

Project - COVID-19 Vaccination Progress

16 April, 2021

Introduction

With the aim of immunizing everyone with a vaccine, COVID-19 vaccination process has now started and progressed in many countries all around the world. There are some vaccines which are much more used than other vaccines in different countries. Some countries are very advanced in the process of vaccinating major percentage of its country population. The economic state is also different in every country due to this global pandemic. In this project, we are using the COVID-19 vaccination data of the world and web-scraping techniques to find answers to some of our key analysis questions related to COVID-19 vaccination progress around the world -

- (1) Which vaccine is mostly used in countries around the world?
- (2) Which countries have the highest daily average vaccinations?
- (3) Where are more people vaccinated per day?
- (4) What is the proportion of fully vaccinated people from entire population of a country?
- (5) Is the association between country GDP and total vaccinations statistically significant?

The purpose of this project is to learn the overall progress of vaccinations in different countries all around the world by analyzing these questions above.

Data and Methods

Primary data source

The primary data source of our project is the COVID-19 Vaccination data from Kaggle. Link - https://www.kaggle.com/gpreda/covid-world-vaccination-progress?select=country_vaccinations.csv

This dataset (country vaccinations) contains information about Country, Country ISO Code, Date, Total number of vaccinations, Total number of people vaccinated, Total number of people fully vaccinated, Daily vaccinations (raw), Daily vaccinations, Total vaccinations per hundred, Total number of people vaccinated per hundred, Total number of people fully vaccinated per hundred, Number of vaccinations per day, Daily vaccinations per million, Vaccines used in the country, Source name, Source website.

Secondary data sources

The secondary sources of data for our project contains information about country population and GDP. Since both these data are absent in our primary data source, we collected these data using web-scraping techniques from Worldometer website.

(1) Country Population - <https://www.worldometers.info/world-population/population-by-country>

(2) Country GDP - <https://www.worldometers.info/gdp/gdp-by-country/>

Methods

For our first three analysis questions, we used the primary data source, data visualization principles and summarized our results in bar charts to answer these questions as accurately as possible.

The last two analysis questions are answered using web-scraping techniques which required fair amount of data cleaning after scraping these data from our secondary data sources. Linear regression is also used to answer the last analysis question.

Data Analysis Results

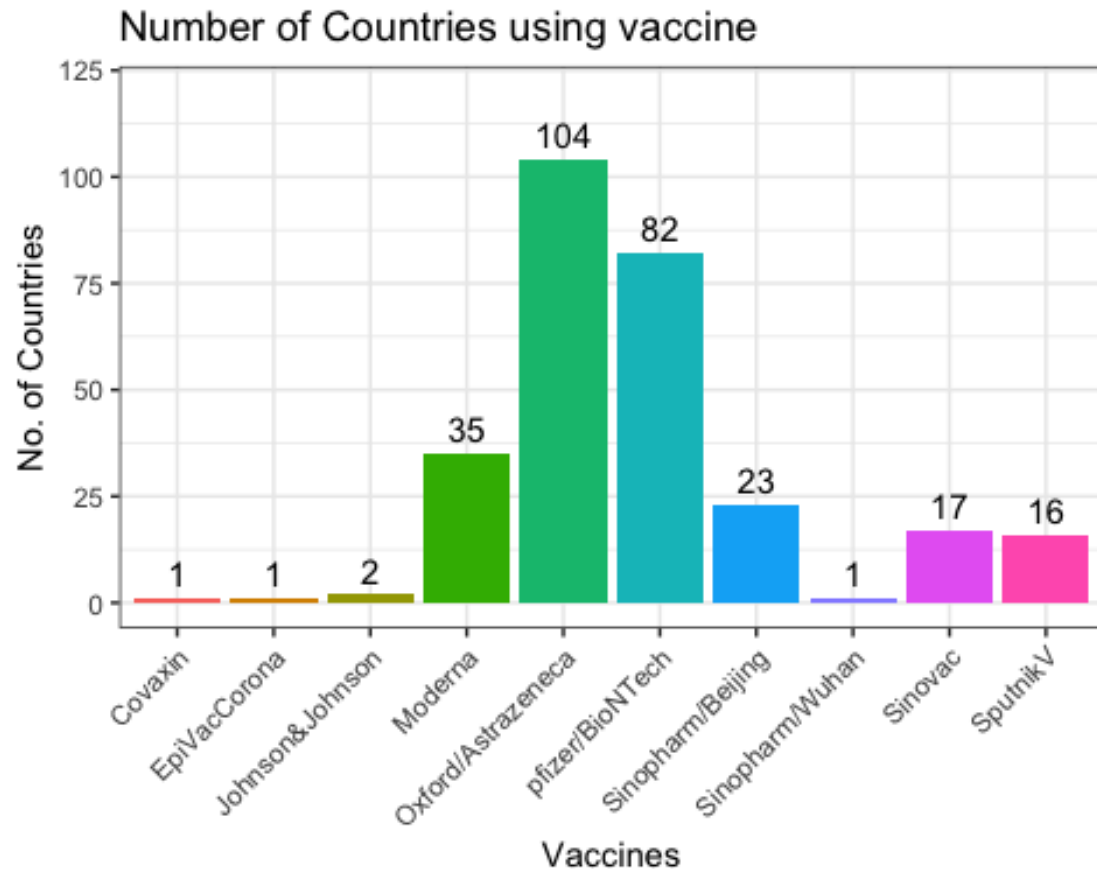
In this section, we will analyze each of our analysis question and present the results of our data analysis separately. We will be using data visualization, web-scraping techniques and linear regression depending on what analysis question we are answering.

