

Agenda

① Central Limit Theorem ✓

② Z-score And Z-stats ✓

③ Z-test, t-test {Solve problems Assignment}

In inferential statistics we draw conclusion for population data using sample data.
Hypothesis testing is one of the means to do that.

Inferential stats

Hypothesis Testing.

If I have a X (which is population) random variable that belongs to normal/gaussian distriⁿ with some mean and std dev and if we create multiple samples out of this (S_1, S_2, S_3, \dots so on) in such a way that each of these sample have fixed sample size(n). Then we determine means for each of these samples and plot it then means of these samples will also follow normal/gaussian distribution. This is called Central limit theorem.

① Central Limit Theorem

$$① X \sim N(\mu, \sigma)$$

1, 2, 3, 4, 5



Note sample size can have any value $n = \text{sample size} \Rightarrow \underline{\text{any value}}$



$$S_1 = \{x_1, x_2, x_3, \dots, x_n\} = \bar{x}_1 //$$

$$S_2 = \{x_2, x_3, x_4, x_5, \dots, x_n\} = \bar{x}_2 //$$

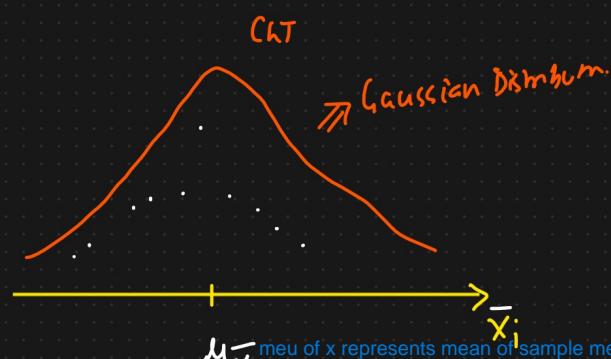
$$S_3 = \{ \dots \} = \bar{x}_3 // \quad \text{means of sample}$$

$$S_4 = \{ \dots \} = \bar{x}_4 //$$

$$S_5 = \{ \dots \} = \bar{x}_5 //$$

$$\vdots \quad \vdots$$

$$S_m = \{ \dots \} = \bar{x}_m //$$



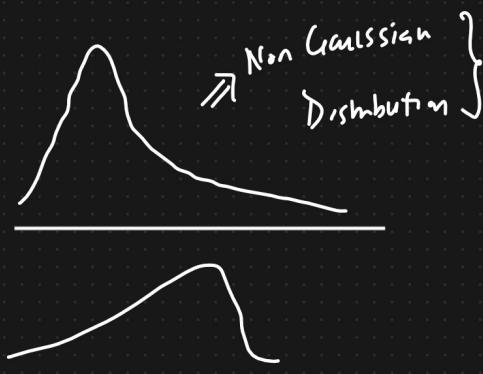
μ_x me of x represents mean of sample means OR 'a mean' for means of different samples calculated above.

$n = \text{Sample size}$

$$n \geq 30$$

In case of random variable X does not follow normal distribution then if plot sample means with sample size greater than or equal to 30 then distribution thus formed for means of sample will follow the normal distribution.

$$② X \not\sim N(\mu, \sigma)$$



$$S_1 = \{ \dots \} = \bar{x}_1 //$$

$$S_2 = \{ \dots \} = \bar{x}_2 //$$

$$S_3 = \{ \dots \} = \bar{x}_3 //$$

$$S_4 = \{ \dots \} = \bar{x}_4 //$$

$$\vdots \quad \vdots$$

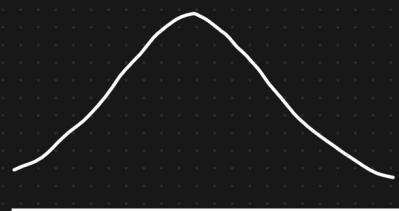
$$\bar{x}_m //$$



① Central Limit Theorem

The central limit theorem relies on the concept of a sampling distribution, which is the probability distribution of a statistic for a large number of samples taken from a population.

The central limit theorem says that the sampling distribution of the mean will always be normally distributed, as long as the sample size is large enough. Regardless of whether the population has a normal, Poisson, binomial, or any other distribution, the sampling distribution of the mean will be normal.



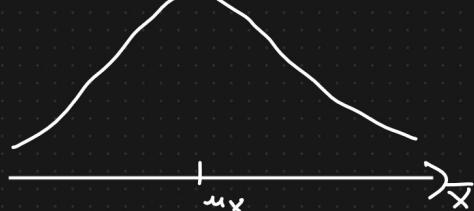
$$X \sim N(\mu, \sigma)$$

This represents population distribution

σ = population std

μ = population mean

Sampling Distribution of the mean



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

standard error.

