



STATISTICS FOR DATA SCIENCE

Central Limit Theorem

D. Uma

Department of Computer Science and Engineering

STATISTICS FOR DATA SCIENCE

Central Limit Theorem

D. Uma

STATISTICS FOR DATA SCIENCE

Topics to be covered...

- ✓ Statistical Inference
- ✓ Sampling Distributions
- ✓ Central Limit Theorem



Statistical Inference



- Statistik



Uniform

$X :$	1	2	3	4	5	6
$P(X=x) :$	$\frac{1}{6}$	$\frac{1}{4}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$

Probability Distribution



Let x_1, x_2, \dots, x_n be n independent r.v.s from a
normal distribution with mean μ and variance σ^2 .
normally distribution

$$\begin{aligned} \mu_{x_1} = \mu & \quad x_1 \sim N(\mu, \sigma^2) \\ \sigma_{x_1}^2 = \sigma^2 & \\ \mu_{x_2} = \mu & \quad x_2 \sim N(\mu, \sigma^2) \\ \sigma_{x_2}^2 = \sigma^2 & \\ \mu_{x_3} = \mu & \quad x_3 \sim N(\mu, \sigma^2) \\ \sigma_{x_3}^2 = \sigma^2 & \end{aligned}$$

Population

Normally
Distributed
 μ, σ

$$\bar{X} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{1}{n}x_1 + \frac{1}{n}x_2 + \dots + \frac{1}{n}x_n$$

Sample
 x_1, \dots, x_n

$n > 30$

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

$$\sigma_{\bar{X}}^2 = \sigma^2 \left(\frac{1}{n}x_1 + \frac{1}{n}x_2 + \dots + \frac{1}{n}x_n \right)^2$$

$$\mu_{\bar{X}} = \mu \left(\frac{1}{n}x_1 + \frac{1}{n}x_2 + \dots + \frac{1}{n}x_n \right)$$

$n > 30$

$$= \frac{1}{n}\mu_{x_1} + \frac{1}{n}\mu_{x_2} + \dots + \frac{1}{n}\mu_{x_n}$$

$$= \frac{1}{n}\mu + \frac{1}{n}\mu + \dots + \frac{1}{n}\mu = \frac{1}{n}(\mu + \mu + \dots + \mu) = \frac{1}{n}(n\mu) = \mu$$

$$\begin{aligned} &= \frac{1}{n^2} \sigma_{x_1}^2 + \frac{1}{n^2} \sigma_{x_2}^2 + \dots + \frac{1}{n^2} \sigma_{x_n}^2 \\ &= \frac{1}{n^2} \left[\sigma^2 + \sigma^2 + \dots + \sigma^2 \right] \\ &= \frac{1}{n^2} (n \sigma^2) = \frac{\sigma^2}{n} \end{aligned}$$

population values = {1,3,5,7} on slips of paper and put them in a box.

List all possible samples of size $n = 2$ and calculate the mean of each.

Find the mean, variance, and standard deviation of the sample means.

Compare your results with the mean ($\mu = 4$) variance ($\sigma^2 = 5$) and standard deviation ($\sigma = 2.236$) of the population.

Predict the sampling distributions.

STATISTICS FOR DATA SCIENCE

Example – Sampling Distribution

List of **all 16 samples** of **size 2** from the **population {1,3,5,7}** and the mean of each sample.

$$M = 4 ; \sigma^2 = 5 ; \sigma = 2.236$$

Sample	Sample Mean (\bar{x})
1, 1	1 ✓
1, 3	2 ✓
1, 5	3 ✓
1, 7	4
3, 1	2 ✓
3, 3	3 ✓
3, 5	4
3, 7	5

Sample	Sample Mean (\bar{x})
5, 1	3 ✓
5, 3	4
5, 5	5
5, 7	6
7, 1	4
7, 3	5
7, 5	6
7, 7	7



PES
UNIVERSITY
ONLINE

Population $N=4$

1, 3,
5, 7

Sample
 $n=2$
?

\bar{x} :
 $P(\bar{x}=x)$:
sample

Population
 \bar{x} ✓ M
Prob
Distribution
? ?
static

Sampling Distribution

Probability Distribution
? a static
 \bar{x} or μ or σ

STATISTICS FOR DATA SCIENCE

Example – Sampling Distribution

List of **all 16 samples** of **size 2** from the **population {1,3,5,7}** and the mean of each sample.

