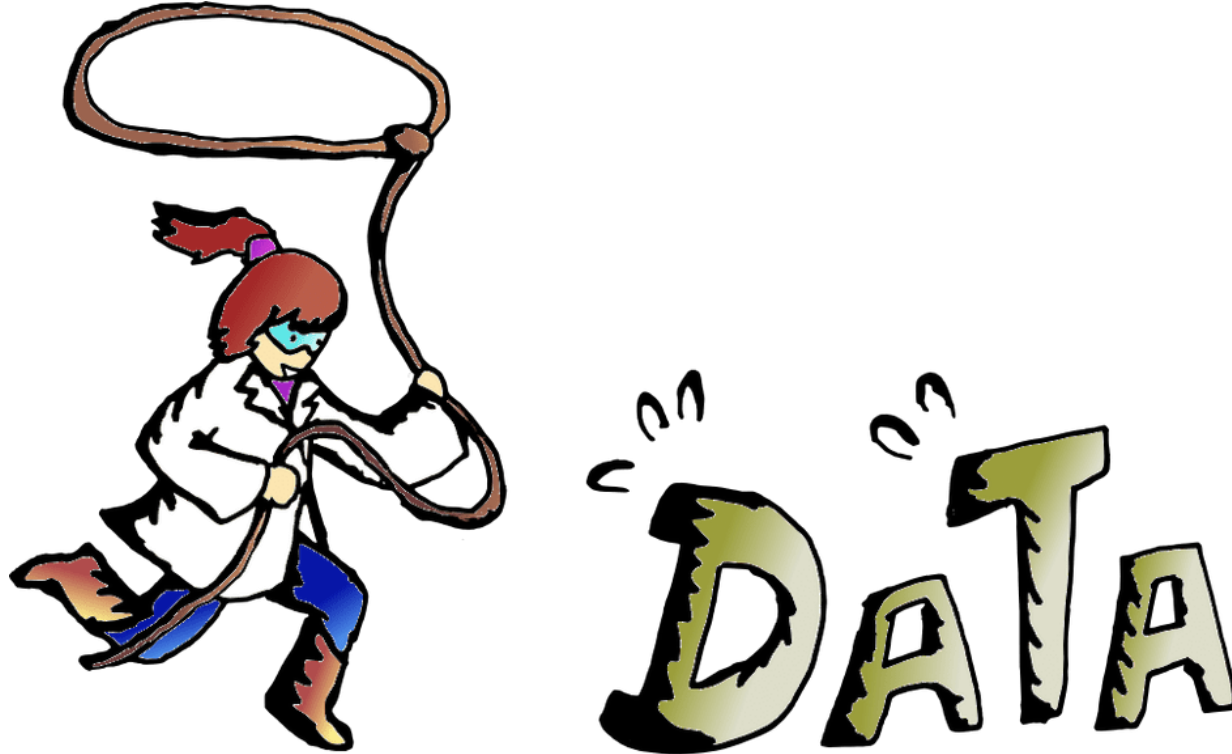# Data wrangling/manipulation

# Overview

Data wrangling/manipulation with dplyr

Brief history of data visualization

# Announcements

A practice midterm exam will be posted by next class

Slides with answers will also be posted soon

Get started on homework 5 early
- I recommend you do the dplyr exercises prior to next class!

Any other questions about class logistics?

# Plan for the semester

**Analysis**                                                              **R**

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
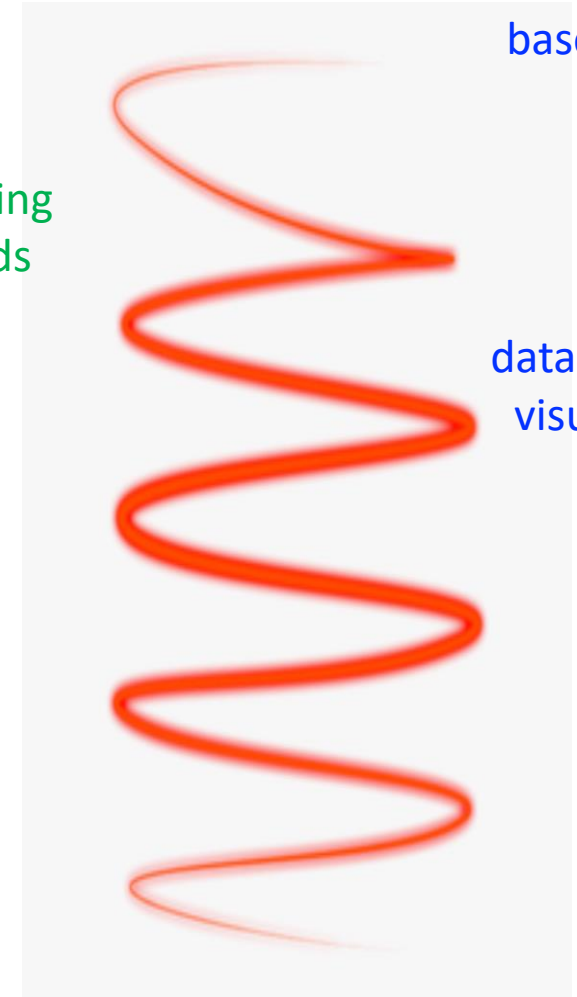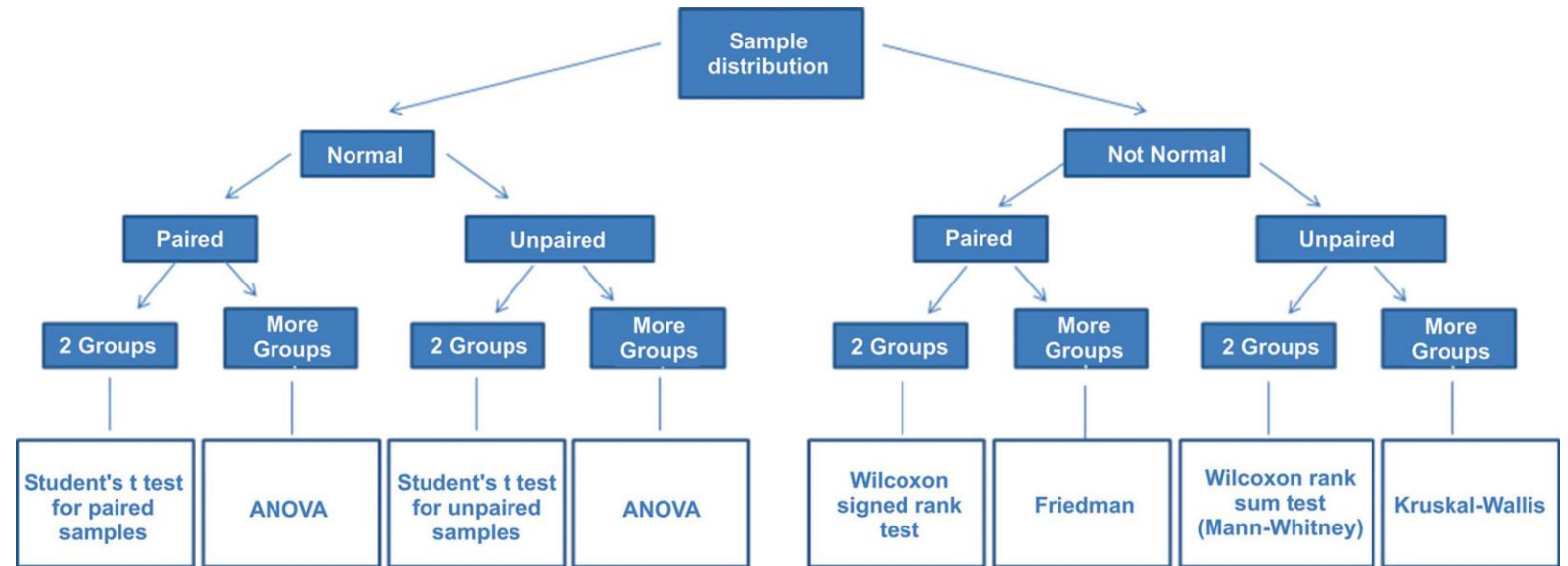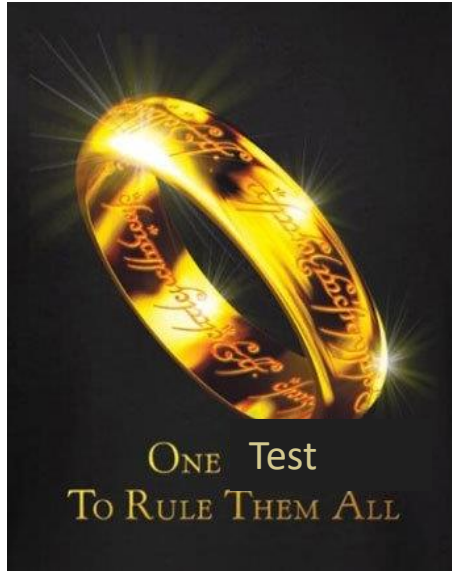
base R

