

# HW06 - DS6306-402

Dan Crouthamel

6/11/2019

Background: Your organization is responsible for building new VA hospitals in the mainland of the United States. You are a little overwhelmed by the prospect, not sure which places need the most help. You decide to begin by exploring healthcare facility data recorded by the U.S. Government.

Disclaimer: While these are real data, the assignment is not an endorsement for any particular position on medical affairs or building hospitals. It is for instructional use only.

## 1 - Mental Health Clinics (40%):

- a) This data set is a survey of every known healthcare facility that offers mental health services in the United States in 2015. Navigate to <https://datafiles.samhsa.gov/study-dataset/nationalmental-health-services-survey-2015-n-mhss-2015-ds0001-nid17098> and select the R download. Look through the codebook PDF for an explanation on certain variables. Upon opening the RDA file, the data set should be inserted into your global environment, which you can then reference.

Use the following command to load the file.

```
# Load .rda file. This will load an object called mh2015_puf
load(file = "N-MHSS-2015-DS0001-data/N-MHSS-2015-DS0001-data-r.rda")
```

- b) Please create code which lists the State abbreviations without their counts, one abbreviation per State value. It does not have to in data frame format. A vector is fine.

```
# Create states vector and mutate LST column to remove whitespace
# Thanks David for tip on %<>%
states <- (mh2015_puf$LST %<>% trimws()) %>% unique()
print(states)
```

```
## [1] "AL" "AK" "AZ" "AR" "CA" "CO" "CT" "DE" "DC" "FL" "GA" "HI" "ID" "IL"
## [15] "IN" "IA" "KS" "KY" "LA" "ME" "MD" "MA" "MI" "MN" "MS" "MO" "MT" "NE"
## [29] "NV" "NH" "NJ" "NM" "NY" "NC" "ND" "OH" "OK" "OR" "PA" "RI" "SC" "SD"
## [43] "TN" "TX" "UT" "VT" "VA" "WA" "WV" "WI" "WY" "AS" "GU" "PR" "VI"
```

- c) Filter the data.frame from 1A. We are only interested in the Veterans Administration (VA) medical centers in the mainland United States—create a listing of counts of these centers by state, including only mainland locations. Alaska, Hawaii, and U.S. territories should be omitted. DC, while not a state, is in the mainland, so it should remain included. Convert this to data.frame()

```
# First, create a vector of states that will be removed from data set.
remove <- c("AK", "AS", "HI", "GU", "PR", "VI")

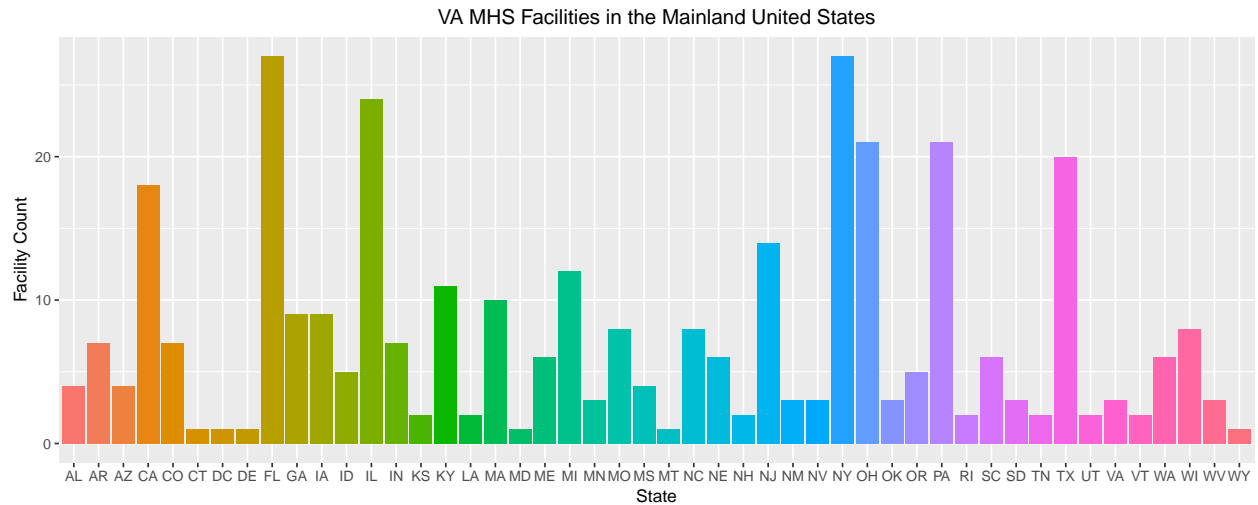
# Create Mainland vector, used to filter below.
mainlandStates <- states[!states %in% remove]

# It seems I have to spell out the VA string? Painful. Codebook says value should be 6 ...
vaMainLand <- mh2015_puf %>% filter(FACILITYTYPE == "Veterans Administration medical center (VAMC) or o

# Create count table for states, give it names. Will use for ggplot below.
stateTable <- vaMainLand %>% count(vaMainLand$LST)
names(stateTable) <- c("State", "Count")
```

- d. Create a ggplot barchart of this filtered data set. Vary the bar's colors by what State it has listed. Give it an appropriately professional title that is centered. Make sure you have informative axis labels. The State axis should be readable, not layered over each other. You're welcome to have a legend or not.

```
statePlot <-ggplot(data=stateTable, aes(x=State, y=Count, fill=State))
statePlot +
  geom_bar(stat="identity") +
  ylab("Facility Count") +
  xlab("State") +
  ggtitle("VA MHS Facilities in the Mainland United States") +
  theme(plot.title = element_text(hjust = 0.5), legend.position = "none")
```