Sample	Sample Mean (\bar{x})
1,1	1
1,3	2
1,5	3
1,7	4
3,1	2
3,3	3
3,5	4
3,7	4

Sample	Sample Mean (\bar{x})
5,1	3
5,3	4
5,5	5
5,7	6
7,1	4
7,3	5
7,5	6
7,7	7

Probability Distribution of all sample means

\bar{x}	frequency	Probability
1	1	$\frac{1}{16}$
2	2	$\frac{2}{16}$
3	3	$\frac{3}{16}$
4		
5		
6		
7		

16

Probability Distribution of all sample means

\bar{x}	frequency	Probability
1	1	1/16
2	2	2/16
3	3	3/16
4	4	4/16
5	3	3/16
6	2	2/16
7	1	1/16

$\Sigma f = 16$

Probability Distribution of all sample means &
Probability Histogram.

\bar{x}	frequency	Probability
1	1	1/16
2	2	2/16
3	3	3/16
4	4	4/16
5	3	3/16
6	2	2/16
7	1	1/16

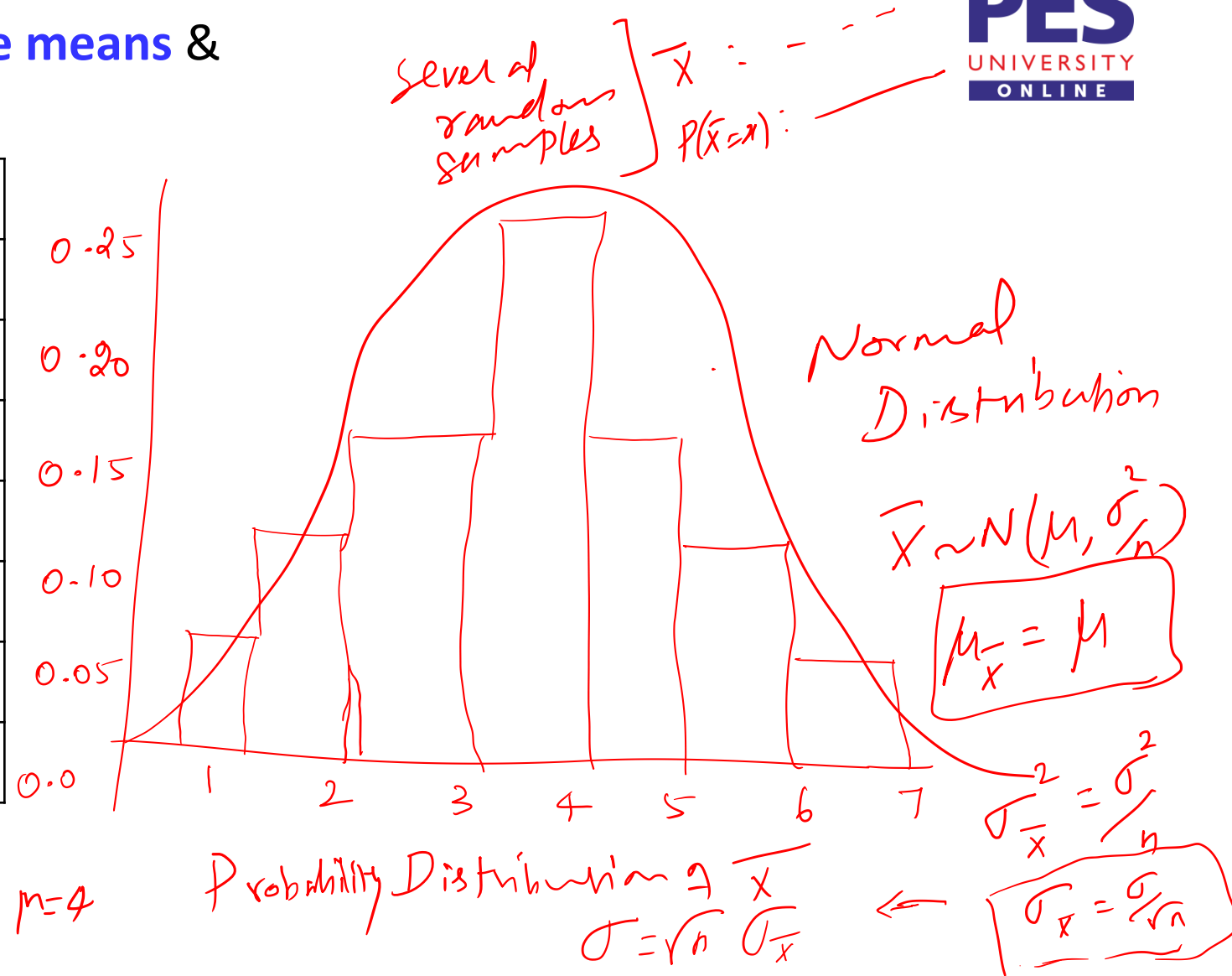
sample mean \bar{x}





Probability Distribution of all sample means & Probability Histogram.

\bar{x}	frequency	Probability
1	1	1/16
2	2	2/16
3	3	3/16
4	4	4/16
5	3	3/16
6	2	2/16
7	1	1/16



$$\mu_{\bar{x}} = 4 \text{ \& } \sigma_{\bar{x}} =$$

$$\mu = 4$$

$$P = \{1, 3, 5, 7\} \Rightarrow n = 4$$

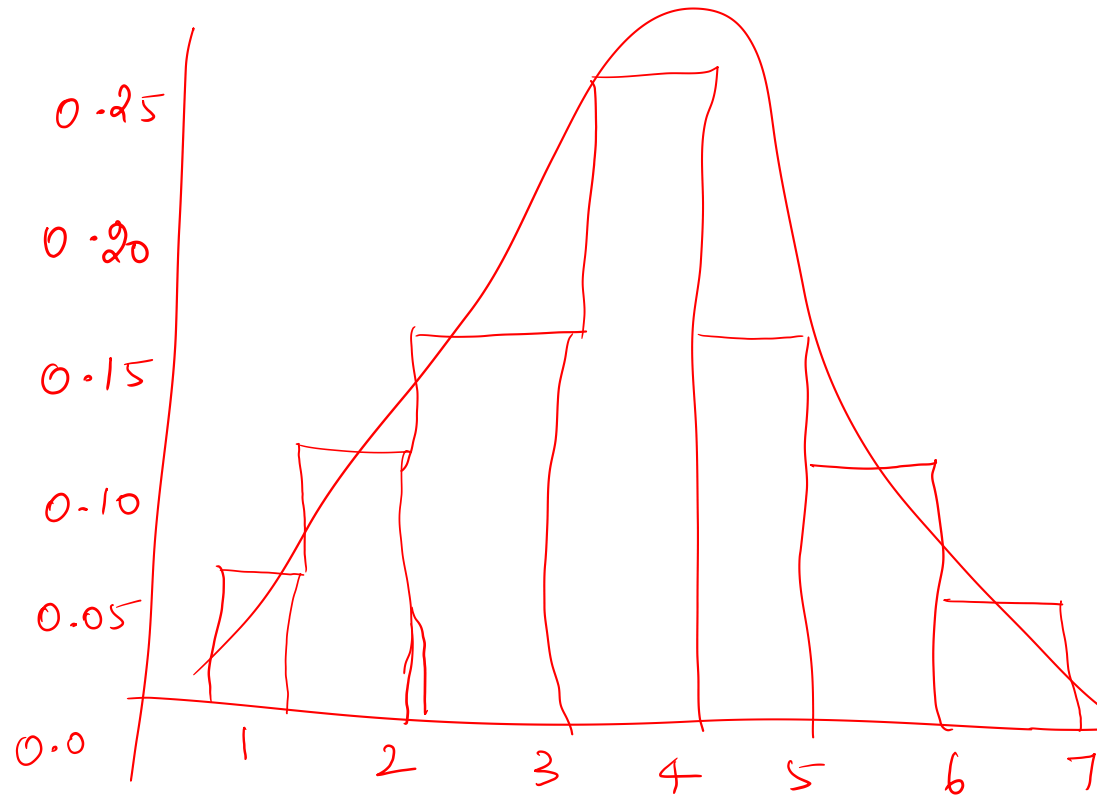
$$\text{Probability Distribution of } \bar{x}$$

$$\sigma = \sqrt{n} \sigma_{\bar{x}}$$

Example – Sampling Distribution & Probability Histogram

Probability Distribution of all sample means & Probability Histogram.

