

# A2 Assignment - Refugee Analysis

MBAN - Group 7

## Team Members:

Abi Joshua GEORG / Eri Yoshimoto / Hakeem GARCIA  
Nattida TAVAROJN / Neha NAGABHUSHAN / Weikang YANG

## Introduction

Refugee migration has been a significant global issue, with various countries experiencing fluctuations in the number of people seeking asylum. This report analyzes refugee data from multiple countries over a span of years to identify trends, patterns, and key insights. The study involves data cleaning, transformation, visualization, and statistical analysis to gain a deeper understanding of global refugee movements.

## Data Preparation and Cleaning

The dataset initially contained raw refugee statistics, which required significant cleaning before analysis.

```
library(tidyverse)
library(readr)
# read data from database
raw_df <- read_csv("data/A2_refugee_status.csv", col_types = cols(.default = "c"))
```

## Handling Missing and Inconsistent Values

The dataset contained placeholders for missing values, such as "D," "X," and "-", which were converted to "0." Numeric values were reformatted by removing commas and converting them into numerical data types.

```
# Set NULL value("D", "X", "-") as "0"
raw_df[raw_df == "D" | raw_df == "X" | raw_df == "-"] <- NA
# set the value column as numbers value and delete comma and transfer to numeric value
raw_df[, -1] <- lapply(raw_df[, -1], function(x) as.numeric(gsub(",", "", x)))
# set cleaned data frame as df and use df in later operation

# format "Congo, Democratic Republic" and "Congo, Republic" to ISO3 official name
raw_df <- raw_df %>%
  mutate(`Continent/Country of Nationality` = case_when(
    `Continent/Country of Nationality` == "Congo, Democratic Republic" ~ "Democratic Republic of the Congo",
    `Continent/Country of Nationality` == "Congo, Republic" ~ "Republic of the Congo",
    TRUE ~ `Continent/Country of Nationality`
  ))
```

## Standardizing Country Names

Countries with alternative naming conventions, such as “China, People’s Republic” and “Korea, North,” were renamed to “China” and “North Korea” for consistency.

```
# Replace the Countries' name with formal format
country_df <- raw_df %>%
  mutate(`Continent/Country of Nationality` = case_when(
    `Continent/Country of Nationality` == "China, People's Republic" ~ "China",
    `Continent/Country of Nationality` == "Korea, North" ~ "North Korea",
    TRUE ~ `Continent/Country of Nationality`
  ))
df <- country_df

# create the new dataframe 'contry_df', and remove the un-country row from the dataframe
# defined the name list that will be removed from the dataframe
non_countries <- c("Africa", "Asia", "Europe", "North America",
  "Oceania", "South America", "Unknown", "Other", "Total")

# use filter to get the row that will be delete
removed_countries <- df %>%
  filter(`Continent/Country of Nationality` %in% non_countries)

# create a new dataframe that only containt the 'countries'
country_df <- df %>%
  filter(!`Continent/Country of Nationality` %in% non_countries)

# print the row that be deleted to make sure all of them are 'non-countries'
print("delete rows with non-contry:")

## [1] "delete rows with non-contry:"

print(removed_countries$`Continent/Country of Nationality`)

## [1] "Africa"      "Asia"        "Europe"      "North America"
## [5] "Oceania"     "South America" "Unknown"     "Other"
## [9] "Unknown"     "Total"
```

```
# check the new data frame
head(country_df)

## # A tibble: 6 x 11
##   Continent/Country of~1 '2006' '2007' '2008' '2009' '2010' '2011' '2012' '2013'
##   <chr>                  <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 Afghanistan           651   441   576   349   515   428   481   661
## 2 Angola                 13     4    NA     8    NA    NA    NA     6
## 3 Armenia                87    29     9     4    NA    15     8     3
## 4 Azerbaijan            77    78    30    38    18    16    10     3
## 5 Belarus               350   219   111   146   103    66    83    10
## 6 Bhutan                  3    NA  5320 13452 12363 14999 15070 9134
## # i abbreviated name: 1: 'Continent/Country of Nationality'
## # i 2 more variables: '2014' <dbl>, '2015' <dbl>
```

```
continent_df <- df %>%
  filter(`Continent/Country of Nationality` %in% c("Africa", "Asia", "Europe",
                                                    "North America", "Oceania", "South America")) %>%
  pivot_longer(cols = -`Continent/Country of Nationality`,
               names_to = "Year", values_to = "Value") %>%
  group_by(`Continent/Country of Nationality`, Year) %>%
  summarise(Total_Refugees = sum(Value, na.rm = TRUE))
```

## 'summarise()' has grouped output by 'Continent/Country of Nationality'. You can  
## override using the '.groups' argument.

```
head(continent_df)
```

```
## # A tibble: 6 x 3
## # Groups:   Continent/Country of Nationality [1]
##   `Continent/Country of Nationality` Year   Total_Refugees
##   <chr>                                <chr>         <dbl>
## 1 Africa                                2006           18129
## 2 Africa                                2007           17486
## 3 Africa                                2008            8943
## 4 Africa                                2009            9678
## 5 Africa                                2010           13325
## 6 Africa                                2011            7693
```

## Trends in Refugee Migration

### Refugee Trends Across Presidential Terms (2006–2015)

The analysis of refugee data from 2006 to 2015 provides valuable insights into global migration patterns and highlights the influence of leadership periods on refugee numbers. By combining stacked bar charts with trend lines and presidential term overlays.

**Key Insights:** Presidential Term Influence:

The data is segmented into two distinct presidential periods:

- **George W. Bush (2006–2008):** During Bush's administration, the total number of Refugees was increased steadily, especially the refugees from Asia. But we can also see, the number of refugees from Africa and Europe was decreased significantly.
- **Barack Obama (2009–2015):** During Obama's administration, the total number of refugees was decreased between 2009 to 2011, but after that, the number was increased back to the steady level. Overall, the number of refugees was increased during Obama's administration with more relaxed immigration policies.
- The period from 2006 to 2015 saw the US experience some of the most substantial variation in the admissions of refugees, with changes attributed to internal policies, external crises, and two presidential administrations: **George W. Bush (2001-2009)** and **Barack Obama (2009-2017)**. Each administration had its own version of reality regarding refugee trends, dictated by both domestic and international developments.

