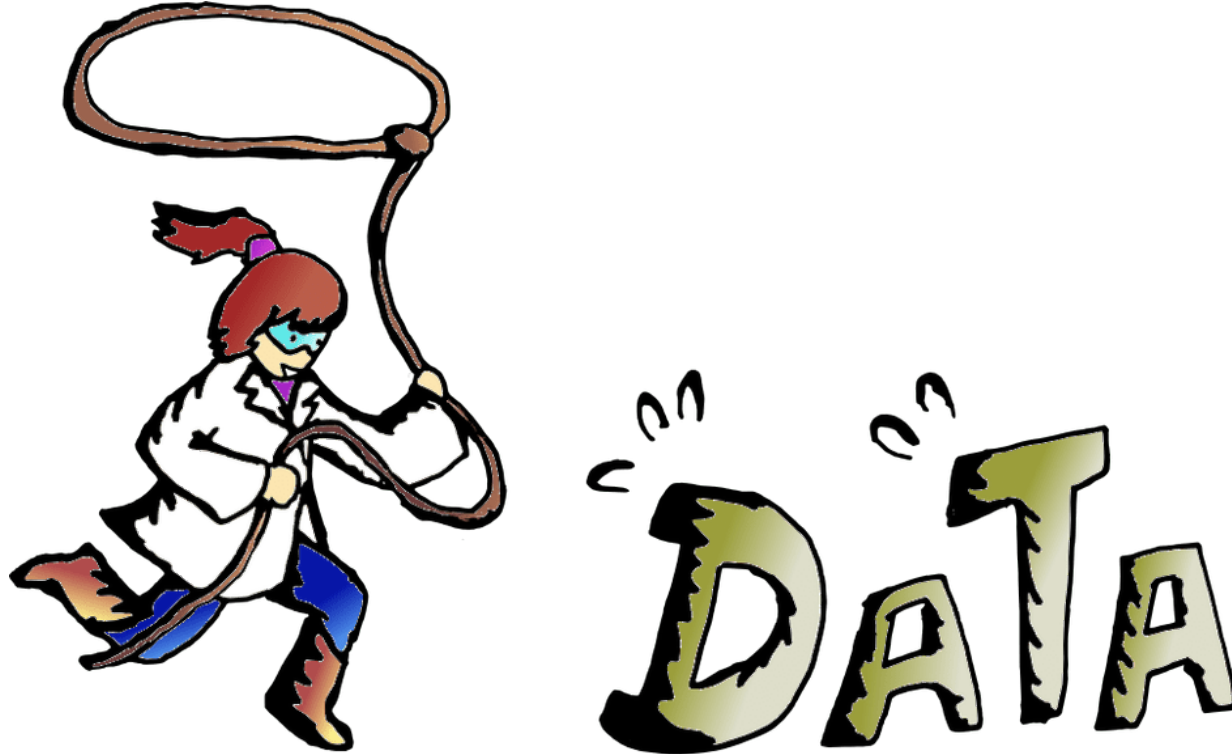# Data wrangling/manipulation

# Overview

Data wrangling/manipulation with dplyr

Brief history of data visualization

# Announcements

A practice midterm exam and slides with the answers will be posted by next class

- Exam format: multiple choice, short essays, short coding
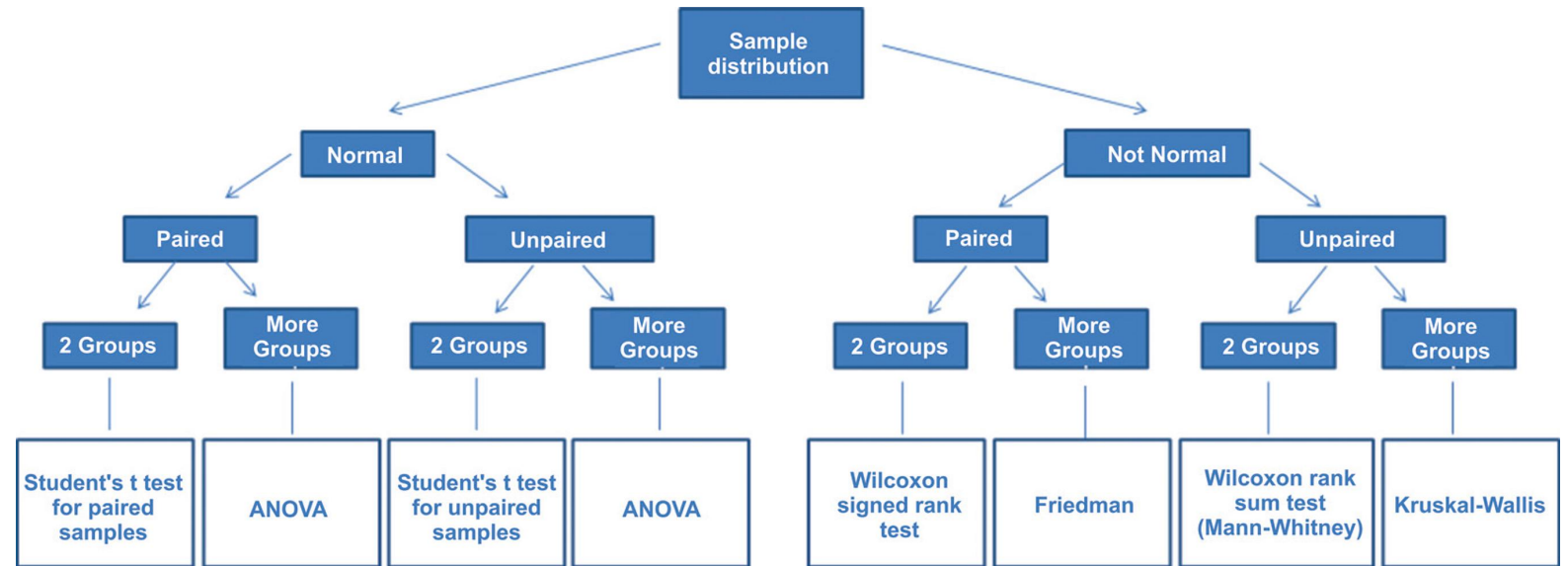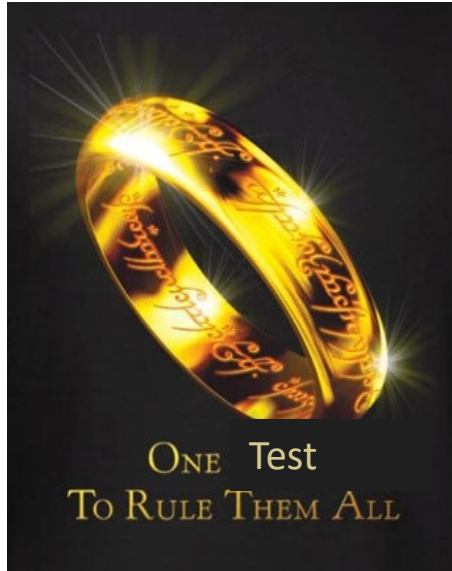
Get started on homework 5 early!!!

