# Economic Forecasting: Predicting World Economic Health

## CMPT 732 – Big Data Programming – I

## Contents

# 1. Introduction

Economic Forecasting deals with trying to predict the future state of economies of countries around the world and the world as a whole using a formulated approach involving hundreds and hundreds of different parameters. These kinds of forecasts generally cater to economists who use these models to predict chief economic parameters like Growth Rate, Gross Domestic Product, Poverty Rates, Development, Human Development Index etc.

The problem with such forecasting is that they are generally meant for economists and are hard to understand in layman terms. We've attempted to resolve that by predicting several indicators using both conventional and unconventional parameters in a way that shall be understandable to most people.

We have obtained our data from myriad sources including but not limited to Web scraping, Big Query accesses, datasets from several international forums and economic websites like United Nations, IMF, World Bank etc. The data at 600MBs contains information of 193 countries and 10 regional groupings for around 100+ features from as far back as 1800s to 2017. The enormous amount of information has been scaled to our use and filtered extensively. We have faced extreme challenges of missing and incorrect data. After an extensive ETL process, the data was made suitable for our use.

We have divided the project into 3 modules:

1. **Global Economic Health Prediction**: Prediction of Growth rates of countries around the world and classification of these rates into 4 categories:
    a. Booming Economies (Growth Rate > 5%) (India, China, Philippines)
    b. Healthy Economies (Growth Rate > 3%) (Ireland, Egypt, Turkey)
    c. Stagnant Economies (Growth Rate > 0%) (Greece, France, Germany)
    d. Declining Economies (Growth Rate < 0%) (Venezuela, North Korea)
2. **GDP Prediction**: The GDP Prediction forecasts for you the global GDP of countries using a host of parameters. The high accuracy that we obtained using GBT Regressor makes our results comparable to that of the World Bank and IMF.
3. **Big Mac Index**: The Big Mac Index is where we present to you a wacky indicator that can be used to state the Global disparity and the strength of certain currencies using mundane day to day items like in this case - A Big Mac Burger
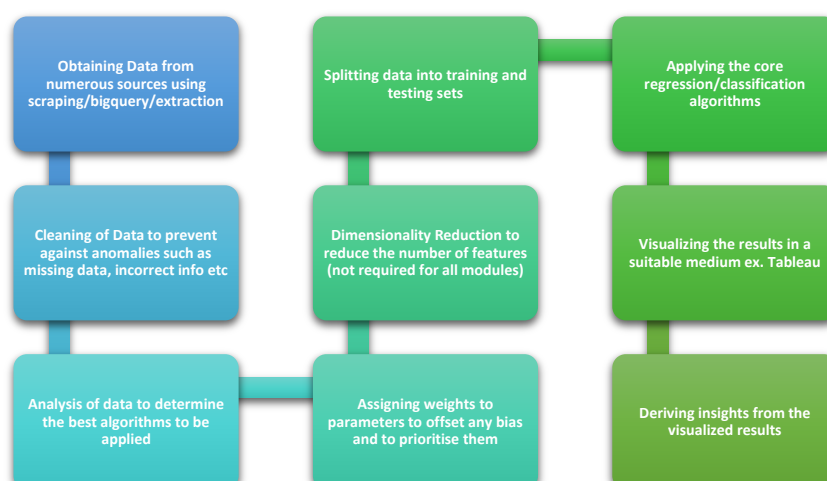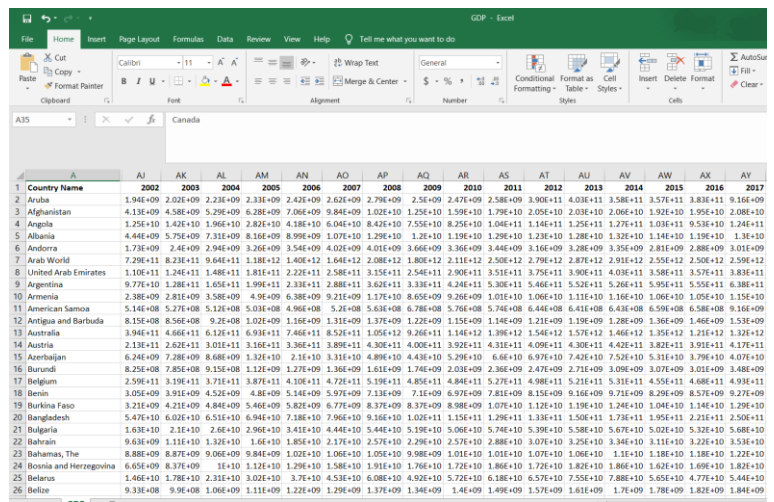
# 2. Methodology



*Fig 2.1 - End-to-End Solution*

## 2.1 Data Processing

The net data that we collected was approximately 600+ MBs. This data ranged from csv files, json files, txt files and other unclassified information. The data was carefully crafted into csv files which were pivoted and edited to obtain a single uniform format of data. The final chosen format was csv files due to their versatility and ease of access using Spark, Python and R.