## Stacked Bar Chart for Continental Trends:

The stacked bar chart provides a breakdown of refugee numbers by continent Over 2 Presidential Terms: - Asia and Africa consistently contributed the highest numbers of refugees, driven by ongoing conflicts and socio-political instability. - Europe and the Americas saw comparatively smaller contributions, reflecting differences in regional refugee patterns. - No refugees was came from Oceania according to the data.

```
library(ggplot2)
library(dplyr)
# Calculate the number of refugees for each state by year
# make sure the 'year' value is in Number data format
continent_df$Year <- as.numeric(continent_df$Year)

presidents <- data.frame(
  President = c("G. W. Bush", "Obama"),
  Start_Year = c(2005, 2009),
  End_Year = c(2009, 2016)
)

continent_df <- continent_df %>%
  mutate(
    Year = as.numeric(Year),
    Total_Refugees = as.numeric(Total_Refugees)
  )

yearly_refugees <- continent_df %>%
  group_by(Year) %>%
  summarise(Total_Refugees = sum(Total_Refugees, na.rm = TRUE))

ggplot() +
  geom_rect(data = presidents, aes(xmin = Start_Year, xmax = End_Year, ymin = 0, ymax = Inf, fill = President),
  geom_bar(data = continent_df, aes(x = Year, y = Total_Refugees, fill = `Continent/Country of Nationality`),
    stat = "identity", position = "stack", width = 0.6) +
  geom_line(data = yearly_refugees, aes(x = Year, y = Total_Refugees, group = 1),
    color = "#9f7abc", size = 2) +
  scale_fill_manual(
    name = "Legend",
    values = c(
      "G. W. Bush" = "red", # color the president
      "Obama" = "blue",
      "Africa" = "#a8b89a",
      "Asia" = "#a29bb3",
      "Europe" = "#d49391",
      "North America" = "#e5d5b2",
      "Oceania" = "#98c4c1",
      "South America" = "#8A3324"
    )
  ) +
  scale_y_continuous(breaks = scales::pretty_breaks(n = 10)) +
  scale_x_continuous(limits = c(2005, 2016), breaks = 2005:2016) +
  theme_minimal() +
  labs(
    title = "Refugee Numbers by Continent (2006-2015) with Presidential Terms",
    x = "Year",
```

```

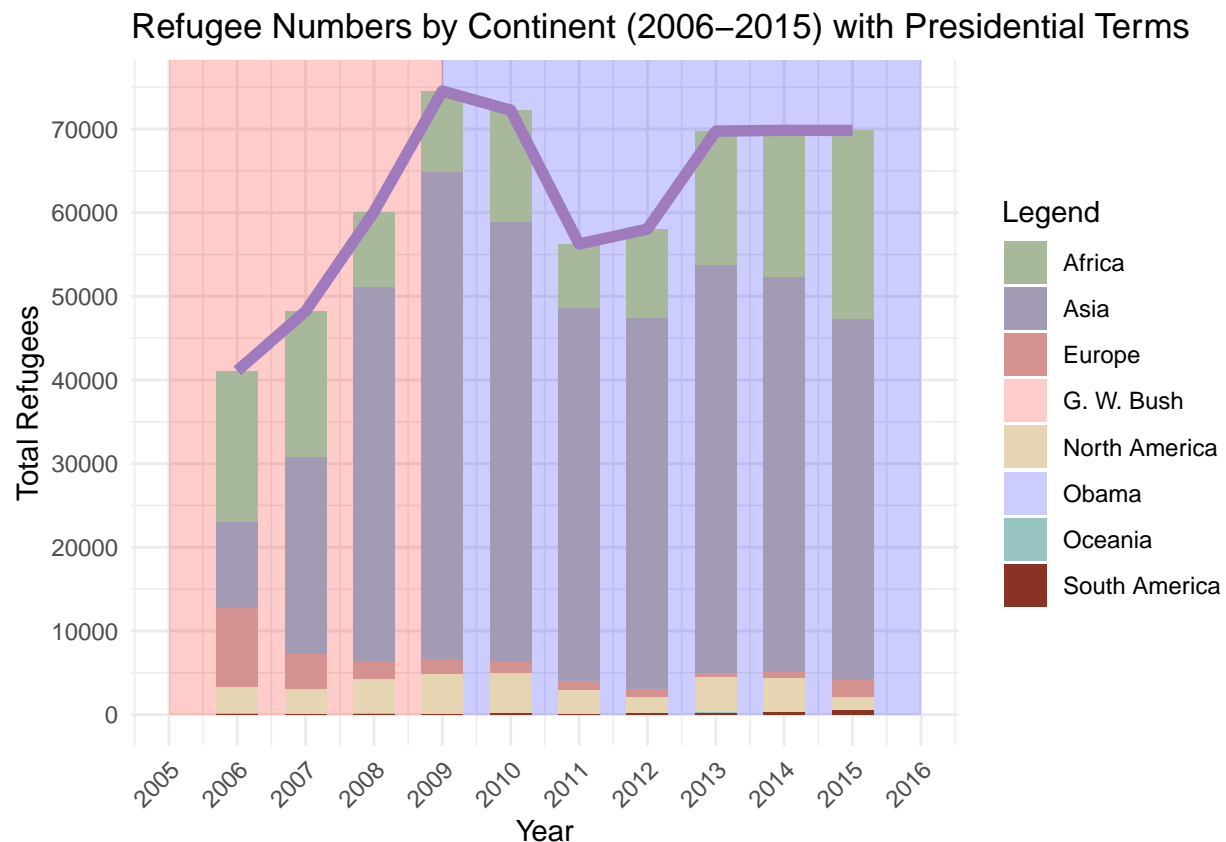
  y = "Total Refugees"
) +
theme(
  axis.text.x = element_text(angle = 45, hjust = 1),
  legend.position = "right"
)

```

```

## Warning: Using 'size' aesthetic for lines was deprecated in ggplot2 3.4.0.
## i Please use 'linewidth' instead.
## This warning is displayed once every 8 hours.
## Call 'lifecycle::last_lifecycle_warnings()' to see where this warning was
## generated.

