

VIRGINIA COMMONWEALTH UNIVERSITY



STATISTICAL ANALYSIS & MODELING

Visualization - Perceptual Mapping for Business

Kirthan Shaker Iyengar
V01108265

Date of Submission: 15/07/2024

CONTENTS

Content:	Page no:
INTRODUCTION	3
OBJECTIVE	3
BUSINESS SIGNIFICANC	3-4
RESULTS AND INTERPRETATIONS	5-13
CODES	13-17

Visualization - Perceptual Mapping for Business

INTRODUCTION

This study applies Visualization techniques, specifically Perceptual Mapping, to the assigned business dataset. The goal is to visualize the underlying relationships and patterns, aiding in strategic decision-making processes.

Perceptual Mapping is a graphical technique used to display the perceptions of customers or respondents. In this analysis, we will preprocess the dataset, handle missing values, and ensure that the visualization assumptions are met.

We will explore the relationships between various business variables, identifying significant dimensions that affect business outcomes. The perceptual map will be utilized to assess the relative positioning of different items or attributes in the minds of respondents. Furthermore, we will compare different visualization methods to determine the most effective way of presenting the data.

By systematically visualizing the data using these methods, we aim to provide comprehensive insights that can guide business strategies and marketing efforts. The findings will highlight the strengths and weaknesses of each visualization approach and contribute to a better understanding of the market structure.

OBJECTIVES

1. Plot a histogram (to show the distribution of total consumption across different districts) and a barplot (To visualize consumption per district with district names) of the data in Assignment A1 to indicate the consumption district-wise for the state assigned to you.
2. Plot {'any variable of your choice'} on the Karnataka (or the state assigned to you) state map using NSSO68.csv data

BUSINESS SIGNIFICANCE

Using these advanced statistical models like Perceptual Mapping for Business can significantly enhance business decision-making. For instance, Perceptual Mapping for Business predict customer churn, enabling targeted retention strategies to maintain revenue streams. Probit regression can analyze consumer behaviour, helping businesses tailor product offerings and marketing campaigns to specific demographics, such as dietary preferences. Tobit regression is crucial for assessing loan default risks, allowing financial institutions to implement risk-based pricing and targeted interventions. These models provide actionable insights, optimizing strategies across customer retention, product development, and risk management.

RESULTS AND INTERPRETATION

1. Plot a histogram (to show the distribution of total consumption across different districts) and a barplot (To visualize consumption per district with district names) of the data in Assignment A1 to indicate the consumption district-wise for the state assigned to you.

#Identifying analysis

Code and Result:

```
KE['total_consumption'] = KE[['ricepds_v', 'Wheatpds_q', 'chicken_q', 'pulsep_q',  
'wheatos_q']].sum(axis=1)
```

```
def summarize_consumption(group_col):
```

```
    summary = KE.groupby(group_col)['total_consumption'].sum().reset_index()
```

```
    summary.sort_values(by='total_consumption', ascending=False, inplace=True)
```

```
    return summary
```

```
district_summary = summarize_consumption('District')
```

```
region_summary = summarize_consumption('Region')
```

```
print("Top Consuming Districts:")
```

```
print(district_summary.head(4))
```

```
print("Region Consumption Summary:")
```

```
print(region_summary)
```

```
district = {'1': 'Kasaragod',
```

```
           '2': 'Kannur',
```

```
           '3': 'Wayand',
```

```
           '4': 'Kozhikode',
```

```
           '5': 'Malappuram',
```

```
           '6': 'Palakkad',
```

```
           '7': 'Thrissur',
```

```
           '8': 'Eranakulam',
```

```
           '9': 'Idukki',
```

```
           '10': 'Kottayam',
```

```
           '11': 'Alappuzha',
```

```
           '12': 'Pathanamthitta',
```

