

Analyzing and Predicting the hospitalization Costs using Data Science - DSC520_FinalProject_Part3

Guruprasad Velikadu Krishnamoorthy

2023-03-04

Contents

1	Introduction	2
2	The problem statement being addressed.	3
3	The Approach to address the problem statement	4
3.1	Exploring the dataset and Description of the columns used:	5
3.2	Identifying the Technical Requirements:	8
3.2.1	R Packages used:	8
3.2.2	Plots and Tables used in the Project:	8
3.3	Importing and Cleansing the datasets:	9
3.3.1	Importing and Cleansing commands in R	11
3.3.2	The Final Datasets after cleansing	22
3.3.3	What information is not self-evident?	23
3.3.4	Approach for slicing and dicing the data:	24
3.4	Summarizing the data and plotting graphs for visualization:	24
3.4.1	Plot 1: Top 8 Medicare Procedures in South Dakota in 2020	24

3.4.2	Plot 2: Top 10 Hospitalizations in Mobile, Alabama in 2020	28
3.4.3	Plot 3: Chemotherapy treatment in West Virginia	30
3.4.4	Plot 4: Hospitalizations in Kansas City, MO vs Kansas City,KS for Liver Ailments . .	31
3.4.5	Plot 5: National Average vs Wyoming for Cost of Heart related Ailments treatment .	33
3.4.6	Plot 6 : Top 10 Hospitalizations in each year across the US	36
3.4.7	Plot 7: Creating a Treemap of Hospitalizations in the US for Heart Transplant Proce- dures from 2018 to 2020	38
3.4.8	Plot 8: Most Profitable hospitals in Nebraska	39
3.4.9	Plot 9: Net Profit margin vs payments per patient ‘for-Profit’ hospital in Anchor- age,Alaska	41
3.4.10	Plot 10: Comparing Cost of Spinal Procedure - National vs State Average	44
3.4.11	Plot 11: Inpatient stays for Chemotherapy between 2018-2020	46
3.4.12	Plot 12: How profitable are hospitals between 2011 and 2020 in the Southern states? .	48
3.4.13	Plot 13: Comparing Average income of “For-Profit” vs “Non-Profit” hospitals in 2019	51
3.4.14	Plot 14: Distribution of hospitals that treats PNEUMONIA across contiguous US states	54
3.4.15	Additional Analysis 1:	57
3.4.16	Additional Analysis 2:	60
3.5	Building a model	61
3.6	Using K-means Clustering	67
3.6.1	Cluster 1- To find clusters based on hospital’s Net income	67
3.6.2	Cluster 2 - To find clusters based on Total hospitalization and cost of Treating Bron- chitis in the US	71
4	Summary of the Interesting insights from the Analysis :	73
5	Implications to the Target audience	76
6	Limitations of the analysis and concluding Remarks:	77

7 Session Info	78
References	80

1 Introduction

Gaps in the current system: With the tremendous advancement that has been made in technology in the last few decades, things that used to be considered luxuries such as smartphones or 50-inch television are not very expensive anymore. However, even with all the technological progress one thing that has not changed much is the cost of medical treatment. On the contrary, many medical procedures have gotten so expensive over the years that some are not affordable, even after having Health Insurance and patients have to go for external sources such as GoFundMe to collect funds. When the doctor suggests a medical treatment and when we call the hospital to ask how much the treatment might cost, there is a good chance the hospital may not be willing to tell the actual amount, considering how complicated the health care system in the US is. Many people get sucked into this mess without completely understanding all the nuances of the healthcare system and end up with a massive hospital bill that they may not be able to pay. An article from CNBC states that two-thirds of all bankruptcies in the US are tied to issues with Medical Bills. With such a huge problem in hand, there is no better solution to get an idea of how much an essential but non-emergency procedure would cost in a hospital. (Konish n.d.)

How does this study help with the problem? Medical bills can be quite expensive, especially for those that don't have insurance or if they get treatment from an 'out of network' hospital. Most hospitals are not transparent about sharing the treatment costs and they redirect to reach out to health insurance instead. However, not many people know that medical bills can be negotiated with hospitals. If we have the right stats and numbers, the hospitals may be willing to negotiate the bills and it can save thousands of dollars and sometimes can even save bankruptcies for some families. This research focuses on studying the pattern behind the medical costs across various states and cities in the US and identifying a pattern in the costs. Patients can use the data to choose care wisely and avoid expensive hospitals if they have a choice. Also, these numbers can help negotiate the bills with the hospitals.

Role of Data science in this research: Data science can be used to study patterns and derive statistics and compare prices over a period. Also, data from different datasets can be combined to derive the relationship between the variables and a model can be built to predict the prices in the coming years.

Focus of this research: This research focuses on the study to identify the disparity in hospitalization costs across the country. Though the costs may differ from state to state, many times they can be vastly different even within the same city between different hospitals. Also, hospitals are not transparent about sharing the medical costs, so it can be very frustrating to undergo a medical procedure not knowing how much it costs and if the hospital is overcharging for the procedure or not. The solution addresses this problem by studying

the trends in medical costs (for a medical procedure as an example) across the state or region and over a period and provides an overall idea about the differences in the cost. The results will be plotted in maps and plots for visualization. The study also provides results of costs that are grouped based on a state or city or a region for a medical procedure. This can also be used by the patients to compare the costs with the explanation of benefits to get a picture if the costs look normal or inflated.

2 The problem statement being addressed.

The problem statement being addressed in this study is to identify the patterns in the medical costs and to predict the costs in the future. This is accomplished by answering or analyzing some of the common questions listed below. These questions can help us discover the hidden patterns in the data.

1. What are the top 10 medical procedures in a state (for example, South Dakota) in a given period.
2. What are the common medical procedures done in a city (for example, Mobile,Alabama) in a given period.
3. What is the total number of inpatient stays across the state for a specific medical procedure during a given period?
4. What is the total number of hospitalizations for a medical procedure across two nearby cities? Do patients prefer hospitals in one city over the other?
5. What is the national average cost of a few medical procedures and how does that compare across the given state?
6. What are the top 10 hospitalizations across US for a given period?
7. Which states have more hospitalizations for a medical procedure, for example Liver ailments?
8. What are the most profitable hospitals in a state, Nebraska for example?
9. Is there a relationship between Profits made by the hospital to the cost of a medical procedure?
10. What is the national average cost of a medical procedure and how does that compare with the state averages?
11. What is the average cost of a medical treatment across various cities in a given state?
12. Given a zip code in a state (Texas, for example), if multiple hospitals offer the same treatment, how different are the costs?
13. Is there a substantial difference in the costs of medical procedures between the For-Profit and Non-Profit hospitals in a given state?
14. How much has the profit/Loss of the hospitals in a region varied over the years?

15. Is there any correlation between the population across the state/ region to the number of hospitalizations in that region?
16. Is there a pattern in the variation in the cost of a medical procedure between different states? Can it predict the cost of treatment of a medical procedure in that state?

3 The Approach to address the problem statement

The approach includes importing the data from different periods and comparing the trends in the costs of medical treatment. The average costs in the state is used as a benchmark to compare the costs and identify how the hospital's costs compare with others in the same region. The occupancy rates of hospitals are studied and compared with the state's average to discover the pattern between inpatient occupancy and their costs. Also, the Tables and plots in R will be used to visualize the results for easy understanding. Also, the costs of treatments may vary between counties and if the difference is substantial, one can choose to undergo treatment from a cheaper hospital from a neighboring county or a neighboring state in some cases.

This study contains two broader types of research :

1. The first category is to study the data at *the national level and compare them with the state averages*. This includes analysis such as comparing the cost of the National average for a medical procedure versus the state averages, Identifying the top 10 medical procedures that had the most hospitalizations across the nation, etc.
2. The second type of study is focused on a relatively *smaller geographic area such as a city or State or Zip code*. For instance, identifying the top 10 profitable hospitals in a state, Identifying the cities that offer treatment for medical procedures in the state. This also includes comparing the costs of medical procedures in a city that shares borders across states.

Below are the high-level steps or approach to address the problem statement:

Identifying the datasets : There are a variety of sources to find the datasets such as www.data.cms.gov, www.cdc.gov, www.who.int, www.eia.gov, Kaggle, etc. However, the real datasets on hospitals and the costs of procedures may be a challenge to find. So, the first step is to invest time to find the right dataset.

Exploring the dataset : Once the datasets are identified, the second step is to study the data and understand each field in the dataset. We may not find all the relevant data in the same dataset. So, we may

have to identify more than one dataset and join them to extract meaningful insights. Hence understanding the data is very important.

Identifying the relevant columns : The next step is to come up with a high-level plan (algorithm) for the approach to address the problem. This step will require identifying the fields that we are planning to join across different datasets and analyzing the number of records in the original data to see if a subset of data is required for analysis.

Gathering the Technical requirements : The next step is technical where we identify the list of R packages and plots required to address the problem. In this step, we will perform a very high-level analysis of the technical requirements as more packages can be added later during the development stage. Further steps of cleansing data and transforming the data will be done during the development phase.

Data Cleansing and exploration : The next steps involve loading the data into R and cleansing and transforming the data. Depending upon the size of the dataset, a subset of data can be used for the research. In this stage, we can identify new fields required to be used in the plots to join the datasets.

Developing a Model In this final stage after exploring the data, the fields that can predict the medical costs will be identified and a regression or classification model will be built on that.

These steps are further explained in detail in this section.

3.1 Exploring the dataset and Description of the columns used:

Seven datasets are used in this research and the details are given below.

- i. MUP_IHP_RY21_P02_V10_DY18_PrvcSvc.csv
- ii. MUP_IHP_RY21_P02_V10_DY19_PrvcSvc.csv
- iii. MUP_IHP_RY22_P02_V10_DY20_PrvcSvc.csv

(“Medicare Inpatient Hospitals - by Provider and Service” 2022)

- iv. MUP_IHP_RY21_P02_V10_Dy18_Geo.csv
- v. MUP_IHP_RY21_P02_V10_Dy19_Geo.csv
- vi. MUP_IHP_RY22_P02_V10_Dy20_Geo.csv

(“Medicare Inpatient Hospitals - by Geography and Service” 2022)

- **Description:** The Datasets 1,2 and 3 contain data about the use, hospital charges, and payments for more than 3000 US hospitals that received IPPS payments (Inpatient Prospective Payment System). Each dataset is for the years 2020, 2019, and 2018 respectively. The Datasets 4,5, and 6 are like the ones above except the data is summarized at state and nation for each category of the medical procedure. The 3 datasets are for the years 2020, 2019, and 2018 respectively.
- **Description of the fields:**
 - *Rndrng_Privr_CCN* : This describes the Rendering provider’s CMS certification number, which is a unique Identifier to identify a hospital.
 - *Rndrng_Privr_Org_Name* : Denotes the Rendering Provider organization name which is the name of the hospital.
 - *Rndrng_Privr_St*, *Rndrng_Privr_City*, *Rndrng_Privr_State_Abrvtn*, and *Rndrng_Privr_Zip5* : These fields indicate the street address, City name, State name, and the Zip codes of the hospital.
 - *DRG_Cd*: Classification code for the Diagnosis Related Group that indicates unique code for each type of medical procedure.
 - *DRG_Desc* : Description of the Diagnosis Related Group code. This indicates the description of the medical procedure.
 - *Tot_Dschrgs* : This indicates the total number of all discharges billed by all the providers.
 - *Avg_Tot_Pymt_Amt* : This indicates the average total payment made to all providers which include the copays, deductibles, etc.
 - *Rndrng_Privr_Geo_Lvl* : This indicates if the data is summarized at the National level or State level.
 - *Rndrng_Privr_Geo_Desc* : This indicates the name of the state.
- **Observations about the data**
 - The data in these datasets do not have any empty or missing values.
 - However, to distinguish the data in the dataset from each other a new column for Year should be added.
 - The Filed CCN has some values that are of length 5 and some are of length 6. We may need to standardize the length by adding a leading 0 to join the datasets that are used in this project.
 - Some of the datasets are huge with up to 200,000 rows and have up to 15 columns each, though only the columns described above will be utilized in this project.

- **Period when the data was collected:** Dec, 2022

vii. **NASHP_HCT_Data_2022_November.csv**

(“Hospital Provider Cost Report” 2022)

- **Description:** This contains the Healthcare Cost report details collected by the NASHP team (National Academy for State Health Policy) in collaboration with rice university. It contains the Healthcare cost report details such as hospital details, Bed size, Occupancy Rates, and lots of financials such as net revenue, net income, etc. It contains data for approximately 5100 hospitals across the US for the period between 2010 to 2021.
- **Description of the fields:**
 - **CCN#** : This indicates the Unique identifier for each hospital (CMS Certification Number). This denotes a Medicare Provider Number.
 - **Facility Type** : This indicates if the hospital is a Critical Access Hospital or a Short-Term Facility.
 - **Year** : Indicates the reporting year.
 - **Hospital Name, Hospital Abbreviated Name, Address, Zip Code, City, State** : Indicates the name and address of the hospital facility.
 - **Hospital Ownership Type** : This indicates the classification of the hospital if it is “For-Profit” or “Non-Profit” or “Governmental”.
 - **Bed Size** : Indicates the number of beds available for Inpatient care and does not include the beds for emergency rooms or Outpatient care.
 - **Inpatient Occupancy** : Percentage of beds occupied by inpatient care during the given year.
 - **Net Income (Loss) and Net Profit Margin** : Indicates the financial metrics that will be used in this project.
- **Observations about the data :**
 - The dataset has 85 fields and 51894 rows of data, though in this project about 12 fields will be used.
 - This dataset requires a lot of cleansing. The field “Hospital Name” and “Hospital Abbreviated name” has invalid numeric values that should be excluded.
 - The Address field and City have invalid values such as “,” “X,” “Y”, etc. should be excluded.

- The field’s Bed Size, Occupancy rates, Net Revenue, and Net Profit margin all have invalid values such as “-99999” and “.” Should be excluded or defaulted to 0.

- *Period when the data was collected:* 2022

3.2 Identifying the Technical Requirements:

3.2.1 R Packages used:

- *ggplot2*: This is one of the most important packages used for data visualization and is used to create Aesthetic visualization plots.
- *lubridate*: The dataset contains date fields and hence lubridate is used to extract individual date components if needed.
- *dplyr and plyr*: This is used for various type of data manipulation.
- *stringr*: This is used for String Manipulations.
- *tidyr*: This is used to tidy up the data.
- *purrr*: This is used to apply the function to a list and also to use the map functions.
- *reshape2*: This package is used for Data munging to convert the rows to columns or for data Melting operations.
- *Magrittr* : This package will be used to utilize the pipe functionality
- *DT*: This is for an interactive tabular experience through DataTables JavaScript library
- *rmarkdown and knitr*: These are required to knit the documents as pdf using RMarkdown.

3.2.2 Plots and Tables used in the Project:

- *Plots*
 - *Histograms*: This is to show the distribution of values for a variable.
 - *Scatterplots*: This is to show two variables in comparison with each other
 - *Boxplots*: This is to identify the outliers in the data.
 - *Density plots*: This is to show the density distribution

- *Facet_wrap*: This takes data in one variable and splits data in it and creates separate panes for each set to fit them in the plot.
- *Facet_grid*: This is like facet wrap but it applies all levels of a variable to a row and column and fits them in a plot.
- *line plot*: Line plot to show the relation between two variables.
- *Heat Map* : Heat map can be used to explain the density in the pattern of data based on the Heat color schemes.

- **Tables:**

- *kableExtra* : Tabular data can be printed using kableExtra
- *Datatable*: From DT package

3.3 Importing and Cleansing the datasets:

Importing data:

This project has seven datasets of which the first six datasets have similar structure and all these datasets are in CSV format. Hence Importing the data and cleansing will be similar for the first 6 datasets. The seventh dataset is in Excel format.

To import the CSV datasets, *read.csv()* in R will be used and for importing Excel datasets, *read_excel()* function from the *readxl* package will be used.

Transforming Datatypes:

- After importing the data, the structure of the data is analyzed to see if the fields are of appropriate data types. Some of the fields that have fewer unique values can be converted to factors.
- Some of the fields such as *total amount, number of beds, number of hospitalizations, etc.* can be numeric fields. If they are not numeric, datatype conversion should be done.
- Fields such as *DRG_Cd* which represents the code for medical procedures need not be numeric instead, they should be a Character field.

Renaming Columns and limiting the number of columns:

- The columns should be renamed to shorter and easy-to-understand names. Some of the fields have spaces in them, which are replaced with an underscore so that they are easy to use.

- Some of the datasets have about 85 columns, and as we may not require all of them, only a few columns will be loaded into our final data frame by using the select function.
- The column names in the final data frame can be validated by using *names()* function.

Cleansing the data:

- Some of the fields like Hospital names have invalid characters such as Double quotes, which should be removed.
- The Medical procedures are identified with the DRG_Cd (Unique code for a medical procedure), however, while combining the datasets it is important to make sure the Description of the procedures are the same across the datasets. For example, the below descriptions of the same procedure code will yield different results while plotting the graphs. The names should be cleansed after combining the datasets.

TRACH W MV >96 HRS OR PDX EXC FACE, MOUTH & NECK W/O MAJ O.R.

TRACH WITH MV >96 HRS OR PDX EXC FACE, MOUTH & NECK WITHOUT MAJ O.R.

- The Hospital code, CCN, which uniquely identifies a medical facility should have the same length across all datasets. As these are read from CSV and excel files, the codes may be interpreted as Numbers. They should be converted to Characters and 0's should be padded in the front if the length is not the same for all values. While joining data frames, this can come in handy.
- Some values in Hospital names and city fields have only numbers that are invalid values, and they should be excluded. Some of the address fields have invalid values such as only numbers and one or two characters. Though these are invalid values, they will not be excluded because the address is not a critical field in this analysis and the invalid values can be ignored.
- Also, some character fields have invalid values such as period (".") in them and can be excluded. This should be done cautiously. For instance, "ST. VINCENTS EAST" is valid, while "9.2" is an invalid hospital name. Hence the length of the field along with the presence of periods should be analyzed before excluding them.
- Some of the important fields such as Bed size, and city names have "NA" and they can be excluded.
- Some numeric fields have invalid values such as -99999 for Bed size and they can be excluded. As the dataset have many thousands of records, excluding a few invalid values will be reasonable

3.3.1 Importing and Cleansing commands in R

```
# Calling the Libraries used
```

```
library(readxl)
```

```
library(dplyr)
```

```
library(lubridate)
```

```
library(readr)
```

```
library(ggplot2)
```

```
library(ggthemes)
```

```
library(tidyr)
```

```
library(DT)
```

```
library(scales)
```

```
library(stringr)
```

```
library(ggalt)
```

```
library(treemap)
```

```
library(knitr)
```

```
library(lemon)
```

```
library(kableExtra)
```

```
## Warning in !is.null(rmarkdown::metadata$output) && rmarkdown::metadata$output
```

```
## %in% : 'length(x) = 2 > 1' in coercion to 'logical(1)'
```

```
library(reshape2)
```

```
library(tidyverse)
```

```
library(usmap)
```

```
library(plotly)
```

```
library(leaflet)
```

```
library(magrittr)
```

```
library(olsrr)
```

```
# library(QuantPsyc) library(relaimpo)
```

```
library(car)
```

```
library(sf)
```

```
library(factoextra)
library(NbClust)
library(maps)
```

```
# Exporting Dataset #1
# MUP_IHP_RY22_P02_V10_DY20_PrsvSvc.csv
medic_prsvsvc_2020_df_orig <- read.csv("MUP_IHP_RY22_P02_V10_DY20_PrsvSvc.csv")

# Exporting the Dataset #2-
# MUP_IHP_RY21_P02_V10_DY19_PrsvSvc.csv
medic_prsvsvc_2019_df_orig <- read.csv("MUP_IHP_RY21_P02_V10_DY19_PrsvSvc.csv")

# Exporting the Dataset #3-
# MUP_IHP_RY21_P02_V10_DY18_PrsvSvc.csv
medic_prsvsvc_2018_df_orig <- read.csv("MUP_IHP_RY21_P02_V10_DY18_PrsvSvc.csv")

# Validating the structure of one of the imported dataset
str(medic_prsvsvc_2020_df_orig)
```

```
## 'data.frame': 165281 obs. of 15 variables:
## $ Rndrng_Privr_CCN : int 10001 10001 10001 10001 10001 10001 10001 10001 10001 10001 ...
## $ Rndrng_Privr_Org_Name : chr "Southeast Health Medical Center" "Southeast Health Medical Center" ...
## $ Rndrng_Privr_St : chr "1108 Ross Clark Circle" "1108 Ross Clark Circle" "1108 Ross Clark Circle" ...
## $ Rndrng_Privr_City : chr "Dothan" "Dothan" "Dothan" "Dothan" ...
## $ Rndrng_Privr_State_Abrvtn: chr "AL" "AL" "AL" "AL" ...
## $ Rndrng_Privr_State_FIPS : int 1 1 1 1 1 1 1 1 1 1 ...
## $ Rndrng_Privr_Zip5 : int 36301 36301 36301 36301 36301 36301 36301 36301 36301 36301 ...
## $ Rndrng_Privr_RUCA : num 1 1 1 1 1 1 1 1 1 1 ...
## $ Rndrng_Privr_RUCA_Desc : chr "Metropolitan area core: primary flow within an urbanized area of ...
## $ DRG_Cd : int 23 24 25 27 38 39 56 57 64 65 ...
## $ DRG_Desc : chr "CRANIOTOMY W MAJOR DEVICE IMPLANT OR ACUTE COMPLEX CNS PDX W MCC ...
## $ Tot_Dschrgs : int 51 16 30 11 25 50 14 23 101 99 ...
## $ Avg_Submtd_Cvrd_Chrg : num 166626 105888 182686 117293 112529 ...
## $ Avg_Tot_Pymt_Amt : num 37033 26209 29492 15641 11399 ...
## $ Avg_Mdcr_Pymt_Amt : num 33908 22531 27776 14357 9559 ...
```

```

# Creating a new dataframe with a limited set of fields
medic_prvsvc_2020_df <- medic_prvsvc_2020_df_orig %>%

  select(CCN = Rndrng_Privr_CCN, hosp_name = Rndrng_Privr_Org_Name,
         street = Rndrng_Privr_St, city = Rndrng_Privr_City,
         state = Rndrng_Privr_State_Abrvtn, zip = Rndrng_Privr_Zip5,
         DRG_Cd, DRG_Desc, Tot_Dschrgs, Avg_Tot_Pymt_Amt)

medic_prvsvc_2019_df <- medic_prvsvc_2019_df_orig %>%

  select(CCN = Rndrng_Privr_CCN, hosp_name = Rndrng_Privr_Org_Name,
         street = Rndrng_Privr_St, city = Rndrng_Privr_City,
         state = Rndrng_Privr_State_Abrvtn, zip = Rndrng_Privr_Zip5,
         DRG_Cd, DRG_Desc, Tot_Dschrgs, Avg_Tot_Pymt_Amt)

medic_prvsvc_2018_df <- medic_prvsvc_2018_df_orig %>%

  select(CCN = Rndrng_Privr_CCN, hosp_name = Rndrng_Privr_Org_Name,
         street = Rndrng_Privr_St, city = Rndrng_Privr_City,
         state = Rndrng_Privr_State_Abrvtn, zip = Rndrng_Privr_Zip5,
         DRG_Cd, DRG_Desc, Tot_Dschrgs, Avg_Tot_Pymt_Amt)

# Adding a new column for the Year to each Dataframe
medic_prvsvc_2020_df$year <- factor(2020)
medic_prvsvc_2019_df$year <- factor(2019)
medic_prvsvc_2018_df$year <- factor(2018)

# Cleansing the dataset by converting a few columns as
# Factors
medic_prvsvc_2020_df$DRG_Cd <- factor(medic_prvsvc_2020_df$DRG_Cd)
medic_prvsvc_2019_df$DRG_Cd <- factor(medic_prvsvc_2019_df$DRG_Cd)
medic_prvsvc_2018_df$DRG_Cd <- factor(medic_prvsvc_2018_df$DRG_Cd)

# Cleansing by removing invalid characters
medic_prvsvc_2020_df$DRG_Desc <- str_replace(medic_prvsvc_2020_df$DRG_Desc,

