

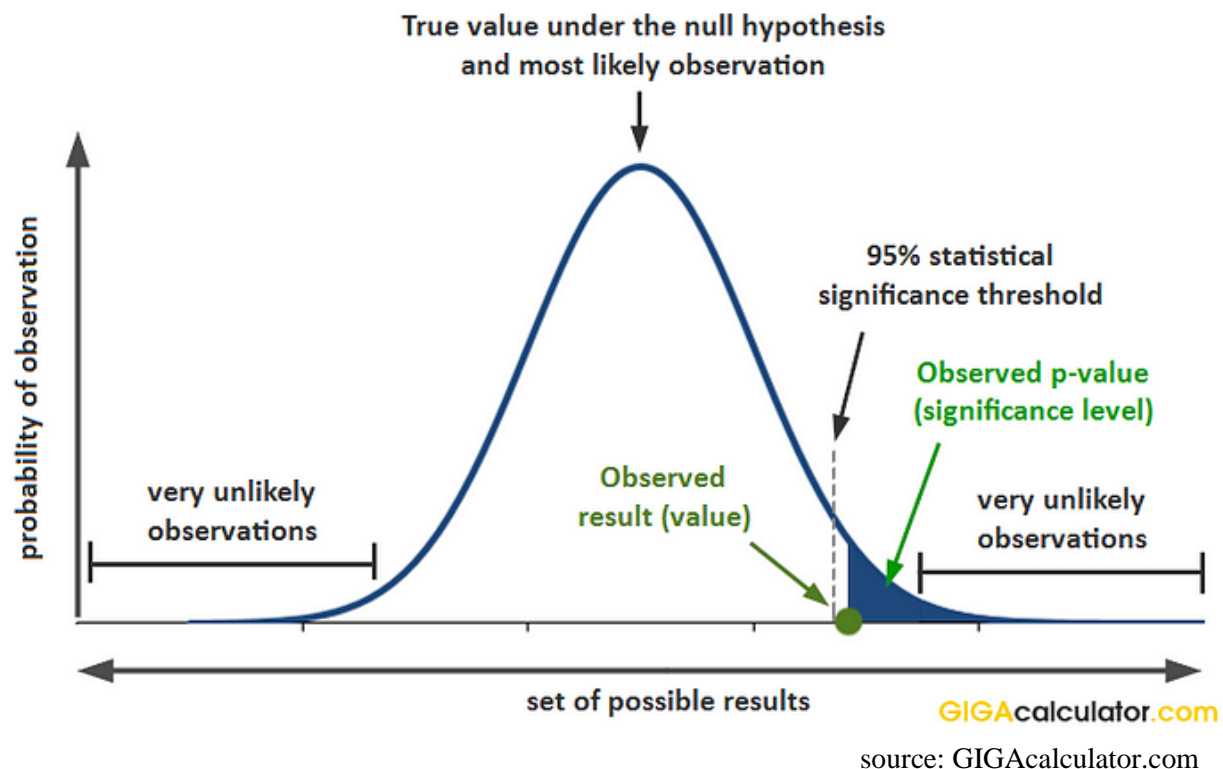
Using Hypothesis Testing such as Z-test, P-Value, T-test, Chi-Square test, ANOVA.

What should be the strategy, when to use what test?

Before proceeding to all the tests, let's have a brief understanding of Hypothesis Testing:

Hypothesis testing is a tool for making statistical inferences about the population data. It is an analysis tool that tests assumptions and determines how likely something is within a given standard of accuracy. Hypothesis testing provides a way to verify whether the results of an experiment are valid.

Probability & Statistical Significance Explained



Null Hypothesis (H₀): The null hypothesis states that the two samples of the population are the same.

Alternate Hypothesis (H_a): It states that there is a significant difference between the two samples of the population.

For the null hypothesis, the same means are assumed to be equal, and we have $H_0: \mu_1 = \mu_2$

& for the alternate hypothesis, the sample means are unequal, and we have $H_a: \mu_1 \neq \mu_2$.

Questions on Null Hypothesis:

Question 1: A medical experiment and trial is conducted to check if a particular drug can serve as the vaccine for Covid-19 and can prevent from occurrence of Corona. Write the null hypothesis and the alternate hypothesis for this situation.

Solution:

The given situation refers to a possible new drug and its effectiveness of being a vaccine for Covid-19 or not. The null hypothesis (H_0) and alternate hypothesis (H_a) for this medical experiment is as follows.

H_0 : The use of the new drug serves as a vaccine and helps for the prevention of Covid-19.

