# Homework 1

**Principles of Data Visualization and Introduction to ggplot2**

I have provided you with data about the 5,000 fastest growing companies in the US, as compiled by Inc. magazine. lets read this in:

```
inc <- read.csv("https://raw.githubusercontent.com/charleyferrari/CUNY_DATA_608/master/module1/Data/inc5
```

And lets preview this data:

```
head(inc)
```

```
##   Rank                        Name Growth_Rate   Revenue
## 1    1                        Fuhu      421.48 1.179e+08
## 2    2         FederalConference.com      248.31 4.960e+07
## 3    3               The HCI Group      245.45 2.550e+07
## 4    4                     Bridger      233.08 1.900e+09
## 5    5                      DataXu      213.37 8.700e+07
## 6    6 MileStone Community Builders      179.38 4.570e+07
##                       Industry Employees         City State
## 1 Consumer Products & Services       104   El Segundo    CA
## 2          Government Services        51     Dumfries    VA
## 3                       Health       132 Jacksonville    FL
## 4                       Energy        50      Addison    TX
## 5        Advertising & Marketing       220       Boston    MA
## 6                  Real Estate        63       Austin    TX
```

```
summary(inc)
```

```
##       Rank           Name           Growth_Rate         Revenue
##  Min.   :   1   Length:5001        Min.   :  0.340   Min.   :2.000e+06
##  1st Qu.:1252   Class :character   1st Qu.:  0.770   1st Qu.:5.100e+06
##  Median :2502   Mode  :character   Median :  1.420   Median :1.090e+07
##  Mean   :2502                      Mean   :  4.612   Mean   :4.822e+07
##  3rd Qu.:3751                      3rd Qu.:  3.290   3rd Qu.:2.860e+07
##  Max.   :5000                      Max.   :421.480   Max.   :1.010e+10
##
##    Industry           Employees          City               State
##  Length:5001        Min.   :    1.0   Length:5001        Length:5001
##  Class :character   1st Qu.:   25.0   Class :character   Class :character
##  Mode  :character   Median :   53.0   Mode  :character   Mode  :character
##                     Mean   :  232.7
##                     3rd Qu.:  132.0
##                     Max.   :66803.0
##                     NA's   :12
```

