

MGT 6203 Team #64 Final Report

Introduction and Problem Statement

As a result of the COVID-19 pandemic, the global economy has faced many unique challenges over the past few years. With the current economic headwinds being faced by the world (supply chain issues, inflation, etc.), many experts are concerned that the direct consequence of the pandemic is a future recession. However, this is not the first recession the world has faced, and every recession presents a different set of challenges and learning opportunities. Governments will look at diverse ways to mitigate pending economic distresses. There are various factors that contribute to a recession that can affect recovery and can be critical to create policies or strategies to handle the current economic crisis.

As determined by (Kose, 2020), there have been 4 global recessions over the last 60 years (1974, 1981, 1990 and 2007), each lasting one year. The paper defines a global recession as a period of negative global economic growth as well as broad weakness in key indicators of global economic activity. It used a weighted average of GDP and purchase-power-parity around the globe to classify periods of global recession statistically and judicially. Our analysis will focus on these four periods. Each year following the global recession is defined as the recovery year where macroeconomic and financial activities rebound.

The goal of this analysis is to determine which factors from previous recessions will aid us in reducing the negative economic impact of a current global recession. By analyzing various economic, social and political impacts on various countries we can determine the predictors that have the greatest impact on a country, we can use these to determine how it would impact a recession. Directly, the response variables will be GDP Growth Rate (%), GDP and Stability.

We will build on the progress report data to find the best models and predictors to use in our conclusion. Based on the progress report, we can already eliminate many predictors that do not have enough data or will not be useful in prediction. We will use these observations to determine the results on GDP and stability.

Data Gathering

Based on the testing we did in the previous stages of this project, for this approach we chose to use the following predictors for Models 1-4.

- FIW (Democratic Indicator) – Aggregate Rating of political rights, civil liberties and freedom of expression in 13 categories. [2013-2022]
- Vaccine Rate - % of children aged 12-23 months [1960-2020]
- Life Expectancy – at birth [1960-2020]
- Population Growth Rate – Annual % [1960-2020]
- Internet Usage - % of population [1990-2020]
- Political Stability Estimate [1996-2020]
- CPIA Transparency, Accountability and Corruption Index – Rating of 1-6 with 1 being low and 6 being high [1995 - 2020]
- Trade – total merchandise imports annually in millions [1980-2019]
- Unemployment Rates – Total % of labor force unemployed [1991-2021]

- Inflation Rates – Annual Inflation Rate % [1960-2021]
- Interest Rates – Real Interest Rates % [1960-2021]

The remaining indicators that were previously identified were removed due to lack of data. For example, murder had almost all missing data points for most years. Migration rates had missing data for the years we were interested in.

Data Cleaning

Imputation and Country Selection:

For each predictor, the data was cleaned to ensure it was easier to use for analysis purposes.

The life expectancy dataset was used as the reference data set as it had the least missing data.

- The columns were renamed to remain consistent across all the predictors.
- Only data from the recession years 1973, 1974, 1975, 1980, 1981, 1982, 1989, 1990, 1991, 2007, 2008 and 2009 were kept.
- The countries that had less than 3 'NA' values across all the years were kept and the rest were removed. This left 198 countries.
- A new predictor called developed was created. Based on some online research, countries that were developed were given a value of 1 and 0 otherwise.

For the rest of the predictors:

- The columns were renamed
- Only data from the recession years 1973, 1974, 1975, 1980, 1981, 1982, 1989, 1990, 1991, 2007, 2008 and 2009 were kept.
- The dataset was merged with life expectancy so only those countries were kept in all the datasets.
- The NA values were imputed with average values from the average values of the developed or non-developed countries. For example, if Aruba was missing data, its value was taken to be the average of the un-developed countries in that same year.

The following predictors were also added to the datasets for analysis:

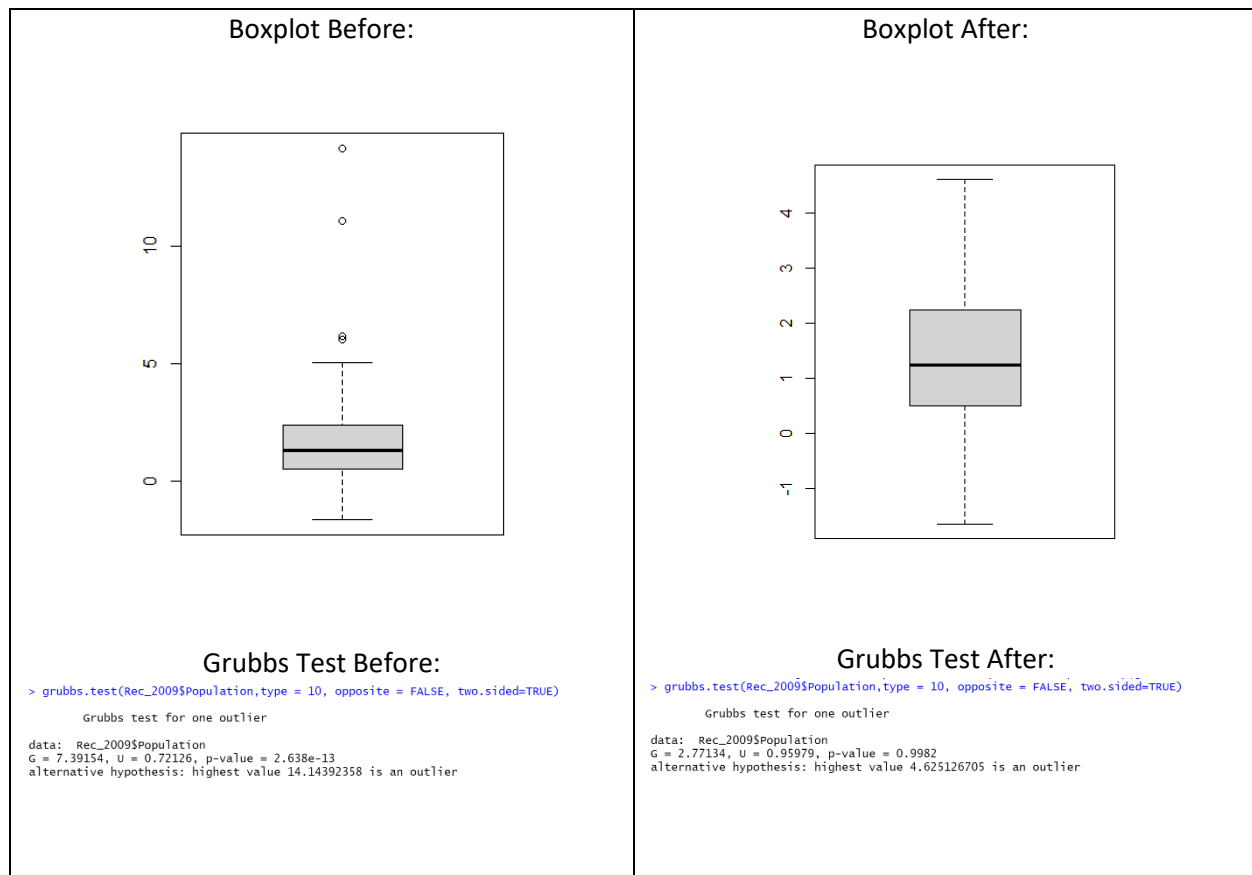
- A year factor was added to each of the data frames with 1, 2, 3 demonstrating before the recession, during the recession and after the recession
- A difference column was added to the recovery years in the year datasets. This contained the difference between the GDP, population, interest rates and inflation rates from before the recession and after. If the values increased, a 1 was assigned to the country indicating a positive sign from the recession, otherwise a 0 was assigned.

