

# Chi-Squared Tests

Nate Wells

Math 141, 4/7/21

## Outline

In this lecture, we will . . .

- Determine whether data follows a certain distribution
- Investigate the chi-squared distribution.
- Use the chi-squared statistic to determine whether two variables are independent

## Section 1

# The Chi-Squared Test for Goodness of Fit

## Inference for Categorical Variables

Suppose we want to investigate either 1 categorical variable or the relationship between 2 categorical variables.

## Inference for Categorical Variables

Suppose we want to investigate either 1 categorical variable or the relationship between 2 categorical variables.

- If the single variable has just 2 levels, we can consider the proportion  $p$  for one level

## Inference for Categorical Variables

Suppose we want to investigate either 1 categorical variable or the relationship between 2 categorical variables.

- If the single variable has just 2 levels, we can consider the proportion  $p$  for one level
- If both response and explanatory variables have 2 levels, we can consider the difference in proportions  $p_1 - p_2$ .

## Inference for Categorical Variables

Suppose we want to investigate either 1 categorical variable or the relationship between 2 categorical variables.

- If the single variable has just 2 levels, we can consider the proportion  $p$  for one level
- If both response and explanatory variables have 2 levels, we can consider the difference in proportions  $p_1 - p_2$ .

What can we do if one or both the variables are categorical with more than 2 levels?

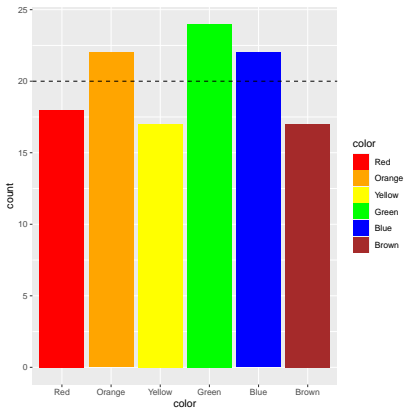
# M&Ms

Suppose we are interested in whether the 6 colors of M&Ms appear with equal frequency. Data from 1 jumbo bag of 120 M&Ms is summarized in the graphic below:



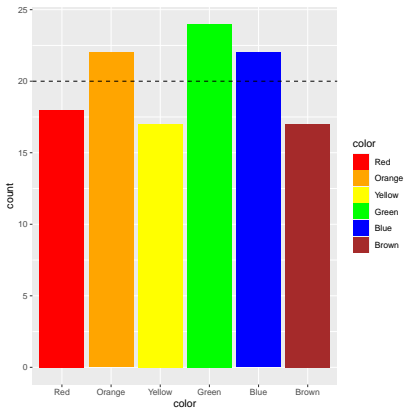
## M&amp;Ms

Suppose we are interested in whether the 6 colors of M&Ms appear with equal frequency. Data from 1 jumbo bag of 120 M&Ms is summarized in the graphic below:



## M&amp;Ms

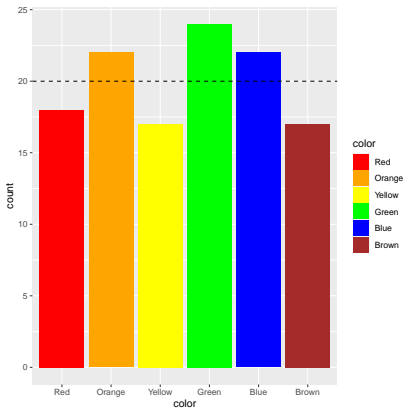
Suppose we are interested in whether the 6 colors of M&Ms appear with equal frequency. Data from 1 jumbo bag of 120 M&Ms is summarized in the graphic below:



- Note that Green M&Ms exceed by the expected count by 20%.

## M&amp;Ms

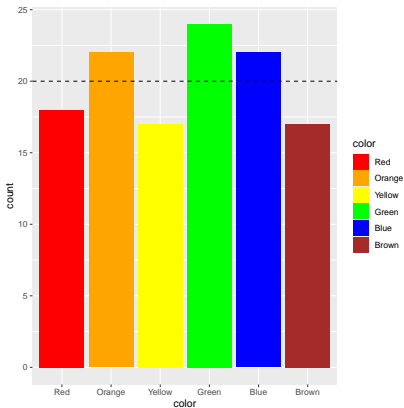
Suppose we are interested in whether the 6 colors of M&Ms appear with equal frequency. Data from 1 jumbo bag of 120 M&Ms is summarized in the graphic below:



- Note that Green M&Ms exceed by the expected count by 20%.
- Does this give good evidence that M&M colors appear at different rates?

## M&amp;Ms

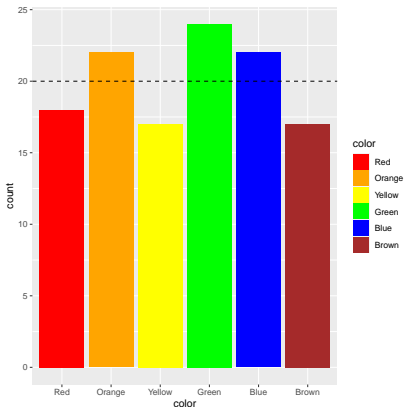
Suppose we are interested in whether the 6 colors of M&Ms appear with equal frequency. Data from 1 jumbo bag of 120 M&Ms is summarized in the graphic below:



- Note that Green M&Ms exceed by the expected count by 20%.
- Does this give good evidence that M&M colors appear at different rates?
  - Suppose we had 20 colors instead of 6...

