

Analysing Count Data

Slide 1:	Introduction.....	3
Slide 2:	Categorical Data	3
Slide 3:	Section 1: Testing Single Proportions	5
Slide 4:	Testing Single Proportions: Soft Drink Example (1)	5
Slide 5:	Testing Single Proportions: Soft Drink Example (2)	6
Slide 6:	Hypothesis Test for p : Soft Drink Example.....	7
Tab 1:	Stating the Null Hypothesis.....	8
Tab 2:	Computing the Test Statistic.....	9
Tab 3:	Evaluating the Test Statistic	9
Slide 7:	Confidence Interval for a Single Proportion.....	10
Slide 8:	Analysing Single Proportions in R.....	11
Slide 9:	Final Remarks on Single Proportions.....	12
Slide 10:	Section 2: Comparing Two Proportions	13
Slide 11:	Fish Oil Example.....	13
Slide 12:	Hypothesis Test for Comparing Two Proportions.....	14
Tab 1:	Stating the Null Hypothesis.....	15
Tab 2:	Computing the Test Statistic.....	15
Tab 3:	Evaluating the Test Statistic	16
Tab 4:	Conclusions.....	17
Slide 13:	Confidence Interval for Comparing Two Proportions.....	17
Slide 14:	Comparing Two Proportions in R.....	18
Slide 15:	Final Remarks on Comparing Two Proportions.....	19
Slide 16:	Section 3: Contingency Tables and Chi-square Tests	19
Slide 17:	Chocolate Example	20
Slide 18:	21
Tab 1:	Expressing the Hypotheses.....	21
Tab 2:	Finding the Expected Values.....	22
Tab 3:	Computing the Test Statistic.....	23
Tab 4:	Identifying the Degrees of Freedom	24
Tab 5:	Evaluating the Test Statistic	24
Tab 6:	Conclusions.....	25
Tab 7:	Chi-square Test for Independence in R.....	26

Slide 19:	26
Tab 1: Assumptions.....	27
Slide 20: Summary	28



Slide 1: **Introduction**

Trinity College Dublin
Coláiste na Tríonóide, Baile Átha Cliath
The University of Dublin

Analysing Count Data

Presenter: Caroline Brophy
Duration: 26:31
School: Computer Science and Statistics

Welcome to this presentation on analysing count or categorical data. My name is Caroline Brophy. In this session we will explore statistical methods for analysing count or categorical data variables.

We will study hypothesis testing and confidence intervals for a single proportion, and for comparing two proportions. We will also explore contingency tables for two categorical variables and a chi-square test of independence for two categorical variables.

Slide 2: **Categorical Data**

Categorical Data

2 of 19

- Categorical data consists of **categorical variables**, where k is the number of categories or levels.
- When the number of levels is two, it is a **binary variable**.
- When a categorical variable has an inherent ordering, it is an **ordinal** categorical variable.

? Which of our examples is ordinal?

Sports Team Supporters	Consumers of Alcohol	Cause of Death
<p>Do you support a sport?</p> <ul style="list-style-type: none">• Possible answers:<ul style="list-style-type: none">• Yes• No• $k = 2$ levels• Binary variable	<p>Do you consume alcohol?</p> <ul style="list-style-type: none">• Possible answers:<ul style="list-style-type: none">• None• Moderate• A lot• $k = 3$ levels• Ordinal variable	<p>What was the cause of death?</p> <ul style="list-style-type: none">• Possible answers:<ul style="list-style-type: none">• Heart disease• Cancer• Accident• Other• $k = 4$ levels

In earlier sessions in this module, the focus was on analysing continuous data. For example, data might have been collected on the height of a particular plant species, at random from the population of the species. In that case, the hypothesis test or confidence interval aims to make statements about the true population mean height of the tree species.

In this session, we will examine categorical data, where k is the number of categories or levels.

Let's take a look at some examples of categorical data variables.

Suppose a random sample of people from a population are asked if they a supporter of a particular sports team. The possible answers are ‘yes’ or ‘no’ and the response for each person is recorded. This is an example of a categorical data variable, with $k =$ two levels. A study such as this might be carried out to find out what proportion of people in the true population are supporters of the sports team.

Another example of a categorical variable is the amount of alcohol a person consumes, with possible answers ‘none’, ‘moderate’, or ‘a lot’. This time the categorical variable has three levels. And of course, the decision as to how to label the categories, or how many categories to have, will often lie with the person conducting the study.

A third example is cause of death where possible answers are ‘heart disease’, ‘cancer’, ‘accident’, and ‘other’. This variable has four levels, so k is equal to 4.

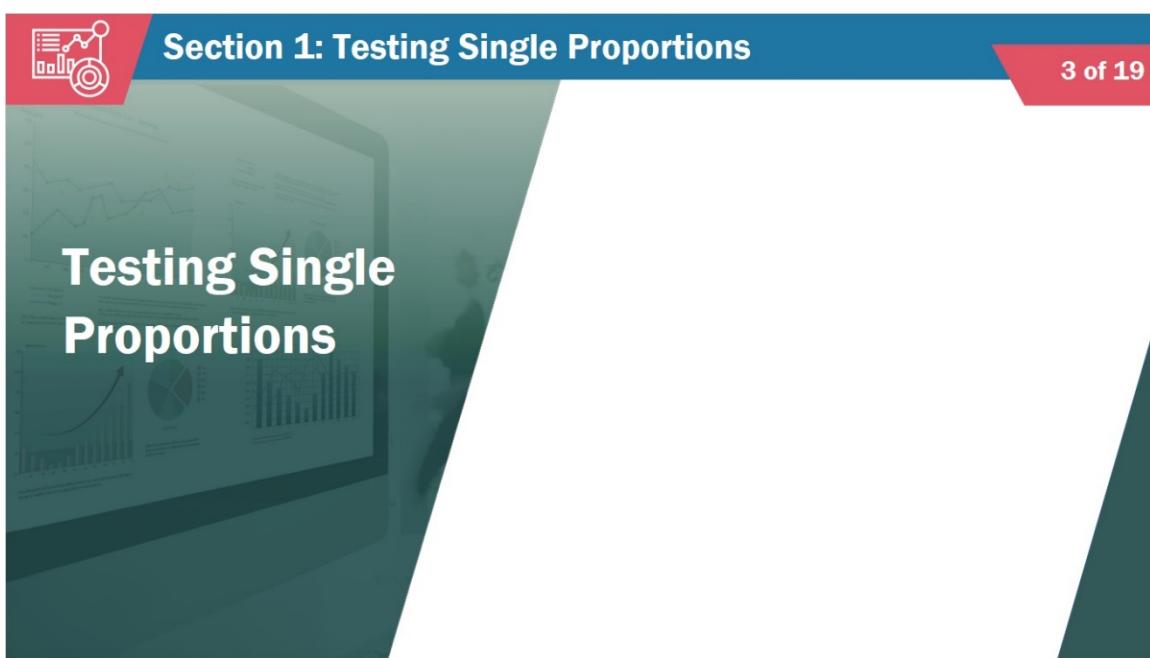
When the number of levels is two, we call this a binary variable. In the examples, the supporter of a sports team is an example of a binary variable.

When the ordering of the levels of a categorical variable has meaning, we call it ‘ordinal’.

Let's consider our three examples, are any of them ordinal? The sports team binary variable is not: it does not matter whether yes or no comes first. The causes of deaths could be listed in any order, for example, it would not matter if ‘accident’ came first instead of third. However, for the alcohol question, the levels of the category are increasing in magnitude of consumption and the order has logic to it. This variable is an example of an ordinal variable.

