

---

# **Welcome to DATA 151**

**I'm so glad you're here!**

---

# DATA 151: CLASS 8A

## INTRODUCTION TO DATA SCIENCE (WITH R)

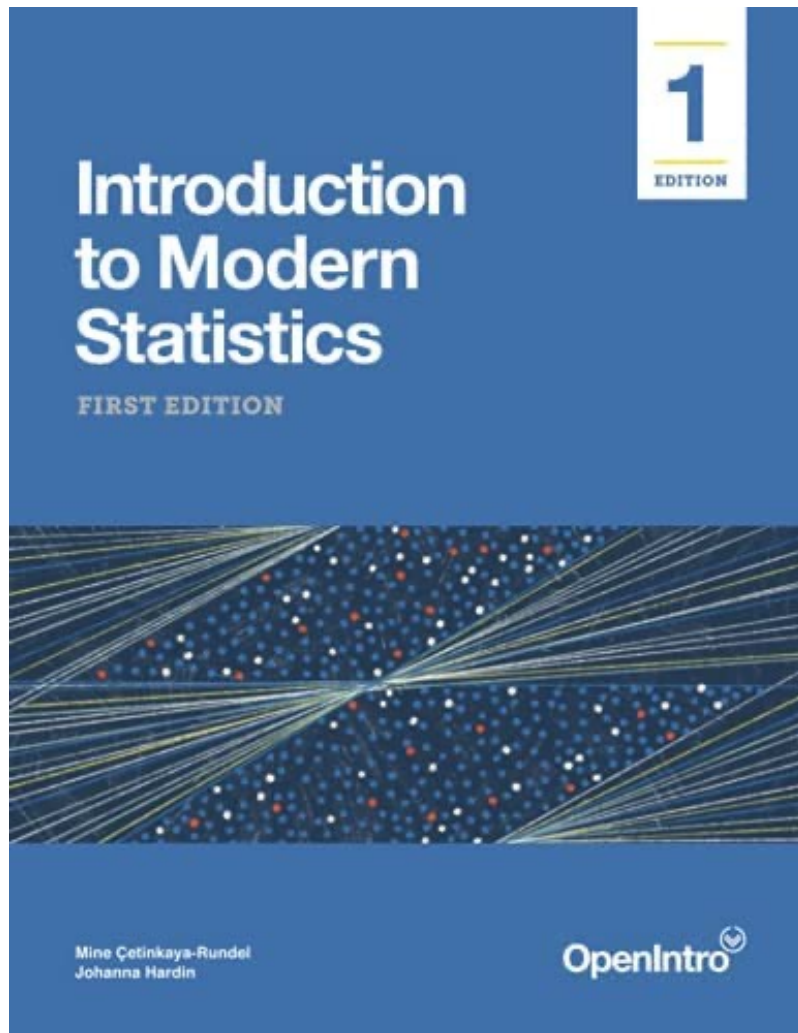
CATEGORICAL DATA ANALYSIS: TABLES AND BARS (PART 2)



# ANNOUNCEMENTS



## RELEVANT READING



## *Introduction to Data Science:*

- Tuesday:
  - Introduction to Modern Statistics
  - Ch 4: Exploring Categorical Data
- Thursday:
  - Introduction to Modern Statistics
  - Ch 5: Exploring Numeric Data

## HOMework REMINDER

***Due this week: (EXTENSION DUE 10/18)***

- *HW #6: DC Introduction to Data Visualization in ggplot2*
  - ***No submission on WISE necessary, do on DataCamp***
- *Project Milestone #3: EDA Step 1*
  - Ask questions and form hypotheses

## HOMEWORK REMINDER

***Due this week: (DUE 10/20)***

- *HW #7: DC Exploratory Data Analysis with Categorical Data*
  - *Just one chapter*
  - ***No submission on WISE necessary, do on DataCamp***
- *Project Milestone #4: EDA Step 2*
  - **Create Tables and Bar Graphs**

## EXTRA CREDIT OPPORTUNITY

### **Data & Computing Tea**

On the Research Experience for Undergraduates,  
Thursday the 20th, 11:30 AM, Ford 201

Meelad Doroodchi, a major in the department who completed a REU (Research Experience for Undergraduates) this summer. Meelad will share some thoughts on the program and his experience, and then we will have time for open discussion.

**Pizza will be provided.**

## EXTRA CREDIT OPPORTUNITY

If you go to the presentation and do a 1-page write up about your take-aways and how this work relates to how we are learning data science in this class, I will give you **4 extra credit points** toward your Midterm #1 grade!



# ANOTHER EXTRA CREDIT OPPORTUNITY

The logo consists of a white rectangular box with a thin black border. Inside the box, the word "FALL" is written in white capital letters on an orange rectangular background, followed by the words "DATA CHALLENGE" in black capital letters on a white background.

FALL DATA  
CHALLENGE

The graphic features the title "AFTER THE BELL" in large, teal, sans-serif capital letters. The word "AFTER" is on the top line, and "THE BELL" is on the bottom line. The letter "O" in "THE" is replaced by a stylized orange bell icon with a white center and a small orange arc above it. The background is a light purple color with a repeating pattern of small, faint white icons representing various educational concepts like books, lightbulbs, and graduation caps.

# AFTER THE BELL

Students will work in teams as they analyze data on the K-12 educational experience “After the Bell.” This year’s theme will require student teams to dive into data on the impacts of school choice and family engagement in school activities and homework. Teams will provide recommendations on factors that best optimize family involvement and support of K-12 students’ academic excellence.

**Submissions will be accepted from October 17 to November 6, at 11:59pm EST.**

## BOARD OF TRUSTEES

On Thursday Oct 20

The Willamette Board of Trustees will be meeting in  
FORD 102 (our classroom space)