## M&amp;Ms

Suppose we are interested in whether the 6 colors of M&Ms appear with equal frequency. Data from 1 jumbo bag of 120 M&Ms is summarized in the graphic below:



- Note that Green M&Ms exceed by the expected count by 20%.
- Does this give good evidence that M&M colors appear at different rates?
  - Suppose we had 20 colors instead of 6. . .
  - Would it really be unusual for 1 color to be over- or under-represented?

## Data

Let's consider some numeric data:

Color	Red	Orange	Yellow	Green	Blue	Brown
Frequency	.15	.183	.142	.2	.183	.142
Counts	18	22	17	24	22	17
Expected Counts	20	20	20	20	20	20
Difference (Obs - Exp)	-2	2	-3	4	2	-3

## Data

Let's consider some numeric data:

Color	Red	Orange	Yellow	Green	Blue	Brown
Frequency	.15	.183	.142	.2	.183	.142
Counts	18	22	17	24	22	17
Expected Counts	20	20	20	20	20	20
Difference (Obs - Exp)	-2	2	-3	4	2	-3

We want to test the following hypotheses:

$$H_0 : p_r = \frac{1}{6} \quad p_o = \frac{1}{6} \quad p_y = \frac{1}{6} \quad p_g = \frac{1}{6} \quad p_b = \frac{1}{6} \quad p_{br} = \frac{1}{6}$$

$$H_a : p_r \neq \frac{1}{6} \quad p_o \neq \frac{1}{6} \quad p_y \neq \frac{1}{6} \quad p_g \neq \frac{1}{6} \quad p_b \neq \frac{1}{6} \quad p_{br} \neq \frac{1}{6}$$

# Randomization

- Since we have theoretical values for each proportion, we can simulate samples



# Randomization

- Since we have theoretical values for each proportion, we can simulate samples

```
## # A tibble: 6 x 8
##   color  `1`  `2`  `3`  `4`  `5`  expected observed
##   <chr> <chr> <chr> <chr> <chr> <chr> <chr>   <chr>
## 1 Blue   22   10   22   13   18   20     22
## 2 Brown  15   25   17   17   24   20     17
## 3 Green  28   17   24   23   18   20     24
## 4 Orange 19   21   23   29   26   20     22
## 5 Red    19   20   23   19   13   20     18
## 6 Yellow 17   27   11   19   21   20     17
```

# Randomization

- Since we have theoretical values for each proportion, we can simulate samples

```
## # A tibble: 6 x 8
##   color  `1`  `2`  `3`  `4`  `5`  expected observed
##   <chr> <chr> <chr> <chr> <chr> <chr> <chr>  <chr>
## 1 Blue   22   10   22   13   18   20     22
## 2 Brown  15   25   17   17   24   20     17
## 3 Green  28   17   24   23   18   20     24
## 4 Orange 19   21   23   29   26   20     22
## 5 Red    19   20   23   19   13   20     18
## 6 Yellow 17   27   11   19   21   20     17
```

- How does the observed data compare?

## A Statistic

- Rather than looking at difference for a single level, we might look at **total difference**  
 $\sum(\text{Observed} - \text{Expected})$

## A Statistic

- Rather than looking at difference for a single level, we might look at **total difference**  
 $\sum(\text{Observed} - \text{Expected})$
- But since any negative difference in one location corresponds to a positive difference elsewhere, we might instead consider **total absolute differences**  
 $\sum|\text{Observed} - \text{Expected}|$

## A Statistic

- Rather than looking at difference for a single level, we might look at **total difference**  
 $\sum(\text{Observed} - \text{Expected})$
- But since any negative difference in one location corresponds to a positive difference elsewhere, we might instead consider **total absolute differences**  
 $\sum|\text{Observed} - \text{Expected}|$
- Since we want one large difference to be more influential than many small differences, we look at the **squared differences**

$$\sum(\text{Observed} - \text{Expected})^2$$

- But we should also anticipate larger observations lead to larger **total squared differences**, so we might normalize  $\sum \frac{(\text{Observed} - \text{Expected})^2}{\text{st dev}}$

## A Statistic

- Rather than looking at difference for a single level, we might look at **total difference**  
 $\sum(\text{Observed} - \text{Expected})$
- But since any negative difference in one location corresponds to a positive difference elsewhere, we might instead consider **total absolute differences**  
 $\sum|\text{Observed} - \text{Expected}|$
- Since we want one large difference to be more influential than many small differences, we look at the **squared differences**

$$\sum(\text{Observed} - \text{Expected})^2$$

- But we should also anticipate larger observations lead to larger **total squared differences**, so we might normalize  $\sum \frac{(\text{Observed} - \text{Expected})^2}{\text{st dev}}$ 
  - The standard deviation for squared differences is approximately  $\sqrt{\text{Expected}}$