Slide 3:

Section 1: Testing Single Proportions

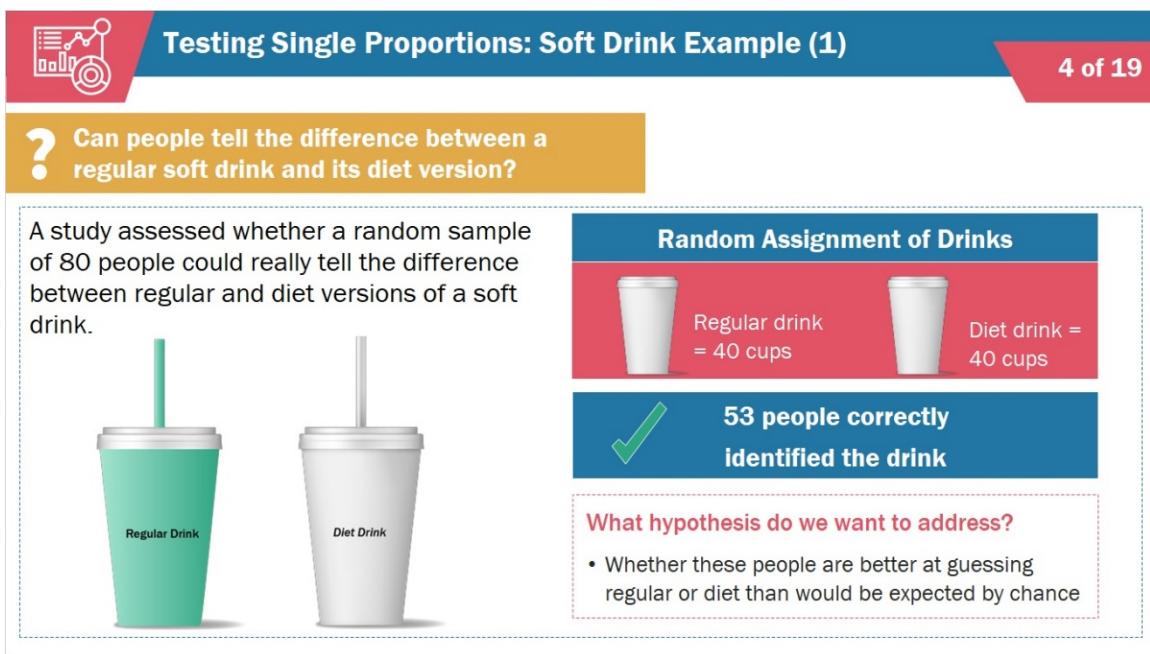


The slide features a red header bar with the title "Section 1: Testing Single Proportions". On the left, there is a small icon of a chart and target. On the right, it says "3 of 19". Below the header, the main title "Testing Single Proportions" is displayed in large white text against a dark green background with faint charts.

We will first explore how to analyse a single proportion.

Slide 4:

Testing Single Proportions: Soft Drink Example (1)



The slide has a red header bar with the title "Testing Single Proportions: Soft Drink Example (1)". On the left, there is a small icon of a chart and target. On the right, it says "4 of 19". Below the header, a yellow box contains the question: "Can people tell the difference between a regular soft drink and its diet version?".

A study assessed whether a random sample of 80 people could really tell the difference between regular and diet versions of a soft drink. Two cups, one green labeled "Regular Drink" and one white labeled "Diet Drink", are shown side-by-side.

Random Assignment of Drinks

 Regular drink = 40 cups	 Diet drink = 40 cups
--	---

53 people correctly identified the drink

What hypothesis do we want to address?

- Whether these people are better at guessing regular or diet than would be expected by chance

We will start with an example which involves a soft drink.

Here is a question for you to consider: can people tell the difference between regular and diet versions of soft drinks? If you were given a soft drink, for example, a glass of 7-up, and asked if it was diet 7-up or not, do you think you could you answer correctly?

A study was carried out to assess whether people who believe they can tell the difference between a regular soft drink and a diet version of it, can actually tell the difference.

Some people believe that they can tell the difference between a regular soft drink and a diet version of it. A researcher wanting to test this claim randomly sampled 80 such people. Eighty plain white cups were filled with the soft drink, half diet and half regular through random assignment. Each person took a sip from their cup and was asked to identify the drink as diet or regular. Of the 80 participants, 53 correctly identified the drink.

What is the question or hypothesis that we want to address here? Really, we want to know if these people are better at guessing regular or diet than would be expected by chance.

And what would be expected by chance?

Well, we would expect that a person would guess correctly 50% of the time.

Slide 5: Testing Single Proportions: Soft Drink Example (2)

 Testing Single Proportions: Soft Drink Example (2)
5 of 19

Soft Drink Study		
Population: <ul style="list-style-type: none"> • People who believe that they can tell the difference between regular and diet soft drinks 	Question: <ul style="list-style-type: none"> • Is there evidence that this population is any better or worse at telling the difference between regular and diet versions than by random guessing? <ul style="list-style-type: none"> • p = proportion of population who can identify this difference 	Data: <ul style="list-style-type: none"> • 53 of 80 participants correctly identified the drink type <ul style="list-style-type: none"> • We can estimate p from the data: $\hat{p} = \frac{x}{n} = \frac{53}{80} = 0.6625$
? How reliable is this estimate and what does it mean with respect to the population?		

The population here is people who believe that they can tell the difference between regular and diet soft drinks.

The question here is really whether or not this population of people are any better or worse than random guessing at telling the difference between a diet and regular soft drink. We use p to denote the proportion of people in the population who can identify the difference between a diet and regular soft drink.

The data arising from the study showed that 53 of 80 participants correctly identified the drink type. We can use this data to find an estimate of p , our parameter of interest.

Here we find that the estimate of p is 53 divided by 80 which equals 0.6625. Or we estimate that 66% of the population can indeed tell the difference.

Note the use of the hat symbol over the p ; this is to indicate that it is an estimate from data. When we say p , we mean the true population parameter, and when we say p_{hat} we mean the estimate of the proportion from the data.

Of course, with any estimate comes uncertainty. We will now look at ways to explore how reliable our estimate is and understand what it means with respect to the population.

Slide 6: Hypothesis Test for p : Soft Drink Example


Hypothesis Test for p : Soft Drink Example
6 of 19

Stating the Hypotheses
Computing the Test Statistic
Evaluating the Test Statistic

Introduction

- We will conduct a hypothesis test for p :
- p is the population of interest, those who can identify the difference between a diet and a regular soft drink.



 Click each tab to learn more. Then, click Next to continue.

We are going to conduct a hypothesis test for p , the proportion of people in the population who can identify the difference between a diet and regular soft drink. The logic we will follow to do this will be similar to that used in previous sessions for testing the population mean. We will specify our hypotheses, compute our test statistic, evaluate the test statistic against the reference distribution, and give our conclusions.

Click each tab in turn to learn more. When you are ready, click next to continue.

Tab 1: Stating the Null Hypothesis


Hypothesis Test for p : Soft Drink Example
6 of 19

Stating the Hypotheses
Computing the Test Statistic
Evaluating the Test Statistic

Stating the Hypotheses

- The null hypothesis is $H_0: p = p_0$
- Applied to the soft drink example, the null hypothesis is $H_0: p = 0.5$
- The alternative hypothesis is $H_A: p \neq 0.5$

$H_0: p = 0.5$ versus $H_A: p \neq 0.5$



! By “under the null hypothesis”, we mean “assuming that the null hypothesis is true”.