Outliers:

Outliers were removed. A boxplot was done on each of the datasets to visually inspect the data, then a Grubbs test was done to analytically determine the outliers. If the outlier was significant, it was removed. However, since the data had so much random variation due to imputation, not all the outliers could be removed without greatly decreasing the size of the dataset. For this reason, the Grubbs test

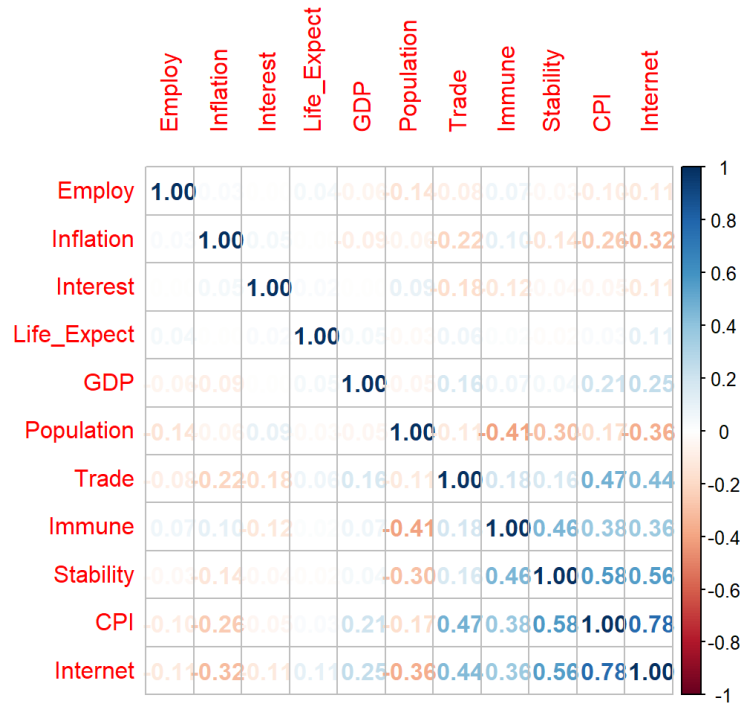
and boxplot were used in conjunction with a scatterplot to determine which datapoints were most impactful.

Below, you can see the outlier removal for the 2009 data.



Correlation:

After outliers were removed, correlation matrices were created for all the data to determine if any should be removed. Below is the correlation matrix for 2009. Only CPI and Internet are moderately correlated. Since they are not strongly correlated, all the predictors are kept in the model.



After the data was cleaned, it left the following data frames for analysis:

a. Year Data:

1. data_1973, data_1974, data_1975: life expectancy, GDP, population, trade, developed status
2. data_1980, data_1981, data_1982: life expectancy, GDP, population, trade, developed status, immunization rate
3. data_1989, data_1990, data_1991: life expectancy, GDP, population, trade, developed status, immunization rate, inflation rate, interest rate
4. data_2007, data_2008, data_2009: life expectancy, GDP, population, trade, developed status, immunization rate, inflation rate, interest rate, political stability, CPI index, internet usage, unemployment rate

b. Recession Data:

1. Rec_1: data from 1973-1975
2. Rec_2: data from 1980-1982
3. Rec_3: data from 1989-1991
4. Rec_4: data from 2007-2009

c. data_all: years 1973-2009 with common predictors - life expectancy, GDP, population, trade, developed status

