

Six Sigma
Prof. Jithesh J Thakkar
Department of Industrial and Systems Engineering
Indian Institute of Technology, Kharagpur

Lecture - 30
Hypothesis Testing: Two Population Test

Hello friends, we are going through our Six Sigma journey. And typically we are discussing the analysed phase of DMAIC cycle. And we have already advanced in our discussion on hypothesis testing. We had seen the concepts; we had seen the hypothesis testing for single population. And this lecture 30; we will basically try to help to appreciate Hypothesis Testing for Two Population Test.

So, the concept will remain same and the general steps that we have follow in the hypothesis testing will remain same. But it would be done for two populations; I am interested to compare or to check the fact for two different population. And this could be let us say you are receiving the material from two different vendor, and you want to check that whether the percentage defective from vendor 1 vendor 2 are the same or they are different.

Or let us say you want to check some fact about the people who are less than 30 years of age, more than 30 years of age, or let us say you want to check some gender specific fact then male versus female. So, you have two different population of interest. And I want to test the hypothesis for the comparison of this two different population and this is called my hypothesis testing for two population test.

(Refer Slide Time: 01:53)

Science is advanced by proposing
and testing hypothesis, not by
declaring questions unsolvable.

— Nick Matzke —

AZ QUOTES

The slide features a yellow background with a dark blue curved shape on the right. On the left, there is a small inset photo of Nick Matzke, a man with glasses and a beard. The quote is centered in white text on a black rectangular background. At the bottom, there are logos for 'swayam' and 'AZ QUOTES'.

So, once again you are reminded that scientific science is advanced by proposing and testing hypothesis, not by declaring questions unsolvable.

(Refer Slide Time: 02:07)


Recap

- ❑ Hypothesis testing for a Single Population Mean using
 - ❖ z statistic
 - ❖ t statistic
 - ❖ For Proportion

The slide has a yellow background with a dark blue curved shape on the right. The title 'Recap' is in bold black font. Below it, there is a list of topics for hypothesis testing, each preceded by a red square icon. The list includes 'z statistic', 't statistic', and 'For Proportion'. At the bottom, there are logos for 'swayam' and 'AZ QUOTES'.

So, we had seen the z statistic, t statistic, and for proportion in the previous lecture.

(Refer Slide Time: 02:15)



CONCEPTS COVERED

Concepts Covered:

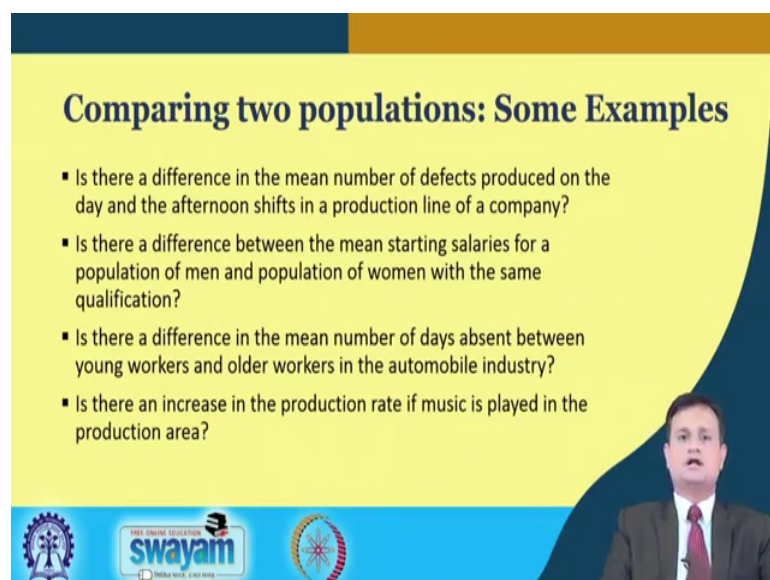
How to use hypothesis testing for comparing the difference between

- ☐ The means of two independent populations
- ☐ The proportions of two independent populations
- ☐ The variances of two independent populations by testing the ratio of the two variances
- ☐ Dependent (Related) Samples (Same group before vs. after treatment)

The slide features a dark blue background on the left with the title 'CONCEPTS COVERED' in yellow. The right side has a yellow background with the text 'Concepts Covered:' and 'How to use hypothesis testing for comparing the difference between'. Below this is a bulleted list of four topics, each preceded by a red square icon. At the bottom right, there is a small video inset of a man in a suit. The bottom of the slide contains logos for 'swayam' and 'INDIAN INSTITUTE OF TECHNOLOGY'.

Now, this lecture we would basically like to focus on two population test where we have four different cases. Case number 1: the means of two independent populations I want to compare. Number 2- the proportion of two independent populations. Number 3- the variance of two independent populations by testing the ratio of two variances. And number 4- dependent related samples same group before versus after treatment.

(Refer Slide Time: 02:51)



Comparing two populations: Some Examples

- Is there a difference in the mean number of defects produced on the day and the afternoon shifts in a production line of a company?
- Is there a difference between the mean starting salaries for a population of men and population of women with the same qualification?
- Is there a difference in the mean number of days absent between young workers and older workers in the automobile industry?
- Is there an increase in the production rate if music is played in the production area?

The slide has a yellow background with a dark blue header and footer. The title 'Comparing two populations: Some Examples' is in bold black text. Below the title is a bulleted list of four questions, each preceded by a red square icon. At the bottom right, there is a small video inset of a man in a suit. The bottom of the slide contains logos for 'swayam' and 'INDIAN INSTITUTE OF TECHNOLOGY'.

So, just try to create your interest in this particular topic. Just see that what could be the questions that may prompt you to go for such kind of hypothesis testing for two

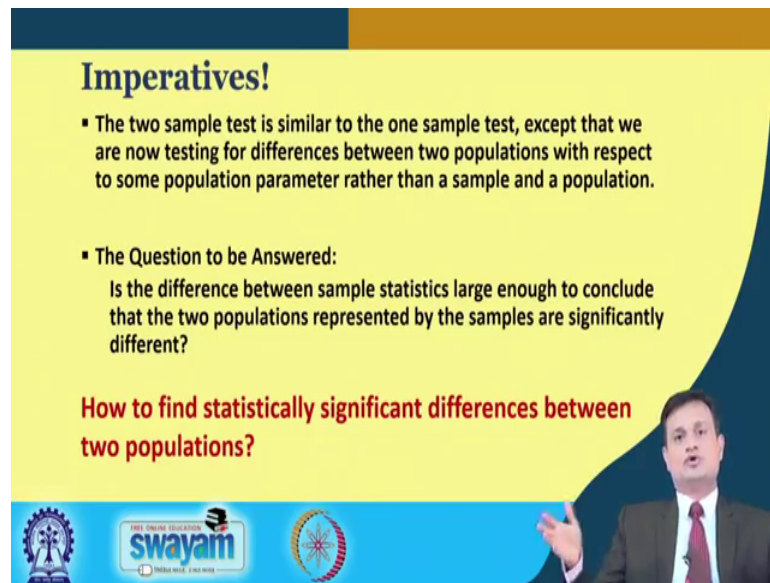
population. Say is there a difference in mean number of defects produced on the day and the afternoon shift in a production line of a company. You want to investigate; your production line is same manufacturing setup is same, but there could be influence of workers skill or some other parameter issues.

And is there really a difference so for the defects are concerned in the morning shift and in the afternoon shift. Is there a difference between the mean starting salaries for a population of man and population of women with the same qualification? You want to check is there a gender bias? Or is there a preference for one particular gender and is there a difference in the average salary. Then is there a difference in mean number of days absent between young workers and the older workers in automobile industry it is also interesting.

That you are the young worker and older worker when you just go by the number you will say fine young worker they were say 27 days absent on an average and the old worker let us say 35 days. But this is just the number descriptive statistics mean I am not making any inferential analysis this could be statistically different this could be same. So, I want to check is there really a statistical evidence available that; average number of absenteeism young worker and the old worker there is a difference.

Is there an increase in the production rate if music is played in the production area very interesting can that have a soothing effect a positive effect on the behaviour skill of the worker and that can help them to improve the say production or the productivity. So, these are some of the interesting questions that you would like to investigate through two population test.

(Refer Slide Time: 05:09)



Imperatives!

- The two sample test is similar to the one sample test, except that we are now testing for differences between two populations with respect to some population parameter rather than a sample and a population.
- The Question to be Answered:
Is the difference between sample statistics large enough to conclude that the two populations represented by the samples are significantly different?

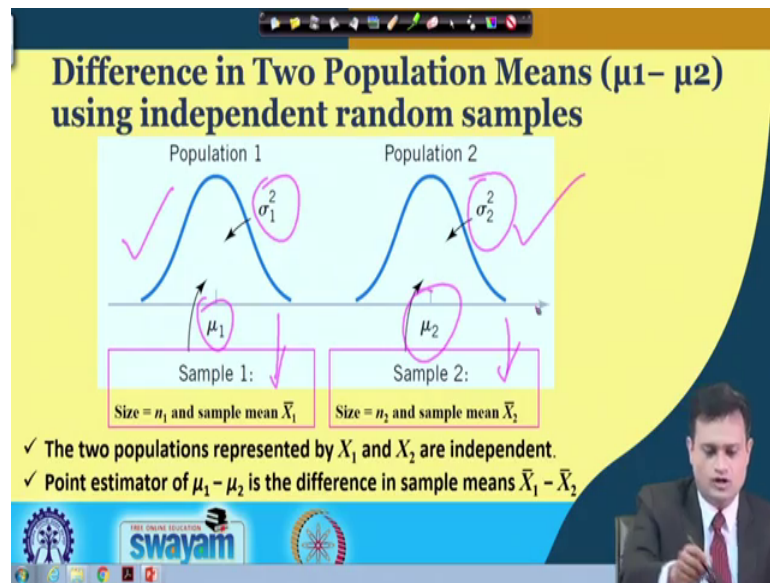
How to find statistically significant differences between two populations?

The slide features a yellow background with a dark blue curved border on the right. At the bottom, there is a blue banner with logos for 'swayam' and 'INDIA WISE, I AND WISE'. A small inset image of a man in a suit is visible in the bottom right corner of the slide.

So, the imperatives are there two sample test is similar to one sample test. As I said there would not be any difference in the procedure and hence I would not spend too much time on explaining the critical value p value concept. I would more focus on the real life application and the interpretation part in this lecture.

So, the question to be answered; is there a difference between sample statistics large enough to conclude that two populations represented by the samples are significantly different. So, we have already read some of the interesting questions that can really help you to appreciate the importance of hypothesis testing for two population.

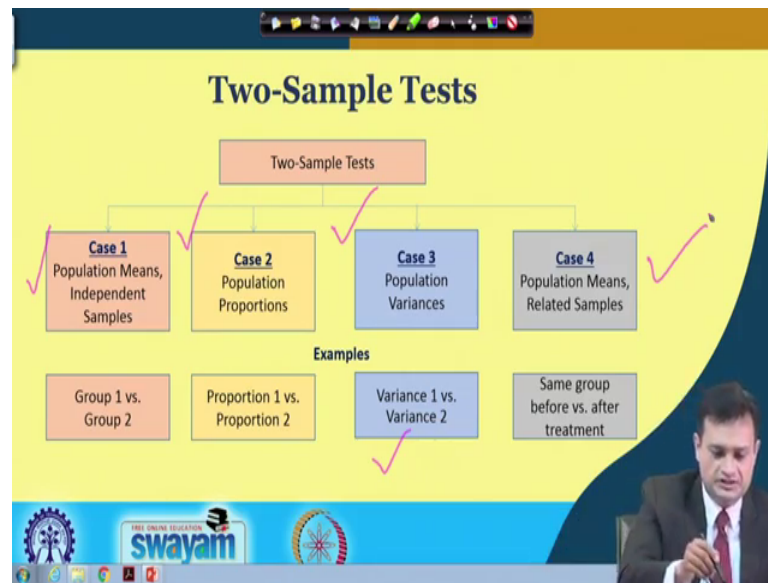
(Refer Slide Time: 06:01)



Just see this to get a better idea, you have basically difference in two population mean μ_1 minus μ_2 . And using independent random sample let us say; we have population 1 and population 2. You have σ_1^2 for population 1 σ_2^2 for population 2. This is my μ_1 , this is my μ_2 and I am drawing this sample 1 for population 1 sample 2 for population 2.