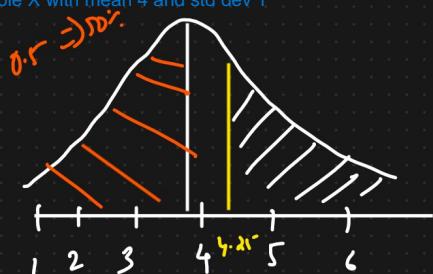
n = Sample size

Acc to CLT sample mean will be approximately equal to population mean but sample standard deviation will be equal to population standard deviation divided by root of sample size which is denoted as standard error.

Below we will discuss about Z score and Z stats

$$(1) \quad X \sim N(4, 1)$$

Eg: Calculating Zscore at point 4.25 for normally distributed random variable X with mean 4 and std dev 1



$$\text{Zscore} = (\text{pt. where Zscore needs to be determined} - \text{mean}) / \text{std dev}$$

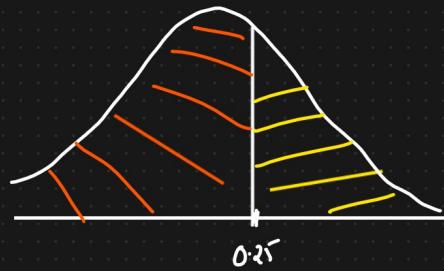
$$\lambda_i = 4.25$$

$$\text{Z-score} = \frac{4.25 - 4}{1} = 0.25$$

Z score is a pt. that basically tells that how much standard deviation away a pt. is from mean. This pt. is then used for calculating the area under the curve using Z table (Z table maps Zscore with area under the curve from left hand side). This curve is basically something whose internal area can be measured in terms of probability(PDF).

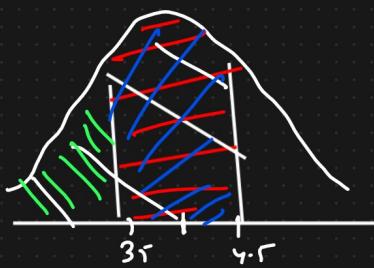
For Eg: Amazon manager asks to find the probability of any item to get delivered in 3 days. So here in this case we will plot histogram with kde=True where x axis represents number of days whereas y axis represents number of times that product get delivered on that particular day. After this we can simply calculate area under 4 th day using Z score and Z table and present to the Amazon manager.

Q) What percentage of score falls above 4.25?



$$1 - 0.59871 = 0.4013 = 40.13\%$$

Q) What percentage of score lies between 3.5 to 4.5?

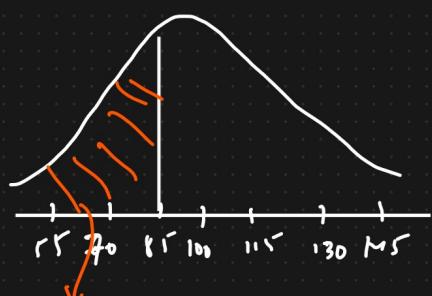


$$Z\text{-Score} = \frac{4.5 - 4}{1} = +0.5 = 0.69146$$

$$\begin{aligned} Z\text{-Score} &= \frac{3.5 - 4}{1} = -0.5 = 0.30854 \\ &= 0.3829 \\ &= 38.29\% \end{aligned}$$

Q) In India the average IQ is 100, with a standard deviation of 15. What is the percentage of the population would you expect to have an IQ lower than 85?

Ans). $\mu = 100$ $\sigma = 15$



$$0.1586 = 15.86\%$$

$$Z\text{-Score} = \frac{85 - 100}{15} = -1$$

Q) $IQ > 85$

$$1 - 0.1586 = 84.13\%$$

$| Z >, IQ \leq 100 | \Rightarrow \text{Internal Assignment}$

① Hypothesis Testing And Statistical Analysis (Statistical analysis wrt inferential statistics.)

① Z-test

② t-test

Below we will try to achieve hypothesis testing using all these(Ztest,Ttest, chi square , Anova)

③ Chi square

④ ANOVA

Note in numerical problems discussed below we are performing "1 sample Ztest and Ttest" since for these only 1 sample with some sample size is taken.

① Z-test

NOTE: Z test and T test always follows for Gaussian distribution.

①) The average heights of all residents in a city is 168cm. A doctor believes the mean to be different. He measured the height of 36 individuals and found the average height to be 169.5 cm.

(a) State null and Alternative Hypothesis

(b) At a 95% confidence level, is there enough evidence to reject the null hypothesis.