H_a : The use of the new drug is not helpful for the prevention of Covid-19.

Question 2: The teacher has prepared a set of important questions and informs the student that preparing these questions helps in scoring more than 60% marks in the board exams. Write the null hypothesis and the alternate hypothesis for this situation.

Solution:

The given situation refers to the teacher who has claimed that her important questions help to score more than 60% marks in the board exams. The null hypothesis (H_0) and alternate hypothesis (H_a) for this situation is as follows.

H_0 : The important questions given by the teacher are helpful for the students to score more than 60% marks in the board exams.

H_a : The important questions given by the teacher do not really help the students to get a score of more than 60% in the board exams.

Hypothesis Testing P Value:

- In hypothesis testing, the p value is used to indicate whether the results obtained after conducting a test are statistically significant or not.
- Basically, it decides whether we should accept our Null Hypothesis or reject it. The lower the p-value, the more surprising the evidence is, the more ridiculous our null hypothesis looks. And when we feel ridiculous about our null hypothesis, we simply reject it and accept our Alternate Hypothesis.
- If we found the p-value is lower than the predetermined significance value (often called alpha or threshold value (The alpha level can be defined as the acceptable risk of incorrectly rejecting the null hypothesis. The alpha level is usually chosen between 1% to 5%.)) then we reject the null hypothesis. The alpha should always be set before an experiment to avoid bias.
- For example, we generally consider a large population data to be in Normal Distribution so while selecting alpha for that distribution we select it as 0.05 (it means we are accepting if it lies in the 95 percent of our distribution). This means that if our p-value is less than 0.05 we will reject the null hypothesis.

Steps1: Calculate the Standard Error of the sample, which is the Population Standard Deviation divided by the square root of sample size (n).

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

SE = Standard Error
 σ = Population Standard Deviation
 n = Sample Size

Step2: After finding the Standard Error, we take a sample and the mean of that sample and then find the Z-score associated with that mean value.

$$Z_c = \frac{(\bar{X} - \mu)}{(\sigma_x)}$$

$$Z = (\text{Sample Mean} - \text{Population Mean}) / (\text{Standard Error})$$

Z_c = Z-Score
 \bar{X} = Sample Mean
 μ = Population Mean
 σ_x = Standard Error associated with the Sample

Step3: After the p-value associated with the Z-score is calculated, we refer the Z-table to find the probability of the Z-score calculated. Then, to find the p-value, we subtract that probability from 1.

$$\text{P-Value} = 1 - \text{Probability}(\text{Z-score})$$

Step4: Finally, we check if the calculated p-value is greater than the significance level or not.

If the ***P-value*** > ***Significance Level***, then we ***Fail to Reject the Null Hypothesis***.

Or else, if the ***P-value*** < ***Significance Level***, we ***Reject the Null Hypothesis***.

Questions on P-Value:

Question 1: A statistician is testing the hypothesis $H_0: \mu = 120$ using the approach of alternative hypothesis $H_a: \mu > 120$ and assuming that $\alpha = 0.05$. The sample values that he took are as $n = 40$, $\sigma = 32.17$ and $\bar{x} = 105.37$. What is the conclusion for this hypothesis?

Solution:

We know that: $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$

On substitution: $32.17/\sqrt{40} = 5.0865$

As per the test static formula, we get.

$$t = (105.37 - 120) / 5.0865$$

$$t = -2.8762$$

Using the Z-Score table, finding the value of $P(t > -2.8762)$, we get,

$$P(t < -2.8762) = P(t > 2.8762) = 0.003$$

If,

$$P(t > -2.8762) = 1 - 0.003 = 0.997$$

$$P\text{-value} = 0.997 > 0.05$$

As the value of $p > 0.05$, the null hypothesis is accepted.

Therefore, the **null hypothesis is accepted**.

Question 2: P-value is 0.3105. If the level of significance is 5%, find if we can reject the null hypothesis.

Solution: Looking at the P-value table, the p-value of 0.3105 is greater than the level of significance of 0.05 (5%), we fail to reject the null hypothesis.

Hypothesis Testing Z-test:

- A z-test is a test that is used to check if the means of two populations are different or not provided the data follows a normal distribution.
- To start with, the null hypothesis and the alternative hypothesis must be set up and the value of the z test statistic must be calculated. The decision criterion is based on the z critical value.

