

Final Presentation

Arthi Vaidyanathan and Denise Shkurovich

11/16/2021

Introduction and Data:

We are interested in looking at the relationship between COVID-19 outcomes and nutrition worldwide. The USDA Center for Nutrition Policy and Promotion suggests a dietary intake which consists of 30% grains, 40% vegetable, 10% fruits, and 20% proteins (dietaryguidelines.gov). Previous studies demonstrate an increased mortality in patients infected with COVID-19 which have chronic inflammatory diseases such as obesity, diabetes, and hypertension.

The prevalence of these chronic inflammatory diseases are known to be correlated with an individual's diet (Onishi 2020). Furthermore, previous studies show that maintaining a healthy diet can decrease risk of severe infection by promoting the immune system (Messina et al. 2020, Iddir et. al 2020). Adequate protein consumption is essential for antibody production and poor nutrient consumption has been shown to increase inflammation and oxidative stress (Iddir et. al 2020). We are ultimately interested in seeing if countries that tend to consume similar diets to those suggested by the USDA show decreased rates of mortality from COVID-19 and how it/if it is related to income levels. We also explored each nutritional group in isolation as well as alcohol consumption, fat consumption, obesity, and undernourishment to analyze if there was a significant correlation with COVID recovery rates across countries.

This dataset, "COVID-19 Healthy Diet Dataset" comes from Kaggle. The dataset provides energy intake (kcal) as percentages of total diet by food group. In addition, it provides percentages of obesity and undernourished individuals. Finally it provides data for total confirmed COVID-19 cases, recovered COVID cases, COVID deaths, and active COVID cases for 170 countries. The food supply quantities in addition to the prevalence of obesity and undernourishment in the populations were obtained from the Food and Agricultural Organization of the United Nations, the population count was taken from the Population Reference Bureau, and the Johns Hopkins Center for Systems Science and Engineering was used for COVID-19 data and all data was last updated in February of 2021. The Food and Agricultural Organization of the United Nations as well as the Johns Hopkins Center for Systems Science and Engineering did not specify the methodology used to collect their data. In addition, data on the income group classification of each country was obtained from the World Bank.

Works Cited:

Iddir M., et al. Strengthening the Immune System and Reducing Inflammation and Oxidative Stress through Diet and Nutrition: Considerations during the COVID-19 Crisis. *Nutrients*. 2020;12(6):1562.

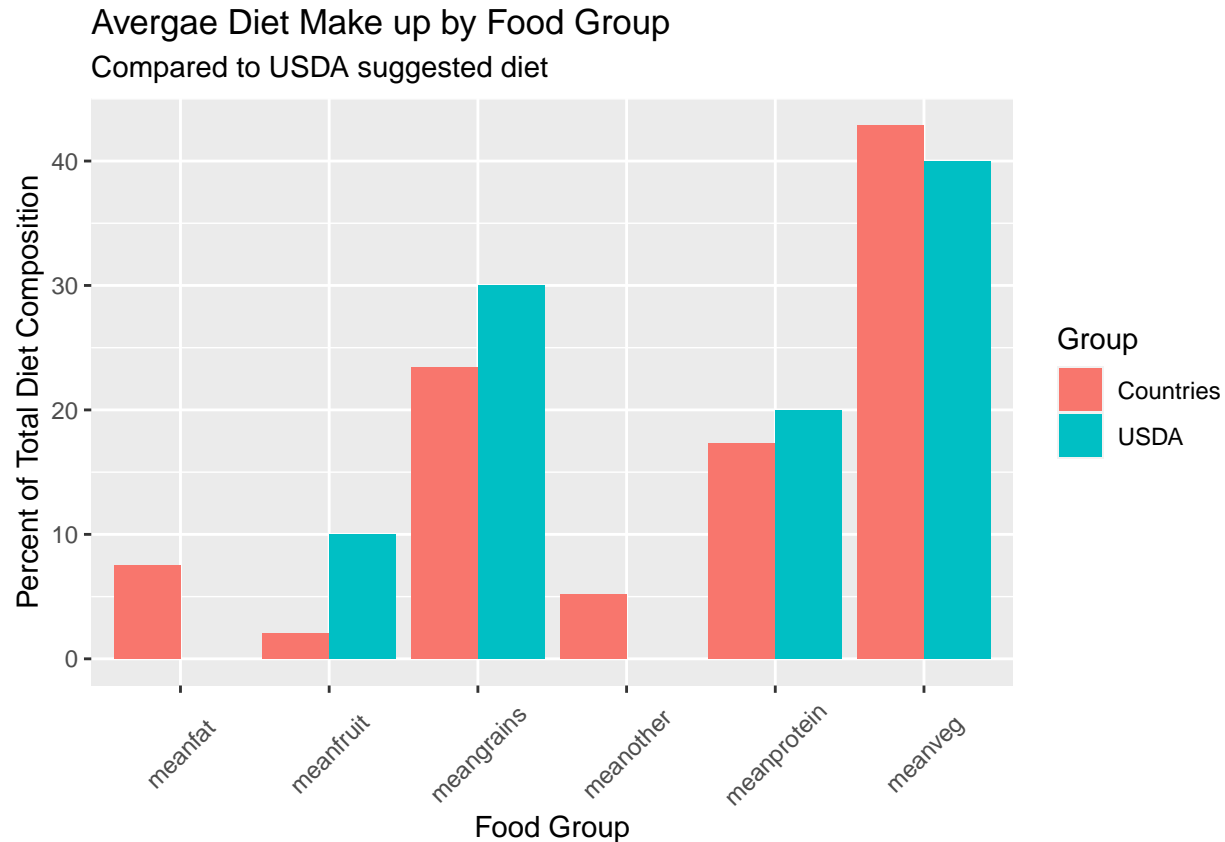
Messina G., et al. Functional Role of Dietary Intervention to Improve the Outcome of COVID-19: A Hypothesis of Work. *International Journal of Molecular Sciences*. 2020; 21(9):3104

