



ETC3250: Data visualisation

Week 6, class 1

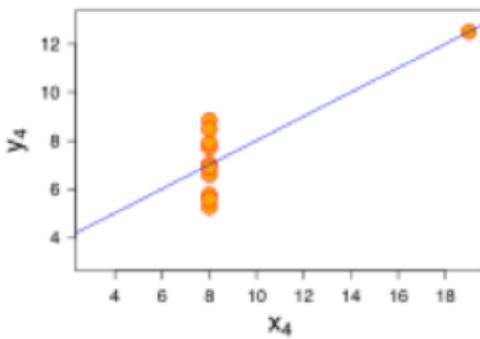
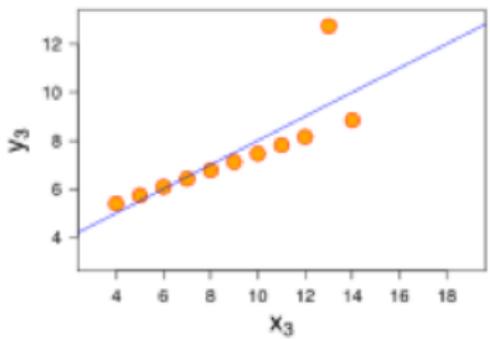
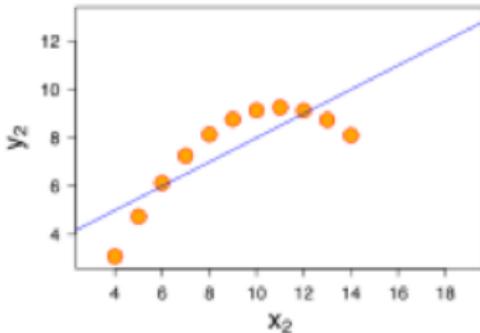
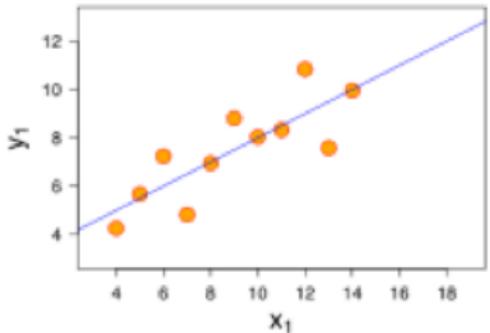
Professor Di Cook, Econometrics and
Business Statistics



Why make a plot of data??

Anscombe's quartet

All of these have the same intercept, slope, correlation.