```

```

    "\\\"", "")
medic_prvsvc_2019_df$DRG_Desc <- str_replace(medic_prvsvc_2019_df$DRG_Desc,
    "\\\"", "")
medic_prvsvc_2018_df$DRG_Desc <- str_replace(medic_prvsvc_2018_df$DRG_Desc,
    "\\\"", "")

# Padding a 0 to standardize the length of CCN field
medic_prvsvc_2020_df$CCN <- str_pad(medic_prvsvc_2020_df$CCN,
    6, pad = "0")
medic_prvsvc_2019_df$CCN <- str_pad(medic_prvsvc_2019_df$CCN,
    6, pad = "0")
medic_prvsvc_2018_df$CCN <- str_pad(medic_prvsvc_2018_df$CCN,
    6, pad = "0")

```

```

# Combining the First 3 datasets into 1 Dataset
medic_prvsvc_combined_df <- rbind(medic_prvsvc_2018_df, medic_prvsvc_2019_df,
    medic_prvsvc_2020_df)

# Validating the counts of individual dataset and the
# combined one
nrow(medic_prvsvc_2018_df)

```

```
## [1] 193003
```

```
nrow(medic_prvsvc_2019_df)
```

```
## [1] 188806
```

```
nrow(medic_prvsvc_2020_df)
```

```
## [1] 165281
```



```
nrow(medic_prvsvc_combined_df)
```

```
## [1] 547090
```

```
identical(nrow(medic_prvsvc_combined_df), nrow(medic_prvsvc_2018_df) +  
          nrow(medic_prvsvc_2019_df) + nrow(medic_prvsvc_2020_df))
```

```
## [1] TRUE
```

```
# Cleansing the combined dataframe by Standardizing the  
# DRG_Desc Values by eliminating the duplicates  
unique_drg_cd_prvsvc <- medic_prvsvc_combined_df %>%  
  count(DRG_Cd, DRG_Desc) %>%  
  select(DRG_Cd, DRG_Desc) %>%  
  distinct(DRG_Cd, .keep_all = TRUE)  
names(unique_drg_cd_prvsvc) <- c("unique_DRG_Cd", "Unique_DRG_Desc")  
# Using Merge function to standardize the Values  
medic_prvsvc_combined_df <- merge(x = medic_prvsvc_combined_df,  
  y = unique_drg_cd_prvsvc, by.x = "DRG_Cd", by.y = "unique_DRG_Cd") %>%  
  select(CCN, hosp_name, street, city, state, zip, DRG_Cd,  
         DRG_Desc = Unique_DRG_Desc, Tot_Dschrgs, Avg_Tot_Pymt_Amt,  
         year)
```

```
# Importing 4th dataset MUP_IHP_RY22_P02_V10_Dy20_Geo.csv  
medic_geo_2020_df_orig <- read.csv("MUP_IHP_RY22_P02_V10_Dy20_Geo.csv")  
# Importing 5th Dataset MUP_IHP_RY21_P02_V10_Dy19_Geo.csv  
medic_geo_2019_df_orig <- read.csv("MUP_IHP_RY21_P02_V10_Dy19_Geo.csv")  
# Importing 6th Dataset MUP_IHP_RY21_P02_V10_Dy18_Geo.csv  
medic_geo_2018_df_orig <- read.csv("MUP_IHP_RY21_P02_V10_Dy18_Geo.csv")  
  
# Creating a new dataframe with a limited set of fields  
medic_geo_2020_df <- medic_geo_2020_df_orig %>%  
  select(Geo_lvl = Rndrng_Prldr_Geo_Lvl, Geo_cd = Rndrng_Prldr_Geo_Cd,
```

```

    Geo_desc = Rndrng_Privr_Geo_Desc, DRG_Cd, DRG_Desc,
    Tot_Dschrgs, Avg_Tot_Pymt_Amt)
medic_geo_2019_df <- medic_geo_2019_df_orig %>%
  select(Geo_lvl = Rndrng_Privr_Geo_Lvl, Geo_cd = Rndrng_Privr_Geo_Cd,
    Geo_desc = Rndrng_Privr_Geo_Desc, DRG_Cd, DRG_Desc,
    Tot_Dschrgs, Avg_Tot_Pymt_Amt)
medic_geo_2018_df <- medic_geo_2019_df_orig %>%
  select(Geo_lvl = Rndrng_Privr_Geo_Lvl, Geo_cd = Rndrng_Privr_Geo_Cd,
    Geo_desc = Rndrng_Privr_Geo_Desc, DRG_Cd, DRG_Desc,
    Tot_Dschrgs, Avg_Tot_Pymt_Amt)

# Adding a new column for the Year to each Dataframe
medic_geo_2020_df$year <- factor(2020)
medic_geo_2019_df$year <- factor(2019)
medic_geo_2018_df$year <- factor(2018)

# Cleansing the dataset by converting a few columns as
# Factors
medic_geo_2020_df$DRG_Cd <- factor(medic_geo_2020_df$DRG_Cd)
medic_geo_2019_df$DRG_Cd <- factor(medic_geo_2019_df$DRG_Cd)
medic_geo_2018_df$DRG_Cd <- factor(medic_geo_2018_df$DRG_Cd)

# Cleansing by removing invalid characters
medic_geo_2020_df$DRG_Desc <- str_replace(medic_geo_2020_df$DRG_Desc,
  "\\\"", "")
medic_geo_2019_df$DRG_Desc <- str_replace(medic_geo_2019_df$DRG_Desc,
  "\\\"", "")
medic_geo_2018_df$DRG_Desc <- str_replace(medic_geo_2018_df$DRG_Desc,
  "\\\"", "")

# Combining the Datasets 4,5 and 6 into one DataFrame
medic_geo_combined_df <- rbind(medic_geo_2018_df, medic_geo_2019_df,
  medic_geo_2020_df)

```

```
nrow(medic_geo_combined_df)
```

```
## [1] 82831
```

```
# Validating the counts of the combined Dataframe
```

```
identical(nrow(medic_geo_combined_df), nrow(medic_geo_2018_df) +  
          nrow(medic_geo_2019_df) + nrow(medic_geo_2020_df))
```

```
## [1] TRUE
```

```
# Cleansing the combined dataframe by Standardizing the
```

```
# DRG_Desc Values by eliminating the duplicates
```

```
unique_drg_cd_geo <- medic_geo_combined_df %>%
```

```
  count(DRG_Cd, DRG_Desc) %>%
```

```
  select(DRG_Cd, DRG_Desc) %>%
```

```
  distinct(DRG_Cd, .keep_all = TRUE)
```

```
names(unique_drg_cd_geo) <- c("unique_DRG_Cd", "Unique_DRG_Desc")
```

```
# Using Merge function to standardize the Values
```

```
medic_geo_combined_df <- merge(x = medic_geo_combined_df, y = unique_drg_cd_geo,
```

```
  by.x = "DRG_Cd", by.y = "unique_DRG_Cd") %>%
```

```
  select(DRG_Cd, DRG_Desc = Unique_DRG_Desc, Geo_lvl1, Geo_cd,
```

```
    Geo_desc, Tot_Dschrgs, Avg_Tot_Pymt_Amt, year)
```

```
# Importing the 7th dataset
```

```
# NASHP_HCT_Data_2022_November.xlsx
```

```
NASHP_excel <- "NASHP_HCT_Data_2022_November.xlsx"
```

```
NASHP_orig_df <- suppressWarnings(read_excel(NASHP_excel, sheet = 2))
```

```
# Creating a new dataframe with limited columns and also
```

```
# renaming the fields
```

```
NASHP_df <- NASHP_orig_df %>%
```

```
  select(CCN = `CCN#`, facility_type = `Facility Type`, year = Year,
```

```
    hosp_name = `Hospital Name`, hosp_abbr_name = `Hospital Abbreviated Name`,
```

```
    address = Address, zip = `Zip Code`, city = City, state = State,
```

```

health_system = `Health System`, hosp_ownership = `Hospital Ownership Type`,
bed_size = `Bed Size`, inpatient_occupancy = `Inpatient Occupancy`,
net_income = `Net Income (Loss)`, net_profit_margin = `Net Profit Margin`)
# Examining the structure of the new dataframe
str(NASHP_df)

```

```

## tibble [51,894 x 15] (S3: tbl_df/tbl/data.frame)
##  $ CCN                : chr [1:51894] "010001" "010001" "010001" "010001" ...
##  $ facility_type       : chr [1:51894] "Short Term (General and Specialty) Hospitals" "Short Term (Gen
##  $ year                : num [1:51894] 2011 2012 2013 2014 2015 ...
##  $ hosp_name           : chr [1:51894] "SOUTHEAST ALABAMA MEDICAL CENTER" "SOUTHEAST ALABAMA MEDICAL C
##  $ hosp_abbr_name      : chr [1:51894] "SE AL MC" "SE AL MC" "SE AL MC" "SE AL MC" ...
##  $ address             : chr [1:51894] "1108 ROSS CLARK CIRCLE" "1108 ROSS CLARK CIRCLE" "1108 ROSS CL
##  $ zip                 : chr [1:51894] "36301" "36301" "36301" "36301" ...
##  $ city                : chr [1:51894] "DOTHAN" "DOTHAN" "DOTHAN" "DOTHAN" ...
##  $ state               : chr [1:51894] "AL" "AL" "AL" "AL" ...
##  $ health_system       : chr [1:51894] "Southeast Health AL" "Southeast Health AL" "Southeast Health AL"
##  $ hosp_ownership      : chr [1:51894] "Governmental" "Governmental" "Governmental" "Governmental" ..
##  $ bed_size            : num [1:51894] 408 408 408 398 400 387 387 327 342 332 ...
##  $ inpatient_occupancy: num [1:51894] 0.607 0.578 0.633 0.683 0.614 ...
##  $ net_income          : chr [1:51894] "4089394" "6672718" "12739814" "1125576" ...
##  $ net_profit_margin   : chr [1:51894] "1.6318108370371089E-2" "2.6184635846616799E-2" "4.86257836126

```

```

# Cleansing the dataset by converting a few columns as
# Factors and some as Numeric
NASHP_df$CCN <- as.factor(NASHP_df$CCN)
NASHP_df$facility_type <- as.factor(NASHP_df$facility_type)
NASHP_df$year <- as.factor(NASHP_df$year)
NASHP_df$state <- as.factor(NASHP_df$state)
NASHP_df$health_system <- as.factor(NASHP_df$health_system)
NASHP_df$hosp_ownership <- as.factor(NASHP_df$hosp_ownership)
NASHP_df$net_income <- as.numeric(NASHP_df$net_income)

```

```
## Warning: NAs introduced by coercion
```

```
NASHP_df$net_profit_margin <- as.numeric(NASHP_df$net_profit_margin)
```

```
## Warning: NAs introduced by coercion
```

```
# Creating a reusable function to check for Number only  
# values  
numbers_only <- function(var1) {  
  !grepl("\\D", var1)  
}  
  
# Checking the Number of rows that have only Numbers in  
# the Hospital name field and filtering them as the count  
# is low.  
nrow(NASHP_df[numbers_only(NASHP_df$hosp_name), ])
```

```
## [1] 5
```

```
NASHP_df <- NASHP_df[!numbers_only(NASHP_df$hosp_name), ]  
# Validating if the cleanup is done correctly  
NASHP_df[numbers_only(NASHP_df$hosp_name), ]
```

```
## # A tibble: 0 x 15  
## # ... with 15 variables: CCN <fct>, facility_type <fct>, year <fct>,  
## #   hosp_name <chr>, hosp_abbr_name <chr>, address <chr>, zip <chr>,  
## #   city <chr>, state <fct>, health_system <fct>, hosp_ownership <fct>,  
## #   bed_size <dbl>, inpatient_occupancy <dbl>, net_income <dbl>,  
## #   net_profit_margin <dbl>
```

```
# Checking other fields such as City and State to see if  
# they have 'only Number' values and there is none.  
sum(numbers_only(NASHP_df$city))
```

```
## [1] 0
```

```
sum(numbers_only(NASHP_df$state))
```

```
## [1] 0
```

```
# Validating To find if any city has name that is lesser  
# than or 2 characters. Results none.  
NASHP_df[nchar(NASHP_df$city) <= 2, ]
```

```
## # A tibble: 0 x 15  
## # ... with 15 variables: CCN <fct>, facility_type <fct>, year <fct>,  
## #   hosp_name <chr>, hosp_abbr_name <chr>, address <chr>, zip <chr>,  
## #   city <chr>, state <fct>, health_system <fct>, hosp_ownership <fct>,  
## #   bed_size <dbl>, inpatient_occupancy <dbl>, net_income <dbl>,  
## #   net_profit_margin <dbl>
```

```
# Creating a function that checks if a Period is present  
contains_period <- function(var1) {  
  grepl("\\.", var1)  
}  
# Checking if there are any periods in the city name and  
# if the length is lesser than or equal to 5. Results 1  
# row  
NASHP_df[contains_period(NASHP_df$city) & nchar(NASHP_df$city) <=  
  5, ]
```

```
## # A tibble: 1 x 15  
##   CCN    facil~1 year hosp_~2 hosp_~3 address zip   city  state healt~4 hosp_~5  
##   <fct> <fct>   <fct> <chr>   <chr>   <chr>   <chr> <chr> <fct> <fct>   <fct>  
## 1 280030 Short ~ 2011 CREIGH~ CREIGH~ 30      68131 9.01  NE    Indepe~ For-Pr~  
## # ... with 4 more variables: bed_size <dbl>, inpatient_occupancy <dbl>,  
## #   net_income <dbl>, net_profit_margin <dbl>, and abbreviated variable names  
## #   1: facility_type, 2: hosp_name, 3: hosp_abbr_name, 4: health_system,  
## #   5: hosp_ownership
```

```
# Excluding 1 row from above
NASHP_df <- NASHP_df[!(contains_period(NASHP_df$city) & nchar(NASHP_df$city) <=
  5), ]
# finding invalid values in bed size and excluding them
nrow(NASHP_df[is.na(NASHP_df$bed_size), ])
```

```
## [1] 71
```

```
NASHP_df <- NASHP_df[!is.na(NASHP_df$bed_size), ]
# Excluding the -99999 values from bed size
NASHP_df <- NASHP_df[-c(which(NASHP_df$bed_size < 0)), ]
# Checking for nulls in other numeric fields
sum(is.na(NASHP_df$net_profit_margin))
```

```
## [1] 1004
```

```
sum(is.na(NASHP_df$net_income))
```

```
## [1] 1005
```

```
NASHP_df <- NASHP_df[!is.na(NASHP_df$net_profit_margin), ]
NASHP_df <- NASHP_df[!is.na(NASHP_df$net_income), ]
# Converting the net_profit_margin into percentage
NASHP_df$net_profit_margin <- round(NASHP_df$net_profit_margin *
  100, digits = 2)
# Finding invalid values in the Net Income and
# eliminating them
sum(NASHP_df$net_income == ".")
```

```
## [1] 0
```

```
NASHP_df <- NASHP_df[!NASHP_df$net_income == ".", ]
# Displaying Final row count in the data frame after all
# Cleanups
nrow(NASHP_df)
```

[1] 50796

3.3.2 The Final Datasets after cleansing

1. The datasets `medic_prsvsvc_2018_df`, `medic_prsvsvc_2019_df`, and `medic_prsvsvc_2020_df` will be combined to form a final combined dataset for Provider Services- ***medic_prsvsvc_combined_df***

```
kbl(head(medic_prsvsvc_combined_df[1:6, c(1:6)]), caption = "Medical Provider Services DataFrame",  
     booktabs = T) %>%  
kable_styling(latex_options = c("striped", "hold_position"))
```

Table 1: Medical Provider Services DataFrame

CCN	hosp_name	street	city	state	zip
220116	Tufts Medical Center	800 Washington Street	Boston	MA	2111
220110	Brigham And Women's Hospital	75 Francis Street	Boston	MA	2115
520098	University Of Wi Hospitals & Clinics Authority	600 Highland Avenue	Madison	WI	53792
060024	University Of Colorado Hospital Authority	12605 E 16th Ave	Aurora	CO	80045
100128	Tampa General Hospital	1 Tampa General Cir	Tampa	FL	33606
490009	University Of Virginia Medical Center	1215 Lee Street	Charlottesville	VA	22908

2. The datasets `medic_geo_2018_df`, `medic_geo_2019_df`, and `medic_geo_2020_df` will be combined to form the final combined dataset- ***medic_geo_combined_df***

```
kbl(head(medic_geo_combined_df[1:6, c(1, 3:8)]), caption = "Medical Geographics DataFrame",  
     booktabs = T) %>%  
kable_styling(latex_options = c("striped", "hold_position"))
```

Table 2: Medical Geographics DataFrame

DRG_Cd	Geo_lvl	Geo_cd	Geo_desc	Tot_Dschrgrs	Avg_Tot_Pymt_Amt	year
1	National	NA	National	2228	276916.9	2018
1	State	20	Kansas	21	186416.0	2020
1	State	36	New York	177	318296.2	2018
1	State	55	Wisconsin	32	284826.5	2018
1	State	17	Illinois	93	303856.4	2020
1	State	1	Alabama	12	208240.9	2019

3. The 7th dataset after cleansing will look like the below: ***NASHP_df***


```
kbl(head(NASHP_df[1:6, c(1, 3, 5:8)]), caption = "NASHP DataFrame",
     booktabs = T) %>%
kable_styling(latex_options = c("striped", "hold_position"))
```

Table 3: NASHP DataFrame

CCN	year	hosp_abbr_name	address	zip	city
010001	2011	SE AL MC	1108 ROSS CLARK CIRCLE	36301	DOTHAN
010001	2012	SE AL MC	1108 ROSS CLARK CIRCLE	36301	DOTHAN
010001	2013	SE AL MC	1108 ROSS CLARK CIRCLE	36301	DOTHAN
010001	2014	SE AL MC	1108 ROSS CLARK CIRCLE	36301	DOTHAN
010001	2015	SE AL MC	1108 ROSS CLARK CIRCLE	36301	DOTHAN
010001	2016	SE AL MC	1108 ROSS CLARK CIRCLE	36301	DOTHAN

3.3.3 What information is not self-evident?

- Each individual dataset (after cleansing) for Medical Provider services- *medic_prvsvc_combined_df* contains data such as Hospital details, Details of the medical procedures, and the Average payment amount for Medicare for each service each year. However, what is not self-evident is the combined metrics of the cost of a medical procedure, such as the average cost of a procedure for a city or state, or zip code.
- Also, each dataset contains the details of medical procedures only for a year, combining the data for each year helps us study the trends in the costs across a geographic location, which is not self-evident.
- All the above points also apply to the datasets *medic_geo_combined_df*. This dataset contains data at the state and National level but not at the city or Zip code. Hence, combining the data for each year helps us study the trends at the state/ National level for changes in the cost of medical procedures over the years, as it is not self-evident.
- The dataset *NASHP_df* contains the details of the hospital such as the address, type of hospital, Ownership type, and a lot of financial information such as Net Revenue, Profit margin, etc. However, it only presents the raw data if the hospital was profitable or not. Combining the data from the other Medicare datasets and studying the patterns in the payments made for Medicare helps us understand why the hospitals were profitable for a certain period. This information is not self-evident unless the datasets are combined, and the underlying patterns are studied.
- In some cities that share borders between states, the medical costs can vary significantly across states. For instance, for someone that lives in *Boise City, OK* (which is in Oklahoma's panhandle), the

states of *Texas, Colorado, New Mexico, and Kansas* are all within 30 minutes of driving distance. This research can help a patient in Boise City compare the costs across the hospitals and choose the appropriate one. This information is not self-evident from the raw data.

3.3.4 Approach for slicing and dicing the data:

- For the most part, the data will be sliced based on the state or City. For some comparison, data will be filtered based on the Zip-code.
- For the study that focuses on the cost of a medical Procedure, the data will be sliced based on the Procedure code (DRG_Cd) or the description (DRG_Desc).
- For the study that involves the medical costs for a given period, filters will be applied to the Year column and required years will be selected.
- In the NASHP_df dataset, apart from the state and city filters, the data will be sliced based on the Hospital ownership type(hosp_ownership). This dataset also has other fields such as health systems and data will be sliced on this field as well.
- The NASHP_df dataset has data for the years 2011 to 2021, though most of the study is focused on the years 2018-2020, hence filter of the year field will be applied to these years.