## A Statistic

- Rather than looking at difference for a single level, we might look at **total difference**  
 $\sum(\text{Observed} - \text{Expected})$
- But since any negative difference in one location corresponds to a positive difference elsewhere, we might instead consider **total absolute differences**  
 $\sum|\text{Observed} - \text{Expected}|$
- Since we want one large difference to be more influential than many small differences, we look at the **squared differences**

$$\sum(\text{Observed} - \text{Expected})^2$$

- But we should also anticipate larger observations lead to larger **total squared differences**, so we might normalize  $\sum \frac{(\text{Observed} - \text{Expected})^2}{\text{st dev}}$ 
  - The standard deviation for squared differences is approximately  $\sqrt{\text{Expected}}$
- This gives a statistic

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

## A Statistic

- Rather than looking at difference for a single level, we might look at **total difference**  
 $\sum(\text{Observed} - \text{Expected})$
- But since any negative difference in one location corresponds to a positive difference elsewhere, we might instead consider **total absolute differences**  
 $\sum|\text{Observed} - \text{Expected}|$
- Since we want one large difference to be more influential than many small differences, we look at the **squared differences**

$$\sum(\text{Observed} - \text{Expected})^2$$

- But we should also anticipate larger observations lead to larger **total squared differences**, so we might normalize  $\sum \frac{(\text{Observed} - \text{Expected})^2}{\text{st dev}}$ 
  - The standard deviation for squared differences is approximately  $\sqrt{\text{Expected}}$
- This gives a statistic

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

- Then large values of  $\chi^2$  should correspond to extreme samples



## Observed Statistic

- What is the  $\chi^2$  statistic for our observed sample?

$$\chi^2 = \frac{(22 - 20)^2}{20} + \frac{(17 - 20)^2}{20} + \frac{(24 - 20)^2}{20} + \frac{(22 - 20)^2}{20} + \frac{(18 - 20)^2}{20} + \frac{(17 - 20)^2}{20} = 2.3$$

## Observed Statistic

- What is the  $\chi^2$  statistic for our observed sample?

$$\chi^2 = \frac{(22 - 20)^2}{20} + \frac{(17 - 20)^2}{20} + \frac{(24 - 20)^2}{20} + \frac{(22 - 20)^2}{20} + \frac{(18 - 20)^2}{20} + \frac{(17 - 20)^2}{20} = 2.3$$

- But what counts as *large*?

## Observed Statistic

- What is the  $\chi^2$  statistic for our observed sample?

$$\chi^2 = \frac{(22 - 20)^2}{20} + \frac{(17 - 20)^2}{20} + \frac{(24 - 20)^2}{20} + \frac{(22 - 20)^2}{20} + \frac{(18 - 20)^2}{20} + \frac{(17 - 20)^2}{20} = 2.3$$

- But what counts as *large*?
- Let's compute the  $\chi^2$  statistic for each of the previous 5 samples from the theoretical population

```
## # A tibble: 5 x 2
##   r   chi2
##   <int> <dbl>
## 1     1   5.2
## 2     2   9.2
## 3     3   6.4
## 4     4   7.5
## 5     5   5.5
```

## Observed Statistic

- What is the  $\chi^2$  statistic for our observed sample?

$$\chi^2 = \frac{(22 - 20)^2}{20} + \frac{(17 - 20)^2}{20} + \frac{(24 - 20)^2}{20} + \frac{(22 - 20)^2}{20} + \frac{(18 - 20)^2}{20} + \frac{(17 - 20)^2}{20} = 2.3$$

- But what counts as *large*?
- Let's compute the  $\chi^2$  statistic for each of the previous 5 samples from the theoretical population

```
## # A tibble: 5 x 2
##   r   chi2
##   <int> <dbl>
## 1     1   5.2
## 2     2   9.2
## 3     3   6.4
## 4     4   7.5
## 5     5   5.5
```

- So our statistic is much smaller than the statistics for these 5 samples.

## Observed Statistic

- What is the  $\chi^2$  statistic for our observed sample?

$$\chi^2 = \frac{(22 - 20)^2}{20} + \frac{(17 - 20)^2}{20} + \frac{(24 - 20)^2}{20} + \frac{(22 - 20)^2}{20} + \frac{(18 - 20)^2}{20} + \frac{(17 - 20)^2}{20} = 2.3$$

- But what counts as *large*?
- Let's compute the  $\chi^2$  statistic for each of the previous 5 samples from the theoretical population

```
## # A tibble: 5 x 2
##   r   chi2
##   <int> <dbl>
## 1     1   5.2
## 2     2   9.2
## 3     3   6.4
## 4     4   7.5
## 5     5   5.5
```

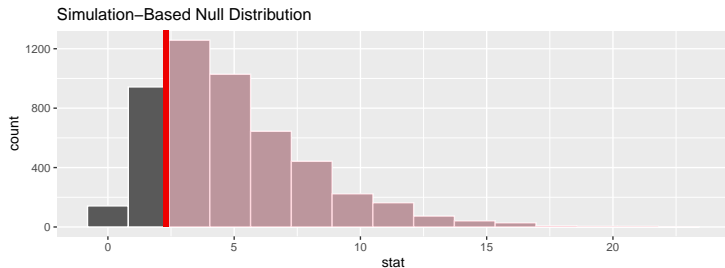
- So our statistic is much smaller than the statistics for these 5 samples.
  - But is this a fluke?

