

Statistics

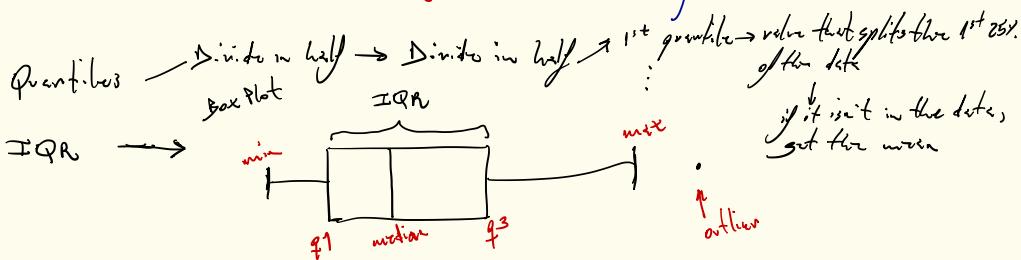
Philippe Fanaro

	<i>heavily influenced by outliers</i>	
mean \rightarrow "average"	$\frac{\sum x_i}{n}$	
median \rightarrow middle value		no spread-out description
mode \rightarrow most common value		no description of the shape

} Central Tendency

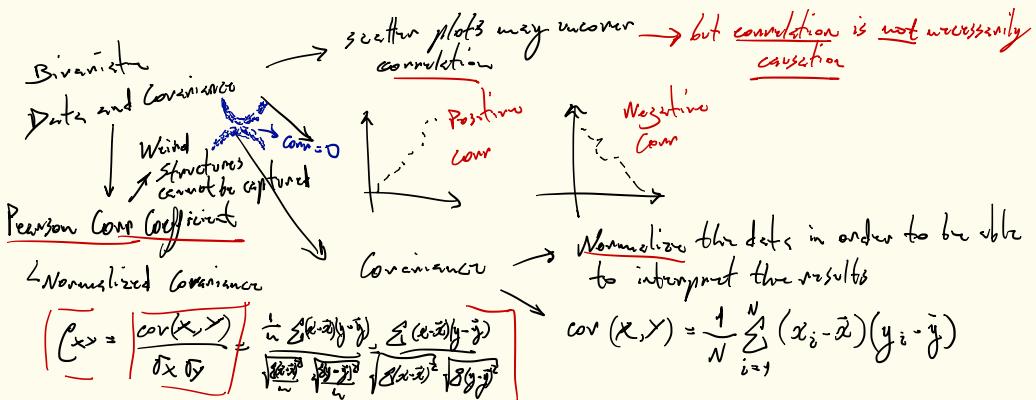
range $\rightarrow (\max - \min)$	sample	population \leftarrow depends on the context it's not always everything
variance $\rightarrow \frac{\sum (x_i - \bar{x})^2}{n-1}$	$\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{N}$	
standard deviation $\rightarrow \sqrt{\sigma^2}$	Degrees of freedom correction	
		use <u>units</u> as the main variable

} Dispersion



Fence and Outliers \rightarrow Commonly, we fence at $1.5 \times \text{IQR}$ beyond q_1 and q_3 \rightarrow outliers are drawn away from the whiskers

Determined by the data, not a prior definition!



Permutation

- $P_{n \times n}! \quad \text{PERMUT}(n, n)$
- Out of n subset → $P_r = \frac{n!}{(n-r)!}$ (no repetition) $\text{PERMUTATION}_R(n, r)$
- with repetition → n^r arrangement → $A_{n, r}$

Combinations → Unordered

$${}^n C_r = C_{n,r} = \frac{n!}{r!(n-r)!}$$

With repetition $\rightarrow C_{w+r-1, r} = \frac{(w+r-1)!}{r!(w-1)!}$ combin A(w,r) could happen

Interventions, Unions and Complements

$$\Rightarrow A \wedge B \rightarrow \text{AND}$$

$$A \cup B \Rightarrow \boxed{\text{OR}} \rightarrow P(A) + P(B) - P(A \cap B)$$

$$\rightarrow \bar{A} = U - A \Rightarrow 1 - P(A)$$

— $A \cup B \neq V$ could happen!