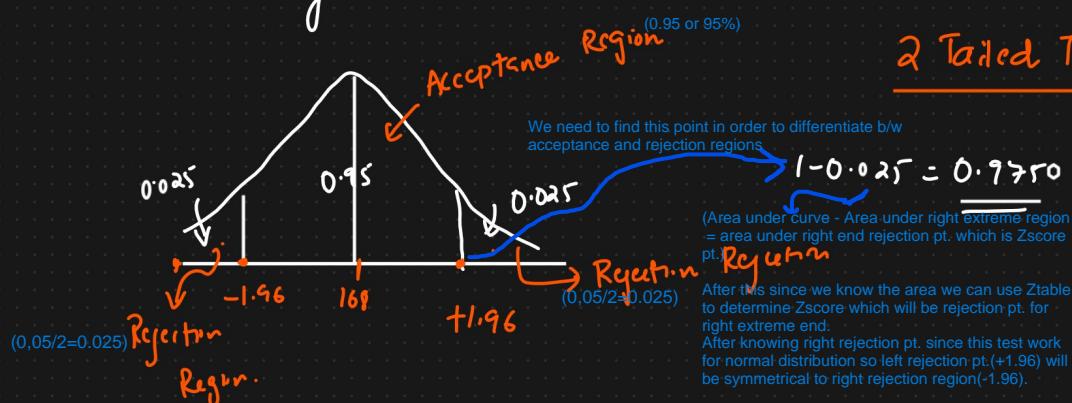
(In short in above question, using sample information we are trying to draw information for population using Z test. This is called inferential statistics).

Ans) $\mu = 168 \text{ cm}$ $\sigma = 3.9$ $n = 36$ $\bar{x} = 169.5$ $(I=0.95)$ $\alpha = 1 - I$
 Population mean Population dev sample size sample mean confidence interval
 $= 0.05\%$

① Null Hypothesis $H_0 : \mu = 168 \text{ cm}$ (Default info given in problem)

Alternate Hypothesis $H_1 : \mu \neq 168 \text{ cm}$

② Decision Boundary based on I .



Since, in question it is said that population mean can be different so it can be either greater than 168 cm or less than 168 cm. That's why this lie under the category of 2 tailed test.

2 Tailed Test

If, doctor would have said that he believes that population mean to be greater than 168 cm then this will cover only 1 extreme region and not both extreme regions as above(as in above case). Then this will fall under the category of 1 tailed test.

Similarly, if doctor says that he believes that population mean to be less than 168 cm then also this case will fall under 1 tailed test(since only 1 extreme region considered).

After this since we know the area we can use Ztable to determine Zscore which will be rejection pt. for right extreme end.
After knowing right rejection pt. since this test work for normal distribution so left rejection pt.(+1.96) will be symmetrical to right rejection region(-1.96).

If Z-test is less than -1.96 or greater than +1.96, Reject the Null Hypothesis

$$\textcircled{4} \quad Z_{\text{test}} = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{169.5 - 168}{\frac{3.9}{\sqrt{36}}} = \boxed{2.31}$$

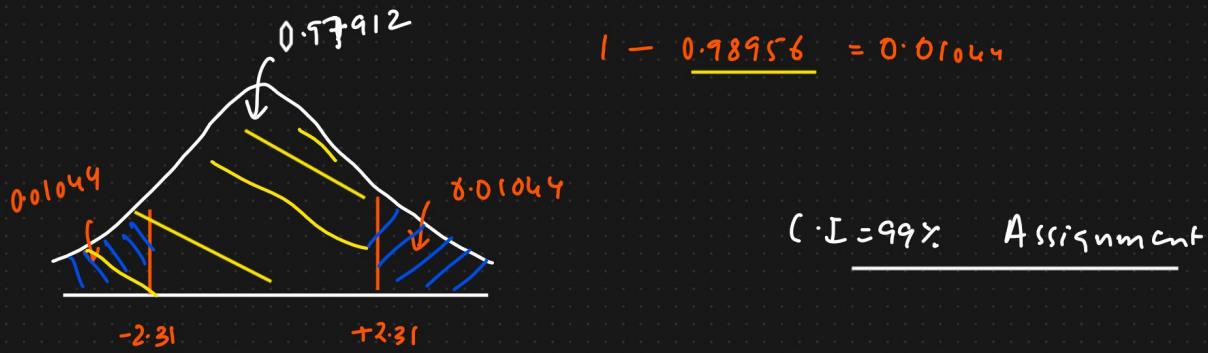
↓
 Sample mean Population mean
 Population std dev $\left\{ \frac{\sigma}{\sqrt{n}} \right\}$ root of sample size
 ↓ (Standard error)
 { CLT }

NOTE: Sample size meaning is number of elements in a sample and NOT Number of samples.

