

Assignment 2. BA

2024-03-10

Summary

My analysis is derived from "usa_00008". This involved examining a comprehensive dataset related to demographics, utility costs, and other variables within the United States. The summary provides insights into the dataset's characteristics and findings:

- Missing Values:** Exploration with glimpse and colSums revealed that there are no missing values ensuring the dataset's readiness for further analysis.
- The highest electricity costs across all 51 states are the same.**
49 states have the highest gas costs.
Hawaii has the highest water cost.
- Most of the states have imbalance in sex.** In every state there's a difference in the number of men and women; none have an equal ratio of males to females.
- California (state 6) had the highest total cost in 2021 and 2022.**
- Maine (state 23) has the oldest, on average, residents**
- In the year 2022 below points are observed,**
 - the average age in Ohio is 43.
 - there are more female than male in Ohio.
 - it is observed that white people are the largest population group in Ohio, followed by Black/African American, with Japanese being the smallest.
 - English was the most widely spoken language, while Yeshiv, Celtic, Aleut, Eskimo, and Ingotruan were the least spoken languages.

To load required libraries and import the data

```
library(readr)
library(tidyverse)
library(tidyverse)

## Warning: package 'tidyverse' was built under R version 4.3.3

## --- Attaching core tidyverse packages --- tidyverse 2.0.0 ---
## # dplyr      1.1.4      ✓ readr      2.1.5
## # # forcats  1.0.0      ✓ stringr  1.5.1
## # # ggplot2   3.4.4      ✓ tibble    3.2.1
## # # lubridate 1.9.3      ✓ tidyqr    1.3.1
## # # purrr     1.0.2
## # --- Conflicts --- tidyverse_conflicts() ---
## # # dplyr::filter() masks stats::filter()
## # # dplyr::lag()   masks stats::lag()
## # Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become errors

library(tidyverse)
library(tidyverse)
library(tidyverse)
Bigdata <- read_csv("D:/Users/vaishnavi/OneDrive - Kent State University/BA/Assignment 2/usa_00008.csv")
```

To confirm that all the data has been properly imported

```
head(Bigdata)

## # YEAR SAMPLE SERIAL CBSERIAL HMMT CLUSTER STATEFIP STRATA GO COSTELEC
## # 1 2021 202101 1 2.02101e+12 13 2.021e+12 1 88881 3 0
## # 2 2021 202101 2 2.02101e+12 51 2.021e+12 1 88881 3 0
## # 3 2021 202101 3 2.02101e+12 17 2.021e+12 1 120801 3 0
## # 4 2021 202101 4 2.02101e+12 61 2.021e+12 1 178801 3 0
## # 5 2021 202101 5 2.02101e+12 15 2.021e+12 1 16881 3 0
## # 6 2021 202101 6 2.02101e+12 46 2.021e+12 1 16881 4 0
## # COSTGAS COSTMATR COSTFUEL PERMUT SEX AGE RACE RACED LANGUAGE LANGUED
## # 1 0 0 0 0 1 15 1 100 1 100
## # 2 0 0 0 0 1 51 2 67 2 200 1 100
## # 3 0 0 0 0 1 17 1 74 1 100 1 100
## # 4 0 0 0 0 1 61 1 16 1 100 1 100
## # 5 0 0 0 0 1 15 1 83 1 100 1 100
## # 6 0 0 0 0 1 46 2 19 1 100 1 100

tail(Bigdata)

## # YEAR SAMPLE SERIAL CBSERIAL HMMT CLUSTER STATEFIP STRATA GO COSTELEC
## # 6625997 2022 202201 1505108 2.02201e+12 12 2.02201e+12 56 38956 1
## # 6625973 2022 202201 1505107 2.02201e+12 119 2.02201e+12 56 48956 1
## # 6625974 2022 202201 1505107 2.02201e+12 119 2.02201e+12 56 48956 1
## # 6625975 2022 202201 1505107 2.02201e+12 119 2.02201e+12 56 48956 1
## # 6625976 2022 202201 1505108 2.02201e+12 126 2.02201e+12 56 20856 1
## # 6625977 2022 202201 1505108 2.02201e+12 126 2.02201e+12 56 20856 1
## # COSTELEC COSTGAS COSTMATR COSTFUEL PERMUT SEX AGE RACE RACED
## # 6625973 848 968 418 9993 1 72 1 55 1 100
## # 6625973 2488 968 388 258 1 119 1 35 1 100
## # 6625974 2488 968 388 258 2 89 2 27 1 100
## # 6625975 44483994 968 388 258 3 177 1 1 100
## # 6625976 3888 1320 78 9993 1 126 1 66 1 100
## # 6625977 3888 1320 78 9993 2 187 2 58 1 100
## # LANGUAGE LANGUED
## # 6625972 1 100
## # 6625973 1 100
## # 6625974 1 100
## # 6625975 0 0
## # 6625976 1 100
## # 6625977 1 100

nrow(Bigdata)

## [1] 6625977
```

