**The University of Texas at Dallas**
Naveen Jindal School of Management

# MIS 6380.002 DATA VISUALIZATION PROJECT

# GLOBAL MIGRATION ANALYSIS

# INSTRUCTOR: DR. JUDD D. BRADBURY
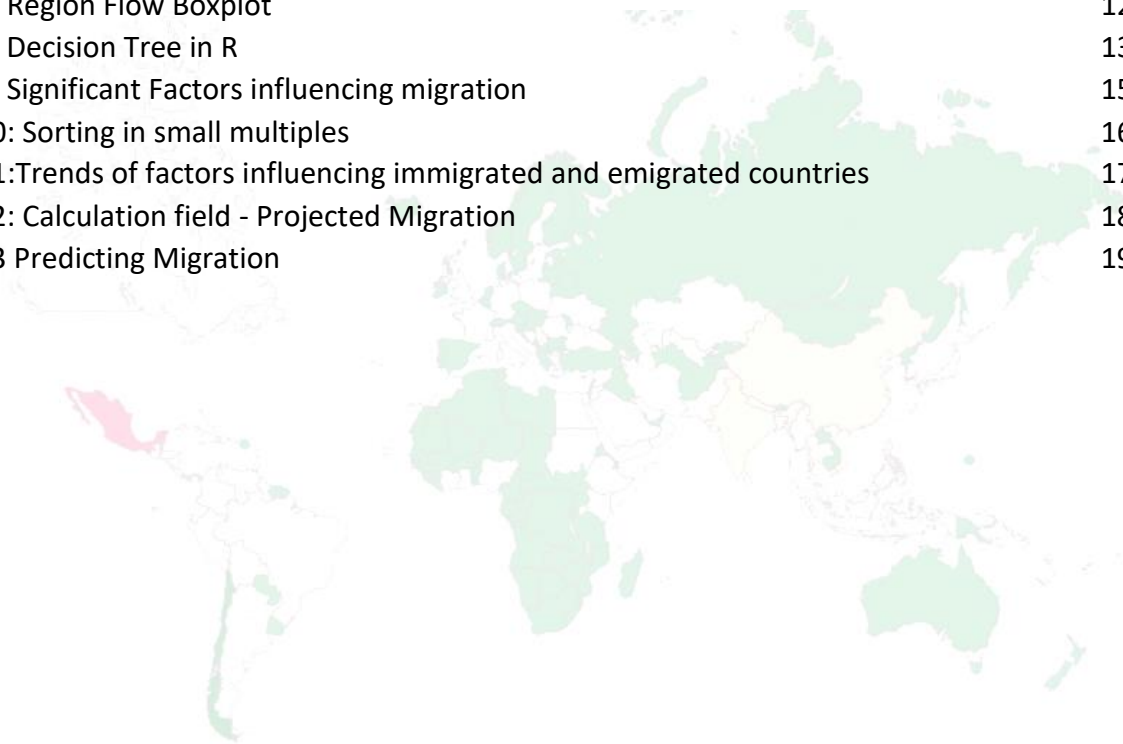
*by*
**GROUP – 12**

*Roobiya Rahamathulla Khan - rxr156930*
*Venkatesh Viswanathan – vxv161130*

# Table of Contents

## Table of Figures

# INTRODUCTION

Migration and development have become a topic of high interest in many countries around the world. Academia as well as policy makers acknowledge that further qualitative as wells as quantitative research is need to address migration and mobility issues to elaborate policy recommendations and implementations. Our team aimed at exploring the "GLOBAL MIGRATIONS' and discovering hidden patterns and relations between migration using the potential of visual representations, mobility and development. Using this data, we have successfully performed a whole series of visualizations for different research questions.

# OVERVIEW

Migration is a key feature of our increasingly interconnected world. It has also become a flashpoint for debate in many countries, which underscores the importance of understanding the patterns of global migration and the economic impact that is created when people move across the world's borders. The past few decades have seen an incredible rise in global immigration, with numbers of immigrants raising three times higher than those recorded in 1960. Currently, 3.3% of the world's population is living in a different country than the one they were born in.

Although media attention has been focused heavily on refugee migration, refugees are only a small share of migrant numbers. Migration is increasingly driven by opportunity-seeking behaviours due to economic disparity. Considering this mass migration movement, it's unsurprising that global migration has created major demographic changes in the world.

Perhaps one of the most interesting observations stemming from this middle-class growth is the massive cleavage that divides the middle class in developed countries versus that of the emerging economies. Although both enjoy the same lifestyle and opportunities, the middle-class of developed countries is chronically stressed and lacks confidence in the future. In contrast, the middle class of emerging economies is increasingly optimistic, opportunistic, and enjoys an 8% growth per year.

As migration flow data are often incomplete and not comparable across nations, the data we obtained has an estimate of several movements by linking changes in migrant stock data over time. Using statistical missing data methods, the data is estimated on five-year migrant flows that are required to meet differences in migrant stock totals. The data collected is from 1990 to 2010.

The objective of this visualization project is to generate various insights from the data set and support it with the interactive and compelling visualization evidence.

## BUSINESS OBJECTIVE

As a research organisation based out of New York, to better understand the global migration flow of people, we have worked to develop various insights and compelling evidences of people migration. As mentioned perviously, economic specific factors analysed include GDP growth, Labor participation rate, population total and employment among other factors. Integrated data analysis is performed using R software and statistical techniques like K-means clustering, Decision trees and Linear Regressions are utilised in the analysis.

## DATA ACQUISITION

Our **Project** focuses on **"Immigration of people from countries all over the world and the factors that influence the migration"**. Our data sets are taken from multiple sources. Global migration data set captures the number of people who change their country of residence over 5-year periods and the other data set on world development indicators that can be viewed as an information on economic and its associated factors of a country to support our first data set.

### 1.Primary Dataset:

Our **Primary data source** has **38,416 records** and **17 attributes**. The data is taken from http://www.global-migration.info and is based on global data flow collected for every five years from 1990 to 2010. The data can be directly pulled from the below link as csv file http://www.global-migration.info/Data%20on%20the%20global%20flow%20of%20people_Version%20March2014.csv. Country_origin indicates the countries people migrate from and Country_destination denotes the countries people migrate to. We also have a region wise origin and destination to get an overall idea of region based migration.

*File Name: Migration Flow in 1995-2010.xlsx*

### 2.Secondary Dataset:

Our **Secondary data source** provides various factors that enable Immigration. The data is taken from http://databank.worldbank.org/data/reports.aspx?source=world-development-indicators Out of 1500 various indicators, we pulled out the factors that are closely related to our analysis. Some of the factors include GDP growth, employment in industries, employment in services, life expectancy, literacy rate, the net enrollment rate of primary school children, health expenditure etc. It **has 264 records** and **134 attributes** collected every five years from 1990 to 2010.

*File Name: Economic Factors Influencing Migration.xlsx*

## DATA MANIPULATION

As the first step, we cleaned the data sources before merging into one dataset. The primary data set has records of countries with same origin and destination. Those records were removed. Thus, we have 38,221 records but the attributes remain same (17). In the secondary data set, we find that there are many missing values. Therefore, the attributes with more than 70% of missing values were deleted. Thereby we get 114 attributes with no change in number of records (264). The missing values denoted by two dots(..) were deleted and left as blank spaces so that tableau considers those null values as numeric values and includes in measures instead of dimensions.

The economic factors of all years 1990-2010 are grouped by average for each country. This is done for the purpose of performing statistical analysis such as decision trees and Regression in R.  Also, the net migration of country is calculated by identifying the inflow and outflow of a country.

## DATA MERGING

The next step is to merge both the datasets. We see that the countries are named differently in both. So, we renamed it for a proper match in their country names. Now we have performed inner join between the two datasets based on country name in Tableau. This helps us to determine the factors which influence the migration.

## INSIGHTS AND DATA VISUALIZATIONS

### I.INSIGHT – MIGRATION PATTERN

Our first insight is based on the below facts that explain the common migration pattern.

- The USA as one of the country having highest number of immigrants followed by United Arab Emirates, Spain and Italy
- Less migration to countries such as Pakistan and Afghanistan

The above insight is drawn upon the evidence captured from the below visualizations.

### II.INSIGHT – FACTORS INFLUENCING MIGRATION PATTERN

On combining the migration data set with the development factors of each country, we obtained rock solid evidence from the visualizations why certain countries were highly immigrated and

others not. Moving more labor to higher-productivity settings boosts global GDP. Migrants of all skill levels contribute to this effect, whether through innovation and entrepreneurship or through freeing up natives for higher-value work There are many development factors considered, including, but not limited to, GDP, employment in Industry, Labor force participation, Life expectancy, net migration and population.

