

* Population vs Sample

* Sample Statistics

Mean

Variance

SD

* Point Estimates

* Sampling Distribution

* Standard Error

* Uniform Distribution

PMF
PDF
CDF

Population vs Sample

Goal: Find average income of
ALL the Professional in
BLR

- * ideally: we would like to collect all the data possible
- * Not feasible
- * Very Costly



Most Practical Scenario

Sample



Calculate statistics

Avg Salary with some error margin

Sample Statistics

Population Stats (m)

$$\bar{X} = \frac{\sum x_i}{n}$$

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

$$\sigma = \sqrt{\sigma^2}$$

Sample Stats (n)

$$\bar{x} = \frac{\sum x_i}{n}$$

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

Sample Variance

Bernie's Correction
(Why?)

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

Why do we need Sample Statistics

Analyze population



Sample Statistics

estimate

Population Parameters

- There are 45 students in a class.
 5 students were randomly selected from this class and their heights (in cm) were recorded as follows:
 [131, 150, 140, 142, 152] Calculate Sample mean and sample variance

143

$$\frac{(131-143)^2 + (150-143)^2 + \dots}{n-1}$$

Estimate

Point Estimate
 ↓
 90 Lakh

Range Estimate
 ↓
 18 - 22 Lakh
 95%

Sampling Technique

① Probabilistic Sampling

Every single member of population will have a probability assigned for it to be part of Sample

⑤ Non-probabilistic Sampling X

↳ Generally based on convenience and hence cheaper.

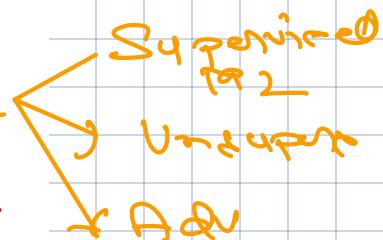
↳ Convenience Sampling

↳ Snowball Sampling

There are four main types of probability samples.

- ✓ 1. Simple random sampling (Today)
- 2. Systematic sampling
- * 3. Stratified sampling
- 4. Cluster sampling

} ML Models



10,000 MLs

Population Size

R

$\frac{1}{10,000}$

Sampling : Margin of Error

Case - 1 : 0% Margin of Error $\rightarrow n = 10000$

Case - 2 : 5% Margin of Error $\rightarrow n < 10000$

Case - 3 : 10% Margin of Error $\rightarrow n_2 < n_1$

- * In Random Sampling all points will have equal probability of being picked

10000 people

Steps for Random Sampling

- ① Define Population and sample-size
- ② Pick Sampling technique and perform Sampling
- ③ Collect Sample Data
 - ML Model --
 - Statistical Analysis --

Multiple Samples

Sample 1
 \bar{x}_1
 s_1^2

Sample 2
 \bar{x}_2
 s_2^2

Sample 3
 \bar{x}_3
 s_3^2

\vdots
 \bar{x}_k
 s_k^2

Sampling Distribution (CLT)
 \bar{x}_i

(Sample of Sample)

- * Estimate Population Parameters with Higher Reliability (CLT)

