

- ⑥ Seaborn KDE plot: Kernel Density Estimation plot.
 - combination of histogram and distribution.
 - `sns.kdeplot(df['sepal length'])`
- ⑦ Distplot: Histogram + KDE Plot.
- ⑧ KDE:
 - for col in ['sepal length', 'sepal width', 'petal length', 'petal width']:
 - `sns.kdeplot(df[col], shade=True)`
- ⑨ pairplot: plots pairwise relationships in a dataset.
 - `sns.pairplot(df, hue='species', size=2.5);`
- ⑩ HeatMap:
 - `flights_long = sns.load_dataset('flights')`
 - `flights = flights_long.pivot("month", "year", "passengers")`
 - `f, ax = plt.subplots(figsize=(9, 6))`
 - `sns.heatmap(flights, annot=True, fmt="d", linewidth=.5, ax=ax)`
- ⑪ boxplot: `sns.boxplot(y=df["Exchange Rate"])`

Business Statistics

- Statistical Analysis is meant to collect and study the information available in large quantities
- Branch of Mathematics where computation is done over a bulk of data.
- Data collected for analysis here is called Measurements
- if we need to measure the data, a sample is taken out a Population
-
- ① Types of Data:
 - ① Descriptive stats describes characteristics of dataset.
 - ② three basic categories of measure:
 - ① Measures of Central Tendency
 - ② Measures of Variability
 - ③ frequency distribution
 - ③ Central Tendency: Center of dataset (Mean, Median, Mode)
 - ④ Variability: Variance, (Standard Deviation)

→ Measures of frequency distribution describe the occurrence of data within a dataset (Count)

① Mean = Avg = $\frac{\sum(i)}{N} = \frac{25+30+15+25+35}{5} = 26$

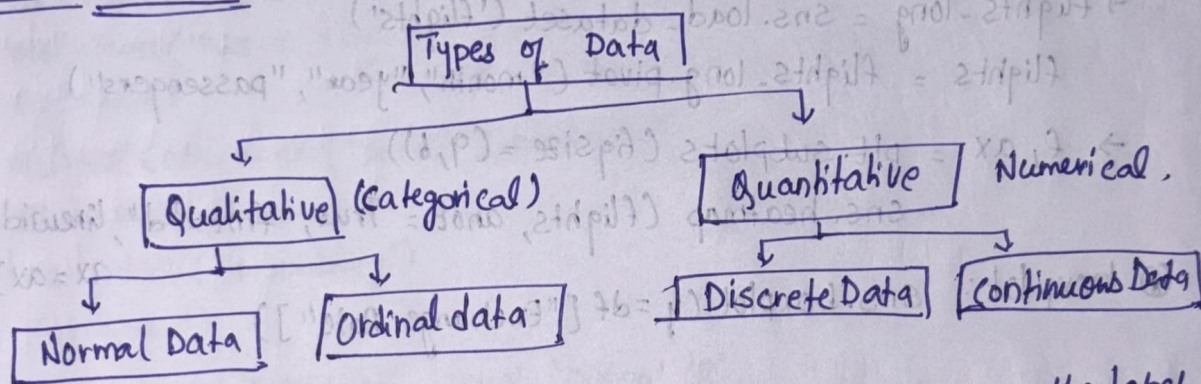
② Median = 15, 25, 25, 30, 35

Median.

③ Mode: 25 → Max no. of occurrences.

→ Standard Deviation (σ) = $\sqrt{\frac{\sum(x_i - \mu)^2}{N}}$
variance.

④ Inferential statistics:



→ Nominal Data: type of qualitative information which helps the label variables without providing numerical value.

→ Ordinal Data: They can be ordered → High, Mid, low

→ Used Mostly in surveys, finance → These data are called ordinal data

⑤ Quantitative: Numerical data represent numerical value

① Discrete Data: Takes only discrete values.

Ex: No. of workers in a Company.

Ex: 40, 45, 50, ...

② Continuous Data: Data that can be counted

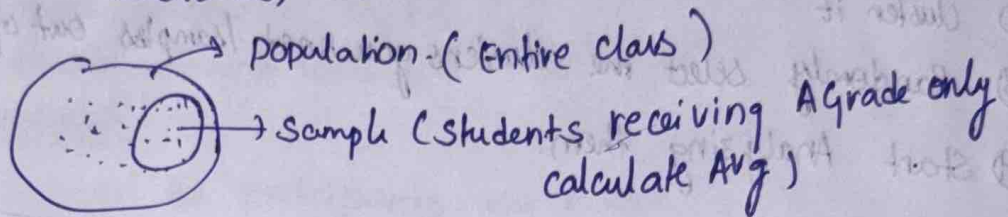
Ex: Temperature Range (97.3, 103.2, ...)

Sampling Techniques

① Population: The Entire group you want to draw conclusions about.

② Sample: Subset of population. The specific group of individuals that you will collect data from.

→ In real world scenarios, whenever work on population data



→ why sampling is important?

- ① Make information faster.
- ② Easily Analyse the data
- ③ Lesser Data Collection Error.

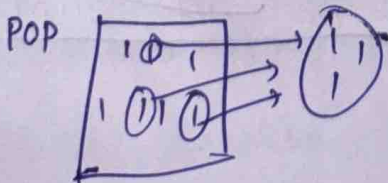
→ Types of sampling:

① Probability sampling

→ Every member of population has an equal chance of being selected.

