

Essential Statistics for Data Science

Topics:

Inferential Statistics, Descriptive Vs Inferential Statistics.

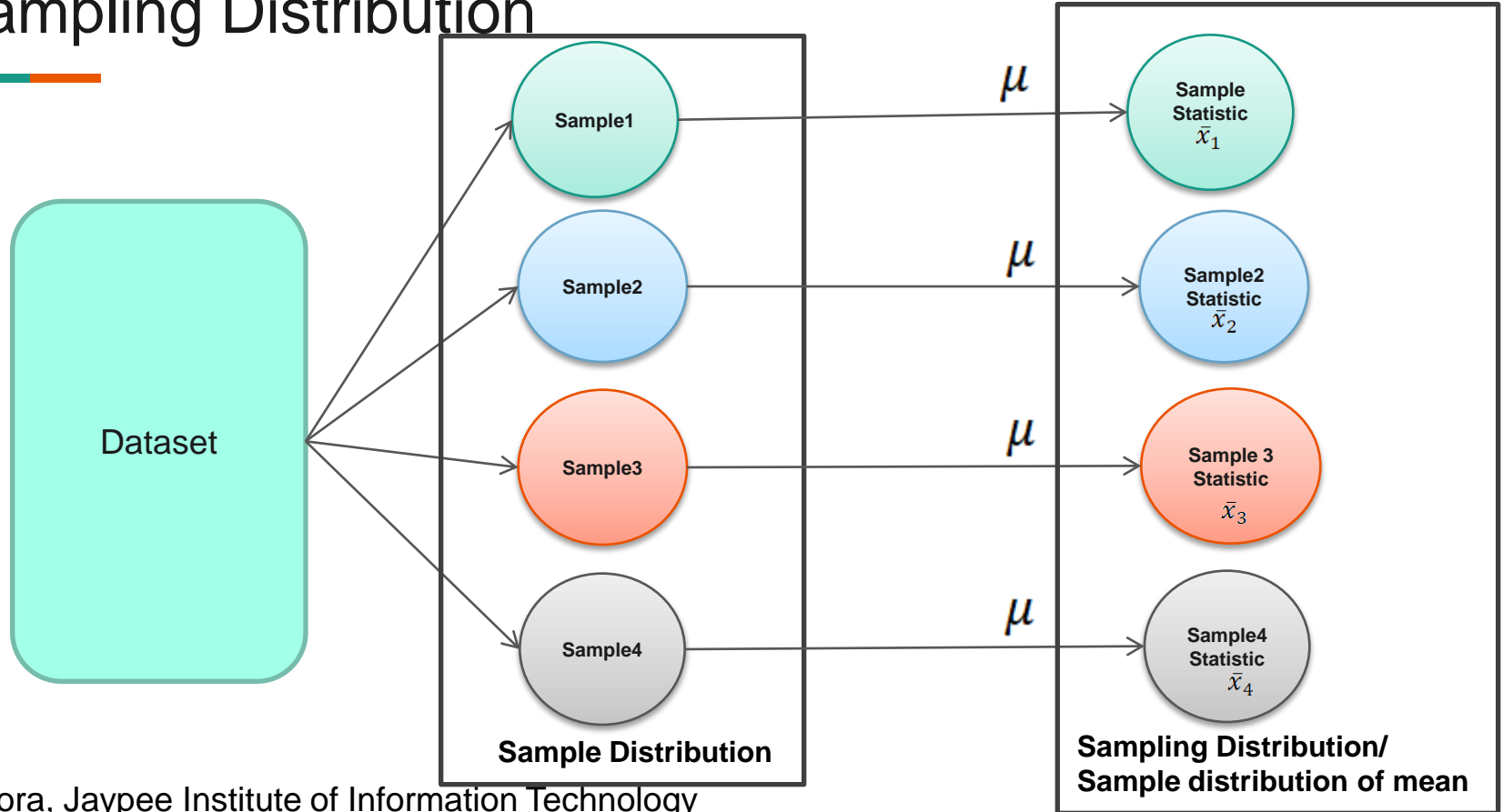
Sampling Distribution, Central Limit Theorem



Deterministic Vs Inferential Statistics

Deterministic Statistics	Inferential Statistics
Measure for definitive measurement	Note the margin of error of research performed
Data Summary: Bar Graph, Histogram, pie chart, etc.	Data Summary: Takes sample data and makes inferences about the sample population. i.e used to make inference or draw a conclusion of population.
Measures of Central Tendency: Mean, Median, Mode Measures of dispersion: Range, Variance, standard deviation	Uses probability to determine how confident we can be that the conclusion we make is correct Normal (Gaussian distribution), Binomial distribution, etc.

Sampling Distribution



Sampling Distribution

Height
of people in
different states of
India

N Population

$$\mu = \frac{x_1 + x_2 + x_3 + \dots + x_n}{N}$$
$$\sigma = \sqrt{\frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N}}$$

Sample Distribution

Delhi: $x_{d1}, x_{d2}, x_{d3}, \dots, x_{dn}$

MP: $x_{mp1}, x_{mp2}, x_{mp3}, \dots, x_{mpn}$

.....

.....

$$\overline{x_d} \quad \sigma_d$$
$$\overline{x_{mp}} \quad \sigma_{mp}$$

Sampling distribution

$$\text{mean}(\bar{x}) \approx \mu$$

$$SD(\bar{x}) < \sigma$$

Central Limit Theorem (CLT)

- Three mean values:

$\mu \rightarrow$ original Population Mean

$\overline{x_{mp}}, \overline{x_d} \rightarrow$ Sample Mean varies for each sample

$\bar{x} \rightarrow$ mean of Sample means

- **CLT Part-I:** Sample is taken from population, So mean of Sample will approximately be same as of original population. Sample size (large/ small) does not matter. Even mean of samples will also be approximately same as shown below

$$\bar{x} \approx \mu \approx \overline{x_{mp}}, \overline{x_d}$$

Central Limit Theorem (CLT)

- **CLT Part-II:** The standard deviation of the sample means decreases as the sample size increases i.e

$$SD(\bar{x}) < \sigma$$

Its formula looks like

$$SD(\bar{x}) = \frac{\sigma}{\sqrt{n}}$$

this is known as standard error as well



n increase



Standard error decrease

- **CLT Part-III:** If the population follows a normal distribution, then the sampling distribution is also normal distribution. The distribution of the sample means approaches a normal distribution, under certain conditions,

$$\bar{x} \sim N(\mu_X, \sigma_X / \sqrt{n})$$

Online simulation Link: <http://www.ltconline.net/green/java/Statistics/clt/cltsimulation.html>

CLT Theorem:

Distribution of sample statistics is nearly normal, centred at the population mean, and with a standard deviation equal to the population standard deviation divided by square root of sample size.

$$\bar{x} \sim N(\mu_X, \sigma_X / \sqrt{n})$$
