

Topics:

"Play" } → experiment

① ANOVA - code

② Chi-square Goodness of fit → χ^2 - test of
indep...

③ Correlation & Covariance

④ Pearson Correlation-coeff — lots of diagrams

⑤ Spearman Rank CG — more diagrams

+ problems

OPS



$$\mu = \sum_{i=1}^n p_i x_i.$$

(When such a discrete [weighted variance](#) is specified by weights whose sum is not 1, then one divides by the sum of the weights.)

The variance of a collection of n equally likely values can be written as

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

where μ is the average value. That is,

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i.$$

The variance of a set of n equally likely values can be equivalently expressed, without directly referring to the mean, in terms of squared deviations of all pairwise squared distances of points from each other:^[2]

$$\text{Var}(X) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \frac{1}{2} (x_i - x_j)^2 = \frac{1}{n^2} \sum_i \sum_{j>i} (x_i - x_j)^2.$$

Absolutely continuous random variable [edit]

If the random variable X has a [probability density function](#) $f(x)$, and $F(x)$ is the corresponding [cumulative distribution function](#), then

$$\text{Var}(X) = \sigma^2 = \int_{\mathbb{R}} (x - \mu)^2 f(x) dx$$

$$= \int_{\mathbb{R}} x^2 f(x) dx - 2\mu \int_{\mathbb{R}} x f(x) dx + \mu^2 \int_{\mathbb{R}} f(x) dx$$

Variance - Wikipedia ANOVA & Chi-Square.ipynb - Covariance - Wikipedia Pearson correlation coefficient Spearman's rank correlation New Tab en.wikipedia.org/wiki/Variance

$\mu = \sum_{i=1}^n p_i x_i.$

(When such a discrete weighted variance is specified by weights whose sum is not 1, then one divides by the sum of the weights.)

μ : pop-mean

The variance of a collection of n equally likely values can be written as

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

$x: x_1, x_2, \dots, x_n$

where μ is the average value. That is,

\bar{x} : sample mean

$$\bar{x} = \mu = \frac{1}{n} \sum_{i=1}^n x_i.$$

The variance of a set of n equally likely values can be equivalently expressed, without directly referring to the mean, in terms of squared deviations of all pairwise squared distances of points from each other:^[2]

$$\text{Var}(X) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \frac{1}{2} (x_i - x_j)^2 = \frac{1}{n^2} \sum_i \sum_{j>i} (x_i - x_j)^2.$$

Absolutely continuous random variable [edit]

If the random variable X has a probability density function $f(x)$, and $F(x)$ is the corresponding cumulative distribution function, then

$$\begin{aligned} \text{Var}(X) &= \sigma^2 = \int_{\mathbb{R}} (x - \mu)^2 f(x) dx \\ &= \int_{\mathbb{R}} x^2 f(x) dx - 2\mu \int_{\mathbb{R}} x f(x) dx + \mu^2 \int_{\mathbb{R}} f(x) dx \end{aligned}$$

4 / 4

The variance of a collection of n equally likely values can be written as

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

where μ is the average value. That is,

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i.$$

The variance of a set of n equally likely values can be equivalently expressed, without directly referring to the mean, in terms of squared deviations of all pairwise squared distances of points from each other:^[2]

$$\left[\text{Var}(X) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \frac{1}{2} (x_i - x_j)^2 = \frac{1}{n^2} \sum_i \sum_{j>i} (x_i - x_j)^2. \right]$$

Absolutely continuous random variable [edit]

If the random variable X has a [probability density function](#) $f(x)$, and $F(x)$ is the corresponding [cumulative distribution function](#), then

$$\begin{aligned} \text{Var}(X) &= \sigma^2 = \int_{\mathbb{R}} (x - \mu)^2 f(x) dx \\ &= \int_{\mathbb{R}} x^2 f(x) dx - 2\mu \int_{\mathbb{R}} x f(x) dx + \mu^2 \int_{\mathbb{R}} f(x) dx \\ &= \int_{\mathbb{R}} x^2 dF(x) - 2\mu \int_{\mathbb{R}} x dF(x) + \mu^2 \int_{\mathbb{R}} dF(x) \\ &= \int_{\mathbb{R}} x^2 dF(x) - 2\mu x dF(x) + \mu^2 dF(x) \end{aligned}$$

variance can also be thought of as the covariance of a random variable with itself

$$\text{Var}(X) = \text{Cov}(X, X).$$

The variance is also equivalent to the second **cumulant** of a probability distribution that generates X . The variance is typically designated as $\text{Var}(X)$, or sometimes as $V(X)$ or $\mathbb{V}(X)$, or symbolically as σ_X^2 or simply σ^2 (pronounced "sigma squared"). The expression for the variance can be expanded as follows:

$$\begin{aligned}\text{Var}(X) &= \mathbb{E}[(X - \mathbb{E}[X])^2] \\ &= \mathbb{E}[X^2 - 2X\mathbb{E}[X] + \mathbb{E}[X]^2] \\ &= \mathbb{E}[X^2] - 2\mathbb{E}[X]\mathbb{E}[X] + \mathbb{E}[X]^2 \\ &= \mathbb{E}[X^2] - \mathbb{E}[X]^2\end{aligned}$$

In other words, the variance of X is equal to the mean of the square of X minus the square of the mean of X . This equation should not be used for computations using floating point arithmetic, because it suffers from catastrophic cancellation if the two components of the equation are similar in magnitude. For other numerically stable alternatives, see [Algorithms for calculating variance](#).

Discrete random variable

If the generator of random variable X is discrete with probability mass function $x_1 \mapsto p_1, x_2 \mapsto p_2, \dots, x_n \mapsto p_n$, then

$$\text{Var}(X) = \sum_{i=1}^n p_i \cdot (x_i - \mu)^2$$

where μ is the expected value. That is

$$\mu = \sum_{i=1}^n p_i x_i.$$

Basic properties [edit]

Variance is non-negative because the squares are positive or zero:

$$\text{Var}(X) \geq 0.$$

The variance of a constant is zero.

$$\text{Var}(a) = 0.$$

Conversely, if the variance of a random variable is 0, then it is [almost surely](#) a constant. That is, it always has the same value:

$$\text{Var}(X) = 0 \iff \exists a : P(X = a) = 1.$$

Variance is [invariant](#) with respect to changes in a [location parameter](#). That is, if a constant is added to all values of the variable, the variance is unchanged:

$$\text{Var}(X + a) = \text{Var}(X).$$

If all values are scaled by a constant, the variance is scaled by the square of that constant:

$$\boxed{\text{Var}(aX) = a^2 \text{Var}(X)}.$$

The variance of a sum of two random variables is given by

$$\text{Var}(aX + bY) = a^2 \text{Var}(X) + b^2 \text{Var}(Y) + 2ab \text{Cov}(X, Y),$$

$$\text{Var}(aX - bY) = a^2 \text{Var}(X) + b^2 \text{Var}(Y) - 2ab \text{Cov}(X, Y),$$

where $\text{Cov}(X, Y)$ is the [covariance](#).

Linear combinations [edit]

In general, for the sum of N random variables $\{X_1, \dots, X_N\}$, the variance becomes:

$$x: x_1, x_2, \dots, x_n$$

$$\text{Var}(x)$$

$$y: a\bar{x}_1, a\bar{x}_2, \dots, a\bar{x}_n$$

$$\text{Var}(y) = \text{Var}(ax)$$

Variance - Wikipedia

ANOVA & Chi-Square.ipynb -

Covariance - Wikipedia

Pearson correlation coefficient

Spearman's rank correlation co

New Tab

colab.research.google.com/drive/1fK356lyrxxvZ-NAHnu0konAvzuktvWaN#scrollTo=luJpH2kXGWpU

Reconnect



+ Code

+ Text

Playing with ANOVA

```
{x} [ ] from scipy.stats import f_oneway  
from scipy.stats import norm  
import seaborn as sns  
import matplotlib.pyplot as plt
```

```
[ ] G1 = norm.rvs(loc=10.0, scale=2.0, size=40)  
G2 = norm.rvs(loc=10.3, scale=2.1, size=40)  
G3 = norm.rvs(loc=9.8, scale=2.3, size=50)
```

```
▶ sns.distplot(G1, hist=False)  
sns.distplot(G2, hist=False)  
sns.distplot(G3, hist=False)  
plt.grid()
```

```
↳ /usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func  
warnings.warn(msg, FutureWarning)  
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func  
warnings.warn(msg, FutureWarning)  
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func  
warnings.warn(msg, FutureWarning)
```

+ Code + Text

Reconnect



Playing with ANOVA

one-way ANOVA

{x}
[] from scipy.stats import f_oneway
from scipy.stats import norm
import seaborn as sns
import matplotlib.pyplot as plt[] G1 = norm.rvs(loc=10.0, scale=2.0, size=40)
G2 = norm.rvs(loc=10.3, scale=2.1, size=40)
G3 = norm.rvs(loc=9.8, scale=2.3, size=50)▶
[] sns.distplot(G1, hist=False)
sns.distplot(G2, hist=False)
sns.distplot(G3, hist=False)
plt.grid()<
[] /usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func
warnings.warn(msg, FutureWarning)
<
[] /usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func
warnings.warn(msg, FutureWarning)
<
[] /usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func
warnings.warn(msg, FutureWarning)



+ Code + Text

Playing with ANOVA

```
[ ] from scipy.stats import f_oneway  
from scipy.stats import norm  
import seaborn as sns  
import matplotlib.pyplot as plt
```

$$\mu_1 \quad \sigma_1$$

```
[ ] G1 = norm.rvs(loc=10.0, scale=2.0, size=40)  
G2 = norm.rvs(loc=10.3, scale=2.1, size=40)  
G3 = norm.rvs(loc=9.8, scale=2.3, size=50)
```

W

```
[ ] sns.distplot(G1, hist=False)  
sns.distplot(G2, hist=False)  
sns.distplot(G3, hist=False)  
plt.grid()
```

Simulation

$$G_1, G_2, G_3$$

$$N(\mu_1, \sigma_1) \quad N(\mu_2, \sigma_2)$$

$$N(\mu_3, \sigma_3)$$

<>
=>
=>

```
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func  
warnings.warn(msg, FutureWarning)  
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func  
warnings.warn(msg, FutureWarning)  
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func  
warnings.warn(msg, FutureWarning)
```

