

COVID Vaccination Project

Adrian Cornejo and Victor Guillera

11/4/2021

Milestone 2

Section 1

This data set will relate to our question if being a homeowner status and available housing impacts vaccination rates per CA county. While not a perfect representation of SES of cases, this data set will help us elucidate if there is a relationship between the homeowner-related to renter rates and the vaccination rate per county.

We have two data sets, one is of COVID-19 vaccination progress throughout CA provided by the CDPH following vaccine dosage by zip code from 1/5/21 to about present 9/14/21, and the other is demographic info in CA across all counties from the US census data for 2014 to 2019. Below, we read in our libraries, data sets, explore our data sets, and then proceed to clean the data sets.

```
library(readr)
library(knitr)
library(tibble)
library(dplyr)
```

```
##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
##
##   filter, lag

## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union
```

```
library(stringr)
countydem <- read.csv("ca_county_demographics.csv")
covidvax <- read.csv("cov_vax_admin.csv")
```

To import the csv files, we used `read_csv` from base R and will inspect our imported files first before we clean our data set. Libraries we loaded in were `readr`, `dplyr`, and `tibble`. We are going to pull variables “persons_fully_vaccinated”, “persons_partially_vaccinated”, “zip_code_tabulation_area”, “age12_plus_population”, and “county” from the `covidvax` csv file. We will pull “name”, “renter_occ”, “owner_occ”, “hse_units” and “vacant” variables.

```
str(covidvax)
head(covidvax)
nrow(covidvax)
ncol(covidvax)
typeof(covidvax)
class(covidvax)
```

```
distinct(covidvax, as_of_date, .keep_all = F)
```

The structure of the “cov_vax_admin” file is a data frame with 11 columns and 65,268 rows. The data frame contains integer, number, and character data types and contains a redacted column that is inaccessible. There are also a total of 27 unique dates, each 7 weeks apart (weekly results).

```
str(countydem)
head(countydem)
nrow(countydem)
ncol(countydem)
typeof(countydem)
class(countydem)
```

The structure of “countydem” reveals a data frame containing 23 columns and 58 rows with integer, character and number data types.

```
#Cleaning the data by removing unnecessary columns, rearranging the data,
#and removing "N/A" data points.
```

```
#COVID Vaccination Rates
```

```
covidvax2 <- select(covidvax, "county","zip_code_tabulation_area", "as_of_date",
                    "age12_plus_population", "persons_fully_vaccinated",
                    "persons_partially_vaccinated")
```

```
names(covidvax2) <- c("county", "zip_code", "date", "total_pop_over_12",
                     "fully_vaccinated", "partially_vaccinated")
```

```
covidvax2[is.na(covidvax2)] <- 0
```

```
#Using the function str() we see that the selected variables in covidvax2
#are chr, int, chr, num, num
#We may need to convert the dates to numeric using the as.numeric() function
#to determine and compare vaccination rate. All other column types should be
#fine as is.
#Converting NA cells to 0 we allow ourselves to manipulate and describe entire
#numeric columns.
```

```
#Descriptions
```

```
distinct(covidvax2, county, .keep_all = F)
distinct(covidvax2, date, .keep_all = F)
distinct(covidvax2, zip_code, .keep_all = F)
```

```
covidvax2 %>% group_by(zip_code) %>%
  summarise(fullyvax_range = range(fully_vaccinated), partiallyvax_range = range(partially_vaccinated))
```

```
covidvax2 %>%
  summarise(total_pop_range = range(total_pop_over_12), total_pop_median = median(total_pop_over_12))
```

```
#County Data
```

```
countydem2 <- select(countydem, "name", "renter_occ","owner_occ",
                     "vacant", "hse_units")
```

```
names(countydem2) <-c("county_name","renters","owners",
                     "vacant_housing", "total_housing_units")
```

```
countydem2[is.na(countydem2)] <- 0
```

Using the function str() we see that the selected variables in countydem2 are chr, int, int, int, int. All column types should be fine as is for our study question. Converting NA cells to 0 we allow ourselves to manipulate

and describe entire numeric columns. We also changed column names to allow for easier/obvious verbiage.

