

# Chapter 11: $\chi^2$ Tests

DSCC 462  
Computational Introduction to Statistics

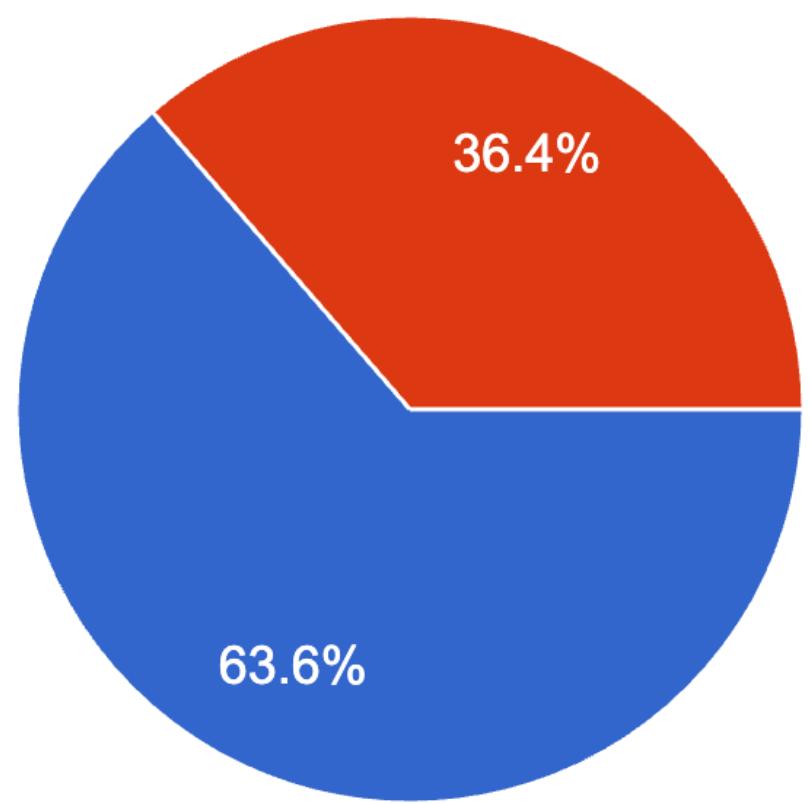
Anson Kahng  
Fall 2022

# Survey Feedback ( $n = 44$ out of 54 so far)

The material is clear

The instructor presents concepts clearly.

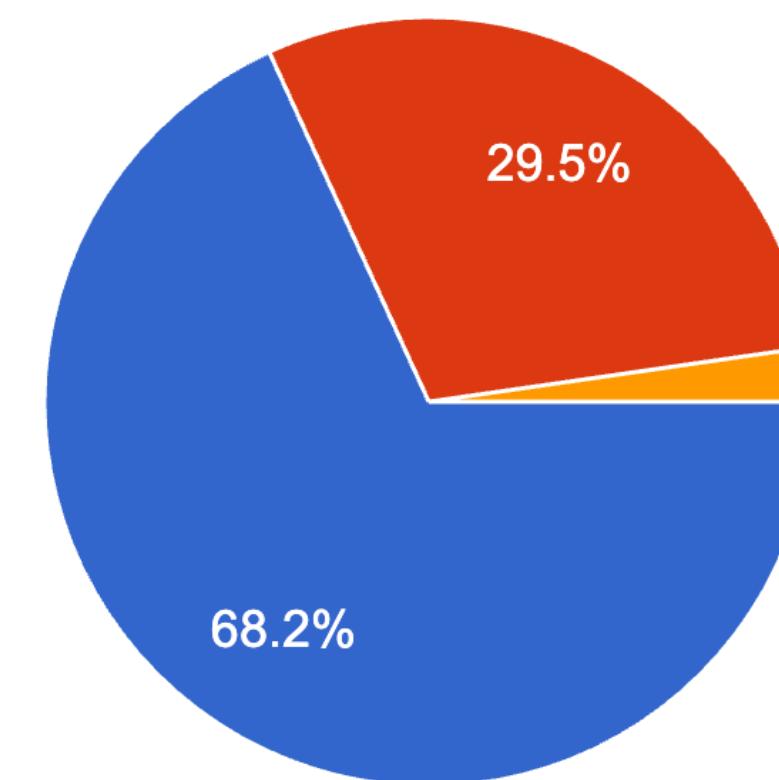
44 responses



The assignments are useful

The assignments help me apply concepts that we learn

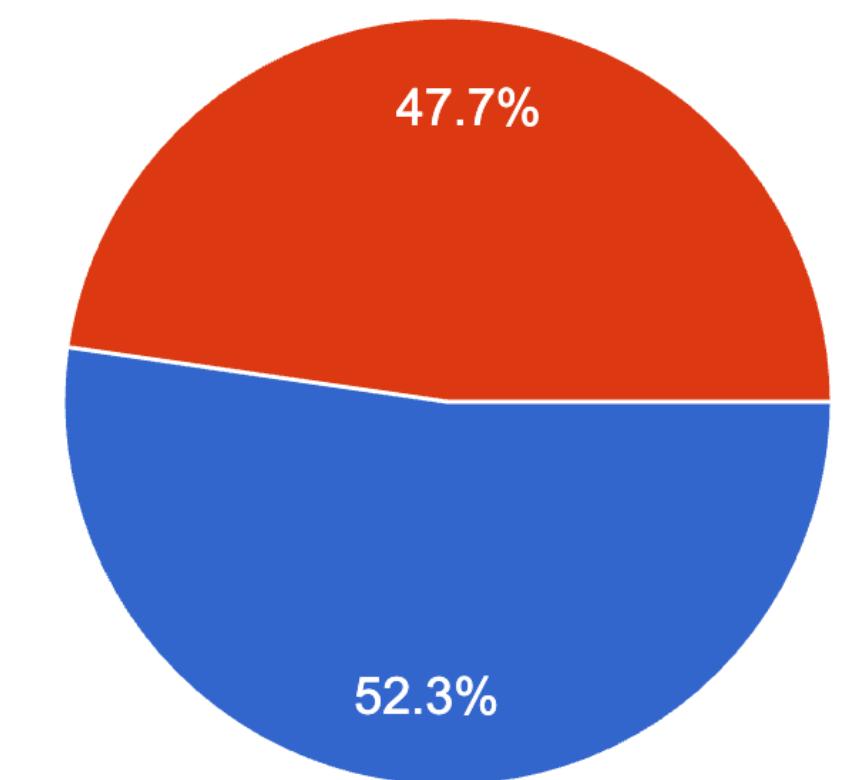
44 responses



You are learning!

I am learning a lot in this class.

44 responses



- Strongly agree
- Agree
- Neutral
- Disagree
- Strongly disagree

# Survey Feedback (n = 44 out of 54 so far)

Some of you mentioned that...

What we'll do going forward

# Survey Feedback (n = 44 out of 54 so far)

START

Some of you mentioned that...

- You like **outlines / summaries** of lecture
- You find **examples** helpful
- You want more **feedback** on homeworks
- You want to see more **R code** on slides

What we'll do going forward

+ *instructor slides!*

- I will add **summary slides** to lecture
- I will use **examples** to frame concepts
- TAs will give more **detailed feedback** on assignments
- I will add **R code** to slides

# Survey Feedback (n = 44 out of 54 so far)

Some of you mentioned that...