```



```

dir.create("output_plot", showWarnings = FALSE)
ggsave("./output_plot/Refugee Numbers by Continent with Presidential Terms.svg", width = 10, height = 6, dpi = 300)
ggsave("./png/Refugee Numbers by Continent with Presidential Terms.png", width = 10, height = 6, dpi = 300)

```

```

# Load required libraries
library(ggplot2)
library(dplyr)
library(tidyr)
library(patchwork)
library(countrycode) # For country to continent mapping
theme_set(

```

```

theme_minimal() + # Base theme
  theme()
)
# Load the dataset
data <- read.csv("country_data.csv")

# Convert the dataset to long format
data_long <- data %>%
  pivot_longer(cols = -Country,
               names_to = "Year",
               values_to = "Value")

# Step 1: Map countries to continents using countrycode package
data_long <- data_long %>%
  mutate(
    Continent = countrycode(Country, origin = "country.name", destination = "continent")
  )

# Step 2: Manually classify certain countries as North America
data_long <- data_long %>%
  mutate(
    Continent = ifelse(Country %in% c("Haiti", "Cuba", "Honduras"), "North America", Continent)
  )

# Step 3: Exclude countries from Oceania and South America
data_long <- data_long %>%
  filter(!Continent %in% c("Oceania", "South America"))

# Step 4: Define a function to get top 5 countries by total refugees for each continent
top_5_per_continent <- function(data, continent, rank) {
  data %>%
    filter(Continent == continent) %>%
    group_by(Country) %>%
    summarise(Total_Refugees = sum(Value, na.rm = TRUE)) %>%
    arrange(desc(Total_Refugees)) %>%
    slice_head(n = rank) %>%
    inner_join(data, by = "Country")
}

# Step 5: Get the top 5 countries for each continent
asia_data <- top_5_per_continent(data_long, "Asia", 5)
africa_data <- top_5_per_continent(data_long, "Africa", 5)
north_america_data <- top_5_per_continent(data_long, "North America", 3) # Now includes manually class
europe_data <- top_5_per_continent(data_long, "Europe", 5)

# Step 6: Define a function to create line plots for each continent
plot_continent <- function(data, title) {
  ggplot(data, aes(x = Year, y = Value, color = Country, group = Country)) +
    geom_line(size = 1) +
    theme_minimal() +
    labs(title = title, x = "Year", y = "Number of Refugees", color = "Countries") +
    theme(
      legend.position = "bottom",