Onishi J., et al. Can Dietary Fatty Acids Affect the COVID-19 Infection Outcome in Vulnerable Populations? *mBio*. 11(4).

Exploratory Data Analysis

Food Group Visualizations

Our dataset included more specific dietary composition than the USDA values. Therefore, to be able to compare them, we grouped the datasets into broader dietary categories including grains, vegetables, fruits, fats, proteins, and a general miscellaneous category. We did so by reading the descriptions of what composes each variable in the dataset.

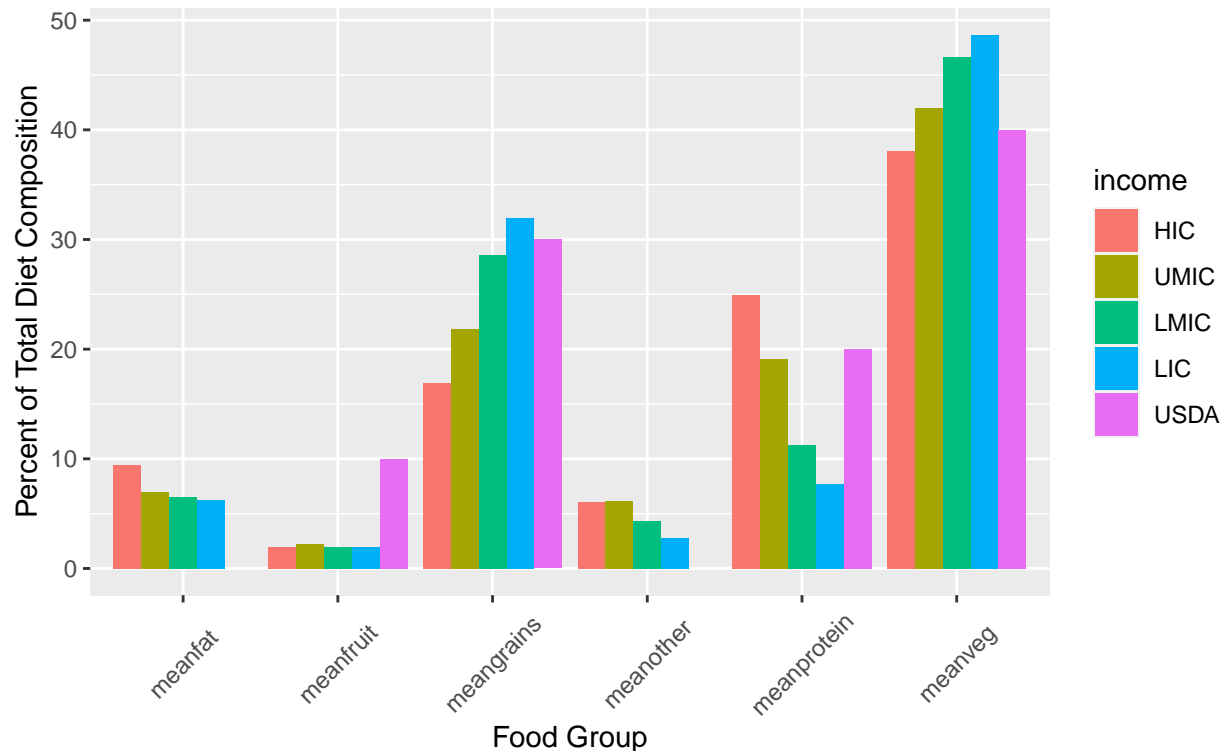


We were then curious to see if the mean percentages of the dietary categories of countries generally aligned with the USDA values. This helps inform us if we are making a valid comparison between the data and USDA values for the duration of the project. In order to do so we calculated the means of each of the values and added an additional column with the USDA values before plotting the bar graph. As shown, the values for grains, proteins, and vegetables are comparable to the USDA values indicating that our grouping of variables likely successfully captured the majority of dietary groups belonging to the broader categories. The value for fruits is significantly lower indicating that either the dataset does not accurately capture fruit consumption or countries on average eat significantly less fruit than recommended by the USDA. Furthermore, while USDA does not provide a value for suggested fat, it is a nutritional group significantly represented in our data set.

One of the major confounding factors to consider in both dietary consumption as well as COVID-19 outcomes is income. Since some foods such as grains are typically cheaper than others such as proteins, income could play a significant role as a confounding factor. In order to account for this discrepancy we created a new income variable and classified each country according to their World Bank income status. HIC refers to high income countries, UMIC refers to upper middle income countries, LMIC refers to lower middle income countries, and LIC refers to low income countries.