Question 1 - Which vaccine is mostly used in countries around the world?

For answering this question, firstly we need to find how many vaccines are in there in this dataset. After exploring the dataset, we found that these are the 10 unique vaccines which are being used by countries all over the world.

```
## [1] "Oxford/AstraZeneca" "Pfizer/BioNTech" "SputnikV"  
## [4] "Moderna"           "Sinovac"          "Sinopharm/Beijing"  
## [7] "Covaxin"           "EpiVacCorona"     "Johnson&Johnson"  
## [10] "Sinopharm/Wuhan"
```

Now, we find the number of countries using the above vaccines and visualize our results in a bar plot.

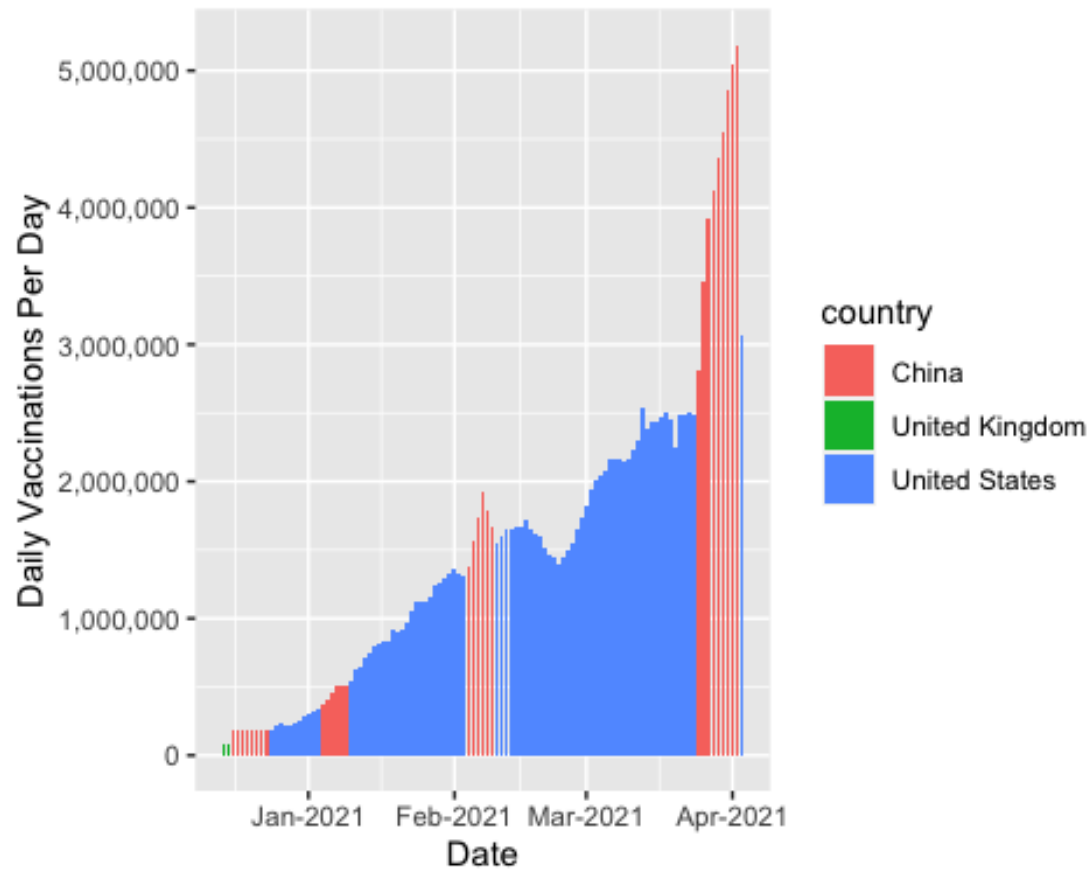


From the above analysis and data visualization, we can see that the most used vaccine around the world is Oxford/Astrazeneca followed by pfizer/BioNTech. A total of 99 countries using Oxford/Astrazeneca vaccine and a total of 81 countries are using the pfizer/BioNTech vaccine.