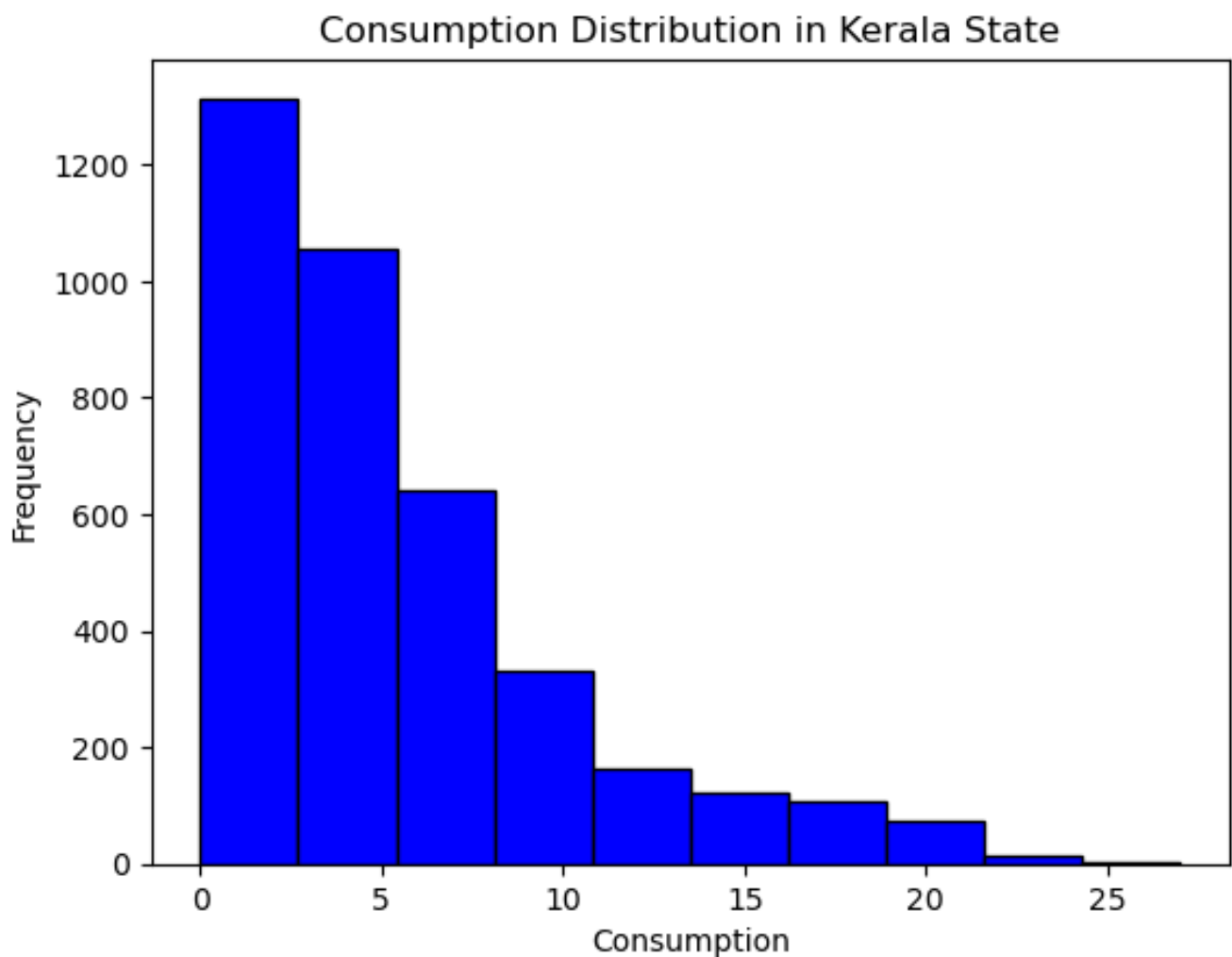
```
           '13': 'Kollam',
```

```
           '14': 'Thiruvananthapuram',
```

```
}
```

```
sector = {
    '2': 'URBAN',
    '1': 'RURAL'
}
plt.hist(KE['total_consumption'], bins=10, color='blue', edgecolor='black')
plt.xlabel("Consumption")
plt.ylabel("Frequency")
plt.title("Consumption Distribution in Kerala State")
plt.show()
```

Result:



Interpretation:

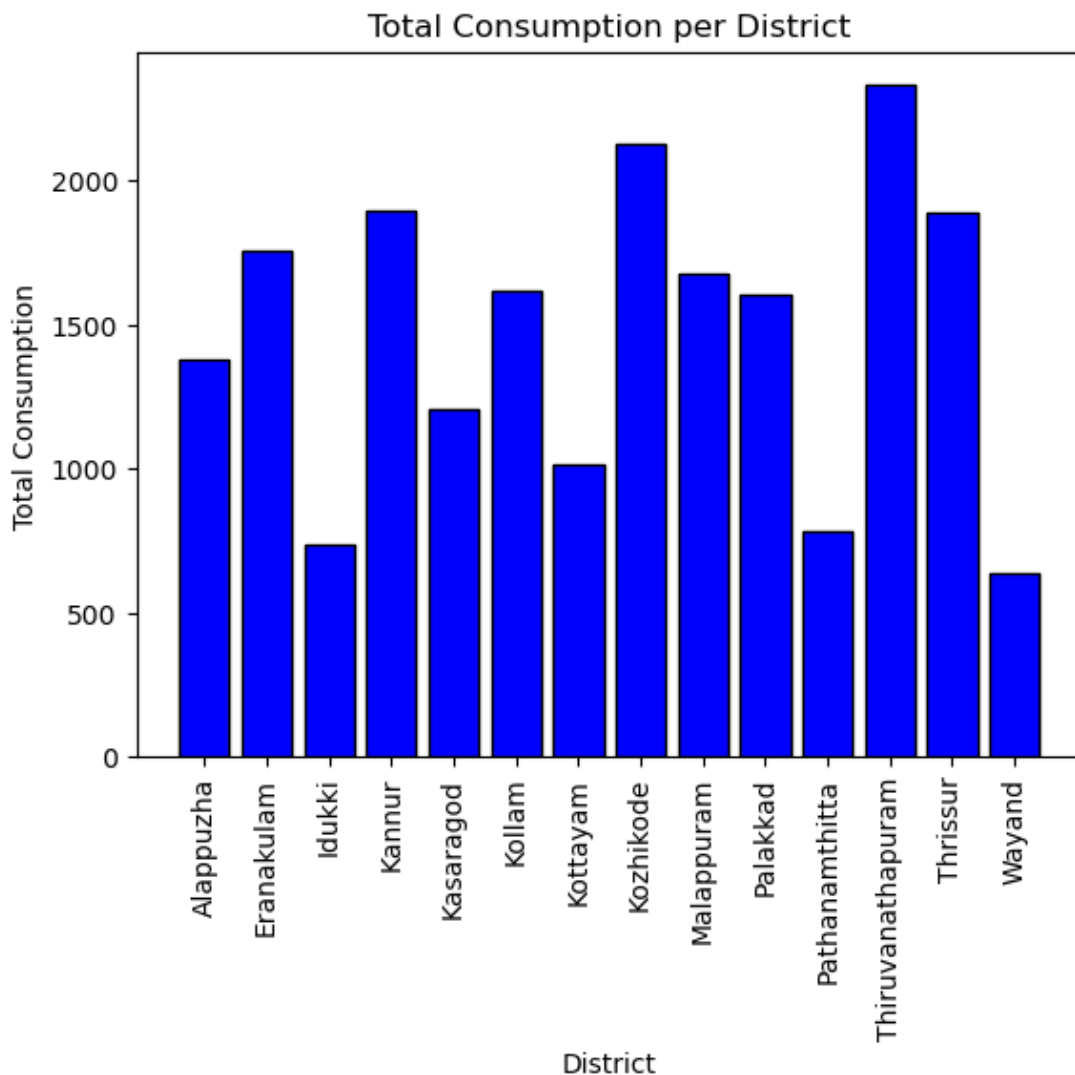
The histogram shows the distribution of total consumption across different districts in Kerala State. Here is an interpretation of the histogram:

1. **Skewness:** The distribution is right-skewed, indicating that most districts have relatively low consumption, with a few districts having higher consumption values.
2. **Frequency:**

- The highest frequency of consumption values falls within the 0-5 range, with more than 1200 districts.
 - The number of districts with consumption values between 5-10 is lower but still significant, with around 800 districts.
 - As the consumption value increases, the frequency of districts decreases. For instance, the range of 10-15 consumption units has fewer than 400 districts.
3. **Outliers:** There are very few districts with consumption values exceeding 20. These can be considered outliers.
 4. **Overall Consumption Trend:** The general trend indicates that a majority of the districts have lower consumption values, while a small number of districts have higher consumption values.

Code :

```
plt.bar(KE_consumption['District'], KE_consumption['total_consumption'], color='blue',
edgecolor='black')
plt.xlabel("District")
plt.ylabel("Total Consumption")
plt.title("Total Consumption per District")
plt.xticks(rotation=90) # Rotate district names for better visibility
plt.show()
```



Interpretation:

- **Districts with Highest Consumption:**

- **Thiruvananthapuram:** This district has the highest total consumption, with a value exceeding 2000.
- **Kozhikode:** Following closely, this district also has a high consumption value, slightly less than Thiruvananthapuram.

- **Districts with Moderate Consumption:**

- **Ernakulam, Malappuram, and Thrissur:** These districts have moderate consumption values, ranging between 1500 and 2000.
- **Kannur, Palakkad, Pathanamthitta, and Kottayam:** These districts have similar moderate values, ranging between 1000 and 1500.

- **Districts with Lower Consumption:**

- **Alappuzha, Idukki, Kasaragod, Kollam, and Wayanad:** These districts have relatively lower consumption values, ranging from around 500 to 1000.

- **Insights:**

- The consumption is not evenly distributed across the districts, with some districts like Thiruvananthapuram and Kozhikode showing significantly higher values.
- There is a noticeable variation in consumption values, indicating potential differences in population size, industrial activity, or other factors contributing to consumption levels.

- **Visualization:**

- The bar chart effectively uses color and edge color to distinguish the bars, and the x-axis labels are rotated for better readability, ensuring that the district names are easily readable.