We start by stating the null hypothesis; p is the population proportion of interest, p_0 is the hypothesised value of p . So, our null hypothesis is that p equals p_0 .

In this example, the null hypothesis is that these people are no better or worse than random guessing at telling the difference between diet and regular soft drinks. With random guessing, there is a 50:50 chance of getting it right, so under the null hypothesis, we hypothesise that the population proportion is equal to 0.5. By saying “under the null hypothesis, we mean “assuming that the null hypothesis is true”.

The alternative hypothesis is that the population proportion does not equal 0.5.

Tab 2: Computing the Test Statistic

Hypothesis Test for p : Soft Drink Example

6 of 19

Stating the Hypotheses **Computing the Test Statistic** **Evaluating the Test Statistic**

Computing the Test Statistic

- Assuming the H_0 is true, we compute the test statistic:

$$Z_{obs} = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} = \frac{0.6625 - 0.5}{\sqrt{\frac{0.5(1-0.5)}{80}}} = 2.907$$

Diagram illustrating the components of the test statistic formula:

- Estimate of \hat{p}** : Points to the value 0.6625.
- Value of p_0** : Points to the value 0.5.
- Number of Observations**: Points to the value 80.
- Test Statistic**: Points to the result 2.907.

- The denominator is the standard error.
 - It uses p_0 under the assumption that the null hypothesis is true.

The test statistic takes a similar form to what you saw in previous sessions: the parameter estimate from the data, minus the value of the parameter under the null hypothesis, divided by the standard error. Since the test statistic is constructed assuming that the null hypothesis is true, p_0 is used in computing the standard error on the denominator.

We fill in the values: our estimate p hat is 0.6625, our value of p under the null hypothesis is 0.5 and n the number of observations in the sample is 80. Our test statistic value works out to be 2.907.

Tab 3: Evaluating the Test Statistic

Hypothesis Test for p : Soft Drink Example

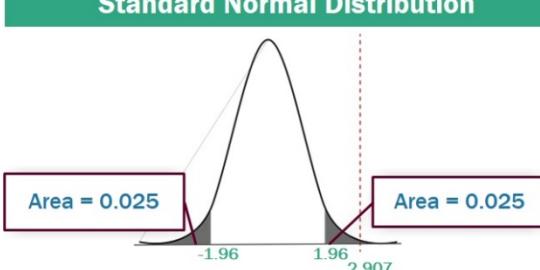
6 of 19

Stating the Hypotheses **Computing the Test Statistic** **Evaluating the Test Statistic**

Evaluating the Test Statistic

- The reference distribution for assessing Z_{obs} is the standard normal distribution.
- We use alpha = 0.05.
- As the test statistic of 2.907 is greater than 1.96, we reject the null hypothesis.
- We conclude that $p \neq 0.5$.
 - The true proportion differs from what would be expected by chance.

Standard Normal Distribution



?

If the null hypothesis is true, is the observed test statistic of 2.907 unusual?

If the null hypothesis is true, the test statistic is a random draw from a standard normal distribution. If the null hypothesis is true, is the observed test statistic of 2.907 unusual? Let's take a look.

Here we can see the standard normal distribution. Is our test statistic value unusual for this distribution? Our observed test statistic, or our observed z value, is equal to 2.907. What defines unusual? Is 2.907 unusual?

We will use alpha = 0.05 here, and since we are doing a two-sided test, we find the critical values such that the area under the curve outside our critical values sum to 0.05, with each one equalling $0.05/2 = 0.025$. So, the shaded grey area to the left of -1.96 is equal to 0.025 and the same for the area to the right of 1.96, it is also equal to 0.025. We say that any test statistic that lies outside these bounds is unusual for the distribution.

We have observed a test statistic of 2.907, which is greater than 1.96, and so we reject the null hypothesis and conclude that p does not equal 0.5. We can say that the true population probability differs from what would be expected based on chance alone.

Slide 7: Confidence Interval for a Single Proportion



Confidence Interval for a Single Proportion

7 of 19

- To construct a confidence interval for a single proportion p , use:

$$\hat{p} \pm z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}$$

- Here, the standard error uses \hat{p} since there is no assumed null hypothesis as there was in the hypothesis test.
- For the soft drink example, the 95% confidence interval is:

$$\begin{aligned} \hat{p} \pm z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}} &= 0.6625 \pm 1.96 \sqrt{\frac{0.6625(1 - 0.6625)}{80}} \\ &= (0.559, 0.766) \end{aligned}$$

- We are 95% confident that the true proportion p lies between 0.559 and 0.766.

We may also wish to construct a confidence interval for p , our population parameter of interest. The structure of the confidence interval is similar to what you saw in previous sessions: the parameter estimate from the data, plus or minus the critical value times the standard error.

We use \hat{p} to compute the standard error in the confidence interval, as opposed to p_0 for the hypothesis test standard error. This is because the test statistic is computed assuming the null hypothesis is true, but for a confidence interval, there is no null hypothesis assumption.

We plug in our estimates, p_{hat} equals 0.6625, $n = 80$ and our critical value is 1.96. Our interval goes from 0.559 up to 0.766.

We are 95% confident that the true population parameter p lies between 0.559 and 0.766. Remember, the interpretation of a confidence interval is always in relation to the population parameter of interest.

Slide 8: Analysing Single Proportions in R


Analysing Single Proportions in R
8 of 19

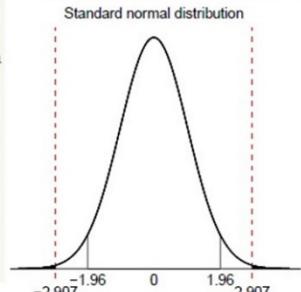
R Code: Hypothesis Test and Confidence Interval Construction

```

1 prop.test(x=53, n=80, p=0.5, correct = FALSE)
2
3     1-sample proportions test without continuity correction
4
5 data: 53 out of 80, null probability 0.5
6 X-squared = 8.45, df = 1, p-value = 0.00365
7 alternative hypothesis: true p is not equal to 0.5
8 95 percent confidence interval:
9 0.5535652 0.7565439
10 sample estimates:
11 p
12 0.6625

```

Standard Normal Distribution



This is the probability of observing this test statistic or a more extreme one, given that the null hypothesis is true.

We can do the hypothesis test and construct the confidence interval using the R code shown. We use the function `prop.test()` and specify the number of successes, the total number of observations and the value of p under the null hypothesis.

We see that our p-value is 0.00365 and this is the probability of observing our test statistic or a more extreme value, assuming that the null hypothesis is true. Graphically, 0.00365 is the area under the curve to the right of 2.907 plus the area under the curve to the left of -2.907. Effectively, we condition on the null hypothesis being true, and examine the probability that our observed test statistic or a more extreme value comes from the null distribution.

Slide 9: Final Remarks on Single Proportions


Final Remarks on Single Proportions
9 of 19

- With any statistical test or model:
 - Always be aware of the assumptions
 - Ensure these assumptions are reasonable

Two assumptions must hold for the hypothesis test and confidence interval to be true for the soft drink example.

1. The sample is a simple random sample from a large population.
2. The expected number of 'successes' and 'failures' is at least 10, that is:
 - $np \geq 10$ and $n(1 - p) \geq 10$

! Assumption validation is a key component of statistical analysis.

Should we carry out a one-sided or two-sided hypothesis test?

- Researchers had no prior belief that people differentiated between regular and diet versions, despite individuals believing they could.
 - Thus, a two-sided hypothesis test was conducted.