+ Code + Text

Reconnect



Playing with ANOVA

{x}

□

```
[ ] from scipy.stats import f_oneway  
from scipy.stats import norm  
import seaborn as sns  
import matplotlib.pyplot as plt
```

$$\mu_1 \quad \sigma_1$$

```
[ ] G1 = norm.rvs(loc=10.0, scale=2.0, size=40)  
G2 = norm.rvs(loc=10.3, scale=2.1, size=40)  
G3 = norm.rvs(loc=9.8, scale=2.3, size=50)
```

$$\mu_2, \sigma_2$$

$$\mu_3, \sigma_3$$

$$n = m \times k$$

```
[ ] sns.distplot(G1, hist=False)  
sns.distplot(G2, hist=False)  
sns.distplot(G3, hist=False)  
plt.grid()
```

```
<> /usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func  
      warnings.warn(msg, FutureWarning)  
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func  
      warnings.warn(msg, FutureWarning)  
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func  
      warnings.warn(msg, FutureWarning)
```

Variance - Wikipedia ANOVA & Chi-Square.ipynb - Covariance - Wikipedia Pearson correlation coefficient Spearman's rank correlation New Tab colab.research.google.com/drive/1fK356lyrxxVZ-NAHnu0konAvzuktvWaN#scrollTo=luJpH2kXGWpU

+ Code + Text Reconnect

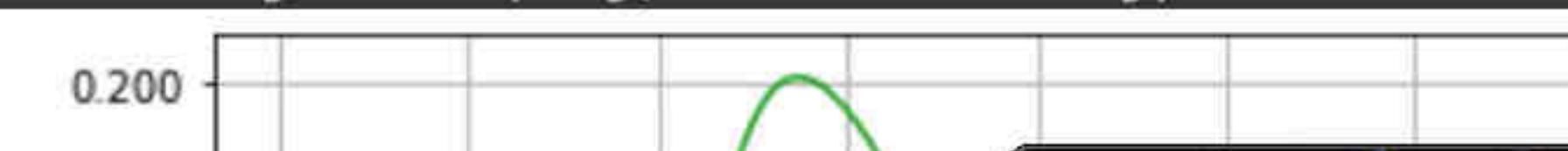
Playing WITH ANOVA

[] `from scipy.stats import f_oneway
from scipy.stats import norm
import seaborn as sns
import matplotlib.pyplot as plt`

[] `G1 = norm.rvs(loc=10.0, scale=2.0, size=40)
G2 = norm.rvs(loc=10.3, scale=2.1, size=40)
G3 = norm.rvs(loc=9.8, scale=2.3, size=50)`

▶ `sns.distplot(G1, hist=False)
sns.distplot(G2, hist=False)
sns.distplot(G3, hist=False)
plt.grid()`

↳ /usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func
 warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func
 warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func
 warnings.warn(msg, FutureWarning)



0.200

12 / 12

+ Code + Text

Reconnect



Playing with ANOVA



```
[ ] from scipy.stats import f_oneway  
from scipy.stats import norm  
import seaborn as sns  
import matplotlib.pyplot as plt
```

```
[ ] G1 = norm.rvs(loc=10.0, scale=2.0, size=40)  
G2 = norm.rvs(loc=10.3, scale=2.1, size=40)  
G3 = norm.rvs(loc=9.8, scale=2.3, size=50)
```

▶

```
sns.distplot(G1, hist=False)  
sns.distplot(G2, hist=False)  
sns.distplot(G3, hist=False)  
plt.grid()
```

```
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func  
warnings.warn(msg, FutureWarning)  
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func  
warnings.warn(msg, FutureWarning)  
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func  
warnings.warn(msg, FutureWarning)
```



Variance - Wikipedia

ANOVA & Chi-Square.ipynb

Covariance - Wikipedia

Pearson correlation coefficient

Spearman's rank correlation coefficient

New Tab



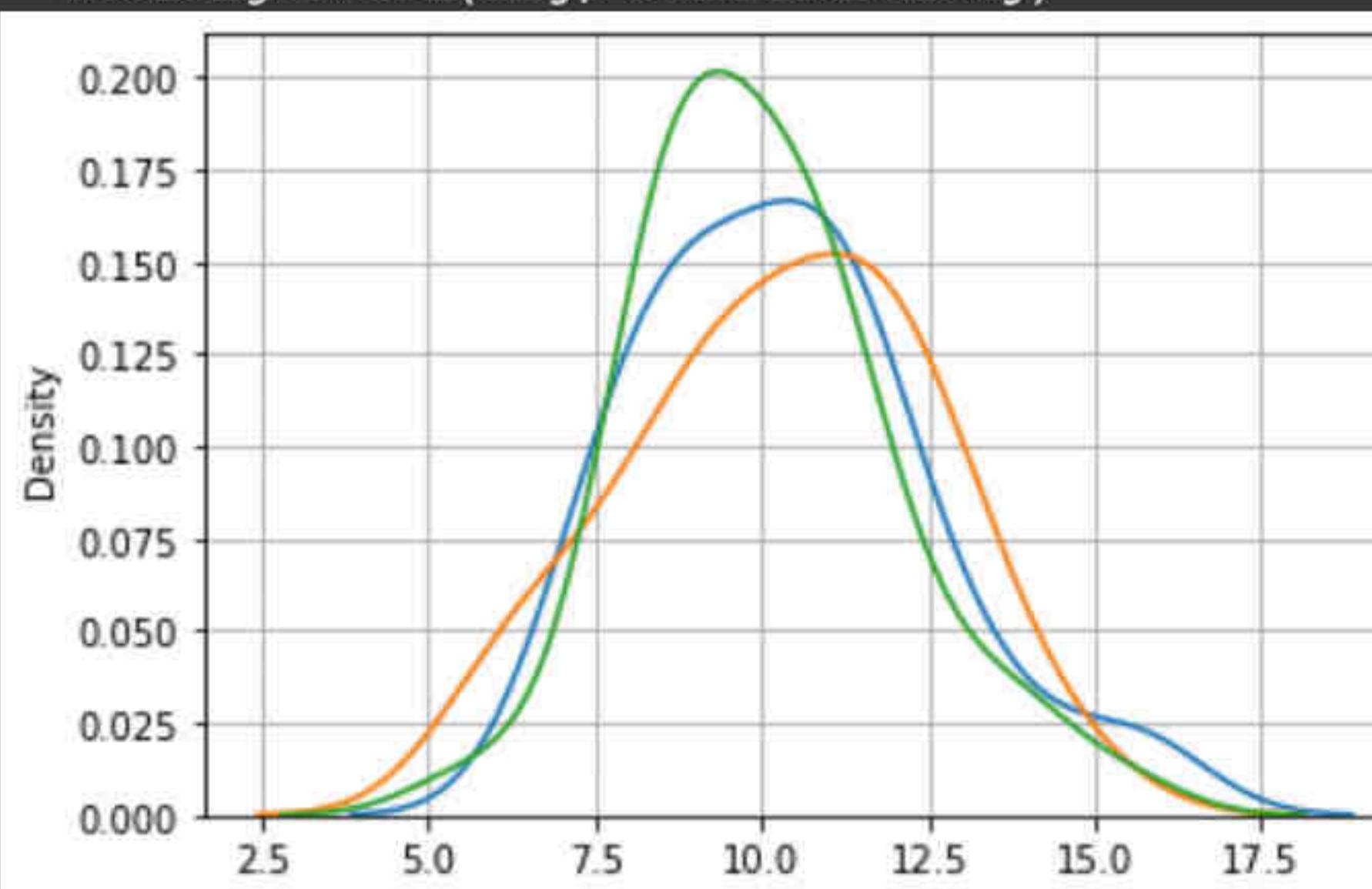
+ Code + Text

```
[ ] sns.distplot(G3, hist=False)  
    plt.grid()
```

Reconnect



```
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func  
  warnings.warn(msg, FutureWarning)  
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func  
  warnings.warn(msg, FutureWarning)  
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func  
  warnings.warn(msg, FutureWarning)
```



(Q)

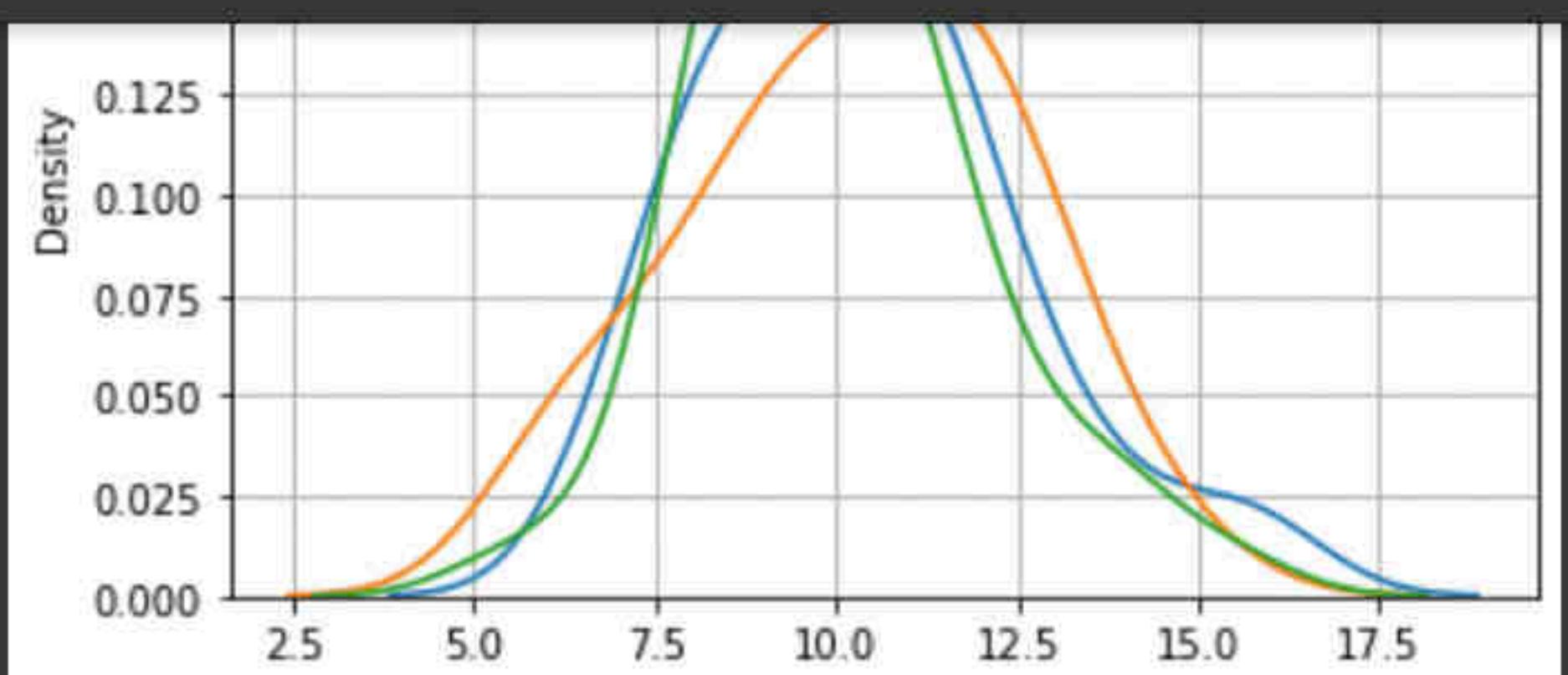
$H_0: \mu_1 = \mu_2 = \mu_3$

$H_a: \text{Otherwise}$



```
[ ] f_oneway(G1, G2, G3)
```





p-val < α

```
[ ] f_oneway(G1, G2, G3)
F_onewayResult(statistic=17.22314794729301, pvalue=2.408534447184806e-07)
```