- I strongly recommend you do the dplyr exercises prior to next class

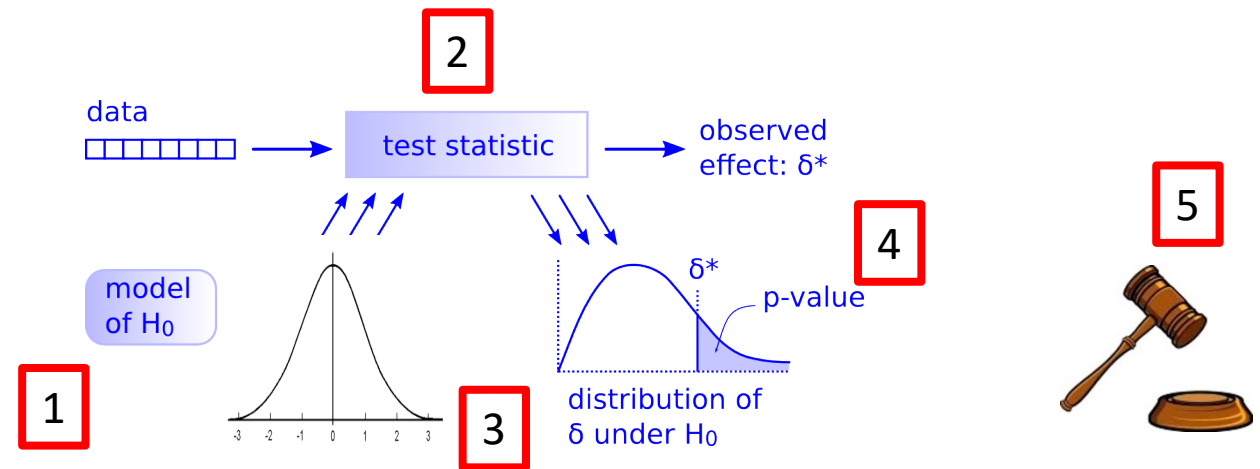Any other questions about class logistics?

# Plan for the semester

1   Sep 2     Course overview, introduction to R, descriptive statistics

2   Sep 7-9     Review of central statistical concepts and exploratory analysis using R

3   Sep 14-16     Confidence Intervals and the bootstrap

4   Sep 21-23     Review of hypothesis tests and permutation tests in R

5   Sep 28-30     Parametric, non-parametric and theories of hypothesis testing

6   Oct 5-7     Data manipulation and visualization

7   Oct 12-14     Mapping, review and midterm exam
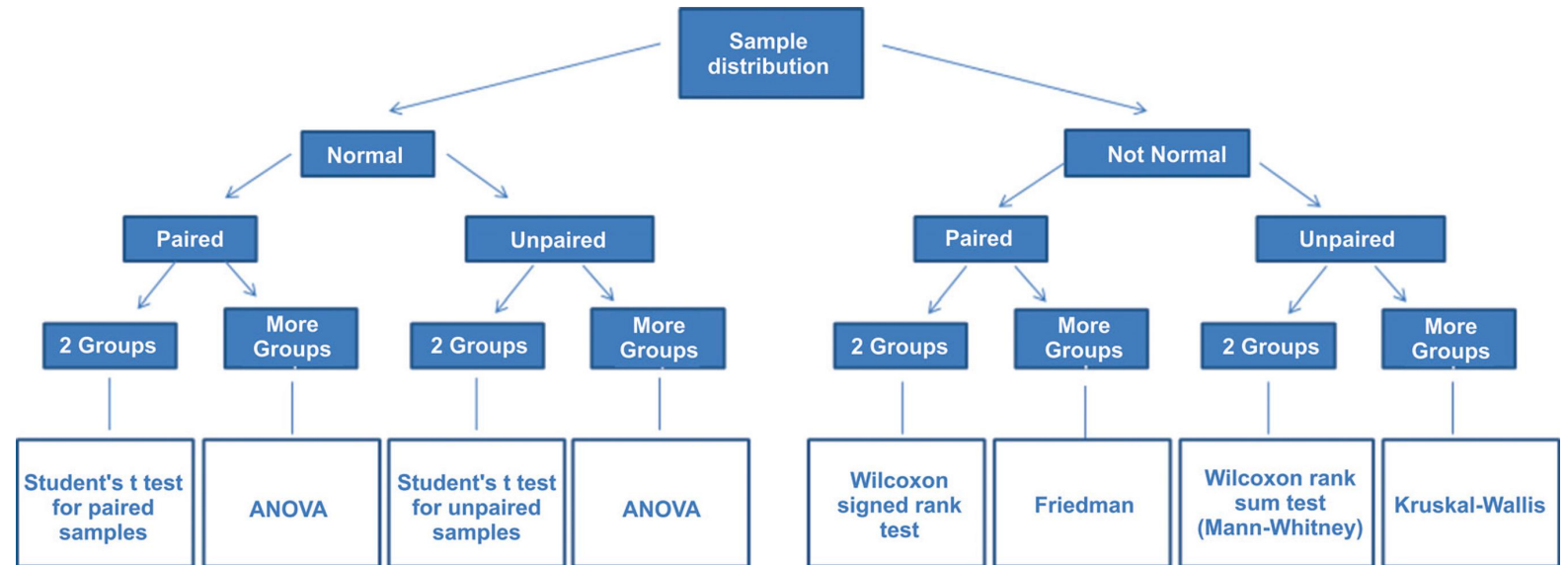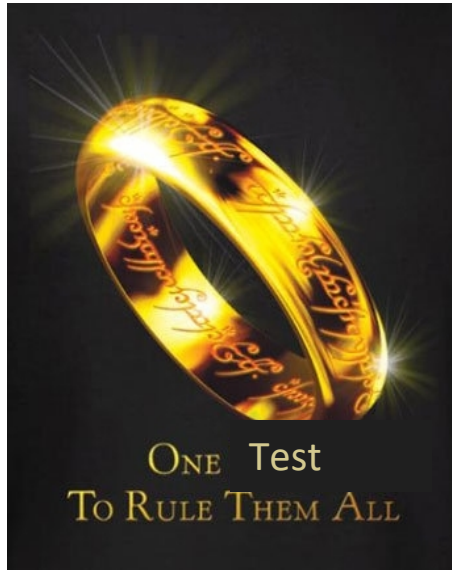
8   Oct 22     October break

base R

resampling methods

data wrangling visualization

# Very quick review



Test

ONE Test TO RULE THEM ALL

Just need to follow 5 steps!

Sample distribution

Normal → Paired → 2 Groups → Student's t test for paired samples

Normal → Paired → More Groups → ANOVA

Normal → Unpaired → 2 Groups → Student's t test for unpaired samples

Normal → Unpaired → More Groups → ANOVA

Not Normal → Paired → 2 Groups → Wilcoxon signed rank test

Not Normal → Paired → More Groups → Friedman

Not Normal → Unpaired → 2 Groups → Wilcoxon rank sum test (Mann-Whitney)

Not Normal → Unpaired → More Groups → Kruskal-Wallis

data → test statistic → observed effect: $\delta^*$

model of $H_0$

distribution of $\delta$ under $H_0$

$\delta^*$

p-value

1

2

3

4

5

# Very quick review



One Test To Rule Them All



To select the appropriate parametric test, focus on the parameters being tested in the null hypothesis
- E.g., $H_0: \pi = 0.5$ $\quad$ $H_0: \mu = 0.5$ $\quad$ $H_0: \mu_T = \mu_C$ $\quad$ $H_0: \mu_1 = \mu_2 = \ldots = \mu_k$

Parametric tests are derived from particular mathematical assumptions
- E.g., data from the two samples comes from normal populations with the same variance
- Some hypothesis tests are "robust" to violations of these assumptions
  - The robustness can be evaluated this through computer simulations