## Distribution of $\chi^2$ statistics

- Let's calculate the  $\chi^2$  statistic for several thousand other samples and plot the distribution

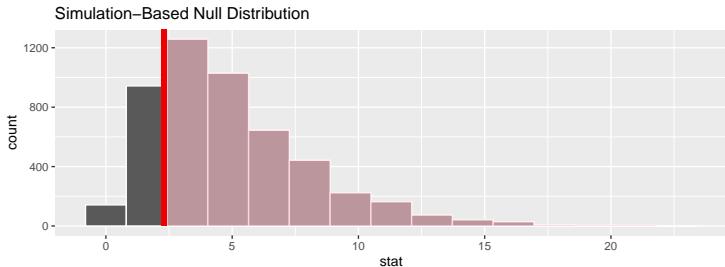
## Distribution of $\chi^2$ statistics

- Let's calculate the  $\chi^2$  statistic for several thousand other samples and plot the distribution



## Distribution of $\chi^2$ statistics

- Let's calculate the  $\chi^2$  statistic for several thousand other samples and plot the distribution

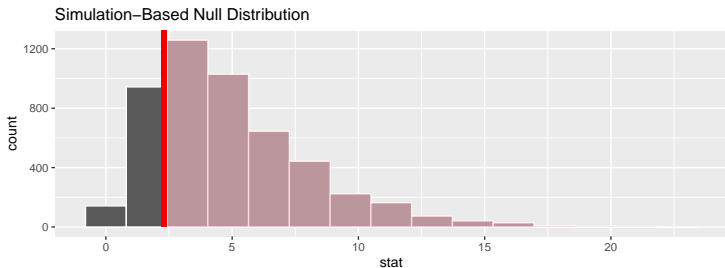


- For this data, it seems that statistics between 0 and 8 are typical.



## Distribution of $\chi^2$ statistics

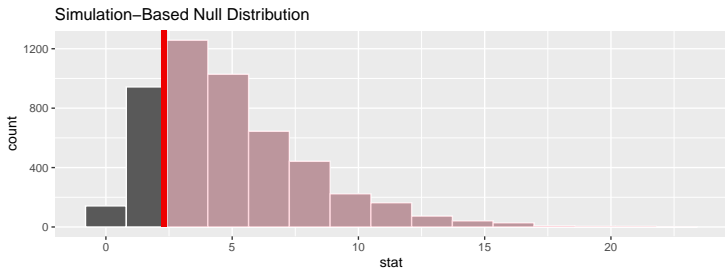
- Let's calculate the  $\chi^2$  statistic for several thousand other samples and plot the distribution



- For this data, it seems that statistics between 0 and 8 are typical.
  - Almost no statistic is greater than 15. And NONE are greater than 20.

## Distribution of $\chi^2$ statistics

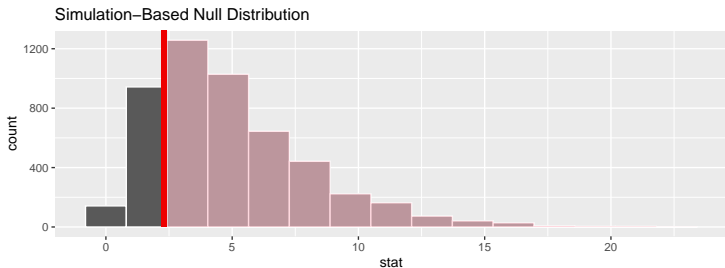
- Let's calculate the  $\chi^2$  statistic for several thousand other samples and plot the distribution



- For this data, it seems that statistics between 0 and 8 are typical.
  - Almost no statistic is greater than 15. And NONE are greater than 20.
- Our observed statistic of  $\chi^2 = 2.3$  is very moderate

## Distribution of $\chi^2$ statistics

- Let's calculate the  $\chi^2$  statistic for several thousand other samples and plot the distribution



- For this data, it seems that statistics between 0 and 8 are typical.
  - Almost no statistic is greater than 15. And NONE are greater than 20.
- Our observed statistic of  $\chi^2 = 2.3$  is very moderate
  - A statistic more extreme would occur about 80% of the time!

## Using infer

- How do we find the probability that a particular  $\chi^2$  value would occur?

## Using infer

- How do we find the probability that a particular  $\chi^2$  value would occur?
  - Use infer!

