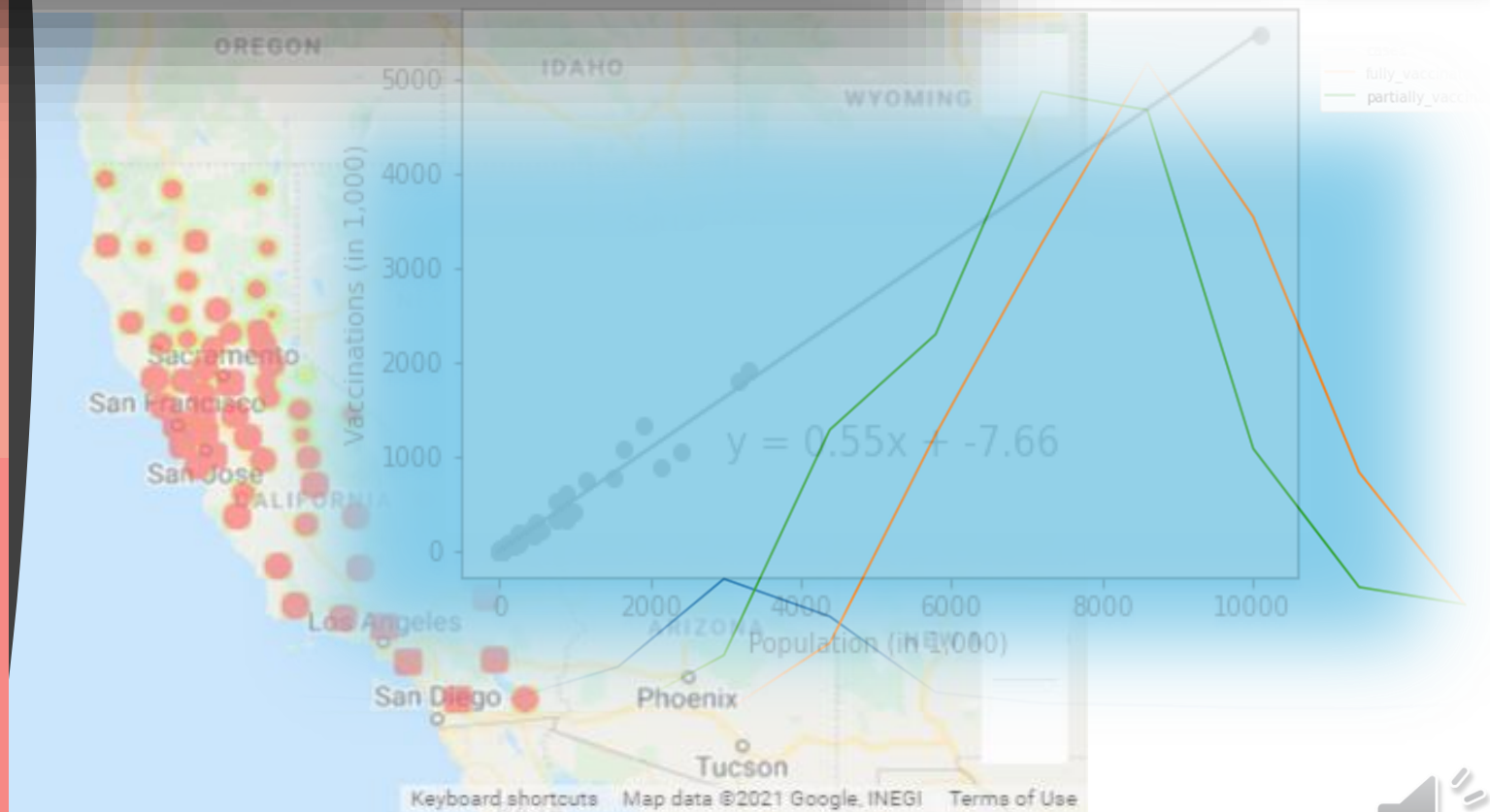


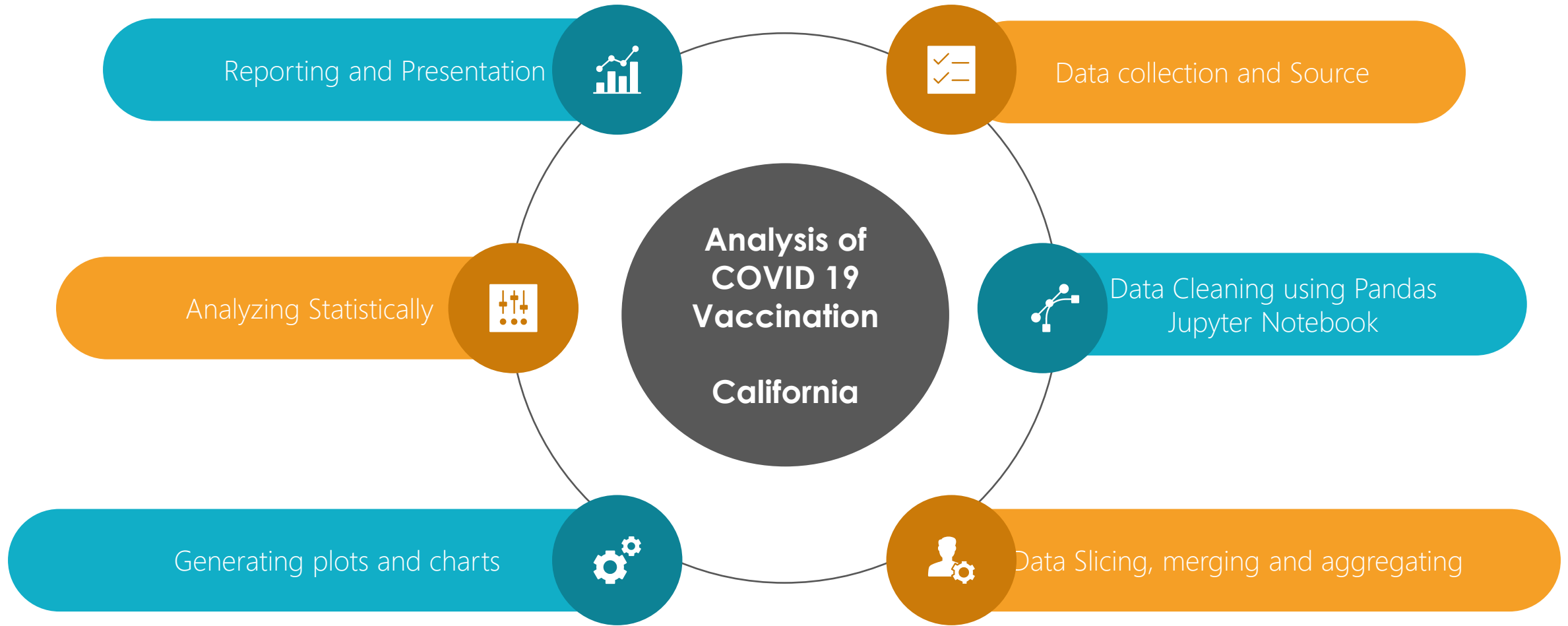
COVID-19 Vaccination California

Team members:

- Tikaram Subedy
- Kristian Hamilton
- David Koski



Project Outline



Project Plan



California Data

Cases, Vaccination and Demographic data



Data Cleaning and Preparation

Data Cleaning and Preparation using Python and Pandas in Jupyter Notebook



ANALYSIS

Performed Analysis on Overall California, and Counties in Cases, Demographics, and Vaccinations



Trends and Charts

Cases and Vaccination trends were analyzed. Fitted regression models (linear)



Report and Presentation

Prepared data analysis report. Recommendations and further



Summarize where and how we found the data used to answer these questions

- 1.Importing COVID19 data and preparing it for the analysis by slicing and aggregating.
- 2.Deciding on and calculating a good measure for analysis.
- 3.Merging datasets and finding correlations among the dataset.
- 4.Visualizing the analysis results using Pandas and Matplotlib.



What Affects Vaccination in CA?

How is the vaccination trend in overall CA as well counties?

Is there a relationship between cases and deaths that has affected in vaccination?

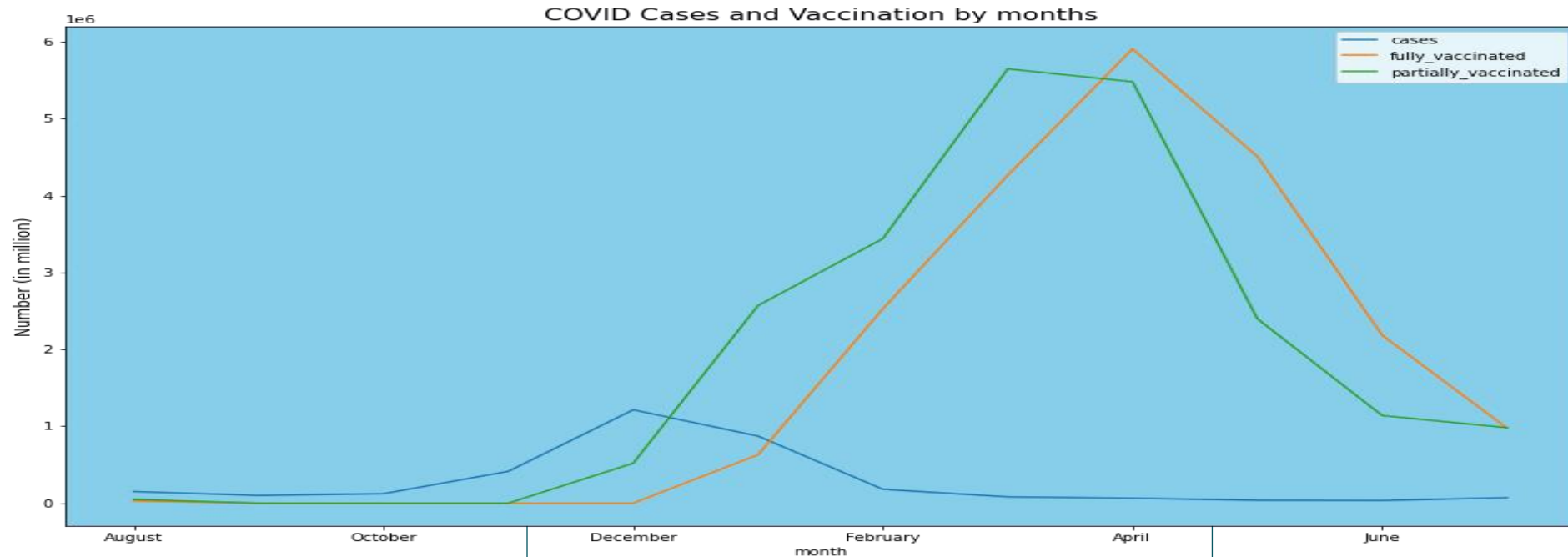
Do demographics affect on vaccination?

Does location and time affect vaccination?



Analysis

After the vaccination started in December and made available to wider range of population, the cases showed a dropping trend and remaining low.