Conditions to apply z-test:

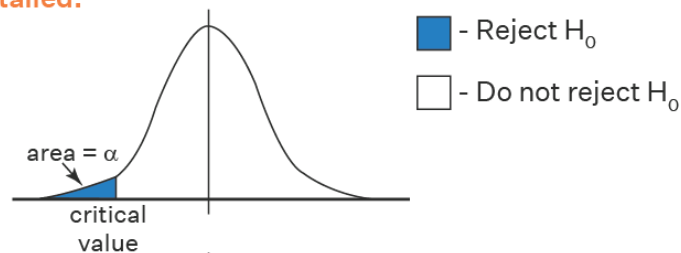
- Z test is a statistical test that is conducted on normally distributed data to check if there is a difference in means of two data sets.
- The sample size should be greater than 30 and the population variance must be known to perform a z test.
- The one-sample z test checks if there is a difference in the sample and population mean,
- The two sample z test checks if the means of two different groups are equal.

Refer: <https://www.z-table.com/>

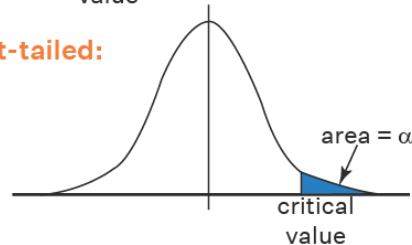
Rejection Region for Null Hypothesis



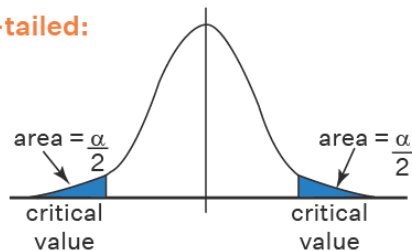
left-tailed:



right-tailed:



two-tailed:



Question 1: A teacher claims that the mean score of students in his class is greater than 82 with a standard deviation of 20. If a sample of 81 students was selected with a mean score of 90 then check if there is enough evidence to support this claim at a 0.05 significance level.

Solution:

As the sample size is 81 and population standard deviation is known, this is an example of a right-tailed one-sample z test.

Ho: $\mu=82$

Ha: $\mu>82$

From the z table the critical value at $\alpha = 1.645$

$$Z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

$\bar{x} = 90, \mu = 82, n = 81, \sigma = 20$

$z = 3.6$

As $3.6 > 1.645$ thus, the null hypothesis is rejected, and it is concluded that there is enough evidence to support the teacher's claim.

Answer: Reject the null hypothesis

Question 2: An online medicine shop claims that the mean delivery time for medicines is less than 120 minutes with a standard deviation of 30 minutes. Is there enough evidence to support this claim at a 0.05 significance level if 49 orders were examined with a mean of 100 minutes?

Solution: As the sample size is 49 and population standard deviation is known, this is an example of a left-tailed one-sample z test.

Ho: $\mu=120$

Ha: $\mu<120$

From the z table the critical value at $\alpha = -1.645$. A negative sign is used as this is a left tailed test.

$$Z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

$$\bar{x} = 100, \mu = 120, n = 49, \sigma = 30$$

$$z = -4.66$$

As $-4.66 < -1.645$ thus, the null hypothesis is rejected, and it is concluded that there is enough evidence to support the medicine shop's claim.

Answer: Reject the null hypothesis

Hypothesis Testing T-test:

- The t-test formula helps us to compare the average values of two data sets and determine if they belong to the same population or are they different.
- The t-score is compared with the critical value obtained from the t-table. The large t-score indicates that the groups are different, and a small t-score indicates that the groups are similar.

There are three different versions of t-tests:

- One sample t-test which tells **whether means of sample and population are different**.

$$t = \frac{\bar{x} - \mu}{s_{\bar{x}}} \quad \text{where} \quad s_{\bar{x}} = \frac{s}{\sqrt{n}}$$

where

μ = Proposed constant for the population mean

\bar{x} = Sample mean

n = Sample size (i.e., number of observations)

s = Sample standard deviation

$s_{\bar{x}}$ = Estimated standard error of the mean (s/\sqrt{n})

- Two sample t-test also is known as independent t-test — it compares the means of two independent groups and determines whether there is statistical evidence that the associated population means are significantly different.

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{s^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

$$s^2 = \frac{\sum_{i=1}^{n_1} (x_i - \bar{x}_1)^2 + \sum_{j=1}^{n_2} (x_j - \bar{x}_2)^2}{n_1 + n_2 - 2}$$

- Paired t-test when you want to compare means of the different samples from the same group or **which compares means from the same group at different times.**

Condition to apply t-test:

- The sample size is smaller than 30 and the population variance is not known to perform a t-test.