# Very quick review: theories of hypothesis testing
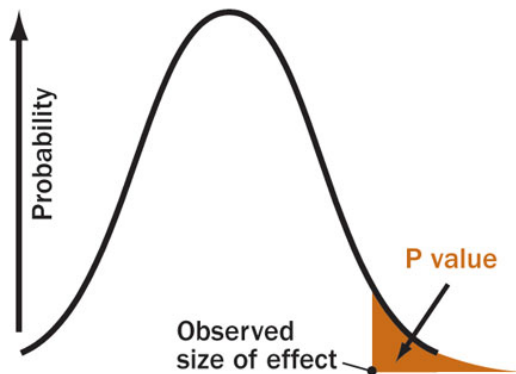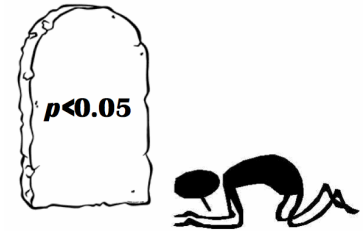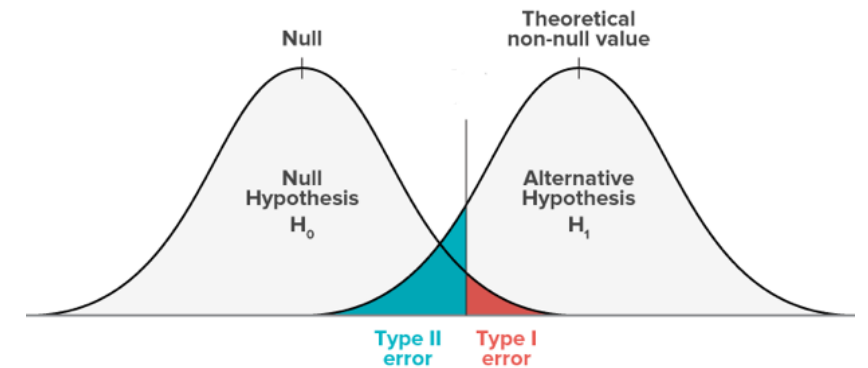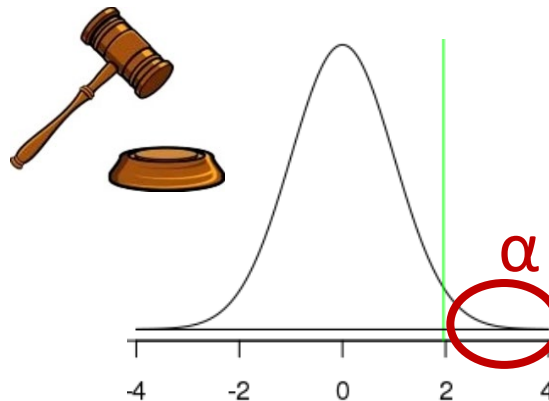


Fisher (1890-1962)

Neyman (1894-1981)    Pearson (1895-1980)

p<0.05

p-value a strength of evidence

Use p-value to make a decision

# Questions?

# The tidyverse and dplyr

# The 'tidyverse'

The tidyverse is set of R packages that operate 'tidy data'
- i.e., that operate on data frames    (or tibbles)

Tidy data is data where:
- Each variable must have its own column
- Each observation must have its own row
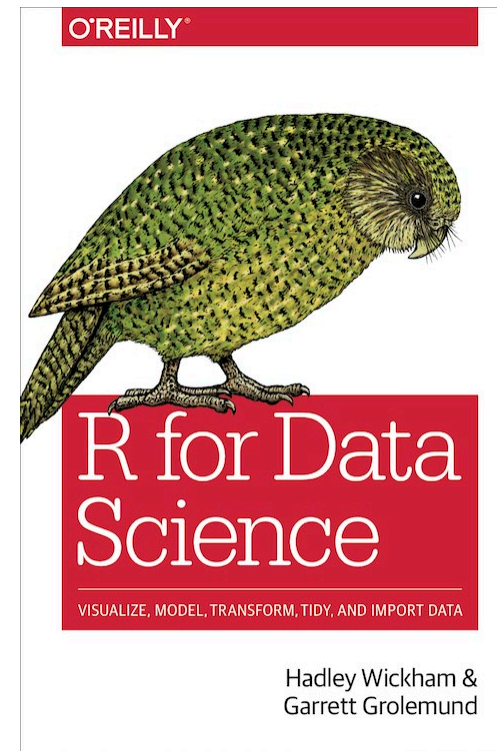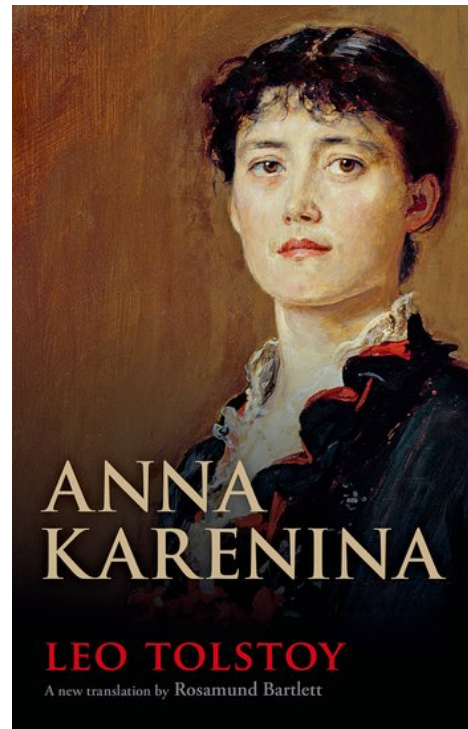- Each value must have its own cell

# Messy data…

What would be an example of data that is not tidy?

| Curve information - Curve quality data | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Name | Formula | Slope at | Intercept | ED-20 | ED-50 | ED-80 | Correlati | Forced through origo | | | | | | | | |
| Standard | Calc 1: C | standard | standard | 3792394 | 27752 | 0.2 | 0.5 | 0.8 | 1 | No | | | | | | |
| | | | | | | | | | | | | | | | | |
| Plate information | | | | | | | | | | | | | | | | |
| Plate | Repeat | Barcode | Measured | Chamber | Chamber | Humidity | Humidity | Ambient | Ambient | Formula | Measurement date | | | | | |
| 1 | 1 | | N/A | N/A | N/A | N/A | N/A | N/A | N/A | Calc 1: C | standard | standard | 10.12.2013 10:23:33 | | | |
| | | | | | | | | | | | | | | | | |
| Background information | | | | | | | | | | | | | | | | |
| Plate | Label | Result | Signal | Flashes/1 | Meastime | MeasInfo | | | | | | | | | | |
| 1 | PicoGree | 0 | 110307 | 10 | 0 | De=1st Ex=Top Em=Top Wdw=N/A | | | | | | | | | | |
| | | | | | | | | | | | | | | | | |
| Calculate | standard | standards on each plate) where Label: PicoGreenFilterTop(1) channel 1 | | | | | | | | | | | | | | |
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | | | | |
| A | -0.0011 | -0.0011 | -0.001 | -0.001 | -0.0011 | -0.0012 | -0.0011 | -0.0011 | -0.0012 | -0.0012 | 0.9973 | 1.0026 | | | | |
| B | 0.0012 | 0.0014 | 0.0013 | 0.0012 | 0.0013 | 0.0012 | 0.0014 | 0.0003 | -0.0011 | -0.0011 | 0.0981 | 0.103 | | | | |
| C | 0.0016 | 0.0013 | 0.0013 | 0.0011 | 0.0012 | 0.0015 | 0.0016 | -0.0004 | -0.0011 | -0.0011 | 0.0104 | 0.0095 | | | | |
| D | 0.0019 | 0.0024 | 0.0018 | 0.0015 | -0.001 | -0.001 | -0.001 | -0.001 | -0.0011 | -0.0011 | 0.0008 | 0.0009 | | | | |
| E | -0.001 | -0.0011 | -0.0011 | -0.0011 | -0.001 | -0.0012 | -0.0011 | -0.001 | -0.0009 | -0.0011 | -0.0001 | -0.0002 | | | | |
| F | -0.001 | -0.0011 | -0.001 | -0.001 | -0.0012 | -0.0011 | -0.0011 | -0.0009 | -0.001 | -0.001 | -0.0003 | -0.0002 | | | | |
| G | -0.0011 | -0.0011 | -0.0011 | -0.001 | -0.001 | -0.0012 | -0.0011 | -0.001 | -0.001 | -0.0011 | -0.0002 | 0.0012 | | | | |
| H | -0.0011 | -0.0012 | -0.0011 | -0.001 | -0.0011 | -0.0011 | -0.0012 | -0.0011 | -0.0011 | -0.001 | -0.0003 | -0.0003 | | | | |

