

Data 608 HW1

Alice Friedman

Sept 6., 2020

```
library(ggplot2)
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
```

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:

```
#class of each
lapply(inc, class)
```

```
## $Rank
## [1] "integer"
##
## $Name
## [1] "factor"
##
## $Growth_Rate
## [1] "numeric"
##
## $Revenue
## [1] "numeric"
##
## $Industry
## [1] "factor"
##
## $Employees
## [1] "integer"
##
## $City
```

```
## [1] "factor"
##
## $State
## [1] "factor"
```

```
sd(inc$Growth_Rate)
```

```
## [1] 14.12369
```

```
sd(inc$Revenue)
```

```
## [1] 240542281
```

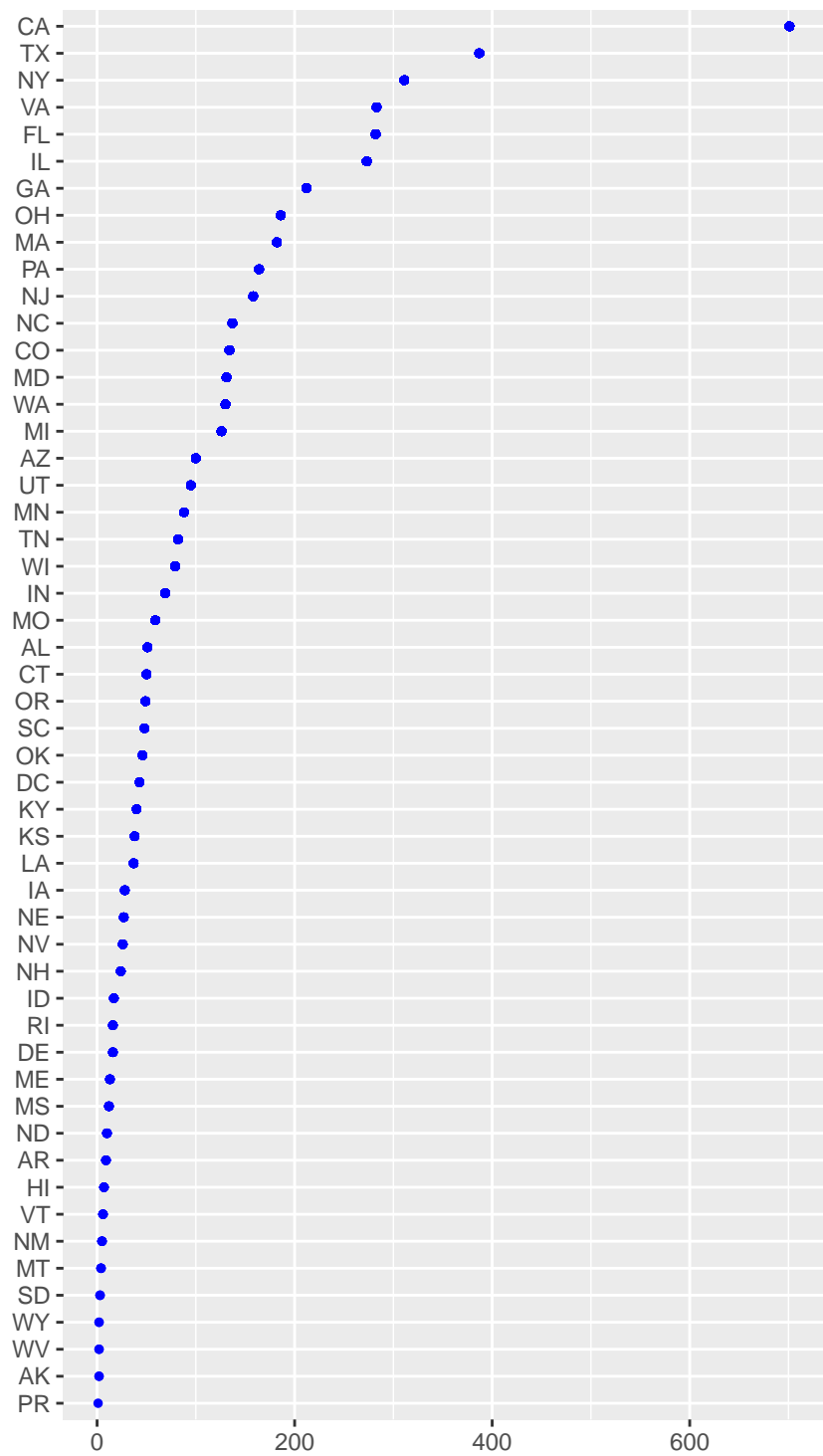
Question 1

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.

```
df <- inc %>% group_by(State) %>% mutate(Companies = n())

ggplot(df, aes(y = reorder(State, Companies), x=Companies)) +
  geom_point(color = "blue", size=1) +
  labs(title='Number of Companies Registered by State',
       caption = "Data source: 5,000 Fastest Growing Companies, Inc. Magazine") +
  theme(axis.title.x=element_blank(),
        axis.title.y=element_blank(),
        plot.margin = margin(0, 0, 0, 0))
```