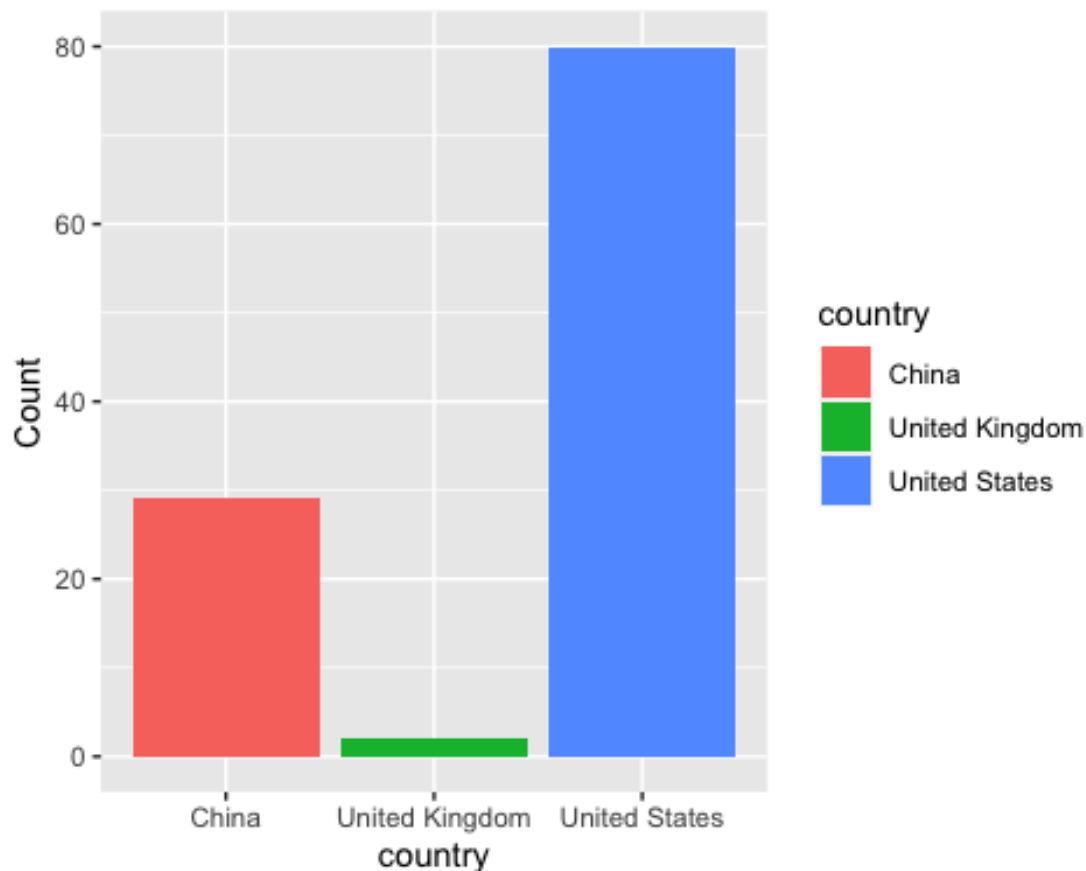
Question 2 - Which countries have the highest daily average vaccination rate?

For this analysis question, we need to find the mean of `daily_vaccinations_per_million` column grouped by country and find the top 10 countries with highest rate of daily average vaccination per million.

Now, we can visualize the above data in a bar-chart to show the countries with the highest number of vaccinated people per day.



```
## # A tibble: 3 × 2
## # Groups:   country [3]
##   country      n
##   <chr>      <int>
## 1 China        29
## 2 United Kingdom    2
## 3 United States   80
```