**START**

- You like **outlines / summaries** of lecture
- You find **examples** helpful
- You want more **feedback** on homeworks
- You want to see more **R code** on slides

What we'll do going forward

- I will add **summary slides** to lecture
- I will use **examples** to frame concepts
- TAs will give more **detailed feedback** on assignments
- I will add **R code** to slides

**STOP**

- Homework questions are **repetitive**
- Homeworks are slightly **long**

- I can't promise homeworks will be shorter, but I'll make sure the questions **cover more concepts** 😊

# Survey Feedback (n = 44 out of 54 so far)

Some of you mentioned that...

**START**

- You like **outlines / summaries** of lecture
- You find **examples** helpful
- You want more **feedback** on homeworks
- You want to see more **R code** on slides

What we'll do going forward

- I will add **summary slides** to lecture
- I will use **examples** to frame concepts
- TAs will give more **detailed feedback** on assignments
- I will add **R code** to slides

**STOP**

- Homework questions are **repetitive**
- Homeworks are slightly **long**

- I can't promise homeworks will be shorter, but I'll make sure the questions **cover more concepts** 😊

**CONTINUE**

- Being **responsive** to students' needs
- Seeking **feedback** from students
- **Answering questions** during class
- **Solving examples** in class

- I'll **continue doing these things!**

# Key Reminders

# Key Reminders

- I have three office hours each week! Only a few of you have used them regularly

# Key Reminders

- I have three office hours each week! Only a few of you have used them regularly
- These are **by appointment** so please email me ahead of time

# Key Reminders

- I have three office hours each week! Only a few of you have used them regularly
  - These are **by appointment** so please email me ahead of time
  - If you cannot make them, **email me so we can find another time!**

# Key Reminders

- I have three office hours each week! Only a few of you have used them regularly
  - These are **by appointment** so please email me ahead of time
  - If you cannot make them, **email me so we can find another time!**
- The TAs are also here to support you

# Key Reminders

- I have three office hours each week! Only a few of you have used them regularly
  - These are **by appointment** so please email me ahead of time
  - If you cannot make them, **email me so we can find another time!**
- The TAs are also here to support you
  - Learning from a group can help; be sure to prepare questions in advance

# Key Reminders

- I have three office hours each week! Only a few of you have used them regularly
  - These are **by appointment** so please email me ahead of time
  - If you cannot make them, **email me so we can find another time!**
- The TAs are also here to support you
  - Learning from a group can help; be sure to prepare questions in advance
- Educational videos:

# Key Reminders

- I have three office hours each week! Only a few of you have used them regularly
  - These are **by appointment** so please email me ahead of time
  - If you cannot make them, **email me so we can find another time!**
- The TAs are also here to support you
  - Learning from a group can help; be sure to prepare questions in advance
- Educational videos:
  - YouTube: zedstatistics, 3blue1brown, jbstatistics

# Community Norms

# Community Norms

- Let's keep side conversations to a minimum

# Community Norms

- Let's keep side conversations to a minimum
- Let's stay awake

# Lecture Plan for Today

# Lecture Plan for Today

- Goodness-of-Fit Test

# Lecture Plan for Today

- Goodness-of-Fit Test
  - True proportion = expected proportion?

# Lecture Plan for Today

- Goodness-of-Fit Test
  - True proportion = expected proportion?
  - Generalization of proportion hypothesis tests

# Lecture Plan for Today

- Goodness-of-Fit Test
  - True proportion = expected proportion?
  - Generalization of proportion hypothesis tests
- Chi-Squared ( $\chi^2$ ) Test of Independence

# Lecture Plan for Today

- Goodness-of-Fit Test
  - True proportion = expected proportion?
  - Generalization of proportion hypothesis tests
- Chi-Squared ( $\chi^2$ ) Test of Independence
  - Are variables related or not?

# Multi-Category Proportions

# Multi-Category Proportions

- Last lecture, we looked at inference for proportions

# Multi-Category Proportions

- Last lecture, we looked at inference for proportions
  - In this setting, a variable could take on one of two values

# Multi-Category Proportions

- Last lecture, we looked at inference for proportions
  - In this setting, a variable could take on one of two values
  - What if the variable had more categories?

# Multi-Category Proportions

- Last lecture, we looked at inference for proportions
  - In this setting, a variable could take on one of two values
  - What if the variable had more categories?
- Example: Let's say we are trying to figure out what proportion of people have each season (winter, spring, summer, fall) as their favorite. How can we do inference in this setting?

# Multi-Category Proportions

- Last lecture, we looked at inference for proportions
  - In this setting, a variable could take on one of two values
  - What if the variable had more categories?
- Example: Let's say we are trying to figure out what proportion of people have each season (winter, spring, summer, fall) as their favorite. How can we do inference in this setting?
  - Instead of one  $p$ , we have to infer values for  $p_1, p_2, p_3$

# Goodness-of-Fit

# Goodness-of-Fit

- Consider a categorical variable with multiple categories

# Goodness-of-Fit

- Consider a categorical variable with multiple categories
  - E.g., eye color: brown, hazel, blue, other

# Goodness-of-Fit

- Consider a categorical variable with multiple categories
  - E.g., eye color: brown, hazel, blue, other
- Perhaps we want to test whether the true proportion of people falling into each category is equal to some value

# Goodness-of-Fit

- Consider a categorical variable with multiple categories
  - E.g., eye color: brown, hazel, blue, other
- Perhaps we want to test whether the true proportion of people falling into each category is equal to some value
- Use the Goodness-of-Fit Test

# Goodness-of-Fit

# Goodness-of-Fit

- Let  $p_i$  be the true proportion of the population that falls into category  $i$

# Goodness-of-Fit

- Let  $p_i$  be the true proportion of the population that falls into category  $i$ 
  - $i = 1, \dots, k$ , where  $k$  is the number of categories

# Goodness-of-Fit

- Let  $p_i$  be the true proportion of the population that falls into category  $i$ 
  - $i = 1, \dots, k$ , where  $k$  is the number of categories
  - Note that  $\sum_{i=1}^k p_i = 1$

# Goodness-of-Fit

- Let  $p_i$  be the true proportion of the population that falls into category  $i$ 
  - $i = 1, \dots, k$ , where  $k$  is the number of categories
- Note that  $\sum_{i=1}^k p_i = 1$
- Our hypotheses are as follows:

# Goodness-of-Fit

- Let  $p_i$  be the true proportion of the population that falls into category  $i$ 
  - $i = 1, \dots, k$ , where  $k$  is the number of categories
  - Note that  $\sum_{i=1}^k p_i = 1$
- Our hypotheses are as follows:
  - $H_0 : p_1 = p_{1_0}, p_2 = p_{2_0}, \dots, p_k = p_{k_0}$

# Goodness-of-Fit

- Let  $p_i$  be the true proportion of the population that falls into category  $i$ 
  - $i = 1, \dots, k$ , where  $k$  is the number of categories
  - Note that  $\sum_{i=1}^k p_i = 1$
- Our hypotheses are as follows:
  - $H_0 : p_1 = p_{1_0}, p_2 = p_{2_0}, \dots, p_k = p_{k_0}$
  - $H_1$  : at least one of these equalities does not hold

$$\sum_{i=1}^k p_{i_0} = 1$$

# Goodness-of-Fit

# Goodness-of-Fit

- Test the hypothesis using a chi-squared ( $\chi^2$ ) test

# Goodness-of-Fit

- Test the hypothesis using a chi-squared ( $\chi^2$ ) test
- The test statistic is

