# INTRODUCTION TO STATISTICS

Pierre Dragicevic

# WHAT YOU WILL LEARN

Statistical theory

Applied statistics

This lecture

# GOALS

- Learn basic intuitions and terminology

- Perform basic statistical inference with R

- Focus on high-level aspects

- Accent on estimation rather than hypothesis testing ("the New Statistics")

# ORGANIZATION

- Part I - Elementary notions

- Part II - Tutorial with R

- Part III - Assignments

# A DEFINITION

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

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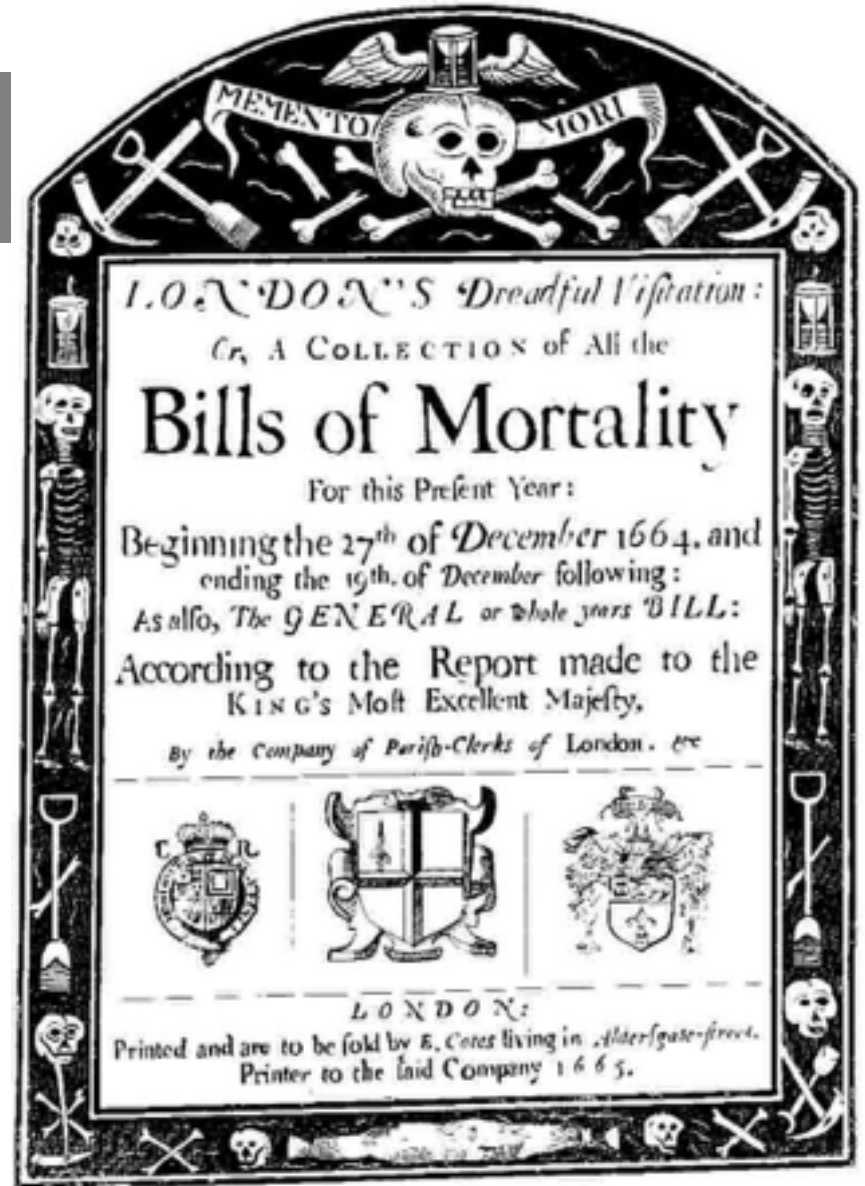
# ORIGINS