# Messy data…

"Happy families are all alike; every unhappy family is unhappy in its own way."
— Leo Tolstoy





"Tidy datasets are all alike, but every messy dataset is messy in its own way." –
– Hadley Wickham

# Messy data...

Messy data can be difficult to deal with

# The 'tidyverse'

The packages share a common design philosophy
- Most written by Hadley Wickham

# dplyr:  A grammar for data wrangling

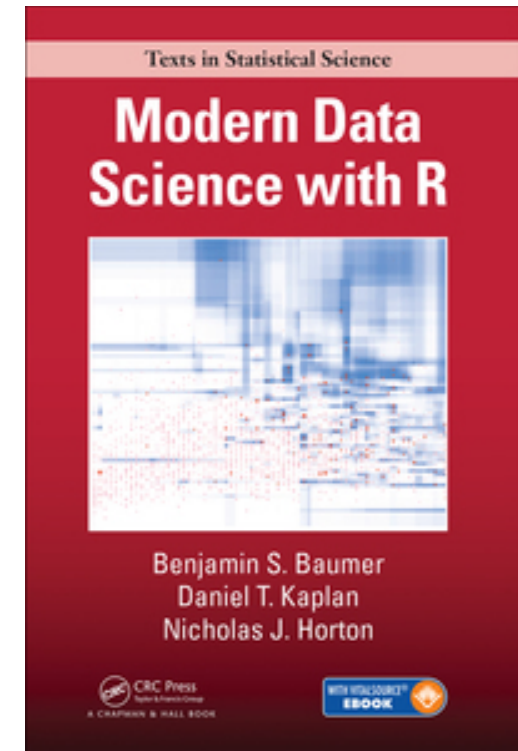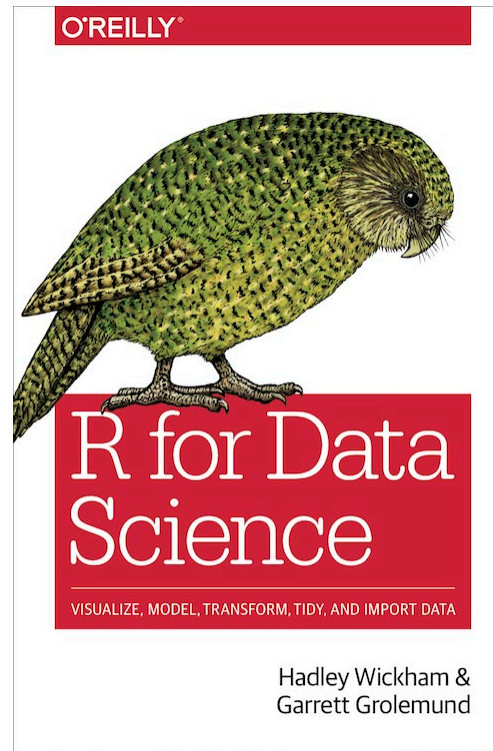**Grammar**:  a set of components that can be combined to achieve a goal

**dplyr** is a package that has a set of verbs that are useful for transformations data:

1. filter()
2. select()
3. mutate()
4. arrange()
5. summarize()
6. group_by()

All these function **take a data frame** and other arguments and **return a data frame**

```
> library(dplyr)   # load the dplyr package
```

# Quick overview of the dplyr functions

# 1. filter()

The filter() function allows you to select a subset of rows in data frame

# 2. select()

The select() function allows you to select a subset of columns

# 3. mutate()

The mutate() function allows you to create new columns that are functions of existing columns

# 4. arrange()

The arrange() function arranges the rows based values in a column
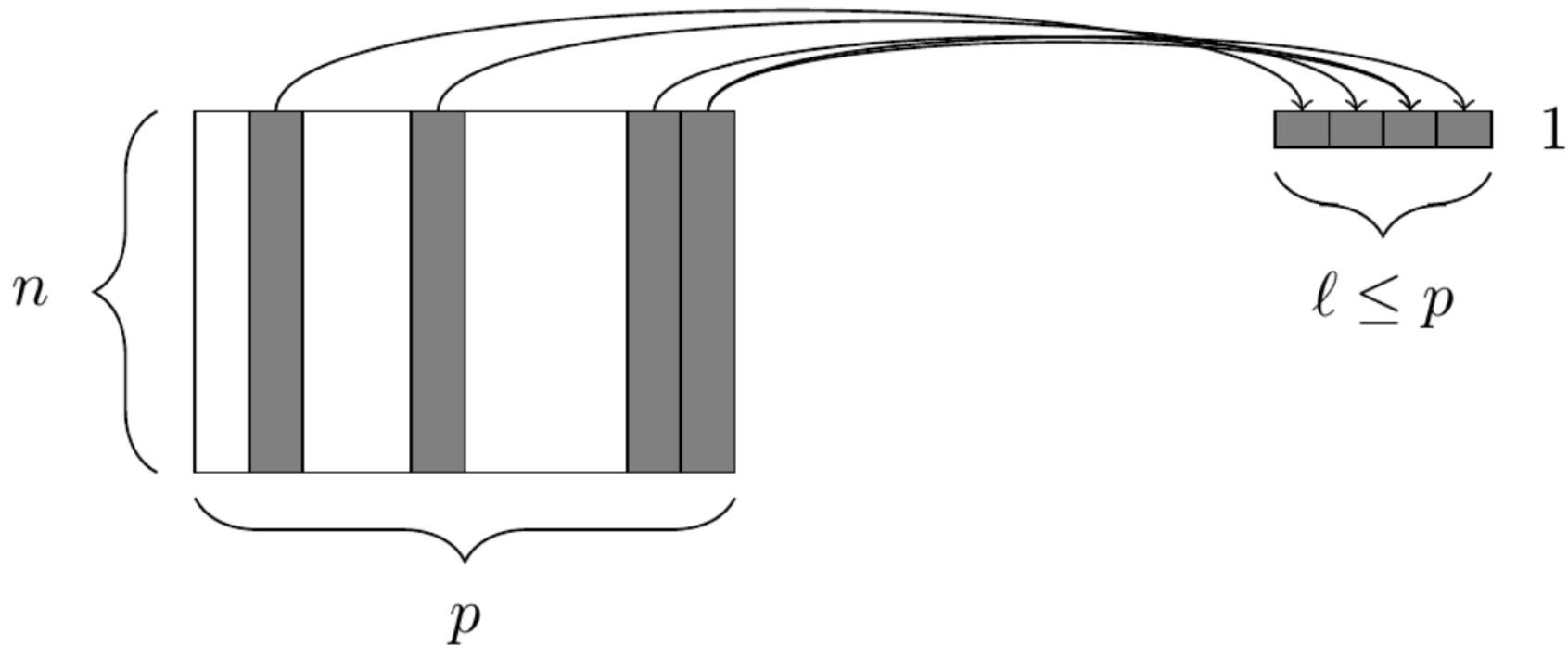  - arrange(desc())  arranges from largest to smallest