• (1) $s_1 \rightarrow 22.8 \quad \bar{x}_1$

(1) $s_2 \rightarrow 23.8 \quad \bar{x}_2$

(1) $s_3 \rightarrow 21.2 \quad \bar{x}_3$

CLT

$\left[\bar{x}_1, \bar{x}_2, \bar{x}_3, \dots, \bar{x}_k \right]$

Sampling Distribution

Standard Error

To quantify Variability across Sample

How far population mean is from Sampling Distribution mean

① Population SD is known

$$SE = \frac{S}{\sqrt{n}}$$

$S \rightarrow \text{Population SD}$
 $n \rightarrow \text{Sample size}$

② Population SD Not known

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

$s \rightarrow \text{Sample SD}$
 $n \rightarrow \text{Sample size}$

* If SE is large : Pop-mean is far away from Sample-mean

* If SE is Low : Pop-mean is Close from Sample-mean

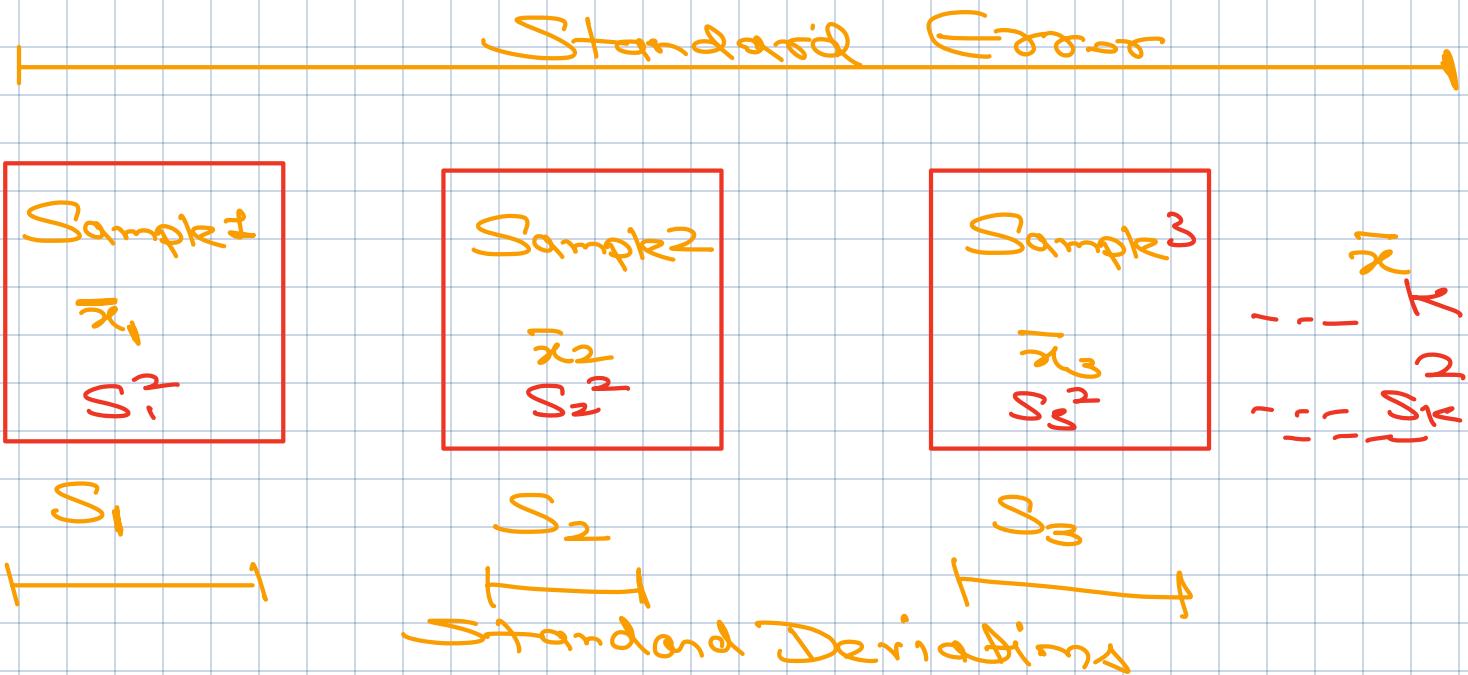
Key Takeaway

Since $SE \propto \frac{1}{\sqrt{\text{Sample size}}}$

Higher Sample Size, Lesser the Standard Error

* Law of Large Numbers

⑤ Higher the Sample Size, the closer Sample mean gets to Population Mean.



Collection of Sample means \rightarrow Sampling Distribution

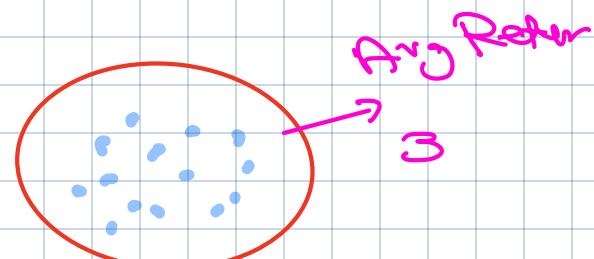
A sample of 30 latest returns on XYZ stock reveals a mean return of 4 with a sample standard deviation of 0.13. Estimate the SE of the sample mean.

