

# DS Visualisation and Analysis

Abbey Waldron

# This Week - Making the right Plot!

1. Show last week's work to the class
2. Introduction to making the right plot
3. Group Challenge - guess the question
4. Individual work - 2014 ebola data

# Random Student Generator

# Problem Set From Last Week: Irises

1. Use `help(iris)` to understand what variables there are - make sure you know what they all mean.
2. Make a histogram (`hist`) with 20 bins of petal width for the Iris Setosa.
3. Make a scatterplot (tip: try `ggplot2` `ggplot`) of sepal length vs petal length. Show each of the three species of Iris on the same plot with a coloured legend to separate them.
4. Make a scatterplot of sepal length vs sepal width for all Irises whose petal width is greater than 1.5.
5. Make one more plot that shows something interesting about the inter-species differences of the Irises.

# What does it mean to make the right plot?

# What does it mean to make the right plot?

- ▶ Ask a good question
- ▶ Answer the question in one plot

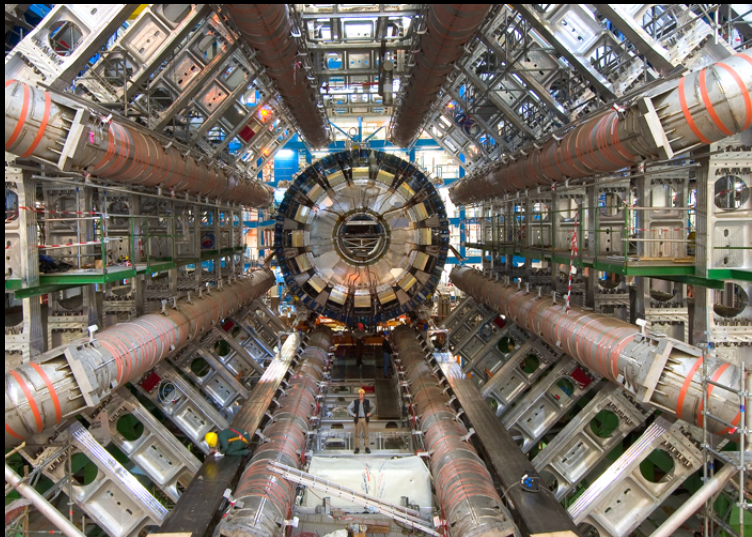
# Practice and THINK

# Calibration (normalisation)

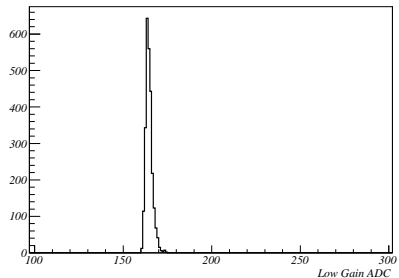
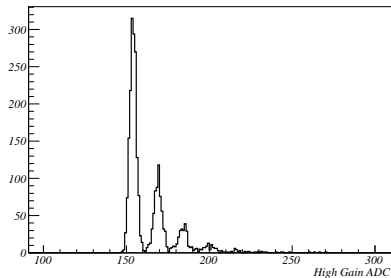
Same input  $\rightarrow$  same output



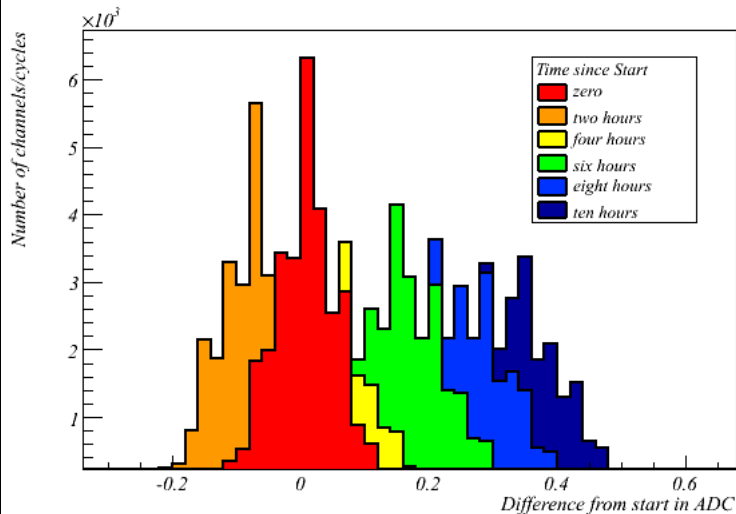
# A little story



# Noise



# What Worked



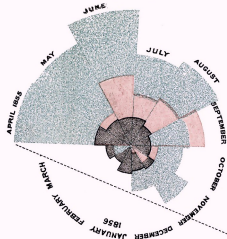
# Group Challenge

Here are some plots: what is the question they answer?