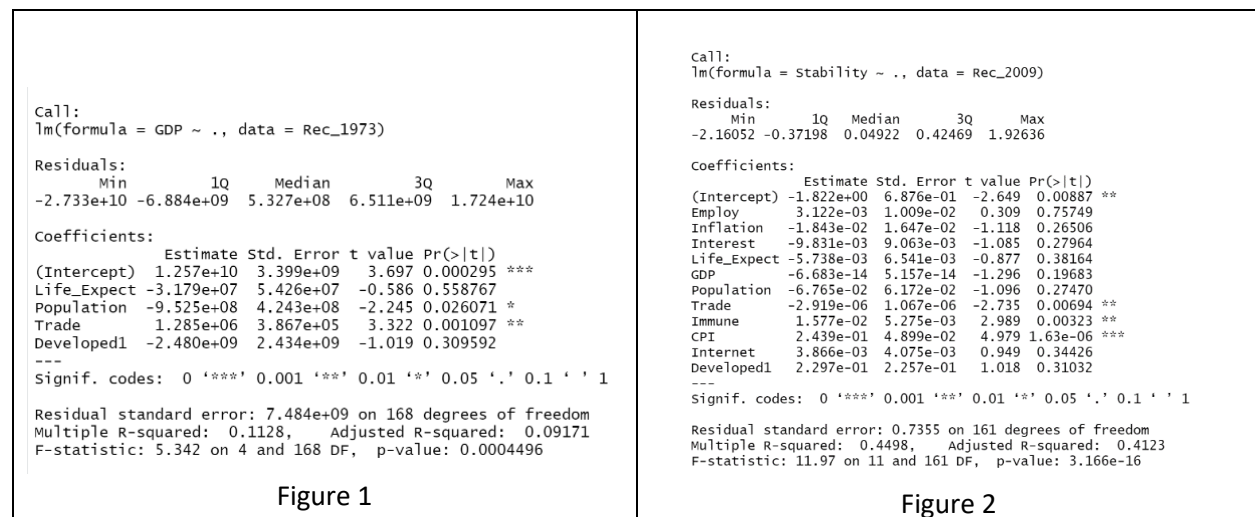
Modelling

We tried various models and modelling techniques to determine the best model to use for future predictions.

Model 1: Linear Regression with no Transformations

Figure 1: Using GDP as a response, this is an example from 1973. Population and trade were significant. Population negatively impacted GDP and Trade positively impacted GDP. However, the R2 value was very low so this model does not have predictive power. 1974, 1975, 1980, 1981 and 1982 showed similar results. 1989-2009 has no significant predictors so it is possible the addition of predictors decreased the predictive power of the model.

Figure 2: This is an example from 2009. Trade, vaccination rates and CPI were significant. Trade was negatively related to Stability. Vaccine and CPI were positively related. The R2 is 0.41. This is better than the model using GDP. 2007 and 2008 had similar results.

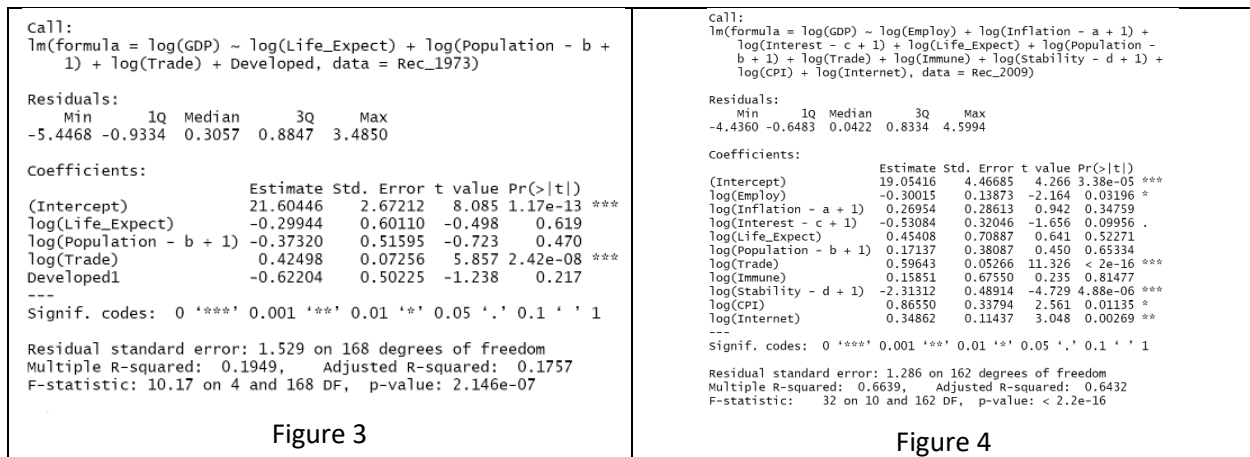


Based on the above, linear regression with no transformations is not a good model to be used on this data. It is likely because the relationships in this data are not linear.

Model 2: Linear Regression with Transformations

Figure 3: This is a sample from 1973. The minimum value of Population was added to the log to avoid 0 values. Only log(Trade) was significant and it was positively related to log(GDP). The R2 was only 0.1757. In 1974, 1980 and 1989 log(Population) was also significant. In 1990, log(interest) became significant.

Figure 4: This is a sample from 2009. log(Employ), log(Trade), log(Stability), log(CPI) and log(Internet) are significant. Log(stability) and log(employment) are negatively related to GDP. The R2 value is 0.6432. This is better than Model 1.



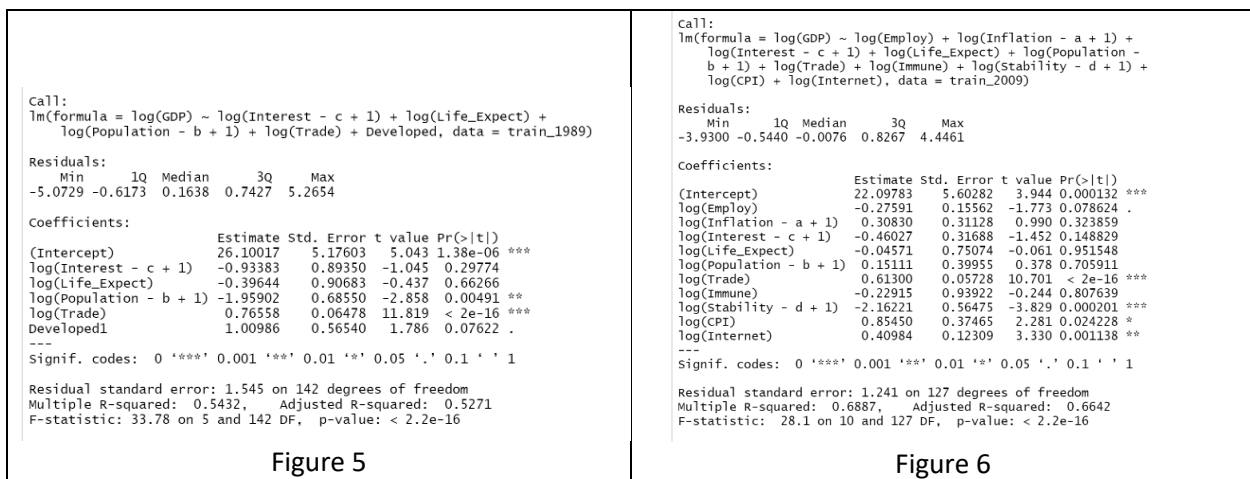
Based on the above, trade, stability, CPI, internet, and unemployment may be valuable predictors.

