

Project 1

Cesar L. Espitia

2/11/2018

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/inc5000_data.csv", header= TRUE)
```

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
## Min.      :    1   (Add)ventures                :    1   Min.      :    0.340
## 1st Qu.:1252   @Properties                        :    1   1st Qu.:    0.770
## Median :2502   1-Stop Translation USA:          :    1   Median :    1.420
## Mean    :2502   110 Consulting                   :    1   Mean    :    4.612
## 3rd Qu.:3751   11thStreetCoffee.com             :    1   3rd Qu.:    3.290
## Max.     :5000   123 Exteriors                   :    1   Max.     :421.480
##                                     (Other)       :4995
##
##           Revenue                               Industry       Employees
## Min.      :2.000e+06   IT Services                : 733   Min.      :    1.0
## 1st Qu.:5.100e+06   Business Products & Services: 482   1st Qu.:    25.0
## Median :1.090e+07   Advertising & Marketing                : 471   Median :    53.0
## Mean    :4.822e+07   Health                                  : 355   Mean    :   232.7
## 3rd Qu.:2.860e+07   Software                                  : 342   3rd Qu.:   132.0
## Max.     :1.010e+10   Financial Services                      : 260   Max.     :66803.0
##                                     (Other)       :2358   NA's     :12
##
##           City                               State
## New York      : 160   CA                : 701
## Chicago       :  90   TX                : 387
## Austin        :  88   NY                : 311
## Houston       :  76   VA                : 283
## San Francisco:  75   FL                : 282
## Atlanta       :  74   IL                : 273
## (Other)       :4438   (Other):2764
```

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:

```
# The following provides additional variables for Revenue and No. of Employees to help further describe the shape of the data.
```

```
library(pastecs)
```

```
## Loading required package: boot
```

```
stat.desc(inc$Revenue)
```

```
##           nbr.val      nbr.null      nbr.na           min           max
## 5.001000e+03  0.000000e+00  0.000000e+00  2.000000e+06  1.010000e+10
##           range           sum      median           mean      SE.mean
## 1.009800e+10  2.411609e+11  1.090000e+07  4.822254e+07  3.401441e+06
## CI.mean.0.95           var      std.dev      coef.var
## 6.668317e+06  5.786059e+16  2.405423e+08  4.988172e+00
```

```
stat.desc(inc$Employees)
```

```
##          nbr.val      nbr.null      nbr.na          min          max
## 4.989000e+03 0.000000e+00 1.200000e+01 1.000000e+00 6.680300e+04
##          range          sum          median          mean          SE.mean
## 6.680200e+04 1.161030e+06 5.300000e+01 2.327180e+02 1.915720e+01
## CI.mean.0.95          var          std.dev          coef.var
## 3.755654e+01 1.830955e+06 1.353128e+03 5.814454e+00
```