## Using infer

- How do we find the probability that a particular  $\chi^2$  value would occur?
  - Use infer!

```
set.seed(1)
chisq_null <- MMs %>%
  specify(response = color) %>%
  hypothesize(null = "point",
              p = c("Red" = 1/6, "Orange" = 1/6, "Yellow" = 1/6,
                    "Green" = 1/6, "Blue" = 1/6, "Brown" = 1/6)) %>%
  generate(reps = 5000, type = "simulate") %>%
  calculate(stat = "Chisq")

chisq_null %>% get_p_value(obs_stat = 2.3, direction = "right")

## # A tibble: 1 x 1
##   p_value
##   <dbl>
## 1    0.813
```

## Conclusions

- We tested the following hypotheses:

$$H_0 : p_r = \frac{1}{6} \quad p_o = \frac{1}{6} \quad p_y = \frac{1}{6} \quad p_g = \frac{1}{6} \quad p_b = \frac{1}{6} \quad p_{br} = \frac{1}{6}$$

$$H_a : p_r \neq \frac{1}{6} \quad p_o \neq \frac{1}{6} \quad p_y \neq \frac{1}{6} \quad p_g \neq \frac{1}{6} \quad p_b \neq \frac{1}{6} \quad p_{br} \neq \frac{1}{6}$$

## Conclusions

- We tested the following hypotheses:

$$H_0 : p_r = \frac{1}{6} \quad p_o = \frac{1}{6} \quad p_y = \frac{1}{6} \quad p_g = \frac{1}{6} \quad p_b = \frac{1}{6} \quad p_{br} = \frac{1}{6}$$

$$H_a : p_r \neq \frac{1}{6} \quad p_o \neq \frac{1}{6} \quad p_y \neq \frac{1}{6} \quad p_g \neq \frac{1}{6} \quad p_b \neq \frac{1}{6} \quad p_{br} \neq \frac{1}{6}$$

- Our observed statistic  $\chi^2 = 2.3$  had a simulated p-value of approximately 0.8



## Conclusions

- We tested the following hypotheses:

$$H_0 : p_r = \frac{1}{6} \quad p_o = \frac{1}{6} \quad p_y = \frac{1}{6} \quad p_g = \frac{1}{6} \quad p_b = \frac{1}{6} \quad p_{br} = \frac{1}{6}$$

$$H_a : p_r \neq \frac{1}{6} \quad p_o \neq \frac{1}{6} \quad p_y \neq \frac{1}{6} \quad p_g \neq \frac{1}{6} \quad p_b \neq \frac{1}{6} \quad p_{br} \neq \frac{1}{6}$$

- Our observed statistic  $\chi^2 = 2.3$  had a simulated p-value of approximately 0.8
- We do not reject  $H_0$  at the  $\alpha = 0.05$  significance level (or at any reasonable level)
  - It is likely that such a difference in counts would arise due to chance, if the null hypothesis were true.

## Conclusions

- We tested the following hypotheses:

$$H_0 : p_r = \frac{1}{6} \quad p_o = \frac{1}{6} \quad p_y = \frac{1}{6} \quad p_g = \frac{1}{6} \quad p_b = \frac{1}{6} \quad p_{br} = \frac{1}{6}$$

$$H_a : p_r \neq \frac{1}{6} \quad p_o \neq \frac{1}{6} \quad p_y \neq \frac{1}{6} \quad p_g \neq \frac{1}{6} \quad p_b \neq \frac{1}{6} \quad p_{br} \neq \frac{1}{6}$$

- Our observed statistic  $\chi^2 = 2.3$  had a simulated p-value of approximately 0.8
- We do not reject  $H_0$  at the  $\alpha = 0.05$  significance level (or at any reasonable level)
  - It is likely that such a difference in counts would arise due to chance, if the null hypothesis were true.
- The test provides inconclusive evidence that frequency differs among colors.

## Conclusions

- We tested the following hypotheses:

$$H_0 : p_r = \frac{1}{6} \quad p_o = \frac{1}{6} \quad p_y = \frac{1}{6} \quad p_g = \frac{1}{6} \quad p_b = \frac{1}{6} \quad p_{br} = \frac{1}{6}$$

$$H_a : p_r \neq \frac{1}{6} \quad p_o \neq \frac{1}{6} \quad p_y \neq \frac{1}{6} \quad p_g \neq \frac{1}{6} \quad p_b \neq \frac{1}{6} \quad p_{br} \neq \frac{1}{6}$$

- Our observed statistic  $\chi^2 = 2.3$  had a simulated p-value of approximately 0.8
- We do not reject  $H_0$  at the  $\alpha = 0.05$  significance level (or at any reasonable level)
  - It is likely that such a difference in counts would arise due to chance, if the null hypothesis were true.
- The test provides inconclusive evidence that frequency differs among colors.
  - Importantly, it does not verify that colors ARE equally distributed.

## Theory-based p-values for $\chi^2$ statistics

If we have independent observations on a categorical variable with  $k$  levels, and each observed count is at least 5,

## Theory-based p-values for $\chi^2$ statistics

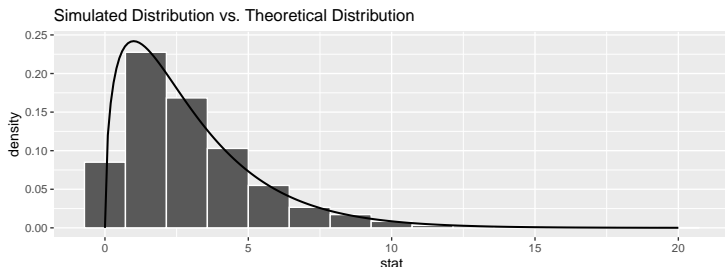
If we have independent observations on a categorical variable with  $k$  levels, and each observed count is at least 5,

- Then  $\chi^2$  is approximately the Chi-Squared distribution with  $k - 1$  degrees of freedom.