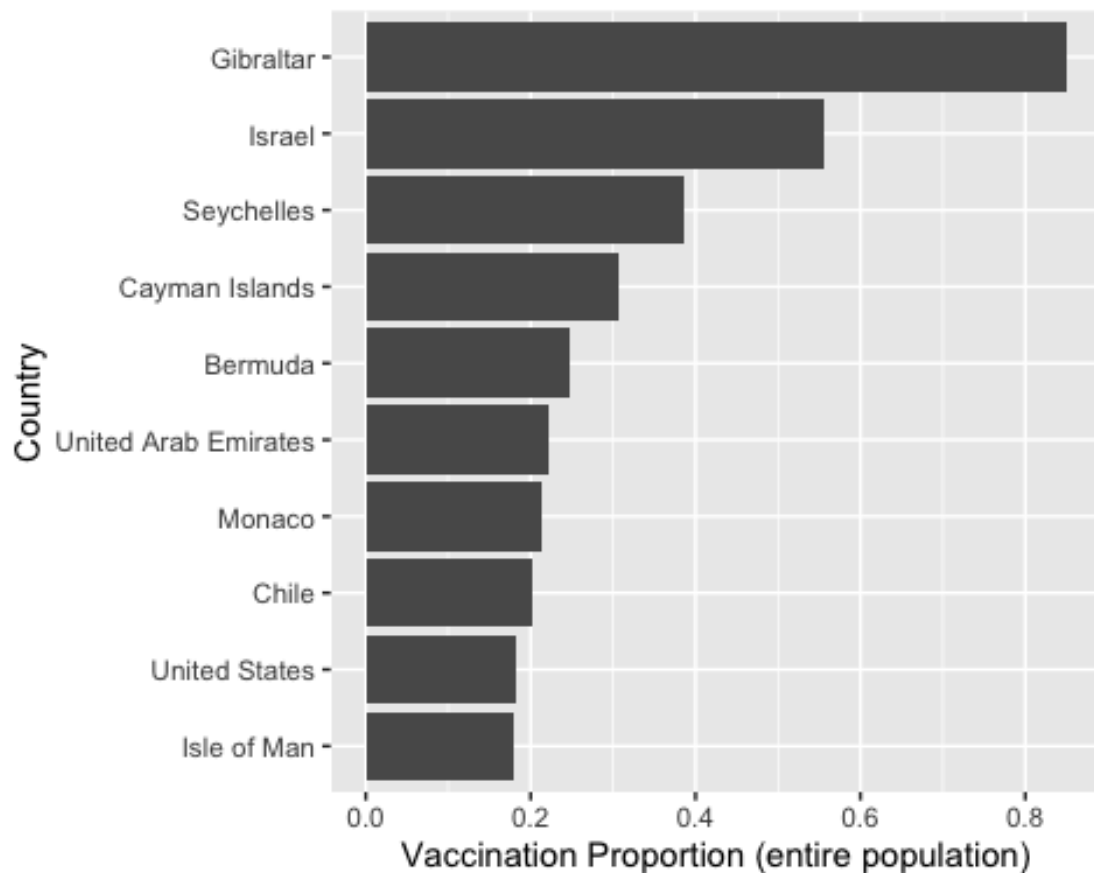
As we can see from the chart above, only 3 countries have the highest number of vaccinated people on 108 different dates. China has appeared 27 times which means that China has the highest number of vaccinated people on 27 different dates. We can also say that United Kingdom has appeared only 2 times which is the lowest. On the other hand, United States has the highest number of vaccinated people on 79 different dates which is the highest among the countries in our chosen data-set.

Question 4 - What is the proportion of fully vaccinated people from entire population of a country?

In this analysis question, we will be using web-scraping to get country population data from worldometer website as described above in data and methods section.

Firstly, we get the data of fully vaccinated people from our primary data source. Then, we use web-scraping to get data of country and it's population from worldometer website. Since we don't need any other info after web-scraping this table, we clean up our data to only keep information about country and population. Since population column has commas and character datatype, we remove the commas from the column and convert it to a numeric datatype.

Secondly, we inner-join our data from primary source and secondary source by country. Then, we can visualize top 10 countries with highest proportion of fully vaccinated people out of their entire population in the following way.



As we can see above from this data visualization, we have a list of top 10 countries of highest proportion of fully vaccinated people out of their entire population. Gibraltar is leading on top by having highest proportion of fully vaccinated people out of its entire population whereas United States is at the bottom of the list.

Question 5 - Is the association between country GDP and total vaccinations statistically significant?

In this analysis question, we will be using web-scraping to get country GDP data from worldometer website as described above in data and methods section.

Firstly, we get the data of total vaccinations of each country from our primary data source. Then, we use web-scraping to get data of country and its GDP from worldometer website. Since we don't need any other info after web-scraping this table, we clean up our data to only keep information about country and GDP. Since GDP column has commas and \$ signs, we remove them from the column and make it to a numeric datatype from character datatype for further analysis later.

Secondly, we inner-join our data from primary source and secondary source by country. Since the values of GDP and total vaccinations are very large, we also need to transform the variables GDP and total vaccination using a log transformation.

Now, we can fit a linear model and find the confidence interval to determine if the association between log_GDP and log_total_vaccinations statistically significant.

