

Group Members:

Wesley Chok

Project 1

1. Data Wrangling

The current steps completed are: a, c, d, f and g. For step a, to read the csv data, we need to ensure to use the `read_csv` function. With this the following csv hyper link or file are placed into the parameters. To save the csv I have used the leftwards assignment operator to save it to a variable. For step b, I have removed the rows that contains an NA value, vaccination rates that are 0 and the confirmed shots that are 0. Touching on step c, the tables that I have tidied are the covid and demographic tables. For the covid data, the reason it is messy is because it contains column names as dates and the variables as vaccination shots. To tidy this, I have filtered out all the column names except for the dates, set the name of the column to "dates" and then set the values to "confirmed". Instead of having an abundance of columns, now dates are formatted to rows and the shots relating to that date are placed to the next column. For the demographics table, YR2015 contains multiple values from other columns. To tidy this, I have pivoted wider to have Series Code as the column name and the values as YR2015. Moving on to step d, to get the vaccination rate I have mutated the data. The variables I have used are the population variable and the confirmed variable that I have tidied. For step e, to get the number of days since the first non-zero vaccination number I had to use the `mutate()`. The value will be whatever the current row number is. In order to reset the days so it does not bleed into other countries, I used the `group_by()` function and grouped by `Combined_Key`. To reset the value, I've used the `cumsum()` function. Going to step f, to get only the number of hospital beds from the most recent year which was 2019 according to the data table. For step g I have used the mutate and select function to combine the male/female population. For each serial code group, I have added the "FE" and "MA"

together. Lastly for step h, I've merged all the tables through the `inner_join()` function. As suggested, I've compared all the country names and merged all the tables together.

2. Data Modeling

For the modeling combinations, I have compared the population with 5 of the demographic series codes. Since there were multiple series code in the demographic file, I thought it would be interesting to see the differences with population. SP.URB.TOTL, SP.POP.TOTL, SP.POP.80UP and SP.POP.1564.IN all had similar R-squared values in comparison to the population. The only outlier that existed was SP.DYN.LE00.IN which had a R-squared value of 0.06027.

