

Topics:

- Types of ANOVA ←
- α, β , power & effect-size - Defns & terms
- Kurtosis ✓
- Covariance ✓
- Pearson correlation coefficient
- Spearman rank corr. coefficient

Search Google or type a URL

Gmail Images



Google

Search Google or type a URL



Colaboratory



My Drive



YouTube



Learning



InterviewBit S...



GitHub



Scaler Academ...



Reduce the file...



InterviewBit



Add shortcut

✓ Programming → Mock-interview

✓ SQL → Mock-interview

[Foundational Math →]

✓ "Framework"

✓ "Framework"

→ Tens of Variations

learnt

One way

✓ indep. variable

drug-type → d_1, d_2, d_3

→ #days to recover
dependent var

✓
2-way

dep-var
= #days to recover

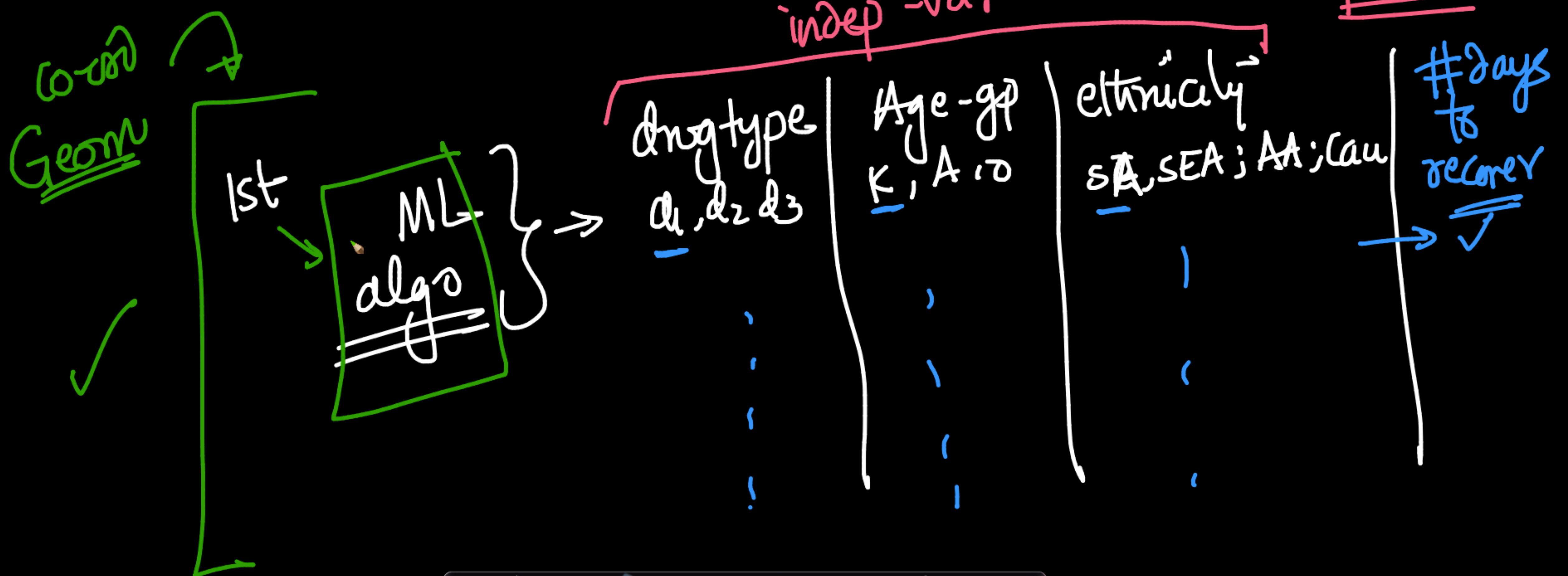
indep-var ✓
d₁, d₂, d₃
2 → { k, A, D }

ANOVA
Variations

✓
Multiway ANOVA
72

dep-var = #days to recover

indep-var = → d_1, d_2, d_3
→ K, A, D
→ SA; SEA
|| AA; Gau



MANOVA

indep-var

college-degree

[HS, B, M, PhD]

one-way Manova

Dep-var:

multiple

Salary

sq.feet of
the home

2-way
=

MANOVA:

$\alpha = \underline{\text{S}y}$.

dep. var:

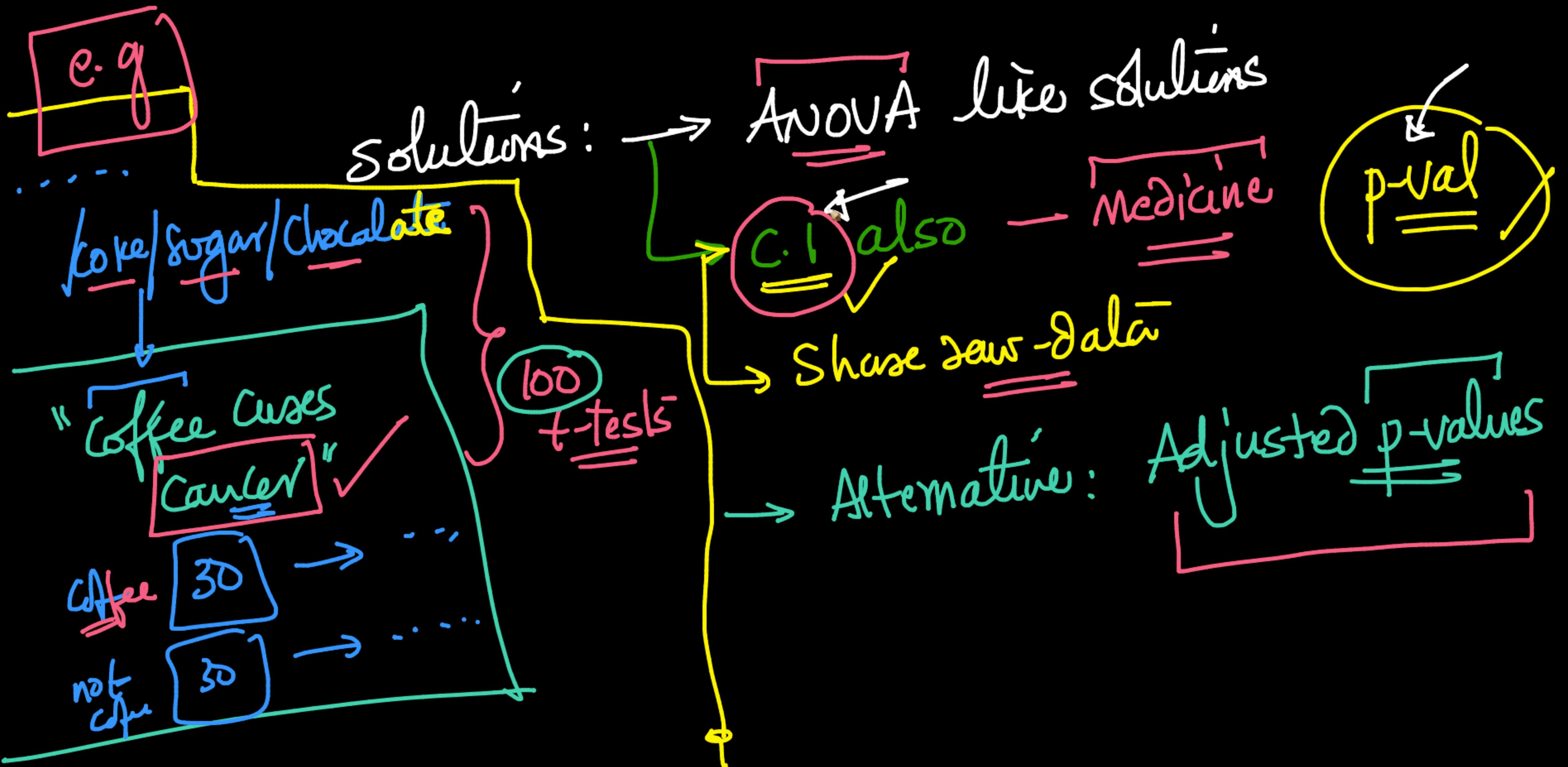
↳ Salary }
↳ sq ft of the
home }

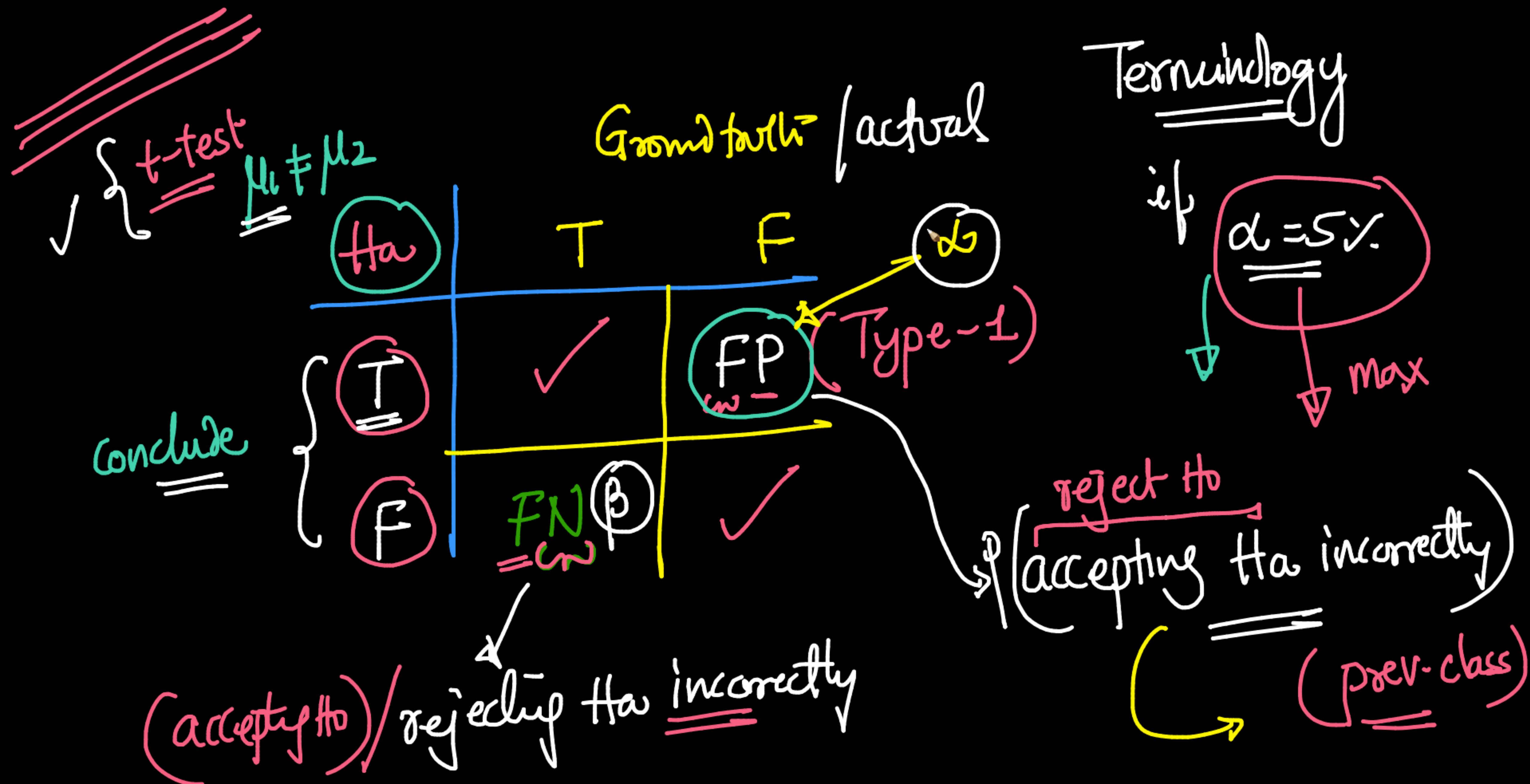
1
2 indep. var
↳ College-degree
[HS, B, M, PhD]
↳ 18-score
[L, M, H, NHT]

multiple-stat-test:
 K_{C_2}

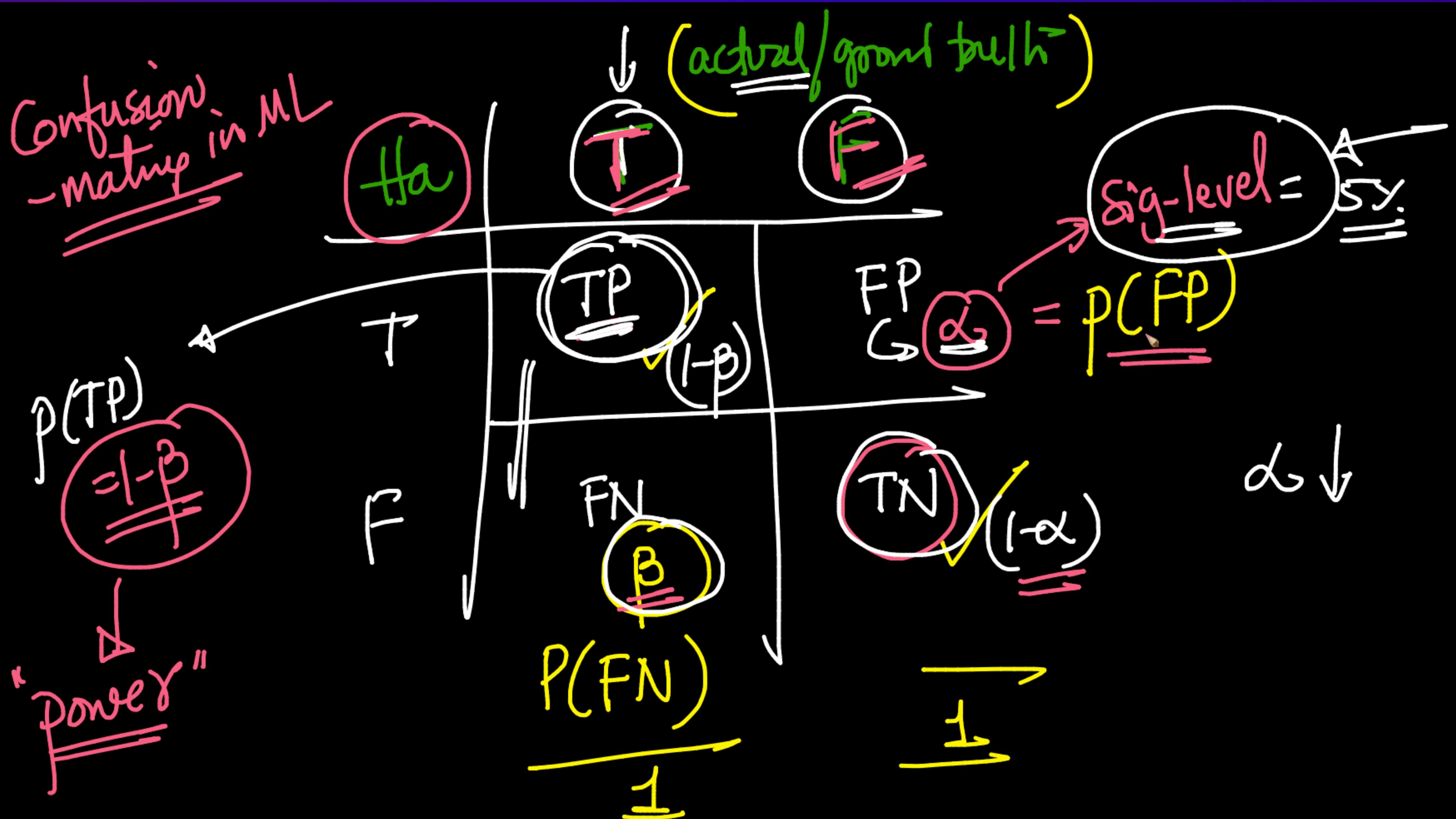
p-hacking

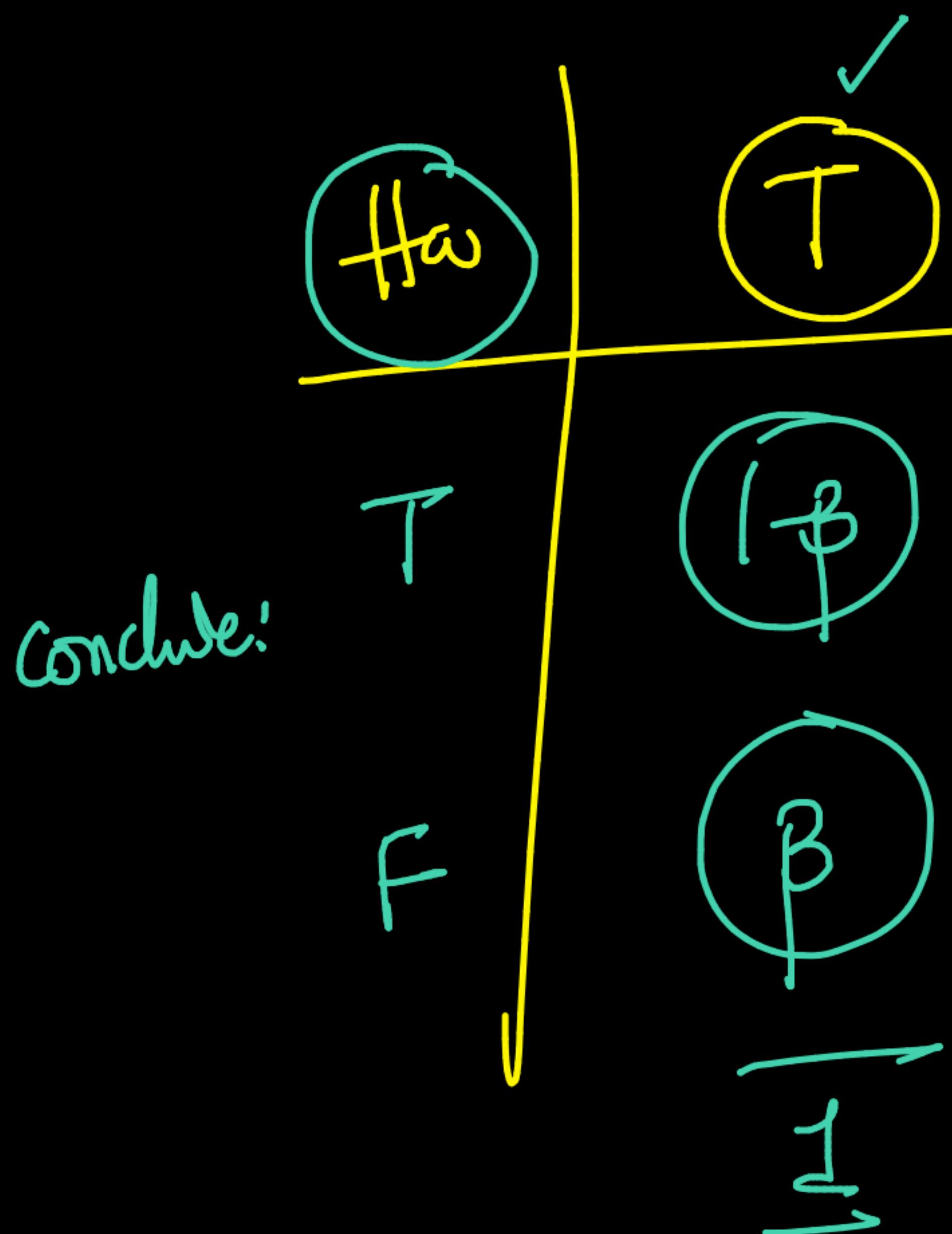
Too Many
Simultaneous 2-sample tests
K-drags $\rightarrow K_{C_2}$ z/t-tests (BAD)
 $\hookrightarrow p$ (accept one one incorrectly) ≈ 0.9
when $t=10$
 \hookrightarrow (prev-class)



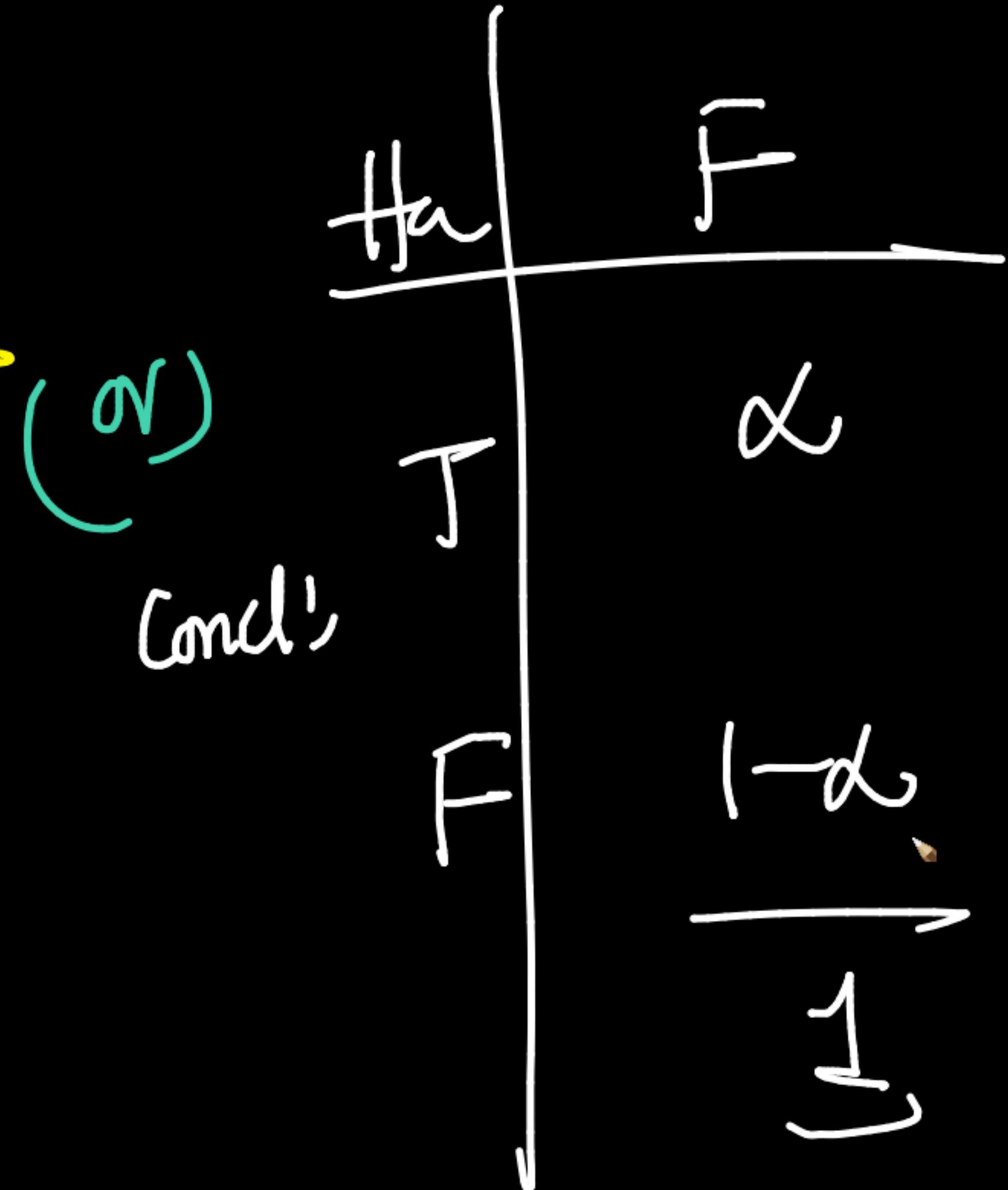


Confusion in ML





Conclude:



$\alpha = \text{sig-level}$

$\beta = \text{power}$

Effect-size

H_s : heights of Swedish people

$\underline{H_t}$: " "

$\underline{H_c}$: " "

Thailand

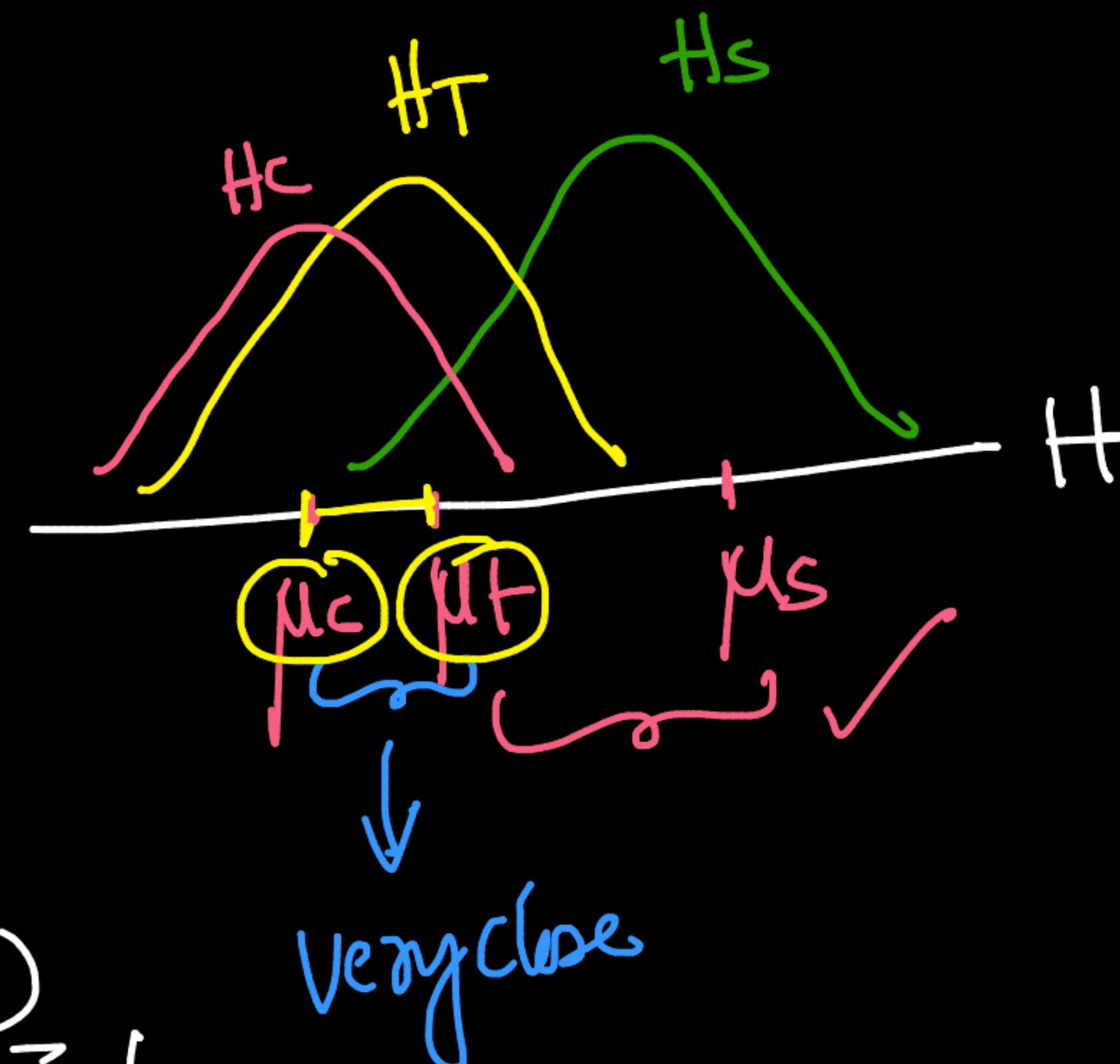
China

t-test:

① $\mu_C = \mu_T$

effect size is small

(TP)
power the test ↓

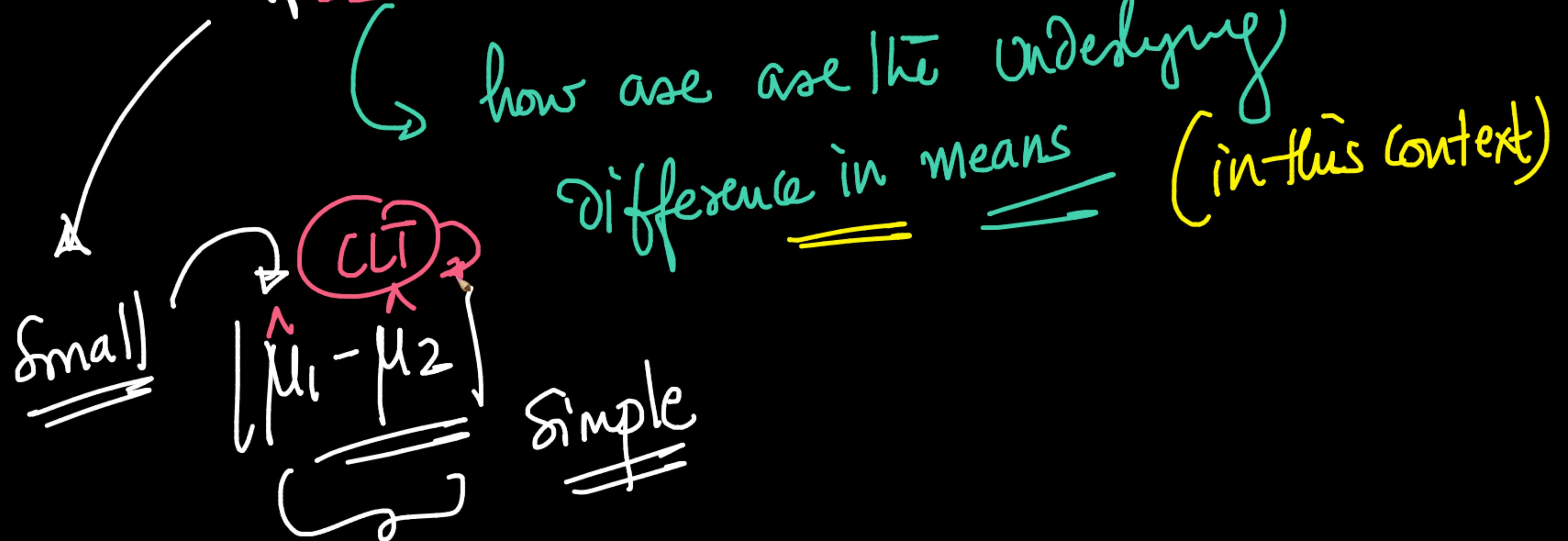


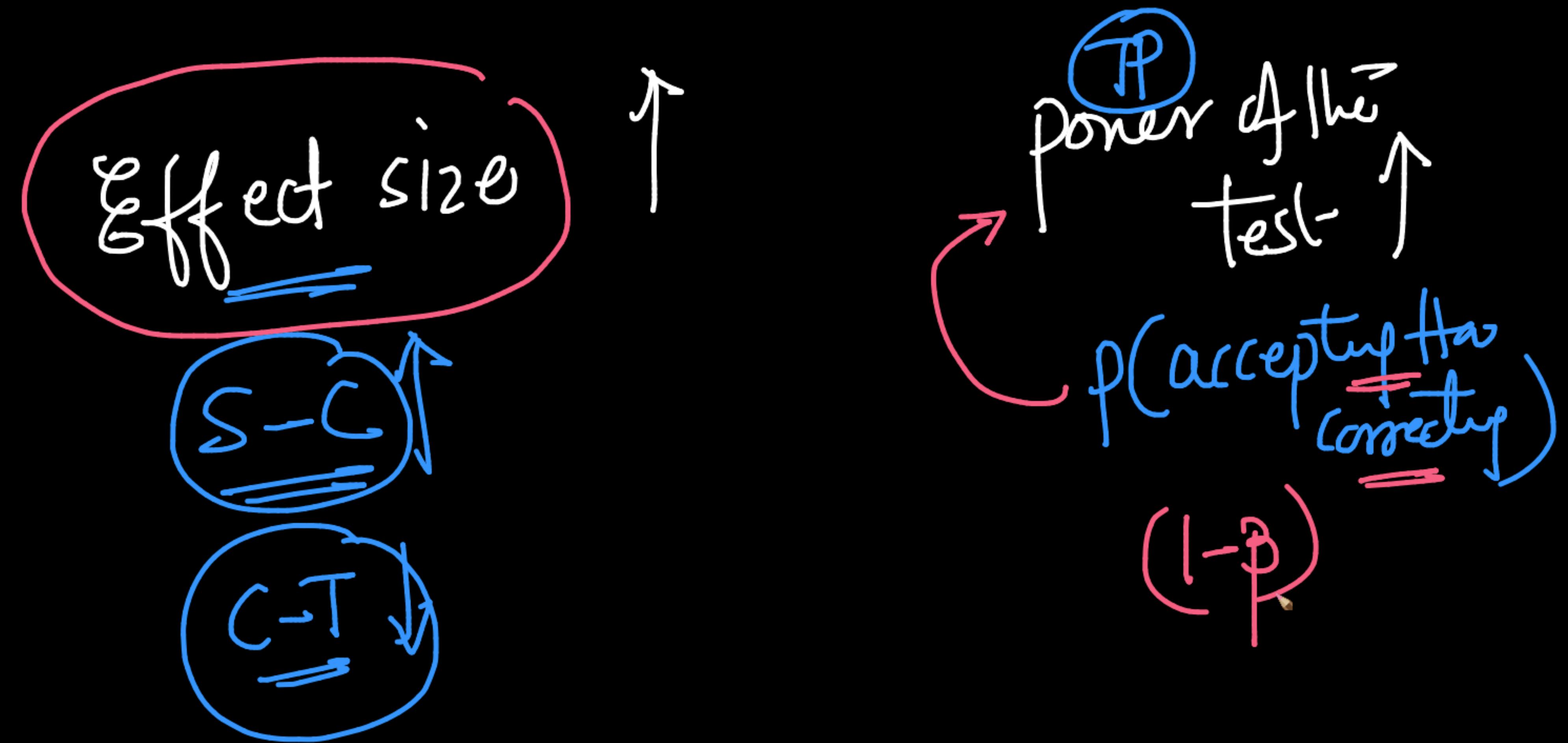
② $\mu_T = \mu_S$

effect size is larger

Power of my test ↑

Effect-size in hypothesis





✓ T-test: $\rightarrow \alpha, \beta$
(problematic)

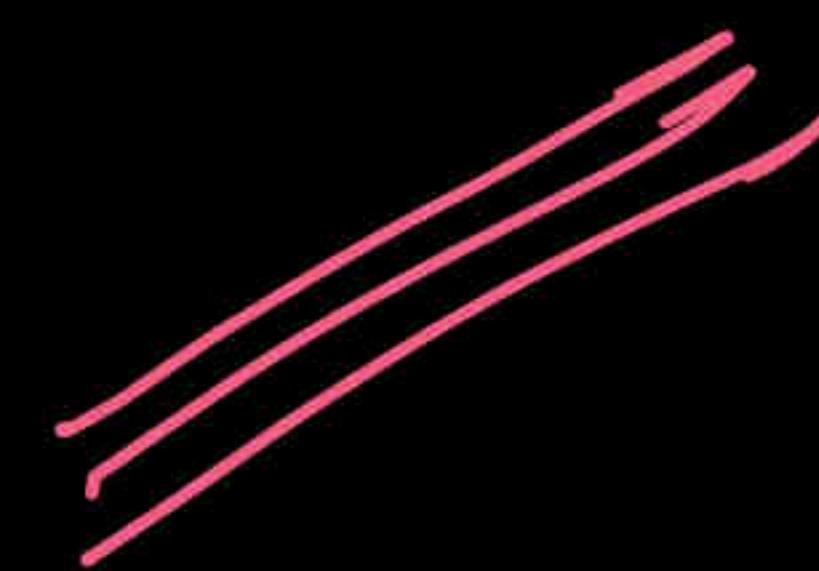
✓ TX

{ Effect-size: how different are the
underlying populations }

Practical
Increase (n_1, n_2) \hookrightarrow Small \rightarrow Power $\downarrow\downarrow$

TP

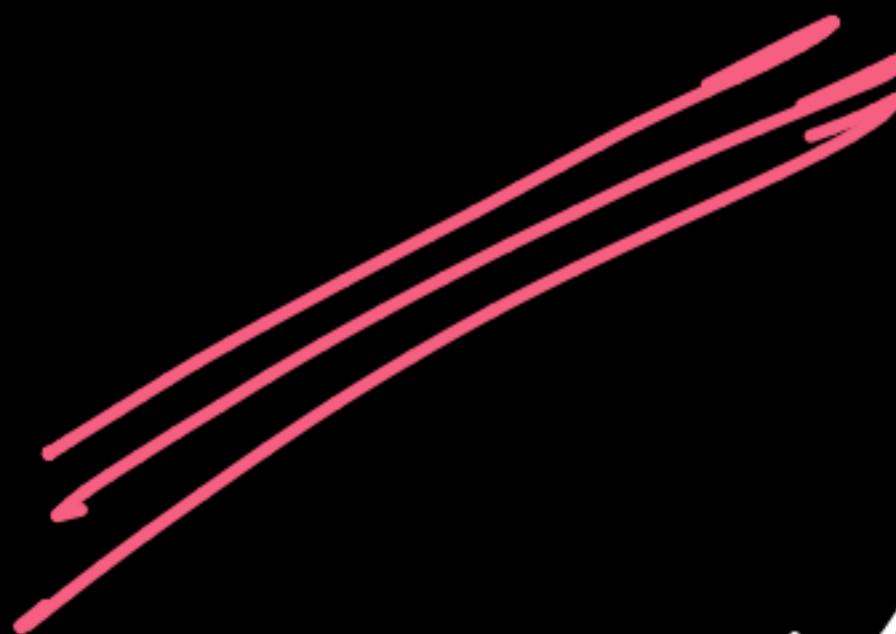
H_a



A/B test

$\alpha = 0.1\%$
 β = function of test; $(n_1, n_2) \uparrow \dots$

power \uparrow
 \equiv
 $1 - \alpha \uparrow$



measure
of tailoredness

$$K = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^4$$

Kurtosis

Mean: measure
of central
value

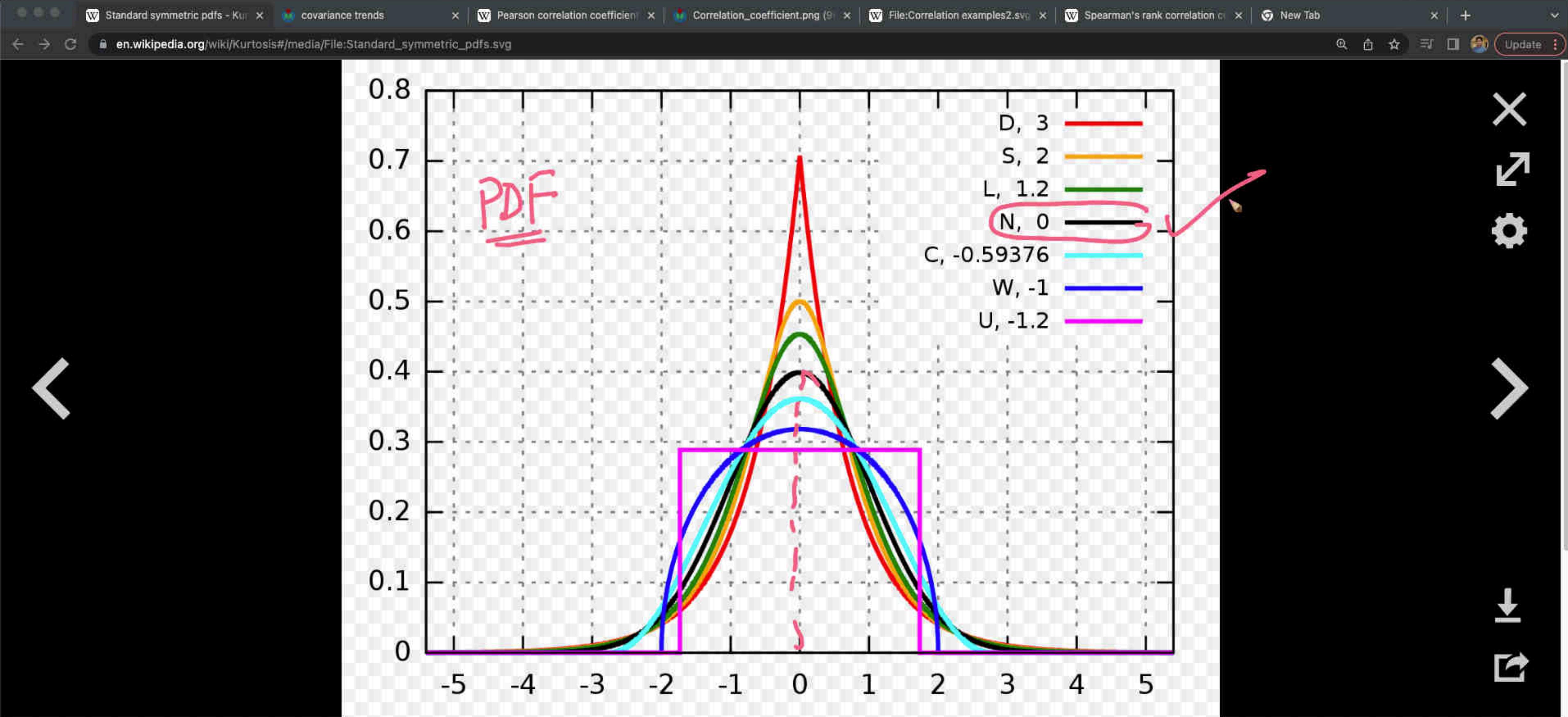
Var: measure
of spread

$K_{\text{normal}} = 3$ (proven...)

"Excess-kurtosis" = $K - 3$

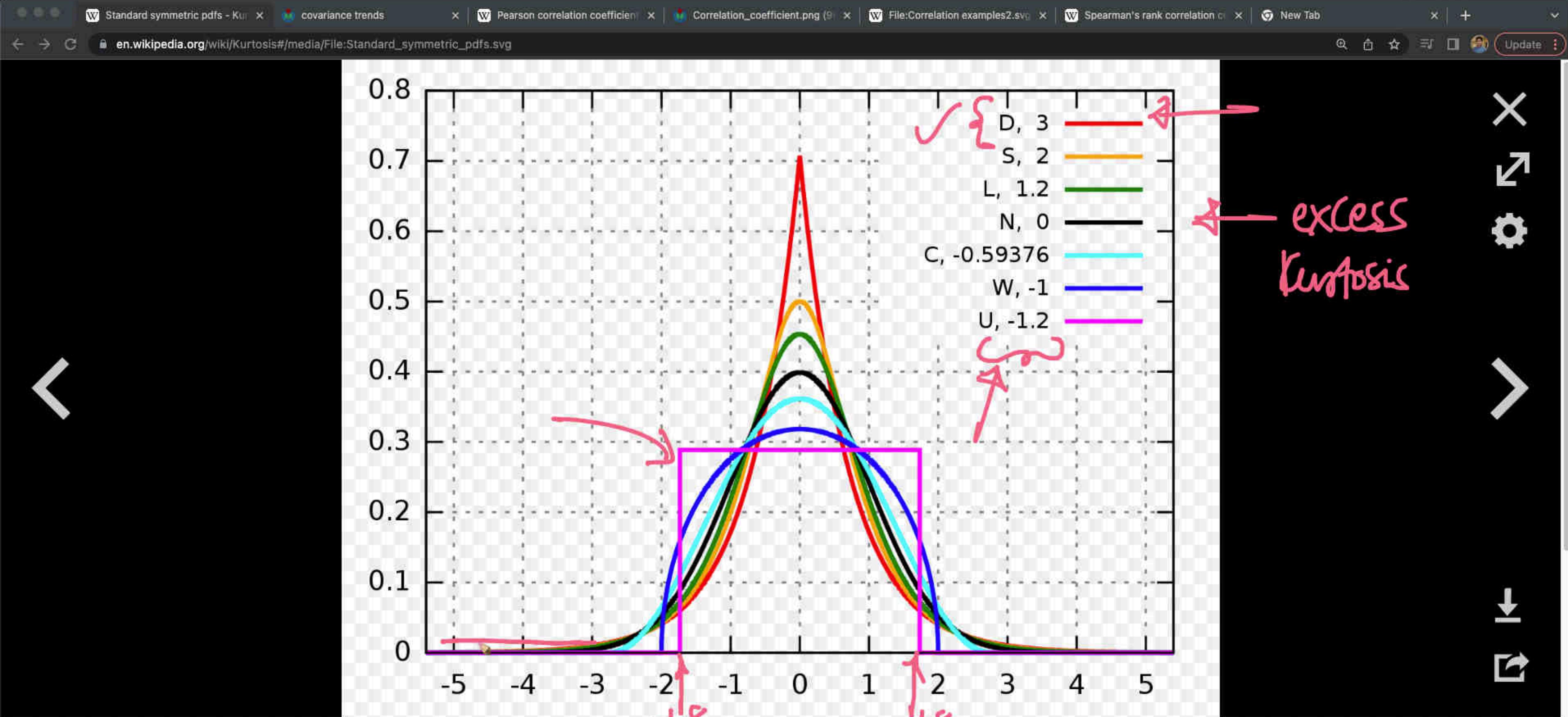
Normal = 0

$K_e > 0$ fatter tails
than Gaussian



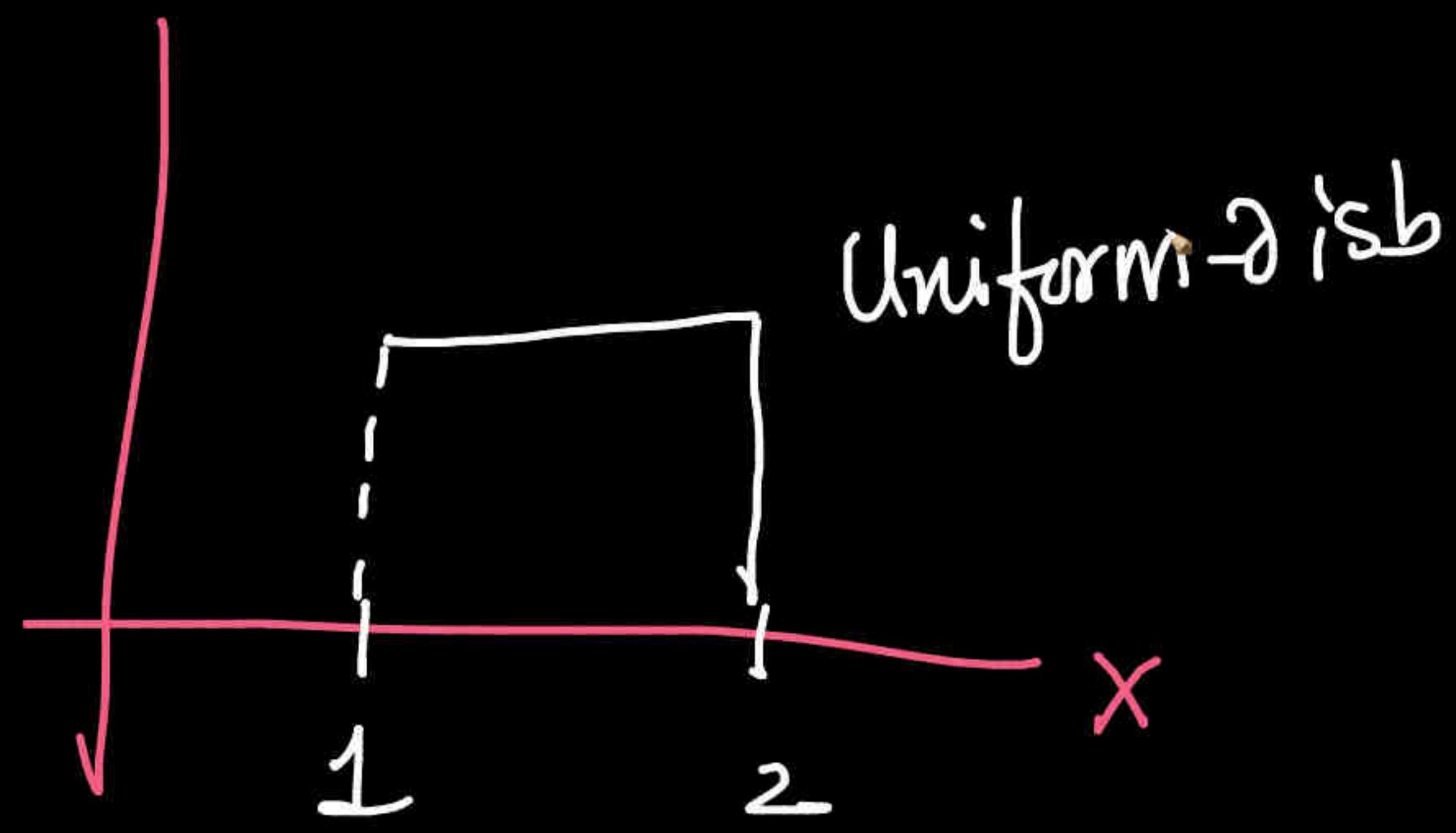
Probability density functions for selected distributions with mean 0, variance 1 and different excess kurtosis

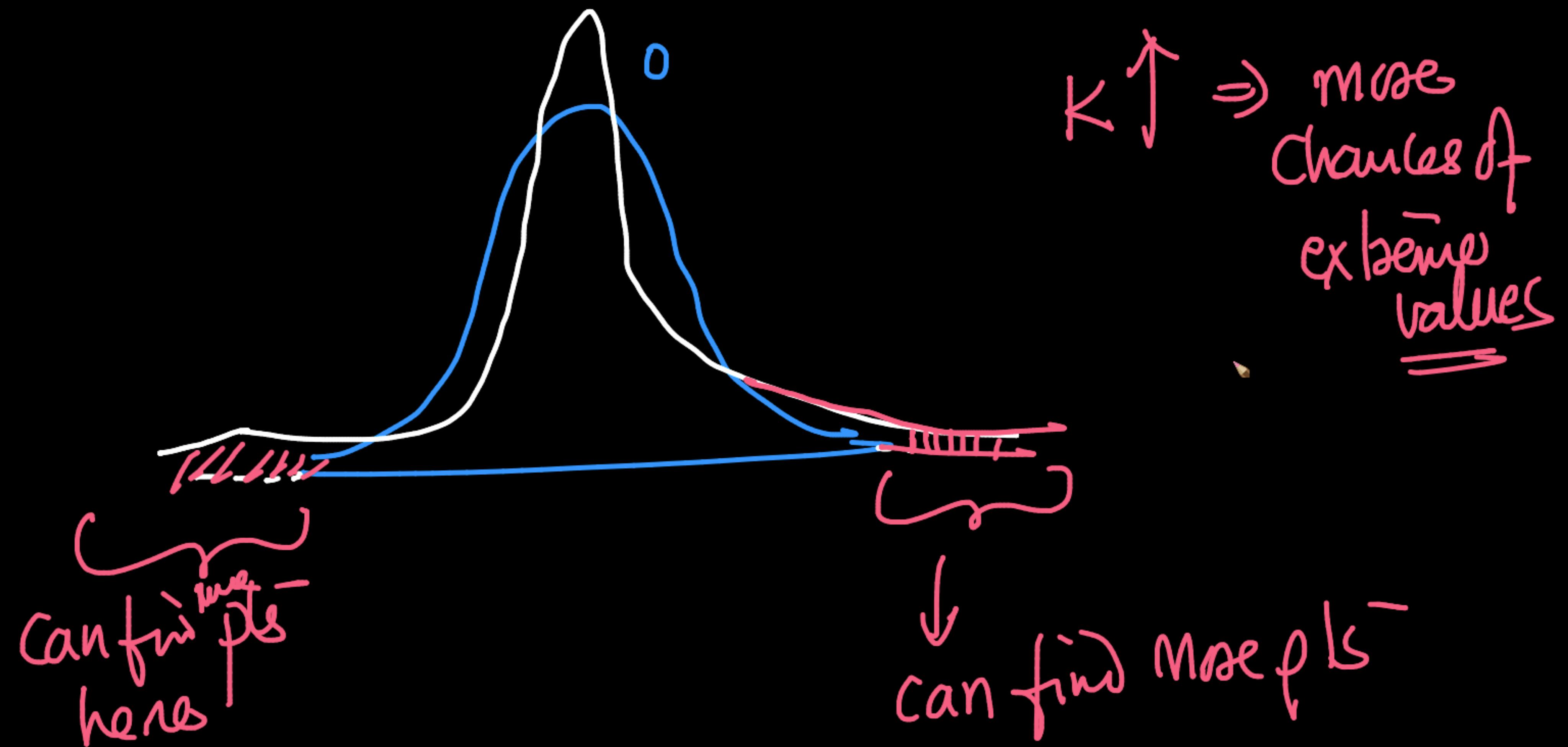
More details

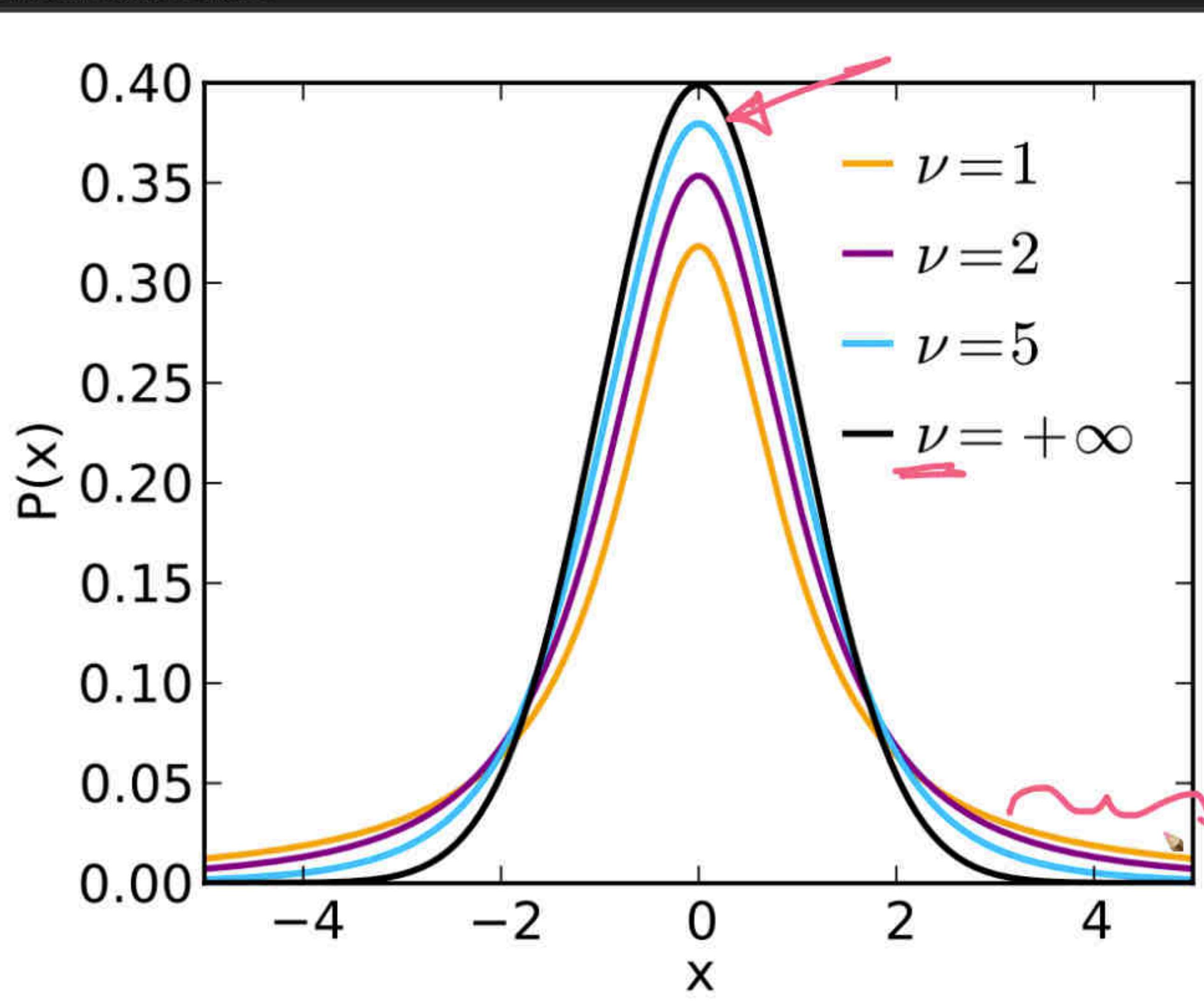


Probability density functions for selected distributions with mean 0, variance 1 and different excess kurtosis

More details







Plot of the density function for several members of the Student t family.