```

```

    legend.text = element_text(size = 8)
  ) +
  guides(color = guide_legend(ncol = 2)) +
  scale_y_continuous(
    limits = c(0, max(data$Value, na.rm = TRUE) * 1.2), # Extend Y-axis by 20%
    breaks = scales::pretty_breaks(n = 10) # Add more breaks for readability
  )
}

# Step 7: Create individual line plots for each continent
plot_asia <- plot_continent(asia_data, "Top 5 Refugee Countries in Asia")
plot_africa <- plot_continent(africa_data, "Top 5 Refugee Countries in Africa")
plot_north_america <- plot_continent(north_america_data, "Top 5 Refugee Countries in North America")
plot_europe <- plot_continent(europe_data, "Top 5 Refugee Countries in Europe")

# Define custom colors for countries
custom_colors <- c(
  # Asia
  "Bhutan" = "#a8b89a", "Iraq" = "#a29bb3", "Burma" = "#d49391",
  "Vietnam" = "#e5d5b2", "Iran" = "#8A3324",

  # Africa
  "Burundi" = "#a8b89a", "Somalia" = "#a29bb3", "Congo" = "#d49391",
  "Sudan" = "#e5d5b2", "Eritrea" = "#8A3324",

  # North America
  "Cuba" = "#a8b89a", "Honduras" = "#a29bb3", "Haiti" = "#d49391",

  # Europe
  "Belarus" = "#a8b89a", "Russia" = "#a29bb3", "Latvia" = "#d49391",
  "Ukraine" = "#e5d5b2", "Moldova" = "#8A3324"
)

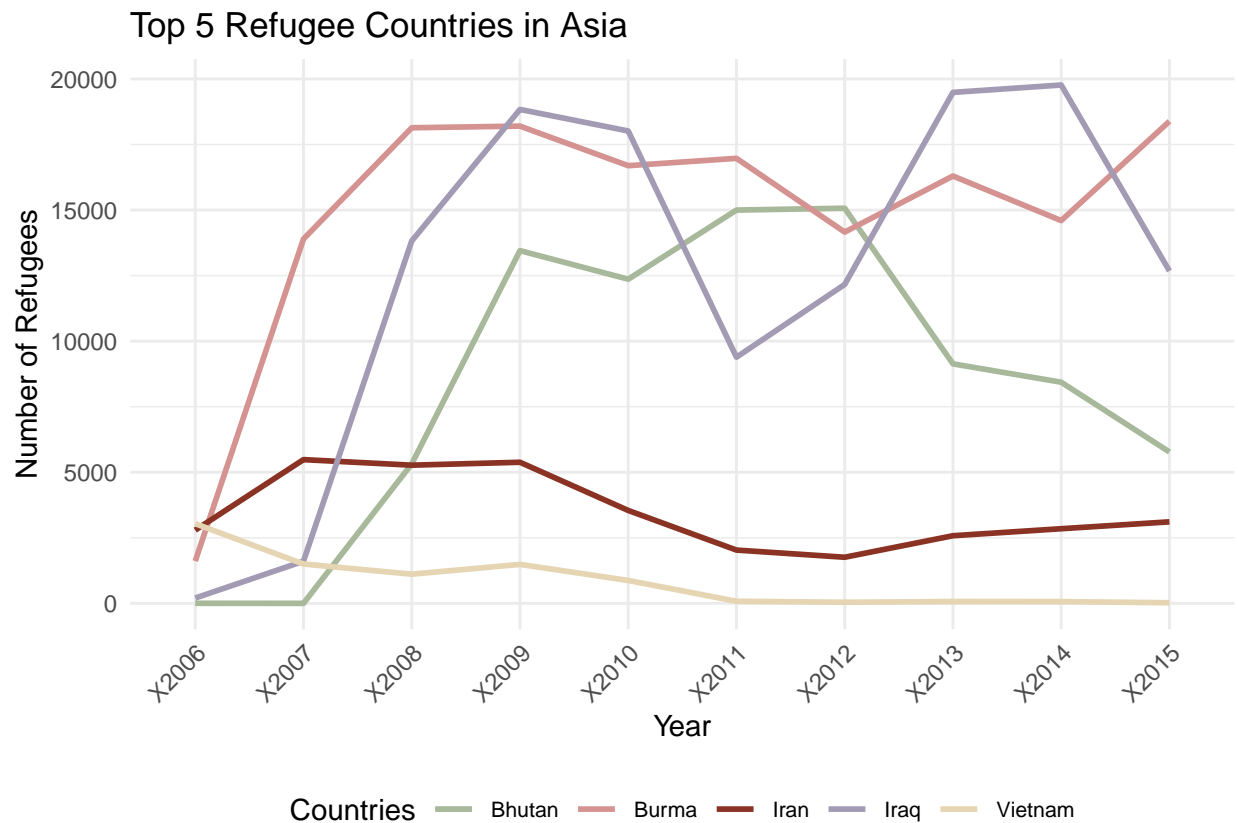