Think a bit on what these summaries mean. Use the space below to add some more relevant non-visual exploratory information you think helps you understand this data:

```
library(dplyr)
```

```
##
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:stats':
##
##     filter, lag
```

```
## The following objects are masked from 'package:base':
##
##     intersect, setdiff, setequal, union
```

```
library(ggplot2)

#showing sample of the table
str(inc)
```

```
## 'data.frame':    5001 obs. of  8 variables:
##  $ Rank       : int  1 2 3 4 5 6 7 8 9 10 ...
##  $ Name       : chr  "Fuhu" "FederalConference.com" "The HCI Group" "Bridger" ...
##  $ Growth_Rate: num  421 248 245 233 213 ...
##  $ Revenue    : num  1.18e+08 4.96e+07 2.55e+07 1.90e+09 8.70e+07 ...
##  $ Industry   : chr  "Consumer Products & Services" "Government Services" "Health" "Energy" ...
##  $ Employees  : int  104 51 132 50 220 63 27 75 97 15 ...
##  $ City       : chr  "El Segundo" "Dumfries" "Jacksonville" "Addison" ...
##  $ State      : chr  "CA" "VA" "FL" "TX" ...
```

```
# selection of industries
with(inc, table(Industry))
```

```
## Industry
##        Advertising & Marketing Business Products & Services
##                            471                           482
##              Computer Hardware                  Construction
##                             44                           187
## Consumer Products & Services                     Education
##                            203                            83
##                         Energy                   Engineering
##                            109                            74
##         Environmental Services            Financial Services
##                             51                           260
##                Food & Beverage           Government Services
##                            131                           202
##                         Health               Human Resources
##                            355                           196
##                      Insurance                   IT Services
##                             50                           733
```

```
##    Logistics & Transportation                    Manufacturing
##                             155                             256
##                           Media                     Real Estate
##                              54                              96
##                          Retail                        Security
##                             203                              73
##                        Software              Telecommunications
##                             342                             129
##           Travel & Hospitality
##                              62
```

```
with(inc, table(State))
```
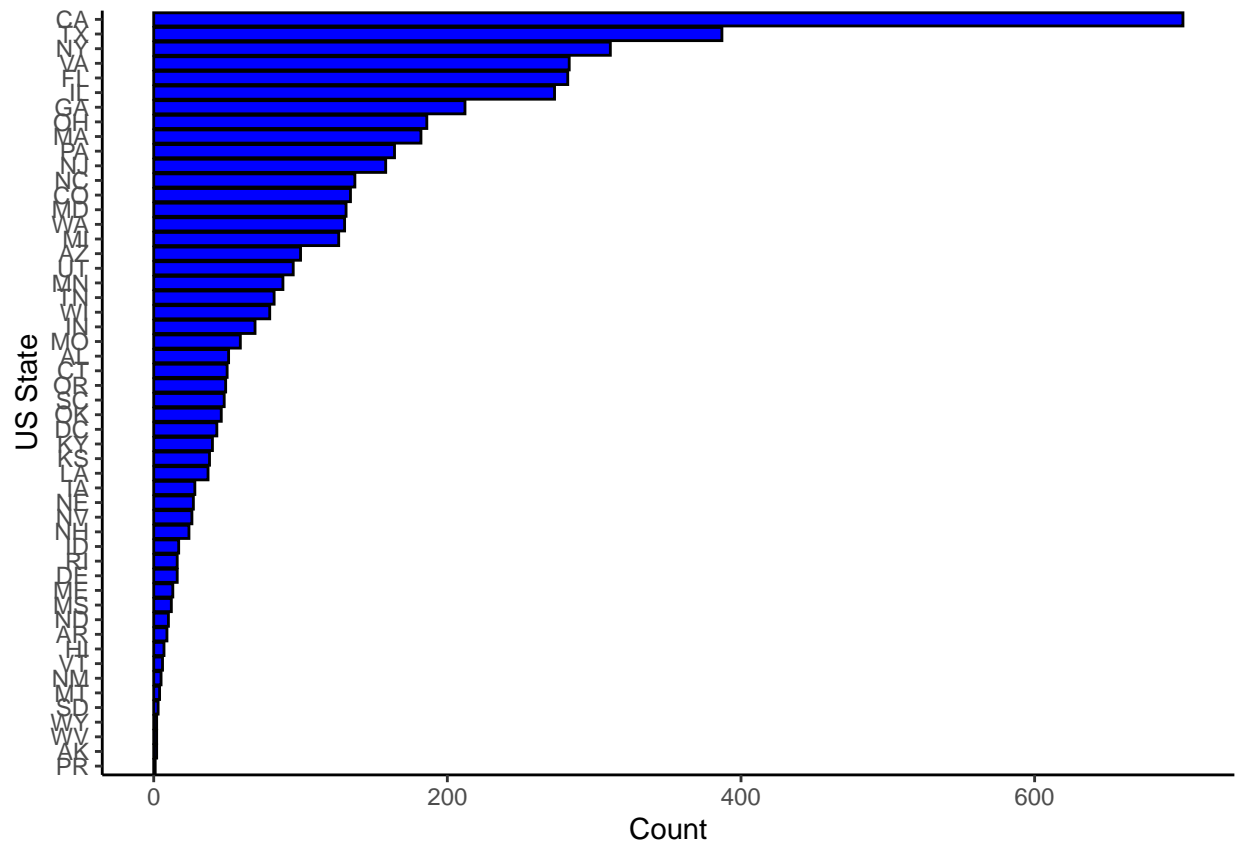
```
## State
##  AK  AL  AR  AZ  CA  CO  CT  DC  DE  FL  GA  HI  IA  ID  IL  IN  KS  KY  LA  MA
##   2  51   9 100 701 134  50  43  16 282 212   7  28  17 273  69  38  40  37 182
##  MD  ME  MI  MN  MO  MS  MT  NC  ND  NE  NH  NJ  NM  NV  NY  OH  OK  OR  PA  PR
## 131  13 126  88  59  12   4 137  10  27  24 158   5  26 311 186  46  49 164   1
##  RI  SC  SD  TN  TX  UT  VA  VT  WA  WI  WV  WY
##  16  48   3  82 387  95 283   6 130  79   2   2
```