resampling methods

data wrangling
visualization

# Very quick review



ONE Test
TO RULE THEM ALL



Just need to follow 5 steps!

# Very quick review



One Test To Rule Them All



```
                              Sample
                            distribution
              ┌──────────────────┴──────────────────┐
           Normal                               Not Normal
        ┌─────┴─────┐                        ┌──────┴──────┐
     Paired      Unpaired                 Paired        Unpaired
    ┌──┴──┐      ┌──┴──┐                 ┌──┴──┐        ┌──┴──┐
 2 Groups  More  2 Groups  More     2 Groups  More  2 Groups  More
          Groups          Groups             Groups          Groups
```

| Student's t test for paired samples | ANOVA | Student's t test for unpaired samples | ANOVA | Wilcoxon signed rank test | Friedman | Wilcoxon rank sum test (Mann-Whitney) | Kruskal-Wallis |

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

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
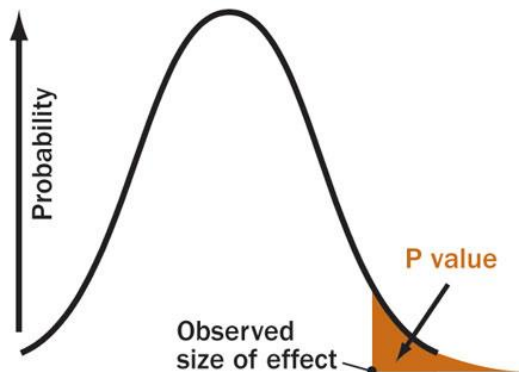


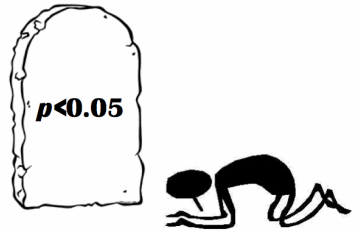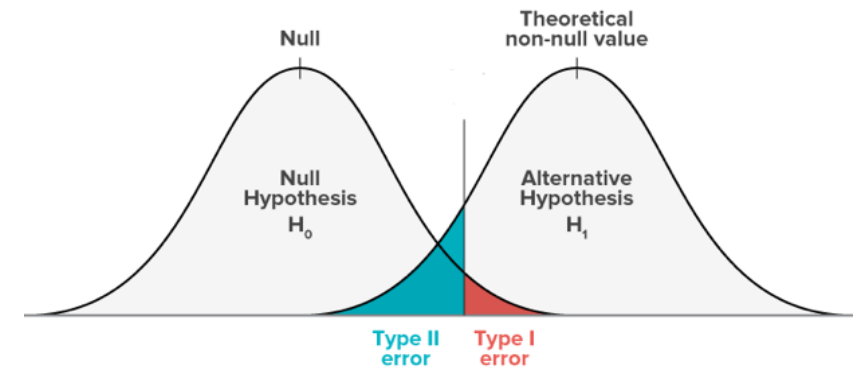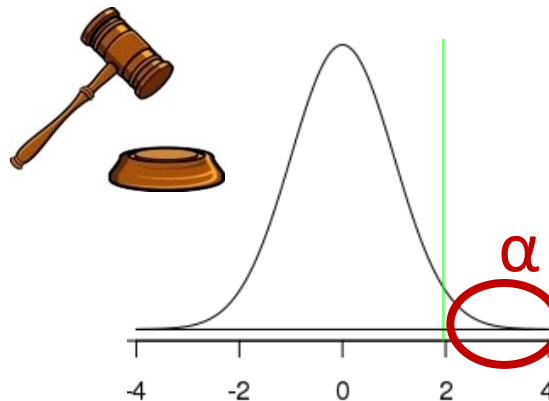Fisher (1890-1962)

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



variables                observations                values

# Messy data...

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

# Messy data…

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

# The 'tidyverse'

The packages share a common design philosophy

- Most written by Hadley Wickham

# dplyr:  A grammar for data wrangling

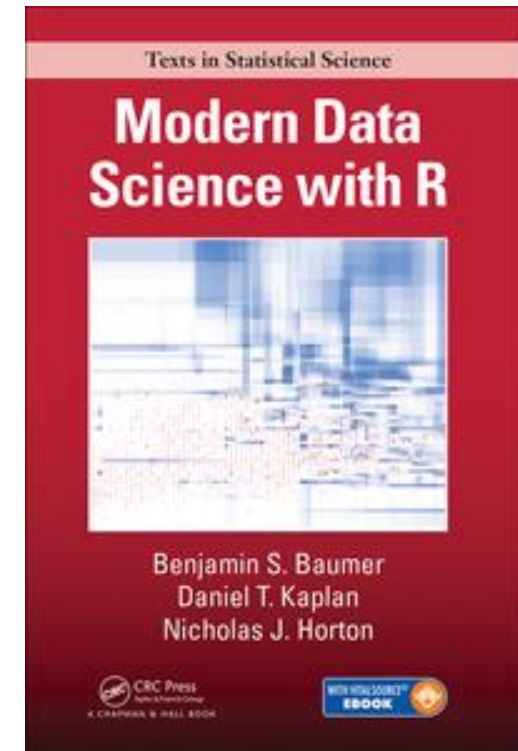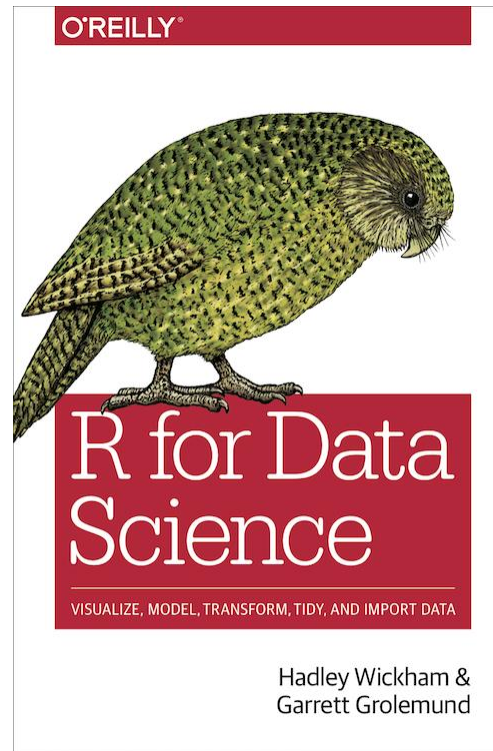**Grammar**:  a set of components that can be combined to achieve a goal

