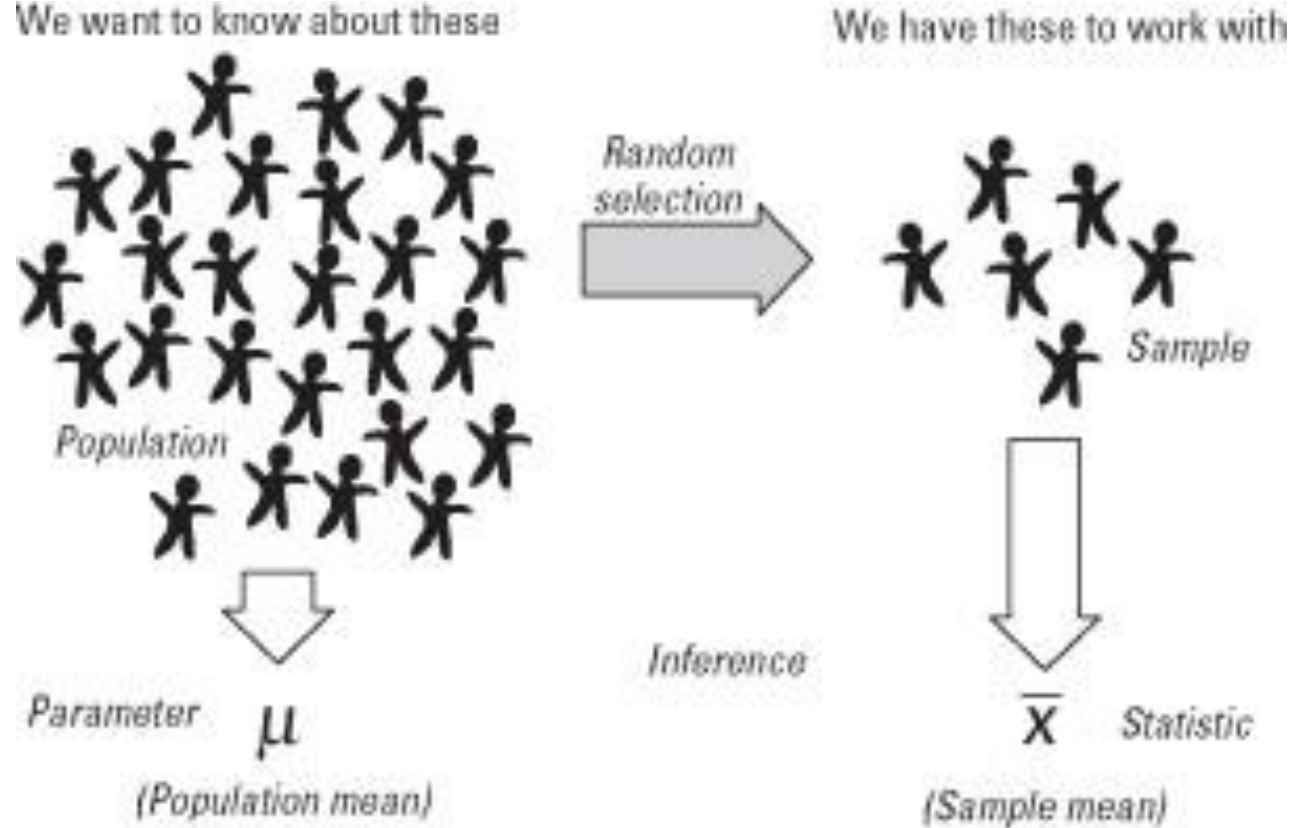


Module 5

Hypothesis Testing-1

Dr. Parvez Alam
Department of Mathematics, SAS
VIT Vellore



Population: A population or universe is a set of similar items or events which is of interest for some question or experiment.

Sample: A finite subset of statistical individuals in a population is called sample.

Example 1: Suppose from a village is having 5000 population and we selected 100 people for a survey. So, here village is population set and 100 peoples are sample set.

Example 2: Suppose in a Godown a huge amount a medicine is stocked item and 30 packets of the medicine taken for the test sample. So, here whole stock is population set and 30 packets are sample set.

Virat Kohli or Sachin Tendulkar!

Who is best? (Claim)



- ☐ Data available.
- ☐ Guess (who is the best?).
- ☐ Test the guess with valid evidence.
- ☐ So, the given statement is a **hypothesis**.

1. Hypothesis Tests

A **hypothesis test** is a process that uses sample statistics to test a claim about the value of a population parameter.

If a manufacturer of rechargeable batteries claims that the batteries they produce are good for an average of at least 1,000 charges, a sample would be taken to test this claim.

A verbal statement, or claim, about a population parameter is called a **statistical hypothesis**.

To test the average of 1000 hours, a pair of hypotheses are stated – one that represents the claim and the other, its complement. When one of these hypotheses is false, the other must be true.

2. Stating a Hypothesis

A **null hypothesis H_0** is a statistical hypothesis that contains a statement of equality such as \leq , $=$, or \geq .

A **alternative hypothesis H_a (or H_1)** is the complement of the null hypothesis. It is a statement that must be true if H_0 is false and contains a statement of inequality such as $>$, \neq , or $<$.

To write the null and alternative hypotheses, translate the claim made about the population parameter from a verbal statement to a mathematical statement.

Note: Claim may come under **H_0** or **H_a** it will depend on the problems symbol \leq , $=$, \geq , $>$, \neq , or $<$.

Cont...

Example 1:

Write the claim as a mathematical sentence. State the null and alternative hypotheses and identify which represents the claim.

A manufacturer claims that its rechargeable batteries have an average life of at least 1,000 charges.

$$\mu \geq 1000$$

└→ Condition of
equality

$H_0: \mu \geq 1000$ (Claim; see the symbol \geq on last slide comes under null)

$H_a: \mu < 1000$

└→ Complement of the null hypothesis, that is Alternate

Cont...

Example 2: VIT claims that 94% of their graduates find employment within six months of graduation.

First write hypothesis, then check under which claims come

$H_0: p = 0.94$ (Claim)

$H_a: p \neq 0.94$

$p = 0.94$

└ Condition of equality, it comes under null

└ Alternate hypothesis

Cont...

Example 4: A cigarette manufacturer claims that **less than one-eighth** of the US adult population smokes cigarettes.

$$H_0: p \geq 0.125$$

$$H_a: p < 0.125 \text{ (Claim)}$$

Example 5: A local telephone company claims that the average length of a phone call is 8 minutes.

$$H_0: \mu = 8 \text{ (Claim)}$$

$$H_a: \mu \neq 8$$

3. P-values

If the null hypothesis is true, a *P*-value (or probability value) of a hypothesis test is the probability of obtaining a sample statistic with a value as extreme or more extreme than the one determined from the sample data.

The *P*-value of a hypothesis test depends on the nature of the test.