More details

$$\uparrow K = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \mu}{\sigma} \right)^4$$

$$\uparrow \text{Var} = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

→ 1

2

:

i

:

n

	$x_1 - \bar{x}_1$	$y_1 - \bar{y}_1$	(pairs)
2	$x_2 - \bar{x}_2$	$y_2 - \bar{y}_2$	
:	:	:	
i	$(x_i - \bar{x}_i)$	$(y_i - \bar{y}_i)$	
:	:	:	
n	$x_n - \bar{x}_n$	$y_n - \bar{y}_n$	

Data

$\mu_x = \text{Sample-mean of } x$

$\mu_y = \text{Sample Mean of } y$

Co-variance

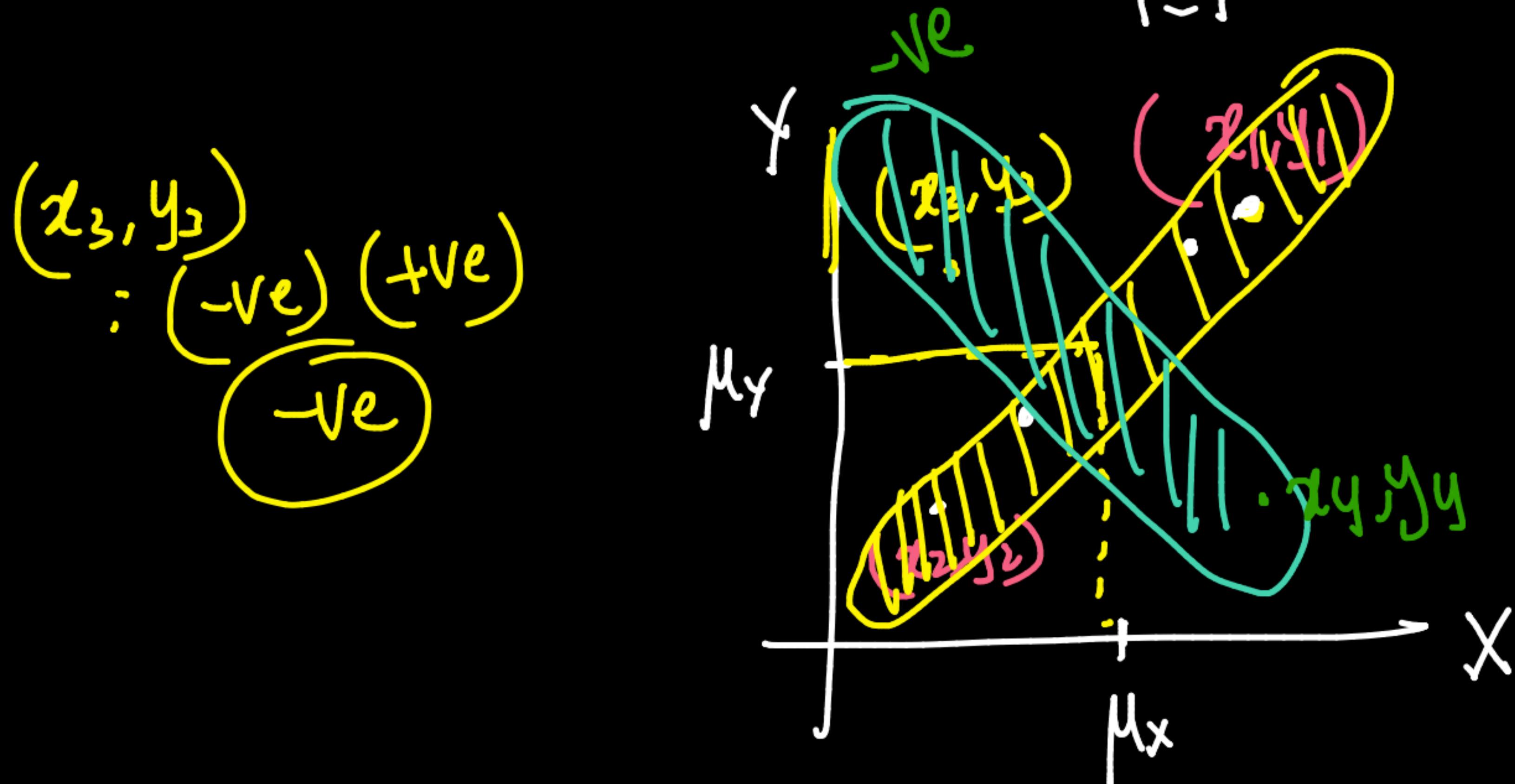
"'in general'"

$$\text{Co-var}(\underline{x}, \underline{y}) = \frac{1}{n} \sum_{i=1}^n (\underline{x}_i - \underline{\mu_x})(\underline{y}_i - \underline{\mu_y})$$

$$\text{Var}(x) = \frac{1}{n} \sum_{i=1}^n (\underline{x}_i - \underline{\mu_x})(\underline{x}_i - \underline{\mu_x})$$

$$\text{Cov}(x, x) = \text{Var}(x)$$

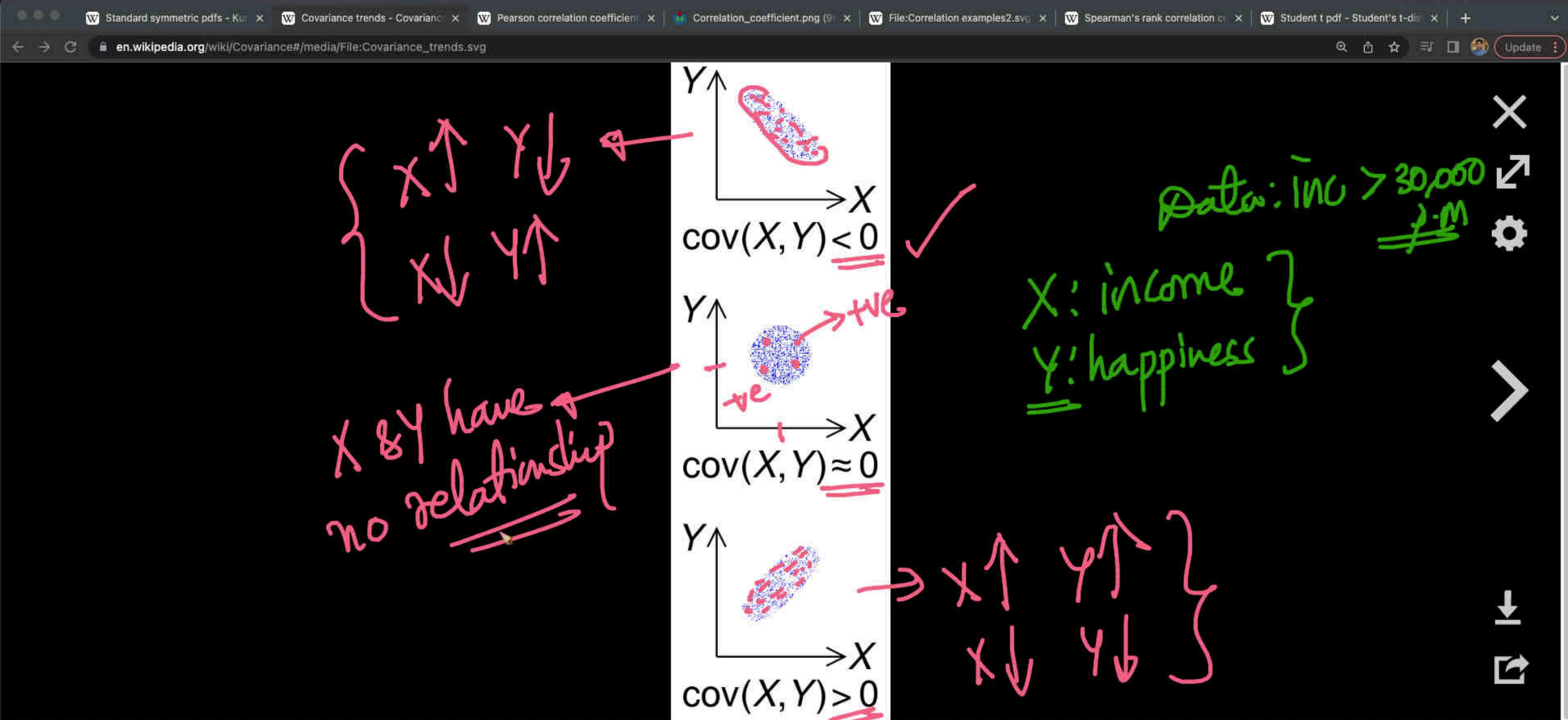
$$\text{Cov}(X,Y) = \frac{1}{n} \sum_{i=1}^n (\underline{x_i - \mu_x}) (\underline{y_i - \mu_y})$$



(x_1, y_1): +ve, +ve
↳ +ve

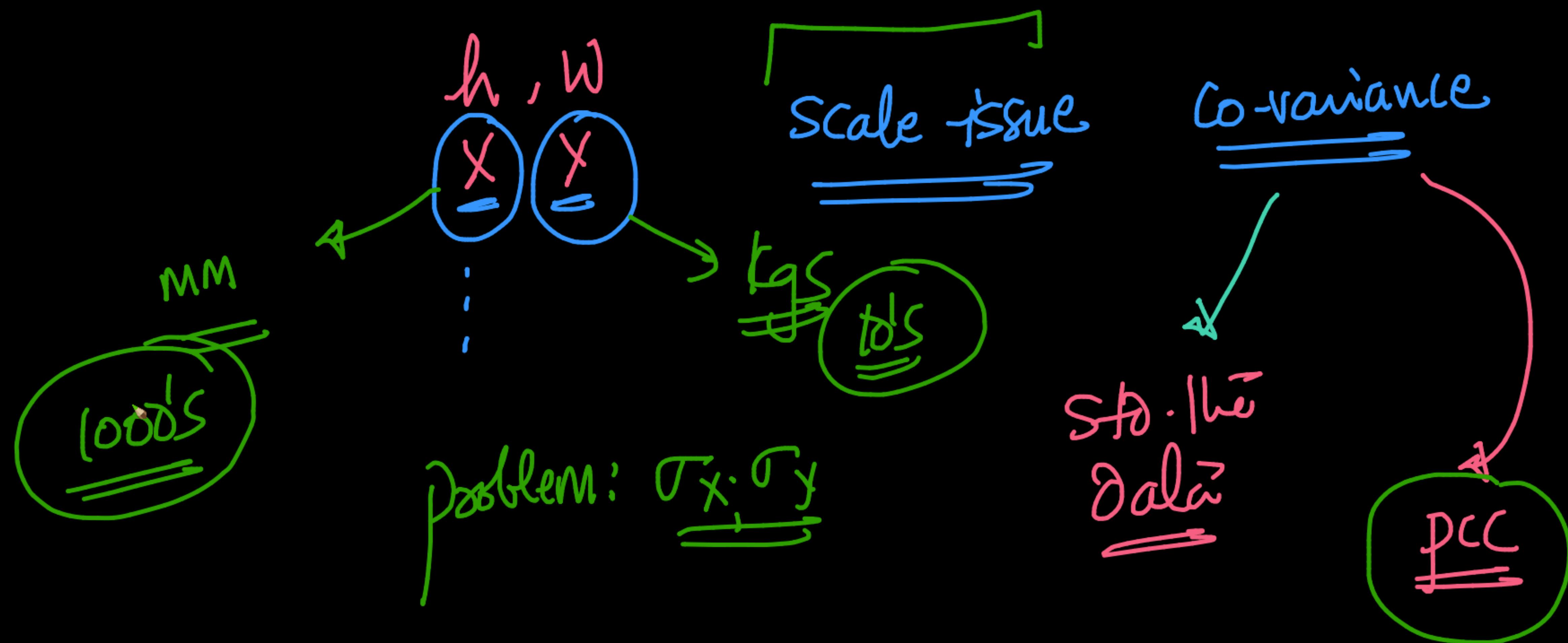
(x_2, y_2): -ve, -ve
(+ve) /

(x_3, y_3): -ve
↳ -ve



The sign of the covariance of two random variables X and Y

More details

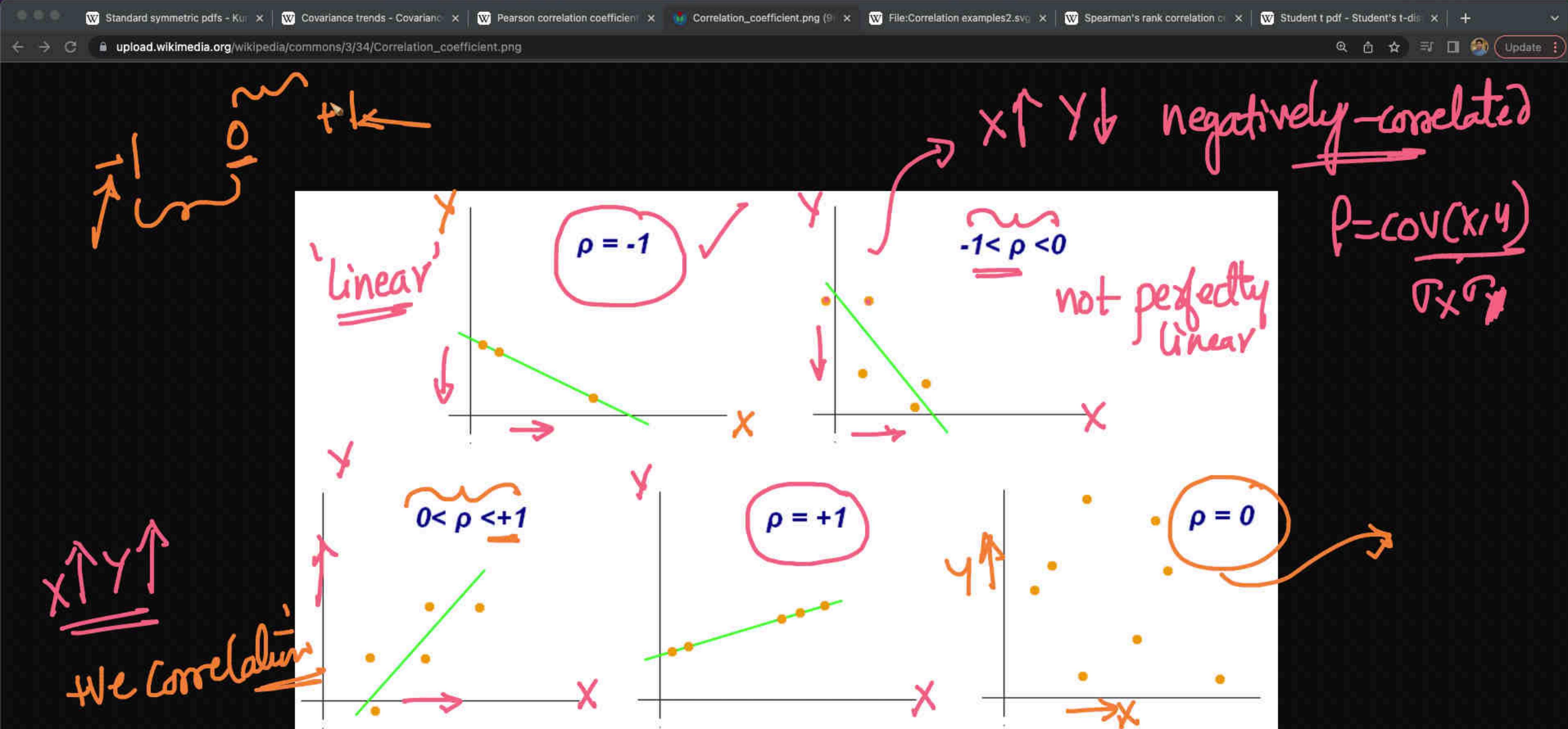


P.C.C
Very widely

$$\rho_{x,y} = \frac{\text{cov}(x,y)}{\sigma_x \sigma_y}$$

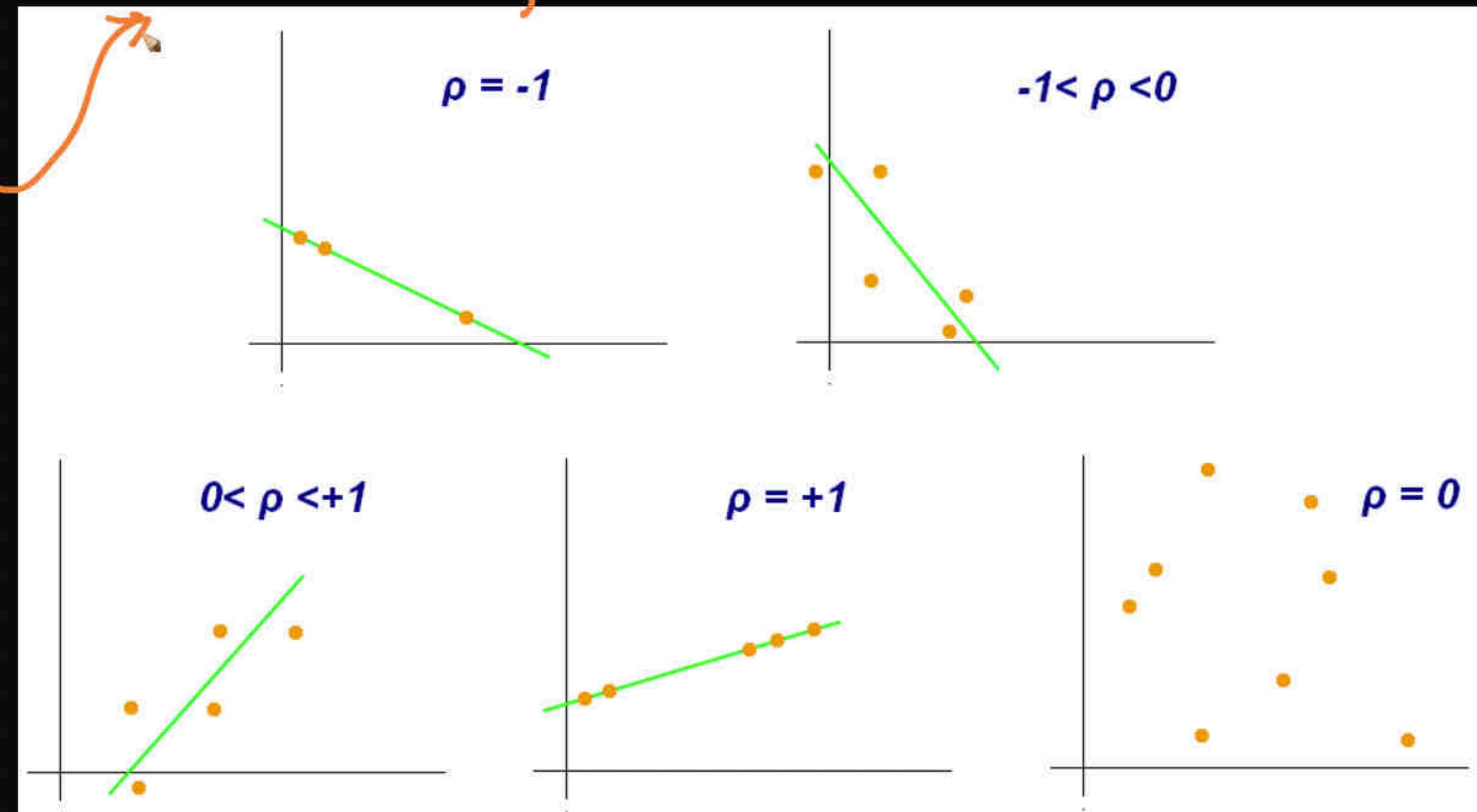
corrects issues with std.dev

[to be continued ...]



P.C.C → relationship is linear

e



WIKIPEDIA

The Free Encyclopedia

Main page

Contents

Current events

Random article

About Wikipedia

Contact us

Donate

Contribute

Help

Learn to edit

Community portal

Recent changes

Upload file

Tools

What links here

Special pages

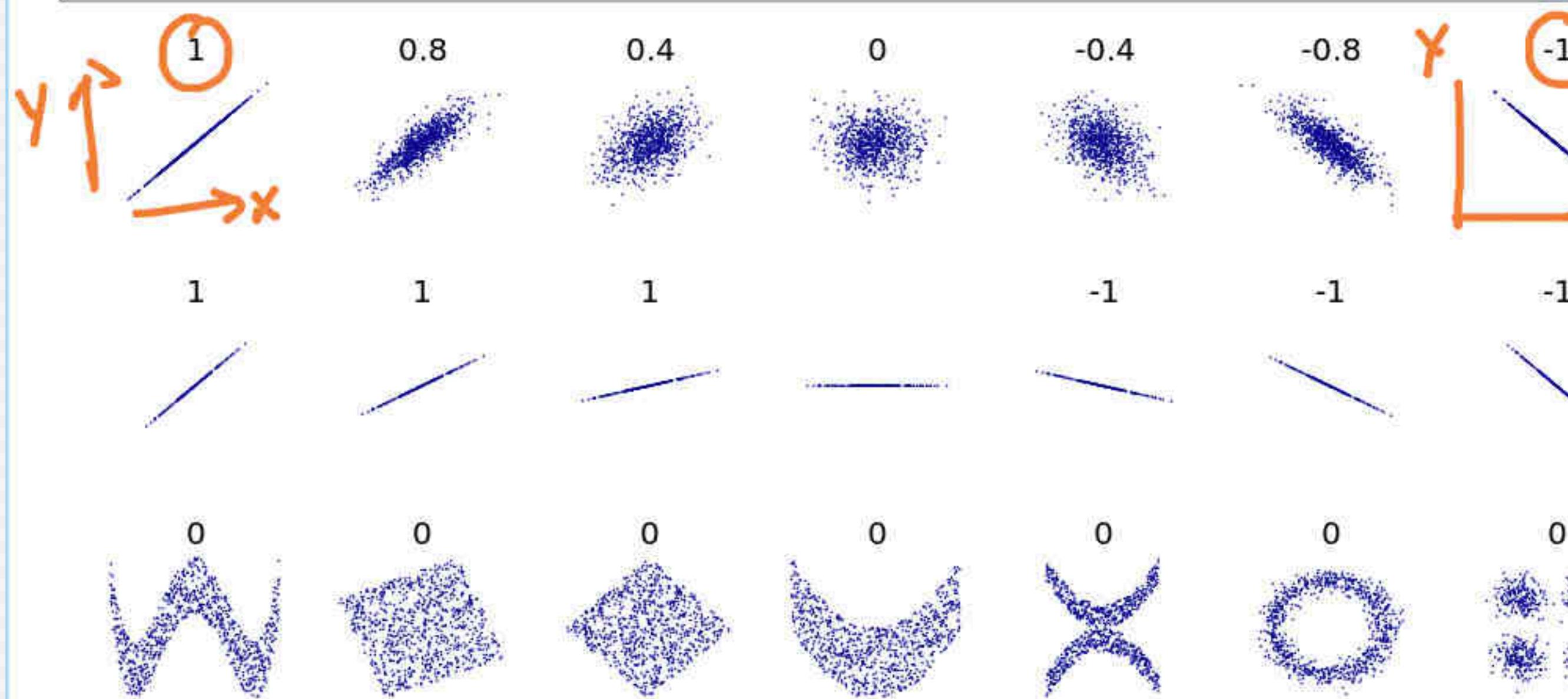
Printable version

Page information

File:Correlation examples2.svg

From Wikipedia, the free encyclopedia

File File history File usage Global file usage Metadata



Size of this PNG preview of this SVG file: 506 × 231 pixels. Other resolutions: 320 × 146 pixels | 640 × 292 pixels | 1,024 × 467 pixels | 1,280 × 584 pixels | 2,560 × 1,169 pixels.

