

## **Data Science Tableau Project**

By

Sidak Pal Singh Woodwal

102003149

COE-7

Thapar University

**Submitted to**

Dr. Kashish Goyal

Thapar University

<b>S.no</b>	<b>Subject</b>	<b>Page</b>
1	Introduction and Datasets	2-3
2	Data Cleaning Process	4-7
3	Dashboard Explanation	8-12

**Objective:** To analyse and try to establish the relationship between Civil liberty and Judicial independence across nations and highlight the importance of an independent judiciary in maintaining and protecting the liberties of its people.

**Datasets used:** The datasets have been sourced from gapminder.org. Gapminder is an independent educational non-profit fighting global misconceptions. 2 datasets have been used in this project, namely:

**ATtribution: FREE DATA FROM WORLD BANK VIA [GAPMINDER.ORG](https://gapminder.org), CC-BY LICENSE**

- 1) **Judicial independence index:** The judicial independence sub attribute denotes the extent to which the courts are not subject to undue influence from the other branches of government, especially the executive. This sub attribute includes five indicators: High Court independence, Lower Court independence, compliance with High Court, compliance with judiciary and law and order. These five indicators were aggregated into the judicial independence sub attribute using IRT.
- 2) **Index of the civil liberties: "This is one of five subindexes of the democracy index composed from the following indicators estimated by experts in each field:**
  1. Is there a free electronic media?
  2. Is there a free print media?
  3. Is there freedom of expression and protest (bar only generally accepted restrictions, such as banning advocacy of violence)?
  4. Is media coverage robust? Is there open and free discussion of public issues, with a reasonable diversity of opinions?
  5. Are there political restrictions on access to the Internet?
  6. Are citizens free to form professional organisations and trade unions?
  7. Do institutions provide citizens with the opportunity to petition government to redress grievances?
  8. The use of torture by the state.
  9. The degree to which the judiciary is independent of government influence. Consider the views of international legal and judicial watchdogs. Have the courts ever issued an important judgement against the government, or a senior government official?
  10. The degree of religious tolerance and freedom of religious expression. Are all religions permitted to operate freely, or are some restricted? Is the right to worship permitted both publicly and privately? Do some religious groups feel intimidated by others, even if the law requires equality and protection?
  11. The degree to which citizens are treated equally under the law. Consider whether favoured groups or individuals are spared prosecution under the law.;
  12. Do citizens enjoy basic security?

## Civil Liberties

[illegible]

## Judicial Independence

[illegible]

## Data Cleaning Process

## Summary:

- 1) Rename variables
- 2) Replace NA values
- 3) Derive new columns from existing ones
- 4) Merge cleaned datasets into one
- 5) Use merged set to create final variables used in visualization

## Civil Liberties Dataset:

The zoo library is used to plot the na\_pareta plots

The dataset was read into RStudio and I carried out an evaluation of the summary and structure of the dataset.

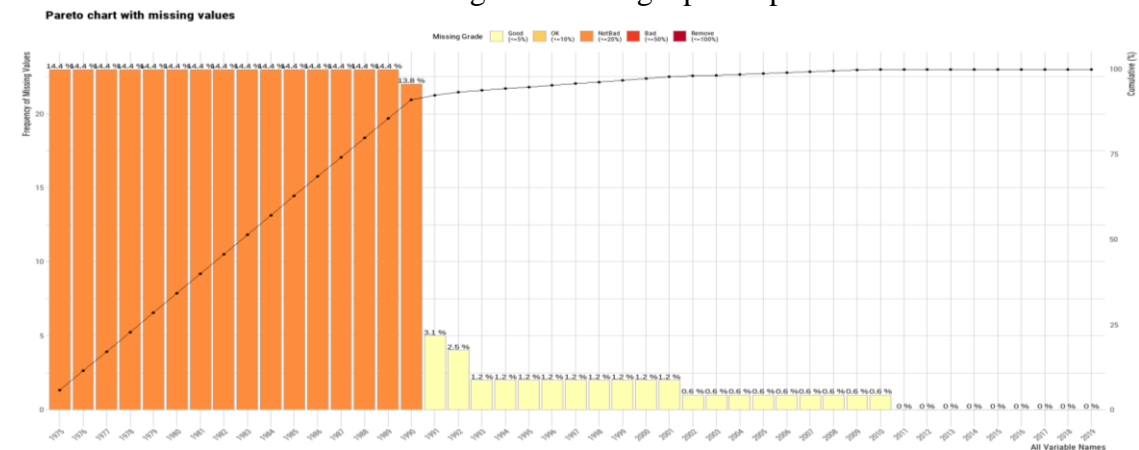
[illegible]

Then the variable names were modified to drop the prefix “X” from the year names.