3.4 Summarizing the data and plotting graphs for visualization:

3.4.1 Plot 1: Top 8 Medicare Procedures in South Dakota in 2020

The data can be grouped based on the *Medical procedure code and/or description* and then *summarized based on Average payments, total payments, and Total and Average hospitalizations*. This is shown below:

```
# creating re-usable function for display purposes
topN <- function(df_name, var, N = 10) {
  df_name %>%
    arrange(desc(var)) %>%
    head(N)
}

# Summarizing based on Medical Procedures done and using
```

```
# the function topN to display 10 rows
medic_prsvsvc_summary_df1 <- medic_prsvsvc_2020_df %>%
  group_by(DRG_Cd, DRG_Desc) %>%
  dplyr::summarise(total_amt_paid = sum(Avg_Tot_Pymt_Amt),
    mean_amt_paid = mean(Avg_Tot_Pymt_Amt), total_discharges = sum(Tot_Dschrgs)) %>%
  topN(., "total_amt_paid", 10)
kbl(head(medic_prsvsvc_summary_df1), caption = "Results of Summary Dataframe 1",
  booktabs = T) %>%
  kable_styling(full_width = F, latex_options = c("striped",
    "hold_position")) %>%
  column_spec(2, width = "20em")
```

Table 4: Results of Summary Dataframe 1

DRG_Cd	DRG_Desc	total_amt_paid	mean_amt_paid	total_discharges
1	HEART TRANSPLANT OR IMPLANT OF HEART ASSIST SYSTEM WITH MCC	20535801.3	297620.31	1343
3	ECMO OR TRACHEOSTOMY WITH MV >96 HOURS OR PRINCIPAL DIAGNOSIS EXCEP	57683736.0	171677.79	9186
4	TRACH W MV >96 HRS OR PDX EXC FACE, MOUTH & NECK W/O MAJ O.R.	38420005.1	103279.58	7219
5	LIVER TRANSPLANT WITH MCC OR INTESTINAL TRANSPLANT	6918481.0	130537.38	1025
6	LIVER TRANSPLANT W/O MCC	371192.7	46399.09	100
7	LUNG TRANSPLANT	3768959.9	144960.00	498

The state information can be included in the grouping above and then filters can be applied to a specific state to study the numbers for that state. In this example, Medicare payments in South Dakota are studied based on the grouping done in the states.

```
# Creating a summary data frame that summarizes the Total
# payments made in each state per medical procedure
prsvsvc_2020_summary_df1 <- medic_prsvsvc_2020_df %>%
  group_by(DRG_Cd, DRG_Desc, state) %>%
  dplyr::summarise(total_amt_paid = sum(Avg_Tot_Pymt_Amt),
    mean_amt_paid = mean(Avg_Tot_Pymt_Amt), total_discharges = sum(Tot_Dschrgs)) %>%
  arrange(desc(total_amt_paid))
```

Table 5: Results of Summary Dataframe 2

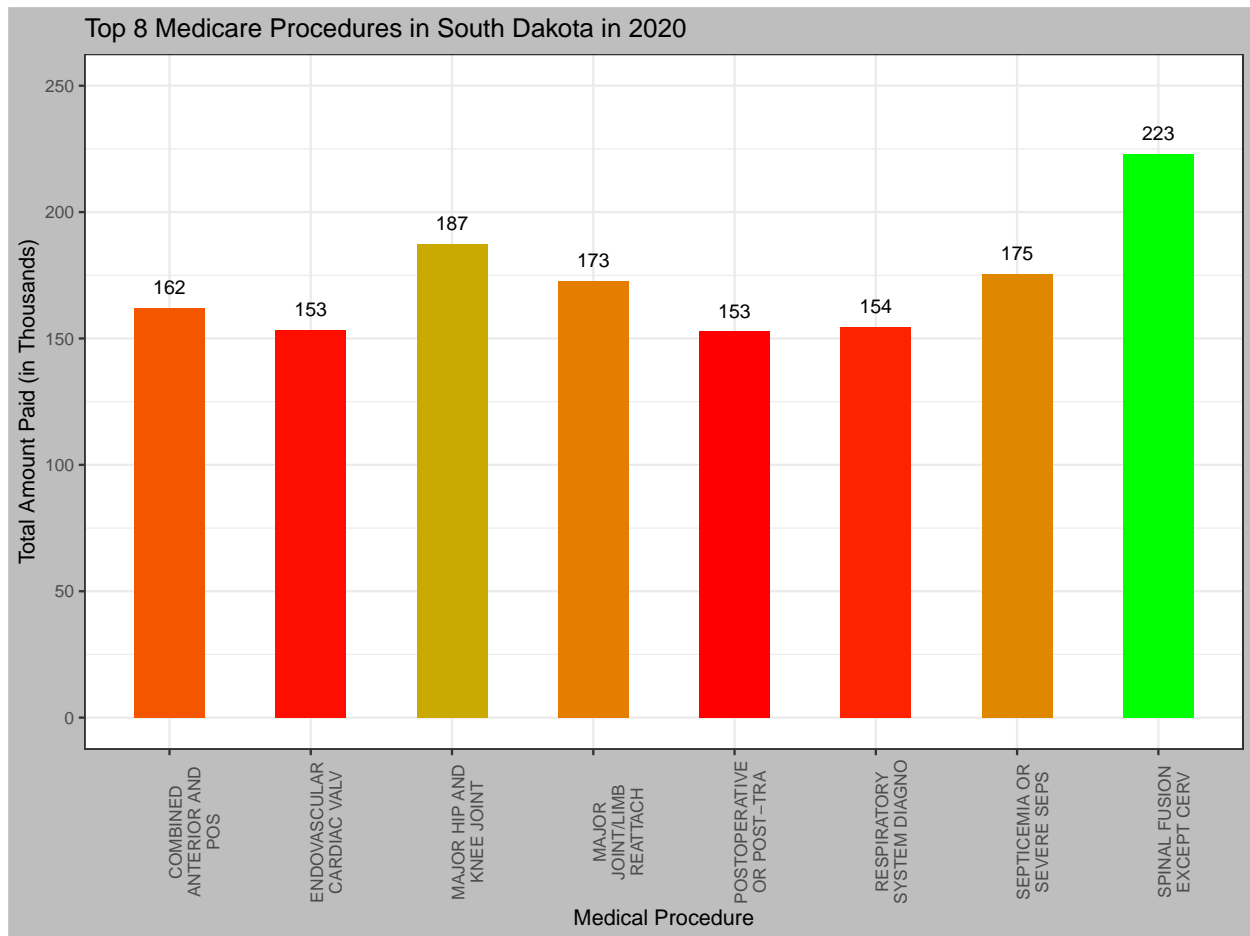
DRG_Cd	DRG_Desc	state	total_amt_paid	mean_amt_paid	total_discharges
460	SPINAL FUSION EXCEPT CERVICAL WITHOUT MCC	SD	223023.9	27877.99	225
470	MAJOR HIP AND KNEE JOINT REPLACEMENT OR REATTACHMENT OF LOWER EXTREM	SD	187349.3	14411.49	2452
870	SEPTICEMIA OR SEVERE SEPSIS WITH MV >96 HOURS	SD	175364.8	58454.92	60
483	MAJOR JOINT/LIMB REATTACHMENT PROCEDURE OF UPPER EXTREMITIES	SD	172643.0	17264.30	563
454	COMBINED ANTERIOR AND POSTERIOR SPINAL FUSION WITH CC	SD	162167.4	40541.86	76
207	RESPIRATORY SYSTEM DIAGNOSIS WITH VENTILATOR SUPPORT >96 HOURS	SD	154435.6	51478.52	47

```

# Top 8 Medicate Payments in the state of South Dakota
# with the number of Hospitalizations.
SD_2020_df <- prvsvc_2020_summary_df1 %>%
  filter(state == "SD") %>%
  topN(., "total_discharges", 8)
kbl(head(SD_2020_df), caption = "Results of Summary Dataframe 2",
  booktabs = T) %>%
  kable_styling(full_width = F, latex_options = c("striped")) %>%
  column_spec(2, width = "10em")

```

```
# Top 8 Medicare Procedures in South Dakota in 2020
ggplot(data = SD_2020_df, aes(str_wrap(str_sub(DRG_Desc, 1,
25), 15), total_amt_paid/1000, fill = total_amt_paid)) +
  geom_bar(stat = "identity", width = 0.5, position = "dodge") +
  geom_text(aes(label = round(total_amt_paid/1000)), vjust = -1,
    size = 3) + scale_fill_gradient(low = "Red", high = "Green") +
  labs(x = "Medical Procedure", y = "Total Amount Paid (in Thousands)",
    title = "Top 8 Medicare Procedures in South Dakota in 2020") +
  scale_y_continuous(labels = comma) + ylim(0, 250) + theme_bw() +
  theme(plot.title = element_text(size = 12), axis.text.x = element_text(size = 8,
    angle = 90), axis.text.y = element_text(size = 8),
    axis.title = element_text(size = 10), plot.background = element_rect(fill = "Grey"),
    legend.position = "none")
```



Analysis of Plot 1 Results: The results indicate Top payments made in South Dakota in 2020 by

Medicare was for *Spinal fusion procedures* followed by *Hip & Knee joint treatments*.

3.4.2 Plot 2: Top 10 Hospitalizations in Mobile, Alabama in 2020

In the example below, the intention is to study the top ten hospitalizations in Mobile, Alabama in the year 2020. Hence the filter for the state and city will be applied first, and then data is grouped based on the Procedure code & its description. Then the summary of total Hospitalization is extracted and top 10 results are displayed.

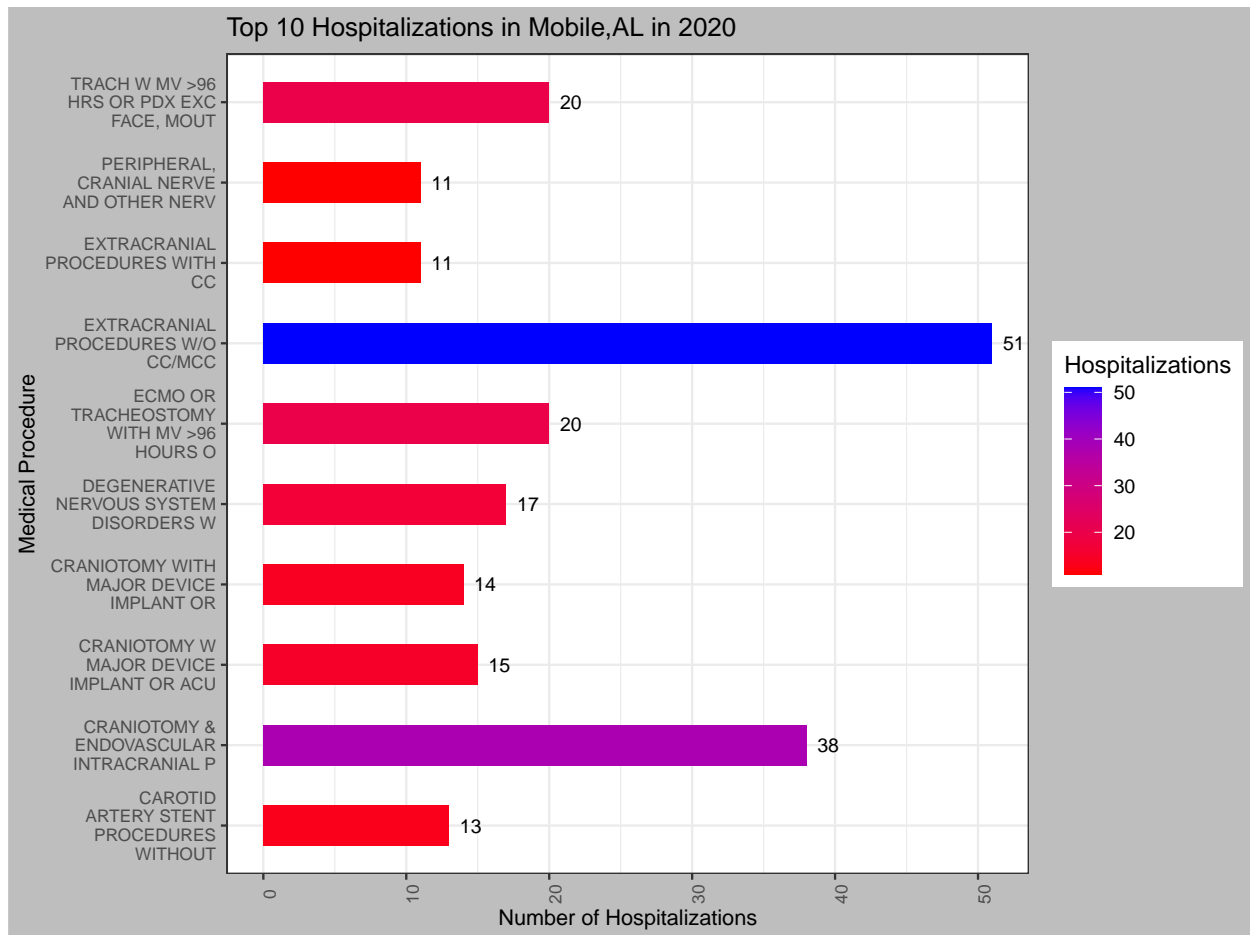
```
# Top 10 Hospitalizations in Mobile, AL in 2020 and using
# Reusable function topN
mobileAL_hosp_df_2020 <- medic_prvsvc_2020_df %>%
  filter(state == "AL" & city == "Mobile") %>%
  group_by(DRG_Cd, DRG_Desc) %>%
  dplyr::summarise(Hospitalizations = sum(Tot_Dschrgs)) %>%
  topN(., "Hospitalizations")

kbl(head(mobileAL_hosp_df_2020), caption = "Results of Summary Dataframe 3",
     booktabs = T) %>%
  kable_styling(full_width = F, latex_options = c("striped",
    "hold_position")) %>%
  column_spec(2, width = "20em")

# Top 10 Hospitalizations in Mobile, Alabama in 2020
ggplot(data = mobileAL_hosp_df_2020, aes(str_wrap(str_sub(DRG_Desc,
  1, 40), 15), Hospitalizations)) + geom_bar(stat = "identity",
  width = 0.5, aes(fill = Hospitalizations)) + scale_fill_gradient(low = "red",
  high = "blue") + coord_flip() + labs(x = "Medical Procedure",
  y = "Number of Hospitalizations", title = "Top 10 Hospitalizations in Mobile,AL in 2020") +
  geom_text(aes(label = Hospitalizations), hjust = -0.5,
    size = 3) + theme_bw() + theme(plot.title = element_text(size = 12),
  axis.text.x = element_text(size = 8, angle = 90), axis.text.y = element_text(size = 8),
  axis.title = element_text(size = 10), plot.background = element_rect(fill = "Grey"))
```

Table 6: Results of Summary Dataframe 3

DRG_Cd	DRG_Desc	Hospitalizations
3	ECMO OR TRACHEOSTOMY WITH MV >96 HOURS OR PRINCIPAL DIAGNOSIS EXCEP	20
4	TRACH W MV >96 HRS OR PDX EXC FACE, MOUTH & NECK W/O MAJ O.R.	20
23	CRANIOTOMY W MAJOR DEVICE IMPLANT OR ACUTE COMPLEX CNS PDX W MCC OR	15
24	CRANIOTOMY WITH MAJOR DEVICE IMPLANT OR ACUTE COMPLEX CNS PRINCIPAL	14
25	CRANIOTOMY & ENDOVASCULAR INTRACRANIAL PROCEDURES W MCC	38
36	CAROTID ARTERY STENT PROCEDURES WITHOUT CC/MCC	13



Analysis of Plot 2 Results: Results indicate that Most hospitalization stays in Mobile,Alabama in 2020 was for *Extra cranial Procedures* followed by *Craniotomy & endovascular procedures*

Table 7: Chemotherapy treatment in West Virginia

year	DRG_Cd	state	city	total_Hospitalization	avg_Hospitalization
2018	847	WV	Charleston	50	50
2018	847	WV	Huntington	11	11
2018	847	WV	Morgantown	34	17
2019	847	WV	Charleston	22	22
2019	847	WV	Huntington	12	12
2019	847	WV	Morgantown	15	15

3.4.3 Plot 3: Chemotherapy treatment in West Virginia

```
# Creating a Summary dataframe from combined data based
# on total & average payments, total and average
# hospitalizations based on the variables Year, Medical
# procedure, State and City

prvsvc_comb_summary_df1 <- medic_prvsvc_combined_df %>%
  group_by(year, DRG_Cd, DRG_Desc, state, city) %>%
  dplyr::summarise(Total_payments = sum(Avg_Tot_Pymt_Amt),
    Average_payment = mean(Avg_Tot_Pymt_Amt), total_Hospitalization = sum(Tot_Dschrgrs),
    avg_Hospitalization = mean(Tot_Dschrgrs))

WV_Chemo_df <- prvsvc_comb_summary_df1 %>%
  filter(DRG_Cd == "847" & state == "WV") %>%
  topN(., "city")

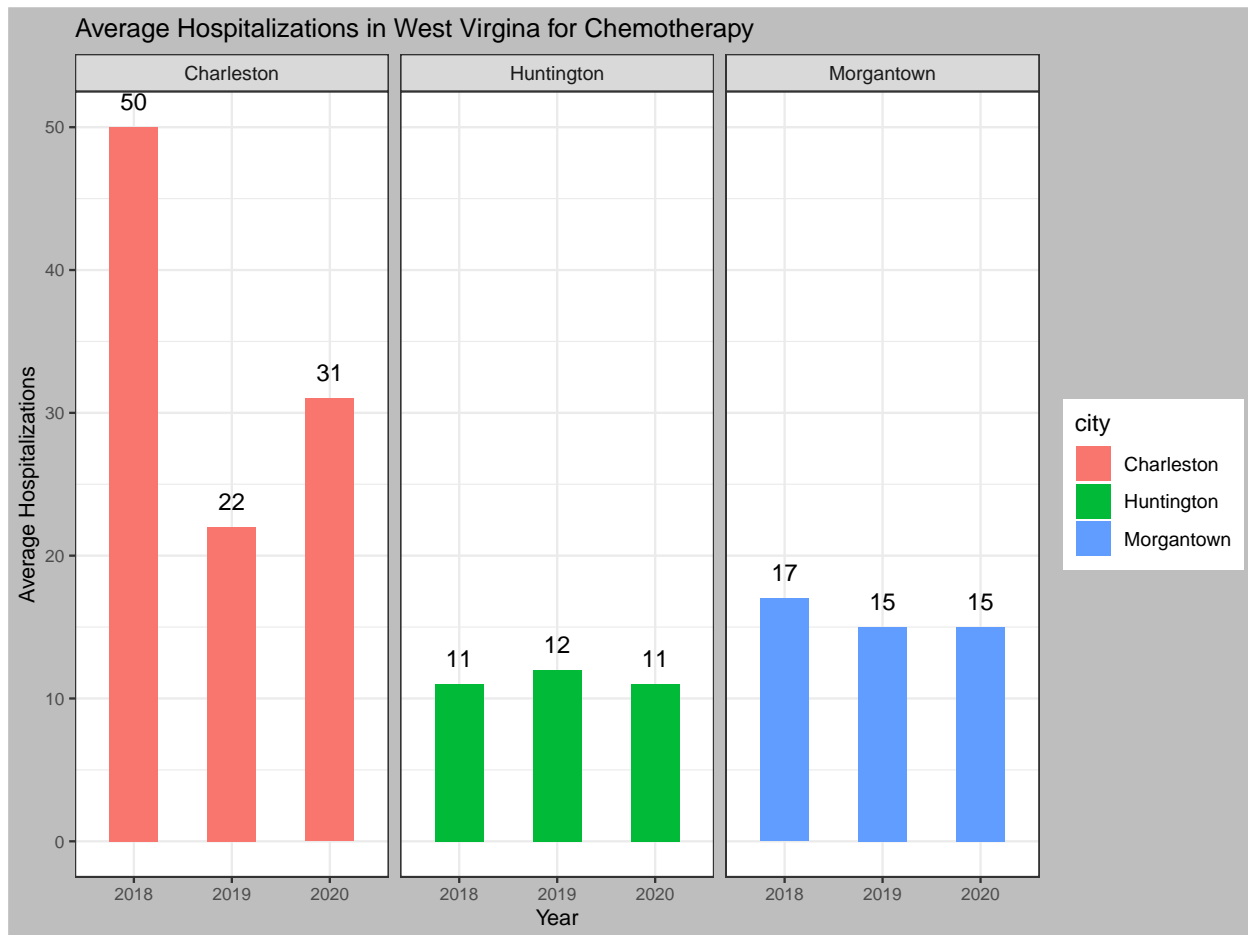
# Printing The results of the WV_Chemo_df
kbl(head(WV_Chemo_df[, c(1:2, 4:5, 8:9)]), caption = "Chemotherapy treatment in West Virginia",
  booktabs = T) %>%
  kable_styling(full_width = F, latex_options = c("striped"))

# Chemotherapy treatment in West Virginia

ggplot(WV_Chemo_df, aes(year, avg_Hospitalization)) + geom_bar(stat = "identity",
  width = 0.5, aes(fill = city)) + geom_text(aes(label = avg_Hospitalization),
  vjust = -1) + facet_wrap(~city) + theme_bw() + labs(x = "Year",
  y = "Average Hospitalizations", title = "Average Hospitalizations in West Virginia for Chemotherapy")
  theme(plot.title = element_text(size = 12), axis.text.x = element_text(size = 8),
```



```
axis.text.y = element_text(size = 8), axis.title = element_text(size = 10),
plot.background = element_rect(fill = "Grey"))
```



Analysis of Plot 3 Results: Results indicate most number of Hospitalization stays for *Chemotherapy treatment* in WV was in *Charleston* between 2018 to 2020 followed by *Morgantown*.

3.4.4 Plot 4: Hospitalizations in Kansas City, MO vs Kansas City,KS for Liver Ailments

```
# Liver related procedures in Kansas City MO and KS
kansasCity_df <- medic_prvsvc_combined_df %>%
  filter(city %in% "Kansas City" & state %in% c("MO", "KS") &
    grepl("Liver", medic_prvsvc_combined_df$DRG_Desc, ignore.case = TRUE)) %>%
  group_by(year, state) %>%
  dplyr::summarise(Tot_payment = sum(Avg_Tot_Pymt_Amt), Avg_payment = mean(Avg_Tot_Pymt_Amt),
    Total_discharges = sum(Tot_Dschrgrs), Avg_discharges = mean(Tot_Dschrgrs)) %>%
```

Table 8: Liver treatment in Kansas City

year	state	Tot_payment	Avg_payment	Total_discharges	Avg_discharges
2018	KS	256876.00	42812.67	265	44.16667
2018	MO	59911.77	19970.59	52	17.33333
2019	KS	226576.14	32368.02	227	32.42857
2019	MO	134934.12	26986.82	90	18.00000
2020	KS	186266.70	31044.45	194	32.33333
2020	MO	65865.98	21955.33	45	15.00000

```

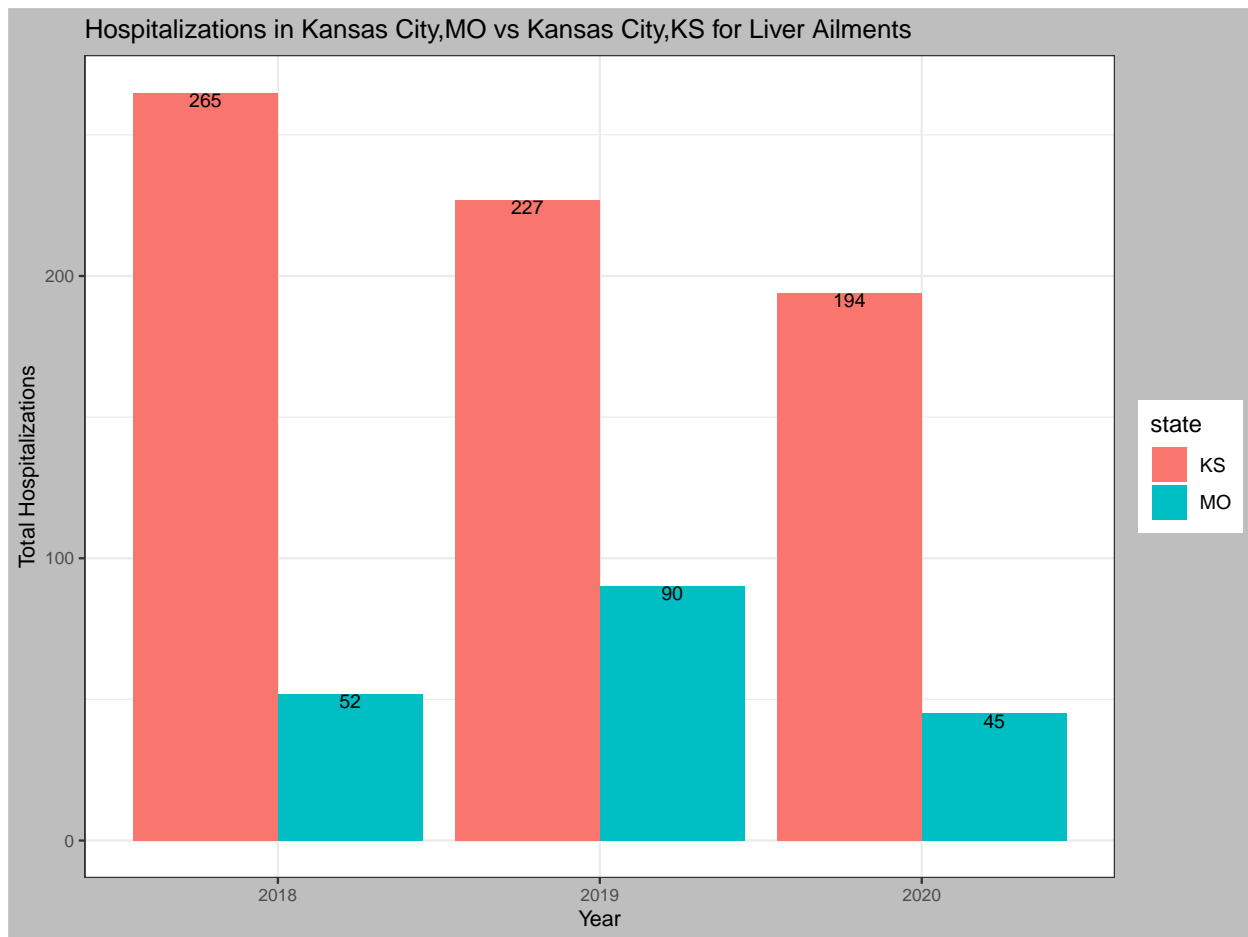
arrange(year)

# Printing The results of the kansasCity_df
kbl(head(kansasCity_df[, 1:6]), caption = "Liver treatment in Kansas City",
     booktabs = T) %>%
  kable_styling(full_width = F, latex_options = c("striped")) %>%
  column_spec(2, width = "15em")

# Hospitalizations in Kansas City, MO vs KS for Liver
# Ailments
ggplot(kansasCity_df, aes(year, Total_discharges, fill = state)) +
  geom_col(position = "dodge") + geom_text(aes(label = Total_discharges,
  fill = state), color = "black", size = 3, vjust = 1, position = position_dodge(0.9)) +
  labs(x = "Year", y = "Total Hospitalizations", title = "Hospitalizations in Kansas City, MO vs Kansas City, KS") +
  theme_bw() + theme(plot.title = element_text(size = 12),
  axis.text.x = element_text(size = 8), axis.text.y = element_text(size = 8),
  axis.title = element_text(size = 10), plot.background = element_rect(fill = "Grey"))

## Warning in geom_text(aes(label = Total_discharges, fill = state), color =
## "black", : Ignoring unknown aesthetics: fill

```



Analysis of Plot 4 Results: Results indicate that most of the patients went *chose KansasCity, KS hospitals than Kansascity, MO hospitals* for inpatient services. Most research to be done to analyze the cost of treatments across these states to understand the pattern.

