* Population vs Sample
* Sample Statistics
  - Mean
  - Variance
  - SD

* Point Estimates
* Sampling Distribution
* Standard Error

* Uniform Distribution, ← PMF
                          ← PDF
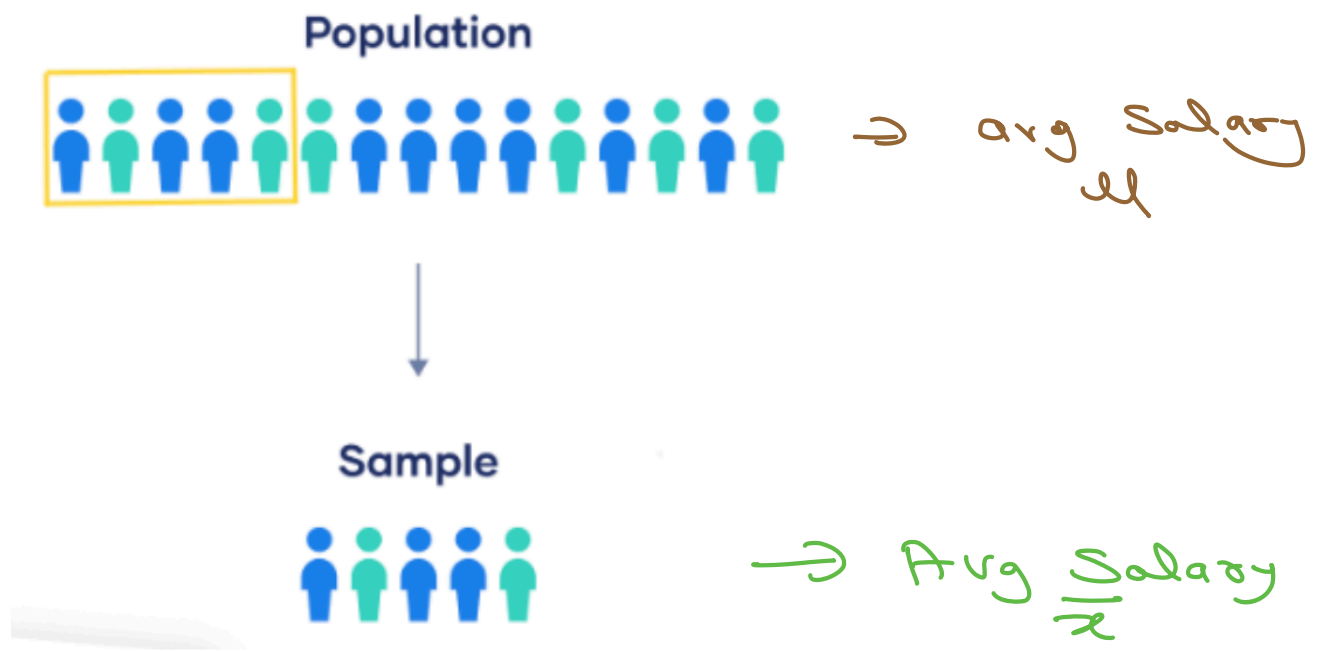                          ← CDF

Population vs Sample

BLR → 1,00,000 Data Scientist

Avg salary and σ of Salary

Survey ⇒

Population ⟹ 1,00,000

⟹ Reaching Everyone for Survey is impractical and Costly

Population



⟹ avg Salary
$\mu$

Sample

→ Avg $\dfrac{Salary}{x}$

* Sampling is the Solution for large populations

# Sample Statistics

| Population Stats | Sample Stats |
|---|---|
| ① $\mu$ ⇒ $\dfrac{\sum_{i=1}^{n} x_i}{n}$ | ① $\bar{X}$ ⇒ $\dfrac{\sum_{i=1}^{n} x_i}{n}$ |
| ② $\sigma^2$ (Pop Variance) ↓ $\sum_{i=1}^{n} \dfrac{(x_i - \mu)^2}{n}$ ↓ $\sqrt{\sigma^2}$ | ② $S^2$ (Sample Var) ↓ $\sum_{i=1}^{n} \dfrac{(x_i - \bar{X})^2}{n-1}$ |
| ③ $\sigma$ ⇒ $\sqrt{\sigma^2}$ | ③ $S$ ⇒ $\sqrt{S^2}$ |

* $n-1$ is called Bessel's correction
and it is use for correcting the
Bias due to sample

$$\bar{X} ⇒ 6/3 ⇒ 2K$$
$$n ⇒ ③$$

| $x_i$ | | $x_i - \bar{X}$ |
|---|---|---|
| 1 | 1000 | -1000 |
| 2 | 2000 | 0 |
| 3 | 3000 | +1000 |