Total Population in California

39,283,497

Completely Vaccinated

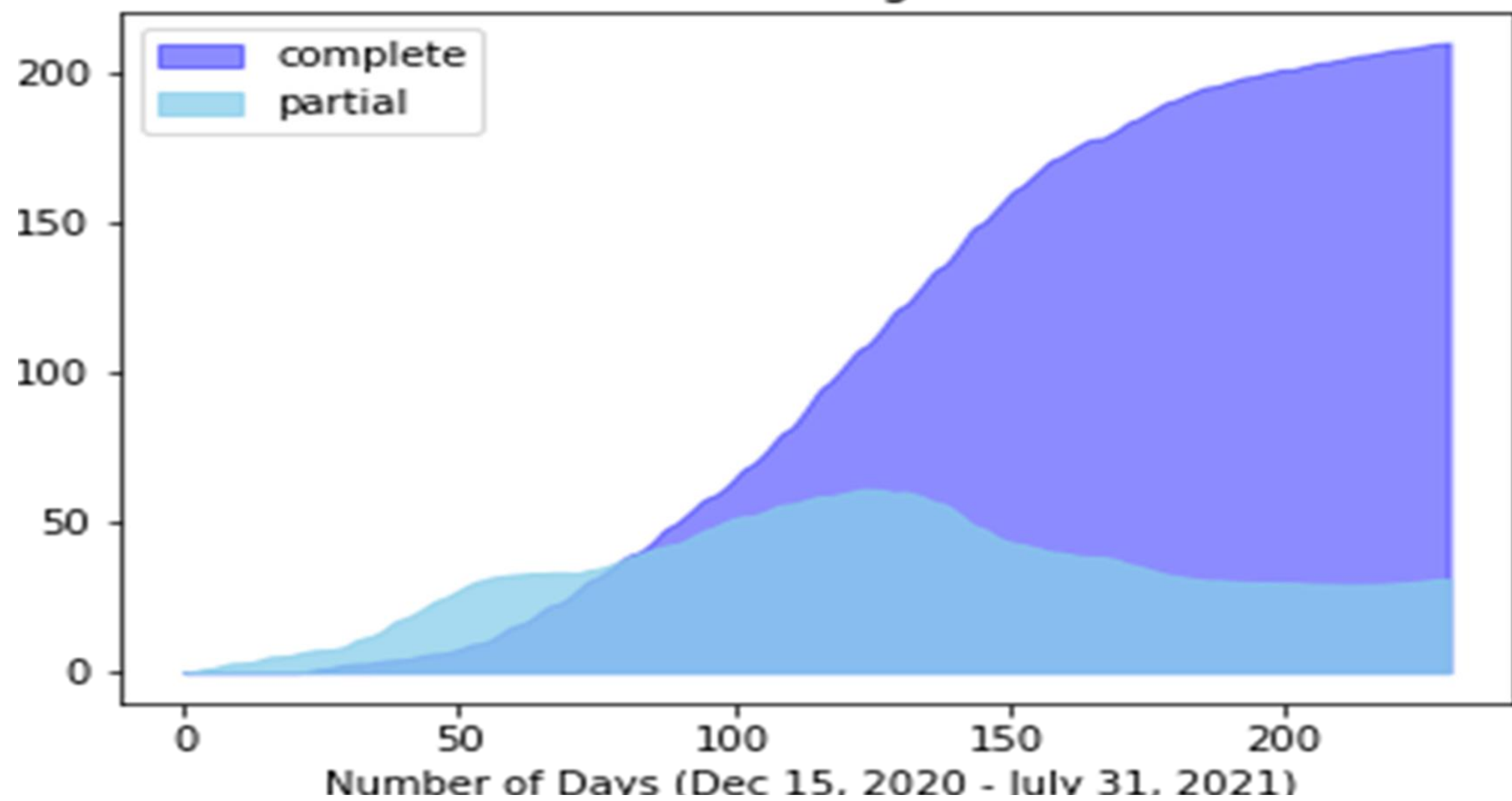
21,025,768
(53.52%)

Partially Vaccinated

22,224,679
(56.57%)

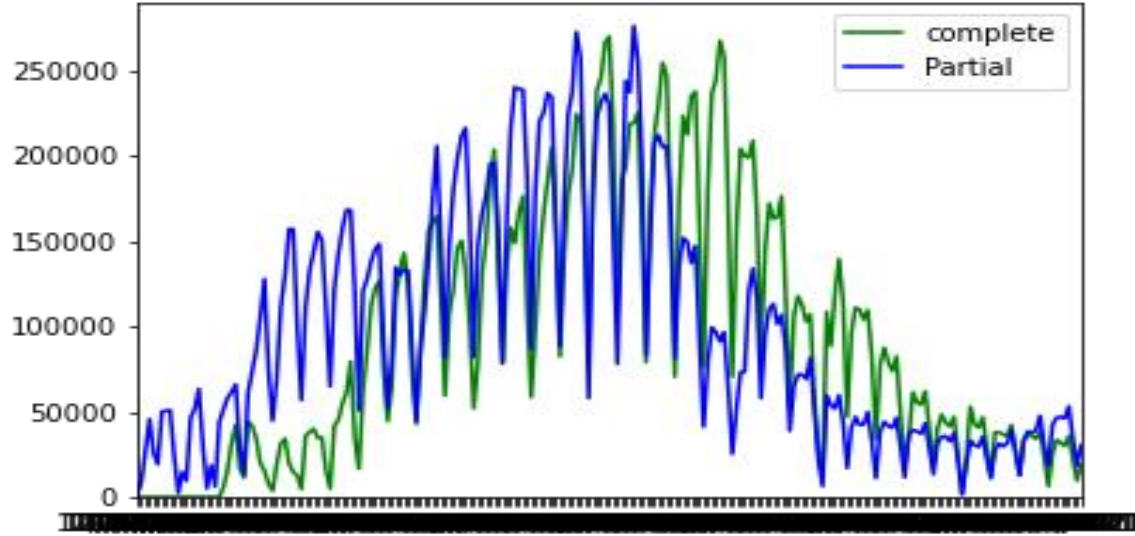


Vaccination Coverage Over Time



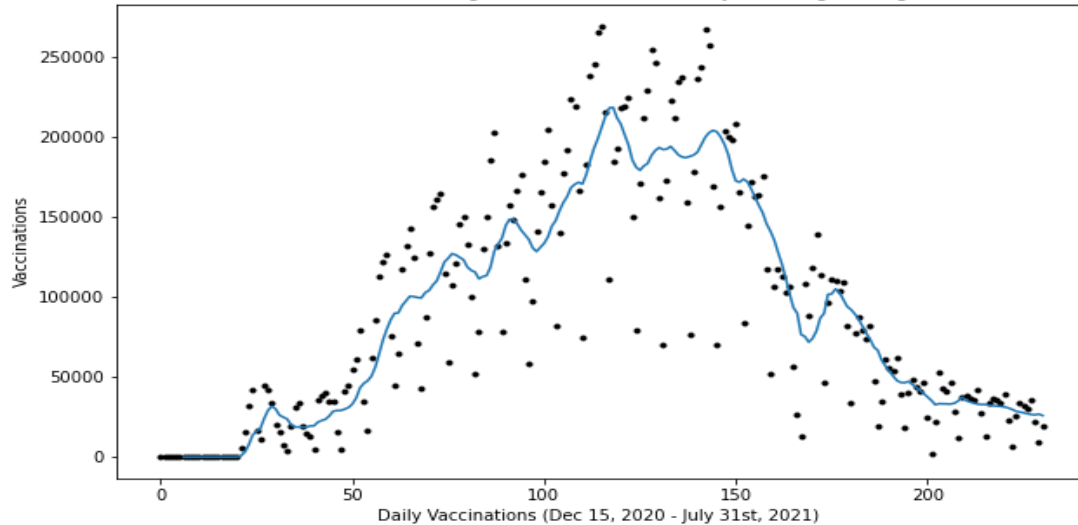
Analysis

Daily Vaccination Over Time

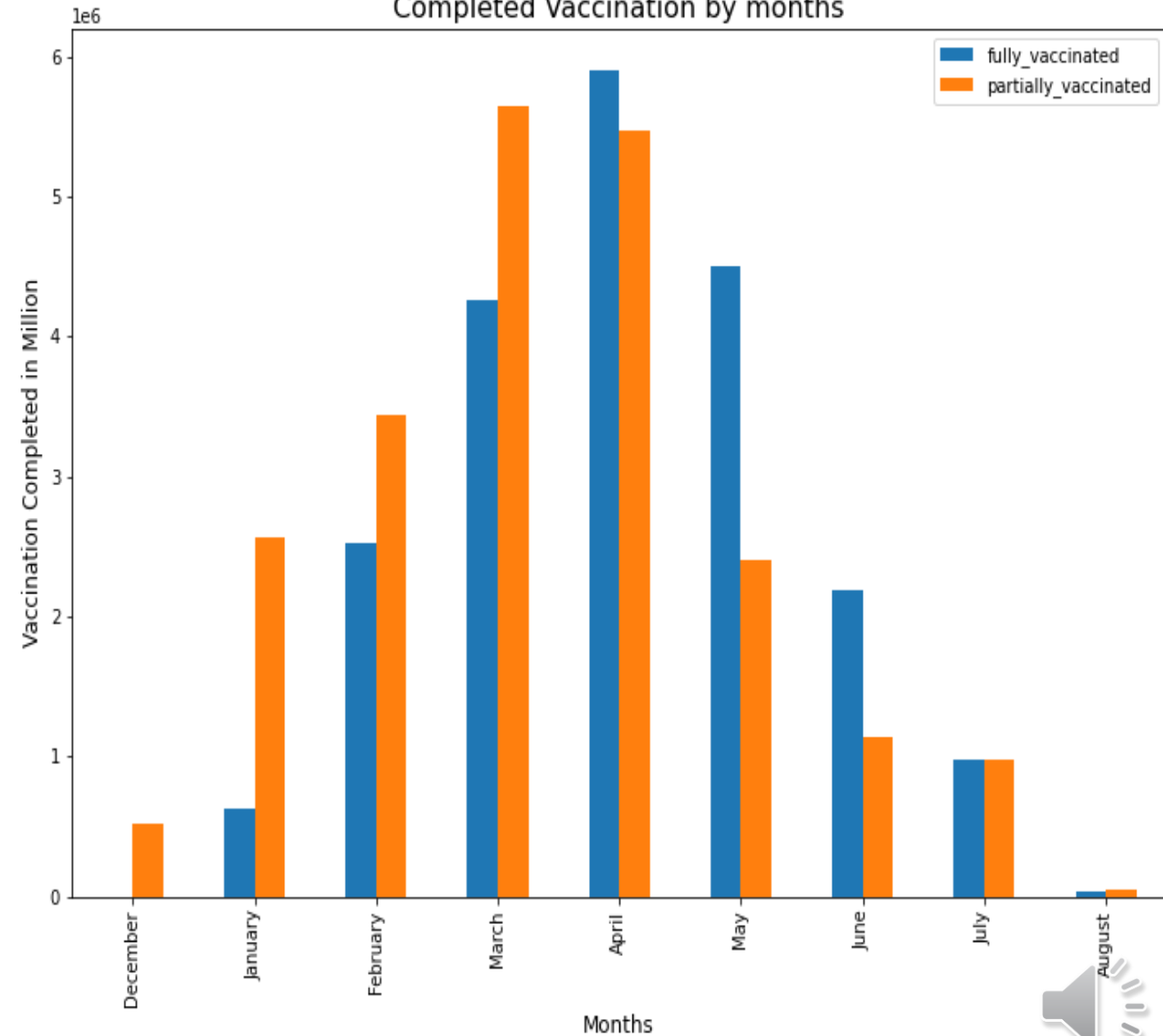


Vaccination Dates

Vaccination Coverage Over Time (with 7 days moving average)