2. Ploted Consumption on the Kerala state map using NSSO68.csv data

Code and Result

```
In [98]: data_map = gpd.read_file("/Users/kirthanshaker/Desktop/SCMA 631 Data Files /KERALA_DISTRICTS (1).geojson")

In [99]: print(data_map.columns)
print(KE_consumption.columns)
Index(['dtname', 'stname', 'stcode11', 'dtcode11', 'year_stat', 'Shape_Length',
       'Shape_Area', 'OBJECTID', 'test', 'Dist_LGD', 'State_LGD', 'geometry'],
      dtype='object')
Index(['District', 'total_consumption'], dtype='object')

In [100]: data_map['District'] = KE_consumption['District']

In [101]: data_map_data = data_map.merge(KE_consumption, left_on='dtname', right_on='District')

In [102]: print(data_map.columns)
Index(['dtname', 'stname', 'stcode11', 'dtcode11', 'year_stat', 'Shape_Length',
       'Shape_Area', 'OBJECTID', 'test', 'Dist_LGD', 'State_LGD', 'geometry',
       'District'],
      dtype='object')

In [103]: import geopandas as gpd
import pandas as pd
import matplotlib.pyplot as plt

In [104]: data_map = gpd.read_file("/Users/kirthanshaker/Desktop/SCMA 631 Data Files /KERALA_DISTRICTS (1).geojson")

In [105]: data_map = data_map.rename(columns={'dtname': 'District'})

In [106]: display(data_map.rename())

In [108]: KE_consumption = pd.read_csv("/Users/kirthanshaker/Desktop/SCMA 631 Data Files /KERALA_DISTRICTS (1).geojson", low_m

In [109]: KE_consumption = KE.groupby('District')['total_consumption'].sum().reset_index()

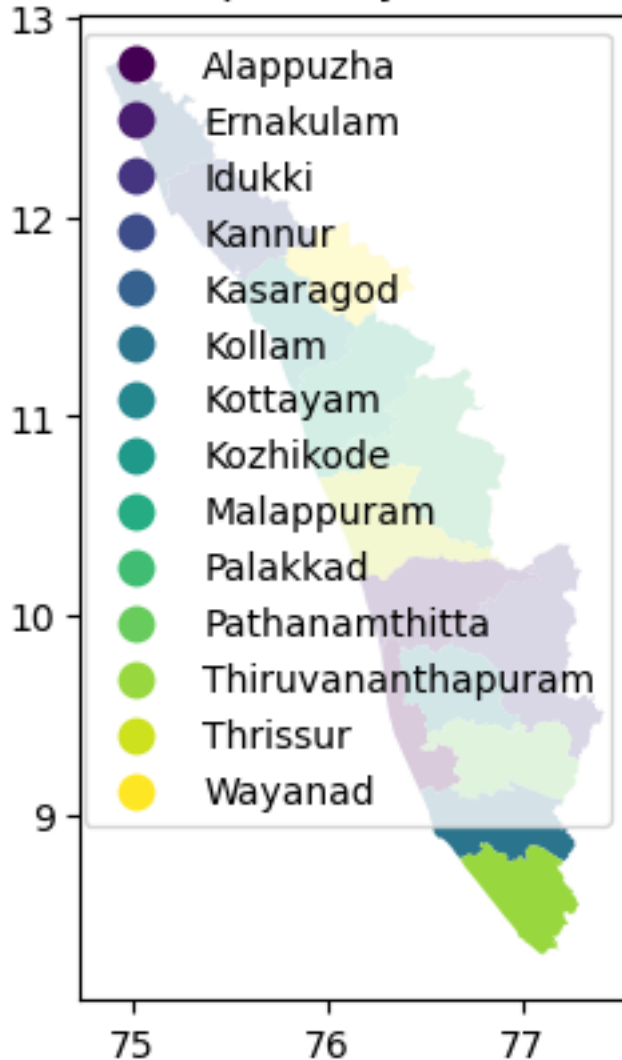
In [110]: print(KE_consumption.head())
   District  total_consumption
0  Alappuzha      1378.539983
1  Ernakulam      1755.480476
2    Idukki       734.516746
3   Kannur       1896.350831
4  Kasaragod      1207.148675

In [113]: data_map = gpd.read_file("/Users/kirthanshaker/Desktop/SCMA 631 Data Files /KERALA_DISTRICTS (1).geojson")
data_map = data_map.rename(columns={'dtname': 'total_consumption'})

In [114]: fig, ax = plt.subplots(1, 1)
data_map.plot(column='total_consumption', cmap='viridis', legend=True, ax=ax)
ax.set_title('Total Consumption by District in Kerala')
plt.show()
```

Result:

Total Consumption by District in Kerala



Interpretation:

is a map visualization of the total consumption by district in Kerala. Each district is color-coded to represent its consumption level. Here is an interpretation of the map:

Interpretation of the Map:

- Title:**
 - The title "Total Consumption by District in Kerala" indicates that the map visualizes consumption data across various districts within the state of Kerala.
- Legend:**
 - The legend on the left provides the color coding for each district, making it easy to identify them on the map. Each color corresponds to a specific district.
- Geographical Distribution:**

- The map displays the geographical boundaries of Kerala and highlights each district according to its total consumption.

Insights:

1. **Regional Patterns:**
 - The map helps in understanding the regional distribution of consumption. Districts with higher or lower consumption can be easily identified by their location on the map and their corresponding color.
2. **Comparative Analysis:**
 - By comparing the colors on the map with the legend, we can discern which districts have higher or lower consumption levels. For example, darker colors might represent higher consumption compared to lighter colors, depending on the specific data representation.

Specific Consumption Values:

1. **Alappuzha:**
 - Total Consumption: **1378.54**
 - Represented by a purple color on the map.
2. **Eranakulam:**
 - Total Consumption: **1755.48**
 - Represented by a dark blue color on the map.
3. **Idukki:**
 - Total Consumption: **734.52**
 - Represented by a light blue color on the map.
4. **Kannur:**
 - Total Consumption: **1896.35**
 - Represented by a blue color on the map.
5. **Kasaragod:**
 - Total Consumption: **1207.15**
 - Represented by a light green color on the map.