Sample Data

[131, 150, 140, 142, 152]

$$\bar{X} \implies \frac{(131 + 150 + 140 + 142 + 152)}{5} ?$$

$$\downarrow$$
143

$$\boxed{S^2} \implies \frac{x_i - \bar{X}}{4}$$

| | |
|---|---|
| A | $\bar{X} = 143, \sigma^2 = 71$ |
| B | $\bar{X} = 142.4, \sigma^2 = 66.3$ |
| C | $\bar{X} = 147, \sigma^2 = 73.2$ |
| D | $\bar{X} = 152, \sigma^2 = 64.5$ |

\* Estimate the Avg Salary of all DS in Bengaluru

$$\downarrow$$

$$\boxed{\text{Point Estimate}}$$

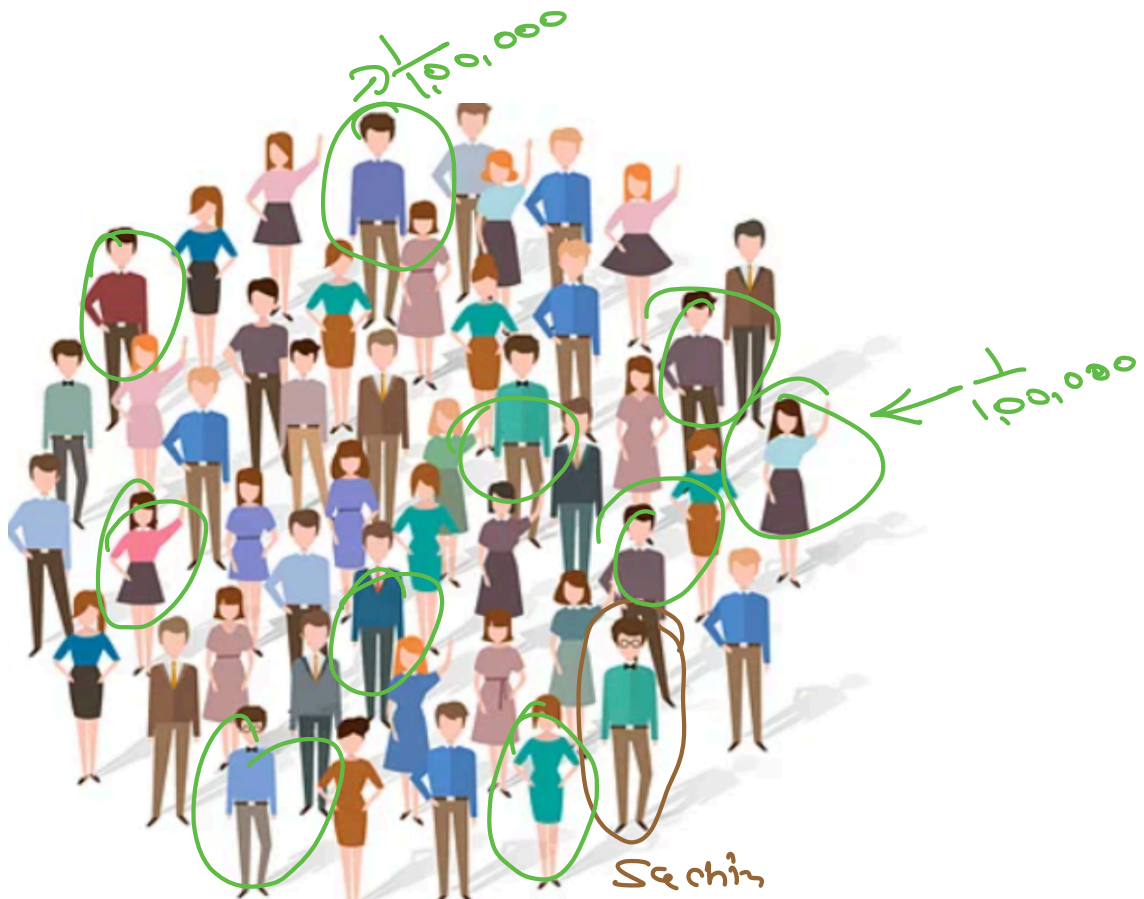① Estimating Population Statistics based on Sample Statistic

Sachin

# BiasNess in Sampling

①  All people we sampled can
    have lower Salary then AVG

②  All people we sampled can
    have Higher Salary then AVG

## Sampling Technique

① Probabilistic Sampling

② Non-Probabilistic Sampling



$2 \frac{1}{1,00,000}$

$\frac{1}{1,00,000}$

Sachin

Pop → 10.0000

Prob of 1 DS → $\frac{1}{1,00,000}$

# Simple Random Sampling

**( Every Entity in population Ras Equal chance of Getting picked )**

There are four main types of probability samples.

