

# R Notebook

## 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/inc.csv")
```

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:

```
##
## 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

# Listing of states and number of companies within each state, sorted from having most companies to least
inc %>% count(State) %>% arrange(desc(n))

## # A tibble: 52 x 2
##   State      n
##   <fct> <int>
## 1 CA       701
## 2 TX       387
## 3 NY       311
## 4 VA       283
## 5 FL       282
## 6 IL       273
## 7 GA       212
## 8 OH       186
## 9 MA       182
## 10 PA      164
## # ... with 42 more rows

# Summary statistics for the number of companies in each state
inc %>% count(State) %>% summary()

##           State           n
## AK      : 1   Min.      : 1.00
## AL      : 1   1st Qu.: 15.25
## AR      : 1   Median : 48.50
## AZ      : 1   Mean     : 96.17
## CA      : 1   3rd Qu.:131.75
## CO      : 1   Max.      :701.00
## (Other):46
```

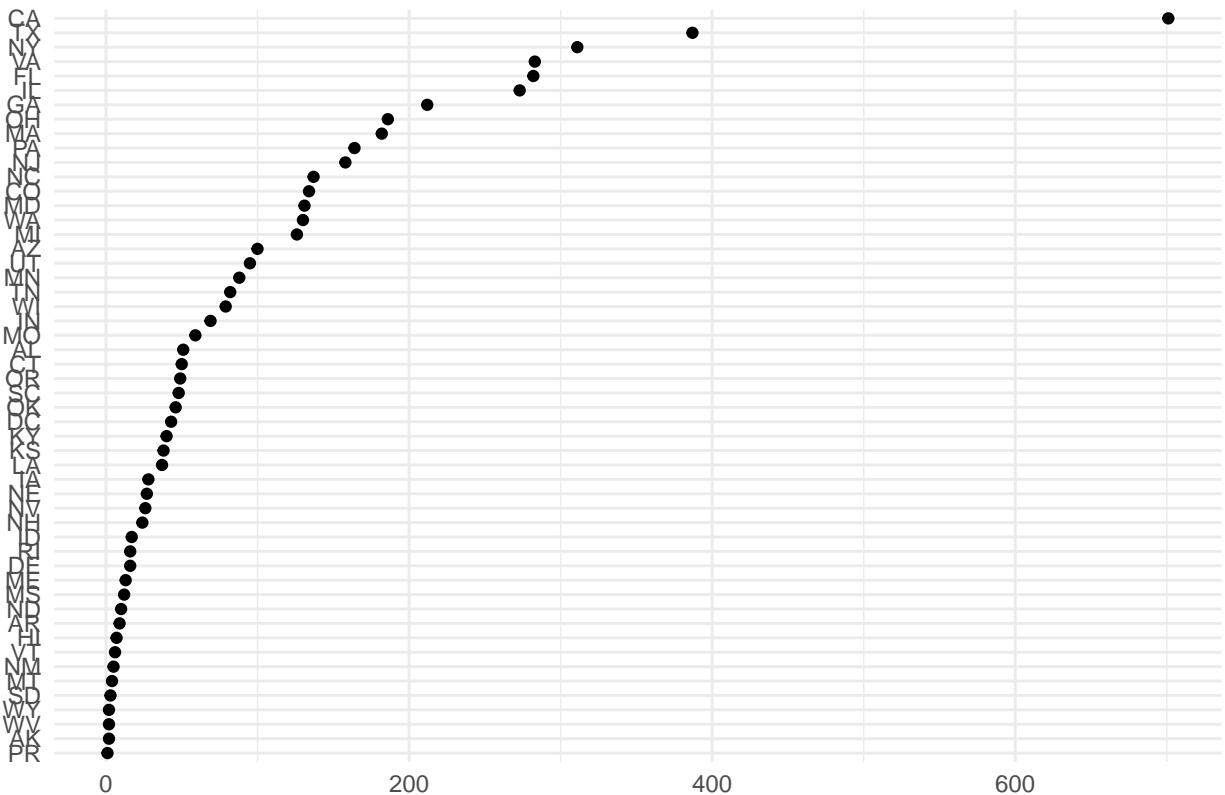
## Question 1

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

inc %>% count(State) %>% ggplot(aes(x = n, y = reorder(State, n))) + geom_point() + ggtitle("Distribution of companies by State")
```