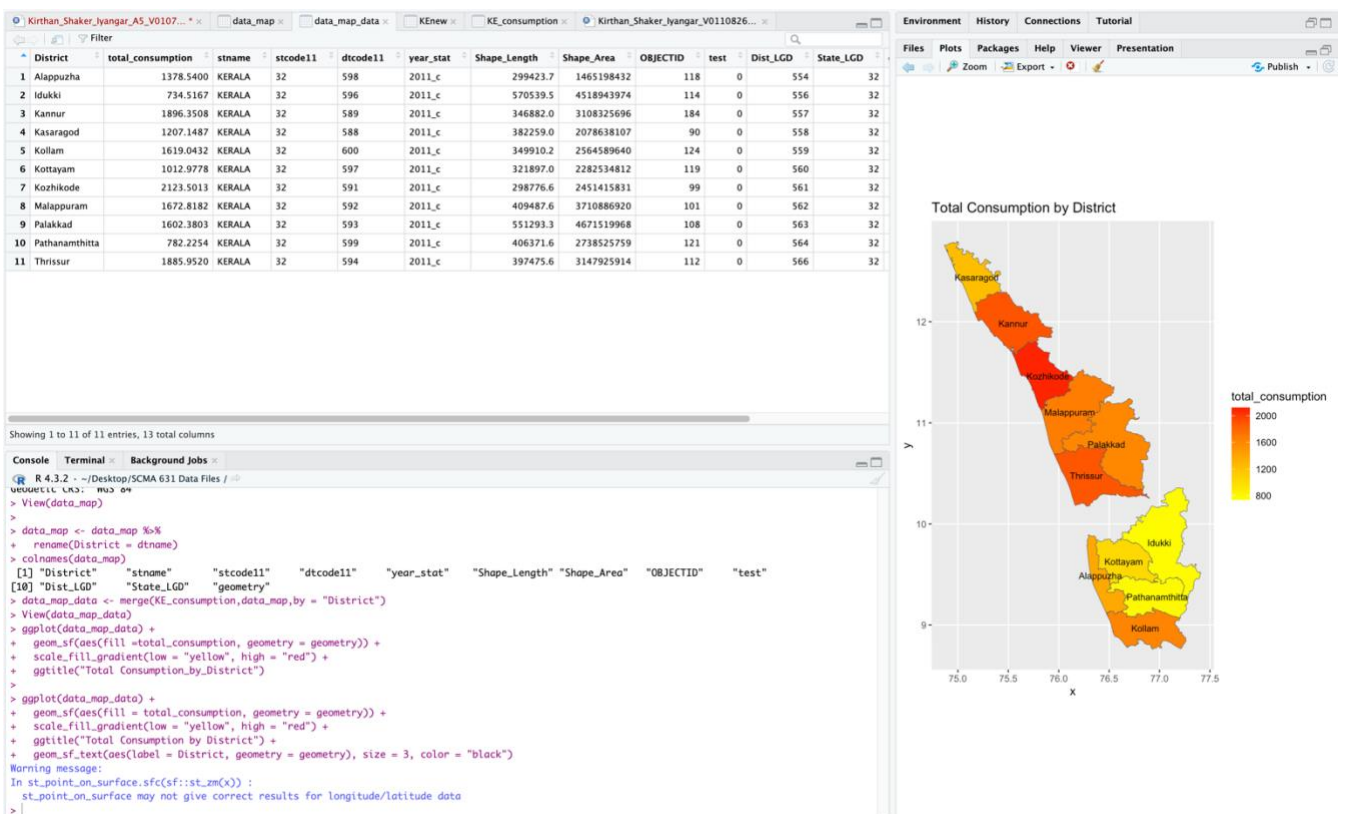
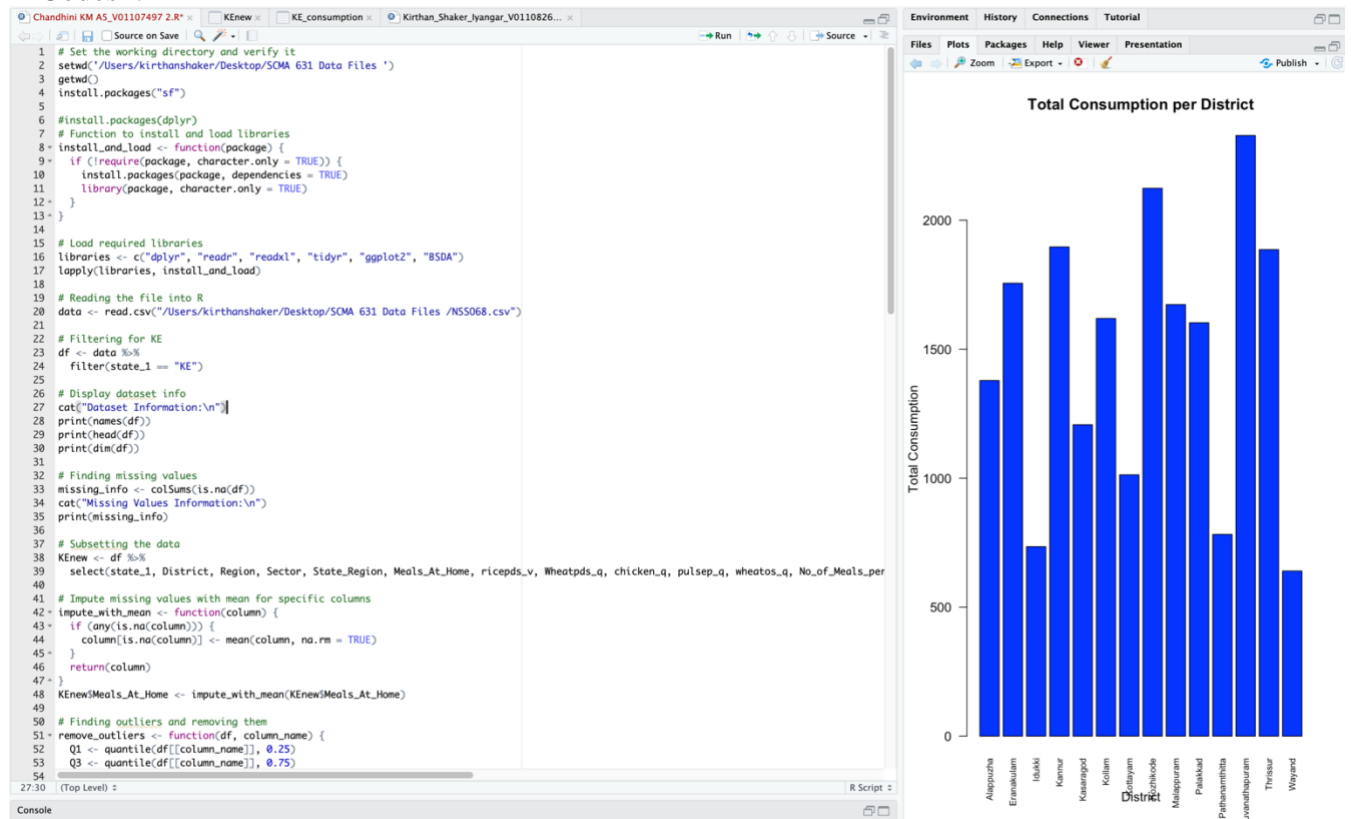
Business Significance