For this one day only we have been asked to  
relocate our class to Eaton 209



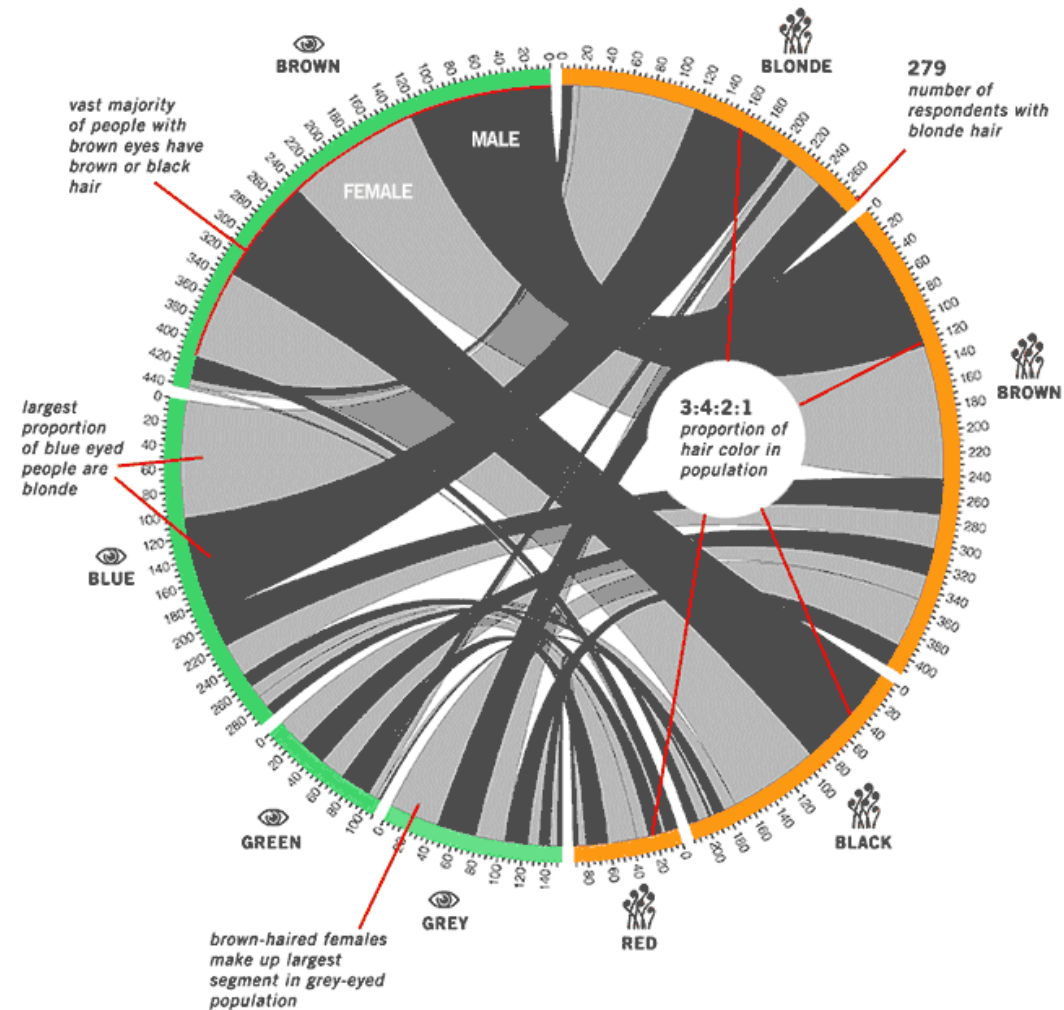
# WORKING WITH CATEGORICAL DATA



## REFRESH: TYPES OF DISTRIBUTIONS

- Joint Distribution
- Marginal Distribution
- Conditional Distribution

## WARM UP EXAMPLE



## WARM UP EXAMPLE

Consider the following survey of 100 students:

	<b>Black</b>	<b>Brunette</b>	<b>Red</b>	<b>Blonde</b>	<b>Total</b>
<b>Black</b>	11	20	4	1	<b>36</b>
<b>Blue</b>	3	14	3	16	<b>36</b>
<b>Hazel</b>	3	9	2	2	<b>16</b>
<b>Green</b>	2	5	2	3	<b>12</b>
<b>Total</b>	<b>19</b>	<b>48</b>	<b>11</b>	<b>22</b>	<b>100</b>

## RECALL: TYPES OF DISTRIBUTIONS

What type of distribution is this?

	<b>Black</b>	<b>Brunette</b>	<b>Red</b>	<b>Blonde</b>	<b>Total</b>
<b>Black</b>	0.11	0.20	0.04	0.01	<b>0.36</b>
<b>Blue</b>	0.03	0.14	0.03	0.16	<b>0.36</b>
<b>Hazel</b>	0.03	0.09	0.02	0.02	<b>0.16</b>
<b>Green</b>	0.02	0.05	0.02	0.03	<b>0.12</b>
<b>Total</b>	<b>0.19</b>	<b>0.48</b>	<b>0.11</b>	<b>0.22</b>	<b>1.00</b>

## RECALL: TYPES OF DISTRIBUTIONS

What type of distribution is this?

	<b>Black</b>	<b>Brunette</b>	<b>Red</b>	<b>Blonde</b>	<b>Total</b>
<b>Black</b>	0.11	0.20	0.04	0.01	<b>0.36</b>
<b>Blue</b>	0.03	0.14	0.03	0.16	<b>0.36</b>
<b>Hazel</b>	0.03	0.09	0.02	0.02	<b>0.16</b>
<b>Green</b>	0.02	0.05	0.02	0.03	<b>0.12</b>
<b>Total</b>	<b>0.19</b>	<b>0.48</b>	<b>0.11</b>	<b>0.22</b>	<b>1.00</b>



## RECALL: TYPES OF DISTRIBUTIONS

What type of distribution is this?

	<b>Black</b>	<b>Brunette</b>	<b>Red</b>	<b>Blonde</b>	<b>Total</b>
<b>Black</b>	0.58	0.42	0.36	0.04	---
<b>Blue</b>	0.16	0.29	0.27	0.73	---
<b>Hazel</b>	0.16	0.19	0.18	0.09	---
<b>Green</b>	0.10	0.10	0.18	0.14	---
<b>Total</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00*</b>	<b>1.00</b>	---

\* rounding error

## RECALL: TYPES OF DISTRIBUTIONS

Where do these numbers come from?

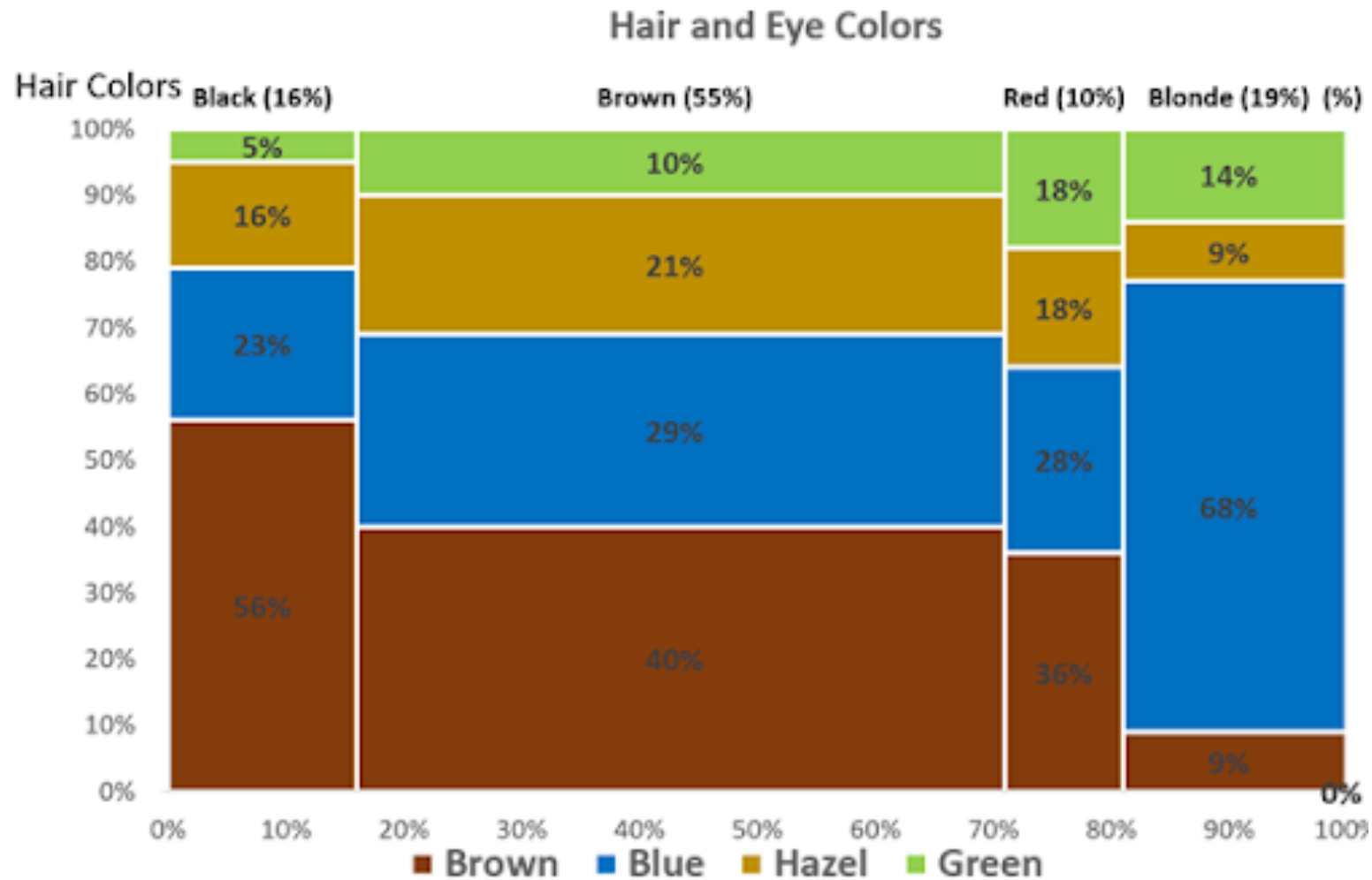
	<b>Black</b>
<b>Black</b>	11
<b>Blue</b>	3
<b>Hazel</b>	3
<b>Green</b>	2
<b>Total</b>	<b>19</b>