#Descriptions

```
distinct(covidvax2, county, .keep_all = F)
distinct(covidvax2, date, .keep_all = F)
distinct(covidvax2, zip_code, .keep_all = F)
```

```
covidvax2 %>% group_by(zip_code) %>%
  summarise(fullyvax_range = range(fully_vaccinated), partiallyvax_range = range(partially_vaccinated))
```

'summarise()' has grouped output by 'zip_code'. You can override using the '.groups' argument.

Population Range and Median for COVID vaccinations for those over 12 years old.

```
covidvax2 %>% summarise(total_pop_range = range(total_pop_over_12),
  total_pop_median = median(total_pop_over_12))
```

*#Using the function str() we see that the selected variables
#in countydem2 are chr, int, int, int, int
#All column types should be fine as is for our study question
#Converting NA cells to 0 we allow ourselves to manipulate and describe
#entire numeric columns.
#Changed column names to easier/obvious verbiage*

Descriptions

```
distinct(countydem2, county_name, .keep_all = F)
```

```
countydem2 %>% summarise(renters_mean = mean(renters),
  owners_mean = mean(owners),
  vacant_mean = mean(vacant_housing), total_mean = mean(total_housing_units))
```

```
countydem2 %>% summarise(renters_range = range(renters),
  owners_range = range(owners),
  vacant_range = range(vacant_housing), total_range = range(total_housing_units))
```

Below, is the minimum, first quartile, median, third quartile, and max values.

```
fivenum(covidvax2$total_pop_over_12)
fivenum(covidvax2$fully_vaccinated)
fivenum(covidvax2$partially_vaccinated)
```

#Descriptions

#We get the mean & range for all numeric and integer vectors in countydem2 below

```
distinct(countydem2, county_name, .keep_all = F)
```

```
countydem2 %>%  
  summarise(renters_mean = mean(renters),  
            owners_mean = mean(owners),  
            vacant_mean = mean(vacant_housing),  
            total_mean = mean(total_housing_units))
```

```
countydem2 %>%  
  summarise(renters_range = range(renters),  
            owners_range = range(owners),  
            vacant_range = range(vacant_housing),  
            total_range = range(total_housing_units))
```

Below, we get the five number summary for data of interest in countydem2.

```
fivenum(countydem2$renters)
```

```
## [1]    140    5878   25140   87737 1696455
```

```
fivenum(countydem2$owners)
```

```
## [1]    357   12629   39306  123646 1544749
```

```
fivenum(countydem2$vacant_housing)
```

```
## [1]    827.0   3328.0   8580.5  18747.0 203872.0
```

```
fivenum(countydem2$total_housing_units)
```

```
## [1]   1760.0  23910.0  76183.5 233755.0 3445076.0
```

Milestone 3

#Covid Vaccination by County data manipulation (covidvax files)

```
unique(covidvax2$county)
```

```
covidvax2 <- covidvax2 %>% filter(county!="0")
```

```
unique(covidvax2$date)
```

```
covidvax3 <- covidvax2 %>%  
  mutate(month = case_when(  
    str_detect(date, "-01-") == T ~ "January",  
    str_detect(date, "-02-") == T ~ "February",  
    str_detect(date, "-03-") == T ~ "March",  
    str_detect(date, "-04-") == T ~ "April",  
    str_detect(date, "-05-") == T ~ "May",  
    str_detect(date, "-06-") == T ~ "June",  
    str_detect(date, "-07-") == T ~ "July",  
    str_detect(date, "-08-") == T ~ "August",  
    str_detect(date, "-09-") == T ~ "September"  
  ))
```

```
unique(covidvax3$month)
```

,

```
covidvax4 <- covidvax3 %>%  
  group_by(county,month)%>%  
  mutate(county_total_pop_over12 = sum(total_pop_over_12)) %>%  
  mutate(county_fully_vaccinated = sum(fully_vaccinated)) %>%  
  mutate(county_partially_vaccinated = sum(partially_vaccinated))
```

```
covidvax4 <- covidvax4 %>%  
  select("-zip_code",-"total_pop_over_12",-"fully_vaccinated",  
    -"partially_vaccinated")
```