So, two populations are represented by X_1 X_2 which are independent and μ_1 minus μ_2 is the difference in the mean \bar{X}_1 and \bar{X}_2 . So, I want to check compare my claim for two different populations which are independent. And this is what exactly we would do in this particular lecture.

(Refer Slide Time: 06:53)



Let me give you the frame work of this particular lecture and what are the different cases which I would like to discuss as a part of hypothesis testing for two population test. So, case 1 we have population means and independent samples group 1 versus group 2. Suppose average performance of the boys in the school may be higher secondary and average performance of the girls in sports in the higher secondary is there a really significance difference or they are same?

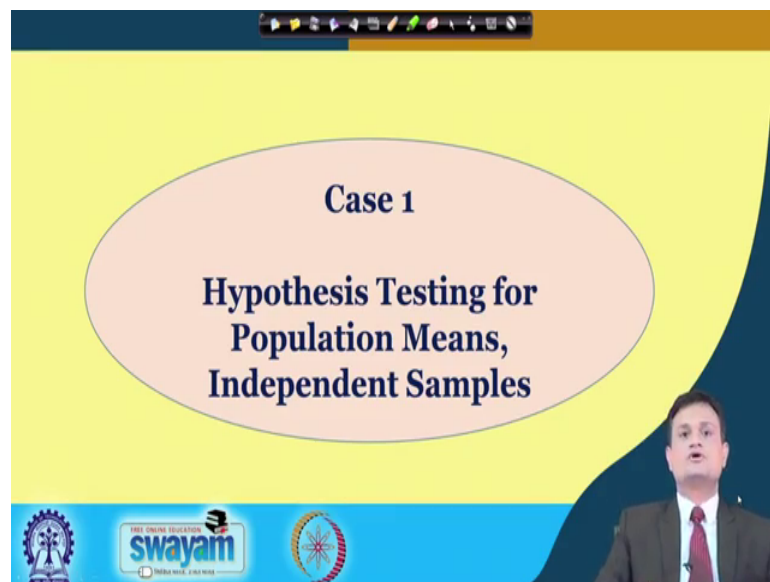
We have seen average salary for male average salary for female production in the morning shift, production in the afternoon shift. So, this would fall typically in the first case that is population means independent samples. Now just see the case 2 population proportion as I have explained previously number of success and what is that proportion that something is happening some number of time and I would like to compare proportion 1 with proportion 2 and not the two means as we are doing in the case 1. Now you have case 3.

So, many a times the variance is of interest and I want to check the variability present in the population 1 and the variability in the population 2. So, I want to compare variance one versus variance 2 by taking an appropriate sample and this is my case 3. When I say case 4 population means related samples same group before or after treatment. Now this would be really interesting suppose you are implementing a six sigma program for a particular function or process.

Now you want to see that is there really a reduction in the defect rate after implementing a six sigma program. So, you have the before data you have the after data and then you would like to check it. Similar way you can think that there are patients suffering from a particular viral disease they are given the antibiotics and you want to see the improvement in their symptom.

So, you have the same group before and after treatment and you are interested to compare this two. And here these two before and after are treated as two different population. So, with this broad classification let us try to discuss each particular case with some real life example.

(Refer Slide Time: 09:39)



Case 1

**Hypothesis Testing for
Population Means,
Independent Samples**

swayam

So, case 1- hypothesis testing for population means independent samples.

(Refer Slide Time: 09:45)

Difference in Two Population Means ($\mu_1 - \mu_2$) using independent random samples

$$E(\bar{X}_1 - \bar{X}_2) = E(\bar{X}_1) - E(\bar{X}_2) = \mu_1 - \mu_2$$

If two samples are independent

$$V(\bar{X}_1 - \bar{X}_2) = V(\bar{X}_1) + (-1)^2 V(\bar{X}_2) = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}$$

standard error = $\sigma_{\bar{X}_1 - \bar{X}_2} = \sqrt{\sigma_{\bar{X}_1}^2 + \sigma_{\bar{X}_2}^2} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$

So, just I would like to present little bit mathematics. So, you have expectation E stands for expectation $\bar{X}_1 - \bar{X}_2$ is expressed as expectation of \bar{X}_1 minus expectation of \bar{X}_2 . And in statistics we say that expectation of \bar{X}_1 is μ_1 . Means expectation of \bar{x} as per central limit theorem is close to μ . So, I have μ_1 minus μ_2 as expectation of $\bar{X}_1 - \bar{X}_2$ so this is $\mu_1 - \mu_2$.

Now if two samples are independent; variance of $\bar{X}_1 - \bar{X}_2$ would be expressed as; variance of \bar{X}_1 plus $(-1)^2$ variance of \bar{X}_2 σ_1^2 square divided by n_1 σ_2^2 square divided by n_2 as per central limit theorem. Typically suppose you have two variable then variance of $y_1 - y_2$ is typically expressed as variance of y_1 plus variance of y_2 minus co variance of y_1 into y_2 . Because you have these two events y_1 and y_2 independent my covariance is basically 0.

So, it is expected that you have some basic knowledge of statistics, you should revise the suggested book. And on this logic I have variance of $\bar{X}_1 - \bar{X}_2$ is equal to σ_1^2 square divided by n_1 plus σ_2^2 square divided by n_2 . And typically you can express your standard error as $\sigma_{\bar{X}_1 - \bar{X}_2}$ which is equal to square root of $\sigma_{\bar{X}_1}^2 + \sigma_{\bar{X}_2}^2$. And this is basically the square root of σ_1^2 square divided n_1 plus σ_2^2 square divided by n_2 . So, this is what we say.

(Refer Slide Time: 11:57)

Difference in Means, Variances Known

- Sampling Distribution of the variable $\bar{X}_1 - \bar{X}_2$ follows normal distribution, i.e.,
- When $(\bar{X}_1 - \bar{X}_2) \sim N\left(\mu_1 - \mu_2, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}\right)$
 - Both populations are normal
- or
- Both $n_1 = 30$ and $n_2 = 30$ (Central Limit Theorem).
- This result will be used for tests of hypotheses and confidence intervals on $\mu_1 - \mu_2$, when σ_1 and σ_2 are known.

The statistic

$$Z = \frac{\bar{X}_1 - \bar{X}_2 - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim N(0,1)$$

Handwritten notes on the slide include: $Z = \frac{\bar{X} - \mu}{\sigma}$ and a normal distribution curve with mean $\mu_1 - \mu_2$ and standard deviation $\sigma_{\bar{X}_1 - \bar{X}_2}$.

Now, you have two populations. So, I want to check the difference in mean and suppose variances are known. So, variances are known means I am aware about the population variance and I want to say that my sampling distribution of the variable $\bar{X}_1 - \bar{X}_2$ follows the normal distribution. So, please understand that a distribution typically represents the behaviour of a particular random variable studied over a period of time.

And here I measure the random variable as $\bar{X}_1 - \bar{X}_2$. So, this is my sampling distribution for $\bar{X}_1 - \bar{X}_2$. And that is why here it is expressed as $\mu_1 - \mu_2$ that is population mean. And $\sigma_{\bar{X}_1 - \bar{X}_2} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$. So, both population are normal both n_1 is equal to 30 n_2 is equal to 30 central limit theorem can be applied.

And your typically standard normal z statistics can be computed as basically z is equal to $\bar{x} - \mu$ divided by σ . Same thing is basically transcribed here as $\bar{X}_1 - \bar{X}_2$. Because I am considering $\bar{X}_1 - \bar{X}_2$ as my random variable so $\bar{X}_1 - \bar{X}_2 - \mu_1 + \mu_2$ equivalent to $\bar{X}_1 - \bar{X}_2 - \mu_1 + \mu_2$ divided by my standard deviation. So, this approaches to standard normal. And this is my test statistics.

(Refer Slide Time: 13:37)

Confidence Interval on a Difference in Means
[Variances (σ_1^2 and σ_2^2) Known]

- The 100(1 - α)% confidence interval on the difference in two means $\mu_1 - \mu_2$
 $P(-z_{\alpha/2} < z < z_{\alpha/2}) = 1 - \alpha$
or

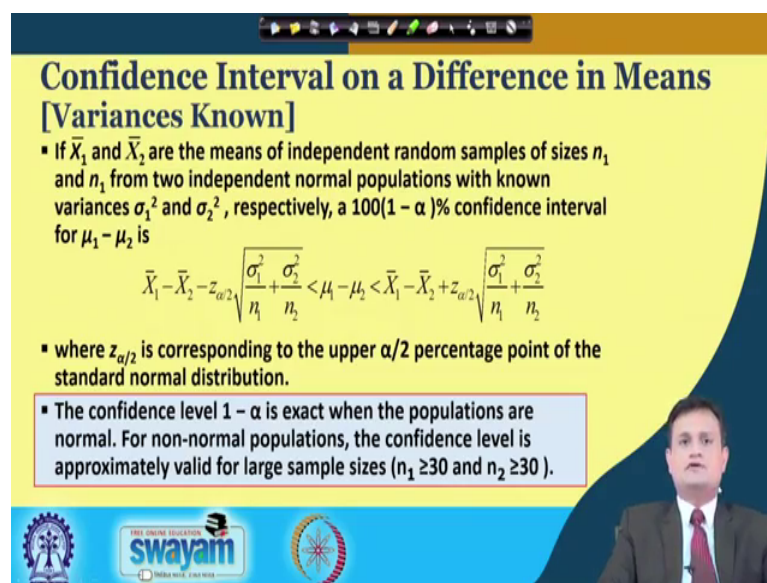
$$P\left(-z_{\alpha/2} < \frac{\bar{X}_1 - \bar{X}_2 - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} < z_{\alpha/2}\right) = 1 - \alpha$$
- This can be rearranged as

$$P\left(\bar{X}_1 - \bar{X}_2 - z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} < \mu_1 - \mu_2 < \bar{X}_1 - \bar{X}_2 + z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}\right) = 1 - \alpha$$

So, now I can also create the confidence interval for plus or minus z sigma or k sigma. And you can see here that I want to find the probability which is probability of minus z alpha by 2 less than z less than z alpha by 2 it should be 1 minus alpha. So, I want to typically find a region in this particular normal distribution. So, I have this alpha by 2 I have this alpha by 2 I am interested to find 1 minus alpha region this particular region. And say probability when I convert this into standard normal you can see this is what we had.

So, this particular quantity is my standard normal z statistics minus z alpha by 2, z alpha by 2 is equal to 1 minus alpha. So, you can just do little bit iteration to convert this expression simply having mu 1 minus mu 2 in the centre. And you will find that this is the expression after rearranging gives me the probability that what is that particular confidence interval which gives me 1 or minus alpha region. So, this is what we can do as confidence interval.

(Refer Slide Time: 15:01)



Confidence Interval on a Difference in Means [Variances Known]

- If \bar{X}_1 and \bar{X}_2 are the means of independent random samples of sizes n_1 and n_2 from two independent normal populations with known variances σ_1^2 and σ_2^2 , respectively, a $100(1 - \alpha)\%$ confidence interval for $\mu_1 - \mu_2$ is

$$\bar{X}_1 - \bar{X}_2 - z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} < \mu_1 - \mu_2 < \bar{X}_1 - \bar{X}_2 + z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

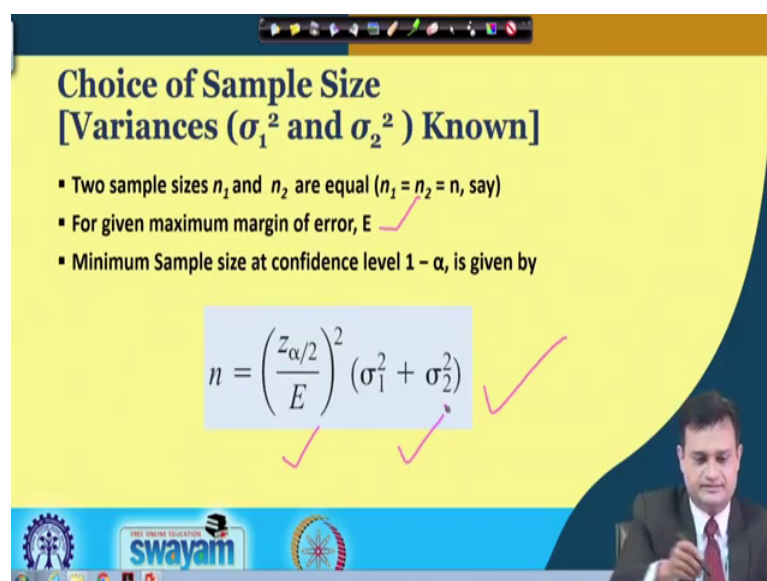
- where $z_{\alpha/2}$ is corresponding to the upper $\alpha/2$ percentage point of the standard normal distribution.
- The confidence level $1 - \alpha$ is exact when the populations are normal. For non-normal populations, the confidence level is approximately valid for large sample sizes ($n_1 \geq 30$ and $n_2 \geq 30$).

swayam

Now, you can have confidence interval in mean for variances known. Now here I assume that my variances are known about the population and if \bar{X}_1 and \bar{X}_2 are the means of independent random sample of size n_1 and typically n_1 from two independent normal population with known variance.