# Adjust plot dimensions for a single column layout
options(repr.plot.width = 15, repr.plot.height = 30)

# Modify the plot_continent function to include custom line colors
plot_continent <- function(data, title) {
  ggplot(data, aes(x = Year, y = Value, color = Country, group = Country)) +
    geom_line(size = 1) +
    theme_minimal() +
    labs(title = title, x = "Year", y = "Number of Refugees", color = "Countries") +
    theme(
      legend.position = "bottom",
      legend.text = element_text(size = 8),
      axis.text.x = element_text(angle = 45, hjust = 1)
    ) +
    scale_y_continuous(
      breaks = scales::pretty_breaks(n = 5) # Adjust Y-axis breaks dynamically
    ) +
    scale_color_manual(values = custom_colors) # Apply custom colors
}

```

```
# Create updated plots for each continent
plot_asia <- plot_continent(asia_data, "Top 5 Refugee Countries in Asia")
plot_africa <- plot_continent(africa_data, "Top 5 Refugee Countries in Africa")
plot_north_america <- plot_continent(north_america_data, "Top 5 Refugee Countries in North America")
plot_europe <- plot_continent(europe_data, "Top 5 Refugee Countries in Europe")

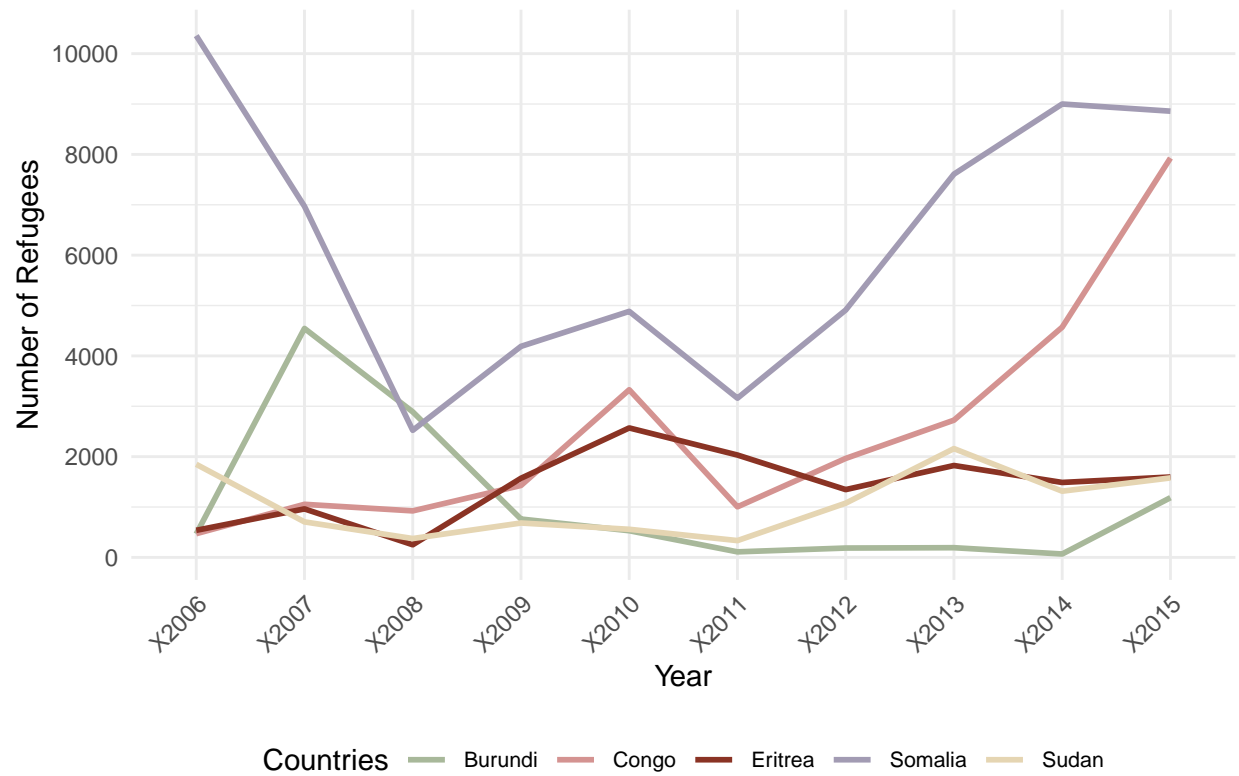
# Combine plots with a single column layout
plot_asia
```