# 5. summarize()

The summarize() function reduces values in many rows into single values

# 6. The group_by() function

The group_by() function groups variables for future operations

# The pipe operator

The pipe operator %>% allows us to chain commands together

# Let's try it out!

# A very brief history of data visualization

**The Golden Age of Statistical Graphics**

**Michael Friendly**

# Data visualization

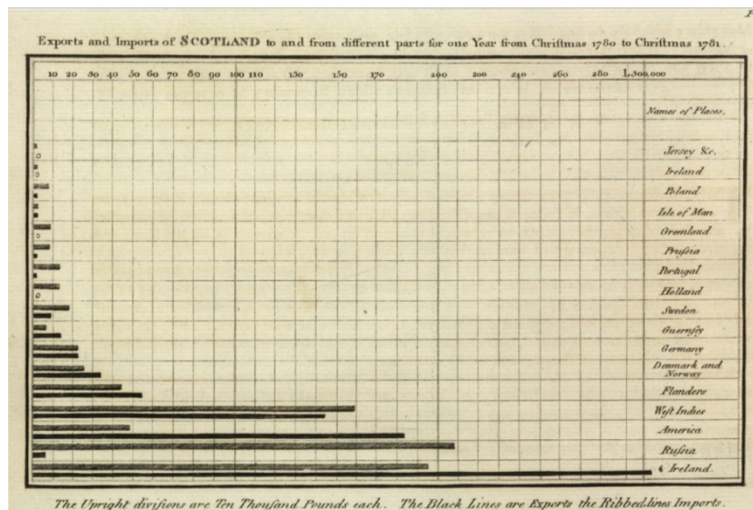What are some reasons we visualize data rather than just reporting statistics?

*Whatever relates to extent and quantity may be represented by geometrical figures. Statistical projections which speak to the senses without fatiguing the mind, possess the advantage of fixing the attention on a great number of important facts.*

*—Alexander von Humboldt, 1811*

# A very brief history of data visualization

The age of modern statistical graphs began around the beginning of the 19th century

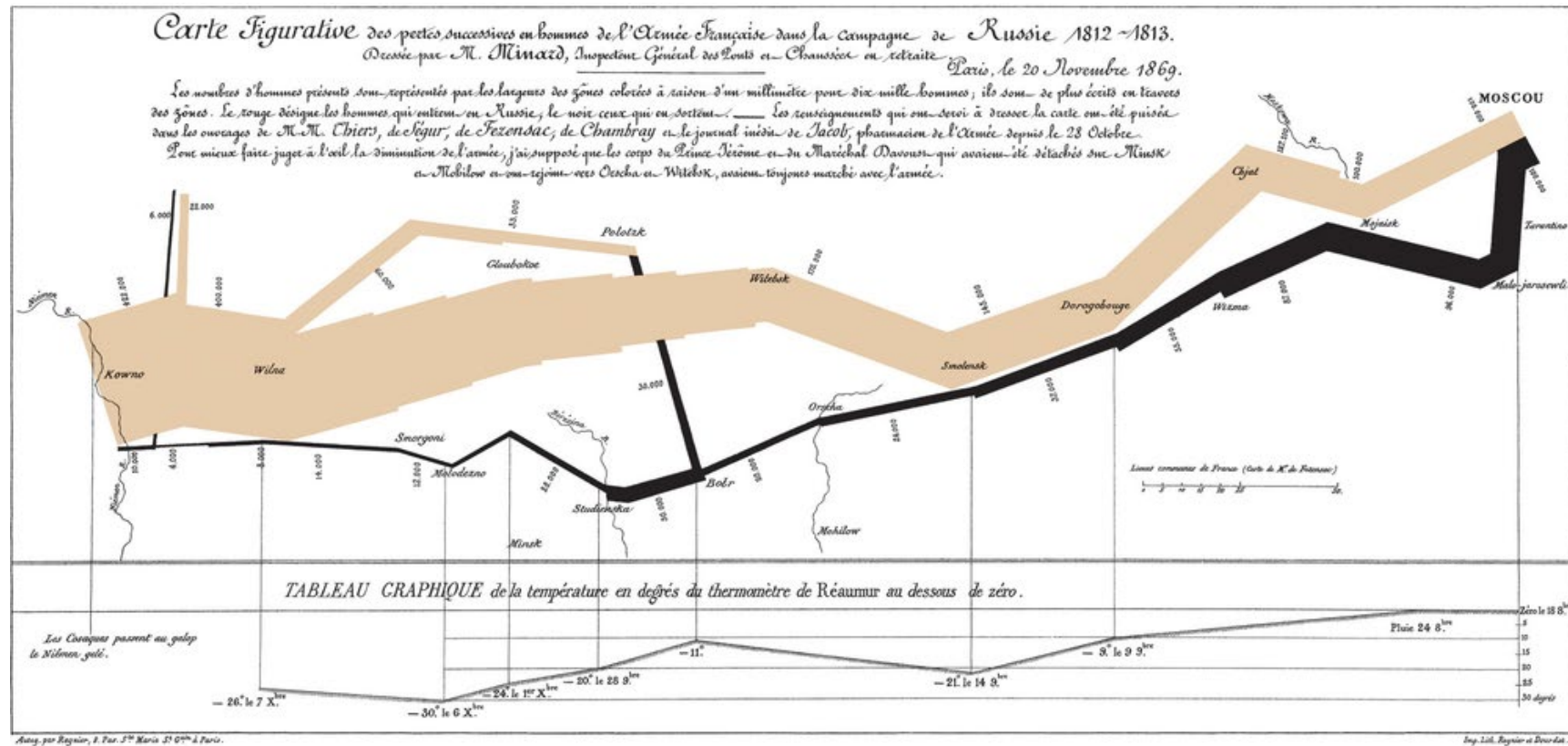[William Playfair](#) (1759-1823) credited with inventing the line graph, bar chart and pie chart



Exports and Imports of SCOTLAND to and from different parts for one Year from Christmas 1780 to Christmas 1781

The Upright divisions are Ten Thousand Pounds each. The Black Lines are Exports the Ribbed lines Imports.