Question 1: A company wants to improve its sales. The previous sales data indicated that the average sale of 25 salesmen was \$50 per transaction. After training, the recent data showed an average sale of \$80 per transaction. If the standard deviation is \$15, find the t-score. Has the training provided improved the sales?

Solution:

Ho: The population mean = the claimed value $\Rightarrow \mu = \mu_0$

Ha: The population mean not equal to the claimed value $\Rightarrow \mu \neq \mu_0$

t-test formula for independent test is $t = \frac{\bar{x} - \mu}{s/\sqrt{n}}$

Mean sale = 80, $\mu = 50$, $s = 15$ and $n = 25$

substituting the values, we get $t = (80-50)/(15/\sqrt{25})$

$$t = (30 \times 5)/15 = 10$$

Looking at the t-table we find $10 > 1.711$. (I.e., CV for $\alpha = 0.05$). \therefore the accepted hypothesis is not true.

Thus, we conclude that the training boosted the sales.

Hypothesis Testing Chi-Square Test:

- The Chi-Square test is used when we perform **hypothesis testing on two categorical variables** from a single population or we can say that to **compare categorical variables** from a single population.
- By this we find is there any significant association between the two categorical variables.

$$\chi^2 = \sum \frac{(O - E)^2}{E}$$

Where:

χ^2 = Chi Square obtained
 \sum = the sum of
 O = observed score
 E = expected score

The hypothesis being tested for chi-square is:

Null: Variable A and Variable B are independent

Alternate: Variable A and Variable B are not independent.

The Chi-Square test gives a P-value to help you know the correlation if any!

A hypothesis is in consideration, that a given condition or statement might be true, which we can test later.

For example

- A very small Chi-Square test statistic indicates that the collected data matches the expected data extremely well.
- A very large Chi-Square test statistic indicates that the data does not match very well. If the chi-square value is large, the null hypothesis is rejected.

Applications:

- used by Biologists to determine if there is a significant association between the two variables, such as the association between two species in a community.
- used by Genetic analysts to interpret the numbers in various phenotypic classes.
- used in various statistical procedures to help to decide if to hold onto or reject the hypothesis.
- used in medical literature to compare the incidence of the same characteristics in two or more groups.

Question 1: Calculate the Chi-square value for the following data of incidences of water-borne diseases in three tropical regions.

	India	Equador	South America	Total
Typhoid	31	14	45	90
Cholera	2	5	53	60
Diarrhoea	53	45	2	100
	86	64	100	250

Solution:

Setting up the following table:

Observed	Expected	$O_i - E_i$	$(O_i - E_i)^2$	$(O_i - E_i)^2/E_i$
31	30.96	0.04	0.0016	0.0000516
14	23.04	9.04	81.72	3.546
45	36	9	81	2.25
2	20.64	18.64	347.45	16.83
5	15.36	10.36	107.33	6.99
53	24	29	841	35.04
53	34.4	18.6	345.96	10.06
45	25.6	19.4	376.36	14.7
2	40	38	1444	36.1
				125.5160516

Answer: Chi Square = 125.516

Hypothesis Testing ANOVA Test:

- It is also called an analysis of variance and is used to compare multiple (three or more) samples with a single test.
- It is used when the categorical feature has more than two categories.

The hypothesis being tested in ANOVA is:

Null: All pairs of samples are same i.e., all sample means are equal

Alternate: At least one pair of samples is significantly different

ANOVA Test Table



Source of Variation	Sum of Squares	Degrees of Freedom	Mean Squares	F Value
Between Groups	$SSB = \sum n_j(\bar{X}_j - \bar{X})^2$	$df_1 = k - 1$	$MSB = SSB / (k - 1)$	$f = MSB / MSE$
Error	$SSE = \sum \sum (X - \bar{X}_j)^2$	$df_2 = N - k$	$MSE = SSE / (N - k)$	
Total	$SST = SSB + SSE$	$df_3 = N - 1$		

One Way ANOVA:

The one way ANOVA test is used to determine whether there is any difference between the means of three or more groups. A one way ANOVA will have only one independent variable. The hypothesis for a one way ANOVA test can be set up as follows:

Null Hypothesis: H_0 : $\mu_1 = \mu_2 = \mu_3 = \dots = \mu_k$

Alternative Hypothesis: H_a : The means are not equal

Decision Rule: If *test statistic* > *critical value* then rejects the null hypothesis and conclude that the means of at least two groups are statistically significant.