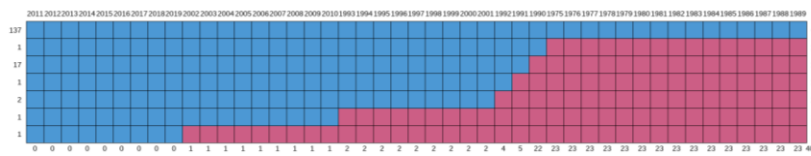
```
> colnames(fd)
[1] "x1977" "x1976" "x1977" "x1978" "x1979" "x1980" "x1981" "x1982" "x1983" "x1984" "x1985" "x1986" "x1987" "x1988" "x1989"
[16] "x1990" "x1991" "x1992" "x1993" "x1994" "x1995" "x1996" "x1997" "x1998" "x1999" "x2000" "x2001" "x2002" "x2003" "x2004"
[31] "x2005" "x2006" "x2007" "x2008" "x2009" "x2010" "x2011" "x2012" "x2013" "x2014" "x2015" "x2016" "x2017" "x2018" "x2019"
> colnames(fd) <- gsub("X", "", colnames(fd))
[1] "1977" "1976" "1977" "1978" "1979" "1980" "1981" "1982" "1983" "1984" "1985" "1986" "1987" "1988" "1989" "1990" "1991"
[16] "1992" "1993" "1994" "1995" "1996" "1997" "1998" "1999" "2000" "2001" "2002" "2003" "2004" "2005" "2006" "2007" "2008"
[35] "2009" "2010" "2011" "2012" "2013" "2014" "2015" "2016" "2017" "2018" "2019"
```

Following this I measured count of missing values country-wise as well as sum of missing values. These values were handled using the `na_locf` which means last observation carried forward as in case of these indices, year by year value for a nation is much more likely to be similar to the previous year rather than an aggregate measure.