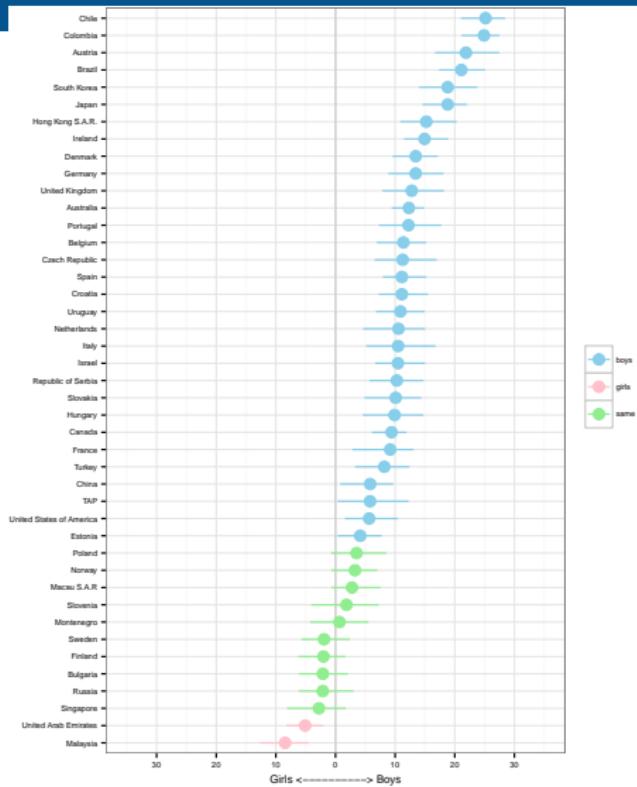


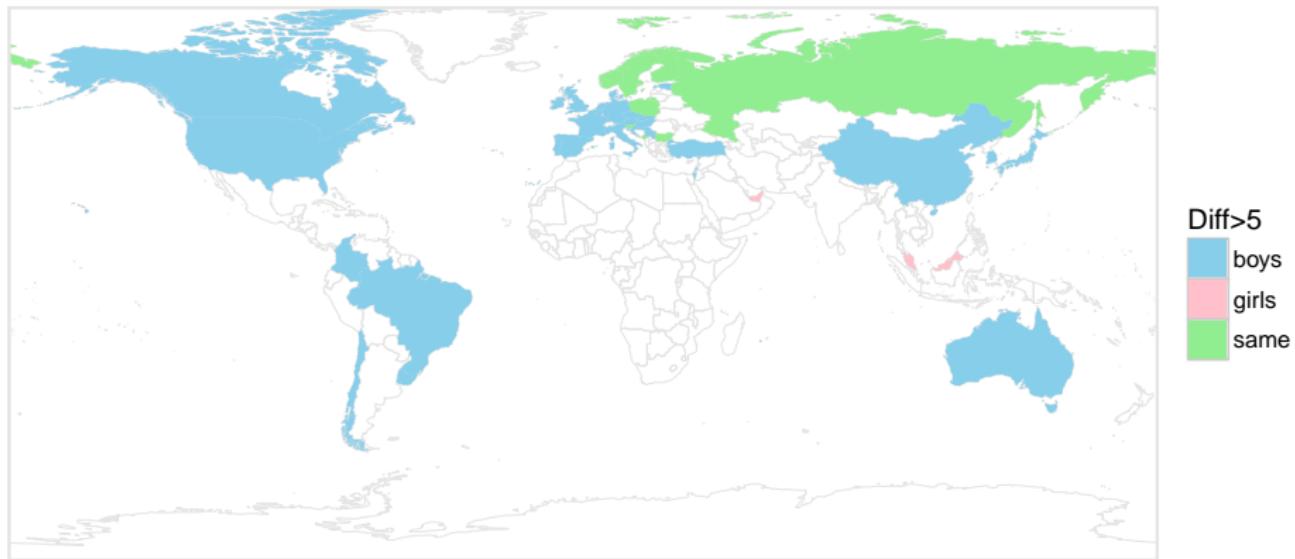
Examples from my own work



- Education
- Climate change
- US Election polls
- Airline traffic patterns
- Wages

- OECD PISA survey “the world's global metric for quality, equity and efficiency in school education”.