# Goodness-of-Fit

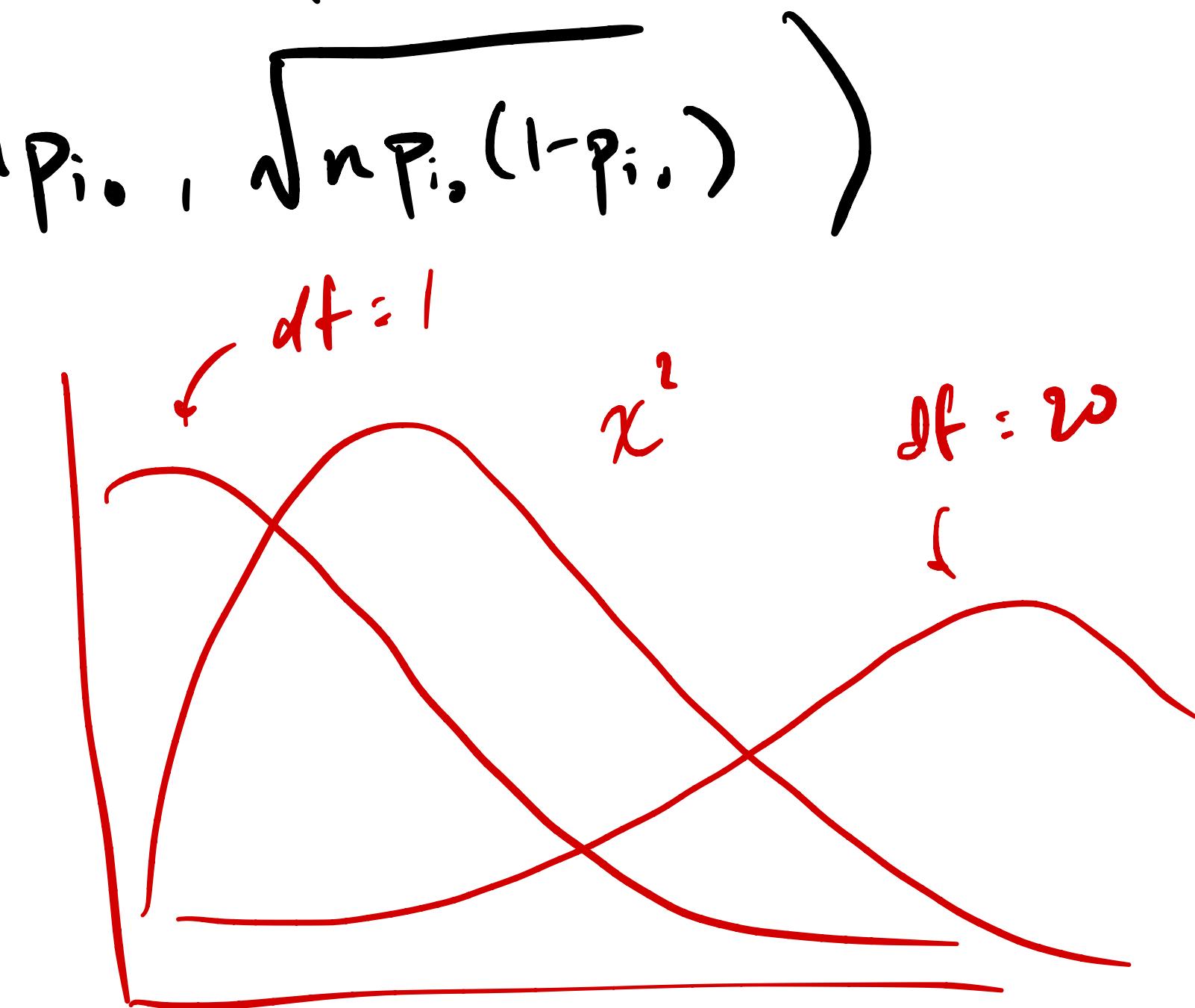
- Test the hypothesis using a chi-squared ( $\chi^2$ ) test
- The test statistic is

$$X^2 = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i} \sim \underline{\chi^2_{k-1}}$$

$\mu_{X_1}$

$$\boxed{(1-p_{i0})Z^2} = \left( \frac{X_i - np_{i0}}{\sqrt{np_{i0}(1-p_{i0})}} \right)^2 = \frac{(X_i - E_i)^2}{E_i}$$

$$X_i \sim \text{Binom}(n, p_{i0}) \\ \simeq N(np_{i0}, \sqrt{np_{i0}(1-p_{i0})})$$



# Goodness-of-Fit

- Test the hypothesis using a chi-squared ( $\chi^2$ ) test
- The test statistic is

$$X^2 = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i} \sim \chi^2_{k-1}$$

- Here,  $O_i$  is the observed number of people who fall in category  $i$

# Goodness-of-Fit

- Test the hypothesis using a chi-squared ( $\chi^2$ ) test
- The test statistic is

$$X^2 = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i} \sim \chi^2_{k-1}$$

- Here,  $O_i$  is the observed number of people who fall in category  $i$
- $E_i$  is the expected number of people who fall into category  $i$  under the null hypothesis

# Goodness-of-Fit

# Goodness-of-Fit

- We calculate the test statistic using the following information

# Goodness-of-Fit

- We calculate the test statistic using the following information
- Recall that the test statistic is  $X^2 = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i}$

# Goodness-of-Fit

- We calculate the test statistic using the following information

- Recall that the test statistic is  $X^2 = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i}$

$p_{1o}$

Category	1	2	...	k
Observed	$O_1$	$O_2$	...	$O_k$
Expected	$E_1 = n^*p_{1o}$	$E_2 = n^*p_{2o}$	...	$E_k = n^*p_{ko}$

# Goodness-of-Fit

# Goodness-of-Fit

- We are interested in the p-value of  $\Pr(\chi^2 > X^2)$

# Goodness-of-Fit

- We are interested in the p-value of  $\Pr(\chi^2 > X^2)$
- We can find this p-value by using a  $\chi^2$  distribution with  $df = k - 1$