Here is a visualization of the missing values using a pareta plot

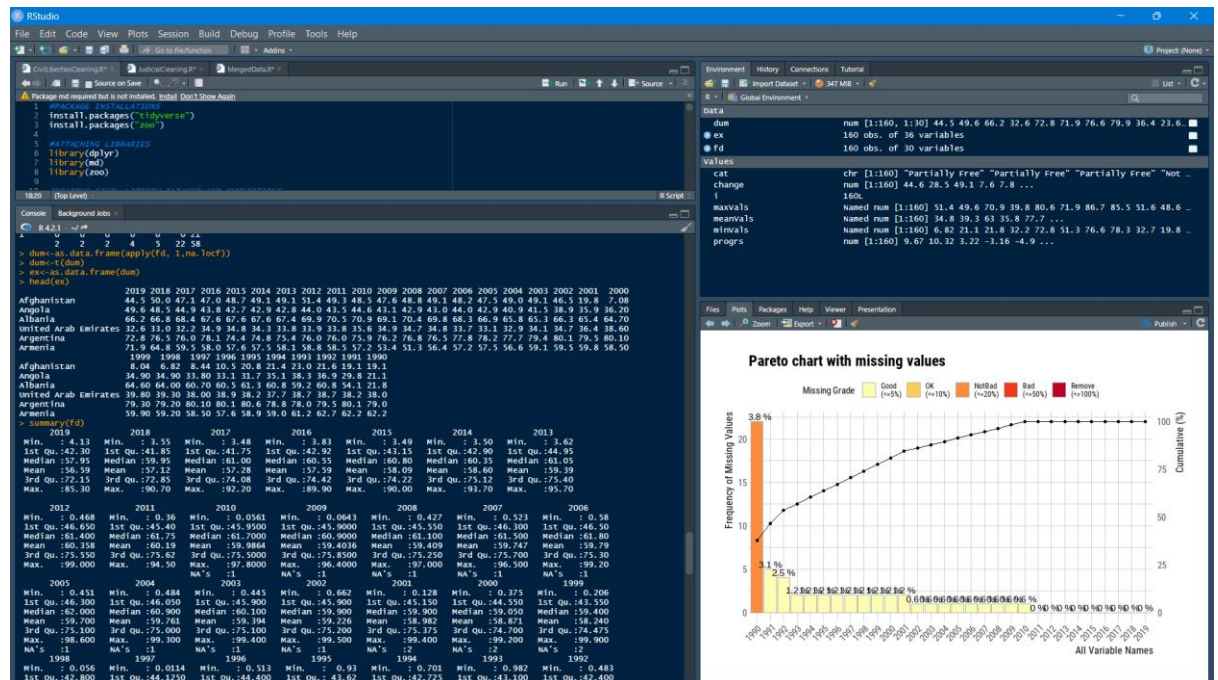






Plot of non-missing(blue) vs missing(red) values

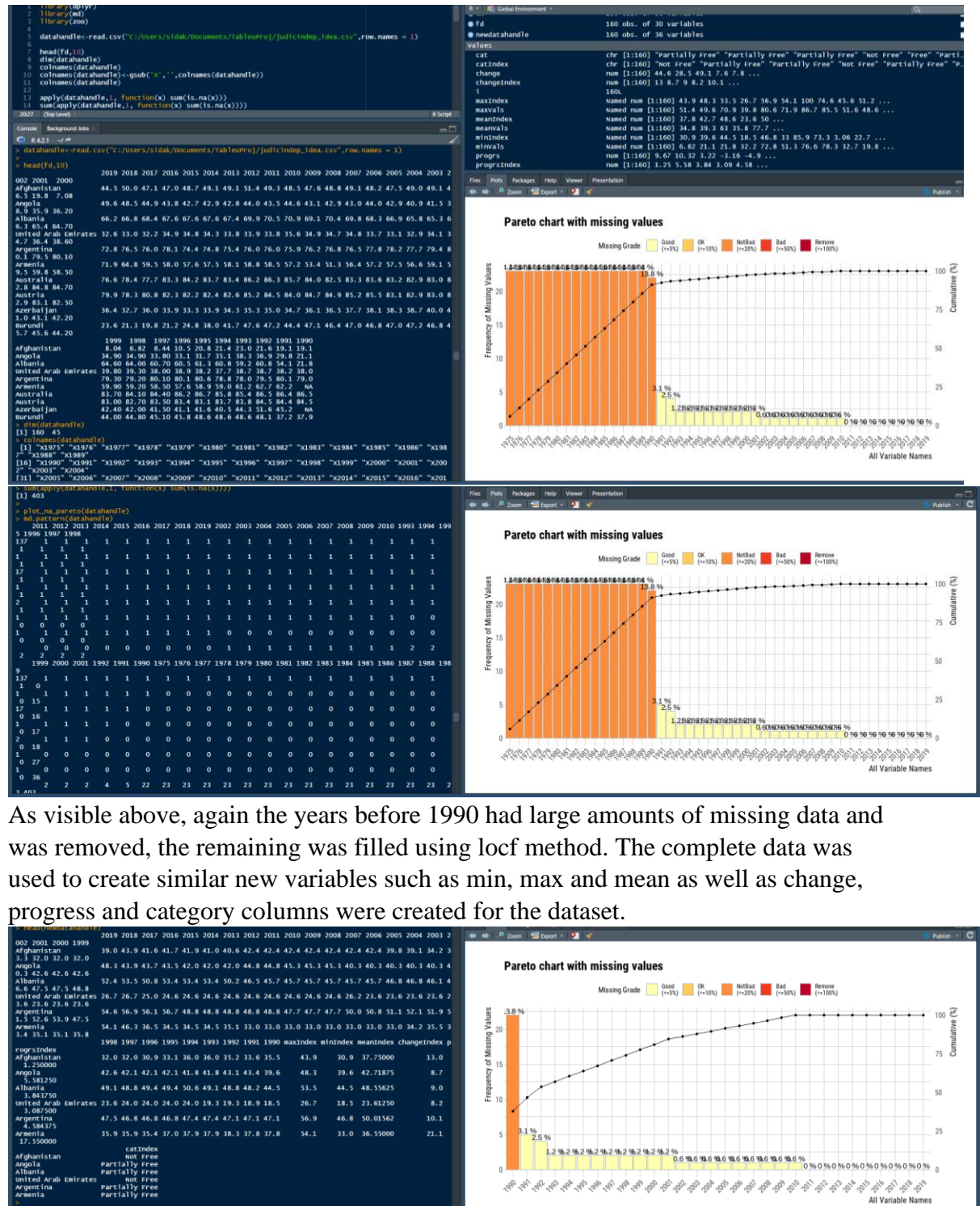
Since there are many missing values before 1990, I truncated the data to only include data from 1990 to 2019.