- This map is useful for policymakers, researchers, and analysts to understand and address regional disparities in consumption.
- It can help in planning resource allocation, implementing targeted interventions, and formulating policies that address the specific needs of each district.

The following codes and analyses done on R Codes too :

- The analysis and result were the same, following below is a screenshot of the codes and analysis.

R Codes :



```

89 KNew$Sector <- as.character(KNew$Sector)
90 KNew$District <- ifelse(KNew$District %in% names(district_mapping), district_mapping[KNew$District], KNew$District)
91 KNew$Sector <- ifelse(KNew$Sector %in% names(sector_mapping), sector_mapping[KNew$Sector], KNew$Sector)
92
93 View(KNew)
94
95 hist(KNew$total_consumption, breaks = 10, col = 'blue', border = 'black',
96      xlab = "Consumption", ylab = "Frequency", main = "Consumption Distribution in Mizoram State")
97
98 KE_consumption <- aggregate(total_consumption ~ District, data = KNew, sum)
99 View(KE_consumption)
100 ??barplot
101 barplot(KE_consumption$total_consumption,
102         names.arg = KE_consumption$District,
103         las = 2, # Makes the district names vertical
104         col = 'blue',
105         border = 'black',
106         xlab = "District",
107         ylab = "Total Consumption",
108         main = "Total Consumption per District",
109         cex.names = 0.7) # Adjust the size of district names if needed
110
111 # b) Plot {'any variable of your choice'} on the Karnataka state map using NSS068.csv data
112 install.packages("sf")
113 install.packages("sf", type = "https://cran.rstudio.com/bin/macosx/big-sur-arm64/contrib/4.3/sf_1.0-16.tgz")
114
115 library(ggplot2)
116 library(sf) # mapping
117 library(dplyr)
118 Sys.setenv("SHAPE_RESTORE_SHX" = "YES")
119
120 data_map <- st_read("/Users/kirthanshaker/Desktop/SCMA 631 Data Files /KERALA_DISTRICTS (1).geojson")
121 View(data_map)
122
123 data_map <- data_map %>%
124   rename(District = dname)
125 colnames(data_map)
126 data_map_data <- merge(KE_consumption, data_map, by = "District")
127 View(data_map_data)
128 ggplot(data_map_data) +
129   geom_sf(aes(fill = total_consumption, geometry = geometry)) +
130   scale_fill_gradient(low = "yellow", high = "red") +
131   ggtitle("Total Consumption by District")
132
133 ggplot(data_map_data) +
134   geom_sf(aes(fill = total_consumption, geometry = geometry)) +
135   scale_fill_gradient(low = "yellow", high = "red") +
136   ggtitle("Total Consumption by District") +
137   geom_sf_text(aes(label = District, geometry = geometry), size = 3, color = "black")
138
139
140
141

```

CODES For Both Python and R :

R Codes :

```

# Set the working directory and verify it
setwd('/Users/kirthanshaker/Desktop/SCMA 631 Data Files ')
getwd()
install.packages("sf")

#install.packages(dplyr)
# Function to install and load libraries
install_and_load <- function(package) {
  if (!require(package, character.only = TRUE)) {
    install.packages(package, dependencies = TRUE)
    library(package, character.only = TRUE)
  }
}