Model 3: Linear Regression with Validation and Transformations

Now that the log transformations have been determined to be more successful than without, we can try using validation to further improve the model.

Figure 5: The sample below is from 1989. log(Population) and log(Trade) are significant and the R2 is 0.5271 for 1989. Population is negatively related and trade is positively related. This is similar to results from model 1. Only trade was significant for 1973, 1975, 1980, 1981 and 1982. In 1990, developing country status was significant.

Figure 6: log(Trade), log(Stability), log(CPI) and log(Internet) are significant and the R2 is 0.6637. Trade, and internet are positively related. Population is negatively related.



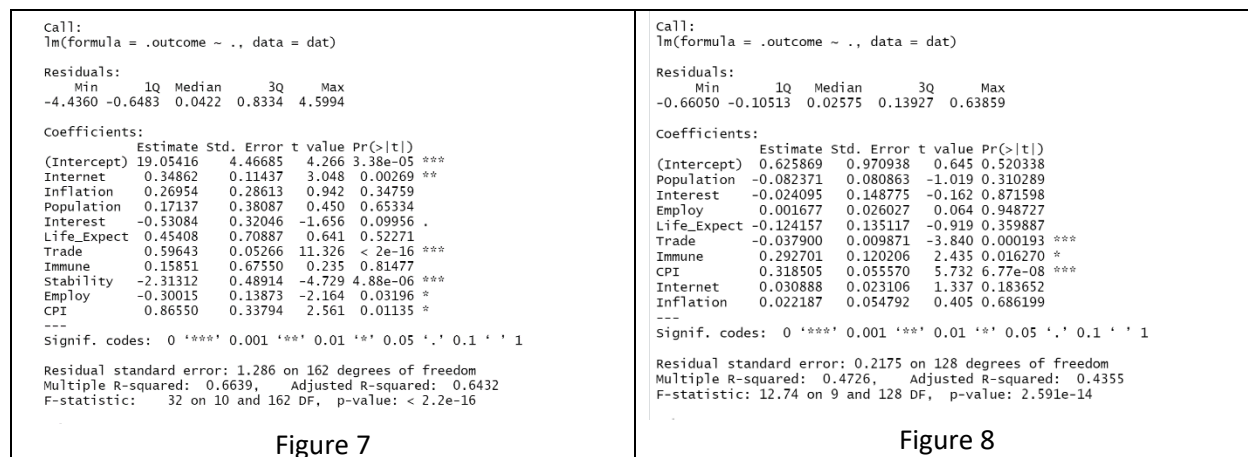
Based on the above trade and population may be valuable predictors.

Model 4: Linear Regression with Cross Validation and Transformations

Given the results from models 1-4, we decided only to use 2009 data for this model.

Figure7: For 2009, with GDP as response, log(Internet), log(Trade), log(Stability), log(CPI) and log(Employ) are significant and the R2 is 0.6432. Internet, trade and CPI are positively related while Stability and unemployment are negatively related.

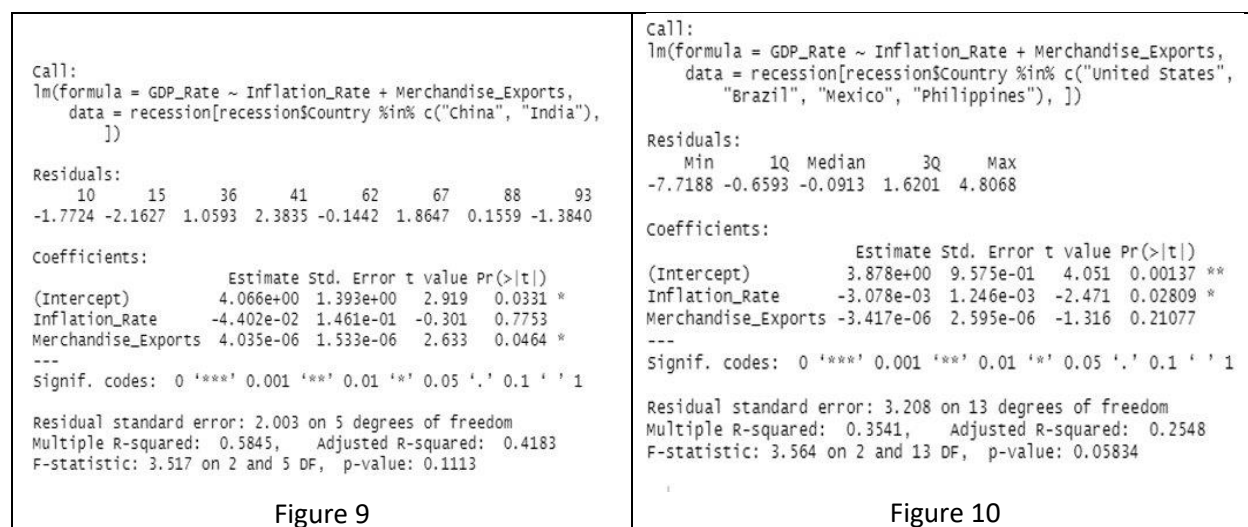
Figure 8: log(CPI), log(Trade), log(Immune) are significant and the R2 is 0.4355. Only trade is negatively related while the rest are positively related.



Based on the above trade, stability and CPI may be valuable predictors.

Model 5: Linear Regression with Trade Data

The goal of this model was to build upon the linear regression models that we created for our progress report where we were examining the relationship between each country's value of merchandise exports and their GDP growth rate during each recession time frame. However, this time we separated the countries into lists by population size to see if countries with larger populations would have similar or different trends to countries with smaller populations. Additionally, we ran these models twice: the first time against the data that we collected from the recession years and the second time against the data that we collected from the recession recovery years. Last, we added inflation rates as a second predictor variable. In summary, the output that we obtained from these models provided us with a better high-level understanding of trade on a country's performance during recessions and in the recovery phases, specifically by examining the effect of inflation rates and values of merchandise exports on the annual GDP growth rates of our countries by population size.