## Insight III – PREDICTING PROJECTION

This insight speaks about the migration prediction of all countries using the regression analysis in R. Migration pattern predicated shows a different pattern of migration on performing the analysis. Regression analysis is done using R software. The statistical method used is the linear regression model. A Linear regression model uses a dependent variable and list of the independent variables. The dependent variable in our case is the Net Migration where Net migration is the net total of migrants during the period, that is, the total number of immigrants less the annual number of emigrants, including both citizens and noncitizens and the independent variables are factors identified in our decision tree.

### 1.TWO DECADES OF WORLD MIGRATION

Before performing an analysis on migration data, opted to visualize the overall flow of people between countries. This would help us to find out the countries with high and low migration. Therefore, we decided to visualize migration data with a heat map and classify it based on migration type (Type 1 - Very Low migration, Type 2 – Low migration, Type 3 – High migration & Type 4 - Very high migration)
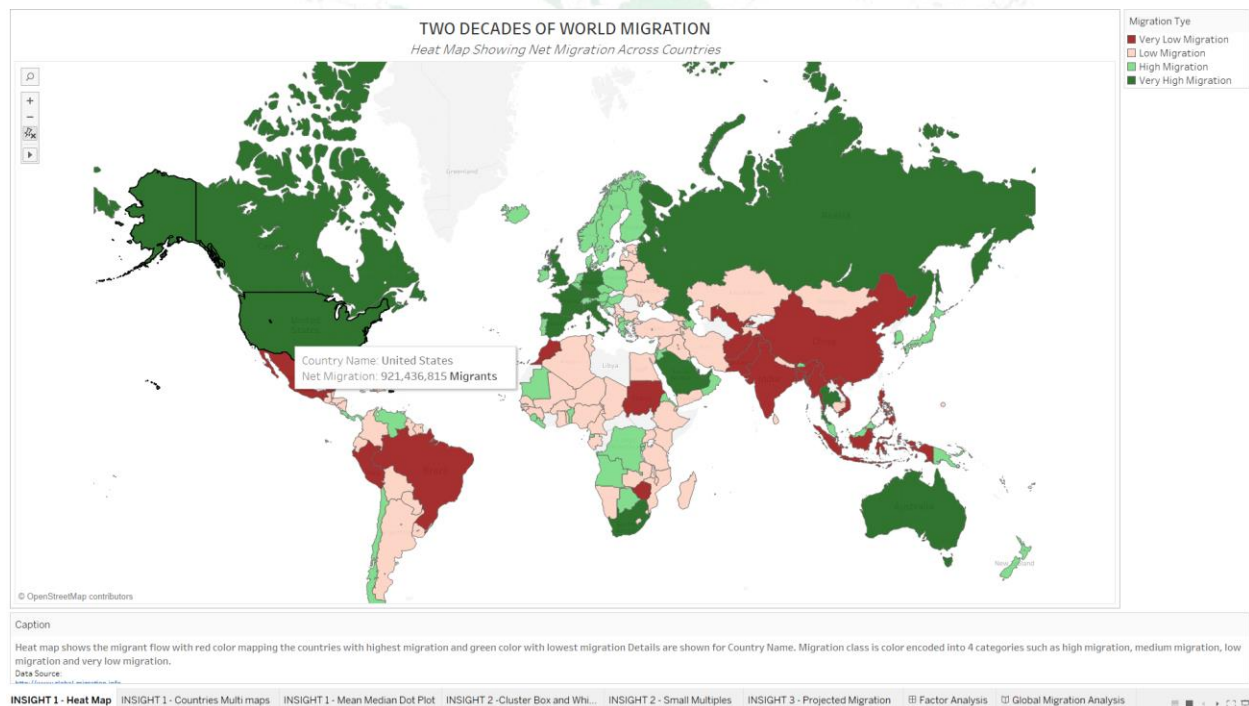


*Figure 1: Heat Flow Map*

The above map shows the global migration at Country level for the past two decades, which allows the target audience (i.e. everyone) to identify migrant countries on the spatial level. Countries are color encoded with Net Migrants. The green color represents the countries with high migrants and red color represents the countries with least migrants.

## 2. FACTORS ANALYSIS- MEAN, MEDIAN, DOT PLOT

We represent all the characters of our analysis in a graphical approach rather than traditional shape representation approach. The characters in the migration flow are countries and factors. These characters are visualized using a Mean-Median-Dot plot.

The Mean-Median-Dot plot gives the important insight on factors influencing migration. The custom sorting, filtering, dynamic titles are used.

The interesting migration classes out of four are class 1 and class 4, as these are the countries where the migration is very low and very high respectively. As per the secondary dataset, there are various factors that affects the migration. At a glance, the factors like GDP Growth, Population Total, Life Expectancy Birth, Employment in Agriculture, Consumer Expenditure, Health Expenditure play a vital role in the migration of people.

Initially, maximum value of each factor is identified and their average is calculated. For example, highest value of GDP is 10, it's average is calculated as below,



Figure 2: Calculation field- Maximum GDP

The following maximum value is calculated for all the factors.

=# Max ConExp
=# Max EmpAgri
=# **Max GDP**
=# Max HealthExp
=# Max LifeExp
=# Max PopTotal

The length is calculated for each factor to give the scope of maximum value range.

| Lenght EmpAgri | EconomicFactors_Linear_Regressi (EconomicFactors_Linear_Regression) (2) | ✕ |

AVG (-40)|

▶

The calculation is valid.                    Sheets Affected ▾    Apply    OK

*Figure 3: Calculation field - Length of employment agriculture*

These calculated fields are used to create the dot plot. The mean and the median reference lines for each factor are added. Then graph is designed in such a dynamic way by using the dynamic sorting method.
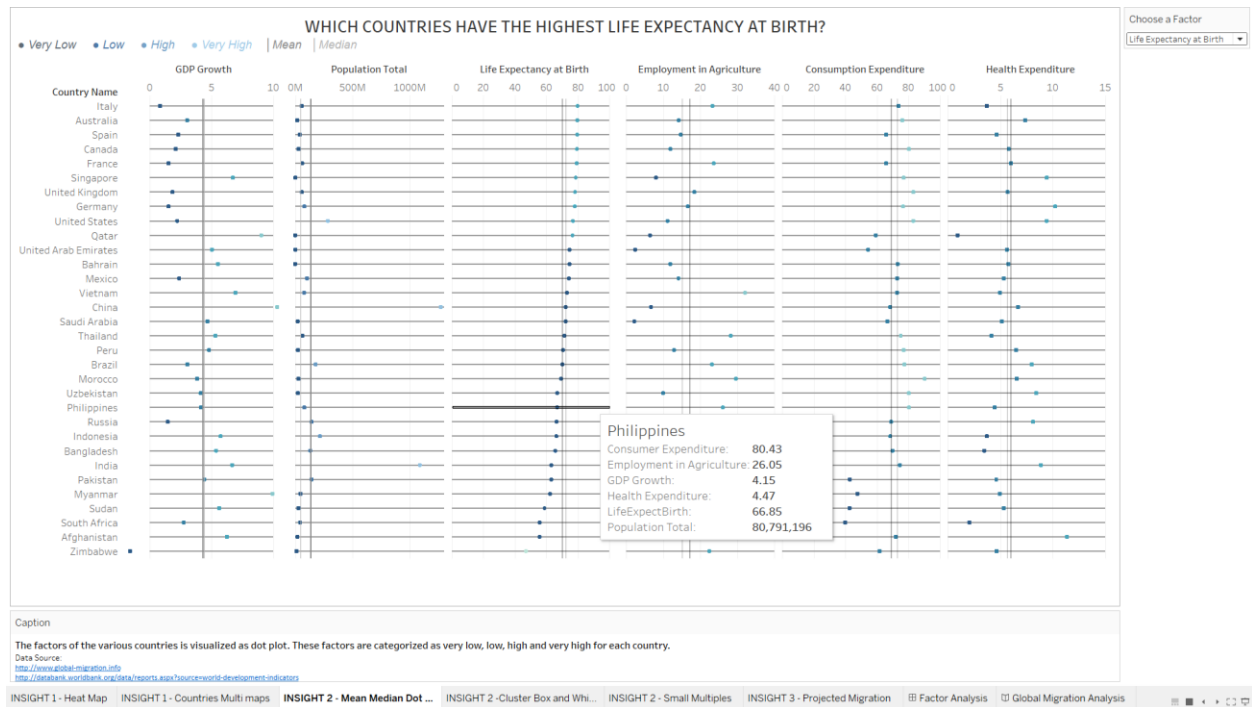
**Inference**:



*Figure 4: Countries having highest life expectancy at birth*

We can clearly see the characters (countries and factors). The dynamic sorting helps us to understand where each country stands in each factor. This helps us to compare factors against all other countries. Italy stands first in the Life Expectancy at Birth, it is also way ahead in the mean and median of the same factor. It also has high employment in agriculture, closely to the median in the consumer expenditure. These factors may stimulate the migration. On the other side, the GDP of the country pretty much on the lower side and the expenditure to health industry is very less, which makes the people to rethink before migrating to Italy.