→

→

<b>Conditional</b>	<b>Given Black</b>
11/19	0.58
3/19	0.16
3/19	0.16
2/19	0.10
<b>19/19</b>	<b>1.00</b>

# MOSAIC PLOTS!





## EXAMPLE 2: IMMIGRATION POLICY



## MOTIVATING EXAMPLE #2: IMMIGRATION POLICY

**Views on immigration.** Nine-hundred and ten (910) randomly sampled registered voters from Tampa, FL were asked if they thought workers who have illegally entered the US should be (i) allowed to keep their jobs and apply for US citizenship, (ii) allowed to keep their jobs as temporary guest workers but not allowed to apply for US citizenship, or (iii) lose their jobs and have to leave the country. The results of the survey by political ideology are shown below.<sup>48</sup>

# QUESTIONS OF INTEREST

- a. What percent of these Tampa, FL voters identify themselves as conservatives?
- b. What percent of these Tampa, FL voters are in favor of the citizenship option?
- c. What percent of these Tampa, FL voters identify themselves as conservatives and are in favor of the citizenship option?
- d. What percent of these Tampa, FL voters who identify themselves as conservatives are also in favor of the citizenship option? What percent of moderates share this view? What percent of liberals share this view?
- e. Do political ideology and views on immigration appear to be associated? Explain your reasoning.



TRANSITION TO R STUDIO  
FOR OUR HANDS-ON ACTIVITY

# GETTING STARTED

## Step 0: Install the package

```
#install.packages("openintro")  
library(openintro)
```

## Step 1: Load the Data

```
data("immigration")  
str(immigration)
```

```
## tibble [910 × 2] (S3: tbl_df/tbl/data.frame)  
## $ response : Factor w/ 4 levels "Apply for citizenship",...:  
1 1 1 1 1 1 1 1 1 1 ...  
## $ political: Factor w/ 3 levels "conservative",...: 1 1 1 1 1  
1 1 1 1 1 ...
```



## STEP 2

### Step 2: Re-level categories

By default R will order a variable alphabetically, but we might not want that.

```
immigration$political<-as.character(immigration$political)
immigration$political<-factor(immigration$political,
                             levels = c("conservative", "moderate", "liberal"))
```

## QUESTION 1

What percent of these Tampa, FL voters identify themselves as conservatives?

We will learn two new functions to work with individual level data:

- `table()`
- `prop.table()`

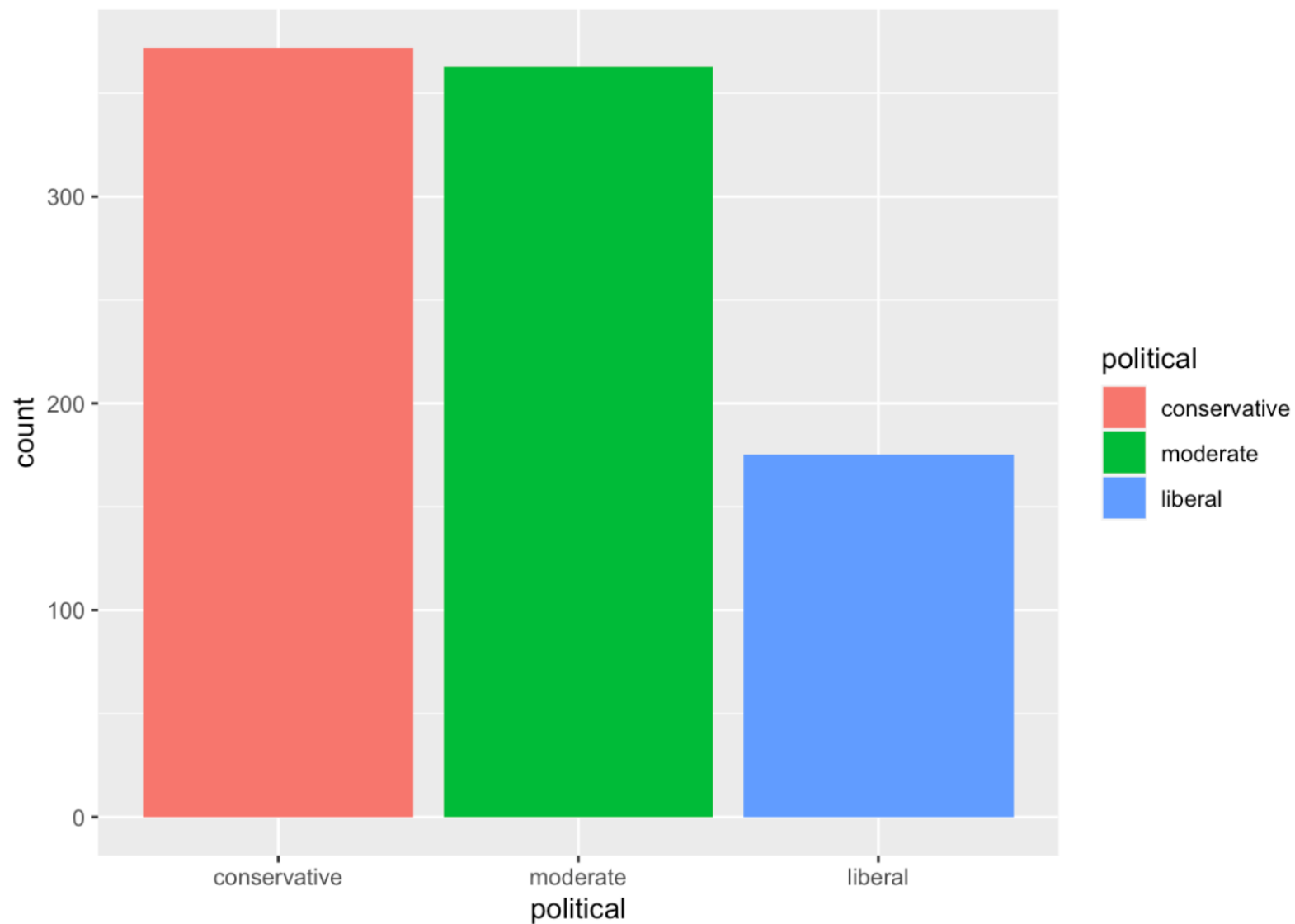
---

```
# Table for Political affiliation  
# use table() function  
tabPol<-table(immigration$political)
```

```
# the prop.table() function must take a table object  
prop.table(tabPol)
```

```
##  
## conservative      moderate      liberal  
##      0.4087912      0.3989011      0.1923077
```

```
# create a graph to display the distribution  
ggplot(immigration, aes(x=political, fill=political))+  
  geom_bar()
```



We can also use `kable` to make tables in R markdown:

```
library(knitr)
kable(tabPol, col.names = c('Party', 'Count'),
      caption = "Distribution of Political Identities")
```

Distribution of Political Identities

Party	Count
conservative	372
moderate	363
liberal	175

## QUESTION 2

What percent of Tampa, FL voters are in favor of the citizenship option?