# Goodness-of-Fit

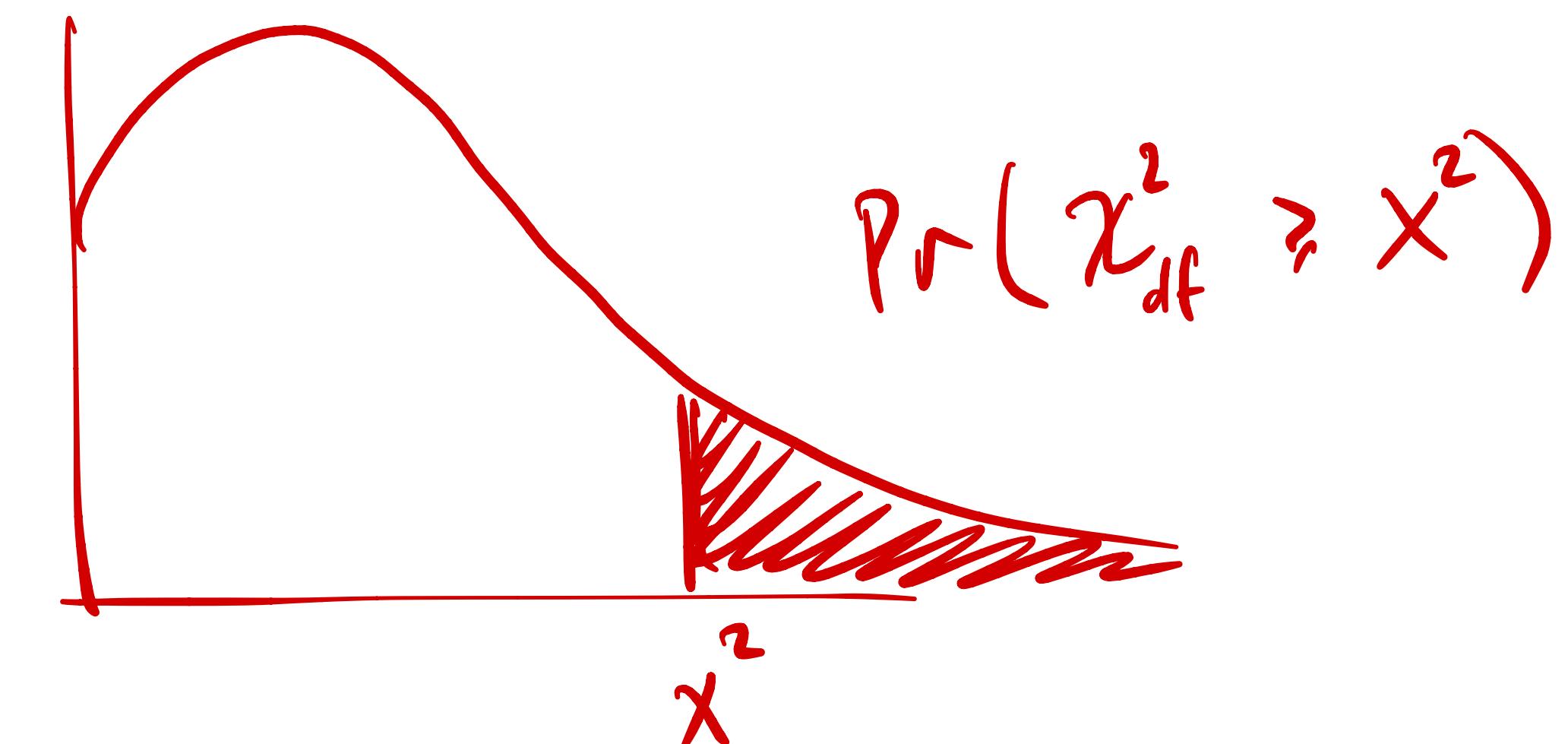
- We are interested in the p-value of  $\Pr(\chi^2 > X^2)$
- We can find this p-value by using a  $\chi^2$  distribution with  $df = k - 1$ 
  - In R:  $p = 1 - \text{pchisq}(X^2, df)$

# Goodness-of-Fit

$$\chi^2 = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i}$$

- We are interested in the p-value of  $\Pr(\chi^2 > X^2)$
- We can find this p-value by using a  $\chi^2$  distribution with  $df = k - 1$ 
  - In R:  $p = 1 - \text{pchisq}(X^2, df)$
  - Always looking for *upper tail probability* (probability of seeing an outcome that is as extreme or more extreme than what we observed)

$$\chi^2 \sim \chi^2_{k-1}$$



# Goodness-of-Fit

- We are interested in the p-value of  $\Pr(\chi^2 > X^2)$
- We can find this p-value by using a  $\chi^2$  distribution with  $df = k - 1$ 
  - In R:  $p = 1 - \text{pchisq}(X^2, df)$
  - Always looking for *upper tail probability* (probability of seeing an outcome that is as extreme or more extreme than what we observed)
- If  $p \leq \alpha$ , we reject  $H_0$

# Goodness-of-Fit

- We are interested in the p-value of  $\Pr(\chi^2 > X^2)$
- We can find this p-value by using a  $\chi^2$  distribution with  $df = k - 1$ 
  - In R:  $p = 1 - \text{pchisq}(X^2, df)$
  - Always looking for *upper tail probability* (probability of seeing an outcome that is as extreme or more extreme than what we observed)
- If  $p \leq \alpha$ , we reject  $H_0$
- If  $p > \alpha$ , we fail to reject  $H_0$

# Goodness-of-Fit

- We are interested in the p-value of  $\Pr(\chi^2 > X^2)$
- We can find this p-value by using a  $\chi^2$  distribution with  $df = k - 1$ 
  - In R:  $p = 1 - \text{pchisq}(X^2, df)$
  - Always looking for *upper tail probability* (probability of seeing an outcome that is as extreme or more extreme than what we observed)
- If  $p \leq \alpha$ , we reject  $H_0$
- If  $p > \alpha$ , we fail to reject  $H_0$
- State your conclusion in the context of the problem

# Goodness-of-Fit

- We are interested in the p-value of  $\Pr(\chi^2 > X^2)$
- We can find this p-value by using a  $\chi^2$  distribution with  $df = k - 1$ 
  - In R:  $p = 1 - \text{pchisq}(X^2, df)$
  - Always looking for *upper tail probability* (probability of seeing an outcome that is as extreme or more extreme than what we observed)
- If  $p \leq \alpha$ , we reject  $H_0$
- If  $p > \alpha$ , we fail to reject  $H_0$
- State your conclusion in the context of the problem
- Note: For this test to be valid, all the expected counts must be at least 5

$$\begin{aligned} n\hat{P}_{1,0}, \dots, n\hat{P}_{k,0} \\ \geq 5 \end{aligned}$$

# Goodness-of-Fit: Example

# Goodness-of-Fit: Example

- Consider eye color, broken down into four categories: brown, hazel, blue, other

# Goodness-of-Fit: Example

- Consider eye color, broken down into four categories: brown, hazel, blue, other
- Based on previous knowledge, we believe that 40% of people have brown eyes, 10% have hazel eyes, 5% have blue eyes, and 45% have some other color eyes

# Goodness-of-Fit: Example

- Consider eye color, broken down into four categories: brown, hazel, blue, other
- Based on previous knowledge, we believe that 40% of people have brown eyes, 10% have hazel eyes, 5% have blue eyes, and 45% have some other color eyes
- To see if this is correct, we go out and take a sample of 200 people and determine their eye color. We get 84 people with brown eyes, 17 people with hazel eyes, 16 people with blue eyes, and 83 people with other color eyes

# Goodness-of-Fit: Example

- Consider eye color, broken down into four categories: brown, hazel, blue, other
- Based on previous knowledge, we believe that 40% of people have brown eyes, 10% have hazel eyes, 5% have blue eyes, and 45% have some other color eyes
- To see if this is correct, we go out and take a sample of 200 people and determine their eye color. We get 84 people with brown eyes, 17 people with hazel eyes, 16 people with blue eyes, and 83 people with other color eyes
- Test at the  $\alpha = 0.05$  significance level

# Goodness-of-Fit: Example

# Goodness-of-Fit: Example

- $H_0 : p_1 = 0.4, p_2 = 0.1, p_3 = 0.05, p_4 = 0.45$  vs.  $H_1$  : at least one of these proportions does not hold

# Goodness-of-Fit: Example

- $H_0 : p_1 = 0.4, p_2 = 0.1, p_3 = 0.05, p_4 = 0.45$  vs.  $H_1$  : at least one of these proportions does not hold
- Calculate the test statistic:

# Goodness-of-Fit: Example

- $H_0 : p_1 = 0.4, p_2 = 0.1, p_3 = 0.05, p_4 = 0.45$  vs.  $H_1$  : at least one of these proportions does not hold
- Calculate the test statistic:  $n = 200$

Category	Brown	Hazel	Blue	Other
Observed	84	17	16	83
Expected	80	20	10	90

$$\chi^2 = \frac{(84-80)^2}{80} + \frac{(17-20)^2}{20} + \frac{(16-10)^2}{10} + \frac{(83-90)^2}{90} = 4.79$$

$$p = \Pr(\chi^2_3 > 4.79) = 1 - \text{pchisq}(4.79, df=3) = 0.188. > \alpha = 0.05$$