## Theory-based p-values for $\chi^2$ statistics

If we have independent observations on a categorical variable with  $k$  levels, and each observed count is at least 5,

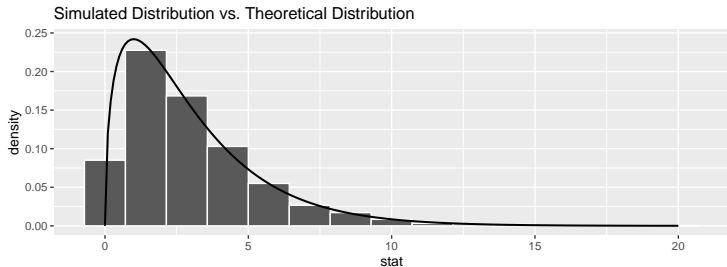
- Then  $\chi^2$  is approximately the Chi-Squared distribution with  $k - 1$  degrees of freedom.



## Theory-based p-values for $\chi^2$ statistics

If we have independent observations on a categorical variable with  $k$  levels, and each observed count is at least 5,

- Then  $\chi^2$  is approximately the Chi-Squared distribution with  $k - 1$  degrees of freedom.

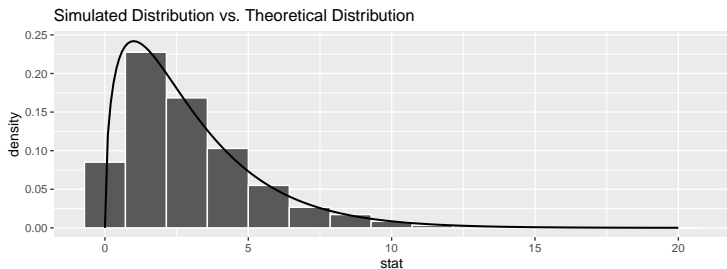


- Use `pchisq(q = ..., df = ..., lower.tail = F)` to find the area right of the observed statistic  $q$ .

## Theory-based p-values for $\chi^2$ statistics

If we have independent observations on a categorical variable with  $k$  levels, and each observed count is at least 5,

- Then  $\chi^2$  is approximately the Chi-Squared distribution with  $k - 1$  degrees of freedom.



- Use `pchisq(q = ..., df = ..., lower.tail = F)` to find the area right of the observed statistic  $q$ .

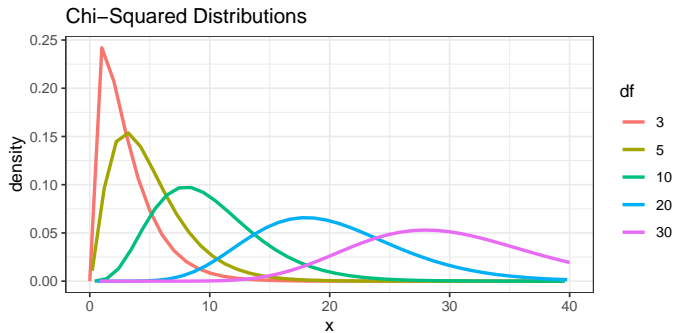
```
pchisq(q = 2.3, df = 5, lower.tail = F)
```

```
## [1] 0.81
```



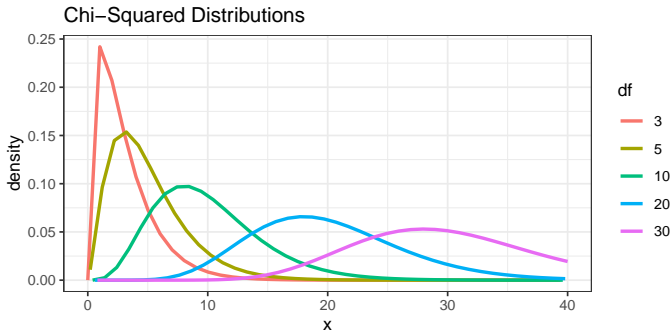
# The Chi-Squared Distribution

Just Normal distributions are described by their mean  $\mu$  and standard deviation  $\sigma$ , the Chi-Square distribution is described by its degrees of freedom  $df$ .



