



STATISTICS FOR DATA SCIENCE

Central Limit Theorem

Prof. Uma D

Prof. Silviya Nancy J

Prof. Suganthi S

Department of Computer Science and Engineering

STATISTICS FOR DATA SCIENCE

Central Limit Theorem

Prof. Uma D
Prof. Silviya Nancy J
Prof. Suganthi S

STATISTICS FOR DATA SCIENCE

Topics to be covered...

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



- It allows us to make accurate decisions from the numerical descriptive measures of parameter from statistic.
- Accuracy can be measured using probability.
- It is entrusted by identifying the distribution of the random variables from the population from where they come from.
- With this the shape and the location of the sample mean, decision about the population parameters can be made.
- To understand the shape of the sample data, sampling distribution can be employed.

Examples of How do we deal with unknown parameters?

- A medical researcher surveys and the responses are based on agreement or disagreement which will follow Binomial distribution, but the proportion (p) of people who agree in the population may be unknown.
- The heights of the male in the city is assumed to be normally distributed. Estimate the mean (μ) and standard deviation (σ) which are unknown.
- So, to interpret the above two cases, and to get reliable information about the population, the sample should be chosen accordingly.
- This can be done by choosing the appropriate sampling methods.

- The probability distribution of statistic is called as sampling distribution.
- This distribution based on how many trials are repeatedly taken of size 'n' from a population.
- If the sample statistic is the sample mean, then the distribution is sampling distribution of sample means.
- Every sample statistic has a sampling distribution.

You write the population values {1,3,5,7} on slips of paper and put them in a box. Then you randomly choose two slips of paper, with replacement. List all possible samples of size $n = 2$ and calculate the mean of each. These means form the sampling distribution of the sample means. 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.

Example – Sampling Distribution

List of all 16 samples of size 2 from the population and the mean of each sample.

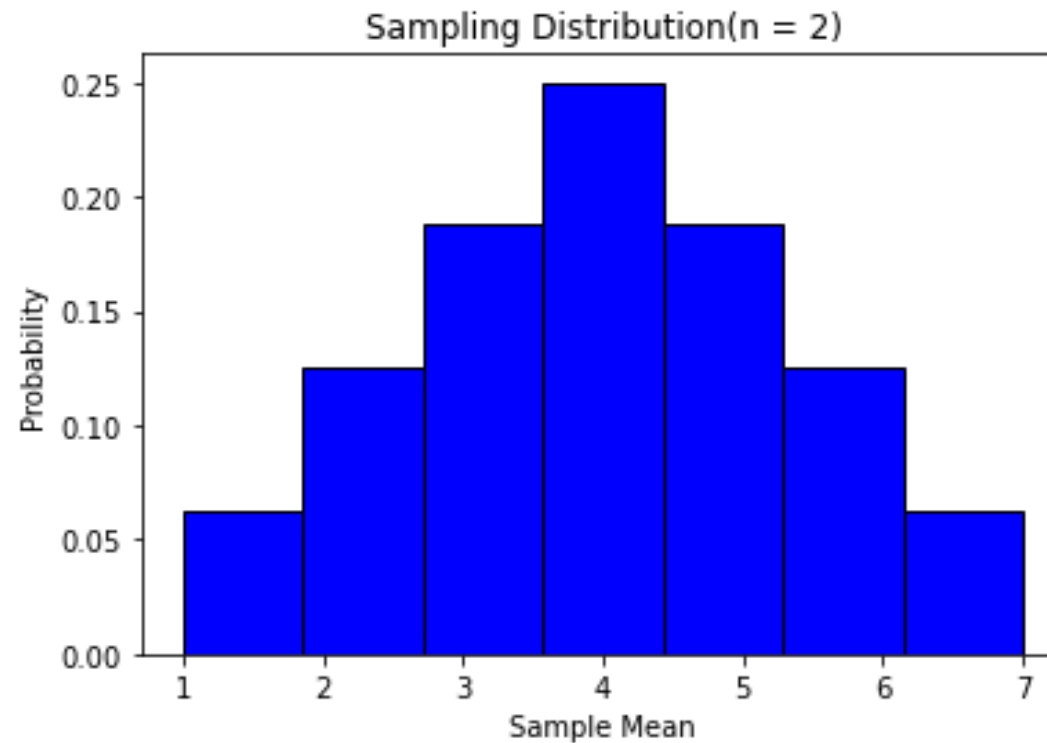
Sample	Sample Mean (\bar{x})
1,1	1
1,3	2
1,5	3
1,7	4
3,1	2
3,3	3
3,5	4
3,7	5