# Goodness-of-Fit: Example

- $H_0 : p_1 = 0.4, p_2 = 0.1, p_3 = 0.05, p_4 = 0.45$  vs.  $H_1$  : at least one of these proportions does not hold
- Calculate the test statistic:

<b>Category</b>	<b>Brown</b>	<b>Hazel</b>	<b>Blue</b>	<b>Other</b>
Observed				
Expected				

- $X^2 =$

# Goodness-of-Fit: Example

- $H_0 : p_1 = 0.4, p_2 = 0.1, p_3 = 0.05, p_4 = 0.45$  vs.  $H_1$  : at least one of these proportions does not hold
- Calculate the test statistic:

<b>Category</b>	<b>Brown</b>	<b>Hazel</b>	<b>Blue</b>	<b>Other</b>
Observed				
Expected				

- $X^2 =$
- $\Pr(\chi^2 > X^2) =$

# Goodness-of-Fit: Example

- $H_0 : p_1 = 0.4, p_2 = 0.1, p_3 = 0.05, p_4 = 0.45$  vs.  $H_1$  : at least one of these proportions does not hold
- Calculate the test statistic:

<b>Category</b>	<b>Brown</b>	<b>Hazel</b>	<b>Blue</b>	<b>Other</b>
Observed				
Expected				

- $X^2 =$
- $\Pr(\chi^2 > X^2) =$
- Conclusion:

# Goodness-of-Fit: Example

$O_i$

$P_{i0}$

```
> chisq.test(c(84,17,16,83),p=c(0.4,0.1,0.05,0.45))
```

Chi-squared test for given probabilities

data: c(84, 17, 16, 83)  
X-squared = 4.7944, df = 3, p-value = 0.1875

# Goodness-of-Fit: Example

- Can do directly in R as well

```
> chisq.test(c(84,17,16,83),p=c(0.4,0.1,0.05,0.45))
```

Chi-squared test for given probabilities

```
data: c(84, 17, 16, 83)
X-squared = 4.7944, df = 3, p-value = 0.1875
```

# Contingency Table

# Contingency Table

- Now, let's consider the case of two categorical variables

# Contingency Table

- Now, let's consider the case of two categorical variables
  - We used a normal approximation to the binomial distribution and formed a two-proportion z-test for binary variables

# Contingency Table

- Now, let's consider the case of two categorical variables
  - We used a normal approximation to the binomial distribution and formed a two-proportion z-test for binary variables
- A generalized technique for testing proportions is through the  $\chi^2$  test of independence for contingency tables

# Contingency Table

$$O_{11} = \frac{r_1 \times c_1}{n}$$

$O_{ij}$   
row  
col

do you drink?  
 $O_{11}$

do you smoke?  
smoke?

		Variable 2		Total
		Yes	No	
Variable 1	Yes	10 $O_{11}$	15 $O_{12}$	25 $r_1$
	No	20 $O_{21}$	5 $O_{22}$	25 $r_2$
Total	30 $c_1$	20 $c_2$		$n$

$$\frac{c_1}{n} = \Pr(\text{drink})$$

$$c_1 = O_{11} + O_{21}$$

$$\frac{r_1}{n} = \Pr(\text{smoker})$$

$$r_1 = O_{11} + O_{12}$$

$A$  and  $B$  are indep iff

$$\Pr(A \cap B) = \Pr(A) \times \Pr(B)$$

# Testing Whether Variables are Independent

# Testing Whether Variables are Independent

- Setting:

# Testing Whether Variables are Independent

- Setting:
  - Consider two categorical variables: favorite season (winter, spring, summer, fall) and whether or not someone has pets (yes, no)

# Testing Whether Variables are Independent

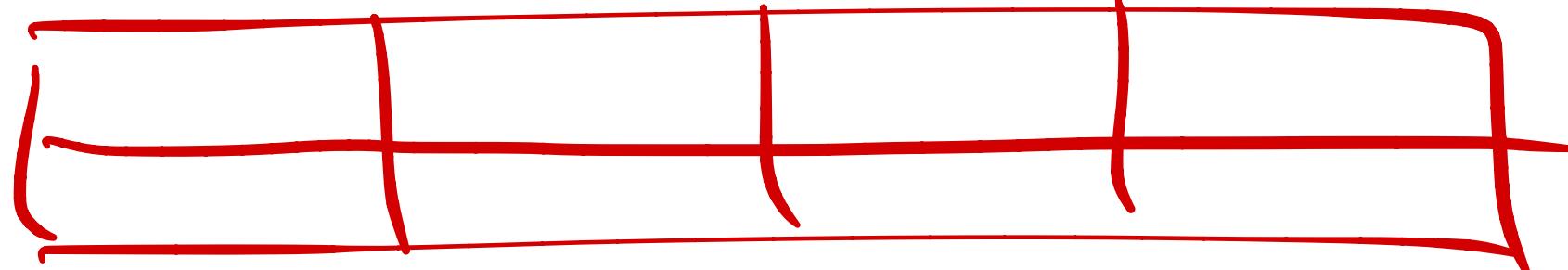
- Setting:
  - Consider two categorical variables: favorite season (winter, spring, summer, fall) and whether or not someone has pets (yes, no)
  - We are interested in whether a population's favorite season is independent of whether or not they are a pet owner

# Testing Whether Variables are Independent

- Setting:
  - Consider two categorical variables: favorite season (winter, spring, summer, fall) and whether or not someone has pets (yes, no)
  - We are interested in whether a population's favorite season is independent of whether or not they are a pet owner
  - How can we test such a hypothesis?

# Testing Whether Variables are Independent

- Setting:
  - Consider two categorical variables: favorite season (winter, spring, summer, fall) and whether or not someone has pets (yes, no)
  - We are interested in whether a population's favorite season is independent of whether or not they are a pet owner
  - How can we test such a hypothesis?
  - Idea: Extend the goodness-of-fit test to multiple dimensions



# $\chi^2$ Test of Independence

# $\chi^2$ Test of Independence

- We are testing the following hypotheses:

# $\chi^2$ Test of Independence

- We are testing the following hypotheses:
  - $H_0$  : the two variables are independent

# $\chi^2$ Test of Independence

- We are testing the following hypotheses:
  - $H_0$  : the two variables are independent
  - $H_1$  : the two variables are associated (i.e., not independent)

# $\chi^2$ Test of Independence

# $\chi^2$ Test of Independence

- Similar to the goodness-of-fit test, the test of independence compares the observed frequencies in each category of the contingency table with the expected frequencies given that the null hypothesis is true

# $\chi^2$ Test of Independence

- Similar to the goodness-of-fit test, the test of independence compares the observed frequencies in each category of the contingency table with the expected frequencies given that the null hypothesis is true
- Let  $O$  be the observed frequencies

# $\chi^2$ Test of Independence

- Similar to the goodness-of-fit test, the test of independence compares the observed frequencies in each category of the contingency table with the expected frequencies given that the null hypothesis is true
  - Let  $O$  be the observed frequencies
  - Let  $E$  be the expected frequencies under the null hypothesis

# $\chi^2$ Test of Independence

- Similar to the goodness-of-fit test, the test of independence compares the observed frequencies in each category of the contingency table with the expected frequencies given that the null hypothesis is true
- Let  $O$  be the observed frequencies
- Let  $\underline{E}$  be the expected frequencies under the null hypothesis
- Use the chi-square test to determine whether the deviations between the observed and expected frequencies are too large to be attributed to chance