Conclusion

2.31 > 1.96 Reject the Null Hypothesis.

Below is another way of determining whether accept or reject the null hypothesis using p value which is basically Probability value.



$$\textcircled{1} \quad p \text{ value} = 0.01044 + 0.01044$$

(Probability of distribution lying in rejection region for the points calculated based Zscore determined using Ztest)

$$= 0.02088$$

$P < 0.05 \Rightarrow$ Reject the Null Hypothesis.

(Note here 0.05 is the default probability of rejection regions calculated based on CI and Significance level defined by domain expert)

Example (2) A factory manufactures bulbs with an average warranty of 5 years

with standard deviation of 0.50. A worker believes that the bulb will malfunction in less than 5 years. He tests a sample of 40 bulbs and finds the average time to be 4.8 years.

(a) State null and alternate hypothesis

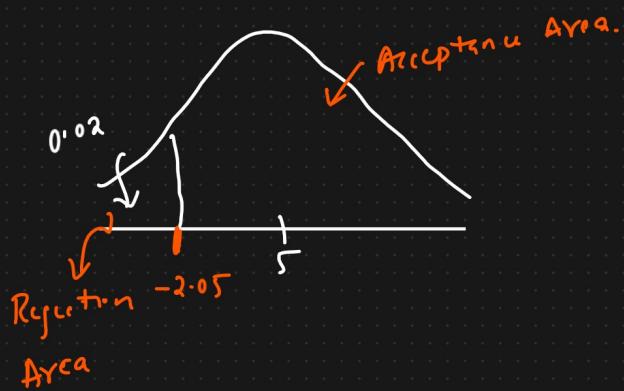
(b) At a 2% significance level, is there enough evidence to support the idea that the warranty should be revised?

$$\text{Ans}) \quad \mu = 5 \quad \sigma = 0.50 \quad n = 40 \quad \bar{x} = 4.8 \text{ years.} \quad C.I = 0.98 \quad \delta = 0.02$$

$$\textcircled{1} \quad H_0 \quad \mu = 5$$

$$H_1 \quad \mu < 5 \quad \{ \text{1 Tail Test} \}.$$

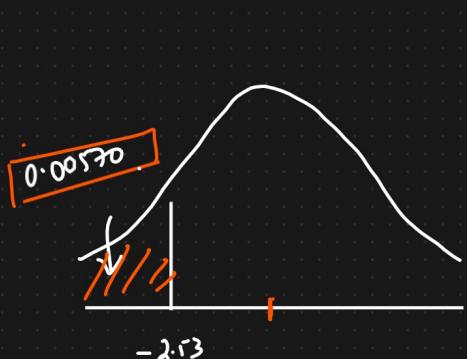
\textcircled{2} Decision Boundary



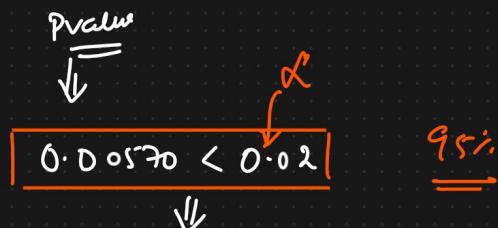
If, $Z_{\text{test}} < -2.05$ then, Reject the Null Hypothesis.

$$\textcircled{*} \quad Z_{\text{test}} = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} = \frac{4.8 - 5}{0.50 / \sqrt{40}} = [-2.53]$$

$-2.53 < -2.05 \Rightarrow \text{True} \Rightarrow \text{Reject the Null Hypothesis.}$



So, to conclude factory needs to revise the warranty of bulbs manufactured.
But what will be new warranty?



Reject the Null Hypothesis.

May be again testing of bulbs needs to done and based on outcome average warranty will be mentioned. Since, worker has already calculated warranty as 4.8 years on a sample of 40 bulbs this can be treated as new warranty. Or one can take even bigger sample of bulbs and perform testing. All this depends on domain expert means the people who are making bulbs or working in industry. Data scientist will not go and perform testing on bulbs to determine new warranty. Data scientist can only suggest whether to revise or not revise the warranty and its on domain experts to again perform testing of their product and come up with revised warranty if required. Data scientist just rely on info given by domain experts.

② T Test

Difference b/w Ztest and Ttest:
In Ztest we are provided with population standard deviation.

Whereas, in Ttest we are not provided population standard deviation. In Ttest we are provided with sample standard deviation.

