

⑦ STATISTICS ⑦

⑦ Inferential statistics ⑦

⑦ Inferential statistics is a part of statistics where we can estimate the value of any parameter of population data based on the suitable or appropriate experiment of sample data.

In one word, here, we can estimate the information about population based on the information of sample.

⑦ The important topics of inferential statistics are —

⑦ Central Limit Theorem ⑦

⑦ Definition : The central limit theorem states that the distribution of the sample means of a large number of independent and identically distributed random variables will approach a normal distribution, regardless of the underlying distribution of the variables.

The conditions required for the CLT to hold are —

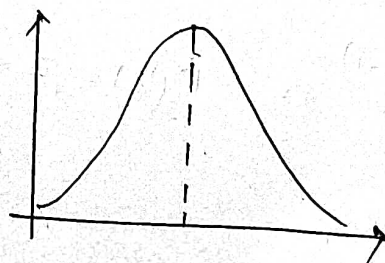
- ① The sample size ≥ 30 .
- ② The sample is drawn from a finite population or infinite population with finite variance.
- ③ The random variables in the sample are independent and identically distributed.

⑦ If the population has the mean μ and variance σ^2 then the empirical mean and variance would be μ and σ^2/n .

$$\bar{X}_i \sim N(\mu, \sigma^2/n) \quad N = \text{sample size.}$$

⑦ Z-score ⑦

for normal distribution —



① Hypothesis Testing :

A statistical hypothesis test is a method of statistical inference used to decide whether the data at hand sufficiently support a particular hypothesis. Hypothesis testing allows us to make probabilistic statements about population parameters.

② Null hypothesis and Alternative Hypothesis :

① Null hypothesis :

In simple word, the null hypothesis is a statement that assumes there is no significant effect or relationship between the variables being studied.

② Alternative hypothesis :

Alternative hypothesis is a statement that contradicts the null hypothesis and claims that there is significant effects or relationship between the variables being studied.

It is also called Research hypothesis.

③ The purpose of hypothesis testing is to gather evidence (data) to either reject or fail to reject the null hypothesis in favour of alternative hypothesis.

④ It is important to note that failing to reject the null hypothesis doesn't necessarily mean that the null hypothesis is true. It just means that there is not enough evidence to support the alternative hypothesis.

⑤ Type-I and Type-II error

In hypothesis testing, there are two types of errors that can occur when making a decision about the null hypothesis.

① Type-I ② Type-II

① Type-I Error :

It occurs when the sample results lead to the rejection of null hypothesis, when it is actually true. (False positive)

② It is the mistake of finding a significant effect or relationship when there is none.

It is also known as significance level, denoted by α . By choosing a significance level researchers can control the risk of making a type I error.

③ Type-II error

It occurs when based on the sample results, the null hypothesis is not rejected when it is actually false. This means that the researchers fail to detect a significant effect or relationship when at least one actually exists. It is denoted by β .

		Truth about the population	
		H_0 true	H_0 false
Decision based on sample	Reject H_0	False positive (Type-I) (α)	Correct Decision
	Accept H_0	Correct Decision	False Negative (Type-II) (β)

④ Trade off between Type-I error and Type-II error :

Let say we are decreasing our α value that is type-I error. So, we reduce the probability of rejecting null hypothesis when it is true. Indirectly we are raising the chance of happening the type-II error. (that we are accepting Null hypothesis may be when it is actually false).

① One Sided and Two Sided Test

① One Sided (one tailed) Test :

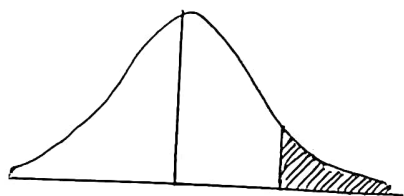
A one sided test is used when the researcher is interested in testing the effect in a specific direction (either greater than or less than)

Example — A researcher wants to test whether a new medication increases the average recovery rate compared to the existing medication.

So, here the Null hypothesis be —

H_0 : There is no significant changes.

H_a : New medication increases the average recovery rates.



have
Here we are interest in the right side of graph (greater than). That is also called Right tailed test. (same logic for left tail test.

① Advantage :

(i) More powerful — Here the entire significance level (α) is allocated to one tail of distribution, that means that the test is more likely to detect an effect in that specified direction.

(ii) Directional Hypothesis.

① Disadvantage :

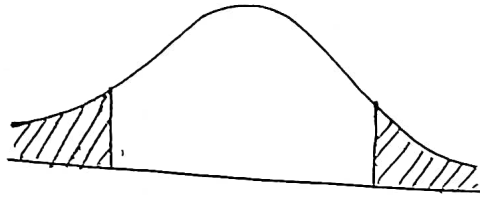
(i) Missed effect in another side.

(ii) Increase the risk of type-I error.

⑩ Two-side or two tailed test :

A two-sided test is used when the researcher is interested in testing the effect in both directions.

Example — A researcher wants to test a new medication has different average recovery rate compare to existing one. So, here the researcher wants to test in both directions. (Maybe it is less than or greater than).