- Workforce readiness of 15-year old students
- 500,000 students were tested across 65 countries and 18,000 schools
- Math, reading and science
- Data available from <http://www.oecd.org/pisa>





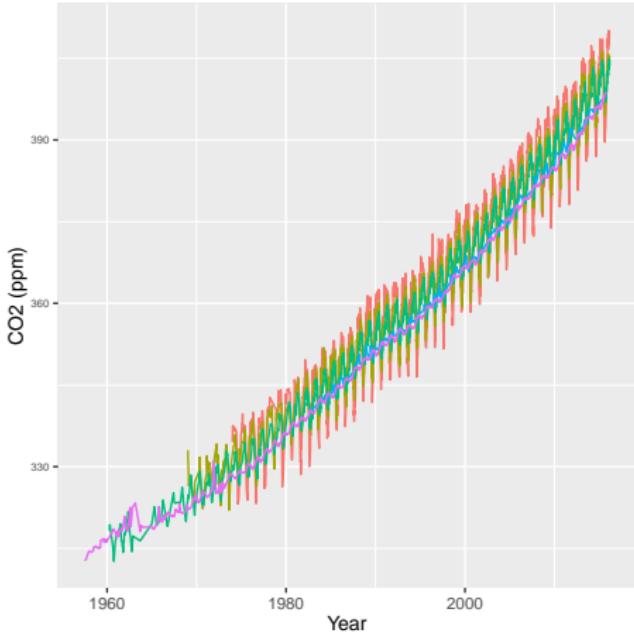
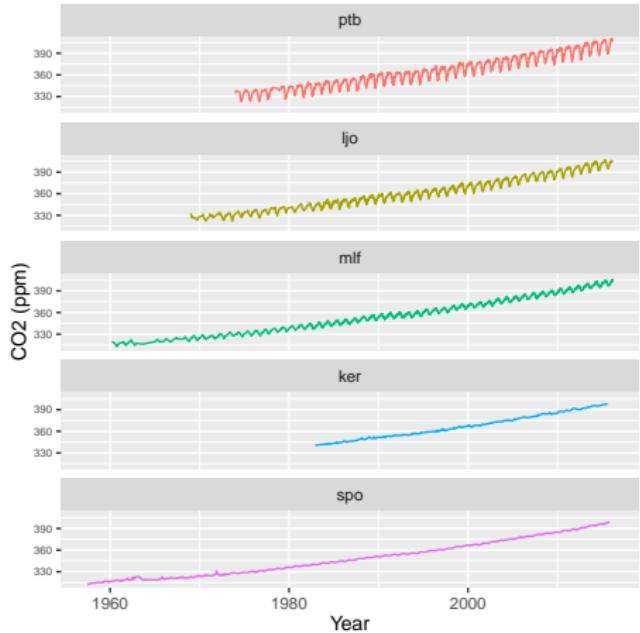
What do we learn?



