$ of the mean will lie here.

The bigger the "Variance" the bigger the spread of data



The less the S.D more the "height" / ^{peak} of bell curve and less the "spread".

Percentiles & Quartiles

is a value below which a certain percentage of observations lies.

Ex If say I got 95 percentile. this means that I have better marks than 95% of other students or 95% of people got marks less than me.

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

What is percentile rank^r of 10?

$$\text{percentile rank}^r = \frac{\# \text{No of elements below } 10}{\text{total number of elements}} \times 100$$

$$= \frac{16}{20} \times 100 = 80 \text{-percentile}$$

So rank of "10" is 80 percentile that means 80% of values in dataset are below 10.

(\Leftarrow) Value = $\frac{\text{Percentile}}{100} \times (n+1)$

? What value of 25 percentile
 $= \frac{25}{100} \times 21 = (5.25)$ \hookrightarrow Index

5th if 6th value - average

Quartile \rightarrow 25th percentile - Ist Quartile
50th percentile - IInd Quartile
75th percentile - IIIrd Quartile

Inter Quartile Range = 3rd Quartile - 1st Quartile

used for finding or detecting outliers.

How to construct Box plots for outliers.

5 Number Summary and Box Plot

- (1) Minimum
- (2) first Quartile (25or.) Q1
- (3) Median
- (4) Third Quartile (Q3)
- (5) Maximum

These 5 component help us plot a Box plot which in turn help us find an "Outlier".

→ any value that is different from dataset.

Technique to find Outlier

[Lower Fence \longleftrightarrow Higher fence]

Any value smaller than lower fence & bigger than higher fence would be considered as an outlier.

How to determine lower fence & higher fence

$$\Rightarrow \text{lower fence} = Q1 - 1.5 (\text{IQR})$$

$$\text{higher fence} = Q3 + 1.5 (\text{IQR})$$

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

$$Q1 = \frac{25}{100} \times (19+1) \quad \begin{matrix} 19 \text{ values} \\ \text{Data should be sorted in order} \end{matrix}$$

$$Q1 = 5^{\text{th}} \text{ Index}$$

$$Q1 = 3$$

$$Q3 = \frac{75}{100} \times (20) = 15^{\text{th}} \text{ Index} = 7$$

$$IQR = Q_3 - Q_1 = 7 - 3 = 4$$

$$\text{lower fence} = Q_1 - 1.5 \times IQR$$

$$= 3 - 6 = -3$$

$$\text{upper fence} = Q_3 + 1.5(IQR) = 7 + 6 = 13$$

$$[-3 \longleftrightarrow 13]$$

outlier
↑

One value is greater than 13 i.e. 27

Similarly our 5 numbers will be

① Minimum = 1

② First Quartile = $Q_1 = 3$

③ Median = 5

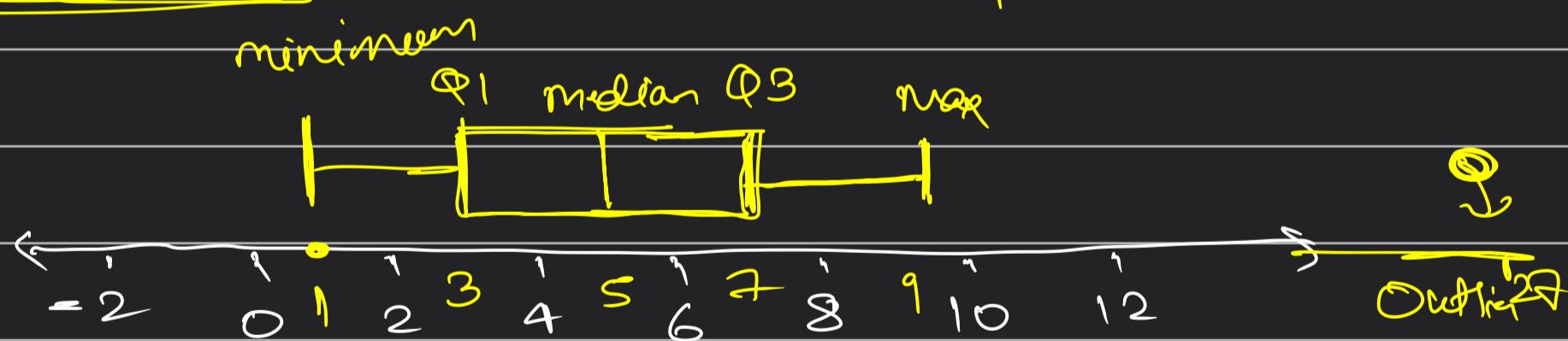
④ Third Quartile = $Q_3 = 7$

⑤ Maximum = 27 = 9

~~outlier removed~~



Using this we plot our Box Plot



- finding outliers using these techniques using python. [Cheek criterias] Sack pccel Akyay Data Science.