Number of Companies Registered by State



Data source: 5,000 Fastest Growing Companies, Inc. Magazine

Question 2

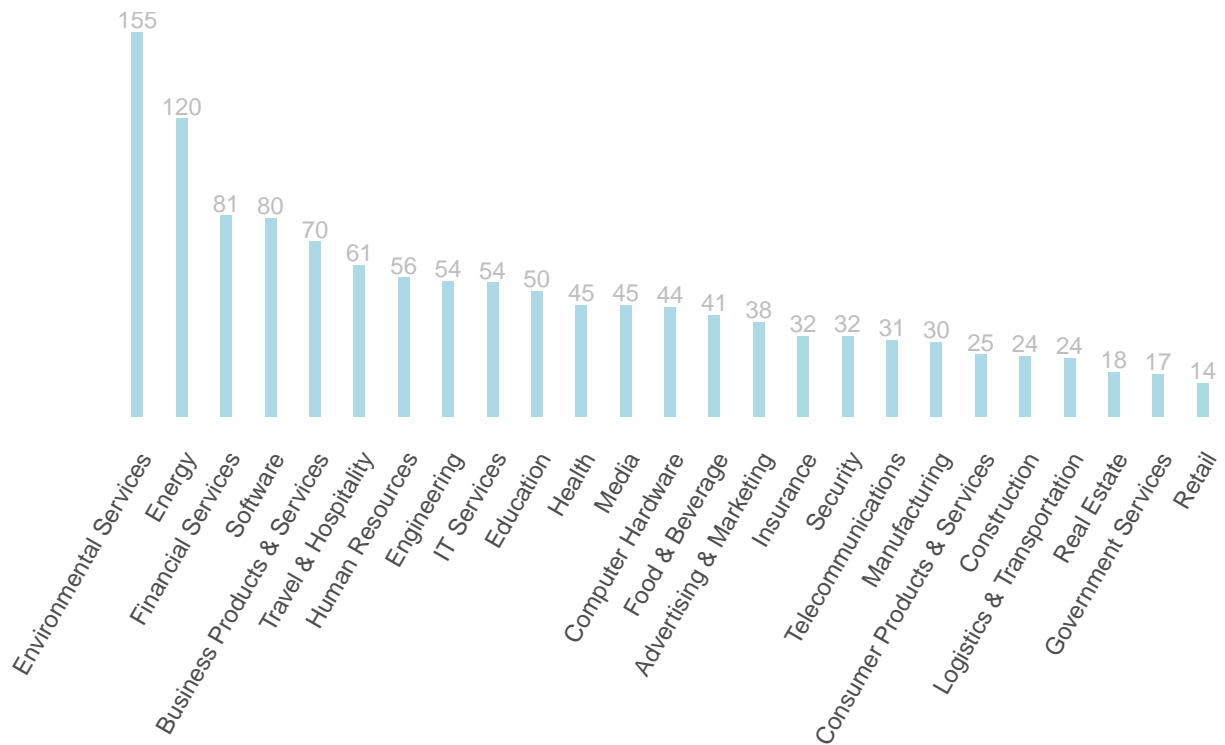
Lets dig in on the state with the 3rd most companies in the data set. Imagine you work for the state and are interested in how many people are employed by companies in different industries. Create a plot that

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 <- inc[complete.cases(inc), ] %>% dplyr::filter(State == 'NY')
```