Use the functions

- `table()`
- `prop.table()`

```

# Table for citizenship response
# use table() function
tabResp<-table(immigration$response)

# use prop.table()
prop.table(tabResp)

```

```

##
## Apply for citizenship          Guest worker          Leave the cou
ntry
##          0.30549451          0.28791209          0.3846
1538
##          Not sure
##          0.02197802

```

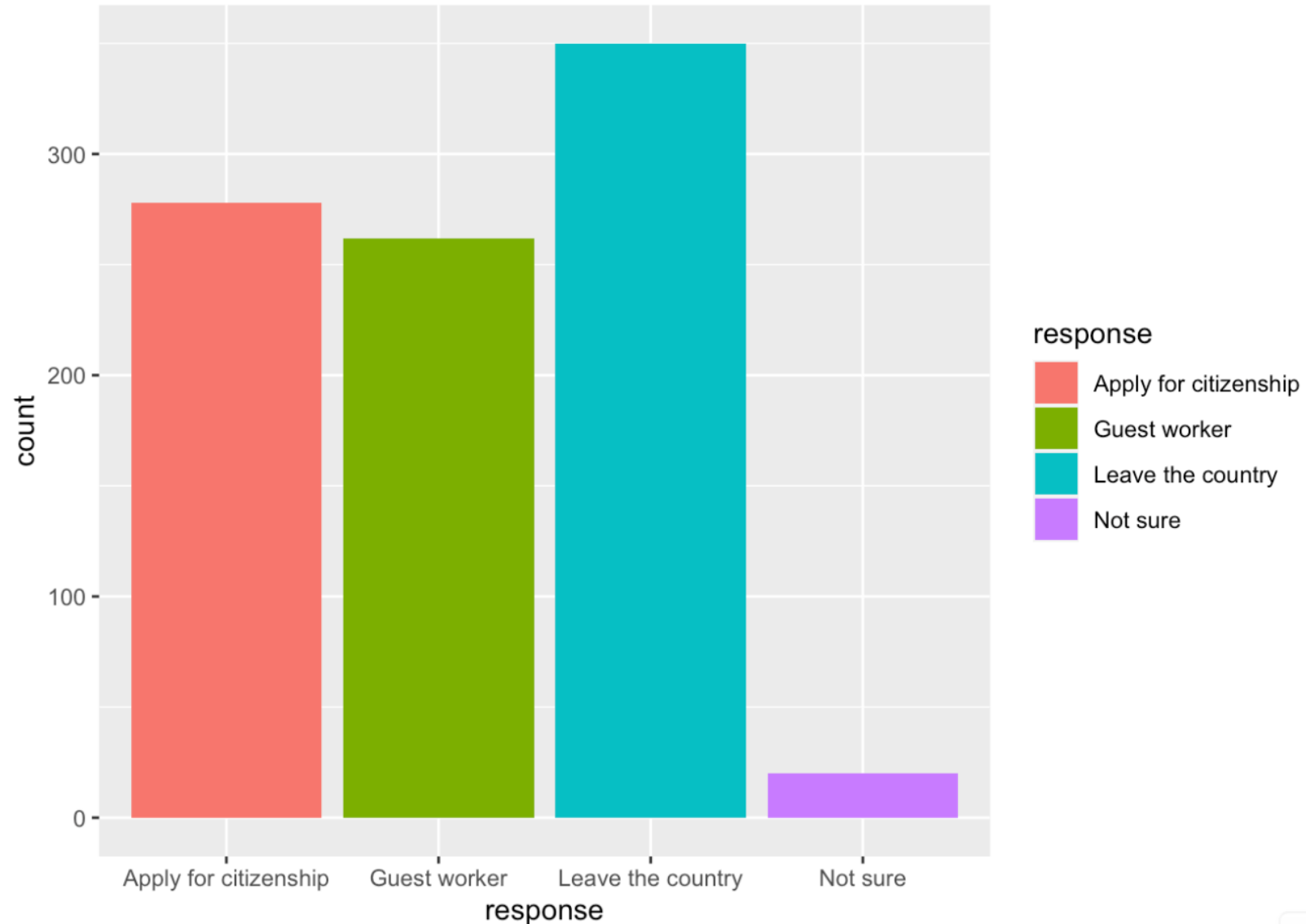
```
# kable
kable(tabResp, col.names = c('Response', 'Count'),
      caption = "Distribution of Response to Citizenship")
```

## Distribution of Response to Citizenship

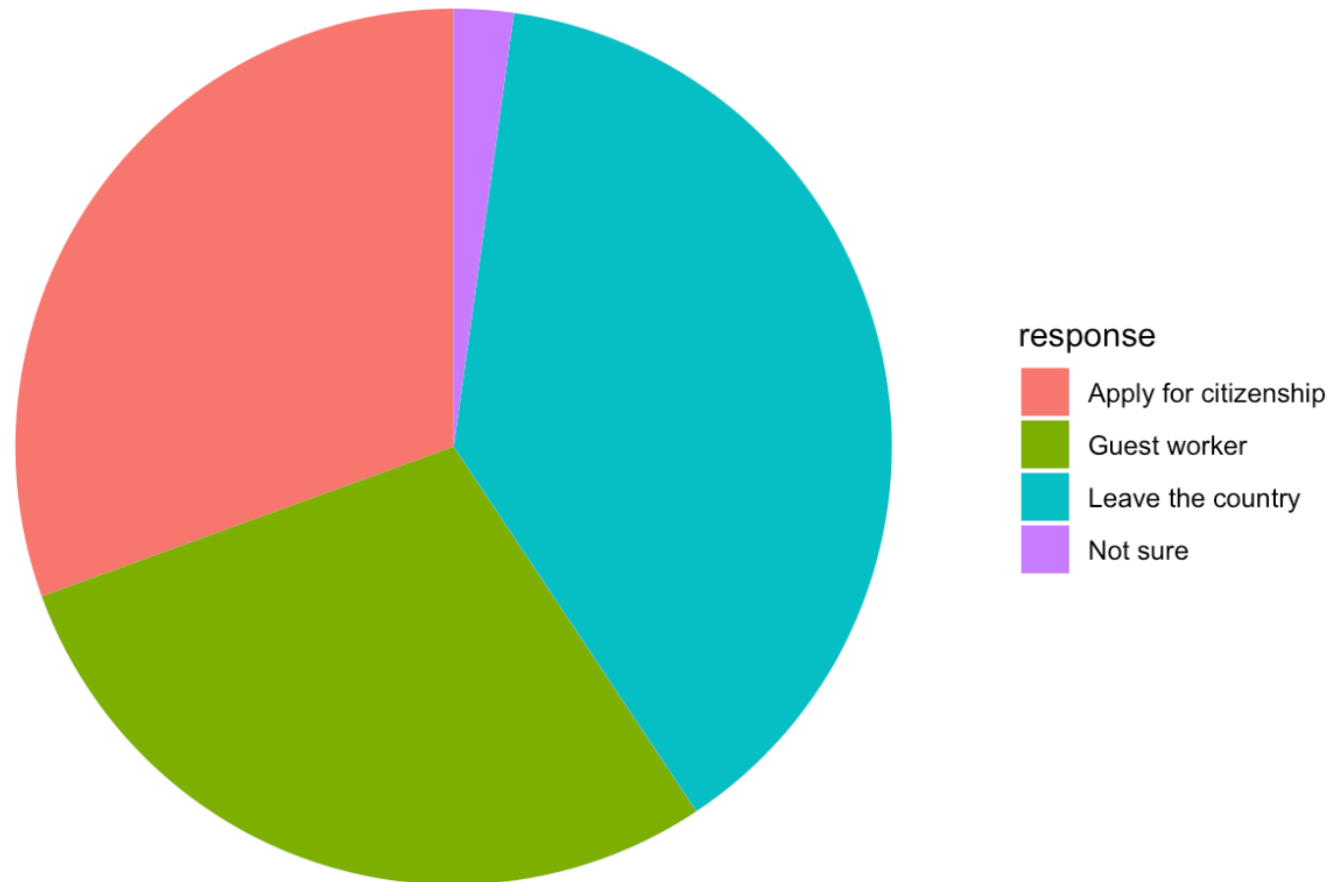
Response	Count
Apply for citizenship	278
Guest worker	262
Leave the country	350
Not sure	20



```
# create a graph to display the distribution  
ggplot(immigration, aes(x=response, fill=response))+  
  geom_bar()
```



```
# pie graph  
ggplot(immigration, aes(x=1, fill=response))+  
  geom_bar()+  
  coord_polar("y", start=0)+  
  theme_void()
```