Similarly, we can choose various factors and see which country ranks best in what factors. Below are the screenshots of the Consumer expenditure and GDP,
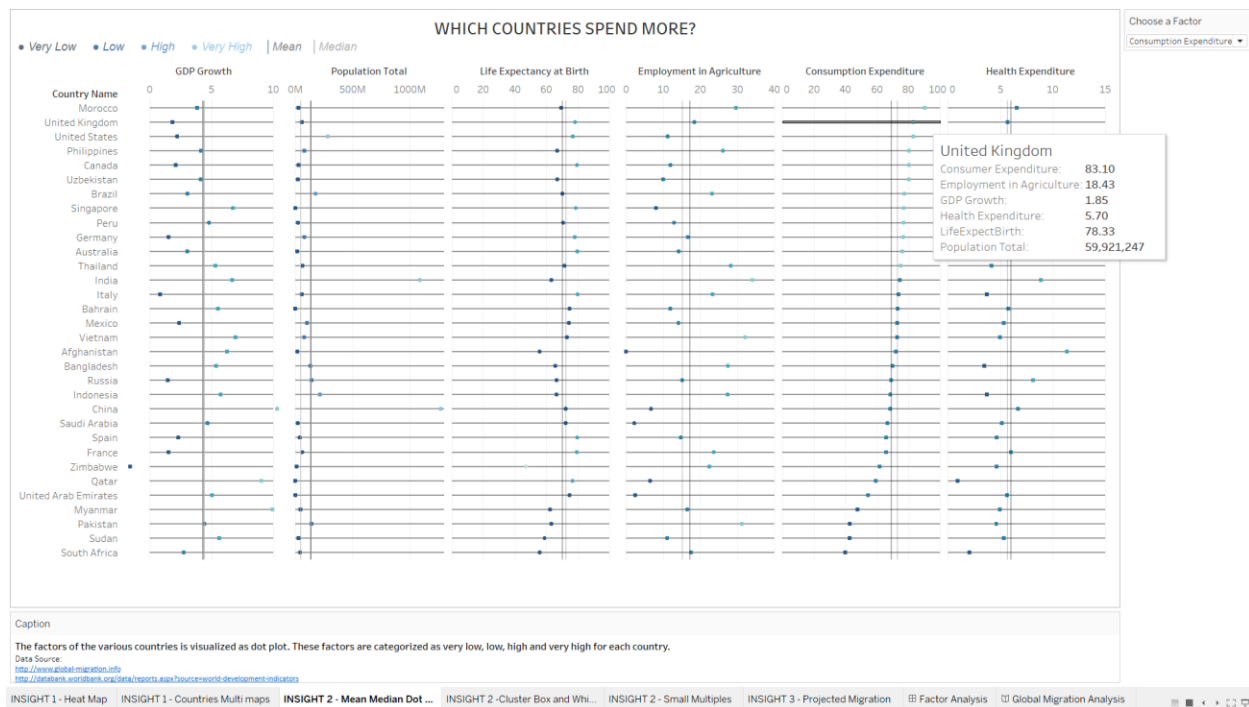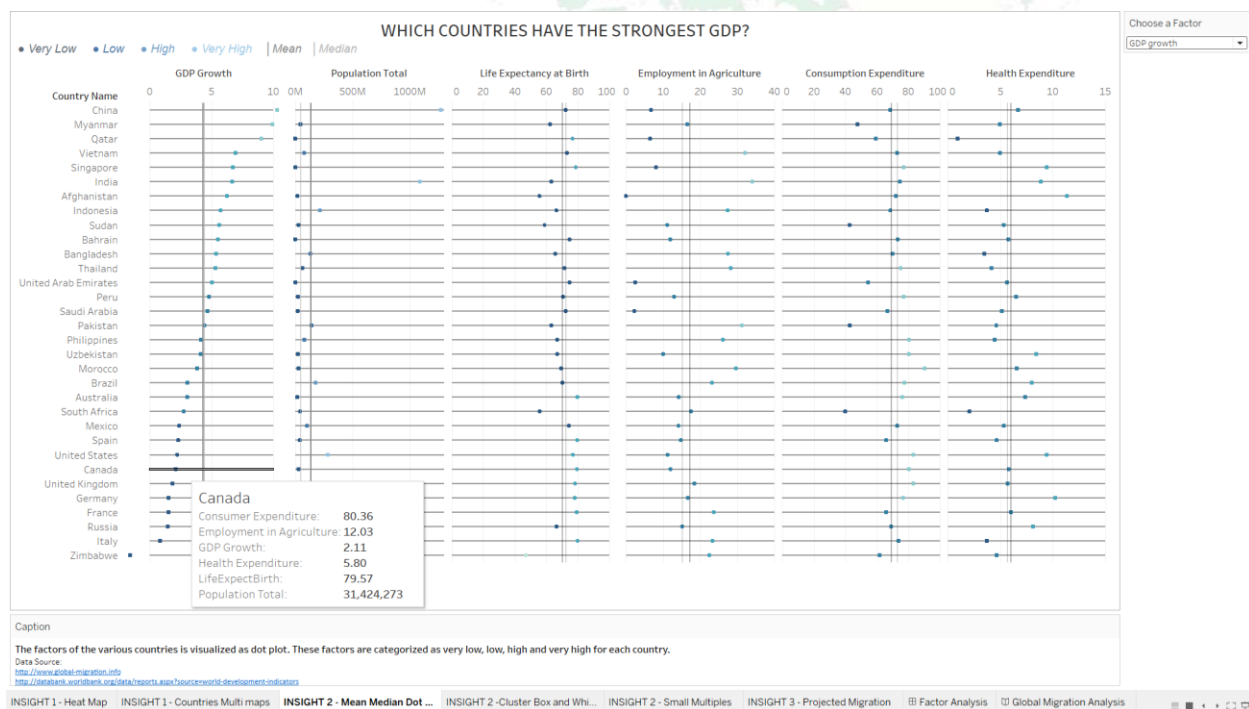
*Figure 5: Countries revenue in expenditure*



*Figure 6: Countries with the strongest GDP*

The mean and median can be used effectively to comprehend how good a particular country in a particular factor.


## INTEGRATED ANALYSIS WITH R

Using the joined dataset, we created a boxplot to determine the distribution of region flow and decision tree to determine the most significant variables influencing the global migrations.


### a. Box Plot:

Box Plot is a most standardized way of displaying the distribution of data, its central value and its variability on quartiles. The below box plot depicts the overall region flow over the two decades from 1990 to 2010. We can infer that Africa has the versatile number of migrants compared to other regions, while its mean remaining less.
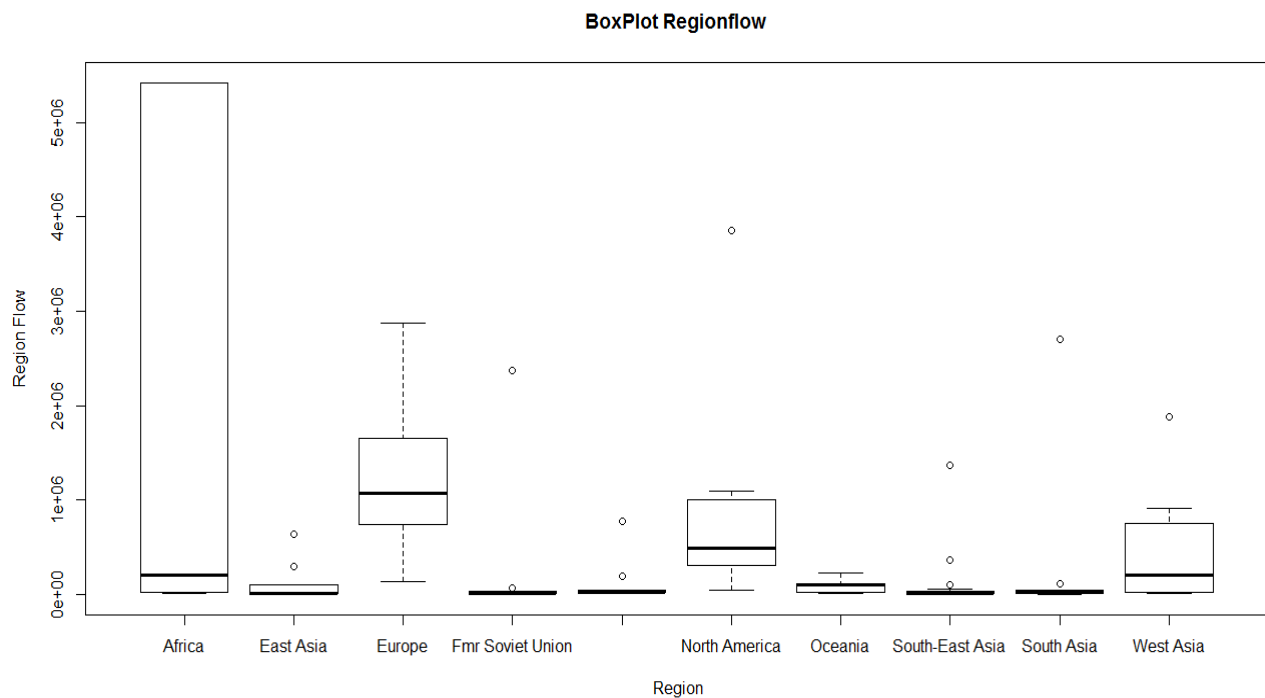


*Figure 7 Region Flow Boxplot*

## b. Decision Tree:

Decision tree is basically a nested structure that uses a branching method to illustrate every possible outcome of a decision, hierarchized from top to bottom on most significant factors. It is the most common predictive analysis technique. The first node is called the Root node and the last nodes as Leaf node.