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

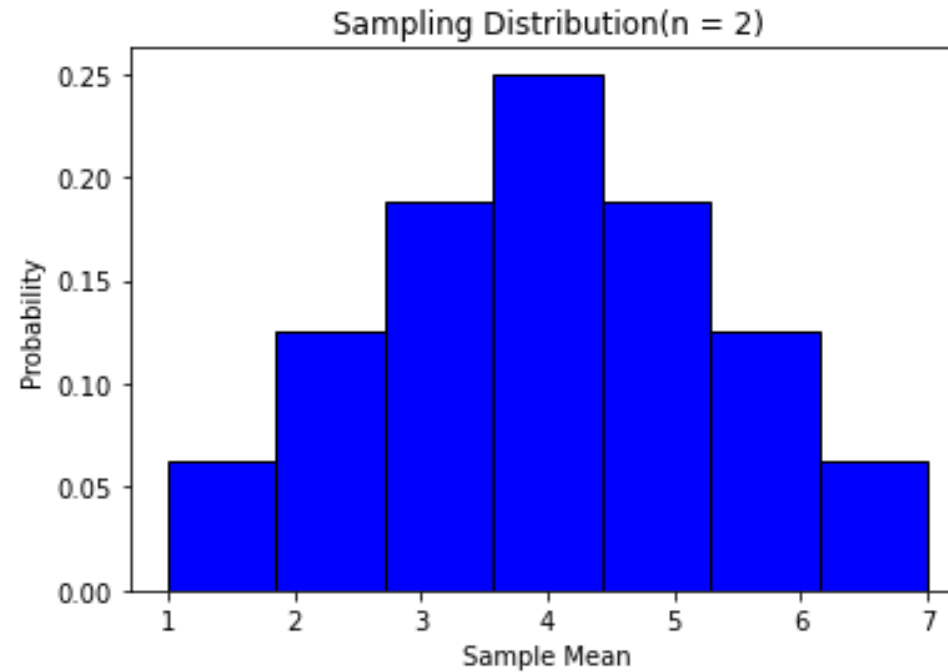
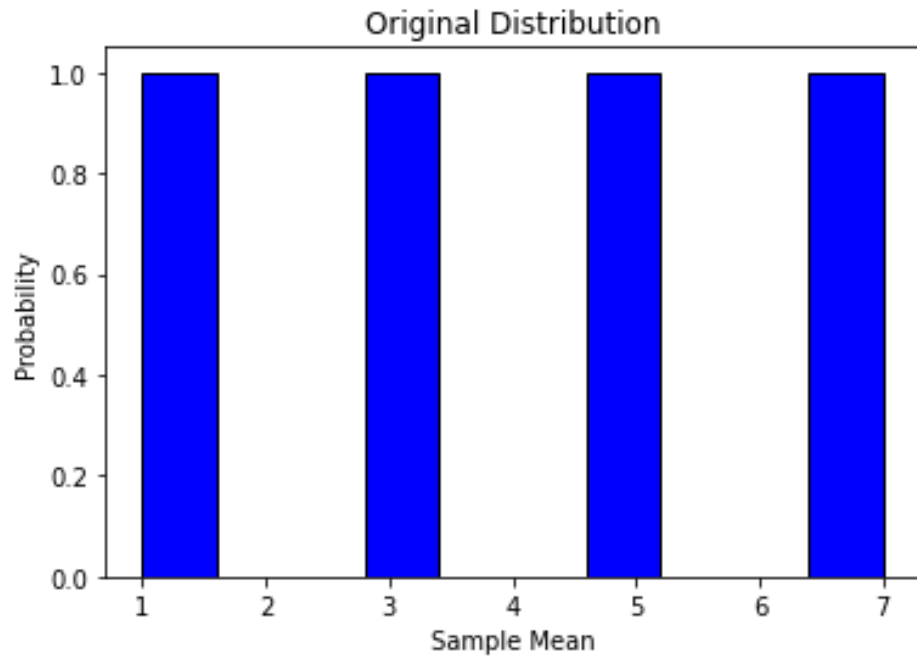
Example – Sampling Distribution & Probability Histogram

Probability Distribution of all sample means & Probability Histogram.