Diet Composition by Food Group

broken down by income category and compared to USDA suggested diet



In this graph, we summarized the break down of diet by income country in comparison to the recommended values from USDA. This visualization shows that low income countries seem to be consuming less than the recommended amount of protein and more grains than recommended. We can also see that higher income countries consume more protein and less grains than lower income countries. This graph suggests that in our analysis we should in fact include income level as a possible predictor of COVID outcomes since countries seem to show trends in their nutrition in most categories by income level.

Food Group Data Analysis

```
##           Df Sum Sq Mean Sq F value Pr(>F)
## income      3   5104  1701.4   86.92 <2e-16 ***
## Residuals 166   3249    19.6
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Using the overall F test, we identified a significant difference in percent of diet that comes from grains among income groups. Our p value is «0.001.

```
##           Df Sum Sq Mean Sq F value Pr(>F)
## income      3   2617   872.4   67.49 <2e-16 ***
## Residuals 166   2146    12.9
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Using the overall F test, we identified a significant difference in percent of diet that comes from vegetables among income groups. Our p value is «0.001.

```
##           Df Sum Sq Mean Sq F value Pr(>F)
## income      3     2.8   0.9269   0.458  0.712
```

```
## Residuals    166   336.1   2.0250
```

There are no significant differences in percent of diet that comes from fruit among income groups. Our p value is great than 0.05.

```
##              Df Sum Sq Mean Sq F value    Pr(>F)
## income         3     288    95.99    14.3 2.48e-08 ***
## Residuals     166    1114     6.71
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Using the overall F test, we identified a significant difference in percent of diet that comes from fats among income groups. Our p value is «0.001.

