# Module 3 Assignment

1. What is hypothesis testing in statistics?

Ans: Hypothesis testing is a structured method used to determine if the findings of a study provide evidence to support a specific theory relevant to a larger population.

Hypothesis Testing is a type of statistical analysis anal in which you put your assumptions about a population parameter to the test. It is used to estimate the relationship between 2 statistical variables. When performing hypothesis testing, you must understand different data types, such as nominal data.

2.  What is the null hypothesis, and how does it differ from the alternative hypothesis?

Ans: Null hypothesis suggests that there is no relationship between the two variables. Null hypothesis is also exactly the opposite of the alternative hypothesis. Null hypothesis is generally what researchers or scientists try to disprove and if the null hypothesis gets accepted then we have to make changes in our opinion i.e. we have to make changes in our original opinion or statement in order to match null hypothesis. Null hypothesis is represented as H0. If my alternative hypothesis is that 55**%** of boys in my town are taller than girls then my alternative hypothesis will be that 55% of boys in my town are not taller than girls.

Alternative hypothesis is a method for reaching a conclusion and making inferences and judgements about certain facts or a statement. This is done on the basis of the data which is available. Usually, the statement which we check regarding the null hypothesis is commonly known as the alternative hypothesis. Most of the times alternative hypothesis is exactly the opposite of the null hypothesis. This is what generally researchers or scientists try to approve. Alternative hypothesis is represented as Ha or H1. If my null hypothesis is that 55% of boys in my town are not taller than girls then my alternative hypothesis will be that 55% of boys in my town are taller than girls.

3. What is the significance level in hypothesis testing, and why is it important?

Ans: Hypothesis testing is a fundamental concept in statistics used to make decisions about a population based on sample data. One of the key components of hypothesis testing is the significance level, often denoted by the Greek letter alpha ($\alpha$). Let's delve into what significance level means and its importance in hypothesis testing.

The significance level ($\alpha$) is a threshold set by the researcher before conducting the test. It represents the probability of rejecting the null hypothesis when it is actually true. In other words, it is the probability of making a Type I error, which is a false positive.

Common significance levels are 0.05 (5%), 0.01 (1%), and 0.10 (10%). For instance, if $\alpha = 0.05$, there is a 5% risk of concluding that there is an effect when there is none.

The significance level is crucial because it helps control the likelihood of making incorrect conclusions. By setting a significance level, researchers can manage the risk of Type I errors. A lower significance level means stricter criteria for rejecting the null hypothesis, reducing the chance of a false positive but increasing the chance of a Type II error (failing to reject a false null hypothesis)

4. What does a P-value represent in hypothesis testing?

Ans: A p value is used in hypothesis testing to help you support and reject the null hypothesis. The p value is the evidence against a null hypothesis. The smaller the p-value, the stronger the evidence that you should reject the null hypothesis.

P values are expressed as decimals although it may be easier to understand what they are if you convert them to a percentage. For example, a p value of 0.0254 is 2.54%. This means there is a 2.54% chance your results could be random (i.e. happened by chance). That's pretty tiny. On the other hand, a large p-value of .9(90%) means your results have a 90% probability of being completely random and *not* due to anything in your experiment. Therefore, the smaller the p-value, the more important your results.

When you run  hypothesis test, you compare the p value from your test to the alpha level you selected when you ran the test. Alpha levels can also be written as percentages.

5. How do you interpret the P-value in hypothesis testing?

Ans: The p-value measures the probability of obtaining results at least as extreme as the observed results, assuming that the null hypothesis is true. Essentially, it helps us gauge whether the observed data could have occurred by random chance.

### Common Significance Levels

- **0.05 (5%)**: A common threshold. If the p-value is less than 0.05, you reject the null hypothesis.

- **0.01 (1%)**: More stringent. If the p-value is less than 0.01, you have stronger evidence against the null hypothesis.

- **0.10 (10%)**: More lenient. If the p-value is less than 0.10, you might still consider rejecting the null hypothesis, but with caution.

6. What are Type 1 and Type 2 errors in hypothesis testing?

 Ans: In hypothesis testing, a Type I error is  false positive while a Type II error is a false negative.

7. What is the difference between a one-tailed and a two-tailed test in hypothesis testing?

Ans: One and Two-Tailed Test**s** are ways to identify the relationship between the statistical variables. For checking the relationship between variables in a single direction (Left or Right direction), we use a one-tailed test. A two-tailed test is used to check whether the relations between variables are in any direction or not.

8. What is the Z-test, and when is it used in hypothesis testing?

Ans: Z test are statistical hypothesis testing techniques that are used to determine whether the null hypothesis relating to comparing sample means or proportions with that of population at a given significance level can be rejected or otherwise based on the z-statistics or z-score. As a data scientist, you must get a good understanding of the z-tests and its applications to test the hypothesis for your statistical models.

9. How do you calculate the Z-score, and what does it represent in hypothesis testing?

Ans: A z-score (also called a *standard score*) gives you an idea of how far from the mean a data point is. More technically, it's a measure of how many standard deviations below or above the population mean a row score.

10. What is the T-distribution, and when should it be used instead of the normal distribution?

Ans: A probability distribution which is used when we are working with small sample sizes or when the population variance is unknown. The t-distribution accounts for the extra uncertainty that comes with these conditions. By understanding the nuances between the t-distribution and the normal distribution, analysts can avoid missteps in data analysis that could lead to incorrect inferences about their data.

when estimating the mean of a normally distributed population in situations where the sample size is small and the population standard deviation is unknown".

11. What is the difference between a Z-test and a T-test?

Ans: Z-tests are used when the population variance is known and the sample size is large, while t-tests are used when the population variance is unknown and the sample size is small.