#Ensuring we get the last date for each month's final total vaccinations

```
unique(covidvax4$date)
```

```
covidvax4 <- covidvax4 %>% filter(date == "2021-01-26" | date == "2021-02-23" |  
                                date == "2021-03-30" |  
                                date == "2021-04-27" |  
                                date == "2021-05-25" |  
                                date == "2021-06-29" |  
                                date == "2021-07-27" |  
                                date == "2021-08-31" |  
                                date == "2021-09-14")
```

*#Removing duplicate columns after combining the zip codes for each county
#and making column totals for each . Then we check to ensure we have 58 counties
#for each month.*

```
covidvax5 <- covidvax4 %>% distinct(county,month, .keep_all=T)
```

```
covidvax5 %>% group_by(month) %>% summarise(n())
```

```
covidvax5 <- covidvax5 %>%  
mutate(county_percent_fullvax =  
(county_fully_vaccinated)/(county_total_pop_over12)) %>%  
mutate(county_percent_partialvax =  
(county_partially_vaccinated)/(county_total_pop_over12))
```

```
#Checking that our percentages make sense.  
summary(covidvax5$county_percent_fullvax)  
summary(covidvax5$county_percent_partialvax)
```

County Demographic Data Editing

```
countydem3 <- countydem2 %>%  
  mutate(total_in_use_housing = total_housing_units - vacant_housing) %>%  
  mutate(renter_housing_prop = (renters/(total_in_use_housing))) %>%  
  mutate(owner_housing_prop = (owners/(total_in_use_housing))) %>%  
  mutate(renter_owner_ratio = renters/owners)
```


#Data Dictionary

```
str(countydem3)
str(covidvax5)
```

#County Demographics Data Set

Name: county_name Data type: character Description: A list of all 58 counties in California.

Name:renters Data type: integer Description: Number of renters in each county in California.

Name:owners Data type: integer Description: Number of owners in each county in California.

Name:vacant_housing Data type: integer Description: Number of empty homes in each county in California.

Name:total_housing_units Data type: integer Description: Number of total homes in each county in California.

Name:total_in_use_housing Data type: integer Description: Number of total homes currently in use in each county in California.

Name:renter_housing_prop Data type: number Description: Proportion of renters living in occupied homes in each county in California.

Name:owner_housing_prop Data type: number Description: Proportion of home owners living in occupied homes in each county in California.

Name:renter_owner_ratio Data type: number Description: Renter to Home owner ratio for each county in California.

#Covid Vaccination Rate Data Set

Name:county Data type: character Description: A list containing the 58 counties in California.

Name:date Data type: character Description:The dates when the data was compiled, ranging from January 5, 2021 to September 14, 2021.

Name:month Data type: character Description: The month in which the testing took place. Created to more easily combine data from specific dates into monthly variables for each county in the future.

Name:county_total_pop_over12 Data type: number Description: Total population over the age of 12 years old in each county in California for a specific month.

Name:county_fully_vaccinated Data type: number Description: Total number of people fully vaccinated over the age of 12 in each county in California for a specific month.

Name:county_fully_vaccinated Data type: number Description: Total number of people fully vaccinated over the age of 12 in each county in California for a specific month.

Name:county_partially_vaccinated Data type: number Description: Total number of people partially vaccinated over the age of 12 in each county in California for a specific month.

Name:county_percent_fullvax Data type: number Description: Percentage of people fully vaccinated over the age of 12 over the total population over 12 years old in each county in California for a specific month.

Name:county_percent_partialvax Data type: number Description: Percentage of people partially vaccinated over the age of 12 over the total population over 12 years old in each county in California for a specific month.