The steps to perform the one way ANOVA test are given below:

Step 1: Calculate the mean for each group.

Step 2: Calculate the total mean. This is done by adding all the means and dividing it by the total number of means.

Step 3: Calculate the SSB.

Step 4: Calculate the between groups degrees of freedom.

Step 5: Calculate the SSE.

Step 6: Calculate the degrees of freedom of errors.

Step 7: Determine the MSB and the MSE.

Step 8: Find the f test statistic.

Step 9: Using the f table for the specified level of significance, α , find the critical value. This is given by $F(\alpha, df1, df2)$.

Step 10: If $f > F$ then reject the null hypothesis.

Limitations of One Way ANOVA Test:

The one way ANOVA is an omnibus test statistic. This implies that the test will determine whether the means of the various groups are statistically significant or not. However, it cannot distinguish the specific groups that have a statistically significant mean. Thus, to find the specific group with a different mean, a post hoc test needs to be conducted.

Two Way ANOVA:

- The two way ANOVA has two independent variables. Thus, it can be thought of as an extension of a one way ANOVA where only one variable affects the dependent variable.
- A two way ANOVA test is used to check the main effect of each independent variable and to see if there is an interaction effect between them.
- To examine the main effect, each factor is considered separately as done in a one way ANOVA. Furthermore, to check the interaction effect, all factors are considered at the same time.

There are certain assumptions made for a two way ANOVA test. These are given as follows:

1. The samples drawn from the population must be independent.
2. The population should be approximately normally distributed.
3. The groups should have the same sample size.
4. The population variances are equal

Suppose in the two way ANOVA example, as mentioned above, the income groups are low, middle, high. The gender groups are female, male, and transgender. Then there will be 9 treatment groups and the three hypotheses can be set up as follows:

H01: All income groups have equal mean anxiety.

H11: All income groups do not have equal mean anxiety.

H02: All gender groups have equal mean anxiety.

H12: All gender groups do not have equal mean anxiety.

H03: Interaction effect does not exist

H13: Interaction effect exists.

Important Notes on ANOVA Test:

- ANOVA test is used to check whether the means of three or more groups are different or not by using estimation parameters such as the variance.
- An ANOVA table is used to summarize the results of an ANOVA test.

- There are two types of ANOVA tests - one way ANOVA and two way ANOVA.
- One way ANOVA has only one independent variable while a two way ANOVA has two independent variables.

Question1: Three types of fertilizers are used on three groups of plants for 5 weeks. We want to check if there is a difference in the mean growth of each group. Using the data given below apply a one way ANOVA test at 0.05 significant level.

Fertilizer 1	Fertilizer 2	Fertilizer 3
6	8	13
8	12	9
4	9	11
5	11	8
3	6	7
4	8	12

Solution:

H₀: $\mu_1 = \mu_2 = \mu_3$

H_a: The means are not equal

Fertilizer 1	Fertilizer 2	Fertilizer 3
6	8	13
8	12	9
4	9	11
5	11	8
3	6	7
4	8	12
$\bar{X}_1 = 5$	$\bar{X}_2 = 9$	$\bar{X}_3 = 10$

Total mean, $\bar{X} = 8$

$n_1 = n_2 = n_3 = 6$ & $k = 3$

$SSB = 6(5 - 8)^2 + 6(9 - 8)^2 + 6(10 - 8)^2$

$= 84$

$df_1 = k - 1 = 2$

Fertilizer 1	$(X - 5)^2$	Fertilizer 2	$(X - 9)^2$	Fertilizer 3	$(X - 10)^2$
6	1	8	1	13	9
8	9	12	9	9	1
4	1	9	0	11	1
5	0	11	4	8	4
3	4	6	9	7	9
4	1	8	1	12	4
$X_1^- = 5$	Total = 16	$X_1^- = 9$	Total = 24	$X_1^- = 10$	Total = 28

$$SSE = 16 + 24 + 28 = 68$$

$$N = 18$$

$$df2 = N - k = 18 - 3 = 15$$

$$MSB = SSB / df1 = 84 / 2 = 42$$

$$MSE = SSE / df2 = 68 / 15 = 4.53$$

$$ANOVA \text{ test statistic, } f = MSB / MSE = 42 / 4.53 = 9.33$$

Using the f table at $\alpha = 0.05$ the critical value is given as $F(0.05, 2, 15) = 3.68$

As $f > F$, thus, the null hypothesis is rejected and it can be concluded that there is a difference in the mean growth of the plants.

Answer: Reject the null hypothesis