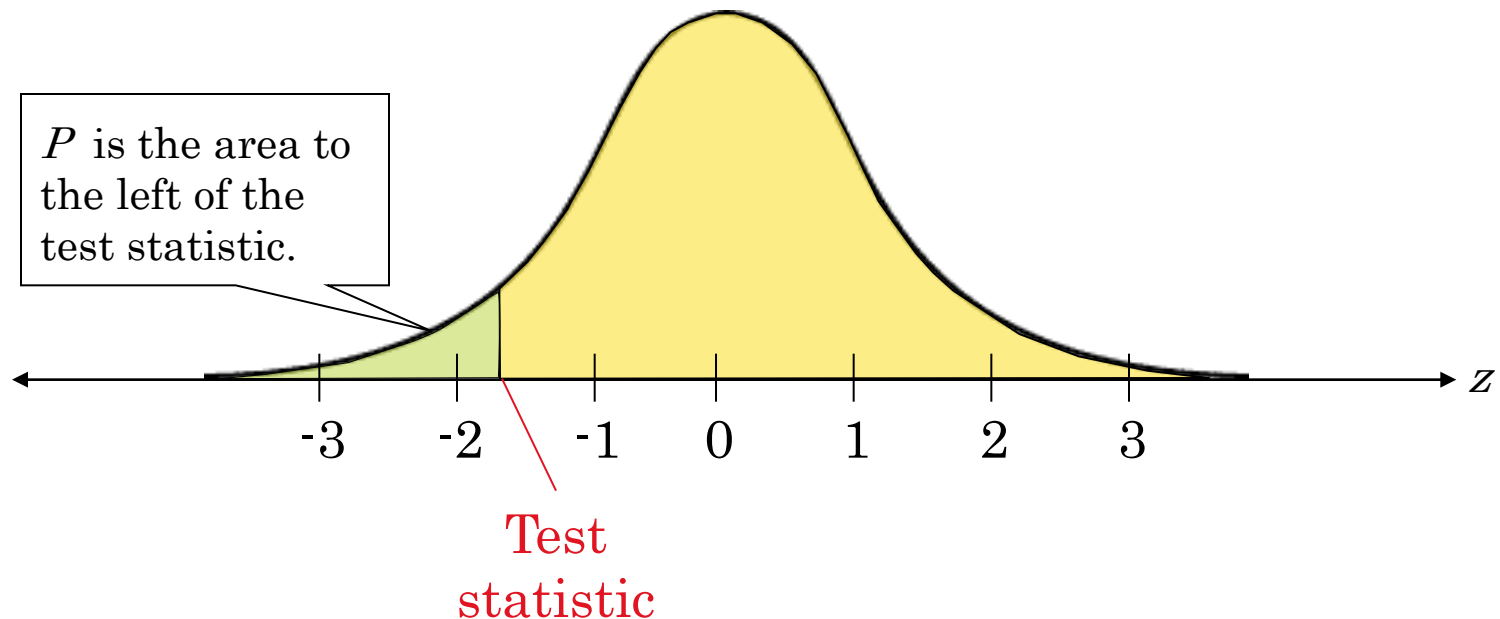
There are three types of hypothesis tests – left-, right-, or two-tailed test. The type of test depends on the region of the sampling distribution that favors a rejection of H_0 . This region is indicated by the alternative hypothesis.

4. Left-tailed Test

1. If the **alternative hypothesis** contains the **less-than inequality symbol** ($<$), the hypothesis test is a **left-tailed test**.

$$H_0: \mu \geq k$$

$$H_a: \mu < k$$

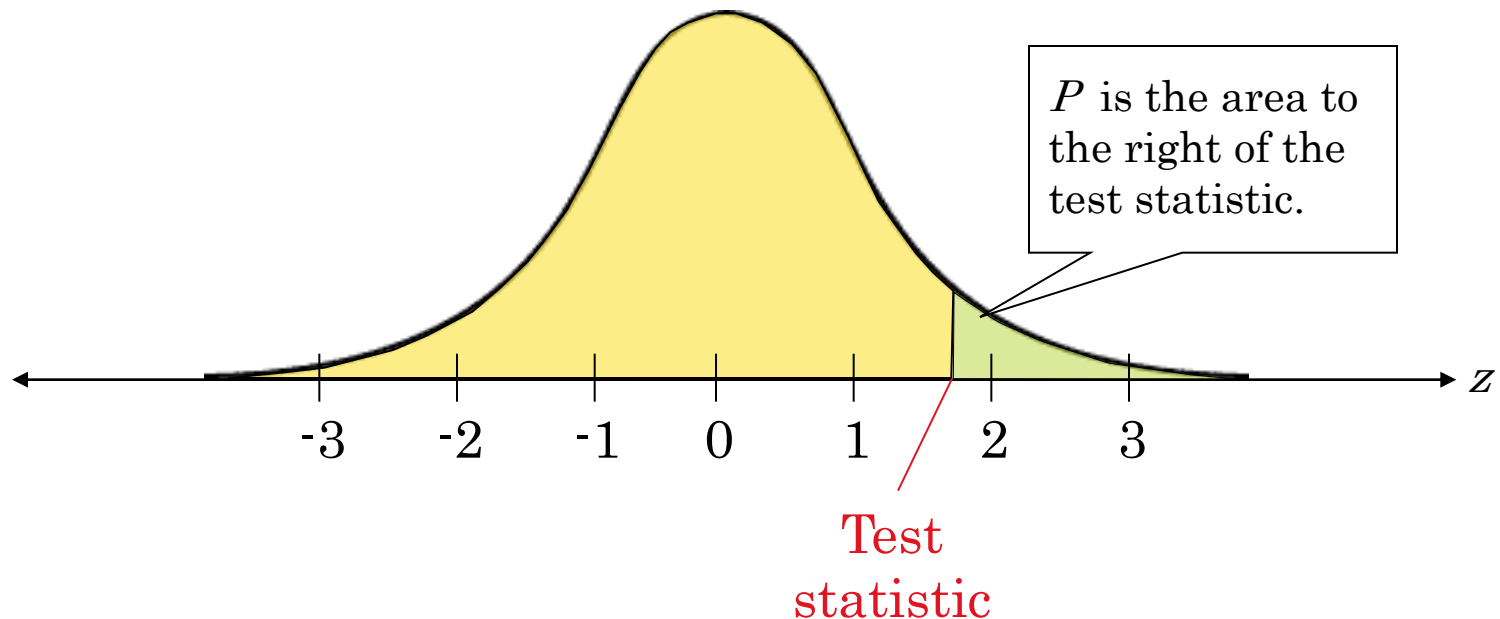


5. Right-tailed Test

2. If the **alternative hypothesis** contains the **greater-than symbol** ($>$), the hypothesis test is a **right-tailed test**.