Question 1: Are there any missing values?

Exploration with glimpse and colSums revealed that there are no missing values ensuring the dataset's readiness for further analysis.

```
# To check if there are any missing values
glimpse(Bigdata)

## Rows: 6,625,977
## # YEAR SAMPLE SERIAL CBSERIAL HMMT CLUSTER STATEFIP STRATA GO COSTELEC
## # 1 2021 202101 1 2.02101e+12 13 2.021e+12 1 88881 3 0
## # 2 2021 202101 2 2.02101e+12 51 2.021e+12 1 88881 3 0
## # 3 2021 202101 3 2.02101e+12 17 2.021e+12 1 120801 3 0
## # 4 2021 202101 4 2.02101e+12 61 2.021e+12 1 178801 3 0
## # 5 2021 202101 5 2.02101e+12 15 2.021e+12 1 16881 3 0
## # 6 2021 202101 6 2.02101e+12 46 2.021e+12 1 16881 4 0
## # COSTGAS COSTMATR COSTFUEL PERMUT SEX AGE RACE RACED LANGUAGE LANGUED
## # 1 0 0 0 0 1 15 1 100 1 100
## # 2 0 0 0 0 1 51 2 67 2 200 1 100
## # 3 0 0 0 0 1 17 1 74 1 100 1 100
## # 4 0 0 0 0 1 61 1 16 1 100 1 100
## # 5 0 0 0 0 1 15 1 83 1 100 1 100
## # 6 0 0 0 0 1 46 2 19 1 100 1 100
## # LANGUAGE LANGUED
## # 6625972 1 100
## # 6625973 1 100
## # 6625974 1 100
## # 6625975 0 0
## # 6625976 1 100
## # 6625977 1 100

colSums(is.na(Bigdata))

## # YEAR SAMPLE SERIAL CBSERIAL HMMT CLUSTER STATEFIP STRATA GO COSTELEC
## # 0 0 0 0 0 0 0 0 0 0
## # 0 0 0 0 0 0 0 0 0 0
## # AGE RACE RACED LANGUAGE LANGUED
## # 0 0 0 0 0
```

Question 2: Identify the states with the highest cost of electricity, gas, and water.

The highest electricity costs across all 51 states are the same.

49 states have the highest gas costs.