3.4.5 Plot 5: National Average vs Wyoming for Cost of Heart related Ailments treatment

```
# National Average vs Treatment costs in Wyoming for
# Heart related Ailments

wyoming_heart_df <- medic_geo_combined_df %>%
  filter(year %in% c("2020", "2019", "2018") & Geo_desc %in%
    c("National", "Wyoming") & DRG_Cd %in% c("269", "291",
    "292", "293")) %>%
  select(DRG_Cd, DRG_Desc, Geo_desc, Avg_Tot_Pymt_Amt, year)
# Using Spread function to create new columns for
```

Table 9: National Average vs Wyoming for Cost of Heart related Ailments treatment

DRG_Cd	DRG_Desc	year	National	Wyoming
269	AORTIC AND HEART ASSIST PROCEDURES EXCEPT PULSATION BALLOON W/O MCC	2018	31366.63	41865.42
269	AORTIC AND HEART ASSIST PROCEDURES EXCEPT PULSATION BALLOON W/O MCC	2019	31366.63	41865.42
269	AORTIC AND HEART ASSIST PROCEDURES EXCEPT PULSATION BALLOON W/O MCC	2020	33384.54	41676.46
291	HEART FAILURE & SHOCK W MCC	2018	10538.62	12875.81
291	HEART FAILURE & SHOCK W MCC	2019	10538.62	12875.81
291	HEART FAILURE & SHOCK W MCC	2020	10918.66	13140.10

```
# National average and Wyoming
wyoming_heart_df_spread <- wyoming_heart_df %>%
  spread(key = Geo_desc, value = Avg_Tot_Pymt_Amt)

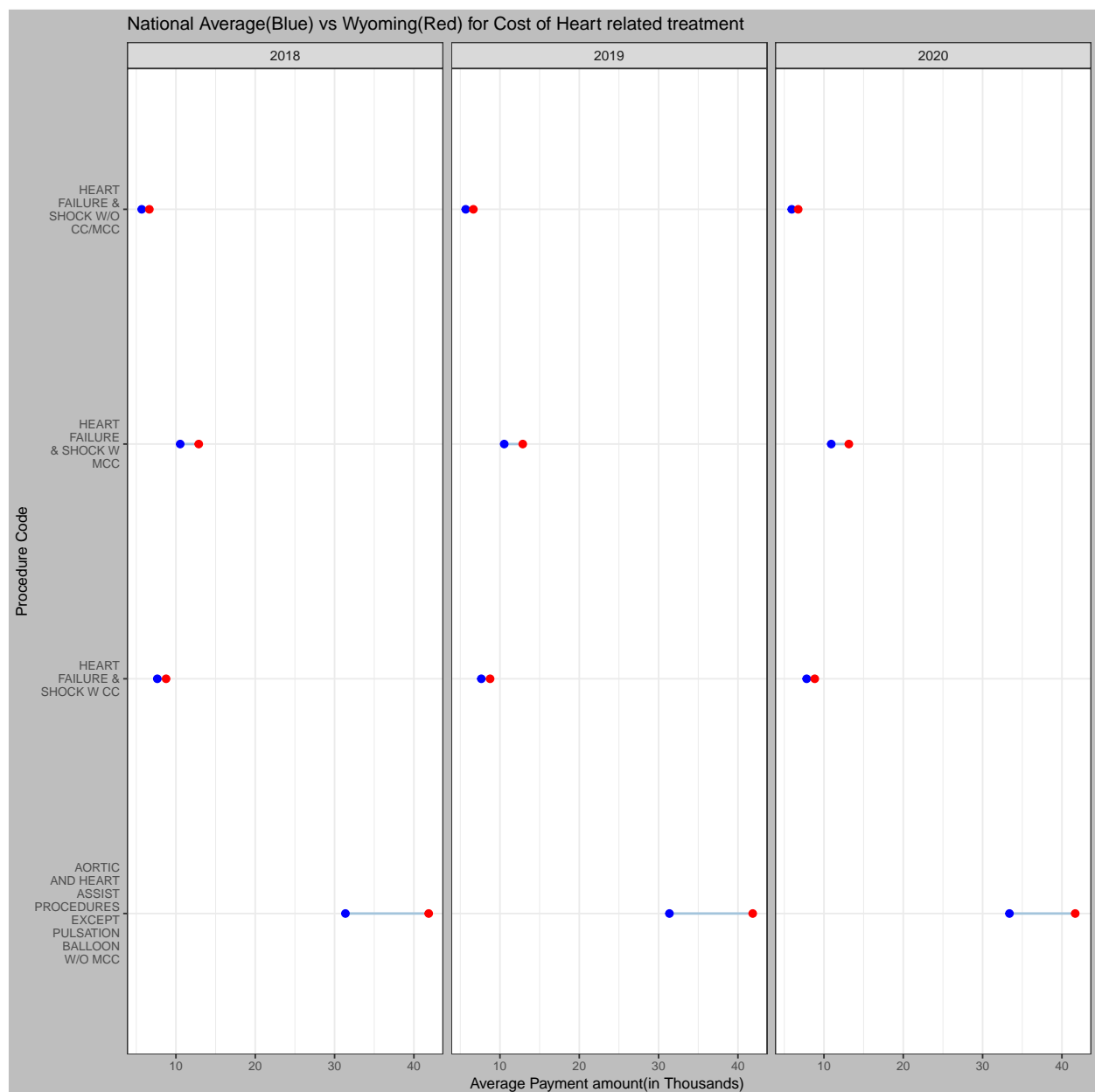
# Printing The results of the wyoming_heart_df_spread
kbl(head(wyoming_heart_df_spread), caption = "National Average vs Wyoming for Cost of Heart related Ailments",
  booktabs = T) %>%
  kable_styling(full_width = F, latex_options = c("striped")) %>%
  column_spec(2, width = "15em")
```

```
# National Average vs Wyoming for Heart related Ailments
ggplot(wyoming_heart_df_spread, aes(x = National/1000, xend = Wyoming/1000,
  y = DRG_Desc, group = DRG_Cd)) + geom_dumbbell(color = "#a3c4dc",
  size = 0.75, colour_x = "Blue", colour_xend = "Red", size_x = 2,
  size_xend = 2) + facet_wrap(~year) + labs(x = "Average Payment amount(in Thousands) ",
  y = "Procedure Code", title = "National Average(Blue) vs Wyoming(Red) for Cost of Heart related treatments",
  color = "Legend") + theme_bw() + theme(plot.title = element_text(size = 12),
```

```
axis.text.x = element_text(size = 8), axis.text.y = element_text(size = 8),
axis.title = element_text(size = 10), plot.background = element_rect(fill = "Grey"),
legend.position = "bottom") + scale_y_discrete(labels = function(x) str_wrap(x,
width = 10)) + scale_color_manual(name = "Regression Model",
breaks = c("Linear", "Quadratic", "Cubic"), values = c(Cubic = "pink",
Quadratic = "blue", Linear = "purple"))
```

Warning: Using the 'size' aesthetic with geom_segment was deprecated in ggplot2 3.4.0.

i Please use the 'linewidth' aesthetic instead.



Analysis of Plot 5 Results: Results indicate that the Average cost of heart related treatment in *Wyoming* is slightly above the national average except for the *Aortic and heart assist* procedures.

3.4.6 Plot 6 : Top 10 Hospitalizations in each year across the US

```
# Top 10 Hospitalizations in each year across the country

Top10_hosp_geo_df_2020 <- medic_geo_combined_df %>%
  filter(Geo_desc == "National" & year %in% c("2020")) %>%
  arrange(desc(Tot_Dschrgs)) %>%
  head(n = 10)

Top10_hosp_geo_df_2019 <- medic_geo_combined_df %>%
  filter(Geo_desc == "National" & year %in% c("2019")) %>%
  arrange(desc(Tot_Dschrgs)) %>%
  head(n = 10)

Top10_hosp_geo_df_2018 <- medic_geo_combined_df %>%
  filter(Geo_desc == "National" & year %in% c("2018")) %>%
  arrange(desc(Tot_Dschrgs)) %>%
  head(n = 10)

# Creating Top10_hosp_geo_df by combining data from above
# 3 Dataframes

Top10_hosp_geo_df <- rbind(Top10_hosp_geo_df_2020, Top10_hosp_geo_df_2019,
  Top10_hosp_geo_df_2018) %>%
  select(DRG_Cd, DRG_Desc, Tot_Dschrgs, year)

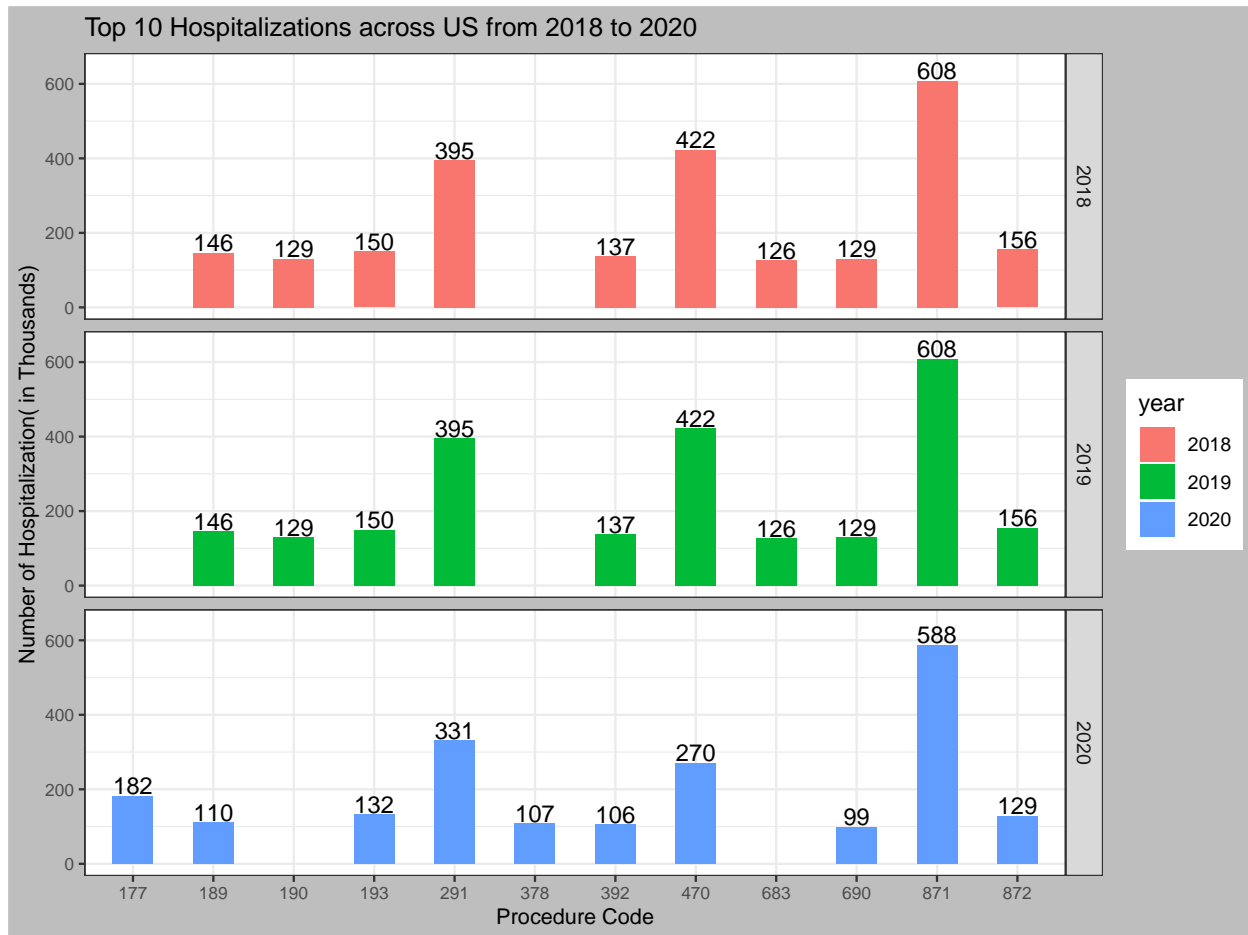
# Printing The results of the Top10_hosp_geo_df

kbl(head(Top10_hosp_geo_df), caption = "Top 10 Hospitalizations in each year across U.S",
  booktabs = T) %>%
  kable_styling(full_width = F, latex_options = c("striped")) %>%
  column_spec(2, width = "15em")
```

Table 10: Top 10 Hospitalizations in each year across U.S

DRG_Cd	DRG_Desc	Tot_Dschrgrs	year
871	SEPTICEMIA OR SEVERE SEPSIS W/O MV >96 HOURS W MCC	587611	2020
291	HEART FAILURE & SHOCK W MCC	331381	2020
470	MAJOR HIP AND KNEE JOINT REPLACEMENT OR REATTACHMENT OF LOWER EXTREM	270402	2020
177	RESPIRATORY INFECTIONS & INFLAMMATIONS W MCC	181981	2020
193	SIMPLE PNEUMONIA & PLEURISY W MCC	131572	2020
872	SEPTICEMIA OR SEVERE SEPSIS W/O MV >96 HOURS W/O MCC	128651	2020

```
# Top 10 Hospitalizations in each year across the U.S
ggplot(Top10_hosp_geo_df, aes(DRG_Cd, Tot_Dschrgrs/1000)) +
  geom_bar(stat = "Identity", width = 0.5, aes(fill = year)) +
  labs(x = "Procedure Code", y = "Number of Hospitalization( in Thousands)",
       title = "Top 10 Hospitalizations across US from 2018 to 2020") +
  facet_grid(year ~ .) + ylim(0, 650) + theme_bw() + theme(plot.title = element_text(size = 12),
axis.text.x = element_text(size = 8), axis.text.y = element_text(size = 8),
axis.title = element_text(size = 10), plot.background = element_rect(fill = "Grey")) +
  geom_text(aes(label = (round(Tot_Dschrgrs/1000)), vjust = -0.1))
```

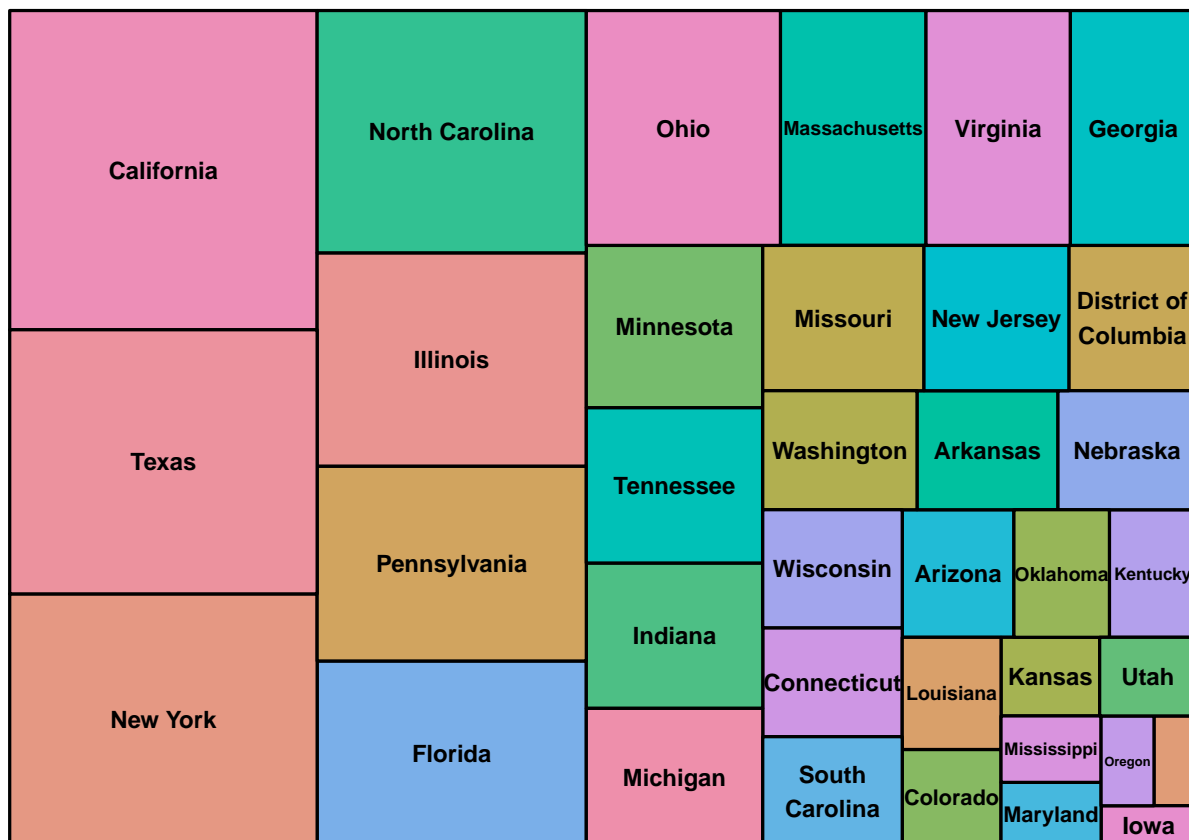


Analysis of Plot 6 Results: Results are shown based on the Procedure code due to space constraints. The table can be used to identify the name of the Procedure description which indicates that maximum hospitalization across US from 2018 to 2020 was for **SEPTICEMIA OR SEVERE SEPSIS** followed by **MAJOR HIP AND KNEE JOINT REPLACEMENT** and **HEART FAILURES**

3.4.7 Plot 7: Creating a Treemap of Hospitalizations in the US for Heart Transplant Procedures from 2018 to 2020

```
# Treemap of Hospitalizations in the US for Heart
# Transplant Procedures from 2018 to 2020
treemap(medic_geo_combined_df %>%
  filter(Geo_lvl == "State" & DRG_Cd == "1"), index = "Geo_desc",
  vSize = "Tot_Dschrgs", type = "index", title = "Treemap of Hospitalizations in US for Heart Transpl",
  fontsize.title = 18)
```


Treemap of Hospitalizations in US for Heart Transplant



Analysis of Plot 7 Results: The Treemap indicates that most hospitalizations in the US for Heart transplant was in *California* followed by *Texas*, *New York*, *North Carolina*, *Illinois*, *Pennsylvania* and *others*. Surprisingly, *Washington DC* with population less than a million has more hospitalizations than states like *Colorado*, *Arizona*, *South Carolina* which are more populated than DC.