## 2 - Cleaning and Bringing in New Features (60%):

- a) This graph (1D) might be somewhat misleading, as bigger states may have more hospitals, but could be more sparsely located. Read `statesize.csv` into your R environment. This contains essentially a vector of square miles for each state. In trying to merge it with your `data.frame()` from 1C, you find that they don't match. Use `paste()` on your LST column in 1C to see what the matter is, and write what you observe in a comment.

```
# Please Note! I already cleaned the LST column in problem 1 above, the whitespace issue.
```

```
# Read file
```

```
stateSize <- read.csv("statesize.csv")
```

- b) Correct the problem with the LST column using any method in R that is programmatic and easily understandable. Once you have made these state abbreviations identical to `statesize.csv`'s `Abbrev` column, merge the `data.frame()` from 1C and `statesize.csv` in order to add size information.

```
# Column correction already done in problem 1 above.
```

```
# Merge data sets, but use inner_join instead of merge
```

```
finalTable <- inner_join(stateTable, stateSize, by = c("State"= "Abbrev"))
```

```
## Warning: Column `State`/`Abbrev` joining character vector and factor,  
## coercing into character vector
```

```
# Could have used the merge command below as well
```

```
# finalTable <- merge(stateTable, stateSize, by.x = "State", by.y = "Abbrev", all=FALSE)
```

- c) Calculate a new variable in your combined `data.frame()` which indicates the VA hospitals per thousand square miles.

```
finalTable$NumPer1000SquareMiles <- finalTable$Count / finalTable$SqMiles * 1000
```

- d) Create another ggplot which considers the VAs per square thousand miles, rather than just frequency.
- Make sure the State axis is readable, like before. Change the title and axes as appropriate.
  - Modify the ggplot syntax to make your bars in descending order (there are StackOverflow topics for this, and I have demonstrated how in Live Coding in prior classes).
  - Color-code the bars based on Region (see the merged `data.frame()`)—however, change the color scheme from the default. Any set of colors is fine, so long as it is readable.
  - Keep the legend—you should have four regions and therefore four colors.

```
sizePlot <- ggplot(data=finalTable, aes(x=State, y=NumPer1000SquareMiles, fill=Region))
```

```
sizePlot +
```

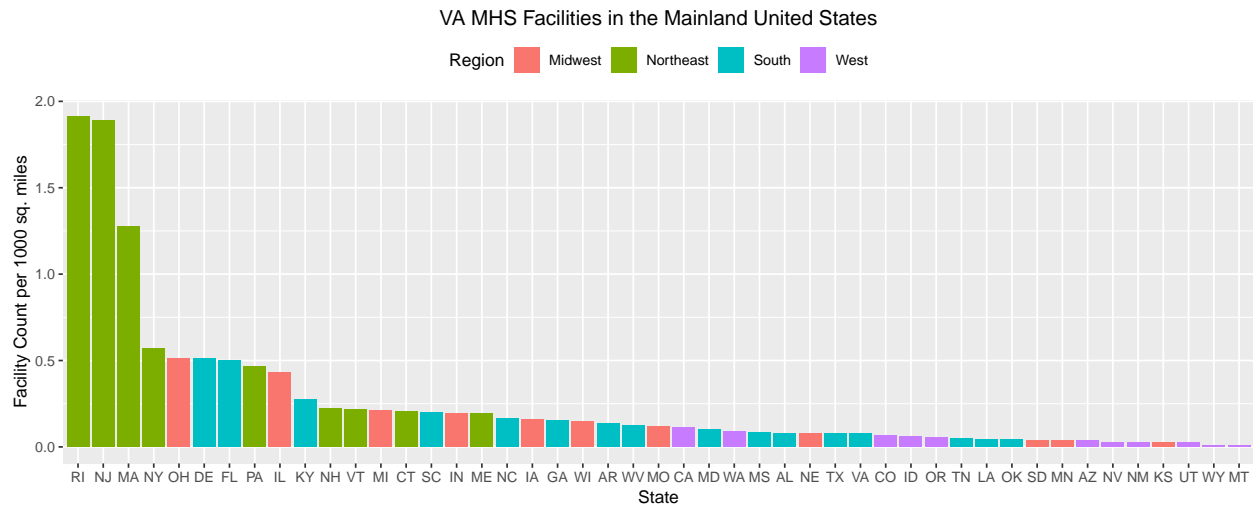
```
  geom_bar(aes(reorder(State,-NumPer1000SquareMiles)), stat="identity", position = position_dodge(1)) +
```

```
  ylab("Facility Count per 1000 sq. miles") +
```

```
  xlab("State") +
```

```
  ggtitle("VA MHS Facilities in the Mainland United States") +
```

```
  theme(plot.title = element_text(hjust = 0.5), legend.position = "top")
```



- e) What patterns do you see? By this metric, is there any region that seems relatively high for VA medical centers per thousand square miles? How about low? Given these data, what advice might you give your boss before you start modeling (and why)?

**Answer:** It's pretty clear that the Northeast region has a high amount and the West region is low. But California is a pretty big state, and Rhode Island is tiny. I think before making any recommendations one needs to consider the population of people in the area (population density). Perhaps look at the county level for a given state, and consider the population AND land size of that county.