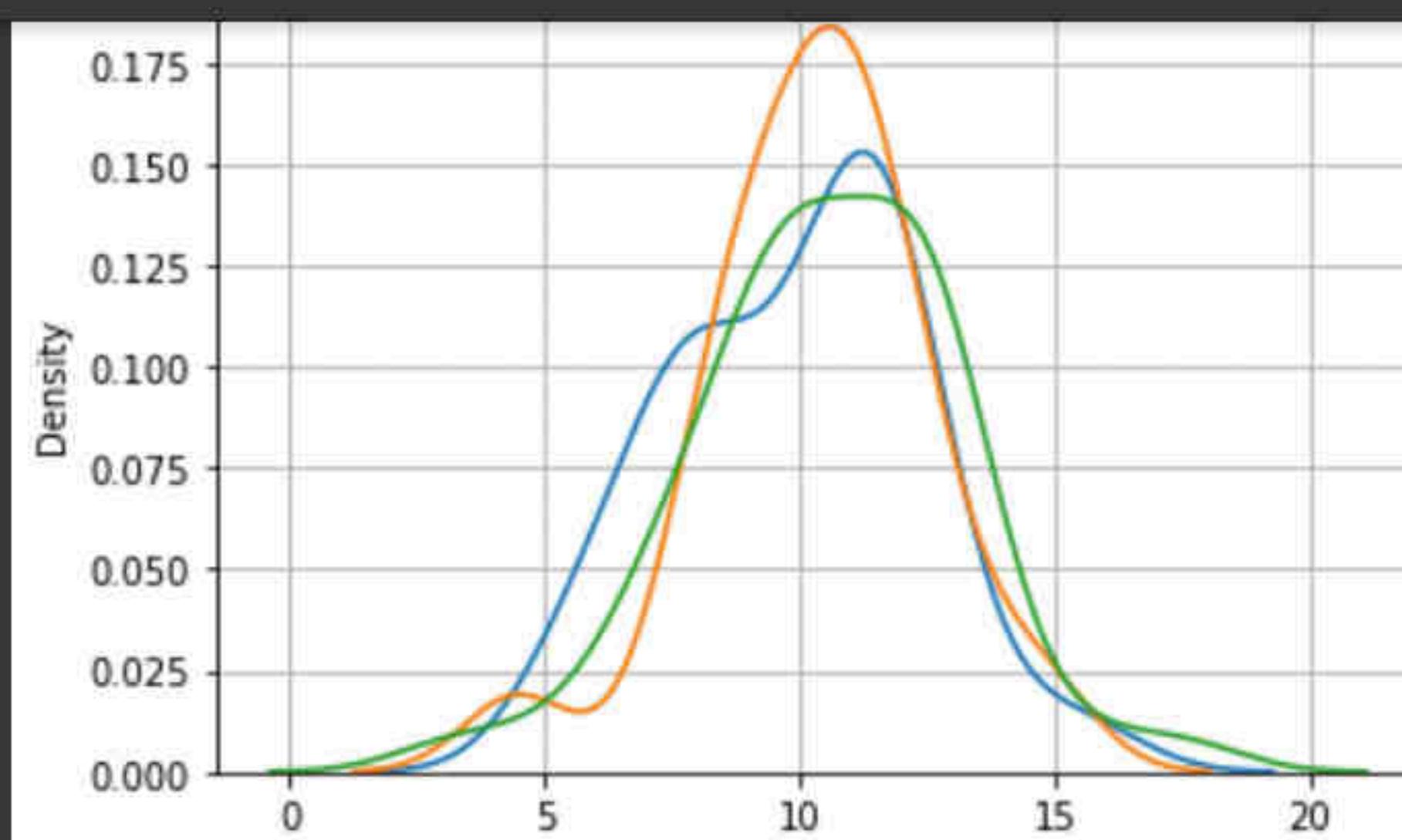
[]

Double-click (or enter) to edit

[]

#

+ Code + Text



49.68 → $\alpha = 5\%$
p-val

[64] f_oneway(G1, G2, G3)

F_onewayResult(statistic=0.7034160770793035, pvalue=0.4968095155468297)

Double-click (or enter) to edit



[]

#

+ Code + Text

✓ RAM Disk



Playing with ANOVA

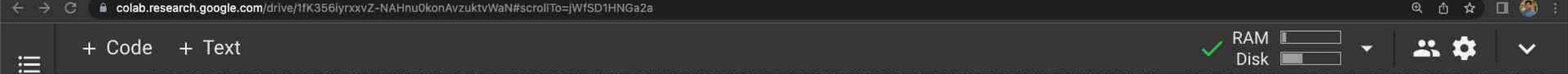
{x}

```
[61] from scipy.stats import f_oneway  
     from scipy.stats import norm  
     import seaborn as sns  
     import matplotlib.pyplot as plt
```

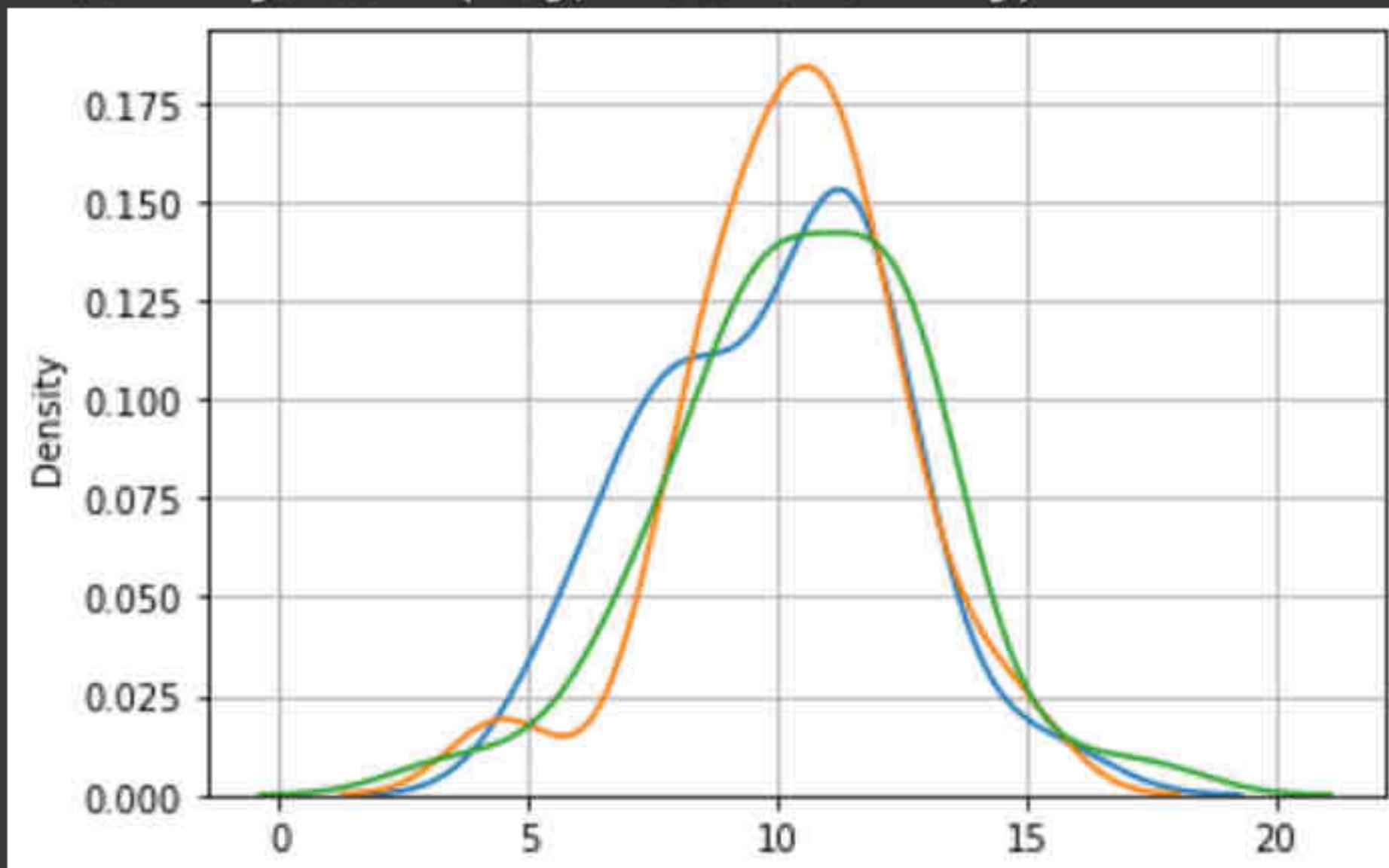
```
[62] G1 = norm.rvs(loc=10.0, scale=2.0, size=40)  
     G2 = norm.rvs(loc=10.3, scale=2.1, size=40)  
     G3 = norm.rvs(loc=9.8, scale=2.3, size=50)
```

```
[63] sns.distplot(G1, hist=False)  
     sns.distplot(G2, hist=False)  
     sns.distplot(G3, hist=False)  
     plt.grid()
```

```
<> /usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func  
      warnings.warn(msg, FutureWarning)  
      /usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func  
      warnings.warn(msg, FutureWarning)  
      /usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func  
      warnings.warn(msg, FutureWarning)
```



```
[63]: warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func
warnings.warn(msg, FutureWarning)
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func
warnings.warn(msg, FutureWarning)
```