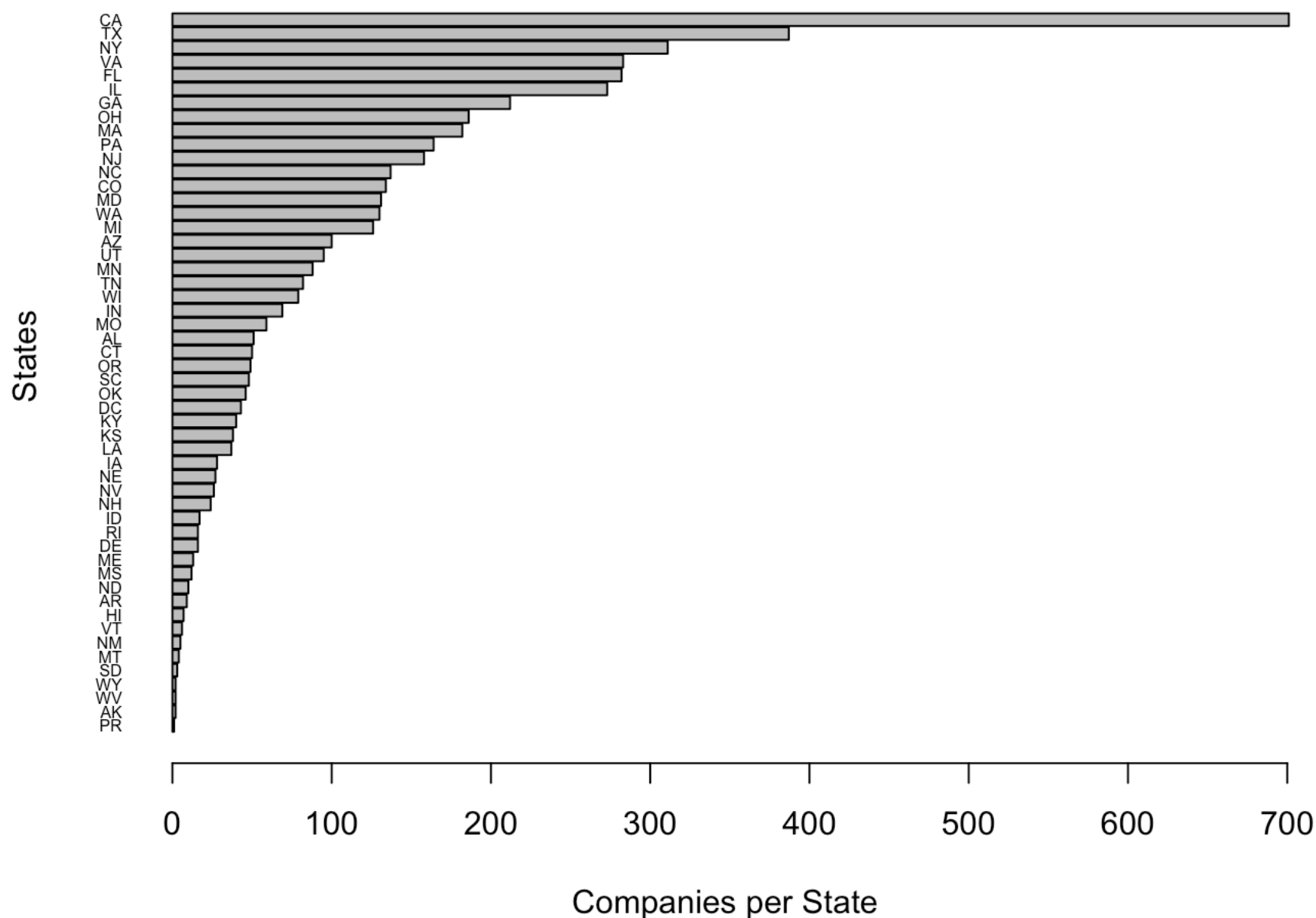
Question 1

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.

```
# Answer Question 1 here
library(plyr)
library(ggplot2)
state <- count(inc, "State")

state<-state[order(state[,2],decreasing=FALSE),]

barplot(state$freq, ylab= "States", xlab="Companies per State", horiz=TRUE, beside =
TRUE, space=0.1, ylim = c(0,60), yaxp=c(0,5,1),
names.arg=state$State, cex.names=0.5, las=1)
```