$$H_0: \mu \leq k$$

$$H_a: \mu > k$$

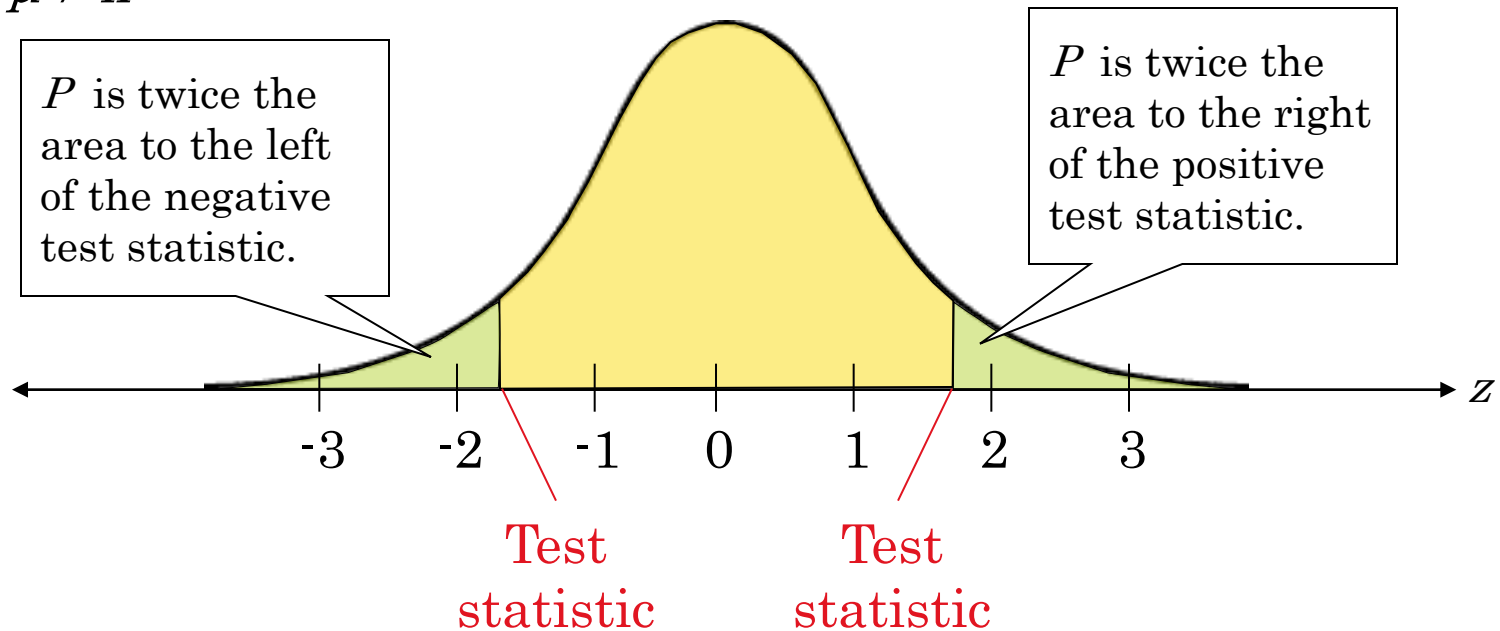


6. Two-tailed Test

3. If the **alternative hypothesis** contains the **not-equal-to symbol (\neq)**, the hypothesis test is a **two-tailed test**. In a two-tailed test, each tail has an area of $\frac{1}{2}P$.

$$H_0: \mu = k$$

$$H_a: \mu \neq k$$



7. Identifying Types of Tests

Example 1:

For each claim, state H_0 and H_a . Then determine whether the hypothesis test is a left-tailed, right-tailed, or two-tailed test.

- a.) A cigarette manufacturer claims that less than one-eighth of the US adult population smokes cigarettes.

$$H_0: p \geq 0.125$$

$$H_a: p < 0.125 \text{ (Claim)}$$

Left-tailed test



- b.) A local telephone company claims that the average length of a phone call is 8 minutes.

$$H_0: \mu = 8 \text{ (Claim)}$$

$$H_a: \mu \neq 8$$

Two-tailed test



8. Types of Errors

No matter which hypothesis represents the claim, always begin the hypothesis test **assuming that the null hypothesis is true.**

At the end of the test, one of two decisions will be made:

1. reject the null hypothesis, **or**
2. fail to reject the null hypothesis.





A **type I error** occurs if the null hypothesis is rejected when it is true.

A **type II error** occurs if the null hypothesis is not rejected when it is false.

Cont...

Actual Truth of H_0		
Decision	H_0 is true	H_0 is false
Do not reject H_0	Correct Decision	Type II Error <i>(Acceptance of False/Alternate hypo.)</i>
Reject H_0	Type I Error <i>(Rejection of truth/Null hypo.)</i>	Correct Decision

Cont...

VERDICT	<i>Set Free</i>		Type 2 Error (β) 
	<i>Jailed</i>	Type 1 Error (α) 	
		<i>Innocent</i>	<i>Guilty</i>
		TRUTH	

Cont...

Example 1:

VIT claims that 94% of their graduates find employment within six months of graduation. What will a type I or type II error be?

H_0 : $p = 0.94$ (Claimed, i.e. Null Hypo.)

H_a : $p \neq 0.94$ (Alt. Hypo.)

A type I error is rejecting the null when it is true.

The population proportion is actually 0.94, but is rejected.
(We believe it is not 0.94.)

A type II error is failing to reject the null when it is false.

The population proportion is not 0.94, but is not rejected.
(We believe it is 0.94.)

9. Level of Significance

In a hypothesis test, the **level of significance** is your **maximum allowable probability** of making a type I error. It is denoted by α (alpha).

└→ Hypothesis tests
are based on α .

The probability of making a type II error is denoted by β (beta).

By setting the level of significance at a small value, you are saying that you want the probability of rejecting a true null hypothesis to be small.

Commonly used levels of significance:

$$\alpha = 0.01$$

$$\alpha = 0.02$$

$$\alpha = 0.05$$

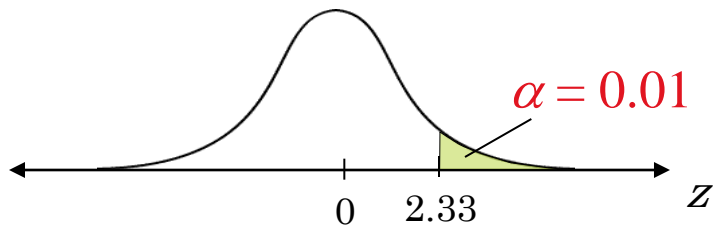
$$\alpha = 0.1$$

10. Rejection Regions and Critical Values

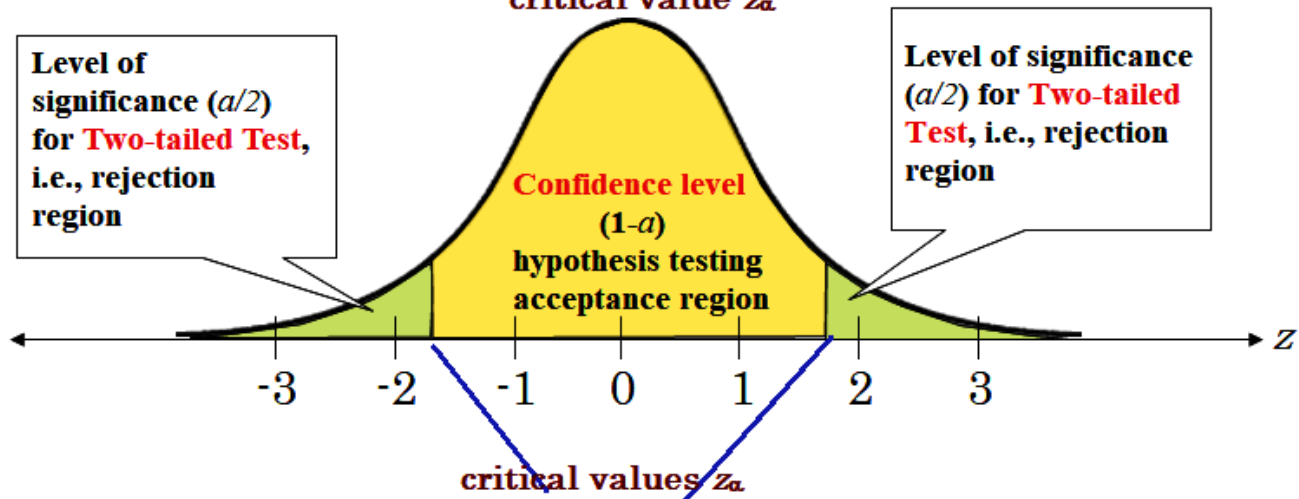
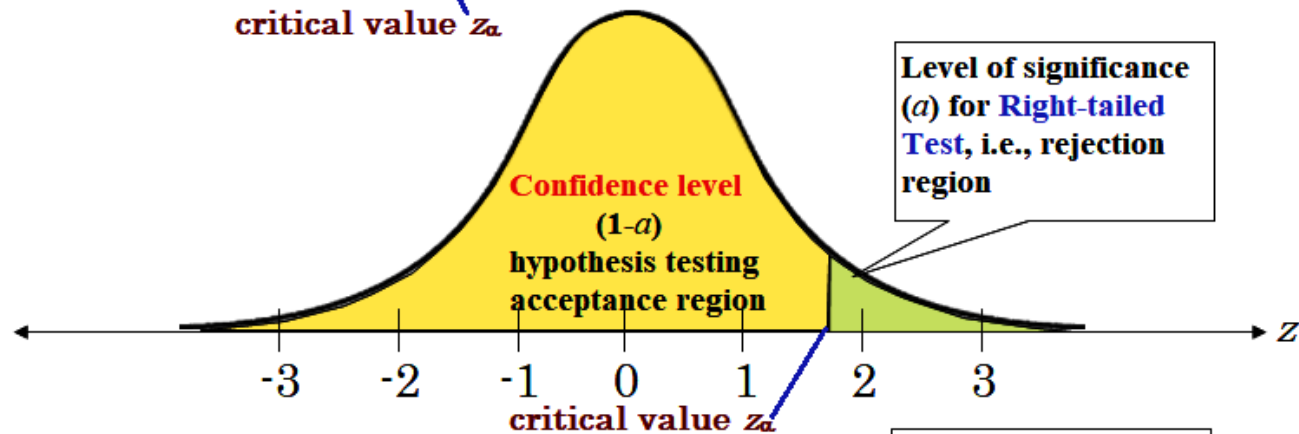
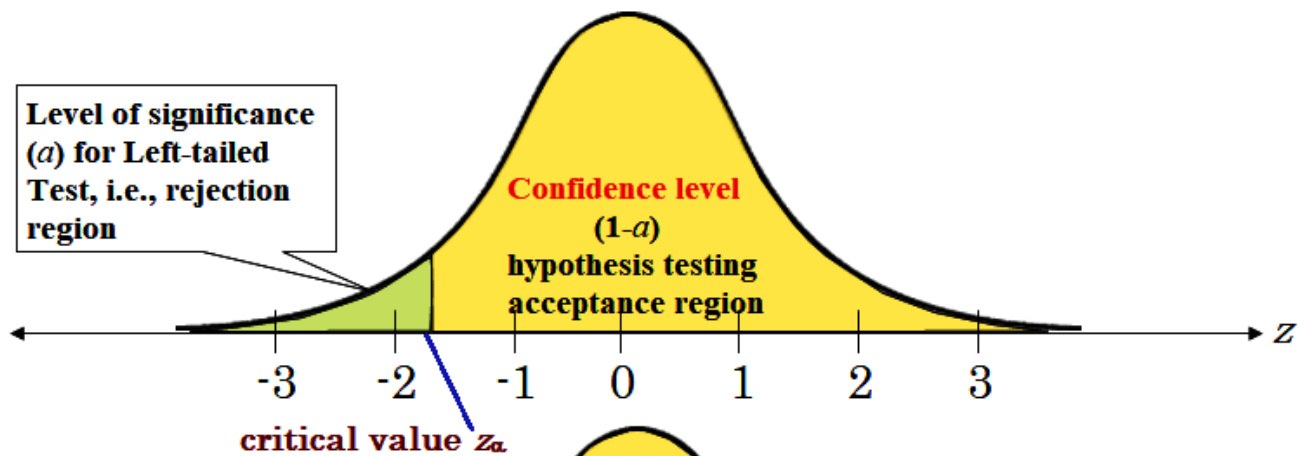
A **rejection region** (or **critical region**) of the sampling distribution is the range of values for which the null hypothesis is not probable. If a test statistic falls in this region, the null hypothesis is rejected. A critical value z_α separates the **rejection region** from the **non-rejection region**.

Example:

Find the critical value and rejection region for a right tailed test with $\alpha = 0.01$.



The rejection region is to the right of $z_\alpha = 2.33$.



Test of Significance: Large Samples (Z-test)

Either the population is normally distributed or $n \geq 30$.

We have following types of z-tests

1. Test of significance for single mean
2. Test of significance for difference of means of two large samples
3. Test of significance for a single proportion
4. Test of significance for difference of proportions

We have two methods for above kind of tests (problems),

1. Rejection region method
2. P-Value method

Type -1 Problem: (If single mean given)

z-Test using sample mean \bar{x} and population mean μ ($n \geq 30$)

Let x_1, x_2, \dots, x_n be a random sample of size n , drawn from a large population with mean μ and variance σ^2 .

Let \bar{x} denote the mean of the sample and s^2 denote the variance of the sample.

We know that $\bar{x} \sim N(\mu, \sigma^2/n)$. The standard normal variate corresponding to \bar{x} is $Z = \frac{\bar{x} - \mu}{S.E.(\bar{x})}$, where $S.E.(\bar{x}) = \sigma/\sqrt{n}$.

We set up the null hypothesis that there is no difference between the sample mean and the population mean. The test statistic is

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

If σ is known.

$$\frac{\sigma}{\sqrt{n}} = \text{standard error} = \sigma_{\bar{x}}$$

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

If σ is not known. Here, s is the standard deviation of the sample.

Method 1: Hypothesis testing using rejection region (using Critical value)

Procedure for Hypothesis Testing

Step 1: State the Null (H_0) and Alternative (H_1) Hypotheses

Step 2: Decide the nature of test (one-tailed or two-tailed based on H_1)

Step 3: Obtain z_α value which depends upon α value and the nature of test.

Step 4: Choose appropriate formula and calculate test statistic, that is, z -value.

Step 5: Comparison and Conclusion.

- If $|z| < |z_\alpha|$, H_0 is accepted or H_1 is rejected, that is, there is no significant difference at $\alpha\%$ LOS.
- If $|z| > |z_\alpha|$, H_0 is rejected or H_1 is accepted,, that is, there is significant difference at $\alpha\%$ LOS.

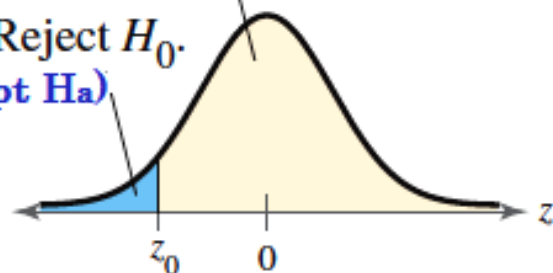
Explanation of last slide rule's

If the standardized
test statistic

1. is in the rejection region, then reject H_0 .
2. is *not* in the rejection region, then fail to reject H_0 .