Conditional Prob and Bayes Theorem

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(B|A) = \frac{P(A \cap B)}{P(A)}$$

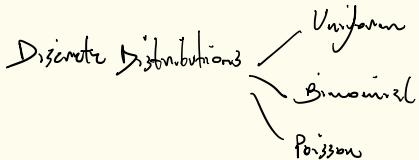
Bayes Theorem:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

$$P(Day | Pos) = \frac{P(Pos | Day) P(Day)}{P(Pos)} =$$

$$P(\text{dry} \mid P_{03}) = \frac{0,99 \cdot 0,165}{0,99 \cdot 0,165 + 0,01}$$

$$\begin{aligned}
 &= \frac{0,99 \cdot 0,002}{P(T \text{ Frob}) + P(F \text{ Frob})} = \frac{0,99 \cdot 0,002}{P(P_{01} | D_{01}) \cdot P(D_{01})} = \\
 &\quad + P(P_{02} | \overline{D}_{01}) \cdot P(\overline{D}_{01}) \\
 P(P_{01} | T) + P(P_{02} | NF) &= \\
 &= \frac{0,99 \cdot 0,002}{0,99 \cdot 0,002 + 0,01 \cdot 0,998} = 0,165
 \end{aligned}$$



Bernoulli Trial

- 2 possible outcomes
- success
- independent trials

$$P(x; n, p) = \binom{n}{x} p^x (1-p)^{n-x}$$

from sci-py.stats import binom
 $\text{binom.pmf}(3, 16, 1/6)$

Normal Distribution

$$\mu = np$$

$$\sigma^2 = np(1-p)$$

Poisson \rightarrow success out of n trials

Poisson \rightarrow successes per unit of time (continuous unit)

$$\lambda = \frac{\# \text{ occurrences}}{\text{interval}} = \mu = E(X) = \sigma^2$$

$P(x) = \frac{\lambda^x e^{-\lambda}}{x!} \rightarrow$ Prob of seeing # occurrences per interval

$$\text{cdf}(x) = P(X \leq x) = \sum_{i=0}^x \frac{\lambda^i e^{-\lambda}}{i!}$$

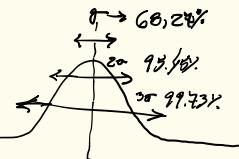
assumes that
 that prob during
 small time interval
 is proportional to the
 entire length, i.e.,
 $\lambda_{\text{min}} = \lambda_{\text{max}}$

Continuous Distributions

- Normal
- Exponential
- Beta

Normal

Standard Normal Distribution $N(0, 1)$

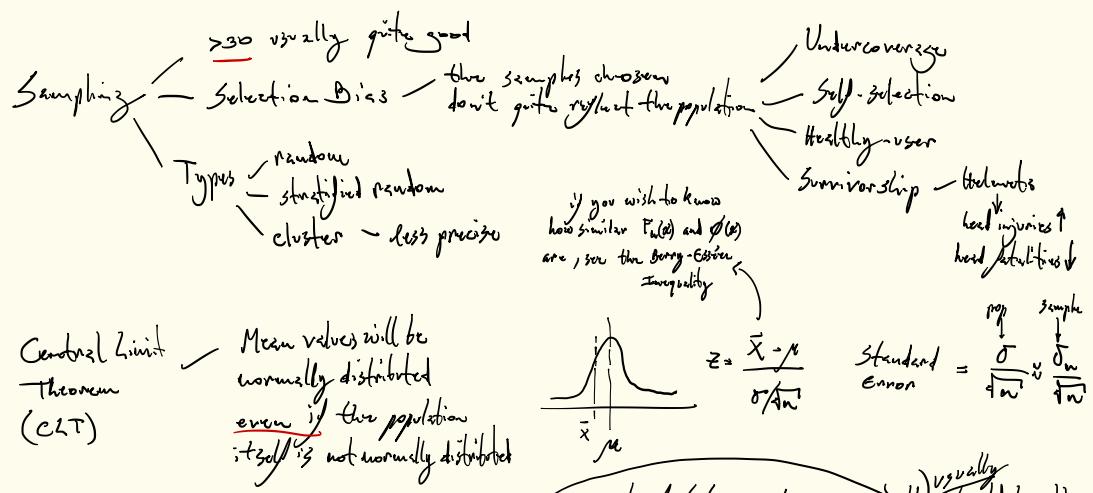


$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Z-Scores

$$z = \frac{x-\mu}{\sigma}$$

from sci-py import stats
 stats.norm.cdf(z=0.4)
 stats.norm.pdf(mu=0.95)



