Introduction to data science & artificial intelligence (INF7100)

Arthur Charpentier

#221 Statistical Inference: Average

été 2020

Statistical Indicators

						_					_								_				16101	14633	1647	1671	1600		In 20
						TI	Чε	2	ΓA	BL	E	01	7 (CA	5 U	A L	T	lE	S.				1630	1634	1648	1652	1656		Years.
The Years of our Lord	1647	1648	1649	1650	1651	1652	1653	1654	1655	656	1657	1658	1659	1660	1629	630	1631	1632	1633	1634	1635	1636	1632	1636	1649 1650	1654	1658	1659	rears,
Abortive, and (b)born	335			351					483	419 801	463	467	421	544	199	439	410 661	445 671	500 704	475 623	507 794	523	1793	2005	1342 3336	1587	1832	1247	8559 15757
Aged Ague, and Fever	1260			970	780 1018	1212	1182		689	875	999	1800	2303	2148	956	712			953	1279	1622		4418	6235	3865	4903	4363	4010	23784
Apoplex, and fodulity Riesch	68	74	64	74	106	111	118	86	92	102	113	138	91	67	22	36	- 1	17	24	35	26		75	85	280	421	445	177	1306 15
Blaffed	4	٠.	1 1	,	6	8			4	'	5	3	3	8	13	8	10	13	6	4		4	54	14	5	12	14	16	99
Bleeding	. 3	176	801	189	811	762	200	386	168	368	362	233	7 346	251	449	438	352	4 348	278	S12	346	330	1587	1466	1422	12 2181	1161	1597	65 7818
Bloudy Flux, Scouring, and Flux Rurat, and Scalded	155	1/6		209	033	8	5	7	100	500	7	4	6	6	3	10	7	5	1	312	12	330	25		24	31	26	19	125
Calenture Cancer, Gangrene, and Fiftula	26	29		. 1	31	53	36	37	73	31	3 24	35	63	52	20	14	23	28	27	30	24	30	85	112	105	157	150	114	609
Wolf	20			19	1	,,	- 1				- 1	"		- 1	- 1	14	*3	-	-/	30	*1			8	1				8
Carker, Spre-mouth, and Thruth Childhed	161				206	213	158	72 192	177	81	19 236	27	73	194	150	157	4	171	132	143	163	74	15 590	79 668	190	244	161 839	133	689 3364
Chrifomes, and Infants	1369		1065	990	1237	1280	1050	1343	1089	1393	1161	1144	858	1123	2596	2378					2113	1895	9277	8453	4678	4910	4788	4519	32106
Colick, and Wind Cold, and Cough	103	71	85	82	76	102	80	101	85	120	113	179	33	167	48	57	51	55	45	54	37	50	174		341	359 77		247	1389
Confemption, and Cough			2388	1988	2350	2410	2286	2868	2606	3184	2757	3610	2982	3414	1827	1910	1713	1797	1754	1955	2030	2477	5157	8266	8999	9914	12157	7197	44487
Convultion Gramp	684	491	530	493	369	653	606	828	702	1017	807	841	742	1031	52	87	18	241	223	386	418	709	498	1734	1198	2656	3377	1324	9073
Cut of the Stone		2	î	3		1		2	4		3	5	46	48			1	5	1	5	2	2	- 5	10	6	4	13	47	38
Dropfy, and Tympany Drowned	185	434			49	556	617 53	704	660 43	706	631	931	646 37	872 48	235 43	252	2.79	280	37	250	329	389 45	199	1734	1538	182	215	1302	9623 827
Exceffive drinking			2							- 1	1				- 1		- 1			- 1	- 1				2	- 1	1	1	2
Executed Fainted in a Bath	8	17	29	43	2.4	12	19	2.1	19	22	20	18	7	18	19	13	12	18	13	13	13	13	62	52	97	76	79	55	384
Falling-Sicknets	3	2		3	1	3	4		4	. 3	. 1		4	5	3	10	- 7	7	2	5	6	8	27	2.1	10	8	8	9	74