# A very brief history of data visualization

According to Friendly, statistical graphics researched its golden age between 1850-1900

# A very brief history of data visualization

Joseph Minard (1781-1870)



Map of Napoleon's march on Russia

# A very brief history of data visualization

John Snow (1813-1858)



Clusters of cholera cases in London epidemic of 1854

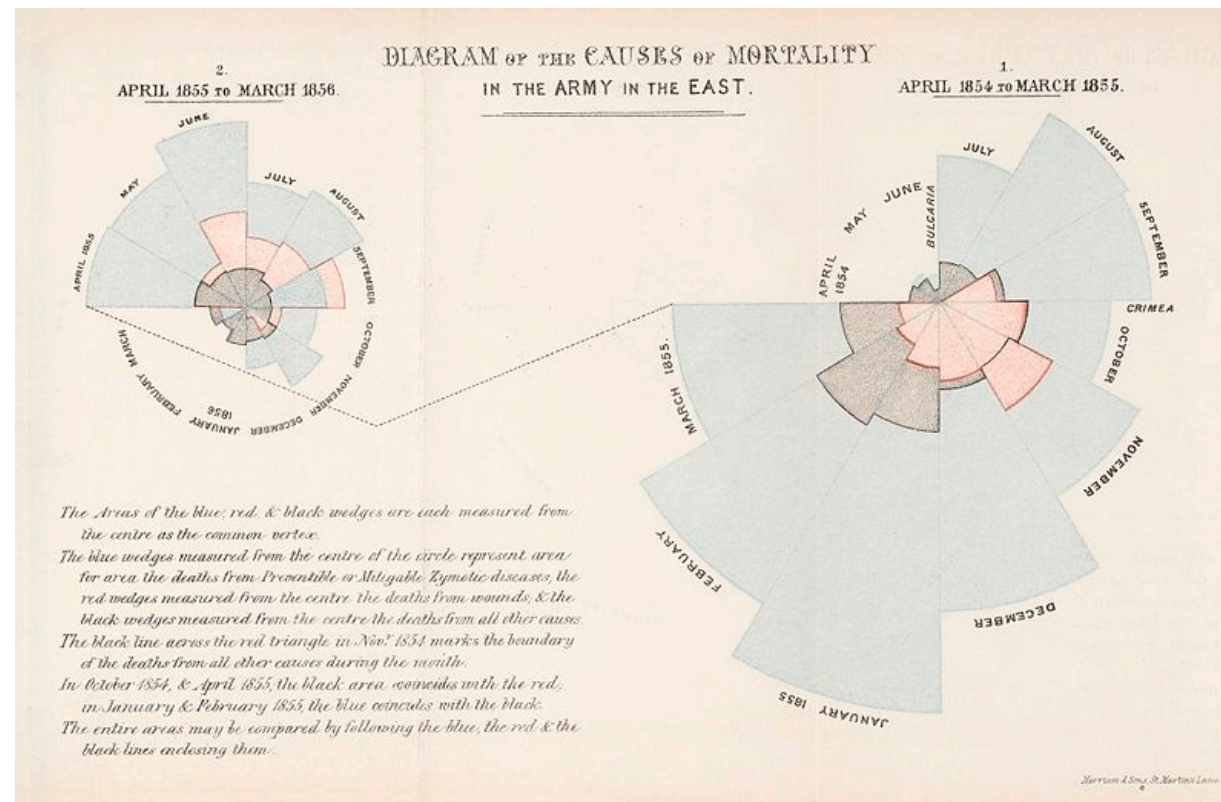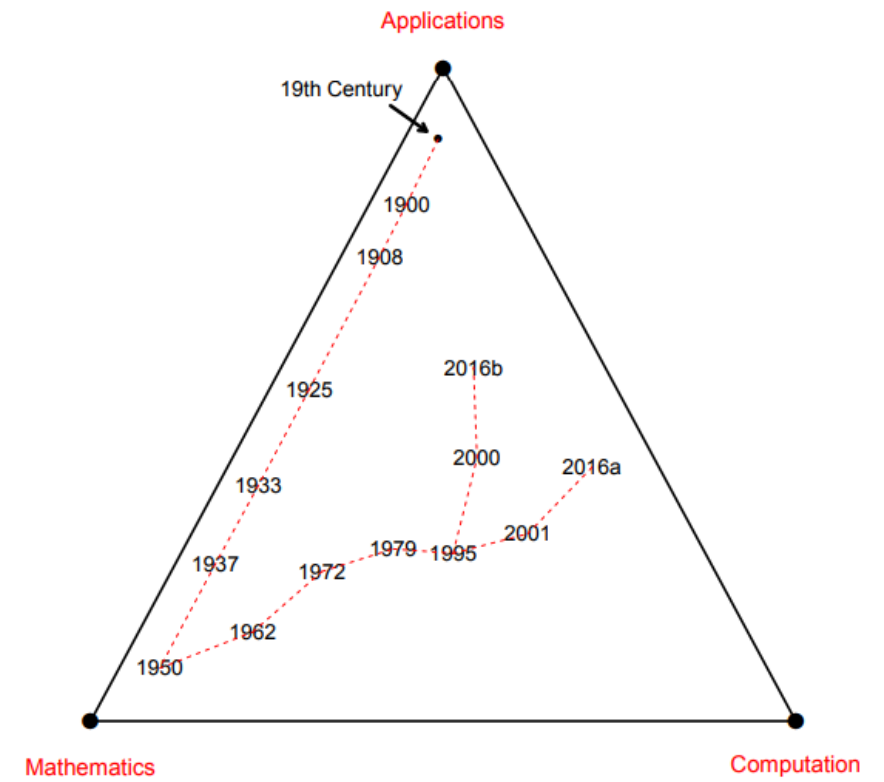# A very brief history of data visualization
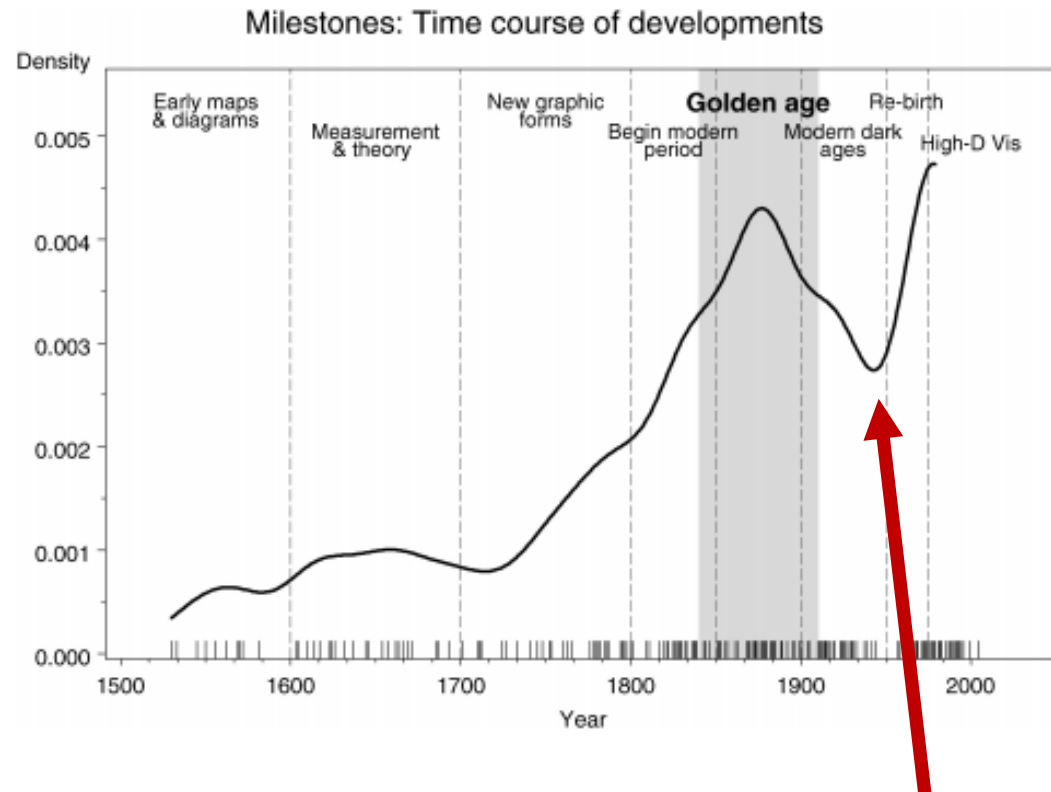
Florence Nightingale (1820-1910)



Diagram of the causes of mortality in the army in the east

# A very brief history of data visualization

[Francis Galton](#) (1822-1911)



WEATHER CHART, MARCH 31, 1875.