```
##              Df Sum Sq Mean Sq F value    Pr(>F)
## income         3    6911    2304   67.69 <2e-16 ***
## Residuals     166    5649     34
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Using the overall F test, we identified a significant difference in percent of diet that comes from proteins among income groups. Our p value is «0.001.

```
##              Df Sum Sq Mean Sq F value    Pr(>F)
## income         3    98.03    32.68   32.79 <2e-16 ***
## Residuals     166   165.44     1.00
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Using the overall F test, we identified a significant difference in percent of diet that comes from Alcohol/Stimulants among income groups. Our p value is «0.001.

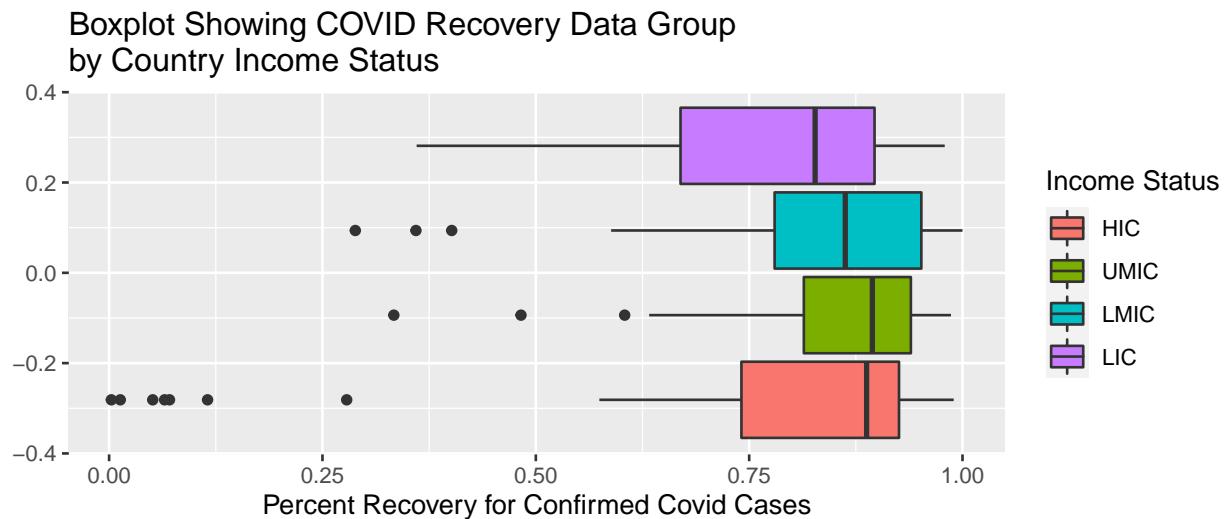
```
##              Df Sum Sq Mean Sq F value    Pr(>F)
## income         3    5888   1962.7   33.61 <2e-16 ***
## Residuals     163    9518     58.4
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 3 observations deleted due to missingness
```

Using the overall F test, we identified a significant difference in percent of the population that is obese among income groups. Our p value is «0.001.

```
##              Df Sum Sq Mean Sq F value    Pr(>F)
## income         3    9328   3109.2   36.36 <2e-16 ***
## Residuals     159   13595     85.5
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 7 observations deleted due to missingness
```

Using the overall F test, we identified a significant difference in percent of the population that is malnourished among income groups. Our p value is «0.001.

Recovery Rate Visualizations



In order to further our analysis, we created a variable that quantifies the COVID recovery rate for each country for comparison. We did this by dividing the percentage people in a country that recovered from COVID over the percentage of total people in that country that were confirmed COVID cases to obtain a rate of what percentage of people in that country that had COVID were able to recover. We also explored if income had an effect on percent recovery for COVID. High income countries hypothetically might have more resources to be able to provide better care for their citizens and improve their recovery rate. However, according to the boxplot, countries in the 4 income groups seem to have comparable recovery rates with the median for UMIC actually being the highest by a slight margin. Furthermore, the HIC group has numerous outliers near around 0.3 which might be important to consider. This seemed unusual to us due to the nature of disease and the resources the countries have. We speculate that more serious cases might be reported at disproportionately higher rates or there might potentially be other sources of error in the data set that we will later discuss.

Recovery Rate Analysis

```
##           Df Sum Sq Mean Sq F value Pr(>F)
## income      3  0.324  0.10785    2.546  0.058 .
## Residuals 156  6.607  0.04236
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