## Question 1

Create a graph that shows the distribution of companies in the dataset by State (ie how many are in each state). There are a lot of States, so consider which axis you should use. This visualization is ultimately going to be consumed on a 'portrait' oriented screen (ie taller than wide), which should further guide your layout choices.

```r
df <- inc$State %>% table() %>% as.data.frame(stringsAsFactors=FALSE)

colnames(df) <- c('State', 'Count')
ggplot(df, aes(x=reorder(State, Count),y=Count, color=State)) +
  geom_bar(stat='identity', color = 'black', fill='blue')+
  coord_flip() +
  xlab('US State')+
  theme_classic()
```
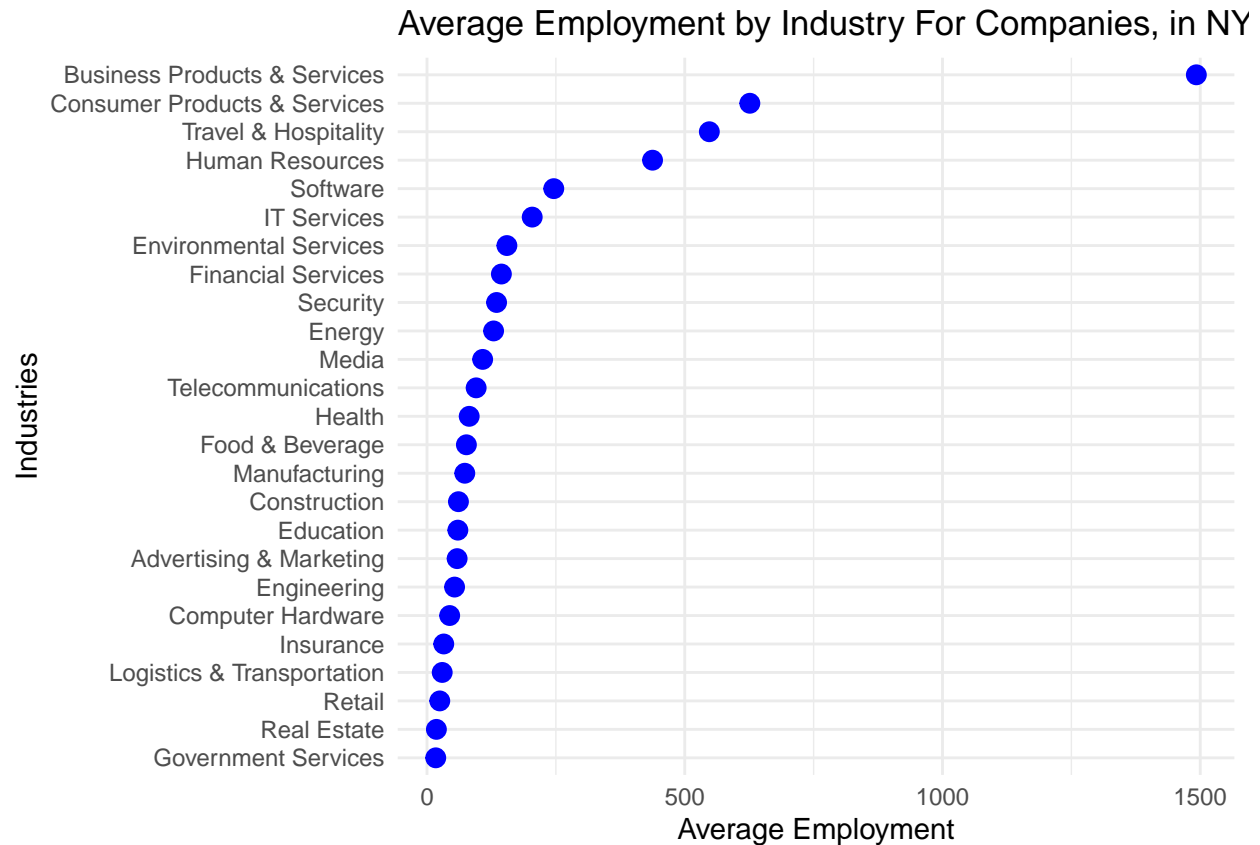
## Quesiton 2

Lets dig in on the state with the 3rd most companies in the data set. Imagine you work for the state and are interested in how many people are employed by companies in different industries. Create a plot that shows the average and/or median employment by industry for companies in this state (only use cases with full data, use R's `complete.cases()` function.) In addition to this, your graph should show how variable the ranges are, and you should deal with outliers.

```
state <- inc %>% count(State) %>% arrange(desc(n))
state3rd <- state$State[3]

df2 <- inc[complete.cases(inc), ]

df2 %>% filter(State == state3rd) %>% group_by(Industry) %>% summarise(avg = mean(Employees)) %>% ggplot
```

## Average Employment by Industry For Companies, in NY



## Question 3

Now imagine you work for an investor and want to see which industries generate the most revenue per employee. Create a chart that makes this information clear. Once again, the distribution per industry should be shown.

```
avgR <-inc[complete.cases(inc),] %>%
                group_by(Industry) %>%
                summarise(R=sum(Revenue),E=sum(Employees)) %>%
                mutate(revenue = R/E)

ggplot(avgR, aes(x =reorder(Industry, revenue), y = revenue)) +
  geom_bar(stat="identity", width=0.5, fill="blue") +coord_flip()+
  ggtitle("Revenue Per Employee ")+
  labs(x="Industry",y="Mean")+  theme(axis.text.x = element_text(angle = 60, hjust = 1))
```

## Revenue Per Employee