The complete decision tree is shown below which was created in RStudio, the root node always denotes the variable(vector) of the highest importance. We can see that, the countries with Life Expectancy Birth Lesser than 74 tend to have more migration.

Therefore, the most significant factors derived from the decision tree are as follows:

1. **Life Expectancy at Birth**
2. **Total Population**
3. **Population in Slum**
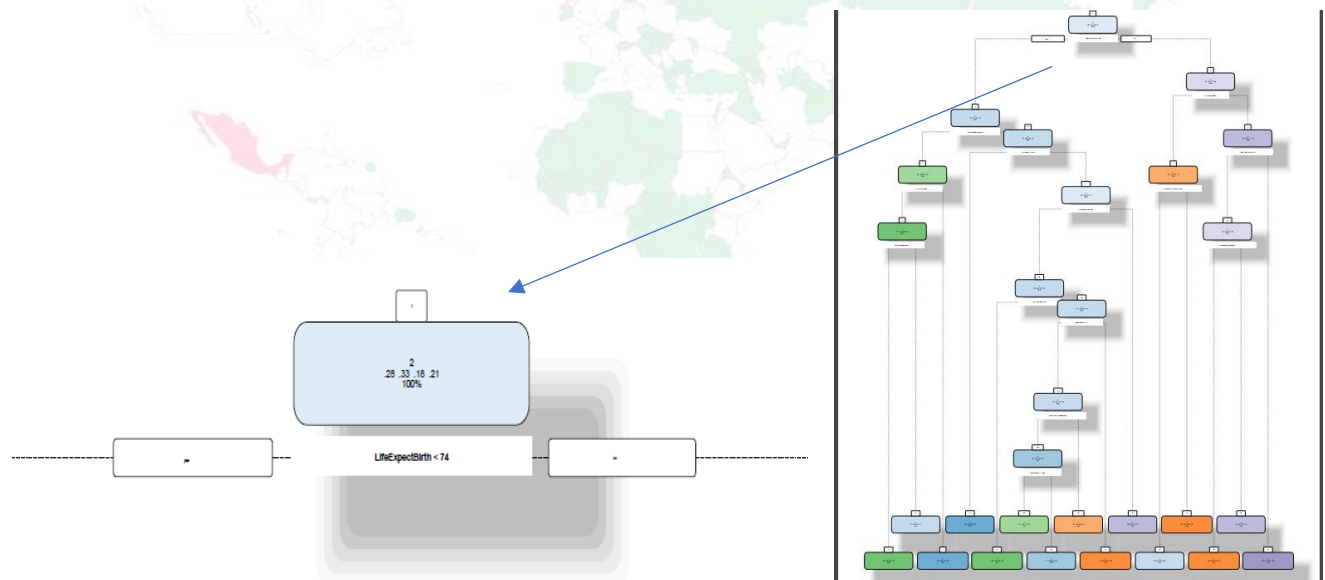4. **High Technology Exports**
5. **Technical Grants**



*Figure 8 Decision Tree in R*

We further performed K-mean Clustering using R based on the above significant factors influencing migration and determined the relationships among them. In addition, we integrated the results of the Clustering analysis to Tableau and visualized our data and found interesting pattern.

## 3. COUNTRIES MULTI MAPS

Country Multi Maps are created 2*2 matrix showing the top 4 emigrant countries ranked based on life expectancy at birth. India, Pakistan, Bangladesh and Mexico are the top 4 countries with the highest people outflow. We see that the health expenditure is high in these countries favoring emigration. We also created layered map showing the population inflow in various states of US which has high amount of immigration population.



*Figure 8:  Top 4 Ranked emigrant countries and their life expectancy at birth*

## 4. SIGNIFICANT FACTORS INCLUENCING MIGRATION

This Boxplot visualization is an integrated analysis with R using decision tree and clustering. R functions and models were invoked through tableau by creating a calculated field named as Cluster. We have performed K means for clustering.

We have created four Migration Bins based on the Net Migrations as tabulated below and visualized them against each cluster. The Cluster's has been color coded

| Migration Bin | Net Migrants |
|---|---|
| 1 | < -3,70,000 Migrants |
| 2 | -3,70,000 < Migrants < 0 |
| 3 | 0 > Migrants < 3,70,00 |
| 4 | 0 > 3,70,000 Migrants |

## 4a. Cluster Analysis:

- Cluster 4 has only the Immigration Bin 1 with highest total population, this shows higher the population corresponds to higher emigration
- Cluster 2 and Cluster 4 are majorly comprised of high technology exports and Cluster 6 has the Primary School enrollment below the average
- Cluster 2 has highest Life Expectancy Birth with values above the average and Population in Slum below the average. This makes evident that life expectancy at birth contradicts the population is slums
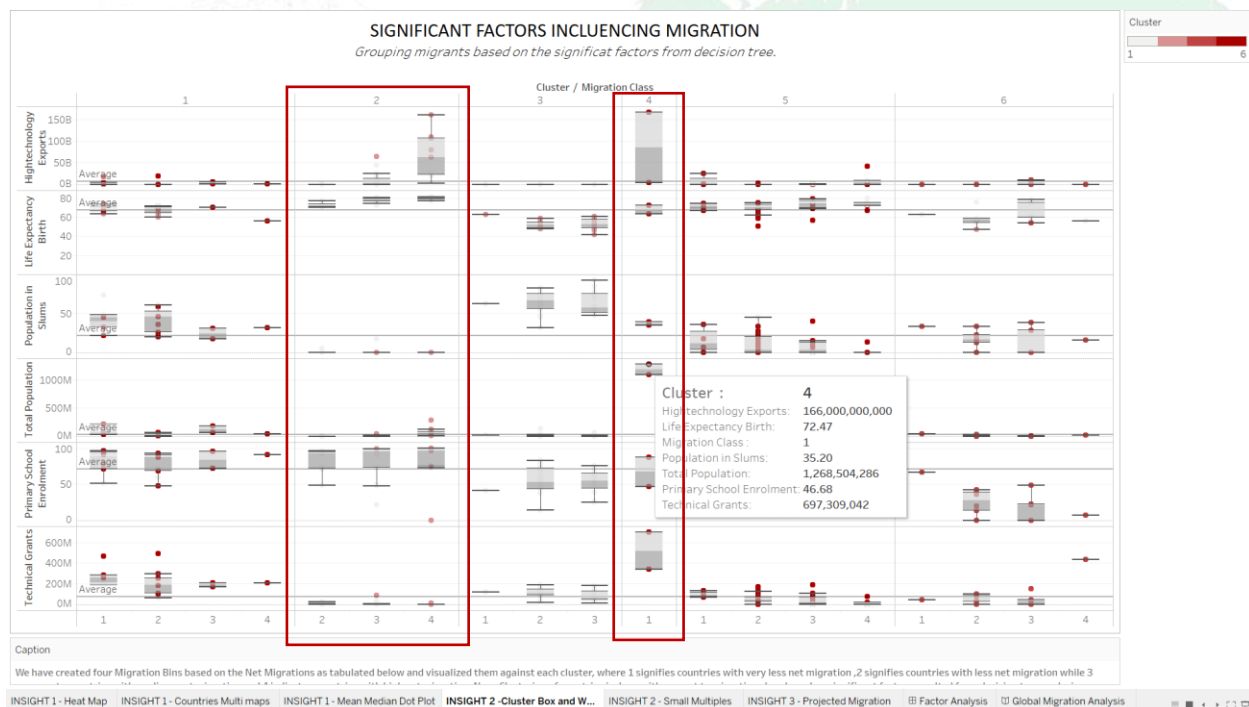


*Figure 9 Significant Factors influencing migration*

## 5. TRENDS OF FACTORS INFLUENCING EMIGRATED AND IMMIGRATED COUNTRIES

Based on the results of clustering performed, we determined the factors that influence migration in a positive and negative way. These factors show a pattern similar to the characteristics of highly emigrated and immigrated countries. A migration class is created to categorize the records with a negative migration and positive migration, which are called Emigration and Immigration. Emigration denotes people moving out and Immigration denotes people moving into one country. A well-known fact is highly immigrating countries are developed economies of the world and highly emigrating countries are developing economies. On analyzing these factors further with the migration class viz., Immigration and Emigration, below are the insights and evidence created. A small multiple graph is created is out of this analysis. Countries of importance are filtered and sorted in this graph.