- Math gender gap is not universal
- Many countries have no substantial difference
- Reverse gap exists in surprising places
- Australia has a 10 point gap (10 points out of 1000 points)
- Individuals show different pattern, highest math score in US is by a girl
- Australia has a huge variation in scores, one of the highest countries, but also one of the lowest countries
- Reading gap is universal in favour of girls

Carbon dioxide data

- Data is collected at a number of locations world wide.
- See Scripps Inst. of Oceanography
- Let's pull the data from the web and take a look . . .
-
- Recordings from South Pole (SPO), Kermadec Islands (KER), Mauna Loa Hawaii (MLF), La Jolla Pier, California (LJO), Point Barrow, Alaska (PTB).



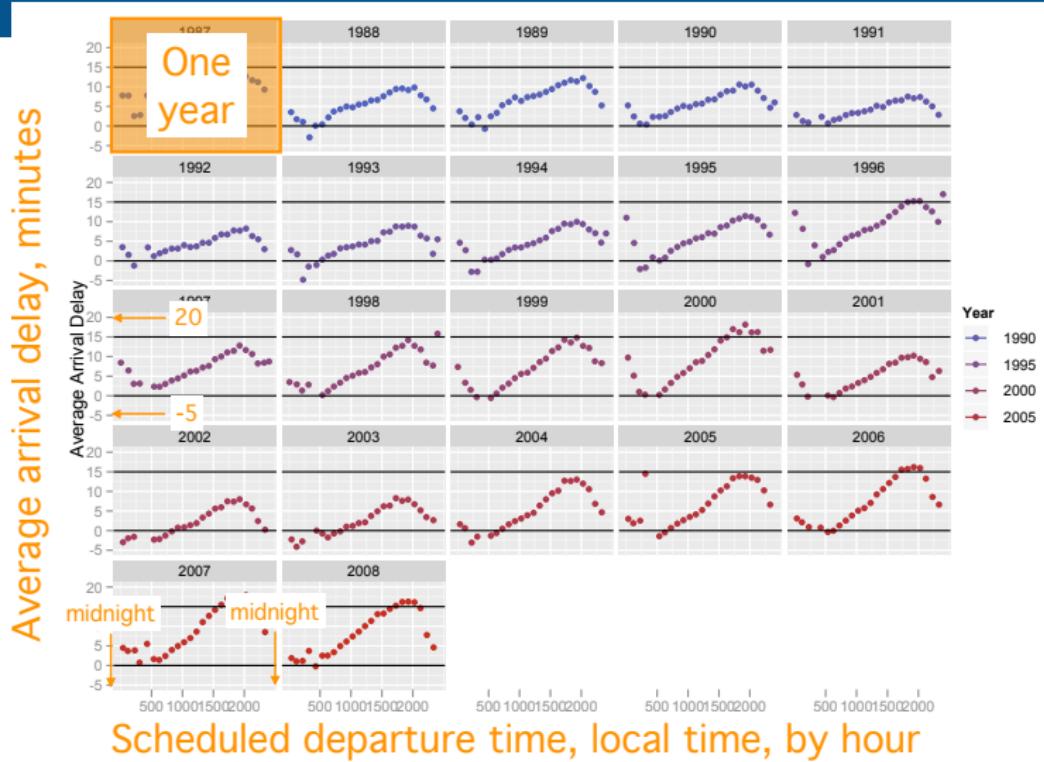


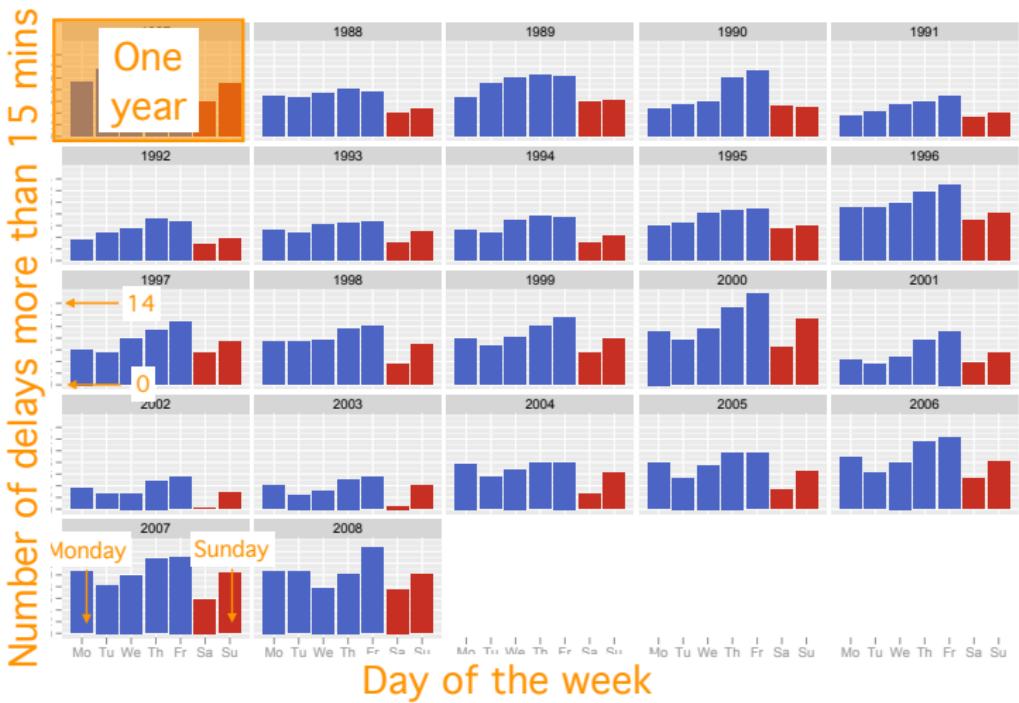
What do we learn?