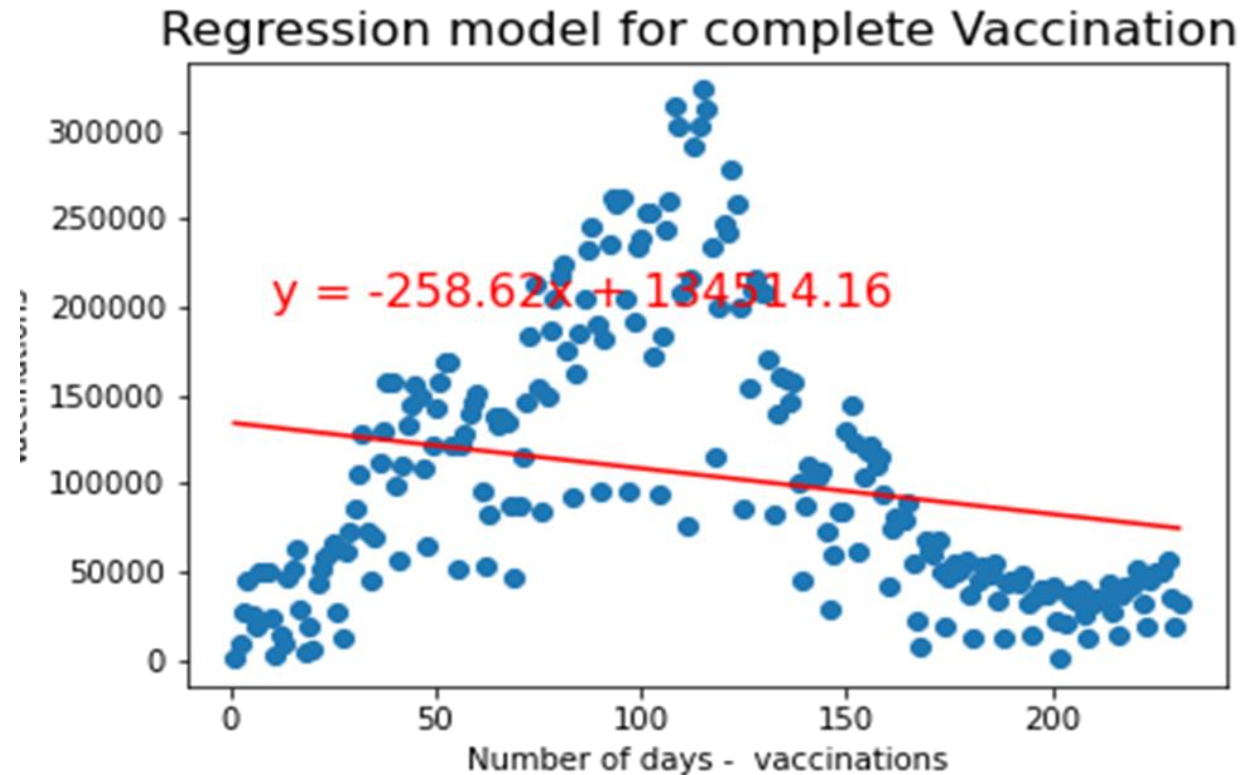


Completed Vaccination by months

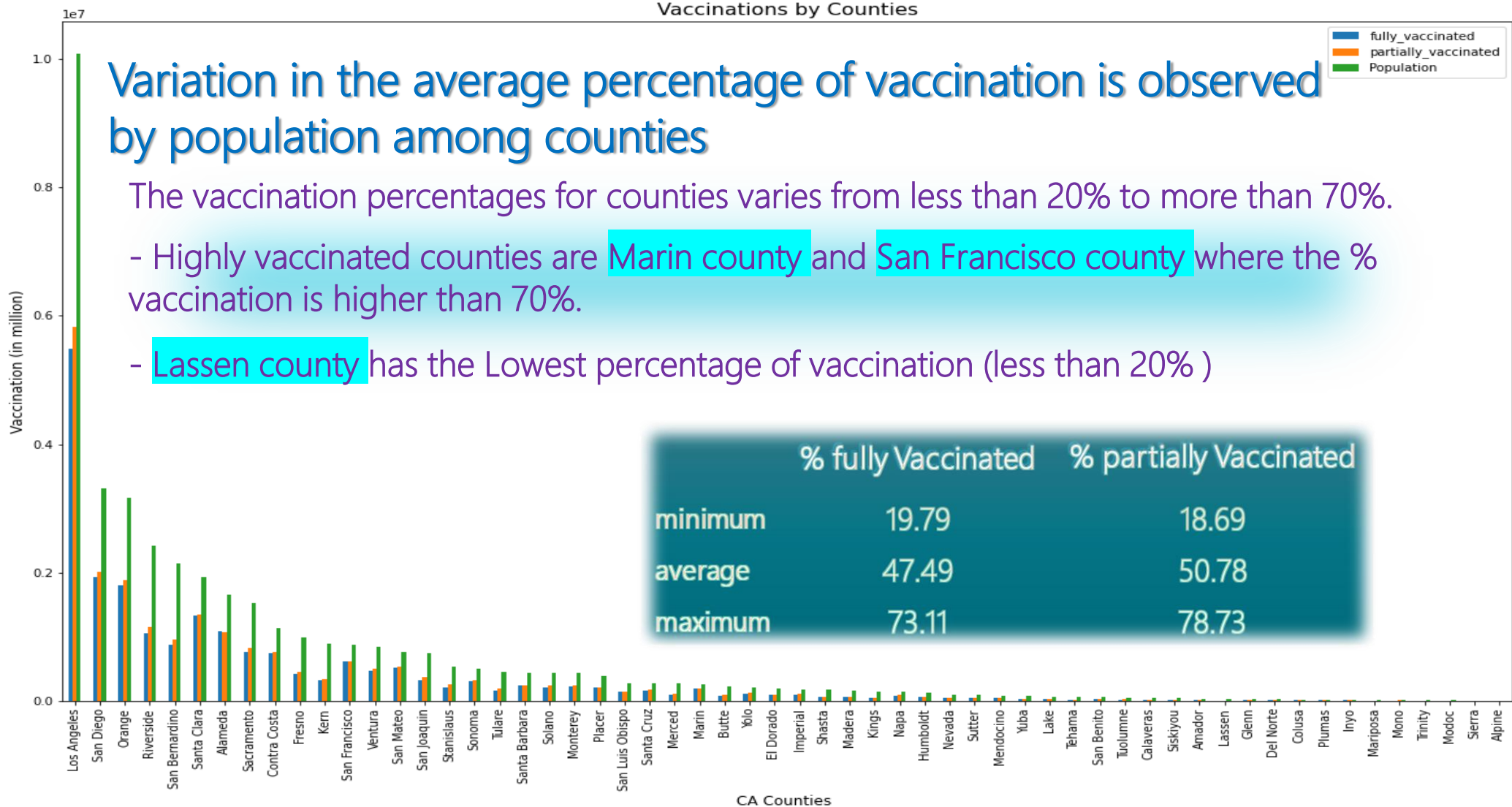


Analysis

The statistics shows, it can be estimated that the vaccination in California will be completed in about 521 days from the beginning, while all other compounding factors remain unchanged.



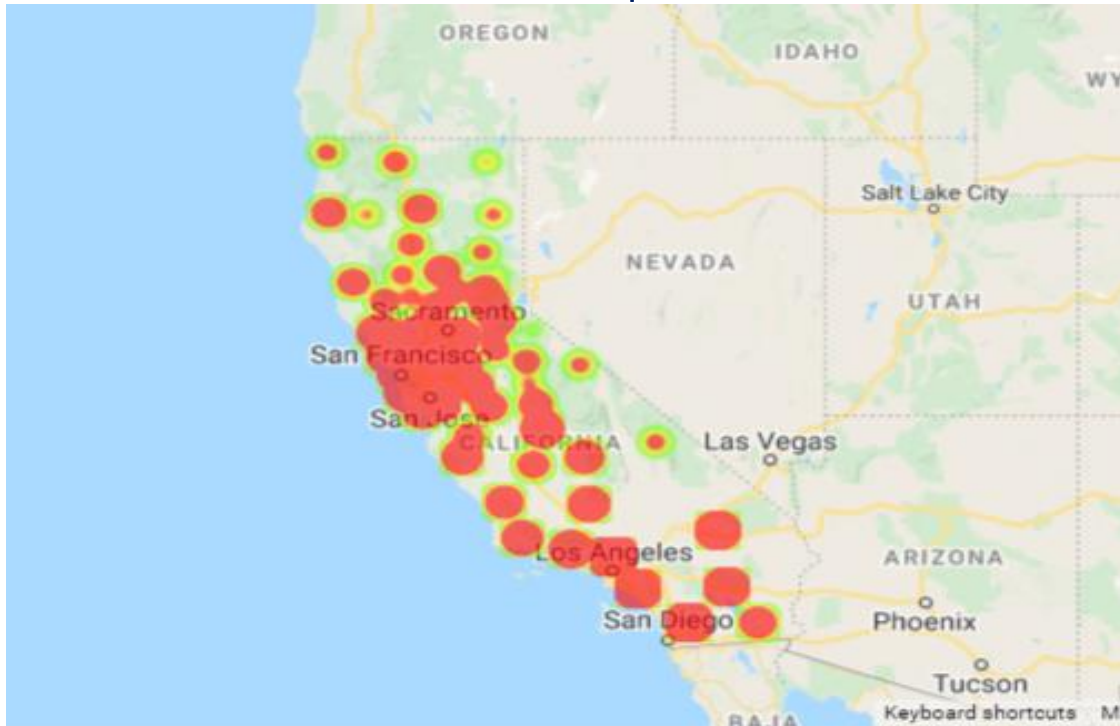
Analysis



Locations and vaccination

Urban and highly populated counties are found to have higher vaccination compared to the rural and less populated counties.

Counties - Complete Vaccination



Counties- higher average vaccination than State average



How strong is the correlation between cases/deaths, and vaccination rates?



CDC

- I downloaded all of the information on COVID cases from the CDC website, and then got to work cleaning it.



USAFacts.org

- I found the data on COVID deaths in California counties from USAFacts.org. I discovered this site after seeing that they were referenced as a trusted data source in an NPR (National Public Radio) article I was reading.