Original file (SVG file, nominally 506 × 231 pixels, file size: 2.18 MB)



Open in Media Viewer



WIKIPEDIA

The Free Encyclopedia

Main page

Contents

Current events

Random article

About Wikipedia

Contact us

Donate

Contribute

Help

Learn to edit

Community portal

Recent changes

Upload file

Tools

What links here

Special pages

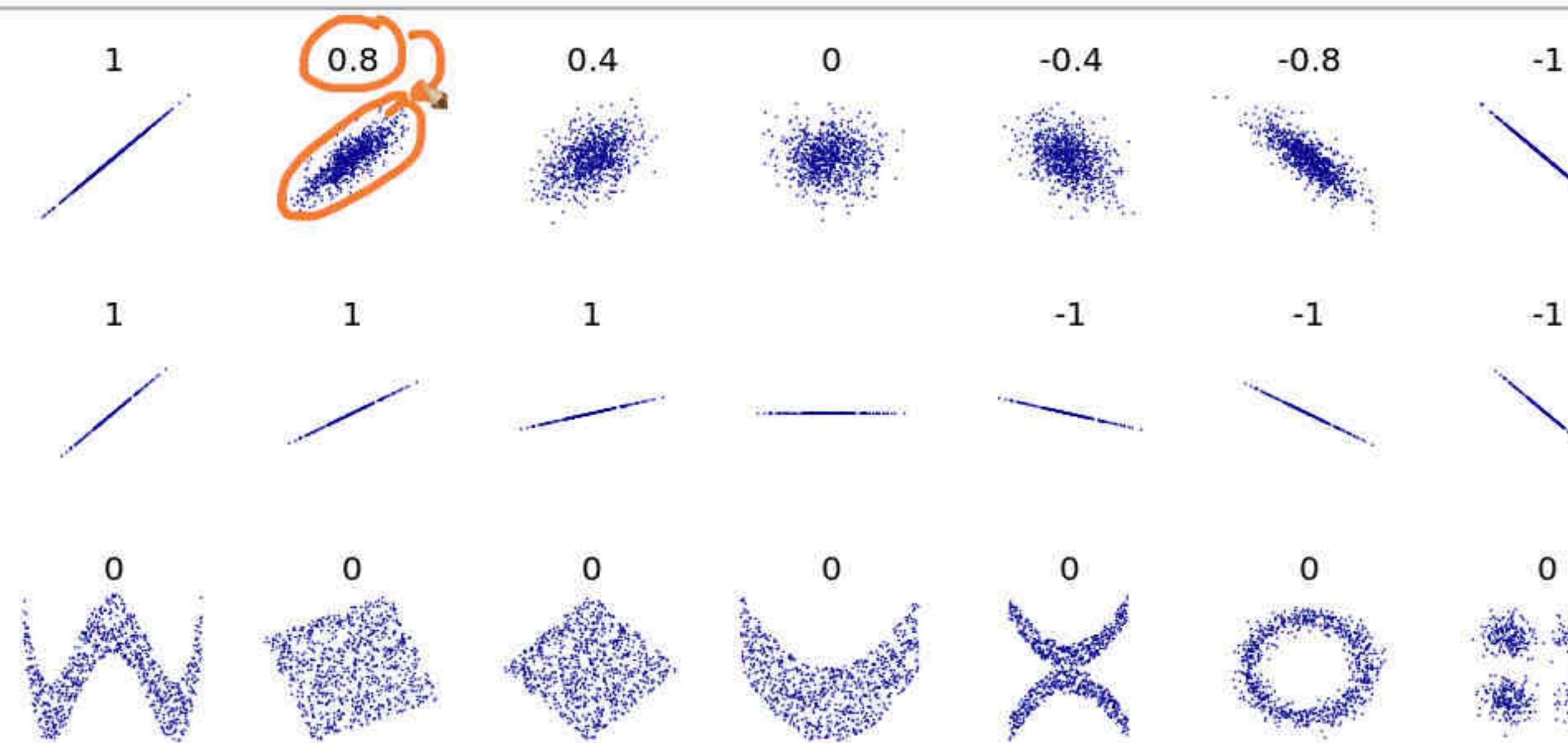
Printable version

Page information

File:Correlation examples2.svg

From Wikipedia, the free encyclopedia

File File history File usage Global file usage Metadata



Size of this PNG preview of this SVG file: 506 × 231 pixels. Other resolutions: 320 × 146 pixels | 640 × 292 pixels | 1,024 × 467 pixels | 1,280 × 584 pixels | 2,560 × 1,169 pixels.

Original file (SVG file, nominally 506 × 231 pixels, file size: 2.18 MB)

Open in Media Viewer



This



38 / 39

description page there is shown below.

WIKIPEDIA

The Free Encyclopedia

Main page

Contents

Current events

Random article

About Wikipedia

Contact us

Donate

Contribute

Help

Learn to edit

Community portal

Recent changes

Upload file

Tools

What links here

Special pages

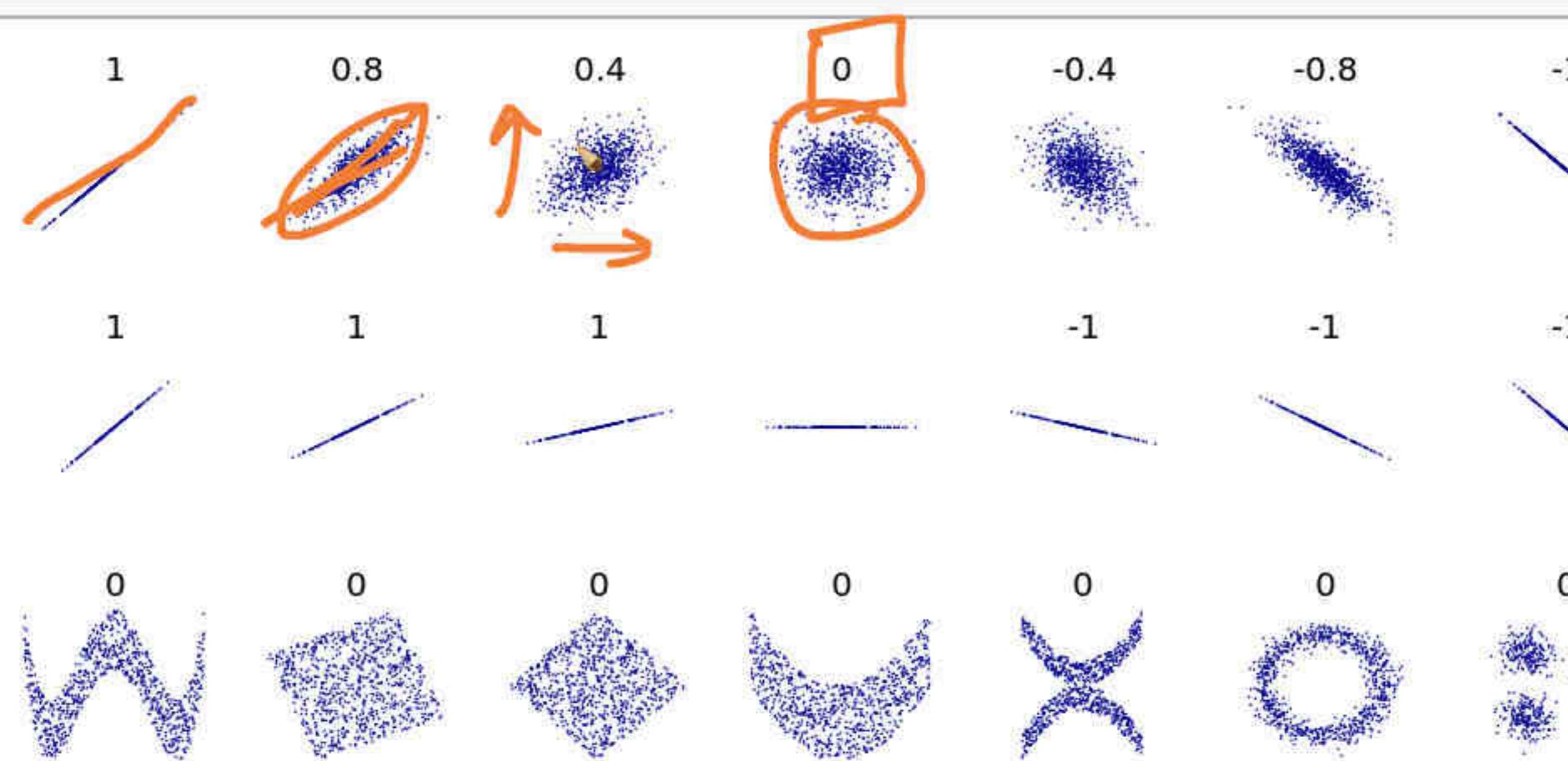
Printable version

Page information

File:Correlation examples2.svg

From Wikipedia, the free encyclopedia

File File history File usage Global file usage Metadata



Size of this PNG preview of this SVG file: 506 × 231 pixels. Other resolutions: 320 × 146 pixels | 640 × 292 pixels | 1,024 × 467 pixels | 1,280 × 584 pixels | 2,560 × 1,169 pixels.

Original file (SVG file, nominally 506 × 231 pixels, file size: 2.18 MB)



Open in Media Viewer



WIKIPEDIA

The Free Encyclopedia

Main page

Contents

Current events

Random article

About Wikipedia

Contact us

Donate

Contribute

Help

Learn to edit

Community portal

Recent changes

Upload file

Tools

What links here

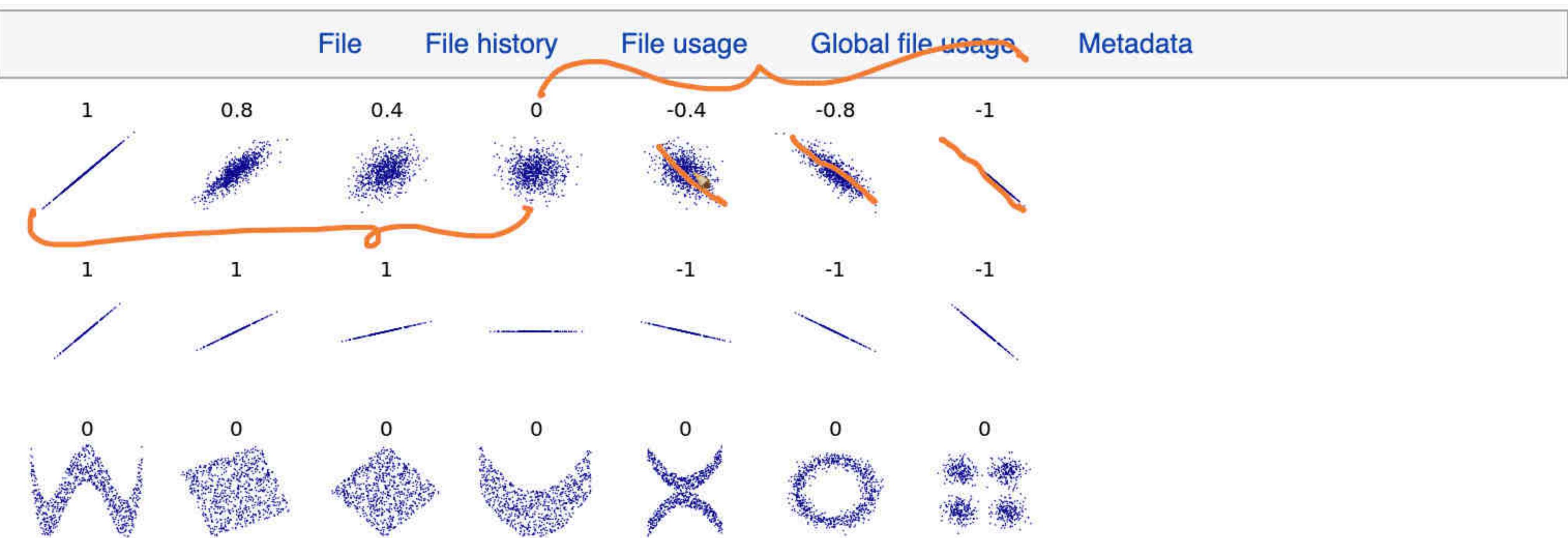
Special pages

Printable version

Page information

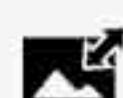
File:Correlation examples2.svg

From Wikipedia, the free encyclopedia



Size of this PNG preview of this SVG file: 506 × 231 pixels. Other resolutions: 320 × 146 pixels | 640 × 292 pixels | 1,024 × 467 pixels | 1,280 × 584 pixels | 2,560 × 1,169 pixels.

Original file (SVG file, nominally 506 × 231 pixels, file size: 2.18 MB)



Open in Media Viewer



WIKIPEDIA

The Free Encyclopedia

Main page

Contents

Current events

Random article

About Wikipedia

Contact us

Donate

Contribute

Help

Learn to edit

Community portal

Recent changes

Upload file

Tools

What links here

Special pages

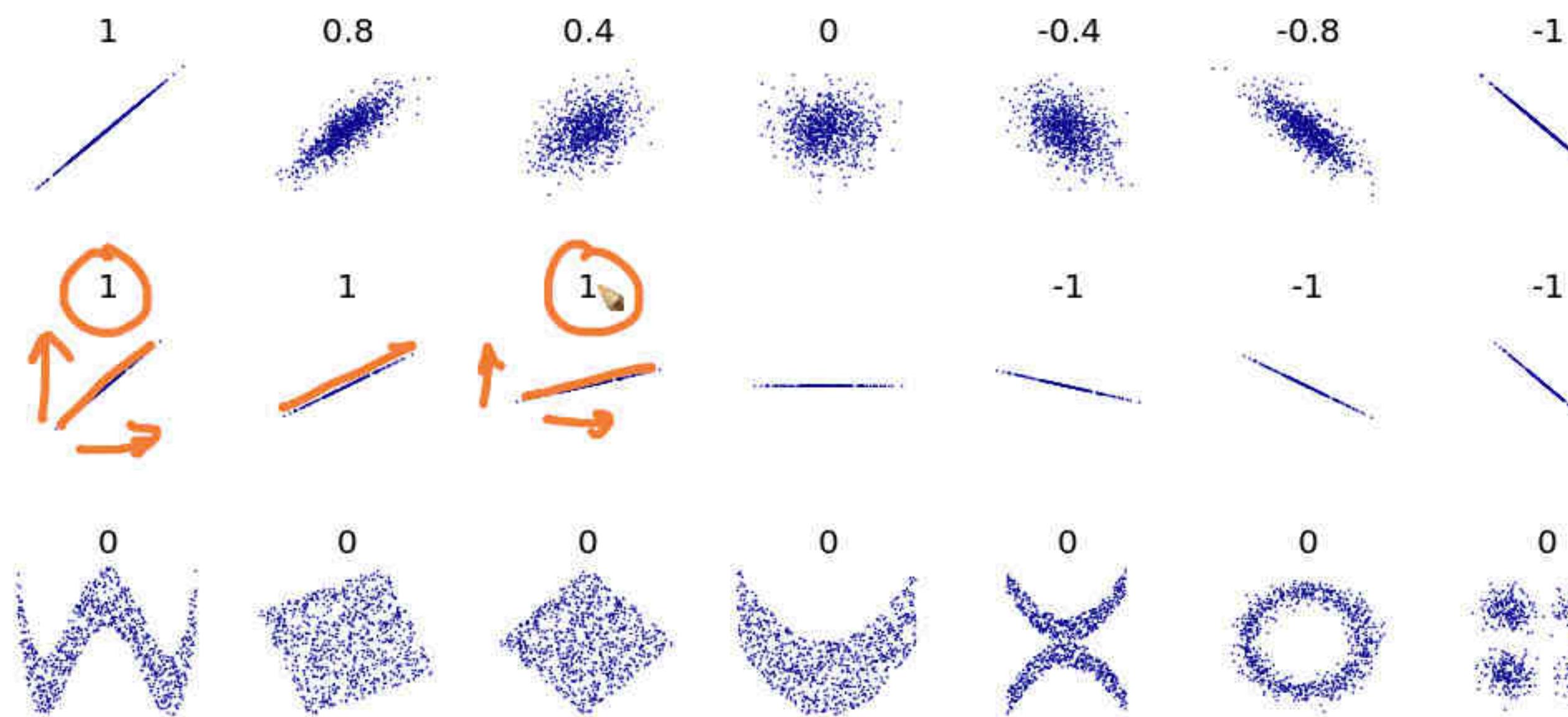
Printable version

Page information

File:Correlation examples2.svg

From Wikipedia, the free encyclopedia

[File](#) [File history](#) [File usage](#) [Global file usage](#) [Metadata](#)



Size of this PNG preview of this SVG file: 506 × 231 pixels. Other resolutions: 320 × 146 pixels | 640 × 292 pixels | 1,024 × 467 pixels | 1,280 × 584 pixels | 2,560 × 1,169 pixels.

Original file (SVG file, nominally 506 × 231 pixels, file size: 2.18 MB)



[Open in Media Viewer](#)



WIKIPEDIA

The Free Encyclopedia

Main page

Contents

Current events

Random article

About Wikipedia

Contact us

Donate

Contribute

Help

Learn to edit

Community portal

Recent changes

Upload file

Tools

What links here

Special pages

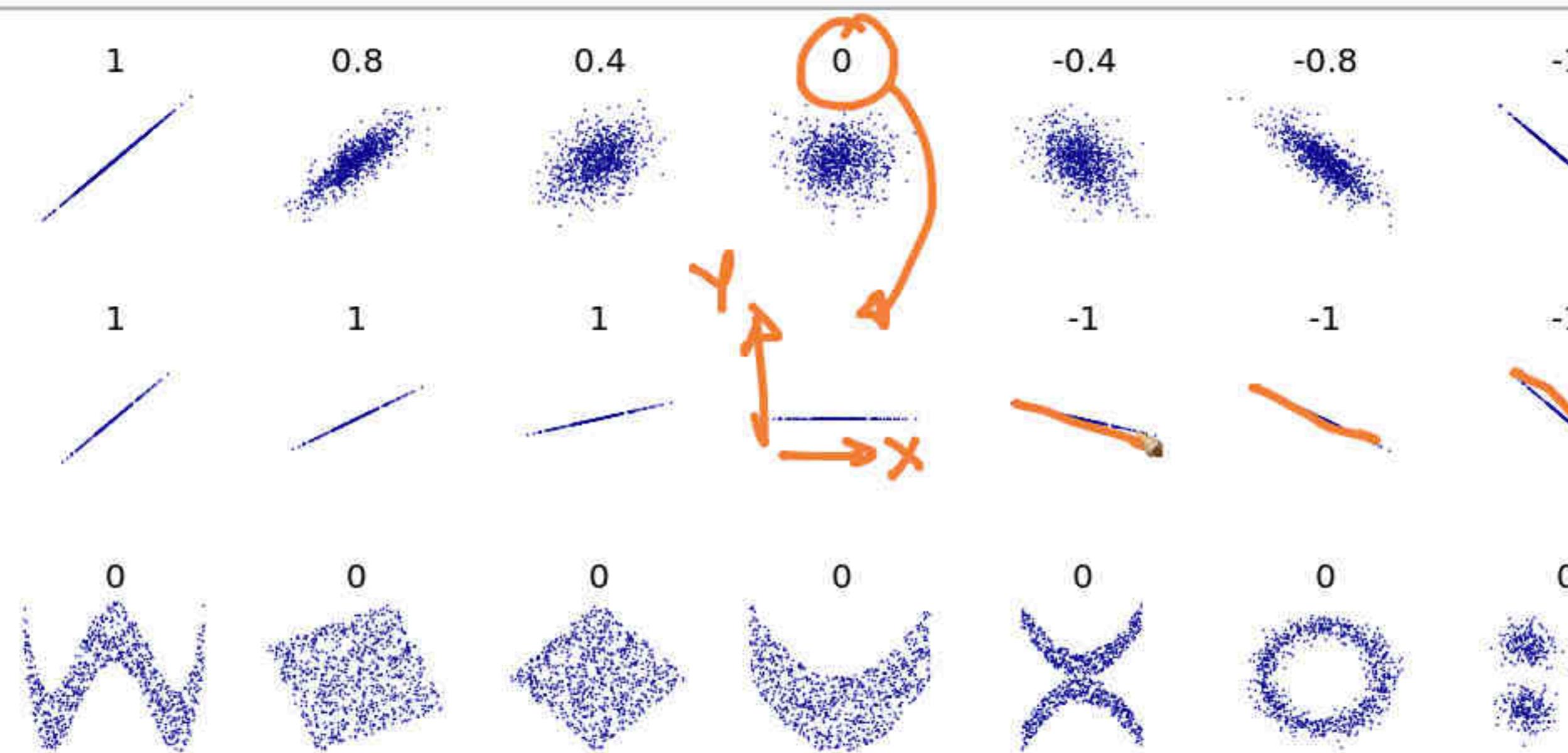
Printable version

Page information

File:Correlation examples2.svg

From Wikipedia, the free encyclopedia

[File](#) [File history](#) [File usage](#) [Global file usage](#) [Metadata](#)



Size of this PNG preview of this SVG file: 506 × 231 pixels. Other resolutions: 320 × 146 pixels | 640 × 292 pixels | 1,024 × 467 pixels | 1,280 × 584 pixels | 2,560 × 1,169 pixels.

Original file (SVG file, nominally 506 × 231 pixels, file size: 2.18 MB)



[Open in Media Viewer](#)



WIKIPEDIA

The Free Encyclopedia

Main page

Contents

Current events

Random article

About Wikipedia

Contact us

Donate

Contribute

Help

Learn to edit

Community portal

Recent changes

Upload file

Tools

What links here

Special pages

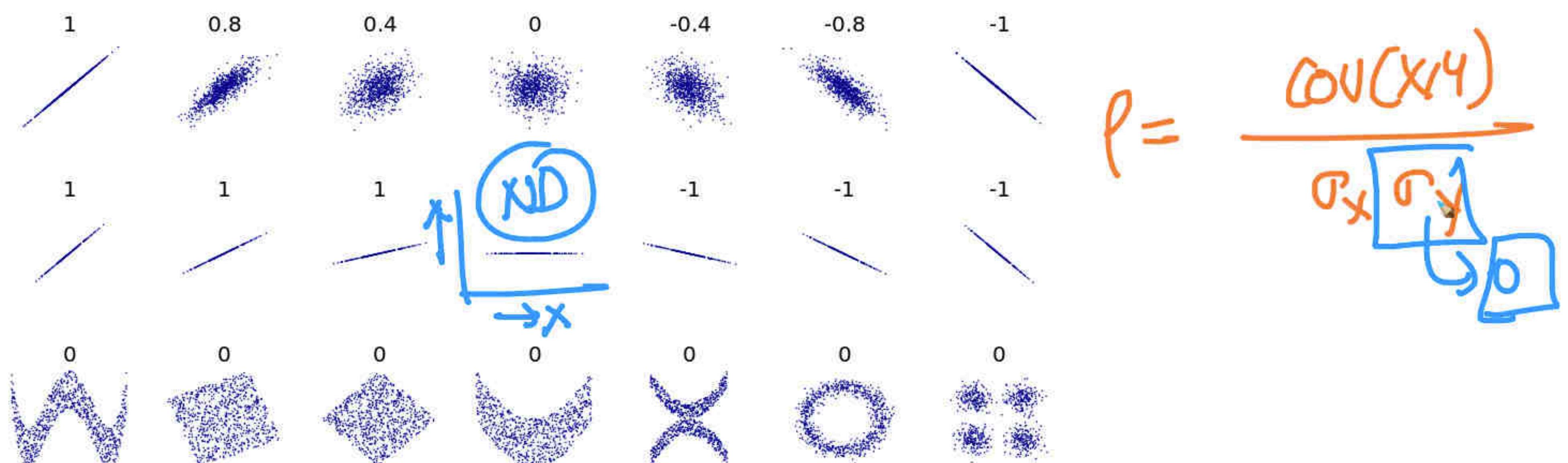
Printable version

Page information

File:Correlation examples2.svg

From Wikipedia, the free encyclopedia

File File history File usage Global file usage Metadata



Size of this PNG preview of this SVG file: 506 × 231 pixels. Other resolutions: 320 × 146 pixels | 640 × 292 pixels | 1,024 × 467 pixels | 1,280 × 584 pixels | 2,560 × 1,169 pixels.

Original file (SVG file, nominally 506 × 231 pixels, file size: 2.18 MB)



Open in Media Viewer



WIKIPEDIA

The Free Encyclopedia

Main page

Contents

Current events

Random article

About Wikipedia

Contact us

Donate

Contribute

Help

Learn to edit

Community portal

Recent changes

Upload file

Tools

What links here

Special pages

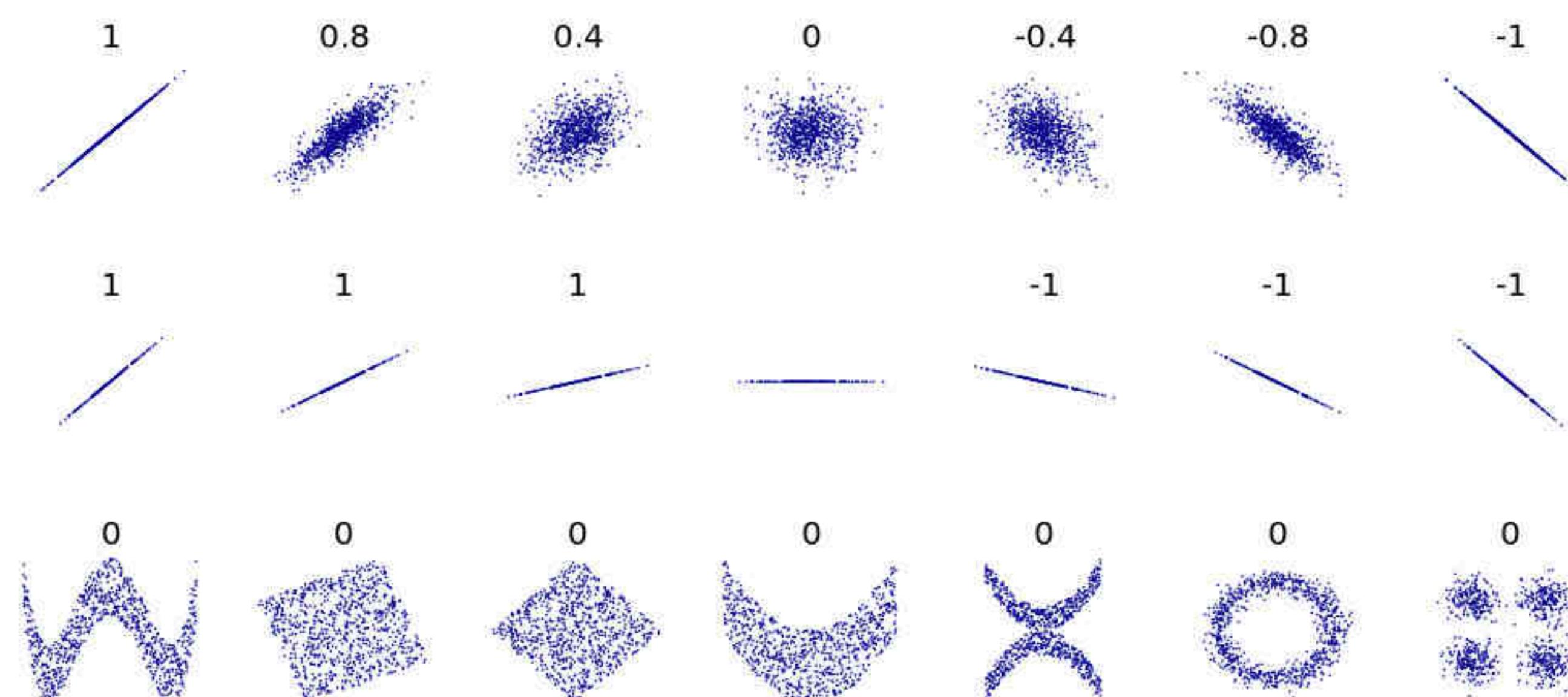
Printable version

Page information

File:Correlation examples2.svg

From Wikipedia, the free encyclopedia

[File](#) [File history](#) [File usage](#) [Global file usage](#) [Metadata](#)



Size of this PNG preview of this SVG file: 506 × 231 pixels. Other resolutions: 320 × 146 pixels | 640 × 292 pixels | 1,024 × 467 pixels | 1,280 × 584 pixels | 2,560 × 1,169 pixels.