Fail to reject H_0 . (Accept H_0)

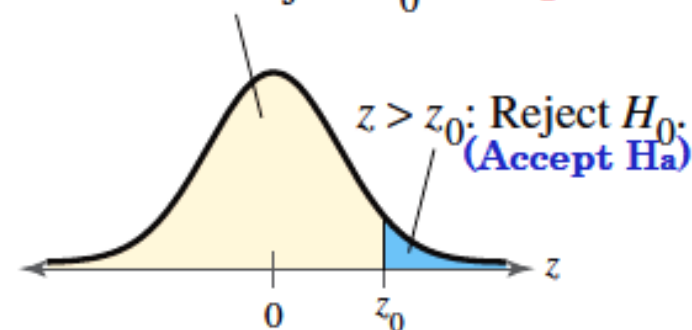
$z < z_0$: Reject H_0 .
(Accept H_a)



Left-Tailed Test

Fail to reject H_0 . (Accept H_0)

$z > z_0$: Reject H_0 .
(Accept H_a)



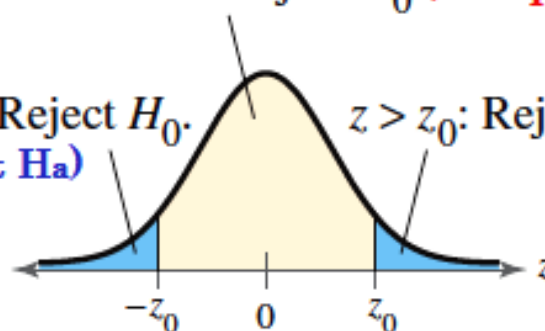
Right-Tailed Test

Two-Tailed Test

Fail to reject H_0 . (Accept H_0)

$z < -z_0$: Reject H_0 .
(Accept H_a)

$z > z_0$: Reject H_0 . (Accept H_a)



Critical Values for some standard LOS's

The critical values for some standard LOS's are given in the following table:

Table 1:

Type of Test	$\alpha = 1\%(0.01)$	$\alpha = 2\%(0.02)$	$\alpha = 5\%(0.05)$	$\alpha = 10\%(0.1)$
Two-Tailed	$ z_\alpha = 2.58$	$ z_\alpha = 2.33$	$ z_\alpha = 1.96$	$ z_\alpha = 1.645$
Right-Tailed	$z_\alpha = 2.33$	$z_\alpha = 2.055$	$z_\alpha = 1.645$	$z_\alpha = 1.28$
Left-Tailed	$z_\alpha = -2.33$	$z_\alpha = -2.055$	$z_\alpha = -1.645$	$z_\alpha = -1.28$

Example: A local telephone company **claims** that the average length of a phone call is 8 minutes. In a random sample of 58 phone calls, the sample mean was 7.8 minutes and the standard deviation was 0.5 minutes. Is there enough evidence to support this claim at $\alpha = 0.05$?

Solution: $\mu = 8$, $n = 58$, $\bar{x} = 7.8$, $\sigma = 0.5$, LOS $\alpha = 0.05$

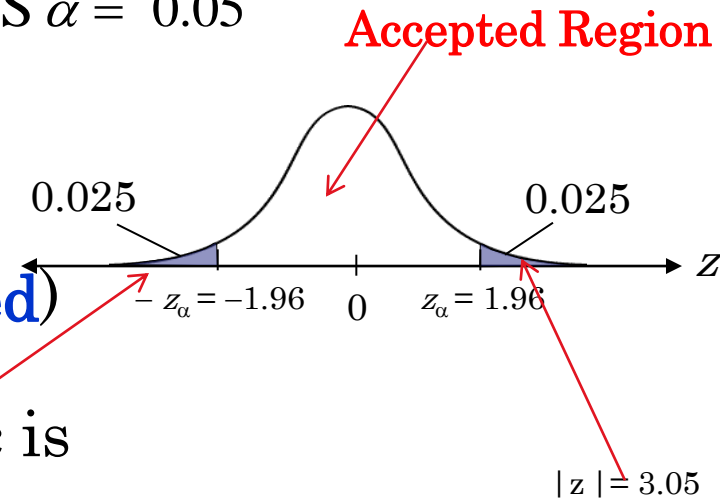
Step 1: $H_0: \mu = 8$ (Claim) $H_a: \mu \neq 8$

Step 2: Two tailed test

Step 3: For $\alpha = 0.05$, $|z_\alpha| = 1.96$ (2 tailed)

Step 4: The standardized test statistic is

$$z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} = \frac{7.8 - 8}{0.5 / \sqrt{58}} \approx -3.05.$$



Step 5: Here $|z_\alpha| < |z|$, value so f z is coming into the rejected region, so the test statistic falls in the rejection region, so H_0 is rejected

At the 5% level of significance, there is enough evidence to reject the claim that the average length of a phone call is 8

Problem 2: A sample of 400 items is taken from a population whose standard deviation is 10. The mean of the sample is 40. Test whether the sample has come from the population with mean 38. Also calculate 95% confidence interval for the population mean.

Solution:

$$H_0: \mu = 38 \quad \text{(Claim)} \quad H_1: \mu \neq 38$$

Level of significance: 5%

We apply the two tailed test.

$$\text{Test statistic is } Z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} = \frac{40 - 38}{10 / \sqrt{400}} = 4.$$

$$\therefore |Z| = 4$$

Table value of Z at 5% level of significance = 1.96. Since calculated value of Z at 5% level of significance is greater than the table value of Z , we reject H_0 at 5% level of significance.

95% confidence interval for the population mean is given by

$$\bar{x} \pm z_{\alpha} \frac{\sigma}{\sqrt{n}} = 40 \pm \left(1.96 \times \frac{10}{\sqrt{400}} \right) = [39.02, 40.98]$$

Problem

A sample of 100 students is taken from a large population. The mean height of the students in this sample is 160 cm. Can it be reasonably regarded that, in the population, the mean height is 165 cm and S.D. is 10 cm? LOS is 1%.

Solution:

Here, $n = 100$, $\bar{x} = 160$, $\mu = 165$ and $\sigma = 10$.

Step 1: $H_0 : \bar{x} = \mu$ (Claim) and $H_1 : \bar{x} \neq \mu$.

Step 2: In this case, we will use two-tailed test based on H_1 .

Step 3: Type of test and LOS value imply that $|z_\alpha| = 2.58$.

Step 4: The test statistic:

$$z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{160 - 165}{10/\sqrt{100}} = 5.$$

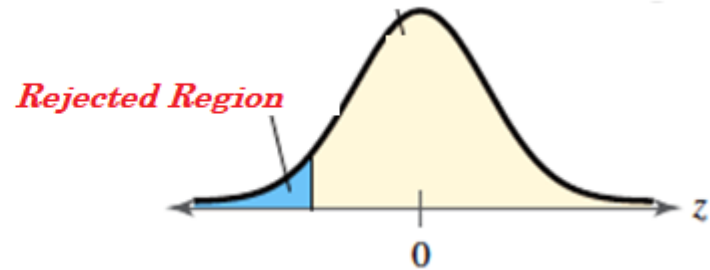
Step 5: Comparing the tabulated and calculated values of z , we have $|z| > |z_\alpha|$. That is, we reject H_0 at $\alpha = 0.01$.

Method 2: Hypothesis testing using P -values

1. State the null H_0 and Alternate H_a hypotheses.
2. Specify the level of significance.
3. Determine the standardized test statistic (z).
4. Find the area that corresponds to z .
5. Find the P -value, according to L, R or 2 tailed tests;
 - a. For a left-tailed test, $P = (\text{Area in left tail})$.
 - b. For a right-tailed test, $P = (\text{Area in right tail})$.
 - c. For a two-tailed test, $P = 2(\text{Area in tail of test statistic})$.
6. Compare the P -value with α for hypothesis test as;
 - a. If $P \leq \alpha$, then reject H_0 (acceptance of Alternate H_a)
 - b. If $P > \alpha$, then reject H_a (acceptance of Null H_0)

Calculation of P Value

1. For Right Tailed: Calculate z value, calculate the right tailed area (Blue area). For this first find area “A” by normal distribution table, then as we know left side curve area is 0.5, then calculate P value as **$P = 0.5 - A$** .