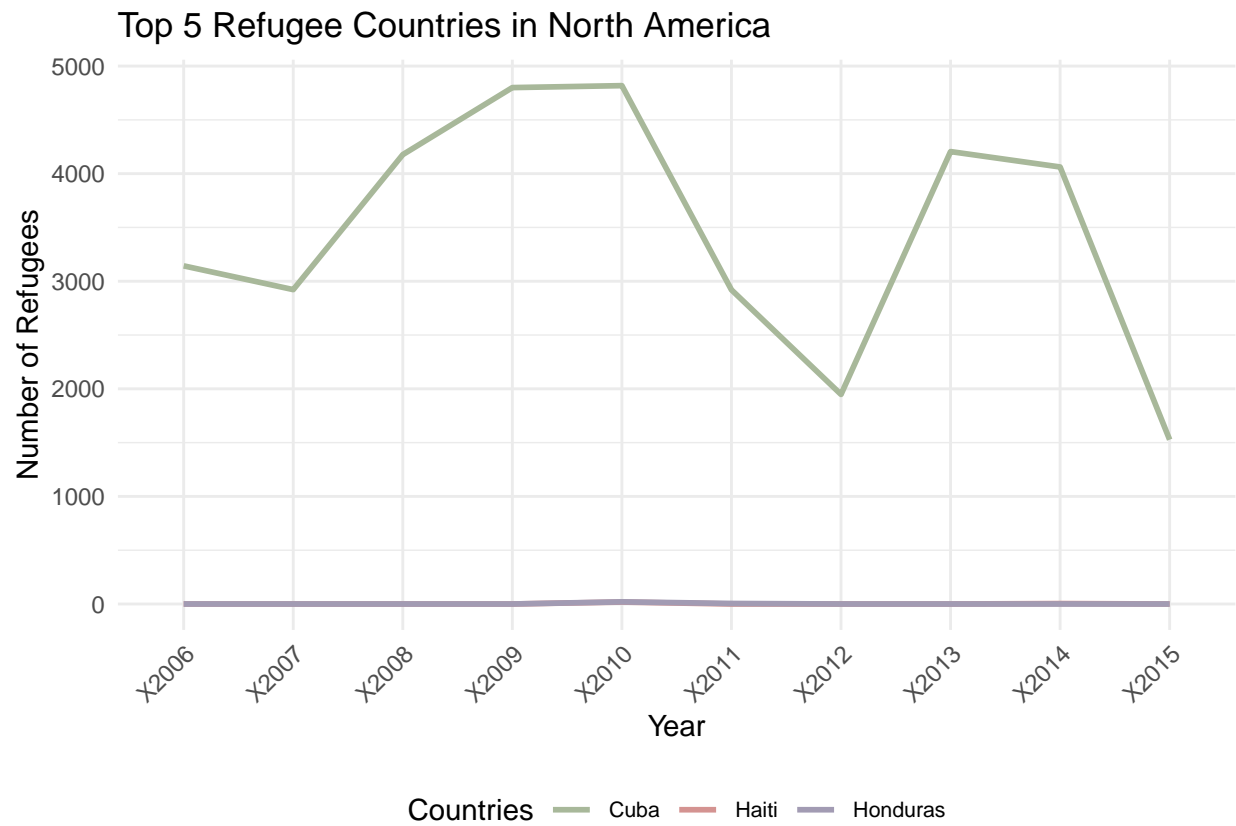
```
plot_africa
```



Top 5 Refugee Countries in Africa

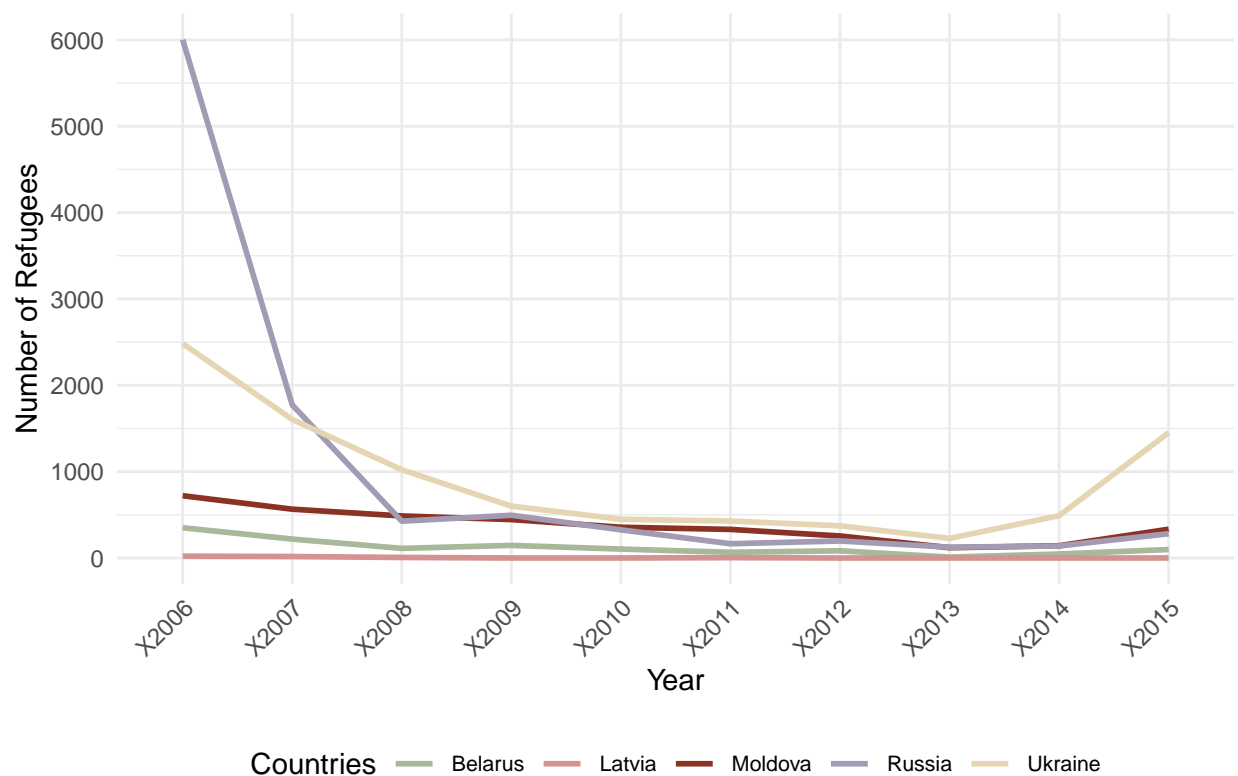


plot\_north\_america



plot\_europe

### Top 5 Refugee Countries in Europe



```
continent_names <- c("asia", "africa", "north_america", "europe")
plots <- list(plot_asia, plot_africa, plot_north_america, plot_europe)
```

```
lapply(seq_along(continent_names), function(i) {
  ggsave(paste0("output_plot/plot_", continent_names[i], ".svg"),
    plots[[i]], width = 10, height = 6, dpi = 300)
})
```

```
## [[1]]
## [1] "output_plot/plot_asia.svg"
##
## [[2]]
## [1] "output_plot/plot_africa.svg"
##
## [[3]]
## [1] "output_plot/plot_north_america.svg"
##
## [[4]]
## [1] "output_plot/plot_europe.svg"
```

```
lapply(seq_along(continent_names), function(i) {
  ggsave(paste0("png/plot_", continent_names[i], ".png"),
    plots[[i]], width = 10, height = 6, dpi = 300)
})
```

```
## [[1]]
## [1] "png/plot_asia.png"
##
## [[2]]
## [1] "png/plot_africa.png"
##
## [[3]]
## [1] "png/plot_north_america.png"
##
## [[4]]
## [1] "png/plot_europe.png"
```

```
# Load required libraries
library(dplyr)
library(tidyr)

# Load the dataset
data <- read.csv("country_data.csv")
# Convert the dataset to long format
data_long <- data %>%
  pivot_longer(cols = -Country, names_to = "Year", values_to = "Refugees")
# Step 1: Find the top 10 countries by refugees for each year
top10_per_year <- data_long %>%
  group_by(Year) %>%
  slice_max(order_by = Refugees, n = 3) %>%
  ungroup()

# Step 2: Extract all years' data for these top countries
top_countries <- unique(top10_per_year$Country) # Get unique top countries
all_data_top_countries <- data_long %>%
  filter(Country %in% top_countries) # Filter all years for top countries

# View the data
print(all_data_top_countries)
```

```
## # A tibble: 70 x 3
##   Country Year  Refugees
##   <chr>   <chr>    <int>
## 1 Bhutan  X2006         3
## 2 Bhutan  X2007         0
## 3 Bhutan  X2008        5320
## 4 Bhutan  X2009       13452
## 5 Bhutan  X2010       12363
## 6 Bhutan  X2011       14999
## 7 Bhutan  X2012       15070
## 8 Bhutan  X2013        9134
## 9 Bhutan  X2014        8434
## 10 Bhutan X2015        5775
## # i 60 more rows
```