The dotted lines indicate the gradations of barometric pressure. The variations of the temperature are marked by figures, the state of the sea and sky by descriptive words, and the direction of the wind by arrows—barbed and feathered according to its force. ⊙ denotes calm.

First weather map published in a newspaper (1875)

# A very brief history of data visualization

"Graphical dark ages" around 1950



Computer Age Statistical Inference, Efron and Hastie

# A very brief history of data visualization
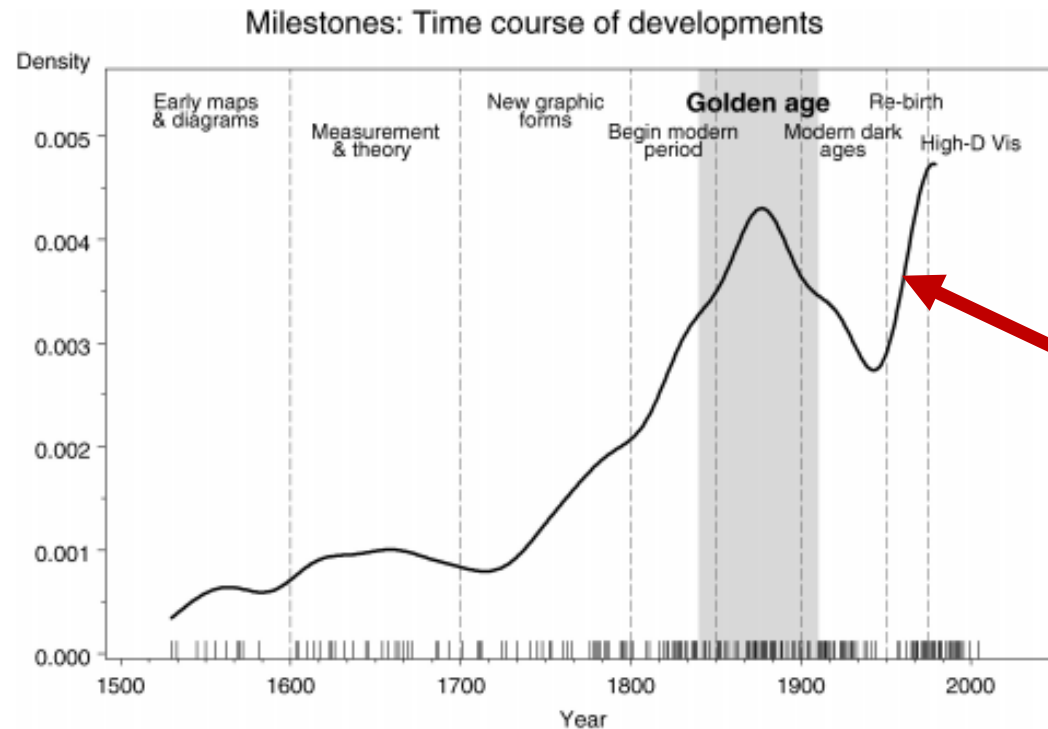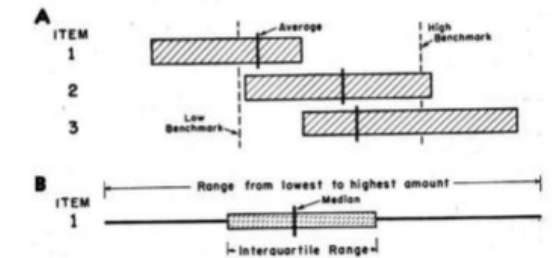
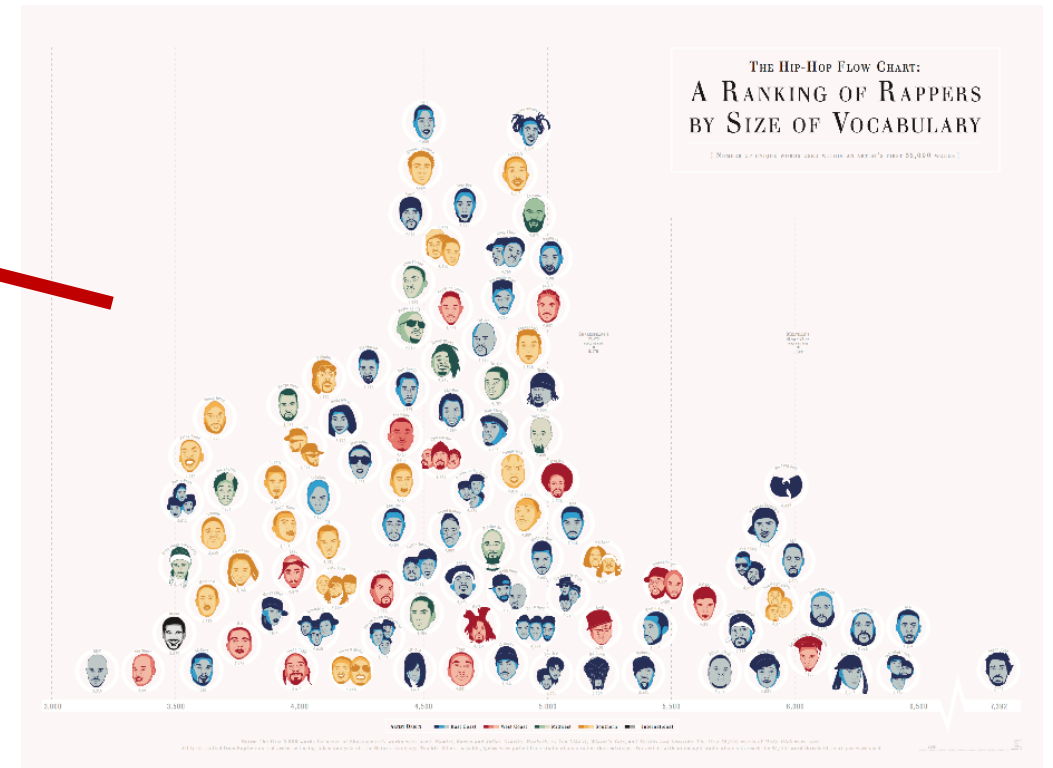Currently undergoing a "Graphical re-birth"

