

DATA608: Module 1

Donald Butler

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      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(tidyverse)
library(scales)
```

Which state are represented in the dataset?

```
inc %>% count(State)
```

```
##      State      n
## 1      AK       2
## 2      AL      51
## 3      AR       9
## 4      AZ     100
## 5      CA    701
## 6      CO     134
## 7      CT      50
## 8      DC      43
## 9      DE      16
## 10     FL     282
## 11     GA     212
## 12     HI       7
## 13     IA      28
## 14     ID      17
## 15     IL     273
## 16     IN      69
## 17     KS      38
## 18     KY      40
## 19     LA      37
## 20     MA     182
## 21     MD     131
## 22     ME      13
## 23     MI     126
## 24     MN      88
## 25     MO      59
## 26     MS      12
## 27     MT       4
## 28     NC     137
## 29     ND      10
## 30     NE      27
## 31     NH      24
## 32     NJ     158
## 33     NM       5
## 34     NV      26
## 35     NY     311
## 36     OH     186
## 37     OK      46
## 38     OR      49
## 39     PA     164
## 40     PR       1
## 41     RI      16
## 42     SC      48
```

```
## 43 SD 3
## 44 TN 82
## 45 TX 387
## 46 UT 95
## 47 VA 283
## 48 VT 6
## 49 WA 130
## 50 WI 79
## 51 WV 2
## 52 WY 2
```

All 50 states are represented along with Washington DC and Puerto Rico.

Which industries are represented?

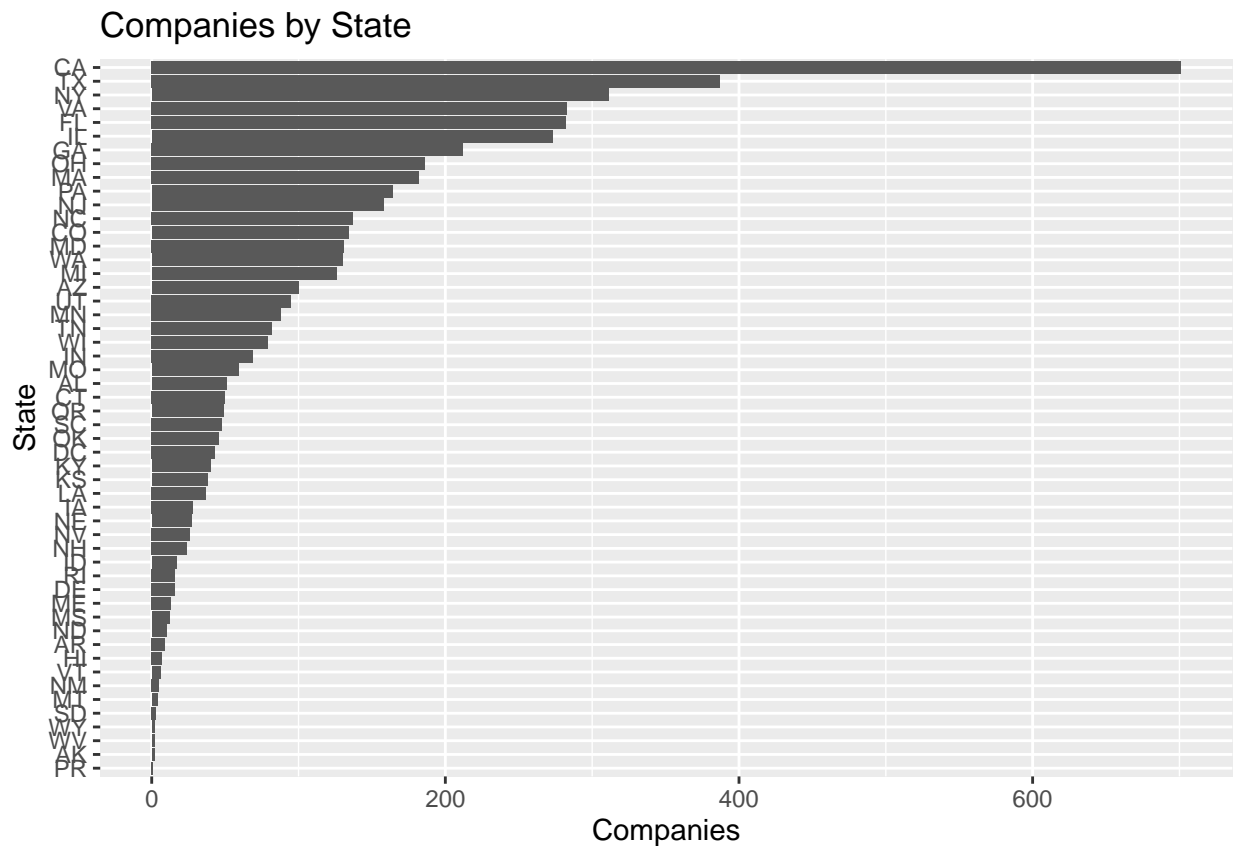
```
inc %>% count(Industry)
```

```
##           Industry  n
## 1 Advertising & Marketing 471
## 2 Business Products & Services 482
## 3 Computer Hardware 44
## 4 Construction 187
## 5 Consumer Products & Services 203
## 6 Education 83
## 7 Energy 109
## 8 Engineering 74
## 9 Environmental Services 51
## 10 Financial Services 260
## 11 Food & Beverage 131
## 12 Government Services 202
## 13 Health 355
## 14 Human Resources 196
## 15 Insurance 50
## 16 IT Services 733
## 17 Logistics & Transportation 155
## 18 Manufacturing 256
## 19 Media 54
## 20 Real Estate 96
## 21 Retail 203
## 22 Security 73
## 23 Software 342
## 24 Telecommunications 129
## 25 Travel & Hospitality 62
```

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.