```

```

}

# Load required libraries
libraries <- c("dplyr", "readr", "readxl", "tidyr", "ggplot2", "BSDA")
lapply(libraries, install_and_load)

# Reading the file into R
data <- read.csv("/Users/kirthanshaker/Desktop/SCMA 631 Data Files /NSSO68.csv")

# Filtering for KE
df <- data %>%
  filter(state_1 == "KE")

# Display dataset info
cat("Dataset Information:\n")
print(names(df))
print(head(df))
print(dim(df))

# Finding missing values
missing_info <- colSums(is.na(df))
cat("Missing Values Information:\n")
print(missing_info)

# Subsetting the data
KEnew <- df %>%
  select(state_1, District, Region, Sector, State_Region, Meals_At_Home, ricepds_v, Wheatpds_q, chicken_q, pulsep_q, wheatos_q, No_of_Meals_per_day)

# Impute missing values with mean for specific columns
impute_with_mean <- function(column) {
  if (any(is.na(column))) {
    column[is.na(column)] <- mean(column, na.rm = TRUE)
  }
  return(column)
}
KEnew$Meals_At_Home <- impute_with_mean(KEnew$Meals_At_Home)

# Finding outliers and removing them
remove_outliers <- function(df, column_name) {
  Q1 <- quantile(df[[column_name]], 0.25)
  Q3 <- quantile(df[[column_name]], 0.75)
  IQR <- Q3 - Q1
  lower_threshold <- Q1 - (1.5 * IQR)
  upper_threshold <- Q3 + (1.5 * IQR)
  df <- subset(df, df[[column_name]] >= lower_threshold & df[[column_name]] <= upper_threshold)
  return(df)
}
outlier_columns <- c("ricepds_v", "chicken_q")
for (col in outlier_columns) {
  KEnew <- remove_outliers(KEnew, col)
}

# Summarize consumption
KEnew$total_consumption <- rowSums(KEnew[, c("ricepds_v", "Wheatpds_q", "chicken_q", "pulsep_q", "wheatos_q")], na.rm = TRUE)

# Summarize and display top consuming districts and regions
summarize_consumption <- function(group_col) {
  summary <- KEnew %>%
    group_by(across(all_of(group_col))) %>%
    summarise(total = sum(total_consumption)) %>%
    arrange(desc(total))
}

```

```

    return(summary)
  }
  district_summary <- summarize_consumption("District")
  region_summary <- summarize_consumption("Region")

  cat("Top Consuming Districts:\n")
  print(head(district_summary, 4))
  cat("Region Consumption Summary:\n")
  print(region_summary)

# Rename districts and sectors
district_mapping <- c("1" = "Kasaragod", "2" = "Kannur", "3" = "Wayand", "4" = "Kozhikode", "5" = "Malappuram", "6" = "Palakkad", "7" = "Thrissur", "8" =
"Ernakulam", "9" = "Idukki", "10" = "Kottayam", "11" = "Alappuzha", "12" = "Pathanamthitta", "13" = "Kollam", "14" = "Thiruvananthapuram")
sector_mapping <- c("2" = "URBAN", "1" = "RURAL")

KEnew$District <- as.character(KEnew$District)
KEnew$Sector <- as.character(KEnew$Sector)
KEnew$District <- ifelse(KEnew$District %in% names(district_mapping), district_mapping[KEnew$District], KEnew$District)
KEnew$Sector <- ifelse(KEnew$Sector %in% names(sector_mapping), sector_mapping[KEnew$Sector], KEnew$Sector)

View(KEnew)

hist(KEnew$total_consumption, breaks = 10, col = 'blue', border = 'black',
     xlab = "Consumption", ylab = "Frequency", main = "Consumption Distribution in Mizoram State")

KE_consumption <- aggregate(total_consumption ~ District, data = KEnew, sum)
View(KE_consumption)
??barplot
barplot(KE_consumption$total_consumption,
       names.arg = KE_consumption$District,
       las = 2, # Makes the district names vertical
       col = 'blue',
       border = 'black',
       xlab = "District",
       ylab = "Total Consumption",
       main = "Total Consumption per District",
       cex.names = 0.7) # Adjust the size of district names if needed

# b) Plot { 'any variable of your choice' } on the Karnataka state map using NSSO68.csv data
install.packages("sf")
install.packages("sf", type = "https://cran.rstudio.com/bin/macosx/big-sur-arm64/contrib/4.3/sf_1.0-16.tgz")

library(ggplot2)
library(sf) # mapping
library(dplyr)
Sys.setenv("SHAPE_RESTORE_SHX" = "YES")

data_map <- st_read("/Users/kirthanshaker/Desktop/SCMA 631 Data Files /KERALA_DISTRICTS (1).geojson")
View(data_map)

data_map <- data_map %>%
  rename(District = dname)
colnames(data_map)
data_map_data <- merge(KE_consumption, data_map, by = "District")
View(data_map_data)
ggplot(data_map_data) +
  geom_sf(aes(fill = total_consumption, geometry = geometry)) +
  scale_fill_gradient(low = "yellow", high = "red") +
  ggtitle("Total Consumption_by_District")

ggplot(data_map_data) +

```

```
geom_sf(aes(fill = total_consumption, geometry = geometry)) +
scale_fill_gradient(low = "yellow", high = "red") +
ggtitle("Total Consumption by District") +
geom_sf_text(aes(label = District, geometry = geometry), size = 3, color = "black")
```