[Original file \(SVG file, nominally 506 × 231 pixels, file size: 2.18 MB\)](#)

[Open in Media Viewer](#)



WIKIPEDIA

The Free Encyclopedia

Main page

Contents

Current events

Random article

About Wikipedia

Contact us

Donate

Contribute

Help

Learn to edit

Community portal

Recent changes

Upload file

Tools

What links here

Special pages

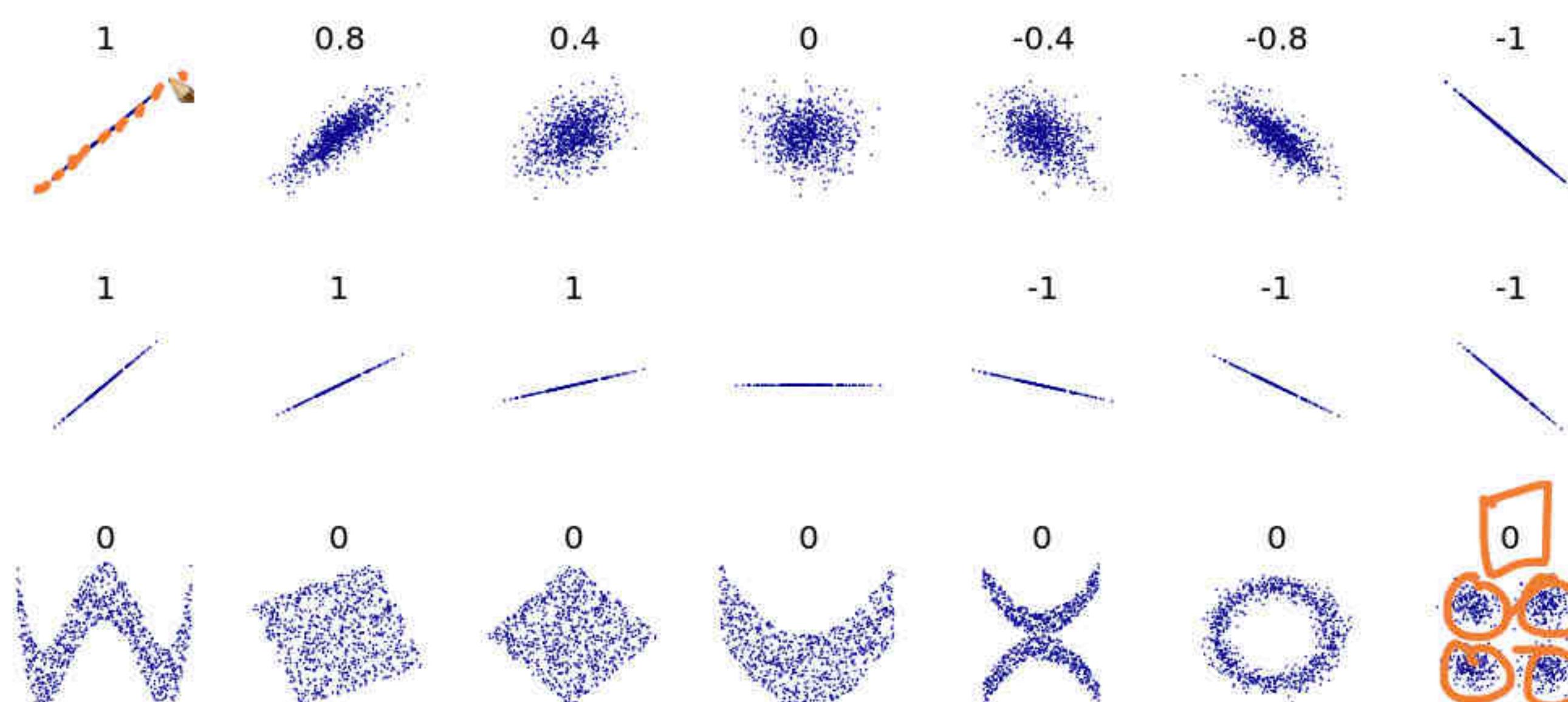
Printable version

Page information

File:Correlation examples2.svg

From Wikipedia, the free encyclopedia

File File history File usage Global file usage Metadata



Size of this PNG preview of this SVG file: 506 × 231 pixels. Other resolutions: 320 × 146 pixels | 640 × 292 pixels | 1,024 × 467 pixels | 1,280 × 584 pixels | 2,560 × 1,169 pixels.

Original file (SVG file, nominally 506 × 231 pixels, file size: 2.18 MB)



Open in Media Viewer



WIKIPEDIA

The Free Encyclopedia

Main page

Contents

Current events

Random article

About Wikipedia

Contact us

Donate

Contribute

Help

Learn to edit

Community portal

Recent changes

Upload file

Tools

What links here

Special pages

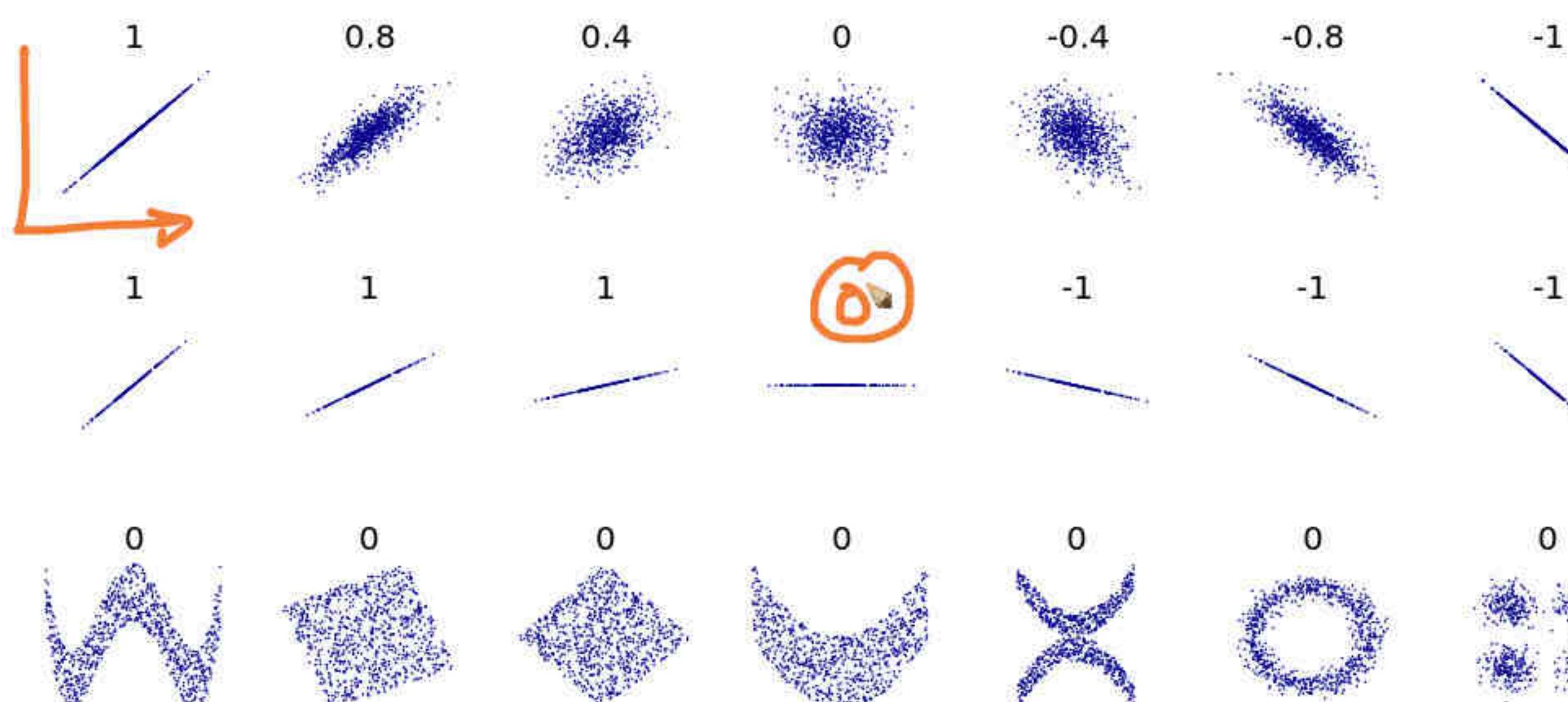
Printable version

Page information

File:Correlation examples2.svg

From Wikipedia, the free encyclopedia

File File history File usage Global file usage Metadata



Size of this PNG preview of this SVG file: 506 × 231 pixels. Other resolutions: 320 × 146 pixels | 640 × 292 pixels | 1,024 × 467 pixels | 1,280 × 584 pixels | 2,560 × 1,169 pixels.

Original file (SVG file, nominally 506 × 231 pixels, file size: 2.18 MB)

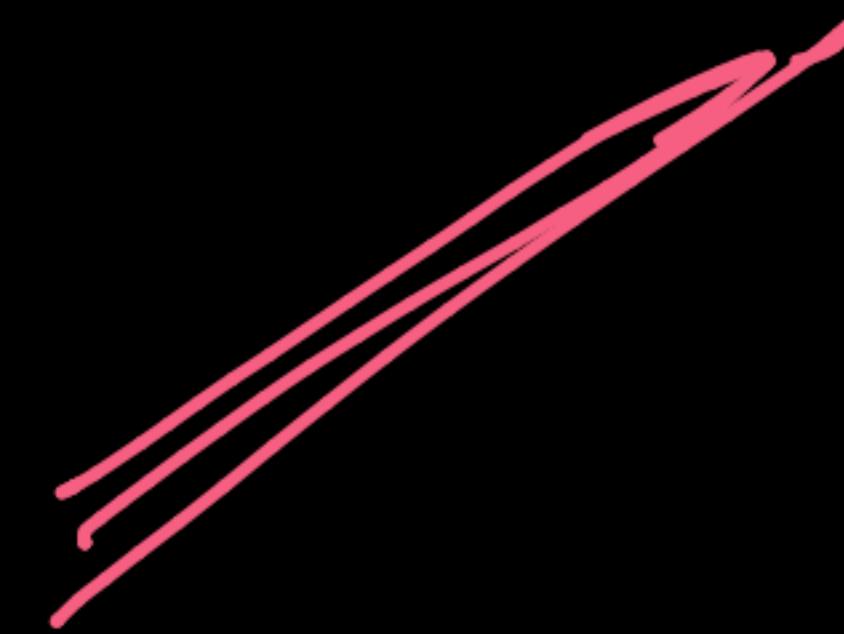


Open in Media Viewer



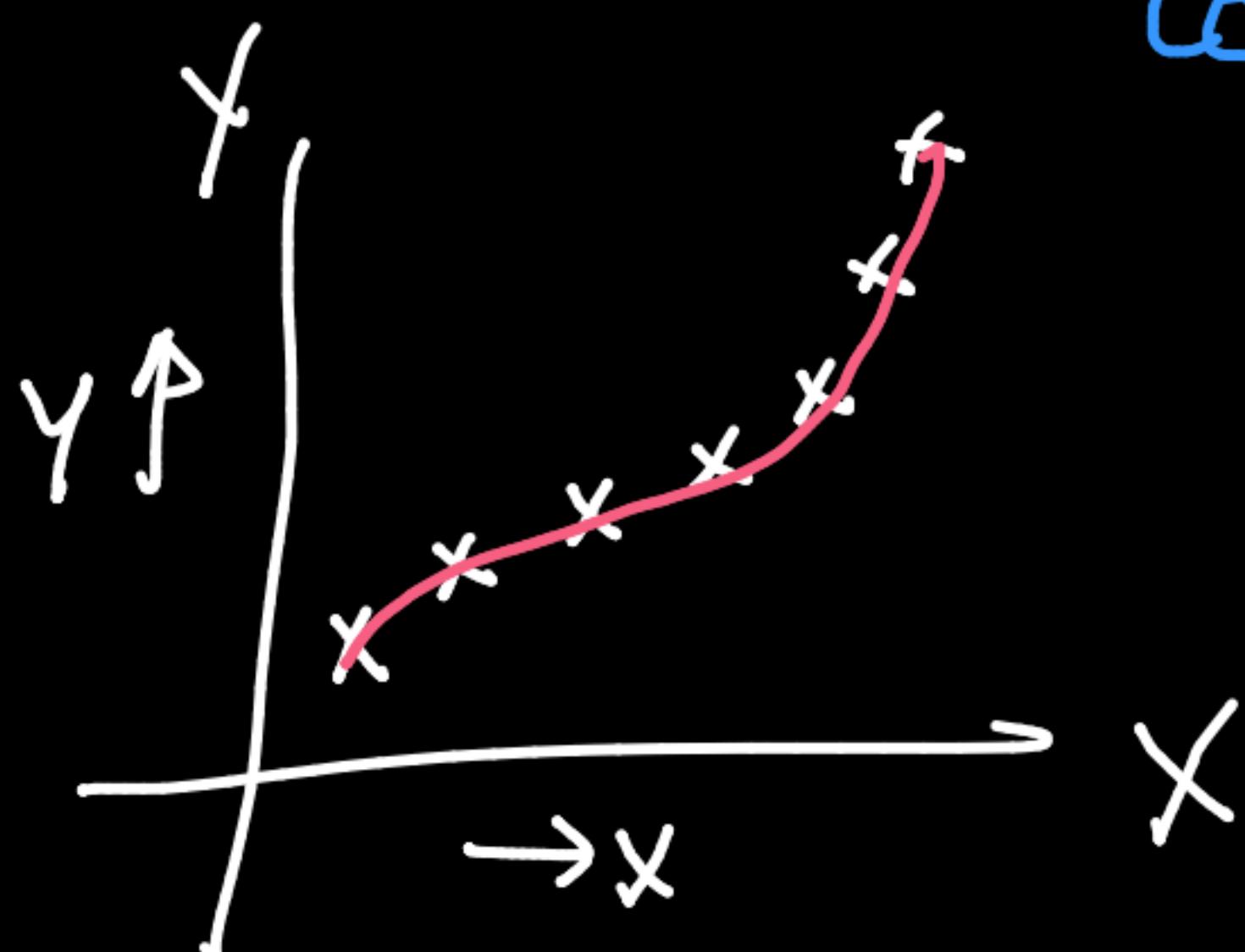




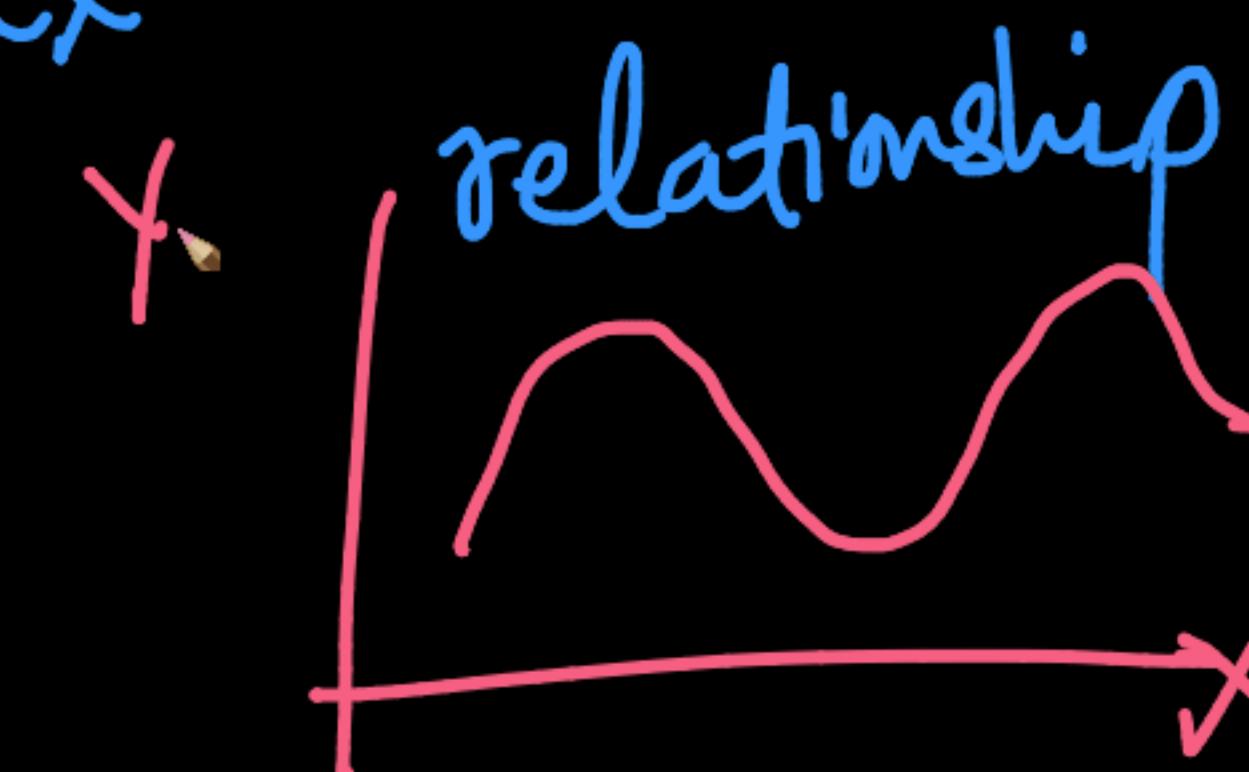


Linear correlation \rightarrow PCC

Monotonic relationship \rightarrow hack in PCC



Complex non-linear relationship



\rightarrow ML-algos

GBDT / RF / DL
(lately)

non-gaussian

boxcox
→ gaussian
log
sqrt
f()

Spearmann- rank- correlation:

D :

X	Y
10	20
12	30
8	12
6	4
5	1

Sort by
Sort X's:
Sort Y's:

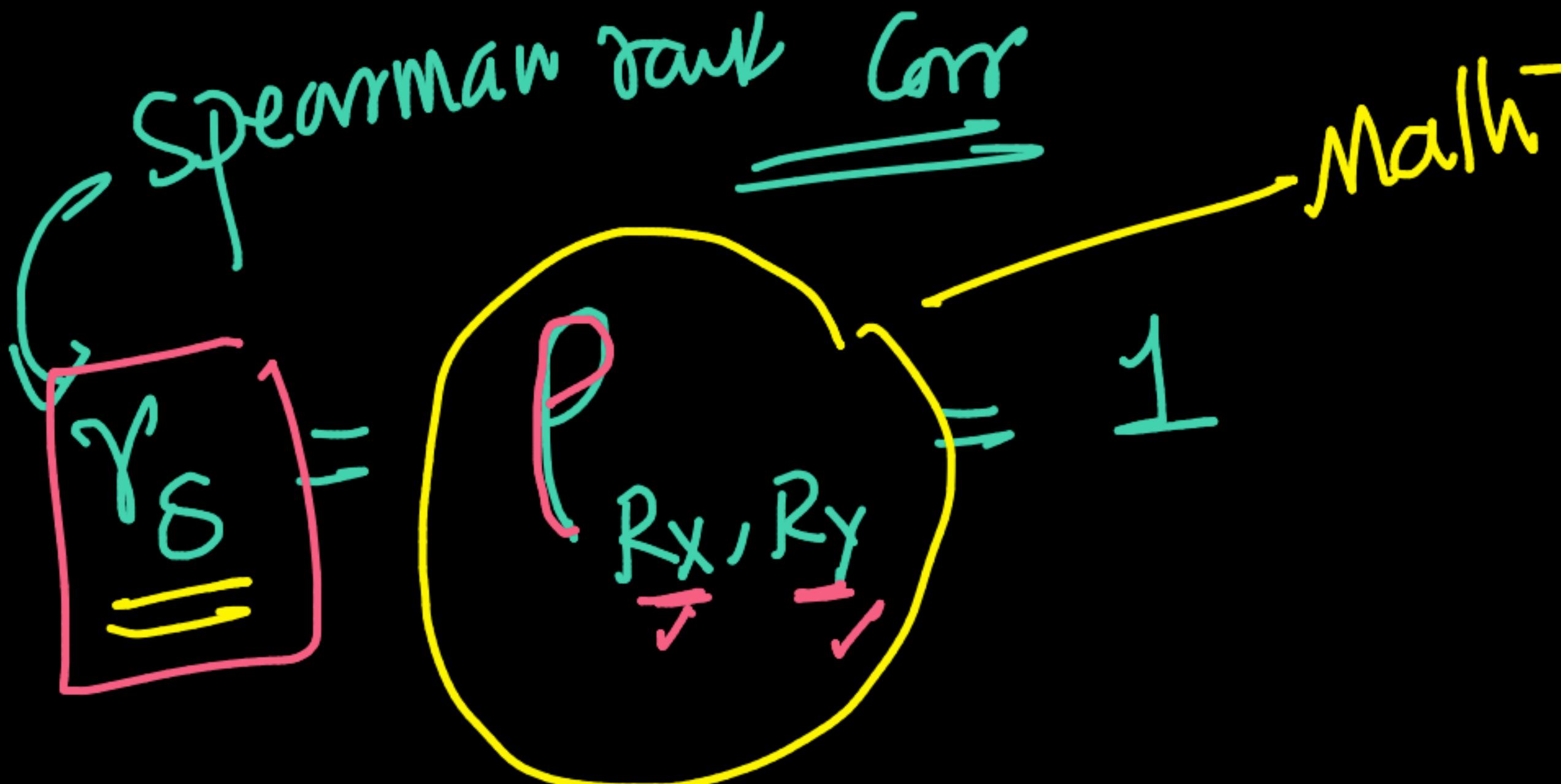
5, 6, 8, 10, 12
1 2 3 4 5

1, 4, 12, 20, 30
1 2 3 4 5

(monotonic relationship)

D'

R_X	R_Y
4	4
5	5
3	3
2	2
1	1



rank-transform

ranks
Converts
monotonicity
to linearity

ranks nullify scale

A hand-drawn diagram showing the process of ranks nullifying scale. It starts with a blue asterisk (*) pointing to a bracket that encloses the words 'Capture monotonicity'. An arrow points from this bracket down to the word 'linearity'. To the right of the bracket is a circled 'P'.

Spearman's rank correlation coefficient

From Wikipedia, the free encyclopedia

In statistics, **Spearman's rank correlation coefficient** or **Spearman's ρ** , named after Charles Spearman and often denoted by the Greek letter ρ (rho) or as r_s , is a nonparametric measure of rank correlation (statistical dependence between the rankings of two variables). It assesses how well the relationship between two variables can be described using a monotonic function.

The Spearman correlation between two variables is equal to the Pearson correlation between the rank values of those two variables; while Pearson's correlation assesses linear relationships, Spearman's correlation assesses monotonic relationships (whether linear or not). If there are no repeated data values, a perfect Spearman correlation of +1 or -1 occurs when each of the variables is a perfect monotone function of the other.

Intuitively, the Spearman correlation between two variables will be high when observations have a similar (or identical for a correlation of 1) rank (i.e. re

rank - transform

A Spearman correlation of 1 results when the two variables being compared are monotonically related, even if their relationship is not linear. This means that all data points with greater x values than that of a given data point will have greater y values as well. In

γ_s over P

Monotonically \rightarrow linearly
(conceptually)

Standard symmetric pdfs - Kur... Covariance trends - Covarian... Pearson correlation coefficient Correlation_coefficient.png (9) File:Correlation examples2.sv Spearman's rank correlation c... Student t pdf - Student's t-d...

Permanent link

Page information

Cite this page

Wikidata item

Print/export

Download as PDF

Printable version

In other projects

Wikimedia Commons



Languages

العربية

Español

Français

한국어

हिन्दी

日本語

★ Polski

Português

中文

文 18 more

Edit links

variable. 1st, 2nd, 3rd, etc.) between the two variables, and how when observations have a dissimilar (or fully opposed for a correlation of -1) rank between the two variables.

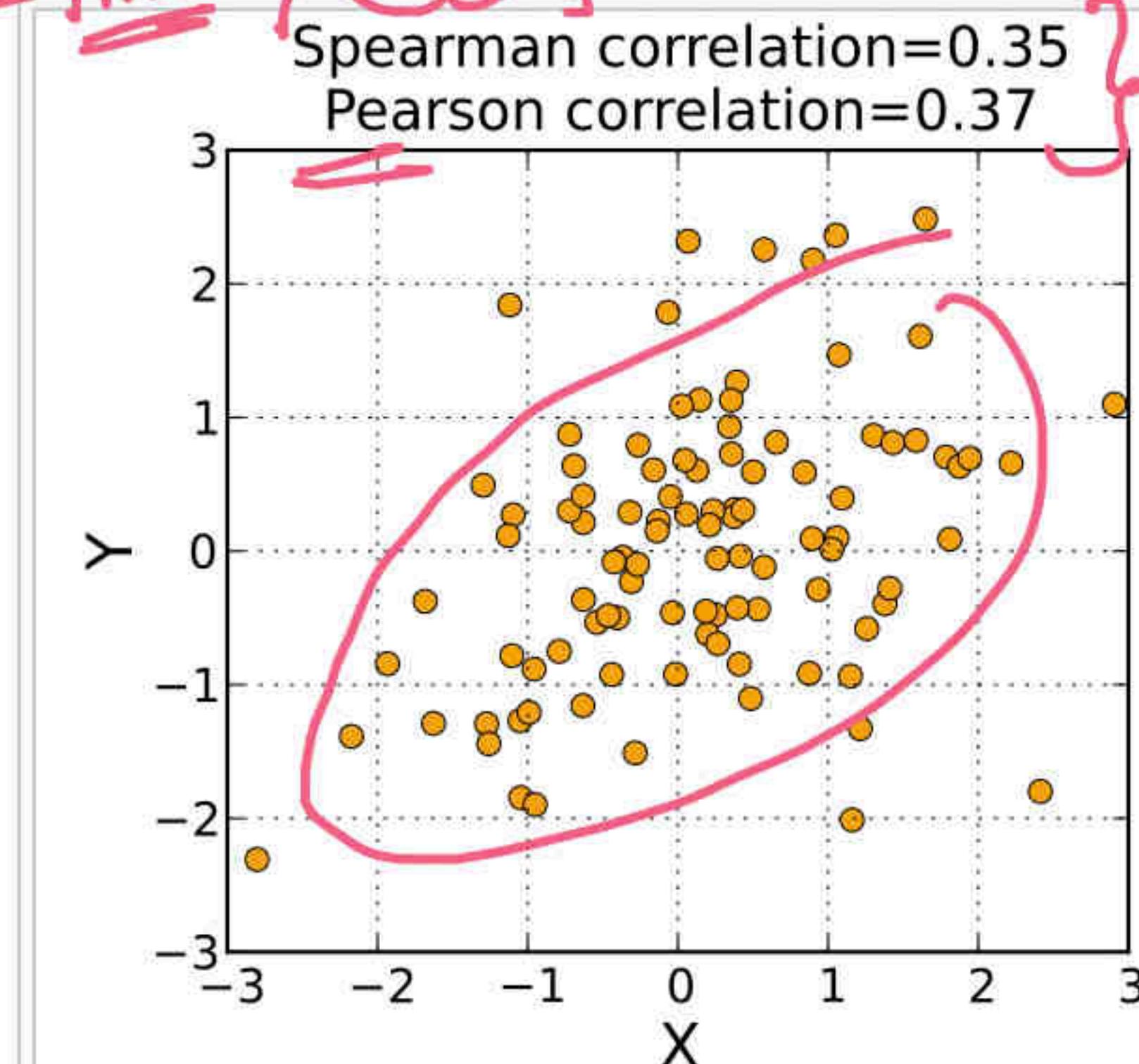
Spearman's coefficient is appropriate for both [continuous](#) and discrete [ordinal variables](#).^{[1][2]} Both Spearman's ρ and Kendall's τ can be formulated as special cases of a more [general correlation coefficient](#).