After truncating data, I removed the NAs by replacing them with the previous year's index for the country. Following this, I created new columns in the data frame including minimum, maximum and mean values of liberty indices of each country . Using these values, I created new columns which calculate the difference between the maximum and minimum value of indices to enumerate the maximum change a nation has been through. I created another column which measures the latest index of a nation (2019) against the mean of values from 1990 to 2019 to show how a nation performed lately. Finally, I created a categorical column on basis of maximum index of a nation and classified the first and last quantile as Free and No Free respectively.

## Judicial Independence Index

This dataset was cleaned similar to the previous dataset where I viewed and visualized the dataset, renamed the variables and calculated the total number of NA values in the dataset.



## Merging the datasets together

I merged the datasets together using natural join and pipelined the output into dplyr select functions to remove the year wise data and only keep the columns that were previously defined by me. The variables were renamed.

The screenshot shows the RStudio interface with a script editor on the left and a console window on the right. The script editor contains R code that reads a CSV file, filters rows based on a date range, and calculates the maximum and minimum values for several variables across different categories. The console window shows the output of the code, including the dimensions of the data frame and the results of the calculations.

```

# RStudio
File Edit View Code Plots Session Build Debug Profile Tools Help

# Load the data
library(readr)
library(dplyr)
library(tibble)
library(stringr)

# Read the data
fdata_csv = read_csv("data/fdata.csv")

# Filter the data
fdata_csv = fdata_csv %>%
  filter(date >= "2019-01-01", date <= "2019-12-31")

# Calculate the maximum and minimum values for each variable
fdata_csv = fdata_csv %>%
  summarise(
    max_idx = which.max(fdata_csv[, "max_idx"]),
    min_idx = which.min(fdata_csv[, "min_idx"])
  )

# Print the results
print(fdata_csv)

```

The console window shows the following output:

```

# A tibble: 1 x 1
  max_idx
  <dbl>
1      1

```

The script editor also contains the following code:

```

# Read the data
fdata_csv = read_csv("data/fdata.csv")

# Filter the data
fdata_csv = fdata_csv %>%
  filter(date >= "2019-01-01", date <= "2019-12-31")

# Calculate the maximum and minimum values for each variable
fdata_csv = fdata_csv %>%
  summarise(
    max_idx = which.max(fdata_csv[, "max_idx"]),
    min_idx = which.min(fdata_csv[, "min_idx"])
  )

# Print the results
print(fdata_csv)

```

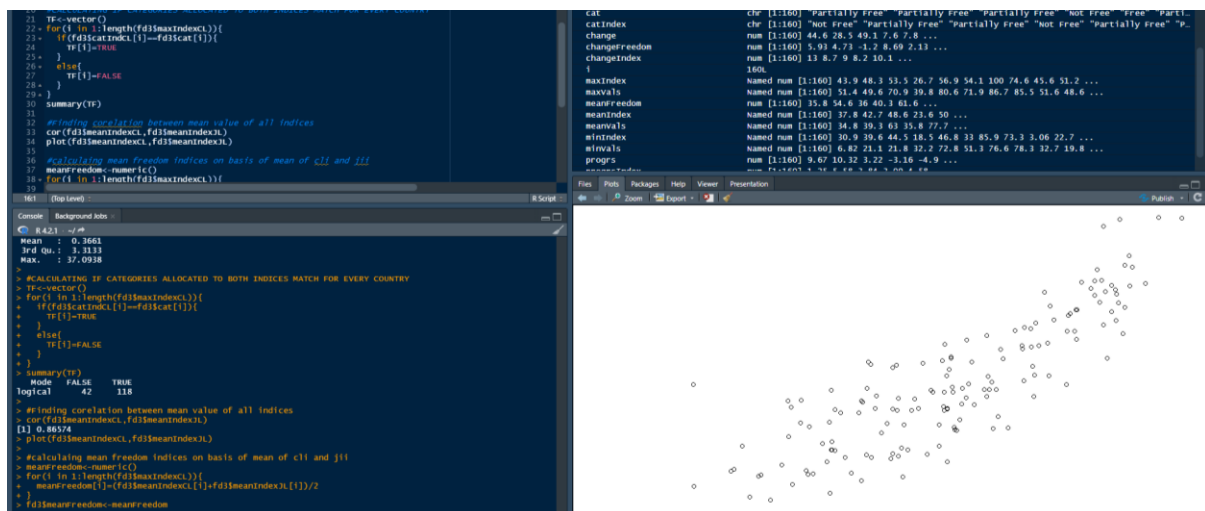
The console window shows the following output:

```

# A tibble: 1 x 1
  max_idx
  <dbl>
1      1