Then $100(1 - \alpha)\%$ confidence interval for $\mu_1 - \mu_2$ can be constructed like this. So, you only need to follow the same procedure in order to get the confidence interval for the given situation.

(Refer Slide Time: 15:51)



Choice of Sample Size [Variances (σ_1^2 and σ_2^2) Known]

- Two sample sizes n_1 and n_2 are equal ($n_1 = n_2 = n$, say)
- For given maximum margin of error, E
- Minimum Sample size at confidence level $1 - \alpha$, is given by

$$n = \left(\frac{z_{\alpha/2}}{E} \right)^2 (\sigma_1^2 + \sigma_2^2)$$

swayam

Now, you have the sample size to be determined and for choice of sample size. Let us say if I am aware about the variance σ_1^2 and σ_2^2 of two different population. And two sample n_1 and n_2 let us say are equal and equal to n_1 is equal to n_2 is equal to n . So, for given maximum margin of error E this can be expressed as n is equal to $z_{\alpha/2}$ divided by E square plus σ_1^2 plus σ_2^2 .

We are not going into the details of derivation part. We accept some of the expressions as it is in the standard form just in order to facilitate our computation and the investigation for a typical real life problem through hypothesis testing.

(Refer Slide Time: 16:49)

Confidence Interval on a Difference in Means
[Variances (σ_1^2 and σ_2^2) Known] large samples

- Both sample sizes are ≥ 30
- Population standard deviations are unknown
- Use sample variances to estimate population variances : Normal approximation remains valid
- The statistic is z value.

$$Z = \frac{\bar{X}_1 - \bar{X}_2 - (\mu_1 - \mu_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

➤ A $100(1 - \alpha)\%$ confidence interval for $\mu_1 - \mu_2$ is:

$$(\bar{X}_1 - \bar{X}_2) \pm z_{\alpha/2} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

So, you have let us say confidence interval on a difference in means variance σ_1^2 and σ_2^2 are known in your large samples. So, as we have seen that your typical say expression would follow $\bar{X}_1 - \bar{X}_2 - \mu_1 - \mu_2$ divided by s_1^2 divided by n_1 plus s_2^2 square divided by n_2 .

So, you have this sample standard sample variance for population 1 and as well as for population 2. And you can create the $100(1 - \alpha)\%$ confidence for $\mu_1 - \mu_2$. So, we have different cases where we tried to.

(Refer Slide Time: 17:37)

**Confidence Interval on a Difference in Means:
Variances unknown, small samples**

- Populations are normally distributed
- The population variances are assumed equal ($\sigma_1^2 = \sigma_2^2 = \sigma^2$), so use the two sample variances and pool them to estimate σ^2 .
- the test statistic is a t value with $(n_1 + n_2 - 2)$ degrees of freedom.

$$t = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{s_p \sqrt{\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$$

The pooled estimator of σ^2 , denoted by s_p^2 , is defined by

$$s_p^2 = \frac{\sum_{i=1}^{n_1} (X_{1i} - \bar{X}_1)^2 + \sum_{i=1}^{n_2} (X_{2i} - \bar{X}_2)^2}{(n_1 - 1) + (n_2 - 1)} = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$$

Say compute the z statistics t statistics and create the confidence interval. Now you want to create the confidence interval on let us say mean for a small sample. So, when I talk about small sample usually we refer to the student t distribution and your expression would become something like this. So, if this is your t statistics $\bar{X}_1 - \bar{X}_2 - \mu_1 - \mu_2$ divided by s_p which is called pooled variance.

So, s_p square root of $\frac{1}{n_1} + \frac{1}{n_2}$ a standard expression of t statistics we are accepting for the checking or testing of the hypothesis. And your pooled variance is given by this particular expression $\frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$. Because I am considering 2 sample so I have to consider the $n_1 + n_2 - 2$ as the denominator.

(Refer Slide Time: 18:53)

**Confidence Interval on a Difference in Means:
Variances unknown, small samples**

A $100(1 - \alpha)\%$ confidence interval for $\mu_1 - \mu_2$ is:

$$(\bar{X}_1 - \bar{X}_2) \pm t_{\alpha/2, n_1 + n_2 - 2} s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

Where $t_{\alpha/2}$ has $(n_1 + n_2 - 2)$ degrees of freedom, and

$$s_p = \sqrt{\frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}}$$

The pooled estimate of the variance is a weighted average of the two individual sample variances, with weights proportional to the sizes of the two samples. That is, larger weight is given to the variance from the larger sample.

swayam

So, you have $100(1 - \alpha)\%$ confidence interval for $\mu_1 - \mu_2$. And basically it is your t because we are referring the small sample case t distribution.

So, basically \bar{X} or μ plus or minus k sigma z sigma in case of t it is t sigma and pooled variance is given by this. So, this is the standard procedure we follow.

(Refer Slide Time: 19:21)

**Confidence Interval on a Difference in Means:
Variances unknown, small samples**

- Populations are normally distributed
- The population variances are assumed not equal ($\sigma_1^2 \neq \sigma_2^2$).
- the test statistic $t' = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$ is approximately distributed as t, with degrees of freedom given by

$$v = \frac{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)^2}{\frac{(s_1^2/n_1)^2}{n_1 - 1} + \frac{(s_2^2/n_2)^2}{n_2 - 1}}$$

called Smith-Satterthwaite degrees of freedom

- If v is not integer then round it to the nearest integer

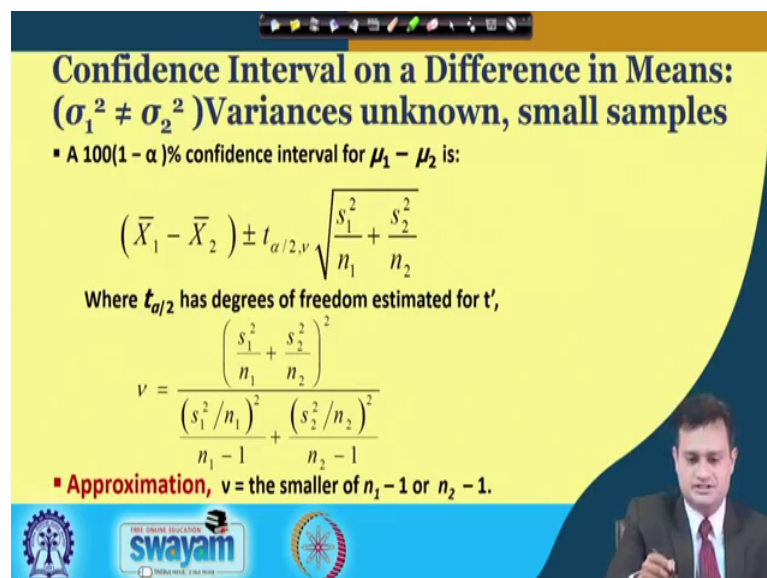
swayam

And you have another situation where; confidence interval on difference in means variances are unknown. So, so far we assumed that I am aware of the variance of the population. But many a times say particular phenomenon you are encountering first time

or little evidences is available little pass data is available. Then you are not very much aware of the population variance in this case you will little bit modify the expression.

And your test statistic will become t dash express like this and you will consider nu typically nu something called s 1 square divided by n 1 plus s 2 square divided by n 2. And you have denominator s 1 square divided by n 1 whole square divided by n 1 minus n 2 and this is also called Smith-Satterthwaite say degree of freedom. And this degree of freedom you will use to find the critical value or the p value from your say t table statistical table.

(Refer Slide Time: 20:35)



Confidence Interval on a Difference in Means:
 $(\sigma_1^2 \neq \sigma_2^2)$ Variances unknown, small samples

- A $100(1 - \alpha)\%$ confidence interval for $\mu_1 - \mu_2$ is:

$$(\bar{X}_1 - \bar{X}_2) \pm t_{\alpha/2, v} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

Where $t_{\alpha/2}$ has degrees of freedom estimated for t',

$$v = \frac{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2} \right)^2}{\frac{(s_1^2/n_1)^2}{n_1 - 1} + \frac{(s_2^2/n_2)^2}{n_2 - 1}}$$

- **Approximation**, $v =$ the smaller of $n_1 - 1$ or $n_2 - 1$.

So, we have this particular 100 1 minus alpha percent confidence for my t when variances are not known.

(Refer Slide Time: 20:47)

Comparisons of Two Population Means:
Hypothesis testing

- Let D_0 be the hypothesized difference between μ_1 and μ_2
- The three forms for a hypothesis test are as follows:

Left tailed test:	Right tailed test:	Two-tailed test:
$H_0: \mu_1 - \mu_2 \geq D_0$	$H_0: \mu_1 - \mu_2 \leq D_0$	$H_0: \mu_1 - \mu_2 = D_0$
$H_1: \mu_1 - \mu_2 < D_0$	$H_1: \mu_1 - \mu_2 > D_0$	$H_1: \mu_1 - \mu_2 \neq D_0$
- In many applications, $D_0 = 0$.

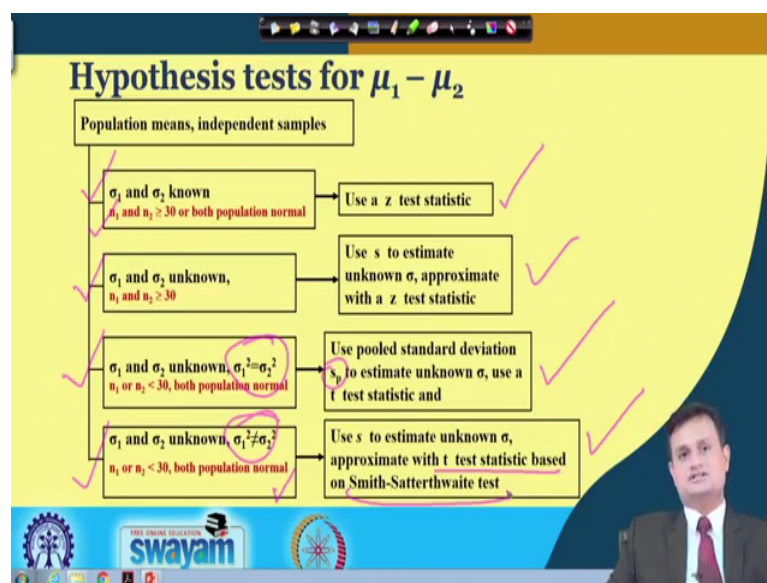
Left tailed test:	Right tailed test:	Two-tailed test:
$H_0: \mu_1 - \mu_2 \geq 0$	$H_0: \mu_1 - \mu_2 \leq 0$	$H_0: \mu_1 - \mu_2 = 0$
$H_1: \mu_1 - \mu_2 < 0$	$H_1: \mu_1 - \mu_2 > 0$	$H_1: \mu_1 - \mu_2 \neq 0$
i.e.,	i.e.,	i.e.,
$H_0: \mu_1 \geq \mu_2$	$H_0: \mu_1 \leq \mu_2$	$H_0: \mu_1 = \mu_2$
$H_1: \mu_1 < \mu_2$	$H_1: \mu_1 > \mu_2$	$H_1: \mu_1 \neq \mu_2$

Handwritten note: $D_0 = \mu_1 - \mu_2$

And I have the situation like this that let D_0 be hypothesised difference between μ_1 and μ_2 so typically $\mu_1 - \mu_2$. So, D_0 is $\mu_1 - \mu_2$ and I can have different forms of hypothesis to be tested. I can say that X_0 is $\mu_1 - \mu_2$ greater than or equal to D_0 or less than equal to D_0 .