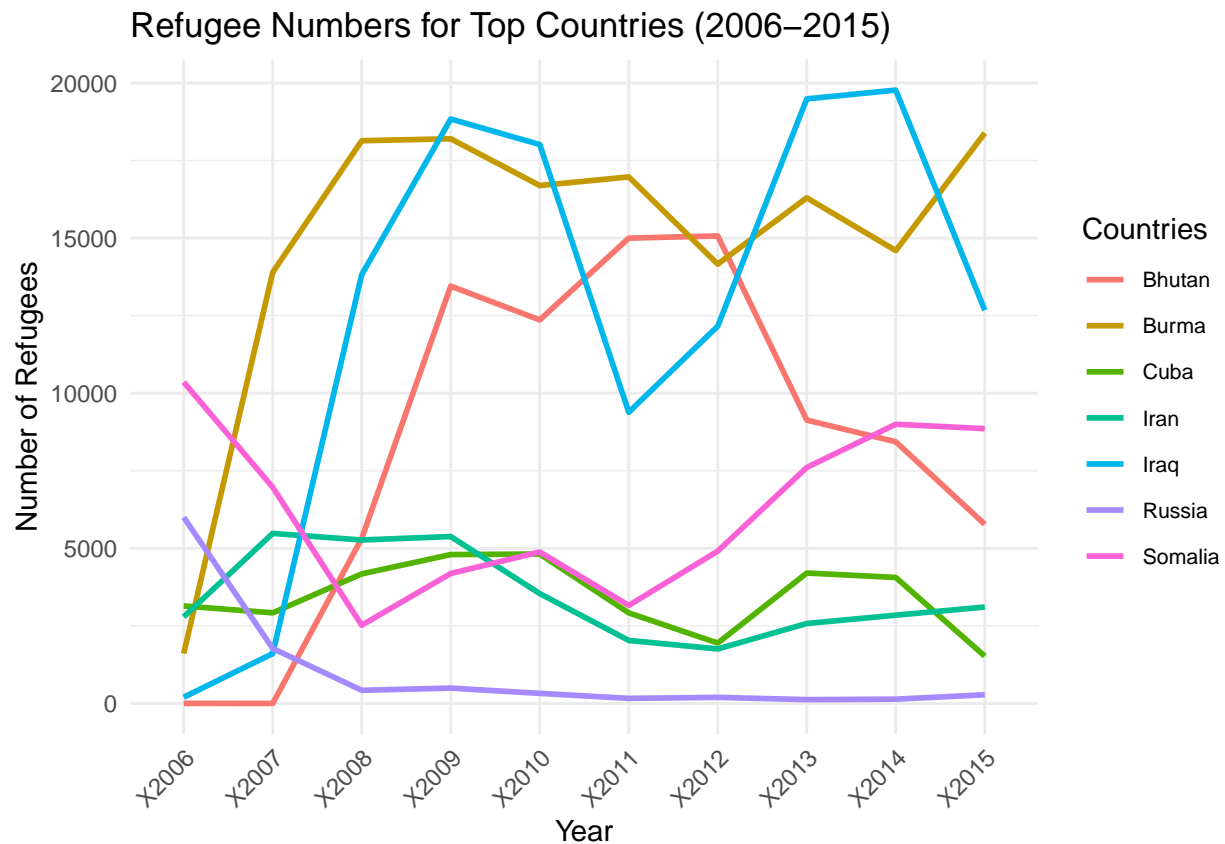
```
# Save to a new CSV file if needed
write.csv(all_data_top_countries, "all_top_countries_refugees.csv", row.names = FALSE)
```

```

# Ensure data is loaded into all_data_top_countries
# If not already loaded, reload the dataset
all_data_top_countries <- read.csv("all_top_countries_refugees.csv")

# Create a line plot for all top countries over the years
p00 <- ggplot(all_data_top_countries, aes(x = Year, y = Refugees, color = Country, group = Country)) +
  geom_line(size = 1) + # Add lines for each country
  theme_minimal() + # Use a minimal theme
  labs(
    title = "Refugee Numbers for Top Countries (2006-2015)",
    x = "Year",
    y = "Number of Refugees",
    color = "Countries"
  ) +
  theme(
    legend.position = "right", # Move legend to the bottom
    legend.text = element_text(size = 8), # Adjust legend text size
    axis.text.x = element_text(angle = 45, hjust = 1) # Rotate x-axis labels
  )
ggsave("output_plot/Refugee Numbers for Top Countries 2006-2015.svg", p00, width = 8, height = 6, dpi =
ggsave("png/Refugee Numbers for Top Countries 2006-2015.png", p00, width = 8, height = 6, dpi = 300)
p00

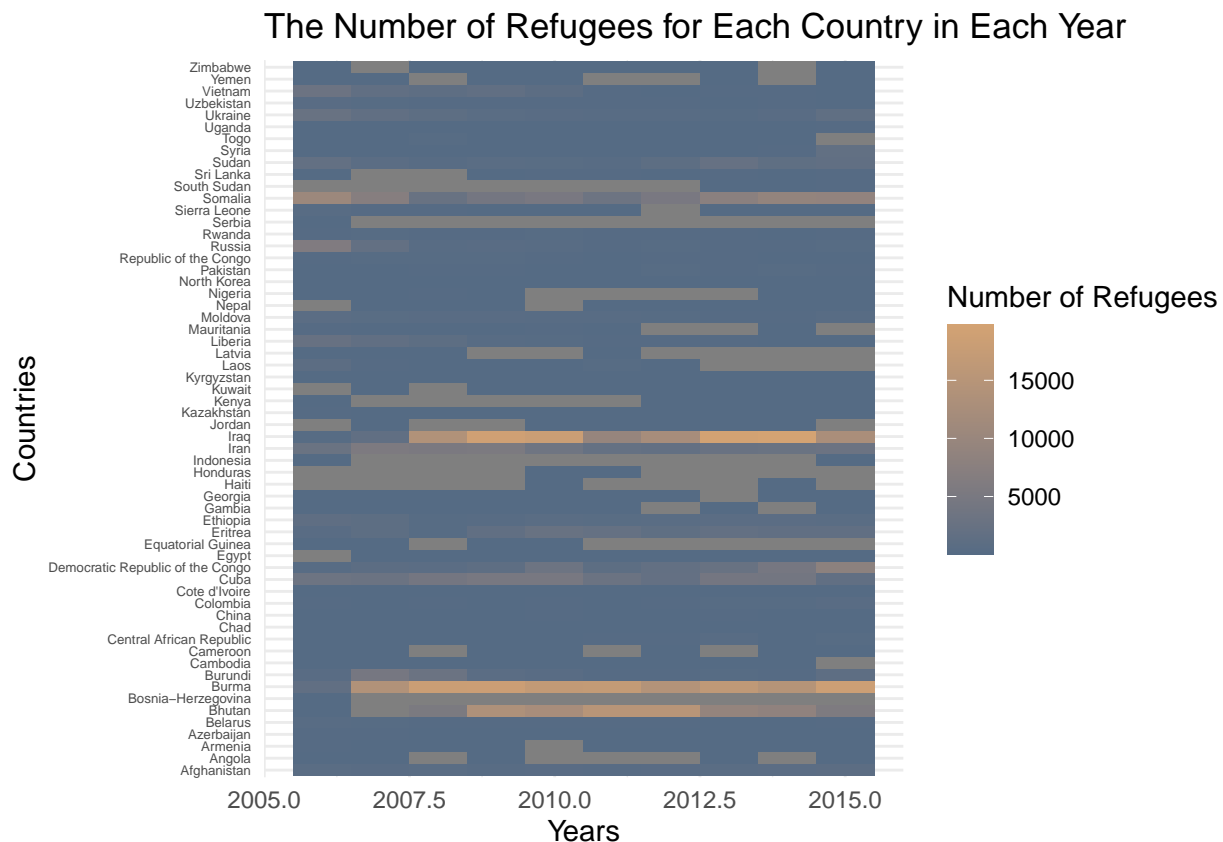
```