1. **Simple random sampling** ✓
2. **Systematic sampling**
3. **Stratified sampling**
4. **Cluster sampling** } ml

## Steps for Point Estimates

**Step 1 :** Define population
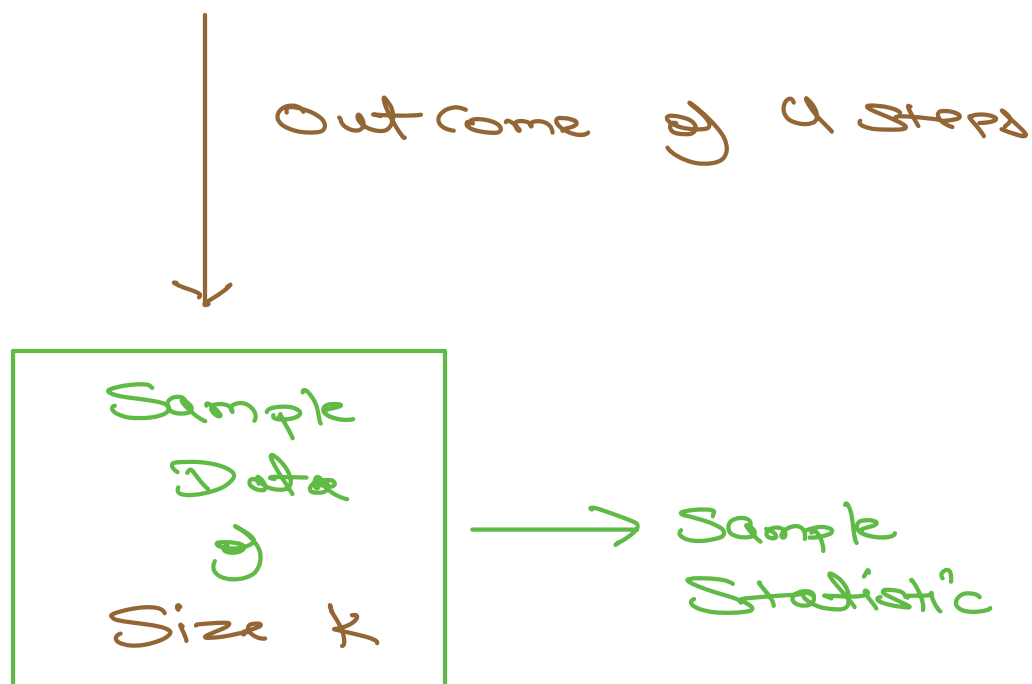
**Step 2 :** Determine Sample Size

    ⊙ certainity ⇒ **95%**

    ⊙ Error of Margin⊙ 5%

Let's say $\xrightarrow{\text{for } 95\%}$ (5,000)