3.4.8 Plot 8: Most Profitable hospitals in Nebraska

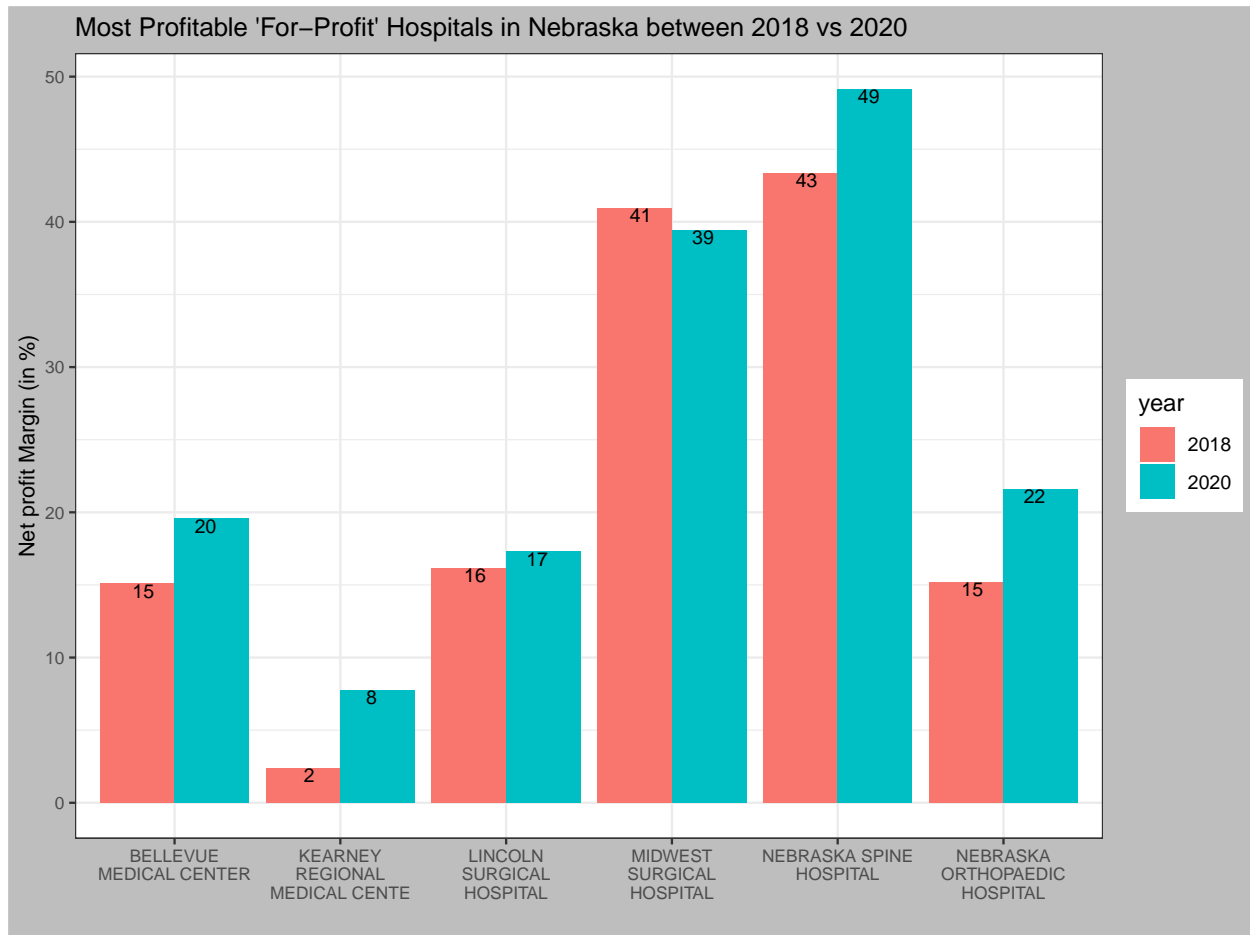
```
nebraska_hosp_df <- NASHP_df %>%
  filter(state == "NE" & hosp_ownership == "For-Profit" &
    year %in% c("2018", "2020") & net_income > 0) %>%
  group_by(year) %>%
  arrange(year, desc(net_profit_margin)) %>%
  select(CCN, year, hosp_name, city, state, health_system,
    net_income, net_profit_margin)
```

Table 11: Most Profitable 'For-Profit' Hospitals in Nebraska between 2018 vs 2020

year	hosp_name	city	state	net_income	net_profit_margin
2018	NEBRASKA SPINE HOSPITAL	OMAHA	NE	19942302	43.33
2018	MIDWEST SURGICAL HOSPITAL	OMAHA	NE	21566392	40.95
2018	LINCOLN SURGICAL HOSPITAL	LINCOLN	NE	10025338	16.14
2018	NEBRASKA ORTHOPAEDIC HOSPITAL LLC	OMAHA	NE	12222977	15.21
2018	BELLEVUE MEDICAL CENTER	BELLEVUE	NE	10999128	15.09
2018	KEARNEY REGIONAL MEDICAL CENTER	KEARNEY	NE	2831353	2.40

```
# Printing The results of the nebraska_hosp_df
kbl(head(nebraska_hosp_df[1:6, c(2:5, 7:8)]), caption = "Most Profitable 'For-Profit' Hospitals in Nebraska",
      booktabs = T) %>%
  kable_styling(full_width = F, latex_options = c("striped"))
```

```
# Most Profitable hospitals in Nebraska
ggplot(nebraska_hosp_df, aes(str_wrap(str_sub(hosp_name, 1,
  30), 15), net_profit_margin, fill = year)) + geom_col(position = "dodge") +
  geom_text(aes(label = round(net_profit_margin)), color = "black",
    size = 3, vjust = 1, position = position_dodge(0.75)) +
  labs(x = "", y = "Net profit Margin (in %)", title = "Most Profitable 'For-Profit' Hospitals in Nebraska",
    theme_bw() + theme(plot.title = element_text(size = 12),
    axis.text.x = element_text(size = 8), axis.text.y = element_text(size = 8),
    axis.title = element_text(size = 10), plot.background = element_rect(fill = "Grey"))
```



Analysis of Plot 8 Results: Results indicate *Nebraska Spine Hospital* in Omaha is the most profitable across the state followed by *Midwest Surgical hospital* in Omaha (again!) for the period 2018 and 2020. Also these results indicate the most profitable hospitals are located on the *eastern part of Nebraska* which is more populated than other regions of the state(except Kearney, NE).

3.4.9 Plot 9: Net Profit margin vs payments per patient 'for-Profit' hospital in Anchorage, Alaska

```
# How many 'For-profit' Hospitals operate in Alaska
NASHP_df %>%
  filter(state == "AK" & hosp_ownership == "For-Profit") %>%
  distinct(CCN, city)
```

```
## # A tibble: 3 x 2
```

```
##   CCN      city
```

```
##   <fct>  <chr>
## 1 020006 PALMER
## 2 020017 ANCHORAGE
## 3 021307 CORDOVA
```

```
# checking the profit margin of 'For-Profit' hospital in
# Anchorage, Alaska
NASHP_Anchorage_df <- NASHP_df %>%
  filter(CCN %in% c("020017") & year %in% c("2018", "2019",
    "2020")) %>%
  select(CCN, hosp_name, city, state, year, net_income, net_profit_margin)

# Extracting the results from other dataframe
# medic_prvsvc_combined_df for the same Anchorage
# hospital 020017
medic_prvsvc_Anchorage_df <- medic_prvsvc_combined_df %>%
  filter(CCN %in% c("020017")) %>%
  group_by(CCN, year) %>%
  dplyr::summarise(total_payment = sum(Avg_Tot_Pymt_Amt),
    total_discharges = sum(Tot_Dschrgs), payment_per_patient = total_payment/total_discharges)
medic_prvsvc_Anchorage_df
```

```
## # A tibble: 3 x 5
## # Groups:   CCN [1]
##   CCN      year total_payment total_discharges payment_per_patient
##   <chr>  <fct>      <dbl>          <int>          <dbl>
## 1 020017 2018      990021.         1465           676.
## 2 020017 2019     1347698.         1931           698.
## 3 020017 2020     1302531.         1526           854.
```

Table 12: Most Profitable Hospitals in Anchorage Alaska between 2018 to 2020

CCN	year	hosp_name	city	state	net_income	total_discharges	Type
020017	2018	ALASKA REGIONAL HOSPITAL	ANCHORAGE	AK	62828656	1465	net_profit_margin
020017	2019	ALASKA REGIONAL HOSPITAL	ANCHORAGE	AK	96080416	1931	net_profit_margin
020017	2020	ALASKA REGIONAL HOSPITAL	ANCHORAGE	AK	123580707	1526	net_profit_margin
020017	2018	ALASKA REGIONAL HOSPITAL	ANCHORAGE	AK	62828656	1465	payment_per_patient
020017	2019	ALASKA REGIONAL HOSPITAL	ANCHORAGE	AK	96080416	1931	payment_per_patient
020017	2020	ALASKA REGIONAL HOSPITAL	ANCHORAGE	AK	123580707	1526	payment_per_patient

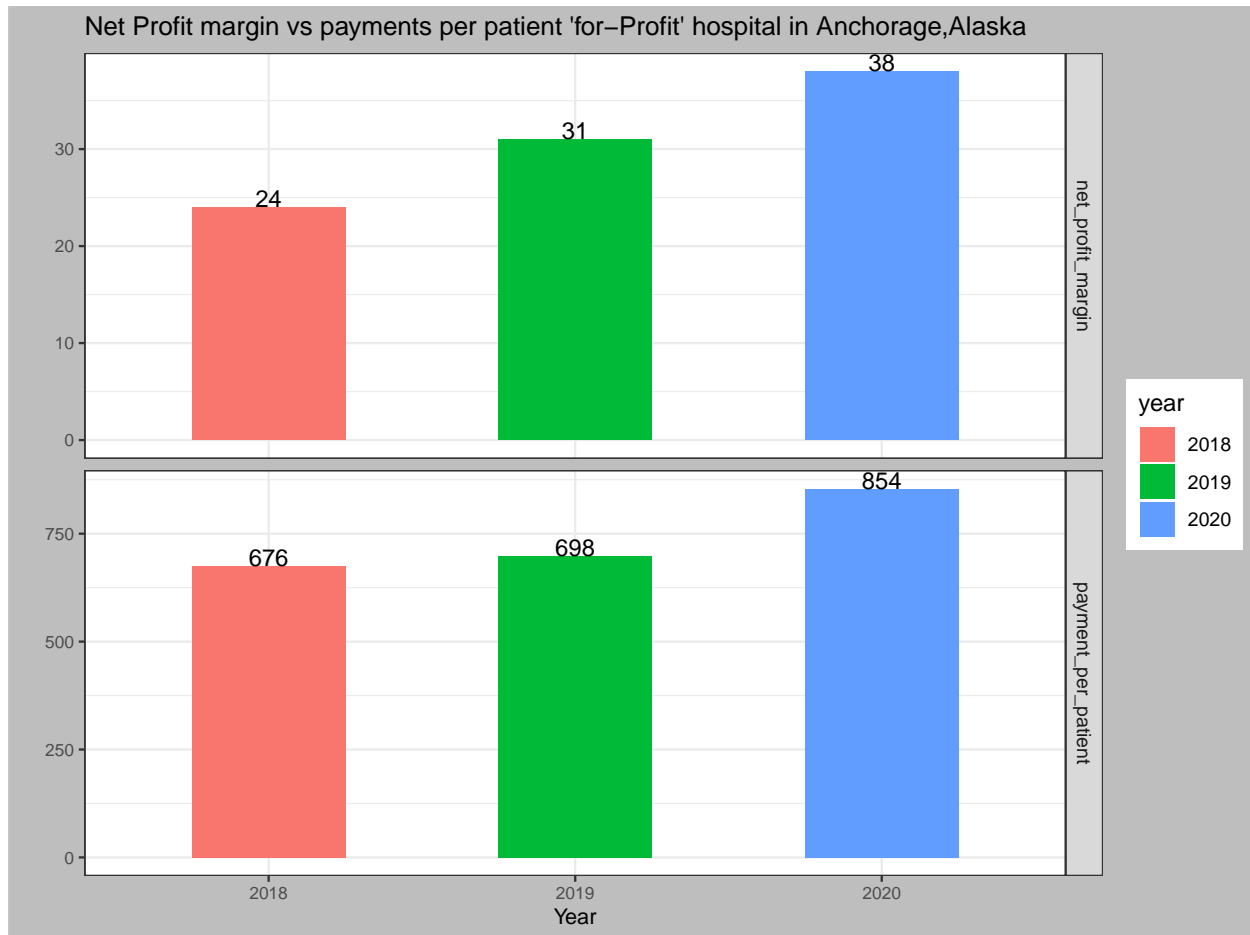
```
# Joining both the dataframes based on the Hospital code
# and year

Anchorage_df <- merge(x = NASHP_Anchorage_df, y = medic_prvsvc_Anchorage_df,
  by = c("CCN", "year")) %>%
  select(-total_payment) %>%
  gather(key = Type, value = Stats, net_profit_margin, payment_per_patient)
Anchorage_df$Stats <- round(Anchorage_df$Stats)

# Printing The results of the Anchorage_df
kbl(head(Anchorage_df), caption = "Most Profitable Hospitals in Anchorage Alaska between 2018 to 2020",
  booktabs = T) %>%
  kable_styling(full_width = F, latex_options = c("striped")) %>%
  column_spec(c(3, 8), width = "10em")
```

```
ggplot(Anchorage_df, aes(x = year, y = Stats, fill = year)) +
  geom_bar(stat = "identity", position = "dodge", width = 0.5) +
  facet_grid(Type ~ ., scales = "free_y") + labs(x = "Year",
  y = "", title = "Net Profit margin vs payments per patient 'for-Profit' hospital in Anchorage, Alaska")
  geom_text(aes(label = Stats), vjust = 0) + theme_bw() +
  theme(plot.title = element_text(size = 12), axis.text.x = element_text(size = 8),
```

```
axis.text.y = element_text(size = 8), axis.title = element_text(size = 10),
plot.background = element_rect(fill = "Grey"))
```



Analysis of Plot 9 Results: The results indicate same pattern in net profit margin vs Payments made by Medicare per patient between 2018 to 2020 for the *Alaska regional Hospital*. Also the table indicates that the hospital *treated lesser patients* in 2020 than in 2019, but *made more profit* than in 2019. More analysis is required to understand the reasoning behind it.

3.4.10 Plot 10: Comparing Cost of Spinal Procedure - National vs State Average

```
spinal_2020_national_avg <- medic_geo_combined_df %>%
  filter(Geo_lvl == "National" & DRG_Cd == "28" & year ==
    "2020") %>%
  select(Avg_Tot_Pymt_Amt)
spinal_2020_df <- medic_geo_combined_df %>%
```

Table 13: Cost of Spinal Procedure- National vs State Average

	Geo_desc	Avg_Tot_Pymt_Amt	difference	comparePrice
7	Arkansas	35743.00	-12830.020	below
24	West Virginia	39427.00	-9146.020	below
38	South Carolina	39695.48	-8877.543	below
15	Mississippi	39941.20	-8631.820	below
42	Iowa	40014.26	-8558.759	below
23	Nevada	40064.00	-8509.020	below

```

filter(Geo_lvl == "State" & DRG_Cd == "28" & year == "2020") %>%
select(Geo_desc, Avg_Tot_Pymt_Amt) %>%
mutate(difference = (Avg_Tot_Pymt_Amt - as.numeric(spinal_2020_national_avg)))
spinal_2020_df$comparePrice <- ifelse(spinal_2020_df$difference <
  0, "below", "above")
spinal_2020_df <- spinal_2020_df[order(spinal_2020_df$difference),
  ]

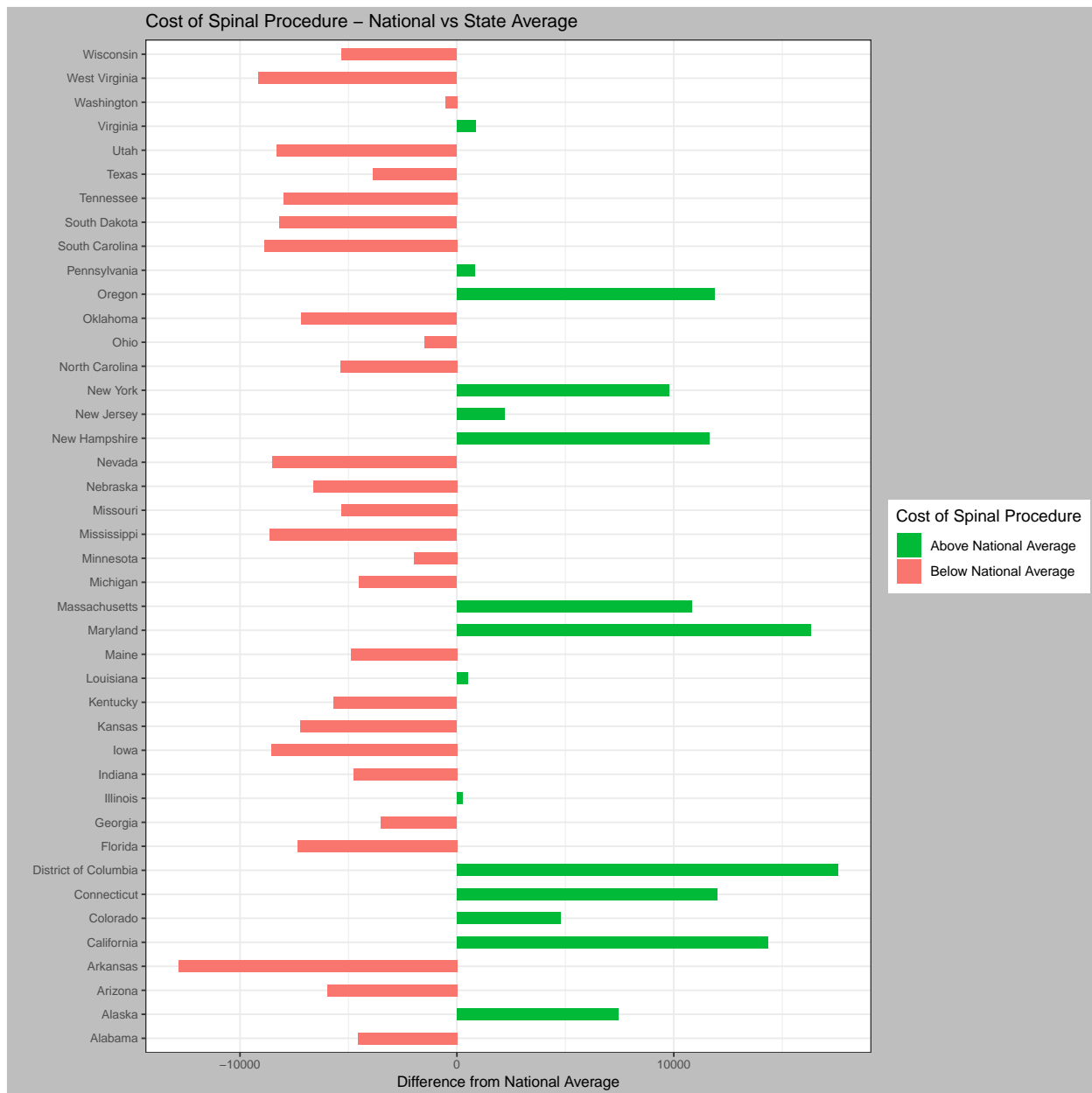
kbl(head(spinal_2020_df), caption = "Cost of Spinal Procedure- National vs State Average",
  booktabs = T) %>%
kable_styling(full_width = F, latex_options = c("striped"))

```

```

ggplot(spinal_2020_df, aes(x = Geo_desc, y = difference)) +
  geom_bar(stat = "identity", aes(fill = comparePrice), width = 0.5) +
  scale_fill_manual(name = "Cost of Spinal Procedure", labels = c("Above National Average",
    "Below National Average"), values = c(above = "#00ba38",
    below = "#f8766d")) + coord_flip() + theme_bw() + theme(plot.title = element_text(size = 12),
  axis.text.x = element_text(size = 8), axis.text.y = element_text(size = 8),
  axis.title = element_text(size = 10), plot.background = element_rect(fill = "Grey")) +
  labs(x = "", y = "Difference from National Average", title = "Cost of Spinal Procedure - National vs State Average")

```



Analysis of Plot 10 Results: The plot indicates that the Spine procedures are *most expensive in Washington D.C, followed by Maryland and California*, while *Arkansas and West Virginia are among the cheapest* in the nation.

3.4.11 Plot 11: Inpatient stays for Chemotherapy between 2018-2020

```
# Inpatient stays for Chemotherapy between 2018-2020
```

```
chemo_states_df <- medic_geo_combined_df %>%
```



```

filter(Geo_lvl == "State" & grepl("Chemo", DRG_Desc, ignore.case = TRUE)) %>%
group_by(Geo_desc) %>%

dplyr::summarise(tot_dschrgrs = sum(Tot_Dschrgrs)) %>%

select(state = Geo_desc, tot_dschrgrs)

chemo_states_df

```

```

## # A tibble: 50 x 2
##   state          tot_dschrgrs
##   <chr>          <int>
## 1 Alabama          691
## 2 Alaska           95
## 3 Arizona        1027
## 4 Arkansas         669
## 5 California     6483
## 6 Colorado         995
## 7 Connecticut     755
## 8 Delaware        146
## 9 District of Columbia 473
## 10 Florida       2461
## # ... with 40 more rows

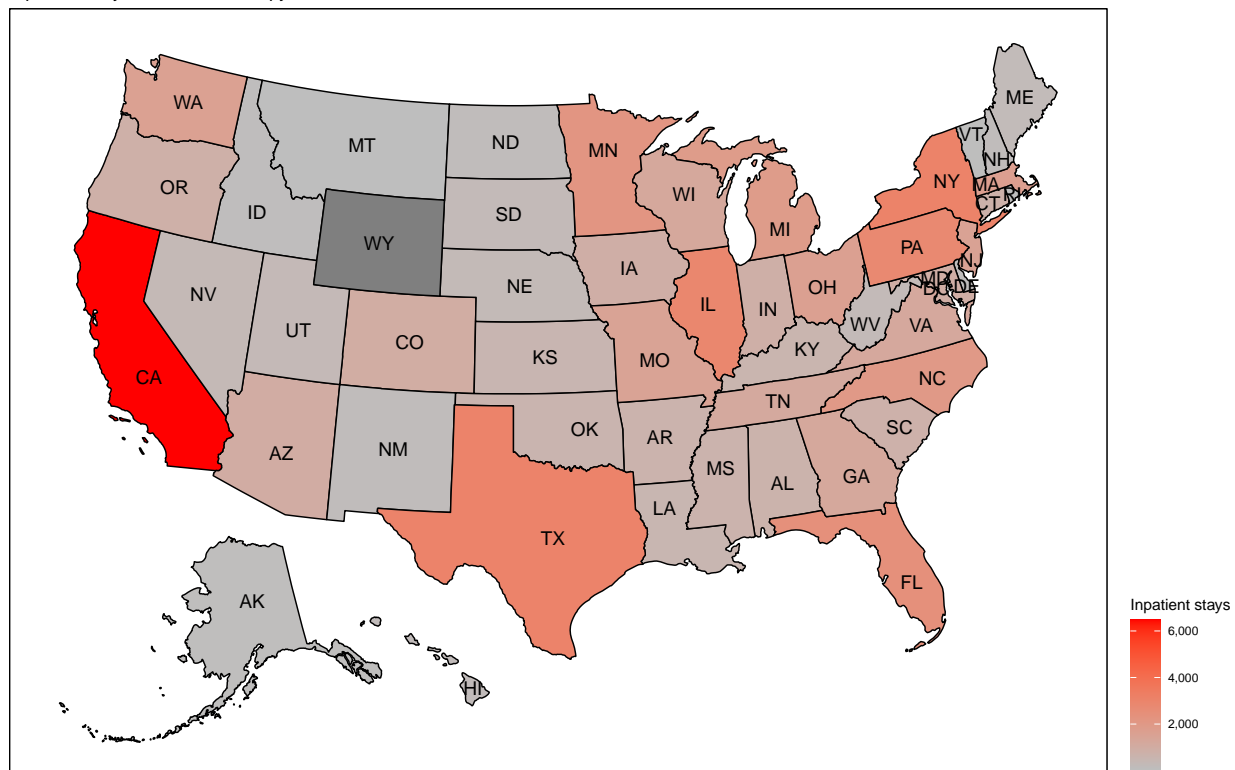
```

```

plot_usmap(data = chemo_states_df, values = "tot_dschrgrs",
  color = "black", labels = TRUE) + scale_fill_continuous(name = "Inpatient stays",
  label = scales::comma, low = "grey", high = "red") + theme(legend.position = "right",
  panel.background = element_rect(colour = "black")) + ggtitle("Inpatient stays for Chemotherapy betw

```

Inpatient stays for Chemotherapy between 2018–2020



Analysis of Plot 11 Results: The Map indicates that *California has maximum number of Inpatient stays* for Chemotherapy *followed by Texas, Illinois, Pennsylvania, New York and Florida*