Descriptive Statistics Tables

```
descstatsdem <- countydem3 %>%
  summarise(total_in_use_housing_range = summary(total_in_use_housing),
            renter_housing_ratio_range = summary(renter_housing_prop),
            owner_housing_ratio_range = summary(owner_housing_prop),
            renter_owner_ratio_range = summary(renter_owner_ratio))
rownames(descstatsdem) <- c("Min", "1st Quartile", "Median", "Mean", "3rd Quartile", "Max")

descstatsdem2 <- countydem3 %>%
  summarise(total_renters = sum(renters),
            total_owners = sum(owners))
rownames(descstatsdem2) <- c("Total Count")

library(kableExtra)

##
## Attaching package: 'kableExtra'
## The following object is masked from 'package:dplyr':
##
##   group_rows

countydem_kable <- kable(descstatsdem, longtable=T, booktabs=T,
                        col.names = c("Total in Use Housing", "Renter/Housing Ratio",
                                      "Owner/Housing Ratio", "Renter/Owner Ratio"),
                        caption = "Summary of CA Housing Across All Counties in 2021")
countydem_kable
```

Table 1: Summary of CA Housing Across All Counties in 2021

	Total in Use Housing	Renter/Housing Ratio	Owner/Housing Ratio	Renter/Owner Ratio
Min	497.00	0.2311765	0.3575537	0.3006887
1st Quartile	19040.75	0.3426035	0.5753113	0.5211578
Median	70284.50	0.3854496	0.6145504	0.6272347
Mean	216853.41	0.3833366	0.6166634	0.6472471
3rd Quartile	207711.50	0.4246887	0.6573965	0.7381916
Max	3241204.00	0.6424463	0.7688235	1.7967828

```
countydem_kable2 <- kable(descstatsdem2, longtable=T, booktabs=T,
                        col.names = c("Renters", "Owners"),
                        caption = "Total Number of Renters and Owners in CA in 2021")
countydem_kable2
```

Table 2: Total Number of Renters and Owners in CA in 2021

	Renters	Owners
Total Count	5542127	7035371

*#does not make sense to visualize covidvax dataset in a kable, will update visualization
#with graph using ggplot in milestone 4*

Milestone 4

```
#Creating a combined data set.

countydem3$county_name <- toupper(countydem3$county_name)

combined <- right_join(countydem3, covidvax5, by=c("county_name" = "county"))

combined <- combined %>% filter(county_name == "SAN FRANCISCO" |
                                county_name == "LOS ANGELES" |
                                county_name == "MONTEREY" |
                                county_name == "SANTA BARBARA" |
                                county_name == "YOLO" |
                                county_name == "SIERRA" |
                                county_name == "NEVADA" |
                                county_name == "EL DORADO" |
                                county_name == "AMADOR" |
                                county_name == "CALAVERAS") %>%
  select(county_name, month, renter_owner_ratio, county_percent_fullvax)

library(tidyr)

combined1 <- combined %>% pivot_wider(names_from = month,
                                       values_from = county_percent_fullvax)

combined1 <- combined1 %>% arrange(desc(renter_owner_ratio))

combined2 <- combined1 %>% mutate_if(is.numeric, round, digits = 3)

combined_kable <- kable(combined2, longtable=T, booktabs=T,
                        col.names = c("County", "Renter/Owner Ratio",
                                      "Jan", "Feb", "March", "April",
                                      "May", "June",
                                      "July", "Aug", "Sep"),
                        caption = "COVID-19 Vaccination Rate by Renter/Owner
                                  Ratio in Top 5 and Bottom 5 Counties in CA in 2021") %>%
  kable_styling(bootstrap_options = c("striped", "hover"), full_width = F,
                font_size = 7) %>%
  row_spec(1:5, bold = TRUE, color = "#CB3155") %>%
  row_spec(6:10, bold = TRUE, color = "#28A69D") %>%
  footnote(general_title = "Note.",
           footnote_as_chunk = TRUE,
           threeparttable = TRUE,
           general = "Pink =Top 5 counties & Blue = Bottom 5 counties."
  )
```

combined_kable

Table 3: COVID-19 Vaccination Rate by Renter/Owner Ratio in Top 5 and Bottom 5 Counties in CA in 2021