*Figure 10: Sorting in small multiples*

- Total Population:

The total population is based on the de facto definition of population, which counts all residents regardless of legal status or citizenship. The values collected were midyear estimates. On visualizing population with the migration class, highly developed countries, such as the USA and the United Kingdom have population lesser than developing countries. Countries like China and India exhibit huge dominance in this factor with 1-2 billion.

This means is that country/place where the population is higher people tend to move out. High Technology Exports High-technology exports are products with high R&D intensity, such as in aerospace, computers, pharmaceuticals, scientific instruments, and electrical machinery. Data    are in

current U.S. dollars. On visualizing this factor with the migration class, countries leading in the high technology export attract more population from developing countries, so these countries highly immigrate. The reverse is the case for developing countries. It can also be interpreted in the view that countries good in this factor, create more jobs in the industry and attracting skilled labor from other countries.

- Population in slums:

Population living in slums is the proportion of the urban population living in slum households. A slum household is defined as a group of individuals living under the same roof lacking one or more of the following conditions: access to improved water, access to improved sanitation, sufficient living area, and durability of housing.

On visualizing this factor with migration class, countries that are dominant in this factor are highly emigrated countries as this explain the people moving out to developed economies escape from poverty.
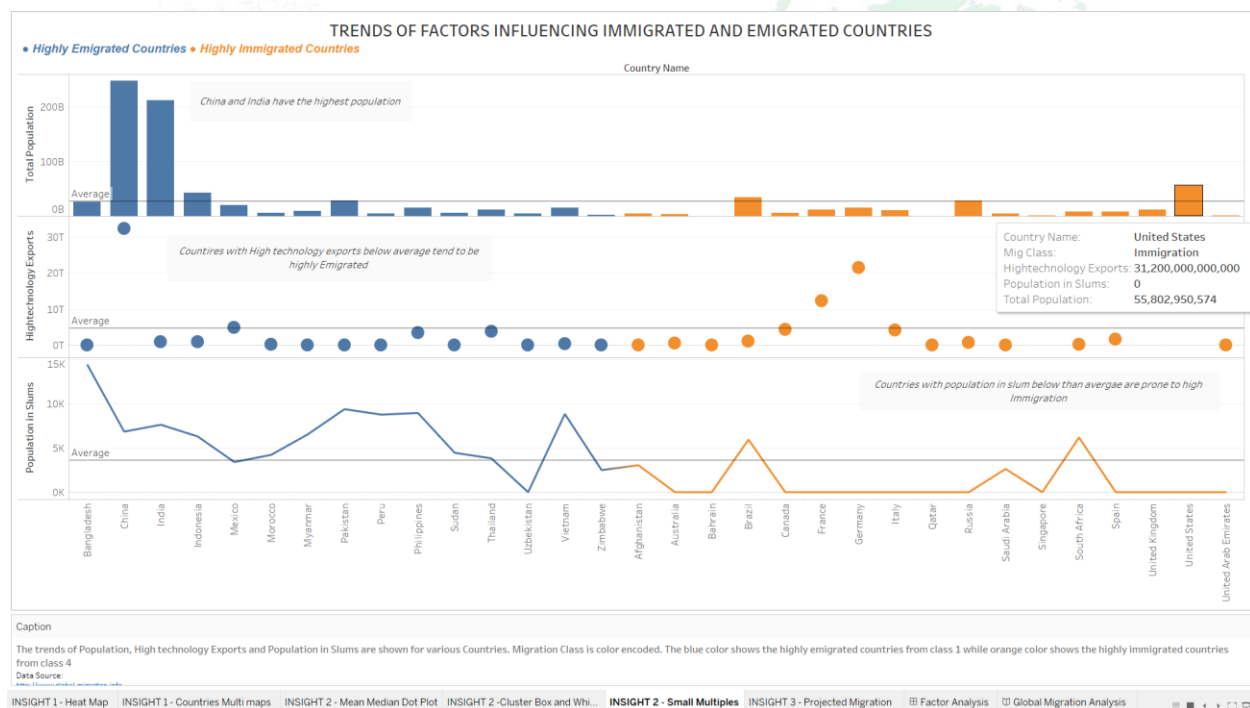


*Figure 11 Trends of factors influencing immigrated and emigrated countries*

Below is the table that shows how a country's migration can affect based on our analysis.

| Factors ( High) | Immigration | Emigration |
|---|---|---|
| Population Total | Less | High |
| High Technology Exports | High | Low |
| Population in Slums | Less | High |

## 6.PREDICTING MIGRATION

Below calculated field is created in Tableau to perform Linear Regression,



```
Projected Migration                    flow+                                                    ×

Results are computed along Table (across).
SCRIPT_REAL("
fit<-lm(.arg1 ~ .arg2+ .arg3+ .arg4+ .arg5+ .arg6+ .arg7+ .arg8+ .arg9)

fit$fitted
"
,
AVG([Net Migration]),AVG([Life Expectancy Birth]),AVG([Hightechnology Exports]),
AVG([Technical Grants]),AVG([Labor Participation Rate]),AVG([National Income]),
AVG([Total Population]),AVG([Population in Slums]),AVG([Registered Business])
)




                                                                    Default Table Calculation
The calculation is valid.                              Sheets Affected ▾    Apply        OK
```

*Figure 12: Calculation field - Projected Migration*

On plotting the Projected Migration calculated field, a surprising migration pattern is obtained for some of the countries. Countries that exhibited high immigration and emigration were just displayed in the visualization.

One of the surprising pattern found is The United States of America, with a two-year-old decade migration of 5,192,065 is now predicated as 1,929,677. This reduction can be attributed to various reasons such as the increase in the population of the country, less GDP growth and other political reasons.  The other pattern is that Germany's projected migration is 1,624,378 which is higher than its two-year-old decade migrant number. The thought process behind this could be the refugee migration from the Arab countries as a result of Arab Springs. Europe too envisions high migration as the reason could be the same as mentioned for Germany.

Another interesting fact to know that Mexico's projected migration has reduced significantly and this could be due to efforts of the USA government to stop illegal migration from Mexico.