```
inc %>%
  count(State) %>%
  ggplot(aes(reorder(State,n),n)) + geom_col() + coord_flip() + labs(title = "Companies by State", x =
```

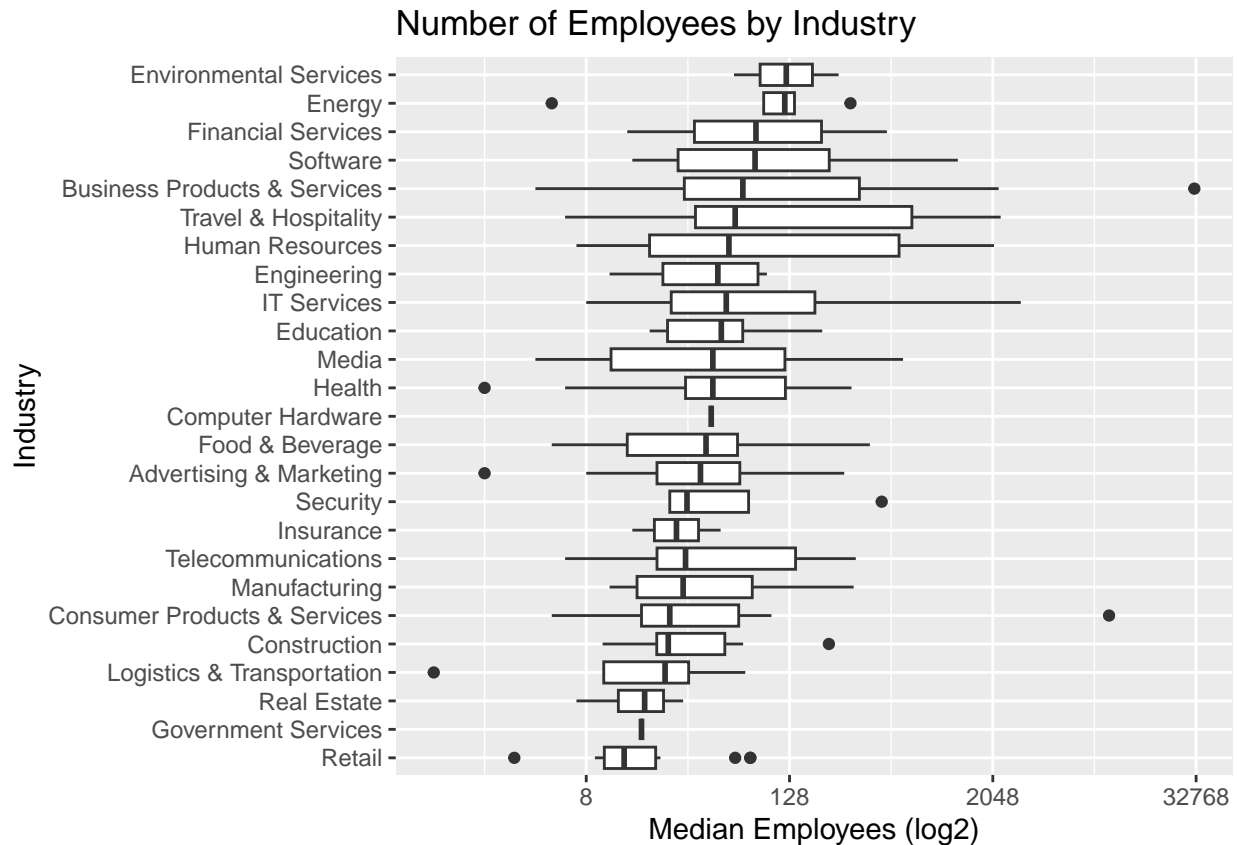


Question 2

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.

```
# NY is the 3rd most represented

inc %>%
  filter(State == "NY" & complete.cases(.)) %>%
  ggplot(aes(x = reorder(Industry, Employees, FUN=median), y = Employees)) +
  geom_boxplot() +
  scale_y_continuous(trans = log2_trans()) +
  labs(title = "Number of Employees by Industry", x = "Industry", y = "Median Employees (log2)") +
  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.

```
inc %>%
  filter(complete.cases(.)) %>%
  group_by(Industry) %>%
  summarise(TotalRevenue = sum(Revenue), TotalEmployees = sum(Employees)) %>%
  mutate(EmployeeRevenue = TotalRevenue / TotalEmployees / 1000) %>%
  ggplot(aes(x = reorder(Industry, EmployeeRevenue), y = EmployeeRevenue)) +
  geom_bar(stat = "identity") +
  labs(title = "Revenue per Employee by Industry", x = "Industry", y = "Revenue per Employee (thousands)") +
  coord_flip()
```

Revenue per Employee by Industry