**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()

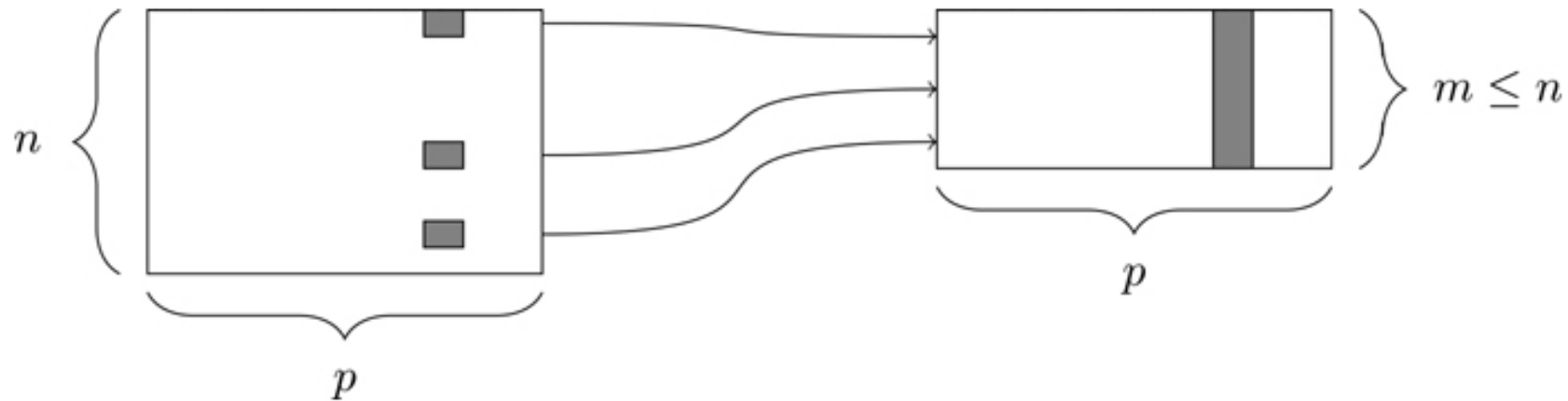All these function **take a data frame** and other arguments and **return a data frame**

> library(dplyr)   # load the dplyr package
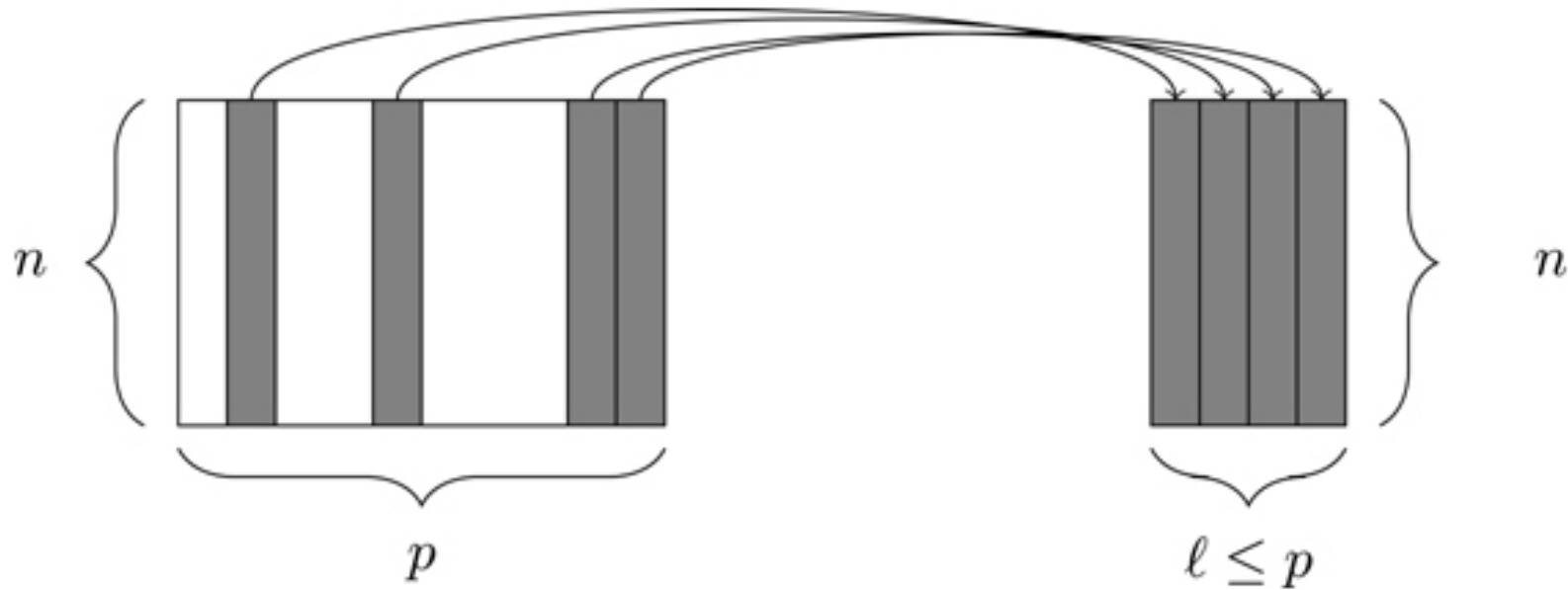
# Quick overview of the dplyr functions