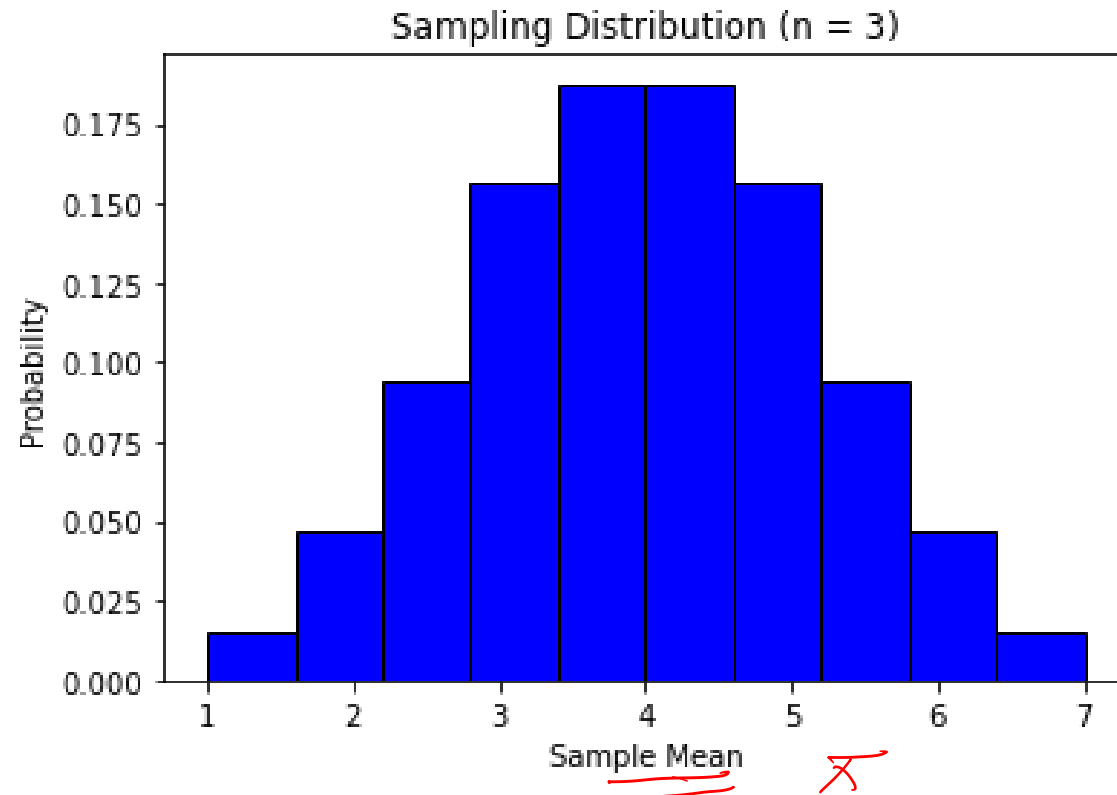
\bar{x}	frequency	Probability
1	1	$1/16 = 0.0625$
2	2	$2/16 = 0.1250$
3	3	$3/16 = 0.1875$
4	4	$4/16 = 0.2500$
5	3	$3/16 = 0.1875$
6	2	$2/16 = 0.1250$
7	1	$1/16 = 0.0625$



Example – Sampling Distribution



Sampling distribution when $n = 3$.



$N = 4$
 $n = 3$
 $\leftarrow n = 4$
 $n \uparrow$
 $n = 5$
 $n = 10$
 $n = 20$
 $n = 50$

Pop
Dish

Sampling distribution
of statistic
 $\sim N(\mu, \sigma)$

It is understood that when the sample size (n) increases, the shape is getting closer and closer to the normal distribution.