With any statistical test or model, there will be assumptions on which the validity of the inference from the analysis depends. It is a statistician's responsibility to firstly, know what these assumptions are, and secondly, to ensure they are reasonable.

Here, there are two assumptions to check for the validity of the hypothesis test and the confidence interval.

The data must be a simple random sample from a large population. It is crucial to determine how a study was carried out and assess its validity. In this case, we are told people were randomly sampled from the population of interest.

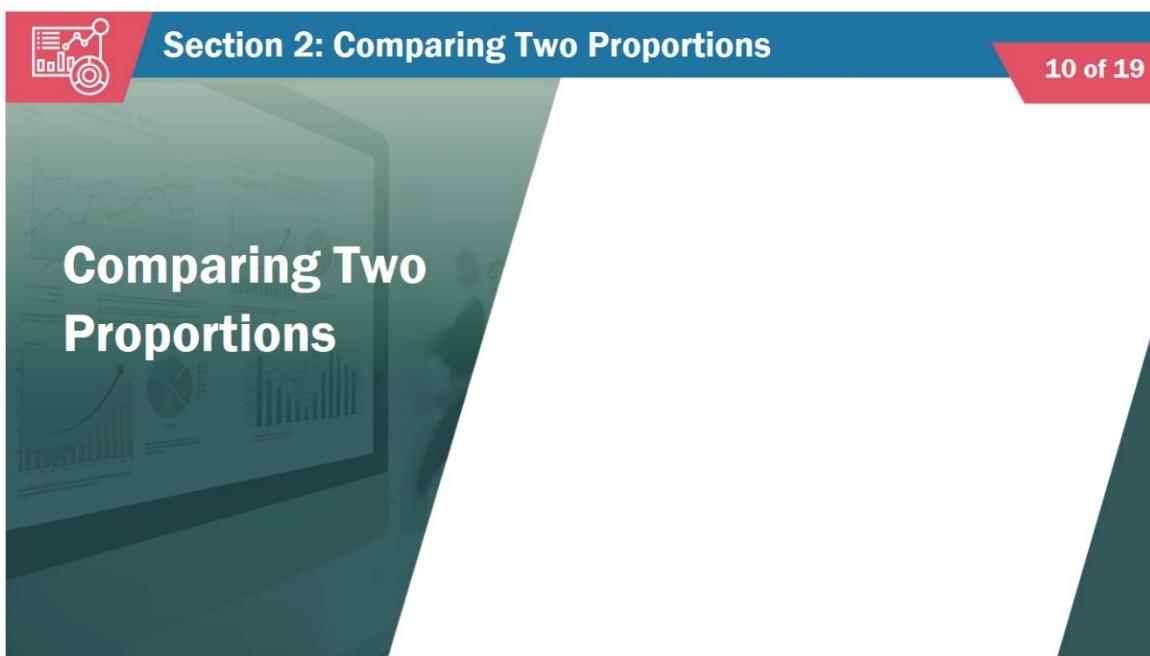
The expected number of successes and failures must be at least 10. We check this for the hypothesis test using p_0 and for the confidence interval using p_{hat} . For the hypothesis test, $np_0 = n(1 - p_0) = 80 * 0.5 = 40$, which is greater than 10. For the confidence interval, $n * p_{\text{hat}} = 53$ and $n * (1 - p_{\text{hat}}) = 27$, which are both greater than 10.

It is crucial that these assumptions are checked since the analysis results may not be valid if the assumptions are not met. Assumption validation is a key component of any statistical analysis.

When performing a hypothesis test, you must decide whether a one-sided test or two-sided test is best. The only reason to do a one-sided test is if the experimenters have an *a priori* belief that the outcome is in one particular direction. When analysing data, we never decide on the direction of the alternative hypothesis, or the value of the parameter under the null hypothesis, based on our observations in the data.

In the soft drink example, the researchers conducting the study had no prior belief that people could tell the difference between regular and diet soft drinks, despite the individuals believing they could. Thus, a two-sided hypothesis test was conducted.

Slide 10: **Section 2: Comparing Two Proportions**



Section 2: Comparing Two Proportions

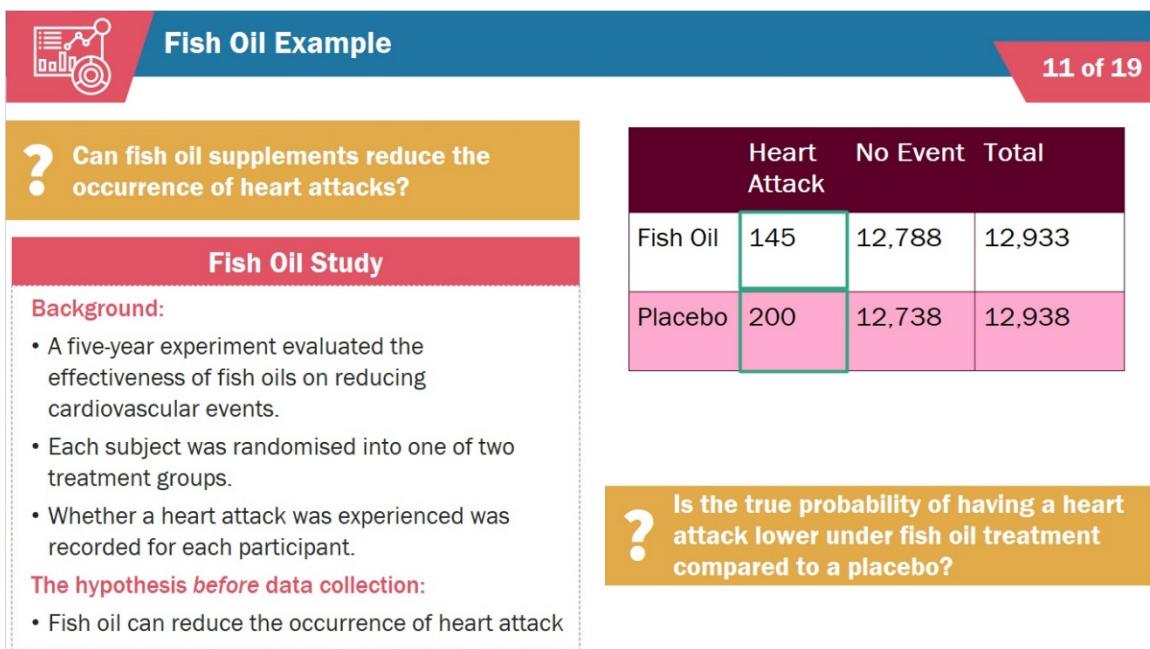
10 of 19

Comparing Two Proportions

This slide is titled "Comparing Two Proportions". It features a background image of a presentation screen displaying various charts and graphs. A red header bar at the top contains the title and page number. Below the header, the main title is displayed prominently.

We will now move on to comparing two proportions.

Slide 11: **Fish Oil Example**



Fish Oil Example

11 of 19

? Can fish oil supplements reduce the occurrence of heart attacks?

	Heart Attack	No Event	Total
Fish Oil	145	12,788	12,933
Placebo	200	12,738	12,938

Fish Oil Study

Background:

- A five-year experiment evaluated the effectiveness of fish oils on reducing cardiovascular events.
- Each subject was randomised into one of two treatment groups.
- Whether a heart attack was experienced was recorded for each participant.

The hypothesis before data collection:

- Fish oil can reduce the occurrence of heart attack

? Is the true probability of having a heart attack lower under fish oil treatment compared to a placebo?