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.

```
library(ggplot2)
library(gridExtra)
state[order(-state$freq), ][3, ]

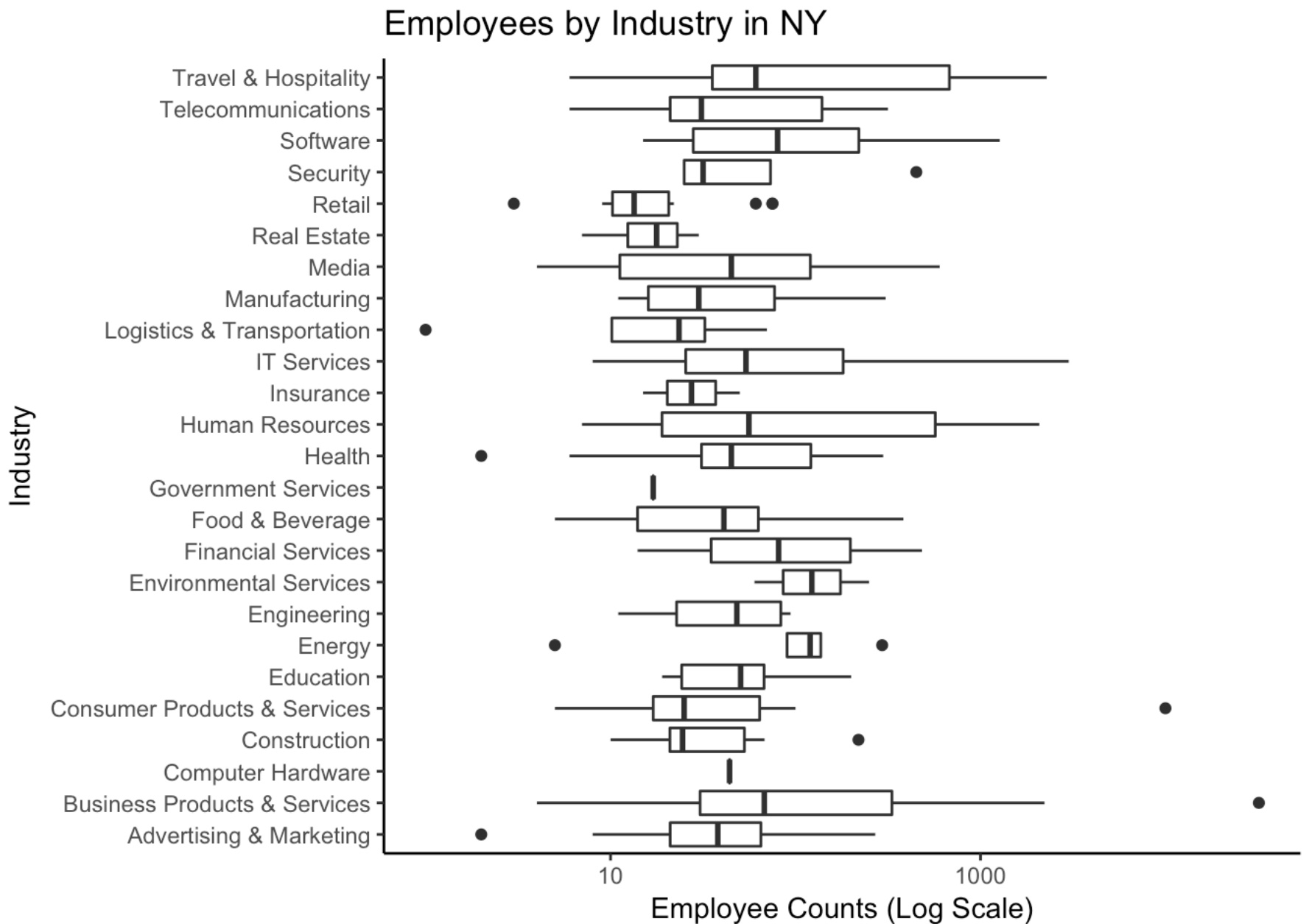
##      State freq
## 35      NY   311
```

```
#NY is third
complete.cases(inc[1:8,])
```