Mean and Standard Deviation of the Sampling Distribution

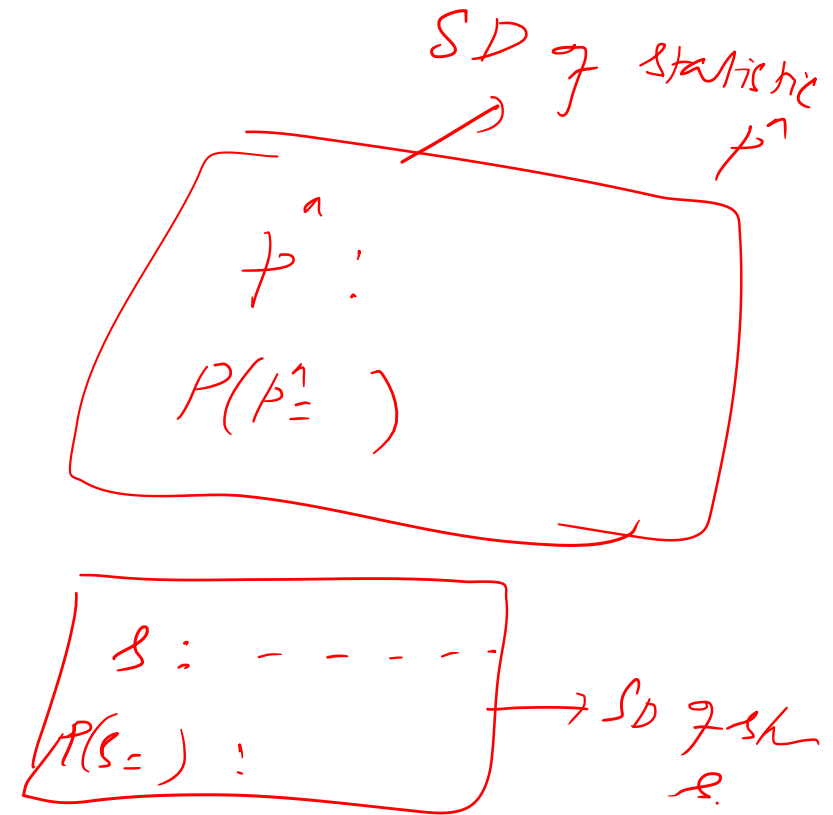
$$\text{If } \bar{X} \sim N(\mu, \sigma^2/n)$$

$$\mu_{\bar{X}} = \mu$$

$$\sigma_{\bar{X}}^2 = \frac{\sigma^2}{n}$$

$$\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$$

- The **probability distribution of statistic** is called as **sampling distribution**.
- Trials are repeated by taking sample size 'n' from a population.
- The distribution is sampling distribution of sample means.
- Every sample statistic has a sampling distribution.