*Fig 2.2 - Sample File*

By the end of cleaning process, we were left with 20 datasets of features from the years 1960-2017 for all countries. These datasets were then used for analysis. In certain cases (like the LifeExpectancy v/s IncomePerCapita graph) we traced back separate datasets to data since 1800s.

## 2.2 Data Analysis

Our Data Analysis has been performed in a scalable and efficient manner. We have used popular algorithms for analysis for all 3 modules. Primarily we have performed classification, regression, dimensionality reduction and several optimization algorithms. We have used R, Python and Spark machine learning libraries. Spark excels at iterative computation, enabling MLlib to run fast. At the same time, we care about algorithmic performance: MLlib contains high-quality algorithms that leverage iteration, and can yield better results.

For the Global Economy Health prediction, the dataset was scraped and obtained from different sources in a largely incoherent format. This data was cropped to obtain relevant parameters and converted to suitable datatypes. Using Spark with casting made this possible. Later the data was imported to Python module where it was run through Pearson Correlation Coefficient to obtain parameters with clear cut correlation to facilitate their removal and choosing appropriate parameters. Next dimensionality reduction using Principal Component Analysis was performed. The PCA analysis rendered the number of significant components at 3. These parameters were then passed though the Train_test_split module of ScikitLearn's model selection module. The distribution was split among Training, Validation and Test data. The data was then trained through Logistic Regression and Decision Trees. For Logisitic Regression we provided a binary classification and for Decision Trees a 4 way classification. The greater accuracy obtained for Decision Trees (65%) made them the modus operandi for choice. With the Growth Rate predictions done for every year(2010-17), we classified the results into categories and displayed them accordingly in a Tableau representation.

The Big Mac dataset was obtained from 'The Economist' which publishes the Big Mac Index every year. This data was filtered according to our needs and preprocessed accordingly. The base currencies selected were the US Dollar, Euro, British Pound, Japanese Yen and Chinese Yuan. The index was calculated by dividing the local price by the price of the base currency. Next, Linear Regression of GDP and Big Mac Index was taken in order to calculated which currency was undervalued and overvalued.

The GDP forecast analysis was done with data obtained from similar sources as the Growth Rate datasets. The data was also converted into a coherent format using typecasting and conversion of datatypes and substituting missing values. The data of the 5 base parameters was then loaded into

spark through read_csv function and a common table constituting relevant factors was constructed using Spark SQL. For the prediction module, we constructed a schema. The data was loaded and taken through this schema and passed through a vector assembler which was imported using Spark Machine Learning library i.e. Spark.ml.feature and the module Vector Assembler along with StringIndexer and SQL Transformer. For the input to the regressor, we converted the country names into string indexer. The converted data which held the string indexed country names and features was then passed to the GBT Regressor.

The above GDP prediction was put into a pipeline and the model was trained upon the training data was again obtained using train_test_split with a ratio of 3:1. The data was then transformed. The output obtained was the GDP prediction with an accuracy of 87%.

## 2.3 Web Interface and Android App

We created a website ([www.ecoforecast.ml](www.ecoforecast.ml)) to provide an interactive interface for our project. The website contains Tableau dashboards embedded in it. These visualizations provided a dynamic outlook to our website and helps the user analyze the results in layman terms. We also created an Android application to better represent our data for the users who are more phone-savvy. The application displays the latest statistics for 47 countries along with predicted GDPs and other statistics.

# 3. Problems & Challenges

## 3.1 Data Processing

As mentioned above, extraction of data was done with great efforts since several underdeveloped countries do not publish their data annually. Moreover obtaining historical data was another cause of concern for most countries. The data quality in several cases was abysmal and had to be thoroughly revamped. Methods like scraping often fail when such varied data must be collected from hundreds of different websites.

## 3.2 Data Analysis

The data is highly diverse. Parameters like slum populations and high-tech exports have no correlation with each other but significantly affect the growth rate. Due to this, the right algorithm to predict the data is very difficult to find and accuracy obtained is seldom high. The same goes for the Decision Trees we used as their presented accuracy goes only till 60% even after using the parameters causing the greatest impact on the growth rate.