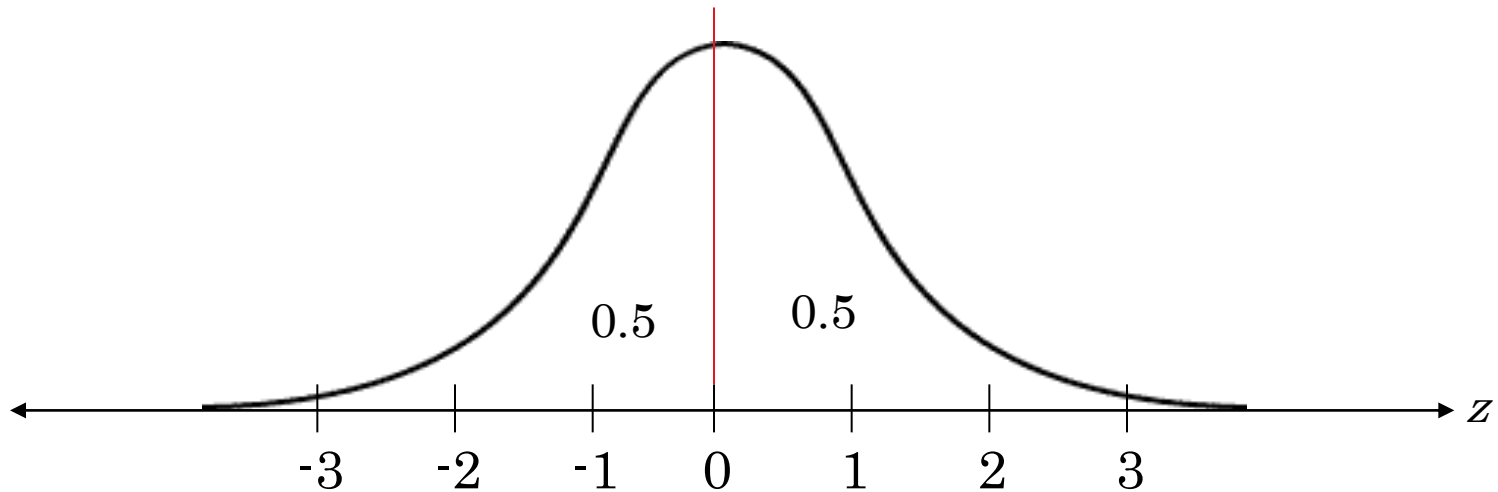
2. For Left Tailed: Calculate z value, calculate the left tailed area (Blue area). For this first find area “A” by normal distribution table, then as we know left side curve area is 0.5, then calculate P value as **$P = 0.5 - A$** .



3. For Right Tailed: Calculate z value, calculate the right and left both tailed area. For this first find area “A” by normal distribution table, then as we know left side curve area is 0.5, then calculate P value as **$P = 2(0.5 - A)$** .

Very Important Note:

At $z = \pm 4$ or $|z| < 4$ area under the curve is 0. and 0.5 (both sides).
It means area $z > +4$ and $z < -4$ are zero.

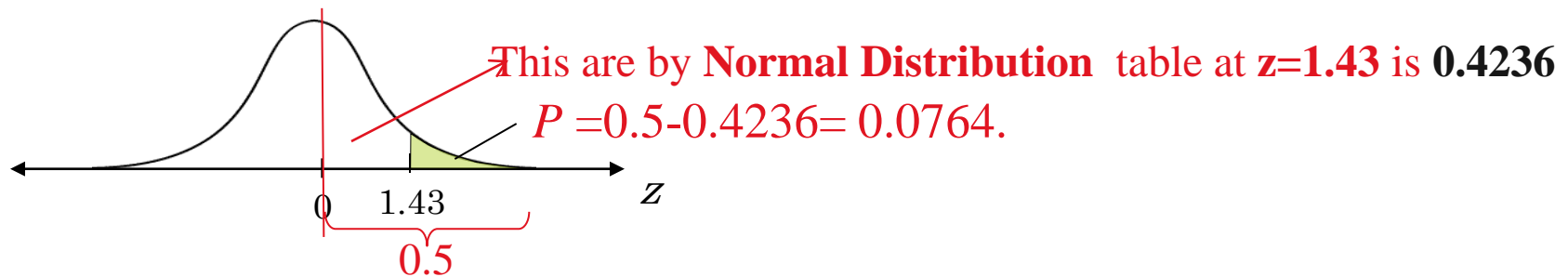


Example: A manufacturer claims that its rechargeable batteries are good for an average of more than 1,000 charges. A random sample of 100 batteries has a mean life of 1002 charges and a standard deviation of 14. Is there enough evidence to support this claim at $\alpha = 0.01$?

Soln: $H_0: \mu \leq 1000$ $H_a: \mu > 1000$ (Claim)

The level of significance is $\alpha = 0.01$.

The standardized test statistic is $z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{1002 - 1000}{14/\sqrt{100}} \approx 1.43$



Here we got that $P > \alpha$, therefore, reject H_a . (Accept H_0), so claim is rejected

At the 1% level of significance, there is not enough evidence to support the claim that the rechargeable battery has an average life of at least 1000 charges.

Problem 2: A sample of 400 items is taken from a population whose standard deviation is 10. The mean of the sample is 40. Test whether the sample has come from the population with mean 38. Also calculate 95% confidence interval for the population mean.

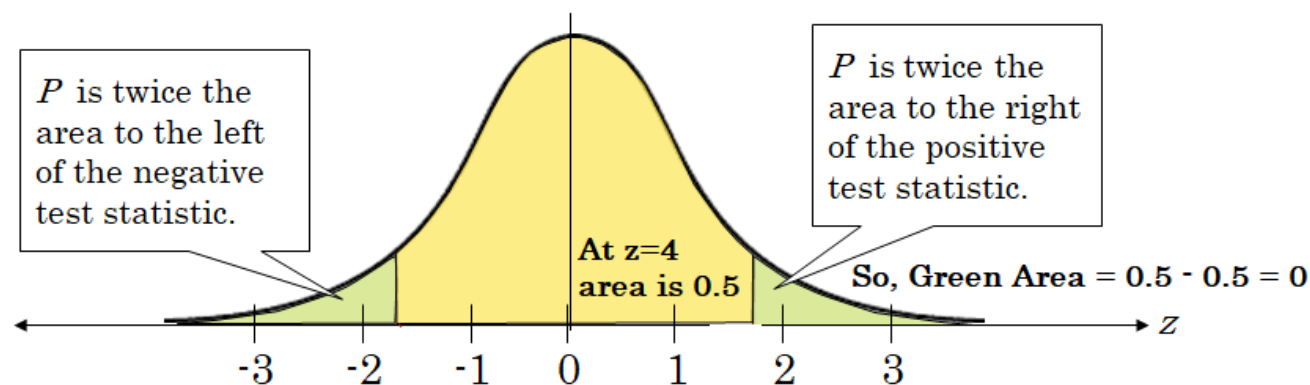
Solution:

$$H_0: \mu = 38 \quad \text{(Claim)} \quad H_1: \mu \neq 38$$

Level of significance: 5%

We apply the two tailed test.

$$\text{Test statistic is } Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{40 - 38}{10/\sqrt{400}} = 4. \quad \therefore |Z| = 4$$



Therefore, for two tailed test P value = $2 * 0 = 0$

Level of confidence is given 0.95.