I can have right tailed test I can have two tailed test I can have say many applications suppose D_0 is equal to 0 then; left tailed test right tailed test and two tailed test. So, you just need to appreciate that depending upon your interest right tailed test two tailed test you would like to set the hypothesis for checking your claim. So, we will see couple of examples; so the idea would be better clear.

(Refer Slide Time: 21:49)



Now, just see that what could be the different situations in which we can refer our case 1. We are with case 1 and we have population mean independent sample this is my case 1. So, I can have a situation where sigma 1 and sigma 2 known means population standard deviations are known and n 1 and n 2 greater than equal to 30. So, both the populations are normal. You compute z statistic the guiding rule is that you compute z statistic.

You have sigma 1 and sigma 2 unknown so n 1 and n 2 greater than equal to 30. Use s to estimate unknown sigma approximate with z test statistic; sigma 1 and sigma 2 unknown sigma 1 square is equal to sigma 2 square, n 1 or n 2 less than 30 both population be assumed to be normal. You go for t statistic because your sample size is less than 30 you have sigma 1 and sigma 2 unknown sigma 1 square is not equal to sigma 2 square n 1 or n 2 they less than 30 both population normal. You go for s estimate of unknown sigma and approximate with t statistic based on Smith-Satterthwaite test.

So, this are the variance things you can see in this last two case that here; this two are equal I would just like to draw your attention. Here this two are not equal so in one case we are using the pooled standard deviation in another case you go with the t test statistics based on Smith-Satterthwaite test. So, these are the standard guidelines to follow for testing the hypothesis for two population independent of each other

(Refer Slide Time: 23:45)

σ_1 and σ_2 known n_1 and $n_2 \geq 30$ or both population normal

Population means, independent samples

- σ_1 and σ_2 known
 n_1 and $n_2 \geq 30$ or both population normal**
- σ_1 and σ_2 unknown,
 n_1 and $n_2 \geq 30$
- σ_1 and σ_2 unknown, $\sigma_1^2 = \sigma_2^2$
 n_1 or $n_2 < 30$, both population normal
- σ_1 and σ_2 unknown, $\sigma_1^2 \neq \sigma_2^2$
 n_1 or $n_2 < 30$, both population normal

* The test statistic for $\mu_1 - \mu_2$ is:

$$z = \frac{(\bar{X}_1 - \bar{X}_2) - D_0}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

And the variance test statistics which I explained before I have just summarized here if you have a case 1; this is the test statistic which needs to be computed for finding the critical value, as well as finding the p value.

(Refer Slide Time: 24:23)

**σ_1 and σ_2 unknown, small samples
 $\sigma_1^2 = \sigma_2^2$, both population normal**

Population means, independent samples

- σ_1 and σ_2 known
 n_1 and $n_2 \geq 30$ or both population normal
- σ_1 and σ_2 unknown,
 n_1 and $n_2 \geq 30$
- σ_1 and σ_2 unknown, $\sigma_1^2 = \sigma_2^2$
 n_1 or $n_2 < 30$, both population normal**
- σ_1 and σ_2 unknown, $\sigma_1^2 \neq \sigma_2^2$
 n_1 or $n_2 < 30$, both population normal

The test statistic for $\mu_1 - \mu_2$ is:

$$t = \frac{(\bar{X}_1 - \bar{X}_2) - D_0}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

* Where t_a has $(n_1 + n_2 - 2)$ d.f., and

$$s_p = \sqrt{\frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}}$$

Now second one is this one so, you have the case sigma 1 sigma 2 known this is your z statistic you will make use of pooled variance. The third one; is the case where you have sigma 1 square is equal to sigma 2 square you will use this as the test statistic.

(Refer Slide Time: 24:39)

**σ_1 and σ_2 unknown, small samples
 $\sigma_1^2 \neq \sigma_2^2$, both population normal**

Population means, independent samples

- σ_1 and σ_2 known
 n_1 and $n_2 \geq 30$ or both population normal
- σ_1 and σ_2 unknown,
 n_1 and $n_2 \geq 30$
- σ_1 and σ_2 unknown, $\sigma_1^2 = \sigma_2^2$
 n_1 or $n_2 < 30$, both population normal
- σ_1 and σ_2 unknown, $\sigma_1^2 \neq \sigma_2^2$ *
 n_1 or $n_2 < 30$, both population normal

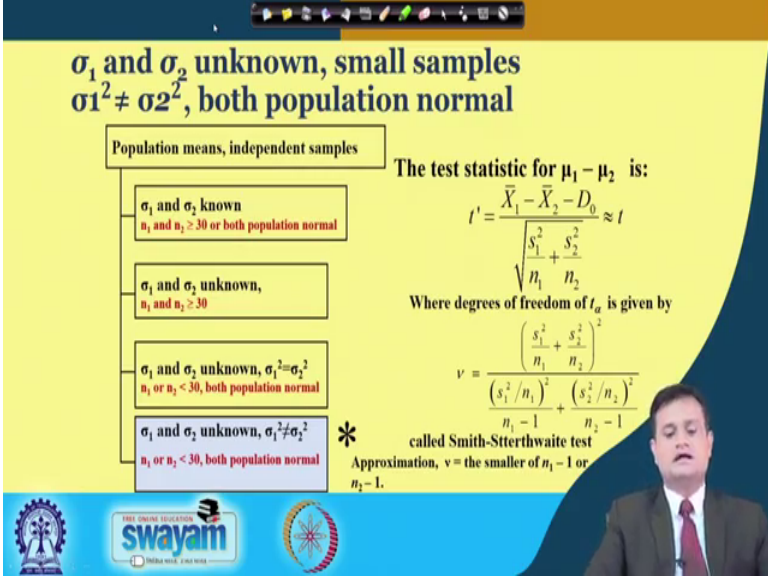
The test statistic for $\mu_1 - \mu_2$ is:

$$t' = \frac{\bar{X}_1 - \bar{X}_2 - D_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \approx t$$

Where degrees of freedom of t_d is given by

$$v = \frac{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)^2}{\frac{(s_1^2/n_1)^2}{n_1 - 1} + \frac{(s_2^2/n_2)^2}{n_2 - 1}}$$

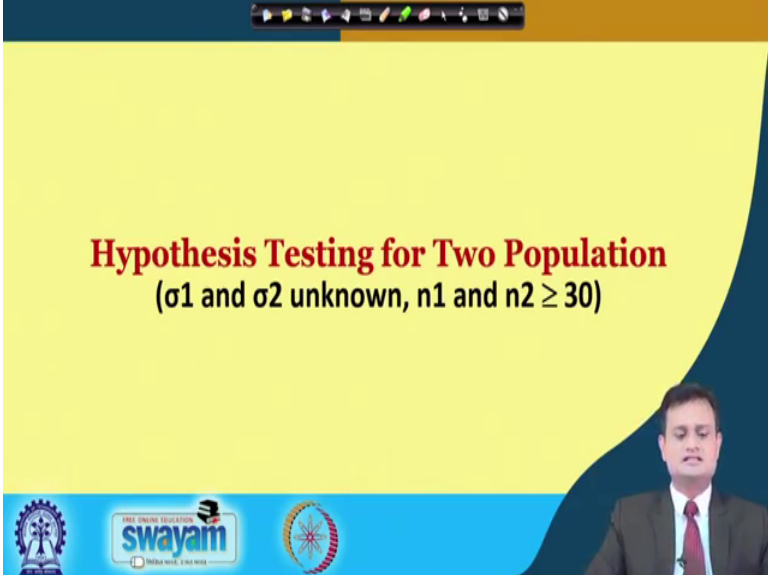
called Smith-Stterthwaite test
Approximation, v = the smaller of $n_1 - 1$ or $n_2 - 1$.



And similar way you have fourth case; where sigma 1 and sigma 1 square and sigma 2 square are basically not equal. And you will go by this Smith-Satterthwaite test to check your hypothesis to test your hypothesis.

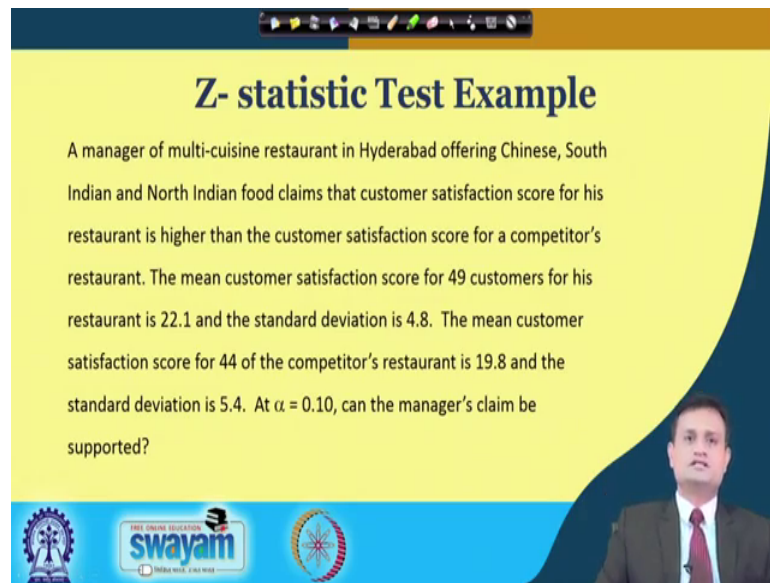
(Refer Slide Time: 24:55)

**Hypothesis Testing for Two Population
(σ_1 and σ_2 unknown, n_1 and $n_2 \geq 30$)**



So, we have some examples to discuss hypothesis testing for two population sigma 1 and sigma 2 unknown n 1 and n 2 greater than or equal to 30.

(Refer Slide Time: 25:05)



Z- statistic Test Example

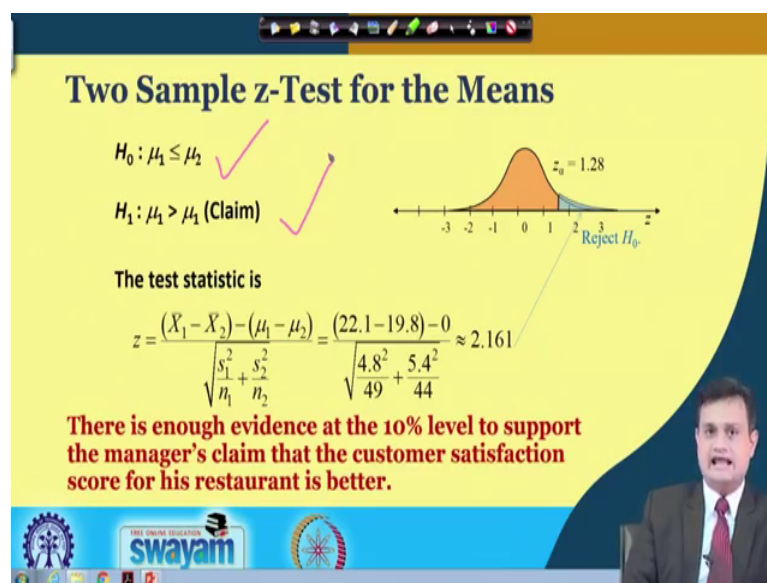
A manager of multi-cuisine restaurant in Hyderabad offering Chinese, South Indian and North Indian food claims that customer satisfaction score for his restaurant is higher than the customer satisfaction score for a competitor's restaurant. The mean customer satisfaction score for 49 customers for his restaurant is 22.1 and the standard deviation is 4.8. The mean customer satisfaction score for 44 of the competitor's restaurant is 19.8 and the standard deviation is 5.4. At $\alpha = 0.10$, can the manager's claim be supported?

swayam
INDIA WIDE, LIFELONG LEARNING

Let us say my situation is like this a manager of multi cuisine restaurant; in Hyderabad offering Chinese, South Indian, North Indian food. And he claims that the customer satisfaction score for his restaurant is higher than the customer satisfaction score for a competitors restaurant.

The mean customer satisfaction score for 49 customer for his restaurant is 22.1 in the standard deviation he found is 4.8. The mean customer satisfaction score for 44 of the competitors restaurant is 19.8 and the standard deviation is 5.4 he wants to check his claim that alpha is equal to 0.1 and he want to see that whether his claim is supported or not.

(Refer Slide Time: 25:59)



So, I will not go too much into detail of explaining; how to find critical value and refer it all this is well explained you have the hypothesis, μ_1 is less than equal to μ_2 it is about customer satisfaction score μ_2 is greater than μ_1 this is the claim of manager. And I want to check whether I am with H_0 or H_1 . So, you will compute z statistic and 2.16 is the value. So, you will say that 2.1 and your z alpha is 1.28.