# 1. filter()

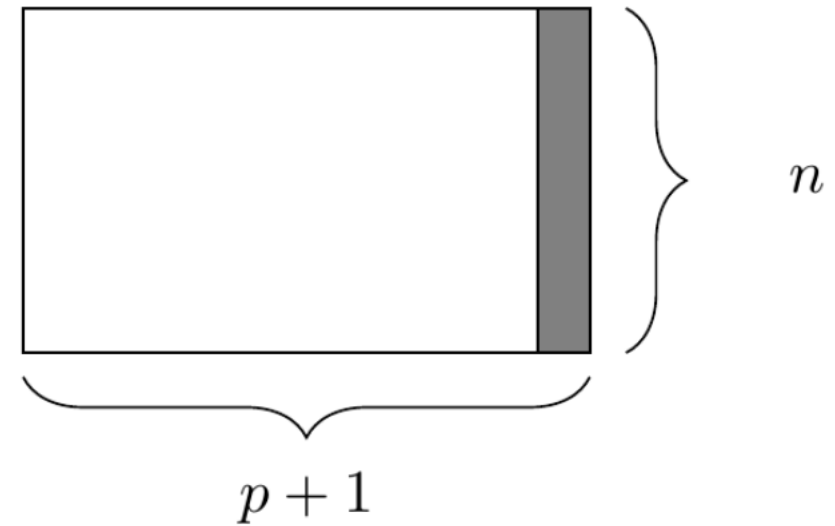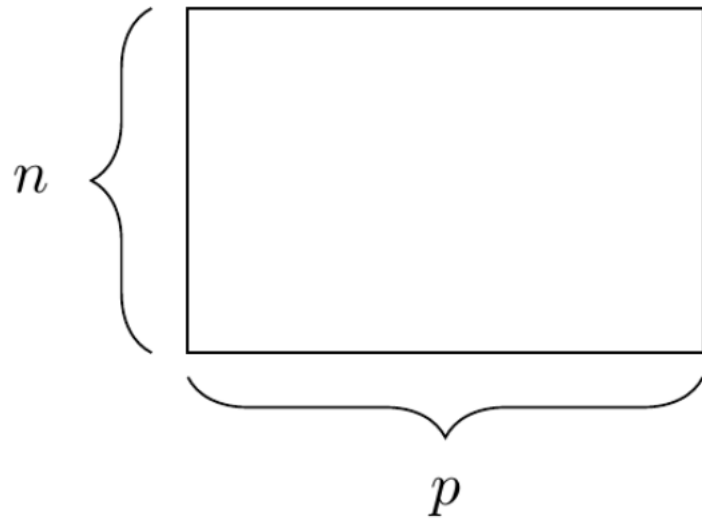The filter() function allows you to select a subset of rows in data frame

# 2. select()

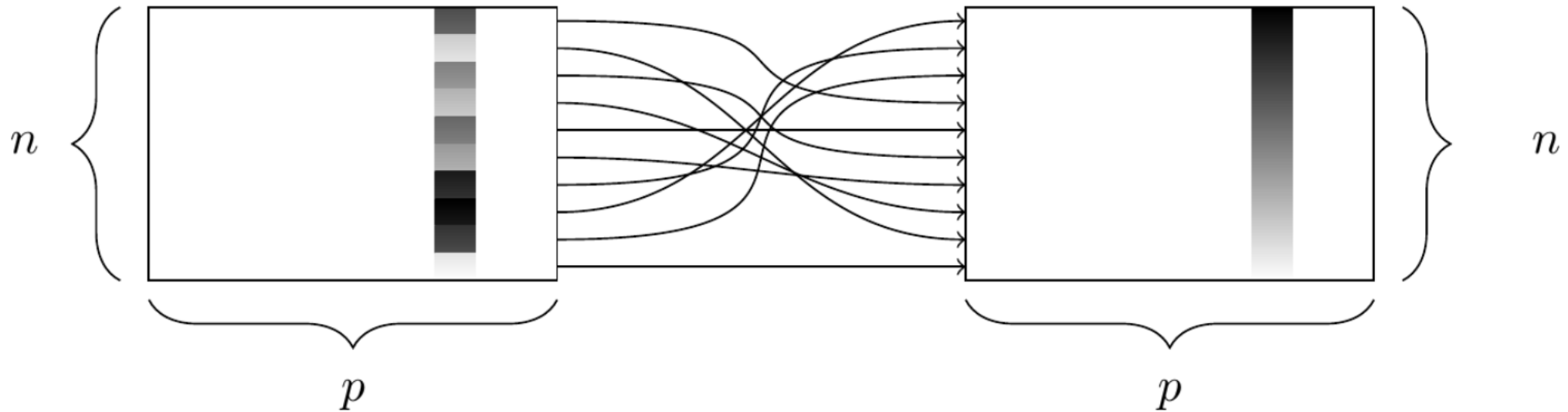The select() function allows you to select a subset of columns

# 3. mutate()

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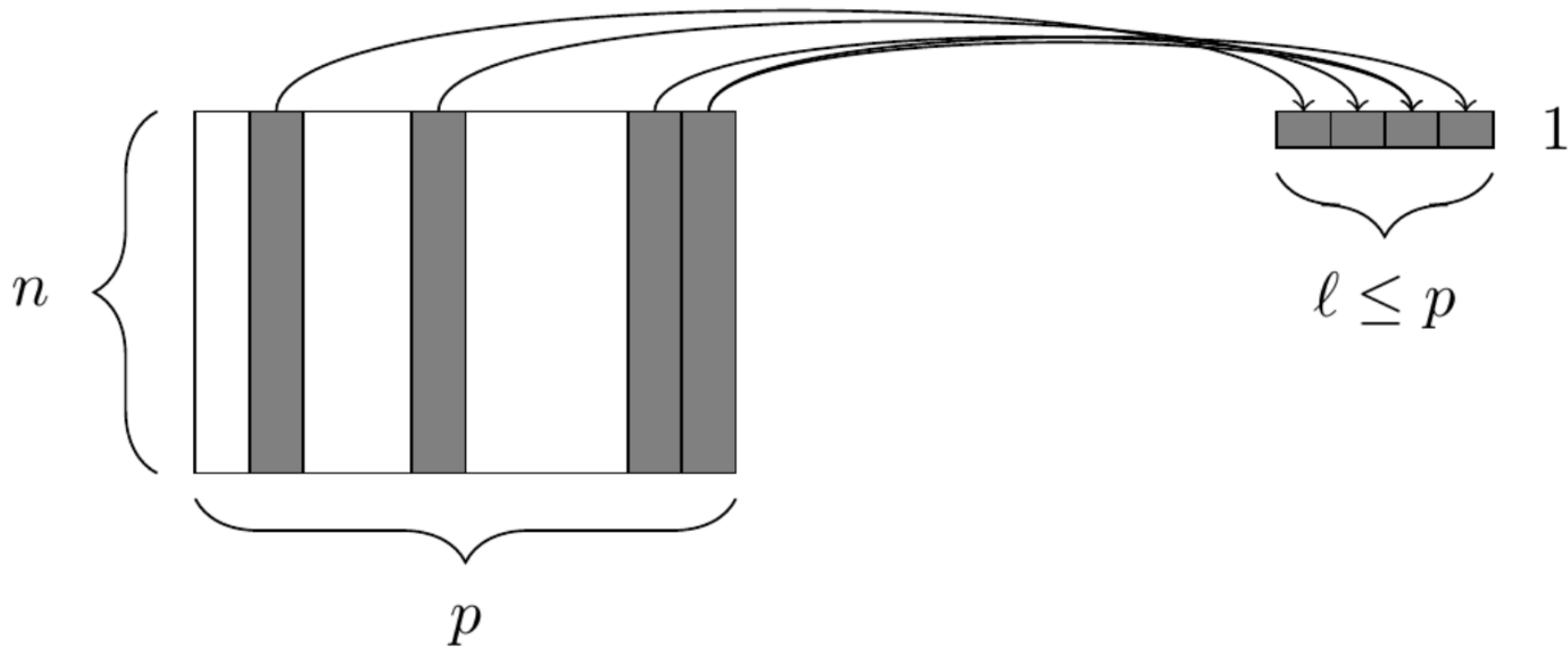
# 4. arrange()