India and China show same migration pattern wherein the projected migration is higher than that of previous two-year decade migration. One or more economic factors of these countries as discussed in small multiples could influence the migration number.
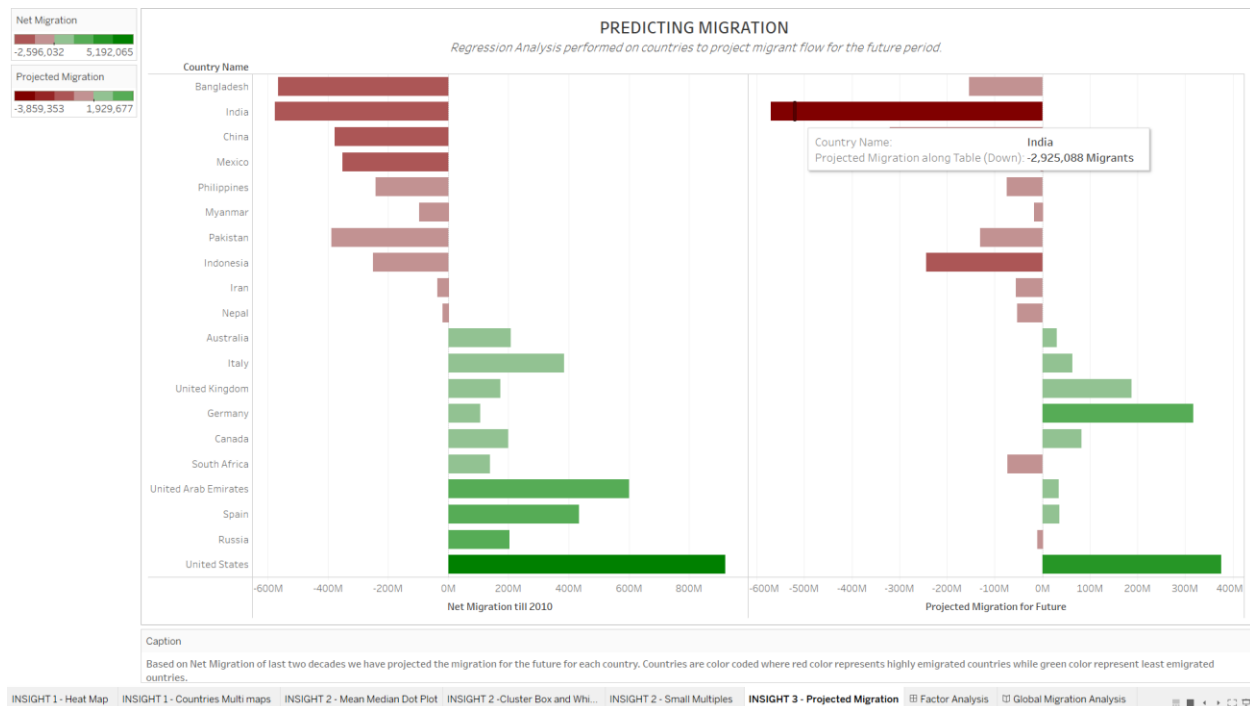
*Figure 13 Predicting Migration*
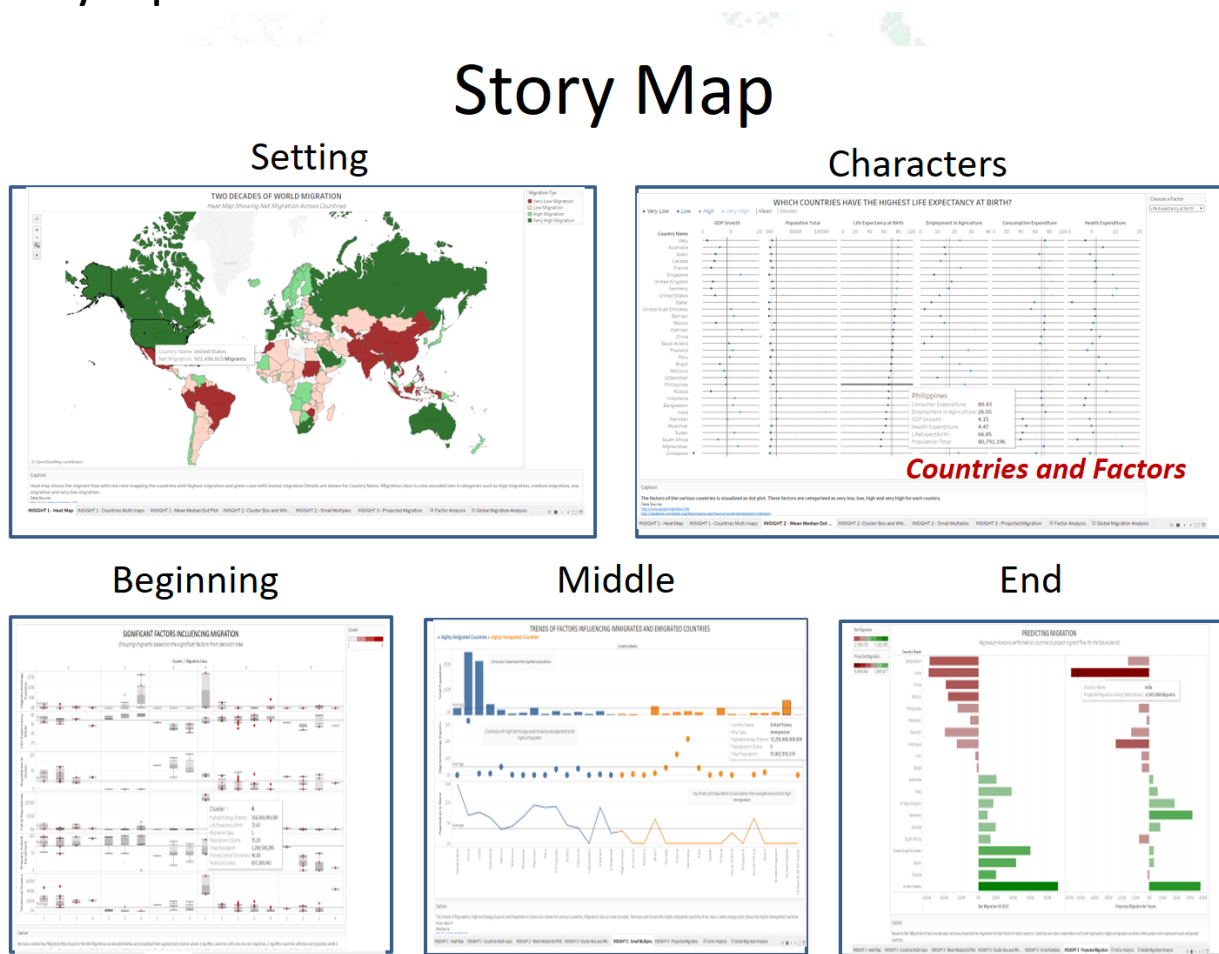
# DATA STORYLINE

## 1.Rhetoric – Logic:

The motive behind choosing this rhetoric was to allow users to understand and draw conclusions from the presented visual evidence. Our project explores different factors contributing to immigration and population movements. Why people migrate?

People migrate for various reasons. Economic, social, political and ecological factors are the main forces driving migration. Through our story, we aim to analyze the various economic factors that influences the country flow from 1990-2010. Push factors like poverty, unemployment, etc encourages people to move out of their home countries while pull factors like opportunity, safety, freedom, etc attract people as their migratory destinations. We see that the incentive to migrate is a lot higher in areas that have a high level of economic inequality. During 2000-2005, the more developed regions of the world gained an estimated 2.6 million migrants annually from the less developed regions. This amounts to about 13.1 million migrants over the whole period. Northern America gained the most from net migration: 1.4 million migrants annually.
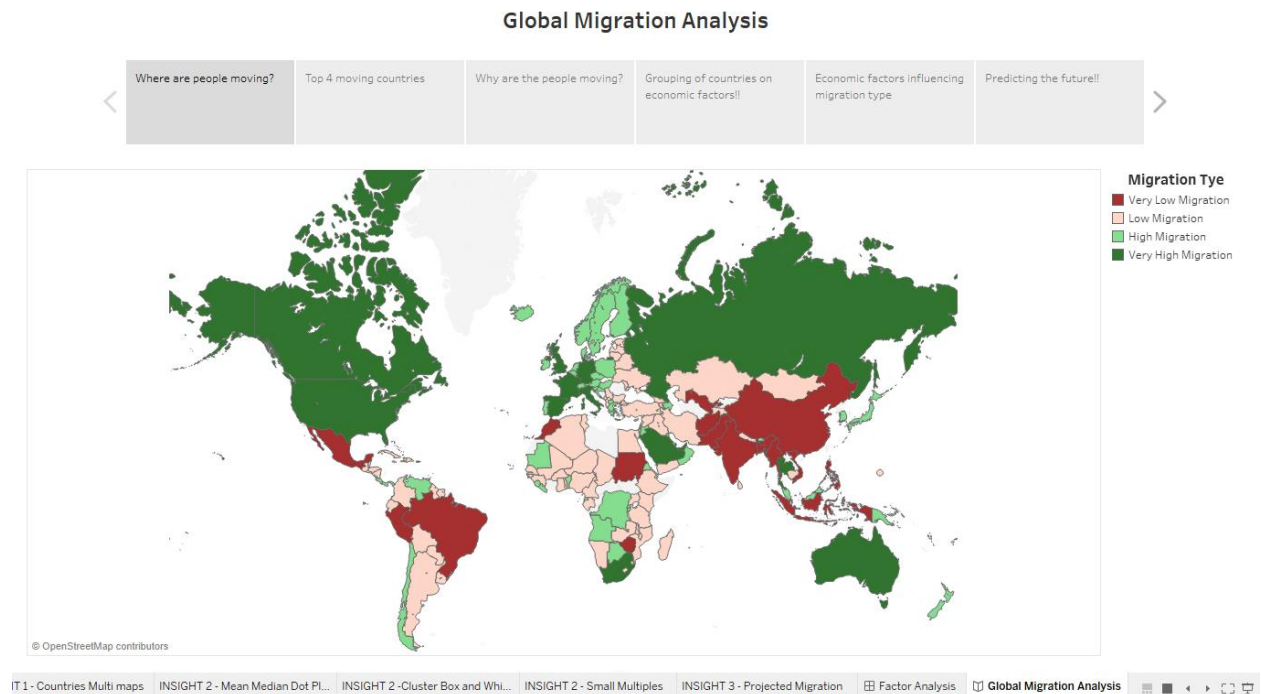
## 2.Audience:

In this project, we are presenting our visual analysis to help people understand on migration. Even though, migration of human population started as early as 2 million years ago, the past few decades have seen an incredible rise in global immigration. Pull factors within the destination country are more likely to influence the decision making process of economic migrants. These major factors influence population movements and immigration. Our visual analysis helps us in this.