So, Level of sign. is 0.05

Finally, we have $P < \alpha$

we reject H_0 at 5% level of significance.

Claim is rejected

95% confidence interval for the population mean is given by

$$\bar{x} \pm z_{\alpha} \frac{\sigma}{\sqrt{n}} = 40 \pm \left(1.96 \times \frac{10}{\sqrt{400}} \right) = [39.02, 40.98]$$

Type -2 Problem: For difference of means of two large samples

Let \bar{x}_1 be the mean of an independent random sample of size n_1 from a population with mean μ_1 and variance σ_1^2 . Again, let \bar{x}_2 be mean of an independent random sample of size n_2 from a population with mean μ_2 and variance σ_2^2 . Here, n_1 and n_2 are large.

Clearly, $\bar{x}_1 \sim N\left(\mu_1, \frac{\sigma_1^2}{n_1}\right)$ and $\bar{x}_2 \sim N\left(\mu_2, \frac{\sigma_2^2}{n_2}\right)$.

Hence, under the null hypothesis $H_0: \mu_1 = \mu_2$, the test statistic is

$$z = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}.$$

Cont...

🔷 If $\sigma_1 = \sigma_2 = \sigma$, then the test statistic is

$$Z = \frac{\bar{X}_1 - \bar{X}_2}{\sigma \left(\sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \right)}.$$

🔷 If σ_1 and σ_2 are not known, then the test statistic is

$$Z = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}}.$$

🔷 If $\sigma_1 = \sigma_2 = \sigma$ and σ is not known, we compute σ^2 by using the formula

$$\sigma^2 = \frac{n_1^2 s_1^2 + n_2^2 s_2^2}{n_1 + n_2}.$$

In this case, the test statistic is

$$Z = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{n_1^2 s_1^2 + n_2^2 s_2^2}{n_1 + n_2} \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}.$$

Problem: Random samples drawn from two places gave the following data relating to the heights of children

	Place A	Place B
Mean height	68.50	68.58
Standard deviation	2.5	3.0
Sample size	1200	1500

Test at 5% level that the mean height is the same for the children at two places.

Solution: $\bar{x}_1 = 68.50$, $\bar{x}_2 = 68.58$, $n_1 = 1200$, $n_2 = 1500$

Level of significance = 0.05 $\sigma_1 = 2.5$, $\sigma_2 = 3.0$

Null hypothesis: $H_0 : \mu_1 = \mu_2$ (Claim)

Alternative hypothesis: $H_1 : \mu_1 \neq \mu_2$ Two tailed test

$$Z = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} = -0.756.$$

$\therefore |z| = 0.756 < z_\alpha = 1.96$. Hence, null hypothesis is accepted.

So, claim is correct.

Type -3 Problem: For the difference between sample proportion and population proportion (single proportion)

Let X be the number of successes in n independent Bernoulli trials in which the probability of success for each trial is a constant $= P$ (say). Then it is known that X follows a binomial distribution with mean $E(X) = nP$ and variance $V(X) = nPQ$

When n is large, X follows $N(nP, \sqrt{nPQ})$, i.e. a normal distribution with mean nP and

$$\text{S.D. } \sqrt{nPQ}, \quad \text{where } Q = 1 - P. \quad \frac{X}{n} \text{ follows } N\left(\frac{nP}{n}, \sqrt{\frac{nPQ}{n^2}}\right)$$

Now $\frac{X}{n}$ is the proportion of successes in the sample consisting of n trials,

that is denoted by p . Thus the sample proportion p follows $N\left(P, \sqrt{\frac{PQ}{n}}\right)$.

$$\text{Test Statistic } z = \frac{p - P}{\sqrt{\frac{PQ}{n}}}, \quad \text{where } Q = 1 - P.$$

Note: When P is not known, the 95 percent confidence limits for P are given by

$$p - 1.96\sqrt{\frac{pq}{n}} \leq P \leq p + 1.96\sqrt{\frac{pq}{n}}$$

Ex. The CEO of a large electric utility claims that at least 80 percent of his 10,00,000 customers are very satisfied with the service they receive. To test this claim, the local newspaper surveyed 100 customers, using simple random sampling. Among the sampled customers, 73 percent say they are very satisfied. Based on these findings, can we reject the CEO's hypothesis that at least 80 percent of the customers are very satisfied? Use a 0.05 level of significance.

Son: In this problem, $P = 0.8$, $p = 0.73$ and $n = 100$.

Step 1: Null hypothesis (H_0) : $P \geq 0.80$, **(Claim)**
at least 80 percent customers are satisfied.

Alternative hypothesis (H_1) : $P < 0.80$

Step 2: Note that these hypotheses constitute a one-tailed (left-tailed) test.

Step 3: As it is given $\alpha = 0.05$ and one-tailed test, we have $z_\alpha = -1.645$.

Step 4: Test Statistic $z = \frac{p - P}{\sqrt{\frac{PQ}{n}}} = \frac{0.73 - 0.8}{\sqrt{\frac{0.8 \times 0.2}{100}}} = -1.75$.

Step 5: It can be viewed that $|z| = 1.75 > 1.645 = |z_\alpha|$. Therefore, we reject the null hypothesis H_0 , that is, H_1 is accepted.

\therefore The CEO's claim is wrong.

Some time p is not given the calculate it by formula $p=x/n$

Problem : If 20 people were attacked by a disease and only 18 survived, will you reject the hypothesis that the survival rate if attacked by this disease is 85% in favor of the hypothesis that is more at 5% level.

Solution: Number of people survived = $x = 18$. Size of the sample = $n = 20$.

$$p = \text{Proportion of the people survived} = \frac{x}{n} = \frac{18}{20} = 0.9$$

$$\text{It is given that } P = 85\% = 0.85. \quad Q = 1 - P = 1 - 0.85 = 0.15$$

Null hypothesis: $H_0: P = 0.85$

Level of significance = $\alpha = 0.05$

Alternative hypothesis: $H_1: P > 0.85$

$$\text{Test statistic: } z = \frac{p-P}{\sqrt{\frac{PQ}{n}}} = 0.6265.$$

Table value of $z_{\alpha} = 1.645$.

Calculated value of z is less than the table value of z_{α} at 5% level of significance. Null hypothesis is accepted.

Ex. Experience has shown that 20 percent of a manufactured product is of top quality. In one day's production 400 articles, only 50 are of top quality. Show that either the production of the day chosen was not a representative sample or the hypothesis 20 percent was strong.

Son. **Step 1:** Null Hypothesis (H_0) : $P = \frac{1}{5}$, that is, 20 percent of the products manufactured is of top quality.

Alternative Hypothesis (H_1) : $P \neq \frac{1}{5}$, that is, 20 percent of the products manufactured is not of top quality.

Step 2: Assume that $\alpha = 5\%$ and one can note that the type of test is two-tailed based on H_1 .

Step 3: Since it is two-tailed test and $\alpha = 5\%$, $z_\alpha = 1.96$.