3.4.12 Plot 12: How profitable are hospitals between 2011 and 2020 in the Southern states?

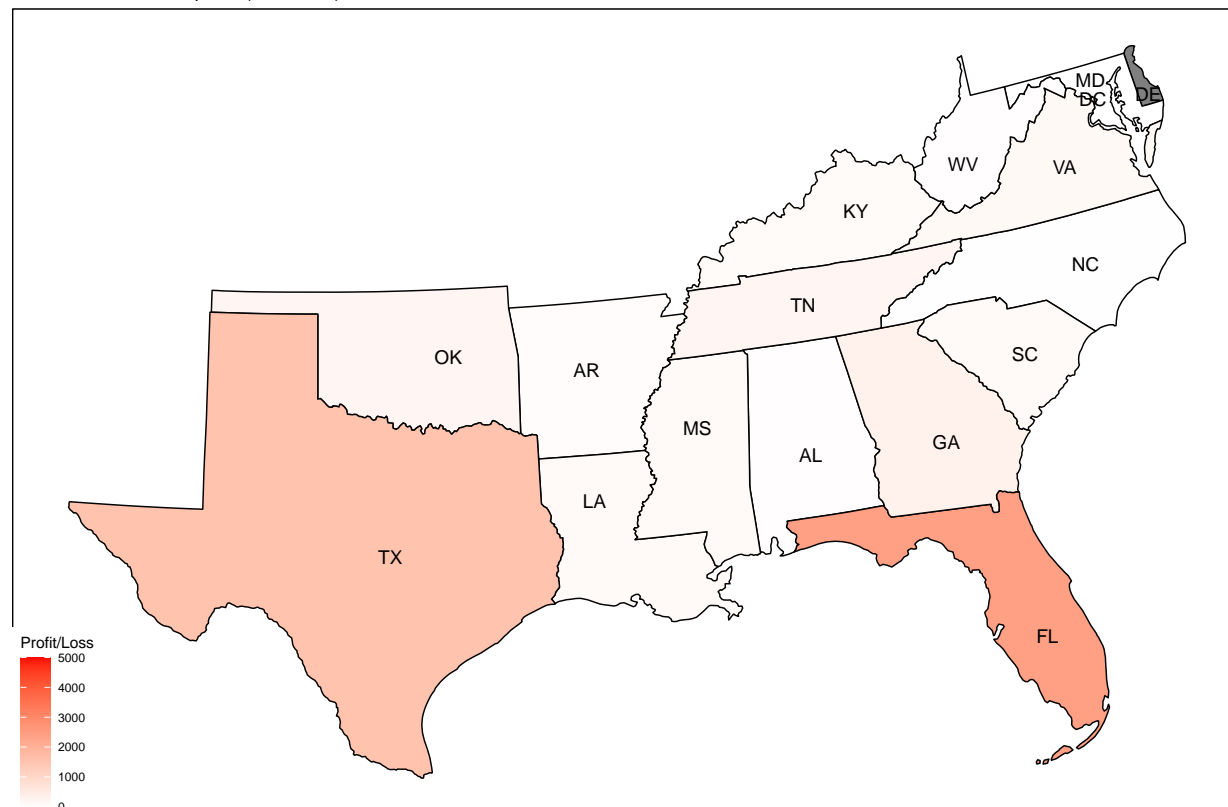
```
NAS_2011_PnL_df <- NASHP_df %>%
  filter(hosp_ownership == "For-Profit" & year == "2011") %>%
  group_by(state) %>%
```

```

    dplyr::summarise(Total_PnL = sum(net_income)/1e+06)
NAS_2011_PnL_df$Total_PnL <- ifelse(NAS_2011_PnL_df$Total_PnL <
    0, 0, NAS_2011_PnL_df$Total_PnL)
NAS_2020_PnL_df <- NASHP_df %>%
    filter(hosp_ownership == "For-Profit" & year == "2020") %>%
    group_by(state) %>%
    dplyr::summarise(Total_PnL = sum(net_income)/1e+06)
NAS_2020_PnL_df$Total_PnL <- ifelse(NAS_2020_PnL_df$Total_PnL <
    0, 0, NAS_2020_PnL_df$Total_PnL)
plot2011 <- plot_usmap(data = NAS_2011_PnL_df, include = .south_region,
    values = "Total_PnL", labels = TRUE) + scale_fill_continuous(low = "white",
    high = "red", name = "Profit/Loss", limits = c(0, 5000)) +
    labs(title = "Profit/Loss of US Hospitals (in Millions) in 2011 in the Southern states") +
    theme(panel.background = element_rect(colour = "black"))
plot2011

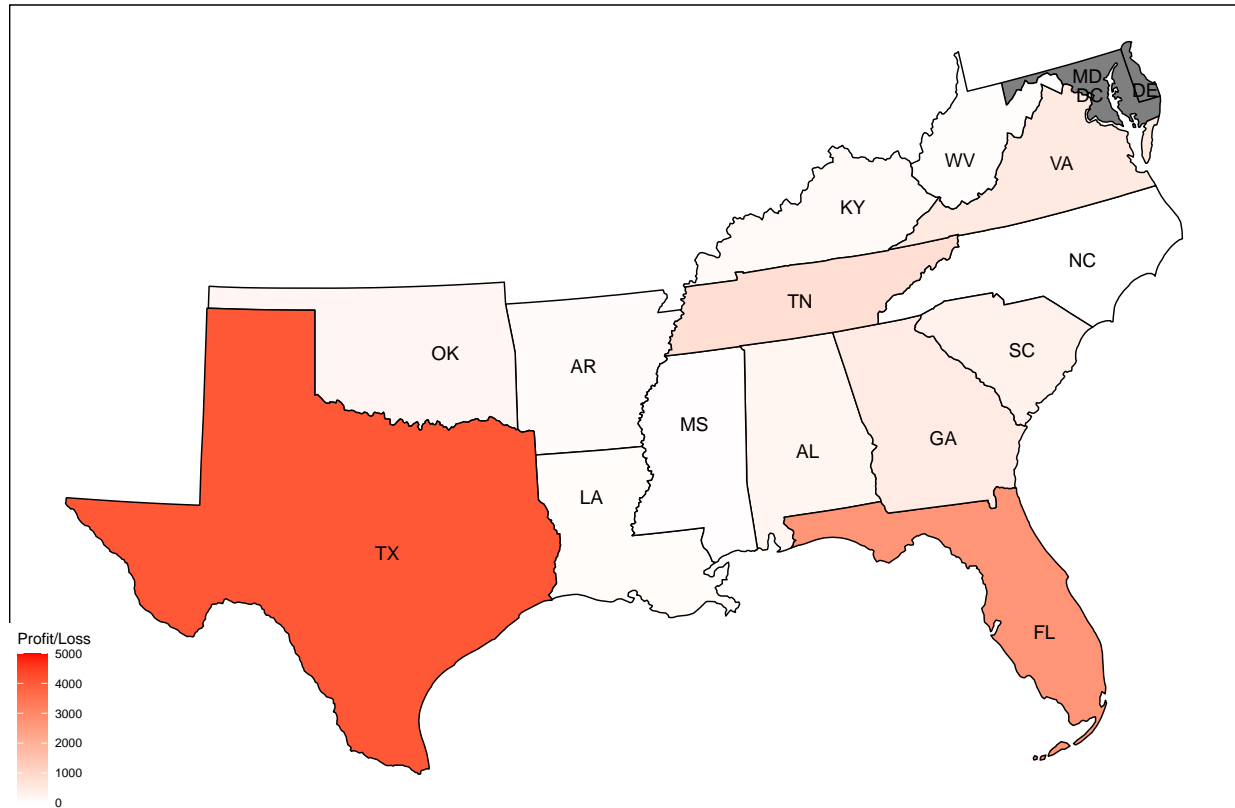
```

Profit/Loss of US Hospitals (in Millions) in 2011 in the Southern states



```
plot2020 <- plot_usmap(data = NAS_2020_PnL_df, include = .south_region,
  values = "Total_PnL", labels = TRUE) + scale_fill_continuous(low = "white",
  high = "red", name = "Profit/Loss", limits = c(0, 5000)) +
  labs(title = "Profit/Loss of US Hospitals (in Millions) in 2020 in the Southern states") +
  theme(panel.background = element_rect(colour = "black"))
plot2020
```

Profit/Loss of US Hospitals (in Millions)) in 2020 in the Southern states



Analysis of Plot 12 Results: The Map shows that, in the southern region, *Florida was most profitable* for Hospitals in 2011 in terms of the total net income earned, but *Texas became the most profitable state in 2020*, followed by Florida ,Tennessee and Virginia.

3.4.13 Plot 13: Comparing Average income of “For-Profit” vs “Non-Profit” hospitals in 2019

```
NASHP_profit_nonprofit_df <- NASHP_df %>%
  filter(hosp_ownership %in% c("For-Profit", "Non-Profit") &
```

```

    year %in% c("2019"))

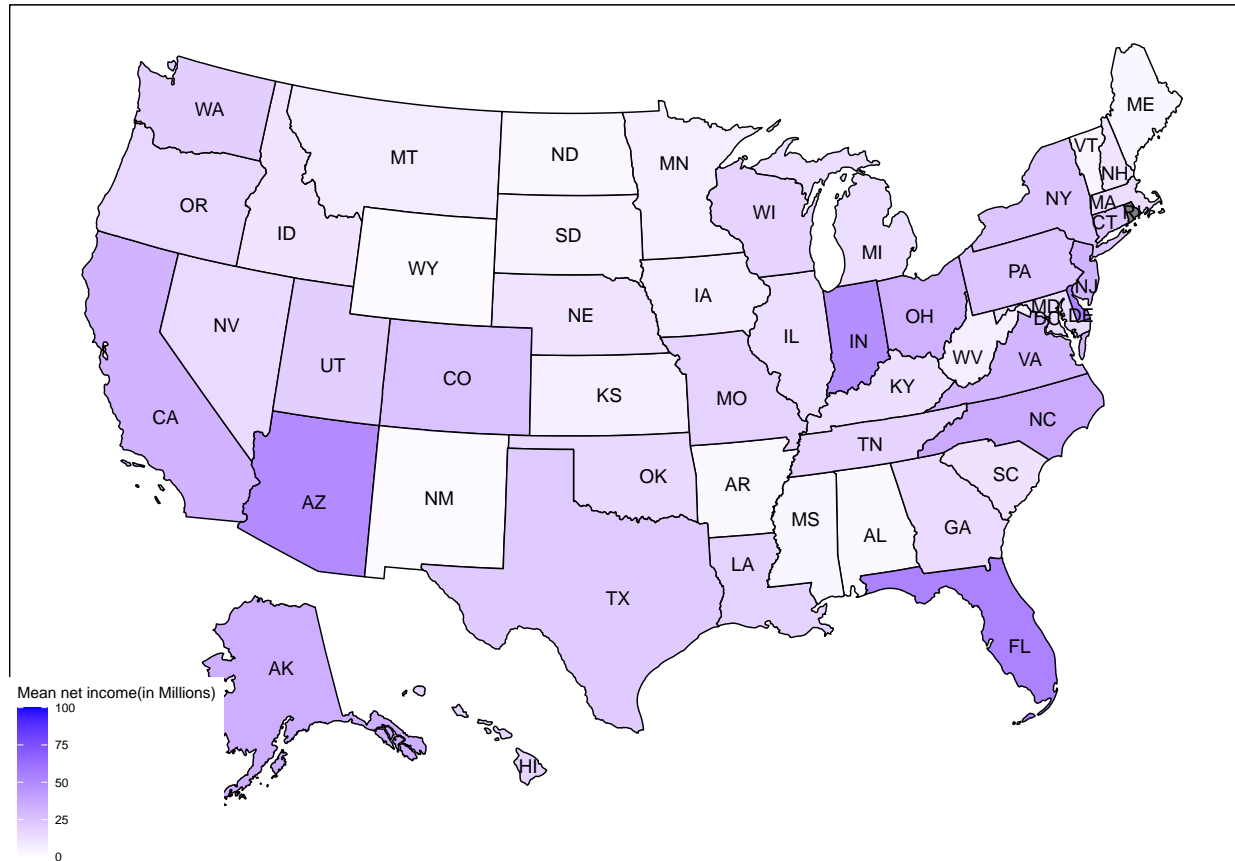
NASHP_profit_nonprofit_sum1_df <- NASHP_profit_nonprofit_df %>%
  group_by(state, hosp_ownership) %>%
  dplyr::summarise(mean_net_income = mean(net_income/1e+06))

non_profit_2019_plot <- plot_usmap(data = NASHP_profit_nonprofit_sum1_df[NASHP_profit_nonprofit_sum1_df$
  "Non-Profit", ], values = "mean_net_income", labels = TRUE) +
  scale_fill_continuous(low = "white", high = "blue", name = "Mean net income(in Millions)",
    limits = c(0, 100)) + labs(title = "Mean Net income of Non-profit hospitals (in Millions) in 20
  theme(panel.background = element_rect(colour = "black"))

non_profit_2019_plot

```

Mean Net income of Non-profit hospitals (in Millions) in 2019



```
profit_2019_plot <- plot_usmap(data = NASHP_profit_nonprofit_sum1_df[NASHP_profit_nonprofit_sum1_df$hosp_type ==
  "For-Profit", ], values = "mean_net_income", labels = TRUE) +
  scale_fill_continuous(low = "white", high = "blue", name = "Mean net income(in Millions)",
    limits = c(0, 100)) + labs(title = "Mean Net income of For-profit hospitals (in Millions) in 2019") +
  theme(panel.background = element_rect(colour = "black"))
profit_2019_plot
```

Mean net income(in Millions)

100
75
50
25
0

3.4.14 Plot 14: Distribution of hospitals that treats PNEUMONIA across contiguous US states


```

medic_prsvsc_combined_df_2019 <- medic_prsvsc_combined_df %>% filter(year=="2019")
PNEUMONIA_df <- medic_prsvsc_combined_df_2019 %>% filter(DRG_Cd=="193" ) %>% select(hosp_name,city,state)
PNEUMONIA_df$country <- "USA"

PNEUMONIA_shp_df <- st_read("C:/Users/Gurup/GURU/Learning/Masters/Term_2/DSC520_T302_Statistics_for_Data_Science/Week_9/data/s

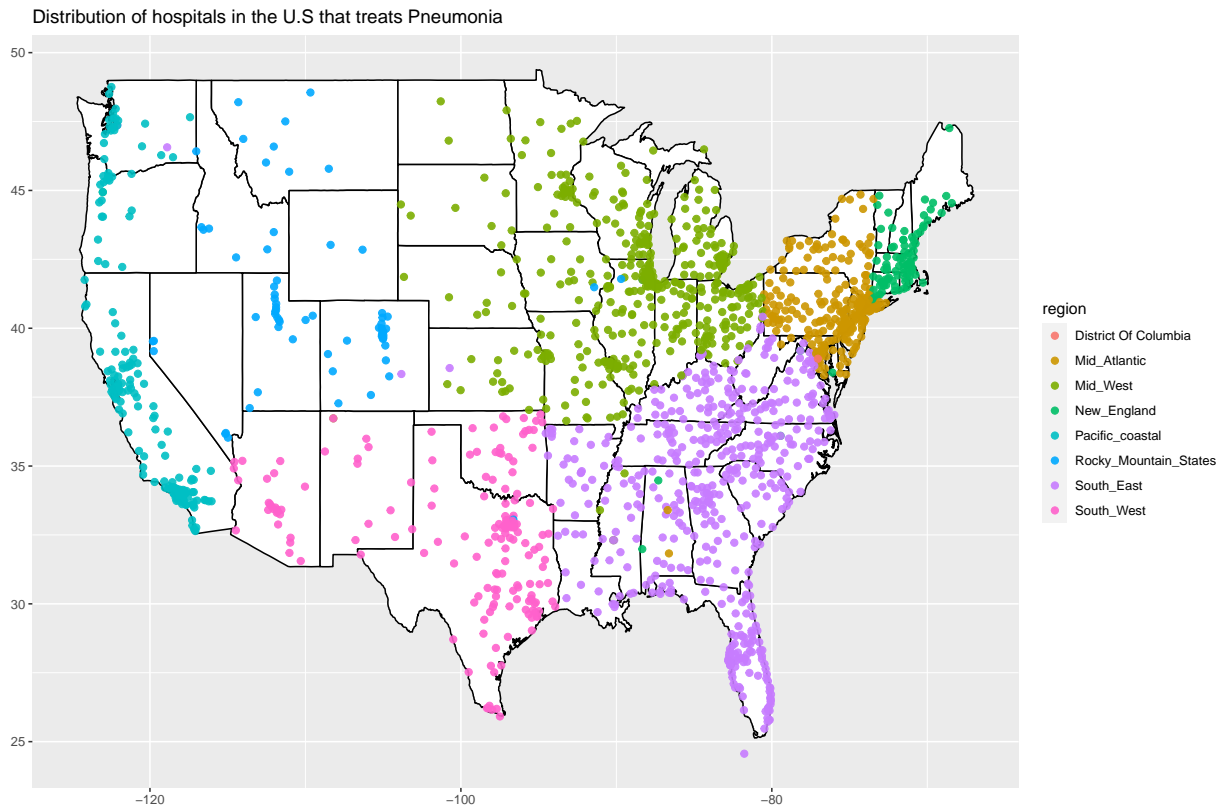
## Reading layer 'PNEUMONIA_df' from data source
##   'C:/Users/Gurup/GURU/Learning/Masters/Term_2/DSC520_T302_Statistics_for_Data_Science/Week_9/data/s
##   using driver 'ESRI Shapefile'
## Simple feature collection with 1701 features and 11 fields
## Geometry type: POINT
## Dimension:      XY
## Bounding box:   xmin: -159.3688 ymin: 20.88947 xmax: -68.42431 ymax: 64.83784
## Geodetic CRS:   WGS 84

PNEUMONIA_shp_df[c('lat', 'lon')] <- str_split_fixed(PNEUMONIA_shp_df$latlong, ',', 2)
PNEUMONIA_shp_df$lon <- as.numeric(PNEUMONIA_shp_df$lon)
PNEUMONIA_shp_df$lat <- as.numeric(PNEUMONIA_shp_df$lat)

MainStates <- map_data("state")

# Using reusable function to assign Regions
assign_regions <- function(state_code)
{
  ifelse(state_code %in% c("ME", "VT", "MA", "RI", "CT", "NH"), "New_England"
    , ifelse(state_code %in% c("NY", "PA", "NJ", "DE", "MD"), "Mid_Atlantic",
      ifelse(state_code %in% c("AR", "LA", "MS", "AL", "GA", "FL", "TN", "KY", "VA", "WV", "NC", "SC",
        ifelse(state_code %in% c("ND", "SD", "NE", "KS", "MO", "IA", "MN", "WI", "IL", "MI", "IN",
          ifelse(state_code %in% c("NV", "UT", "CO", "WY", "ID", "MT"), "Rocky_Mountain",
            ifelse(state_code %in% c("WA", "CA", "OR", "AK", "HI"), "Pacific_coast",
              ifelse(state_code %in% c("AZ", "NM", "OK", "TX"), "South_West",
                ifelse(state_code %in% c("DC"), "District Of Columbia",

```

Analysis of Plot 14 Results: The results indicate the hospitals that treats Pneumonia were located in most urban cities in the US. Not a surprise there! Also, there are *no hospitals* that treats the disease *in West texas except Elpaso* region, almost no hospitals in Nevada except the *Las Vegas and Reno area*, Also no hospitals in the *Northern Maine*

3.4.15 Additional Analysis 1:

```
# This is to study the cost of treatment for
# Gastrointestinal Hemorrhage in Arkansas and Tennessee
# states and compare the costs. These two are neighboring
# states.
medic_prsvsc_combined_df_2019 <- medic_prsvsc_combined_df %>%
  filter(year == "2019")

profit_nonprofit_merged_df <- merge(x = NASHP_profit_nonprofit_df,
  y = medic_prsvsc_combined_df_2019, by = c("CCN", "year")) %>%
```

```

select(CCN, year, hosp_name = hosp_name.x, city = city.x,
       state = state.x, zip = zip.x, health_system, hosp_ownership,
       bed_size, inpatient_occupancy, net_income, net_profit_margin,
       DRG_Desc, DRG_Cd, Tot_Dschrgs, Avg_Tot_Pymt_Amt)

TN_For_Proft_df <- profit_nonprofit_merged_df %>%
  filter(DRG_Cd == "377" & hosp_ownership == "For-Profit" &
         state == "TN" & city %in% c("CHATTANOOGA", "NASHVILLE",
                                       "MEMPHIS")) %>%
  group_by(city) %>%
  summarize(Avg_Tot_Pymt_Amt = mean(Avg_Tot_Pymt_Amt), Tot_Dschrgs = mean(Tot_Dschrgs),
            net_income = mean(net_income))

kbl(TN_For_Proft_df, caption = "Gastrointestinal Hemorrhage treatment cost in Tennessee's For-Profit hospitals",
     booktabs = T) %>%
  kable_styling(latex_options = c("striped", "hold_position"))

```

Table 14: Gastrointestinal Hemorrhage treatment cost in Tennessee's For-Profit hospitals

city	Avg_Tot_Pymt_Amt	Tot_Dschrgs	net_income
CHATTANOOGA	11752.25	20	78062717
MEMPHIS	11493.45	22	-2416698
NASHVILLE	11580.89	38	89932690

```

TN_Non_Proft_df <- profit_nonprofit_merged_df %>%
  filter(DRG_Cd == "377" & hosp_ownership == "Non-Profit" &
         state == "TN" & city %in% c("CHATTANOOGA", "NASHVILLE",
                                       "MEMPHIS")) %>%
  group_by(city) %>%
  summarize(Avg_Tot_Pymt_Amt = mean(Avg_Tot_Pymt_Amt), Tot_Dschrgs = mean(Tot_Dschrgs),
            net_income = mean(net_income))

```

```
kbl(TN_Non_Proft_df, caption = "Gastrointestinal Hemorrhage treatment cost in Tennessee's Non-Profit hospitals",
    booktabs = T) %>%
    kable_styling(latex_options = c("striped", "hold_position"))
```

Table 15: Gastrointestinal Hemorrhage treatment cost in Tennessee's Non-Profit hospitals

city	Avg_Tot_Pymt_Amt	Tot_Dschrgs	net_income
CHATTANOOGA	10167.74	90.00000	53409000
MEMPHIS	12835.86	171.50000	77632184
NASHVILLE	13567.80	54.66667	101988193

```
AR_Non_Proft_df <- profit_nonprofit_merged_df %>%
    filter(DRG_Cd == "377" & hosp_ownership == "For-Profit" &
           state == "AR") %>%
    group_by(city) %>%
    summarize(Avg_Tot_Pymt_Amt = mean(Avg_Tot_Pymt_Amt), Tot_Dschrgs = mean(Tot_Dschrgs),
              net_income = mean(net_income))

kbl(AR_Non_Proft_df, caption = "Gastrointestinal Hemorrhage treatment cost in Arkansas's Non-Profit hospitals",
    booktabs = T) %>%
    kable_styling(latex_options = c("striped", "hold_position"))
```

Table 16: Gastrointestinal Hemorrhage treatment cost in Arkansas's Non-Profit hospitals

city	Avg_Tot_Pymt_Amt	Tot_Dschrgs	net_income
BENTON	9783.80	15	922285
EL DORADO	11554.50	14	3104989
HOT SPRINGS	10264.03	30	13299474
RUSSELLVILLE	10208.33	15	14337858
SPRINGDALE	10666.15	20	22068676

```
AR_For_Proft_df <- profit_nonprofit_merged_df %>%
    filter(DRG_Cd == "377" & hosp_ownership == "Non-Profit" &
           state == "AR") %>%
    group_by(city) %>%
    summarize(Avg_Tot_Pymt_Amt = mean(Avg_Tot_Pymt_Amt), Tot_Dschrgs = mean(Tot_Dschrgs),
```

```

net_income = mean(net_income))

kbl(AR_For_Proft_df, caption = "Gastrointestinal Hemorrhage treatment cost in Arkansas's For-Profit hospitals",
    booktabs = T) %>%
    kable_styling(latex_options = c("striped", "hold_position"))

```

Table 17: Gastrointestinal Hemorrhage treatment cost in Arkansas's For-Profit hospitals

city	Avg_Tot_Pymt_Amt	Tot_Dschrgrs	net_income
BATESVILLE	10206.452	31.0	-181726.0
CONWAY	9647.467	17.5	-97924.5
FAYETTEVILLE	10390.350	40.0	75075105.0
FORT SMITH	10000.257	50.0	23101802.0
HOT SPRINGS	10351.214	14.0	17245958.0
JONESBORO	9991.909	50.0	13636842.5
LITTLE ROCK	10468.948	87.0	-30326879.0
MOUNTAIN HOME	9618.606	33.0	11380461.0
NORTH LITTLE ROCK	9439.283	53.0	11040244.0
PINE BLUFF	12851.192	26.0	2952809.0
ROGERS	11348.500	18.0	32228982.0
SEARCY	11041.111	27.0	17918859.0

Analysis of the above results indicates that *“For-Profit” hospitals costs lesser than “Non-Profit” hospitals* for Gastrointestinal Hemorrhage in *Memphis and Nashville, Tennessee, but not in Chattanooga, TN*. Also, many cities in the neighboring state of *Arkansas offers the same treatment for lesser cost*. There is almost a *\$1000 difference between the costs in Little rock and North Little rock* for treating Gastrointestinal Hemorrhage

3.4.16 Additional Analysis 2:

```

# Compare costs within a city between For-Profit and
# Non-Profit hospitals in Houston, Texas

houston_tx_df <- profit_nonprofit_merged_df %>%
  filter(state == "TX" & DRG_Cd == "305" & city == "HOUSTON" &
    zip %in% c("77090", "77082", "77070", "77094")) %>%
  select(hosp_name, city, state, zip, hosp_ownership, bed_size,
    Avg_Tot_Pymt_Amt, DRG_Desc)

```

```
kbl(houston_tx_df, caption = "Costs of treatments for HyperTension related disorders in Houston,TX",
    booktabs = T) %>%
    kable_styling(latex_options = c("striped", "hold_position"))
```

Table 18: Costs of treatments for HyperTension related disorders in Houston,TX

hosp_name	city	state	zip	hosp_ownership	bed_size	Avg_Tot
HOUSTON NORTHWEST MEDICAL CENTER	HOUSTON	TX	77090	For-Profit	288	
HCA HOUSTON HEALTHCARE WEST	HOUSTON	TX	77082	For-Profit	230	
METHODIST WILLOWBROOK HOSPITAL	HOUSTON	TX	77070	Non-Profit	301	
METHODIST WEST HOUSTON HOSPITAL	HOUSTON	TX	77094	Non-Profit	238	

Analysis by studying the costs of Hypertension treatments in Houston,TX. There is a *difference of about \$1,500 in the costs* between different hospitals in the same city.

