

SESSION 13

DATA VISUALIZATION 1

R FOR SOCIAL DATA SCIENCE

JEFFREY ZIEGLER, PHD

ASSISTANT PROFESSOR IN POLITICAL SCIENCE & DATA SCIENCE
TRINITY COLLEGE DUBLIN

FALL 2022

ROAD MAP FOR TODAY

Last week:

- Data input and output
- Data frames and alternatives
- 'tidyverse' packages
- Working with tabular data
- Summary statistics

This time:

- Effectively tell a story with visuals
- Plotting in base R

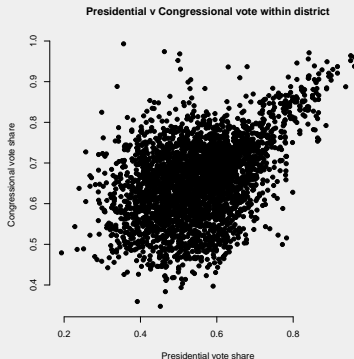
PRINCIPALS OF "GOOD" DATA VIS - TUFTE

- "Show" data
- Think about substance rather than about methodology, graphic design, technology, etc.
- Avoid distorting data
- Present many numbers (info) in a small space
- Make large data sets coherent
- Encourage eye to compare different pieces of data
- Reveal data at several levels of detail
- Integrate statistical and verbal descriptions of a data set

SCATTERPLOTS

■ Scatter plot can be created using `plot(x, y)`

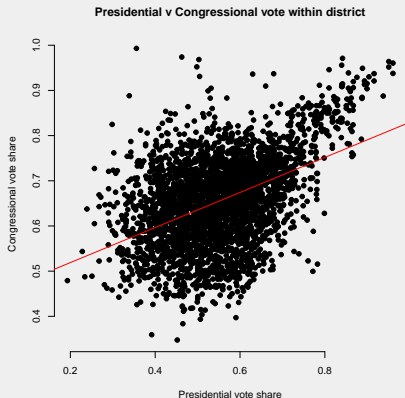
```
1 # Plot with main and axis titles
2 # Change point shape (pch = 19) and remove frame
3 plot(inc_local$presvote, inc_local$voteshare,
4       main = "Presidential v Congressional vote within district",
5       xlab = "Presidential vote share", ylab = "Congressional vote
share", pch = 19, frame = F)
```



SCATTERPLOTS

- `lm()` can be used to fit linear models between `y` and `x`
- Regression line can be added on the plot using `abline()`, which takes output of `lm()` as an argument

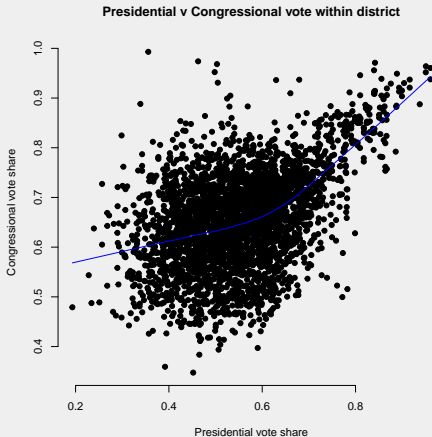
```
1 abline(lm(voteshare ~ presvote, data = inc_local), col = "red")
```



SCATTERPLOTS

■ Can also add a smoothing line using `loess()`

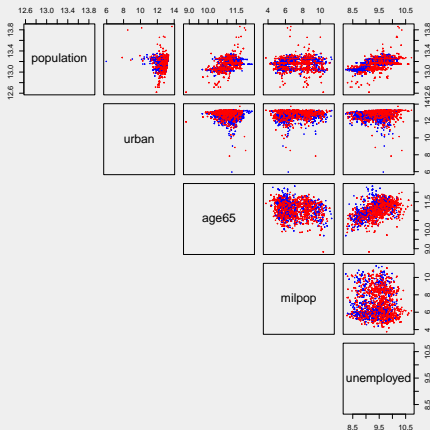
```
1 lines(lowess(inc_local$presvote, inc_local$voteshare), col = "blue")
```