3. Conclusion

Most of the vaccination rates seem to have a very close margin to 0 vaccinations throughout all the countries. These rates are most likely to be the most impactful data as it's all clumped together the most noticeable outlier out of all of the countries had a vaccination rate of approximately 0.24. The outlier is likely to be the least significant factor as it does not accurately represent the general vaccination rates per country,

```
1 library(tidyverse)
2
3 covid <- read_csv("https://raw.githubusercontent.com/govex/COVID-
19/master/data_tables/vaccine_data/global_data/time_series_covid19_vaccine_doses_admin_global.csv")
4 hospital <- read_csv("data.csv")
5 demo <- read_csv("demographics.csv")
6
7 # Data wrangling for vaccinations data. Created column for dates and
8 # confirmed vaccinations shots
9 covid %>%
10   pivot_longer(cols = -one_of('UID', 'iso2', 'iso3', 'code3', 'Admin2',
11                               'Province_State', 'Country_Region', 'Lat',
12                               'Long_', 'Combined_Key', 'Population'),
13               names_to='date', values_to= 'confirmed') -> covid
14 covid
15
16 # New variable to calculate the vaccination rate (confirmed/ Population)
17 vacRate <- covid %>% mutate(vacRate = confirmed / Population)
18
19 # Removed unnecessary data
20 hello <- vacRate %>% filter(!is.na(confirmed), !is.na(vacRate), vacRate > 0,
21                             confirmed > 0)
22
23 # New variable to calculate number of days since last non-zero vaccination shot
24 value <- hello %>%
25   group_by(Combined_Key,
26           group = cumsum(date = lag(date, default = "12/31/2021"))) %>%
27   mutate(daysSinceStart = row_number()) %>% ungroup %>% select(-group)
28
29 # Completed data wrangling vaccination data
30 value
31
32 demo_tidy <- demo %>% select(-`Series Name`) %>%
33   pivot_wider(names_from = "Series Code", values_from = YR2015)
34
35 demo_merged <- demo_tidy %>%
36   mutate(SP.POP.80UP=SP.POP.80UP.FE+SP.POP.80UP.MA) %>%
37   mutate(SP.POP.1564.IN=SP.POP.1564.MA.IN+SP.POP.1564.FE.IN) %>%
38   mutate(SP.POP.0014.IN=SP.POP.0014.MA.IN+SP.POP.0014.FE.IN) %>%
39   mutate(SP.DYN.AMRT=SP.DYN.AMRT.MA+SP.DYN.AMRT.FE) %>%
40   mutate(SP.POP.TOTL.IN=SP.POP.TOTL.FE.IN+SP.POP.TOTL.MA.IN) %>%
41   mutate(SP.POP.65UP.IN=SP.POP.65UP.FE.IN+SP.POP.65UP.MA.IN) %>%
42   select(-contains(".FE")) %>% select(-contains(".MA"))
43
44 names(demo_merged) <- gsub(" ", "_", names(demo_merged))
45
46 # Completed demographics wrangling vaccination data
47 demo_merged
48
49 #hospital <- filter(hospital, Year == 2019)
50
51 # Completed hospital wrangling vaccination data
52 hospital
53
54
55 # Changed name to South Korea
56 demo_merged <- demo_merged %>%
```

```
57 mutate(Country_Name = replace(Country_Name , Country_Name ==
58                                "Korea, Dem. People's Rep.", "South Korea"))
59 demo_merged <- demo_merged %>%
60   mutate(Country_Name = replace(Country_Name, Country_Name == "Korea, Rep.",
61                                "South Korea"))
62
63 # Changed name to Bahamas
64 demo_merged <- demo_merged %>%
65   mutate(Country_Name = replace(Country_Name, Country_Name == "Bahamas, The",
66                                "Bahamas"))
67
68 # Changed name to Iran
69 demo_merged <- demo_merged %>%
70   mutate(Country_Name = replace(Country_Name,
71                                Country_Name == "Iran, Islamic Rep.", "Iran"))
72
73 # Changed name to Hong Kong
74 demo_merged <- demo_merged %>%
75   mutate(Country_Name = replace(Country_Name,
76                                Country_Name == "Hong Kong SAR, China",
77                                "Hong Kong"))
78
79 # Merged demographics data into value
80 value_joined <- value %>%
81   inner_join(demo_merged, by=c(Country_Region = "Country_Name"))
82
83 # Merged hospital data into value
84 value_joined <- value_joined %>%
85   inner_join(hospital, by=c(Country_Region = "Country"))
86
87 # Completed merge of all 3 tables
88 value_joined
89
90 # Linear modeling data, Population vs SP.DYN.LE00.IN
91 modPopulationDyn <- lm(data=value_joined, formula=Population~SP.DYN.LE00.IN)
92 summary(modPopulationDyn)
93
94 # Linear modeling data, Population vs SP.URB.TOTL
95 modPopulationUrb <- lm(data=value_joined, formula=Population~SP.URB.TOTL)
96 summary(modPopulationUrb)
97
98 # Linear modeling data, Population vs SP.POP.TOTL
99 modPopulationUrbPop <- lm(data=value_joined, formula=Population~SP.POP.TOTL)
100 summary(modPopulationUrbPop)
101
102 # Linear modeling data, Population vs SP.POP.80UP
103 modPopulationUrb80UP <- lm(data=value_joined, formula=Population~SP.POP.80UP)
104 summary(modPopulationUrb80UP)
105
106 # Linear modeling data, Population vs SP.POP.1564.IN
107 modPopulationUrb80In <- lm(data=value_joined, formula=Population~SP.POP.1564.IN)
108 summary(modPopulationUrb80In)
109
110
111 # Calculates the most recent vaccination rate per country
112 valueTable <- value_joined %>% group_by(Country_Region) %>%
113   summarize(daysSinceStart = max(daysSinceStart), vacRate = first(vacRate))
```

```
114  
115 valueTable  
116  
117  
118 # Scatterplot  
119 ggplot(data=valueTable) + geom_point(mapping = aes(x=daysSinceStart, y=vacRate))  
120  
121 ggsave("scatterplot.pdf")  
122  
123 Model <- c(0.06027, 0.3993, 0.5123, 0.2982, 0.4891)  
124  
125 df <- data.frame(Model)  
126  
127 # Bar graph  
128 ggplot(data=df) + geom_bar(mapping = aes(x=Model))  
129  
130 ggsave("bar_graph.pdf")
```