Data Cleanup David

- Describe the data exploration and cleanup process



CDC Cleaning

```
In [ ]: 1 import pandas as pd
        2 import json
        3 import matplotlib.pyplot as plt
        4 import numpy as np
        5 import os
        6 import csv
```

```
In [ ]: 1 #read in all CDC COVID data for the US as of early August
        2 csv_cdcpath = "./covid19_cdc_us_data.csv"
        3 cdc_data = pd.read_csv(csv_cdcpath, encoding="utf-8")
        4 cdc_data_df = pd.DataFrame(cdc_data)
        5 cdc_data_df.head()
```

```
In [ ]: 1 #filter to only include data for CA
        2 ca_cdc_data = cdc_data_df.loc[cdc_data_df['res_state'] == 'CA' ]
        3 ca_cdc_data.head()
```

```
In [ ]: 1 #filter again to only include month, state, county, age group, current status, sex, race, and ethnicity
        2 ca_cdc_data_short = ca_cdc_data[['case_month', 'res_state', 'res_county',
        3                                     'age_group', 'current_status', 'sex', 'race', 'ethnicity' ]]
        4 ca_cdc_data_short.head(2)
        5 #write the narrowed results to a csv and start a new notebook so that the computer
        6 #doesn't have to read a 4GB file every time the kernel is restarted
        7
        8 ca_output_data = 'ca_output_data.csv'
        9 ca_cdc_data_short.to_csv(ca_output_data, index = False)
```



CA CDC Data Cleaning

```
1 #extract only the data for August 2020 on, was originally just for 2021
2 ca_cases_aug20_on = ca_covid.loc[(ca_covid['case_month'] == '2020-08') |(ca_covid['case_month'] == '2020-09') |
3                                   (ca_covid['case_month'] == '2020-10') |(ca_covid['case_month'] == '2020-11') |
4                                   (ca_covid['case_month'] == '2020-12') |
5                                   (ca_covid['case_month'] == '2021-01') | (ca_covid['case_month'] == '2021-02') |
6                                   (ca_covid['case_month'] == '2021-03') | (ca_covid['case_month'] == '2021-04')
7                                   | (ca_covid['case_month'] == '2021-05') | (ca_covid['case_month'] == '2021-06')
8                                   | (ca_covid['case_month'] == '2021-07') ]
9
10 #extract just Dec 2020 on for group members to potentially use
11 ca_cases_2021 = ca_cases_aug20_on.copy()
12 ca_cases_2021 = ca_cases_2021.loc[(ca_cases_2021['case_month'] == '2020-12') |
13                                   (ca_cases_2021['case_month'] == '2021-01') | (ca_cases_2021['case_month'] == '2021-02') |
14                                   (ca_cases_2021['case_month'] == '2021-03') | (ca_cases_2021['case_month'] == '2021-04')
15                                   | (ca_cases_2021['case_month'] == '2021-05') | (ca_cases_2021['case_month'] == '2021-06')
16                                   | (ca_cases_2021['case_month'] == '2021-07') ]
17 ##print(ca_cases_aug20_on.head())
18 #print(ca_cases_2021.head())
```

```
1 #write out the modified dataframe
2 ca_cases_aug20_on.head()
3 ca_cases_aug20_on.to_csv('./Resources/CA_cases_by_county/ca_aug2020_on_case_data.csv',
4                           index=False,header=True)
```

```
1 #convert the case_month to month names so that it can merge with Tikaram
2 #use DatetimeIndex so that computer does not time out in for loop
3 ca_cases_aug20_on['month'] = pd.DatetimeIndex(ca_cases_aug20_on['case_month']).month_name()
4 ca_cases_2021['month'] = pd.DatetimeIndex(ca_cases_2021['case_month']).month_name()
```



CategoricalDtype Sorting

*#This cell is will organize the months we are analyzing by calendar placement
instead of alphabetical order when I plot the data*

```
month_order = CategoricalDtype(['December', 'January', 'February', 'March', 'April',  
                                'May', 'June', 'July'], ordered=True)
```

```
ca_merged_df['month'] = ca_merged_df['month'].astype(month_order)  
ca_merged_df = ca_merged_df.rename(columns={'month': 'Month', 'fully_vaccinated': 'Full Vaccinations',  
                                             'partially_vaccinated': 'Partial Vaccinations', 'cumulative_fully_vaccinated': 'Cumulative Fully Va  
                                             'cumulative_at_least_one_dose': 'People with at Least One Dose'})
```



Manually Code Cases per Month

```
1 #since since using groupby followed by .sum()  
2 #is taking so long, make variables for the case count each month  
3  
4 counts = ca_cases_aug20_on['month'].value_counts()  
5  
6 aug = counts['August']  
7 sept = counts['September']  
8 octo = counts['October']  
9 nov = counts['November']  
10 dec = counts['December']  
11 jan = counts['January']  
12 feb = counts['February']  
13 mar = counts['March']  
14 apr = counts['April']  
15 may = counts['May']  
16 june = counts['June']  
17 july = counts['July']  
18 month_list = counts.tolist()  
19 ordered_month_list = [151345,100714,123809,415153,1216142,873586,181288,82660,65037,37661,34899,71623]  
20 #print(month_list)  
21 print(counts)
```



Finding Monthly Death Totals

```
In [6]: 1 #create a dataframe for deaths by month for the entire state
2 #from December 2020 on so that it can be compared with vaccination/cases.
3 #since the original df is aggregates it is necessary to subtract from the sum of
4 #each month the sum of the month before it
5 dec = deaths_by_county['December'].sum() - deaths_by_county['November'].sum()
6 jan = deaths_by_county['January'].sum() - deaths_by_county['December'].sum()
7 feb = deaths_by_county['February'].sum() - deaths_by_county['January'].sum()
8 mar = deaths_by_county['March'].sum() - deaths_by_county['February'].sum()
9 apr = deaths_by_county['April'].sum() - deaths_by_county['March'].sum()
10 may = deaths_by_county['May'].sum() - deaths_by_county['April'].sum()
11 june = deaths_by_county['June'].sum() - deaths_by_county['May'].sum()
12 july = deaths_by_county['July'].sum() - deaths_by_county['June'].sum()
13
14 #print(dec)
15 #print(jan)
16 #print(feb)
17 #print(mar)
18 #print(apr)
19 #print(may)
20 #print(june)
21
22 ca_deaths_by_month = pd.DataFrame({
23     'month': ['December', 'January', 'February',
24             'March', 'April', 'May', 'June', 'July'],
25     'Deaths': [dec, jan, feb, mar, apr, may, june, july]
26 })
27
28 ca_deaths_by_month
```

Out[6]:

	month	Deaths
0	December	6171
1	January	15311
2	February	11282
3	March	5799



Merging the Cleaned Dataframes

```
2 vac_cleaned_grouped = vac_cleaned.groupby('month').sum()[['fully_vaccinated','partially_vaccinated']]
3 vac_cleaned_grouped_df = pd.DataFrame(vac_cleaned_grouped[['fully_vaccinated','partially_vaccinated']])
4
5 ca_merged_df = pd.merge(ca_deaths_by_month,vac_cleaned_grouped_df, on='month')
6 ca_merged_df = pd.merge(ca_merged_df,cases_by_county,on='month')
7 ca_merged_df = pd.merge(ca_merged_df,cumulative_fully_vac_grouped_df,on='month')
8 ca_merged_df
```

Out[11]:

	month	Deaths	fully_vaccinated	partially_vaccinated	Cases	cumulative_fully_vaccinated	cumulative_at_least_one_dose
0	December	6171	75	520210	1216142	77	539074
1	January	15311	631193	2571205	873586	648991	3188197
2	February	11282	2530930	3441366	181288	3237951	6783529
3	March	5799	4278021	5655057	82660	7592057	13100499
4	April	2689	5917068	5482020	65037	13603297	19420686
5	May	1329	4517693	2401439	37661	18173799	22118245
6	June	1078	2187579	1139761	34899	20388307	23475751
7	July	917	981594	986489	71623	21386589	24590667

In []:

```
1
```

In [12]:

```
1 #This cell is will organize the months we are analyzing by calendar placement
2 # instead of alphabetical order when I plot the data
3
4 month_order = CategoricalDtype(['December', 'January', 'February', 'March', 'April',
5                                'May', 'June', 'July'], ordered=True)
6
7 ca_merged_df['month'] = ca_merged_df['month'].astype(month_order)
8 ca_merged_df = ca_merged_df.rename(columns={'month':'Month', 'fully_vaccinated':'Full Vaccinations',
9                                             'partially_vaccinated':'Partial Vaccinations','cumulative_fully_vaccinated':'Cumulative Fully Va
10                                             'cumulative_at_least_one_dose':'People with at Least One Dose'})
11
```