So, when you start from 0 your 2.161 is somewhere here and this falls in the blue region; which is the rejection region. So, you will say there is enough evidence that ten percent level alpha is equal to 0.1. Support the claim of the manager that customer satisfaction score for his restaurant is better. So, he can feel satisfied that yes customers they are more happy with my strategy, my offerings compared to the competitor. So, this is something that can really give the confidence to the business and the managers.

(Refer Slide Time: 27:13)

Hypothesis Testing for Two Population

(σ_1 and σ_2 unknown, $\sigma_1^2 = \sigma_2^2$, n_1 or $n_2 < 30$, both population normal)

We can have another case; σ_1 and σ_2 unknown σ_1^2 is equal to σ_2^2 n_1 or n_2 is less than 30 both population normal.

(Refer Slide Time: 27:23)

Pooled-Variance t-Test Example

You are a financial analyst for a brokerage firm. Is there a difference in dividend yield between stocks listed on the BSE & NSE? You collect the following data:

	BSE	NSE
Number	21	25
Sample mean	3.27	2.53
Sample standard dev.	1.30	1.16

- Assuming both populations are approximately normal with equal variances, is there a difference in mean yield ($\alpha = 0.05$)?


So, let us say another interesting case you have the stock exchange BSE and NSE. Suppose I am interested to see or an analyst is interested to see that is there a difference in dividend yield between stocks listed on BSE and NSE. So, he just said some data crunching found let us say; 21 BSE and 25 say NSE share. Sample mean is 3.27 for dividend and standard deviation is 1.3, 1.16 wants to check at alpha is equal to 0.05.

(Refer Slide Time: 28:07)

Pooled-Variance t Test Example: Calculating the Test Statistic

$H_0: \mu_1 - \mu_2 = 0$ i.e. $(\mu_1 = \mu_2)$
 $H_1: \mu_1 - \mu_2 \neq 0$ i.e. $(\mu_1 \neq \mu_2)$

The test statistic is:

$$t = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{S_p^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}} = \frac{(3.27 - 2.53) - 0}{\sqrt{1.5021 \left(\frac{1}{21} + \frac{1}{25} \right)}} = 2.040$$
$$S_p^2 = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{(n_1 - 1) + (n_2 - 1)} = \frac{(21 - 1)1.30^2 + (25 - 1)1.16^2}{(21 - 1) + (25 - 1)} = 1.5021$$


So, now we can just compute the t statistic. So, this is my null hypothesis $\mu_1 - \mu_2$ is equal to 0; $\mu_1 - \mu_2$ is not equal to 0 it means there is not much difference in the dividend yield by the BSE or NSE, null hypothesis says this and alternate says there is a difference. So, I will compute here t for this case and comes out to be 2.040. My pooled variance is computed with this expression which is already given to you and I can take the decision.

(Refer Slide Time: 28:47)

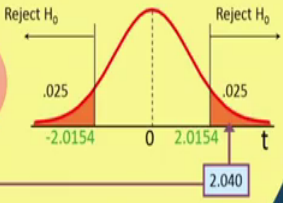
Pooled-Variance t Test Example: Hypothesis Test Solution

$H_0: \mu_1 - \mu_2 = 0$ i.e. $(\mu_1 = \mu_2)$
 $H_1: \mu_1 - \mu_2 \neq 0$ i.e. $(\mu_1 \neq \mu_2)$
 $\alpha = 0.05$
 $df = 21 + 25 - 2 = 44$
Critical Values: $t = \pm 2.0154$


Test Statistic:

$$t = \frac{3.27 - 2.53}{\sqrt{1.5021 \left(\frac{1}{21} + \frac{1}{25} \right)}} = 2.040$$

Decision:
Reject H_0 at $\alpha = 0.05$

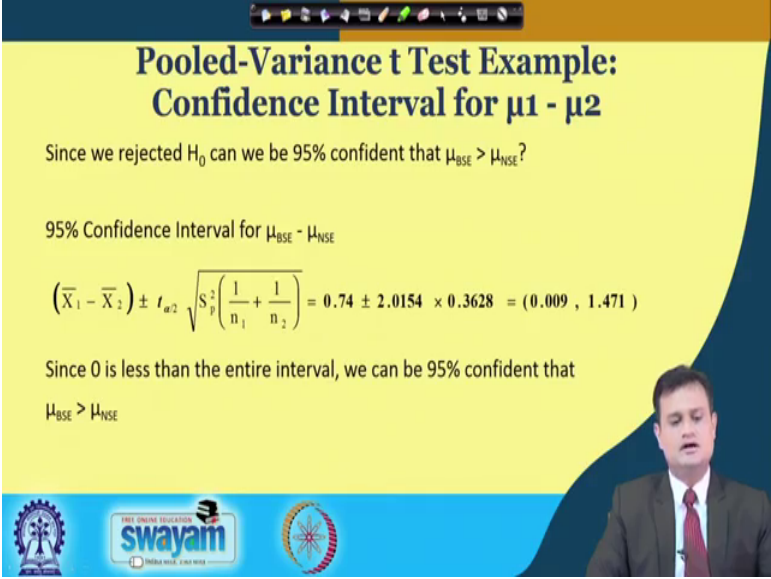


Conclusion:
There is evidence of a difference in means.



So, here you will see that alpha is equal to 0.05 degree of freedom is 21 plus 25 n 1 plus n 2 minus 2 is 44. So, my computed statistics 2.40 falls in this there is a evidence of difference in mean. So, yes the dividend is different when you look at the stocks listed on the BSE and NSE that is what my financial analysis says. And I reject the null hypothesis so my alternate claim is true that is mu 1 is not equal to mu 2. It means there is a difference in the dividend if I refer BSE and if I refer a stock on the NSE. So, this is something really interesting.

(Refer Slide Time: 29:37)



**Pooled-Variance t Test Example:
Confidence Interval for $\mu_1 - \mu_2$**

Since we rejected H_0 can we be 95% confident that $\mu_{BSE} > \mu_{NSE}$?

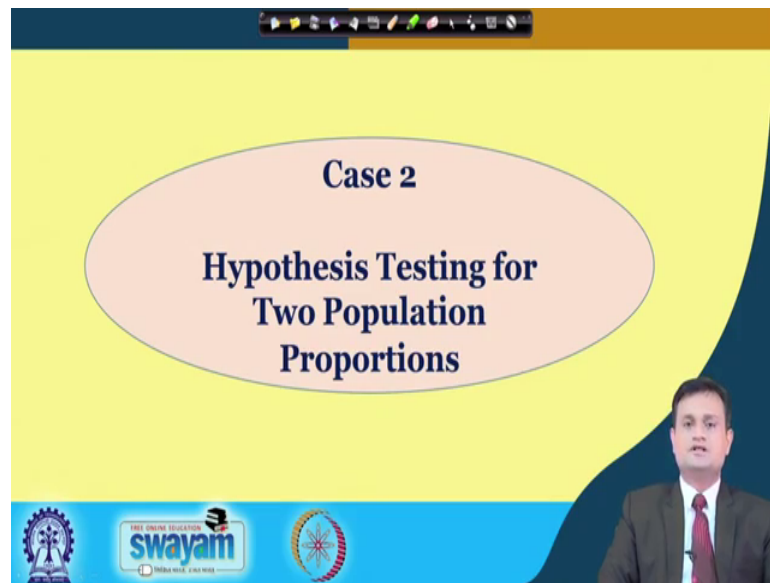
95% Confidence Interval for $\mu_{BSE} - \mu_{NSE}$

$$(\bar{X}_1 - \bar{X}_2) \pm t_{\alpha/2} \sqrt{S_p^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)} = 0.74 \pm 2.0154 \times 0.3628 = (0.009, 1.471)$$

Since 0 is less than the entire interval, we can be 95% confident that $\mu_{BSE} > \mu_{NSE}$

And I can also have confidence interval on my \bar{X}_1 minus \bar{X}_2 for BSE and NSE. So, you can use this particular expression \bar{X}_1 minus \bar{X}_2 plus or minus t alpha by 2. And you can see that whether my fall value falls within this will again yield the same result for 95 percent confidence interval. That mu of BSE is greater than mu of NSE.

(Refer Slide Time: 30:09)



I have case 2 so this is about two population proportion we have already discussed; the case for single population proportion case here there is only a difference that it would be two population proportion.

(Refer Slide Time: 30:25)

A presentation slide with a yellow background and a blue wavy border on the right. The text on the slide includes: "Goal: test a hypothesis or form a confidence interval for the difference between two population proportions," followed by a box containing $\pi_1 - \pi_2$. Below this, it says "Assumptions:" followed by two lines: $n_1 \pi_1 \geq 5, n_1(1 - \pi_1) \geq 5$ and $n_2 \pi_2 \geq 5, n_2(1 - \pi_2) \geq 5$. Then it says "The point estimate for the difference is" followed by $p_1 - p_2$. On the left, there is a blue box labeled "Population proportions". At the bottom left, there are logos for "swayam" and "INDIA WIDE, FREE WIDE". A presenter in a suit is visible in the bottom right corner.

So, we will just little bit do it quickly. So, I have the proportion $\pi_1 - \pi_2$ and I have $n_1 \pi_1 \geq 5$, $n_1(1 - \pi_1) \geq 5$. Let us say $n_2 \pi_2 \geq 5$, $n_2(1 - \pi_2) \geq 5$. The point estimate for the difference is $p_1 - p_2$ and I want to check the difference in population proportions.

(Refer Slide Time: 30:49)

Two Population Proportions

Population proportions

In the null hypothesis we assume the null hypothesis is true, so we assume $\pi_1 = \pi_2$ and pool the two sample estimates

The pooled estimate for the overall proportion is:

$$\bar{p} = \frac{X_1 + X_2}{n_1 + n_2}$$

where X_1 and X_2 are the number of items of interest in samples 1 and 2

The slide features a yellow background with a blue header and footer. The footer includes the Swayam logo and the text 'FREE ONLINE EDUCATION swayam'. A small video inset of a man in a suit is visible in the bottom right corner.

So, this is how I can easily compute my $\bar{p} = \frac{X_1 + X_2}{n_1 + n_2}$. And X_1 and X_2 are the numbers of items of interest in sample 1 and sample 2. And I have the sample size n_1 and n_2 so you can easily compute this.

(Refer Slide Time: 31:09)

Two Population Proportions

Population proportions

The test statistic for $\pi_1 - \pi_2$ is a Z statistic:

$$Z_{STAT} = \frac{(p_1 - p_2) - (\pi_1 - \pi_2)}{\sqrt{\bar{p}(1 - \bar{p})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$$

where $\bar{p} = \frac{X_1 + X_2}{n_1 + n_2}$, $p_1 = \frac{X_1}{n_1}$, $p_2 = \frac{X_2}{n_2}$

The slide features a yellow background with a blue header and footer. The footer includes the Swayam logo and the text 'FREE ONLINE EDUCATION swayam'. A small video inset of a man in a suit is visible in the bottom right corner.


And here I would follow the Z distribution Z statistics for this and various values.