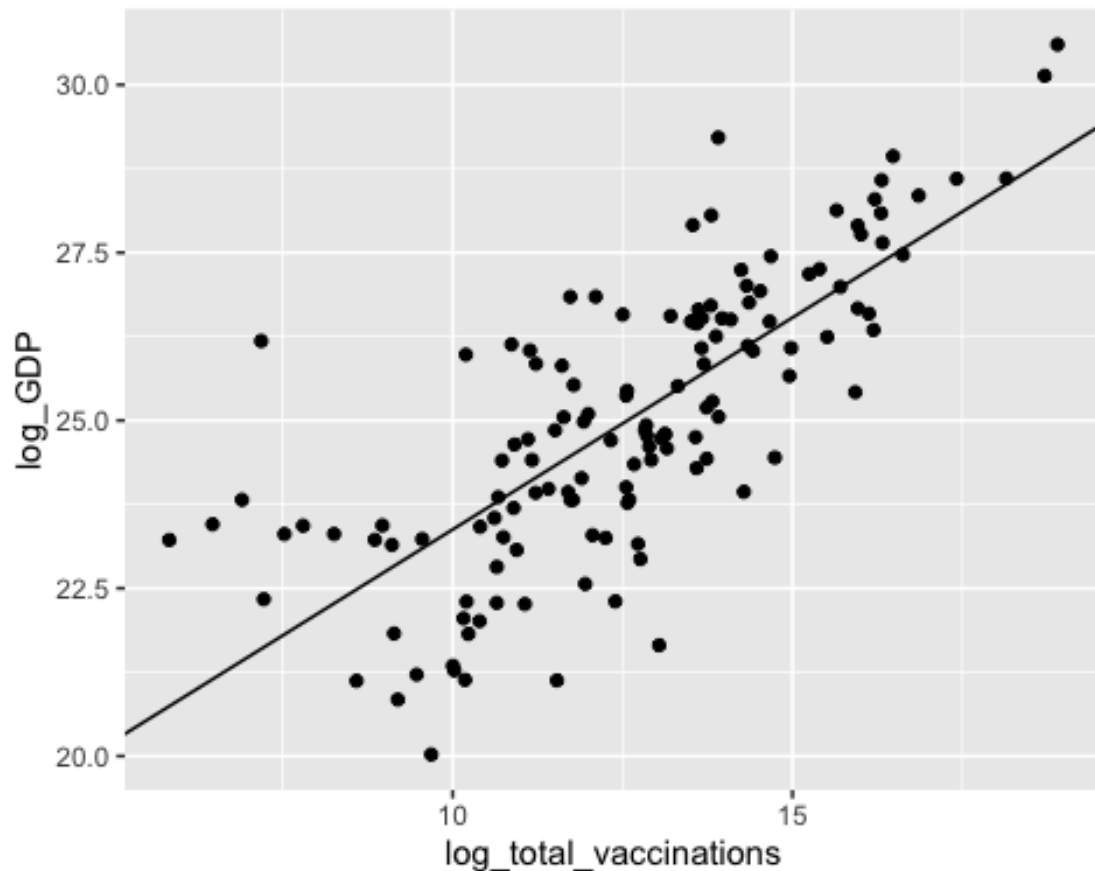
```
##
## Call:
## lm(formula = log_GDP ~ log_total_vaccinations, data =
vaccination_summary2_log)
##
## Coefficients:
##              (Intercept)  log_total_vaccinations
##              17.0640              0.6312

##              2.5 %      97.5 %
## (Intercept)      15.884726 18.2432248
## log_total_vaccinations 0.539404 0.7230014
```

To determine whether this difference is statistically significant, we can look at 95% confidence intervals for the regression coefficients.

Because 0, which corresponds to no difference in average value, is not in the confidence interval for the regression coefficient log_total_vaccinations, we conclude that the association between log_GDP and log_total_vaccinations is statistically significant.

The following scatterplot shows the relationship between log_total_vaccinations and log_GDP and the linear line provides a good description of the dataset.



Conclusion

In summary, this data analysis project gave us the chance to explore and understand the overall progress of COVID-19 vaccinations all around the world. We learned about the most used vaccines in different countries. Additionally, we learned about the vaccination progress of different countries in different aspects. Lastly, we investigated to find association between country GDP and their total vaccinations count using linear regression model. In short, we are able to achieve the real purpose of our project which we originally intended in the beginning through extensive data analysis of all of our analysis questions.

Project Tweet (280 characters)

The most used vaccine around the world is Oxford/Astrazeneca currently. Bhutan has the highest rate of daily average vaccinations per million. People are getting vaccinated more daily in USA. The association between country GDP and total vaccinations is statistically significant.