This slide details a study on fish oil and heart attacks. It includes a question about the study's purpose, a 2x4 table of study results, and sections on study background and hypothesis. A separate section asks if the true probability of a heart attack is lower under fish oil treatment compared to a placebo.

Let's start with a motivating example: can fish oil supplements help reduce the occurrence of a heart attack? An experiment was carried to address this question.

Here are the details of the experiment:

A 5-year experiment was conducted to evaluate the effectiveness of fish oils on reducing cardiovascular events, where each subject was randomised into one of two treatment

groups. Whether a heart attack was experienced was recorded for each participant and shown in the following table. Before the data was collected, it was hypothesised that fish oil could reduce the occurrence of heart attack.

The data are displayed here in a contingency table. Participants are categorised according to whether they received the fish oil supplement or a placebo and whether or not they had a heart attack during the five years of the study. For example, there were 145 people who had heart attacks and were on the fish oil treatment, while there were 200 on the placebo treatment who had a heart attack. There were approximately 13,000 participants in each group.

We want to know if the probability of having a heart attack is lower under the fish oil supplement treatment compared to placebo.

Slide 12: Hypothesis Test for Comparing Two Proportions


Hypothesis Test for Comparing Two Proportions
12 of 19

Stating the Null Hypothesis
Computing the Test Statistic
Evaluating the Test Statistic
Conclusions

Introduction

- We will conduct a hypothesis test to compare the probability of having a heart attack under the fish oil treatment compared to a placebo.



 Click each tab to learn more. Then, click Next to continue.

We will conduct a hypothesis test to compare the probability of having a heart attack under the fish oil treatment compared to placebo. Once again, we will specify our hypotheses, compute our test statistic, evaluate the test statistic against the reference distribution, and give our conclusions.

Click each tab in turn to learn more. When you are ready, click next to continue.

Tab 1: Stating the Null Hypothesis

Hypothesis Test for Comparing Two Proportions 12 of 19

Stating the Null Hypothesis Computing the Test Statistic Evaluating the Test Statistic Conclusions

Stating the Null Hypothesis

- Let p_1 and p_2 be the true probability of heart attack under fish oil and placebo respectively.
- $H_0: p_1 - p_2 = 0$ versus $H_A: p_1 - p_2 < 0$
- We conduct a lower-tailed test due to the a priori belief that:
 - Fish oil would reduce the probability of a heart attack, compared to the placebo



We start off by defining our notation. We let p_1 be the probability of a heart attack under fish oil, and similarly p_2 is the probability under placebo. Under the null hypothesis, the two probabilities are the same, or their difference is equal to 0. The alternative hypothesis is that p_1 is lower than p_2 , or that $p_1 - p_2$ is less than 0. Remember, we were told in the example that it was hypothesised prior to doing the experiment that taking the fish oil would reduce the probability of a heart attack, compared to the placebo group. This *a priori* belief is the reason that we are doing a lower tailed test (as opposed to a two-sided test).

Tab 2: Computing the Test Statistic

Hypothesis Test for Comparing Two Proportions 12 of 19

Stating the Null Hypothesis Computing the Test Statistic Evaluating the Test Statistic Conclusions

Computing the Test Statistic

- The test statistic is:

?

Is this an unusual value for a standard normal distribution?

$$z_{obs} = \frac{(\hat{p}_1 - \hat{p}_2) - 0}{\sqrt{\hat{p}(1 - \hat{p})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}} = \frac{(0.01121 - 0.01546) - 0}{\sqrt{0.01334(1 - 0.01334)\left(\frac{1}{12933} + \frac{1}{12938}\right)}} = -2.977$$

Test Statistic Value

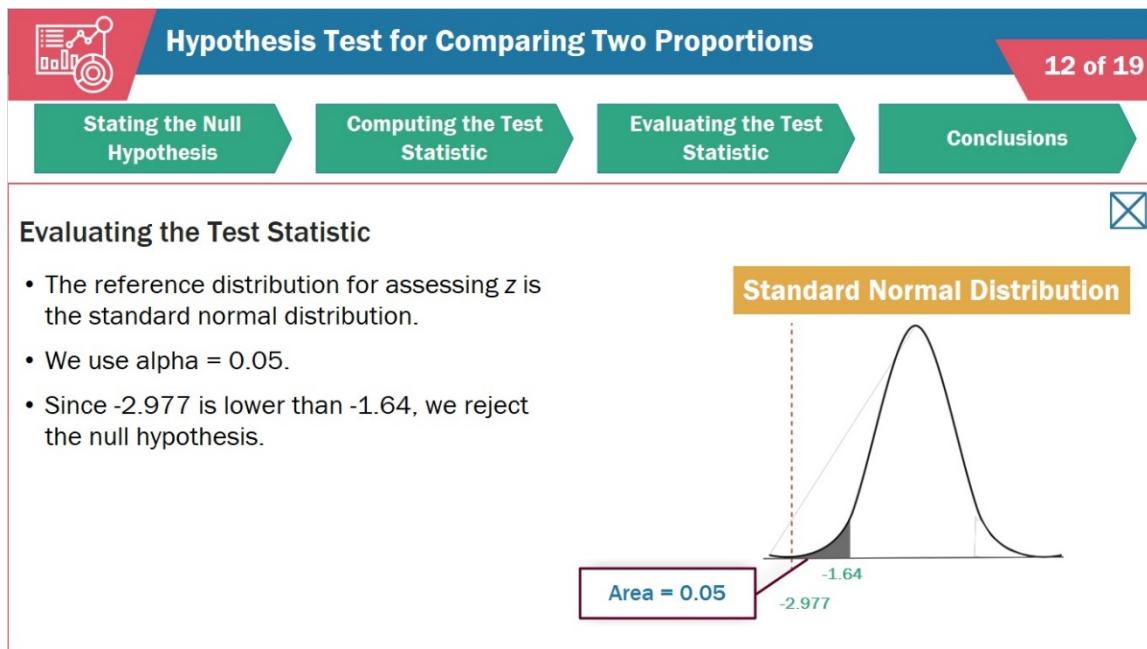
- The standard error is on the denominator.
 - \hat{p} is the pooled estimate of p , which is used because $p_1 = p_2$ under the null hypothesis.
- Under the null hypothesis the test statistic is a random draw from a standard normal distribution.

This is the formula for the test statistic. On the numerator, we have the difference between $p_1\hat{}$ and $p_2\hat{}$, which are the estimates of each probability from the data, minus 0, the difference under the null hypothesis. The standard error is on the denominator, where $p\hat{}$ is the pooled estimate of p , assuming that there is no difference between the two groups (which is what is assumed under the null hypothesis). n_1 and n_2 are the number of participants in the fish oil and placebo groups respectively.

The test statistics value works out to be -2.977.

If the null hypothesis is true, the test statistic is a random draw from a standard normal distribution. We have observed -2.977, and the question now is whether or not this is an unusual value for a standard normal distribution.

Tab 3: Evaluating the Test Statistic



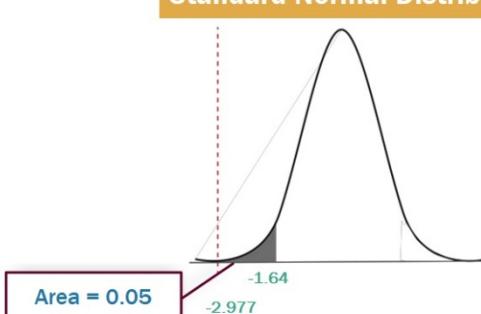
Hypothesis Test for Comparing Two Proportions

12 of 19

Evaluating the Test Statistic

- The reference distribution for assessing z is the standard normal distribution.
- We use $\alpha = 0.05$.
- Since -2.977 is lower than -1.64, we reject the null hypothesis.