Examples:

- Simple random sampling
- Stratified random sampling
- Cluster sampling
- Systematic sampling



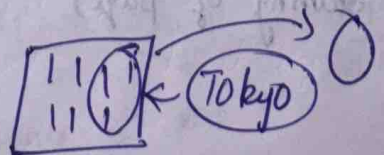
② Non probability sampling

→ selection basis of judgement (or) the convenience of accessing data.

→ Largely depends on a researchers sample selection skills

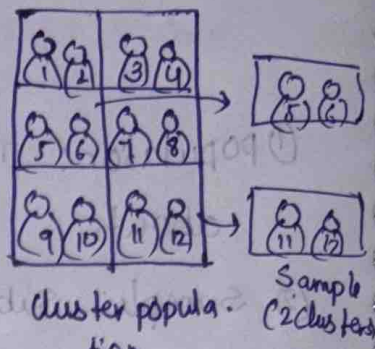
→ Examples:

- convenience sampling
- purposive sampling
- voluntary response
- snow ball sampling



* Cluster random Sampling:

- ① Divide population into groups
- ② Then randomly select group from all the groups

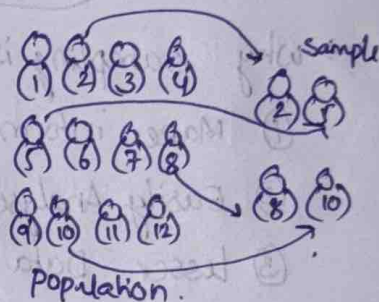


* Steps to form clusters:

- ① Population → diverse, No repetition in clusters
→ cluster should cover the entire population Members
- ② cluster it
- ③ Randomly select the clusters you need / samples out of it.
- ④ Start Analyzing them.

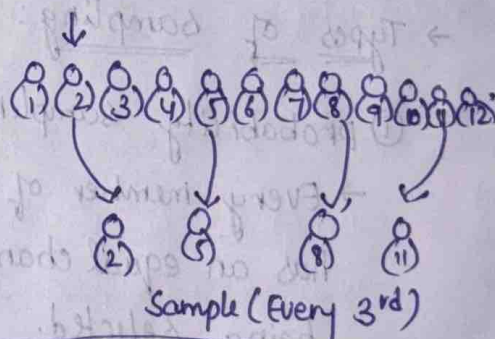
* Random Sampling:

- Randomly chose a member from the population
- Every member and set of member has an equal chance of being selected



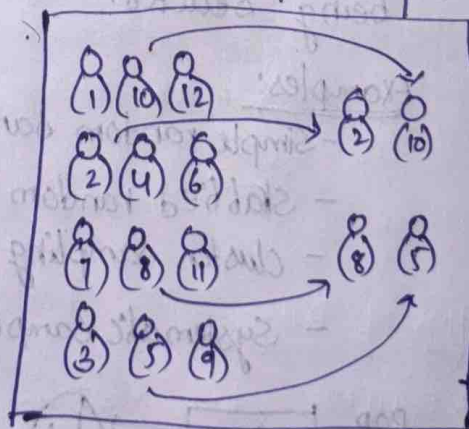
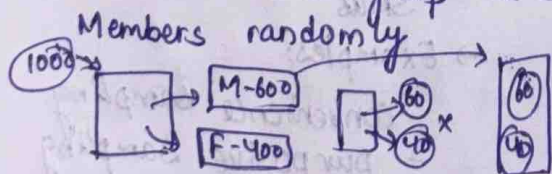
* Systematic Sampling:

- put a member of the population in some order and a starting point is choose as random the every "nth" member is selected to be in a sample.



* Stratified Sampling

- first divide population into groups
- Then from each group we select Members randomly



* Cluster Random sampling (explained starting of page)

* Non probabilistic Sampling

- ① convenience Sampling: sampling based on the survey form. Include respondents/member who are easy to reach for researcher.
- ② purposive Sampling: select sample based on purpose of research
→ Researcher select sample by using their expertise

③ voluntary Response and Knowledge Sampling:

- Based on ease of access
- people volunteer to it, ex: LinkedIn polls.

④ Snowball Sampling:

- Recruit the participants via research participants for test (or) study.
- used where its hard to find potential population for research.
- like a subcontract.

* population Sampling

- Analysing or testing entire population is impossible and also a cost & time taking. we use sample
- subset of population which represent entire population.
- Errors can lead to inaccurate and misleading result

⑤ Standard Deviation: Measure of how spread the numbers are. \sqrt{v}
(σ)

⑥ Variance: $\frac{\sum (x_i - \mu)^2}{N}$ $\mu \rightarrow$ Avg/Mean

Example: x_i

$$\mu = \frac{\sum x_i}{N} = \frac{80}{5} = 16$$

⑦ for sampling: it is going to be (N-1) instead of N.

x_i	$x_i - \mu$	$(x_i - \mu)^2$
5	5 - 16 = -11	121
10	-6	36
15	-1	1
20	4	16
30	14	196
		$\sum (x_i - \mu)^2 = 370$

$$\sigma = \sqrt{\frac{\sum (x_i - \mu)^2}{N}}$$

$$= \sqrt{\frac{370}{5}}$$

$$\sigma = \sqrt{74}$$

Var

But why N-1 and Not N?

→ Bessels correction.

- corrects bias in estimation of population variance
- partially corrects bias in estimation of population SD