Z-test is a statistical test used to determine whether there is a significant difference between sample and population means or between the means of two samples.

It is typically used when the sample size is large (generally $n > 30$) and the population standard deviation is known. The Z-test is based on the standard normal distribution (Z-distribution).

The Z-test compares the means of two populations with a large sample size (typically $\geq 30$) and known population standard deviation. It assesses whether the difference between the means is statistically significant.

T-test is a statistical test used to determine whether there is a significant difference between the means of two groups.

It is particularly useful when the sample size is small (typically $n < 30$) and the population standard deviation is unknown. The T-test relies on the t-distribution, which is similar to the normal distribution but has heavier tails.

12. What is the T-test, and how is it used in hypothesis testing?

Ans: A t test is a statistical test that is used to compare the means of two groups. It is often used in hypothesis testing to determine whether a process or treatment actually has an effect on the population of interest, or whether two groups are different from one another.

13. What is the relationship between Z-test and T-test in hypothesis testing?

Ans: Z-tests are used when the population variance is known and the sample size is large, while t-tests are used when the population variance is unknown and the sample size is small.

This article explains the differences between Z-tests and T-tests, detailing their purposes, assumptions, sample size requirements, and applications in statistical hypothesis testing.


T-test is a statistical test used to determine whether there is a significant difference between the means of two groups.

It is particularly useful when the sample size is small (typically $n < 30$) and the population standard deviation is unknown. The T-test relies on the t-distribution, which is similar to the normal distribution but has heavier tails.

14. What is a confidence interval, and how is it used to interpret statistical results?

Ans: Confidence intervals are a fundamental concept in general statistics and are widely used to quantify uncertainty in an estimate. They have a wide range of applications, from evaluating the effectiveness of a drug, predicting election results, or analyzing sales data. A confidence interval provides the range of values, calculated from the sample, in which we have confidence that the true population parameter lies.

15. What is the margin of error, and how does it affect the confidence interval?

Ans: The margin of error is a range of uncertainty in a statistical estimate. It represents how much the sample results are expected to differ from the actual population value due to random sampling.

A confidence interval is calculated as:

Confidence Interval=Point Estimate±Margin of Error\text{Confidence Interval} = \text{Point Estimate} \pm \text{Margin of Error}Confidence Interval=Point Estimate±Margin of Error

For example, if the average height from a sample is 170 cm and the margin of error is 3 cm, the 95% confidence interval would be:

170±3=[167,173]170 \pm 3 = [167, 173]170±3=[167,173]

So, we can say with 95% confidence that the true population mean lies between 167 cm and 173 cm.

16. How is Bayes' Theorem used in statistics, and what is its significance?

Ans: Bayes' theorem (also known as the Bayes Rule or Bayes Law) is used to determine the conditional probability of event A when event B has already occurred.

The general statement of Bayes' theorem is "The conditional probability of an event A, given the occurrence of another event B, is equal to the product of the event of B, given A, and the probability of A divided by the probability of event B."

Let E1, E2,…, En be a set of events associated with the sample space S, in which all the events E1, E2,…, En have a non-zero probability of occurrence. All the events E1, E2,…, E form a partition of S. Let A be an event from space S for which we have to find the probability, then according to Bayes theorem,

P(Ei|A)=P(Ei)·P(A|Ei)∑k=1nP(Ek)·P(A|Ek)P(Ei|A)=∑k=1nP(Ek)·P(A|Ek)P(Ei)·P(A|Ei)

for k = 1, 2, 3, …., n

17. What is the Chi-square distribution, and when is it used?

Ans: Chi-square ($X^2$) distributions are a family of continuous probability distributions. They're widely used in hypothesis test, including the chi-square goodness of fit test and the chi-square test of independence.

The shape of a chi-square distribution is determined by the parameter $k$, which represents the degrees of freedom. Very few real-world observations follow a chi-square distribution. The main purpose of chi-square distributions is hypothesis testing, not describing real-world distributions. In contrast, most other widely used distributions, like normal

distributions or passion distributions, can describe useful things such as newborns' birth weights or disease cases per year, respectively.

18. What is the Chi-square goodness of fit test, and how is it applied?

Ans: A chi-square ($X^2$) goodness of fit test is a goodness of fit test for a categorical variables. Goodness of fit is a measure of how well a statistical model fits a set of observations.

- When goodness of fit is high, the values expected based on the model are close to the observed values.

- When goodness of fit is low, the values expected based on the model are far from the observed values.

The statistical models that are analysed by chi-square goodness of fit tests are distributions. They can be any distribution, from as simple as equal probability for all groups, to as complex as a probability distributions with many parameters.

19. What is the F-distribution, and when is it used in hypothesis testing?

Ans: F test is a statistical test that is used in hypothesis testing that determines whether the variances of two samples are equal or not. The article will provide detailed information on f test, f statistic, its calculation, critical value and how to use it to test hypotheses. To understand F test firstly we need to have some basic understanding of F-distribution.

20. What is an ANOVA test, and what are its assumptions?

Ans: ANOVA is a statistical test used to examine differences among the means of three or more groups. Unlike a t-test, which only compares two groups, ANOVA can handle multiple groups in a single analysis.

21. What are the different types of ANOVA tests?

Ans:

a) One-Way ANOVA

b) Two-Way ANOVA

c) Repeated Measures ANOVA

d) MANOVA (Multivariate Analysis of Variance)

22. What is the F-test, and how does it relate to hypothesis testing?

Ans: F test is a statistical test that is used in hypothesis testing that determines whether the variances of two samples are equal or not. The article will provide detailed information on f test, f statistic, its calculation, critical value and how to use it to test hypotheses.