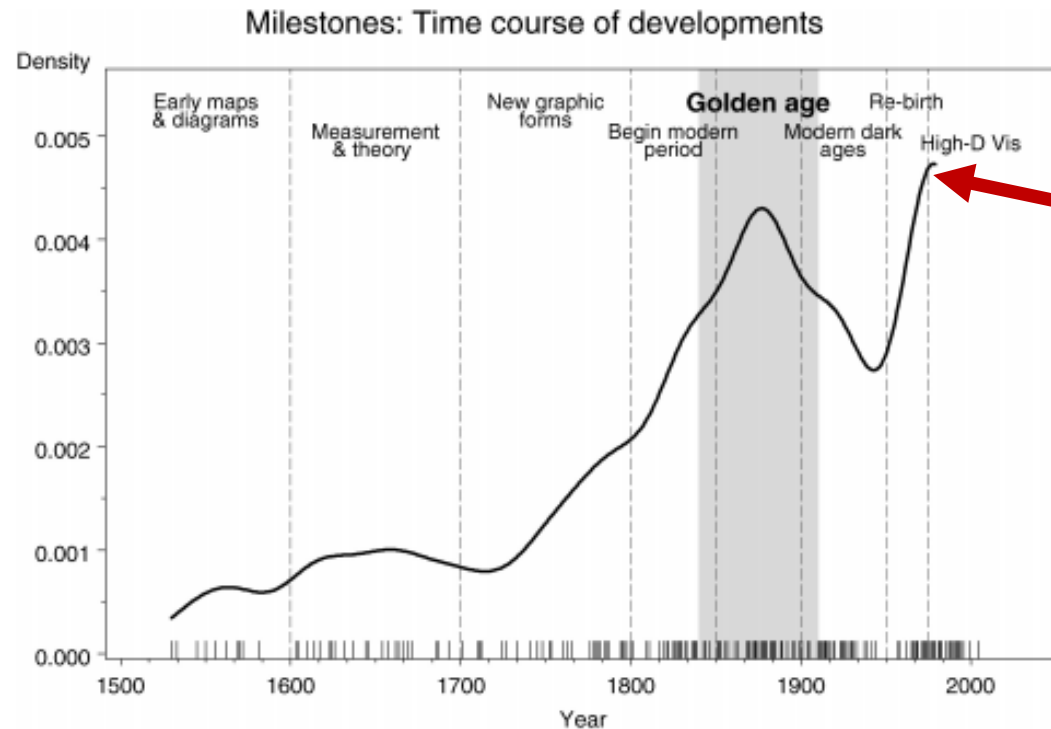Box plot



Spear 1952, Tukey 1970

# A very brief history of data visualization

Currently undergoing a "Graphical re-birth"

# A very brief history of data visualization

Currently undergoing a "Graphical re-birth"



Milestones: Time course of developments

Hans Rosling's gapminder
- Simple version
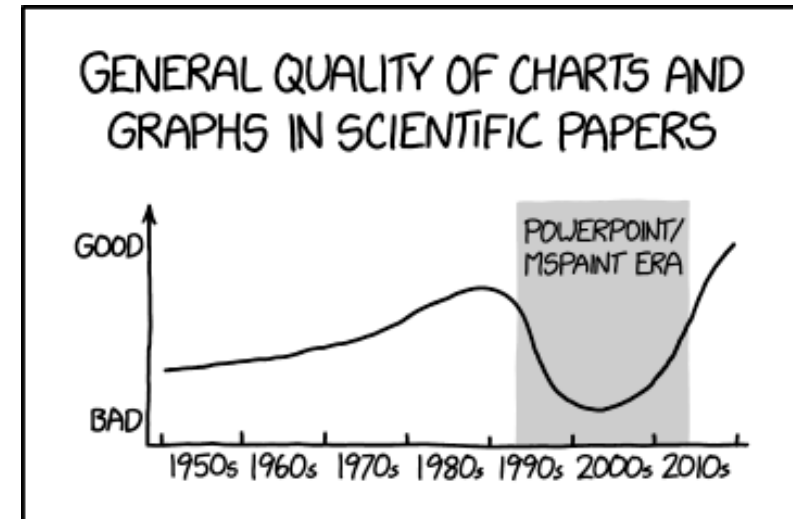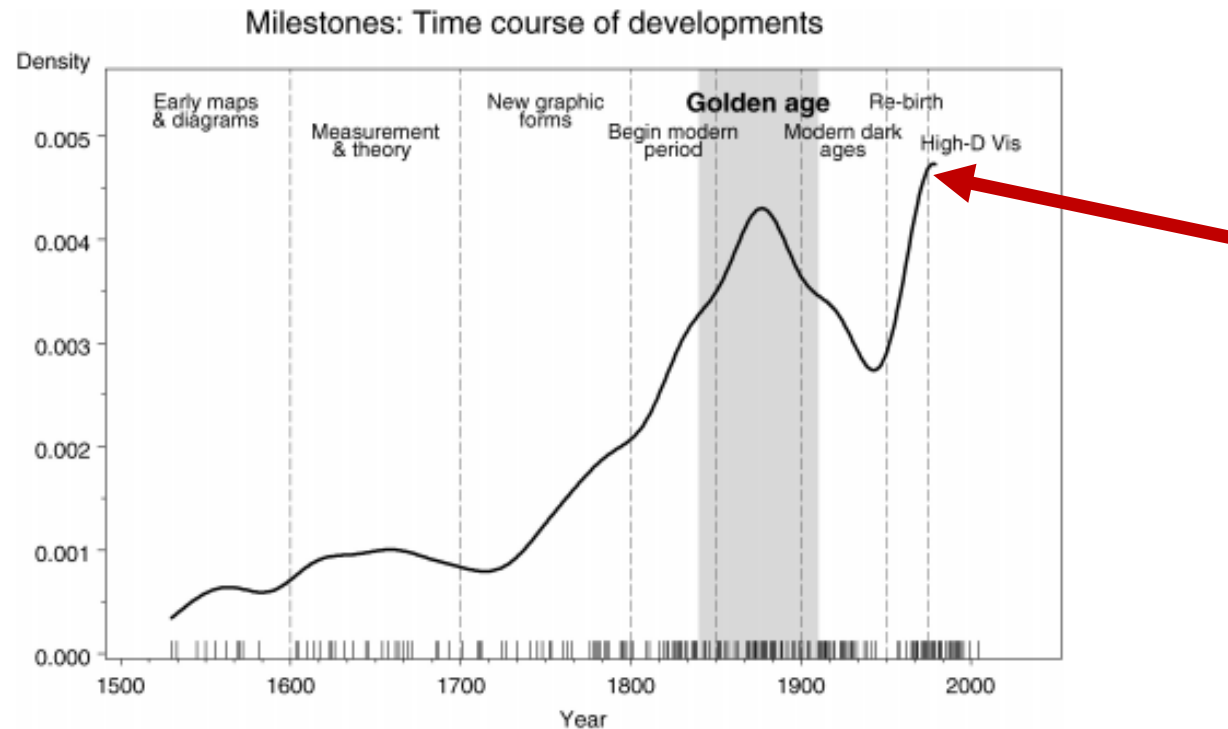- TV special effects
- Ted Talk

Gapminder tools:
https://www.gapminder.org/tools

> library('gapminder')

# A very brief history of data visualization

Currently undergoing a "Graphical re-birth"

# Next class: a grammar of graphics and ggplot

Start on homework 5 early!

Question : Find an interesting data visualization
- https://www.reddit.com/r/dataisbeautiful/
- https://flowingdata.com/