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

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

# ORIGINS

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

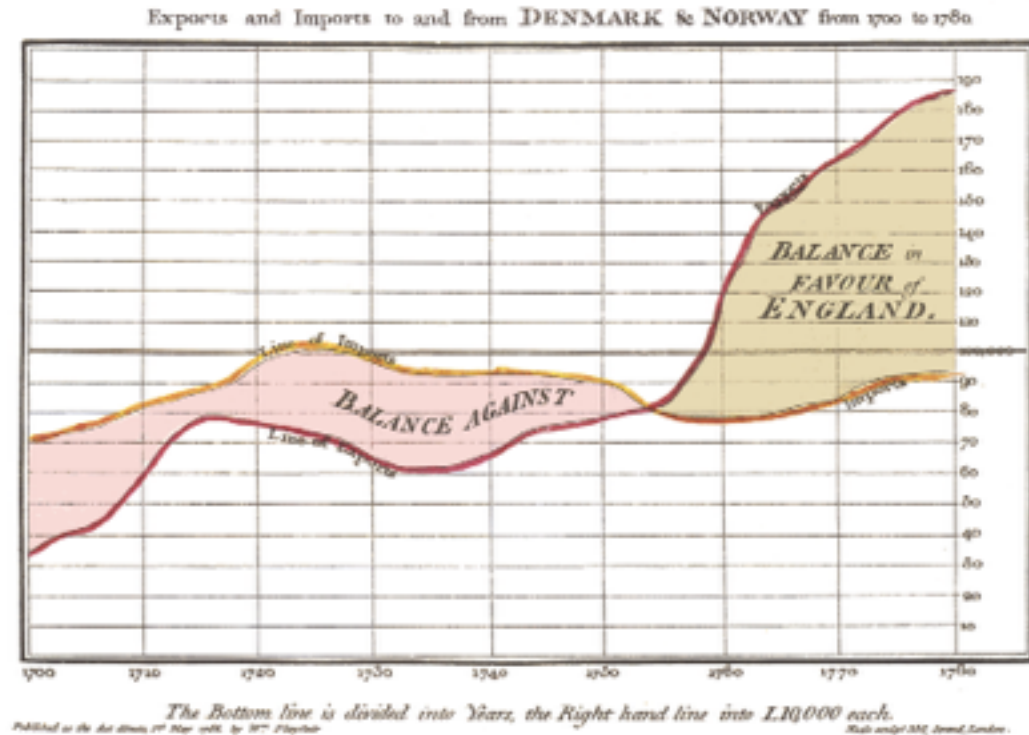| The Years of our Lord | 1647 | 1648 | 1649 | 1650 | 1651 | 1652 | 1653 | 1654 | 1655 | 1656 | 1657 | 1658 | 1659 | 1660 | 1629 | 1630 | 1631 | 1632 | 1633 | 1634 | 1635 | 1636 | 1629–32 | 1633–36 | 1647–50 | 1651–54 | 1655–58 | 1619 1640 1659 | In 20 Years |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 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 Fever | 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 | 105 | 111 | 118 | 86 | 91 | 102 | 113 | 138 | 91 | 67 | 22 | 36 | | 17 | | 35 | 25 | 26 | 75 | 85 | 280 | 421 | 445 | 177 | 1306 |
| Bleach | | | | 1 | 3 | 7 | | | | 1 | | | | | | | | | | | | | | | 4 | 9 | 1 | 1 | 15 |
| Blasted | 4 | 1 | | | 6 | 6 | | 4 | | 5 | 5 | 3 | 7 | 8 | 13 | 8 | 10 | 13 | | 4 | | 4 | 54 | 14 | 5 | 12 | 14 | 16 | 99 |
| Bleeding | 3 | 2 | 5 | 1 | 3 | 4 | 3 | 2 | 7 | 3 | 5 | 4 | 5 | 4 | 5 | 2 | 5 | 4 | | 3 | | | 16 | 7 | 11 | 12 | 19 | 17 | 65 |
| Bloudy Flux, Scouring, and Flux | 155 | 176 | 802 | 289 | 833 | 762 | 200 | 386 | 168 | 368 | 362 | 233 | 346 | 251 | 449 | 418 | 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 | 16 | 19 | 125 |
| Calenture | 1 | | | 1 | | 2 | 1 | 1 | | 3 | | | | | | | | 1 | 3 | | | | 1 | 3 | 4 | 2 | 4 | 3 | 13 |
| Cancer, Gangrene, and Fistula | 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 |
| Wolf | | | 8 | | | | | | | | | | | | | | | | | | | | | 8 | | | | | |
| Canker, Sore-mouth, and Thrush | 66 | 28 | 54 | 42 | 68 | 51 | 73 | 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 | 709 | 838 | 499 | 3364 |
| Chrisomes, and Infants | 1369 | 1254 | 1065 | 990 | 1237 | 1280 | 1050 | 1343 | 1089 | 1393 | 1162 | 1144 | 858 | 1123 | 2590 | 2378 | 2035 | 2258 | 2130 | 2315 | 2113 | 1895 | 9277 | 8453 | 4678 | 4910 | 4788 | 4519 | 32106 |
| Colick, and Wind | 103 | 71 | 85 | 82 | 76 | 102 | 2 | 101 | 85 | 120 | 113 | 179 | 116 | 167 | 48 | 57 | | | 37 | 50 | 105 | 87 | 341 | 359 | 497 | 147 | | 43 | 1389 |
| Cold, and Cough | | | | | | | 41 | 36 | 21 | 58 | 30 | 31 | 33 | 24 | 10 | 58 | 51 | 55 | 45 | 54 | 50 | 57 | 174 | 207 | 60 | 77 | 140 | 43 | 598 |
| Consumpcion, 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 | 8266 | 8999 | 9914 | 12157 | 7197 | 44487 |
| Convulsion | 684 | 491 | 510 | 493 | 569 | 653 | 606 | 828 | 702 | 1027 | 807 | 841 | 742 | 1031 | 52 | 87 | 18 | 21 | 221 | 386 | 418 | 709 | 498 | 1734 | 2198 | 2656 | 3377 | 1324 | 9073 |
| Cramp | | | 1 | | | | | 1 | | 1 | | | | | | | | | | | 1 | | | 1 | | | 1 | | |
| Cut of the Stone | | 2 | 1 | 3 | | 1 | | 2 | 4 | 1 | 3 | 5 | 46 | 48 | | 5 | 1 | | | | 1 | | | | 5 | 4 | 13 | 47 | 38 |
| Dropsy, and Tympany | 185 | 434 | 421 | 508 | 444 | 556 | 617 | 704 | 660 | 706 | 631 | 931 | 646 | 872 | 235 | 252 | 279 | 250 | 266 | 250 | 329 | 385 | 1048 | 1734 | 1535 | 1321 | 2982 | 1302 | 9623 |
| Drowned | 47 | 40 | 30 | 27 | 49 | 50 | 3 | 30 | 43 | 46 | 63 | 60 | 57 | 48 | 43 | 33 | 29 | 14 | 37 | 32 | 32 | 45 | 139 | 147 | 144 | 182 | 215 | 130 | 827 |
| Excessive drinking | | | 2 | | | | | | | | | | | | | | | | | | | | | | | 2 | | 2 | |
| Executed | 8 | 17 | 29 | 43 | 24 | 12 | 19 | 21 | 19 | 22 | | 18 | 7 | 18 | 19 | 13 | 12 | 13 | 13 | 13 | 13 | 13 | 62 | 52 | 97 | 76 | 79 | 55 | 384 |
| Fainted in a Bath | | | | | | | | | | 1 | | | | | | | | | | | | 1 | | | | | | | |
| Falling-Sickness | 3 | 2 | 2 | 3 | | 3 | 1 | | 1 | | 4 | 5 | 5 | 3 | | 10 | 7 | 7 | 5 | 6 | 8 | | 27 | 21 | 10 | 8 | 8 | 9 | 74 |
| Flox, and small Pox | 139 | 400 | 1190 | 184 | 525 | 1272 | 119 | 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 | 6 | 6 | 18 | 33 | 26 | 6 | 13 | 8 | 24 | 24 | 83 | 69 | 29 | 34 | 27 | 29 | 243 |
| French-Pox | 18 | 29 | 15 | 18 | 21 | 20 | 20 | 20 | 29 | 23 | 25 | 53 | 51 | 31 | 17 | 12 | 12 | 12 | 7 | 17 | 12 | 22 | 53 | 45 | 80 | 81 | 130 | 83 | 392 |
| Frighted | 4 | 4 | 1 | | 3 | | 2 | | 1 | | | 9 | | 2 | 1 | | | | | | 3 | 2 | 3 | 9 | 5 | 2 | | 2 | 11 |
| Gout | 9 | 5 | 11 | 9 | 7 | 7 | 5 | 7 | 8 | 13 | 14 | 2 | 2 | 2 | 3 | 1 | 5 | 7 | 8 | 14 | 24 | 35 | 25 | 36 | 28 | 134 |
| Grief | 12 | 13 | 16 | 7 | 17 | 14 | 11 | 17 | 10 | 13 | 24 | 18 | 11 | 36 | 2 | 8 | 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 | 5 | 15 | | 3 | 8 | 7 | 37 | 18 | 48 | 47 | 72 | 32 |
| 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 | | | | | 2 | 2 | | | | 10 | 16 | 13 | 8 | 10 | 8 | 11 | 47 | 35 | 02 | 5 | 6 | 10 | 11 |
| Impostume | 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 | | | | | | | | | | | | | | | 10 | | | | | | 00 | 10 | 01 | | | |
| Killed by several Accidents | 27 | 57 | 39 | 94 | 47 | 45 | 57 | 58 | 52 | 43 | 52 | 47 | 55 | 47 | 54 | 55 | 47 | 46 | 41 | 51 | 60 | 69 | 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 | 20 | 26 | 69 | 97 | 150 | 94 | 94 | 103 | 66 | 533 |
| Lethargy | 3 | 2 | 4 | 3 | 10 | 2 | 4 | 3 | 2 | 4 | 6 | 4 | | 2 | 6 | 3 | 2 | 2 | 2 | 3 | 11 | 5 | 7 | 11 | 11 | 9 | 77 |
| Leprosy | | | 1 | | | | | | | | | | | | | | | | | | | | 2 | 1 | 1 | | | |
| Livergrown, Spleen, and Rickets | 53 | 46 | 50 | 59 | 65 | 72 | 6 | 65 | 52 | 50 | 38 | 51 | 9 | 15 | 94 | 112 | 99 | 87 | 82 | 77 | 98 | 99 | 392 | 356 | 213 | 269 | 191 | 158 | 1178 |

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Cancer, Gangrene, and Fistula Wolf | 26 | 29 | 31 | 19 8 | 31 | 53 | 36 | 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 1 |
| Colick, and Wind | 103 | 71 | 85 | 82 | 76 | 102 | 8 | 101 | 85 | 120 |
| Cold, and Cough | | | | | | | 41 | 36 | 21 | 58 |
| Consumption, and Cough | 2423 | 2200 | 2388 | 1938 | 2350 | 2410 | 2216 | 2868 | 2606 | 3184 2 |
| Convulsion | 684 | 491 | 530 | 493 | 569 | 653 | 666 | 818 | 702 | 1027 |
| Cramp | | | 1 | | | | | | | |
| 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 | 49 |
| 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 | 1190 | 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 | |
| Gout | 9 | 5 | 11 | 9 | 7 | 7 | 5 | 6 | 8 | 7 |
| Grief | 12 | 13 | 16 | 7 | 17 | 14 | 11 | 17 | 10 | 13 |

# ORIGINS

- 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

# ORIGINS

- 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
  - ...

# 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.

BALANCE AGAINST

The Bottom line is divided into Years, the Right hand line into £10,000 each.

# STATS & VISUALIZATION

- Exploratory Data Analysis
  - Tukey, 1977

# Box-and-whisker plots with end values identified

## A) HEIGHTS of 50 STATES

Alaska ○

California
Colorado ○ ● ● Washington
Hawaii ○ ● ● Wyoming

Delaware ● ○ Louisiana
○ Florida

## B) HEIGHTS of 219 VOLCANOS

Height
(feet)
↑

20,000

Guallatiri
Lascar ● ● Cotapaxl
Kilimanjaro ● ● Misti
Tupungatito

15,000

10,000

5,000

0

Ilha Nova ○ ● Anak Krakatau

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

- # Statistical Graphics

  - – AT&T Bell Labs Video, 1985

Baby Name > cr ✕

◉ Both ◯ Boys ◯ Girls

Names starting with 'CR' per million babies

5,000

Cristina

4,000

3,000

Craig          Crystal

2,000

Cristian

1,000

Cruz

1880s  1890s  1900s  1910s  1920s  1930s  1940s  1950s  1960s  1970s  1980s  1990s  2000s  2009

**Last letter of boys' names in 1950** — *Percentage of boys born* vs last letter (a–z)

**Last letter of boys' names in 2010** — *Percentage of boys born* vs last letter (a–z)

| 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

Regent's Park

Cumberland

River Thames

Scale: one-half mile

German bombings in London during WWII

# STATS & VISUALIZATION

- **Confirmatory data analysis**

  - For answering questions rigorously

  - Example: is this new drug effective?

  - Strong focus on automatic procedures, computation and objectivity

  - Looking at data can impair objectivity:

    - Cherry picking, snooping, fishing, data mining

German bombings in London during WWII

# 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

# WHAT ARE STATS?

- A set of tools and methods

- With an old tradition:

  - Origins in demographics

  - Anchored in mathematics & probability theory

  - Visual representations play a role

  - A generally strong focus on (computationally cheap) numerical calculations

# WHAT ARE STATS?

- Good for:

    - Summarizing data for presentation

    - Answering questions rigorously

    - Making predictions

    - Making rational, evidence-based decisions

    - A long accumulated experience!

# STATISTICAL TOOLS

# STATISTICAL TOOLS

## DESCRIPTIVE STATISTICS

# AN EXAMPLE

- Selling encyclopedias

Robert  Steve  Paul  Roger  Geoffrey  Dan

| 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

- When are the mean and the median equal? When do they differ?

# 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

What is the best measure of central tendency?



Income

# DISPERSION

## Standard Deviation

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



Image from Wikipedia

# DEPENDENCE

- Correlation



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 |

How much can we trust this chart?

# LET US TRAVEL TO THE FUTURE

# BACK TO THE PRESENT

| 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 |

How much can we trust this chart?

# STATISTICAL TOOLS

## INFERENTIAL STATISTICS

# STATISTICAL INFERENCE



Sample

$\overline{x}$  Statistic

(Sample mean)

# STATISTICAL INFERENCE



We want to know about these

We have these to work with

Random selection

Population

Sample

Parameter $\mu$

(Population mean)

Inference

$\bar{x}$  Statistic

(Sample mean)

# STATISTICAL INFERENCE

- Terminology:

  - **Sample** vs. **population**

  - Mean, median, standard deviation, correlation, etc:
    - A sample **statistic**
    - A population **parameter**

# STATISTICAL INFERENCE

- ## Unit of statistical analysis



We want to know about these

We have these to work with

Random selection

Population

Sample

*= "the thing that I'm sampling from a larger population"*

# STATISTICAL INFERENCE

- Unit of statistical analysis

| 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 |

# STATISTICAL INFERENCE

- Unit of statistical analysis

| day | Seller 1 |
| --- | --- |
| 1 | €320 |
| 2 | €74 |
| 3 | €340 |
| 4 | €322 |
| 5 | €146 |
| 6 | €24 |
| 7 | €42 |
| 8 | €76 |
| 9 | €99 |
| 10 | €915 |

# STATISTICAL INFERENCE

- Unit of statistical analysis

| 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 |

# STATISTICAL INFERENCE

- Unit of statistical analysis

### Average Sales

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

# STATISTICAL INFERENCE

- Unit of statistical analysis

| 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 |

# SAMPLING DISTRIBUTION

# SAMPLING DISTRIBUTION

- *"The **sampling distribution** of a statistic is the distribution of that statistic, considered as a random variable, when derived from a random sample of size n."*
*[…]*
*"It may be considered as the **distribution of the statistic for all possible samples** from the same population of a given size"*

# SAMPLING DISTRIBUTION



Diastolic Blood Pressure?

Mean = 78 mm Hg

Samples

Mean = 75

Mean = 67

Mean = 71.3

# SAMPLING DISTRIBUTION

- ## Demo http://onlinestatbook.com/stat_sim/sampling_dist/

# SAMPLING DISTRIBUTION



eans, N=5

# SAMPLING DISTRIBUTION



Standard error

95% confidence interval

# SAMPLING DISTRIBUTION

# SAMPLING DISTRIBUTION

- Resampling techniques
  - Bootstrapping

Bootstraps

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

From Germain Salvato-Vallverdu

Complete element space

initial sample with N elements

$N_b$ bootstrap samples

1
2
3
4
$\vdots$
$N_b$

From Germain Salvato-Vallverdu

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 DISTRIBUTION

- Bootstrapping video


Confidence Intervals Using Bootstrapping

# SAMPLING DISTRIBUTION

- 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

- Abraham De Moivre
  1667 - 1754

# MORE HISTORY



**Number of individuals** (y-axis)

**Height in inches** (x-axis)

# 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

# CONFIDENCE INTERVALS

# CONFIDENCE INTERVALS



95% confidence interval

# CONFIDENCE INTERVALS



95% confidence interval
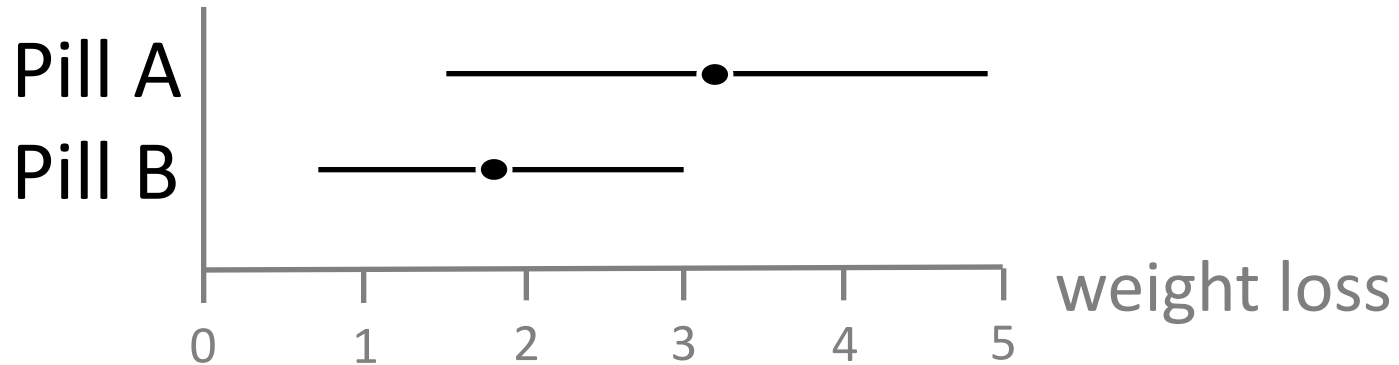
μ

Different random samples

tinyurl.com/danceptrial2

# 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:

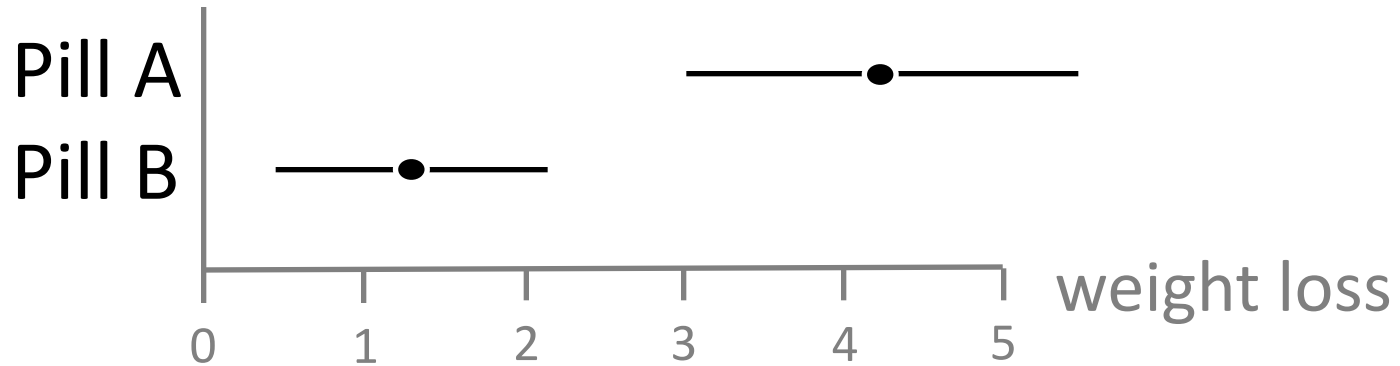  2.2m, 95% CI [1.6m, 2.8m]
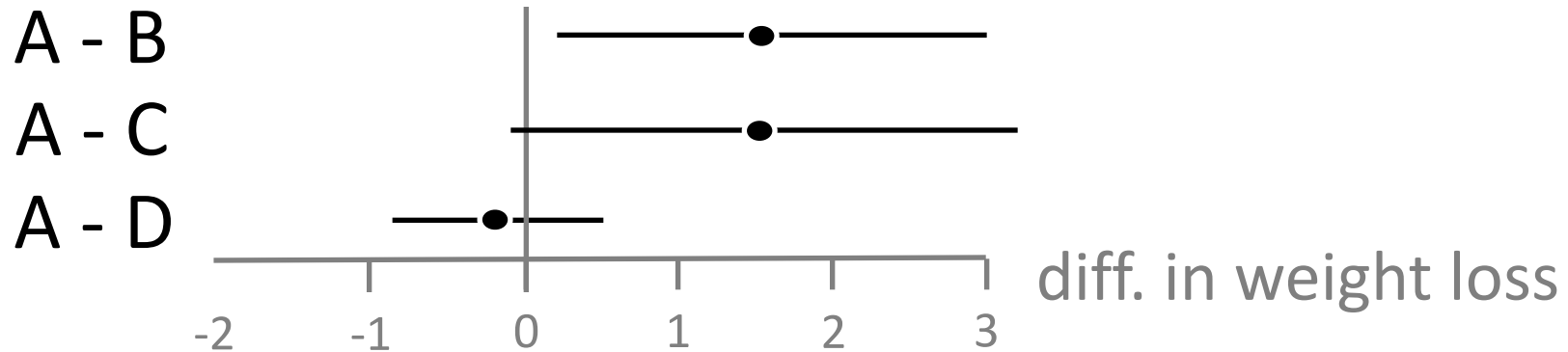
  2.2m +/- 0.6m

  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

- *"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."*
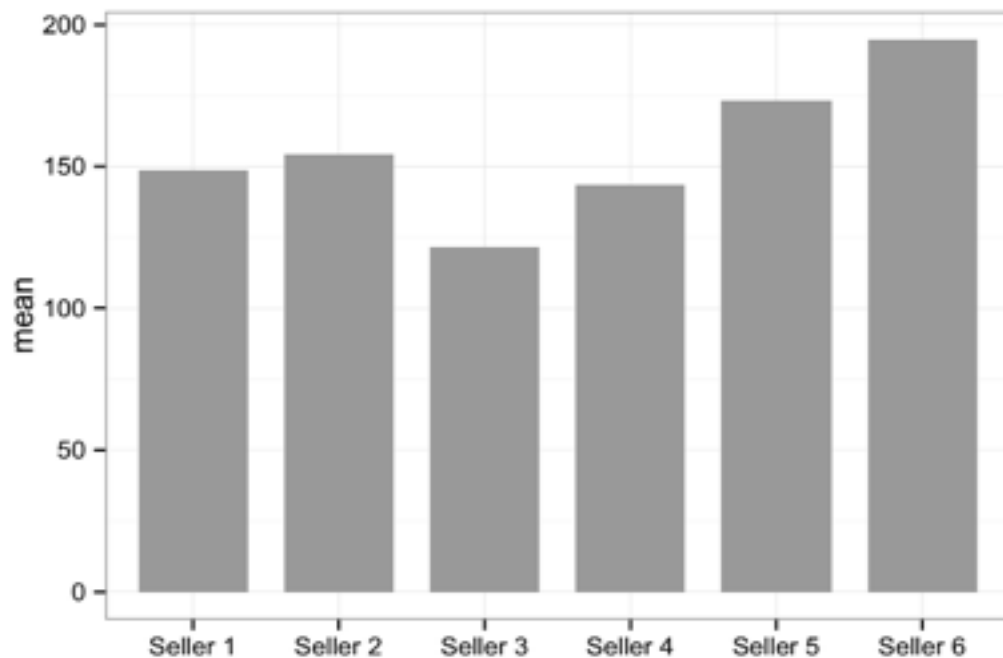
(Cumming, 2013)

# 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 |

http://tinyurl.com/stats-va2015