Step 3: Randomly Sample the Population Based on Size

Step 4: Collect Data of all the Samples

Outcome of 4 Steps

Sample Data of Size k → Sample Statistic
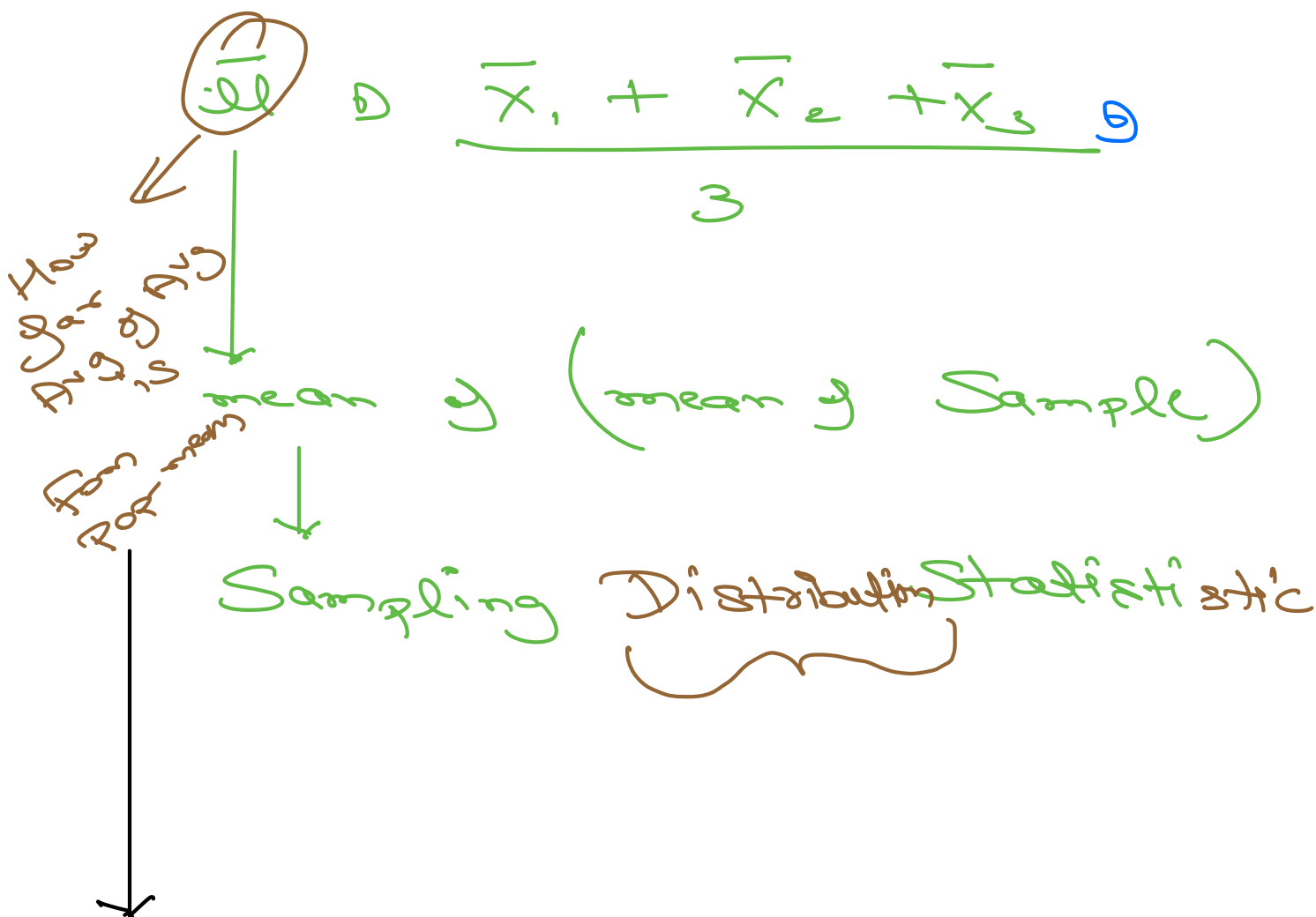
Ⓠ What if we do all 4 steps multiple times

* Later

CLT ⟹ Central Limit Theorem

| Sample 1 $\bar{X}_1$ $S_1^2$ | Sample 2 $\bar{X}_1$ $S_2^2$ | Sample 3 $\bar{X}_1$ $S_3^2$ |
|---|---|---|

$\boxed{\bar{\mu}}$

How far Avg of Avg is from Pop mean

① $\dfrac{\bar{X}_1 + \bar{X}_2 + \bar{X}_3}{3}$ ②

mean of (mean of Sample)

Sampling Distribution Statistic

* **Standard Error**

① Quantifies Variability among multiple Sample

② How much Sample Mean is expected to Deviate from Population Mean

① If population SD ($\sigma$) is given

$$SE \oplus \frac{\sigma}{\sqrt{n}}$$

n : Sample Size

② If Sample SD (s) is given

$$SE \oplus \frac{s}{\sqrt{n}}$$

(T) 0

n : Sample Size

Key Take-away

$$SE \propto \frac{1}{\sqrt{n}}$$

# * Law of Large Numbers

① As Sample Size increases, the Sample mean gets closer to Population Mean

Diff b/w

SE                    and                    S

Standard Error (Sampling Distribution)

| Sample 1 | Sample 2 | Sample 3 |
|----------|----------|----------|

$S_1$                    $S_2$                    $S_3$

## Quiz

① n ⟹ 30

② $\bar{X}$ ⟹ 4

③ SD ⟹ 0.13 (σ)