# Expected Contingency Table

# Expected Contingency Table

- We compare what we observe with what we expect to see if the null hypothesis is true

# Expected Contingency Table

- We compare what we observe with what we expect to see if the null hypothesis is true
- Calculate the expected counts as follows:

# Expected Contingency Table

- We compare what we observe with what we expect to see if the null hypothesis is true
- Calculate the expected counts as follows:

$$\Pr(Y_{ij} \mid V_1) \quad \Pr(Y_{ij} \mid V_2)$$

$$\left(\frac{r_i}{n}\right) \left(\frac{c_j}{n}\right) \cdot n = \frac{r_i \times c_j}{n}$$

Variable 1	Variable 2		Total
	Yes	No	
Yes	$O_{11}$ $E_{11} = \frac{r_1 \times c_1}{n}$	$O_{12}$ $E_{12} = \frac{r_1 \times c_2}{n}$	$r_1$
No	$O_{21}$ $E_{21} = \frac{r_2 \times c_1}{n}$	$O_{22}$ $E_{22} = \frac{r_2 \times c_2}{n}$	$r_2$
Total	$c_1$	$c_2$	$n$

# $\chi^2$ Test of Independence

# $\chi^2$ Test of Independence

- For a contingency table with  $r$  rows and  $c$  columns (for a total of  $rc$  cells), the chi-square test statistic is as follows:

# $\chi^2$ Test of Independence

- For a contingency table with  $r$  rows and  $c$  columns (for a total of  $rc$  cells), the chi-square test statistic is as follows:

$$X^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

# $\chi^2$ Test of Independence

given  $\nwarrow$

- For a contingency table with  $r$  rows and  $c$  columns (for a total of  $rc$  cells), the chi-square test statistic is as follows:

$$X^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

- $X^2$  approximately follows a  $\chi^2$  distribution with  $(r - 1)(c - 1)$  degrees of freedom


Handwritten annotations:

- Row labels:  $r_1, r_2, r_3, r_4$  are written vertically to the right of the first four rows.
- Column labels:  $c_1, c_2, c_3, c_4$  are written horizontally below the first four columns.
- A handwritten arrow points from  $r_1$  to the top row of the grid.
- A handwritten arrow points from  $c_1$  to the left column of the grid.
- A handwritten arrow points from  $\nwarrow$  to the bottom-right corner of the grid.
- A handwritten arrow points from  $\rightarrow n$  to the bottom-right corner of the grid.
- A handwritten arrow points from  $\underline{(r-1)(c-1)}$  to the formula  $(r - 1)(c - 1)$ .

# $\chi^2$ Test of Independence

- For a contingency table with  $r$  rows and  $c$  columns (for a total of  $rc$  cells), the chi-square test statistic is as follows:

$$X^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

- $X^2$  approximately follows a  $\chi^2$  distribution with  $(r - 1)(c - 1)$  degrees of freedom
- Find the p-value  $p = 1 - \text{pchisq}(X^2, \text{df})$

# $\chi^2$ Test of Independence

- For a contingency table with  $r$  rows and  $c$  columns (for a total of  $rc$  cells), the chi-square test statistic is as follows:

$$X^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

- $X^2$  approximately follows a  $\chi^2$  distribution with  $(r - 1)(c - 1)$  degrees of freedom
- Find the p-value  $p = 1 - \text{pchisq}(X^2, \text{ df})$
- If  $p \leq \alpha$ , then reject  $H_0$

# $\chi^2$ Test of Independence

- For a contingency table with  $r$  rows and  $c$  columns (for a total of  $rc$  cells), the chi-square test statistic is as follows:

$$X^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

- $X^2$  approximately follows a  $\chi^2$  distribution with  $(r - 1)(c - 1)$  degrees of freedom
- Find the p-value  $p = 1 - \text{pchisq}(X^2, \text{ df})$
- If  $p \leq \alpha$ , then reject  $H_0$  
- If  $p > \alpha$ , then fail to reject  $H_0$

# $\chi^2$ Test of Independence

- For a contingency table with  $r$  rows and  $c$  columns (for a total of  $rc$  cells), the chi-square test statistic is as follows:

$$X^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

- $X^2$  approximately follows a  $\chi^2$  distribution with  $(r - 1)(c - 1)$  degrees of freedom
- Find the p-value  $p = 1 - \text{pchisq}(X^2, \text{ df})$
- If  $p \leq \alpha$ , then reject  $H_0$
- If  $p > \alpha$ , then fail to reject  $H_0$
- In order for the  $\chi^2$  distribution to be appropriate, no cell should have an expected or observed frequency less than 5

# $\chi^2$ Test of Independence: Example

# $\chi^2$ Test of Independence: Example

- To examine the effectiveness of flossing, we wish to know whether there is an association between the occurrence of gum disease and the use of floss

# $\chi^2$ Test of Independence: Example

- To examine the effectiveness of flossing, we wish to know whether there is an association between the occurrence of gum disease and the use of floss
- Test the following hypotheses:

# $\chi^2$ Test of Independence: Example

- To examine the effectiveness of flossing, we wish to know whether there is an association between the occurrence of gum disease and the use of floss
- Test the following hypotheses:
  - $H_0$  : Flossing and gum disease are independent

# $\chi^2$ Test of Independence: Example

- To examine the effectiveness of flossing, we wish to know whether there is an association between the occurrence of gum disease and the use of floss
- Test the following hypotheses:
  - $H_0$  : Flossing and gum disease are independent
  - $H_1$  : Flossing and gum disease are associated

# $\chi^2$ Test of Independence: Example

- To examine the effectiveness of flossing, we wish to know whether there is an association between the occurrence of gum disease and the use of floss
- Test the following hypotheses:
  - $H_0$  : Flossing and gum disease are independent
  - $H_1$  : Flossing and gum disease are associated
- Perform this test at the  $\alpha = 0.05$  significance level

# $\chi^2$ Test of Independence: Example

# $\chi^2$ Test of Independence: Example

- 350 dental patients were examined to determine whether flossing daily reduced their risk for gum disease

# $\chi^2$ Test of Independence: Example

- 350 dental patients were examined to determine whether flossing daily reduced their risk for gum disease
- Observed contingency table:

# $\chi^2$ Test of Independence: Example

- 350 dental patients were examined to determine whether flossing daily reduced their risk for gum disease
- Observed contingency table:

Daily Flossing	Gum Disease		Total
	Yes	No	
Yes	50	127	177
No	82	91	173
Total	132	218	350

Annotations for the contingency table:

- Blue arrows point to the column totals:  $c_1$  points to 132, and  $c_2$  points to 218.
- Blue arrows point to the row totals:  $r_1$  points to 177, and  $r_2$  points to 173.
- The cell value 50 is underlined.
- The cell value 127 is underlined.

# $\chi^2$ Test of Independence: Example

# $\chi^2$ Test of Independence: Example

- Given the observed contingency table, what are expected counts?

# $\chi^2$ Test of Independence: Example

- Given the observed contingency table, what are expected counts?

<b>Daily Flossing</b>	<b>Gum Disease</b>		<b>Total</b>
	Yes	No	
Yes	50	127	177
No	82	91	173
Total	132	218	350

# $\chi^2$ Test of Independence: Example

- Given the observed contingency table, what are expected counts?

<b>Daily Flossing</b>	<b>Gum Disease</b>		<b>Total</b>
	Yes	No	
Yes	50	127	177
No	82	91	173
Total	132	218	350

<b>Daily Flossing</b>	<b>Gum Disease</b>		<b>Total</b>
	Yes	No	
Yes			177
No			173
Total	132	218	350

# $\chi^2$ Test of Independence: Example

- Given the observed contingency table, what are expected counts?