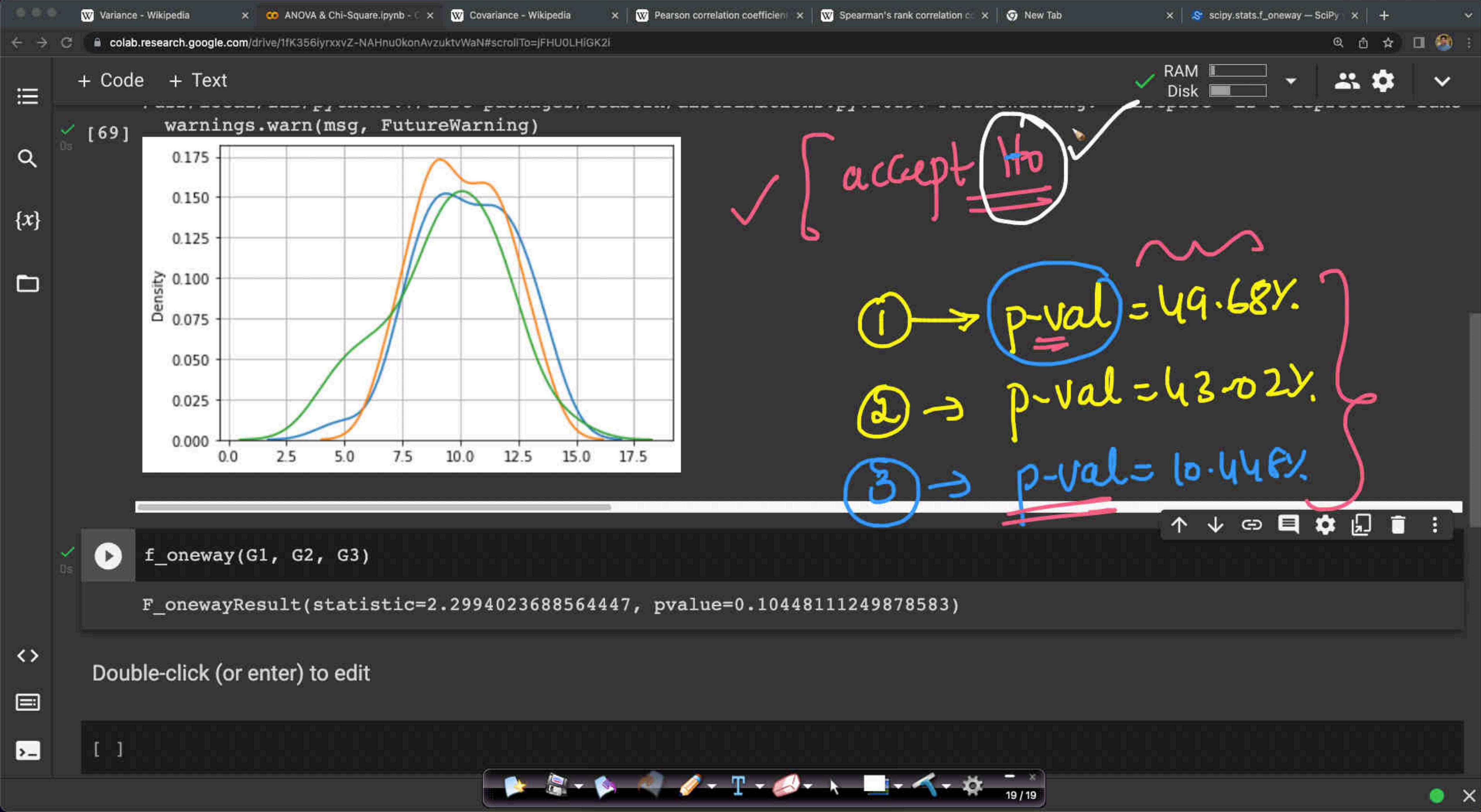


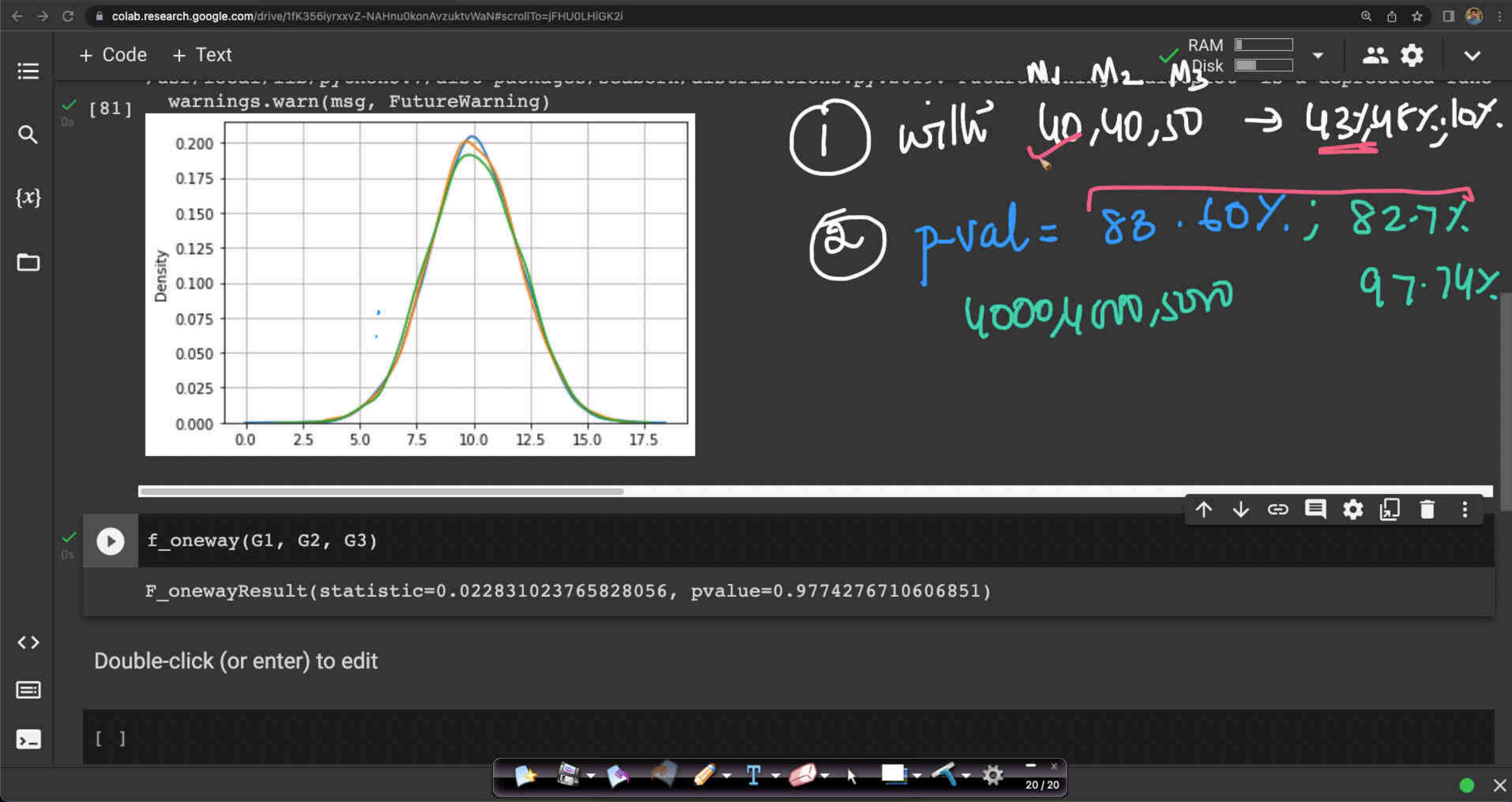
accept H₀

```
[64]: f_oneway(G1, G2, G3)
```

```
F_onewayResult(statistic=0.7034160770793035, pvalue=0.4968095155468297)
```







+ Code + Text

✓ RAM Disk

Playing with ANOVA

{x}

```
[61] from scipy.stats import f_oneway  
from scipy.stats import norm  
import seaborn as sns  
import matplotlib.pyplot as plt
```

```
[80] G1 = norm.rvs(loc=10.0, scale=2.0, size=4000)  
G2 = norm.rvs(loc=10.0, scale=2.01, size=4000)  
G3 = norm.rvs(loc=10.0, scale=2.03, size=5000)
```

```
▶ sns.distplot(G1, hist=False)  
sns.distplot(G2, hist=False)  
sns.distplot(G3, hist=False)  
plt.grid()
```

① as sample sizes ↑
ANOVA can detect small differences in μ_1, μ_2 & μ_3

② Variance in p-values depends on sample sizes.

```
<> /usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func  
warnings.warn(msg, FutureWarning)  
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func  
warnings.warn(msg, FutureWarning)  
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func  
warnings.warn(msg, FutureWarning)
```

+ Code + Text

✓ RAM Disk



Playing with ANOVA

```
[61] from scipy.stats import f_oneway  
     from scipy.stats import norm  
     import seaborn as sns  
     import matplotlib.pyplot as plt
```

```
[80] G1 = norm.rvs(loc=10.0, scale=2.0, size=4000)  
     G2 = norm.rvs(loc=10.0, scale=2.01, size=4000)  
     G3 = norm.rvs(loc=10.0, scale=2.03, size=5000)
```

```
[81] sns.distplot(G1, hist=False)  
     sns.distplot(G2, hist=False)  
     sns.distplot(G3, hist=False)  
     plt.grid()
```

```
<> /usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func  
      warnings.warn(msg, FutureWarning)  
      /usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func  
      warnings.warn(msg, FutureWarning)  
      /usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated func  
      warnings.warn(msg, FutureWarning)
```

together

{

Theory → concrete 'proofs'

=

✓ Simulation → help you build intuition

=

Likelihood 105.357

1

0.000

The resulting chi-square statistic is 102.596 with a p-value of .000. The 2X2 table also includes the expected values.