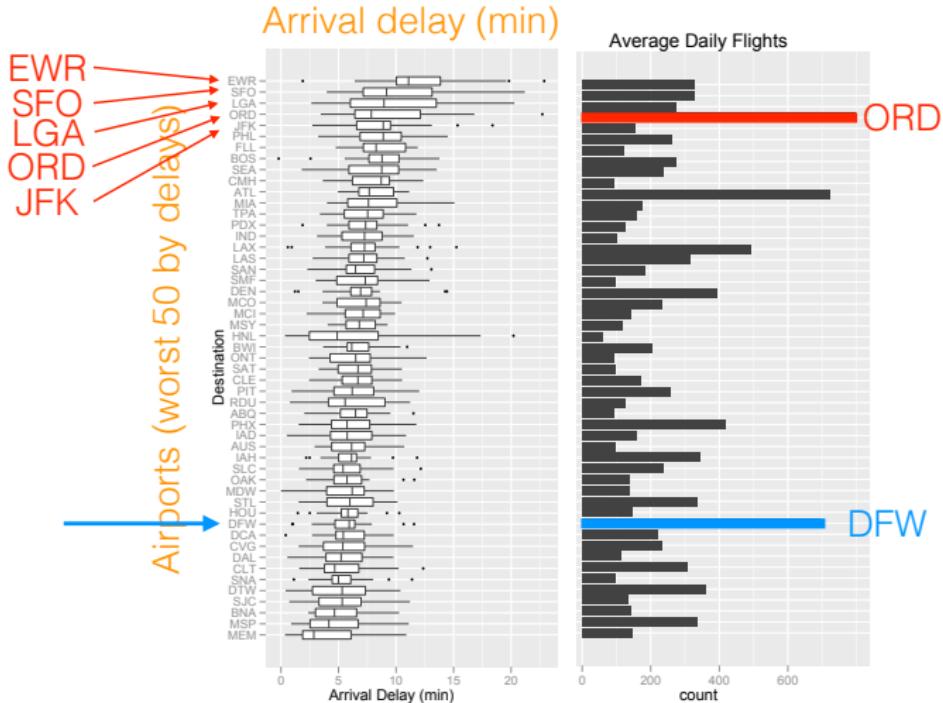


- CO₂ is increasing, and it looks like it is exponential increase. **I really expected that the concentration would have flattened out with all of the efforts to reduce carbon emissions.**
- The same trend is seen at every location - REALLY? Need some physics to understand this.
- Some stations show seasonal pattern - actually the more north the more seasonality - WHY?

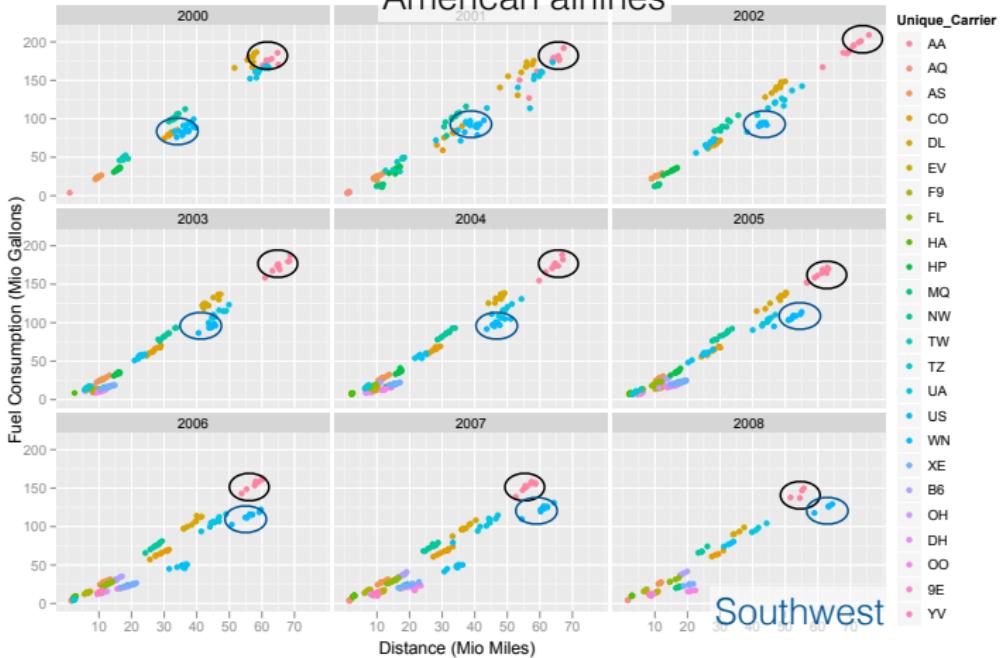
- ~15,000 flights a day
- April 1986 - present (2008)
- RITA - Research and Innovative Technology Administration (flight information, arrival delay, airline, plane id, ...)
- On time performance database - <http://www.transtats.bts.gov/> - yes, you can download this yourself
- Analysis code examples on
<https://github.com/heike/data-technologies>







American airlines



What did we learn?