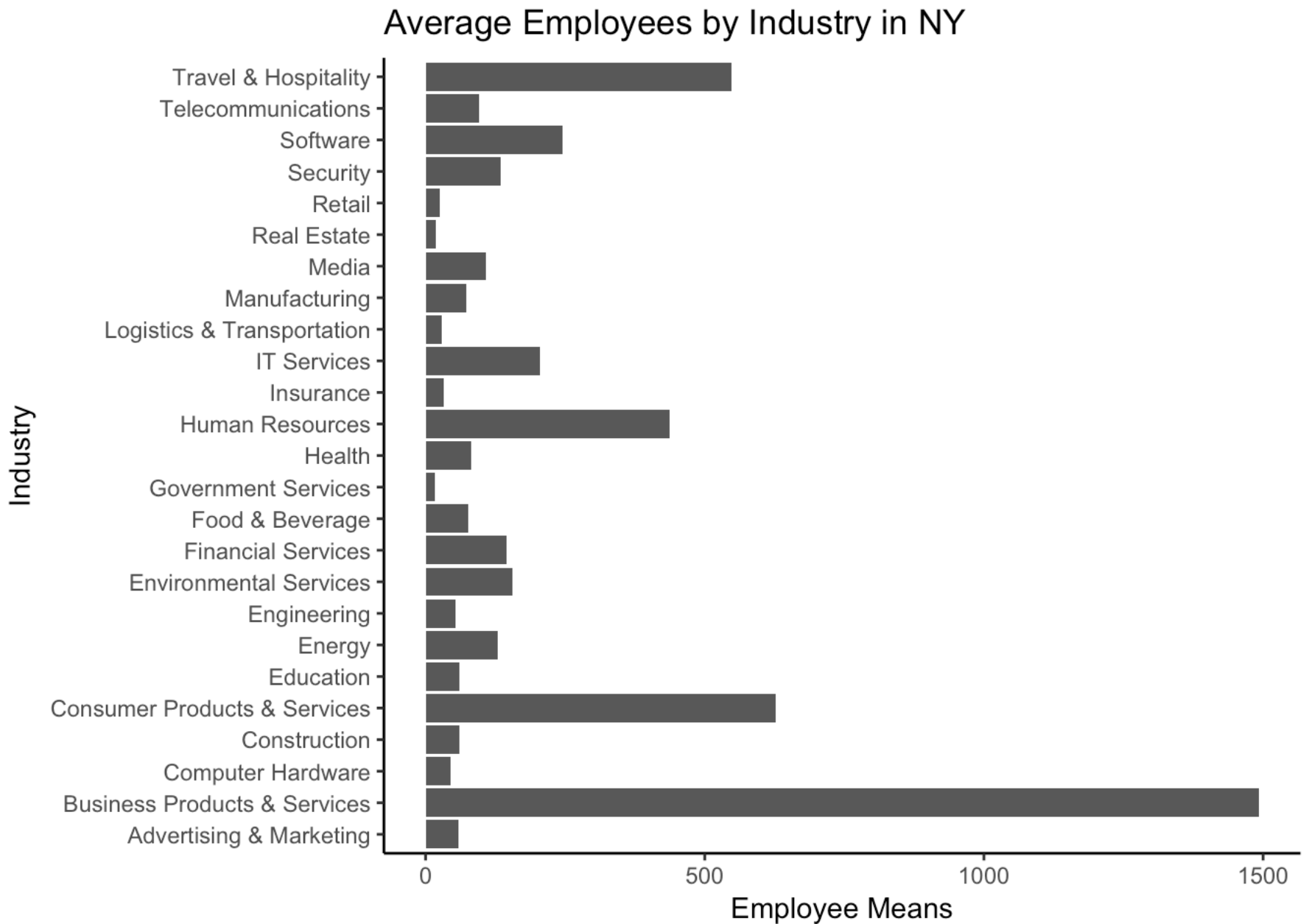
```
## [1] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
```

```
#all cases are complete
NY <- inc[ which(inc$State=="NY"), ]
NYc <- NY[c("Industry","Employees")]
IM <- aggregate(NYc$Employees, by=list(NY$Industry),
  FUN=mean, na.rm=TRUE)
colnames(IM) <- c("Industry","EmployeeMean")

p <- ggplot(NY,aes(NY$Industry, NY$Employees))+geom_boxplot()+scale_y_log10()+theme_c
lassic()+labs(title="Employees by Industry in NY", x="Industry", y="Employee Counts (
Log Scale)")
p+coord_flip()
```



```
e <- ggplot(IM,aes(IM$Industry, IM$EmployeeMean))+geom_bar(stat="identity")+theme_classic()+labs(title="Average Employees by Industry in NY", x="Industry", y="Employee Means")
e+coord_flip()
```



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.

```

# Answer Question 3 here
library(ggplot2)
library(gridExtra)

incc <- inc[c("Industry", "Employees", "Revenue")]
IM <- aggregate(incc[c("Employees", "Revenue")], by=list(incc$Industry),
  FUN=sum, na.rm=TRUE)
colnames(IM) <- c("Industry", "EmployeeSum", "RevenueSum")
IM$RevPerEmp <- with(IM, IM$RevenueSum/IM$EmployeeSum/1000)

IM <- IM[order(IM[,4], decreasing=TRUE), ]

e <- ggplot(IM, aes(IM$Industry, IM$RevPerEmp))+geom_bar(stat="identity")+theme_classic()+labs(title="Revenue per Employee by Industry in NY", x="Industry", y="In 000's")

e+coord_flip()+geom_hline(yintercept=mean(IM$RevPerEmp), linetype="dashed",
  color = "black", size=0.5)+ geom_text(aes(0, mean(IM$RevPerEmp), label = paste(round(mean(IM$RevPerEmp), 0), "Nationwide Average")), hjust=-0.1, vjust = -1, color = "blue")

```