Contents [hide]

- 1 Definition and calculation
- 2 Related quantities
- 3 Interpretation
- 4 Example
- 5 Determining significance
- 6 Correspondence analysis based on Spearman's ρ
- 7 Approximating Spearman's ρ from a stream
- 8 Software implementations
- 9 See also
- 10 References
- 11 Further reading
- 12 External links

contrast, this does not give a perfect Pearson correlation.

More time



When the data are roughly elliptically distributed and there are no prominent outliers, the Spearman correlation and Pearson correlation give similar values.

Spearman correlation=0.84

Português

中文

文 A 18 more

Edit links

[9 See also](#)[10 References](#)[11 Further reading](#)[12 External links](#)

Definition and calculation [edit]

The Spearman correlation coefficient is defined as the Pearson correlation coefficient between the rank variables.^[3]

For a sample of size n , the n raw scores X_i, Y_i are converted to ranks $R(X_i), R(Y_i)$, and r_s is computed as

$$r_s = \rho_{R(X), R(Y)} = \frac{\text{cov}(R(X), R(Y))}{\sigma_{R(X)} \sigma_{R(Y)}},$$

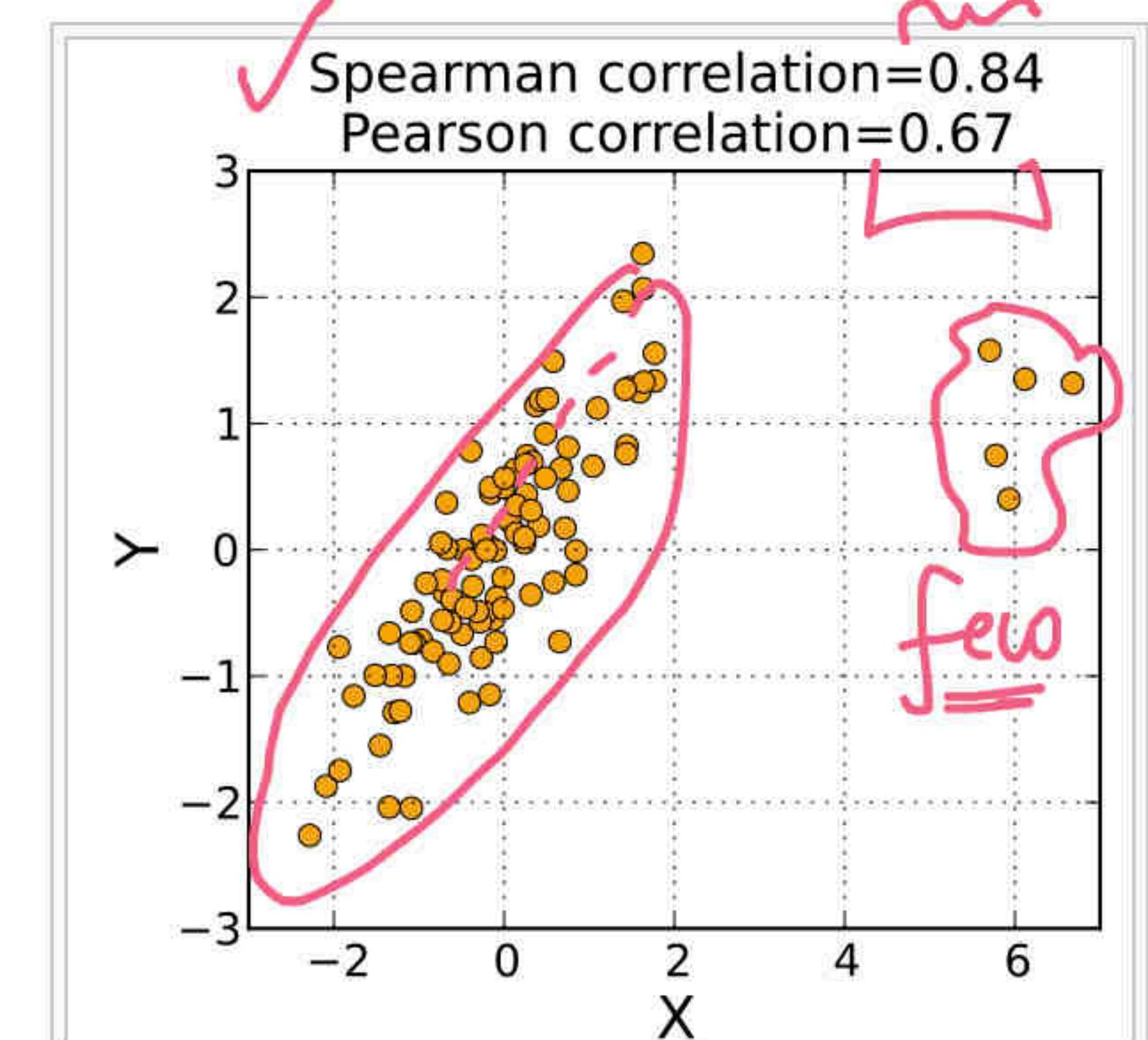
where

ρ denotes the usual Pearson correlation coefficient, but applied to the rank variables,

$\text{cov}(R(X), R(Y))$ is the covariance of the rank variables,

$\sigma_{R(X)}$ and $\sigma_{R(Y)}$ are the standard deviations of the rank variables.

and there are no prominent outliers, the Spearman correlation and Pearson correlation give similar values.

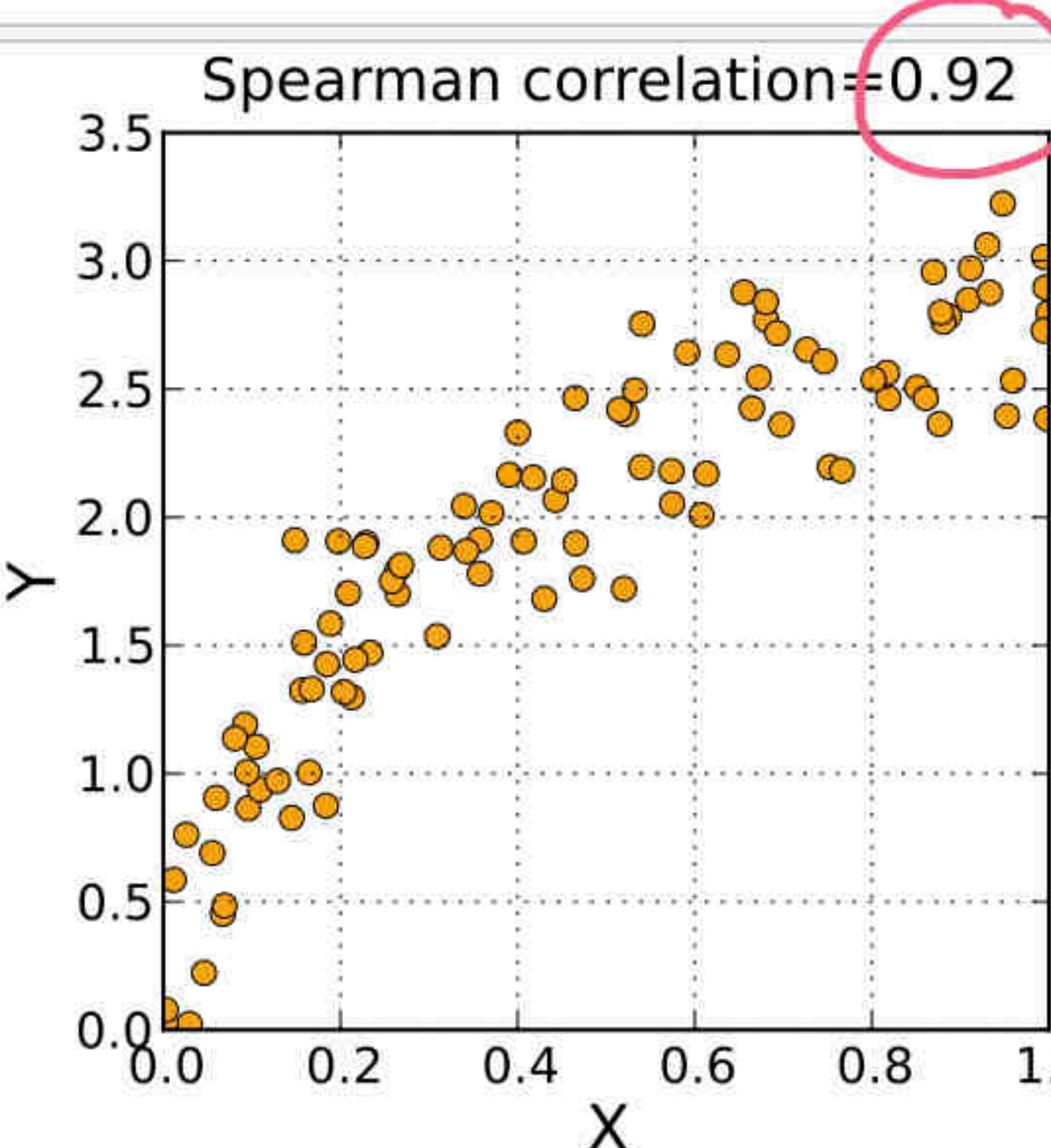


The Spearman correlation is less sensitive than the Pearson correlation to strong outliers that are in the tails of both samples. That is because Spearman's ρ limits the outlier to the value of its rank.

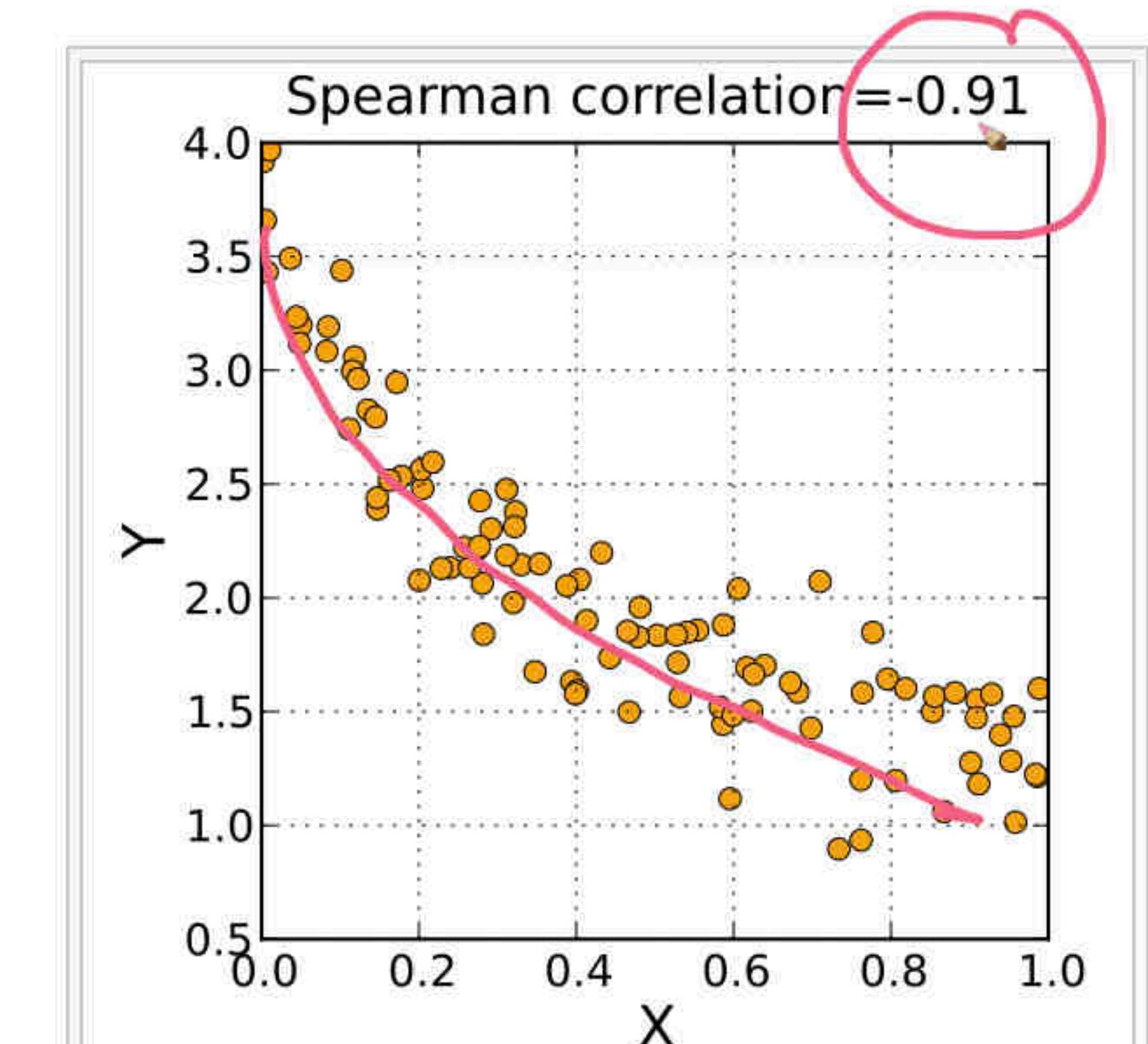
Interpretation [edit]

The sign of the Spearman correlation indicates the direction of association between X (the independent variable) and Y (the dependent variable). If Y tends to increase when X increases, the Spearman correlation coefficient is positive. If Y tends to decrease when X increases, the Spearman correlation coefficient is negative.

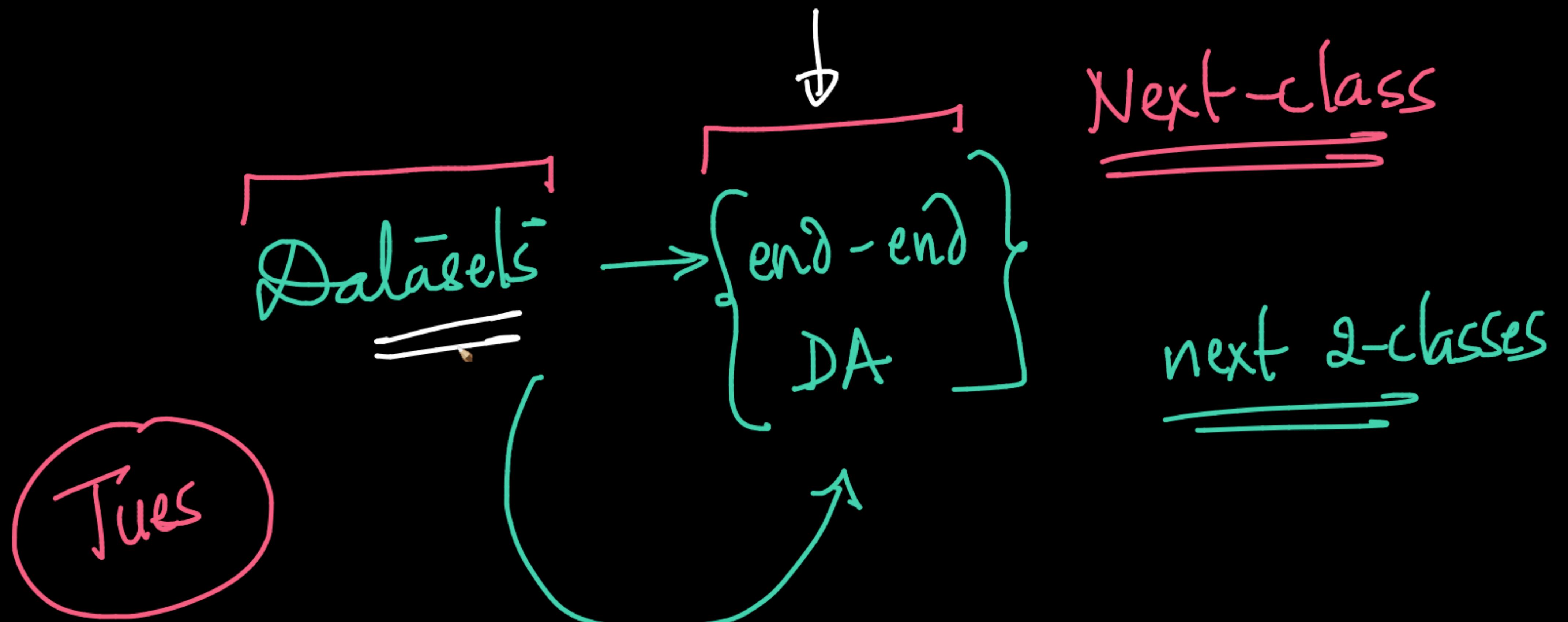
Positive and negative Spearman rank correlations



A positive Spearman correlation coefficient corresponds to an increasing monotonic trend between X and Y .



A negative Spearman correlation coefficient corresponds to a decreasing monotonic trend between X and Y .



~~repeating elements:~~

$x_i: \overbrace{2, 3, 3, 6, 1}^{\sim}$

↓ ↓ ↓ ↓ ↓

2 3 4 5 1

===== ✓

[Print/export](#)[Download as PDF](#)[Printable version](#)[In other projects](#)[Wikimedia Commons](#)[Languages](#)[Deutsch](#)[Español](#)[Français](#)[한국어](#)[日本語](#)[Português](#)[Русский](#)[Tiếng Việt](#)[中文](#)[文 11 more](#)[Edit links](#)

Data [\[edit \]](#)

For all four datasets:

Property	Value	Accuracy
Mean of x	9	exact
Sample variance of $x : s_x^2$	11	exact
Mean of y	7.50	to 2 decimal places
Sample variance of $y : s_y^2$	4.125	±0.003
Correlation between x and y	0.816	to 3 decimal places
Linear regression line	$y = 3.00 + 0.500x$	to 2 and 3 decimal places, respectively
Coefficient of determination of the linear regression : R^2	0.67	to 2 decimal places

- The first [scatter plot](#) (top left) appears to be a simple [linear relationship](#), corresponding to two [variables correlated](#) where y could be modelled as [gaussian](#) with mean linearly dependent on x .
- The second graph (top right); while a relationship between the two variables is obvious, it is not linear, and the [Pearson correlation coefficient](#) is not relevant. A more general regression and the corresponding [coefficient of determination](#) would be more appropriate.
- In the third g... could have a different regression line (a [robust](#)

en.wikipedia.org/wiki/Anscombe%27s_quartet#/media/File:Anscombe%27s_quartet_3.svg

Print/export

Download as PDF

Printable version

In other projects

Wikimedia Commons

Languages 

Deutsch

Español

Français

한국어

日本語

Português

Русский

Tiếng Việt

中文

文 11 more

Edit links

Data [edit]

For all four datasets:

Property	Value	Accuracy
Mean of x	9	exact
Sample variance of $x : s_x^2$	11	exact
Mean of y	7.50	to 2 decimal places
Sample variance of $y : s_y^2$	4.125	±0.003
Correlation between x and y	0.816	to 3 decimal places
Linear regression line	$y = 3.00 + 0.500x$	to 2 and 3 decimal places, respectively
Coefficient of determination of the linear regression : R^2	0.67	to 2 decimal places

• The first scatter plot (top left) appears to be a simple linear relationship, corresponding to two variables correlated where y could be modelled as gaussian with mean linearly dependent on x .

• The second graph (top right); while a relationship between the two variables is obvious, it is not linear, and the Pearson correlation coefficient is not relevant. A more general regression and the corresponding coefficient of determination would be more appropriate.

• In the third g... could have a different regression line (a robust

WIKIPEDIA

The Free Encyclopedia

[Main page](#)
[Contents](#)
[Current events](#)
[Random article](#)
[About Wikipedia](#)
[Contact us](#)
[Donate](#)

[Contribute](#)

[Help](#)
[Learn to edit](#)
[Community portal](#)
[Recent changes](#)
[Upload file](#)

[Tools](#)

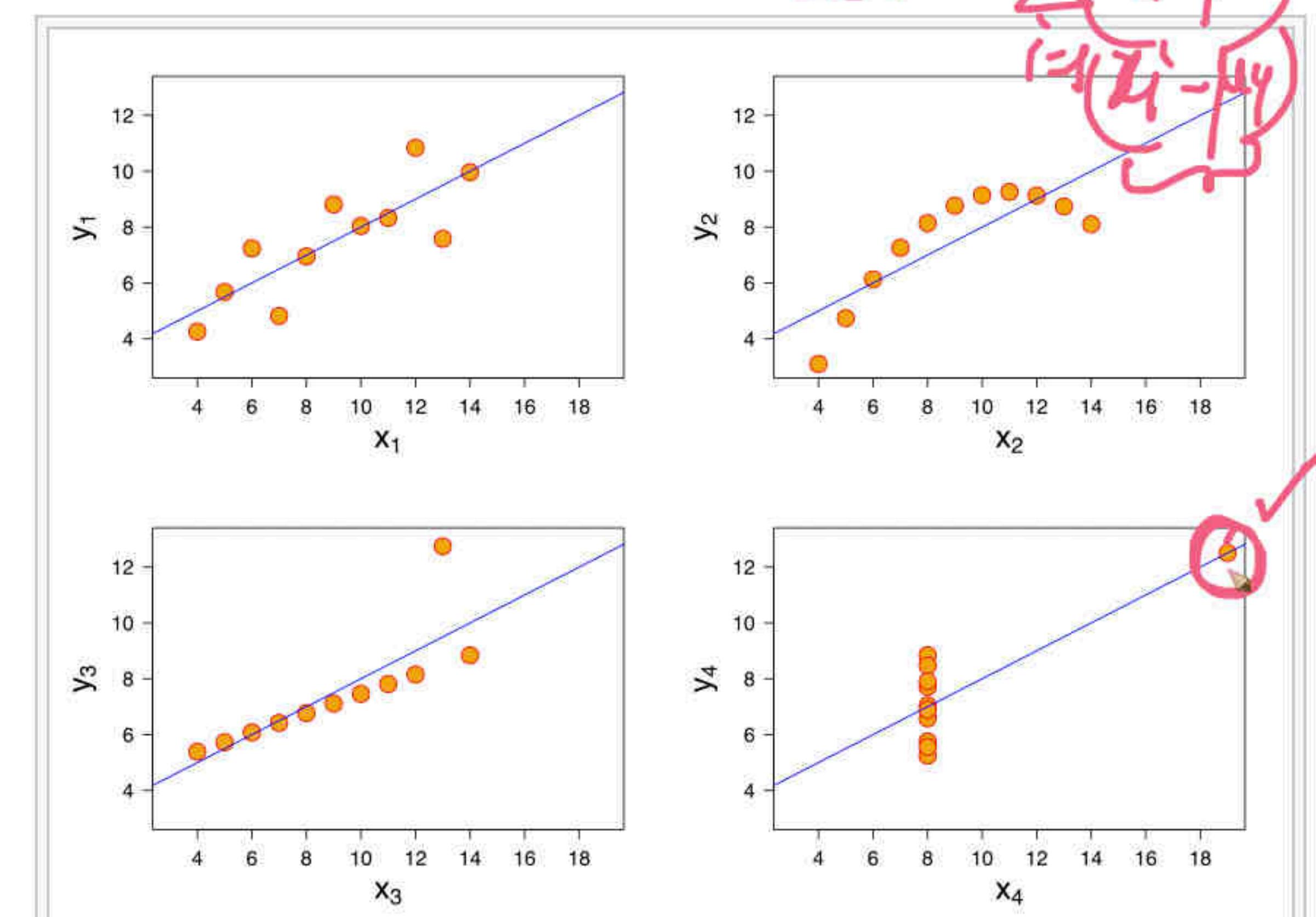
[What links here](#)
[Related changes](#)
[Special pages](#)
[Permanent link](#)
[Page information](#)
[Cite this page](#)

From Wikipedia, the free encyclopedia

Anscombe's quartet comprises four [data sets](#) that have nearly identical simple [descriptive statistics](#), yet have very different [distributions](#) and appear very different when [graphed](#). Each dataset consists of eleven (x,y) points. They were constructed in 1973 by the [statistician Francis Anscombe](#) to demonstrate both the importance of graphing data when analyzing it, and the effect of [outliers](#) and other [influential observations](#) on statistical properties. He described the article as being intended to counter the impression among statisticians that "numerical calculations are exact, but graphs are rough."^[1]

Contents [hide]