Looking at the output from the recession years-only data in Figure 9, for countries with populations over 1 billion people, inflation rates decreased and values of merchandise exports increased as GDP growth rates increased. We see in Figure 10 the same relationship for countries with populations between 100 million and 320 million people. During recovery years, our outputs indicate that the same trend occurred for countries with populations over 1 billion people (Figure 11). However, during recovery years for countries with populations between 100 million and 320 million people, we see that both inflation rates and the value of merchandise exports trended downward as GDP growth rates increased (Figure 12). We can see that the relationship between these variables appears to continue in the same direction as countries decrease in population size, as is noted in Figures 13 and 14 (output of the models for countries with less than 100 million people).

```
Call:
lm(formula = GDP_Rate ~ Inflation_Rate + Merchandise_Exports,
    data = recovery[recovery$Country %in% c("china", "India"),
    ])

Residuals:
    10     15     36     41     62     67     88     93 
-0.796239 -0.606944 -0.007954 -1.714945  3.371728 -1.486884 -0.381790  1.623028 

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)   8.981e+00  1.031e+00   8.714  0.00033 ***
Inflation_Rate -4.689e-01  1.426e-01  -3.288  0.02176 *
Merchandise_Exports 6.130e-07  1.913e-06   0.320  0.76156
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2.015 on 5 degrees of freedom
Multiple R-squared:  0.7122,    Adjusted R-squared:  0.597
F-statistic: 6.185 on 2 and 5 DF,  p-value: 0.04445
```

Figure 11

```
Call:
lm(formula = GDP_Rate ~ Inflation_Rate + Merchandise_Exports,
    data = recession[recession$Country %in% c("Vietnam", "Germany",
    "France", "United Kingdom", "Italy", "Spain", "Argentina",
    "Canada", "Morocco", "Peru", "Afghanistan", "Angola",
    "Australia", "Ecuador", "Cambodia", "Cuba", "Israel",
    "Norway", "Barbados", "Aruba"),
    ])

Residuals:
      Min       1Q   Median       3Q      Max
-8.6313 -2.3432  0.0776  2.0047 16.0656

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)   3.640e+00  5.911e-01   6.158  8.94e-08 ***
Inflation_Rate -1.533e-03  6.128e-04  -2.502  0.0153 *
Merchandise_Exports -3.830e-06  2.292e-06  -1.671  0.1004
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 3.963 on 55 degrees of freedom
(22 observations deleted due to missingness)
Multiple R-squared:  0.1327,    Adjusted R-squared:  0.1012
F-statistic: 4.209 on 2 and 55 DF,  p-value: 0.01992
```

Figure 13

```
Call:
lm(formula = GDP_Rate ~ Inflation_Rate + Merchandise_Exports,
    data = recovery[recovery$Country %in% c("United States",
    "Brazil", "Mexico", "Philippines"),
    ])

Residuals:
      Min       1Q   Median       3Q      Max
-5.9066 -1.7427 -0.1153  2.0410  3.7457

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)   2.053e+00  9.204e-01   2.230  0.0440 *
Inflation_Rate -2.302e-03  7.370e-03  -0.312  0.7597
Merchandise_Exports -6.198e-06  2.802e-06  -2.212  0.0455 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2.857 on 13 degrees of freedom
Multiple R-squared:  0.2741,    Adjusted R-squared:  0.1625
F-statistic: 2.455 on 2 and 13 DF,  p-value: 0.1246
```