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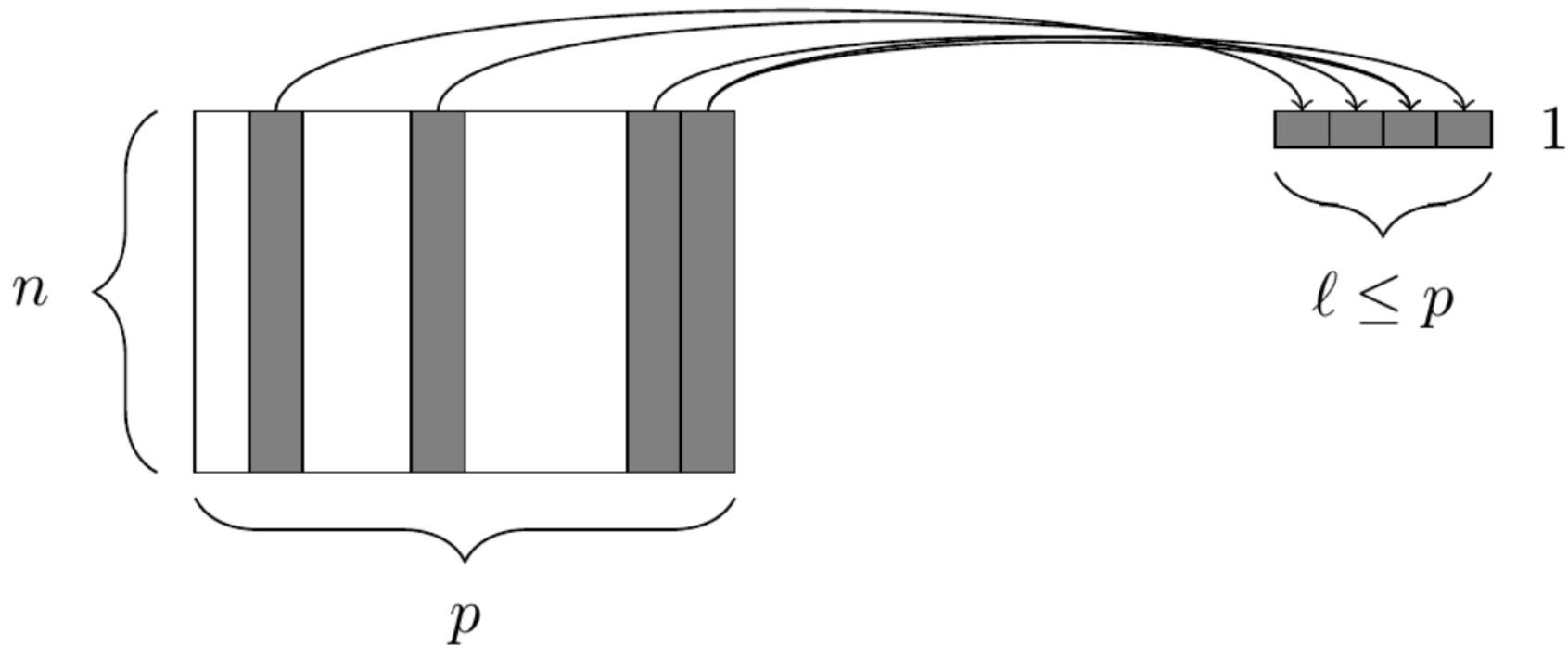- arrange(dec())  arranges from largest to smallest

# 5. summarize()

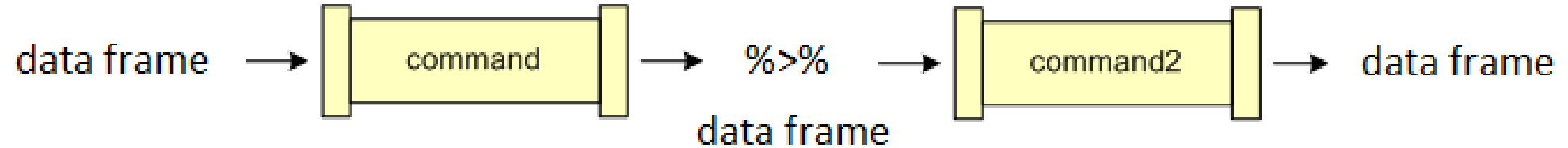The summarize() function reduces values in many rows into single values

# 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!

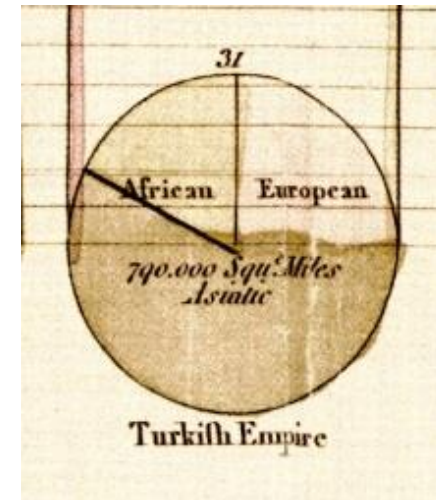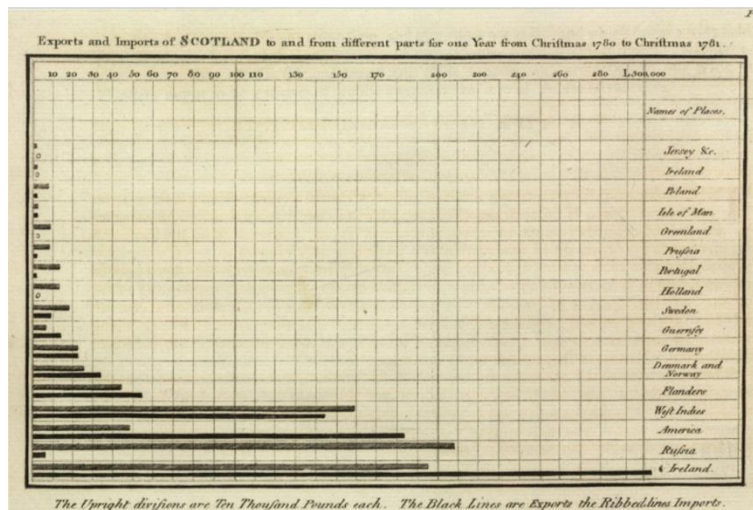# A very brief history of data visualization