Hawaii has the highest water cost.

```
# Highest cost of the electricity
VD=>Bigdata %>%
  filter(COSTELEC < 9993)%>%
  summarise(VD)

## # STATEFIP COSTELEC
## # Min.: 1.00 Min.: 0
## # 1st Qu.:12.00 1st Qu.:1000
## # Median :27.00 Median :1800
## # Mean :27.72 Mean :2110
## # 3rd Qu.:42.00 3rd Qu.:2760
## # Max.:56.00 Max.:9990

# The maximum cost of the electricity is 9990
High_cost_Electricity <- VD %>%
  group_by(STATEFIP)%>%
  summarise(ELECTRICITY = max(COSTELEC))%>%
  slice_max(ELECTRICITY,n=1)

# Highest cost of the gas
VD=>Bigdata %>%
  filter(COSTGAS < 9993)%>%
  summarise(VD)

## # STATEFIP COSTGAS
## # Min.: 1.00 Min.: 0
## # 1st Qu.:13.00 1st Qu.: 360
## # Median :27.00 Median : 720
## # Mean :27.55 Mean :1171
## # 3rd Qu.:41.00 3rd Qu.:1440
## # Max.:56.00 Max.:9990

# The maximum cost of the gas is 9990
High_cost_Gas <- VD %>%
  group_by(STATEFIP)%>%
  summarise(GAS = max(COSTGAS))%>%
  slice_max(GAS,n=1)

# Highest cost of Water
VD=>Bigdata %>%
  filter(COSTMATR < 9993)%>%
  summarise(VD)

## # STATEFIP COSTMATR
## # Min.: 1.00 Min.: 0
## # 1st Qu.:12.00 1st Qu.: 100
## # Median :27.00 Median : 500
## # Mean :27.47 Mean : 680
## # 3rd Qu.:42.00 3rd Qu.:1000
## # Max.:56.00 Max.:7100

# The maximum cost of the Water is 7100
High_cost_Water <- VD %>%
  group_by(STATEFIP)%>%
  summarise(WATER = max(COSTMATR))%>%
  slice_max(WATER,n=1)
```

Question 3: Are there any states with an imbalance in Sex?

Yes, all the states have imbalance in sex. In every state there's a difference in the number of men and women. None have an equal ratio of males to females.

```
Bigdata %>%
  group_by(STATEFIP)%>%
  summarise(Female = sum(PERWT[SEX == 2], na.rm = TRUE),
            Male = sum(PERWT[SEX == 1], na.rm = TRUE))%>%
  mutate(Imbalance = (Female - Male))%>%
  arrange(desc(Imbalance))

## # A tibble: 51 x 4
## # STATEFIP Female Male Imbalance
## # <int> <dbl> <dbl> <dbl>
## # 1 36 29198466 19317018 8780208
## # 2 12 22367476 21584875 769901
## # 3 13 11212206 10903806 52010
## # 4 37 18866687 10384048 482039
## # 5 42 13139111 12769553 342158
## # 6 24 6321838 6087951 313867
## # 7 25 7137816 6828881 308935
## # 8 39 11921517 11613918 307629
## # 9 1 5287913 4986260 301653
## # 10 34 9484501 9124328 280173
## # 41 more rows
```

Question 4: Create a new variable that indicates the Total Annual cost that is the sum of the cost of Electricity, Gas, and Water. Which states have the highest total cost?