Describe the Analysis Process

- Screen shots of code and output
- Describe the analysis process



Fully Vaccinated per Month vs. Fully Vaccinated Total Code

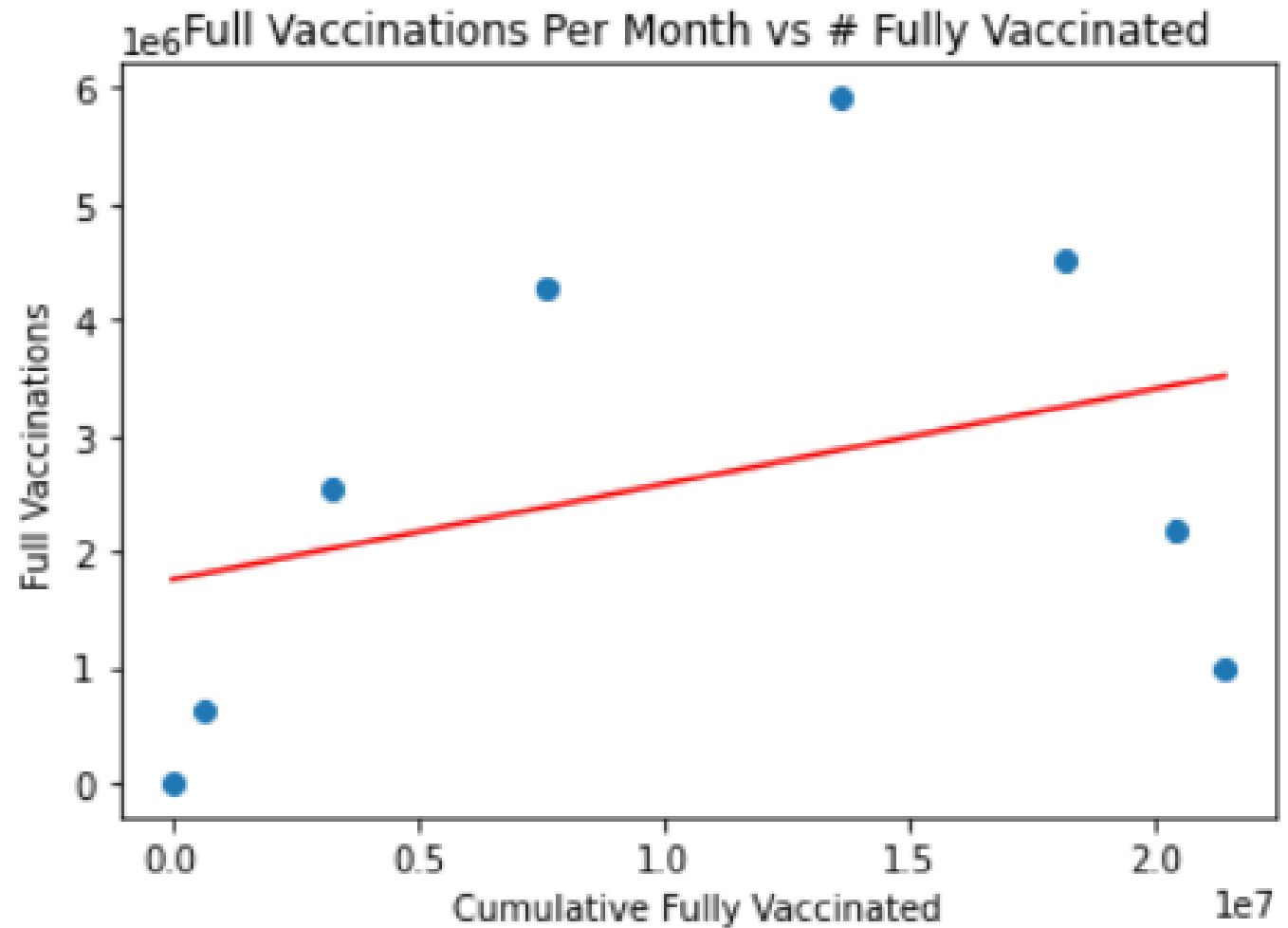
```
1 #show linear regression of vaccinations per month and total full vaccinations in California
2
3 x_values = ca_merged_df['Cumulative Fully Vaccinated']
4 y_values = ca_merged_df['Full Vaccinations']
5 (slope, intercept, rvalue, pvalue, stderr) = linregress(x_values, y_values)
6 regress_values = x_values * slope + intercept
7 line_eq = "y = " + str(round(slope,2)) + "x + " + str(round(intercept,2))
8 plt.scatter(x_values,y_values)|
9 plt.plot(x_values,regress_values,"r-")
10 #plt.annotate(line_eq,(6,10),fontsize=15,color="red")
11 plt.title("Full Vaccinations Per Month vs # Fully Vaccinated")
12 plt.xlabel('Cumulative Fully Vaccinated')
13 plt.ylabel('Full Vaccinations')
14 print(f"The r-value is: {rvalue}")
15 print(f"The r-squared is: {rvalue**2}")
16 plt.savefig('./Resources/CA_covid_deaths_cases_vaccines_combo/Images/linreg_fullper_month_#vaccded.png')
17 plt.show()
```



Fully Vaccinated Per Month vs. Fully Vaccinated Total

The r-value is: 0.3486571944487115

The r-squared is: 0.12156183924084664



Cases vs. Total With One or More Shots

Code

```
1 x_values = ca_merged_df['People with at Least One Dose']
2 y_values = ca_merged_df['Cases']
3 (slope, intercept, rvalue, pvalue, stderr) = linregress(x_values, y_values)
4 regress_values = x_values * slope + intercept
5 line_eq = "y = " + str(round(slope,2)) + "x + " + str(round(intercept,2))
6 plt.scatter(x_values,y_values)
7 plt.plot(x_values,regress_values,"r-")
8 #plt.annotate(line_eq,(6,10),fontsize=15,color="red")
9 plt.title('COVID Cases vs. People With at Least One Dose')
10 plt.xlabel('People with at Least One Dose')
11 plt.ylabel('COVID Cases')
12 print(f"The r-value is: {rvalue}")
13 print(f"The r-squared is: {rvalue**2}")
14 plt.savefig('./Resources/CA_covid_deaths_cases_vaccines_combo/Images/linreg_cases_#dosed.png')
15 plt.show()
```



Deaths vs. Total With One or More Shots

Code

```
1  #show linear regression of deaths and people with at least one shot
2
3  x_values = ca_merged_df['People with at Least One Dose']
4  y_values = ca_merged_df['Deaths']
5  (slope, intercept, rvalue, pvalue, stderr) = linregress(x_values, y_values)
6  regress_values = x_values * slope + intercept
7  line_eq = "y = " + str(round(slope,2)) + "x + " + str(round(intercept,2))
8  plt.scatter(x_values,y_values)
9  plt.plot(x_values,regress_values,"r-")
10 #plt.annotate(line_eq,(6,10),fontsize=15,color="red")
11 plt.title("COVID Deaths vs. People With at Least One Dose")
12 plt.xlabel('People with at Least One Dose')
13 plt.ylabel('COVID Deaths')
14 print(f"The r-value is: {rvalue}")
15 print(f"The r-squared is: {rvalue**2}")
16 plt.savefig('./Resources/CA_covid_deaths_cases_vaccines_combo/Images/linreg_deaths_#dosed.png')
17 plt.show()
```

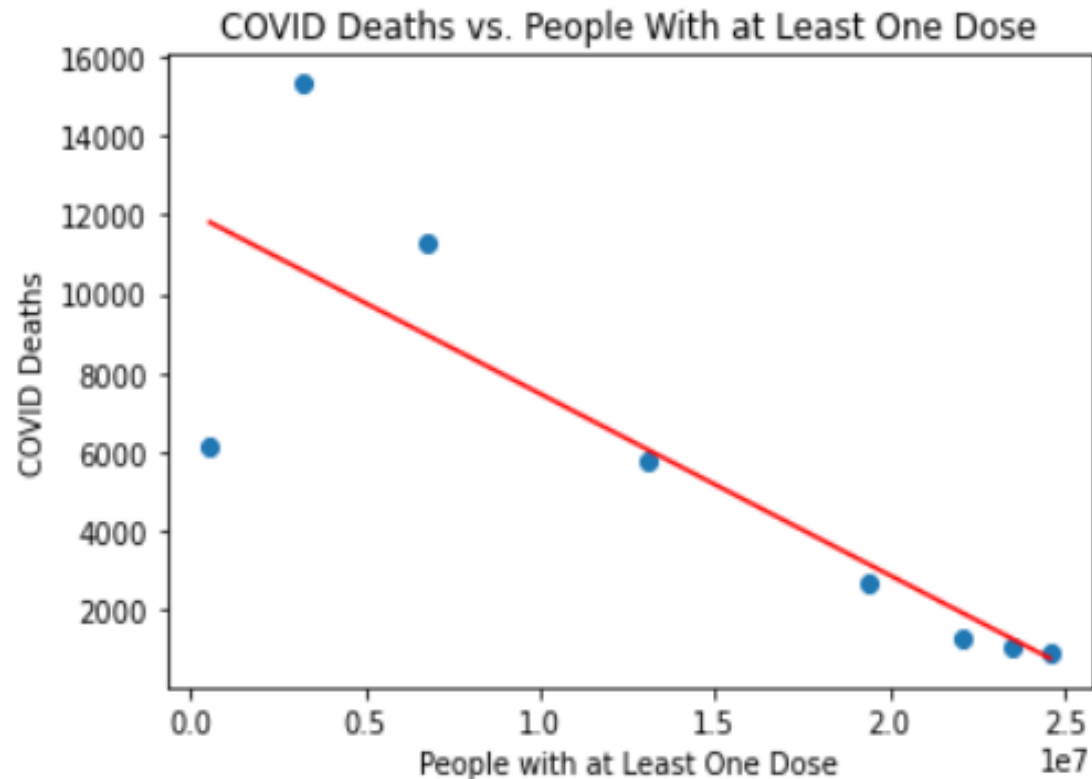
The r-value is: -0.8324963312710234

The r-squared is: 0.6930501415797136

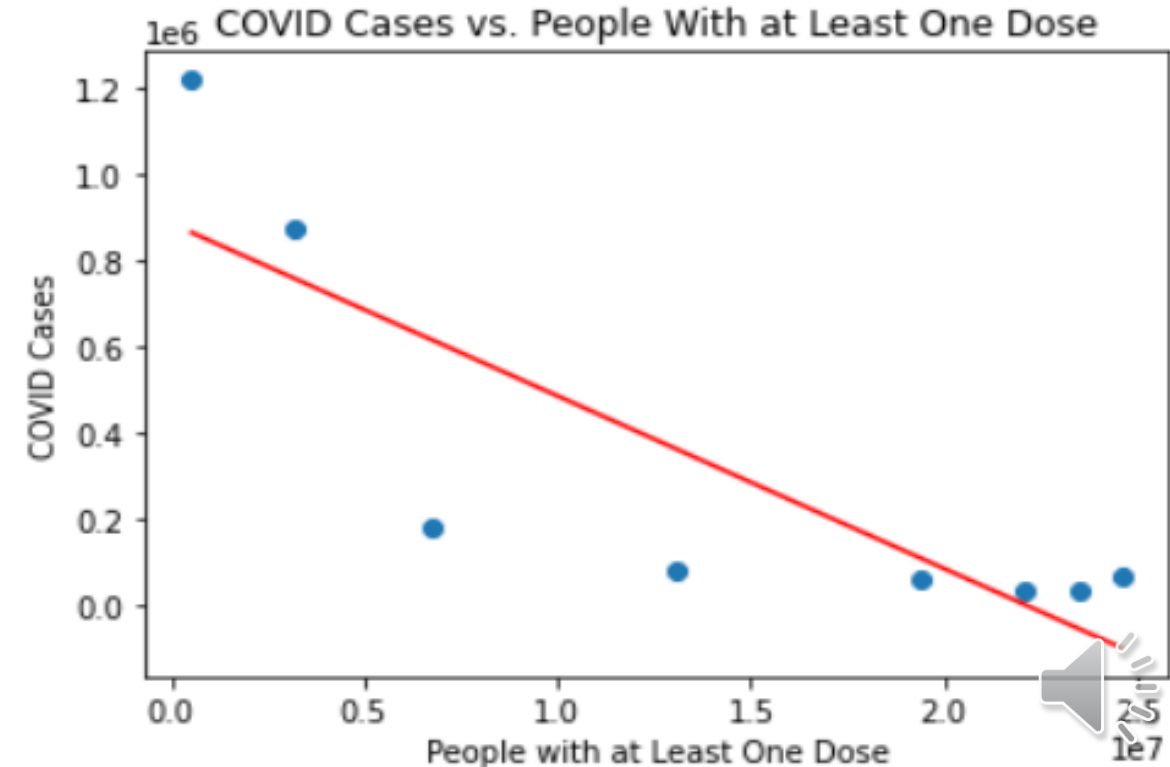


Cases and Deaths vs. Total With One or More Shots

The r-value is: -0.8324963312710234
The r-squared is: 0.6930501415797136




The r-value is: -0.8364258552960373
The r-squared is: 0.6996082114077076



Do demographics affect vaccination rates?