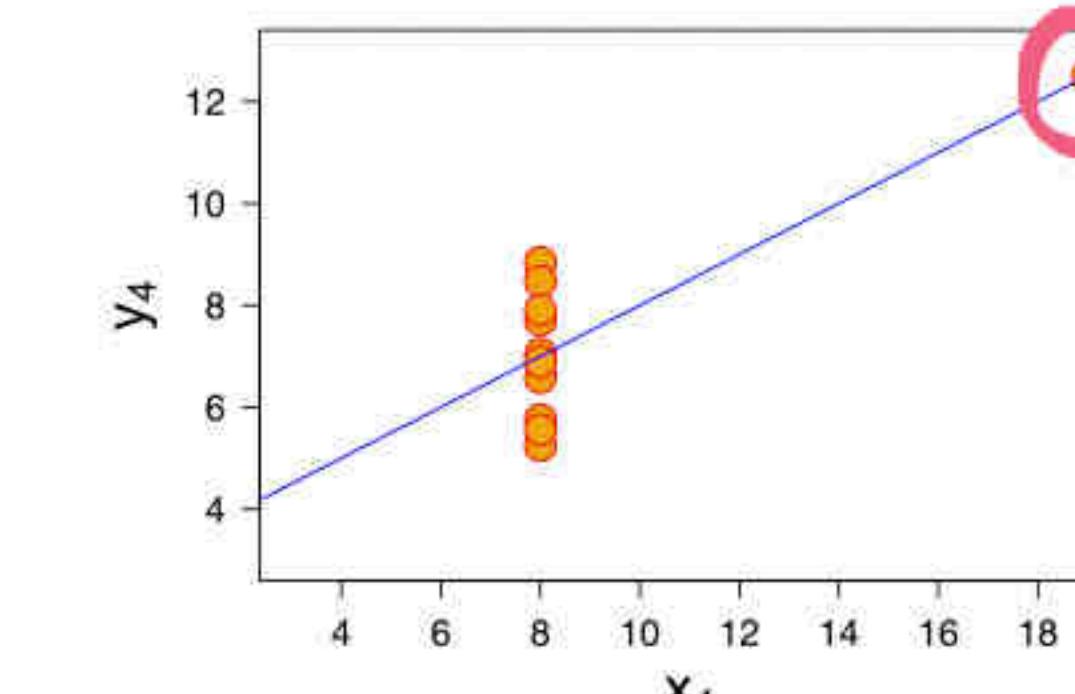
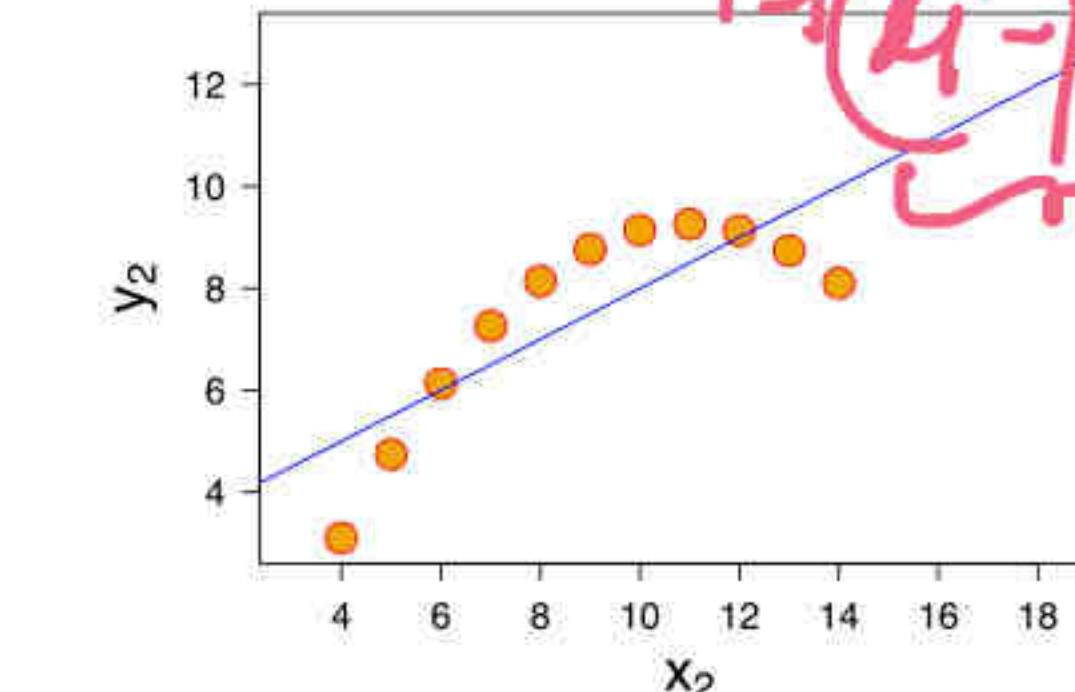
- 1 Data
- 2 See also
- 3 References
- 4 External links



All four sets are identical when examined using simple summary statistics, but vary considerably when graphed

$$\text{COV} = \frac{n}{\sum_{i=1}^n (z_i - \bar{x})(z_i - \bar{y})}$$

$$\sum_{i=1}^n (z_i - \bar{y})$$



WIKIPEDIA

The Free Encyclopedia

Main page

Contents

Current events

Random article

About Wikipedia

Contact us

Donate

Contribute

Help

Learn to edit

Community portal

Recent changes

Upload file

Tools

What links here

Related changes

Special pages

Permanent link

Page information

Cite this page

From Wikipedia, the free encyclopedia

Anscombe's quartet comprises four [data sets](#) that have nearly identical simple [descriptive statistics](#), yet have very different [distributions](#) and appear very different when [graphed](#). Each dataset consists of eleven (x,y) points. They were constructed in 1973 by the [statistician Francis Anscombe](#) to demonstrate both the importance of graphing data when analyzing it, and the effect of [outliers](#) and other [influential observations](#) on statistical properties. He described the article as being intended to counter the impression among statisticians that "numerical calculations are exact, but graphs are rough."^[1]

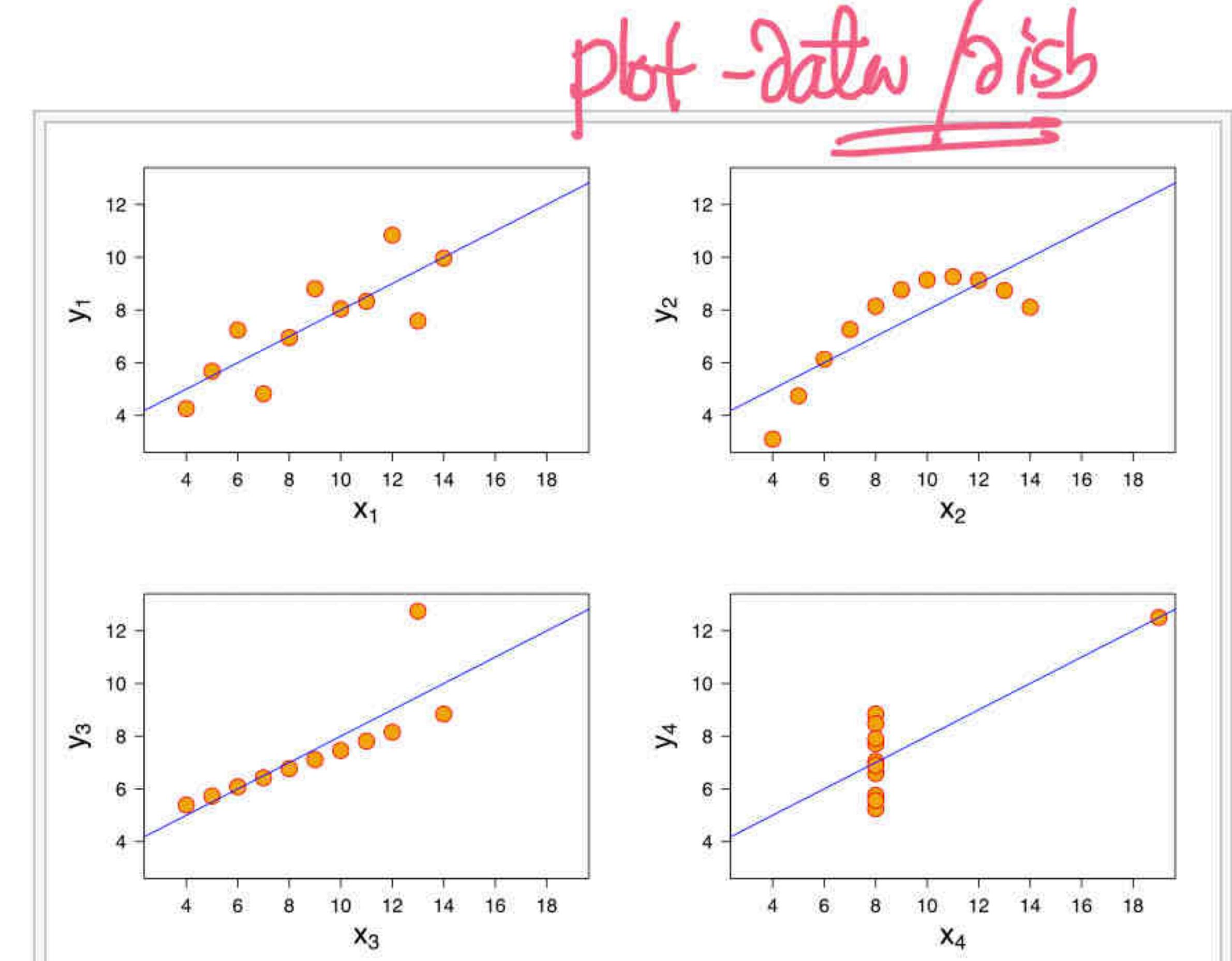
Contents [hide]

1 Data

2 See also

3 References

4 External links



All four sets are identical when examined using simple summary statistics, but vary considerably when graphed

WIKIPEDIA

The Free Encyclopedia

Main page

Contents

Current events

Random article

About Wikipedia

Contact us

Donate

Contribute

Help

Learn to edit

Community portal

Recent changes

Upload file

Tools

What links here

Related changes

Special pages

Permanent link

Page information

Cite this page

From Wikipedia, the free encyclopedia

Anscombe's quartet comprises four data sets that have nearly identical simple descriptive statistics, yet have very different distributions and appear very different when graphed. Each dataset consists of eleven (x,y) points. They were constructed in 1973 by the statistician Francis Anscombe to demonstrate both the importance of graphing data when analyzing it, and the effect of outliers and other influential observations on statistical properties. He described the article as being intended to counter the impression among statisticians that "numerical calculations are exact, but graphs are rough."^[1]

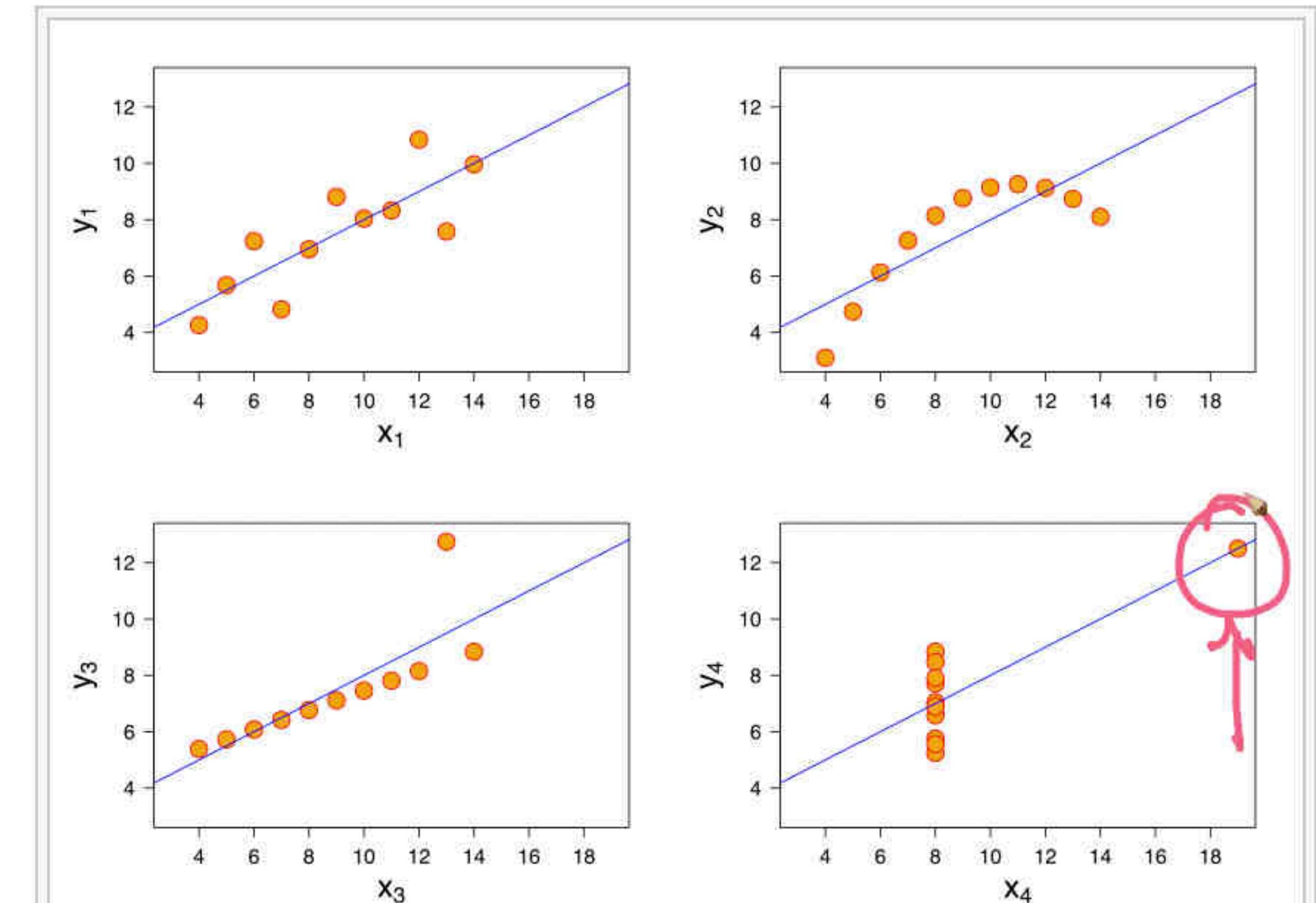
Contents [hide]

1 Data

2 See also

3 References

4 External links



All four sets are identical when examined using simple summary statistics, but vary considerably when graphed

SRCC = $\gamma_s \rightarrow$ test / check for Monotonicity
 $f \rightarrow$ " " " linearity -

multi-collinearity \rightarrow linear regression
(ML)

Chrome File Edit View History Bookmarks Profiles Tab Window Help 11:16

Standard symmetric Covariance trends Pearson correlation Correlation_coefficient File:Correlation_exam Spearman's rank Anscombe's quartet Student t pdf - Study Simpson's paradox + en.wikipedia.org/wiki/Simpson%27s_paradox

Article Talk Read Edit View history Search Wikipedia

SIMPSON'S PARADOX

From Wikipedia, the free encyclopedia

See also: Misuse of statistics

Simpson's paradox, which also goes by several other names, is a phenomenon in probability and statistics in which a trend appears in several groups of data but disappears or reverses when the groups are combined. This result is often encountered in social-science and medical-science statistics,^{[1][2][3]} and is particularly problematic when frequency data is unduly given causal interpretations.^[4] The paradox can be resolved when confounding variables and causal relations are appropriately addressed in the statistical modeling.^{[4][5]} Simpson's paradox has been used to illustrate the kind of misleading results that the misuse of statistics can generate.^{[6][7]}

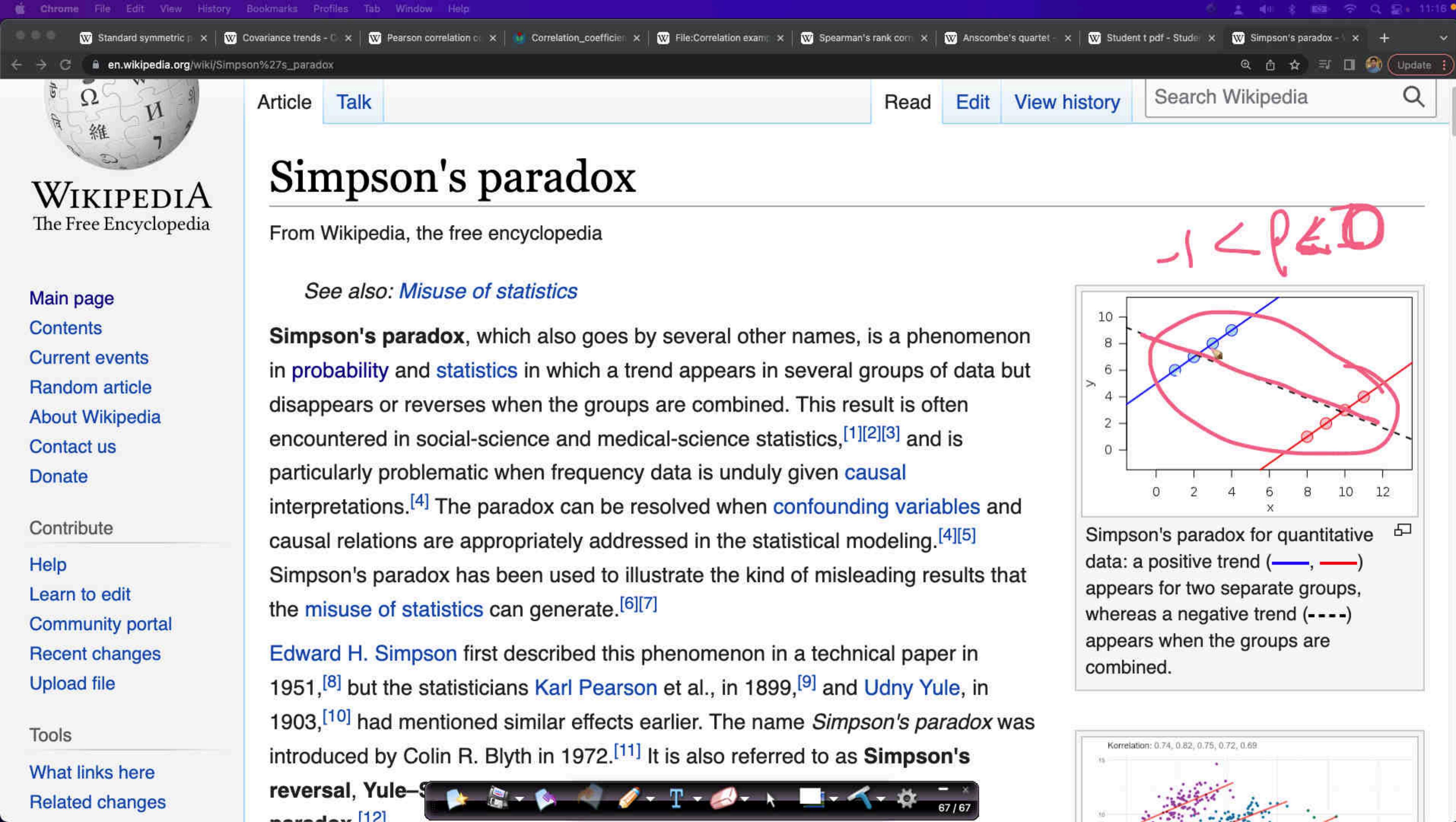
Edward H. Simpson first described this phenomenon in a technical paper in 1951,^[8] but the statisticians Karl Pearson et al., in 1899,^[9] and Udny Yule, in 1903,^[10] had mentioned similar effects earlier. The name *Simpson's paradox* was introduced by Colin R. Blyth in 1972.^[11] It is also referred to as **Simpson's reversal, Yule-Simpson effect, or Simpson-Yule effect**.^[12]

W **M** **W** **OK**

Simpson's paradox for quantitative data: a positive trend (—, —) appears for two separate groups, whereas a negative trend (----) appears when the groups are combined.

Korrelation: -0.74

66 / 66





WIKIPEDIA
The Free Encyclopedia

Main page

Contents

Current events

Random article

About Wikipedia

Contact us

Donate

Contribute

Help

Learn to edit

Community portal

Recent changes

Upload file

Tools

What links here

Related changes

Article Talk

Read

Edit

View history

Search Wikipedia



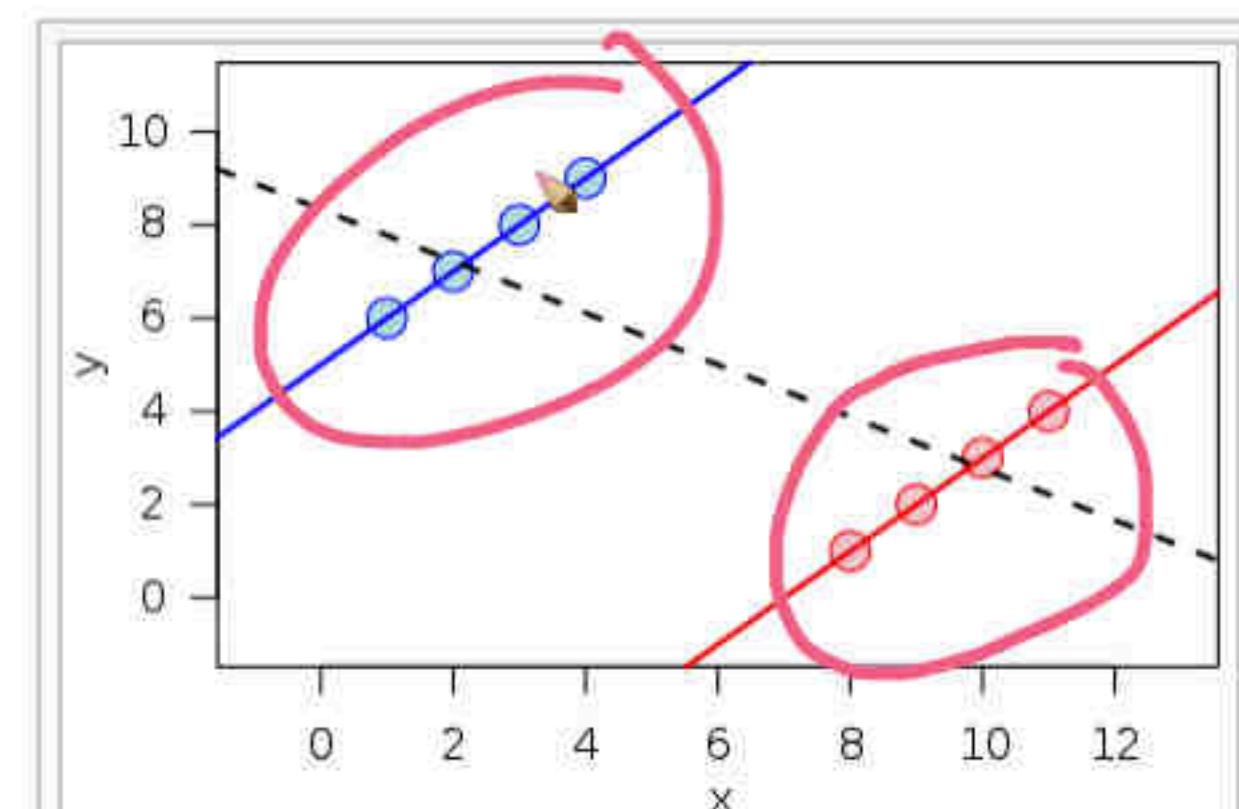
Simpson's paradox

From Wikipedia, the free encyclopedia

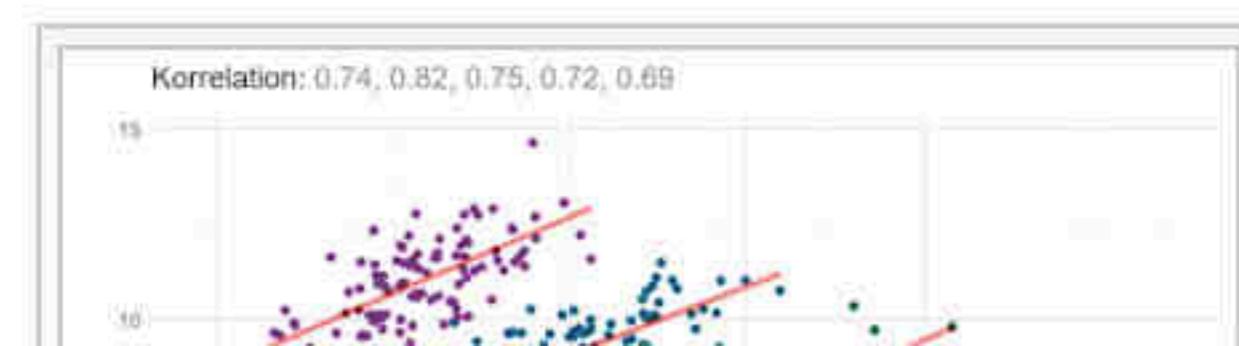
See also: [Misuse of statistics](#)

Simpson's paradox, which also goes by several other names, is a phenomenon in probability and statistics in which a trend appears in several groups of data but disappears or reverses when the groups are combined. This result is often encountered in social-science and medical-science statistics,^{[1][2][3]} and is particularly problematic when frequency data is unduly given causal interpretations.^[4] The paradox can be resolved when confounding variables and causal relations are appropriately addressed in the statistical modeling.^{[4][5]} Simpson's paradox has been used to illustrate the kind of misleading results that the misuse of statistics can generate.^{[6][7]}

Edward H. Simpson first described this phenomenon in a technical paper in 1951,^[8] but the statisticians Karl Pearson et al., in 1899,^[9] and Udny Yule, in 1903,^[10] had mentioned similar effects earlier. The name *Simpson's paradox* was introduced by Colin R. Blyth in 1972.^[11] It is also referred to as **Simpson's reversal**, **Yule-Simpson's paradox**,^[12]



Simpson's paradox for quantitative data: a positive trend (—, —) appears for two separate groups, whereas a negative trend (----) appears when the groups are combined.



Chrome File Edit View History Bookmarks Profiles Tab Window Help

Standard symmetric ... Covariance trends - Pearson correlation co... Correlation_coefficien... File:Correlation exam... Spearman's rank com... Anscombe's quartet - Student t pdf - Stud... Simpson's paradox -

en.wikipedia.org/wiki/Simpson%27s_paradox

Article Talk Read Edit View history Search Wikipedia

WIKIPEDIA The Free Encyclopedia

Main page Contents Current events Random article About Wikipedia Contact us Donate Contribute Help Learn to edit Community portal Recent changes Upload file Tools What links here Related changes

Simpson's paradox

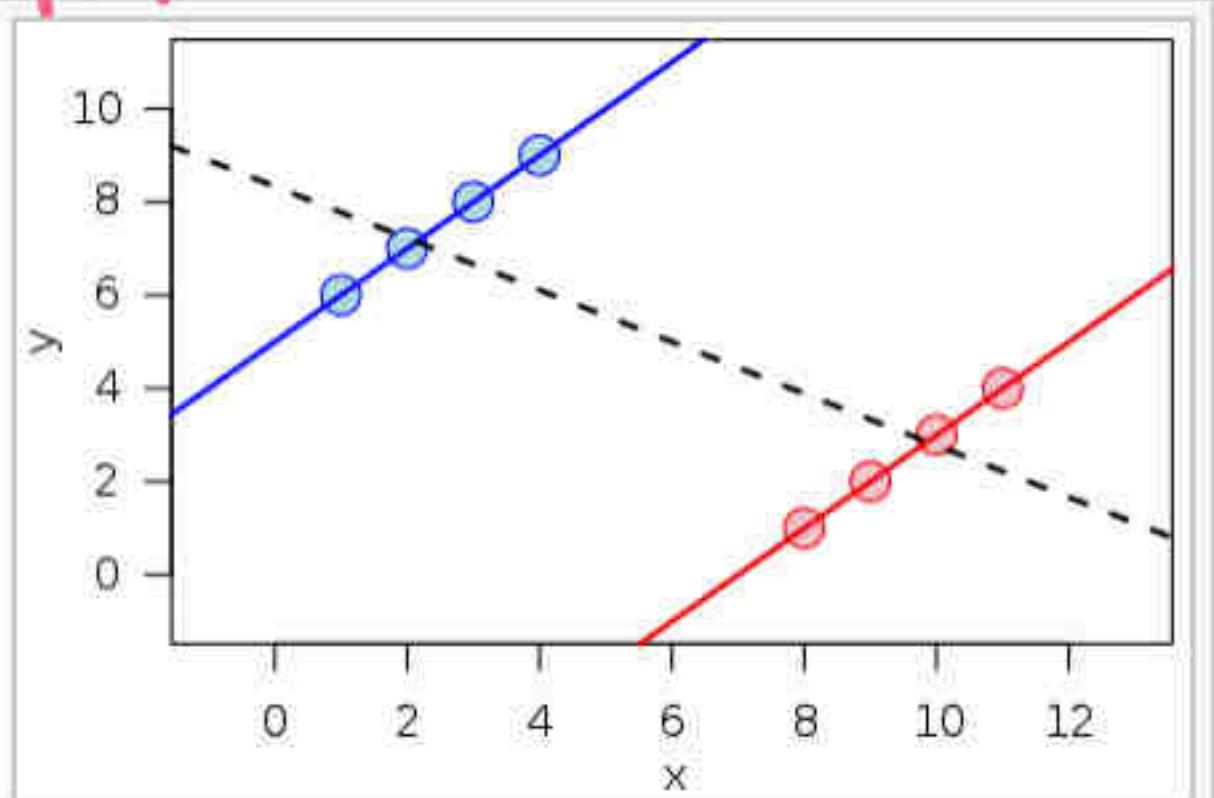
From Wikipedia, the free encyclopedia

See also: *Misuse of statistics*

Simpson's paradox, which also goes by several other names, is a phenomenon in probability and **statistics** in which a trend appears in several groups of data but disappears or reverses when the groups are combined. This result is often encountered in social-science and medical-science statistics,^{[1][2][3]} and is particularly problematic when frequency data is unduly given **causal** interpretations.^[4] The paradox can be resolved when **confounding variables** and causal relations are appropriately addressed in the statistical modeling.^{[4][5]} Simpson's paradox has been used to illustrate the kind of misleading results that the **misuse of statistics** can generate.^{[6][7]}

Edward H. Simpson first described this phenomenon in a technical paper in 1951,^[8] but the statisticians Karl Pearson et al., in 1899,^[9] and Udny Yule, in 1903,^[10] had mentioned similar effects earlier. The name *Simpson's paradox* was introduced by Colin R. Blyth in 1972.^[11] It is also referred to as **Simpson's reversal, Yule-Simpson effect**.^[12]

each gp may show diff beh
as compared to total data



Simpson's paradox for quantitative data: a positive trend (—, —) appears for two separate groups, whereas a negative trend (----) appears when the groups are combined.