3.5 Building a model

A *Linear Regrssion Model* will be built to predict the Cost of treatment based on various fields used in this project. As the dataset has over half a million records, for the purpose of the model a subset of data will be taken.

```
# Creating a subset of dataframe by filtering only Kidney
# related ailments (Nephrology)
Nephrology_prvsvc_df <- medic_prvsvc_combined_df %>%
    filter(grepl("Kidney", DRG_Desc, ignore.case = TRUE))
nrow(Nephrology_prvsvc_df)
```

```
## [1] 26882
```

```
# filtering the years
Nephrology_NASHP_df <- NASHP_df %>%
    filter(year %in% c("2020", "2019", "2018"))
nrow(Nephrology_NASHP_df)
```

```
## [1] 13698
```

```
# building a basic model to predict Avg_Tot_Pymt_Amt
Nephrology_lm1 <- lm(Avg_Tot_Pymt_Amt ~ DRG_Cd, data = Nephrology_prvsvc_df,
  na.action = na.omit)
# The Model predicts 71.89% of the Avg_Tot_Pymt_Amt
summary(Nephrology_lm1)$r.squared
```

```
## [1] 0.7189391
```

```
summary(Nephrology_lm1)$adj.r.squared
```

```
## [1] 0.7187508
```

```
# building a model based on the fields from
# Nephrology_prvsvc_df
Nephrology_lm2 <- lm(Avg_Tot_Pymt_Amt ~ DRG_Cd + year + state,
  data = Nephrology_prvsvc_df, na.action = na.omit)
# The Model predicts 75.53% of the Avg_Tot_Pymt_Amt
summary(Nephrology_lm2)$r.squared
```

```
## [1] 0.7553096
```

```
summary(Nephrology_lm2)$adj.r.squared
```

```
## [1] 0.7546708
```

To improve the efficiency of the model, we combine the data from other dataset *Nephrology_NASHP_df* as shown below.

```
# Merging both the datasets based on CCN and year
Nephrology_merged_df <- merge(x = Nephrology_NASHP_df, y = Nephrology_prvsvc_df,
  by = c("CCN", "year")) %>%
  select(CCN, year = year, hosp_name = hosp_name.x, address,
    zip = zip.x, state = state.x, hosp_ownership, bed_size,
    inpatient_occupancy, net_income, net_profit_margin,
```



```

DRG_Cd, DRG_Desc, Tot_Dschrgs, Avg_Tot_Pymt_Amt)

# building a model including the hospital code(CCN)
# Sampling the data to take 2000 samples
set.seed(42)
Nephrology_merged_df_sampled <- Nephrology_merged_df[sample(nrow(Nephrology_merged_df),
  size = 10000), ]
nrow(Nephrology_merged_df_sampled)

```

```
## [1] 10000
```

```

Nephrology_lm3 <- lm(Avg_Tot_Pymt_Amt ~ DRG_Cd + CCN + bed_size +
  hosp_ownership + net_income + year, data = Nephrology_merged_df_sampled,
  na.action = na.omit)
# By combining datasets and including CCN, the model can
# predict upto 90%
summary(Nephrology_lm3)$r.squared

```

```
## [1] 0.901343
```

```
summary(Nephrology_lm3)$adj.r.squared
```

```
## [1] 0.8688378
```

```

# building a model with the state code instead of CCN
Nephrology_lm4 <- lm(Avg_Tot_Pymt_Amt ~ DRG_Cd + state + year +
  bed_size + hosp_ownership + net_income + Tot_Dschrgs +
  inpatient_occupancy, data = Nephrology_merged_df_sampled,
  na.action = na.omit)
# This model can predict upto 78% of the Treatment costs
summary(Nephrology_lm4)$r.squared

```

```
## [1] 0.7863473
```

```
summary(Nephrology_lm4)$adj.r.squared
```

```
## [1] 0.7847108
```

```
# Identifying the best model and using the step_backward
# method to identify the predictors
best_mod_state <- ols_step_best_subset(model = Nephrology_lm4)
best_mod_state
```

```
##                               Best Subsets Regression
```

```
## -----
## Model Index    Predictors
## -----
##      1         DRG_Cd
##      2         DRG_Cd state
##      3         DRG_Cd state bed_size
##      4         DRG_Cd state bed_size hosp_ownership
##      5         DRG_Cd state bed_size hosp_ownership Tot_Dschrgrs
##      6         DRG_Cd state bed_size hosp_ownership Tot_Dschrgrs inpatient_occupancy
##      7         DRG_Cd state year bed_size hosp_ownership Tot_Dschrgrs inpatient_occupancy
##      8         DRG_Cd state year bed_size hosp_ownership net_income Tot_Dschrgrs inpatient_occupancy
## -----
```

```
##
##                               Subsets Regression Summary
```

```
## -----
##                               Adj.      Pred
## Model    R-Square  R-Square  R-Square  C(p)      AIC      SBIC      SBC
## -----
##      1      0.7250   0.7245   0.7138   2777.6153  196346.7287  167932.9952  196490.9355
##      2      0.7593   0.7577   0.7469   1184.3550  195112.7046  166601.2792  195617.4284
##      3      0.7755   0.7739   0.7633   434.0874   194418.0221  165906.9211  194929.9562
##      4      0.7828   0.7813   0.7705    96.1110   194090.5047  165577.6485  194616.8595
##      5      0.7841   0.7825   0.7716    39.2144   194033.9397  165521.1364  194567.5048
```

```
##      6      0.7856      0.7841      -Inf      -30.0390      193945.2401      165435.3692      194486.0081
##      7      0.7863      0.7847      -Inf      -60.1856      193916.8991      165405.0890      194472.0876
##      8      0.7863      0.7847      -Inf      -59.0000      193918.0784      165406.2851      194480.4771
## -----

## AIC: Akaike Information Criteria
## SBIC: Sawa's Bayesian Information Criteria
## SBC: Schwarz Bayesian Criteria
## MSEP: Estimated error of prediction, assuming multivariate normality
## FPE: Final Prediction Error
## HSP: Hocking's Sp
## APC: Amemiya Prediction Criteria
```

```
# using the Step Backward method to identify the
# predictors to be excluded. Results shows Net_income can
# be excluded.
ols_step_backward_p(model = Nephrology_lm4, details = FALSE)
```

```
##
##
## Elimination Summary
## -----
##      Variable      Adj.
## Step   Removed   R-Square   R-Square   C(p)      AIC      RMSE
## -----
##      1   net_income   0.7863     0.7847   -60.1856   193916.8991   3920.5618
## -----
```

```
# Building a new model based on the recommendation
Nephrology_lm5 <- lm(Avg_Tot_Pymt_Amt ~ DRG_Cd + state + year +
  bed_size + hosp_ownership + Tot_Dschrgs + inpatient_occupancy,
  data = Nephrology_merged_df_sampled, na.action = na.omit)

# This model can predict upto 78% accuracy of the
```

```
# Treatment costs
summary(Nephrology_lm5)$r.squared
```

```
## [1] 0.7863298
```

```
summary(Nephrology_lm5)$adj.r.squared
```

```
## [1] 0.7847148
```

```
# The results indicate that DWT results are almost close
# to 2(which is optimal), So the model is a good
# predictor for the costs.
durbinWatsonTest(Nephrology_lm5)
```

```
## lag Autocorrelation D-W Statistic p-value
## 1 0.002452214 1.995086 0.776
## Alternative hypothesis: rho != 0
```

The Model *Nephrology_lm5* which has year as one of the predictors can be used to predict the Treatment costs in the future. This model can predict the prices overall at state level.

The model *Nephrology_lm3* can predict the costs when CCN (the hospital identifier) is one of the inputs. This can yield much accurate results than *Nephrology_lm5*.

```
# Predict the Treatment costs using Predict function

Nephrology_merged_df_sampled$predicted_TreatmentCost_lm3 <- predict(Nephrology_lm3,
  predicted_TreatmentCost = Nephrology_merged_df_sampled)

Nephrology_predicted <- Nephrology_merged_df_sampled %>%
  select(predicted_TreatmentCost_lm3, Avg_Tot_Pymt_Amt)

kbl(head(Nephrology_predicted), caption = "Predicting Medical costs using Linear Regression model",
  booktabs = T) %>%
  kable_styling(latex_options = c("striped", "hold_position"))
```

Table 19: Predicting Medical costs using Linear Regression model

	predicted_TreatmentCost_lm3	Avg_Tot_Pymt_Amt
18753	8980.000	8980.000
21657	5940.113	5958.717
9290	3933.688	5108.800
1252	6430.602	6393.902
15506	9401.951	9616.760
8826	10290.243	10360.042

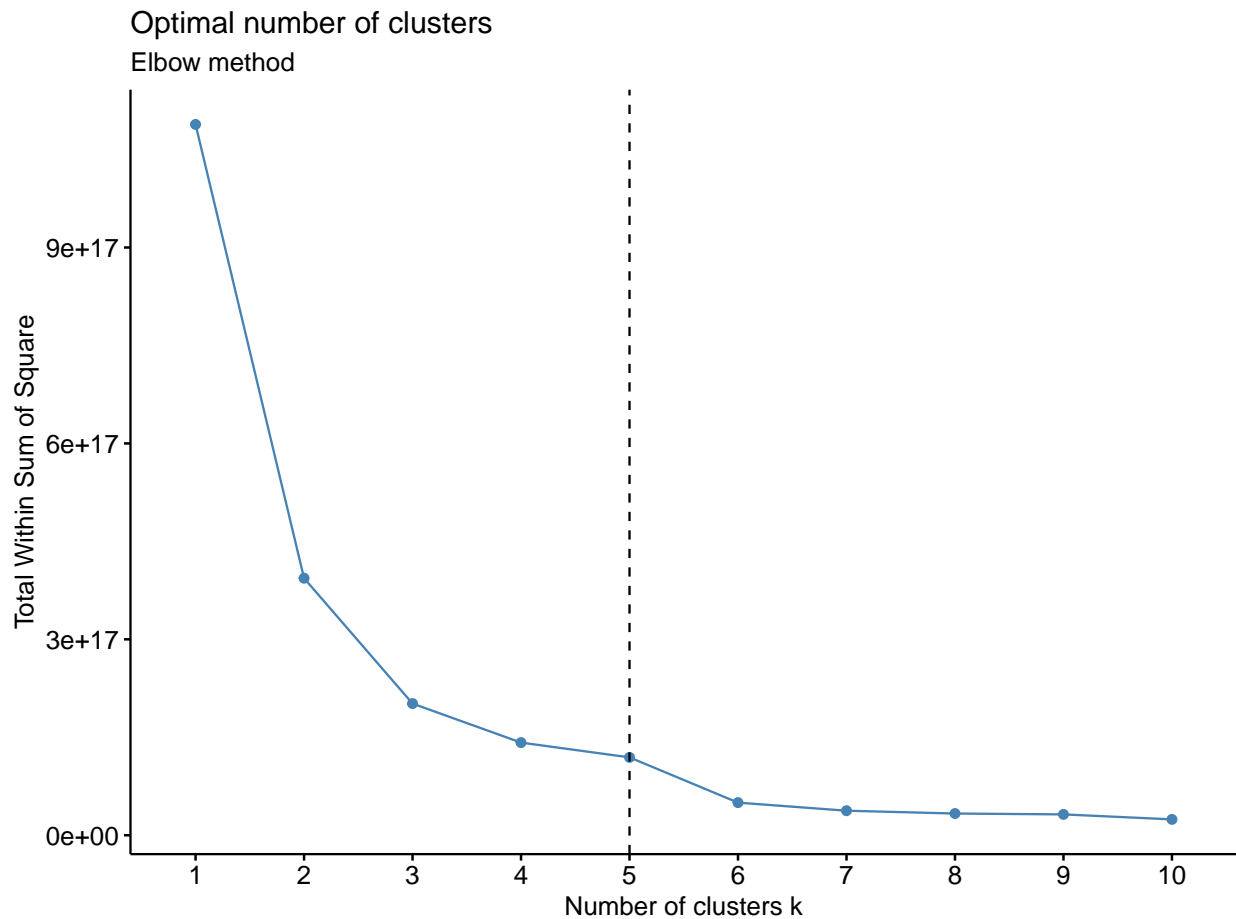
3.6 Using K-means Clustering

3.6.1 Cluster 1- To find clusters based on hospital's Net income

```
# Identifying clusters based on the profits made by the
# 'For-Profit' hospitals in the US
NAS_2020_profit_df <- NASHP_df %>%
  filter(hosp_ownership == "For-Profit" & year == "2020" &
    net_income > 0)
nrow(NAS_2020_profit_df)
```

```
## [1] 610
```

```
# Identifying the optimal clusters which comes to 5
fviz_nbclust(NAS_2020_profit_df[, c("net_income")], kmeans,
  method = "wss") + geom_vline(xintercept = 5, linetype = 2) +
  labs(subtitle = "Elbow method")
```



```
# Applying K-means function on the 'net_income'
NAS_2020_profit_df_km <- kmeans(NAS_2020_profit_df[, c("net_income")],
  5)

# Extracting and saving the cluster details
NAS_2020_profit_df$cluster <- NAS_2020_profit_df_km$cluster
NAS_2020_profit_df$country <- "USA"

# Reading the data from .shp file which contains the
# latitude and longitude data for our dataset. This
# latitude and longitude is obtained from an external
# source.
NAS_2020_profit_shp_df <- st_read("C:/Users/Gurup/GURU/Learning/Masters/Term_2/DSC520_T302_Statistics_for_Data_Science/Week_9/data/s

## Reading layer 'NAS_2020_profit_df_output3' from data source
## 'C:\Users\Gurup\GURU\Learning\Masters\Term_2\DSC520_T302_Statistics_for_Data_Science\Week_9\data\s
```

```

## using driver 'ESRI Shapefile'
## Simple feature collection with 464 features and 25 fields
## Geometry type: POINT
## Dimension: XY
## Bounding box: xmin: -149.2579 ymin: 24.5809 xmax: -70.78932 ymax: 61.56276
## Geodetic CRS: WGS 84

NAS_2020_profit_map_df <- NAS_2020_profit_shp_df %>%
  select(city, state, region, latlong, bed_size, net_income,
         net_profit, cluster)

NAS_2020_profit_map_df[c("lat", "lon")] <- str_split_fixed(NAS_2020_profit_map_df$latlong,
  ",", 2)

# Subsetting the data only for a few states
NAS_2020_profit_map_df <- NAS_2020_profit_map_df %>%
  filter((state %in% c("ID", "NV", "UT", "CO", "NM", "MT",
    "WY", "TN", "KY", "IA", "MN", "OK", "ND", "SD", "NE")))

NAS_2020_profit_map_df$lon <- as.numeric(NAS_2020_profit_map_df$lon)
NAS_2020_profit_map_df$lat <- as.numeric(NAS_2020_profit_map_df$lat)

icon.pop <- awesomeIcons(icon = "users", markerColor = ifelse(NAS_2020_profit_map_df$cluster ==
  "4", "blue", ifelse(NAS_2020_profit_map_df$cluster == "3",
  "orange", ifelse(NAS_2020_profit_map_df$cluster == "2",
  "green", ifelse(NAS_2020_profit_map_df$cluster == "1",
  "red", "grey")))), library = "fa", iconColor = "black")

img1 <- leaflet(NAS_2020_profit_map_df) %>%
  addTiles() %>%
  addAwesomeMarkers(lng = ~lon, lat = ~lat, label = ~cluster,
    icon = icon.pop)

# The leaflet package does not render to PDF, hence the
# image created by the leaflet will be attached here.

```

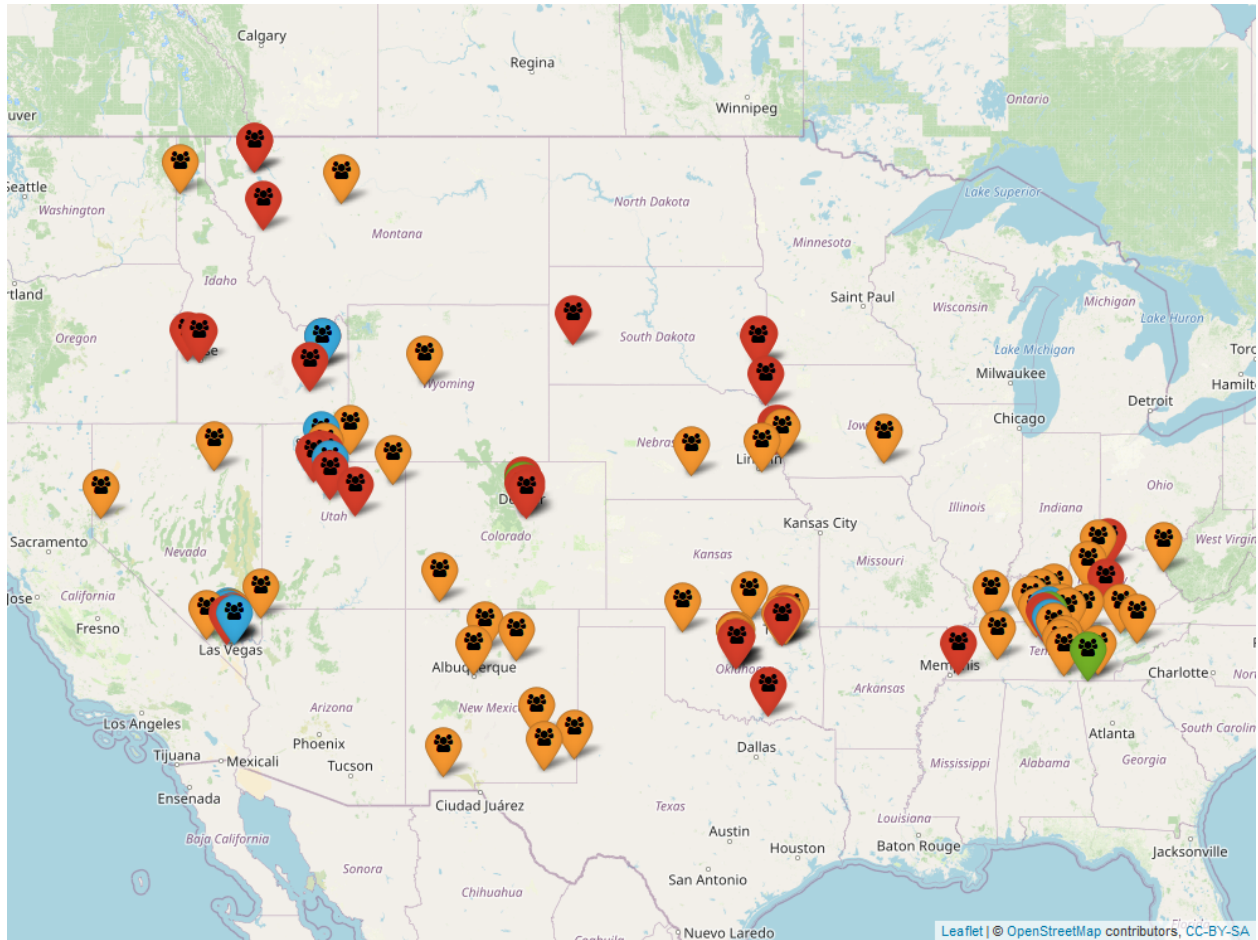


Figure 1: K-means Clustering based on the Net Income of the hospitals (Displaying limited data)