```

After this I took the freedom category columns from both the datasets and evaluated whether the data agreed on the category allotted to each country. Here are the results:

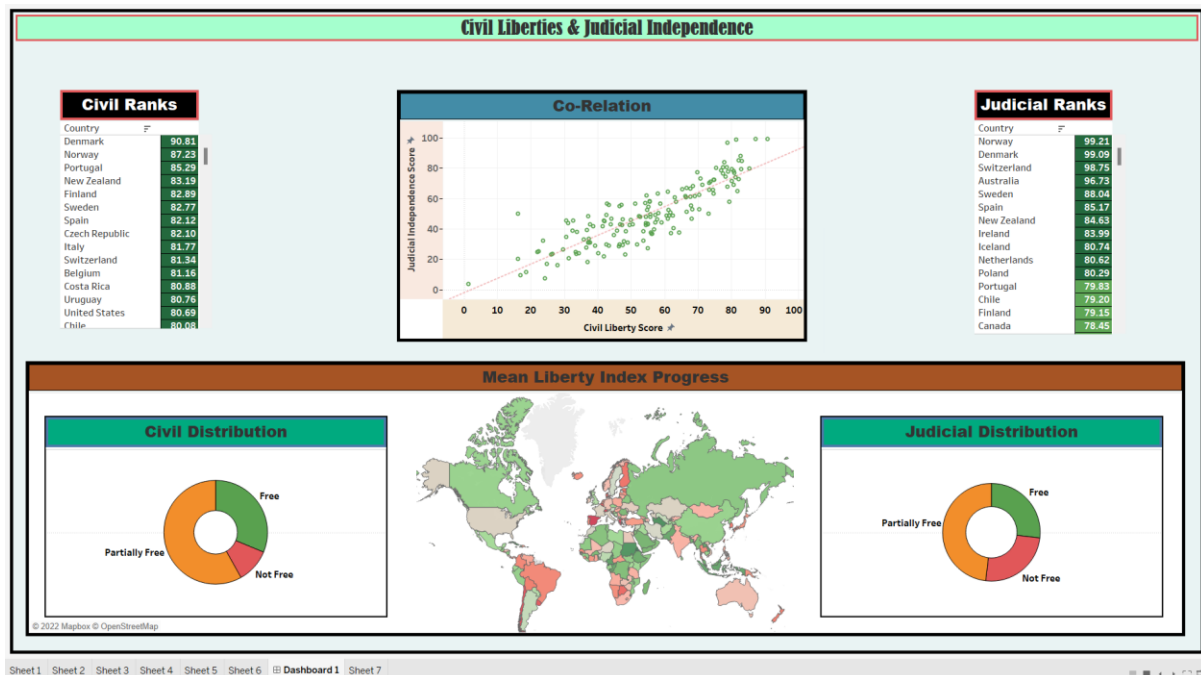


118 instances agree upon the categorisation while 42 disagree. Also the correlation between civil liberty and judicial independence is +0.86574 meaning a strong positive correlation .