 The 2020 Redistricting Data will be available on data.census.gov no later than September 30th. The redistricting legacy format summary files are now available on the [FTP](#) site. Support materials and additional information are available on the [Redistricting Data Program's summary file webpage](#)

// Search / Tables / DP02

SELECTED SOCIAL CHARACTERISTICS IN THE UNITED STATES

Survey/Program: American Community Survey TableID: DP02 Product: 2019: ACS 5-Year Estimates Data Profiles 



Notes



Selections



1 Geo



Years



6 Topics



Surveys



Codes



Filter



Transpose

Margin of Error



Restore



Excel



Download



Print

More Data



Map

	California			
Label	Estimate	Margin of Error	Percent	Percent Margin of Error
▼ HOUSEHOLDS BY TYPE				
▼ Total households	13,044,266	±20,333	13,044,266	(X)
▼ Married-couple family	6,491,236	±30,996	49.8%	±0.2
With own children of the householder under 18 years	2,801,975	±22,451	21.5%	±0.1
▼ Cohabiting couple household	874,868	±6,012	6.7%	±0.1
With own children of the householder under 18 years	331,408	±4,216	2.5%	±0.1
▼ Male householder, no spouse/partner present	2,260,535	±12,994	17.3%	±0.1
With own children of the householder under 18 years	171,733	±3,524	1.3%	±0.1
▼ Householder living alone	1,390,613	±8,068	10.7%	±0.1
65 years and over	403,763	±4,528	3.1%	±0.1
▼ Female householder, no spouse/partner present	3,417,627	±10,396	26.2%	±0.1
With own children of the householder under 18 years	631,664	±6,030	4.8%	±0.1
▼ Householder living alone	1,715,491	±8,335	13.2%	±0.1
65 years and over	836,525	±5,822	6.4%	±0.1
Households with one or more people under 18 years	4,482,879	±21,313	34.4%	±0.1
Households with one or more people 65 years and over	3,803,822	±9,028	29.2%	±0.1
Average household size	2.95	±0.01	(X)	(X)
Average family size	3.53	±0.01	(X)	(X)
▼ RELATIONSHIP				
▼ Population in households	38,462,235	*****	38,462,235	(X)



Imported .csv as DataFrame:

]:

	GEO_ID	NAME	DP02_0001E	DP02_0001M	DP02_0001PE	DP02_0001PM	DP02_0002E	D
0	id	Geographic Area Name	Estimate!!HOUSEHOLDS BY TYPE!!Total households	Margin of Error!!HOUSEHOLDS BY TYPE!!Total hou...	Percent!!HOUSEHOLDS BY TYPE!!Total households	Percent Margin of Error!!HOUSEHOLDS BY TYPE!!T...	Estimate!!HOUSEHOLDS BY TYPE!!Total households...	Error!!HOI BY TYPE!
1	0500000US06001	Alameda County, California	577177	1744	577177	(X)	292079	
2	0500000US06003	Alpine County, California	350	69	350	(X)	194	
3	0500000US06005	Amador County, California	14594	448	14594	(X)	7954	
4	0500000US06007	Butte County, California	85320	891	85320	(X)	37211	

```
#Fixes header for the census dataframe
newHeader = censusDemo_df.iloc[0]
censusDemo_df = censusDemo_df[1:]
censusDemo_df.columns = newHeader
```

	id	Geographic Area Name	Estimate!!HOUSEHOLDS BY TYPE!!Total households	Margin of Error!!HOUSEHOLDS BY TYPE!!Total households	Percent!!HOUSEHOLDS BY TYPE!!Total households	Percent Margin of Error!!HOUSEHOLDS BY TYPE!!Total households	Estimate!!HOUSEHOLDS BY TYPE!!Total households!!Married-couple family	Errc hou
1	0500000US06001	Alameda County,	577177	1744	577177	(X)	292079	



```
#compile dataframe for use with plots
for x in countyList:
```

```
    county_df.loc[x, "Total Population"] =
        censusDemo_df.loc[
            censusDemo_df['Geographic Area Name'] == (x + ' County, California'), 'Estimate!!ANCESTRY!!Total population'
        ].sum()*
```

* **ValueError:** Incompatible indexer with Series

```
#data type for Total Pop needs to be cast as float instead of string to easily calc percentages
county_df['Total Population'] = county_df["Total Population"].apply(pd.to_numeric, downcast = 'float')
county_df.loc[x, "Total Fully Vaccinated"] = CHHS_df.loc[CHHS_df['county'] == x, "fully_vaccinated"].sum()
county_df["% of Pop Fully Vaccinated"] = (county_df["Total Fully Vaccinated"] / county_df["Total Population"]) * 100
```

county_df

	Total Population	Total Fully Vaccinated	% of Pop Fully Vaccinated	Avg People per Household	Total Households
County					
Alameda	1656754.0	1085847.0	65.540629	2.82	577177
Alpine	1039.0	689.0	66.313763	2.87	350
Amador	38429.0	15481.0	40.284681	2.38	14594
Butte	225817.0	88630.0	39.248595	2.57	85320
Calaveras	45514.0	18304.0	40.216197	2.66	16942
Colusa	21454.0	8758.0	40.822224	2.94	7227



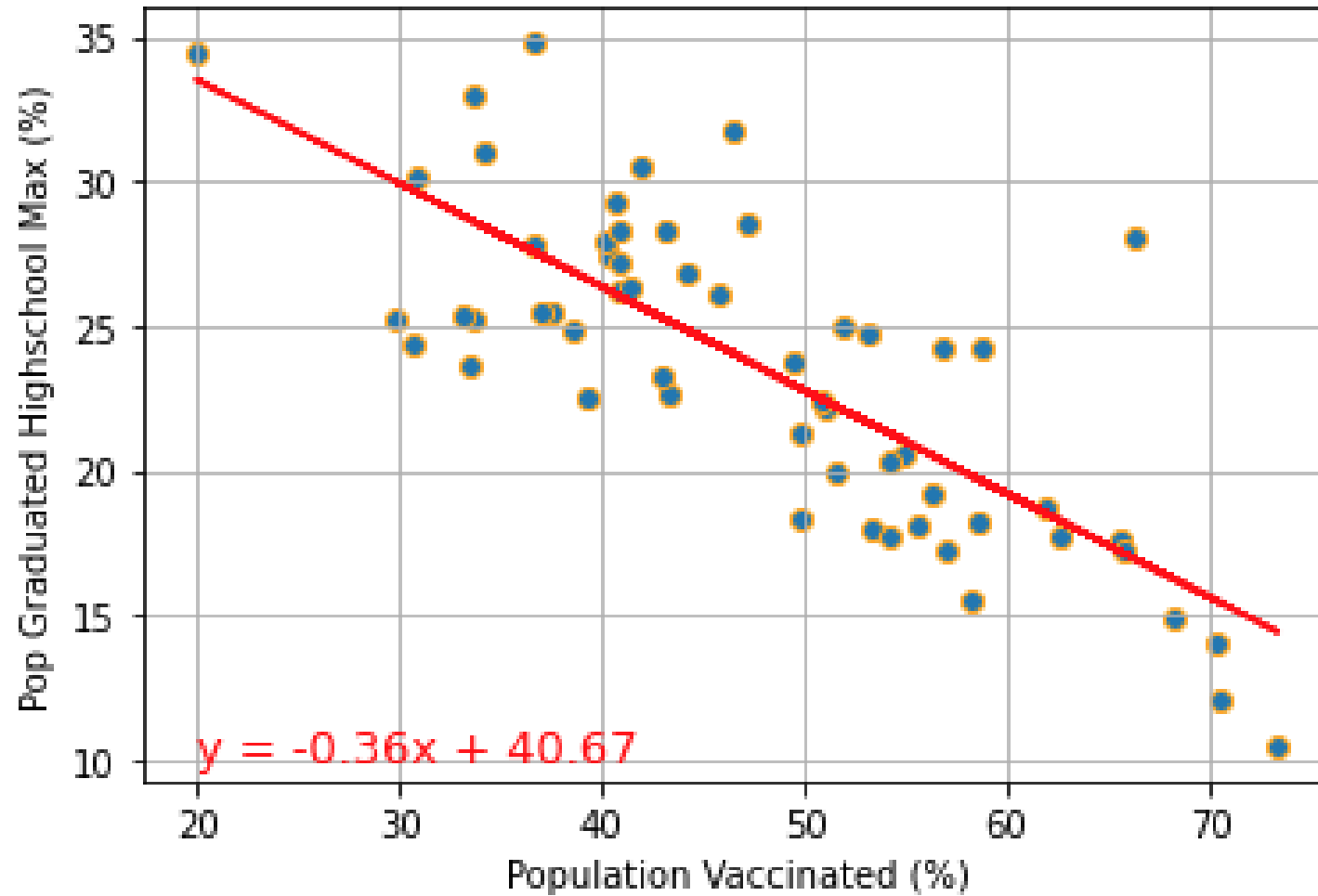
```
#function for scatter plot and linear regression
def plotScatterLinearRegression(xValues, yValues, title, xLabel, yLabel, textCoordinates):

    #regression formula
    (slope, intercept, rvalue, pvalue, stderr) = linregress(xValues, yValues)
    regressValues = xValues * slope + intercept
    lineEquation = "y = " + str(round(slope, 2)) + "x + " + str(round(intercept, 2))

    #scatter plot and regression line plot
    plt.scatter(xValues, yValues, edgecolor = "orange")
    plt.plot(xValues, regressValues, "r-", color = "red")
    plt.annotate(lineEquation, textCoordinates, fontsize = 13, color = "red")
    plt.xlabel(xLabel)
    plt.ylabel(yLabel)
    plt.title(title)
    print(f"R-squared = {round(rvalue,2)}")
    plt.grid()
    plt.show()
```