California (state 6) had the highest total cost in 2021 and 2022.

```
ELE21<-Bigdata %>%
  filter(COSTELEC < 9993 & YEAR==2021)%>%
  select(STATEFIP,COSTELEC)%>%
  group_by(STATEFIP)%>%
  summarise(E_21 = sum(COSTELEC))
head(ELE21)

## # A tibble: 6 x 2
## # STATEFIP E_21
## # <int> <int>
## # 1 119768288
## # 2 12687944
## # 3 4150361536
## # 4 516809752
## # 5 6744658060
## # 6 96558618

ELE22<-Bigdata %>%
  filter(COSTELEC < 9993 & YEAR==2022)%>%
  select(STATEFIP,COSTELEC)%>%
  group_by(STATEFIP)%>%
  summarise(E_22 = sum(COSTELEC))
head(ELE22)

## # A tibble: 6 x 2
## # STATEFIP E_22
## # <int> <int>
## # 1 119768288
## # 2 12687944
## # 3 4150361536
## # 4 516809752
## # 5 6744658060
## # 6 96558618

GAS21<-Bigdata %>%
  filter(COSTGAS < 9993 & YEAR==2021)%>%
  select(STATEFIP,COSTGAS)%>%
  group_by(STATEFIP)%>%
  summarise(G_21 = sum(COSTGAS))
head(GAS21)

## # A tibble: 6 x 2
## # STATEFIP G_21
## # <int> <int>
## # 1 12687944
## # 2 2725424
## # 3 4150361536
## # 4 516809752
## # 5 6744658060
## # 6 96558618

GAS22<-Bigdata %>%
  filter(COSTGAS < 9993 & YEAR==2022)%>%
  select(STATEFIP,COSTGAS)%>%
  group_by(STATEFIP)%>%
  summarise(G_22 = sum(COSTGAS))
head(GAS22)

## # A tibble: 6 x 2
## # STATEFIP G_22
## # <int> <int>
## # 1 12687944
## # 2 2725424
## # 3 4150361536
## # 4 516809752
## # 5 6744658060
## # 6 96558618

MATR21<-Bigdata %>%
  filter(COSTMATR < 9993 & YEAR==2021)%>%
  select(STATEFIP,COSTMATR)%>%
  group_by(STATEFIP)%>%
  summarise(M_21 = sum(COSTMATR))
head(MATR21)

## # A tibble: 6 x 2
## # STATEFIP M_21
## # <int> <int>
## # 1 12687944
## # 2 2725424
## # 3 4150361536
## # 4 516809752
## # 5 6744658060
## # 6 96558618

MATR22<-Bigdata %>%
  filter(COSTMATR < 9993 & YEAR==2022)%>%
  select(STATEFIP,COSTMATR)%>%
  group_by(STATEFIP)%>%
  summarise(M_22 = sum(COSTMATR))
head(MATR22)

## # A tibble: 6 x 2
## # STATEFIP M_22
## # <int> <int>
## # 1 12687944
## # 2 2725424
## # 3 4150361536
## # 4 516809752
## # 5 6744658060
## # 6 96558618

# Merging the datasets for each year
total_cost_2021 <- cbind(ELE21,GAS21,MATR21) %>%
  select(1:5) %>%
  mutate(TOTAL_21 = (E_21+G_21+M_21))
head(total_cost_2021)

## # STATEFIP E_21 G_21 M_21 TOTAL_21
## # 1 119768288 20684952 12687988 162529048
## # 2 12687944 4924608 2725424 28257976
## # 3 4150361536 25762884 44410994 22855514
## # 4 516809752 19762888 44855744 92754884
## # 5 6744658060 23149366 268516166 1248323668
## # 6 96558618 34687424 33807512 144173614

total_cost_2022 <- cbind(ELE22,GAS22,MATR22) %>%
  select(1:5) %>%
  mutate(TOTAL_22 = (E_22+G_22+M_22))
head(total_cost_2022)

## # STATEFIP E_22 G_22 M_22 TOTAL_22
## # 1 119768288 20684952 12687988 162529048
## # 2 12687944 4924608 2725424 28257976
## # 3 4150361536 25762884 44410994 22855514
## # 4 516809752 19762888 44855744 92754884
## # 5 6744658060 23149366 268516166 1248323668
## # 6 96558618 34687424 33807512 144173614
```

Question 5: Which state has the oldest, on average, residents?

Maine (state 23) has the oldest, on average, residents.

```
Bigdata %>%
  group_by(STATEFIP)%>%
  summarise(AvgAge = mean(AGE, na.rm = TRUE))%>%
  arrange(desc(AvgAge))%>%
  head(1)

## # A tibble: 1 x 2
## # STATEFIP AvgAge
## # <int> <dbl>
## # 1 23 47.0
```

Question 6: What can you say about the residents of Ohio based on their age, sex, race, and language. Use only the most recent data.

```
ohio_data <- Bigdata %>%
  filter(STATEFIP == 39 & YEAR == 2022)

# Analyze Age for Ohio
In the year 2022, the average age in Ohio is 43.

ohio_age_summary <- ohio_data %>%
  summarise(AverageAge = mean(AGE, na.rm = TRUE))%>%
  ohio_age_summary

## # AverageAge
## # 1 43.2016
```

Analyze Sex for Ohio

In the year 2022, there are more female than male in Ohio.

```
ohio_sex_distribution <- ohio_data %>%
  group_by(SEX)%>%
  summarise(Count = n())%>%
  mutate(SexLabel = ifelse(SEX == 1, "Male", ifelse(SEX == 2, "Female", "Missing")))
ohio_sex_distribution

## # A tibble: 2 x 3
## # SEX Count SexLabel
## # <int> <int> <chr>
## # 1 50902 Male
## # 2 61724 Female
```