Figure 12

```
Call:
lm(formula = GDP_Rate ~ Inflation_Rate + Merchandise_Exports,
    data = recession[recession$Country %in% c("Vietnam", "Germany",
    "France", "United Kingdom", "Italy", "Spain", "Argentina",
    "Canada", "Morocco", "Peru", "Afghanistan", "Angola",
    "Australia", "Ecuador", "Cambodia", "Cuba", "Israel",
    "Norway", "Barbados", "Aruba"),
    ])

Residuals:
      Min       1Q   Median       3Q      Max
-13.7988 -2.1762 -0.6974  2.1842 19.2615

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)   2.132e+00  8.100e-01   2.633  0.0109 *
Inflation_Rate -1.130e-04  1.029e-02  -0.011  0.9913
Merchandise_Exports -9.125e-06  3.796e-06  -2.404  0.0196 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

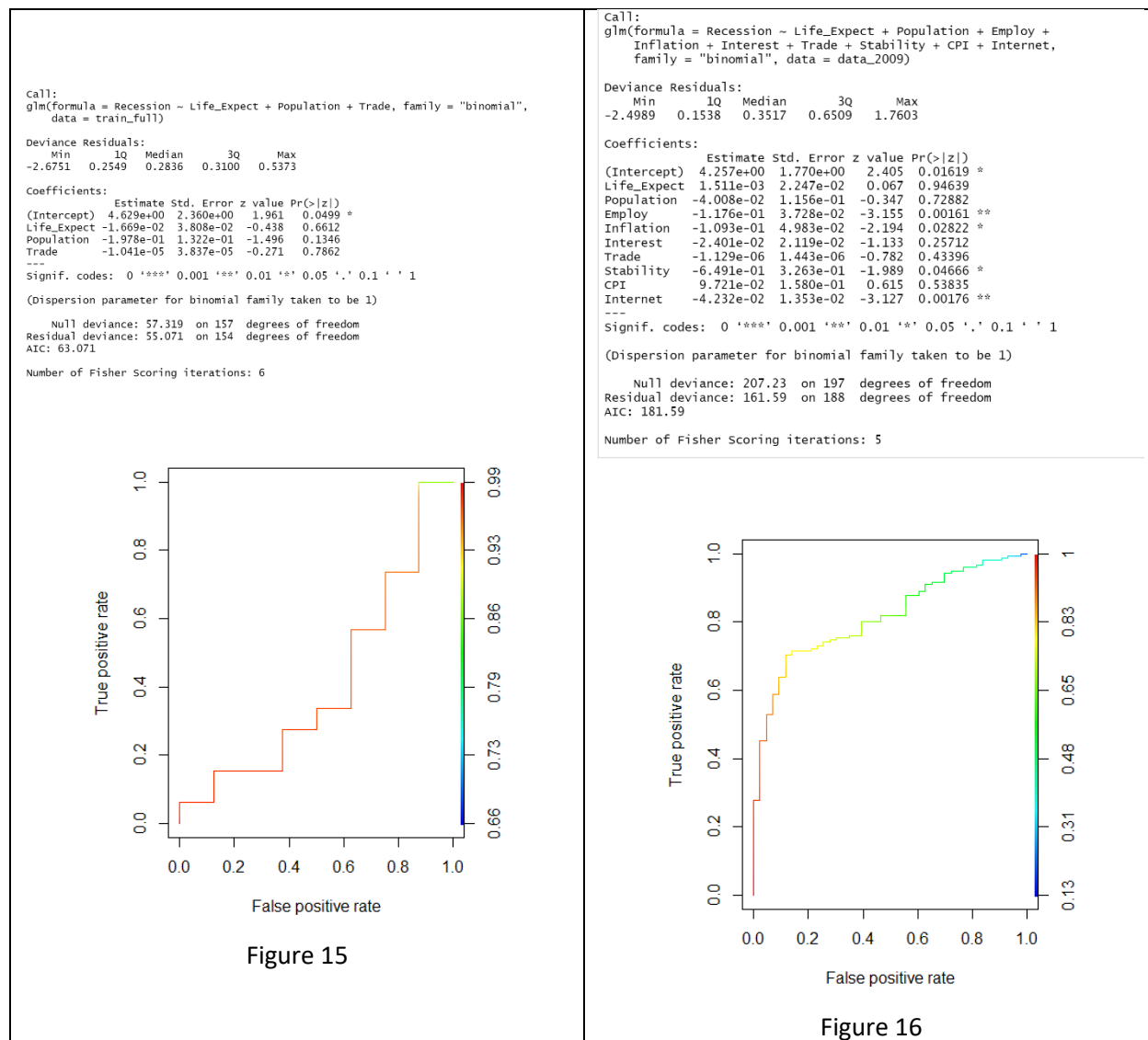
Residual standard error: 5.083 on 56 degrees of freedom
(21 observations deleted due to missingness)
Multiple R-squared:  0.09583,    Adjusted R-squared:  0.06354
F-statistic: 2.968 on 2 and 56 DF,  p-value: 0.05957
```

Figure 14

Model 6: Logistic Regression

Figure 15: Using the recession predictor as the response, there were no significant predictors in 1975 and the AUC was 0.3677. This means the model is not good at showing the variation in data.

Figure 16: Using the recession predictor as the response, unemployment, inflation, stability and the internet were related to recession success or failure. The AUC was 0.8144. This means the model is better than average and good at showing the variation in data. AIC is 181.593. This is lower than the linear regression model so it could be considered better.



In summary, from models 1-5, stability, trade, population and unemployment predictors could be useful for future prediction. Logistic regression seems to be the best kind of model to use for this data as it is not linear.

Model 7: Logistic Regression using aggregate data

For this model, we used Aggregate data from all the Recession Years (1973, 1974, 1980, 1981, 1989, 1990, 2007, 2008) and the Recovery Years (1975, 1982, 1991, 2009). We combined some of the economic and social variables. A binary variable was added to indicate years with recessions and years with no recessions. The variables used were:

Response Variable	Predicting Variables
1. Y = Recession (Binary Variable)	1. GDP Growth Rate
	2. Inflation Rate
	3. Population Growth
	4. Fuel Exports
	5. Life Expectancy
	6. Merchandise Exports

We used data from the 26 previously selected countries:

Afghanistan, Angola, Argentina, Aruba, Australia, Barbados, Brazil, Cambodia, Canada, China, Cuba, Ecuador, France, Germany, India, Israel, Italy, Mexico, Morocco, Norway, Peru, Philippines, Spain, United Kingdom, United States, Vietnam

70% of the data was used to train the model and 30% to test it. We used the MICE package to impute the missing values. Since all the variables were numeric, the imputation method used was Predictive Mean Matching with up to 40 iterations.