SCATTERPLOTS OF MANY VARIABLES

- Can produce a matrix of scatter plots with `pairs()`

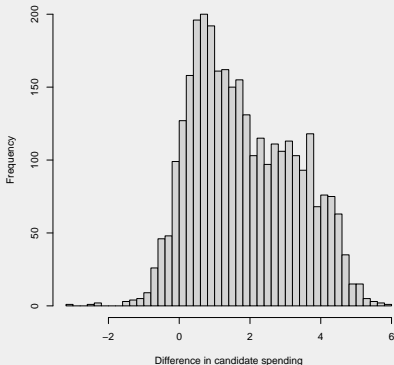
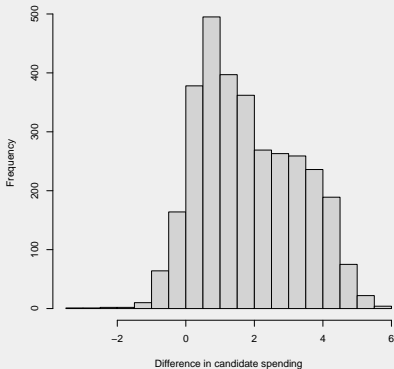
```
1 my_cols <- c("red", "blue")
2 pairs(inc_local[,15:ncol(inc_local)], pch = 19, cex = 0.15,
3       col = my_cols[as.factor(inc_local$south)],
4       lower.panel=NULL)
```



HISTOGRAMS

Summarize distribution of variable as count

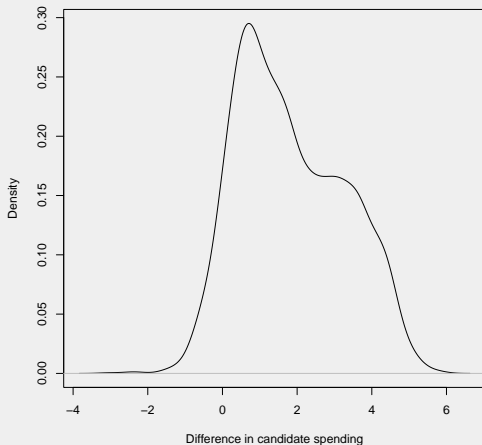
```
1 hist(inc_local$difflog, main = "", xlab = "Difference in candidate spending")  
1 hist(inc_local$difflog, main = "", xlab = "Difference in candidate spending",  
breaks=50)
```



DENSITY PLOT

Summarize distribution of variable as proportion

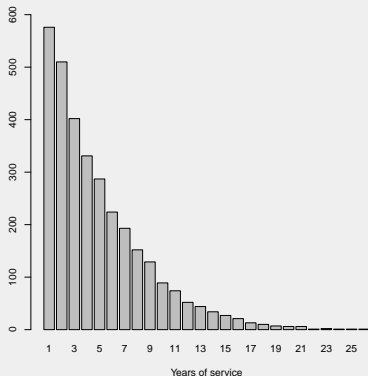
```
1 plot(density(inc_local$difflog), main = "", xlab = "Difference in  
candidate spending")
```



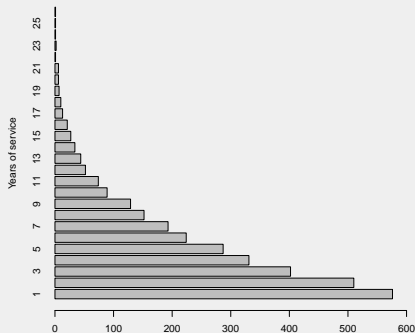
BAR CHART

Summarize count or proportion that you've already calculated

```
1 barplot(table(inc_local$seniority), ylim=c  
(0, 600), xlab="Years of service")
```