## QUESTION 3

What percent of these Tampa, FL voters identify themselves as conservatives and are in favor of the citizenship option?

Use the functions

- `table()`
- `prop.table()`

```
## conservative and citizen
```

```
# Row then col
```

```
tabPolResp<-table(immigration$political, immigration$response)
```

```
tabPolResp
```

```
##
```

```
##           Apply for citizenship Guest worker Leave the country Not sure
```

```
## conservative           57           121           179           15
```

```
## moderate           120           113           126           4
```

```
## liberal           101           28           45           1
```

```
## kable
```

```
kable(tabPolResp)
```

	<b>Apply for citizenship</b>	<b>Guest worker</b>	<b>Leave the country</b>	<b>Not sure</b>
conservative	57	121	179	15
moderate	120	113	126	4
liberal	101	28	45	1

```
## joint
```

```
prop.table(tabPolResp)
```

```
##
```

```
##           Apply for citizenship Guest worker Leave the country    Not sure
## conservative      0.062637363   0.132967033           0.196703297 0.016483516
## moderate          0.131868132   0.124175824           0.138461538 0.004395604
## liberal           0.110989011   0.030769231           0.049450549 0.001098901
```

```
sum(prop.table(tabPolResp))
```

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

```
## kable
```

```
kable(round(prop.table(tabPolResp), 2))
```

	<b>Apply for citizenship</b>	<b>Guest worker</b>	<b>Leave the country</b>	<b>Not sure</b>
conservative	0.06	0.13	0.20	0.02
moderate	0.13	0.12	0.14	0.00
liberal	0.11	0.03	0.05	0.00

## QUESTION 4

What percent of these Tampa, FL voters who identify themselves as conservatives are also in favor of the citizenship option? What percent of moderates share this view? What percent of liberals share this view?



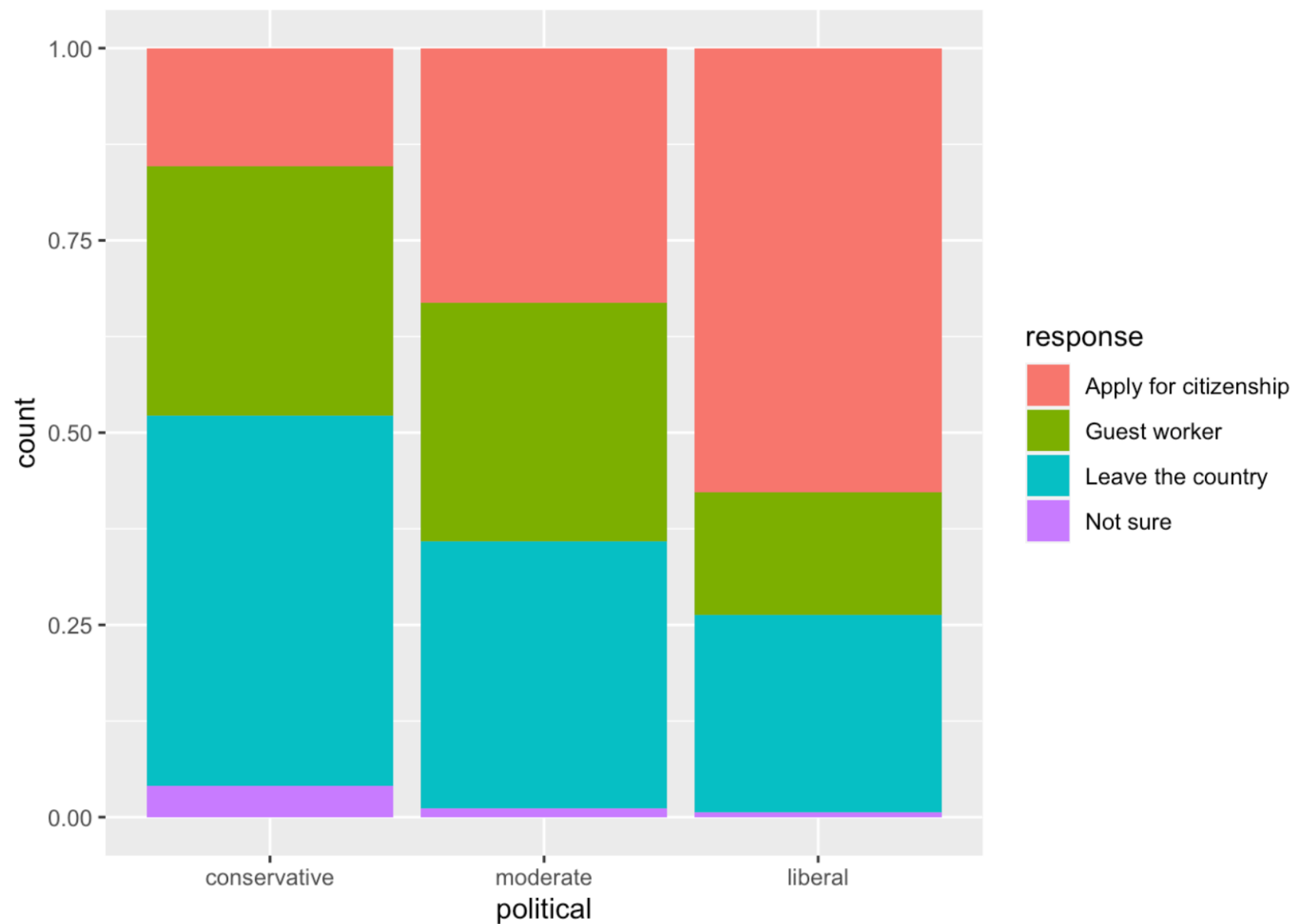
```
## marginal prop
```

```
prop.table(tabPolResp, 1) | = marginal on the row dimension
```

```
##
```