Flox, and finall pox Found dead in the Streets	139	100	1190	184	525	1279	139	812	1294	823	835	409	1523	354	72	40 33	58 26	531	72 13	2354	293	127	701	1846	1913	2755	3361	2785 19	10576
French-Pox	18	29	15	18	21	10	20	20	29	23	25	53	51	31	17	12	12	12	7	17	12	2.2	53	48	80	81	130	83	392
Frighted Gout	9	4	12	,	3	7	5	6	8	7	8		14	9	2	5	3	1	4	5	7	2	14	3 24	9 35	25		28	21 134
Grief	1.2	13	16	7	17	14	11	17	10	13	10	12	13	4	18	20	22	12	14	17	5	20	21	56	48 48	59	45	47	279
Hanged, and trade-away themselves Head-Ach	11	10	13	14	9	14	15	6	14	16	24	18	35	36	8	8	6	15		3	8	7	37	18	14	47	72 17	31 46	222 051
Jaundice	57		39	49	41	43	57	71	61	41	46	77	102	76	47	59	35	43	35	45	54	63	t84	197	180	212	225	188	998
Jaw-fain Impollume	75	61	65	59	80	105	79	90	92	122	80	134	105	96	10	16 76	13 73	74	10	62	73	130	47 282	35 315	260	354	428	228	1639
Inch		1	1										- 1				- 1		10		- 1		00	10	01		194	148	1021
Killed by Several Accidents Kong's Evil	27			94	47	45 20	57	58 26	32	43	52 23	47	55 28	47	16	55	18	18	49 35	41	51 26	69	202	150	217	207	102	66	537
Lethargy	3		2		4	4	3	10	9	4	6	2	6	4	1 2		2	2	3		2	2	5	7	13	21	21	9	67 06
Leprofy Livergrown, Spleen, and Rickets	53	46	1 16	59	65	72	67	65	52	50	38	51	8	15	94	112	99	87	82	77	2 98	99	392		218	269	101		1421
Lunatique	12	18	6	11	7 8	11	9	12	6	7	13	5	14	14	6	11	6	5	4	2	2	5	28	13	42	39	31		158
Meagroin Mealles	12			33	33	62	6	52	11	153	15	80	2	74	42	2	24	80	21	33	27	12	137		133	34 155			132 757
Mother	1 2	1			7	ī	- 1	2	1	3	- 1	3	ŧ	8	- 1			7	- i	6	5	3	01		17	13	- 8	77	18
Murdered Overlayd, and flarved at Nurfe	25	2 2 2 2 2		28	28	3	30	36	58	53	44	50	70 46	43	4	10	13	2	8	14	10	14	34	46		123	211		529
Palfy	27	21	19	20	2.3	20	29	18	2.2	23	20	2.2	17	21	17	23	17	25	14	21	25	17	81		87	90	87	53	423
Plague Plague in the Guts	3597	611	67	15	2.3	16	6	16	9 37	315	446	14	36 253	14 402		1317	274	. 8		-		10400	1599	10401		61	33 844	103	16384
Pleurify	30		13	20	23	19	17	23	10	9	17	16	12	10	26	24	26	36	21		45	24	112	90	89	72	52	51	415
Poyloned Purples, and spotted Fever	145	47	43	7 65	54	60	75	89	56	52	56	126	368	146	32	58	58	38	24	125	245	397	186	791	300	278	200	243	1845
Quanty, and Sore-throat	14	11	12	17	24	20	18	9	15	13	7	10	21	14	01	8	6	7	24	04	- 5	2.2	22	55	54	71	45	34	247
Rickets Mother, rifing of the Lights	150	224 91			134	329	135	372 178	347	458	317	476 228	110	249	44	73	99	98	60	14 84	49 72	50 104	309		780	1190	809	657 969	3681
Rupture	16	7		6	7	16	7	15	11	20	19	18	12	28	2	6	4	9	4	3	10	13	21	30	36	45	68		201
Scal'd-head Scurvy	72	20	23	21	1 29	43	41	44	101	71	82	82	95	12	5	7	9		9		00	35	33	34	94	131	300	115	95 593
Smothered, and flifted			2					- 1		- 1	- 1			- 1		24						1	24		2			- 2	26
Sores, Ulcers, broken and bruifed Shot (Limbs	15	17	17	16	26	32	25	32	2.3	34	40	47	7	48	23		2.0	48	19	19	22	29	91	89	65	115	144	141	504 27
Spleen	12	17					13	13		6	2	- 5	2	2											19	16	13		68