Analyze Race for Ohio

In the year 2022, it is observed that white people are the largest population group in Ohio, followed by Black/African American, with Japanese being the smallest.

```
ohio_race_distribution <- ohio_data %>%
  group_by(RACE)%>%
  summarise(Count = n())%>%
  mutate(RaceLabel = case_when(
    RACE == 1 ~ "White",
    RACE == 2 ~ "Black/African American",
    RACE == 3 ~ "American Indian or Alaska Native",
    RACE == 4 ~ "Chinese",
    RACE == 5 ~ "Japanese",
    RACE == 6 ~ "Other Asian or Pacific Islander",
    RACE == 7 ~ "Other race, nec",
    RACE == 8 ~ "Two major races",
    RACE == 9 ~ "Three or more major races",
    TRUE ~ "Unknown"
  ))
ohio_race_distribution

## # A tibble: 9 x 3
## # RACE Count RaceLabel
## # <int> <int> <chr>
## # 1 98911 White
## # 2 29981 Black/African American
## # 3 278 American Indian or Alaska Native
## # 4 507 Chinese
## # 5 82 Japanese
## # 6 2069 Other Asian or Pacific Islander
## # 7 1598 Other race, nec
## # 8 6712 Two major races
## # 9 519 Three or more major races
```

Analyze Language for Ohio

In the year 2022, English was the most widely spoken language, while Yeshiv, Celtic, Aleut, Eskimo, and Ingotruan were the least spoken languages.

```
ohio_language_distribution <- ohio_data %>%
  group_by(LANGUAGE)%>%
  summarise(Total = n())%>%
  mutate(after = LANGUAGE, LABEL = case_when(
    LANGUAGE == 0 ~ "N/A or blank",
    LANGUAGE == 01 ~ "English",
    LANGUAGE == 02 ~ "German",
    LANGUAGE == 03 ~ "Russian, Jewish",
    LANGUAGE == 04 ~ "Dutch",
    LANGUAGE == 05 ~ "Swedish",
    LANGUAGE == 06 ~ "Danish",
    LANGUAGE == 07 ~ "Norwegian",
    LANGUAGE == 08 ~ "Finnish",
    LANGUAGE == 09 ~ "Scandinavian",
    LANGUAGE == 10 ~ "Italian",
    LANGUAGE == 11 ~ "French",
    LANGUAGE == 12 ~ "Spanish",
    LANGUAGE == 13 ~ "Portuguese",
    LANGUAGE == 14 ~ "Rumanian",
    LANGUAGE == 15 ~ "Celtic",
    LANGUAGE == 16 ~ "Greek",
    LANGUAGE == 17 ~ "Albanian",
    LANGUAGE == 18 ~ "Yiddish, Yiddish",
    LANGUAGE == 19 ~ "Russian",
    LANGUAGE == 20 ~ "Ukrainian, Ruthenian, Little Russian",
    LANGUAGE == 21 ~ "Polish",
    LANGUAGE == 22 ~ "Czech",
    LANGUAGE == 23 ~ "Slovak",
    LANGUAGE == 24 ~ "Slovene",
    LANGUAGE == 25 ~ "Lithuanian",
    LANGUAGE == 26 ~ "Other Baltic Slavic",
    LANGUAGE == 27 ~ "Slavic unknown",
    LANGUAGE == 28 ~ "Persian, Iranian, Farsi",
    LANGUAGE == 29 ~ "Other Persian dialects",
    LANGUAGE == 30 ~ "Hindi and related",
    LANGUAGE == 31 ~ "Romany, Gypsy",
    LANGUAGE == 32 ~ "Arabic",
    LANGUAGE == 33 ~ "Hebrew, Yiddish",