Standard Normal Distribution



Area = 0.05

-1.64

-2.977

Let's examine the standard normal distribution. What you can see here, is the probability density function for the standard normal distribution. The area under the curve to the left of -1.64 has been marked out and this shaded grey area is equal to 0.05 or 5% of the total area under the curve. We say that the critical value for $\alpha = 0.05$ is equal to -1.64. In this case, any value lower than -1.64 is said to be unusual for the distribution, while anywhere else along the axis to the right of -1.64 is considered 'typical' for the distribution.

We have observed -2.977, which is considered unusual for this distribution and not consistent with the null hypothesis.

Since -2.977 is lower than -1.64, we reject the null hypothesis.

Tab 4: Conclusions

Hypothesis Test for Comparing Two Proportions

12 of 19

Stating the Null Hypothesis Computing the Test Statistic Evaluating the Test Statistic Conclusions

Conclusions

p₁ - p₂ is lower than 0:

- The true probability of having a heart attack is lower under the fish oil treatment, compared to placebo



We conclude that $p_1 - p_2$ is lower than 0, meaning that the true probability of having a heart attack is lower under the fish oil treatment, compared to placebo.

Slide 13: Confidence Interval for Comparing Two Proportions

Confidence Interval for Comparing Two Proportions

13 of 19

- To construct a confidence interval for the difference between two proportions:

$$\hat{p}_1 - \hat{p}_2 \pm z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{p}_1(1 - \hat{p}_1)}{n_1} + \frac{\hat{p}_2(1 - \hat{p}_2)}{n_2}}$$

- The standard error uses \hat{p}_1 and \hat{p}_2 , since there is no assumed null hypothesis that the two proportions are equal.
- For the fish oil example, the 90% confidence interval is:

$$\begin{aligned} &= 0.01546 - 0.01121 \pm 1.645 \sqrt{\frac{0.01546(1 - 0.01546)}{12933} + \frac{0.01121(1 - 0.01121)}{12938}} \\ &= (-0.0066, -0.0019) \end{aligned}$$

Interpretation

- We are 90% confident that the true probability of a heart attack under fish oil treatment minus the true probability under placebo treatment lies between -0.0066 and -0.0019.

Note

- 0 is not in the confidence interval.

Let's now construct a confidence interval for $p_1 - p_2$. The structure of the confidence interval is as before: estimate plus or minus the critical value, times the standard error. You'll see that this time we use p_1_{hat} and p_2_{hat} individually in the standard error, as opposed to p_{hat} the pooled estimate that we used in the hypothesis test. This is because in the hypothesis test, we assume the null hypothesis is true and if p_1 equals p_2 , the best estimate of it is the pooled estimate. For a confidence interval, there is no

null hypothesis, and so we use $p1_hat$ and $p2_hat$ separately and make no assumption that $p1$ and $p2$ are equal.

Going back to the fish example, we fill in the values: $p1_hat = 0.01546$, $p2_hat = 0.01121$, the critical value is 1.645, which is for alpha equal to $0.1/2 = 0.05$ for a 90% confidence interval, and $n1$ and $n2$ are 12933 and 12938 respectively. Working out the values, the 90% confidence interval goes from -0.0066 to -0.0019. What is the interpretation of this confidence interval?

We are 90% confident that the true probability of a heart attack under the fish oil treatment minus the true probability under the placebo treatment lies by between -0.0066 and -0.0019. We note that 0 is not in the confidence interval.

Slide 14: Comparing Two Proportions in R



Comparing Two Proportions in R

14 of 19

R Code (prop.test Function): Hypothesis Test for Comparing Two Proportions

```

1 prop.test(x = c(145, 200), n = c(12933, 12938), correct = FALSE,
2           alternative = "less")
3
4   2-sample test for equality of proportions without continuity correction
5
6 data: c(145, 200) out of c(12933, 12938)
7 X-squared = 8.8651, df = 1, p-value = 0.001453
8 alternative hypothesis: less
9 95 percent confidence interval:
10 -1.000000000 -0.001901128
11 sample estimates:
12 prop 1    prop 2
13 0.01121163 0.01545834

1 prop.test(x = c(145, 200), n = c(12933, 12938), correct = FALSE, conf.level = 0.9)
2
3   2-sample test for equality of proportions without continuity correction
4
5 data: c(145, 200) out of c(12933, 12938)
6 X-squared = 8.8651, df = 1, p-value = 0.002907
7 alternative hypothesis: two.sided
8 90 percent confidence interval:
9 -0.006592293 -0.001901128
10 sample estimates:
11 prop 1    prop 2
12 0.01121163 0.01545834

```

We can use R software to perform a hypothesis test for comparing two proportions. We use the `prop.test()` function. The `x` vector shows the number of heart attacks under the two treatments, while the `n` vector shows the number of participants in each group. In the first segment of code, a lower-tailed test is performed. The p-value for the test is 0.001453; this is the probability of observing this test statistic or more extreme, conditioning on the null hypothesis being true. Since this is less than $\alpha = 0.05$, we reject the null hypothesis and conclude that the true probability of having a heart attack is lower under the fish oil treatment, compared to placebo.

The second example uses a two-sided alternative, only so that we can see the full confidence interval.

Slide 15: Final Remarks on Comparing Two Proportions



Final Remarks on Comparing Two Proportions

15 of 19

The following three assumptions must hold for the validity of the hypothesis test and confidence interval to be valid:

1. Each sample must be a random sample from a large population
2. The two samples must be independent of each other
3. Within each group, both the expected number of 'successes' and 'failures' must be at least 10

We conclude that these conditions are all satisfied for the fish oil example.



The validity of the hypothesis test and confidence interval for comparing two proportions relies on three assumptions being valid. First, we need each sample to be a simple random sample from a large population. Second, we need the two samples to be independent of each other. Third, we need both the expected number of successes and failures to be at least 10.

These conditions are all satisfied for the fish oil example.

Slide 16: Section 3: Contingency Tables and Chi-square Tests



Section 3: Contingency Tables and Chi-square Tests

16 of 19

Contingency Tables and Chi-square Tests



Up until now, we have analysed binary variables, that is, categorical data variables with just two levels and we assessed the probability of a ‘success’. We will now explore how to analyse categorical data variables with more than two levels.

Slide 17: Chocolate Example



Chocolate Example

17 of 19

Chocolate Study		Age		
	Chocolate Consumption	18-21	22-25	26-29
(a)	Most days (6-7 days per week)	19	26	20
(b)	Regularly (2-5 days per week)	26	39	28
(c)	Once a week or less often	15	22	29

? Is there an association between age and chocolate consumption?

We'll start with a motivating example. A study was carried out to see if there was an association between consumption of chocolate and the age of adults across the range 18 to 29. Each participant was categorised by their age: 18-21, 22-25 or 26-29, and by their chocolate consumption: (a) Most days (6-7 days per week), (b) Regularly (2-5 days per week) or (c) Once a week or less often.

The data from the study is displayed here in a two-way contingency. For example, there were 19 participants in the 18-21 age category that eat chocolate most days.

Here we have two categorical variables: chocolate consumption and age, each with $k = 3$ levels. We want to know if there is an association between age and chocolate consumption.



Slide 18: Chocolate Example

Chocolate Example

18 of 19

- Expressing the Hypotheses
- Finding the Expected Values
- Computing the Test Statistic
- Identifying the Degrees of Freedom
- Evaluating the Test Statistic
- Conclusions
- Chi-square Test for Independence in R
- Assumptions

Introduction

- We will test for an association between age and chocolate consumption using a chi-square test for independence.