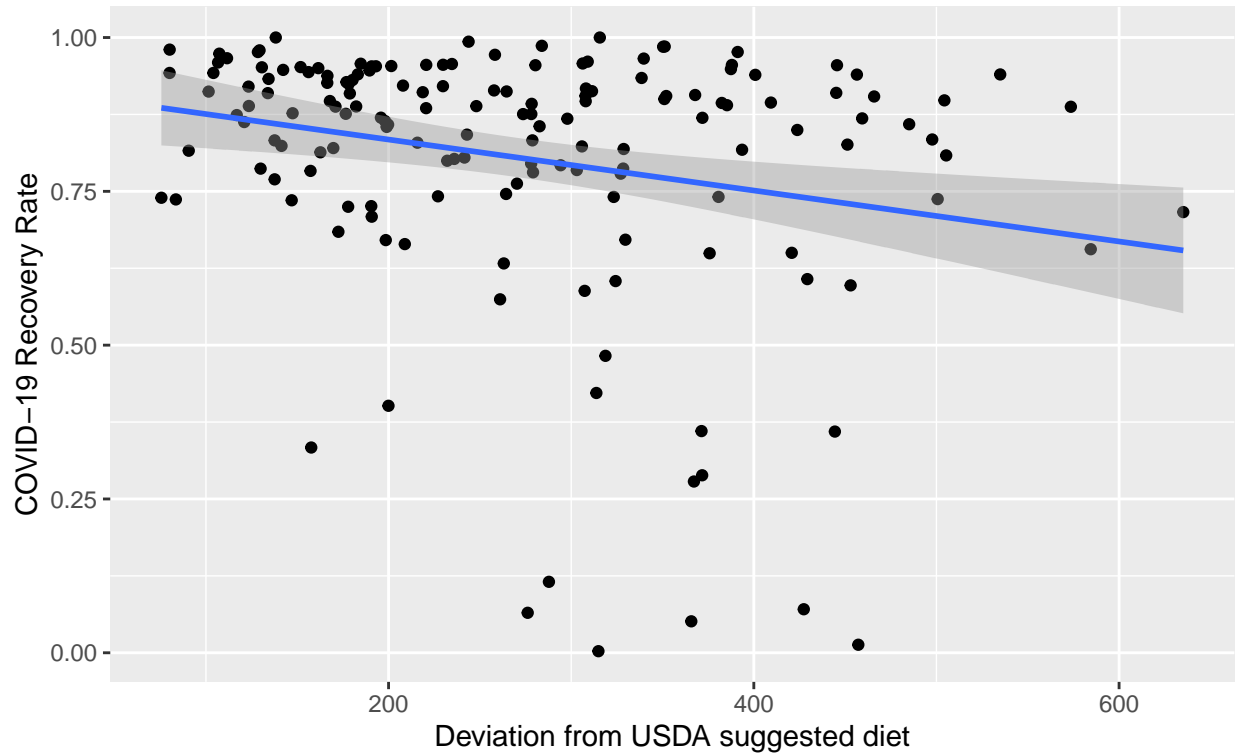
Based on the graphs, we wanted to see whether the mean recovery rate were significantly different between any of the income categories. ANOVA depends on variation between groups/within groups. From our graphs we can see that there is a decent amount of variation in recovery rate within each income category. This is probably due to the fact that these categories are very broad and there is different medical care/amount of COVID cases in each country even within income groups. However, the results of the ANOVA are not significant so therefore we cannot reject the null hypothesis that there are no significant differences due to income in recovery rate from COVID-19.