Average, mean, median, mode, etc

Given a sample
$$\mathbf{x} = \{x_1, \dots, x_n\}$$
, the average is $\overline{x} = \frac{1}{n} \sum_{i=1}^n x_i$

Note that
$$\overline{x} = \operatorname*{argmin}_{m \in \mathbb{R}} \left\{ \sum_{i=1}^n (x_i - m)^2 \right\}$$

In python

- 1 > import statistics
- 2 > x = [1, 2, 3, 4, 5, 6]
- 3 > statistics.mean(x)
- 4 3.5

and in R

- 1 > x = 1:6
- 2 > mean(x)
- 3 [1] 3.5

Average, mean, median, mode, etc.

The average (mean) is the empirical version of the expected value of a random variable,

$$\mathbb{E}(X) = \sum_{x} x \mathbb{P}[X = x] \text{ or } \int x f(x) dx$$

Example: a coin has heads with probability p. Let x = 1 (heads),

$$\mathbb{E}(X) = 1 \cdot p + 0 \cdot (1 - p) = p$$

Linearity:

$$\mathbb{E}(aX+b)=a\mathbb{E}(X)+b, \ \forall a,b\in\mathbb{R},X$$

$$\mathbb{E}(X_1+\cdots+X_k)=\mathbb{E}(X_1)+\cdots\mathbb{E}(X_k), \ \forall X_1,\cdots,X_k$$

Example: toss n coins, of bias p, X is the number of heads

$$\mathbb{E}(X) = \mathbb{E}(X_1 + \dots + X_n) = \mathbb{E}(X_1) + \dots \mathbb{E}(X_n) = np$$













St Petersburg's Paradox

As we will see (law of large numbers) if x_i are realizations of random variables X_i (with identical expected value μ), $\overline{X} \to \mu$ as $n \to \infty$.

A fair coin is tossed at each stage. The initial stake begins at 2 dollars and is doubled every time heads appears. The first time tails appears, the game ends and the player wins whatever is in the pot. Let X denote the gain.

$$\mathbb{E}(X) = \frac{1}{2} \cdot 2 + \frac{1}{4} \cdot 4 + \frac{1}{8} \cdot 8 + \frac{1}{16} \cdot 16 + \dots = 1 + 1 + 1 + 1 + \dots = +\infty$$

the expected value is infinite (but the average always exists)

Average, mean, median, mode, etc

The average is very sensitive to outliers and extremal values.

Denote $\{x_{(i)}\}$ the ordered version of $\{x_i\}$, $x_{(1)} \leq x_{(2)} \leq \cdots \leq x_{(n)}$

Trimmed mean =
$$\frac{1}{n-2p} \sum_{i=1+p}^{n-p} x_{(i)}$$

Weighted mean =
$$\sum_{i=1}^{n} \omega_i x_i$$
 where $\omega_i = \frac{w_i}{w_1 + \dots + w_n}$

(classical on surveys with underrepresentation, see fao or Stat Can) Median $md(x) = x_{(n/2)}$ (50% observations are smaller/larger)

Note that
$$md(x) = \operatorname*{argmin}_{m \in \mathbb{R}} \left\{ \sum_{i=1}^{n} |x_i - m| \right\}$$
 (see #224)

Average ratio

See Citation Analysis and Dynamics of Citation

Networks by Michael Golosovsky.

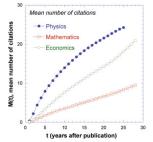
	0	1	2	3	4	5
2014	122	146	163	176	182	186
2015	142	174	198	214	228	
2016	185	215	246	278		'
2017	214	265	312		,	
2018	245	326		,		
2019	261		'			

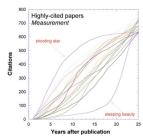
The average ratio is

$$\frac{326 + 265 + 215 + 174 + 146}{245 + 214 + 185 + 142 + 122}$$

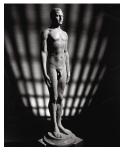
and not the average of ratios

$$\frac{1}{5} \left(\frac{326}{245} + \frac{265}{214} + \frac{215}{185} + \frac{174}{142} + \frac{146}{122} \right)$$

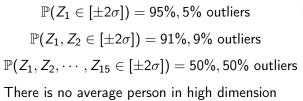




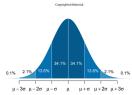
Norms and Averages See The End of Average, by Todd Rose







OUR FUNDAMENTAL ASSUMPTIONS ABOUT TALENT THE END OF AVERAGE UNLOCKING OUR POTENTIAL BY EMBRACING WHAT MAKES **US DIFFERENT** TODD ROSE



(see also curse of dimensionality)

What is an Average

"Take any stock in the United States. The average time in which you hold a stock is - it's gone up from 20 seconds to 22 seconds in the last year" The Telegraph, by Michael Hudson

"Even though most are judged by performance over three-year horizons, their average holding period was about 17 months, and 19% of the managers held the typical stock for one year or less" The Wall Street Journal, by Jason Zweig