Lastly, I have calculated the mean of civil and judicial liberties to find the mean change in index and the mean civil-judicial freedom value

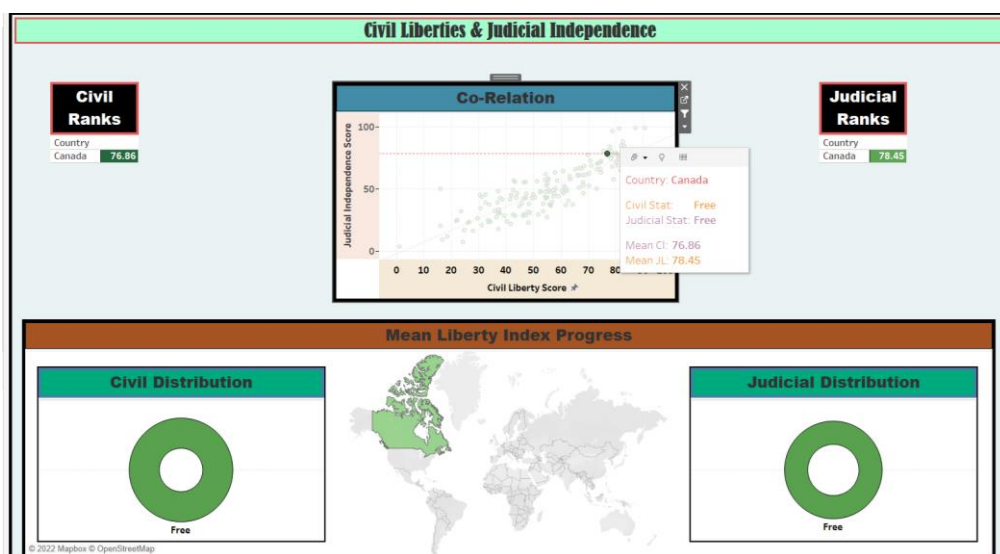
## Tableau Dashboard

Data visualisation is the final aspect of the data science project. To facilitate visualization in my project, I have made use of tableau software. Here is the created dashboard:

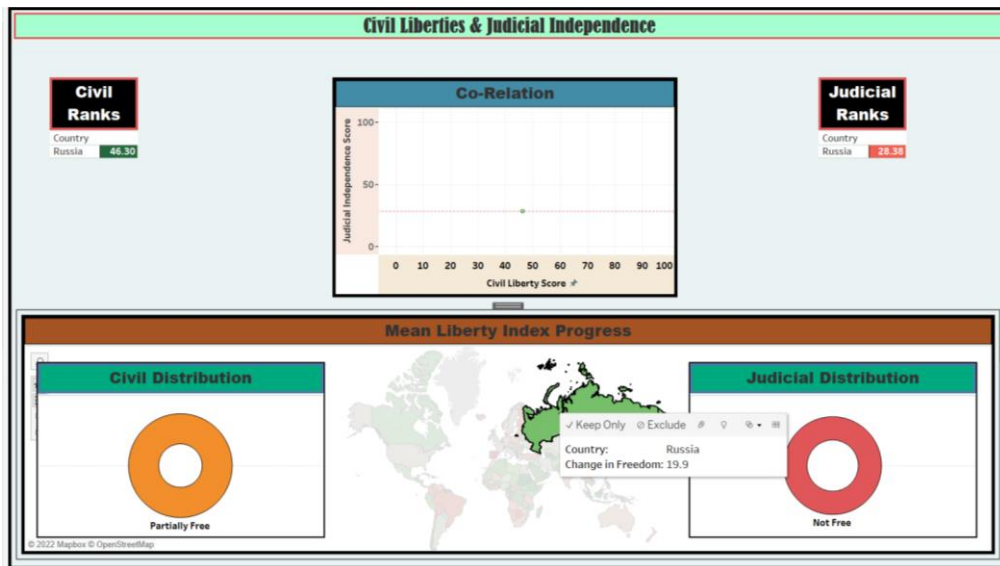


In the centre on the top row I have shown the graph of correlation between mean civil liberty and mean judicial index. To the left and right are highlight tables in which nations are ranked on the maximum score achieved by them in these rankings. The bottom row consists of a map which shows the progress of different nations in 2019 compared to the mean of all years from 2019 to 1990. Beside the map are 2 pie-charts depicting the proportion of categories the nations have been divided into.

Filters have been used to provide capability to select a country from the map or the graph to bring up all details relevant to that country.







Here are the component worksheets:

