# INTRODUCTION TO STATISTICS

## LECTURE 4
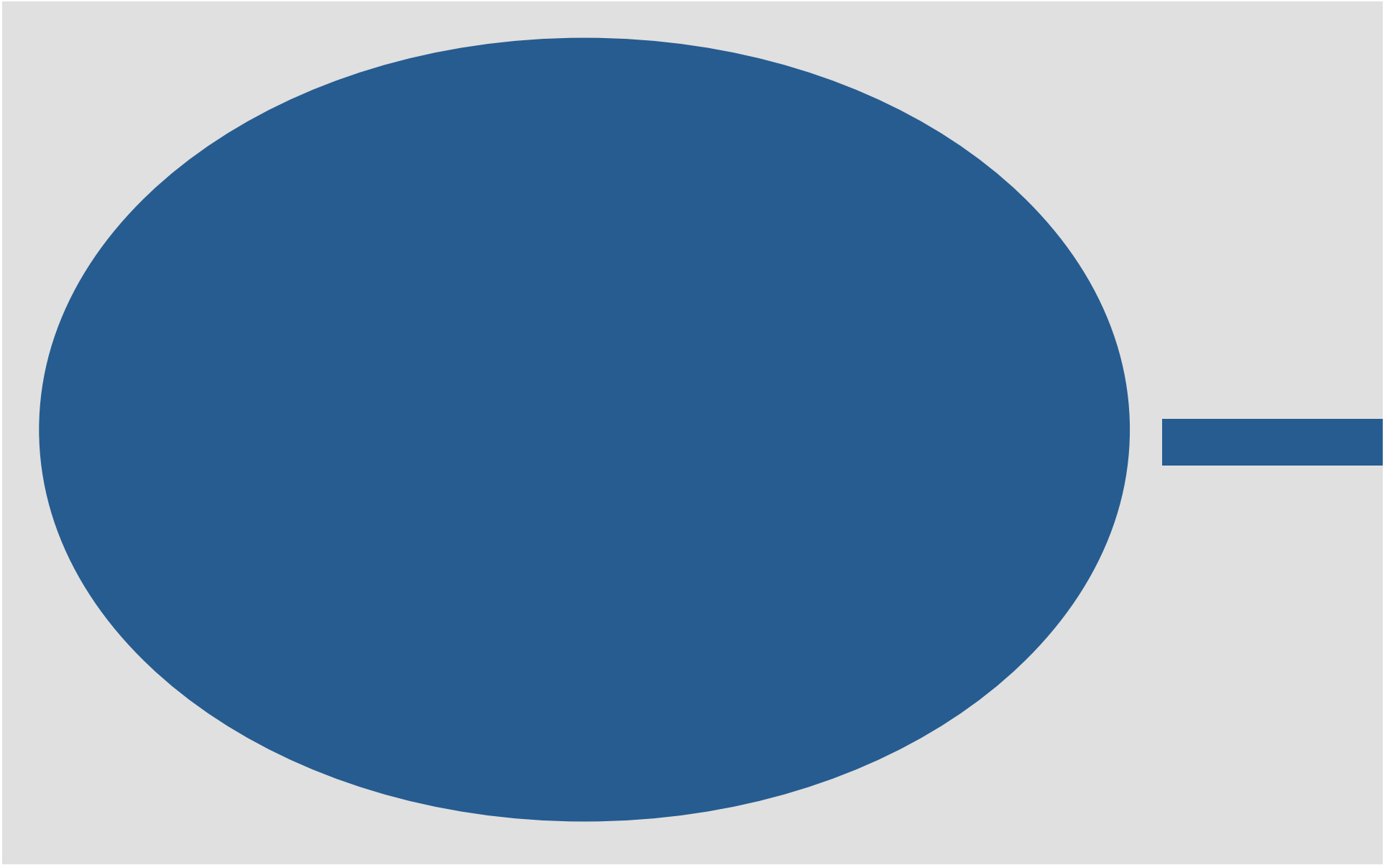
Pierre Dragicevic

# WHAT YOU WILL LEARN



Stats

This class

# STATISTICS

- **Statistics** is the study of the collection, analysis, interpretation, presentation and organization of data.

  Dodge, Y. (2006) The Oxford Dictionary of Statistical Terms, OUP.
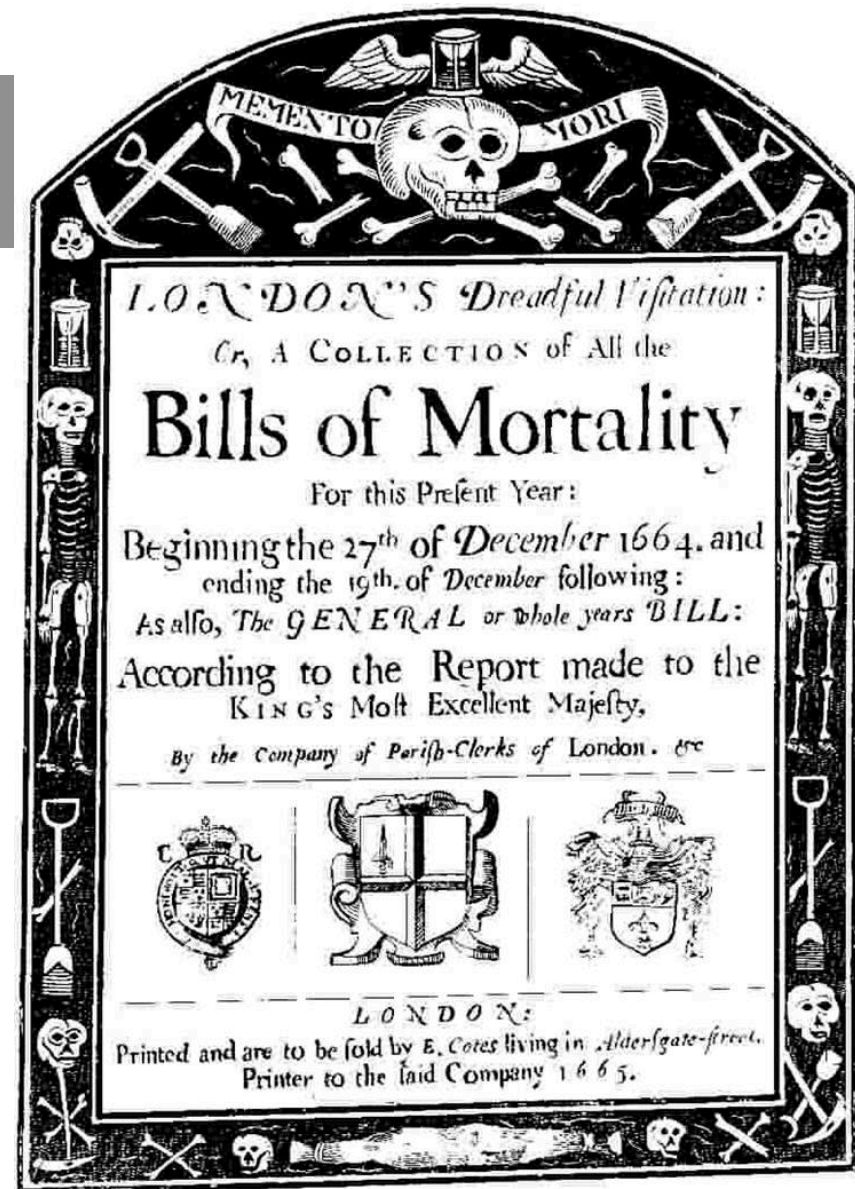
# STATISTICS

- 1750s German "*Statistik*"
  *"analysis of data about the state"*

- Quickly adopted in England
  (previously called "*political arithmetics*")

# STATISTICS

- **John Graunt, 1662**
  *Observations on the bills of mortality*


CAPTAIN JOHN GRAUNT

# THE TABLE OF CASUALTIES

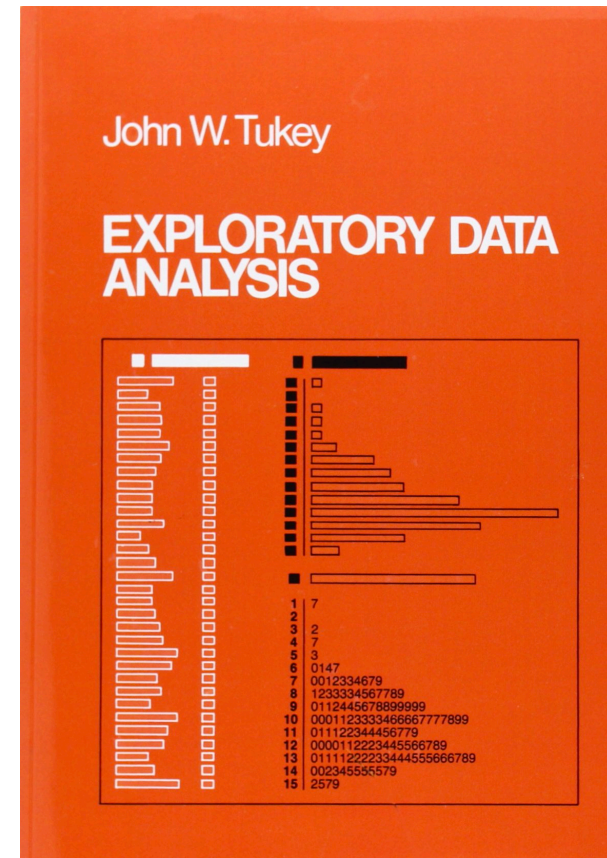| The Years of our Lord | 1647 | 1648 | 1649 | 1650 | 1651 | 1652 | 1653 | 1654 | 1655 | 1656 | 1657 | 1658 | 1659 | 1660 | 1629 | 1630 | 1631 | 1632 | 1633 | 1634 | 1635 | 1636 | 1629 1630 1631 1632 | 1633 1634 1635 1636 | 1647 1648 1649 1650 | 1651 1652 1653 1654 | 1655 1656 1657 1658 | 1629 1649 1659 | In 20 Years 1659 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Abortive, and stilborn | 335 | 329 | 327 | 351 | 389 | 381 | 384 | 433 | 483 | 419 | 463 | 467 | 421 | 544 | 499 | 439 | 410 | 445 | 500 | 475 | 507 | 523 | 1793 | 2005 | 1342 | 1587 | 1832 | 1247 | 8559 |
| Aged | 916 | 835 | 889 | 696 | 780 | 834 | 864 | 974 | 743 | 892 | 869 | 1176 | 909 | 1095 | 579 | 712 | 661 | 671 | 704 | 623 | 794 | 714 | 2473 | 2814 | 3336 | 3452 | 3680 | 2377 | 15757 |
| Ague, and Fevers | 1260 | 884 | 751 | 970 | 1038 | 1212 | 1282 | 1371 | 689 | 875 | 999 | 1800 | 2303 | 2148 | 956 | 1091 | 1115 | 1108 | 953 | 1279 | 1622 | 2360 | 4418 | 6235 | 3865 | 4903 | 4363 | 4010 | 23784 |
| Apoplex, and sodainly | 68 | 74 | 64 | 74 | 106 | 111 | 118 | 86 | 92 | 102 | 113 | 138 | 91 | 67 | 22 | 36 | | 17 | 24 | 35 | 26 | | | 75 | 85 | 280 | 421 | 445 | 177 | 1306 |
| Bleach | | | 1 | 3 | 7 | 2 | | | | | 1 | | | | | | | | | | | | | | 4 | 9 | 1 | 1 | 15 |
| Blasted | 4 | 1 | | | 6 | | | 4 | | 5 | 5 | 3 | 8 | | 13 | 8 | 10 | 13 | 6 | 4 | | 4 | 54 | 14 | 5 | 12 | 14 | 16 | 99 |
| Bleeding | 3 | 2 | 5 | 1 | 3 | 4 | 3 | 2 | 7 | 3 | 5 | 4 | 7 | 2 | 5 | 2 | 5 | 4 | 4 | 3 | | | 16 | 7 | 5 | 12 | 19 | 17 | 65 |
| Bloudy Flux, Scouring, and Flux | 155 | 176 | 802 | 289 | 833 | 762 | 200 | 386 | 168 | 368 | 362 | 233 | 346 | 251 | 449 | 438 | 352 | 348 | 278 | 512 | 346 | 330 | 1587 | 1466 | 1422 | 2181 | 1161 | 1597 | 7818 |
| Burnt, and Scalded | 3 | 6 | 10 | 5 | 11 | 8 | 5 | 7 | 10 | 5 | 7 | 4 | 6 | 6 | 3 | 10 | 7 | 5 | 1 | 3 | 12 | 3 | 25 | 19 | 24 | 31 | 26 | 19 | 125 |
| Calenture | 1 | | | 1 | | 2 | 1 | 1 | 1 | | 3 | | | | | | | | | 1 | 3 | | | 1 | 3 | 4 | 2 | 4 | 3 | 13 |
| Cancer, Gangrene, and Fistula Wolf | 26 | 29 | 31 | 19 | 31 | 53 | 36 | 37 | 73 | 31 | 24 | 35 | 63 | 52 | 20 | 14 | 23 | 28 | 27 | 30 | 24 | 30 | 85 | 112 | 105 | 157 | 150 | 114 | 609 |
| Canker, Sore-mouth, and Thrush | 66 | 28 | 54 | 42 | 68 | 51 | 53 | 72 | 44 | 81 | 19 | 27 | 73 | 68 | 6 | 4 | 4 | 1 | | | 5 | 74 | 15 | 79 | 190 | 244 | 161 | 133 | 689 |
| Childbed | 161 | 106 | 114 | 117 | 206 | 213 | 158 | 192 | 177 | 201 | 236 | 225 | 226 | 194 | 150 | 157 | 112 | 171 | 132 | 143 | 163 | 230 | 590 | 608 | 498 | 769 | 839 | 490 | 3364 |
| Chrisomes, and Infants | 1369 | 1254 | 1065 | 990 | 1237 | 1280 | 1050 | 1343 | 1089 | 1393 | 1162 | 1144 | 858 | 1123 | 2596 | 2378 | 2035 | 2208 | 2130 | 2315 | 2113 | 1895 | 9277 | 8453 | 4678 | 4910 | 4788 | 4519 | 32106 |
| Colick, and Wind | 103 | 71 | 85 | 82 | 76 | 102 | | 101 | 85 | 120 | 113 | 179 | 116 | 167 | 48 | 57 | | | | 37 | 50 | 105 | 87 | 341 | 359 | 497 | 147 | 1389 | |
| Cold, and Cough | | | | 41 | 36 | 21 | 58 | 30 | 31 | 33 | 24 | 10 | 58 | 51 | 55 | 45 | 54 | 50 | 57 | 174 | 207 | 00 | 77 | 140 | 43 | 598 | | | |
| Consumption, and Cough | 2423 | 2200 | 2388 | 1988 | 2350 | 2410 | 2286 | 2868 | 2606 | 3184 | 2757 | 3610 | 2982 | 3414 | 1827 | 1910 | 1713 | 1797 | 1754 | 1955 | 2080 | 2477 | 5157 | 8260 | 8999 | 9914 | 12157 | 7197 | 44487 |
| Convulsion | 684 | 491 | 530 | 493 | 569 | 653 | 666 | 828 | 702 | 1027 | 807 | 841 | 742 | 1031 | 52 | 87 | 18 | 2.1 | 221 | 386 | 418 | 700 | 498 | 1734 | 2198 | 2656 | 3377 | 1324 | 9073 |
| Cramp | | | 1 | | | | | | | | | | | | | | | | | | | | 01 | 00 | 01 | 0 | | | 2 |
| Cut of the Stone | | 2 | 1 | 3 | | 1 | 1 | 2 | 4 | 1 | 3 | 5 | 46 | 48 | | | | | 5 | 1 | 5 | 2 | 5 | 10 | 6 | 4 | 13 | 47 | 38 |
| Dropsy, and Tympany | 185 | 434 | 421 | 508 | 444 | 556 | 617 | 704 | 660 | 706 | 631 | 931 | 646 | 872 | 235 | 252 | 279 | 280 | 266 | 250 | 329 | 389 | 1048 | 1734 | 1533 | 1321 | 2982 | 1302 | 9623 |
| Drowned | 47 | 40 | 30 | 27 | 49 | 50 | 3 | 30 | 43 | 4 | 63 | 60 | 57 | 48 | 43 | 33 | 29 | 14 | 37 | 32 | 32 | 45 | 139 | 147 | 144 | 182 | 215 | 130 | 827 |
| Excessive drinking | | | 2 | | | | | | | | | | | | | | | | | | | | | 2 | | | | | |
| Executed | 8 | 17 | 29 | 43 | 24 | 12 | 19 | 21 | 19 | 22 | 20 | 18 | 7 | 18 | 19 | 13 | 12 | 18 | 13 | 13 | 13 | 13 | 62 | 52 | 97 | 76 | 79 | 55 | 384 |
| Fainted in a Bath | | | | 1 | | | | | | | | | | | | | | | | | | | | | | | | 1 | 1 |
| Falling-Sickness | 3 | 2 | 2 | 3 | | 3 | 4 | 1 | 4 | 3 | 1 | | 4 | 5 | 3 | 10 | 7 | 7 | 2 | 5 | 6 | 8 | 27 | 21 | 10 | 8 | 8 | 9 | 74 |
| Flox, and small Pox | 139 | 400 | 1190 | 184 | 525 | 1279 | 139 | 812 | 1294 | 823 | 835 | 409 | 1523 | 354 | 72 | 40 | 58 | 531 | 72 | 1354 | 293 | 127 | 701 | 1840 | 1913 | 2755 | 3361 | 2785 | 10576 |
| Found dead in the Streets | 6 | 6 | 9 | 8 | 7 | 9 | 14 | 4 | 3 | 4 | 9 | 11 | 2 | 6 | 18 | 33 | 26 | 6 | 13 | 8 | 24 | 24 | 83 | 69 | 39 | 34 | 27 | 29 | 243 |
| French-Pox | 18 | 29 | 15 | 18 | 21 | 20 | 10 | 20 | 29 | 23 | 25 | 53 | 51 | 31 | 17 | 12 | 12 | 12 | 7 | 17 | 12 | 22 | 53 | 48 | 80 | 81 | 130 | 83 | 392 |
| Frighted | 4 | 4 | 1 | | | 3 | | 2 | | 1 | | | 9 | | 1 | | | | 1 | | 3 | 2 | 3 | 9 | 5 | 2 | 2 | 21 | |
| Gout | 9 | 5 | 11 | 9 | 7 | 7 | 5 | 6 | 8 | 7 | 8 | 13 | 14 | 2 | 2 | 5 | 3 | 4 | 4 | 5 | 7 | 8 | 14 | 24 | 35 | 25 | 36 | 28 | 134 |
| Grief | 12 | 13 | 16 | 7 | 17 | 14 | 11 | 17 | 10 | 13 | 10 | 12 | 13 | 4 | 18 | 20 | 22 | 11 | 14 | 17 | 5 | 20 | 71 | 50 | 48 | 59 | 45 | 47 | 279 |
| Hanged, and made-away themselves | 11 | 10 | 13 | 14 | 9 | 14 | 15 | 9 | 14 | 16 | 24 | 18 | 11 | 36 | 8 | 8 | 6 | 15 | | | | 7 | 37 | 18 | 48 | 47 | 72 | 32 | 222 |
| Jaundice | 57 | 35 | 39 | 49 | 41 | 43 | 57 | 71 | 61 | 41 | 46 | 77 | 102 | 76 | 47 | 59 | 35 | 43 | 35 | 45 | 54 | 63 | 184 | 197 | 180 | 212 | 225 | 188 | 998 |
| Jaw-faln | 1 | 1 | | | 3 | | | | | | | 3 | 1 | | 10 | 16 | 13 | 8 | 10 | 10 | 4 | 11 | 47 | 35 | 02 | 5 | 6 | 10 | 95 |
| Imposthume | 75 | 61 | 65 | 59 | 80 | 105 | 79 | 90 | 92 | 122 | 80 | 134 | 105 | 96 | 58 | 76 | 73 | 74 | 50 | 62 | 73 | 130 | 282 | 315 | 260 | 354 | 428 | 228 | 1639 |
| Itch | | 1 | | | | | | | | | | | | | | | | | | | | 00 | 10 | 01 | | 11 | | | |
| Killed by several Accidents | 27 | 57 | 39 | 94 | 47 | 45 | 57 | 58 | 52 | 43 | 52 | 47 | 55 | 47 | 54 | 55 | 47 | 46 | 49 | 41 | 51 | 60 | 202 | 201 | 217 | 207 | 194 | 148 | 1021 |
| King's Evil | 27 | 26 | 22 | 19 | 22 | 20 | 26 | 26 | 27 | 24 | 23 | 28 | 28 | 54 | 16 | 25 | 18 | 38 | 35 | 20 | 26 | 69 | 97 | 150 | 94 | 94 | 102 | 66 | 537 |
| Lethargy | 3 | 4 | 2 | 4 | 4 | 4 | 3 | 10 | 9 | 4 | 6 | 2 | 6 | 4 | P | | | 2 | 3 | | 2 | 2 | 5 | 7 | 13 | 21 | 21 | 9 | 62 |
| Leprosy | | | 1 | | | | | | | | | | 1 | | 2 | | | | | | | 2 | 2 | 1 | | 1 | | 06 | |
| Livergrown, Spleen, and Rickets | 53 | 46 | 56 | 59 | 65 | 72 | 67 | 65 | 52 | 50 | 38 | 51 | 8 | 15 | 94 | 112 | 99 | 87 | 82 | 77 | 98 | 99 | 392 | 356 | 213 | 269 | 191 | 158 | 1421 |
| Lunatique | 12 | 18 | 6 | 11 | 5 | 11 | 5 | 12 | 6 | 7 | 13 | 5 | 14 | 14 | 0 | 11 | 0 | 51 | 4 | 2 | 5 | 28 | 13 | 47 | 39 | 31 | 26 | | |

| Disease | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Cancer, Gangrene, and Fistula, Wolf | 26 | 29 | 31 | 19 8 | 31 | 53 | 16 | 37 | 73 | 31 |
| Canker, Sore-mouth, and Thrush | 66 | 28 | 54 | 42 | 68 | 51 | 53 | 72 | 44 | 81 |
| Childbed | 161 | 106 | 114 | 117 | 206 | 213 | 158 | 192 | 177 | 201 |
| Chrisomes, and Infants | 1369 | 1254 | 1065 | 990 | 1237 | 1280 | 1050 | 1343 | 1089 | 1393 |
| Colick, and Wind | 103 | 71 | 85 | 82 | 76 | 102 | | 101 | 85 | 120 |
| Cold, and Cough | | | | | | | 41 | 36 | 21 | 58 |
| Consumption, and Cough | 2423 | 2200 | 2388 | 1988 | 2350 | 2410 | 2286 | 2868 | 2606 | 3184 |
| Convulsion | | | | 493 | 569 | 653 | 666 | 828 | 702 | 1027 |
| Cramp | | | | | | | | | | |
| Cut of the Stone | | | 2 | 1 | 3 | | 1 | 1 | 2 | 4 | 1 |
| Dropsy, and Tympany | 185 | 434 | 421 | 508 | 444 | 556 | 617 | 704 | 660 | 706 |
| Drowned | 47 | 40 | 30 | 27 | 49 | 50 | 3 | 30 | 43 | 4 |
| Excessive drinking | | | 2 | | | | | | | |
| Executed | 8 | 17 | 29 | 43 | 24 | 12 | 19 | 21 | 19 | 22 |
| Fainted in a Bath | | | | | 1 | | | | | |
| Falling-Sickness | 3 | 2 | 2 | 3 | | 3 | 4 | 1 | 4 | 3 |
| Flox, and small Pox | 139 | 400 | 1150 | 184 | 525 | 1279 | 139 | 812 | 1294 | 823 |
| Found dead in the Streets | 6 | 6 | 9 | 8 | 7 | 9 | 14 | 4 | 3 | 4 |
| French-Pox | 18 | 29 | 15 | 18 | 21 | 20 | 20 | 20 | 29 | 23 |
| Frighted | 4 | 4 | 1 | | 3 | | 2 | | 1 | 1 |
| Gout | 9 | 5 | 12 | 9 | 7 | 7 | 5 | 6 | 8 | 7 |
| Grief | 12 | 13 | 16 | 7 | 17 | 14 | 11 | 17 | 10 | 13 |

# STATISTICS

- ## John Graunt, 1662
  *Observations on the bills of mortality*
  - First "life tables"
  - Dispelled several myths about the plague
  - First analysis of sex ratio
  - First realistic estimate of the population in London

# STATISTICS

- Prompted collection of more data
- Parallel developments in probability theory
- Statistics then developed into a more rigorous discipline and was applied to:
  - Business & industry
  - Medicine
  - Science
  - ...

# WHAT ARE STATS?

- **Statistics** is the study of the collection, analysis, interpretation, presentation and organization of data in numerical form.

  Dodge, Y. (2006) The Oxford Dictionary of Statistical Terms, OUP.

# STATS & VISUALIZATION

• Statistical Charts
  – William Playfair (1759 – 1823)



Exports and Imports to and from DENMARK & NORWAY from 1700 to 1780.

BALANCE in FAVOUR of ENGLAND.

Line of Imports

BALANCE AGAINST

Line of Exports

The Bottom line is divided into Years, the Right hand line into L10,000 each.

Published as the Act directs, 1st May 1786, by Wm Playfair

# STATS & VISUALIZATION

- ## Exploratory Data Analysis
  - – Tukey, 1977

# Box-and-whisker plots with end values identified

**A)** HEIGHTS of 50 STATES      **B)** HEIGHTS of 219 VOLCANOS

**Figure 5.14** Generalized draftsman's display of the four-dimensional iris data (like Figure 5.11), with one flower plotted as an asterisk.

# STATS & VISUALIZATION

- Statistical Graphics

# STATS & VISUALIZATION

- Statistical Graphics

  Gelman and Unwin (2012) *Infovis and Statistical Graphics: Different Goals, Different Looks*

**Last letter of boys' names in 1950**

**Last letter of boys' names in 2010**

# 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 Preventible 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.ʳ 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 1856, the blue coincides with the black
The entire areas may be compared by following the blue, the red & the
    black lines enclosing them

**Mortality rates in the Crimean War from April 1854 to March 1856**

Annualized mortality per thousand

Sanitary commission arrives

- Zymotic diseases
- Wounds, injuries
- All other causes

1000

500

0

A M J J A S O N D J F M A M J J A S O N D J F M

1854          1855          1856

**British Army Size in the Crimean War from April 1854 to March 1856**

Average Army Size

40000

20000

0

A M J J A S O N D J F M A M J J A S O N D J F M

1854          1855          1856

46 64 54 77 67 68 62 56 38 — Population N = 9

$$\mu_x = \frac{\sum x}{N} = \frac{532}{9} = 59.11$$

The Mean of this Population ($\mu_x$) equals 59.11 (i.e. $\mu_x = 59.11$)

Random Sample n = 4

38 62 67 62

$$\overline{X} = \frac{\sum x}{n} = \frac{229}{4} = 57.25$$

The mean of this Random Sample equals 57.25 (i.e. $\overline{X} = 57.25$)

The Central Limit Theorem tells us that $\overline{X}$ is an unbiased estimate of $\mu_x$. ( i.e. $\overline{X} \longrightarrow \mu_x$)

In short, with only one random sample to go on, the mean of the sample ($\overline{X} = 57.25$) is our best estimate of the population mean ($\mu_x$)

German bombings in London during WWII

German bombings in London during WWII

# STATS & VISUALIZATION

- **Confirmatory Analysis**
  - Testing hypotheses
  - Example: is this new drug effective?
  - Strong focus on automatic procedures, computation and objectivity
  - Looking at data can impair objectivity:
    - Data dredging, snooping, fishing, mining

# STATS & VISUALIZATION

**Exploratory data analysis** is sometimes compared to detective work: it is the process of gathering evidence.

**Confirmatory data analysis** is comparable to a court trial: it is the process of evaluating evidence.

Exploratory analysis and confirmatory analysis *"can —and should—proceed side by side"* (Tukey; 1977).

Quoted from the SAS Institute

# STATS & VISUALIZATION

## Workflows

# STATS & VISUALIZATION

**Workflows**

Regent's Park

Cumberland

River Thames

Scale: one-half mile

German bombings in London during WWII

German bombings in London during WWII

# STATS & VISUALIZATION

**Workflows**

# STATS & VISUALIZATION

## Workflows

# STATS & VISUALIZATION

## Workflows

# STATS & VISUALIZATION

[...] the type of "atheoretical" search for patterns that we are sometimes warned against in graduate school **can save us from the humiliation of having to retract conclusions we might ultimately make on the basis of contaminated data**.

We are warned against fishing expeditions for understandable reasons, but **blind application of models without screening our data is a far graver error**.

(Wilkinson, 1999)

# STATS & VISUALIZATION

## Workflows

# STATS & VISUALIZATION

## Workflows

# STATS & VISUALIZATION

## Workflows

# STATS & VISUALIZATION



Example from Jenny Weaver

# STATS & VISUALIZATION

## Workflows

# STATS & VISUALIZATION

**Odds of Coronary Heart Disease by Baldness**

Vertex Baldness

Frontal Baldness

Odds Ratio

None   Mild   Moderate   Severe       None   Present

**Baldness**

From agrippastake.blogspot.dk

# STATS & VISUALIZATION

## Workflows

# STATS & VISUALIZATION

## Workflows

# STATS & VISUALIZATION

- In Teaching



a  *t* and normal distributions

b  *P* values of *t* statistic

# WHAT ARE STATS?

- A set of tools and methods

- Old tradition:
  - Draws from mathematics & probability theory
  - A (generally) strong focus on (computationally cheap) numerical calculations

- Good for:
  - Summarizing data for presentation
  - Rigorously testing hypotheses
  - Making predictions
  - Helping make rational decisions

# STATISTICAL TOOLS

## DESCRIPTIVE STATISTICS

# AN EXAMPLE

- Selling encyclopedias

| day | Seller 1 | Seller 2 | Seller 3 | Seller 4 | Seller 5 | Seller 6 |
|-----|----------|----------|----------|----------|----------|----------|
| 1 | €320 | €80 | €139 | €330 | €133 | €387 |
| 2 | €74 | €60 | €98 | €44 | €182 | €29 |
| 3 | €340 | €67 | €42 | €100 | €51 | €91 |
| 4 | €322 | €54 | €89 | €44 | €67 | €886 |
| 5 | €146 | €195 | €47 | €173 | €49 | €227 |
| 6 | €24 | €288 | €124 | €111 | €730 | €79 |
| 7 | €42 | €249 | €26 | €77 | €672 | €45 |
| 8 | €76 | €67 | €140 | €382 | €195 | €171 |
| 9 | €99 | €312 | €125 | €123 | €43 | €98 |
| 10 | €915 | €77 | €106 | €250 | €149 | €70 |
| 11 | €202 | €504 | €101 | €205 | €682 | €134 |
| 12 | €47 | €167 | €126 | €48 | €93 | €63 |
| 13 | €34 | €65 | €55 | €56 | €333 | €1,157 |
| 14 | €76 | €46 | €89 | €104 | €56 | €470 |
| 15 | €75 | €34 | €184 | €35 | €299 | €205 |
| 16 | €68 | €37 | €275 | €170 | €57 | €192 |

| day | Seller 1 | Seller 2 | Seller 3 | Seller 4 | Seller 5 | Seller 6 |
|---|---|---|---|---|---|---|
| 1 | €320 | €80 | €139 | €330 | €133 | €387 |
| 2 | €74 | €60 | €98 | €44 | €182 | €29 |
| 3 | €340 | €67 | €42 | €100 | €51 | €91 |
| 4 | €322 | €54 | €89 | €44 | €67 | €886 |
| 5 | €146 | €195 | €47 | €173 | €49 | €227 |
| 6 | €24 | €288 | €124 | €111 | €730 | €79 |
| 7 | €42 | €249 | €26 | €77 | €672 | €45 |
| 8 | €76 | €67 | €140 | €382 | €195 | €171 |
| 9 | €99 | €312 | €125 | €123 | €43 | €98 |
| 10 | €915 | €77 | €106 | €250 | €149 | €70 |
| 11 | €202 | €504 | €101 | €205 | €682 | €134 |
| 12 | €47 | €167 | €126 | €48 | €93 | €63 |
| 13 | €34 | €65 | €55 | €56 | €333 | €1,157 |
| 14 | €76 | €46 | €89 | €104 | €56 | €470 |
| 15 | €75 | €34 | €184 | €35 | €299 | €205 |
| 16 | €68 | €37 | €275 | €170 | €57 | €192 |
| 17 | €126 | €23 | €114 | €30 | €43 | €60 |
| 18 | €43 | €290 | €89 | €446 | €57 | €226 |
| 19 | €149 | €215 | €43 | €63 | €62 | €72 |
| 20 | €31 | €81 | €26 | €469 | €60 | €39 |
| 21 | €81 | €127 | €47 | €68 | €315 | €566 |
| 22 | €141 | €70 | €317 | €40 | €160 | €42 |
| 23 | €113 | €947 | €203 | €102 | €108 | €76 |
| 24 | €209 | €48 | €81 | €102 | €50 | €56 |
| 25 | €94 | €95 | €67 | €21 | €54 | €41 |
| 26 | €159 | €125 | €67 | €263 | €69 | €173 |
| 27 | €271 | €176 | €250 | €35 | €48 | €24 |
| 28 | €52 | €85 | €77 | €136 | €95 | €82 |
| 29 | €30 | €12 | €317 | €157 | €240 | €58 |
| 30 | €104 | €31 | €181 | €113 | €45 | €27 |

# CENTRAL TENDENCY

| Name & Meaning | Formula / Example | Used for |
|---|---|---|
| **Arithmetic Mean** [average] | $\dfrac{sum}{size} = \dfrac{a+b+c}{3}$ | Most situations ("average item") |
| **Median** [middle value] | Middle of sorted list (2 middles? Average 'em) | Wildly varying samples (houses, incomes) |
| **Mode** [most popular] | Most popular value | No compromises (winner takes all) |
| **Geometric Mean** [average factor] | $\sqrt[3]{abc}$ | Investments, growth, area, volume |
| **Harmonic Mean** [average rate] | $\dfrac{3}{\dfrac{1}{a}+\dfrac{1}{b}+\dfrac{1}{c}}$ | Speed, production, cost |

# CENTRAL TENDENCY

## Mode (Most Popular)

# CENTRAL TENDENCY

negative skew          symmetric          positive skew

# CENTRAL TENDENCY

# CENTRAL TENDENCY

# CENTRAL TENDENCY



(a) Negatively skewed

Mode
Median
Mean
Frequency
Negative Direction

(b) Normal (no skew)

Mean
Median
Mode
Perfectly Symmetrical
Distribution

(c) Positively skewed

Mode
Median
Mean
Positive Direction

# CENTRAL TENDENCY

**THE UK INCOME DISTRIBUTION IN 2006 / 7**

Number of individuals (millions)

Median, £377

Mean, £463

2.7 million individuals with income above £1,000 per week

Income, £ per week, 2006/07 prices

SOURCE: HBAI data

From Michael Blastland

# DISPERSION

## Standard Deviation

$$\sigma = \sqrt{\frac{1}{N} \sum_{i=1}^{N} (x_i - \mu)^2}$$



Image from Wikipedia

# DEPENDENCE

- Correlation



exam results

hours revising

POSITIVE CORRELATION
- people who do more revision get higher exam results.

# DEPENDENCE

- Correlation

# DEPENDENCE

- Correlation

$r = -0.08$

## Average Sales

| Seller 1 | Seller 2 | Seller 3 | Seller 4 | Seller 5 | Seller 6 |
|----------|----------|----------|----------|----------|----------|
| €149     | €154     | €122     | €143     | €173     | €195     |

# Average Sales

| Seller 1 | Seller 2 | Seller 3 | Seller 4 | Seller 5 | Seller 6 |
| --- | --- | --- | --- | --- | --- |
| €149 | €154 | €122 | €143 | €173 | €195 |

# LOOKING INTO THE FUTURE

September 2014

October 2014

November 2014

december 2014

September 2014

October 2014

November 2014

December 2014

| day | Seller 1 | Seller 2 | Seller 3 | Seller 4 | Seller 5 | Seller 6 |
|-----|----------|----------|----------|----------|----------|----------|
| 1 | €320 | €80 | €139 | €330 | €133 | €387 |
| 2 | €74 | €60 | €98 | €44 | €182 | €29 |
| 3 | €340 | €67 | €42 | €100 | €51 | €91 |
| 4 | €322 | €54 | €89 | €44 | €67 | €886 |
| 5 | €146 | €195 | €47 | €173 | €49 | €227 |
| 6 | €24 | €288 | €124 | €111 | €730 | €79 |
| 7 | €42 | €249 | €26 | €77 | €672 | €45 |
| 8 | €76 | €67 | €140 | €382 | €195 | €171 |
| 9 | €99 | €312 | €125 | €123 | €43 | €98 |
| 10 | €915 | €77 | €106 | €250 | €149 | €70 |
| 11 | €202 | €504 | €101 | €205 | €682 | €134 |
| 12 | €47 | €167 | €126 | €48 | €93 | €63 |
| 13 | €34 | €65 | €55 | €56 | €333 | €1,157 |
| 14 | €76 | €46 | €89 | €104 | €56 | €470 |
| 15 | €75 | €34 | €184 | €35 | €299 | €205 |
| 16 | €68 | €37 | €275 | €170 | €57 | €192 |

September 2014



How much can we trust this chart?

# STATISTICAL TOOLS

## INFERENTIAL STATISTICS

# SAMPLING ERROR

# SAMPLING ERROR

- Terminology:

  – Population vs. sample
  – Sample **statistic** (mean, median, etc.)
  – Population **parameter** (mean, median, etc.)

SAMPLING ERROR

Diastolic Blood Pressure?

Mean = 78 mm Hg

Samples

Mean = 75

Mean = 67

Mean = 71.3

From Lisa Sullivan

# SAMPLING ERROR

- ## Sampling distribution of a statistic
  - Demo

# SAMPLING ERROR

- Bootstrapping



Bootstraps

# SAMPLING ERROR

- Bootstrapping

Complete element space

Complete element space

initial sample with N elements

Complete element space

initial sample with N elements

Complete element space

initial sample with N elements

resample with replacement 1

Complete element space

initial sample with N elements

$N_b$ bootstrap samples

1
2
3
4
⋮
$N_b$

Complete element space

initial sample with N elements

$N_b$ bootstrap samples

1
2
3
4
⋮
$N_b$

**Theorem** (B. Efron, Ann. Statist. 1979)

When N tend to infinity, the distribution of average values computed from bootstrap samples is equal to the distribution of average values obtained from ALL samples with N elements which can be constructed from the complete space. Thus the width of the distribution gives an evaluation of the sample quality.

# SAMPLING ERROR

- Bootstrapping video

# SAMPLING ERROR

- Bootstrapping correlations



$r = 0.78$

# SAMPLING ERROR

- Bootstrapping correlations



bootstrapping on a correlation coeficient

# SAMPLING ERROR

- How did people do before computers?

# MORE HISTORY

- Abraham De Moivre
  1667 - 1754

# MORE HISTORY

- Abraham De Moivre
  1667 - 1754

# MORE HISTORY

- ## Abraham De Moivre
  ### 1667 - 1754

# MORE HISTORY

- ## Abraham De Moivre
  1667 - 1754

# MORE HISTORY

- Abraham De Moivre
  1667 - 1754

# MORE HISTORY



Number of individuals

Height in inches

# NORMAL DISTRIBUTION

# NORMAL DISTRIBUTION

- **Sir Francis Galton**
  1822 – 1911

  Bean Machine
  or Galton Board:

# NORMAL DISTRIBUTION

## Central Limit Theorem

Given certain conditions, the arithmetic mean of a sufficiently large number of iterates of independent random variables, each with a well-defined expected value and well-defined variance, will be approximately normally distributed

# NORMAL DISTRIBUTION

"Exact" Confidence Intervals

$$\overline{X} \pm t\frac{s}{\sqrt{n}}$$

t ~ 1.96 for large samples

# SAMPLING ERROR

- Terminology:
  - **Point estimate** = sample statistic = best guess
  - **Interval estimate** = other good guesses
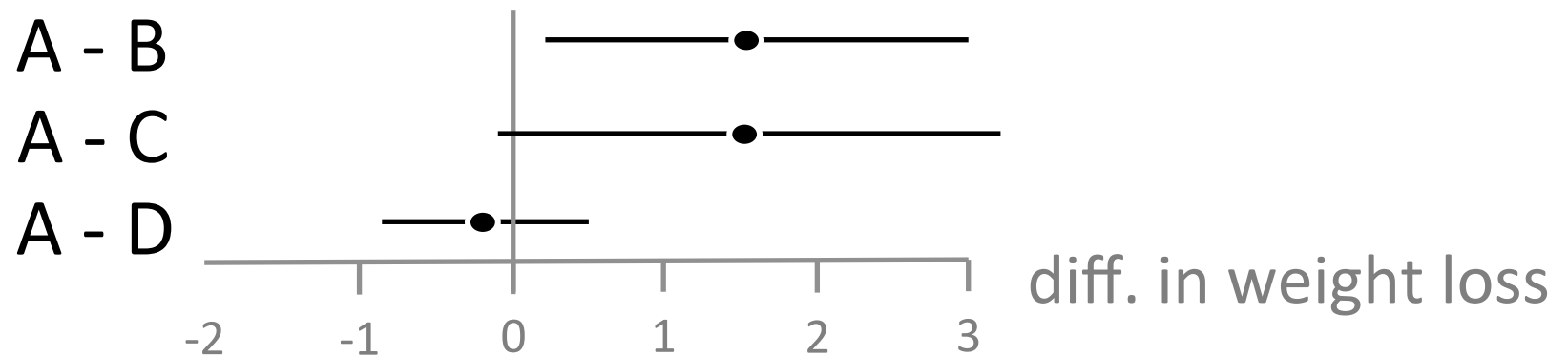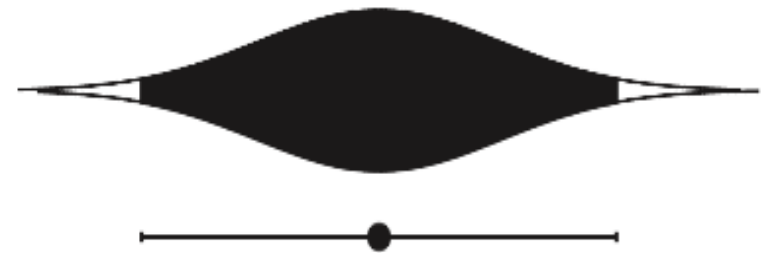
# CONFIDENCE INTERVALS

# CONFIDENCE INTERVALS

- Several interpretations
- « *a range of plausible values for µ. Values outside the CI are relatively implausible.* »
  (Cumming and Finch, 2005)
- Examples of presentation formats:

  2.2m, 95% CI [1.6m, 2.8m]

  2.2m +/- 0.6m

  from 1.6m to 2.8m

# CONFIDENCE INTERVALS

- « *a range of plausible values for μ. Values outside the CI are relatively implausible.* »
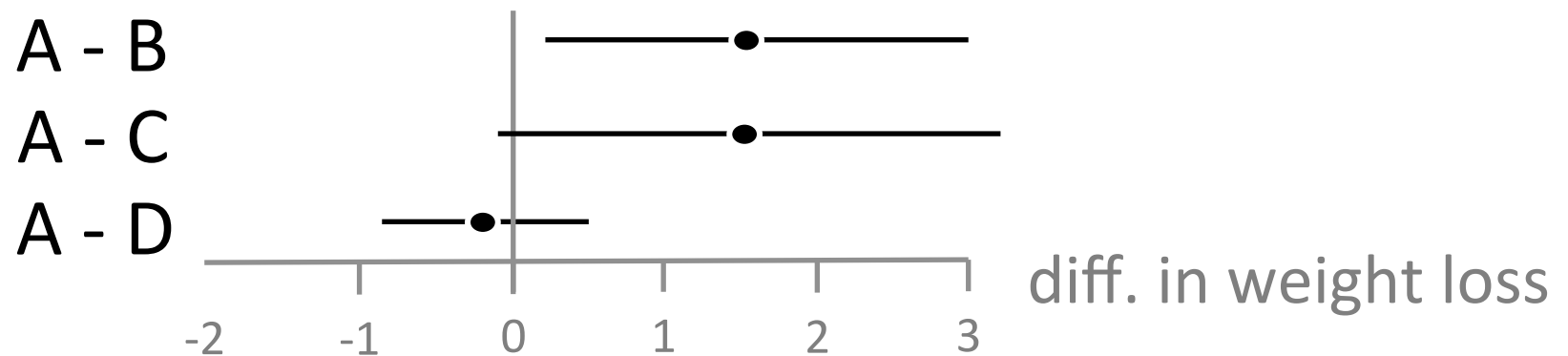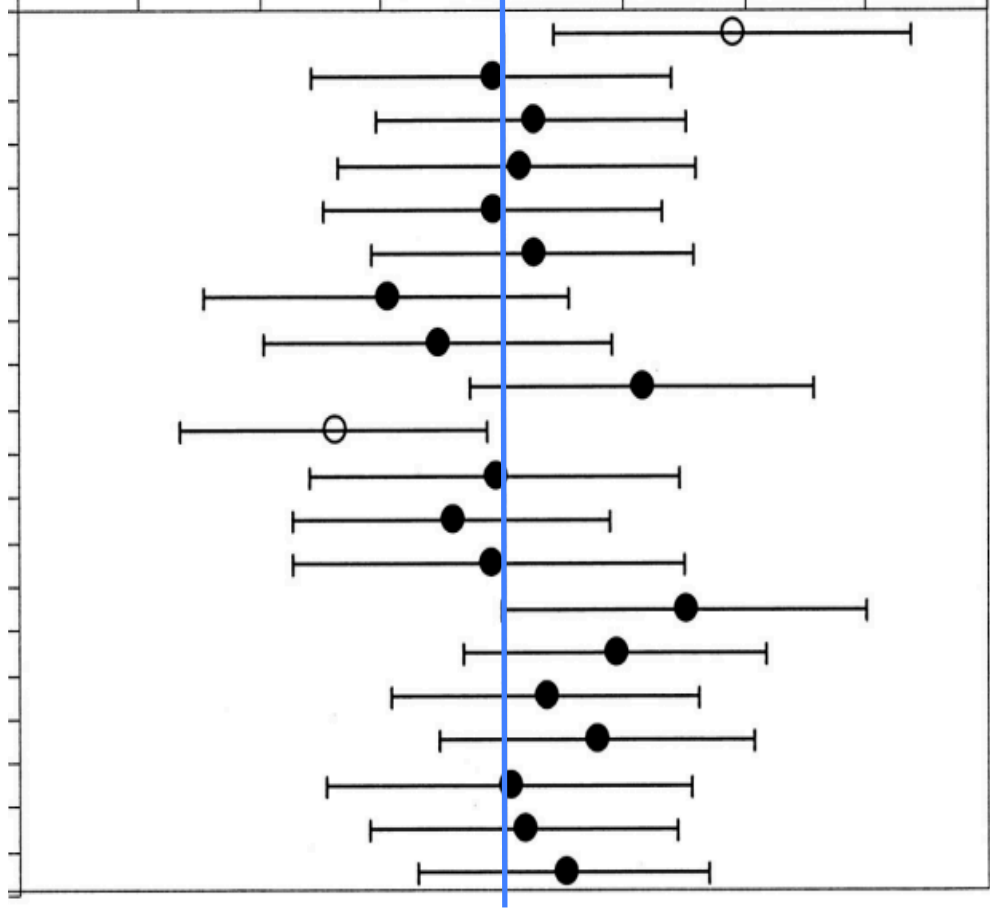  (Cumming and Finch, 2005)

# CONFIDENCE INTERVALS

- « *a range of plausible values for µ. Values outside the CI are relatively implausible.* »
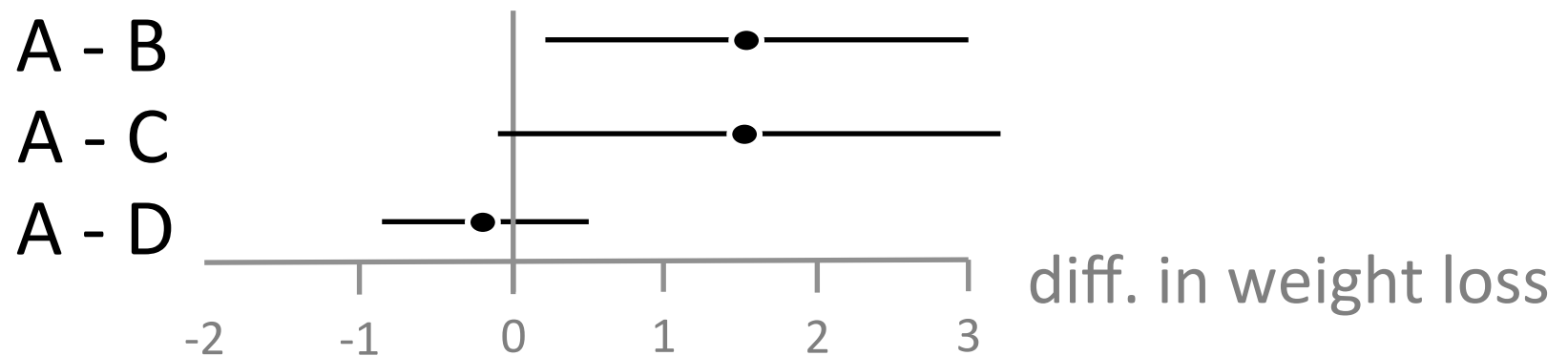  (Cumming and Finch, 2005)

# CONFIDENCE INTERVALS

- « *a range of plausible values for µ. Values outside the CI are relatively implausible.* »
  (Cumming and Finch, 2005)

# CONFIDENCE INTERVALS

- « *a range of plausible values for μ. Values outside the CI are relatively implausible.* »
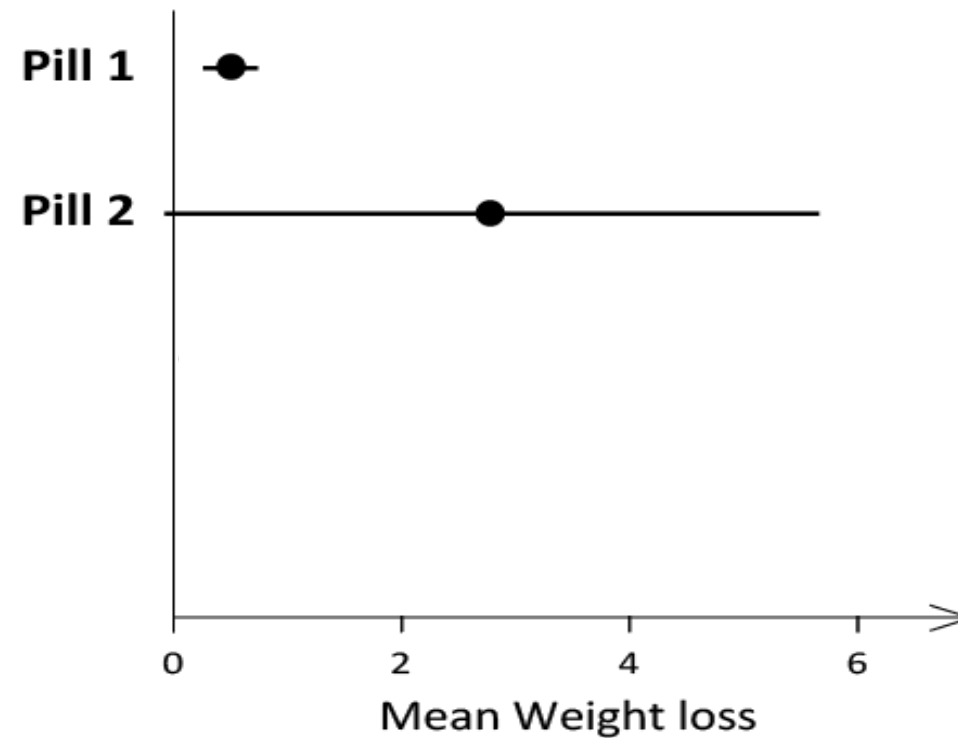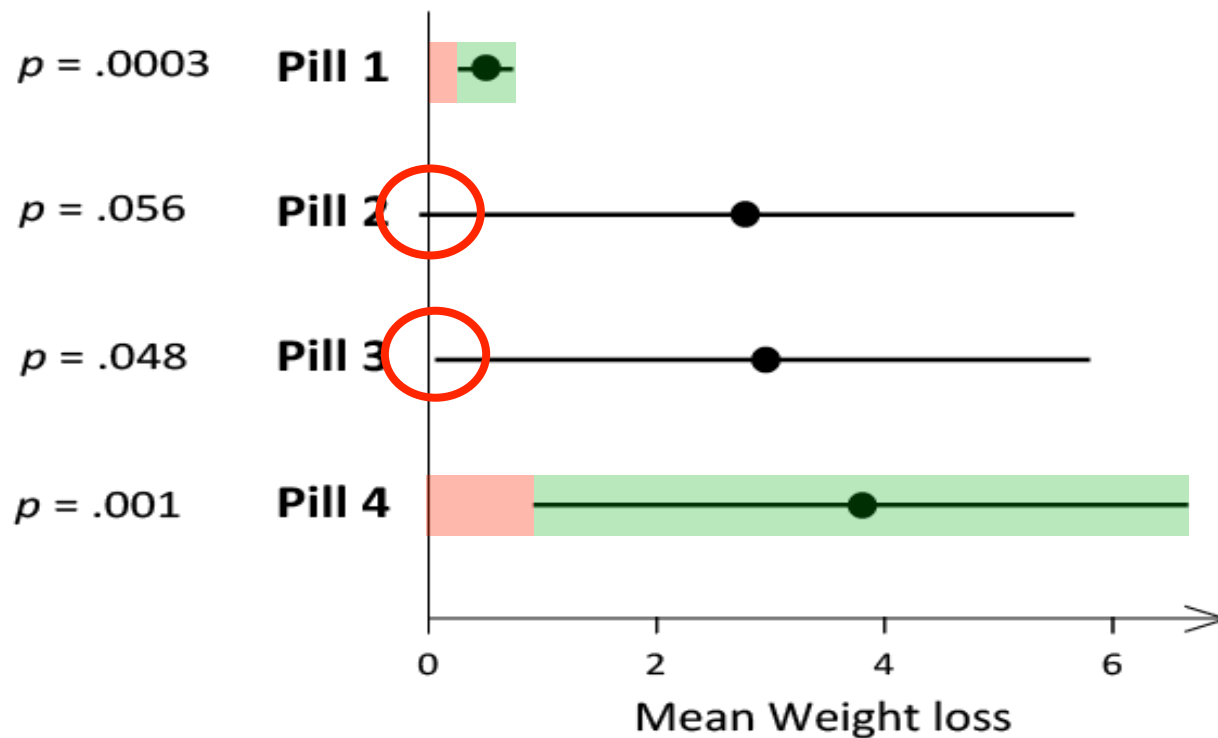  (Cumming and Finch, 2005)

# CONFIDENCE INTERVALS

- "*values close to our M are the best bet for µ, and values closer to the limits of our CI are successively less good bets.*"
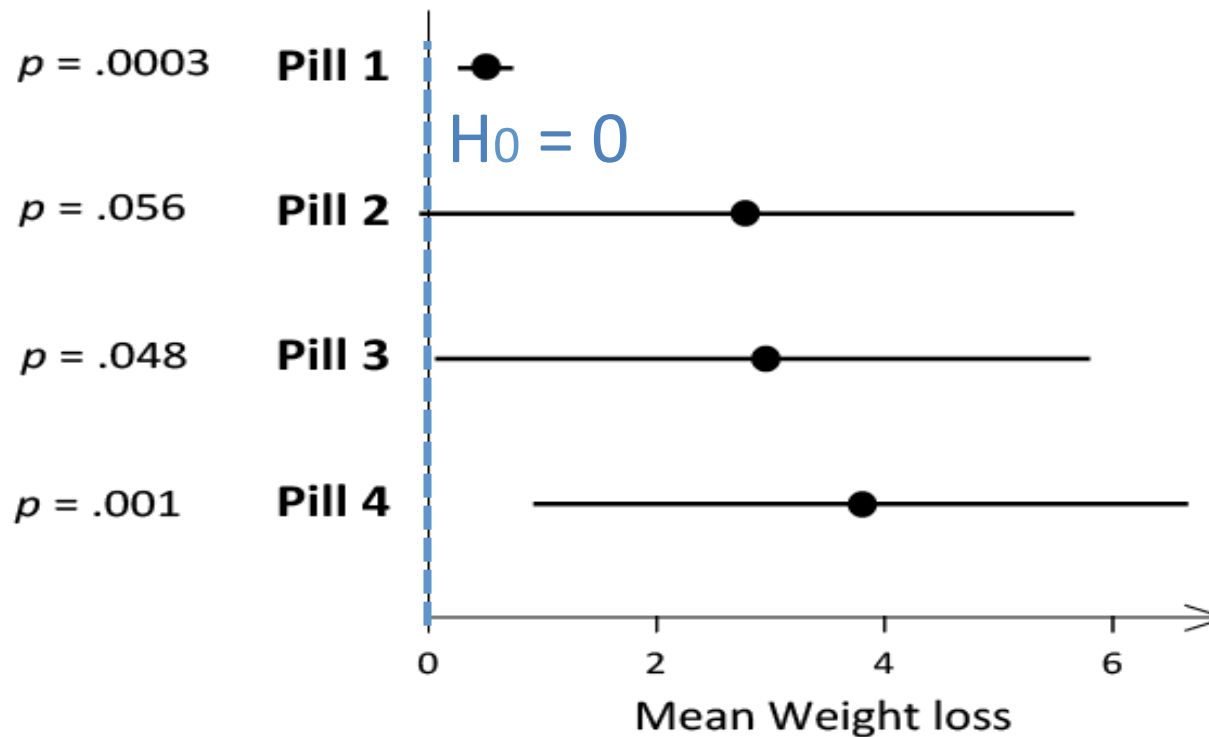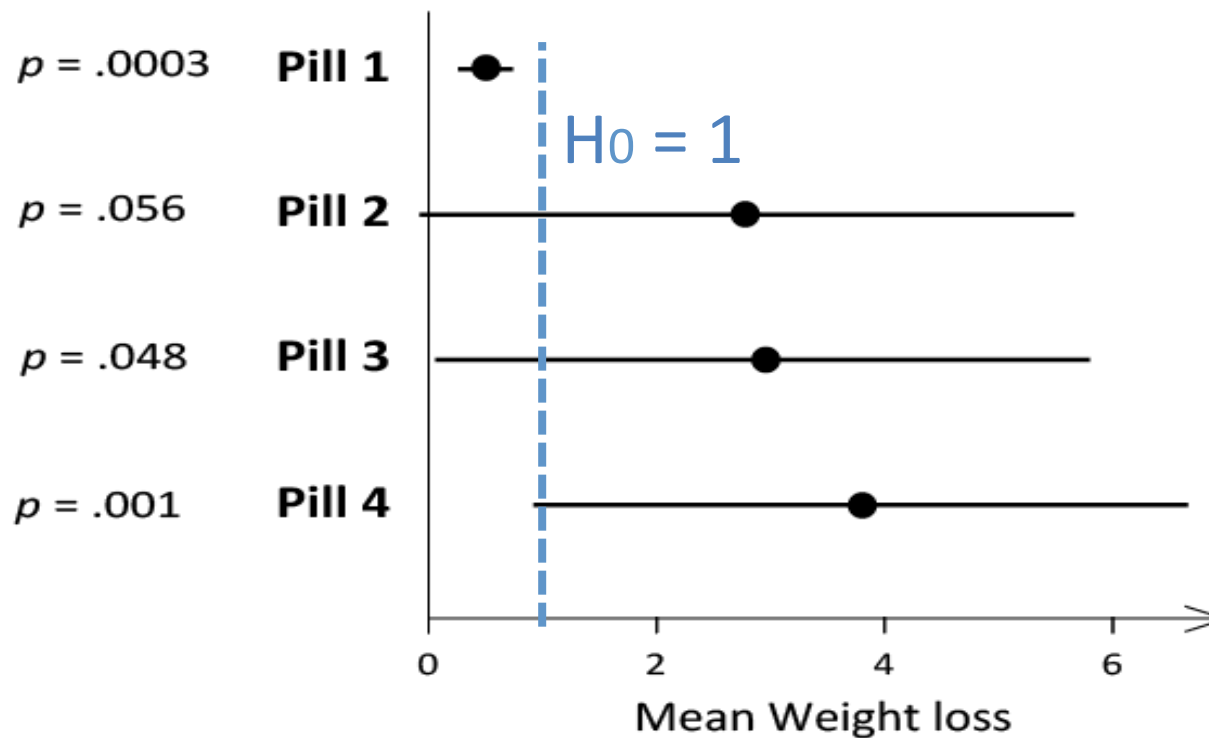
(Cumming, 2013)

# STATISTICAL SIGNIFICANCE



Error bars are 95% CIs
*p*-values are based on a null hypothesis of no effect

# STATISTICAL SIGNIFICANCE



Error bars are 95% CIs
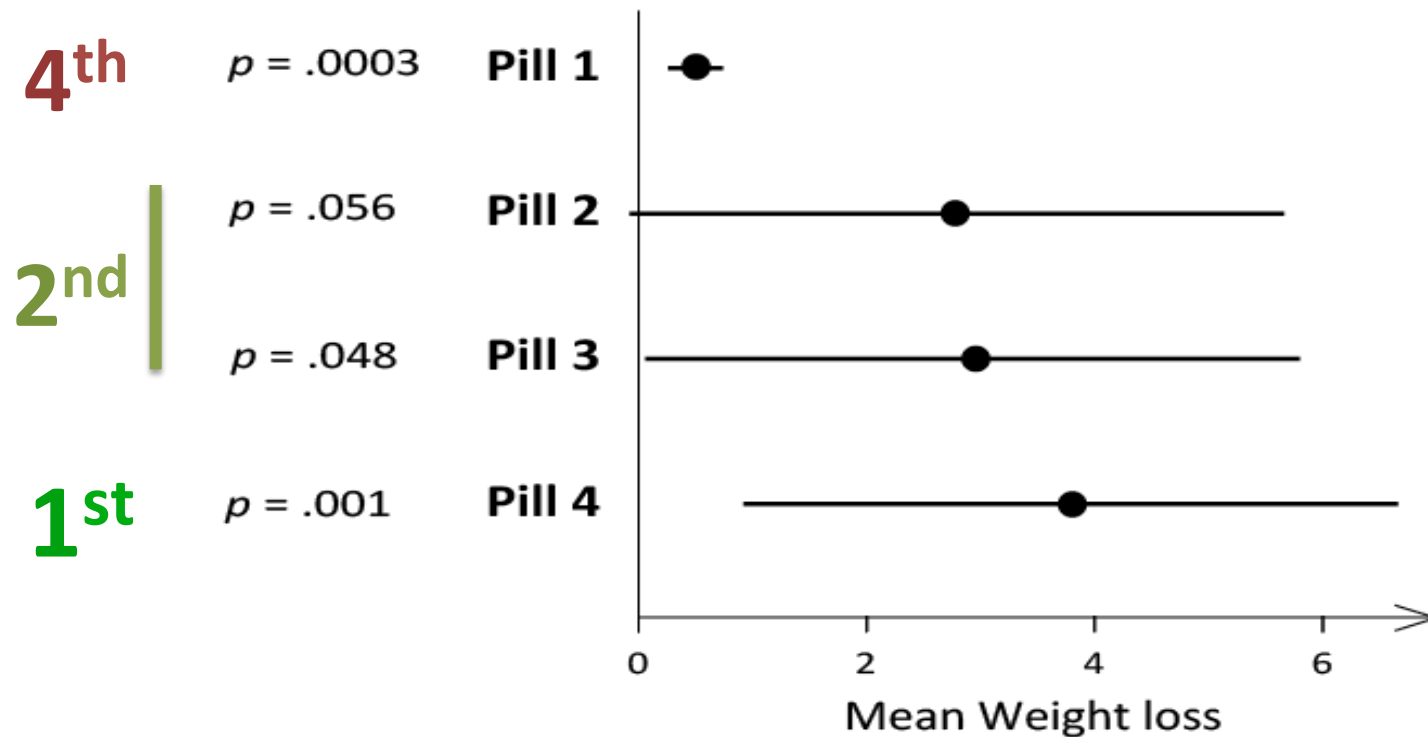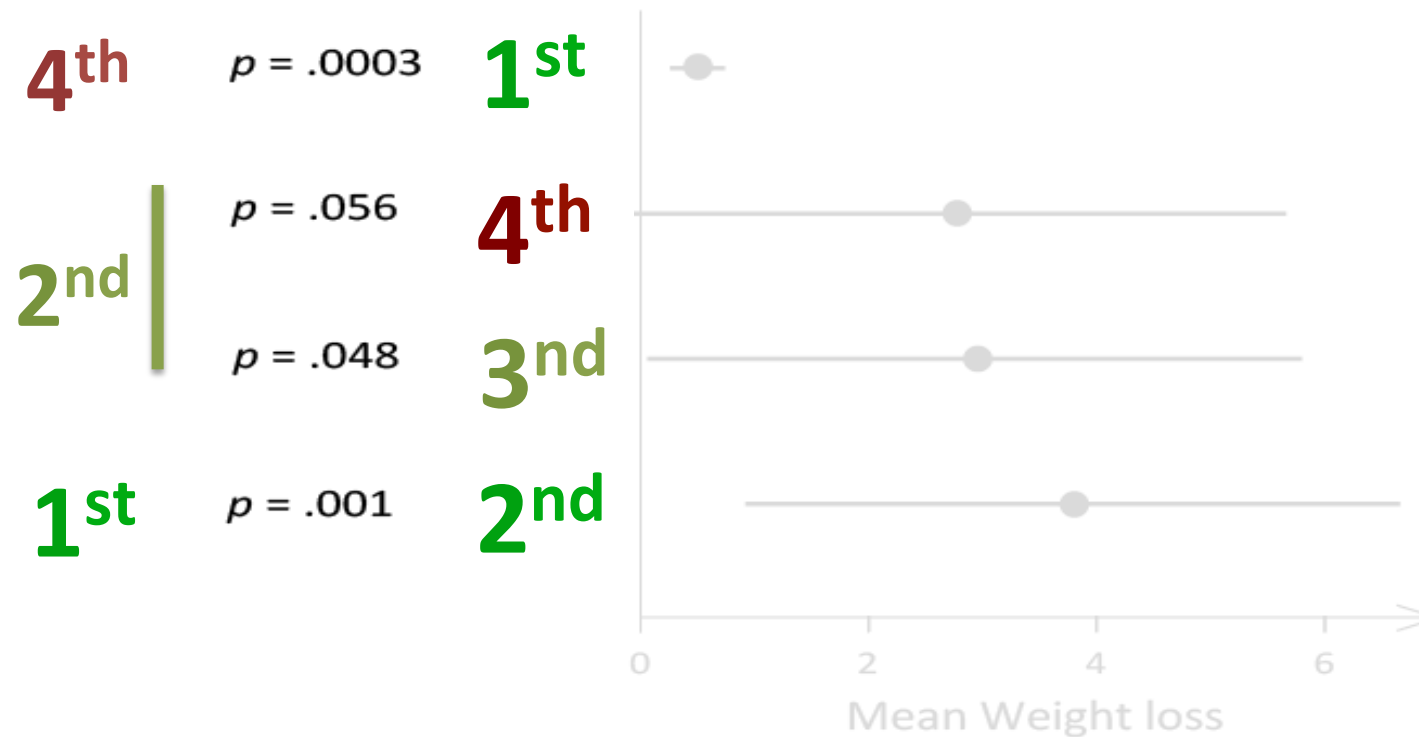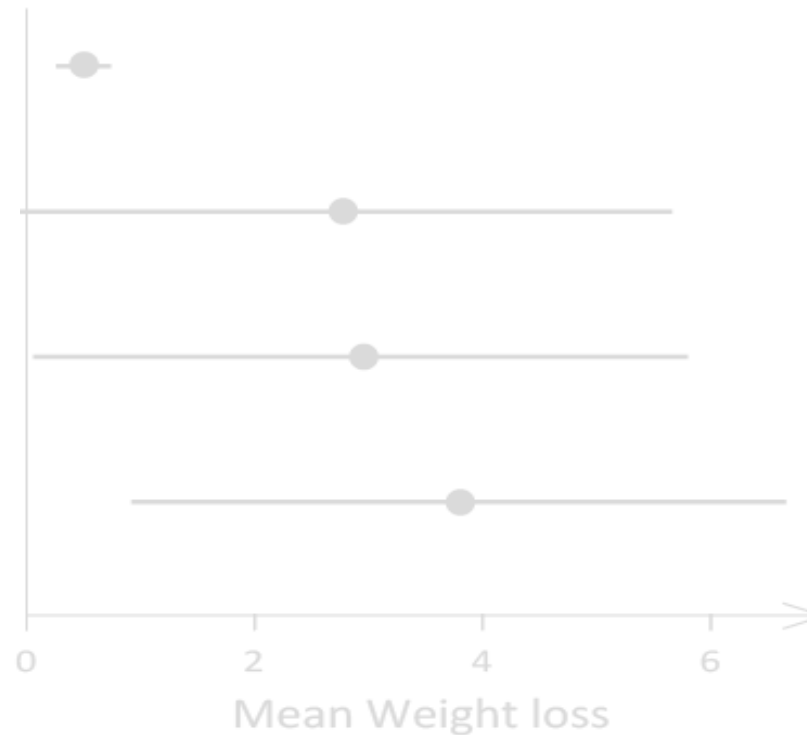p-values are based on a null hypothesis of no effect

STATISTICAL SIGNIFICANCE

4th    p = .0003    Pill 1

2nd    p = .056     Pill 2

        p = .048     Pill 3
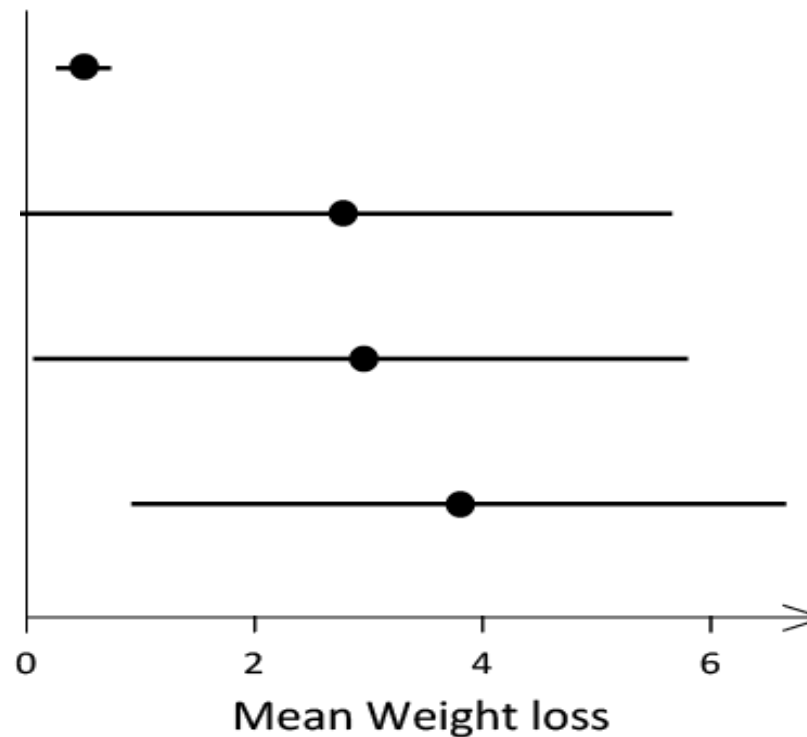
1st    p = .001     Pill 4

Mean Weight loss

Error bars are 95% CIs
p-values are based on a null hypothesis of no effect

# STATISTICAL SIGNIFICANCE

**4th** **Effective**

**2nd** **n.s.**
**(Ineffective?)**
**Effective**

**1st** **Effective**

Mean Weight loss

0     2     4     6

Error bars are 95% CIs
*p*-values are based on a null hypothesis of no effect

Make sure you check _the dance of p-values_ on youtube

# STATISTICAL SIGNIFICANCE

effect of METHOD ($F_{4,44} = 10.1$, $p < 0.0001$
$F_{3,33} = 49.1$, $p < 0.0001$) for both datasets
4) and a significant effect of SCALE for the data
t not for SCALE$\geq$ 4 ($F_{2,22} = 2.7$, $p = 0.0885$).
$= 0.1116$ and $F_{1,11} = 3.9$, $p = 0.0718$).
ractions of METHOD $\times$ W ($F_{12,132} = 6.1$, $p <$
$p < 0.0001$ and $F_{6,66} = 10.6$, $p < 0.0001$) for
SCALE $= 1$ in particular, we have a higher error
his difference vanishes as W increases. The
Mag with other methods. For the remaining
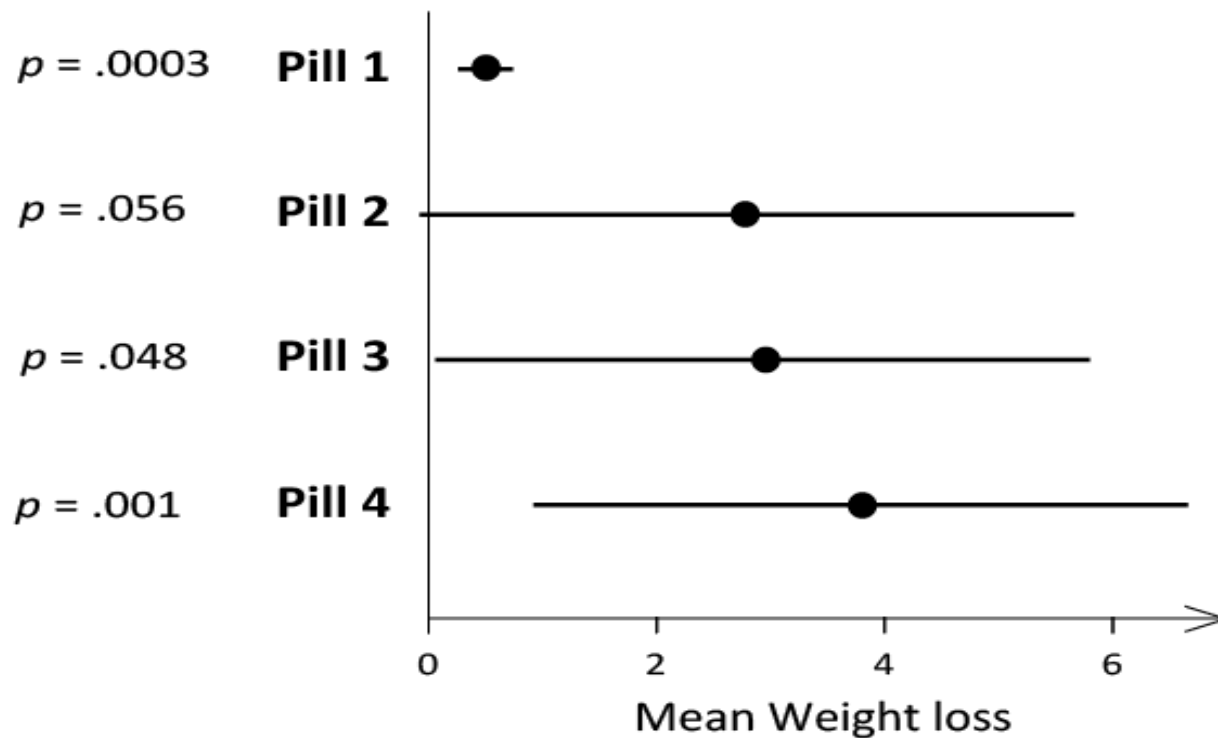s in the error rates.

# STATISTICAL SIGNIFICANCE

*" [NHST] is based upon a fundamental misunderstanding of the nature of rational inference, and is seldom if ever appropriate to the aims of scientific research. "*

Rozeboom (1960)

# STATISTICAL SIGNIFICANCE



Error bars are 95% CIs
*p*-values are based on a null hypothesis of no effect

# STATISTICAL SIGNIFICANCE

" *It seems clear that **no** confidence interval should be interpreted as a a significance test.*"
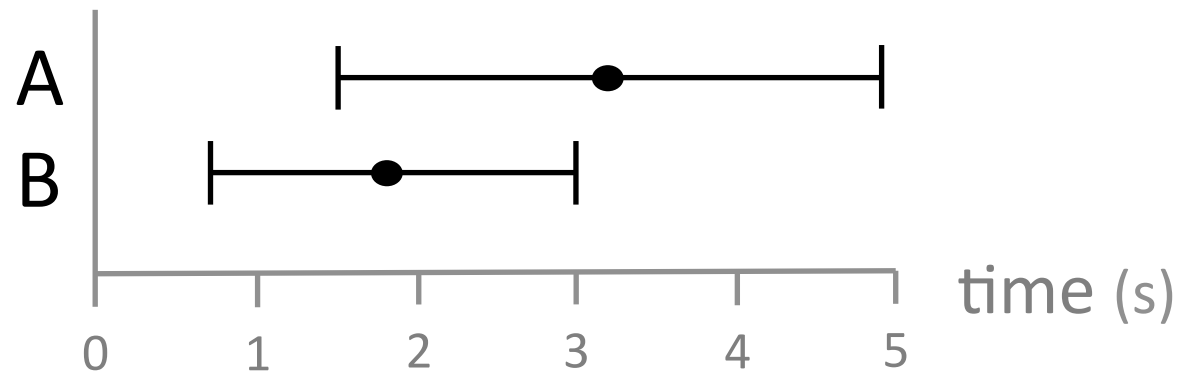
(Schmidt and Hunter, 1997)

## STATISTICAL SIGNIFICANCE

*"It is best for individual researchers to present point estimates and confidence intervals and* **refrain from attempting to draw final conclusions** *about research hypotheses ."*
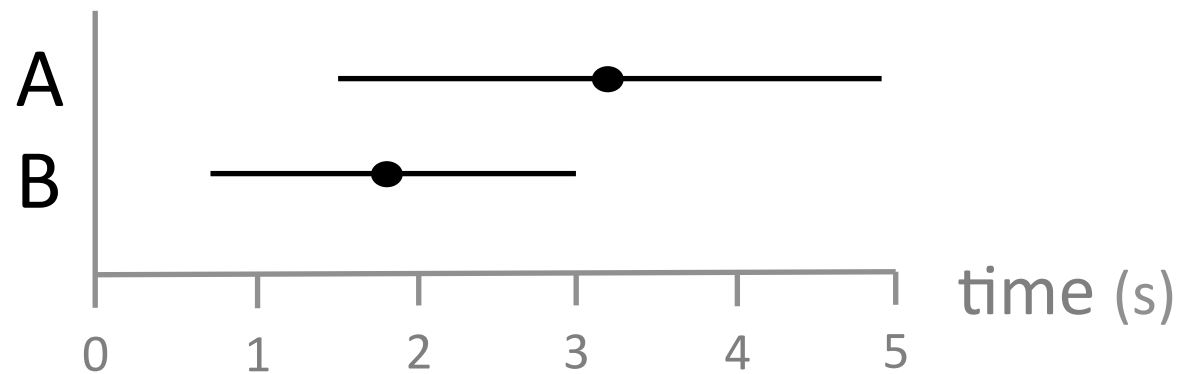
(Schmidt and Hunter, 1997)