$$n = 30$$

$$\bar{x} = 4$$

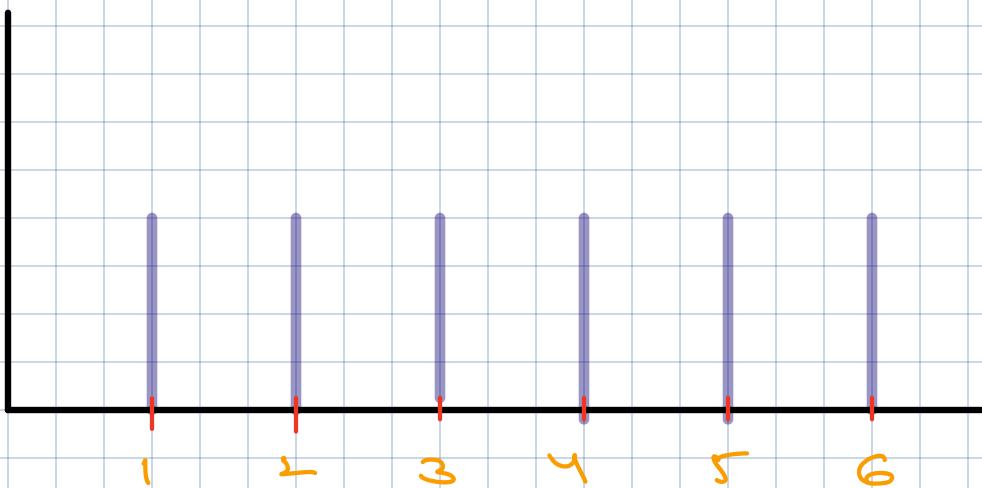
$$s_s = 0.13$$

$$SE \rightarrow \frac{0.13}{\sqrt{30}}$$



Uniform Distribution

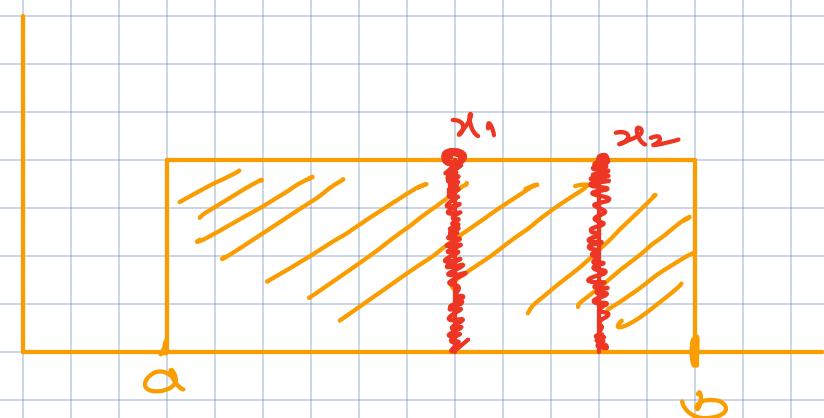
A distribution where all points are equally likely to occur



① Discrete Uniform Distribution

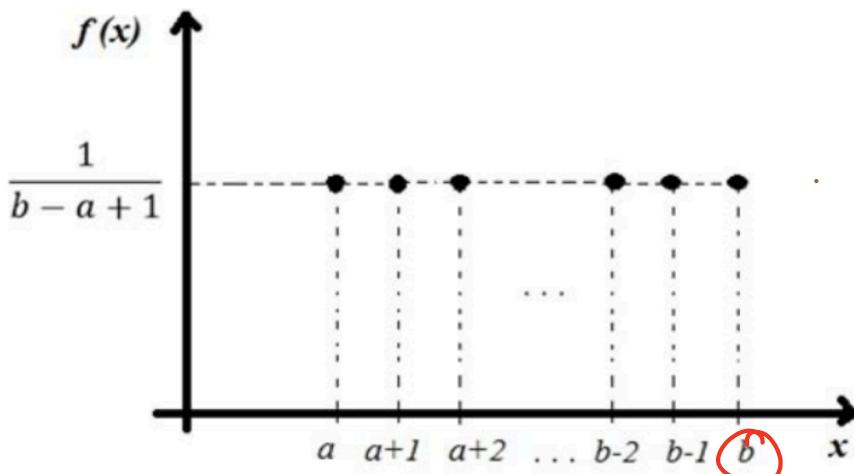


② Continuous Uniform Distribution



PMF θ

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



* PDF θ

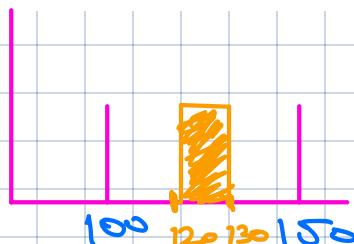
$$\frac{1}{b-a}$$



Suppose the weight of dolphins is uniformly distributed between 100 pounds and 150 pounds. If we randomly select a dolphin at random, then determine the probability that the chosen dolphin will weigh between 120 and 130 pounds.

$$130 - 120 = 10$$

4 options



$$\left(\frac{1}{b-a}\right)$$

$$\frac{1}{50} \times 10$$

Check if random.uniform actually generates data from UD

```
import numpy as np
import matplotlib.pyplot as plt

# Generate random data following a uniform distribution
data = np.random.uniform(1, 7, 10000) # 10,000 samples from a uniform distribution between 1 and 6

# Create a histogram
plt.hist(data, bins=6, edgecolor='black', align='mid', rwidth=0.8, density=True, color='skyblue')

# Set labels and title
plt.xlabel('Outcome')
plt.ylabel('Probability')
plt.title('Histogram of a Uniform Distribution (Die Roll Simulation)')

# Show the plot
plt.grid(True)
plt.show()
```

PDF

=

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

$$\frac{1}{150-100}$$

$$\frac{1}{50}$$