```
ggplot(NY %>% group_by(Industry) %>%
  summarise(`Median Employees` = median(Employees))) +
  geom_col(
    aes(x=reorder(Industry, -`Median Employees`), y = `Median Employees`),
    fill = "light blue",
    width = 0.25) +
  geom_text(
    aes(x = Industry, y = `Median Employees`, label=round(`Median Employees`, digits = 0)),
    vjust=-0.25,
    size=3,
    color="gray") +
  theme(axis.text.x = element_text(angle = 60, hjust = 1),
        axis.text.y = element_blank(),
        axis.title=element_blank(),
        axis.ticks = element_blank(),
        panel.grid = element_blank(),
        panel.background = element_blank(),
        plot.margin = margin(1, 1, 15, 45)
  ) +
  labs(title = "Median Number of Employees per Company by Industry, NY",
       caption = "Data source: 5,000 Fastest Growing Companies, Inc. Magazine")
```

Median Number of Employees per Company by Industry, NY



Data source: 5,000 Fastest Growing Companies, Inc. Magazine

Question 3

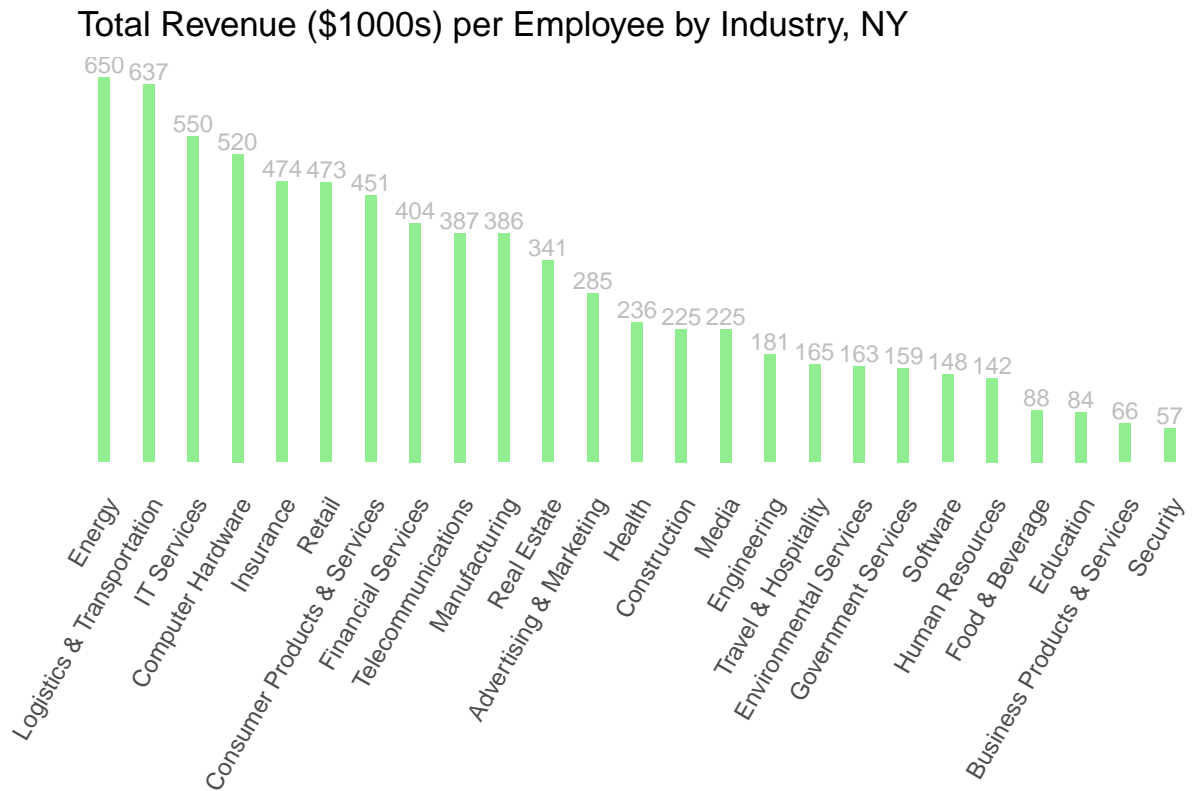
Now imagine you work for an investor and want to see which industries generate the most revenue per employee. Create a chart that makes this information clear. Once again, the distribution per industry should be shown.

```
ggplot(
  NY %>%
    group_by(Industry) %>%
    summarise(`Revenue Per Employee` = sum(Revenue)/sum(Employees))) +
  geom_col(
    aes(x=reorder(Industry, -`Revenue Per Employee`), y = `Revenue Per Employee`),
    fill = "light green",
    width = 0.25) +
  geom_text(
    aes(x = Industry,
        y = `Revenue Per Employee`,
        label=round(`Revenue Per Employee`, digits = -3)/1000),
    vjust=-0.25,
    size=3,
    color="gray") +
  theme(axis.text.x = element_text(angle = 60, hjust = 1),
        axis.text.y = element_blank(),
        axis.title=element_blank(),
        axis.ticks = element_blank(),
```

```

panel.grid = element_blank(),
panel.background = element_blank(),
plot.margin = margin(1, 1, 15, 45)
) +
labs(title = "Total Revenue ($1000s) per Employee by Industry, NY",
caption = "Data source: 5,000 Fastest Growing Companies, Inc. Magazine")

```



Data source: 5,000 Fastest Growing Companies, Inc. Magazine