

Introduction to data science & artificial intelligence (INF7100)

Arthur Charpentier

#201 Data Science or Statistics ?

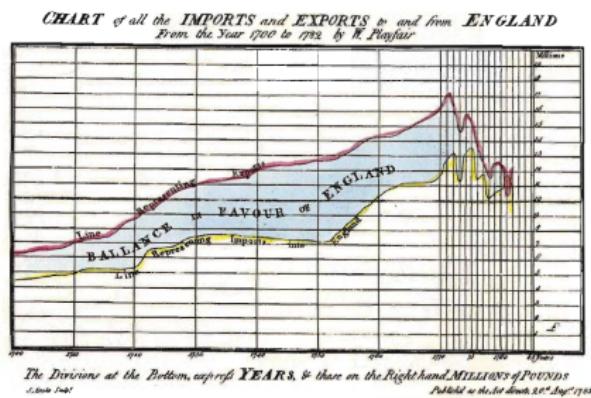
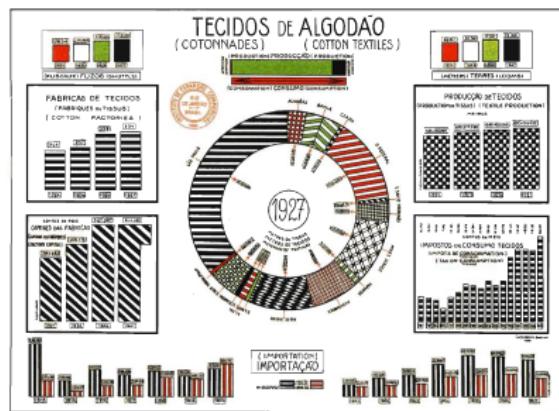
été 2020

Statistics ?

The image shows the front page of The New York Times from May 9, 2020. The top half features a large chart titled "U.S. UNEMPLOYMENT IS WORST SINCE DEPRESSION" showing the percentage of unemployed workers from 1930 to 2020. Below the chart is a headline "Georgia Killing Puts Spotlight on a Police Force's Troubled History". The middle section contains several news stories with headlines like "In Flynn Case, Justice Inquiry Is Barr's Target", "Mexico's Fall Ignores Reality, Mexico's Hospitals Are Overrun", and "A Good Walk, Unfinished". The bottom section includes a "Deaths in America" graphic and a "Deaths in the U.S." chart. The right side has a sidebar for "Late Edition" and a "Deaths in the U.S." box.

Statistics ?

German **Statistik**, introduced by Gottfried Achenwall (1749), originally designated the analysis of data about the state, signifying the "**science of state**" (then called political arithmetic).



Statistik was data to be used by governmental and (often centralized) administrative bodies

Statistics ?

Age.	Per- sons. curr.	Age.	Per- sons.								
Age.	Per- sons. curr.	Age.	Per- sons.								
1	1000	8	580	15	623	22	585	29	539	35	481
2	855	9	570	16	612	23	579	30	531	37	472
3	798	10	651	17	615	24	573	31	523	38	453
4	750	11	553	18	610	25	567	32	515	39	454
5	732	12	546	19	604	26	560	33	507	40	445
6	710	13	640	20	598	27	553	34	495	41	436
7	652	14	634	21	592	28	546	35	490	42	427
Age.	Per- sons. curr.	Age.	Per- sons.								
Age.	Per- sons. curr.	Age.	Per- sons.								
43	417	50	345	57	272	64	262	71	131	78	58
44	407	51	335	58	262	65	192	72	120	79	49
45	397	52	324	59	252	65	182	73	105	80	41
46	387	53	313	60	242	67	172	74	58	81	34
47	377	54	302	61	232	68	162	75	88	82	28
48	367	55	292	62	222	69	152	76	78	83	23
49	357	56	282	63	212	70	142	77	68	84	20
											Sum Total.

Graunt's' (1662) Observations on the Bills of Mortality

Statistics ?

Graunt concluded that the plague was caused by person-to-person infection rather than the competing theory of “infectious air” based on the pattern of infections through time, “thought as we think today [...] they reasoned about their data” *Where Shall the History of Statistics Begin?*

“People use statistics as the drunken man uses lamp posts - for support rather than illumination”,
Andrew Lang (*or not*)



Statistics ? Data ?

Big Data context Data, data everywhere

(source Edward Tufte)

Agenda

Data are x_i 's, or $\mathbf{x}_i = (x_{1,i}, \dots, x_{k,i})$'s.

- ▶ functions, c.d.f. $F(x) = \mathbb{P}[X \leq x]$ and density (histogram)
- ▶ statistical indicators
 - ▶ central value(s): average(s), median
 - ▶ dispersion: variance, standard deviation, inequalities
 - ▶ approximations
 - ▶ quantiles
- ▶ inference: from frequentist to Bayesian
- ▶ testing, significance, p -value
- ▶ bivariate analysis (x_i, y_i)
- ▶ multivariate analysis $(\mathbf{x}_i = (x_{1,i}, \dots, x_{k,i}))$
 - ▶ projection & dimension reduction
 - ▶ cluster analysis
- ▶ networks
- ▶ time series