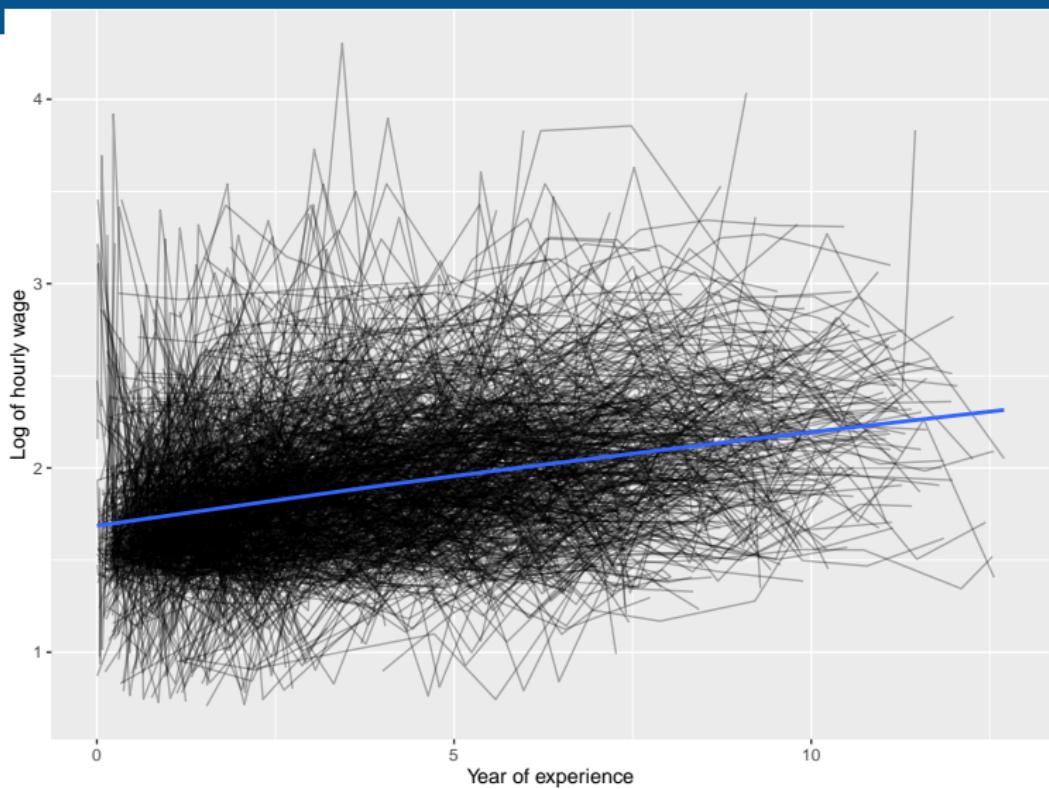
- Fly early in the day, early in the week or weekends (Saturday)
- Avoid ORD, JFK, LGA, EWR
- American Airlines filed for bankruptcy Nov 29, 2013. Mining publicly available data could have sounded the alarms several years in advance