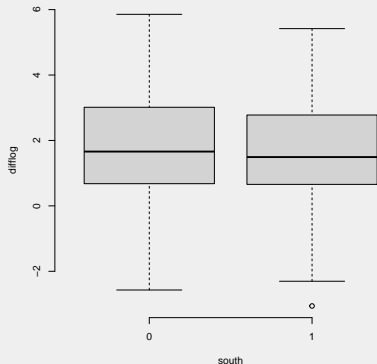
```
1 barplot(table(inc_local$seniority), xlim=c  
(0, 600), ylab="Years of service",  
horiz = T)
```



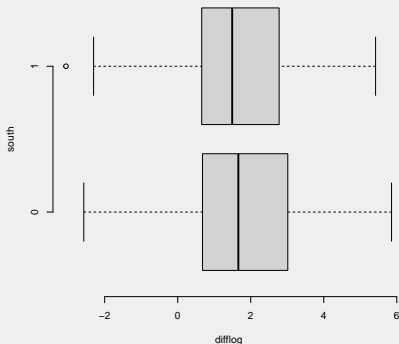
BOX PLOT

Summarize distribution (mean, quartiles, outliers) of variable by group

```
1 boxplot(difflog ~ south, data = inc_local,  
          frame = F)
```



```
1 boxplot(difflog ~ south, data = inc_local,  
          frame = F, horizontal = T)
```



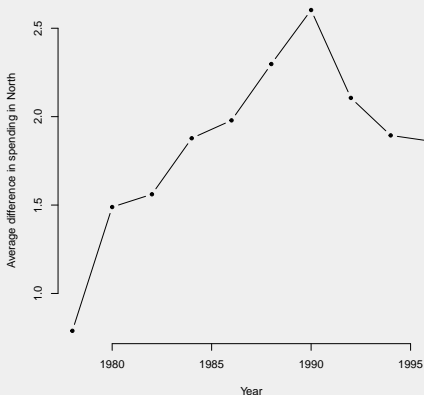
LINE CHART

Type: character indicating type of plotting

- “p” for points
- “l” for lines
- “b” for both points and lines
- “c” for empty points joined by lines
- “o” for overplotted points and lines
- “s” and “S” for stair steps
- “n” does not produce any points or lines
- lty: line types
 - ▶ Line types can either be specified as
 - Integer (0=blank, 1=solid (default), 2=dashed, 3=dotted, 4=dotdash, 5=longdash, 6=twodash), or
 - Character strings “blank”, “solid”, “dashed”, “dotted”, “dotdash”, “longdash”, or “twodash”

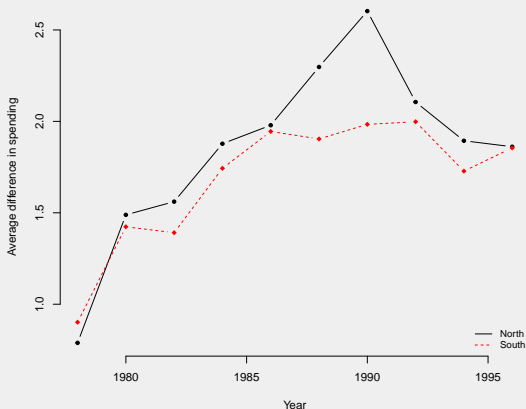
LINE CHART

```
1 south_year_avg <- inc_local %>%  
2   group_by(south, year) %>%  
3   summarize(avg = mean(difflog))  
4 plot(south_year_avg[1:10,]$year, south_year_avg[1:10,]$avg, type = "b", frame = F, pch =  
    20, xlab = "Year", ylab = "Average difference in spending in North")
```