##		Apply for citizenship	Guest worker	Leave the country	Not sure
##	conservative	0.153225806	0.325268817	0.481182796	0.040322581
##	moderate	0.330578512	0.311294766	0.347107438	0.011019284
##	liberal	0.577142857	0.160000000	0.257142857	0.005714286

```
ggplot(immigration, aes(x=political, fill=response))+  
  geom_bar(position="fill")
```





## EXAMPLE 3: SIMPSON'S PARADOX



## SIMPSON'S PARADOX

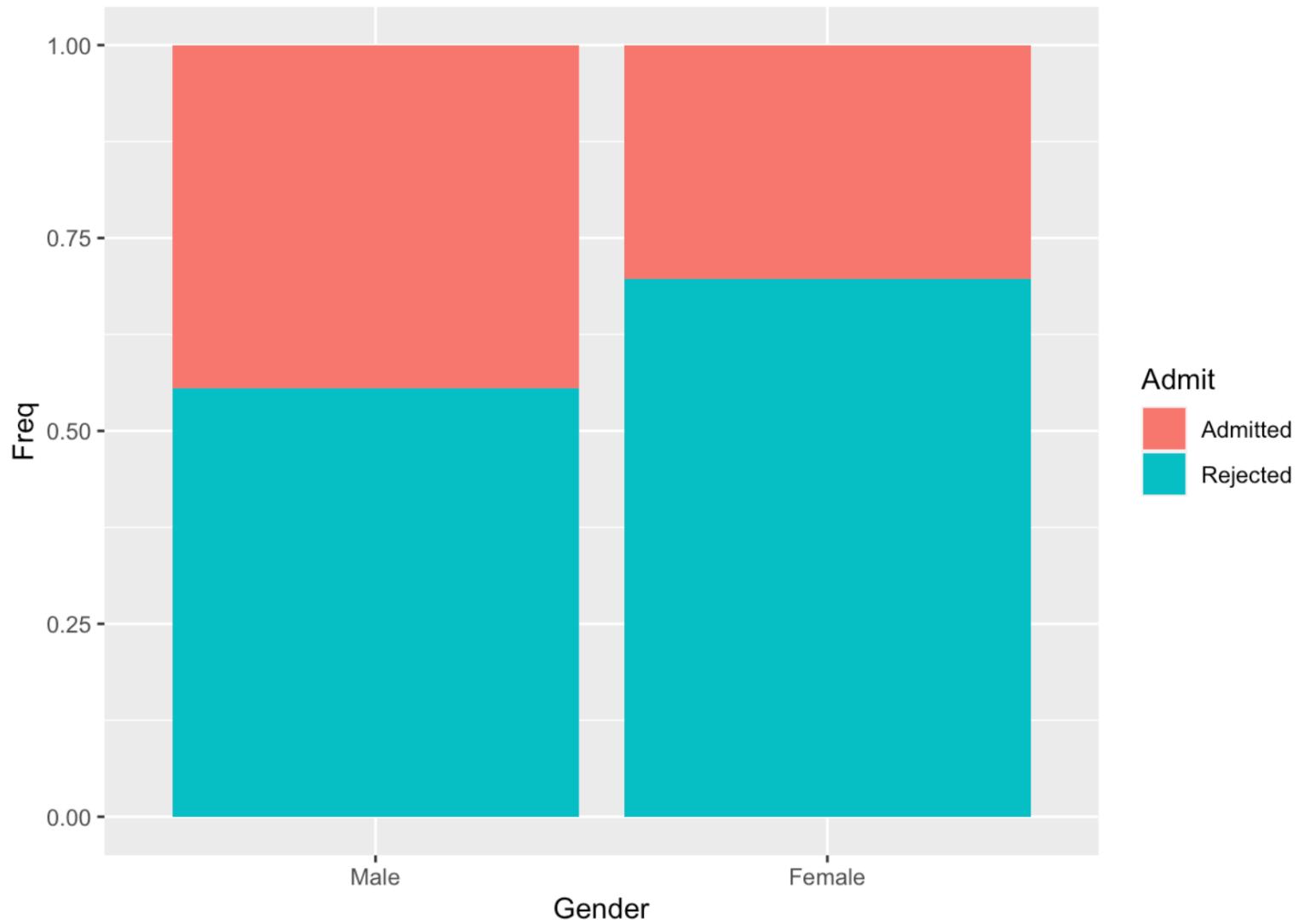
**Simpson's Paradox** (aka The Ecological Fallacy):

A phenomenon in which a trend appears in several different groups but disappears or reverses when the groups are combined.

## EXAMPLE #3: SIMPSON'S PARADOX

- 1973 UC Berkeley Gender Bias in Admissions
- “One of the first universities to be sued for sexual discrimination” (with a statistically significant difference)