## 3.3 UI & Visualization

Android Studio does not possess Big Data libraries and hence we it becomes exceedingly difficult to embed interactive large scale visualizations onto the app. The Neural network libraries that are often not found compatible with popular platforms. Moreover handling large scale data using Tableau is efficient however heavy visualizations make websites heavier requiring excessive usage of already limited resources.
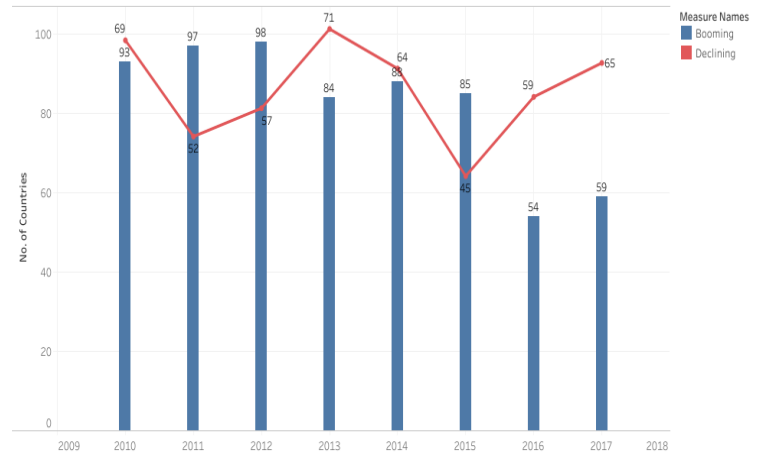
# 4.  Results

The results obtained were displayed graphically as shown below:

1. Health of Global Economy: The graphs below represent how many countries were booming / Declining / Stagnant / Healthy economies. These graphs correspond to the actual growth estimates over the years for ex. The 2011 stock market crash.



*(a)*



*(b)*

*Fig 4.1 – (a) No. of Countries under Healthy & stagnant economy (b) No. of Countries under Booming & Declining economy*
Interactive animated graph made from parameters of Global Economic Health.
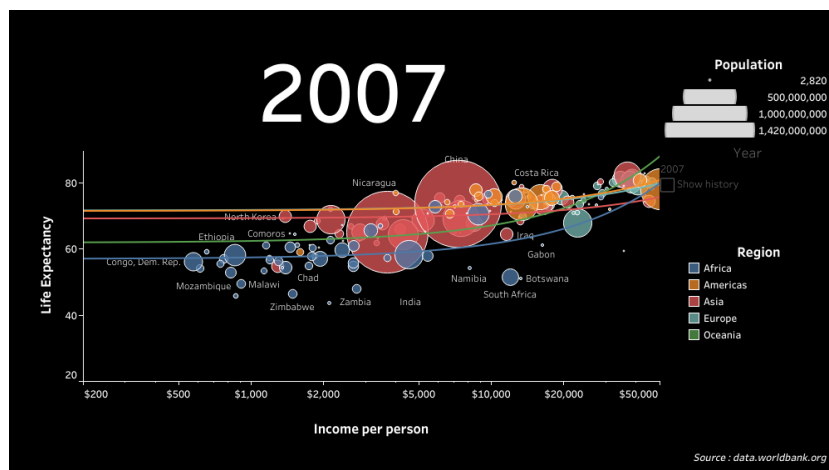


*Fig 4.2 – Income per person v/s Life expectancy*

2. The Big Mac Index: The below graph shows the currency valuations against dollar as the base rate. The graph also shows how these currencies have fared over the years.
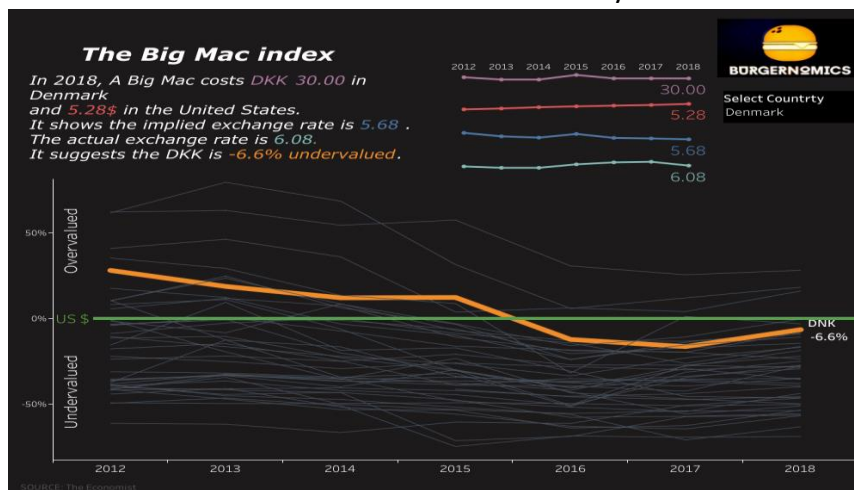


*Fig 4.3 – Currency valuations against dollar*

In the visualization above, we can see how Denmark's currency was valued above dollar till mid of 2015 following which it dropped and was since undervalued.