Korrelation: 0.74, 0.82, 0.75, 0.72, 0.69



WIKIPEDIA
The Free Encyclopedia

Main page

Contents

Current events

Random article

About Wikipedia

Contact us

Donate

Contribute

Help

Learn to edit

Community portal

Recent changes

Upload file

Tools

What links here

Related changes

Article Talk

Read

Edit

View history

Search Wikipedia



Simpson's paradox

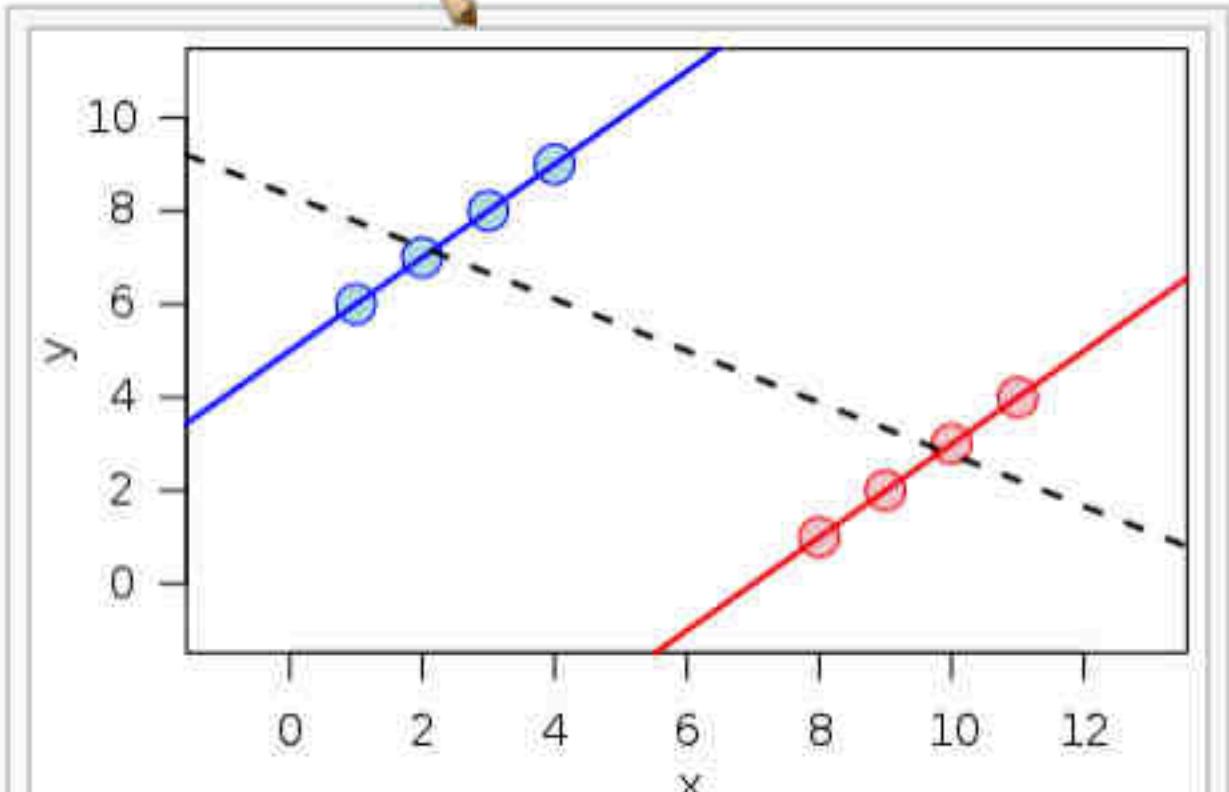
From Wikipedia, the free encyclopedia

#bt/the data

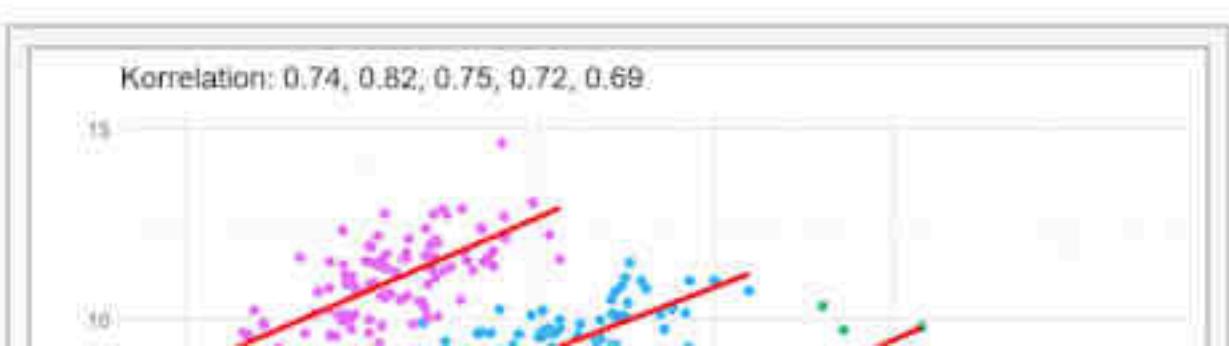
See also: [Misuse of statistics](#)

Simpson's paradox, which also goes by several other names, is a phenomenon in probability and [statistics](#) in which a trend appears in several groups of data but disappears or reverses when the groups are combined. This result is often encountered in social-science and medical-science statistics, [\[1\]](#)[\[2\]](#)[\[3\]](#) and is particularly problematic when frequency data is unduly given [causal](#) interpretations.[\[4\]](#) The paradox can be resolved when [confounding variables](#) and causal relations are appropriately addressed in the statistical modeling.[\[4\]](#)[\[5\]](#) Simpson's paradox has been used to illustrate the kind of misleading results that the [misuse of statistics](#) can generate.[\[6\]](#)[\[7\]](#)

Edward H. Simpson first described this phenomenon in a technical paper in 1951,[\[8\]](#) but the statisticians Karl Pearson et al., in 1899,[\[9\]](#) and Udny Yule, in 1903,[\[10\]](#) had mentioned similar effects earlier. The name *Simpson's paradox* was introduced by Colin R. Blyth in 1972.[\[11\]](#) It is also referred to as **Simpson's reversal**, **Yule–Simpson effect**, **Yule–Simpson paradox**,[\[12\]](#)



Simpson's paradox for quantitative data: a positive trend (—, —) appears for two separate groups, whereas a negative trend (----) appears when the groups are combined.



[Recent changes](#)[Upload file](#)[Tools](#)[What links here](#)[Related changes](#)[Special pages](#)[Permanent link](#)[Page information](#)[Cite this page](#)[Wikidata item](#)[Print/export](#)[Download as PDF](#)[Printable version](#)[In other projects](#)[Wikimedia Commons](#)[Languages](#)[العربية](#)[Deutsch](#)[Español](#)[Français](#)[日本語](#)[Português](#)

Edward H. Simpson first described this phenomenon in a technical paper in 1951,^[8] but the statisticians Karl Pearson et al., in 1899,^[9] and Udny Yule, in 1903,^[10] had mentioned similar effects earlier. The name *Simpson's paradox* was introduced by Colin R. Blyth in 1972.^[11] It is also referred to as **Simpson's reversal, Yule–Simpson effect, amalgamation paradox, or reversal paradox.**^[12]

combined.



Visualization of Simpson's paradox on data resembling real-world variability indicates that risk of misjudgment of true causal relationship can be hard to spot

Contents [hide]

1 Examples

[1.1 UC Berkeley gender bias](#)[1.2 Kidney stone treatment](#)[1.3 Batting averages](#)

2 Vector interpretation

3 Correlation between variables

4 Psychology

5 Probability

6 Simpson's second paradox

7 See also

8 References

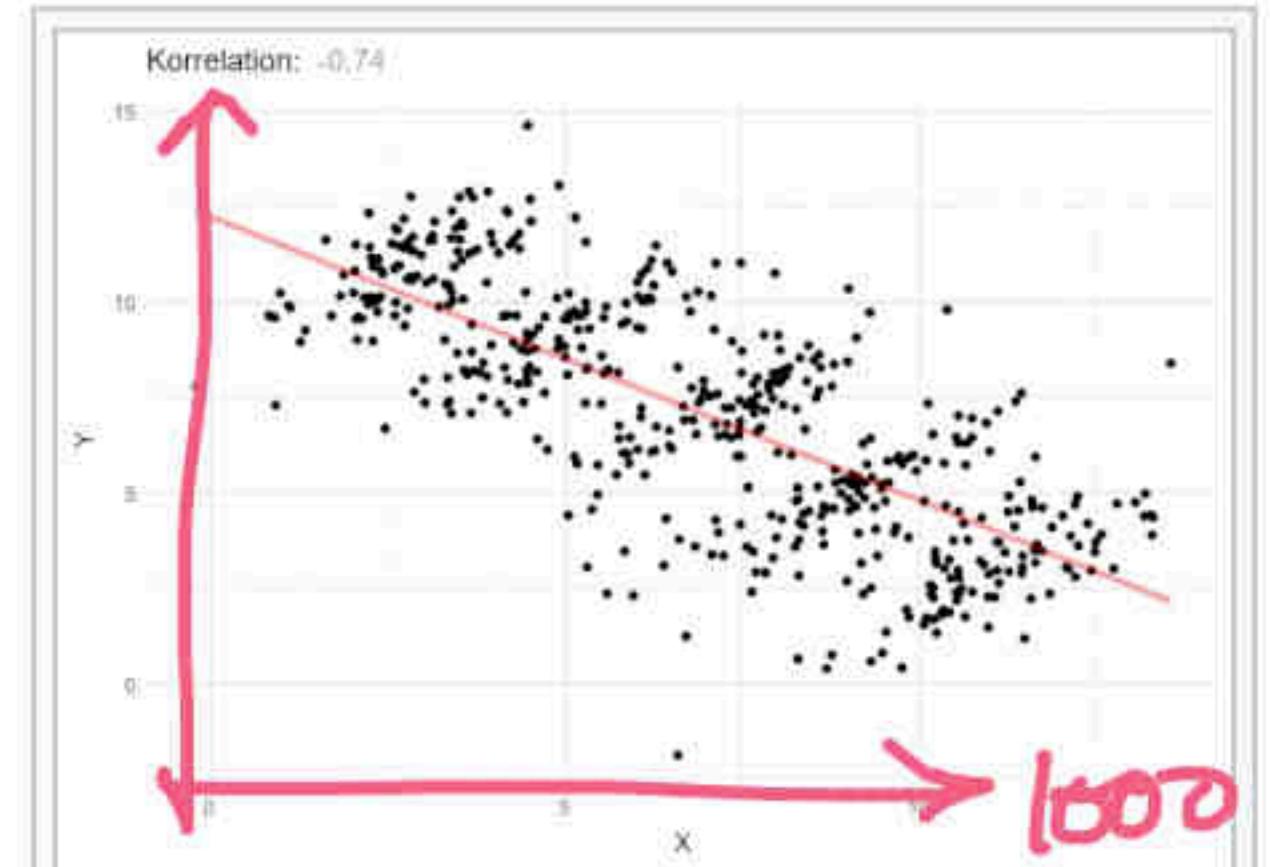
9 Bibliography

10 External links

Edward H. Simpson first described this phenomenon in a technical paper in 1951,^[8] but the statisticians Karl Pearson et al., in 1899,^[9] and Udny Yule, in 1903,^[10] had mentioned similar effects earlier. The name *Simpson's paradox* was introduced by Colin R. Blyth in 1972.^[11] It is also referred to as **Simpson's reversal, Yule–Simpson effect, amalgamation paradox, or reversal paradox.**^[12]

combined.

$$\mu_1 = \mu_2 = \dots = \mu_k$$



Visualization of Simpson's paradox on data resembling real-world variability indicates that risk of misjudgment of true causal relationship can be hard to spot

Contents [hide]

- 1 Examples
 - 1.1 UC Berkeley gender bias
 - 1.2 Kidney stone treatment
 - 1.3 Batting averages
- 2 Vector interpretation
- 3 Correlation between variables
- 4 Psychology
- 5 Probability
- 6 Simpson's second paradox
- 7 See also
- 8 References
- 9 Bibliography
- 10 External links

العربية

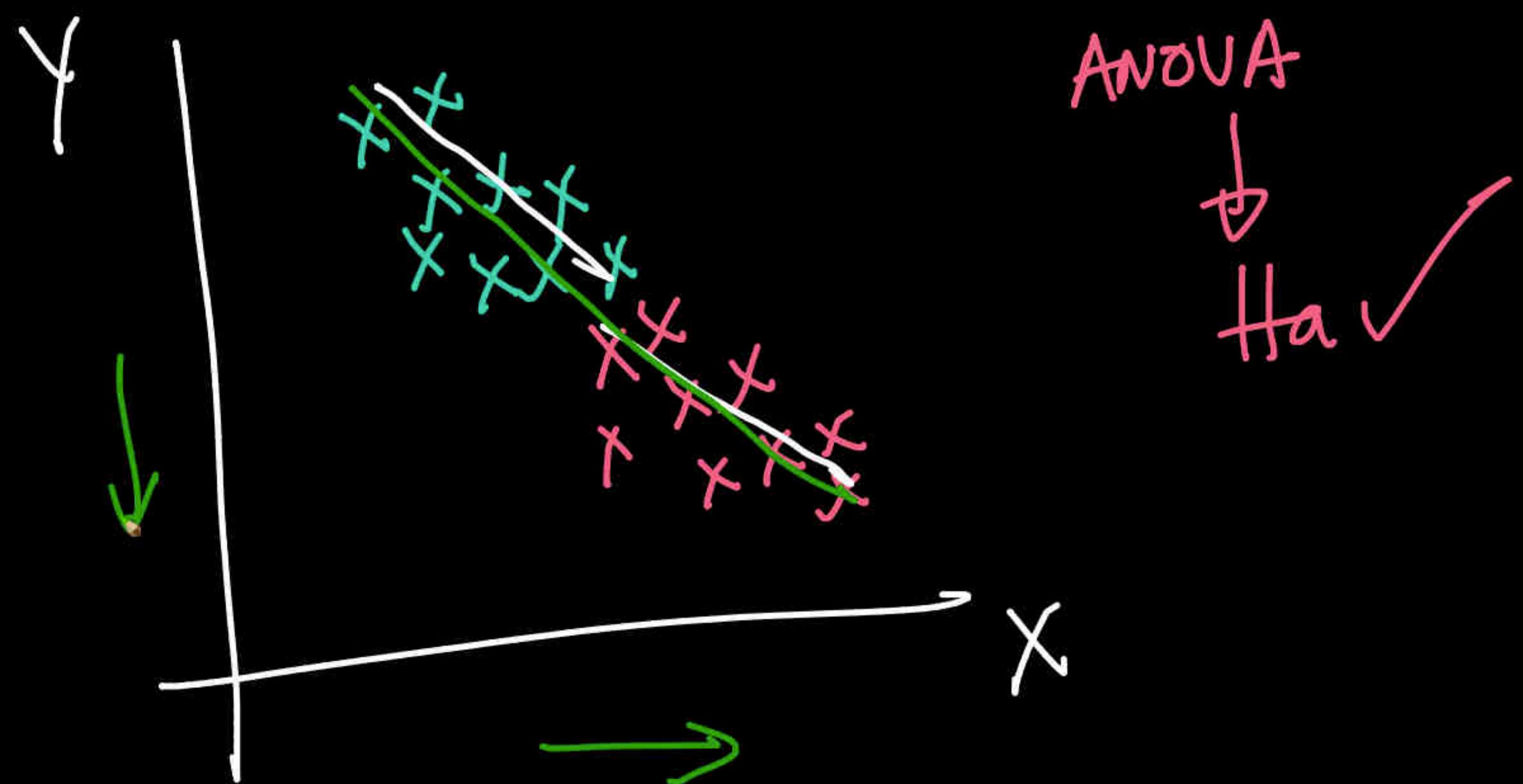
Deutsch

Español

Français

日本語

Português



ANOVA

H_0

H_a

✓

(plot the data)

Sizes: M, S, L, XL (gp)

Task: ...

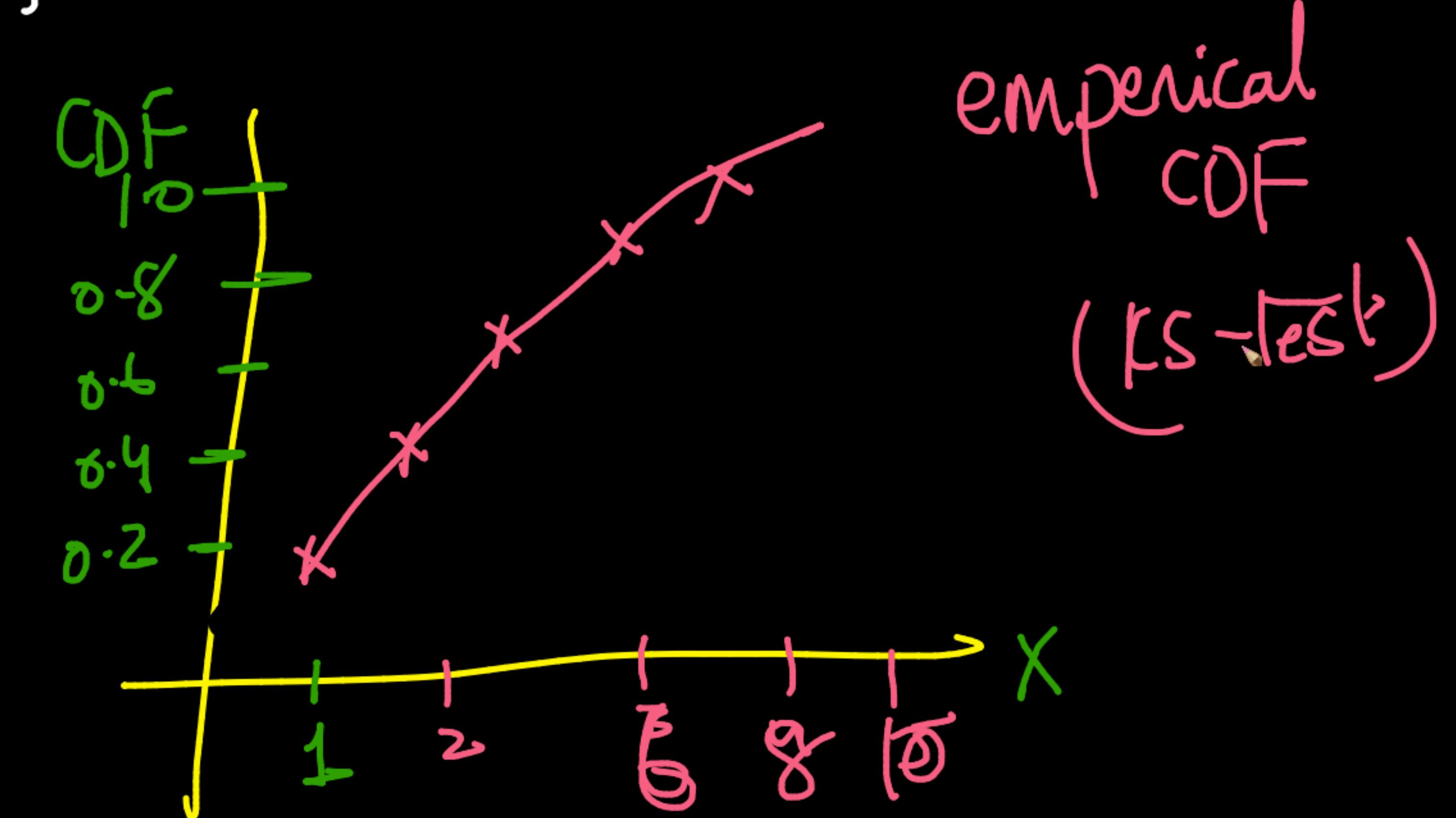
height

Q&A

$10, 8, 6, 2, 1$
 $x_1, x_2, \dots, x_n \rightarrow \text{Sort} \rightarrow$

$n=5$

(1) $2, 6, 8, 10$



Artifacts

$P(H) = 0.5$

as $n \rightarrow \infty$

so $\frac{1}{n} \rightarrow 0$

so $\frac{1}{n}$ heads

$X \sim \text{Bin}(n=100; p=0.9)$

$P(X=10) = \dots$

\approx - Many vs

$X = 10$ (y-a few times)

1000 r.v. of X

≈ 10 values will be 10

105	108
-----	-----

WIKIPEDIA

The Free Encyclopedia

[Main page](#)[Contents](#)[Current events](#)[Random article](#)[About Wikipedia](#)[Contact us](#)[Donate](#)[Contribute](#)[Help](#)[Learn to edit](#)[Community portal](#)[Recent changes](#)[Upload file](#)[Tools](#)[What links here](#)[Related changes](#)[Special pages](#)[Permanent link](#)

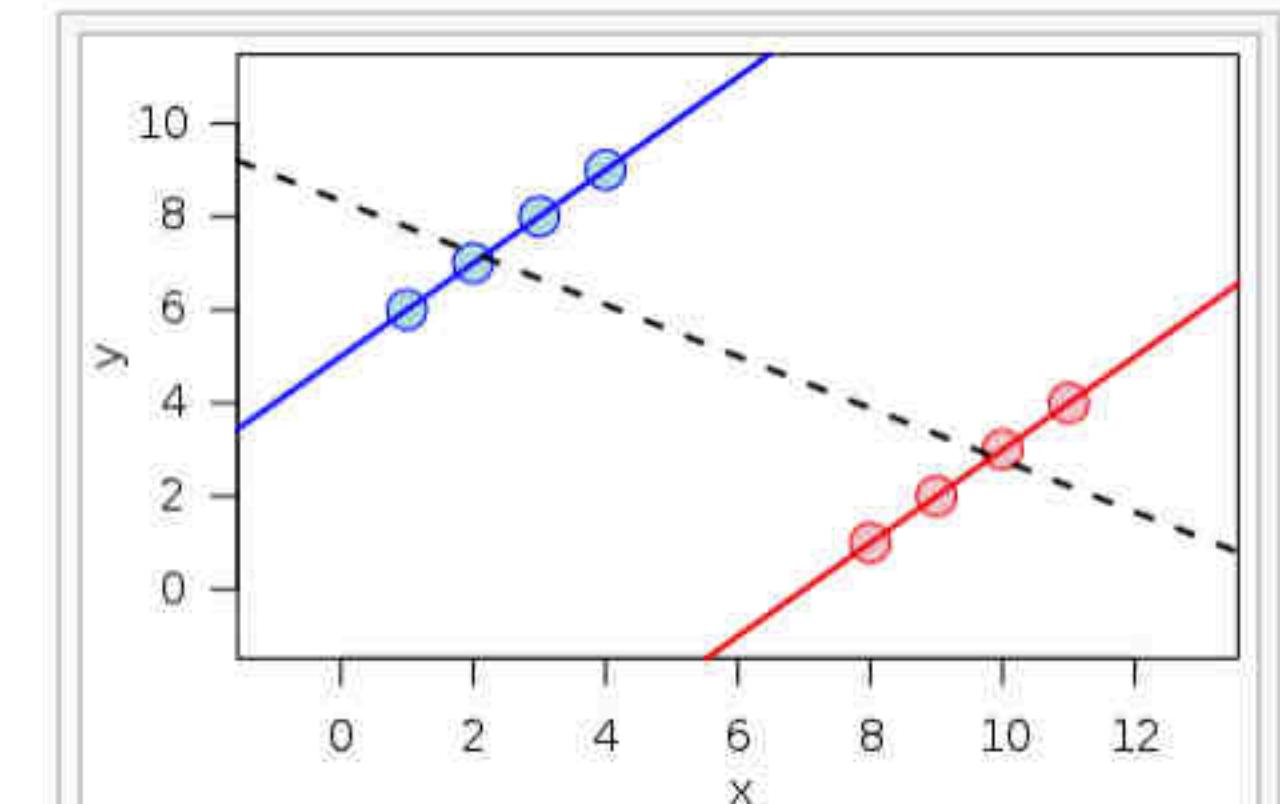
Simpson's paradox

From Wikipedia, the free encyclopedia

See also: Misuse of statistics

Simpson's paradox, which also goes by several other names, is a phenomenon in probability and statistics in which a trend appears in several groups of data but disappears or reverses when the groups are combined. This result is often encountered in social-science and medical-science statistics,^{[1][2][3]} and is particularly problematic when frequency data is unduly given causal interpretations.^[4] The paradox can be resolved when confounding variables and causal relations are appropriately addressed in the statistical modeling.^{[4][5]} Simpson's paradox has been used to illustrate the kind of misleading results that the misuse of statistics can generate.^{[6][7]}

Edward H. Simpson first described this phenomenon in a technical paper in 1951,^[8] but the statisticians Karl Pearson et al., in 1899,^[9] and Udny Yule, in 1903,^[10] had mentioned similar effects earlier. The name *Simpson's paradox* was introduced by Colin R. Blyth in 1972.^[11] It is also referred to as **Simpson's reversal**, **Yule–Simpson effect**, **amalgamation paradox**, or **reversal paradox**.^[12]



Simpson's paradox for quantitative data: a positive trend (—, —) appears for two separate groups, whereas a negative trend (----) appears when the groups are combined.