# HOW TO GRAPH CIS?

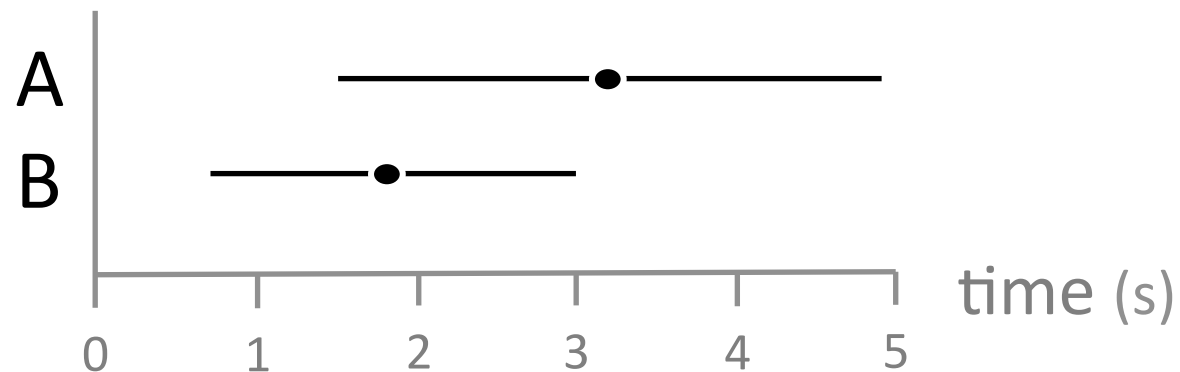- As error bars

# HOW TO GRAPH CIS?
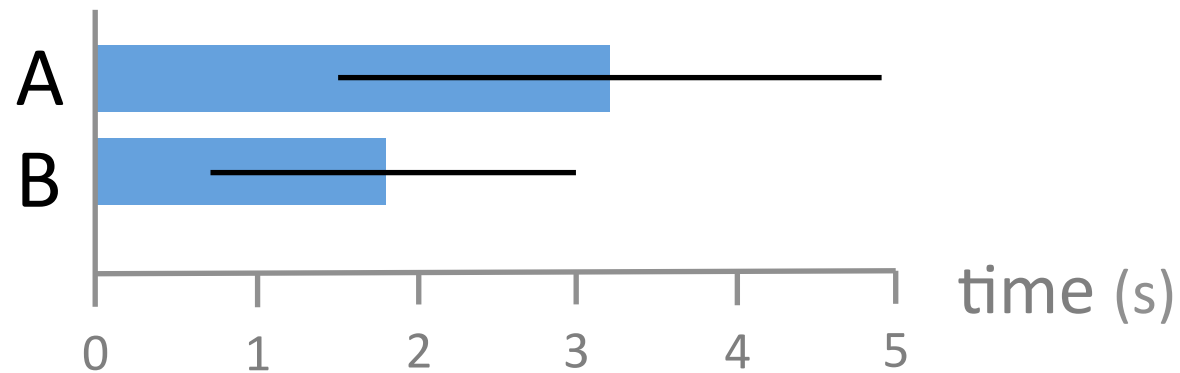
- As error bars
  - Better way:

# HOW TO GRAPH CIS?
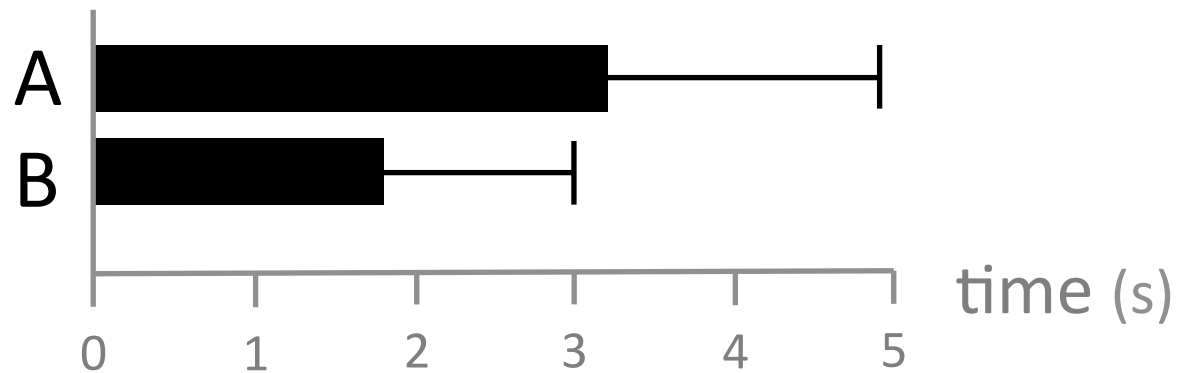
- As error bars
  - Slightly nicer:

# HOW TO GRAPH CIS?

- As error bars
  - With bar charts:
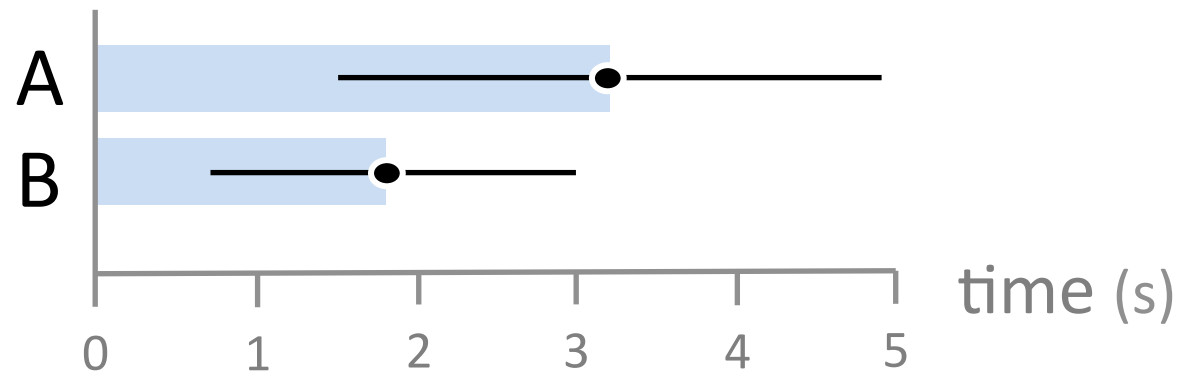
# HOW TO GRAPH CIS?

- ## As error bars
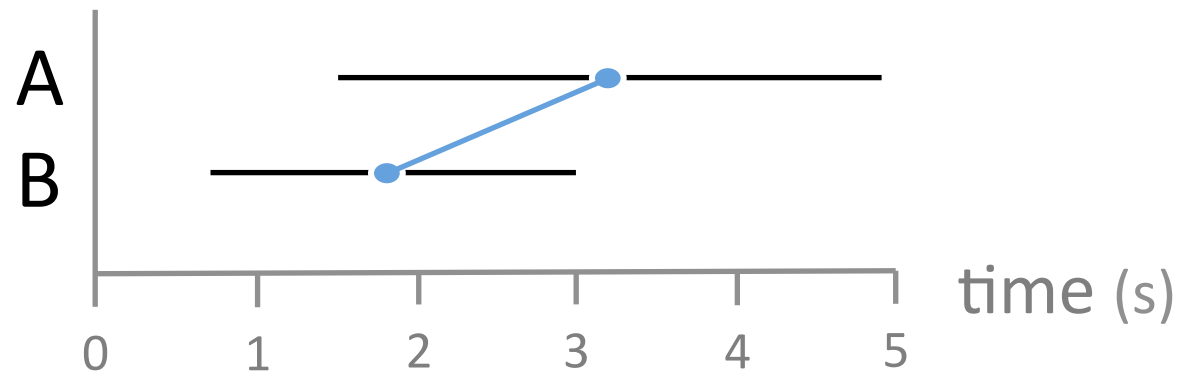  - Dynamite plots:

# HOW TO GRAPH CIS?

- As error bars
  - Perhaps a better approach:

# HOW TO GRAPH CIS?

- As error bars
  - With line charts:
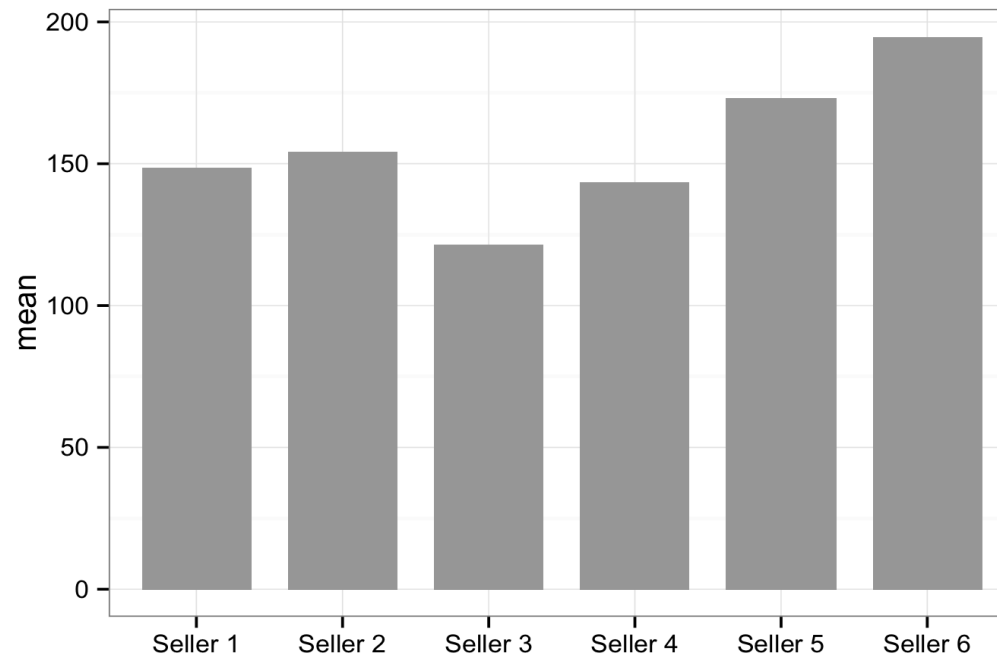
# BACK TO OUR EXAMPLE

- Selling encyclopedias

# Average Sales

| Seller 1 | Seller 2 | Seller 3 | Seller 4 | Seller 5 | Seller 6 |
|----------|----------|----------|----------|----------|----------|
| €149 | €154 | €122 | €143 | €173 | €195 |

| day | Seller 1 | Seller 2 | Seller 3 | Seller 4 | Seller 5 | Seller 6 |
|---|---|---|---|---|---|---|
| 1 | €320 | €80 | €139 | €330 | €133 | €387 |
| 2 | €74 | €60 | €98 | €44 | €182 | €29 |
| 3 | €340 | €67 | €42 | €100 | €51 | €91 |
| 4 | €322 | €54 | €89 | €44 | €67 | €886 |
| 5 | €146 | €195 | €47 | €173 | €49 | €227 |
| 6 | €24 | €288 | €124 | €111 | €730 | €79 |
| 7 | €42 | €249 | €26 | €77 | €672 | €45 |
| 8 | €76 | €67 | €140 | €382 | €195 | €171 |
| 9 | €99 | €312 | €125 | €123 | €43 | €98 |
| 10 | €915 | €77 | €106 | €250 | €149 | €70 |
| 11 | €202 | €504 | €101 | €205 | €682 | €134 |
| 12 | €47 | €167 | €126 | €48 | €93 | €63 |
| 13 | €34 | €65 | €55 | €56 | €333 | €1,157 |
| 14 | €76 | €46 | €89 | €104 | €56 | €470 |
| 15 | €75 | €34 | €184 | €35 | €299 | €205 |
| 16 | €68 | €37 | €275 | €170 | €57 | €192 |