Applicants Admitted		
Men	8442	44%
Women	4321	35%

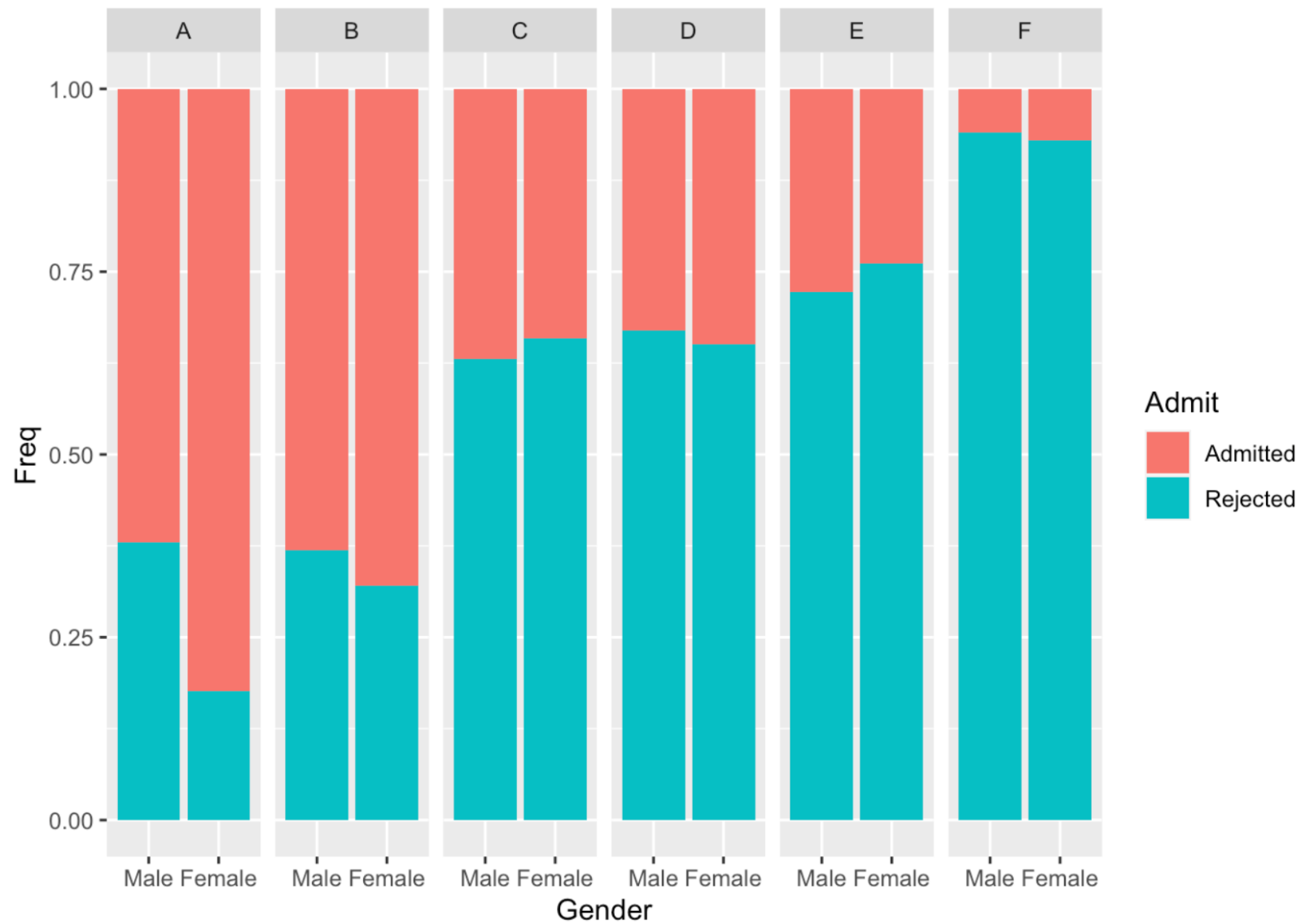


## EXAMPLE #3: SIMPSON'S PARADOX

# SIMPSON'S PARADOX

- But, when you dig into the data...

Department	# of Men	# of Women	Men Accepted	Women Accepted
A	825	108	62%	82%
B	560	25	63%	68%
C	325	593	37%	34%
D	417	375	33%	35%
E	191	393	28%	24%
F	373	341	6%	7%
Total	8442	4321		



## EXAMPLE #3: SIMPSON'S PARADOX



## SIMPSON'S PARADOX

### How does this happen?

*“The simple explanation is that **women tended to apply to the departments that are the hardest to get into**, and men tended to apply to departments that were easier to get into. (Humanities departments tended to have less research funding to support graduate students, while science and engineer departments were awash with money.) So women were rejected more than men. **Presumably, the bias wasn't at Berkeley but earlier in women's education, when other biases led them to different fields of study than men.**”*



TIME FOR GROUP WORK

## MILESTONE #4

- **Due 10/20 - Milestone #4:** Exploratory Data Analysis Step #2
  - Create Tables and Bar Graphs
    - **Goal:** Work to answer at least one of your questions of interest for categorical variables of interest
    - Must have at least one one-way table, and at least one two-way table. Describe interesting distributions (joint, condition, marginal) and what they tell you about your categorical data
    - Must have at least one bar graph with one variable and at least one bar graph with two variables. Describe your graphical voices for how you are presenting your data.
- Please submit using Rmarkdown



ANOTHER ACTIVITY...  
*IF WE HAVE TIME*





# FIVETHIRTYEIGHT ACTIVITY



READ THE ARTICLE

JUN. 26, 2020, AT 7:00 AM

## Voter Registrations Are Way, Way Down During The Pandemic

By [Kaleigh Rogers](#) and [Nathaniel Rakich](#)

Graphics by [Elena Mejia Lutz](#)

Filed under [2020 Election](#)

Get the data on [GitHub](#)



Volunteers instruct citizens how to register to vote during a protest march. ROBERT NICKELSBURG / GETTY IMAGES

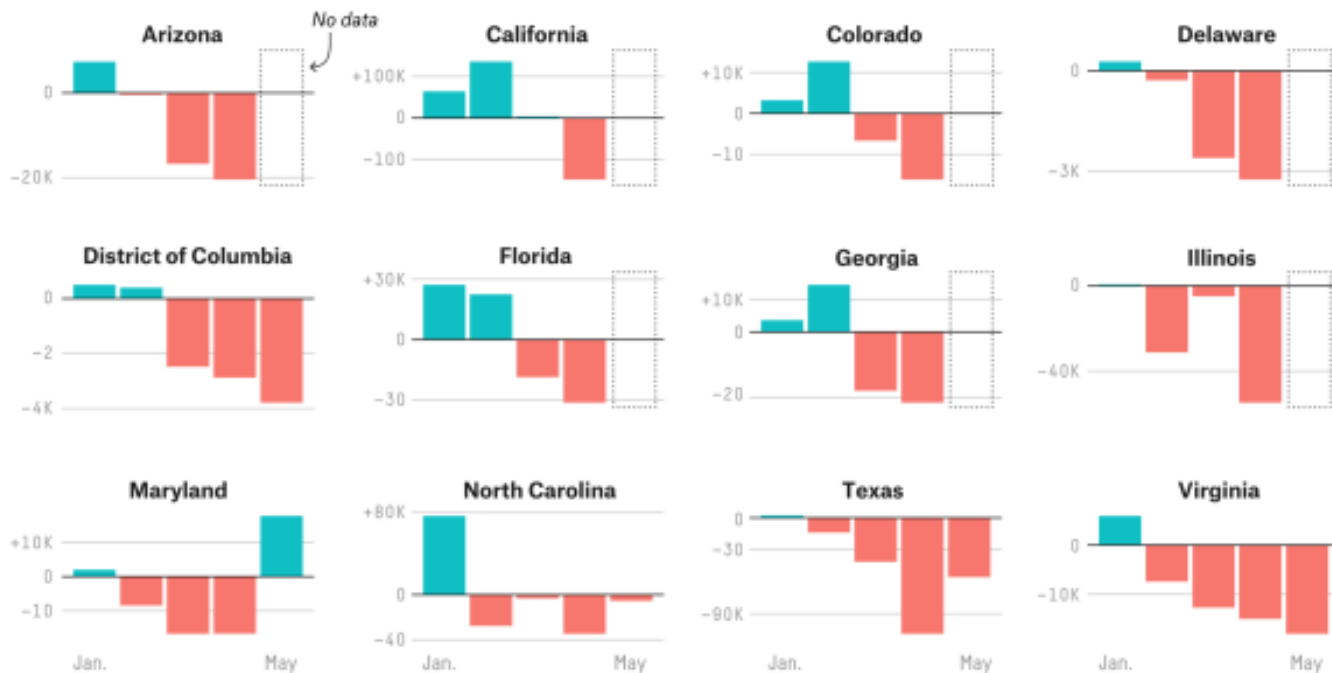
[Poll](#) after [poll](#) showed [a high level of enthusiasm](#) for voting in the general

# DISCUSS IN SMALL GROUPS

1. How are graphics used to tell the author's story?
2. What geometries are used?

## Voter registration dropped dramatically during the pandemic

Difference in the number of newly registered voters for each month in 2020 compared to the same month in 2016