Distribution of Companies by State



## Quesiton 2

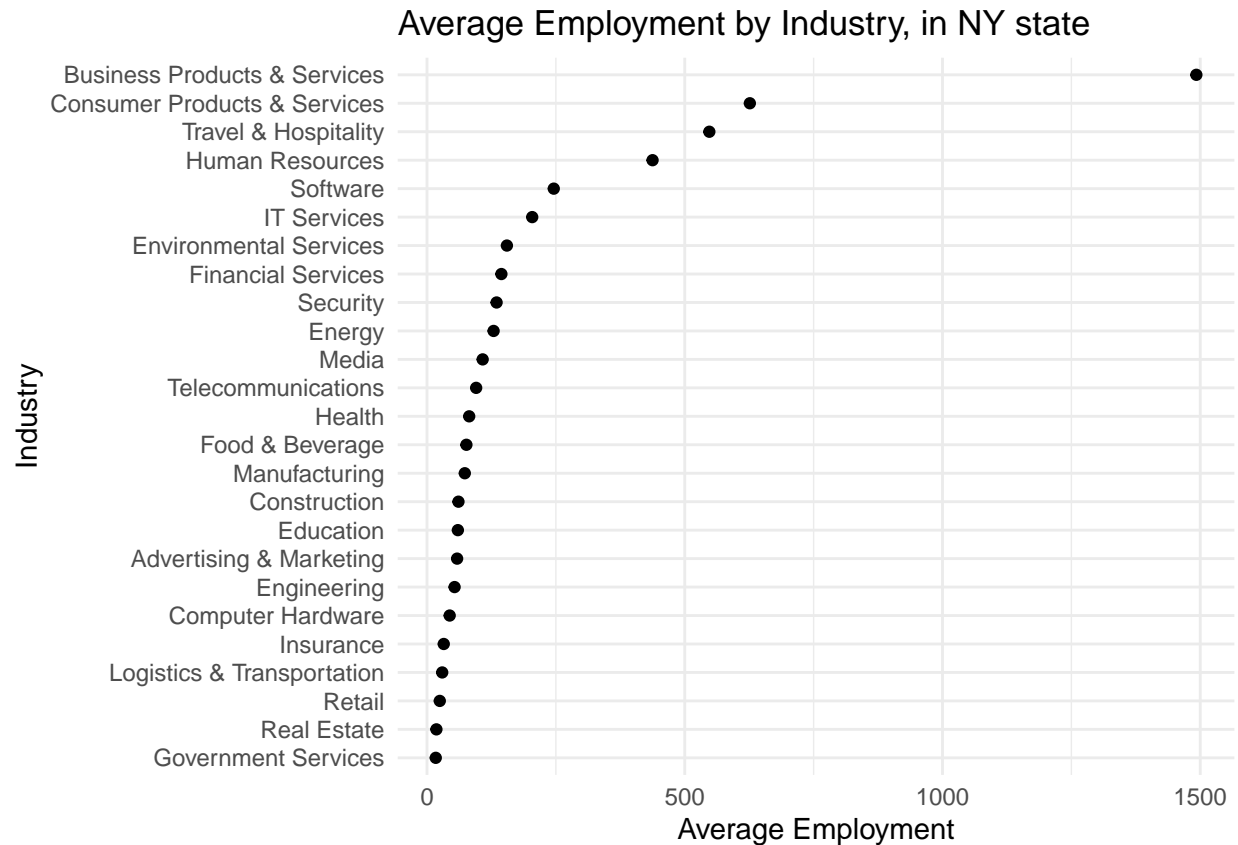
Lets dig in on the state with the 3rd most companies in the data set. Imagine you work for the state and are interested in how many people are employed by companies in different industries. Create a plot that shows the average and/or median employment by industry for companies in this state (only use 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.

*# Answer Question 2 here*

```
# Find the state with the 3rd most companies
s <- inc %>% count(State) %>% arrange(desc(n))
st <- s$State[3] # 3rd state
```

```
# Get full data, using the R's `complete.cases()` function
inc2 <- inc[complete.cases(inc), ]
```

```
inc2 %>% filter(State == st) %>% group_by(Industry) %>% summarise(avg = mean(Employees)) %>% ggplot(aes
```



### 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*

```
inc2 %>% group_by(Industry) %>% summarise(rpe = sum(Revenue) / sum(Employees)) %>% ggplot(aes(x = rpe,
```

# Revenue per Employee in each Industry