Remember the chi-square statistic is comparing the expected values to the observed values from Donna's study. The results of the chi-square indicate this difference (observed - expected is large). Thus, Donna can reject the null hypothesis that entrepreneurialism and geographic location are independent and she can conclude that Entrepreneurialism levels depend on geographic location.

	O ₁	O ₂
O ₁	1	0
O ₂	0	1

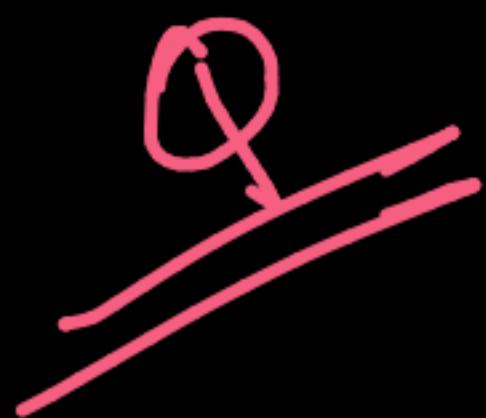
Conditions for Using the Chi-Square Test

Exercise caution when there are small expected counts. Minitab will give a count of the number of cells that have expected frequencies less than five. Some statisticians hesitate to use the chi-square test if more than 20% of the cells have expected frequencies below five, especially if the p-value is small and these cells give a large contribution to the total chi-square value.

Caution!

Sometimes researchers will categorize quantitative data (e.g., take height measurements and categorize as 'below average,' 'average,' and 'above average.') Doing so results in a loss of information - one cannot do the reverse of taking the categories and reproducing the raw quantitative measurements. Instead of categorizing, the data should be analyzed using quantitative methods.

	b ₁	b ₂
a ₁	0	0
a ₂	0	0



more the sample size ↑



more confident we are about H_0 or H_a accepting

less committed to errors

α : FP-error



$1 - \beta$ = power ↑

Φ

{ Maxmin ineq
Chebyshev ineq

P

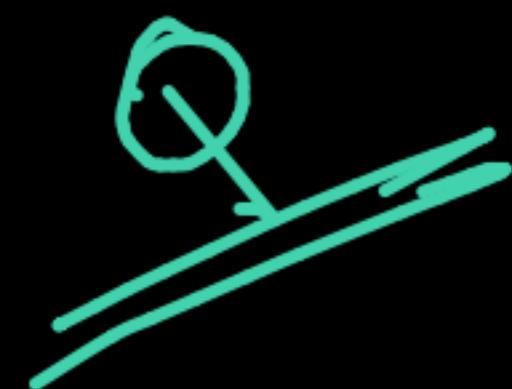
p-val > $\alpha = 5\%$

Threshold

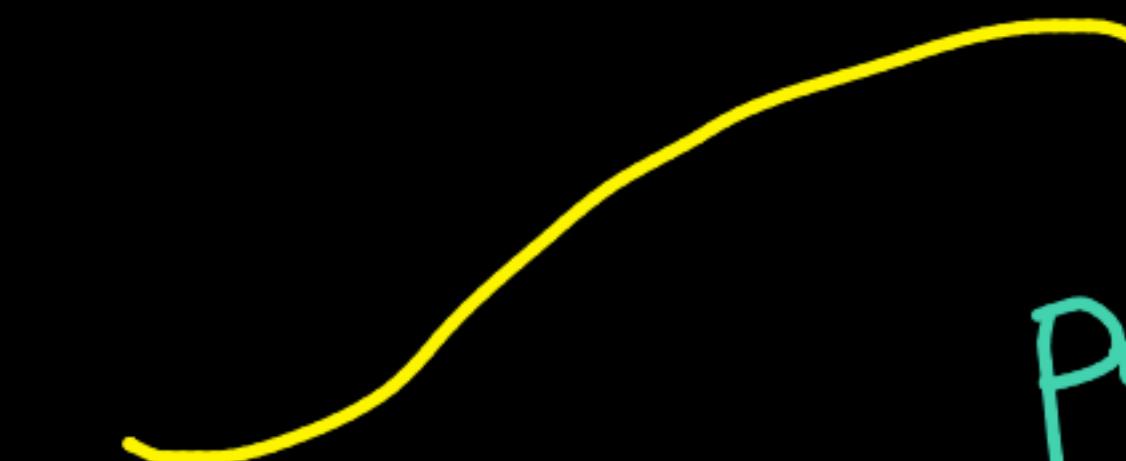
accept H_0

$\hookrightarrow P(T \text{ as extreme as } T_{\text{obs}} | H_0)$

P



①

 G_1 G_2 G_3 G_4 

~~Populations~~ \rightarrow normal

Samples \sim v-close Normal

④

② ~~This~~

$$\sigma_1^2 = \sigma_2^2 = \sigma_3^2$$

~~Py:~~

$$\underline{s_1^2} \approx \underline{s_2^2} \approx \underline{s_3^2}$$



Q

$$T_{\text{ANOVA}} = f = \frac{\sum_{i=1}^k (\bar{x}_i - \bar{x})^2}{\frac{1}{n-k} \sum_{i=1}^k \sum_{j=1}^m (x_{ij} - \bar{x}_i)^2}$$

Annotations:

- A red bracket on the right side of the equation groups the two terms in the denominator.
- A blue circle highlights the term $\sum_{i=1}^k$.
- A blue circle highlights the term $(\bar{x}_i - \bar{x})^2$.
- A blue circle highlights the term $n-k$.
- A blue circle highlights the term m .
- A blue circle highlights the term \bar{x}_i .
- A blue circle highlights the term x_{ij} .
- A blue circle highlights the term \bar{x}_i .

$f \sim F_{df}(k-1, n-k)$

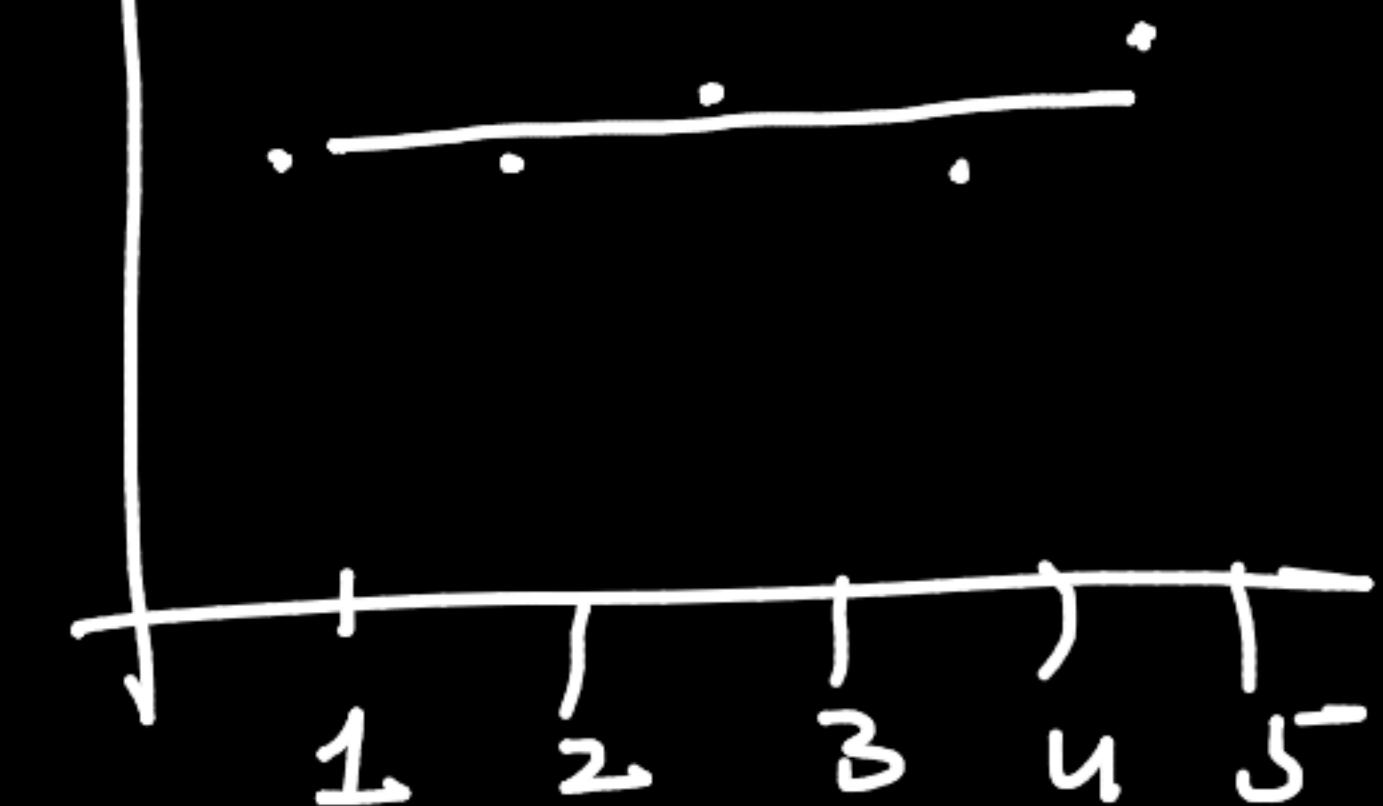
{ Task!

Hyp 2: # learners attending classes is "roughly" the same across the 5 days.