Null hypothesis $H_0 : \mu_{\text{before}} = \mu_{\text{after}}$

Alternative hypothesis $H_0 : \mu_{\text{before}} \neq \mu_{\text{after}}$

(it would be
 $\mu_{\text{before}} > \mu_{\text{after}}$ or
 $\mu_{\text{before}} < \mu_{\text{after}}$)
that's why it is called
two-tailed or two sided
test.

⑪ Advantage :

① Detects effects in both directions.

(ii) More conservative — because the significance level is split between both tails, this reduces the type-I errors in case where the direction of effect is uncertain.

⑫ Disadvantage :

① less powerful : It increase the type-I error as the α is split into two direction of distribution.

(ii) Not appropriate for directional hypothesis.

⑩ P-Value :

P-Value is the probability of getting a sample as or more extreme (having more evidence against (H_0)) than our own sample given the Null hypothesis (H_0) is true.

⑩ Suppose we have done an experiment of tossing a coin 100 times, and we state that our null hypothesis is that the coin is fair. ($N_{OH} = N_{OT}$) and Alternative hypothesis is the ~~coin is not~~ number of getting head is more than 50.

Now, we did the experiment and got that it gives the result as 53 heads and 47 tails.

Let suppose the corresponding P-Value of this experiment = 0.15

It means that there is 15% probability of getting 53 or more heads when we do the experiment, when our null-hypothesis is true.

In simple words, P-Value is a measure of the strength of the evidence against Null hypothesis that is provided by our sample data.

⑩ Interpreting P-Value with significance of level :

if $P\text{-Value} < 0.01 \Rightarrow$ we have strong evidence against the Null hypothesis.

if $0.01 \leq P < 0.05 \Rightarrow$ moderate evidence against the null hypothesis.

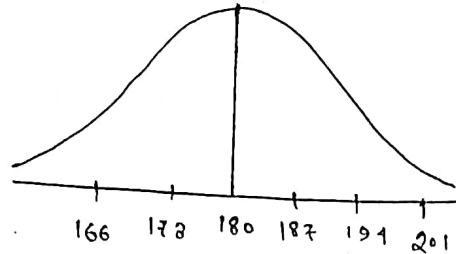
if $0.05 \leq P < 0.1 \Rightarrow$ indicates weak evidence against the null hypothesis. (but there is still some level of uncertainty).

if $P \geq 0.1 \Rightarrow$ indicates weak or no evidence against the null hypothesis.

② Z-score and its Application :

Z-score is a statistic by which we can find how much far it is from mean. For a data point

For example let say we have a dataset whose mean = 180 and s.d = 7.



Now, if I want to know that ~~189~~ How much far 189 from its mean, we apply Z score on it.

So, Formula of

$$Z\text{-score} = \frac{\text{point} - \text{mean}}{\text{s.d}}$$

Here — $Z\text{-score} = \frac{189 - 180}{7} = 1.28$

So, 189 is +1.28 s.d far from its mean (180). (Ans)

③ Application :

① Standardization :

Let suppose we have the following two data columns having different units value.

Age (Y)	Weight (kg)
25	78
20	92
35	65
14	75
26	80

- ① Standardization is a common processing steps that aims to make the input data more suitable for certain machine learning algorithms, specially those that are sensitive to the scale and distribution of the features.

Here we can see that for age column mean = 24.

$$\begin{aligned}\text{Var}(\text{Age}) &= \frac{1}{n} \sum (x_i - \bar{x})^2 \\ &= \frac{1}{5} [1^2 + 1^2 + 11^2 + 10^2 + 2^2] = 18.4\end{aligned}$$

$$\text{S.D} = \sqrt{\text{Var}(\text{Age})} = 6.95$$

in similarly we can find mean = 78 and SD = 8.6948 for weight column.

Now, if we do scaling with respect to these two column our columns looks like —

Age	Weight
0.144	0
-0.575	1.611
1.583	-1.195
-1.139	-0.345
0.288	0.230

It is only done by the formula of Z-score (For each datapoint (x_i))

$$Z\text{-score} = \frac{x_i - \bar{x}}{s}$$

② Compare any score:

Z score is very useful in comparison of two score. (anything).

For example let say in 2020 India's avg run in T20 was 181 with S.D 12 and for 2021, it was 182 with S.D 5. Now, India scored 187 and 185 in the final of 2020 and 2021 respectively. Now, if we compare in which year India score well in final, we can't say directly. So Z score helps us in this situation.

For 2020 —

$$Z_{\text{score}} (187) = \frac{187 - 181}{12} = 0.5$$

For 2021

$$Z_{\text{score}} (185) = \frac{185 - 182}{5} = 0.6$$

so, From here we see that $Z_{\text{score}} (\text{for } 2021) > Z_{\text{score}} (2020)$.
so, although India scored more run in Binal-2020, but
the better & Binal performance was in 2021. (Proved)

① Significance level :

Significance level (α), is a predetermined threshold used in hypothesis testing to determine whether the null hypothesis should be rejected or not. ~~when it actually~~ It represents the probability of rejecting the null hypothesis when it is actually true. (type-I error)

② Type-I error and Type-II error

③ Application of Hypothesis Testing :

- ① testing the effectiveness of interventions or treatment.
- ② Comparing means or proportions
- ③ Analysing relationship between variables
- ④ Evaluating the goodness of fit
- ⑤ Testing the independence of categorical variables
- ⑥ A/B testing