DATA ANALYST

Also, in Ttest we use Z table for determining cutoff points for rejection regions and area under the curve/p value.

Whereas, in case of T test we use T table for the same.

① In the population the average IQ is 100. A team of researchers want to test a new medication to see if it has either a positive or negative effect on intelligence, or no effect at all. A sample of 30 participants who have taken the medication has a mean of 140 with a standard deviation of 20. Did the medication affect intelligence? C.I = 95% $\alpha = 0.05$

$$\text{H}_0: \mu = 100 \quad n=30 \quad \bar{x}=140 \quad s=20 \quad C.I=95\% \\ \text{H}_1: \mu \neq 100 \quad \alpha = 0.05$$

② $\alpha = 0.05$

In Ttest as part of second steps after stating null and alternate hypothesis, we need to calculate Degree of Freedom. Reason why we calculate this(dof) is we are provided with sample std dev and not population std dev.

Degree of freedom

$$dof = n - 1 = 30 - 1 = 29, \dots$$

(dof = sample_size - 1. Dof will be used in T table to find cutoff value for rejection regions)



Explaining dof:
Suppose there are 3 people and 3 chairs in a room.

Choice 1: 1st person can choose b/w 3 chair to sit and will sit on one of the seat.

Choice 2: 2nd person can choose b/w 2 chairs and will sit on one of the seat.

For third person there will be no choice left and he/she has to sit on one remaining seat.

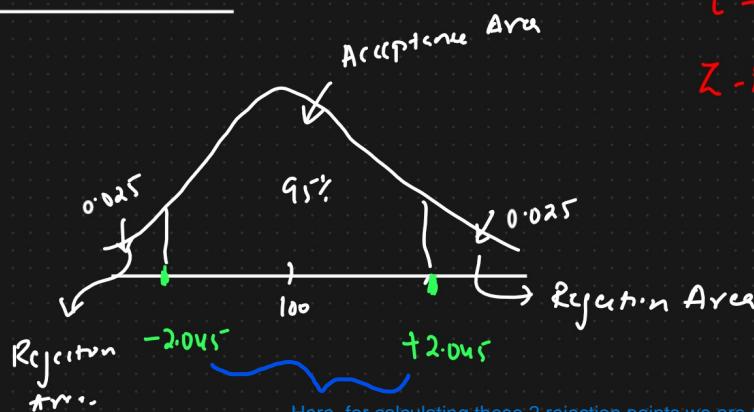
$t-tst \Rightarrow \text{sample std}$

Since, there are 2 choices in total dof in this case will be 2.

$Z-tst \Rightarrow \text{population std}$

From google/byjus:

The degrees of freedom of a system is the number of parameters of the system that may vary independently. For example, a point in the plane has two degrees of freedom for translation: its two coordinates i.e; x and y coordinates.



Here, for calculating these 2 rejection points we are going to user T table and not a Z table.

If t_{tst} is less than -2.045 and greater than $+2.045$, Reject the Null Hypothesis.

In T table check;

- dof
- whether 1 or 2 tail
- p value/area under 1 or 2 table
- C.I

Based on above we calculate Tscore or cutoff value for rejection region in case of T test.

④ T Test statistics

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{140 - 100}{20/\sqrt{30}} = \frac{40}{3.65} = 10.96$$

$t > 2.045$ Reject the Null Hypothesis.

Conclusion : Medication has a true effect on intelligence.



2000

200 - 300 \Rightarrow

There is "A" branch for some bank which has atm installed. As, a data analyst you are asked to perform hypothesis testing to determine whether it would be good to install another atm at a nearby bank's branch "B".

So with the help of domain experts(tellers, cashiers, bank managers) we will collect data like mean transaction happening each day. And based on that will ask domain experts like what should be min and max transaction within which atm installation seems good and beyond which atm installation seems bad(CI and significance value). Using all these data analyst will perform hypotheses testing and tell whether to install or not to install atm at the near by location.

1. Before distributing covid vaccine one will perform hypothesis testing on its efficiency.
2. As a event manager one will perform hypothesis testing on number of plates of food to be ordered for all guests. (Event manager tells take 50 plates extra food. +50 -50 this can be used to place CI which is calculated from domain expert's knowledge).
3. HR wants to order T shirts for whole company. HR will not ask each employee about t shirt size. They will simply take a sample of people and based on that will decide t shirt size of population and give order to manufacturer. But is this assumption going to be true or not? For this we need to perform hypothesis testing. For this HR needs to order first time and then based on feedback received need to judge that how much problem is there and then we need to perform hypothesis testing to determine new size of t shirts that needs to be placed with manufacturer from next time.