3. GDP forecast: The GDP forecast was calculated in a tabular format and fed to Tableau's interactive and animated visualization feature to calculate the variance over the years.

```
+-------------+-------------+----+
|     country|          gdp|year|
+-------------+-------------+----+
|      Canada|1.40454533E10|2016|
|       India|2.76999989E12|2016|
|United States|1.85999997E13|2016|
|       China|1.11999996E13|2016|
|       Japan| 4.9500001E12|2016|
+-------------+-------------+----+
```

```
+-------------+-------------+----+
|     country|          gdp|year|
+-------------+-------------+----+
|      Canada|1.69137828E10|2017|
|       India|2.81999992E12|2017|
|United States|1.94000003E13|2017|
|       China|1.22000004E13|2017|
|       Japan| 4.8700001E12|2017|
+-------------+-------------+----+
```

*Fig 4.4 – GDP prediction for 2016 & 2017*

The animated graph shows the increase in GDP gradually over several decades the last of which have been predicted and displayed.
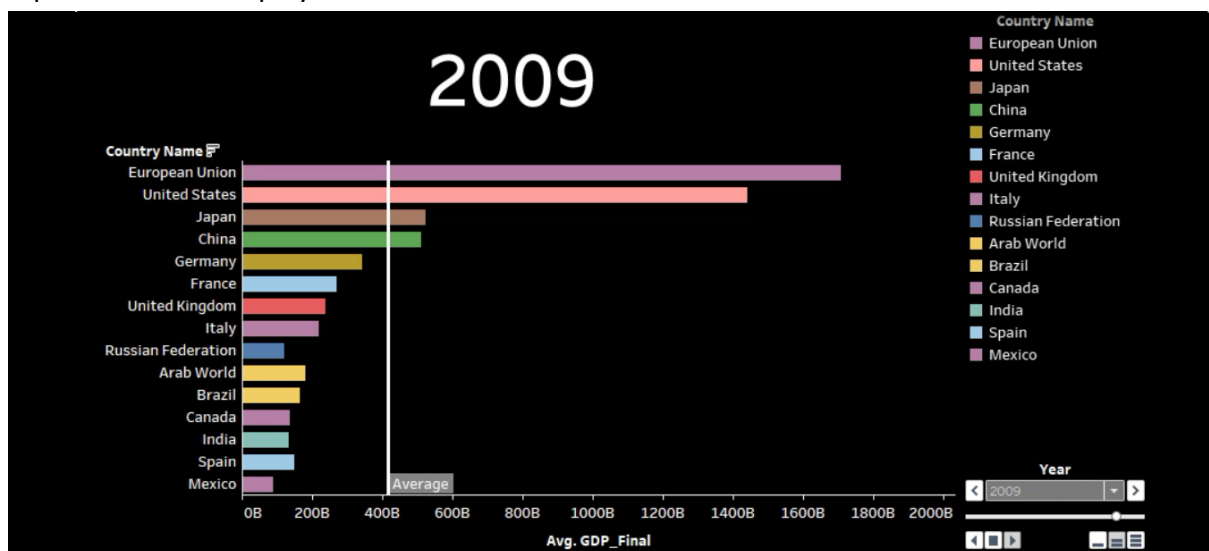


*Fig 4.5 – GDP over the years*

# APPENDIX – Project Summary

| Category | Points |
|---|---|
| **Getting the data**: Scraping data from various websites, using BigQuery to fetch data from Google cloud, datasets from UN,IMF,WB | 2 |
| **ETL:**  **C**leaning the dataset and perform Extract-Transform-Load using Apache Spark and R. | 2 |
| **Algorithmic Work:** MLLib, Scikit to perform analysis and predict GDP. | 5 |
| **Bigness/Parallelization:** 600+ MB of data for 200 countries and groupings with 150+ features and 300,000+ values | 3 |
| **UI:** Developed an interactive website to display the results and embedded it with Tableau dashboards. | 2 |
| **Visualization:** Interactive and animated dashboards created for three features in Tableau and Bar Charts for classification of health of the economy. | 4 |
| **New Technology:** Using Android Studio developed an app to create an interactive medium for the user to see the results. | 2 |
| **Total** | 20 |