Data:

M	1	→	100
T	2	→	98
W	3	→	101
Th	4	→	96
Fr	5	→	103

hist



Obs

1	2	3	4	5
O ₁	O ₂	O ₃	O ₄	O ₅
100	98	101	96	103

exp under H₀

1	2	3	4	5
E ₁	E ₂	E ₃	E ₄	E ₅
100	100	100	100	100

mean

$$H_0: \alpha_1 = \alpha_2 = \alpha_3 = \alpha_4 = \alpha_5$$

Uniform dist

 $H_A:$ otherwise

$$\checkmark \quad \chi^2 = \sum_{i=1}^S \frac{(O_i - E_i)^2}{E_i}$$

$$\sim \chi^2 \text{ dist}(S-1)$$

χ^2 -test Goodness of fit dist

+ Code + Text

l J → 4 cells running

✓ RAM Disk



Chi-Square Goodness of fit

```
from scipy.stats import chisquare
from scipy.stats import poisson
import math
```

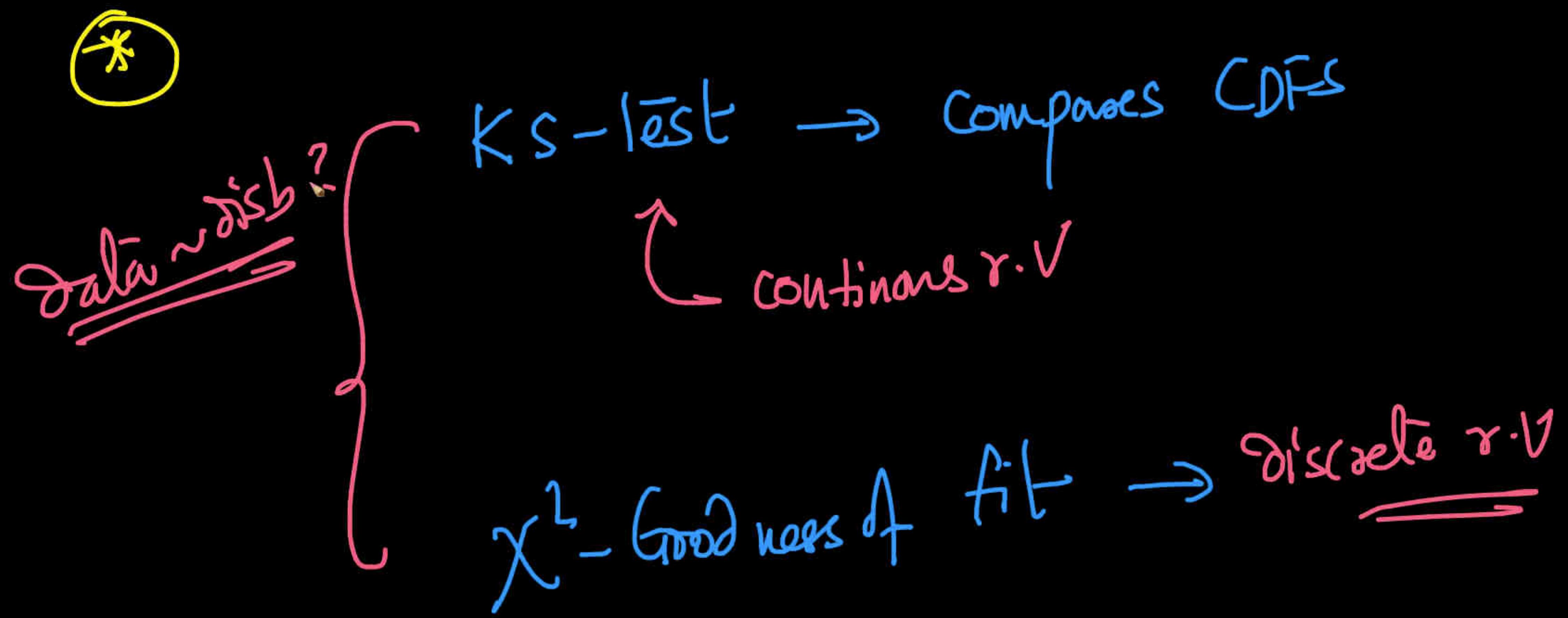
```
[41] p_x_i = poisson.pmf([0,1,2,3,4,5,6,7,8,9,10], mu=2.5)
```

```
[42] freq_exp= (p_x_i*100).round()
freq_exp
array([ 8., 21., 26., 21., 13., 7., 3., 1., 0., 0., 0.])
```

```
[47] freq_obs = [7, 22, 27, 20, 12, 7, 3, 1, 1, 0., 0.]
```

```
[46] chisquare(f_exp=freq_exp, f_obs=freq_obs)
```

```
/usr/local/lib/python3.7/dist-packages/scipy/stats/stats.py:6707: RuntimeWarning: invalid value encountered in true_
terms = (f_obs_float - f_exp)**2 / f_exp
```



Variance - Wikipedia x ANOVA & Chi-Square - Wikipedia x Covariance - Wikipedia x Pearson correlation coefficient x Spearman's rank correlation coefficient x New Tab x scipy.stats.f_onedf x Chi-Square test: Chi-Square Analysis x 6.2 - Chi-Square Test x + colab.research.google.com/drive/1fK356lyrxxVZ-NAHnu0konAvzuktvWaN#scrollTo=So7E5K2_K7MA

+ Code + Text RAM Disk

Chi-Square Goodness of fit

Football goals

$\lambda = 2.5$

```
from scipy.stats import chisquare
from scipy.stats import poisson
import math
```

[41] p_x_i = poisson.pmf([0,1,2,3,4,5,6,7,8,9,10], mu=2.5)

[42] freq_exp= (p_x_i*100).round()
freq_exp

array([8., 21., 26., 21., 13., 7., 3., 1., 0., 0., 0.])

[47] freq_obs = [7, 22, 27, 20, 12, 7, 3, 1, 1, 0., 0.]

[46] chisquare(f_exp=freq_exp, f_obs=freq_obs)

/usr/local/lib/python3.7/dist-packages/scipy/stats/stats.py:6707: RuntimeWarning: invalid value encountered in true_
terms = (f_obs_float - f_exp)**2 / f_exp
Power_divergenceResult(statistic=nan, pvalue=nan)

RAM Disk

35 / 35

W Variance - Wikipedia x ANOVA & Chi-Square - Wikipedia x W Covariance - Wikipedia x W Pearson correlation coefficient x W Spearman's rank correlation coefficient x New Tab x | S scipy.stats.f_onedf x | C Chi-Square test: Chi-Square Analysis x | 6.2 - Chi-Square Test x +

RAM Disk

Chi-Square Goodness of fit

```
from scipy.stats import chisquare
from scipy.stats import poisson
import math
```

```
[41] p_x_i = poisson.pmf([0,1,2,3,4,5,6,7,8,9,10], mu=2.5)
```

```
[42] freq_exp= (p_x_i*100).round()
freq_exp
```

```
array([ 8., 21., 26., 21., 13., 7., 3., 1., 0., 0., 0.])
```

```
[47] freq_obs = [7, 22, 27, 20, 12, 7, 3, 1, 1, 0., 0.]
```

```
[46] chisquare(f_exp=freq_exp, f_obs=freq_obs)
```

```
/usr/local/lib/python3.7/dist-packages/scipy/stats/stats.py:6707: RuntimeWarning: invalid value encountered in true_
    terms = (f_obs_float - f_exp)**2 / f_exp
Power_divergenceResult(statistic=nan, pvalue=nan)
```

$$X \sim \text{Poisson}(\lambda=2.5)$$

$$P(X=0) \rightarrow$$

$$P(X=1) \rightarrow$$

'.'

Variance - Wikipedia x ANOVA & Chi-Square - Wikipedia x Covariance - Wikipedia x Pearson correlation coefficient x Spearman's rank correlation coefficient x New Tab x scipy.stats.f_onedf x Chi-Square test: Chi-Square Analysis x 6.2 - Chi-Square Test x + colab.research.google.com/drive/1fK356lyrxvZ-NAHnu0konAvzuktvWaN#scrollTo=M9BPzgqPL2N0

+ Code + Text RAM Disk

[83] from scipy.stats import chisquare
from scipy.stats import poisson
import math

[85] p_x_i = poisson.pmf([0,1,2,3,4,5,6,7,8,9,10], mu=2.5)
p_x_i

array([8.20849986e-02, 2.05212497e-01, 2.56515621e-01, 2.13763017e-01,
1.33601886e-01, 6.68009429e-02, 2.78337262e-02, 9.94061650e-03,
3.10644266e-03, 8.62900738e-04, 2.15725184e-04])

freq_exp = (p_x_i * 100).round()
freq_exp

array([8., 21., 26., 21., 13., 7., 3., 1., 0., 0., 0.])

[47] freq_obs = [7, 22, 27, 20, 12, 7, 3, 1, 1, 0., 0.]

[46] chisquare(f_exp=freq_exp, f_obs=freq_obs)

/usr/local/lib/python3.7/dist-packages/scipy/stats/stats.py:6707: RuntimeWarning: invalid value encountered in true_

terms = (f_obs_float - f_exp)**2 / f_exp

Power_divergenceResult(statistic=nan, pvalue=nan)

$P(X=2) \times 100 =$

freq_exp = (p_x_i * 100).round()

freq_exp

array([8., 21., 26., 21., 13., 7., 3., 1., 0., 0., 0.])

freq_obs = [7, 22, 27, 20, 12, 7, 3, 1, 1, 0., 0.]

chisquare(f_exp=freq_exp, f_obs=freq_obs)

/usr/local/lib/python3.7/dist-packages/scipy/stats/stats.py:6707: RuntimeWarning: invalid value encountered in true_

terms = (f_obs_float - f_exp)**2 / f_exp

Power_divergenceResult(statistic=nan, pvalue=nan)

Variance - Wikipedia x ANOVA & Chi-Square - Wikipedia x Covariance - Wikipedia x Pearson correlation coefficient x Spearman's rank correlation coefficient x New Tab x scipy.stats.f_onedf x Chi-Square test: Chi-Square Analysis x Chi Square Analysis x 6.2 - Chi-Square Test for Independence x + colab.research.google.com/drive/1fK356lyrxvZ-NAHnu0konAvzuktvWaN#scrollTo=M9BPzgqPL2N0

+ Code + Text RAM Disk

[83] from scipy.stats import chisquare
from scipy.stats import poisson
import math

{x}

[85] p_x_i = poisson.pmf([0,1,2,3,4,5,6,7,8,9,10], mu=2.5)

p_x_i

array([8.20849986e-02, 2.05212497e-01, 2.56515621e-01, 2.13763017e-01, 1.33601886e-01, 6.68009429e-02, 2.78337262e-02, 9.94061650e-03, 3.10644266e-03, 8.62900738e-04, 2.15725184e-04])

freq_exp= (p_x_i*100).round()
freq_exp

array([8., 21., 26., 21., 13., 7., 3., 1., 0., 0., 0.])

freq-UP

[47] freq_obs = [7, 22, 27, 20, 12, 7, 3, 1, 1, 0., 0.]