Wages



- 6402 observations on 888 high school dropouts, 1990-2002

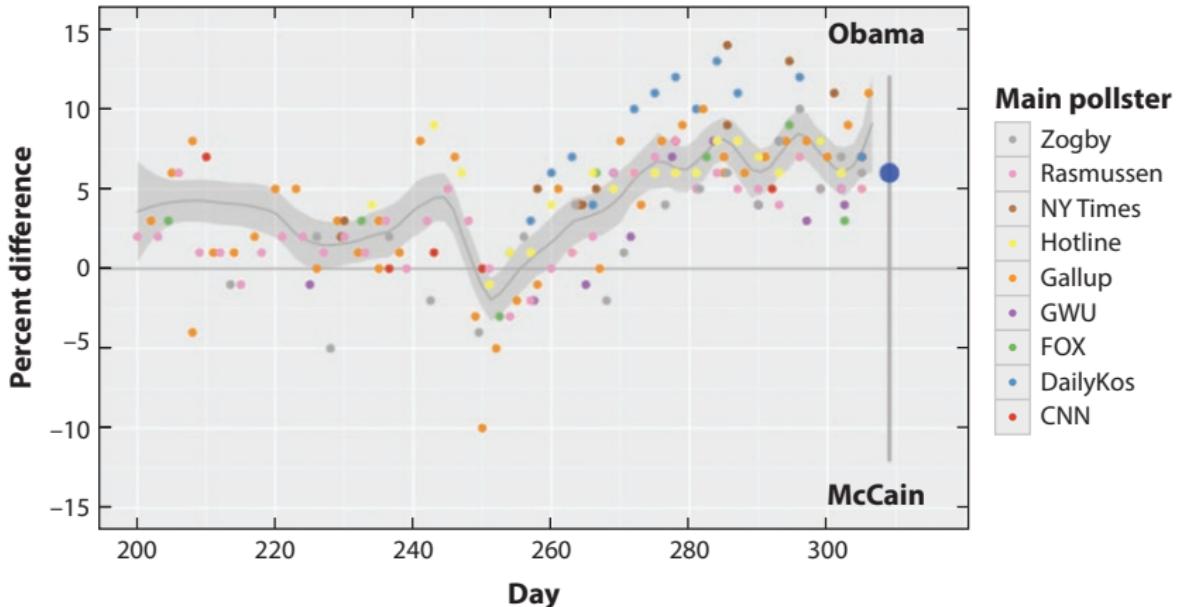
```
# Observations: 6,402
# Variables: 15
# $ id             <fctr> 31, 31, 31, 31, 31, 31, 31, 31, 36, 3
# $ lnw            <dbl> 1.491, 1.433, 1.469, 1.749, 1.931, 1.7
# $ exper           <dbl> 0.015, 0.715, 1.734, 2.773, 3.927, 4.9
# $ ged             <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1
# $ postexp         <dbl> 0.015, 0.715, 1.734, 2.773, 3.927, 4.9
# $ black           <fctr> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0
# $ hispanic        <fctr> 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0
# $ hgc              <int> 8, 8, 8, 8, 8, 8, 8, 8, 9, 9, 9, 9
# $ hgc.9            <int> -1, -1, -1, -1, -1, -1, -1, -1, 0, 0,
# $ uerate           <dbl> 3.21, 3.21, 3.21, 3.30, 2.89, 2.49, 2.
# $ ue.7              <dbl> -3.785, -3.785, -3.785, -3.705, -4.105
# $ ue.centert1      <dbl> 0.000, 0.000, 0.000, 0.080, -0.320, -0.
# $ ue.mean           <dbl> 3.21, 3.21, 3.21, 3.21, 3.21, 3.21, 3.
```



(wages-increasing.mov)

(wages-decreasing.mov)

- They are in the middle of another LONG election
- It's coming to a close
- There is a lot of information about how people might vote
- We looked at how things progressed in 2008 election, in the months leading up to the vote
- We used web scrapers to pull polling data off web sites



Pollsters operating in the US are not all impartial.

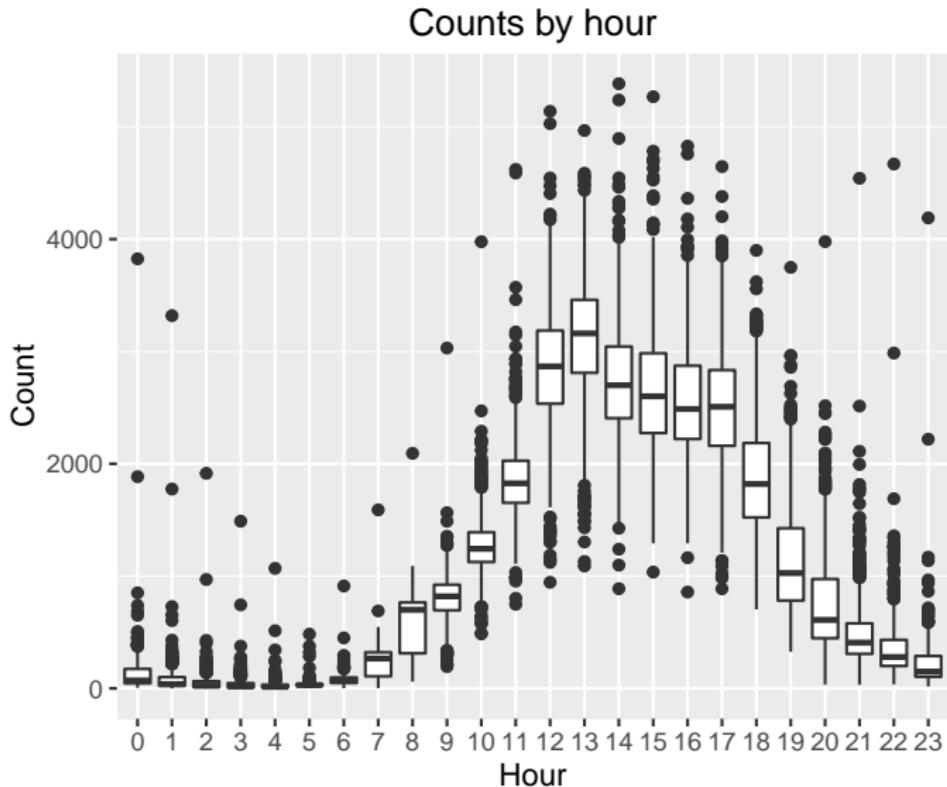
Why make a plot of data??