Click each tab to learn more. Then, click Next to continue.

We will test for an association between age and chocolate consumption using a chi-square test for independence. We want to know, for example, if as age increases does chocolate consumption change. We will once again specify our hypotheses, compute our test statistic, evaluate the test statistic against the reference distribution, and give our conclusions.

Click each tab to learn more. When you are ready, click next to continue.

Tab 1: Expressing the Hypotheses

Chocolate Example

18 of 19

- Expressing the Hypotheses
- Finding the Expected Values
- Computing the Test Statistic
- Identifying the Degrees of Freedom
- Evaluating the Test Statistic
- Conclusions
- Chi-square Test for Independence in R
- Assumptions

Expressing the Hypotheses

Hypotheses

- The hypotheses are:

- H_0 : Age and chocolate consumption are independent

Versus

- H_A : Age and chocolate consumption are not independent

Let's start by expressing the hypotheses. The null hypothesis is that age and chocolate consumption are independent of each other, and the alternative is that they are not independent.

Tab 2: Finding the Expected Values


Chocolate Example
18 of 19

Expressing the Hypotheses Finding the Expected Values Computing the Test Statistic Identifying the Degrees of Freedom Evaluating the Test Statistic Conclusions Chi-square Test for Independence in R Assumptions	<h3>Finding the Expected Values</h3> <p>Expected Values for Each Cell</p> <ul style="list-style-type: none"> We find the expected values for each cell by computing the row total by the column total, divided by the overall total. <table border="1" style="width: 100%; border-collapse: collapse; text-align: center;"> <thead> <tr> <th colspan="2"></th> <th colspan="3">Age</th> </tr> <tr> <th colspan="2">Chocolate consumption</th> <th>18-21</th> <th>22-25</th> <th>26-29</th> </tr> </thead> <tbody> <tr> <td>(a)</td> <td>Most days (6-7 days per week)</td> <td>19</td> <td>26</td> <td>20</td> </tr> <tr> <td>(b)</td> <td>Regularly (2-5 days per week)</td> <td>26</td> <td>39</td> <td>28</td> </tr> <tr> <td>(c)</td> <td>Once a week or less often</td> <td>15</td> <td>22</td> <td>29</td> </tr> <tr> <td>Total</td> <td></td> <td>60</td> <td>87</td> <td>77</td> </tr> <tr> <td></td> <td></td> <td></td> <td></td> <td>225</td> </tr> </tbody> </table> <table border="1" style="width: 100%; border-collapse: collapse; text-align: center;"> <thead> <tr> <th colspan="2"></th> <th colspan="3">Age</th> </tr> <tr> <th colspan="2">Chocolate Consumption</th> <th>18-21</th> <th>22-25</th> <th>26-29</th> </tr> </thead> <tbody> <tr> <td>(a)</td> <td>Most days (6-7 days per week)</td> <td>17.411</td> <td>25.246</td> <td>22.344</td> </tr> <tr> <td>(b)</td> <td>Regularly (2-5 days per week)</td> <td>24.911</td> <td>36.121</td> <td>31.969</td> </tr> <tr> <td>(c)</td> <td>Once a week or less often</td> <td>17.679</td> <td>25.634</td> <td>22.688</td> </tr> </tbody> </table>			Age			Chocolate consumption		18-21	22-25	26-29	(a)	Most days (6-7 days per week)	19	26	20	(b)	Regularly (2-5 days per week)	26	39	28	(c)	Once a week or less often	15	22	29	Total		60	87	77					225			Age			Chocolate Consumption		18-21	22-25	26-29	(a)	Most days (6-7 days per week)	17.411	25.246	22.344	(b)	Regularly (2-5 days per week)	24.911	36.121	31.969	(c)	Once a week or less often	17.679	25.634	22.688
		Age																																																											
Chocolate consumption		18-21	22-25	26-29																																																									
(a)	Most days (6-7 days per week)	19	26	20																																																									
(b)	Regularly (2-5 days per week)	26	39	28																																																									
(c)	Once a week or less often	15	22	29																																																									
Total		60	87	77																																																									
				225																																																									
		Age																																																											
Chocolate Consumption		18-21	22-25	26-29																																																									
(a)	Most days (6-7 days per week)	17.411	25.246	22.344																																																									
(b)	Regularly (2-5 days per week)	24.911	36.121	31.969																																																									
(c)	Once a week or less often	17.679	25.634	22.688																																																									

Before we can compute our test statistic, we need to find expected values for each cell, whereby expected value we mean the expected value assuming that the null hypothesis is true. We do this by multiplying the row total by the column total and dividing by the overall total, where each total is from the raw data table. For example, to find the expected value in the first cell we multiply 65 by 60 and divide by 225 giving 17.411.

Here are the expected values for each cell. Each value is computed in a similar way from the raw data table: row total by column total divided by overall total.

Tab 3: Computing the Test Statistic

Chocolate Example

18 of 19

Expressing the Hypotheses

Finding the Expected Values

Computing the Test Statistic

Identifying the Degrees of Freedom

Evaluating the Test Statistic

Conclusions

Chi-square Test for Independence in R

Assumptions

Computing the Test Statistic

- The test statistic is where the summation is over all cells in the table.

$$\chi_{obs} = \sum \frac{(\text{Observed} - \text{Expected})^2}{\text{Expected}}$$

$$\chi_{obs} = \frac{(19 - 17.411)^2}{17.411} + \frac{(26 - 25.246)^2}{25.246} + \frac{(20 - 22.344)^2}{22.344} + \dots + \frac{(29 - 22.688)^2}{22.688} = 3.86.$$

		Age		
Chocolate consumption		18-21	22-25	26-29
(a)	Most days (6-7 days per week)	19	26	20
(b)	Regularly (2-5 days per week)	26	39	28
(c)	Once a week or less often	15	22	29
Total		60	87	77

		Age		
Chocolate Consumption		18-21	22-25	26-29
(a)	Most days (6-7 days per week)	17.411	25.246	22.344
(b)	Regularly (2-5 days per week)	24.911	36.121	31.969
(c)	Once a week or less often	17.679	25.634	22.688

The test statistic is calculated by computing observed value minus expected value all to be squared, divided by expected for each cell, and then summing over all cells. Observed counts are the raw observed data counts. ‘Expected values’ are the values that would be expected under the null hypothesis, that is, assuming that the null hypothesis is true.

Going back to our example; in the first cell the observed cell count is 19 and the expected value is 17.411. We compute $19 - 17.411$, square it, and divide by 17.411. For the second cell, it is $26 - 25.246$ squared, divided by 25.246, and we continue on for each cell down to the last cell which is $29 - 22.688$ squared, divided by 22.688. When we add up the values for each cell, we get 3.86.

Remember that the expected values are what is expected under the assumption that the null hypothesis is true. The numerator of each cell calculation in the test statistic is effectively measuring how far each observed cell count is from what is expected under the null hypothesis. If the data are consistent with the null hypothesis, we would expect a small test statistic value, because observed and expected will be close together; and if the data depart from what is expected under the null hypothesis, then the test statistic will be large.