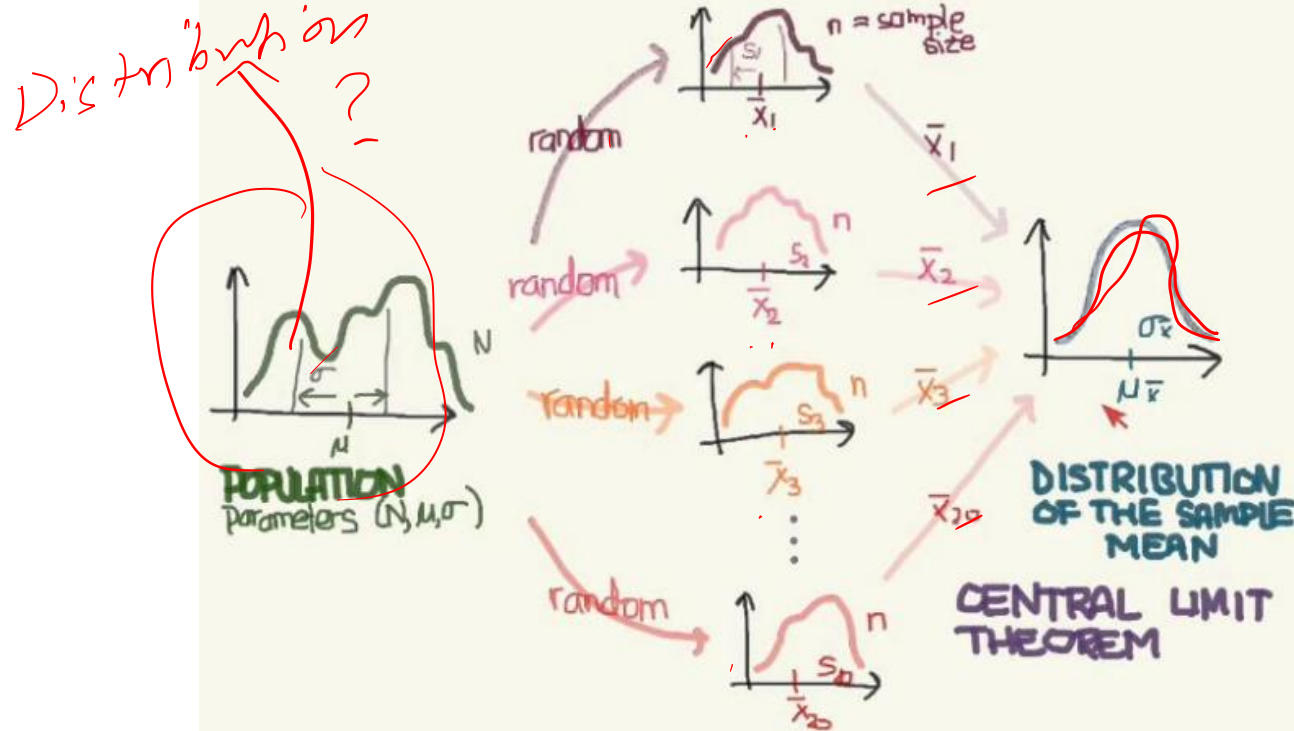


$n > 30$

- “Law of Large numbers” - The **mean of the sample distribution** will be **same as** the **mean of the population** distribution when the size of the **sample increases**.
- Random Selection of samples and independent of each other.
- **Sample size of 30** is mandatory.
- When the sampling is done without replacement, the sample size should **not be more than 10% of the population**.

What is Central Limit Theorem?

Central Limit Theorem states that the distribution of sample means that is calculated from sampling will follow normal distribution as the size of 'n' increases regardless of the samples that may be drawn from any population distribution.



The Central Limit Theorem says that if we draw a large enough sample from a population, then the distribution of the sample mean is approximately normal, no matter what population the sample was drawn.

Let X_1, X_2, \dots, X_n be a random sample from a population with mean μ and variance σ^2 .

Let $\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$ be the sample mean.

Let $S_n = X_1 + X_2 + \dots + X_n$ be the sum of the sample observations.

Then, if n is sufficiently large

✓ $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$ approximately.

$S_n \sim N(\underline{n\mu}, \underline{n\sigma^2})$ approximately.

$$\bar{X} \leftarrow \frac{X_1 + \dots + X_n}{n}$$

$$\bar{X} \sim N(\underline{\underline{\mu}}, \underline{\underline{\sigma^2/n}})$$

$$S_n = n \bar{X}$$

$$\mu_{S_n} = \mu_{n\bar{X}} = n \mu_{\bar{X}} = n(\mu) = n\mu$$

$$\sigma_{S_n}^2 = \sigma_{n\bar{X}}^2 = n^2 \sigma_{\bar{X}}^2 = n^2 \left(\frac{\sigma^2}{n} \right) = n\sigma^2$$

$$\therefore S_n \sim N(n\mu, n\sigma^2)$$

A business client of FedEx wants to deliver urgently a large freight from Denver to Salt Lake City.

When asked about the weight of the cargo they could not supply the exact weight, however they have specified that there are total of 36 boxes.

You are working as a **Business analyst** for FedEx.

And you have been challenged to tell the executives quickly whether or not they can do certain delivery.

Given : Mean of $\mu = 32.66$ kg

standard deviation of $\sigma = 1.36$ kg.

The plane you have can carry the max cargo weight up to 1193 kg.