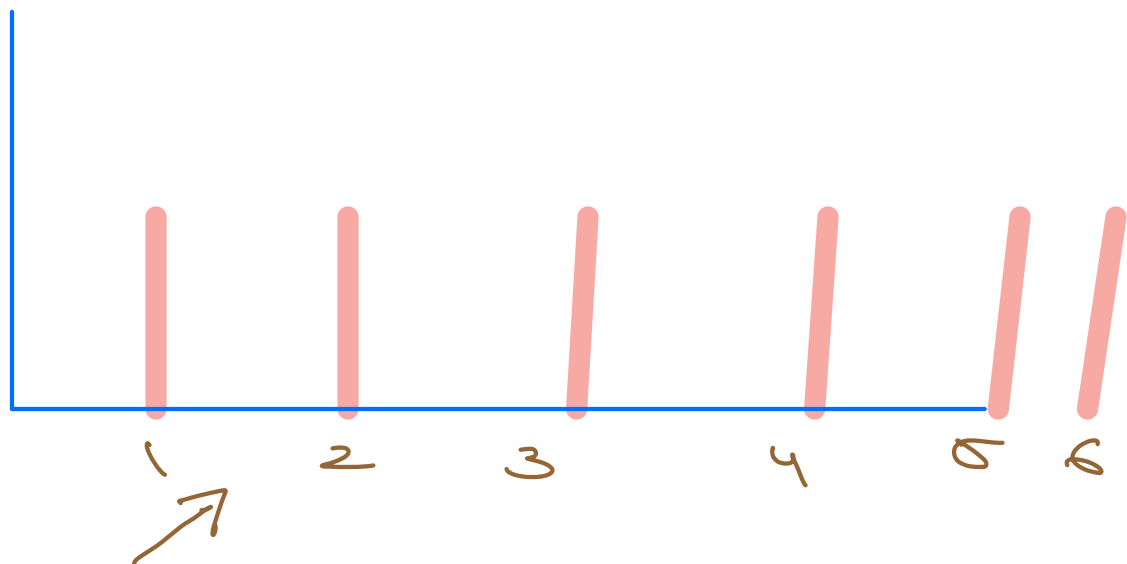② SE ⟹ $\frac{\sqrt{6}}{}$   ⟹ $\frac{\sigma}{\sqrt{n}}$

$$\boxed{\frac{0.13}{\sqrt{30}}}$$

# Uniform Distribution

⑤ Probability of all Outcomes is Equal

Ex→ Die roll → $\{1, 2, 3, 4, 5, 6\}$
↓
Uniform Discrete Distribution



$$\boxed{PMF ⑤ \frac{1}{b-a+1}}$$ for $a \le x \le b$
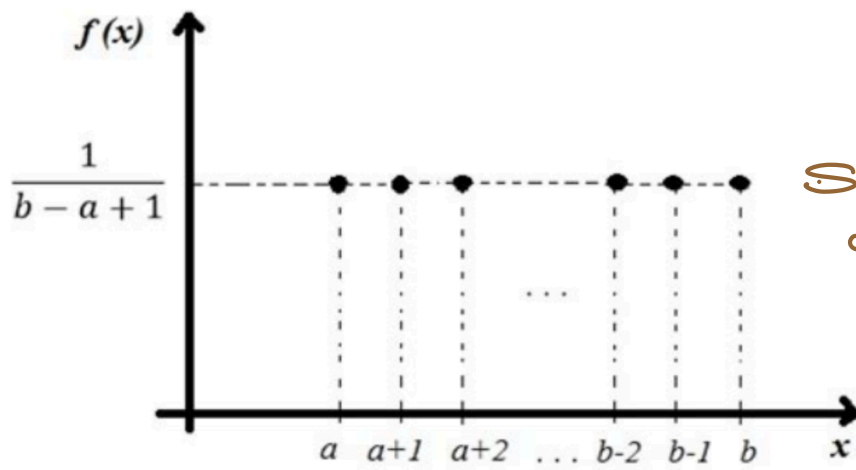↓ ↓
min val max. val

Dice Roll : $\boxed{\frac{1}{6-1+1}} ⇒ \boxed{\frac{1}{6}}$

Coin Toss : $\frac{1}{1-0+1} ⇒ \frac{1}{2} ⑤ 0.5$

$$f(x)$$

$$\frac{1}{b-a+1}$$

a   a+1   a+2   ...   b-2   b-1   b

Same prob across each Outcome

$$a = 1$$
$$b = 6$$ } Same as Dice

$$a = 0$$
$$b = 1$$ } Same as Coin Toss

② $$0 - 10 \leftarrow$$ Random Num Generator
↓
floating

$$PDF = \boxed{\frac{1}{b-a}} \quad (a \leq x \leq b) \quad 9.9$$
$$9.8$$