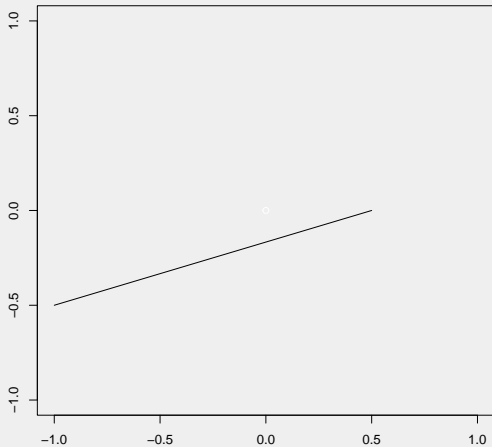
LINE CHART

```
1 lines(south_year_avg[11:20,]$year, south_year_avg[11:20,]$avg, pch = 18, col = "red",  
    type = "b", lty = 2)  
2 # Add a legend to the plot  
3 legend("bottomright", legend=c("North", "South"),  
4       col=c("black", "red"), lty = 1:2, cex=0.8, box.lty=0)
```



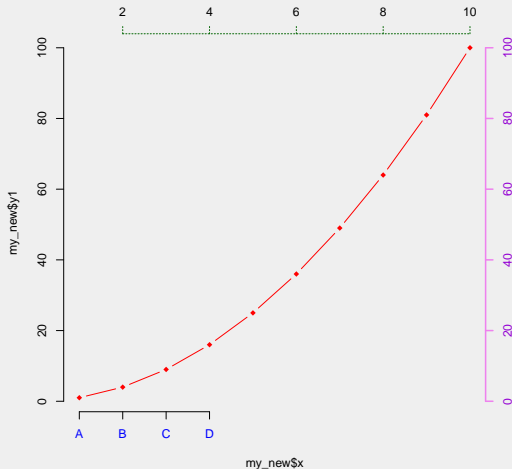
SEGMENTS

```
1 plot(0, 0, col = "white", xlab = "", ylab = "")  
2 segments(x0 = - 1, y0 = - 0.5, x1 = 0.5, y1 = 0)
```



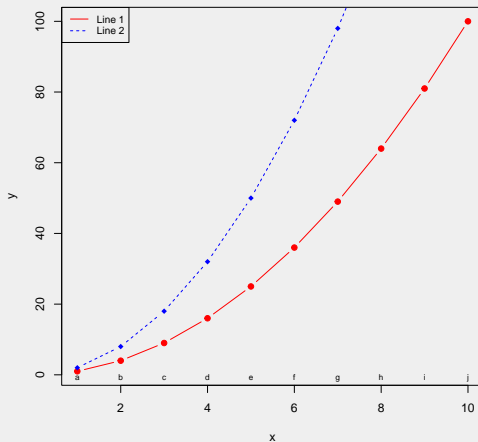
X AND Y AXES

```
1 plot(my_new$x, my_new$y1, pch=18, col="red", type="b",  
2       frame=FALSE, xaxt="n") # Remove x axis  
3 axis(1, 1:4, LETTERS[1:4], col.axis="blue")  
4 axis(3, col = "darkgreen", lty = 2, lwd = 0.5)  
5 axis(4, col = "violet", col.axis = "dark violet", lwd = 2)
```



PLOT LABELS

```
1 text(my_new$x, 2, my_new$x2,  
2      cex=0.65, pos=1, col="black")
```



SAVE PLOTS

There are two ways to think about saving your plots:

- If you're working in RMarkdown, just "knit" your file and your plots will show up as part of your HTML, Word, or PDF document
- If need to save an individual plot for some other purpose (e.g. putting it in a report created in Latex, Powerpoint, Word), use `pdf()` and `dev.off()` functions

```
1 pdf("../graphics/histogram.pdf")
2 hist(inc_local$difflog, main = "", xlab = "Difference in candidate spending")
3 dev.off()
```

TUTORIAL - PLOTTING GROUPS DIFFERENCES

We're interested in whether women promote different policies than men?

- Load in this **dataset** on drinking water facilities and a randomized policy experiment in India, where since the mid-1990s, 13 of village council heads have been randomly reserved for women
- Estimate the effect of the reservation policy on the number of new or repaired drinking water facilities in the villages
 - ▶ Hint: `t.test(y ~ x, alternative="two.sided")`
- Effectively plot this difference

OVERVIEW

This time:

- Plotting in base R

Next time:

- Plotting in ggplot