## The Chi-Squared Distribution

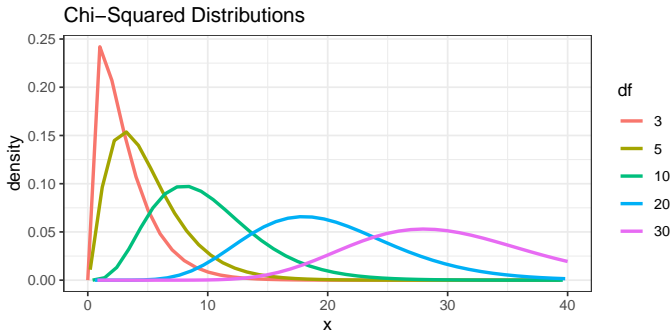
Just Normal distributions are described by their mean  $\mu$  and standard deviation  $\sigma$ , the Chi-Square distribution is described by its degrees of freedom  $df$ .



- Larger degrees of freedom correspond to larger means and larger standard deviations.

## The Chi-Squared Distribution

Just Normal distributions are described by their mean  $\mu$  and standard deviation  $\sigma$ , the Chi-Square distribution is described by its degrees of freedom  $df$ .



- Larger degrees of freedom correspond to larger means and larger standard deviations.
- For Chi-Squared tests, larger degrees of freedom require larger  $\chi^2$  statistic to reject  $H_0$ .

## Section 2

# Chi-Square Test for Independence

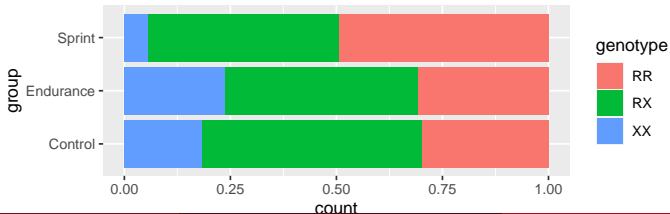
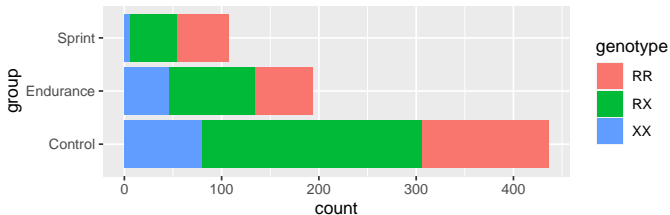
## Genetic Basis for Fast Twitch Muscles

A study on genetics and fast-twitch muscles includes a sample of sprinters, endurance athletes, and a control group of non-athletes.

## Genetic Basis for Fast Twitch Muscles

A study on genetics and fast-twitch muscles includes a sample of sprinters, endurance athletes, and a control group of non-athletes.

- Is there an association between a genotype classification (RR, RX, or XX) and group?



# Contingency Table

Consider the contingency table for group and genotype

```
table(twitch$group, twitch$genotype) %>%  
  addmargins()
```

```
##  
##           RR  RX  XX Sum  
## Control   130 226  80 436  
## Endurance  60  88  46 194  
## Sprint    53  48   6 107  
## Sum       243 362 132 737
```

```
table(twitch$group, twitch$genotype) %>%  
  prop.table( 1)
```

```
##  
##           RR      RX      XX  
## Control   0.298 0.518 0.183  
## Endurance 0.309 0.454 0.237  
## Sprint    0.495 0.449 0.056
```

## Contingency Table

Consider the contingency table for group and genotype

```
table(twitch$group, twitch$genotype) %>%  
  addmargins()
```

```
##  
##           RR  RX  XX Sum  
## Control   130 226  80 436  
## Endurance  60  88  46 194  
## Sprint     53  48   6 107  
## Sum       243 362 132 737
```

```
table(twitch$group, twitch$genotype) %>%  
  prop.table( 1)
```

```
##  
##           RR      RX      XX  
## Control   0.298 0.518 0.183  
## Endurance 0.309 0.454 0.237  
## Sprint    0.495 0.449 0.056
```

- If group and genotype were independent, we would expect proportions to all be equal to the marginal proportions for genotype:

```
table(twitch$genotype) %>% prop.table()
```

```
##  
##  RR  RX  XX  
## 0.33 0.49 0.18
```



## Expected Counts

If the null hypothesis is true, we can multiply the marginal proportions of genotype by the observed counts for group to get expected counts for each genotype-group pair:

	RR	RX	XX
Control	(0.33)(436)	(0.49)(436)	(0.18)(436)
Endurance	(0.33)(194)	(0.49)(194)	(0.18)(194)
Sprint	(0.33)(107)	(0.49)(107)	(0.18)(107)

## Expected Counts

If the null hypothesis is true, we can multiply the marginal proportions of genotype by the observed counts for group to get expected counts for each genotype-group pair:

	RR	RX	XX
Control	(0.33)(436)	(0.49)(436)	(0.18)(436)
Endurance	(0.33)(194)	(0.49)(194)	(0.18)(194)
Sprint	(0.33)(107)	(0.49)(107)	(0.18)(107)

	RR	RX	XX
Control	144	214	78
Endurance	64	95	35
Sprint	35	52	19

## Expected Counts

If the null hypothesis is true, we can multiply the marginal proportions of genotype by the observed counts for group to get expected counts for each genotype-group pair:

	RR	RX	XX
Control	(0.33)(436)	(0.49)(436)	(0.18)(436)
Endurance	(0.33)(194)	(0.49)(194)	(0.18)(194)
Sprint	(0.33)(107)	(0.49)(107)	(0.18)(107)