    LANGUAGE == 34 ~ "Other Afro-Asiatic languages",
    LANGUAGE == 35 ~ "Uraltic",
    LANGUAGE == 36 ~ "Turkic",
    LANGUAGE == 37 ~ "Other Altaic",
    LANGUAGE == 38 ~ "Caucasian, Georgian, Avar",
    LANGUAGE == 39 ~ "Basque",
    LANGUAGE == 40 ~ "Dravidian",
    LANGUAGE == 41 ~ "Indo-European",
    LANGUAGE == 42 ~ "Burmese, Khmer",
    LANGUAGE == 43 ~ "Chinese",
    LANGUAGE == 44 ~ "Tibetan",
    LANGUAGE == 45 ~ "Burmese, Lisu, Lolo",
    LANGUAGE == 46 ~ "Korean",
    LANGUAGE == 47 ~ "Thai, Siamese, Lao",
    LANGUAGE == 48 ~ "Japanese",
    LANGUAGE == 49 ~ "Korean",
    LANGUAGE == 50 ~ "Vietnamese",
    LANGUAGE == 51 ~ "Other East/Southeast Asian",
    LANGUAGE == 52 ~ "Indonesian",
    LANGUAGE == 53 ~ "Other Malay",
    LANGUAGE == 54 ~ "Filipino, Tagalog",
    LANGUAGE == 55 ~ "Micronesian, Polynesian",
    LANGUAGE == 56 ~ "Hawaiian",
    LANGUAGE == 57 ~ "Arabic",
    LANGUAGE == 58 ~ "Near East Arabic dialect",
    LANGUAGE == 59 ~ "Hebrew, Israeli",
    LANGUAGE == 60 ~ "African, Ethiopian, etc.",
    LANGUAGE == 61 ~ "Malay",
    LANGUAGE == 62 ~ "Other Afro-Asiatic languages",
    LANGUAGE == 63 ~ "Sub-Saharan Africa",
    LANGUAGE == 64 ~ "African, n.s.",
    LANGUAGE == 65 ~ "American Indian (all)",
    LANGUAGE == 66 ~ "Aleut, Eskimo",
    LANGUAGE == 67 ~ "Algonquian",
    LANGUAGE == 68 ~ "Yukon",
    LANGUAGE == 69 ~ "Navajo",
    LANGUAGE == 70 ~ "Puebloan",
    LANGUAGE == 71 ~ "Other Penutian",
    LANGUAGE == 72 ~ "Zuni",
    LANGUAGE == 73 ~ "Hokan",
    LANGUAGE == 74 ~ "Other Hokan languages",
    LANGUAGE == 75 ~ "Siouan languages",
    LANGUAGE == 76 ~ "Muskogean",
    LANGUAGE == 77 ~ "Keres",
    LANGUAGE == 78 ~ "Iroquoian",
    LANGUAGE == 79 ~ "Caddoan",
    LANGUAGE == 80 ~ "Shoshonian/Nepi",
    LANGUAGE == 81 ~ "Pima, Papago",
    LANGUAGE == 82 ~ "Yuki and other Sonoran, nec",
    LANGUAGE == 83 ~ "Aztec, Nahuatl, Uto-Aztecan",
    LANGUAGE == 84 ~ "Tanoan languages",
    LANGUAGE == 85 ~ "Other Indian languages",
    LANGUAGE == 86 ~ "Mayan languages",
    LANGUAGE == 87 ~ "Arawakan",
    LANGUAGE == 88 ~ "Carib",
    LANGUAGE == 89 ~ "No language",
    LANGUAGE == 90 ~ "Other not reported",
    LANGUAGE == 91 ~ "Not reported, blank",
    TRUE ~ "Unknown"
  ))
ohio_language_distribution
```

```
## # A tibble: 57 x 3
## # LANGUAGE LABEL Total
## # <int> <chr> <int>
## # 1 0 N/A or blank 519
## # 2 1 English 107913
## # 3 2 German 1116
## # 4 3 Yiddish, Jewish 1
## # 5 4 Dutch 173
## # 6 5 Danish 2
## # 7 6 Norwegian 4
## # 8 7 Norwegian 4
## # 9 8 Italian 148
## # 10 9 French 271
## # 47 more rows
```