Dietary Deviation Variable Creation

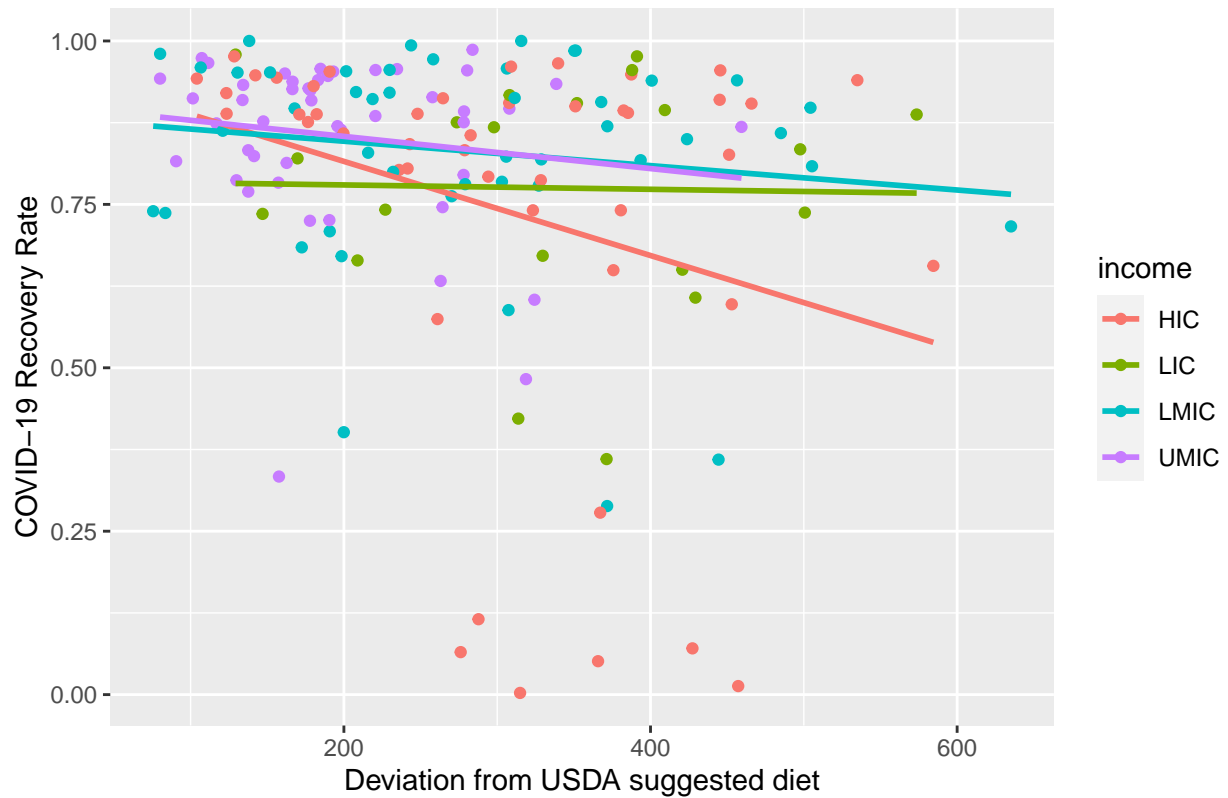
We created a variable for how much the average diet in each country differs from the suggested USDA value. We did so by subtract the mean value from each food group from the USDA value and then squaring that value. We then summed up this value from each food group to create deviance total variable. We squared them to account for negative values if the average diet was less than the USDA value for a certain category as well as increase the spread of deviance values for better comparison and visualization.

Methodology and Results

COVID-19 Recovery Rate in relationship to deviation from USDA suggested linear model by income category



COVID-19 Recovery Rate in relationship to deviation from USDA suggested



Here we compared the relationship between deviation from USDA suggested intake with recovery rates from COVID-19. The recovery rates were calculated by calculated the percent of the population that recovered from COVID-19 with the percent of the population that was diagnosed with COVID-19.

In the first graph, we see that there is a negative correlation between recovery rates and deviation from USDA suggested dietary intake. This suggests that as ones dietary intake begins to differ from the USDA suggested intake, they on average seem to have a lower change of recovering from COVID-19.

To further look at this relationship, we compared this relationship by income level. We can see that while lower middle income and upper middle income countries have a slightly negative relationship, lower income countries almost show no relationship between deviation rate and recovery rate. In contrast, high income countries show the strongest correlation between deviation rate and recovery rate.