(Refer Slide Time: 31:23)

Hypothesis Tests for Two Population Proportions

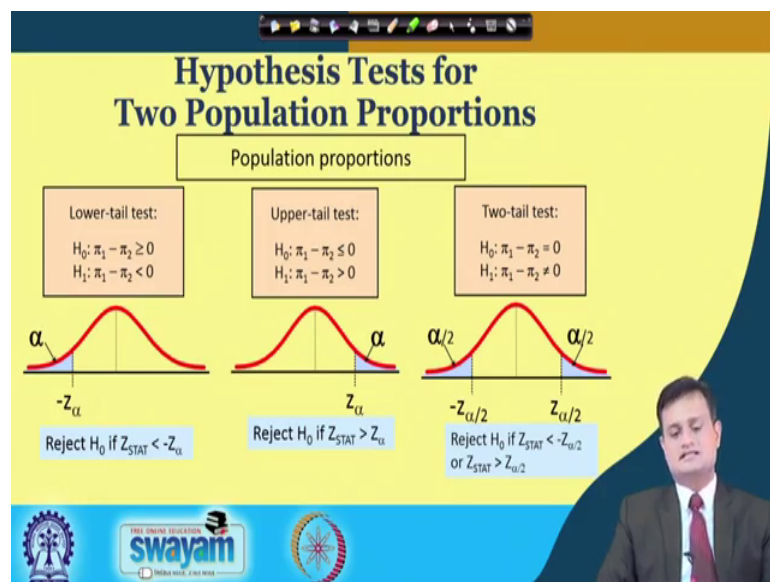
Population proportions

<p>Lower-tail test:</p> $H_0: \pi_1 \geq \pi_2$ $H_1: \pi_1 < \pi_2$ <p style="text-align: center;">i.e.,</p> $H_0: \pi_1 - \pi_2 \geq 0$ $H_1: \pi_1 - \pi_2 < 0$	<p>Upper-tail test:</p> $H_0: \pi_1 \leq \pi_2$ $H_1: \pi_1 > \pi_2$ <p style="text-align: center;">i.e.,</p> $H_0: \pi_1 - \pi_2 \leq 0$ $H_1: \pi_1 - \pi_2 > 0$	<p>Two-tail test:</p> $H_0: \pi_1 = \pi_2$ $H_1: \pi_1 \neq \pi_2$ <p style="text-align: center;">i.e.,</p> $H_0: \pi_1 - \pi_2 = 0$ $H_1: \pi_1 - \pi_2 \neq 0$
--	--	--



You can compute using these expressions. So, for a typical case maybe you can compute the lower tailed test, you can compute the upper tailed test, you can compute the two tailed test and would like to verify your hypothesis.

(Refer Slide Time: 31:37)



So, the decision rule we already discussed this is just kind of revisit. That if it is a lower tailed test this will be your alpha and this would be your rejection region. This would be alpha this would be rejection region this is divided by alpha by 2 and alpha by 2 and same way the rule follows for proportion population hypothesis testing.

(Refer Slide Time: 31:55)

Hypothesis Test Example: Two population Proportions

Is there a significant difference between the proportion of men and the proportion of women who will vote Yes on Proposition A?

- In a random sample, 36 of 72 men and 35 of 50 women indicated they would vote Yes
- Test at the 0.05 level of significance

swamy

So, let us see the example is there a significant difference between the proportion of men and the proportion of women who will vote yes on proposition a. Suppose there is some policy and you want to check that will there be a difference on the proportion of yes and proportion of no from men and women two different gender. I am taking a random sample of 36 of 72 men and 35 of 52 women indicated that they would say yes. I want to check my claim at level of significance 0.05.

(Refer Slide Time: 32:45)

Hypothesis Test Example: Two population Proportions

The test statistic for $\pi_1 - \pi_2$ is:

$$Z_{STAT} = \frac{(p_1 - p_2) - (\pi_1 - \pi_2)}{\sqrt{\bar{p}(1 - \bar{p})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$$
$$= \frac{(.50 - .70) - (0)}{\sqrt{.582(1 - .582)\left(\frac{1}{72} + \frac{1}{50}\right)}} = -2.20$$

Critical Values = ± 1.96
For $\alpha = .05$

Decision: Do not reject H_0

Conclusion: There is not significant evidence of a difference in proportions who will vote yes between men and women.

So, just compute the various values and then when you plug in these values into z statistic; then to check it. And you will see that for this minus point; 2.20 which is my computed value and 1.96 minus and plus it is a two tailed test it is my critical value. So, this particular value minus 2.20 is less than minus 1.96.

So, it is in the rejection region so there is no significant evidence of a difference in proportion who will vote yes between men and women. So, my null hypothesis is basically there is not significant difference and I would say my null hypothesis is basically not supported. So, I will say there is not significant evidence of a difference in proportion.

(Refer Slide Time: 33:43)

Confidence Interval for Two Population Proportions

Population proportions

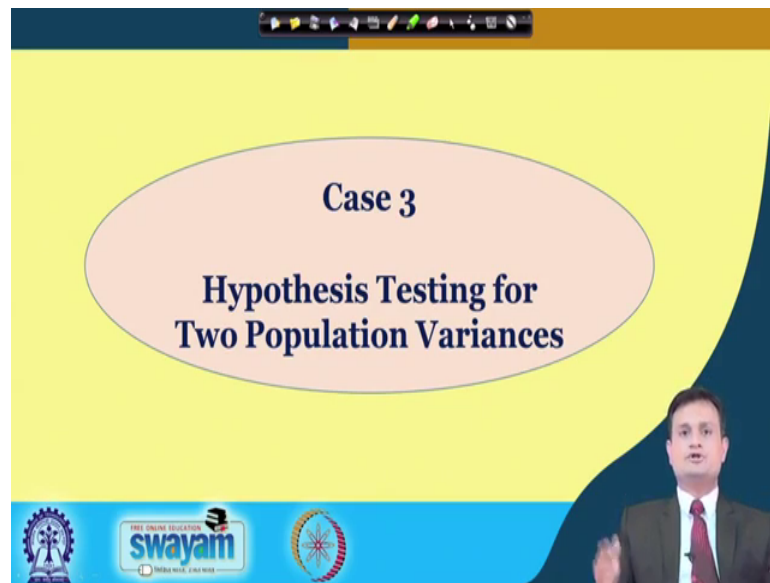
The confidence interval for $\pi_1 - \pi_2$ is:

$$(p_1 - p_2) \pm Z_{\alpha/2} \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}$$

swamyam

This is the confidence interval as usual you can compute it with Z alpha by 2 here and use this expression. So, this will help you to compute the confidence interval.

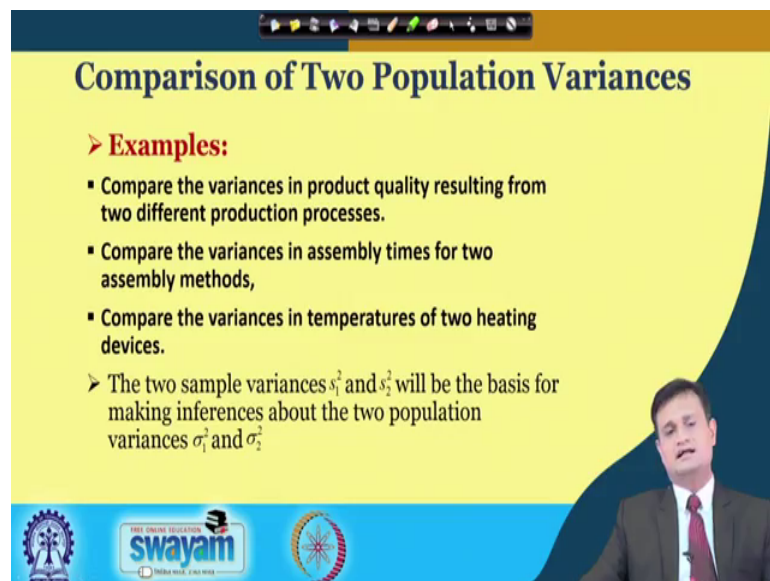
(Refer Slide Time: 33:55)



The slide features a yellow background with a large pink oval in the center. Inside the oval, the text reads "Case 3" followed by "Hypothesis Testing for Two Population Variances" in a bold, dark blue font. At the bottom of the slide, there is a blue banner containing logos for "swayam" and other educational institutions. A small inset video of a male presenter in a suit is located in the bottom right corner.

Now, I have another case 3 this is hypothesis testing for two population variances.

(Refer Slide Time: 34:03)



The slide has a yellow background with a blue banner at the bottom. The title "Comparison of Two Population Variances" is at the top in a bold, dark blue font. Below the title, there is a red heading "Examples:" followed by a list of three bullet points. A fourth point, starting with a red arrow, explains the relationship between sample and population variances. The blue banner at the bottom contains logos for "swayam" and other educational institutions. A small inset video of a male presenter in a suit is located in the bottom right corner.

Comparison of Two Population Variances

➤ **Examples:**

- Compare the variances in product quality resulting from two different production processes.
- Compare the variances in assembly times for two assembly methods,
- Compare the variances in temperatures of two heating devices.

➤ The two sample variances s_1^2 and s_2^2 will be the basis for making inferences about the two population variances σ_1^2 and σ_2^2

So, let us see the practical part that what could be the situations in which I would like to go for this. Suppose compare the variance in product quality; so the important parameter for quality is to check the variability and this is resulting from two different production process. I want to see that the variability in quality from process 1 and the process 2 is there really a significant difference or they are same in terms of variability.

Compare the variance in assembly times for two assembly method compare the variances in temperature of 218 devices. And these are some of the phenomena that could be of interest as a part of case 3. So, two sample variance S_1^2 and S_2^2 will be the basis for making inferences about two population variance σ_1^2 and σ_2^2 square.

(Refer Slide Time: 34:59)

Testing for the Ratio Of Two Population Variances

Hypotheses	F_{STAT}
$H_0: \sigma_1^2 = \sigma_2^2$	S_1^2 / S_2^2
$H_1: \sigma_1^2 \neq \sigma_2^2$	

F test statistic

Where:

S_1^2 = Variance of sample 1 (the larger sample variance) n_2 = sample size of sample 2
 n_1 = sample size of sample 1 $n_1 - 1$ = numerator degrees of freedom
 S_2^2 = Variance of sample 2 (the smaller sample variance) $n_2 - 1$ = denominator degrees of freedom

So, this is something explained like this yes the variability is same or it is different express as null hypothesis alternate. And this is my F statistic which is basically the ratio of my sample variance S_1^2 divided by S_2^2 square and this is your numbers of degrees of freedom that you need to refer the particular table F table and find the F statistics. So, you have S_1^2 / S_2^2 $n_1 - 1$ and $n_2 - 1$.

(Refer Slide Time: 35:37)

The F Distribution

- The F critical value is found from the F table
- There are two degrees of freedom required: numerator and denominator
- The larger sample variance is always the numerator
- When $F_{STAT} = \frac{S_1^2}{S_2^2}$ $df_1 = n_1 - 1$; $df_2 = n_2 - 1$
- In the F table,
 - numerator degrees of freedom determine the column
 - denominator degrees of freedom determine the row

Now, this is my F statistic and these are my degrees of freedom.

(Refer Slide Time: 35:47)

Finding the Rejection Region

$H_0: \sigma_1^2 = \sigma_2^2$
 $H_1: \sigma_1^2 \neq \sigma_2^2$

Reject H_0 if $F_{STAT} > F_{\alpha/2}$

$H_0: \sigma_1^2 \leq \sigma_2^2$
 $H_1: \sigma_1^2 > \sigma_2^2$

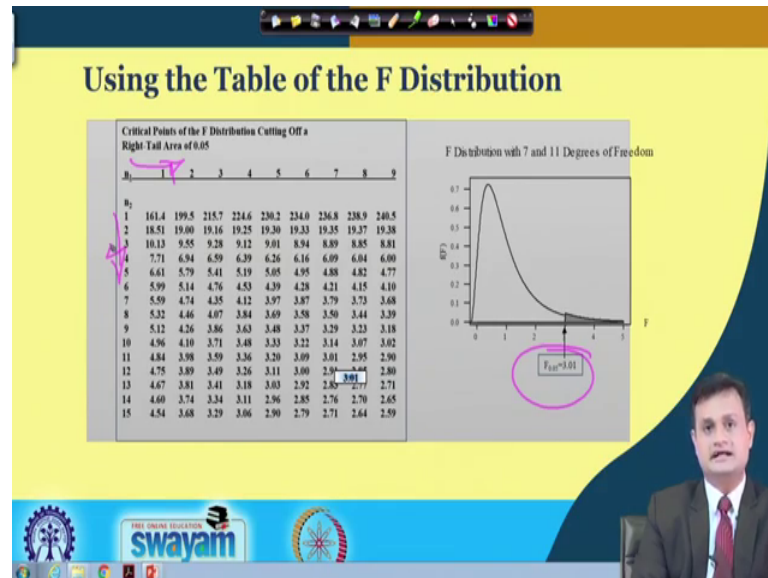
Reject H_0 if $F_{STAT} > F_{\alpha}$

So, I will have a very simple rule to follow this is my alpha by two this is my alpha and this is my do not reject null hypothesis region this is reject null hypothesis. So, if F stat is greater than F alpha by 2. So, calculated value is greater than F alpha by 2 reject null hypothesis only if stat F stat calculated is greater than F alpha by 2.