Daily Flossing	Gum Disease		Total
	Yes	No	
Yes	50 ↓	127 ↑	177
No	82 ↑	91 ↓	173
Total	132	218	350

Daily Flossing	Gum Disease		Total
	Yes	No	
Yes	66.8	110.2	177
No	65.2	107.8	173
Total	132	218	350

- What is the  $X^2$  statistic?

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{ij} - E_{ij})^2}{E_{ij}} = \frac{(50 - 66.8)^2}{66.8} + \dots + \frac{(91 - 107.8)^2}{107.8} = \underline{13.7}.$$

# $\chi^2$ Test of Independence: Example

# $\chi^2$ Test of Independence: Example

- How many degrees of freedom do we have?  $(2 - 1)(2 - 1) = 1$  degree of freedom

# $\chi^2$ Test of Independence: Example      $\chi^2 = 13.7$

- How many degrees of freedom do we have?  $(2 - 1)(2 - 1) = 1$  degree of freedom
- What is the p value?

$$\Pr(\chi^2_1 \geq 13.7) = 1 - \text{pchisq}(13.7, df=1) = 0.00022$$

# $\chi^2$ Test of Independence: Example

# $\chi^2$ Test of Independence: Example

- How many degrees of freedom do we have?

# $\chi^2$ Test of Independence: Example

- How many degrees of freedom do we have?
  - $(2 - 1)(2 - 1) = 1$  degree of freedom

# $\chi^2$ Test of Independence: Example

- How many degrees of freedom do we have?
  - $(2 - 1)(2 - 1) = 1$  degree of freedom
- What is the p-value?

# $\chi^2$ Test of Independence: Example

- How many degrees of freedom do we have?
  - $(2 - 1)(2 - 1) = 1$  degree of freedom
- What is the p-value?
  - $P(X^2 \geq 13.659) = 1 - \text{pchisq}(13.659, 1) = 0.00022$

# $\chi^2$ Test of Independence: Example

- How many degrees of freedom do we have?
  - $(2 - 1)(2 - 1) = 1$  degree of freedom
- What is the p-value?
  - $P(X^2 \geq 13.659) = 1 - \text{pchisq}(13.659, 1) = 0.00022$
- Since the p-value is less than  $\alpha = 0.05$ , we reject the null hypothesis and conclude that there is an association between flossing and gum disease

# $\chi^2$ Test of Independence: Example

- How many degrees of freedom do we have?
  - $(2 - 1)(2 - 1) = 1$  degree of freedom
- What is the p-value?
  - $P(X^2 \geq 13.659) = 1 - \text{pchisq}(13.659, 1) = 0.00022$
- Since the p-value is less than  $\alpha = 0.05$ , we reject the null hypothesis and conclude that there is an association between flossing and gum disease
- Could also use `chisq.test()` in R

# $\chi^2$ Test of Independence: R Code

```
> x <- matrix(c(50, 82, 127, 91), nrow=2, ncol=2)
> x
      [,1] [,2]
[1,]    50   127
[2,]    82    91
> chisq.test(x, correct=F)
```

Pearson's Chi-squared test

data: x  
X-squared = 13.659, df = 1, p-value = 0.0002192

Yates'

"don't do corrections"

$$\left( O_{ij} - E_{ij} - \frac{1}{2} \right)^2$$

# Larger Contingency Tables

# Larger Contingency Tables

- The  $2 \times 2$  tables that we have talked about thus far are characterized by each variable having only two possible outcomes

# Larger Contingency Tables

- The  $2 \times 2$  tables that we have talked about thus far are characterized by each variable having only two possible outcomes
  - This is the case comparable to the two-proportion z-test

# Larger Contingency Tables

- The  $2 \times 2$  tables that we have talked about thus far are characterized by each variable having only two possible outcomes
  - This is the case comparable to the two-proportion z-test
- We can extend the  $\chi^2$  test to accommodate comparison of more than two proportions

# Larger Contingency Tables

- The  $2 \times 2$  tables that we have talked about thus far are characterized by each variable having only two possible outcomes
  - This is the case comparable to the two-proportion z-test
- We can extend the  $\chi^2$  test to accommodate comparison of more than two proportions
  - $r \times c$  tables

# Larger Contingency Tables

- The  $2 \times 2$  tables that we have talked about thus far are characterized by each variable having only two possible outcomes
  - This is the case comparable to the two-proportion z-test
- We can extend the  $\chi^2$  test to accommodate comparison of more than two proportions
  - $r \times c$  tables
  - $r$  categories for the row variable

# Larger Contingency Tables

- The  $2 \times 2$  tables that we have talked about thus far are characterized by each variable having only two possible outcomes
  - This is the case comparable to the two-proportion z-test
- We can extend the  $\chi^2$  test to accommodate comparison of more than two proportions
  - $r \times c$  tables
  - $r$  categories for the row variable
  - $c$  categories for the column variable

# Larger Contingency Tables

- The  $2 \times 2$  tables that we have talked about thus far are characterized by each variable having only two possible outcomes
  - This is the case comparable to the two-proportion z-test
- We can extend the  $\chi^2$  test to accommodate comparison of more than two proportions
  - $r \times c$  tables
  - $r$  categories for the row variable
  - $c$  categories for the column variable
- The inferential procedures are the same for  $r \times c$  tables as for  $2 \times 2$  tables

# Larger Contingency Tables: Example

# Larger Contingency Tables: Example

- Instead of yes/no flossing status, let's investigate different brands of floss

# Larger Contingency Tables: Example

- Instead of yes/no flossing status, let's investigate different brands of floss
  - Oral-B, Colgate, and Reach

# Larger Contingency Tables: Example

- Instead of yes/no flossing status, let's investigate different brands of floss
  - Oral-B, Colgate, and Reach ← 3
  - Investigate gum disease prevalence for users of each floss type

# Larger Contingency Tables: Example

- Instead of yes/no flossing status, let's investigate different brands of floss
  - Oral-B, Colgate, and Reach
- Investigate gum disease prevalence for users of each floss type
- Conduct a study with 260 people who floss

# Larger Contingency Tables: Example

- Instead of yes/no flossing status, let's investigate different brands of floss
  - Oral-B, Colgate, and Reach
- Investigate gum disease prevalence for users of each floss type
- Conduct a study with 260 people who floss
  - Record the floss brand they use and whether they have gum disease

# Larger Contingency Tables: Example

- Instead of yes/no flossing status, let's investigate different brands of floss
  - Oral-B, Colgate, and Reach
- Investigate gum disease prevalence for users of each floss type
- Conduct a study with 260 people who floss
  - Record the floss brand they use and whether they have gum disease
- Hypotheses:

# Larger Contingency Tables: Example

- Instead of yes/no flossing status, let's investigate different brands of floss
  - Oral-B, Colgate, and Reach
- Investigate gum disease prevalence for users of each floss type
- Conduct a study with 260 people who floss
  - Record the floss brand they use and whether they have gum disease
- Hypotheses:
  - $H_0$  : Floss brand and gum disease are independent

# Larger Contingency Tables: Example

- Instead of yes/no flossing status, let's investigate different brands of floss
  - Oral-B, Colgate, and Reach
- Investigate gum disease prevalence for users of each floss type
- Conduct a study with 260 people who floss
  - Record the floss brand they use and whether they have gum disease
- Hypotheses:
  - $H_0$  : Floss brand and gum disease are independent
  - $H_1$  : There is an association between floss brand and gum disease