Appendix

Analysis Question 1 Code

```
library(tidyverse)
covid<-read.csv("country_vaccinations.csv")

data<-covid%>%
  separate(col="vaccines",into=c("vaccine1","vaccine2","vaccine3"),sep=",")

data1<-unique(data$vaccine1)
data2<-unique(data$vaccine2)
data3<-unique(data$vaccine3)
data4<-c(data1,data2,data3)

#data4
data4<-data4%>%str_replace_all(" ","")
data5<-unique(data4)
data5<-data5[!is.na(data5)]
data5

## [1] "Oxford/AstraZeneca" "Pfizer/BioNTech" "SputnikV"
## [4] "Moderna" "Sinovac" "Sinopharm/Beijing"
## [7] "Covaxin" "EpiVacCorona" "Johnson&Johnson"
## [10] "Sinopharm/Wuhan"

library(tidyverse)
library(janitor)
library(stringr)

#Oxford/AstraZeneca
data1<-data%>%
  filter(vaccine1=="Oxford/AstraZeneca"|vaccine2==" Oxford/AstraZeneca"
         |vaccine3==" Oxford/AstraZeneca")
data1<-unique(data1$country)
data1<-as.data.frame(data1)

#Pfizer/BioNTech
data2<-data%>%
  filter(vaccine1=="Pfizer/BioNTech"|vaccine2==" Pfizer/BioNTech"
         |vaccine3==" Pfizer/BioNTech")
data2<-unique(data2$country)
data2<-as.data.frame(data2)
names(data2)[names(data2) == "data2"] <- "Pfizer/BioNTech"

#Sputnik
sputnik<-data%>%
  filter(vaccine1=="Sputnik V"|vaccine2==" Sputnik V"|vaccine3==" Sputnik V")
sputnik<-unique(sputnik$country)
sputnik<-as.data.frame(sputnik)
```

```

#Moderna
moderna<-data%>%
  filter(vaccine1=="Moderna"|vaccine2==" Moderna"|vaccine3==" Moderna")
moderna<-unique(moderna$country)
moderna<-as.data.frame(moderna)

#Sinovac
sinovac<-data%>%
  filter(vaccine1=="Sinovac"|vaccine2==" Sinovac"|vaccine3==" Sinovac")
sinovac<-unique(sinovac$country)
sinovac<-as.data.frame(sinovac)

#Sinobei
sinoBei<-data%>%
  filter(vaccine1=="Sinopharm/Beijing"|vaccine2==" Sinopharm/Beijing"
          |vaccine3==" Sinopharm/Beijing")
sinoBei<-unique(sinoBei$country)
sinoBei<-as.data.frame(sinoBei)

#Covaxin
covaxin<-data%>%
  filter(vaccine1=="Covaxin"|vaccine2==" Covaxin"|vaccine3==" Covaxin")
covaxin<-unique(covaxin$country)
covaxin<-as.data.frame(covaxin)

#EpivacCorona
epi<-data%>%
  filter(vaccine1=="EpiVacCorona"|vaccine2==" EpiVacCorona"
          |vaccine3==" EpiVacCorona")
epi<-unique(epi$country)
epi<-as.data.frame(epi)

#Johnson&Johnson
john<-data%>%
  filter(vaccine1=="Johnson&Johnson"|vaccine2==" Johnson&Johnson"
          |vaccine3==" Johnson&Johnson")
john<-unique(john$country)
john<-as.data.frame(john)

#Sinopharm/Wuhan
sinoWuhan<-data%>%
  filter(vaccine1=="Sinopharm/Wuhan"|vaccine2=="Sinopharm/Wuhan"
          |vaccine3=="Sinopharm/Wuhan")
sinoWuhan<-unique(sinoWuhan$country)
sinoWuhan<-as.data.frame(sinoWuhan)

#Vaccine data
vaccines<-c("Oxford/Astrazeneca","pfizer/BioNTech","SputnikV",

```

```

      "Moderna", "Sinovac", "Sinopharm/Beijing", "Covaxin",
      "EpiVacCorona", "Johnson&Johnson", "Sinopharm/Wuhan")
num_countries<-c(nrow(data1),nrow(data2),nrow(sputnik),nrow(moderna),
               nrow(sinovac),nrow(sinoBei),nrow(covaxin),nrow(epi),
               nrow(john),nrow(sinowuhan))
df<-data.frame(vaccines,num_countries)

#Data visualization
country_vaccine <- ggplot(mapping=aes(x=vaccines, y=num_countries,
                                     fill = vaccines))+
  geom_col() +
  labs(x = "Vaccines", y = "No. of Countries",
       title = "Number of Countries using vaccine")+
  geom_text(aes(label =num_countries ), vjust=-0.5)+
  theme_bw()+
  theme(axis.text.x = element_text(angle = 45, hjust = 1),
        legend.position = "None") +
  expand_limits(y = 120)

```

Analysis Question 2 Code

```

vac_rate<-covid%>%
  group_by(country)%>%
  summarise(mean=mean(daily_vaccinations_per_million,na.rm=TRUE))

top_ten<-vac_rate%>%
  top_n(n=10,wt=mean)

daily_avg_vaccinations <-top_ten%>%
  ggplot(aes(x = reorder(country, mean), y = mean)) +
  geom_bar(stat="identity") +
  coord_flip() +
  labs(x = "Country", y="Daily Average Vaccinations")

```

Analysis Question 3 Code

```

#Find how many dates are there in actual dataset for the daily_vaccination
column
library(tidyverse)

data <- read.csv("country_vaccinations.csv")

total_dates_with_data <- data %>%
  group_by(date) %>%
  filter(!is.na(daily_vaccinations)) %>%
  distinct(date) %>%
  nrow()

total_dates_with_data

## [1] 111