So, it will fall basically in the rejection region. And same way this is reject H 0 this is do not reject H 0 and this is my reject H 0 if F stat is greater than F alpha. Only difference is

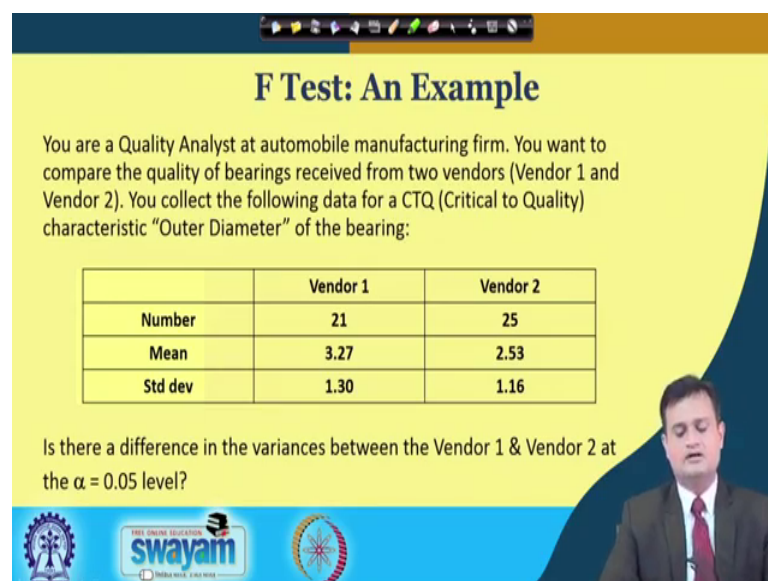
that here I am using alpha by 2 here I am using alpha. So, this is something that we can do for the for checking the variability.

(Refer Slide Time: 36:39)



This is my F table you can see that you can find the critical value from the F table. Let us say $F_{0.05}$ this is the 3.01 and you have n_1 and you have n_2 . So, with this n_1 and n_2 you can figure out that what could be the critical value for your F test.

(Refer Slide Time: 37:07)



Now, let us say I want to see the example that there are two vendors and quality analyst want to check. That quality of bearing coming from vendor one and coming from vendor

two are they same or there is a huge variability in terms of the quality. So, data is vendor 1 vendor 2 number is 21 that is my sample, 25 mean is 3.27, 2.53, 1.30, 1.16 I want to check it at alpha 0.05.


(Refer Slide Time: 37:43)

F Test: Example Solution

- Form the hypothesis test:

$H_0: \sigma_1^2 = \sigma_2^2$ (there is no difference between variances)

$H_1: \sigma_1^2 \neq \sigma_2^2$ (there is a difference between variances)
- Find the F critical value for $\alpha = 0.05$:
- Numerator d.f. = $n_1 - 1 = 21 - 1 = 20$ ✓
- Denominator d.f. = $n_2 - 1 = 25 - 1 = 24$ ✓
- $F_{\alpha/2} = F_{0.025, 20, 24} = 2.33$ ✓



So, we can just compute couple of statistics and we can see that my F alpha by 2 at F 0.025 which is my level of significance alpha by 2. And 20 24 which is my degree of freedom mu 1 and mu 2 it comes out to be 2.33. So, now, what does it mean ?

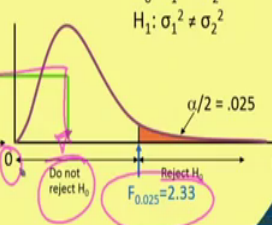
(Refer Slide Time: 38:05)

F Test: Example Solution


- The test statistic is:

$$F_{STAT} = \frac{S_1^2}{S_2^2} = \frac{1.30^2}{1.16^2} = 1.256$$
- $F_{STAT} = 1.256$ is not in the rejection region, so we do not reject H_0 .

$H_0: \sigma_1^2 = \sigma_2^2$
 $H_1: \sigma_1^2 \neq \sigma_2^2$

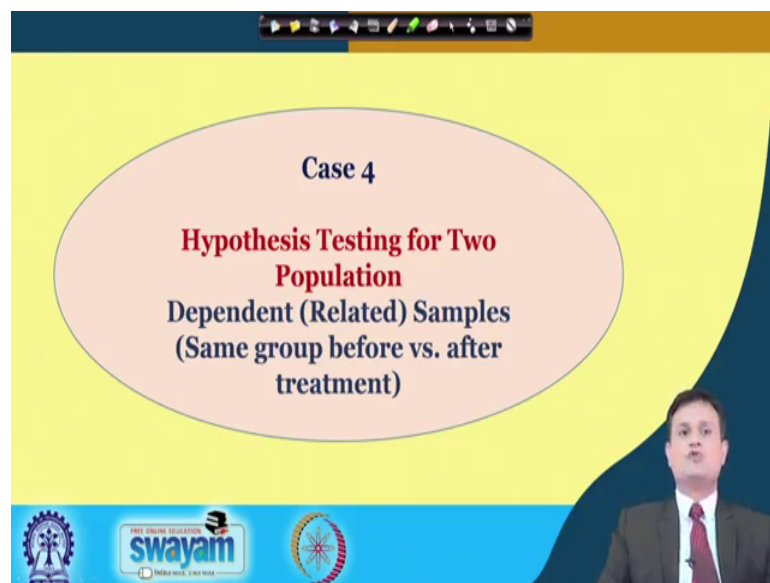


Conclusion: There is not sufficient evidence of a difference in variances at $\alpha = 0.05$



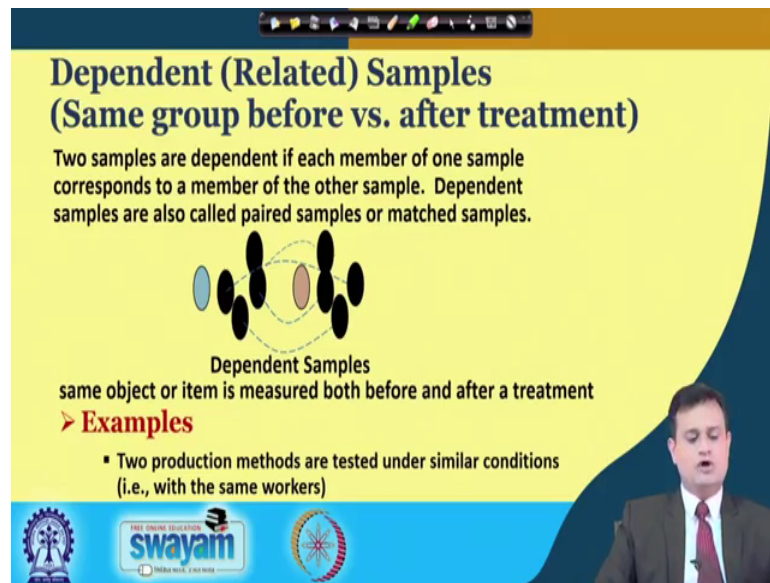
It means that this particular value 2.33 this is my $F_{0.025}$. And F stat when I compute it is 1.256. So, 1.256 is not in the rejection region so we do not reject H_0 . Because 1.256 will fall somewhere here this is the 0. So, I do not reject null hypothesis it means I do not have sufficient evidence to say that quality coming from or variability in quality coming from vendor 1 and the variability quality coming from vendor 2 they are different more or less they are same. So, I do not have sufficient evidence at $\alpha 0.05$ to say that the quality of vendor 1 and quality of vendor 2 is different. So, fine we have seen the case 3.

(Refer Slide Time: 39:09)

The image shows a presentation slide with a yellow background and a blue border. In the center, there is a light pink oval containing the text: "Case 4", "Hypothesis Testing for Two Population", "Dependent (Related) Samples", and "(Same group before vs. after treatment)". The text "Hypothesis Testing for Two Population" is in red, while the rest is in black. In the bottom right corner, there is a small video inset of a man in a suit and tie. At the bottom of the slide, there are logos for "swayam" and "INDIA RISE, YOUNG RISE".


And let us try to end up with one more case that is the case four. So, hypothesis testing for two population dependent related sample. You have the sample patients give them a treatment you are interested to see what is the effect of the treatment before there where a symptoms after treatment what is the reduction in the symptom and you want to check this two situation.

(Refer Slide Time: 39:35)



Dependent (Related) Samples (Same group before vs. after treatment)

Two samples are dependent if each member of one sample corresponds to a member of the other sample. Dependent samples are also called paired samples or matched samples.



Dependent Samples
same object or item is measured both before and after a treatment

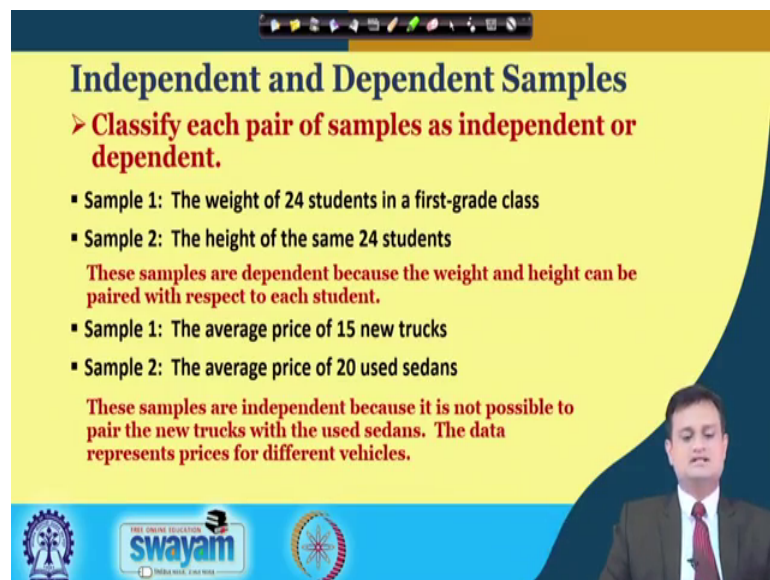
► **Examples**

- Two production methods are tested under similar conditions (i.e., with the same workers)

swayam

So, this is what it is that two production methods may be are tested under similar condition are they same or they are different.

(Refer Slide Time: 39:43)



Independent and Dependent Samples

► **Classify each pair of samples as independent or dependent.**

- Sample 1: The weight of 24 students in a first-grade class
- Sample 2: The height of the same 24 students

These samples are dependent because the weight and height can be paired with respect to each student.

- Sample 1: The average price of 15 new trucks
- Sample 2: The average price of 20 used sedans

These samples are independent because it is not possible to pair the new trucks with the used sedans. The data represents prices for different vehicles.

swayam

We may have many examples that independent and dependent sample. Suppose sample 1 the weight of 24 students in a first grade class student through the height of the same 24 students. The sample are dependent because weight and height can be paired there could be some relationship. But if you see the second one sample 1 the average price of 15 new

trucks sample 2 average price of 20 used sedan these are independent there could not be any dependency or relationship.

(Refer Slide Time: 40:19)

Comparing Two Dependent (Related) Samples

- **Test the Means of Two Related Samples**
 - Paired or matched
 - Repeated measures (before and after)
 - Use difference between i^{th} pair
- **Eliminates variation between subjects**

$$D_i = X_{1i} - X_{2i}$$

The slide features a yellow background with a blue header and footer. The footer includes the Swayam logo and a small video inset of a man in a suit.

So, basically you have paired or matched two related samples repeated measures. And use difference between i^{th} pair; D_i is equal to X_{1i} minus X_{2i} . So, observation specific to sample 1 or before treatment and X_{2i} after treatment and you will ; obviously, take number of sample to test your claim.