# Data visualization

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

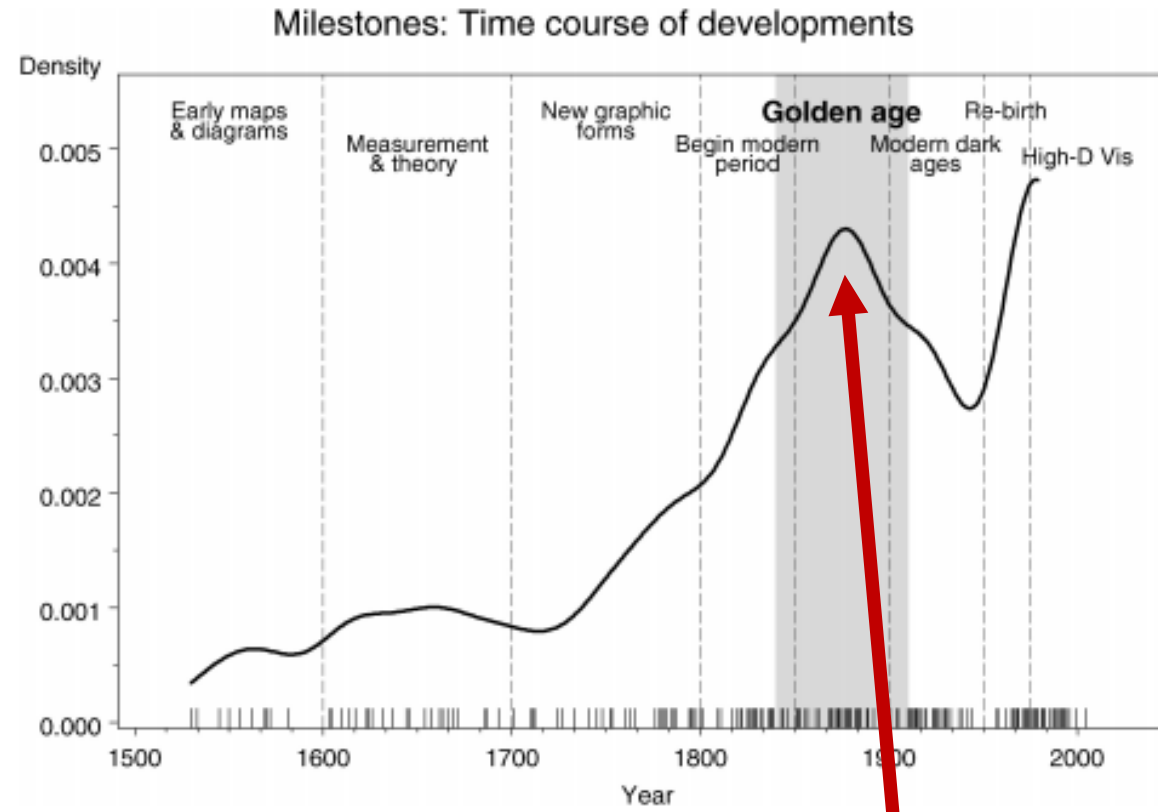# A very brief history of data visualization

The age of modern statistical graphs began around the beginning of the 19<sup>th</sup> century

[William Playfair](#) (1759-1823) credited with inventing the line graph, bar chart and pie chart
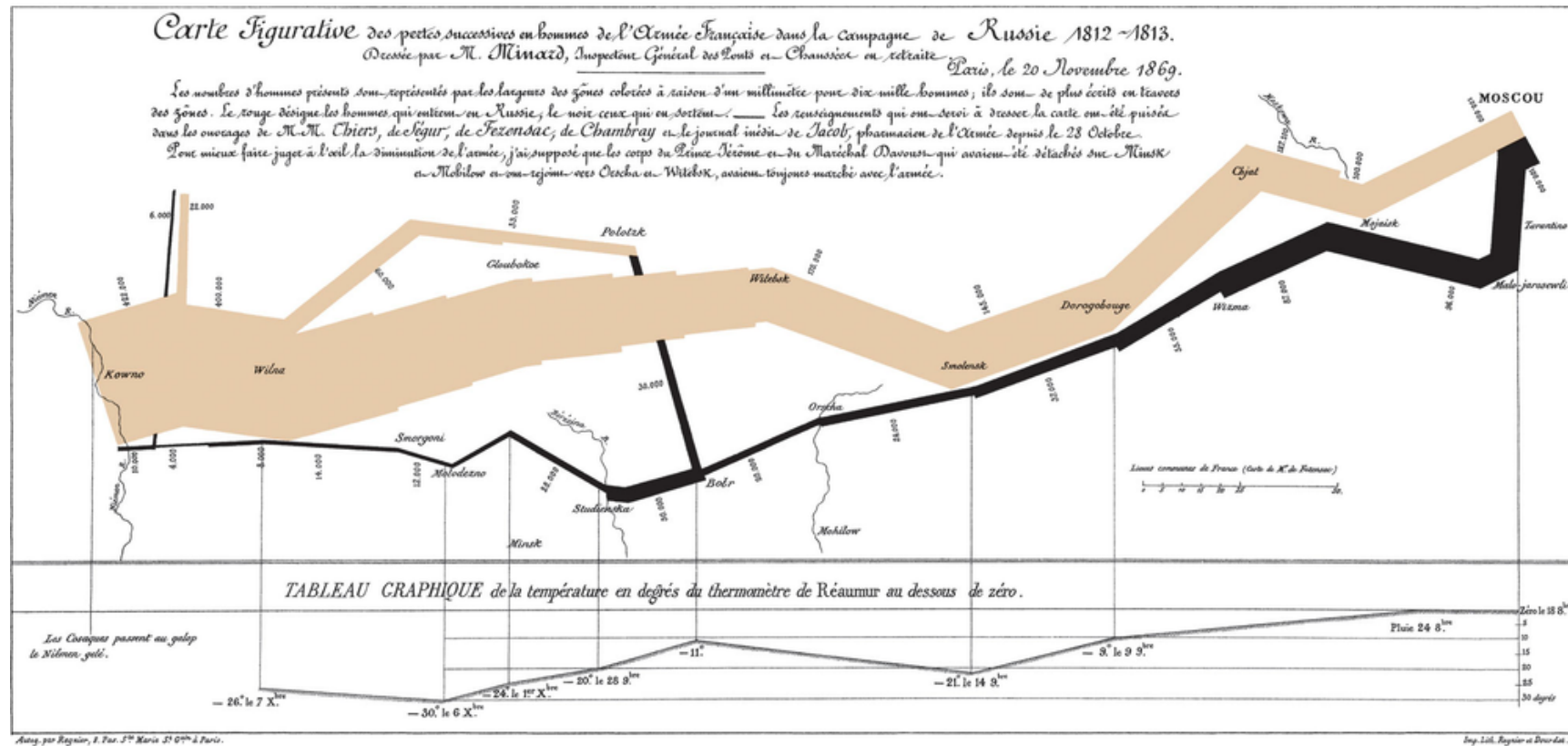
# A very brief history of data visualization