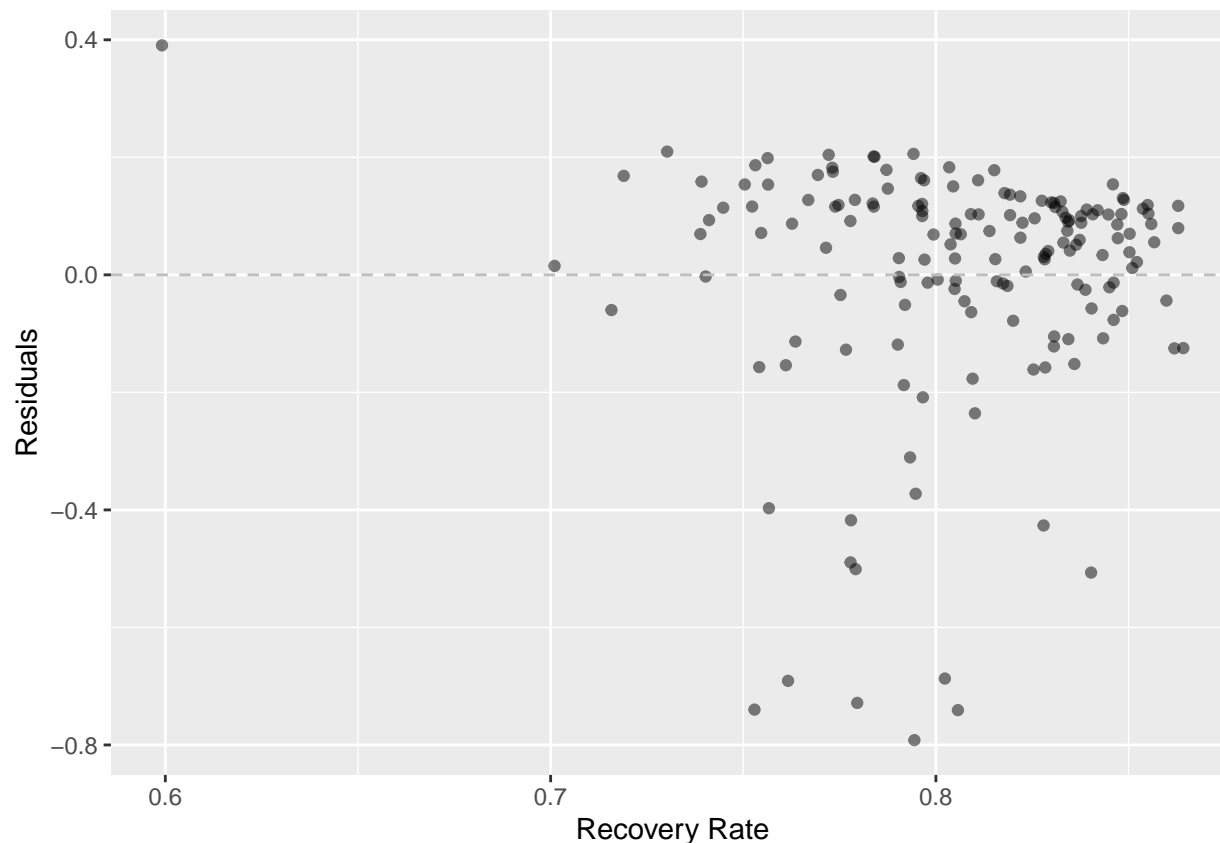
WHY?

Linear model of correlation between dietary deviance and covid recovery rate

```
## # A tibble: 2 x 5
##   term          estimate std.error statistic  p.value
##   <chr>          <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)    0.886      0.0377     23.5 1.83e-53
## 2 deviance_tot -0.000292  0.000123     -2.36 1.94e- 2
## [1] 0.02800442
```

This linear model shows that there is a significant relationship between recovery rate and deviation rate since our p value is less than the alpha level of 0.05. For every integer increase in deviance rate, the chances of recovering are on average 0.0002 lower.

However, the adjusted r squared value is quite low. It indicates that our model only accounts for 2.8% of the variance in the data.



In this model, the residuals do not seem to be even distributed around 0. The graph of the model along the adjusted r square value, suggest that we are perhaps some other predictors of the outcome that we are not including in our model.