(Refer Slide Time: 40:45)

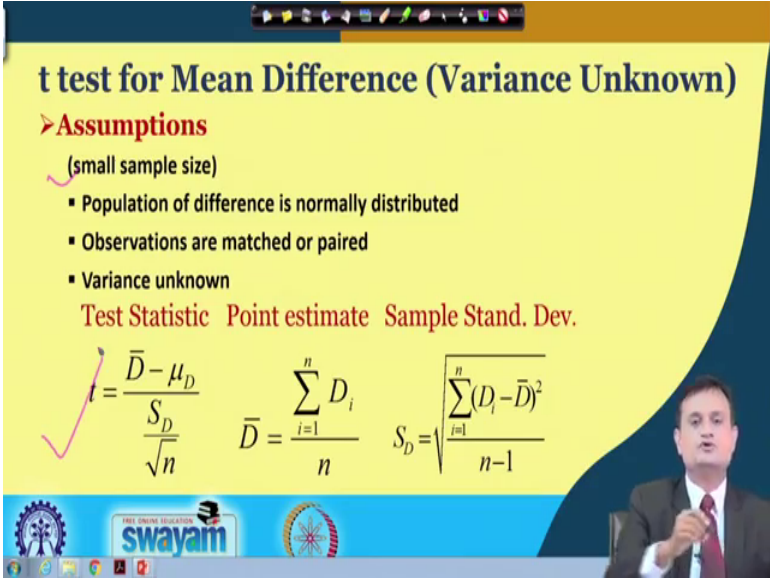
Z Test for Mean Difference (Variance Known)

- **Assumptions**
 - Population of difference is normally distributed or large samples
 - Observations are paired or matched
 - Variance known
- **The point estimate for the population mean paired difference is**
Mean of differences $\bar{D} = \frac{\sum_{i=1}^n D_i}{n}$
 - number of pairs in the paired sample
- **Test Statistic**
$$Z = \frac{\bar{D} - \mu_D}{\frac{\sigma_D}{\sqrt{n}}}$$
 - mean of the difference in values for the population
 - SD of the difference in values for the population

The slide features a yellow background with a blue header and footer. The footer includes the Swayam logo and a small video inset of a man in a suit.

So, you can compute couple of things like \bar{D} this is basically \bar{D} i. So, number of pairs in the paired sample then \bar{D} minus μ_D divided by this. Same procedure same computation as we do for z statistics.

(Refer Slide Time: 41:11)



t test for Mean Difference (Variance Unknown)

➤ **Assumptions**

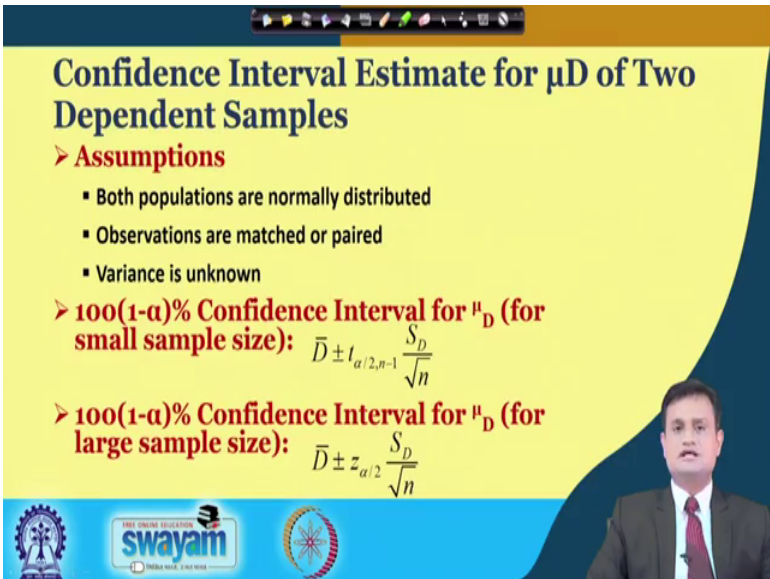
- (small sample size)
 - Population of difference is normally distributed
 - Observations are matched or paired
 - Variance unknown

Test Statistic Point estimate Sample Stand. Dev.

$$t = \frac{\bar{D} - \mu_D}{\frac{S_D}{\sqrt{n}}} \quad \bar{D} = \frac{\sum_{i=1}^n D_i}{n} \quad S_D = \sqrt{\frac{\sum_{i=1}^n (D_i - \bar{D})^2}{n-1}}$$

So, you will suppose you have small sample. Then definitely you will use the logic that you will use student t distribution by plugging in the values of \bar{X} minus μ divided by sigma. Here it is \bar{D} minus μ_D divided by sigma that is S_D divided by square root of n and \bar{D} is this your S_D is this. So, these are the standard expressions.

(Refer Slide Time: 41:37)



Confidence Interval Estimate for μ_D of Two Dependent Samples

➤ **Assumptions**

- Both populations are normally distributed
- Observations are matched or paired
- Variance is unknown

➤ **100(1- α)% Confidence Interval for μ_D (for small sample size):** $\bar{D} \pm t_{\alpha/2, n-1} \frac{S_D}{\sqrt{n}}$

➤ **100(1- α)% Confidence Interval for μ_D (for large sample size):** $\bar{D} \pm z_{\alpha/2} \frac{S_D}{\sqrt{n}}$

We have the assumptions that both populations are normally distributed observations are matched or paired variance is unknown. In this case you can again do the confidence interval analysis and construct it and you can find 100 1 minus alpha percent confidence interval.

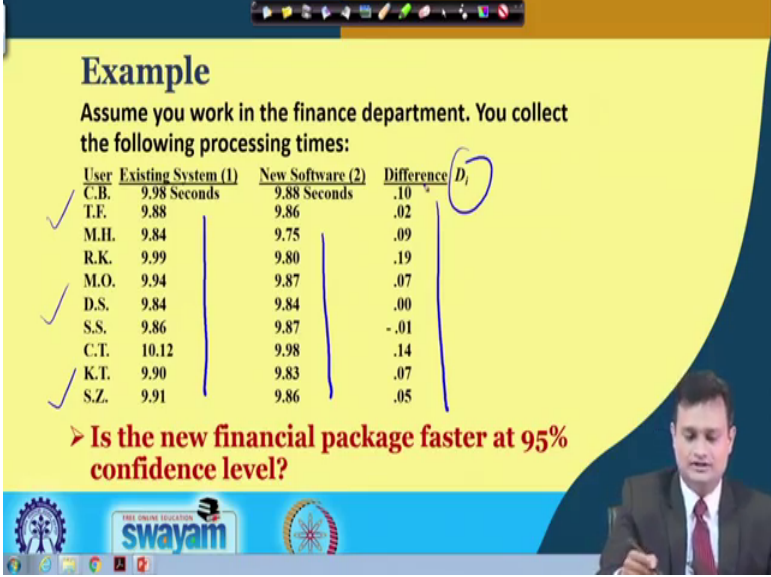
(Refer Slide Time: 42:03)

Example

Assume you work in the finance department. You collect the following processing times:

User	Existing System (1)	New Software (2)	Difference D_i
C.B.	9.98 Seconds	9.88 Seconds	.10
T.F.	9.88	9.86	.02
M.H.	9.84	9.75	.09
R.K.	9.99	9.80	.19
M.O.	9.94	9.87	.07
D.S.	9.84	9.84	.00
S.S.	9.86	9.87	-.01
C.T.	10.12	9.98	.14
K.T.	9.90	9.83	.07
S.Z.	9.91	9.86	.05

➤ Is the new financial package faster at 95% confidence level?



Now, let us say I have the data that you work in the finance department and collect some data regarding the existing system and the new software. And suppose I have let us say some user preference some user rating and let us say this is the existing system in terms of seconds experience or the running time.

This is the new software running time I am finding the D_i value that is my difference. I want to check that is the new financial package faster that 95 percent confidence level. It is a huge investment in software you would like to see that whether the new system purchased is really worthy or not.

(Refer Slide Time: 42:55)

Solution

$$\bar{D} = \frac{\sum D_i}{n} = 0.072 \quad S_D = \sqrt{\frac{\sum (D_i - \bar{D})^2}{n-1}} = 0.06215$$

$H_0: \mu_D \leq 0$
 $H_1: \mu_D > 0$ (Claim)

At $\alpha = 0.05$ and $df = n - 1 = 9$

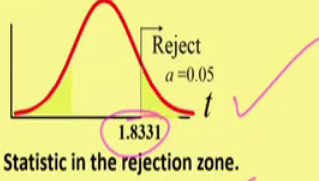
$$t_{9,0.05} = 1.8331$$

Test Statistic

$$t = \frac{\bar{D} - \mu_D}{S_D / \sqrt{n}} = \frac{0.072 - 0}{0.06215 / \sqrt{10}} = 3.66$$

Decision: Reject H_0

Conclusion: The new software package is faster.



t Statistic in the rejection zone.

swayam

So, I would test the hypothesis find the t distribution t statistics. And for this t statistics I will compute the value and this would be compared with the critical value the procedure remains same. And I will say reject H_0 because this is 3.66. So, this falls in the rejection region and the new software package is faster.

(Refer Slide Time: 43:27)

Find a 95% confidence interval

$$t_{\alpha/2, n-1} = t_{0.025, 9} = 2.2622$$

$$\bar{D} \pm t_{\alpha/2, n-1} \frac{S_D}{\sqrt{n}}$$

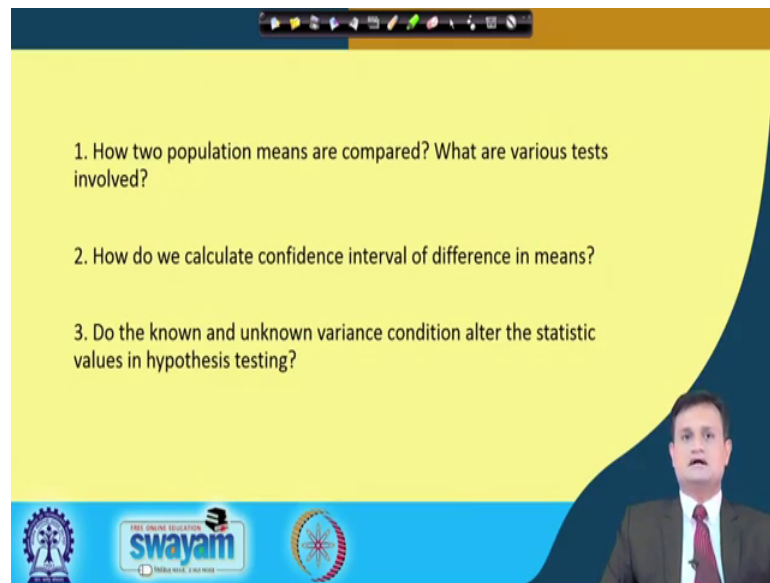
$$0.072 \pm 2.2622 \left(\frac{0.06215}{\sqrt{10}} \right)$$

$$0.0275 < \mu_D < 0.1165$$

swayam

So, with this is what we basically try to do. and you can construct the confidence interval same way by considering plus or minus t alpha by 2 at n minus 1 degree of freedom.

(Refer Slide Time: 43:37)



1. How two population means are compared? What are various tests involved?

2. How do we calculate confidence interval of difference in means?

3. Do the known and unknown variance condition alter the statistic values in hypothesis testing?

The slide features a yellow background with a blue wavy border on the right. At the bottom, there are logos for Anna University, Swayam, and the Ministry of Education, India. A presenter is visible in the bottom right corner.

Let us end this section with couple of think it question. How two population means are compared? What are various tests involved? How do we calculate confidence interval for difference in means and do the known and the unknown variance condition alter the statistics value in hypothesis testing?

(Refer Slide Time: 43:57)



References:

- ❑ Aczel, A., Sounderpandian, J. and Saravanan, P. , Complete Business Statistics, McGraw Hill Publication.
- ❑ David M. Levine, Timothy C. Krehbiel, Mark L. Berenson and P. K. Vishwanathan, Business Statistics, Pearson Publication.
- ❑ T. M. Kubiak, Donald W. Benbow, The Certified Six Sigma Black Belt Handbook, Pearson Publication.
- ❑ Forrest W. Breyfogle III, Implementing Six Sigma, John Wiley & Sons, INC.

The slide features a yellow background with a blue wavy border on the left. The word 'References' is written in a large, stylized font. At the bottom, there are logos for Anna University, Swayam, and the Ministry of Education, India. A presenter is visible in the bottom right corner.

Please refer this reference if you find this session little bit difficult to understand. And I hope this session would provide a good conceptual understanding to digest the concept and see the importance of two population test.

(Refer Slide Time: 44:13)

Conclusion:

- ❖ Compared two independent samples
 - ✓ Performed pooled-variance t test for the difference in two means
 - ✓ Performed separate-variance t test for difference in two means
 - ✓ Formed confidence intervals for the difference between two means
- ❖ Compared two population proportions
 - ✓ Formed confidence intervals for the difference between two population proportions
 - ✓ Performed Z-test for two population proportions
- ❖ Performed F test for the ratio of two population variances

swayam

So, these are the cases basically we have discussed for two population test, two independent sample, compared two population proportion variances and before treatment and after treatment. So, thank you very much for your interest in learning two population hypothesis testing; keep revising, be with me, enjoy.