Hypothesis Testing

Start with a Null hypothesis

We never prove a hypothesis

Null hypothesis should confirm an equality (=) level of significance α → tail(s) of the null hypothesis

→ left-tail
→ right-tail
= → two-tail

Accept / fail to reject

Reject → Assume another mutually exclusive hypothesis

H₀ usually should be the opposite of what you want to "prove"

Testing Means vs Proportions

$$Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

$$Z = \frac{\hat{P} - P}{\sqrt{\frac{P(1-P)}{n}}} = \frac{\hat{P} - P}{\sqrt{P(1-P)/n}}$$

Example 1: Load time $\mu = 3,125$, $\sigma = 0,49$

Desired confidence on improvement: 99% $\rightarrow \alpha = 0,01$
we do

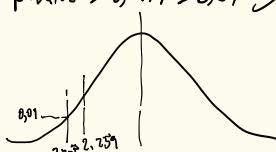
New average $\rightarrow \bar{x} = 2,875$

H₀: $\mu \geq 3,125$ → left-tail

$$Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = -2,259$$

$\rightarrow Z = -2,325 \rightarrow$ fail to reject H₀

p-value $\rightarrow 0,0119 > 0,01 \rightarrow$ fail to reject H₀



- P-value test

- take test statistic
- use it to determine the P-value
- compare the P-value to the level of significance α

P-value is low
Null must go → reject the

P-value is high
Null must stay → fail to reject the

Example 2:

n=100

50% → turners
are most customers
turners?

$$\text{H}_0: P \leq 0,5 \quad \alpha = 0,05 \quad (\text{right-tail})$$

$$Z = \frac{0,50 - 0,5}{\sqrt{\frac{0,5(1-0,5)}{100}}} = \frac{0,00}{\sqrt{0,005}} = 0,00 = 0,00$$

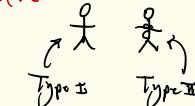
$> 1,645 \rightarrow$ reject H₀

With 99% confidence
the new average
is not an improvement!

Type I and II Errors

		H_0 is	
		Type I $\alpha = P(\text{Type I})$	Type II $\beta = P(\text{Type II})$
Decision about H_0	Fail to Project	Type I $\alpha = P(\text{Type I})$	Type II $\beta = P(\text{Type II})$
	Project	Type I $\alpha = P(\text{Type I})$	Type II $\beta = P(\text{Type II})$

False Positive



"You're pregnant"

"You're not pregnant"

Student's T-Distribution

L William Gossett → Student
Guinness Brewery

→ Select the best barley from small samples where the standard deviation was unknown.

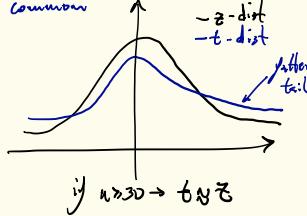
t-table → t-test determines if there is significant difference between t-statistic two sets of data

$$(1): \quad t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \quad \begin{cases} \text{t-table} \\ \text{Degrees of freedom} \\ \text{Chosen significance level} \\ \text{t} \rightarrow t\text{-statistic} \\ t_{n-1, \alpha} = t\text{-critical} \\ n-1 = \text{degrees of freedom} \\ \alpha = \text{significance level} \end{cases}$$

t-test → signals noise

- $$(2): \quad \begin{cases} \text{equal sample sizes, equal variance} \\ \text{unequal sample sizes, equal variance} \\ \text{unequal or equal sample sizes, unequal variance} \end{cases}$$

Most common



$$t = \frac{\text{difference in means}}{\text{sample variability}} = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

degrees of freedom:

$$\hookrightarrow df = \left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2} \right)^{-1} = \frac{1}{\frac{(s_1^2)^2}{n_1-1}} + \frac{1}{\frac{(s_2^2)^2}{n_2-1}}$$

if $s_1 = s_2$:

$$df = n_1 + n_2 - 2 = (n_1 - 1) + (n_2 - 1)$$

Example:

Two plants, same car
Which are to close?

$$n_A = 10 \quad \bar{x}_A = 1222 \quad \text{Is A statistically different from B?}$$

$$n_B = 10 \quad \bar{x}_B = 1186$$

$$H_0: x_A = x_B \quad \text{one-tailed test} \quad S_A^2 = 1248$$

$$H_1: x_A > x_B \quad df = (10+10-3) = 17 \quad S_B^2 = 1246$$

t-table	
over.prob	t.1%
outtail	0.1
two tail	0.2
1%	-
10%	-
19%	-
Critical Value: 1.733	-
2.28	1.733

Reject H_0 → Plant A produces more than B

ANOVA

(analysis of variance)

— Previously: Z- and t-distributions

↓
Now: F-distributions



If you have 3 samples, you can make 3 pairings and t-test them.