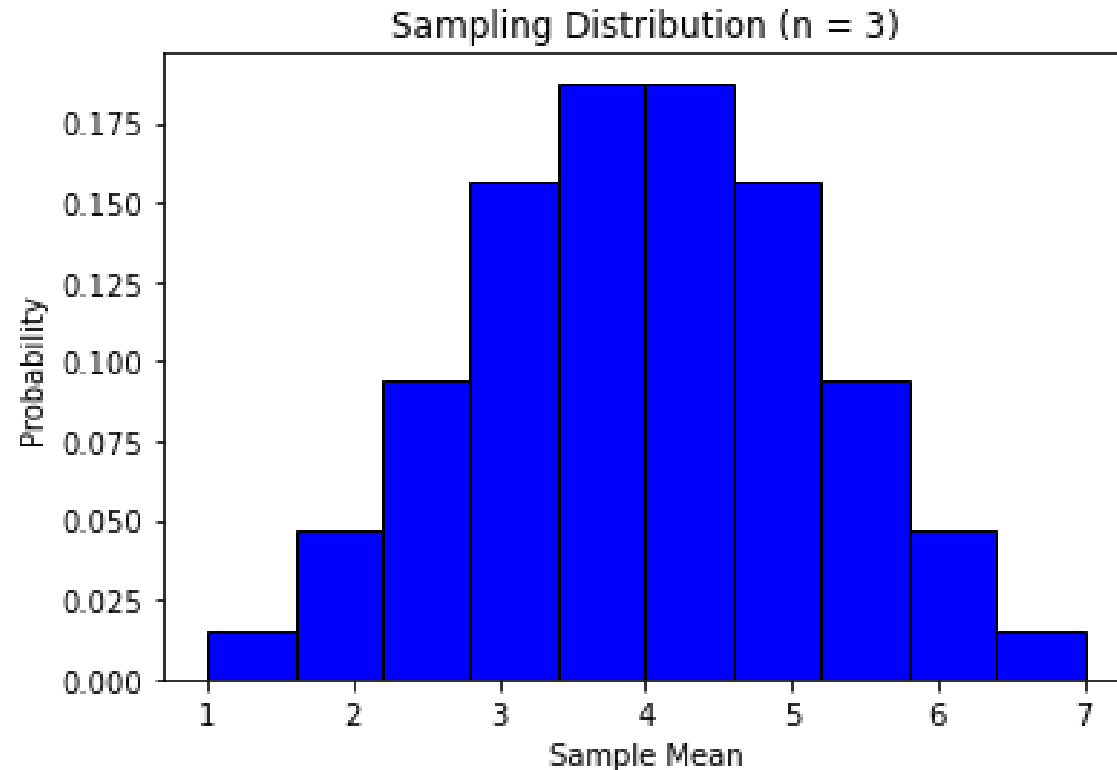
\bar{x}	frequency	Probability
1	1	$1/16 = 0.0625$
2	2	$2/16 = 0.1250$
3	3	$3/16 = 0.1875$
4	4	$4/16 = 0.2500$
5	3	$3/16 = 0.1875$
6	2	$2/16 = 0.1250$
7	1	$1/16 = 0.0625$



The original distribution and the sampling distribution when $n = 2$.



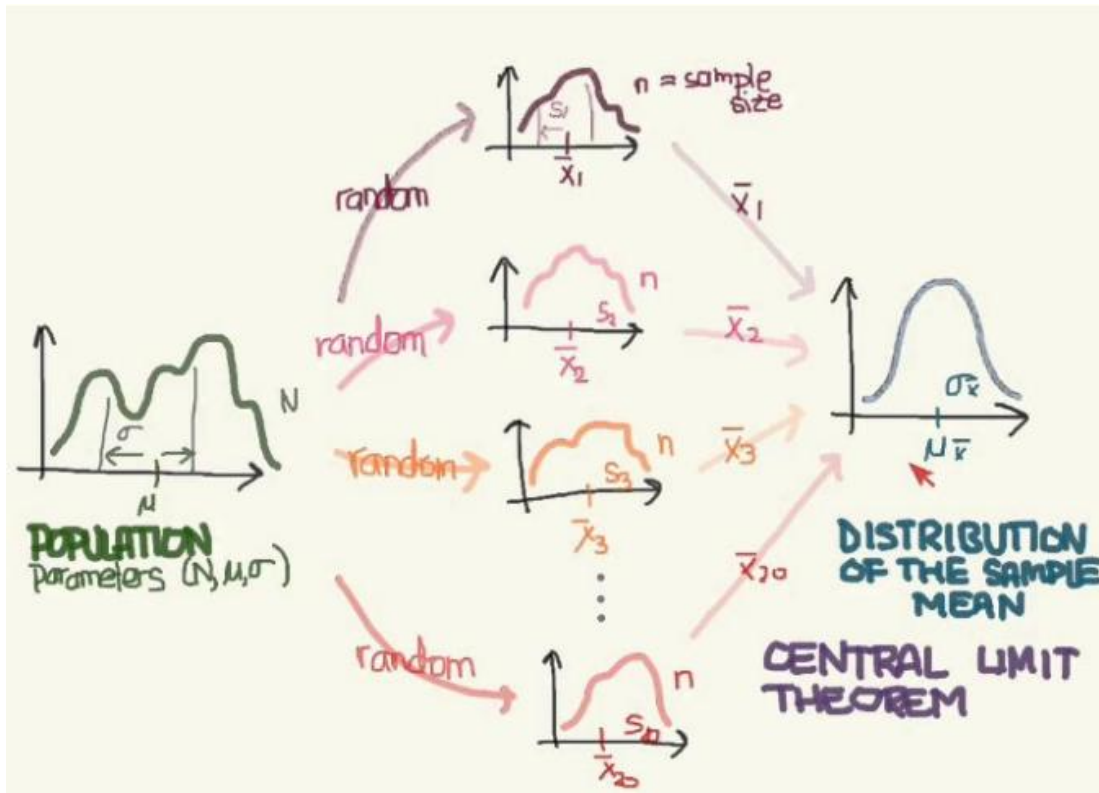
Sampling distribution when $n = 3$.