- What is a (data) plot?
- What are the three most important data plots?
- Is a plot a statistic?

Your turn

M

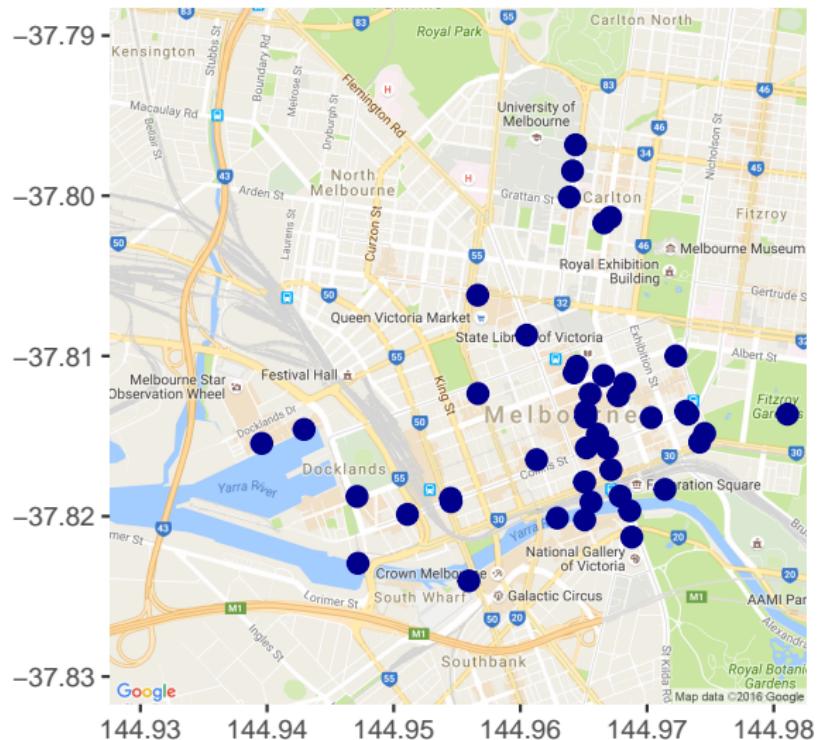
How would you describe this plot?



What about this one?

M

Sensor locations



Using the package ggplot2

Elements of a plot

- data
- aesthetics: mapping of variables to graphical elements
- geom: type of plot structure to use
- transformations: log scale, ...

Additional components

- layers: multiple geoms, multiple data sets, annotation
- facets: show subsets in different plots
- themes: modifying style

Why use a grammar of graphics?

M

Variable in the data is directly mapped to an element in the plot

- Cheat sheet
- ggplot2: Elegant Graphics for Data Analysis, Hadley Wickham, web site
- R Graphics Cookbook, Winston Chang
- Naomi Robbins, Creating More Effective Graphs
- Antony Unwin, Graphical Data Analysis with R

This work is licensed under the Creative Commons Attribution-Noncommercial 3.0 United States License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc/3.0/us/> or send a letter to Creative Commons, 171 Second Street, Suite 300, San Francisco, California, 94105, USA.