## 3.Story Map:

# Story Map

### Setting



### Characters



*Countries and Factors*

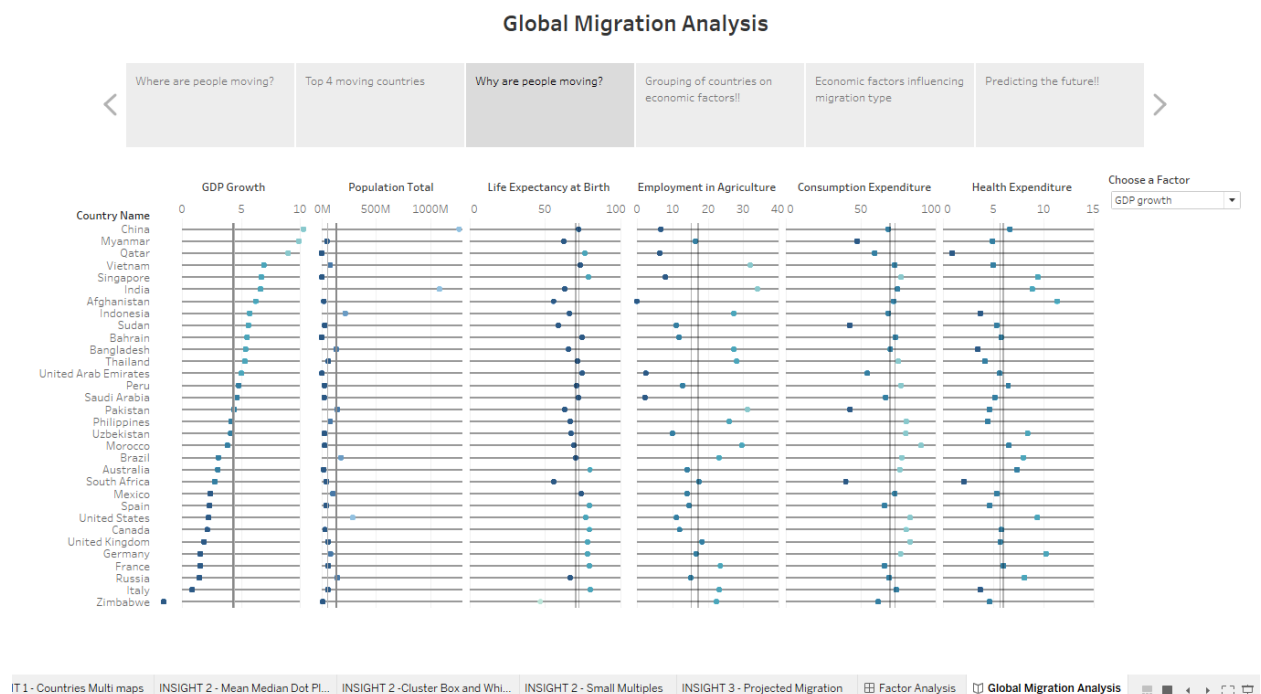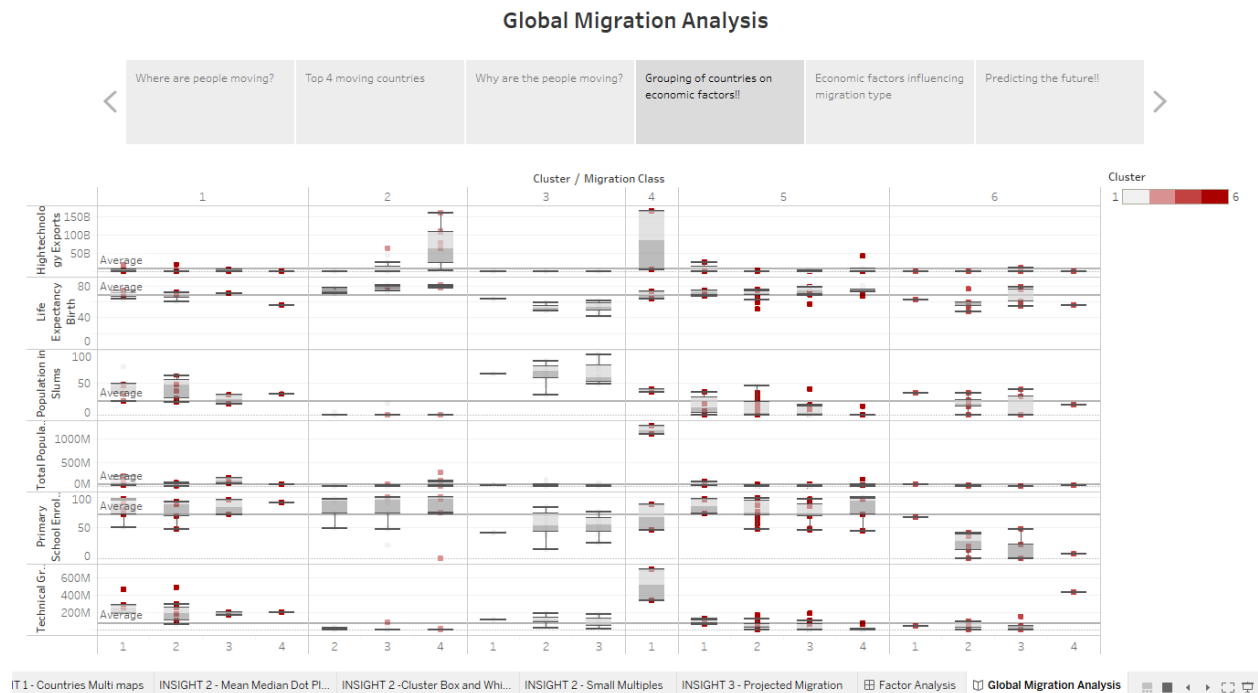### Beginning



### Middle



### End

## 3a. Setting:



The heat map shows the global migration across countries from the past two decades (1990-2010). The green color represents the countries with pull factors. These countries have more immigrant population. The red color represents the countries with push factors. These countries have high emigrant population.
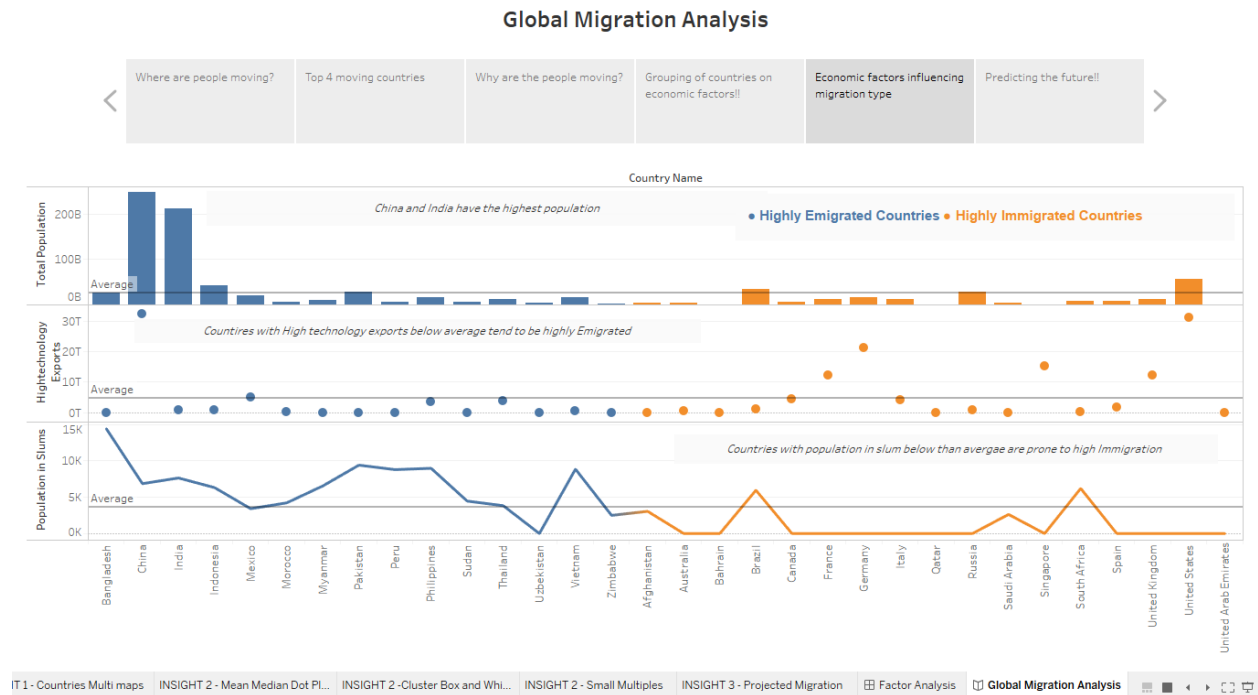
## 3b. Characters:



The main characters for our story are countries and factors.
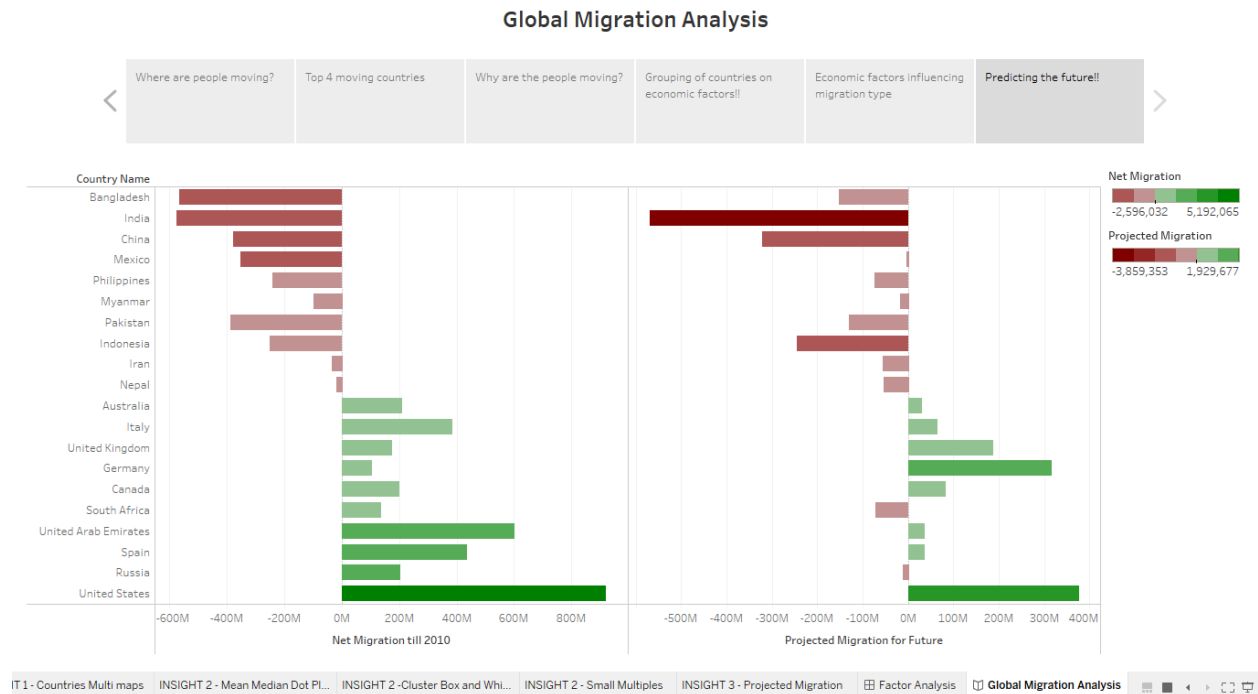
## 3c. Beginning:



We see that factors like life expectancy, primary school, technology, population play major role in immigration. If economic conditions are not favorable and appear to be at risk of declining further, people tend to emigrate to a better economy. Economic migrants are drawn towards international migration because of the prospect of higher wages, better employment opportunities, technological advancement, communications and transport and thereby improving their financial circumstances. These migrants are most likely to come from middle-income countries where the population is becoming increasingly well educated. Salaries and wages, however, are likely to remain relatively low compared to those of individuals with a similar educational background in other, higher-income countries. This disparity has the potential to lead to some highly-skilled individuals from developing countries migrating to more developed countries.

**3d. Middle:**

**Global Migration Analysis**

| Where are people moving? | Top 4 moving countries | Why are the people moving? | Grouping of countries on economic factors!! | Economic factors influencing migration type | Predicting the future!! |



IT 1 - Countries Multi maps   INSIGHT 2 - Mean Median Dot Pl...   INSIGHT 2 -Cluster Box and Whi...   INSIGHT 2 - Small Multiples   INSIGHT 3 - Projected Migration   ⊞ Factor Analysis   ⬢ Global Migration Analysis

We compare factors like high technology exports, population in slums, and total population to understand the huge economic disparities between class1 and class4 countries. Since the mid-twentieth century, the nature of migration has also become largely influenced by globalization by increasing the demand for workers from other countries. Economic disparities between developing and developed nations have accelerated with globalization. In 1900, the ratio of the average income of the five richest countries in the world to the 5-10 poorest countries was about 9:1. Today that ratio is 100:1. These disparities among countries combined with limited opportunities for employment has stimulated increased migration from developing to developed nations.
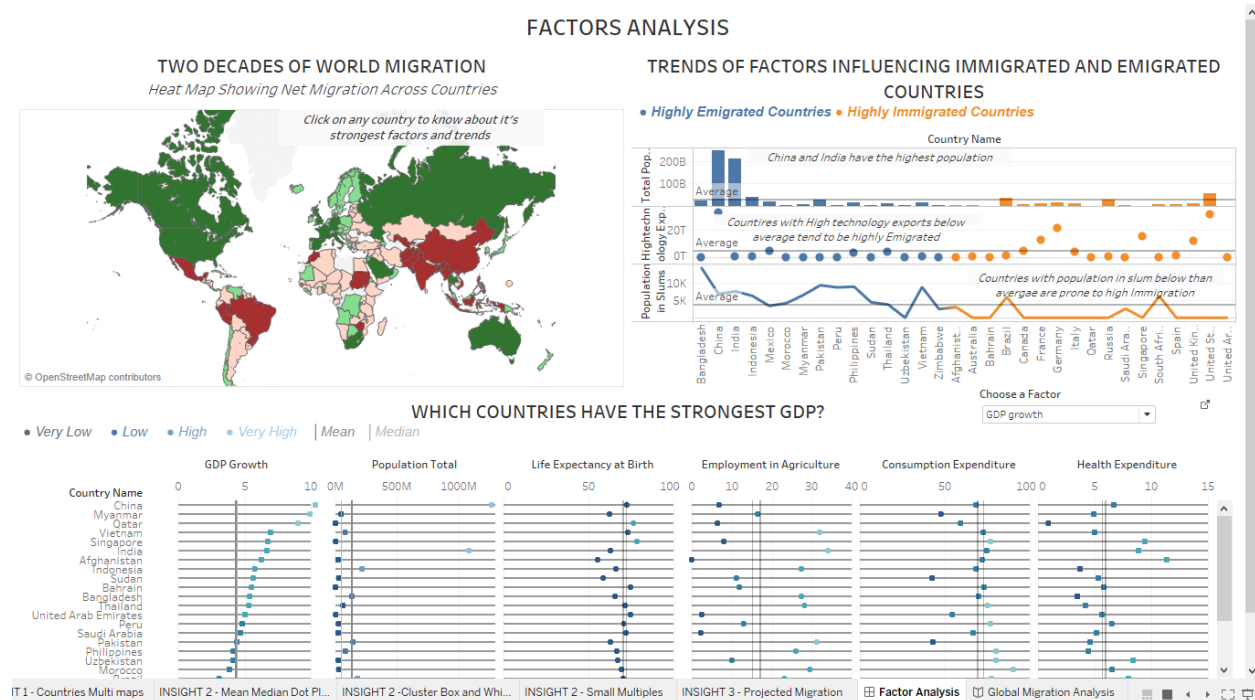
**3e. End:**



As the low- and middle-income countries of today continue to develop and the high-income countries experience slower economic growth, migration from the former could decline. This is seen in projected analysis where countries like Bangladesh, china, mexico seem to have lesser migration compared to previous years. While developed countries like south Africa and Russia reveal a shocking trend with higher emigrant population. Aging populations and low fertility rates in industrialized countries like Germany, United Kingdom have resulted in a greater demand for service-sector jobs and employment opportunities. Developed countries like the U.S. have come to rely on immigrant labor to fulfill their labor needs and will need to so even more in the future as the country faces a mass retirement of baby boomers.

## DASHBOARD - FACTORS ANALYSIS

To perform an analysis based on significant factors, we need to refer the countries and their migration class. We can analyze the countries which are immigrated or emigrated and explore their significant factors that influenced migration.

For instance, choosing Mexico, we can see the various factors that enabled people of Mexico for migration. The mean-median dot plot shows how far it is from mean and median of the respective factors.

## REFERENCES

- Guy J. Abel and Nikola Sander (2014). Quantifying Global International Migration Flows. Science, 343 (6178).
- World Data Bank.