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

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.



Data in Sample Distribution should be

- As “Law of Large numbers” states that the mean of the sample distribution will be same as the mean of the population distribution when the size of the sample increases.
- The samples should be randomly selected.
- The samples must be independent of each other.
- The sample size should be large enough for the distribution to be normal, a sample size of 30 is mandatory for getting more representative sample.
- When the sampling is done without replacement, the sample size should not be more than 10% of the population.

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

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

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

Then if n is sufficiently large,

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

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

- The central limit theorem specifies that $\mu_{\bar{x}} = \mu$ and $\sigma_{\bar{x}}^2 = \frac{\sigma^2}{n}$ which hold for any sample mean.
- The sum of the sample items is equal to the mean multiplied by the sample size, that is $S_n = n\bar{X}$
- It follows that $\mu_{S_n} = n\mu$ and $\sigma_{S_n}^2 = \frac{n^2 \sigma^2}{n} = n\sigma^2$

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.

Since, we have worked with them for so many years and have seen so many freights from them we can confidently say that the type of cargo they follow is a distribution with a

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?

- Using Central Limit Theorem, find the mean and standard deviation of the sample mean.
- Calculate the critical point of each box by dividing the allowable capacity of the plane to carry weight with total number of boxes.
- So, to safely takeoff the plane, the average weight of the each box should not exceed 33.14kg/box.
- Finally, calculate the Z-score.

$$\mu_{\bar{x}} = \mu_x = 32.66 \quad \sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = \frac{1.36}{\sqrt{36}} = 0.227$$

Plane Capacity = 1193 kg

$$x_{critPoint} = \frac{1193}{36} = 33.14 \text{ kg/box}$$

$$Z = \frac{x_{critPoint} - \mu_{\bar{x}}}{\sigma_{\bar{x}}} = \frac{33.14 - 32.66}{0.227} = 2.11$$

$$P(x < x_{critPoint}) = 0.9826 = 98.26\%$$

So, The plan can safely takeoff is 98.26% and 1.7% chance it cannot takeoff.

Problem



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?

Solution for part (a)

Let $S = X_1, \dots, X_{50}$ be the amounts of solution in 50 drums.

$$\mu_x = 30.01 \text{ and } \sigma_x = 0.1$$

Assuming that S is approximately normally distributed, then

$$\mu_S = 50(30.01) = 1500.5 \text{ and } \sigma_S = 0.1\sqrt{50} = 0.7071$$

By calculating z - score,
$$z = \frac{1500 - 1500.5}{0.7071} = -0.71$$

The area to the right of $z = -0.71$ is $1 - 0.2389 = 0.7611$

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

Solution for part (b)

Let $D = X_1, \dots, X_{80}$ be the amounts of solution in 80 drums.

$$\mu_x = 30.01 \text{ and } \sigma_x = 0.1$$

Assuming that D is approximately normally distributed, then

$$\mu_D = 80(30.01) = 2400.8 \text{ and } \sigma_D = 0.1\sqrt{80} = 0.8944$$

By calculating z - score,
$$z = \frac{2401 - 2400.8}{0.8944} = 0.22$$

The area to the left of $z = 0.22$ is 0.5871

$$P(S < 2401) = 0.5871$$

The z – score for 90th percentile is $z = 1.28$

From part (b), D is approximately normally distributed with

$$\mu_D = 80(30.01) = 2400.8 \text{ and } \sigma_D = 0.1\sqrt{80} = 0.8944$$

$$1.28 = \frac{x - 2400.8}{0.8944}$$

$$x = 2401.9L$$



THANK YOU

Prof. Uma D

Prof. Silviya Nancy J

Prof. Suganthi S

Department of Computer Science and Engineering