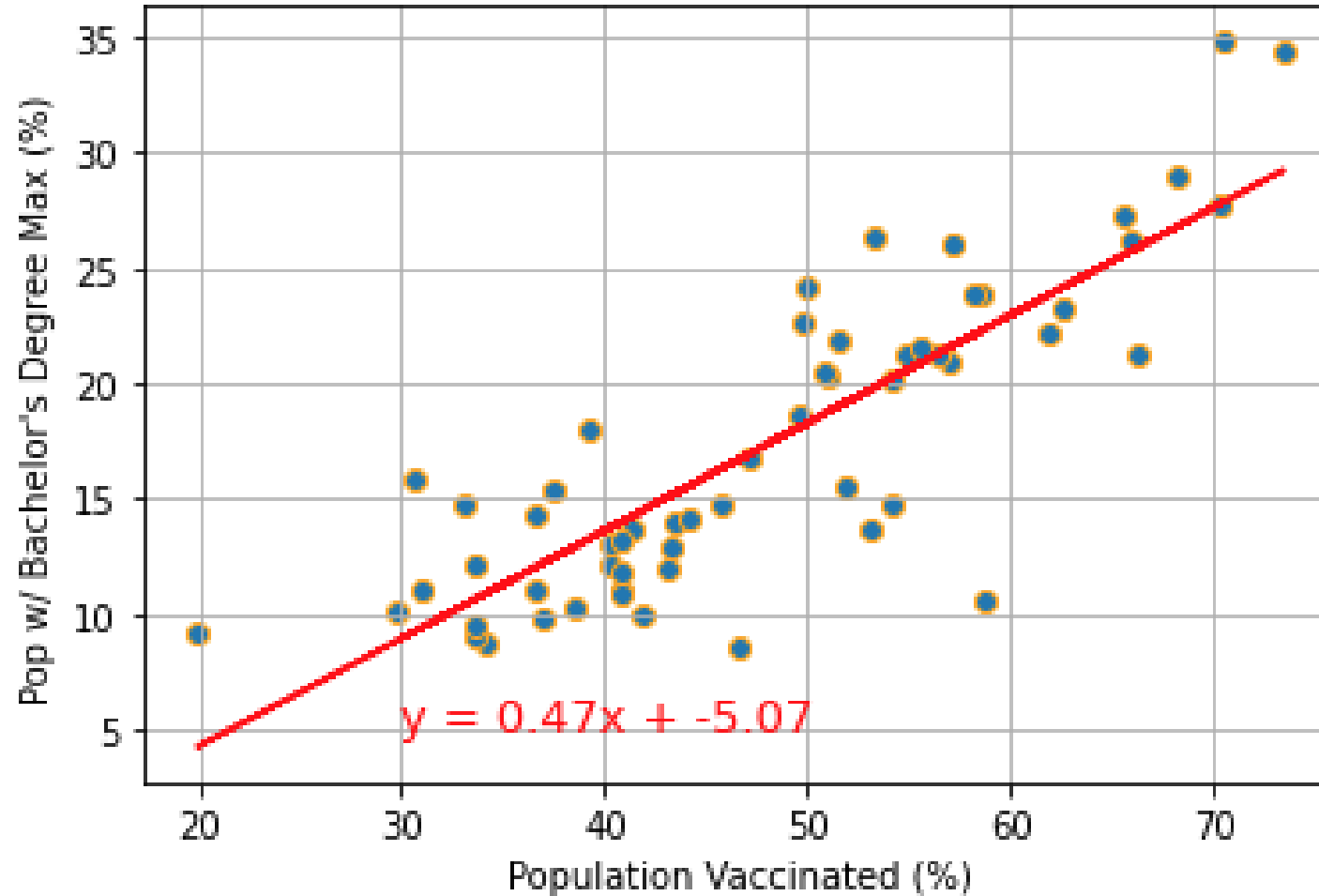


% of Pop, Over 25, Highest Education: High School Graduate by Vaccination Rate

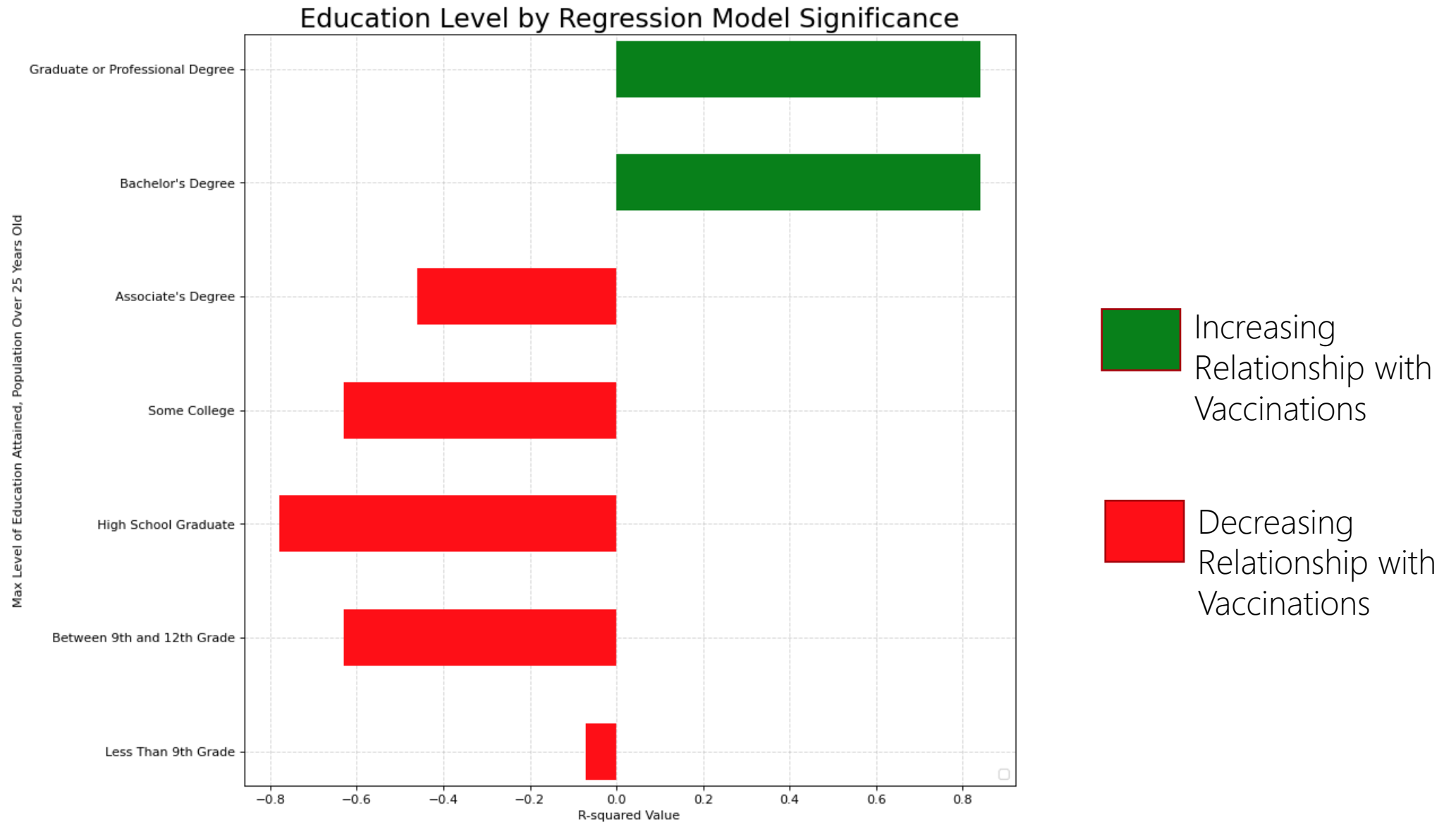


$R^2 = -0.78$

% of Pop, Over 25, Highest Education: Bachelor's Degree by Vaccination Rate



$R^2 = 0.84$



.xls available from website



Save As .csv and import as DataFrame



Report of Registration as of July 16, 2021 Registration by County

County	Eligible	Total Registered	Democratic	Republican	American Independent	Green
Alameda	1,078,848	944,570	566,482	103,575	20,359	5,281
Percent		87.55%	59.97%	10.97%	2.16%	0.56%
Alpine	921	907	394	218	47	7
Percent		98.48%	43.44%	24.04%	5.18%	0.77%
Amador	26,828	26,268	7,336	12,220	1,294	103
Percent		97.91%	27.93%	46.52%	4.93%	0.39%
Butte	150,098	125,414	44,403	44,757	5,490	718
Percent		83.55%	35.41%	35.69%	4.38%	0.57%
Calaveras	36,029	32,406	8,711	14,975	1,688	156
Percent		89.94%	26.88%	46.21%	5.21%	0.48%
Colusa	12,705	10,008	3,204	4,080	334	22
Percent		78.77%	32.01%	40.77%	3.34%	0.22%
Contra Costa	759,452	706,597	375,896	131,228	22,414	2,793
Percent		93.04%	53.20%	18.57%	3.17%	0.40%
Del Norte	18,126	14,879	4,471	5,862	779	96
Percent		82.09%	30.05%	39.40%	5.24%	0.65%

	County	Eligible	Total Registered	Democratic	Republican	American Independent	Green	Libertarian	Peace and Freedom	Unknown	Other	No Party Preference
0	Alameda	1,078,848	944,570	566,482	103,575	20,359	5,281	5,969	3,748	1	5,495	233,660
1	Percent	NaN	87.55%	59.97%	10.97%	2.16%	0.56%	0.63%	0.40%	0.00%	0.58%	24.74%
2	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
3	Alpine	921	907	394	218	47	7	18	2	1	3	217
4	Percent	NaN	98.48%	43.44%	24.04%	5.18%	0.77%	1.98%	0.22%	0.11%	0.33%	23.93%
...
173	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
174	State Total	24,834,943	22,078,290	10,264,695	5,309,040	710,259	88,011	212,686	109,963	118,865	132,121	5,132,650
175	Percent	NaN	88.90%	46.49%	24.05%	3.22%	0.40%	0.96%	0.50%	0.54%	0.60%	23.25%
176	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
177	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN



```
#compile dataframe for use with plots
for x in countyList:
```

```
    #Voter Registration Pull by County
```

```
    county_df.loc[x, "Total Registered Voters"] = caSOS_df.loc[caSOS_df['County'] == x, 'Total Registered'].sum()
    county_df.loc[x, "Democratic"] = caSOS_df.loc[caSOS_df['County'] == x, 'Democratic'].sum()
    county_df.loc[x, "Republican"] = caSOS_df.loc[caSOS_df['County'] == x, 'Republican'].sum()
    county_df.loc[x, "American Independent"] = caSOS_df.loc[caSOS_df['County'] == x, 'American Independent'].sum()
    county_df.loc[x, "Green"] = caSOS_df.loc[caSOS_df['County'] == x, 'Green'].sum()
    county_df.loc[x, "Libertarian"] = caSOS_df.loc[caSOS_df['County'] == x, 'Libertarian'].sum()
    county_df.loc[x, "Peace and Freedom"] = caSOS_df.loc[caSOS_df['County'] == x, 'Peace and Freedom'].sum()
    county_df.loc[x, "No Party Preference"] = caSOS_df.loc[caSOS_df['County'] == x, 'No Party Preference'].sum()
```