County	Renter/Owner Ratio	Jan	Feb	March	April	May	June	July	Aug	Sep
SAN FRANCISCO	1.797	0.009	0.048	0.188	0.419	0.618	0.733	0.768	0.789	0.806
LOS ANGELES	1.098	0.008	0.048	0.152	0.318	0.477	0.576	0.622	0.654	0.682
MONTEREY	0.966	0.006	0.033	0.124	0.309	0.489	0.587	0.633	0.664	0.693
SANTA BARBARA	0.899	0.006	0.044	0.139	0.300	0.483	0.569	0.605	0.631	0.656
YOLO	0.894	0.010	0.051	0.177	0.356	0.497	0.595	0.631	0.654	0.676
SIERRA	0.392	0.006	0.082	0.253	0.365	0.276	0.295	0.272	0.333	0.363
NEVADA	0.389	0.005	0.035	0.163	0.311	0.434	0.501	0.529	0.549	0.569
EL DORADO	0.366	0.007	0.050	0.169	0.316	0.426	0.507	0.540	0.563	0.584
AMADOR	0.339	0.006	0.034	0.140	0.292	0.382	0.429	0.449	0.465	0.484
CALAVERAS	0.301	0.006	0.035	0.148	0.278	0.369	0.414	0.435	0.454	0.474

Note. Pink = Top 5 counties & Blue = Bottom 5 counties.

This table generated via the kable package shows the vaccination rate at the end of each month for the time period provided in the dataset for the top 5 and bottom 5 renter-to-owner ratios and their counties. The table is in descending order of the ratios and as a visual guide, we have color coded specific rows to help split the top 5 from the bottom 5 counties and all of their corresponding data.

```
#Using our "combined" data set, we identify the top 5 and bottom 5 renter to
#owner ratios and their counties using a slice of the original dataset.
#Then we chose 0.75 as a somewhat "arbitrary" ratio to split the two separate
#groups based off of a simple visual analysis.

countydem4 <- countydem3 %>% arrange(desc(renter_owner_ratio))

countydem4subset <- countydem4 %>% slice(1:5, 54:58)

countydem4subset$county_name <- factor(countydem4subset$county_name,
                                       levels = countydem4subset$county_name[order(10:1)])

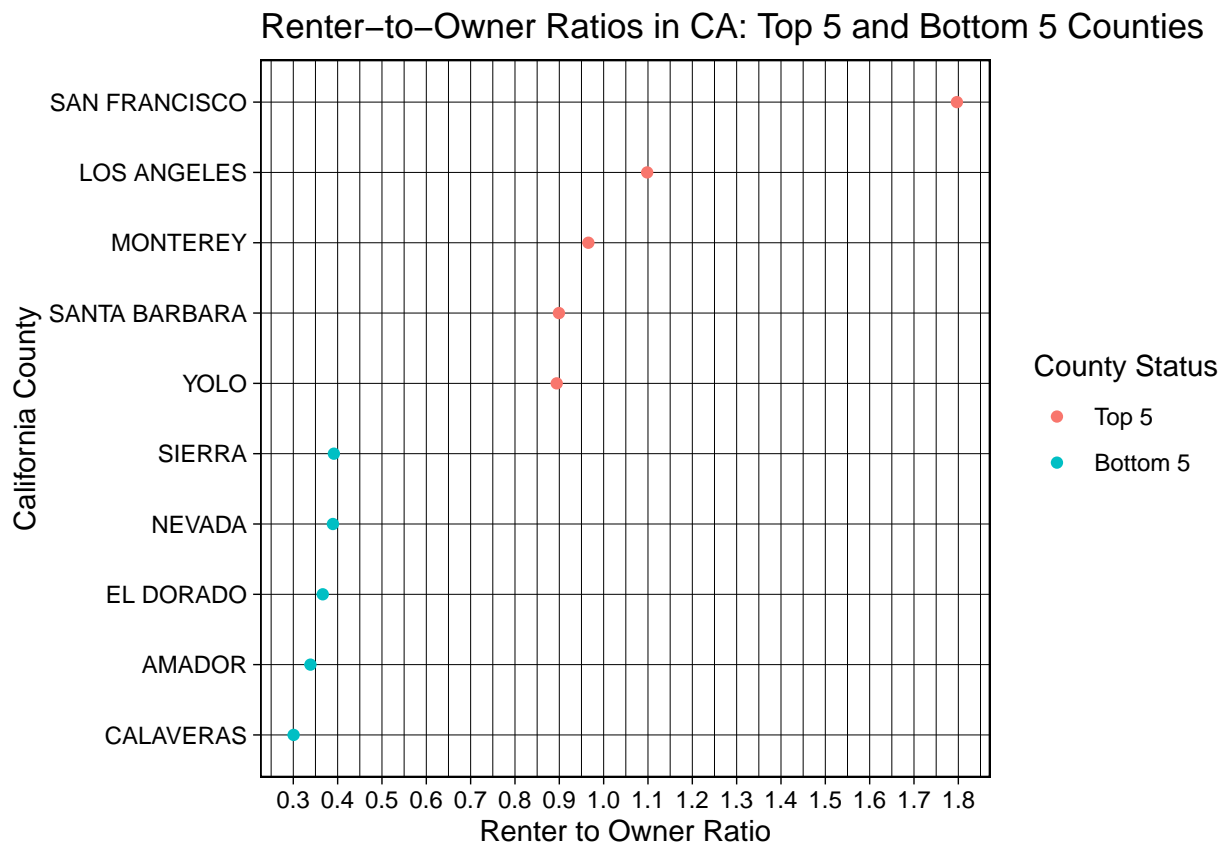
countydem4subset <- countydem4subset %>%
  mutate(status =
    case_when(
      renter_owner_ratio >= 0.75 ~ "Top 5",
      renter_owner_ratio < 0.75 ~ "Bottom 5"
    )
  )

countydem4subset$status <- factor(countydem4subset$status,
                                 levels = c("Top 5", "Bottom 5"))
```

```
library(ggplot2)

#Graph 1 creation

countydem4subset %>% ggplot(aes(y= county_name, x=renter_owner_ratio)) +
  geom_point(aes(color=status)) +
  labs(x = "Renter to Owner Ratio",
       title= "Renter-to-Owner Ratios in CA: Top 5 and Bottom 5 Counties ",
       y ="California County", fill ="Rank",
       color="County Status") +
  scale_x_continuous(breaks=seq(0,2,0.10)) +
  theme_linedraw()
```