Based on this information **what is the probability that all of the cargo can be safely loaded onto the planes and transported?**

$$n = 36$$

$$n > 30$$

CLT

36 boxes

$$\frac{1193}{36} = \underline{\underline{33.14}}$$

$P($

$$\begin{aligned}\bar{X} &\sim N(\mu, \sigma^2/n) \\ \bar{X} &\sim N(32.66, 0.227^2)\end{aligned}$$

$$\begin{aligned}\mu_{\bar{X}} &= \mu \\ \sigma_{\bar{X}}^2 &= \frac{\sigma^2}{n} = \frac{(1.36)^2}{36} \\ \sigma_{\bar{X}} &= 0.227\end{aligned}$$



$$\bar{X} = \frac{1193}{36} = 33.14$$

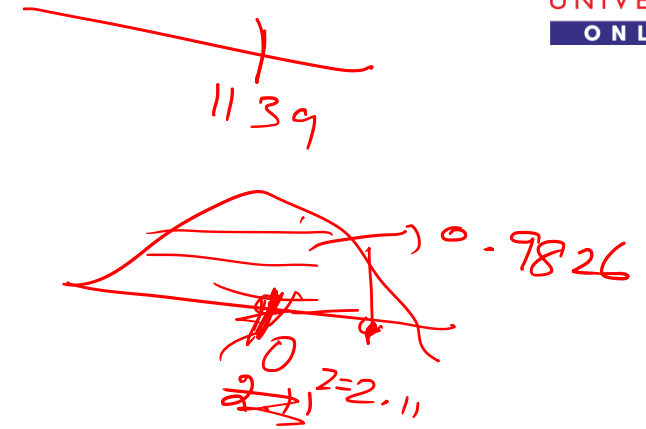
$$Z = \frac{\bar{X} - \mu}{s} = \frac{33.14 - 33.66}{0.227} = \underline{\underline{2.11}}$$

$$P(Z < 2.11) = 0.9826$$

$P(T < 1139)$

can take 98.26 %

1.7 % chance that
the plane can't take
it





We know that $\bar{X} \sim N(\mu_{\bar{X}}, \sigma_{\bar{X}}^2)$ where $\mu_{\bar{X}} = \mu = 32.66 \text{ kg}$

$$\therefore \bar{X} \sim N(32.66, 0.227^2) \quad \sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}} = \frac{1.36}{\sqrt{36}} = 0.227$$

Allowable capacity for the plane is 1193 kg

$$X_{\text{critical point}} = \frac{1193}{36} = 33.14 \text{ kg/box}$$

The corresponding z-score is $z = \frac{x - \mu}{\sigma_{\bar{X}}} = \frac{33.14 - 32.66}{0.227} = 2.11$

Area to the left of $z = 2.11$ is 0.9826



\therefore Probability that plane can take-off safely is 98.26 % or \approx 98.3%.

\therefore Probability that plane can't take off safely is 1.7%.

Example 2

Drums labeled 30 L are filled with a solution from a large vat. The amount of solution put into each drum is random with mean 30.01 L and standard deviation 0.1 L.

- a) What is the probability that the total amount of solution contained in 50 drums is more than 1500 L?
- b) If the total amount of solution in the vat is 2401 L, what is the probability that 80 drums can be filled without running out?
- c) How much solution should the vat contain so that the probability is 0.9 that 80 drums can be filled without running out?



$$n = 30$$
$$\bar{x} = 30.01, \quad s = 0.1$$

Problem

- a) What is the probability that the total amount of solution contained in 50 drums is more than 1500 L?

T = Total amt. of soln in 50 drums

$$\mu_{\bar{x}} = 30.01 \text{ L}$$

$$s_{\bar{x}} = 0.1$$

$$T \text{ or } S = n \bar{x} \quad S \sim N(n\mu, n\sigma^2)$$

$$P(\$ > 1500) = ?$$

$$\mu_S = \mu_{n\bar{x}} = n \mu_{\bar{x}} = 50(30.01)$$

$$= 1500.5$$

$$Z = \frac{\bar{x} - \mu_S}{s_S} = \frac{1500 - 1500.5}{0.7071}$$

$$\sigma_S = \sqrt{n} s = \sqrt{50} * 0.1 = 0.7071$$

$$= -0.71$$

$$P(Z > -0.71) = 1 - 0.2389 = 0.7611$$

$$P(Z < -0.71) = 0.2389$$

$$S \sim N(1500.5, 0.7071)$$

$$\therefore P(S > 1500) = 0.7611$$

Problem

- b) If the total amount of solution in the vat is 2401 L, what is the probability that 80 drums can be filled without running out?