[46] chisquare(f_exp=freq_exp, f_obs=freq_obs)

/usr/local/lib/python3.7/dist-packages/scipy/stats/stats.py:6707: RuntimeWarning: invalid value encountered in true_divide
terms = (f_obs_float - f_exp)**2 / f_exp
Power_divergenceResult(statistic=nan, pvalue=nan)

38 / 38

+ Code + Text

Chi-Square Goodness of fit

discrete r.v

{x}

```
[83] from scipy.stats import chisquare
     from scipy.stats import poisson
     import math
```

□

```
[85] p_x_i = poisson.pmf([0,1,2,3,4,5,6,7,8,9,10], mu=2.5)
     p_x_i
```

```
array([8.20849986e-02, 2.05212497e-01, 2.56515621e-01, 2.13763017e-01,
       1.33601886e-01, 6.68009429e-02, 2.78337262e-02, 9.94061650e-03,
       3.10644266e-03, 8.62900738e-04, 2.15725184e-04])
```



Ds



```
freq_exp= (p_x_i*100).round()
freq_exp
```

```
array([ 8., 21., 26., 21., 13., 7., 3., 1., 0., 0., 0.])
```

<>

Ds

```
[47] freq_obs = [7, 22, 27, 20, 12, 7, 3, 1, 1, 0., 0.]
```

>-

```
[46] chisquare(f_exp=freq_exp, f_obs=freq_obs)
```

Variance - Wikipedia x ANOVA & Chi-Square - Wikipedia x Covariance - Wikipedia x Pearson correlation coefficient x Spearman's rank correlation coefficient x New Tab x scipy.stats.f_onedf x Chi-Square test: Chi-Square Analysis x 6.2 - Chi-Square Test x + colab.research.google.com/drive/1fK356lyrxxvZ-NAHnu0konAvzuktvWaN#scrollTo=M9BPzgqPL2N0

+ Code + Text

[83] import math

[85] p_x_i = poisson.pmf([0,1,2,3,4,5,6,7,8,9,10], mu=2.5)
p_x_i

array([8.20849986e-02, 2.05212497e-01, 2.56515621e-01, 2.13763017e-01, 1.33601886e-01, 6.68009429e-02, 2.78337262e-02, 9.94061650e-03, 3.10644266e-03, 8.62900738e-04, 2.15725184e-04])

[86] freq_exp= (p_x_i*100).round()
freq_exp

array([8., 21., 26., 21., 13., 7., 3., 1., 0., 0., 0.])
0 1 2 3 4 5 6 7 8 9 10

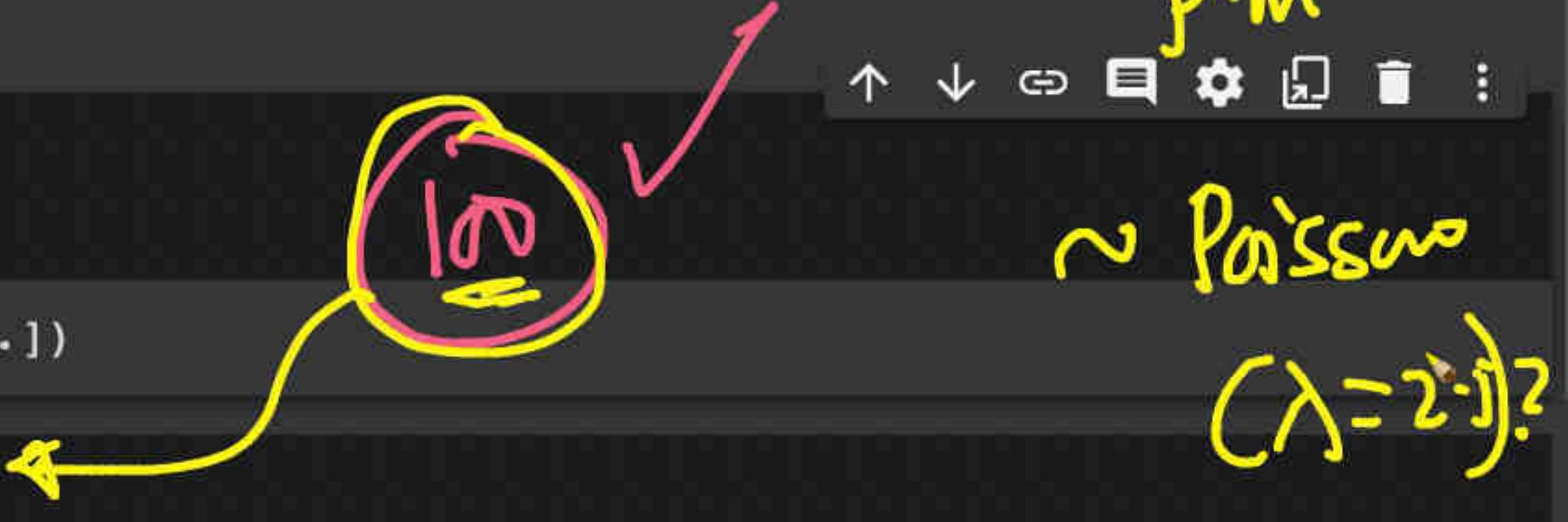
[47] freq_obs = [7, 22, 27, 20, 12, 7, 3, 1, 1, 0., 0.]
7 22 27 20 12 7 3 1 1 0 0

[46] chisquare(f_exp=freq_exp, f_obs=freq_obs)

/usr/local/lib/python3.7/dist-packages/scipy/stats/stats.py:6707: RuntimeWarning: invalid value encountered in true_divide
terms = (f_obs_float - f_exp)**2 / f_exp
Power_divergenceResult(statistic=nan, pvalue=nan)

[52] freq_exp = [8., 21., 26., 21., 13., 7., 3., 1.]

Task: Can we conclude that #goals ~ Poisson ($\lambda = 2.5$)?



RAM Disk

40/40

Variance - Wikipedia x ANOVA & Chi-Square - Wikipedia x Covariance - Wikipedia x Pearson correlation coefficient x Spearman's rank correlation coefficient x New Tab x scipy.stats.f_onedf x Chi-Square test: Chi-Square Analysis x Chi Square Analysis x 6.2 - Chi-Square Test x + colab.research.google.com/drive/1fK356lyrxvZ-NAHnu0konAvzuktvWaN#scrollTo=M9BPzgqPL2N0

+ Code + Text RAM Disk

[83] from scipy.stats import chisquare
from scipy.stats import poisson
import math

[85] p_x_i = poisson.pmf([0,1,2,3,4,5,6,7,8,9,10], mu=2.5)

$P(X=0) \propto 1/\theta$

$P(X=1) \propto 1/\theta$

array([8.20849986e-02, 2.05212497e-01, 2.56515621e-01, 2.13763017e-01, 1.33601886e-01, 6.68009429e-02, 2.78337262e-02, 9.94061650e-03, 3.10644266e-03, 8.62900738e-04, 2.15725184e-04])

freq_exp = (p_x_i*100).round()
freq_exp

array([8., 21., 26., 21., 13., 7., 3., 1., 0., 0., 0.])

$P(X=i) \propto 1/\theta$

[47] freq_obs = [7, 22, 27, 20, 12, 7, 3, 1, 1, 0., 0.]

[46] chisquare(f_exp=freq_exp, f_obs=freq_obs)

/usr/local/lib/python3.7/dist-packages/scipy/stats/stats.py:6707: RuntimeWarning: invalid value encountered in true_

terms = (f_obs_float - f_exp)**2 / f_exp

Power divergenceResult + / statistics

41/41

Variance - Wikipedia x ANOVA & Chi-Square - Wikipedia x Covariance - Wikipedia x Pearson correlation coefficient x Spearman's rank correlation coefficient x New Tab x scipy.stats.f_onedf x Chi-Square test: Chi-Square Analysis x Chi Square Analysis x 6.2 - Chi-Square Test x + colab.research.google.com/drive/1fK356lyrxvZ-NAHnu0konAvzuktvWaN#scrollTo=M9BPzgqPL2N0

+ Code + Text RAM Disk

[83] from scipy.stats import chisquare
from scipy.stats import poisson
import math

[85] p_x_i = poisson.pmf([0,1,2,3,4,5,6,7,8,9,10], mu=2.5)
p_x_i

array([8.20849986e-02, 2.05212497e-01, 2.56515621e-01, 2.13763017e-01,
1.33601886e-01, 6.68009429e-02, 2.78337262e-02, 9.94061650e-03,
3.10644266e-03, 8.62900738e-04, 2.15725184e-04])

freq_exp=(p_x_i*100).round()
freq_exp

array([8., 21., 26., 21., 13., 7., 3., 1., 0., 0., 0.])

freq_obs = [7, 22, 27, 20, 12, 7, 3, 1, 1, 0., 0.]

chisquare(f_exp=freq_exp, f_obs=freq_obs)

/usr/local/lib/python3.7/dist-packages/scipy/stats/stats.py:6707: RuntimeWarning: invalid value encountered in true_

terms = (f_obs_float - f_exp)**2 / f_exp

Power divergenceResult + / statistics

42/42

+ Code + Text

✓ RAM Disk



[83] import math

[85] p_x_i = poisson.pmf([0,1,2,3,4,5,6,7,8,9,10], mu=2.5)
p_x_iarray([8.20849986e-02, 2.05212497e-01, 2.56515621e-01, 2.13763017e-01,
1.33601886e-01, 6.68009429e-02, 2.78337262e-02, 9.94061650e-03,
3.10644266e-03, 8.62900738e-04, 2.15725184e-04])

$$T = \sum_{i=0}^{10} \frac{(O_i - E_i)^2}{E_i}$$

[86] freq_exp= (p_x_i*100).round()
freq_exp

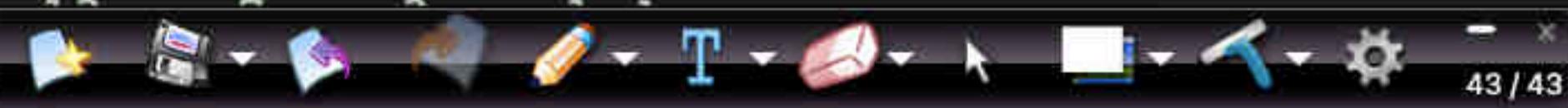
array([8., 21., 26., 21., 13., 7., 3., 1., 0., 0., 0.])