### 3. Heatmap of Refugee Numbers

```
country_long <- country_df %>%
  pivot_longer(cols = -`Continent/Country of Nationality`,
               names_to = "Year", values_to = "Value") %>%
  mutate(
    Year = as.numeric(Year),
    Value = as.numeric(Value)
  )

p0 <- ggplot(country_long, aes(x = as.numeric(Year), y = `Continent/Country of Nationality`, fill = Value)) +
  geom_tile() +
  scale_fill_gradient(low = "#556b84", high = "#d4a373") +
  theme_minimal() +
  labs(title = "The Number of Refugees for Each Country in Each Year",
       x = "Years",
       y = "Countries",
       fill = "Number of Refugees") +
  theme(legend.position = "right",
        axis.text.y = element_text(size = 5))
ggsave("output_plot/The Number of Refugees for Each Country in Each Year.svg",p0, width = 10, height = 6, dpi = 300)
ggsave("png/The Number of Refugees for Each Country in Each Year.png",p0, width = 10, height = 6, dpi = 300)
p0
```



## The World Map with Refugees number

```
library(gganimate)
library(sf)
library(rnaturalearth)
library(rnaturalearthdata)
library(gifski)
library(transformr)

country_long <- country_df %>%
  pivot_longer(cols = -`Continent/Country of Nationality`,
               names_to = "Year", values_to = "Value")
country_long$Year <- as.numeric(country_long$Year)

# load the world map
world_map <- ne_countries(scale = "medium", returnclass = "sf")

# adjust the name of countries
country_long <- country_long %>%
  rename(country = `Continent/Country of Nationality`)

# combine the map data and refugees data
map_data <- world_map %>%
  left_join(country_long, by = c("name" = "country"))

country_long$Year <- as.numeric(country_long$Year)

#
plot_yearly_maps <- function(year) {
  yearly_data <- map_data %>% filter(Year == year)

  ggplot() +
    geom_sf(data = world_map, fill = "gray90", color = "white") +
    geom_sf(data = yearly_data, aes(fill = Value), color = "black") +
    scale_fill_gradient(low = "#556b84", high = "#d4a373", na.value = "gray90") +
    theme_minimal() +
    labs(title = paste("Global Refugee Map - Year", year), fill = "Number of Refugees")
}

p2006 <- plot_yearly_maps(2006)
p2007 <- plot_yearly_maps(2007)
p2008 <- plot_yearly_maps(2008)
p2009 <- plot_yearly_maps(2009)
p2010 <- plot_yearly_maps(2010)
p2011 <- plot_yearly_maps(2011)
p2012 <- plot_yearly_maps(2012)
p2013 <- plot_yearly_maps(2013)
p2014 <- plot_yearly_maps(2014)
p2015 <- plot_yearly_maps(2015)

years <- 2006:2015
```

```
plots <- list(p2006, p2007, p2008, p2009, p2010, p2011, p2012, p2013, p2014, p2015)

lapply(seq_along(years), function(i) {
  ggsave(paste0("output_plot/p", years[i], ".svg"), plots[[i]], width = 8, height = 6, dpi = 300)
})
```

```
## [[1]]
## [1] "output_plot/p2006.svg"
##
## [[2]]
## [1] "output_plot/p2007.svg"
##
## [[3]]
## [1] "output_plot/p2008.svg"
##
## [[4]]
## [1] "output_plot/p2009.svg"
##
## [[5]]
## [1] "output_plot/p2010.svg"
##
## [[6]]
## [1] "output_plot/p2011.svg"
##
## [[7]]
## [1] "output_plot/p2012.svg"
##
## [[8]]
## [1] "output_plot/p2013.svg"
##
## [[9]]
## [1] "output_plot/p2014.svg"
##
## [[10]]
## [1] "output_plot/p2015.svg"
```

```
lapply(seq_along(years), function(i) {
  ggsave(paste0("png/p", years[i], ".png"), plots[[i]], width = 8, height = 6, dpi = 300)
})
```

```
## [[1]]
## [1] "png/p2006.png"
##
## [[2]]
## [1] "png/p2007.png"
##
## [[3]]
## [1] "png/p2008.png"
##
## [[4]]
## [1] "png/p2009.png"
##
## [[5]]
```



```
## [1] "png/p2010.png"
##
## [[6]]
## [1] "png/p2011.png"
##
## [[7]]
## [1] "png/p2012.png"
##
## [[8]]
## [1] "png/p2013.png"
##
## [[9]]
## [1] "png/p2014.png"
##
## [[10]]
## [1] "png/p2015.png"
```

**below code is for generate the GIF** library(gganimate)

```
p <- ggplot() + geom_sf(data = world_map, fill = "gray90", color = "white") + # the world map
geom_sf(data = map_data, aes(fill = Value), color = "black") + # data refugees scale_fill_gradient(low
= "blue", high = "yellow", na.value = "gray90") + theme_minimal() + labs(title = "Global Refugees data
Map (Year: {frame_time})", fill = "Number of Refugees", x = "", y = "") + transition_time(Year) + #
change the plot by year ease_aes('linear') # change smoothly
```

generate the GIF

```
anim <- animate(p, duration = 10, fps = 40, width = 800, height = 500, renderer = gifski_renderer())
anim_save("refugees_map_smooth.gif", animation = anim)
```

## Conclusion

This analysis provides a comprehensive look at global refugee movements, highlighting key trends, high-risk regions, and changes over time. The combination of data cleaning, visualization, and statistical analysis allows for an in-depth understanding of the factors driving refugee migration, which can support future policy-making and humanitarian efforts.