3.6.2 Cluster 2 - To find clusters based on Total hospitalization and cost of Treating Bronchitis in the US

```
# creating a data frame for Bronchitis related ailments
# and excluding the hospitals from Alaska and Hawii for
# the analysis

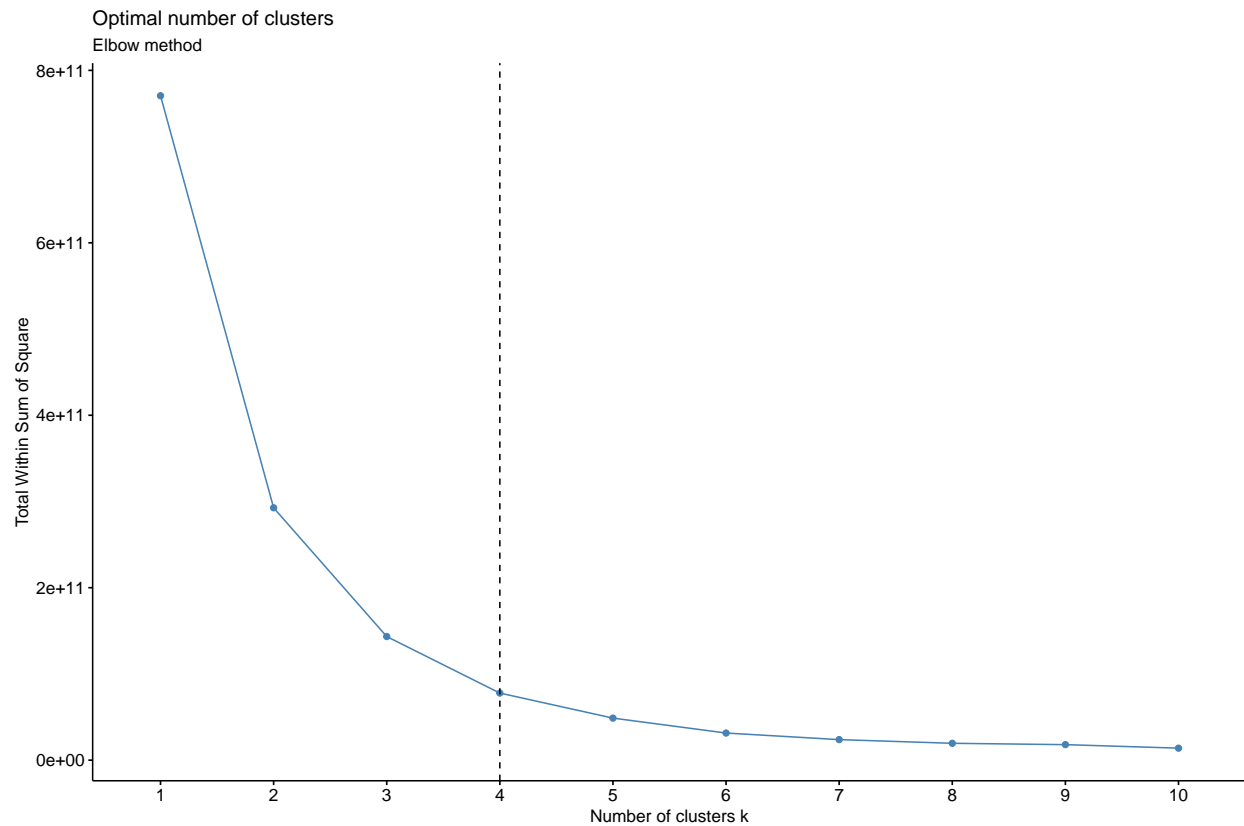
bronchi_summary_df <- medic_prvsvc_combined_df %>%
  filter(grepl("BRONCHITIS", DRG_Desc, ignore.case = TRUE) &
         !(state %in% c("HI", "AK"))) %>%
  select(state, city, Tot_Dschrgrs, Avg_Tot_Pymt_Amt) %>%
  group_by(state, city) %>%
  dplyr::summarise(Tot_Dschrgrs = sum(Tot_Dschrgrs), Avg_Tot_Pymt_Amt = sum(Avg_Tot_Pymt_Amt))

bronchi_summary_df$city_state <- paste(bronchi_summary_df$city,
                                       bronchi_summary_df$state)

# Merging with the us.cities dataset to extract the
# latitude and longitude details

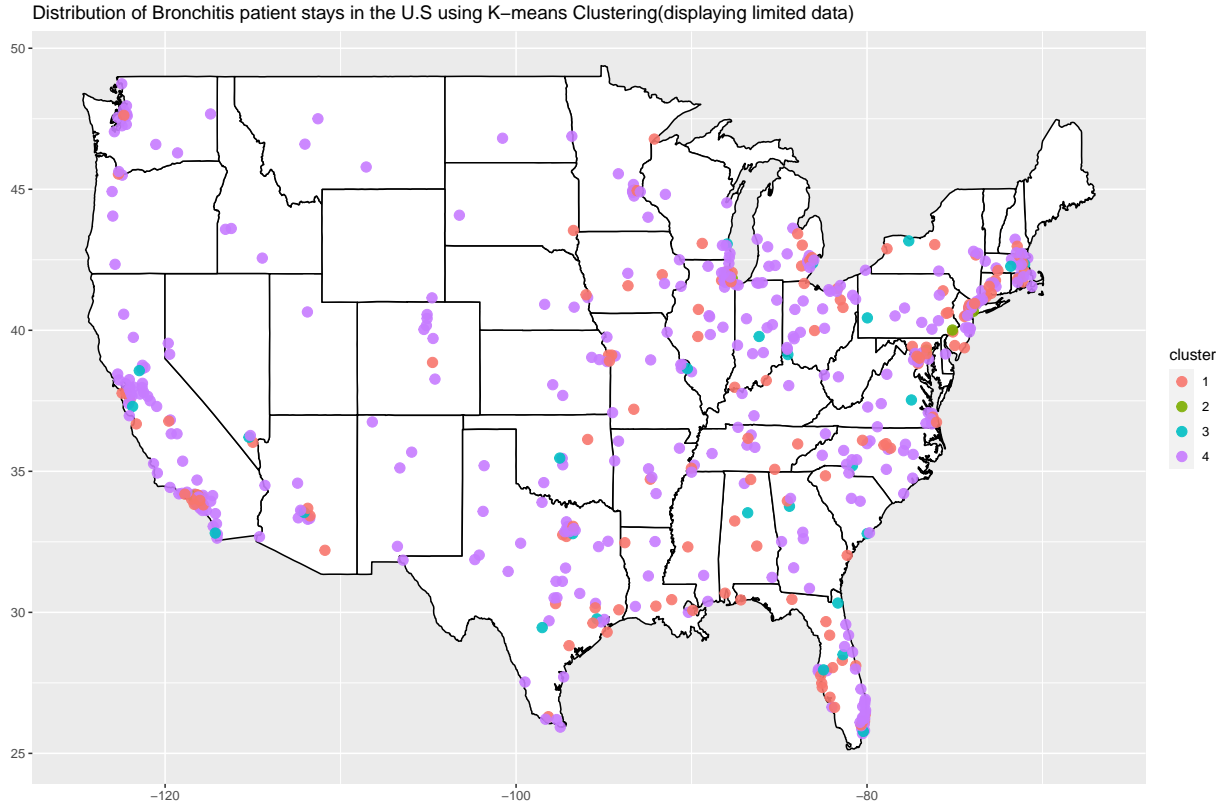
bronchi_merged_df <- merge(x = bronchi_summary_df, y = us.cities,
                           by.x = c("city_state"), by.y = c("name")) %>%
  select(state, city, Tot_Dschrgrs, Avg_Tot_Pymt_Amt, lat,
         long, population = pop)

# find optimal nbr of clusters which comes out to be 4
fviz_nbclust(bronchi_merged_df[, c("Tot_Dschrgrs", "Avg_Tot_Pymt_Amt")],
             kmeans, method = "wss") + geom_vline(xintercept = 4, linetype = 2) +
  labs(subtitle = "Elbow method")
```



```
# Creating a k-means clustering model based on Total
# inpatient stays and Cost of treatment of Bronchitis
bronchi_km <- kmeans(bronchi_merged_df[, c("Tot_Dschrgrs", "Avg_Tot_Pymt_Amt")],
  4)
bronchi_merged_df$cluster <- as.factor(bronchi_km$cluster)

# Plotting on US map
MainStates <- map_data("state")
ggplot() + geom_polygon(data = MainStates, aes(x = long, y = lat,
  group = group), color = "black", fill = "white") + geom_point(data = bronchi_merged_df,
  aes(x = long, y = lat, color = cluster), size = 3, alpha = 0.9) +
  scale_size(name = "Clusters") + labs(title = "Distribution of Bronchitis patient stays in the U.S u
  x = "", y = "") + theme(legend.position = "right")
```



```
# selecting only 100 cities for display purposes
img2 = leaflet(bronchi_merged_df[1:100, ]) %>%
  addTiles() %>%
  addAwesomeMarkers(lng = ~long, lat = ~lat, label = ~cluster,
    icon = icon.pop)
```

4 Summary of the Interesting insights from the Analysis :

A summary analysis of each plot has been included in the study. Below are the high level insights:

1. The *hospitalizations* for the diseases included in this study indicates that over the years, the number of inpatient stays remained consistent and *had a strong correlation with the population* of the city/state. For instance, the number of chemotherapy inpatient stays was higher in Charleston, WV than any other city in the state which is in line with the city population. However, *Charleston WV, had significantly more number of inpatient stays for chemotherapy than Huntington, WV*

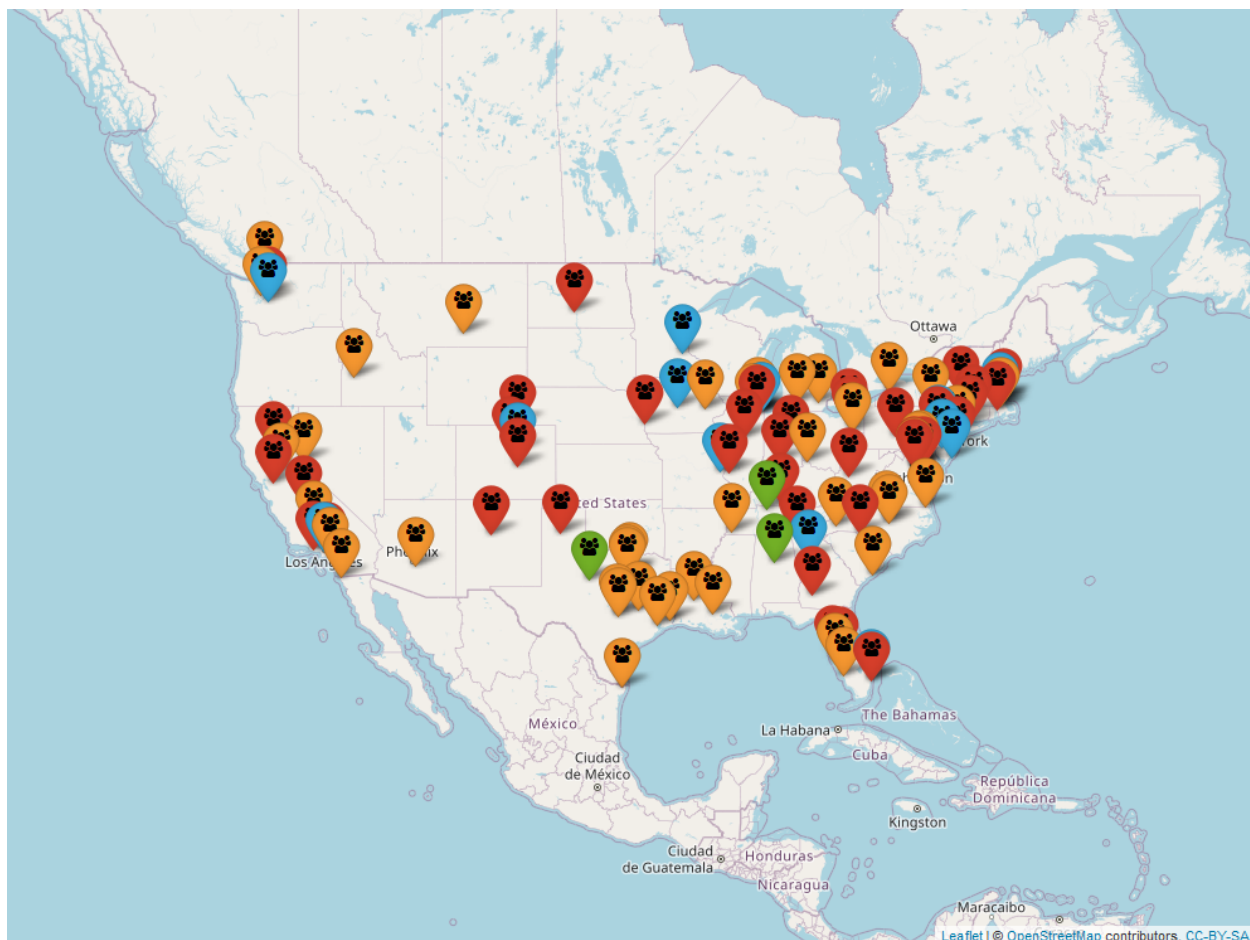


Figure 2: Distribution of Bronchitis patient stays in the U.S using K-means Clustering (Displaying limited data)

though their population is not vastly different. This may need additional research to understand the reasoning behind it.

2. The study helped us *identify the hospitalization rates of same medical treatment between two neighboring cities* across different states and analyze the results. In the comparison of Kansas City MO vs KS, it turns out more patients preferred staying in Kansas hospitals than in Missouri.
3. The study helped us understand how much the *cost of treatment in a state is different from the national average*. In the Wyoming example, the costs of heart related treatments were compared across National average and it was discovered that except for one medical procedure, all other treatment costs for Hear related ailments were close to national average, which means the costs are reasonable in Wyoming.
4. The study included various graphs about *top ten hospitalizations* in different states and cities. This information can be *analyzed over the period of years to understand the shift in the behavior*.
5. In plot 10, the cost of Spine procedure treatments were compared between state and National average and the results indicate that *hospitals in Washington DC was the most expensive followed by Maryland, California, Oregon and New Hampshire*. Similar analysis can be performed for various procedures to understand how the costs vary across each state.
6. While analyzing the *most profitable hospitals in Nebraska*, it was discovered that most of the hospitals were *located in the eastern end of the state* which happens to be more populated region in the state. Similar studies can be done for other states to understand the reasoning why more profitable hospitals prefer a certain portion of the state. This can be compared with *Texas*, which despite being *most-profitable state for the hospitals, there are not many hospitals in the southwest region of the state*.
7. While analyzing the net *profit margin of "For-Profit" hospitals* located in *Anchorage, Alaska* and comparing it with the cost of treatment over the same period, we can find perfect correlation. Such insights can help us reveal the reasoning behind how a hospital is profitable while others lose money.
8. In the map that compares the Profit/loss of southern states *between 2011 and 2020*, we can find how *hospitals in Texas turned out to be a Profit machine over this period, followed by Florida*.
9. The map that includes the results of hospitals that treats Pneumonia indicates that most of the hospitals are located in the *urban areas* and there are hardly a few hospitals in the Rocky mountain states. There are almost *no hospitals in the Northern Maine* except for one. In Nevada, except for regions around *Las Vegas and Reno*, there are no hospitals that treats pneumonia anywhere else in the state.

10. The linear models built in this project indicates that the *Disease code(DRG_Cd) is a strong predictor of the Medical costs*. Adding the Hospital code(CCN) makes the prediction even stronger, though it may lead to over-fitting of the model. *Surprisingly the “net-income” of the hospitals is not a good predictor* for the treatment cost, based on the model results.
11. The results from K-means clustering of Bronchitis hospitalizations indicates, most patient stays are in Urban cities. *Cities such as Boston, NewYork, Dallas, Houston, Southern California, South western Florida are in a cluster*. Interestingly, cities like *Omaha, Nebraska and Sioux Falls,SD* are also in the same cluster. This needs more research to find why there are more Bronchitis hospitalizations in these cities that are comparable to big metropolitan cities .

5 Implications to the Target audience

One of the main goals of this project is to *study the pattern of treatment costs* across the country and analyze the results that can provide recommendations to choose the care providers wisely. Below are some of the implications from the research:

1. While choosing the hospitals, one can *choose “Non-profit” hospitals over “For-profit” hospitals*, if cost is an issue, which is the case for most of the people. The study found that, in most of the scenarios the “For-Profit” hospitals costs more than “Non-Profit” hospitals for the same treatment. Not very surprising, but at least we have numbers to prove that!
2. The study found that the *medical costs in some states are more expensive than others*. For instance, the cost of Spinal procedure is quite expensive in Washington D.C, Maryland, California and Oregon, whilst, the treatment is cheapest in Arkansas and West Virginia. For example, someone that lives in Memphis, Tennessee, can opt to choose treatment in Arkansas hospitals if the costs are lesser which is just a few miles away.
3. The *prediction model build can be used to predict the costs* in the years to come. Two models were build, one with state details and the other one with Hospital details. Hence, either of the models can be used a prediction tool to predict the costs.
4. Treatment costs for same medical procedure can vary over 10-15% between hospitals in neighboring cities. In the Arkansas example, the results indicated that Treatments costs in North Little rock is cheaper than Little rock. Hence *analyzing the costs between nearby cities before choosing the care* using the tool is not a bad idea.

5. Example from treatment costs in Houston,TX indicates that the *treatment costs can be much cheaper in some areas than others within the same city*. Hence, price shopping before choosing the hospitals is a great idea to save money.
6. While searching for doctors for a specific medical procedure we usually google to search for the hospitals or doctors. Sometimes the results can be biased while searching on the internet sites.Hence *this tool which shows unbiased results* can be used to search for doctors/hospitals based on the average costs and then make a decision after that. *This tool can be a starting point while searching for care providers and can help narrow down the search.*
7. This research included bed size of the hospitals and the occupancy rates in the study. Hence knowing the *occupancy rates of hospitals* can come in handy especially at times of pandemic such as Covid-19, *when there is shortage of beds across the hospitals.*
8. The research can *help with the Negotiations of the medical bills* with the hospitals. Not many know that hospitals may be willing to negotiate the bills or come up with a payment plan if we are able to provide right statistics at time of negotiation. Especially with the “For-profit” hospitals, the costs can be compared across various hospitals around the city and facts & numbers can help with the negotiations and sometimes can help save families go bankrupt.

6 Limitations of the analysis and concluding Remarks:

Medical system in the U.S is very complicated as there are so many moving parts associated with the billing. This project may partially address the problem statement of identifying the affordable hospitals in the region as the study does not include the health insurance data in the research. Also this study uses Medicare data for research and at times hospitals may bill a variable price for treatments different from the medicare costs. Hence this study can be used a baseline for further research. Also only the basic Medicare costs are considered for the research and programs such as “Medicare-Advantage” are not included in this study. The costs are compared across hospitals with the assumption that they are all within the same network, which may not be true in all cases.

How someone can enhance the research A more complex model can be built including the Health insurance information, which can provide more accurate results. Visualization of the data using tools such as Tableau or Power BI can be a great enhancement for this research which is not included in this project. Visualization can add more features such as the ability to select multiple hospitals and select medical procedures in the hospitals and compare the costs.Also enhancing the tool with Geo-spatial visualization and

displaying the results in the Map can provide broader use of the project.

Due to the complexity involved with the medical system in the U.S, a single research may not provide a comprehensive solution for all problems. Hence this tool can be a great starting point for further medical research by building complex models including all the bells and whistles in our medical system.

7 Session Info

```
sessionInfo()
```

```
## R version 4.2.2 (2022-10-31 ucrt)
## Platform: x86_64-w64-mingw32/x64 (64-bit)
## Running under: Windows 10 x64 (build 22621)
##
## Matrix products: default
##
## locale:
## [1] LC_COLLATE=English_United States.utf8
## [2] LC_CTYPE=English_United States.utf8
## [3] LC_MONETARY=English_United States.utf8
## [4] LC_NUMERIC=C
## [5] LC_TIME=English_United States.utf8
##
## attached base packages:
## [1] stats      graphics  grDevices  utils      datasets  methods   base
##
## other attached packages:
## [1] maps_3.4.1      NbClust_3.0.1    factoextra_1.0.7 sf_1.0-9
## [5] car_3.1-1       carData_3.0-5    olsrr_0.5.3      magrittr_2.0.3
## [9] leaflet_2.1.1    plotly_4.10.1    usmap_0.6.1      forcats_0.5.2
## [13] purrr_1.0.1      tibble_3.1.8     tidyverse_1.3.2  reshape2_1.4.4
## [17] kableExtra_1.3.4 lemon_0.4.6       knitr_1.41       treemap_2.4-3
## [21] ggalt_0.4.0      stringr_1.5.0    scales_1.2.1     DT_0.27
```



```

## [25] tidyr_1.3.0      ggthemes_4.2.4    ggplot2_3.4.0     readr_2.1.3
## [29] lubridate_1.9.0  timechange_0.1.1  dplyr_1.0.10      readxl_1.4.1
##
## loaded via a namespace (and not attached):
## [1] googledrive_2.0.0  colorspace_2.0-3   ggsignif_0.6.4
## [4] class_7.3-20       ellipsis_0.3.2     fs_1.5.2
## [7] proxy_0.4-27       rstudioapi_0.14    ggpubr_0.6.0
## [10] farver_2.1.1       ggrepel_0.9.2      fansi_1.0.3
## [13] xml2_1.3.3         extrafont_0.19     jsonlite_1.8.4
## [16] broom_1.0.1        gridBase_0.4-7     Rttf2pt1_1.3.12
## [19] dbplyr_2.2.1       shiny_1.7.3        compiler_4.2.2
## [22] httr_1.4.4         backports_1.4.1    assertthat_0.2.1
## [25] fastmap_1.1.0      lazyeval_0.2.2     gargle_1.2.1
## [28] cli_3.4.1          later_1.3.0        formatR_1.12
## [31] htmltools_0.5.3    tools_4.2.2        igraph_1.3.5
## [34] gtable_0.3.1       glue_1.6.2         Rcpp_1.0.10
## [37] cellranger_1.1.0   vctrs_0.5.2        svglite_2.1.1
## [40] extrafontdb_1.0    crosstalk_1.2.0    xfun_0.34
## [43] rvest_1.0.3        mime_0.12           lifecycle_1.0.3
## [46] rstatix_0.7.2      googlesheets4_1.0.1 goftest_1.2-3
## [49] MASS_7.3-58.1      hms_1.1.2          promises_1.2.0.1
## [52] proj4_1.0-12       RColorBrewer_1.1-3 yaml_2.3.6
## [55] gridExtra_2.3      stringi_1.7.8      highr_0.9
## [58] nortest_1.0-4      e1071_1.7-12       rlang_1.0.6
## [61] pkgconfig_2.0.3    systemfonts_1.0.4  evaluate_0.18
## [64] lattice_0.20-45    labeling_0.4.2     htmlwidgets_1.5.4
## [67] tidyselect_1.2.0   usmapdata_0.1.0    plyr_1.8.8
## [70] R6_2.5.1           generics_0.1.3     DBI_1.1.3
## [73] pillar_1.8.1       haven_2.5.1        withr_2.5.0
## [76] units_0.8-1        abind_1.4-5        ash_1.0-15
## [79] modelr_0.1.10      crayon_1.5.2       KernSmooth_2.23-20
## [82] utf8_1.2.2         tzdb_0.3.0         rmarkdown_2.18
## [85] grid_4.2.2         data.table_1.14.4  reprex_2.0.2

```

```
## [88] digest_0.6.30      classInt_0.4-8      webshot_0.5.4
## [91] xtable_1.8-4         httpuv_1.6.6        munsell_0.5.0
## [94] viridisLite_0.4.1
```

References

- “Hospital Provider Cost Report.” 2022. Centers for Medicare & Medicaid Services. <https://catalog.data.gov/dataset/hospital-provider-cost-report>.
- Konish, Lorie. n.d. “This Is the Real Reason Most Americans File for Bankruptcy.” *CNBC*. Accessed February 11, 2023. <https://www.cnbc.com/2019/02/11/this-is-the-real-reason-most-americans-file-for-bankruptcy.html>.
- “Medicare Inpatient Hospitals - by Geography and Service.” 2022. Centers for Medicare & Medicaid Services. <https://catalog.data.gov/dataset/medicare-inpatient-hospitals-by-geography-and-service-1bbc8>.
- “Medicare Inpatient Hospitals - by Provider and Service.” 2022. Centers for Medicare & Medicaid Services. <https://catalog.data.gov/dataset/medicare-inpatient-hospitals-by-provider-and-service-9af02>.