But the overall confidence will drop: 0.95 · 0.95 · 0.95 = 0.854

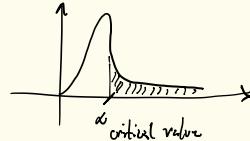
$$\bar{M}_{\text{Total}} \\ M_A, M_B, M_C, \dots$$

Variance

between groups
within groups

$$F = \frac{\text{Var between groups}}{\text{Var within groups}} = \frac{\frac{\sum_{\text{groups}} (M_i - \bar{M})^2}{\text{groups}}}{\frac{\sum_{\text{groups}} \sum_{\text{samples}} (x_{ij} - \bar{x}_i)^2}{(\text{groups} - 1) \cdot \text{samples}}} = \frac{\sum (x_i - \bar{x})^2}{n - 1}$$

F-Distribution



	F-table Upper Tail Area of 0.05			
	1	2	3	...
denominator	2.78			
numerator	4.60			
	:	:	:	:

Example

Give Discounts to different customers → does it affect them? → H₀: discounts yield the same behavior (distribution)

days delayed
with discount

	2% Disc	1% Disc	No Disc
11	21	11	
16	15	11	
9	23	10	
13	10	16	
10	16	21	

$$F = \frac{\frac{11 \cdot 3}{3-1}}{\frac{198}{(5-1) \cdot 3}} = 8,121 \Rightarrow F_{\text{critical}} = 3,885$$

$$F = 8,121 > 3,885 = F_{\text{critical}}$$

↓
Fail to reject H₀

	\bar{M}	12	17	16
M_{Total}	15			
$(M_i - M_{\text{Total}})^2$	9	9	1	$\rightarrow \sum (M_i - M_{\text{Total}})^2 = 18$
$(x_{ij} - \bar{x}_i)^2$	58	106	58	$\rightarrow \sum (x_{ij} - \bar{x}_i)^2 = 198$

Two-Way ANOVA → More variables → e.g. 3 invoices for \$50, 3 invoices for \$100
 different amounts may yield different conversion rates

	2%	1%	0%	Block
\$50	16	23	21	20
\$100	18	21	16	18
\$150	11	16	18	15
\$200	10	15	18	13
\$250	9	10	11	10
Group	12	17	16	15
(Group %)	9	8	1	19
				Total

each row is a block

You want to isolate and remove any variance contributed by the blocks, to better understand the variance in the groups.

↳ Block, SS Block

$$SSE_{\text{error}} = SST_{\text{total}} - SSG_{\text{groups}} - SS_{\text{Blocks}}$$

$$df_{\text{groups}} = n_{\text{groups}} - 1$$

$$df_{\text{error}} = (n_{\text{blocks}} - 1)(n_{\text{groups}} - 1)$$

$$\begin{aligned} F &= \frac{SSB}{df_{\text{groups}}} \\ &= \frac{SSB}{df_{\text{error}}} \end{aligned}$$

$$SSG = 19 \cdot 5 = 95 \quad \text{samples in group}$$

$$df_{\text{groups}} = 3 - 1 = 2$$

$$SSB = 50 \cdot 5 = 250 \quad \text{samples in block}$$

$$F_{\text{groups}} = \frac{250}{25} = \frac{25}{3} = 11.67$$

$$SST = 260$$

$$SSE = SST - SSG - SSB = 8$$

$$df_{\text{error}} = 0.05 \quad df_{\text{groups}} = 2 \quad df_{\text{error}} = 8$$

$$df_{\text{error}} = (5-1)(3-1) = 8$$

$$\text{F critical} = 3.86$$

$$F_{\text{critical}} = 3.86 < 11.67 = F_{\text{groups}}$$

reject H₀ → giving discordant results statistically different results

Two-Way ANOVA with Replication

heights of plants ↘ types of fertilizers A, B, C (groups)
 ↘ temperatures warm and cold (blocks)