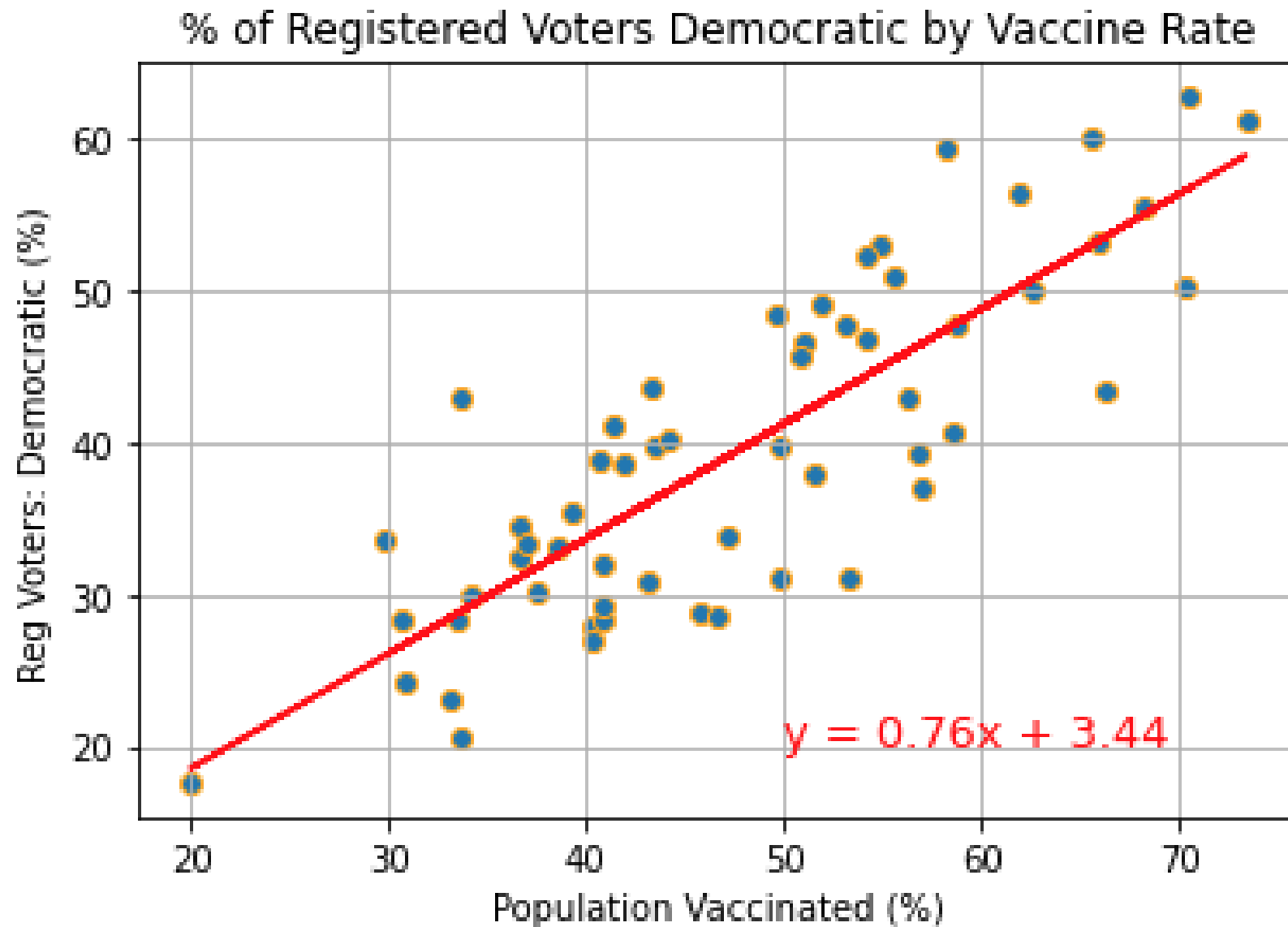
```
#Remove commas from dataframe so that it can be converted from string to float
county_df = county_df.replace(',', '', regex=True)
```

```
county_df["% of Reg Voters: No Party Preference"] =
    (county_df["No Party Preference"].apply(pd.to_numeric, downcast = 'float')
     / county_df["Total Registered Voters"].apply(pd.to_numeric, downcast = 'float')) * 100
```

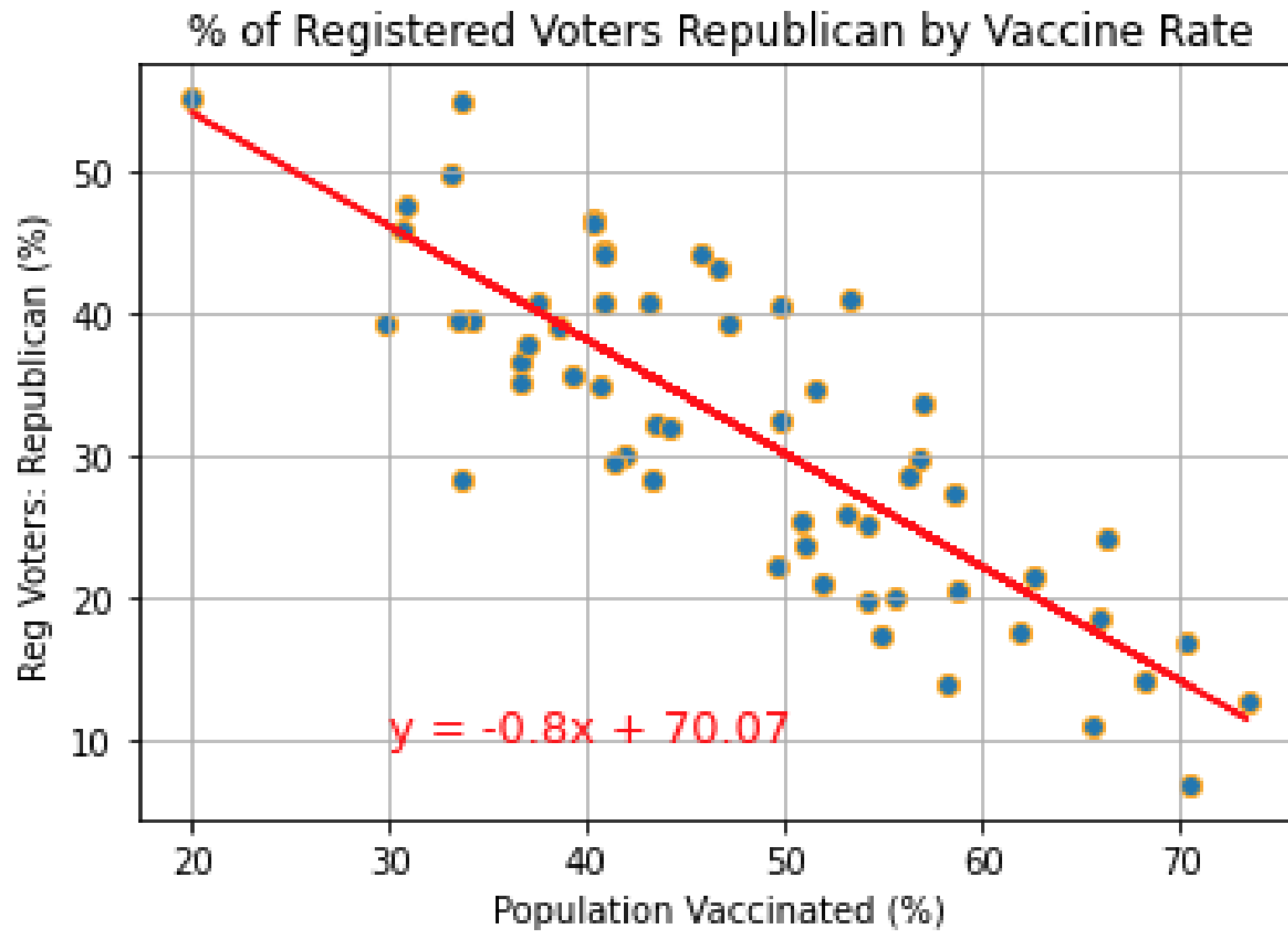
```
county_df["% of Reg Voters: No Party Preference"]
```


County	
Alameda	24.737182
Alpine	23.925028
Amador	17.740216
Butte	20.602166
Calaveras	17.728198
Colusa	21.772582
Contra Costa	22.667518

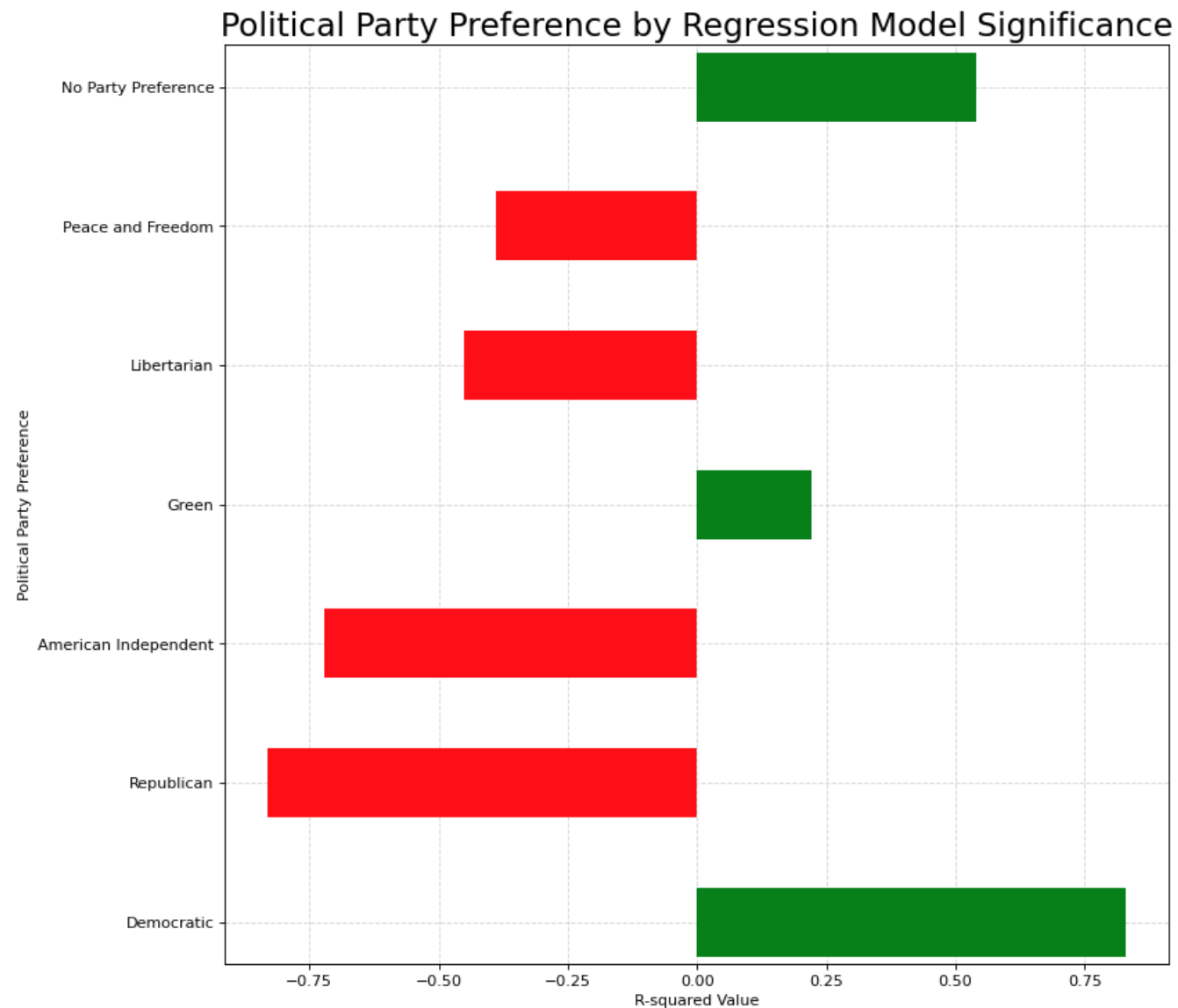




$R^2 = 0.83$



$$R^2 = -0.83$$




Increasing Relationship with Vaccinations

Decreasing Relationship with Vaccinations



Do demographics affect vaccination rates?

Strong R-squared values indicate yes.

Limitations:

- Census data is from 2019.

- Small sample size with only 58 California Counties.

Recommendations:

- It seems like policy encouraging people to take remote college courses may increase vaccination rates slightly in a county. Correlation obviously does not equal causation, however.





Thank You