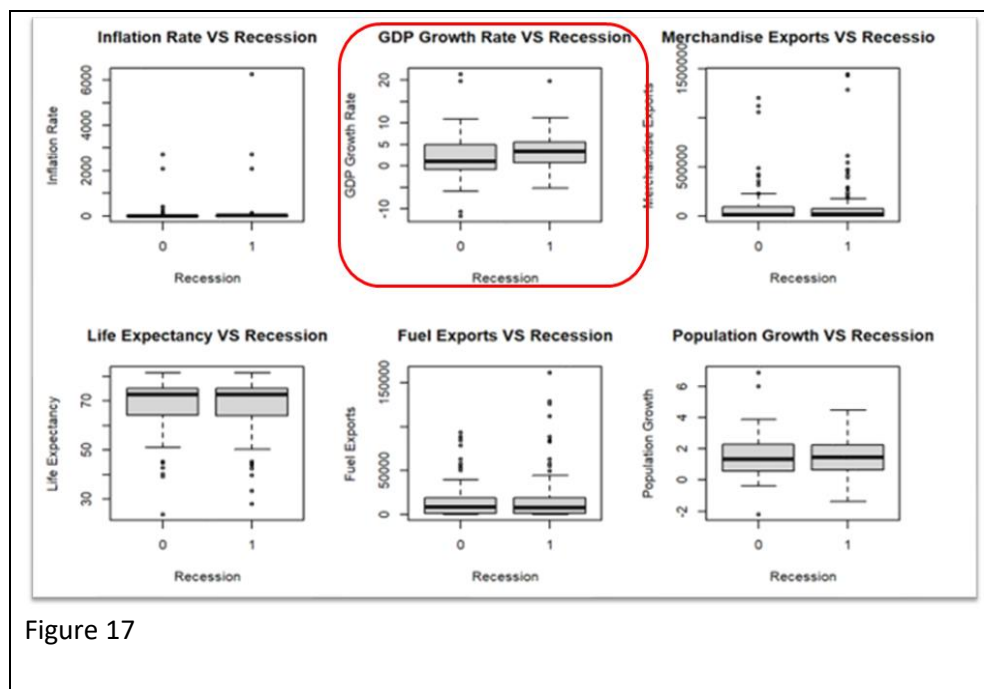
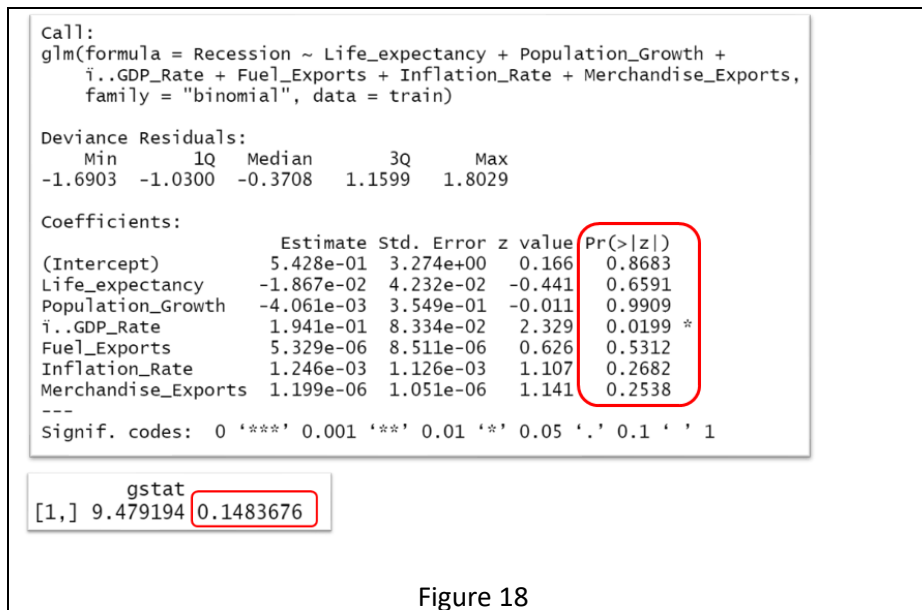


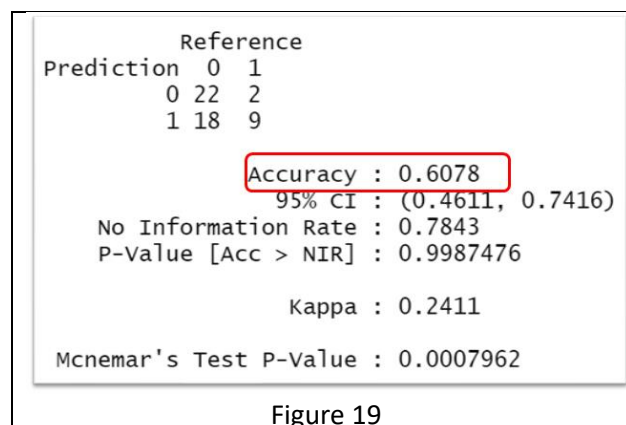
Figure 17 above shows the relationship between the predicting and the response variables. As we can see, there isn't a clear separation between Recession and Non-recession years for any of the variables

except for GDP Growth Rate. We can also see that most of the variables have some outliers which were removed using the Grubbs' Test.

After removing the outliers and fitting the logistic regression model, the output (Figure 18) points out that the only statistically significant coefficient is GDP Growth rate. This was expected after the box-plot analysis. The test for the overall regression shows that the p-value is large (0.14), therefore we can conclude that this model has no predictive power.



We can confirm this by running this model using the test data. In Figure 19 we can see the confusion matrix. The accuracy of this model is about 61% so the results weren't very good. The reason might be that for some of the variables over 25% of the values were missing and had to be imputed.



Model 8: Regression Tree

Usage of a tree model could allow recognition of patterns under different regimes. By dividing all the data into two datasets, Recession and Recovery, the goal of the tree model is to discover which thresholds of different predictors could act as differentiators between countries. Data is aggregated for all countries for all recession years and by predictors.

Prior to fitting a tree model, a straight linear regression model is fitted for both datasets. This model will provide a baseline accuracy to which we can compare the tree model. The outputs for each model are:

- Recession dataset: $R^2 = 0.28$
- Recovery dataset: $R^2 = 0.38$

```
Call:
lm(formula = GDP_Rate ~ ., data = recession)

Residuals:
    Min       1Q   Median       3Q      Max
-7.3928 -1.9122  0.0053  1.8519  6.4071

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  -8.74201    8.35130   -1.047  0.3009
Life_expectancy  0.09142    0.07218    1.267  0.2120
Murder        -0.14753    0.07791   -1.894  0.0649 .
Population_Growth 1.10370    0.57025    1.935  0.0594 .
Inflation_Rate  0.12988    0.05664    2.293  0.0267 *
Tot_Export      1.31128    3.36144    0.390  0.6983
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2.954 on 44 degrees of freedom
Multiple R-squared:  0.2795,    Adjusted R-squared:  0.1976
F-statistic: 3.413 on 5 and 44 DF, p-value: 0.01083
```

Figure 20

```
Call:
lm(formula = GDP_Rate ~ ., data = recovery)

Residuals:
    Min       1Q   Median       3Q      Max
-12.3418 -3.1461 -0.1542  2.6950 11.8931

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  14.5530819 10.2453830   1.420  0.162524
Life_expectancy -0.1385636  0.1275447   -1.086  0.283222
Murder        -0.3080922  0.1386062   -2.223  0.031416 *
Population_Growth 2.7384900  0.7291839    3.756  0.000504 ***
Inflation_Rate  0.0002523  0.0101944    0.025  0.980370
Tot_Export     -0.3984818  0.5412961   -0.736  0.465539
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 5.391 on 44 degrees of freedom
Multiple R-squared:  0.3819,    Adjusted R-squared:  0.3117
F-statistic: 5.438 on 5 and 44 DF, p-value: 0.0005581
```

Figure 21

The first tree model fitted uses the Recession years dataset and the tree is pruned. The tree model uses Inflation Rate and Murder Rate as differentiators. The R^2 for the Recession Regression Tree is approximately 0.41, a substantial jump compared to the baseline.