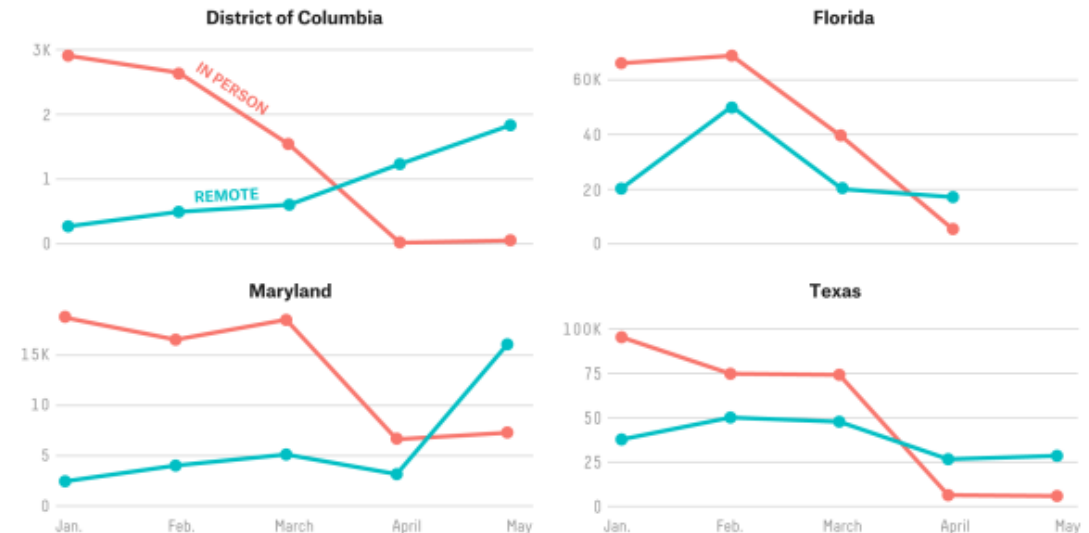
Some states treat voters who move between counties within a state as new registrants because they're unregistered from their old county and newly registered in the new one.

FiveThirtyEight

SOURCE: CENTER FOR ELECTION INNOVATION AND RESEARCH

## In-person registrations dropped as states shut down

Number of new voter registrations submitted in person or remotely, by month, January through May 2020



Does not include voters whose registration method was unclear in state or district data. Some states treat voters who move between counties within a state as new registrants because they're unregistered from their old county and newly registered in the new one.

FiveThirtyEight

SOURCES: TEXAS SECRETARY OF STATE OFFICE, D.C. BOARD OF ELECTIONS, FLORIDA DEPARTMENT OF STATE, MARYLAND STATE BOARD OF ELECTIONS

# WHAT DOES THE RAW DATA LOOK LIKE?

## How to access the data:

```
# Load the tidyverse
library(tidyverse)

# Import data
vreg<-read.csv("https://raw.githubusercontent.com/fivethirtyeight/data/master/voter-registration/new-voter-registrations.csv",
              header=TRUE)
```



# WE MIGHT WANT TO RELEVEL

## 1. Relevel the data so that its in the right order:

```
# Level the Month variable so that its in the right order (ie not alphabetical)
vreg$Month<-factor(vreg$Month,
                  levels=c("Jan", "Feb", "Mar", "Apr", "May"))
```

# ARE WE GOING TO NEED TO TIDY THE DATA?

## 2. Tidy the data:

```
### USE spread() FROM tidyr
vregYear<-vreg%>%
  spread(Year, New.registered.voters)

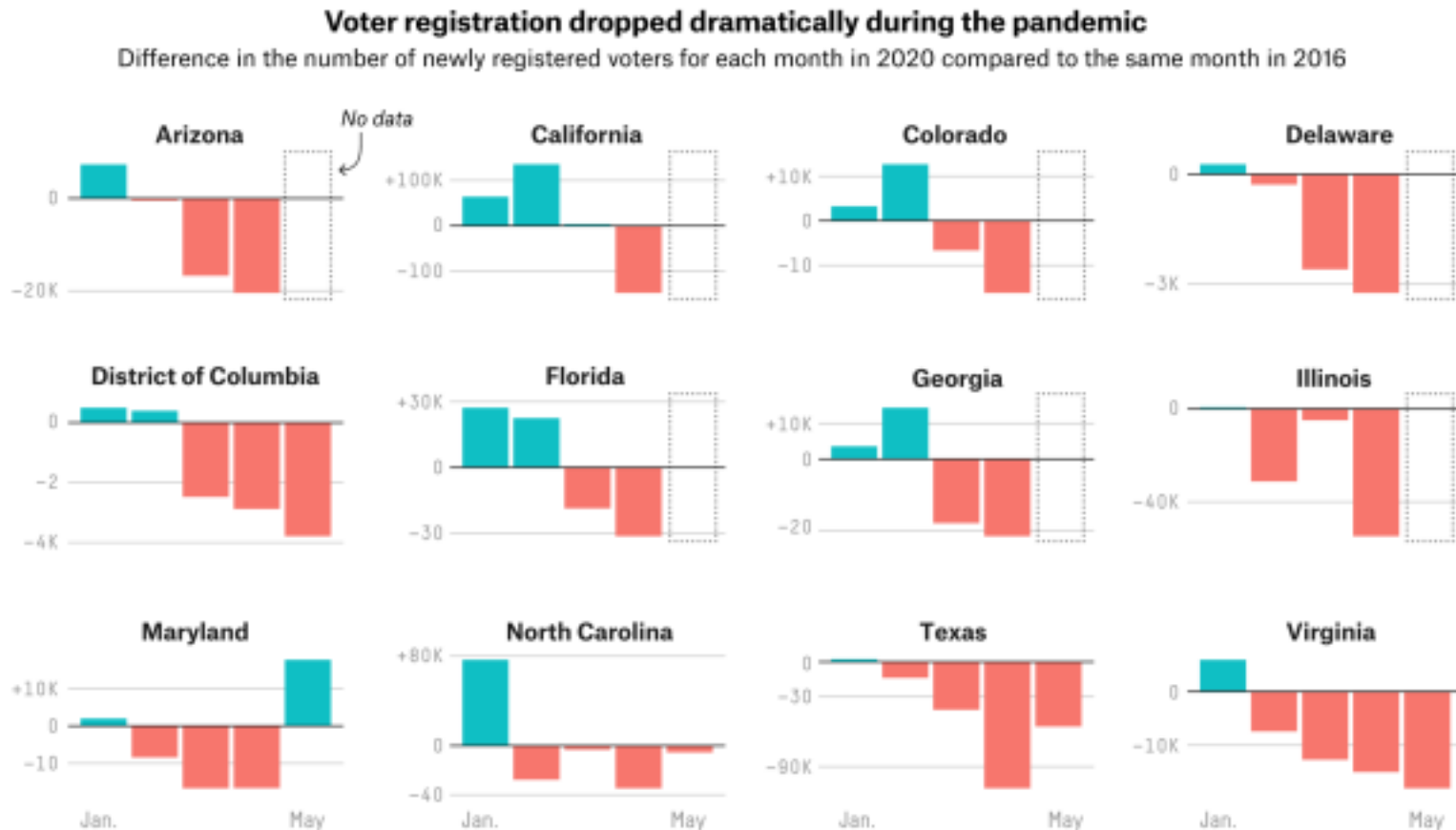
### RENAME THE COLUMNS
colnames(vregYear)<-c("Jurisdiction", "Month", "Y2016", "Y2020")
```

# MUTATE!

## 3. Mutate to add the change:

```
### mutate() FROM dplyr()  
vregChange<-vregYear%>%  
  mutate(change=Y2020-Y2016)
```

# RECREATE THIS GRAPH IN SMALL GROUPS



**Task:** Using the tools we have covered so far, recreate this graph.

Some states treat voters who move between counties within a state as new registrants because they're unregistered from their old county and newly registered in the new one.