Tab 4: Identifying the Degrees of Freedom

Chocolate Example

18 of 19

Expressing the Hypotheses	Identifying the Degrees of Freedom
Finding the Expected Values	Degrees of Freedom
Computing the Test Statistic	<ul style="list-style-type: none"> The degrees of freedom (df) associated with this test are computed as: (# rows - 1) x (# columns - 1)
Identifying the Degrees of Freedom	<ul style="list-style-type: none"> For the chocolate example: (3 - 1) x (3 - 1) = 4 df
Evaluating the Test Statistic	
Conclusions	
Chi-square Test for Independence in R	
Assumptions	

Before we can evaluate the test statistic, we must identify the degrees of freedom. The degrees of freedom (df) associated with this test are computed as the # rows – 1 times the # columns – 1. In this example, there are 3 – 1 times 3-1 equal to 2 by 2 equals to 4 degrees of freedom.

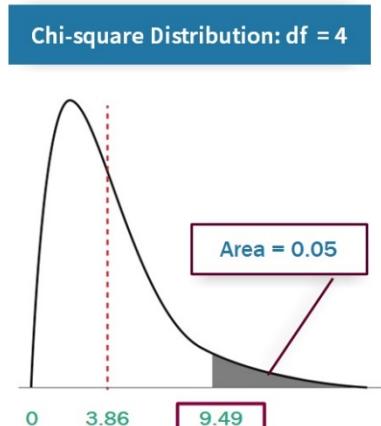
Tab 5: Evaluating the Test Statistic

Chocolate Example

18 of 19

Expressing the Hypotheses	Evaluating the Test Statistic
Finding the Expected Values	Chi-square Distribution
Computing the Test Statistic	<ul style="list-style-type: none"> The reference distribution for assessing the observed test statistic is the chi-square distribution. Its shape is determined by the df. Using alpha = 0.05: We fail to reject the null hypothesis
Identifying the Degrees of Freedom	
Evaluating the Test Statistic	
Conclusions	
Chi-square Test for Independence in R	
Assumptions	

Chi-square Distribution: df = 4



The reference distribution for assessing the observed test statistic is the chi-square distribution, the shape of which is determine by the degrees of freedom. In our example, we have four degrees of freedom and here we can see the shape of the chi-square distribution with 4 degrees of freedom. The value 9.49 marked in on the graph is the

critical value such that the area under the curve to the right of it is equal to 0.05. So, the area in grey is 5% of the total area under the curve.

We use this cut-off of 9.49 to determine if our observed test statistic is extreme or not, anything larger than the cut-off is considered extreme for the distribution. Since our test statistic is equal to 3.86, which is less than 9.49, it is not extreme for this distribution, using alpha = 0.05, we fail to reject the null hypothesis.

Tab 6: Conclusions

📊 Chocolate Example 18 of 19

Expressing the Hypotheses Finding the Expected Values Computing the Test Statistic Identifying the Degrees of Freedom Evaluating the Test Statistic Conclusions Chi-square Test for Independence in R Assumptions	<h3>Conclusions</h3> <p>Chi-square Distribution</p> <ul style="list-style-type: none"> • There is no evidence of an association between age and chocolate consumption. <div style="text-align: center; margin-top: 20px;">  ≠  </div>
--	--

Tab 7: Chi-square Test for Independence in R

Slide 19:



Chocolate Example

- [Expressing the Hypotheses](#)
- [Finding the Expected Values](#)
- [Computing the Test Statistic](#)
- [Identifying the Degrees of Freedom](#)
- [Evaluating the Test Statistic](#)
- [Conclusions](#)
- [Chi-square Test for Independence in R](#)
- [Assumptions](#)

Chi-square Test for Independence in R

Analysis in R

- As the p-value of 0.4252 is greater than 0.05,
 - We fail to reject the null hypothesis
- We conclude that there is no evidence of an association between age and chocolate consumption.

```

1 # Create the table of data
2 choc <- matrix(c(19,26,20,26,39,28,15,22,29), ncol = 3, byrow = TRUE)
3 rownames(choc) <- c("a","b","c")
4 colnames(choc) <- c("18-21","22-25","26-29")
5 choc <- as.table(choc)
```

```

6
7 # Perform the chi-square test
8 test <- chisq.test(choc)
9 test # View the test output
```

```

10 test$observed # observed counts (same as choc)
11 test$expected # expected counts under the null hypothesis
12
13 Pearson's Chi-squared test
14
15 data: choc
16 X-squared = 3.8607, df = 4, p-value = 0.4252
17
18 > test$observed
19   18-21 22-25 26-29
20 a    19    26    20
21 b    26    39    28
22 c    15    22    29
23
24 > test$expected
25   18-21 22-25 26-29
26 a 17.41071 25.24554 22.34379
27 b 24.91071 36.12054 31.96875
28 c 17.67857 25.63393 22.68750
```

We can do this analysis in R. First, we set up a table in R containing the data. We can label the rows and columns appropriately using the row names and col names functions. The first four lines of code are for reading in the data. We can then carry out a chi-square test for independence using the chisq.test function.

In the output, we can see the value of our test statistic equals 3.8607 with p-value 0.4252. The p-value is the probability of observing this test statistic or a more extreme one, assuming that the null hypothesis is true. Since the p-value 0.4252 is greater than 0.05, we fail to reject the null hypothesis, and conclude that we have no evidence of an association between age and chocolate consumption.

Tab 1: Assumptions


Chocolate Example
18 of 19

Expressing the Hypotheses Finding the Expected Values Computing the Test Statistic Identifying the Degrees of Freedom Evaluating the Test Statistic Conclusions Chi-square Test for Independence in R Assumptions	<div style="display: flex; justify-content: space-between;"> Assumptions <input type="checkbox"/> </div> <div style="border: 1px dashed #ccc; padding: 5px; margin-top: 5px;"> Final Remarks <p style="color: #FFB703; margin: 0;">For the chi-square test for independence:</p> <ul style="list-style-type: none"> • All expected cell count should be at least five. • If they are not, the validity of the test may be questionable. <p style="color: #FFB703; margin: 0;">For the chocolate example, all expected counts are above five.</p> <div style="border: 1px dashed #ccc; padding: 5px; margin-top: 10px;"> A chi-square test: <ul style="list-style-type: none"> • Is useful <ul style="list-style-type: none"> • It can test whether two categorical variables are dependent. • Is limited <ul style="list-style-type: none"> • It may show an association between two categorical values, but not tell us anything about the nature of that association </div> </div>
--	---

For the chi-square test for independence, all expected cell counts should be at least 5, if they are not, the validity of the test may be questionable. In the chocolate example, the lowest expected cell count is 17.411 and so this condition is satisfied.

We have seen how a chi-square test can test whether or not two categorical variables are independent; however, there is a limitation to this test. If the result of the test identifies an association between two categorical variables, it will not tell us anything about the nature of that association. But it is a useful starting point for jointly analysing two categorical variables.

Slide 20: **Summary**



Summary

19 of 19

- Having completed this presentation, you should now be able to:
 - Perform a hypothesis test for a single proportion, p , to compare two proportions, $p_1 - p_2$, and to test for independence between two categorical variables
 - Construct and interpret a confidence interval for a single proportion and for comparing two proportions
 - Identify when these analyses are reasonable for your data and implement them in R software
 - Outline the assumptions and limitations associated with these analyses



Developed by Trinity Online Services CLG with the School of Computer Science and Statistics, Trinity College Dublin, The University of Dublin

Having completed this presentation, you should now be able to:

- Perform a hypothesis test for a single proportion, p , to compare two proportions, $p_1 - p_2$, and to test for independence between two categorical variables.
- Construct and interpret a confidence interval for a single proportion and for comparing two proportions.
- Identify when these analyses are reasonable for your data and implement them in R software.
- Outline the assumptions and limitations associated with these analyses.