Python Codes :

```
pip install geopandas
```

```
import geopandas as gpd
```

```
!pip install geojson
```

```
import pandas as pd
import numpy as np
from scipy import stats
import matplotlib.pyplot as plt
import seaborn as sns
import geopandas as gpd
```

```
data = pd.read_csv("/Users/kirthanshaker/Desktop/SCMA 631 Data Files /NSSO68.csv", low_memory=False)
```

```
display(data)
```

```
Kerala_data = data[data['state_1'] == 'KE']
```

```
missing_values = Kerala_data.isna().sum()
print("Missing values in each column:")
print(missing_values)
```

```
KE = Kerala_data[['state_1', 'District', 'Region', 'Sector', 'State_Region', 'Meals_At_Home', 'ricepds_v', 'Wheatpds_q', 'chicken_q', 'pulsep_q', 'wheatos_q', 'No_of_Meals_per_day']]
```

```
def impute_with_mean(column):
    if column.hasnans:
        column.fillna(column.mean(), inplace=True)
    return column
```

```
KE['Meals_At_Home'] = impute_with_mean(KE['Meals_At_Home'])
```

```
def remove_outliers(df, column_name):
    Q1 = df[column_name].quantile(0.25)
    Q3 = df[column_name].quantile(0.75)
    IQR = Q3 - Q1
    lower_threshold = Q1 - (1.5 * IQR)
    upper_threshold = Q3 + (1.5 * IQR)
    df = df[(df[column_name] >= lower_threshold) & (df[column_name] <= upper_threshold)]
    return df
```

```
outlier_columns = ['ricepds_v', 'chicken_q']
for col in outlier_columns:
    KE = remove_outliers(KE, col)
```

```
KE['total_consumption'] = KE[['ricepds_v', 'Wheatpds_q', 'chicken_q', 'pulsep_q', 'wheatos_q']].sum(axis=1)
```

```
def summarize_consumption(group_col):
    summary = KE.groupby(group_col)['total_consumption'].sum().reset_index()
    summary.sort_values(by='total_consumption', ascending=False, inplace=True)
    return summary
```

```
district_summary = summarize_consumption('District')
```



```

region_summary = summarize_consumption('Region')

print("Top Consuming Districts:")
print(district_summary.head(4))
print("Region Consumption Summary:")
print(region_summary)

district = {'1': 'Kasaragod',
            '2': 'Kannur',
            '3': 'Wayand',
            '4': 'Kozhikode',
            '5': 'Malappuram',
            '6': 'Palakkad',
            '7': 'Thrissur',
            '8': 'Eranakulam',
            '9': 'Idukki',
            '10': 'Kottayam',
            '11': 'Alappuzha',
            '12': 'Pathanamthitta',
            '13': 'Kollam',
            '14': 'Thiruvananthapuram',
            }

sector = {
            '2': 'URBAN',
            '1': 'RURAL'
            }

KE.loc[:, 'District'] = KE['District'].astype(str)
KE.loc[:, 'Sector'] = KE['Sector'].astype(str)
KE.loc[:, 'District'] = KE['District'].map(district).fillna(KE['District'])
KE.loc[:, 'Sector'] = KE['Sector'].map(sector).fillna(KE['Sector'])

KE['District'] = KE['District'].astype(str)
KE['Sector'] = KE['Sector'].astype(str)

KE['District'] = KE['District'].map(district).fillna(KE['District'])
KE['Sector'] = KE['Sector'].map(sector).fillna(KE['Sector'])

print(KE.head())

print(KE.columns)

plt.hist(KE['total_consumption'], bins=10, color='blue', edgecolor='black')
plt.xlabel("Consumption")
plt.ylabel("Frequency")
plt.title("Consumption Distribution in Kerala State")
plt.show()

KE_consumption = KE.groupby('District')['total_consumption'].sum().reset_index()

print(KE_consumption.head())

plt.bar(KE_consumption['District'], KE_consumption['total_consumption'], color='blue', edgecolor='black')
plt.xlabel("District")
plt.ylabel("Total Consumption")
plt.title("Total Consumption per District")
plt.xticks(rotation=90) # Rotate district names for better visibility
plt.show()

data_map = gpd.read_file("/Users/kirthanshaker/Desktop/SCMA 631 Data Files /KERALA_DISTRICTS (1).geojson")

```

```

print(data_map.columns)
print(KE_consumption.columns)

data_map['District'] = KE_consumption['District']

data_map_data = data_map.merge(KE_consumption, left_on='dtname', right_on='District')

print(data_map.columns)

import geopandas as gpd
import pandas as pd
import matplotlib.pyplot as plt

data_map = gpd.read_file("/Users/kirthanshaker/Desktop/SCMA 631 Data Files /KERALA_DISTRICTS (1).geojson")

data_map = data_map.rename(columns={'dtname': 'District'})

display(data_map.rename)

KE_consumption = pd.read_csv("/Users/kirthanshaker/Desktop/SCMA 631 Data Files /KERALA_DISTRICTS (1).geojson", low_memory=False)

KE_consumption = KE_consumption.groupby('District')['total_consumption'].sum().reset_index()

print(KE_consumption.head())

data_map = gpd.read_file("/Users/kirthanshaker/Desktop/SCMA 631 Data Files /KERALA_DISTRICTS (1).geojson")
data_map = data_map.rename(columns={'dtname': 'total_consumption'})

```