Fertilizer	A	B	C	Block
Warm	13	21	18	16
	18	19	15	
	12	17	16	
Cold	16	18	15	18
	18	19	13	
	14	18	8	
M _W	13	19	16	Mod
M _C	18	13	12	18

$$\sum (M_B - M_T)^2 = (16 - 15)^2 + (18 - 15)^2 + 2 \quad \text{samples in block}$$

$$SSB = 2 \cdot 5 = 10$$

$$\sum (M_G - M_T)^2 = (15 - 15)^2 + (16 - 15)^2 + (18 - 15)^2 = 2$$

$$SSG = 2 \cdot 6 = 12$$

$$(M_T - M_B)$$

$$df_{\text{groups}} = 3 - 1 = 2$$

$$df_{\text{blocks}} = 2 - 1 = 1$$

$$SST = 168$$

$$SSE = 168 - 18 - 12 - 8 = 50$$

$$df_{\text{error}} = \text{blocks} \cdot \text{groups} (r_{\text{blocks}} - 1) = 2 \cdot 3 \cdot (3 - 1) = 12$$

$$F_{\text{critical}} = 3.86$$

$$F = \frac{18}{50/12} = 1.68 < 3.86 \quad \alpha = 0.05$$

M_W
 sample means

Regression

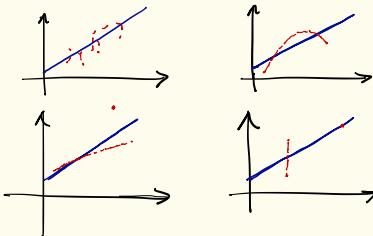
$$\hat{y} = b_0 + b_1 x$$

$$b_1 = \rho_{x,y} \frac{\bar{y}}{\bar{x}} \quad (\rho_{x,y} = \text{Pearson Corr Coefficient} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2}})$$

$$b_0 = \bar{y} - b_1 \bar{x}$$

Limitations — Anscombe Quartet
 illustrates the pitfalls
 of relying on pure calculation
 (1973)

L
 1. simple Regression
 2. faltally different
 underlying functions



Multivariable Regression

$$\hat{y} = b_0 + b_1 x_1 + b_2 x_2 + \dots$$

$$b_1 = \frac{\sum (x_1 - \bar{x}_1)^2 \sum (x_1 - \bar{x}_1)(y - \bar{y}) - \sum (x_1 - \bar{x}_1) \sum (x_2 - \bar{x}_2)(y - \bar{y})}{\sum (x_1 - \bar{x}_1)^2 \sum (x_2 - \bar{x}_2)^2 - (\sum (x_1 - \bar{x}_1)(x_2 - \bar{x}_2))^2}$$

$$b_2 = \dots$$

- Avoid using factors without convolution \rightarrow they will just be noise.

Chi-Squared analysis

χ^2 / Karl Pearson (1900)

how much our observations diverge from the expected

$$\chi^2 = \sum \frac{(O-E)^2}{E} = \frac{(12-9)^2}{9} + \frac{(6-9)^2}{9} = 2$$

↓
expected

$$Q = \sum_{i=1}^{k-1} Z_i^2 \sim \chi^2(k-1)$$

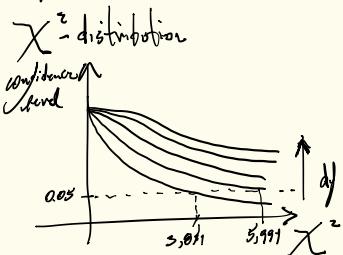
)
the product of two independent chi-squared distributions is a chi-squared distribution

$$df_{\text{coins}} = 2-1 = 1$$

$$\chi^2 = 2 < 3,831 = \chi^2_{\text{crit}}$$

↓
fail to reject H₀

(H₀: 12 heads out of 18 flips was reasonable)



Example:

Based on the following data, can we assume that servers fail at the same rate?

Assumptions:
 - failures are independent.
 - no "degrees of failure", either
 $\rightarrow \frac{\text{observed}}{6} \rightarrow 10$ fail or not.

servers	observed	expected	$(O-E)^2/E$
A	46	40	0,9
B	36	40	0,6
C	52	40	3,6
D	26	40	3,9
E	42	40	0,1
F	38	40	0,1
G	28		

$$\chi^2 = \sum \frac{(O-E)^2}{E} = 10$$

$$\alpha = 0,05$$

$$df = G-1 = 5$$

$$\chi^2_{\text{crit}} = 11,07$$

$$\chi^2 = 10 < 11,07 = \chi^2_{\text{crit}}$$

↓
fail to reject H₀

↳ 95% confidence that they converge to the expected