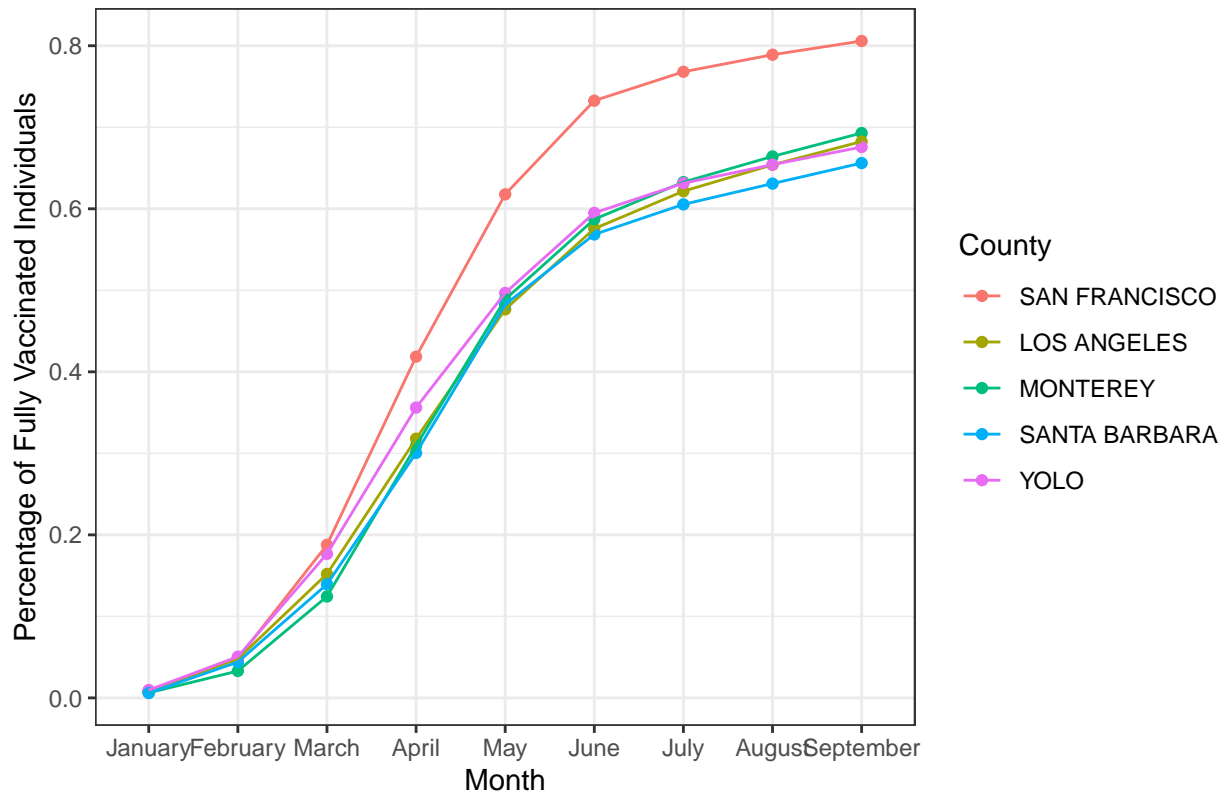


This graph provides a greater visual aid for the disparity of owning a home and the amount of renters for the top 5 and bottom 5 counties. As we would expect, we see a very drastic divide between Sierra county and Yolo county and that San Francisco is an outlier in this graph, with an extremely high amount of people renting housing compared to who actually owns homes.

#Graph 2 creation

```
covidvax6top5 <- covidvax5 %>% filter(county == "SAN FRANCISCO" |  
                                     county == "LOS ANGELES" |  
                                     county == "MONTEREY" |  
                                     county == "SANTA BARBARA" |  
                                     county == "YOLO")  
  
covidvax6top5 <- covidvax6top5 %>% select(-c(county_total_pop_over12,  
                                             county_fully_vaccinated,  
                                             county_partially_vaccinated))  
  
covidvax6bottom5 <- covidvax5 %>% filter(county == "SIERRA" |  
                                         county == "NEVADA" |  
                                         county == "EL DORADO" |  
                                         county == "AMADOR" |  
                                         county == "CALAVERAS")  
  
covidvax6bottom5 <- covidvax6bottom5 %>% select(-c(county_total_pop_over12,  
                                                    county_fully_vaccinated,  
                                                    county_partially_vaccinated))  
  
covidvax6top5$month <- factor(covidvax6top5$month,  
                             levels = c("January", "February", "March",  
                                         "April", "May", "June", "July",  
                                         "August", "September", "October",  
                                         "November", "December"))  
  
covidvax6top5$county <- factor(covidvax6top5$county,  
                              levels = c("SAN FRANCISCO", "LOS ANGELES",  
                                          "MONTEREY", "SANTA BARBARA",  
                                          "YOLO"))  
  
ggplot(data = covidvax6top5, aes(x=month, y=county_percent_fullvax,  
                                color = county, group = county)) +  
  geom_point() +  
  geom_line() +  
  labs(x="Month", y="Percentage of Fully Vaccinated Individuals",  
       title="COVID-19 Vaccination Rates: Top 5 Counties by Renter/Owner Ratio",  
       color="County") +  
  theme_bw()
```