# Larger Contingency Tables: Example

- Instead of yes/no flossing status, let's investigate different brands of floss
  - Oral-B, Colgate, and Reach
- Investigate gum disease prevalence for users of each floss type
- Conduct a study with 260 people who floss
  - Record the floss brand they use and whether they have gum disease
- Hypotheses:
  - $H_0$  : Floss brand and gum disease are independent
  - $H_1$  : There is an association between floss brand and gum disease
- Test at the  $\alpha = 0.05$  significance level

# Larger Contingency Tables: Example

## Observed

Daily Flossing	Gum Disease		Total
	Yes	No	
Oral-B	14	70	84
Colgate	25	71	96
Reach	21	59	80
Total	60	200	260

# Larger Contingency Tables: Example

**Observed**

Daily Flossing	Gum Disease		Total
	Yes	No	
Oral-B	14	70	84
Colgate	25	71	96
Reach	21	59	80
Total	60	200	260

**Expected**

Daily Flossing	Gum Disease		Total
	Yes	No	
Oral-B	19.4	64.6	84
Colgate	22.2	73.8	96
Reach	18.5	61.5	80
Total	60	200	260

$c_1$

$c_2$

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

$$\frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

# Larger Contingency Tables: Example

# Larger Contingency Tables: Example

- What is the test statistic?

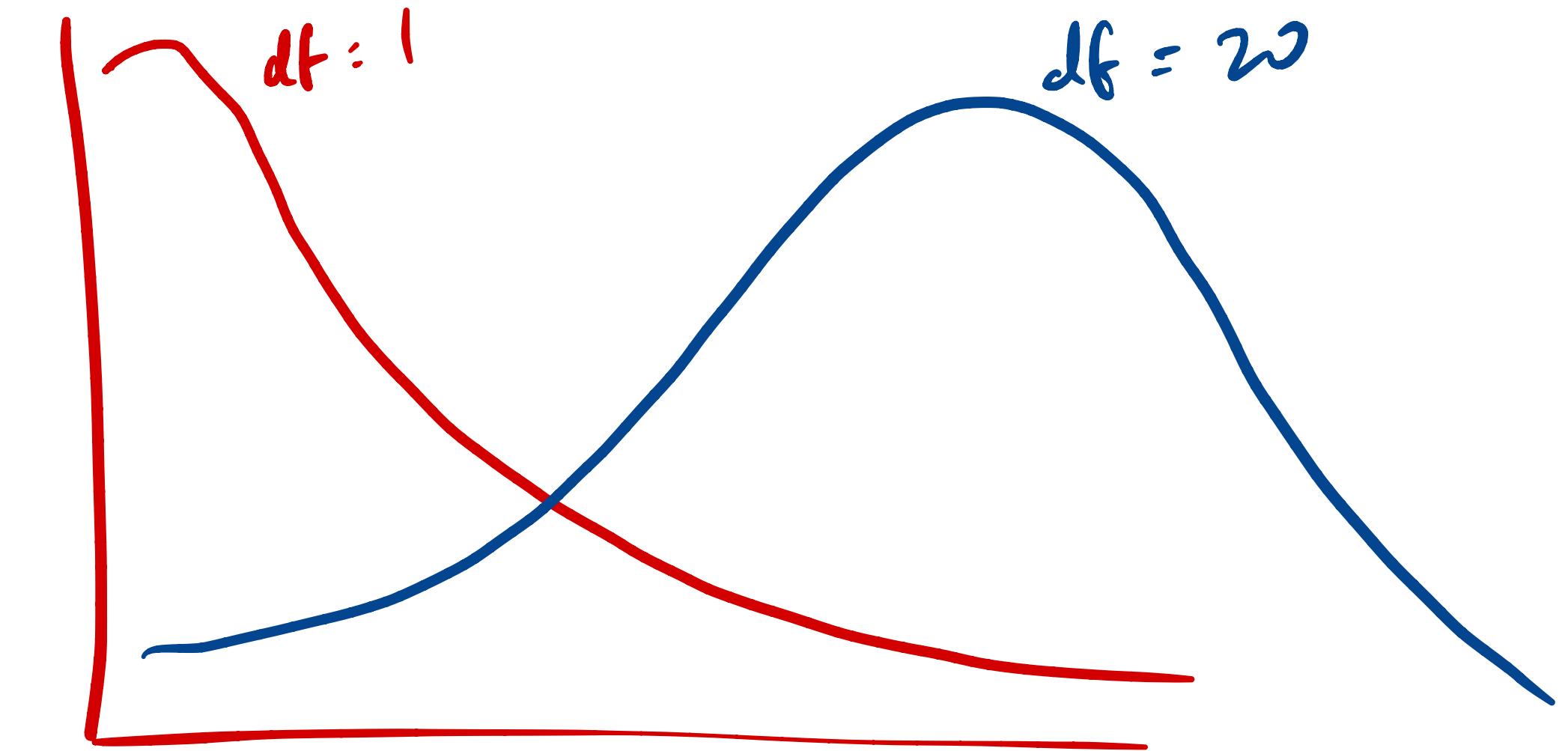
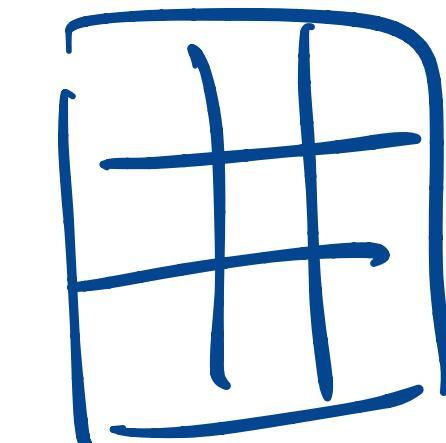
# Larger Contingency Tables: Example

- What is the test statistic?
- How many degrees of freedom?

# Larger Contingency Tables: Example

- What is the test statistic?

$$\underline{\chi^2} = 2.874$$



- How many degrees of freedom?

$$(r-1)(c-1) = 2 \cdot 1 = 2$$

- What is the p-value?

$$\Pr(\chi_2^2 > 2.874) = 0.238 > \alpha = 0.05$$

Fail to reject.

# Larger Contingency Tables: Example

- What is the test statistic?
- How many degrees of freedom?
- What is the p-value?
- Conclusion:

# Summary

# Summary

- Goodness-of-Fit Test

# Summary

- Goodness-of-Fit Test
  - Use this to test if the true proportion = expected proportion

# Summary

- Goodness-of-Fit Test
  - Use this to test if the true proportion = expected proportion
  - Generalization of proportion hypothesis tests

# Summary

- Goodness-of-Fit Test
  - Use this to test if the true proportion = expected proportion
  - Generalization of proportion hypothesis tests
- Chi-Squared ( $\chi^2$ ) Test of Independence

# Summary

- Goodness-of-Fit Test
  - Use this to test if the true proportion = expected proportion
  - Generalization of proportion hypothesis tests
- Chi-Squared ( $\chi^2$ ) Test of Independence
  - Use this to check if variables are independent or not

# Summary

- Goodness-of-Fit Test
  - Use this to test if the true proportion = expected proportion
  - Generalization of proportion hypothesis tests
- Chi-Squared ( $\chi^2$ ) Test of Independence
  - Use this to check if variables are independent or not
- Both rely on the  $\chi^2$  distribution!

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i} \sim \chi^2_{df}$$

## Practise Midterm :

- NOT exhaustive
- 2x as long as real exam (it will be doable)

Review session next week?



Thus.