According to Friendly, statistical graphics researched its golden age between 1850-1900
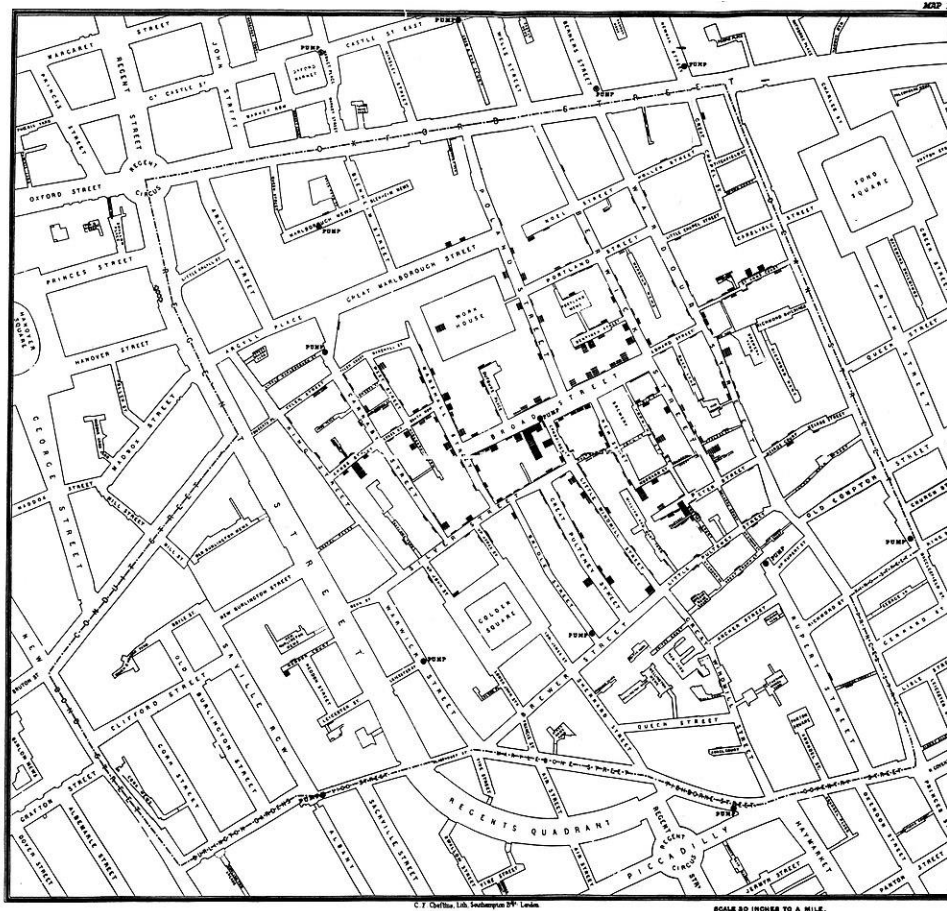
# A very brief history of data visualization

[Joseph Minard](...) (1781-1870)



Map of Napoleon's march on Russia

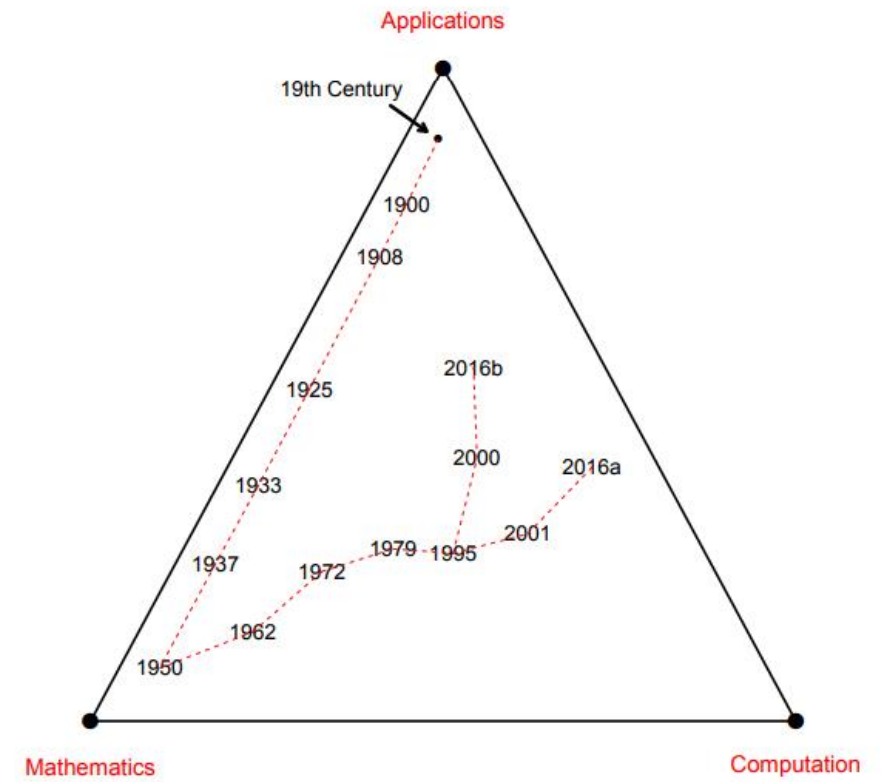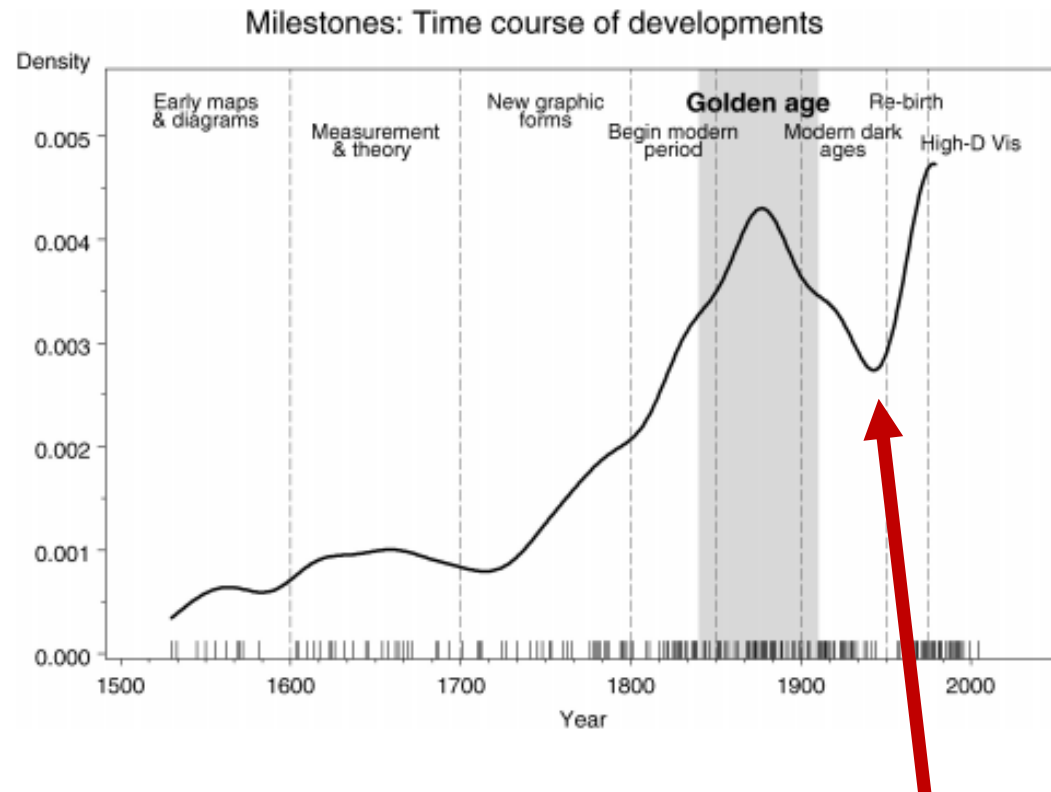# A very brief history of data visualization