	RR	RX	XX
Control	144	214	78
Endurance	64	95	35
Sprint	35	52	19

- We can compare to the observed data:

	RR	RX	XX	Sum
Control	130	226	80	436
Endurance	60	88	46	194
Sprint	53	48	6	107
Sum	243	362	132	737

## Expected Counts

If the null hypothesis is true, we can multiply the marginal proportions of genotype by the observed counts for group to get expected counts for each genotype-group pair:

	RR	RX	XX
Control	(0.33)(436)	(0.49)(436)	(0.18)(436)
Endurance	(0.33)(194)	(0.49)(194)	(0.18)(194)
Sprint	(0.33)(107)	(0.49)(107)	(0.18)(107)

	RR	RX	XX
Control	144	214	78
Endurance	64	95	35
Sprint	35	52	19

- We can compare to the observed data:

	RR	RX	XX	Sum
Control	130	226	80	436
Endurance	60	88	46	194
Sprint	53	48	6	107
Sum	243	362	132	737

- As before, we compute the chi-square statistic

$$\chi^2 = \sum \frac{(\text{Observed} - \text{Expected})^2}{\text{Expected}} = 25$$

## The Null Distribution

Under the null hypothesis, group and genotype are independent.

## The Null Distribution

Under the null hypothesis, group and genotype are independent.

- We can simulate data under  $H_0$  by permuting the group labels among individuals. (Just like we did for hypothesis tests for 2 proportions)
  - After each permutation, we compute a new  $\chi^2$  statistic.
  - The distribution of these statistics gives the null distribution.

# The Null Distribution

Under the null hypothesis, group and genotype are independent.

- We can simulate data under  $H_0$  by permuting the group labels among individuals. (Just like we did for hypothesis tests for 2 proportions)
  - After each permutation, we compute a new  $\chi^2$  statistic.
  - The distribution of these statistics gives the null distribution.

##	ID	group	genotype	##	ID	group	genotype
## 1	1	Endurance	RX	## 1	1	Endurance	RX
## 2	2	Sprint	XX	## 2	2	Sprint	RX
## 3	3	Control	XX	## 3	3	Control	XX
## 4	4	Sprint	RX	## 4	4	Sprint	RR
## 5	5	Control	RX	## 5	5	Control	XX
## 6	6	Sprint	RR	## 6	6	Sprint	RX

## Chi-Square Statistic in `infer`

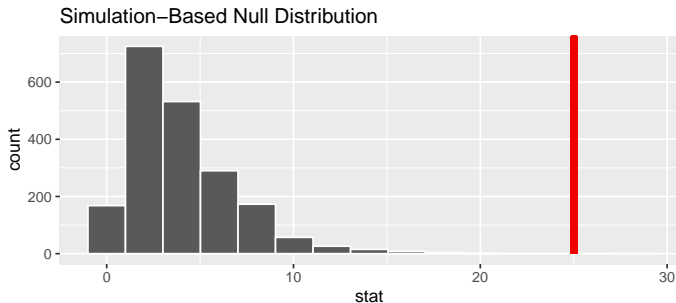
Using `infer...`



# Chi-Square Statistic in infer

Using infer...

```
set.seed(49)
twitch_null <- twitch %>%
  specify(genotype ~ group) %>%
  hypothesize(null = "independence") %>%
  generate(reps = 2000, type = "permute") %>%
  calculate(stat="Chisq")
twitch_null %>% visualize()+shade_p_value(obs_stat = 25, direction = "right")
```



## P-value and conclusions

Using `infer`, the approximate p-value is

```
twitch_null %>% get_p_value(obs_stat = 25, direction = "right")
```

```
## # A tibble: 1 x 1
##   p_value
##   <dbl>
## 1 0.0005
```

## P-value and conclusions

Using `infer`, the approximate p-value is

```
twitch_null %>% get_p_value(obs_stat = 25, direction = "right")
```

```
## # A tibble: 1 x 1
##   p_value
##   <dbl>
## 1 0.0005
```

- At significance  $\alpha = 0.01$ , we reject  $H_0$  in favor the alternative:
  - This sample gives good evidence that group and genotype are associated.

## P-value and conclusions

Using `infer`, the approximate p-value is

```
twitch_null %>% get_p_value(obs_stat = 25, direction = "right")
```

```
## # A tibble: 1 x 1
##   p_value
##   <dbl>
## 1 0.0005
```

- At significance  $\alpha = 0.01$ , we reject  $H_0$  in favor the alternative:
  - This sample gives good evidence that group and genotype are associated.
- What association is there?

## P-value and conclusions

Using `infer`, the approximate p-value is

```
twitch_null %>% get_p_value(obs_stat = 25, direction = "right")
```

```
## # A tibble: 1 x 1
##   p_value
##   <dbl>
## 1 0.0005
```

- At significance  $\alpha = 0.01$ , we reject  $H_0$  in favor the alternative:
  - This sample gives good evidence that group and genotype are associated.
- What association is there?
  - We'll need to further study and experiment to find out.