```

```
#Find the countries with highest number of people vaccinated per day.
total_vaccinations <- data %>%
  group_by(date) %>%
  filter(!is.na(daily_vaccinations)) %>%
  filter(daily_vaccinations == max(daily_vaccinations)) %>%
  summarise(country, daily_vaccinations) %>%
  distinct()

glimpse(total_vaccinations)

## Rows: 111
## Columns: 3
## $ date          <chr> "2020-12-14", "2020-12-15", "2020-12-16",
"2020-12-17", "2020-12-18", "2020-12-19", "2020-12-20", "2020-12-21", "2020-12-22", "2020-12-23", "2020-12-24", "2020-12-25", "2020-12-26", "2020-12-27", "2020-12-28", "2020-12-29", "2020-12-30", "2020-12-31", "2021-01-01", "2021-01-02", "2021-01-03", "2021-01-04", "2021-01-05", "2021-01-06", "2021-01-07", "2021-01-08", "2021-01-09", "2021-01-10", "2021-01-11", "2021-01-12", "2021-01-13", "2021-01-14", "2021-01-15", "2021-01-16", "2021-01-17", "2021-01-18", "2021-01-19", "2021-01-20", "2021-01-21", "2021-01-22", "2021-01-23", "2021-01-24", "2021-01-25", "2021-01-26", "2021-01-27", "2021-01-28", "2021-01-29", "2021-01-30", "2021-01-31", "2021-02-01", "2021-02-02", "2021-02-03", "2021-02-04", "2021-02-05", "2021-02-06", "2021-02-07", "2021-02-08", "2021-02-09", "2021-02-10", "2021-02-11", "2021-02-12", "2021-02-13", "2021-02-14", "2021-02-15", "2021-02-16", "2021-02-17", "2021-02-18", "2021-02-19", "2021-02-20", "2021-02-21", "2021-02-22", "2021-02-23", "2021-02-24", "2021-02-25", "2021-02-26", "2021-02-27", "2021-02-28", "2021-03-01", "2021-03-02", "2021-03-03", "2021-03-04", "2021-03-05", "2021-03-06", "2021-03-07", "2021-03-08", "2021-03-09", "2021-03-10", "2021-03-11", "2021-03-12", "2021-03-13", "2021-03-14", "2021-03-15", "2021-03-16", "2021-03-17", "2021-03-18", "2021-03-19", "2021-03-20", "2021-03-21", "2021-03-22", "2021-03-23", "2021-03-24", "2021-03-25", "2021-03-26", "2021-03-27", "2021-03-28", "2021-03-29", "2021-03-30", "2021-03-31", "2021-04-01", "2021-04-02", "2021-04-03", "2021-04-04", "2021-04-05", "2021-04-06", "2021-04-07", "2021-04-08", "2021-04-09", "2021-04-10", "2021-04-11", "2021-04-12", "2021-04-13", "2021-04-14", "2021-04-15", "2021-04-16", "2021-04-17", "2021-04-18", "2021-04-19", "2021-04-20", "2021-04-21", "2021-04-22", "2021-04-23", "2021-04-24", "2021-04-25", "2021-04-26", "2021-04-27", "2021-04-28", "2021-04-29", "2021-04-30", "2021-05-01", "2021-05-02", "2021-05-03", "2021-05-04", "2021-05-05", "2021-05-06", "2021-05-07", "2021-05-08", "2021-05-09", "2021-05-10", "2021-05-11", "2021-05-12", "2021-05-13", "2021-05-14", "2021-05-15", "2021-05-16", "2021-05-17", "2021-05-18", "2021-05-19", "2021-05-20", "2021-05-21", "2021-05-22", "2021-05-23", "2021-05-24", "2021-05-25", "2021-05-26", "2021-05-27", "2021-05-28", "2021-05-29", "2021-05-30", "2021-05-31", "2021-06-01", "2021-06-02", "2021-06-03", "2021-06-04", "2021-06-05", "2021-06-06", "2021-06-07", "2021-06-08", "2021-06-09", "2021-06-10", "2021-06-11", "2021-06-12", "2021-06-13", "2021-06-14", "2021-06-15", "2021-06-16", "2021-06-17", "2021-06-18", "2021-06-19", "2021-06-20", "2021-06-21", "2021-06-22", "2021-06-23", "2021-06-24", "2021-06-25", "2021-06-26", "2021-06-27", "2021-06-28", "2021-06-29", "2021-06-30", "2021-07-01", "2021-07-02", "2021-07-03", "2021-07-04", "2021-07-05", "2021-07-06", "2021-07-07", "2021-07-08", "2021-07-09", "2021-07-10", "2021-07-11", "2021-07-12", "2021-07-13", "2021-07-14", "2021-07-15", "2021-07-16", "2021-07-17", "2021-07-18", "2021-07-19", "2021-07-20", "2021-07-21", "2021-07-22", "2021-07-23", "2021-07-24", "2021-07-25", "2021-07-26", "2021-07-27", "2021-07-28", "2021-07-29", "2021-07-30", "2021-07-31", "2021-08-01", "2021-08-02", "2021-08-03", "2021-08-04", "2021-08-05", "2021-08-06", "2021-08-07", "2021-08-08", "2021-08-09", "2021-08-10", "2021-08-11", "2021-08-12", "2021-08-13", "2021-08-14", "2021-08-15", "2021-08-16", "2021-08-17", "2021-08-18", "2021-08-19", "2021-08-20", "2021-08-21", "2021-08-22", "2021-08-23", "2021-08-24", "2021-08-25", "2021-08-26", "2021-08-27", "2021-08-28", "2021-08-29", "2021-08-30", "2021-08-31", "2021-09-01", "2021-09-02", "2021-09-03", "2021-09-04", "2021-09-05", "2021-09-06", "2021-09-07", "2021-09-08", "2021-09-09", "2021-09-10", "2021-09-11", "2021-09-12", "2021-09-13", "2021-09-14", "2021-09-15", "2021-09-16", "2021-09-17", "2021-09-18", "2021-09-19", "2021-09-20", "2021-09-21", "2021-09-22", "2021-09-23", "2021-09-24", "2021-09-25", "2021-09-26", "2021-09-27", "2021-09-28", "2021-09-29", "2021-09-30", "2021-10-01", "2021-10-02", "2021-10-03", "2021-10-04", "2021-10-05", "2021-10-06", "2021-10-07", "2021-10-08", "2021-10-09", "2021-10-10", "2021-10-11", "2021-10-12", "2021-10-13", "2021-10-14", "2021-10-15", "2021-10-16", "2021-10-17", "2021-10-18", "2021-10-19", "2021-10-20", "2021-10-21", "2021-10-22", "2021-10-23", "2021-10-24", "2021-10-25", "2021-10-26", "2021-10-27", "2021-10-28", "2021-10-29", "2021-10-30", "2021-10-31", "2021-11-01", "2021-11-02", "2021-11-03", "2021-11-04", "2021-11-05", "2021-11-06", "2021-11-07", "2021-11-08", "2021-1
```

```
geom_col()+  
labs(y = "Count")
```