[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)

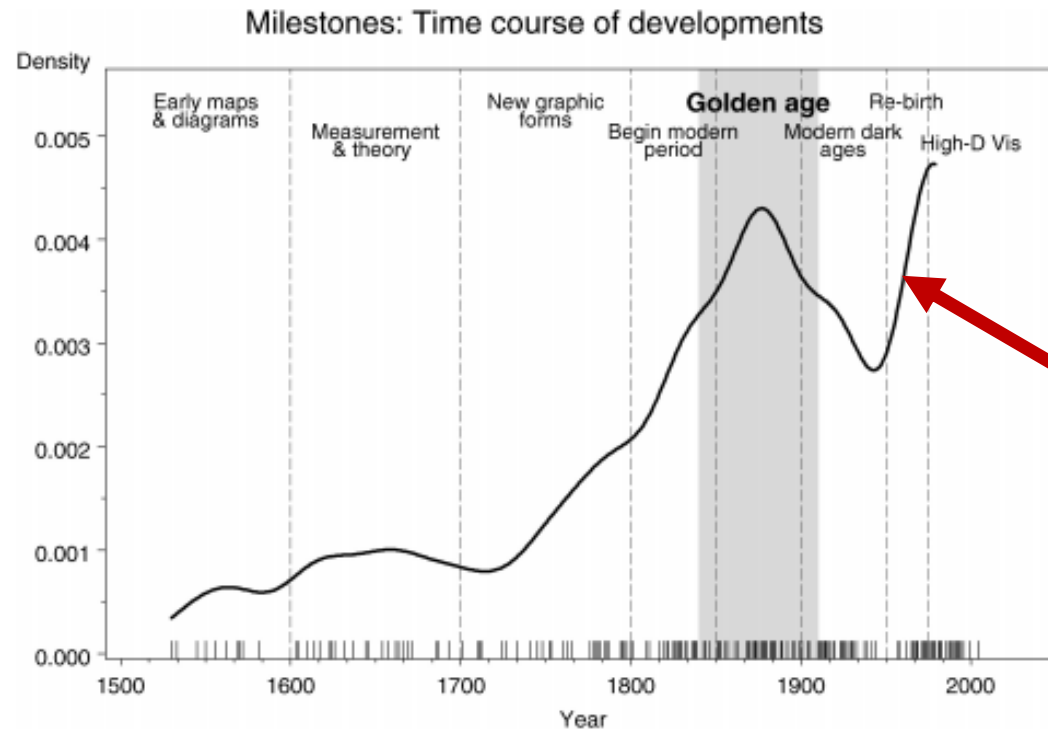# A very brief history of data visualization

"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"





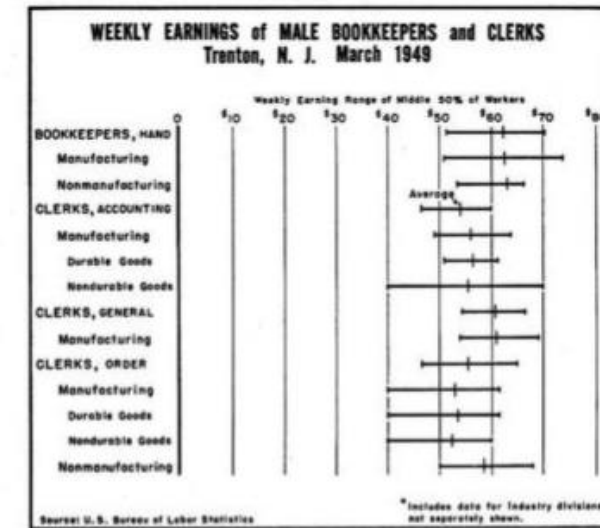[Spear](#) [1952](#), [Tukey](#) 1970

# A very brief history of data visualization

Currently undergoing a "Graphical re-birth"

# A very brief history of data visualization

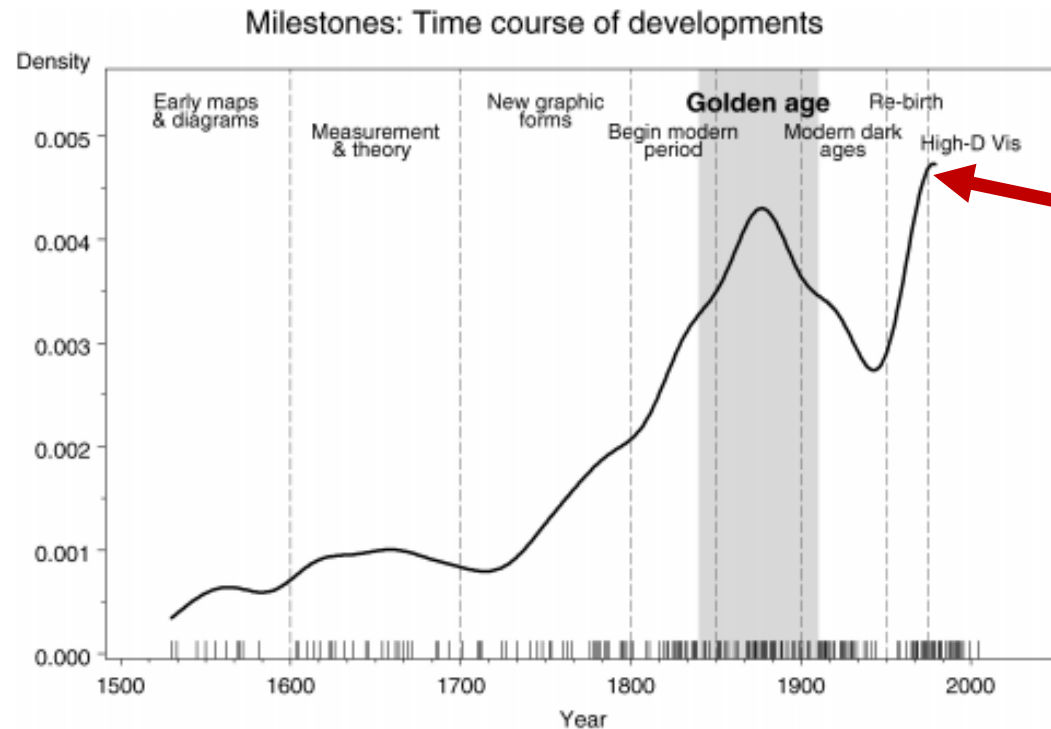Currently undergoing a "Graphical re-birth"



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

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

> library('gapminder')

# 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/