Let T : Total amount of solution in 80 drums.

$$T = X_1 + X_2 + \dots + X_{80} \quad \text{where } X_i: \text{Amt. of soln in drum } X_i$$

As per CLT, $T \sim N(\mu_T, \sigma_T^2)$ where $\mu_T = n\mu = 80(30.01) = 2400.8$

$$\sigma_T = \sqrt{n}\sigma = 0.1 * \sqrt{80} = 0.8944$$

To find: $P(T < 2401)$

$$= P\left(\frac{T - \mu}{\sigma_T} < \frac{2401 - 2400.8}{0.8944}\right)$$
$$= P(Z < 0.22)$$

$$P(T < 2401) = \underline{\underline{0.5871}}$$



Problem

- c) How much solution should the vat contain so that the probability is 0.9 that 80 drums can be filled without running out?

T : Total Amt. of solution in 80 drums.

We know that $T \sim N(\mu_T, \sigma_T^2) \Rightarrow T \sim N(\overset{n\mu}{80(30.01)}, (\overset{(s\sqrt{n})^2}{(0.1\sqrt{80})^2})$

$$T \sim N(2400.8, 0.8944^2)$$

Given probability is 0.9

The corresponding z-score is $Z=1.28$

$$\therefore Z = \frac{X - \mu_T}{\sigma_T} \Rightarrow 1.28 = \frac{X - 2400.8}{0.8944}$$

$$\Rightarrow X = 2401.9 \text{ l}$$

\therefore Vat should contain 2401.9 l solution.

Do It Yourself!!!

A simple random sample of 100 men is chosen from a population with mean height 70 in. & standard deviation 2.5 in.
What is the probability that the average height of the sample men is greater than 69.5 in?

The manufacture of a certain part requires two different machine operations.

- The time on machine 1 has mean 0.5 hours and standard deviation 0.4 hours.
- The time on machine 2 has mean 0.6 hours and standard deviation 0.5 hours.
- The times needed on the machines are independent.
- Suppose 100 parts are manufactured.

- 1) What is the probability that the total time used by machine 1 is greater than 55 hours?
- 2) What is the probability that the total time used by machine 2 is less than 55 hours?
- 3) What is the probability that the total time used by both the machines together is greater than 115 hours.
- 4) What is the probability that the total time used by machine 1 is greater than the total time used by machine 2?

The manufacture of a certain part requires two different machine operations.

- The time on machine 1 has mean 0.5 hours and standard deviation 0.4 hours.
- The time on machine 2 has mean 0.6 hours and standard deviation 0.5 hours.
- The times needed on the machines are independent.
- Suppose 100 parts are manufactured.

1) What is the probability that the total time used by machine 1 is greater than 55 hours?

The Central Limit Theorem: Problems



The manufacture of a certain part requires two different machine operations.

- The time on machine 1 has mean 0.5 hours and standard deviation 0.4 hours.
- The time on machine 2 has mean 0.6 hours and standard deviation 0.5 hours.
- The times needed on the machines are independent.
- Suppose 100 parts are manufactured.

2) What is the probability that the total time used by machine 2 is less than 55 hours?

The Central Limit Theorem: Problems



The manufacture of a certain part requires two different machine operations.

- The time on machine 1 has mean 0.5 hours and standard deviation 0.4 hours.
- The time on machine 2 has mean 0.6 hours and standard deviation 0.5 hours.
- The times needed on the machines are independent.
- Suppose 100 parts are manufactured.

3) What is the probability that the total time used by both the machines together is greater than 115 hours.

The Central Limit Theorem: Problems

The manufacture of a certain part requires two different machine operations.

- The time on machine 1 has mean 0.5 hours and standard deviation 0.4 hours.
- The time on machine 2 has mean 0.6 hours and standard deviation 0.5 hours.
- The times needed on the machines are independent.
- Suppose 100 parts are manufactured.

4) What is the probability that the total time used by machine 1 is greater than the total time used by machine 2?



THANK YOU

D. Uma

Computer Science and Engineering

umaprabha@pes.edu

+91 99 7251 5335