Step 4: Test Statistic = $z = \frac{p - P}{\sqrt{\frac{PQ}{n}}} = \frac{\frac{1}{8} - \frac{1}{5}}{\sqrt{\frac{1}{5} \times \frac{4}{5} \times \frac{1}{400}}} = -3.75$.

Step 5: Now, $|z| = 3.75 > 1.96 = |z_\alpha|$ which implies that H_0 is rejected (or H_1 is accepted).

\therefore The production of the particular day chosen was not a representative sample.

Ex. A recent article in a weekly magazine reported that a job awaits 33% of new college graduates. A survey of 200 recent graduates from your college revealed that 80 students had jobs. At a 99% level of confidence, can we conclude that a larger proportion of students at your college have jobs?

Soln: **Step 1:** $H_0 : P = 0.33$ and $H_1 : P > 0.33$

Step 2: It is one-tailed (right-tailed) test.

Step 3: Here, $\alpha = 1\%$. So $z_\alpha = 2.33$.

Step 4:

$$z = \frac{p - P}{\sqrt{\frac{PQ}{n}}} = \frac{0.4 - 0.33}{\sqrt{0.33 \times 0.67 \times \frac{1}{200}}} = 2.1021.$$

Step 5: $|z| = 2.1021 < 2.33 = |z_\alpha|$ implies that H_0 is not rejected. Therefore, we do not have enough evidence to state that a larger proportion of students at our college have jobs.

Exercise

- ① A salesman in a departmental store claims that at most 60 percent of the shoppers entering the store leaves without making a purchase. A random sample of 50 shoppers showed that 35 percent of them left without making a purchase. Are these sample results consistent with the claim of the salesman at a level of significance of 0.05?
- ② A cubical die is thrown 900 times and a throw of three or four is observed 3240 times. Show that the die cannot be regarded as unbiased one and find the extreme limits between which the probability of a throw of three or four lies.

Type -4 Problem: For the difference of proportions

Suppose two samples of sizes n_1 and n_2 are drawn from two different populations. To test the significance of difference between the two proportions, we consider the following cases.

Case-I When the population proportions P_1 and P_2 are known:

In this case $Q_1 = 1 - P_1$ and $Q_2 = 1 - P_2$. The test statistic is

$$Z = \frac{P_1 - P_2}{\sqrt{\frac{P_1 Q_1}{n_1} + \frac{P_2 Q_2}{n_2}}}$$

Case-II When the population proportions P_1 and P_2 are not known but sample proportions p_1 and p_2 are known:

$$Z = \frac{p_1 - p_2}{\sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}}$$

Cont..

Case-III Method of pooling:

In this method, the sample proportions P_1 and P_2 are pooled into a single proportion P , by using the formula:

$$P = \frac{n_1 P_1 + n_2 P_2}{n_1 + n_2}$$

The test statistic in this case is

$$Z = \frac{P_1 - P_2}{\sqrt{PQ \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

Problem: A machine puts out 16 imperfect articles in sample of 500. After the machine is overhauled, it puts out 3 imperfect articles in a batch of 100. Has the machine improved?

Solution: Here we have $p_1 = \frac{16}{500} = 0.032$, $p_2 = \frac{3}{100} = 0.03$,
 $q_1 = 1 - p_1 = 0.968$ and $q_2 = 1 - p_2 = 0.970$.

Null hypothesis: $H : P_1 = P_2$

Alternative hypothesis: $H : P_1 \neq P_2$

Level of significance=0.05

The test statistic is $z = \frac{p_1 - p_2}{\sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}} = 0.107$

Table value of $z = 1.645$. Since the calculated value of z at 5% level of significance is less than the table value of z , we accept the null hypothesis.

Ex: Suppose the RK Drug Company develops a new drug, designed to prevent Covid19. The company states that the drug is more effective for women than for men. To test this claim, they choose a simple random sample of 100 women and 200 men from a population of 100,000 volunteers.

At the end of the study, 38% of the women caught a Covid19 and 51% of the men caught a Covid19. Based on these findings, can we conclude that the drug is more effective for women than for men? Use a 0.01 level of significance.

Soln: Assume that p_1 represents the effectiveness of drug on women and p_2 represents the effectiveness of drug on men. In this problem, $n_1 = 100$, $p_1 = 0.38$, $n_2 = 200$ and $p_2 = 0.51$.

Step 1: Null hypothesis (H_0): $p_1 \geq p_2$ **(Claim)**

Alternative hypothesis (H_1): $p_1 < p_2$.

Step 2: Note that these hypotheses constitute a left-tailed test.

Step 3: As we have left-tailed test and the significance level is 0.01, $z_\alpha = -2.33$.

Step 4: Test Statistic $(z) = \frac{(p_1 - p_2)}{\sqrt{PQ \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$, where $P = \frac{n_1 p_1 + n_2 p_2}{n_1 + n_2}$.

By simple calculation, one can note that $P = 0.467$ and $Q = 0.533$ which yields $z = -2.13$.

Step 5: Comparing z and z_α values, we obtain $|z| = 2.13 < 2.33 = |z_\alpha|$ which shows that H_0 is accepted. That is, the claim is true.

Exercise: (1) Consider the previous problem (ie, problem 4) with same data and test the claim that the drug is equally effective for men and women.

Ex: In a large city A, 20 percent of a random sample of 900 school days had a slight physical defect. In another large city B, 18.5 percent of a random sample of 1600 school boys had the same defect. Is the difference between the proportions significant? Use the level of significance 5%.

Soln: Given: $n_1 = 900, p_1 = 0.2, n_2 = 1600, p_2 = 0.185,$

Step 1: $H_0 : p_1 = p_2$

$H_1 : p_1 \neq p_2.$

Step 2: It is a two-tailed test.

Step 3: Since it is two-tailed and $\alpha = 0.05, z_\alpha = 1.96.$

Step 4: Test Statistic:

$$z = \frac{(p_1 - p_2)}{\sqrt{PQ \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}, \text{ where } P = \frac{n_1 p_1 + n_2 p_2}{n_1 + n_2}.$$

This implies that $z = \frac{(0.2 - 0.185)}{\sqrt{(0.1904)(0.8096) \left(\frac{1}{900} + \frac{1}{1600} \right)}} = 0.92.$

Step 5: Clearly, $|z| < |z_\alpha|$. Thus, H_0 is accepted, that is, H_1 is rejected.

Ex: 956 children were born in a city A in one year out of which 52.5% were male, while 1406 children were born in cities A and B both out of which proportion of male was 0.496. Is the difference in the proportion of male children in two cities significant?

Soln: Here, $n_1 = 956$, $p_1 = 52.5\%$, $n_2 = 1406 - 956 = 450$ and $P = 0.496 \Rightarrow Q = 1 - P = 0.504$.

We have,

$$P = \frac{n_1 p_1 + n_2 p_2}{n_1 + n_2} \Rightarrow 0.496 = \frac{(956)(0.525) + (450)(p_2)}{956 + 450} \Rightarrow p_2 = 0.432$$

Step 1: $H_0 : p_1 = p_2$ and $H_1 : p_1 \neq p_2$.

Step 2: Two-tailed Test will be used in this case.

Step 3: Let $\alpha = 5\%$, then $z_\alpha = 1.96$. as the type of the test is two-tailed.

Step 4: Test Statistic:

$$z = \frac{(p_1 - p_2)}{\sqrt{PQ \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}} = 3.268.$$

Step 5: Note that $|z| > |z_\alpha|$ which shows that H_1 is accepted.