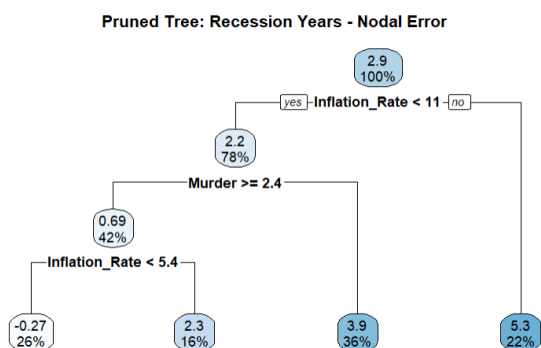


Figure 22

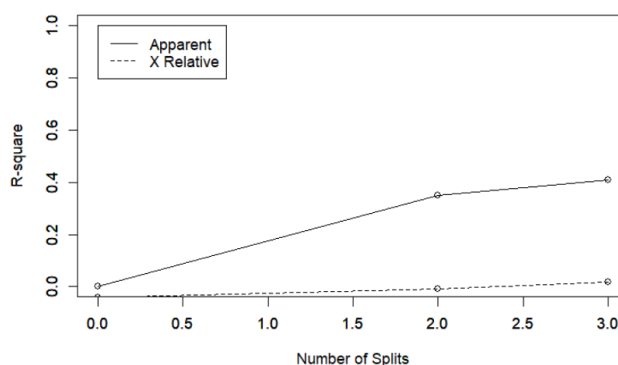


Figure 23

The second tree model was fitted for the Recovery years dataset. In contrast to the first model, this tree model uses Murder Rate and Population Growth Rate as differentiators. The R^2 for the Recovery Regression Tree is approximately 0.42, which is a slight improvement compared to the baseline.

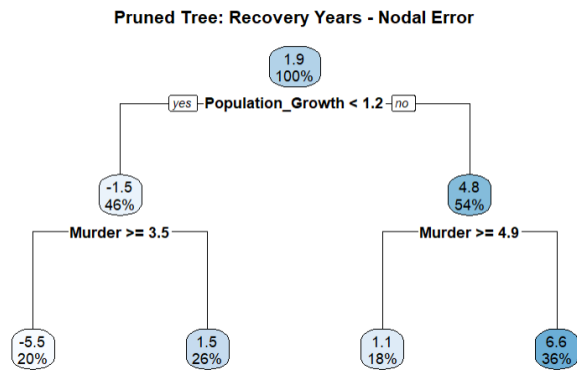


Figure 24

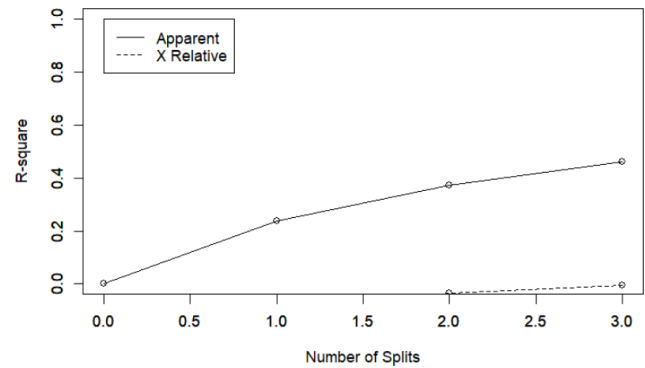


Figure 25

Comparing the trees to the baseline, it is interesting to note that Population Growth, Murder Rate, and Inflation Rate are deemed to be significant predictors that can determine the the response of a country to recession and its recovery after. While the R^2 value for the Recession dataset did indeed see an increas

Conclusion

The different models were definitively able to predict the likelihood of recovery from a recession or GDP Growth Rate with great degree of certainty. However, all models identified similar statistically significant predictors. Based on these models, the most useful predictors are:

- Trade
- Population Growth Rate
- GDP Growth Rate
- Inflation Rate
- Stability Rate

Some Important predictors that require more data are:

- Interest Rate
- Social Data (Crime, CPI Index)
- Unemployment rate

We would recommend that governments use these as a starting point to help limit the negative effects of a potential recession. However, there needs to be more work done on these datasets to get better results from them.

Future Improvements

Overall, the models did not perform as well as they could have. There could be a few reasons for this. The first is regarding the amount of missing data. There was a large quantity of missing data in various years for the recessions. The imputations will not be as accurate as the actual data. Also, there is fewer missing data for more recent years. If we had used more recent data, this problem may have been less prevalent.

The next reason is using more data analysis techniques. We might consider using different models and techniques such as PCA, and variable selection to better determine which predictors are valuable.

References

(n.d.). Retrieved from Our World in Data: <https://ourworldindata.org/>

Carlson, B. (1999, August). *Social Dimensions of Economic Development and Productivity: Inequality and Social Performance*. Retrieved from https://repositorio.cepal.org/bitstream/handle/11362/4659/S99117_en.pdf

Corruption Perceptions Index. (n.d.). Retrieved from Transparency International: <https://www.transparency.org/en/cpi/2020>

Fornari, F. (2010, October). *Predicting Recession Probabilities with Financial Variables Over Multiple Horizons*. Retrieved from European Central Bank: <https://www.ecb.europa.eu/pub/pdf/scpwps/ecbwp1255.pdf>

Freedom in the World. (n.d.). Retrieved from Freedom House: <https://freedomhouse.org/report/freedom-world>

Kose, A. M. (2020, March). *Global Recessions*. Retrieved from <https://documents1.worldbank.org/curated/en/185391583249079464/pdf/Global-Recessions.pdf>

UIS.Stat. (n.d.). Retrieved from UNESCO Institute for Statistics: <http://data.uis.unesco.org/>

World Bank Open Data. (n.d.). Retrieved from The World Bank: <https://data.worldbank.org/>