Analysis Question 4 Code

```
library(rvest)
library(tidyverse)
library(purrr)
library(tidytext)

#Get total vaccinations done in each country from vaccine dataset
#Dataset source: https://www.kaggle.com/gpreda/covid-world-vaccination-progress
vaccine_dataset <- read_csv("country_vaccinations.csv")

total_vaccinations <- vaccine_dataset %>%
  group_by(country) %>%
  filter(!is.na(people_fully_vaccinated)) %>%
  mutate(people_fully_vaccinated =
    max(people_fully_vaccinated, na.rm = TRUE)) %>%
  summarise(country, people_fully_vaccinated) %>%
  distinct()

#Web scraping to collect population of each country in 2020
url_population <-
  "https://www.worldometers.info/world-population/population-by-country"

resource <- read_html(url_population)

country_populations <- resource %>%
  html_node("table") %>%
  html_table() %>%
  rename(
    country = `Country (or dependency)`,
    population = `Population (2020)`
  ) %>%
  select(country, population)

#Clean data - remove commas from population column,
#then convert datatype from chr to numeric
country_populations$population <-
  as.numeric(gsub(",", "", country_populations$population))

#Inner-Join total vaccinations data and country populations by country
#proportion of fully vaccinated people out of
#entire population - data visualization of top 10 countries
vaccination_summary1 <- inner_join(total_vaccinations,
  country_populations,
  by = "country") %>%
  summarise(vaccination_proportion = people_fully_vaccinated/population) %>%
  top_n(n = 10) %>%
  ggplot(aes(x = reorder(country, vaccination_proportion),
    y = vaccination_proportion)) +
```

```
geom_bar(stat="identity") +
coord_flip() +
xlab("Country")
```

Analysis Question 5 Code

```
#Get total vaccinations in countries
vaccinations <- vaccine_dataset %>%
  group_by(country) %>%
  filter(!is.na(total_vaccinations)) %>%
  mutate(total_vaccinations = max(total_vaccinations, na.rm = TRUE)) %>%
  summarise(country, total_vaccinations) %>%
  distinct()

#Web scraping to collect GDP of each country
url_gdp <- "https://www.worldometers.info/gdp/gdp-by-country/"

resource_gdp <- read_html(url_gdp)

country_gdp <- resource_gdp %>%
  html_node("table") %>%
  html_table() %>%
  rename(
    country = `Country`,
    GDP = `GDP (nominal, 2017)`
  ) %>%
  select(country, GDP)

#data cleaning - remove comma and $ from GDP
country_gdp$GDP <- gsub(",", "", country_gdp$GDP)
country_gdp$GDP <- as.numeric(gsub("\\$", "", country_gdp$GDP))

#Inner-Join vaccinations data and country gdp by country
#Data analysis to find if there is any corelation
#between gdp and total vaccinations
vaccination_summary2 <- inner_join(vaccinations,
                                   country_gdp,
                                   by = "country")

vaccination_summary2_log <- mutate(vaccination_summary2,
                                   log_GDP = log(GDP),
                                   log_total_vaccinations =
                                     log(total_vaccinations))

#Fit a linear model and find confidence interval
fit <- lm(log_GDP ~ log_total_vaccinations, data = vaccination_summary2_log)

fit
```



```
##
## Call:
## lm(formula = log_GDP ~ log_total_vaccinations, data =
vaccination_summary2_log)
##
## Coefficients:
##              (Intercept)  log_total_vaccinations
##              17.0640              0.6312

confint(fit)

##              2.5 %      97.5 %
## (Intercept)    15.884726 18.2432248
## log_total_vaccinations 0.539404 0.7230014

#Scatterplot for showing relationship between log_total_vaccinations and
log_GDP
scatter_plot <- ggplot(vaccination_summary2_log, aes(x =
log_total_vaccinations,
              y = log_GDP)) +
  geom_point() +
  geom_abline(intercept = coef(fit)[1],
              slope = coef(fit)[2])
```