COVID-19 Vaccination Rates: Top 5 Counties by Renter/Owner Ratio

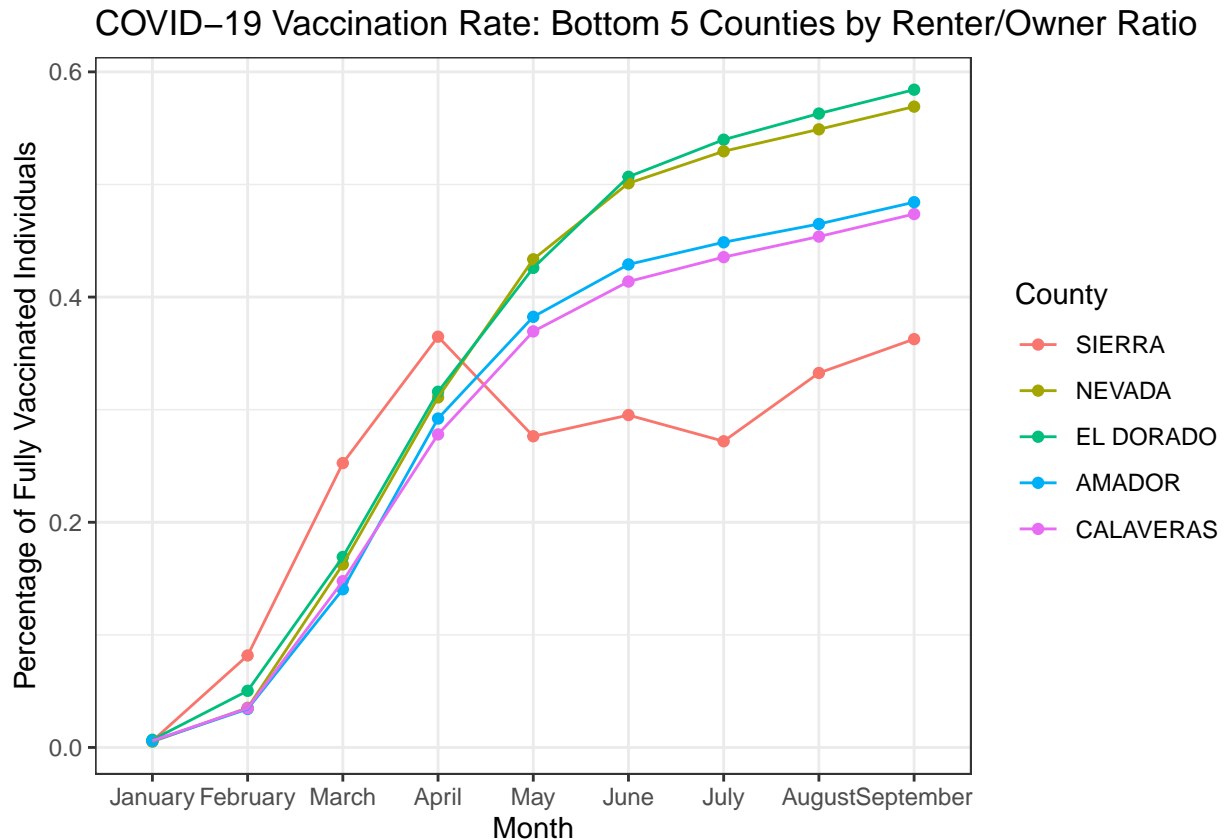


This graph shows the vaccination rates for the top 5 renter-to-owner ratio counties in California over January to September. San Francisco consistently had the highest percentage of fully vaccinated individuals starting in March and reached an astonishing 80% coverage of its “over 12” population. Other counties had a very similar pattern of coverage during this period.


```
covidvax6bottom5$month <- factor(covidvax6bottom5$month,
                                levels = c("January", "February", "March",
                                             "April", "May", "June", "July",
                                             "August", "September", "October",
                                             "November", "December"))

covidvax6bottom5$county <- factor(covidvax6bottom5$county,
                                  levels = c("SIERRA", "NEVADA",
                                               "EL DORADO", "AMADOR",
                                               "CALAVERAS"))

ggplot(data = covidvax6bottom5, aes(x=month, y=county_percent_fullvax,
                                     color = county, group = county)) +
  geom_point() +
  geom_line() +
  labs(x="Month", y="Percentage of Fully Vaccinated Individuals",
       title="COVID-19 Vaccination Rate: Bottom 5 Counties by Renter/Owner Ratio",
       color="County") +
  theme_bw()
```



This graph shows the vaccination rates for the bottom 5 renter-to-owner ratio counties over January to September. The Sierra line looks very strange, as we lose coverage. This could potentially be due to low population density and movement out of the county perhaps. We can only speculate with the given data.