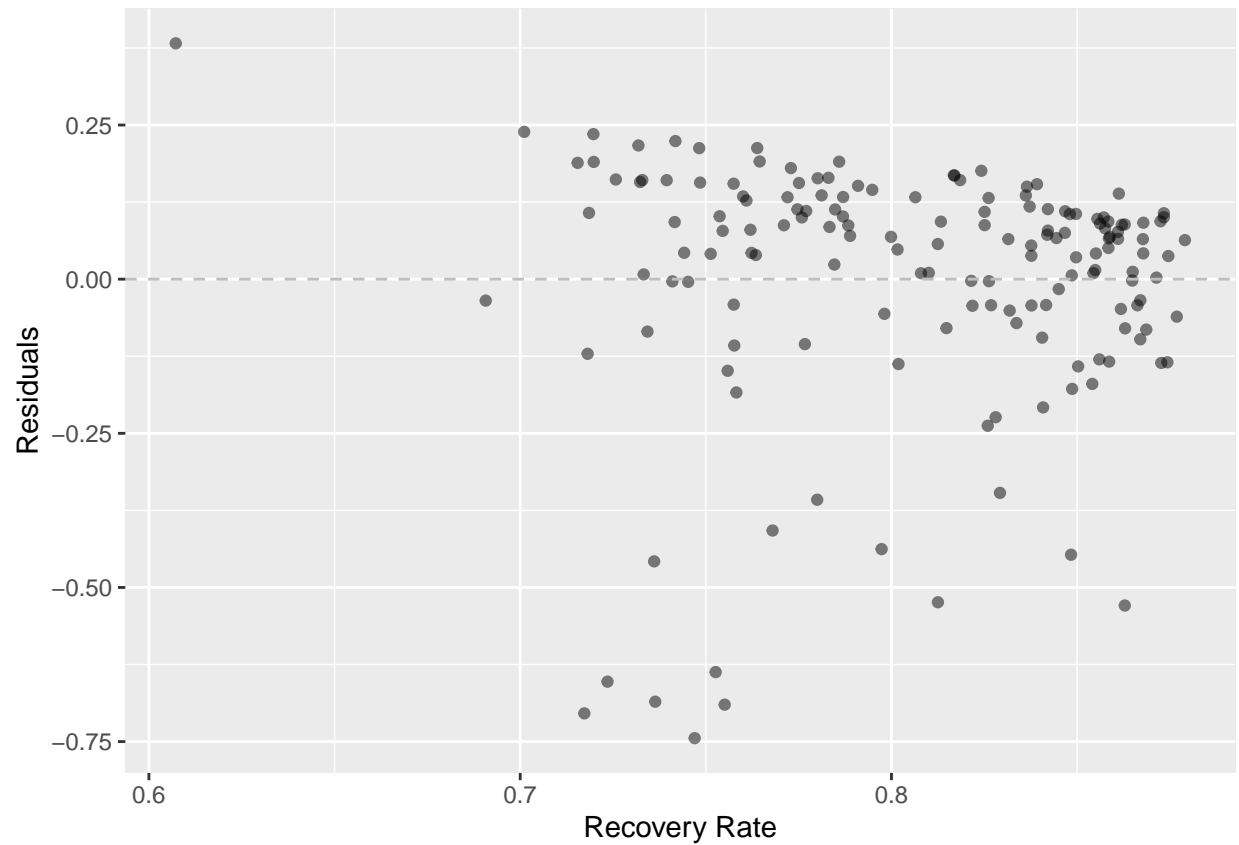
Linear model of correlation between dietary deviance and covid recovery rate accounting for income

```
## # A tibble: 5 x 5
##   term          estimate std.error statistic  p.value
##   <chr>         <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)   0.813      0.0519     15.6 1.05e-33
## 2 deviance_tot -0.000209  0.000133    -1.56 1.20e- 1
## 3 incomeLIC     0.0328     0.0547     0.599 5.50e- 1
## 4 incomeLMIC    0.0774     0.0426     1.82 7.10e- 2
## 5 incomeUMIC    0.0830     0.0453     1.83 6.89e- 2

## [1] 0.03726604
```

In this model, we cannot reject the null hypothesis that there is no relationship between covid recovery rate and dietary deviation from USDA suggested since our p value is greater than 0.05. However, compared to high income countries, lower middle and upper middle income countries on average have a .077 and .083 higher recovery rate holding deviance total constant.

Furthermore, this model has a larger adjusted r squared value, (3.7% compared to 2.8%) indicated that it accounts for more variance in the data than the first model that did not include income.



This graph also indicates that the model that includes income as well as dietary deviance is a better predictor of covid outcomes than a model that solely relies on dietary deviance. The residuals are more evenly spread along 0.

SUMMARY

LIMITATIONS

CONCLUSIONS