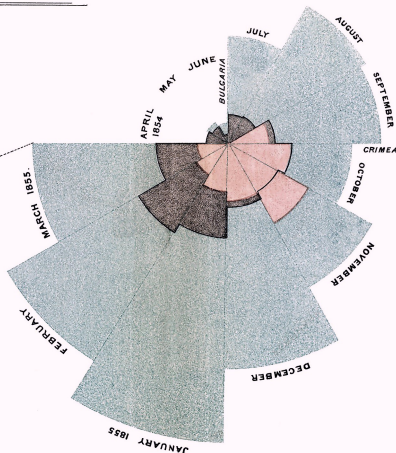
# Florence Nightingale

## DIAGRAM OF THE CAUSES OF MORTALITY IN THE ARMY IN THE EAST.

2.  
APRIL 1855 TO MARCH 1856.



1.  
APRIL 1854 TO MARCH 1855.



*The Areas of the blue, red, & black wedges are each measured from the centre as the common vertex.*

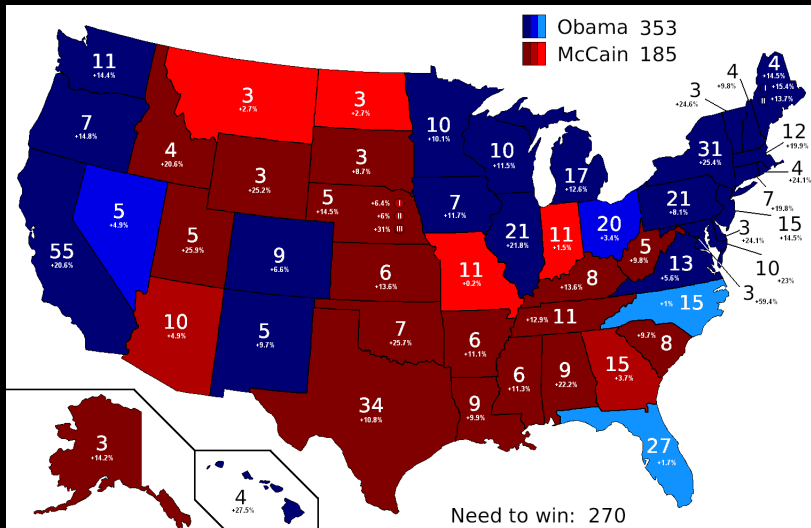
*The blue wedges measured from the centre of the circle represent area for area, the deaths from Preventable or Mitigable Zymotic diseases, the red wedges measured from the centre the deaths from wounds, & the black wedges measured from the centre the deaths from all other causes.*

*The black line across the red triangle in Nov<sup>r</sup> 1854 marks the boundary of the deaths from all other causes during the month.*

*In October 1854, & April 1855, the black area coincides with the red; in January & February 1855, the blue coincides with the black.*

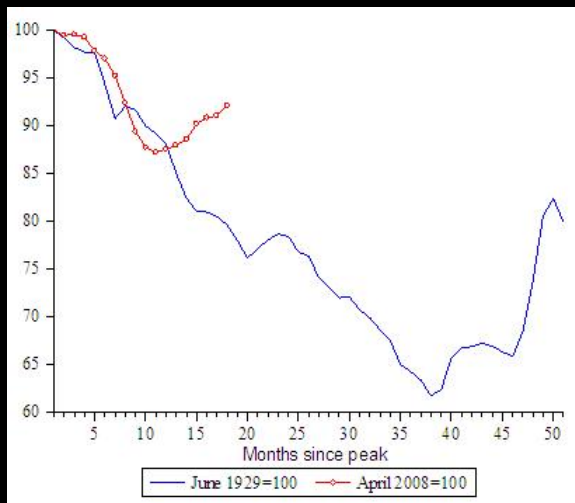
*The entire areas may be compared by following the blue, the red & the black lines enclosing them.*

# US Poll Results pre 2008 election



Wikipedia commons

# Guess



voxeu.org

# This week's problems

We will work with data from the West Africa ebola outbreak of 2013 to present.

Still ongoing, so far 28 000 cases are suspected to have occurred.



# Notes on the data

Data is from the WHO, cite it! The "total" number is only given when data of confirmed, probable and suspected cases or deaths is not known.

Think carefully about which numbers you want to use for the different questions:

- ▶ Confirmed: laboratory tests confirm ebola
- ▶ Probable: health workers see symptoms consistent with ebola
- ▶ Suspected: other cases that may be ebola

# R Code Hints

Some things you may need:

- ▶ `read.csv()`
- ▶ `lapply`, `tapply` etc.

Check out the swirl stats tutorial if you don't know the last ones.

# Backup Slides