[47] freq_obs = [7, 22, 27, 20, 12, 7, 3, 1, 1, 0., 0.]

[46] chisquare(f_exp=freq_exp, f_obs=freq_obs)

{ /usr/local/lib/python3.7/dist-packages/scipy/stats/stats.py:6707: RuntimeWarning: invalid value encountered in true_ terms = (f_obs_float - f_exp)**2 / f_exp
Power_divergenceResult(statistic=nan, pvalue=nan)}

[52] freq_exp = [8., 21., 26., 21.



+ Code + Text

✓ RAM Disk

[52] freq_exp = [8., 21., 26., 21., 13., 7., 3., 1.]
freq_obs = [7, 22, 27, 20, 12, 8, 3, 1]

→ ≈ 17.5% = 0.05

{x}

□

[53] chisquare(f_exp=freq_exp, f_obs=freq_obs)

Power_divergenceResult(statistic=0.47847985347985345, pvalue=0.9995215264692844)

[]

Double-click (or enter) to edit

H_0 ✓

#goals P.M ~ Pois($\lambda=2.5$)

action \Rightarrow

accept $H_0 \}$ ✓
↙ ~~H_0~~

✓ failed to reject $H_0 \} \xrightarrow{P\text{-val}=0.999999}$
↙ ~~H_0~~



Cont - $\tau \cdot V$

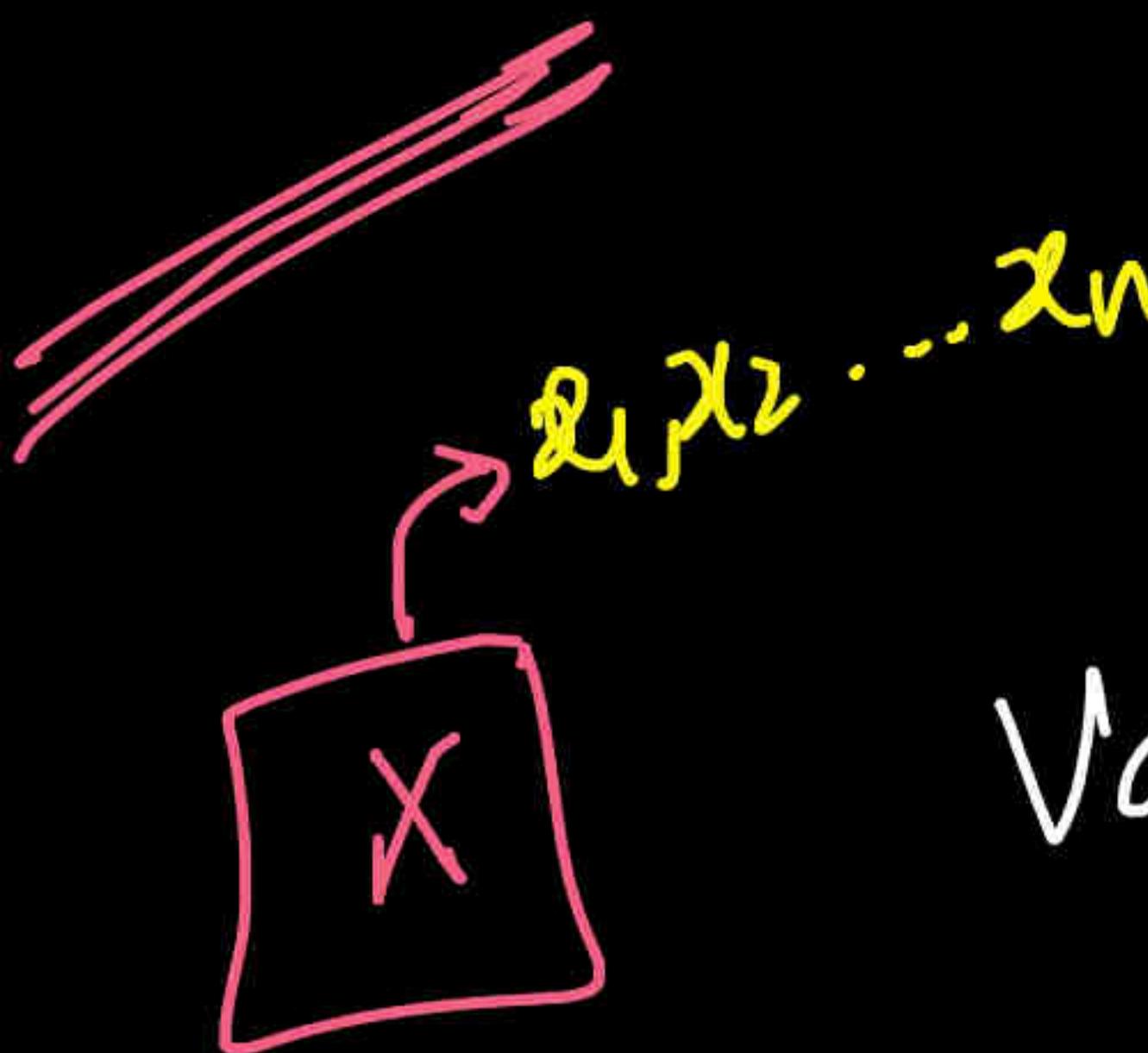
$\chi^2_{\text{good fit}}$

$\rightarrow \underline{\text{CS-test}}$



$\underline{\text{AD-test}}$





Variance

$$\text{Variance} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

Covariance



	x_i	y_i
1	x_1	y_1
2	x_2	y_2
3	x_3	y_3
:		
n	x_n	y_n

$X \downarrow \bar{x}$ $Y \downarrow \bar{y}$



$$\text{Cov}(X, Y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

Annotations:

- $h_i - \bar{h}$ and $w_i - \bar{w}$ are shown in blue, indicating the deviation from the regression line.
- +ve indicates a positive correlation.
- avg indicates the average value of $x_i - \bar{x}$ and $y_i - \bar{y}$.

Case 1:

$$1 \rightarrow (h_1, \omega_1)$$

$$h_1 > \bar{h}$$

$$\omega_1 > \bar{\omega}$$

 Cov ↑

Case 2:

$$h_2 < \bar{h}$$

-ve

$$\omega_2 > \bar{\omega}$$

+ve

 Cov ↓

Case 3:

$$h_2 > \bar{h}$$

+ve

$$\omega_2 < \bar{\omega}$$

-ve

 Cov ↓

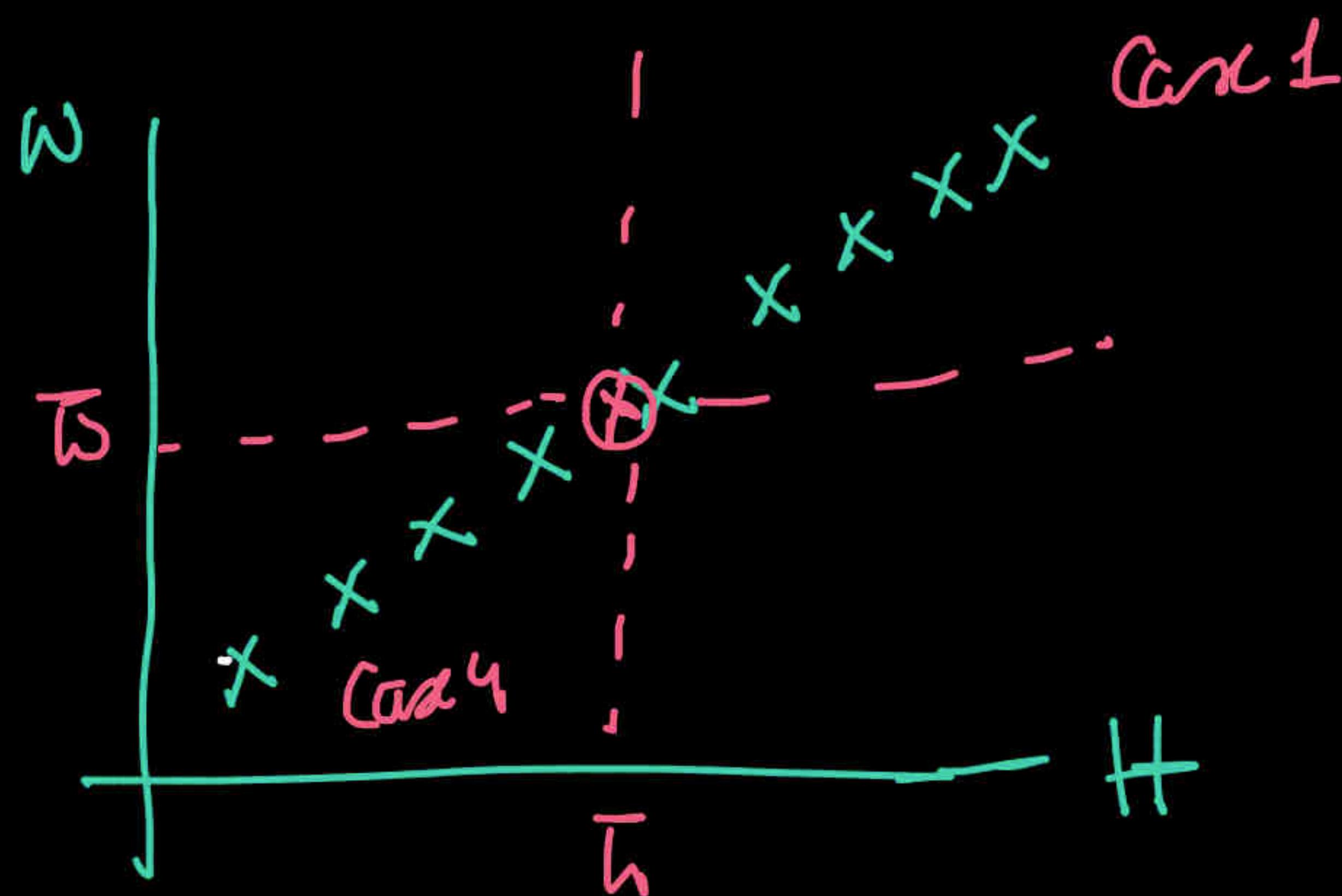
Cash u:

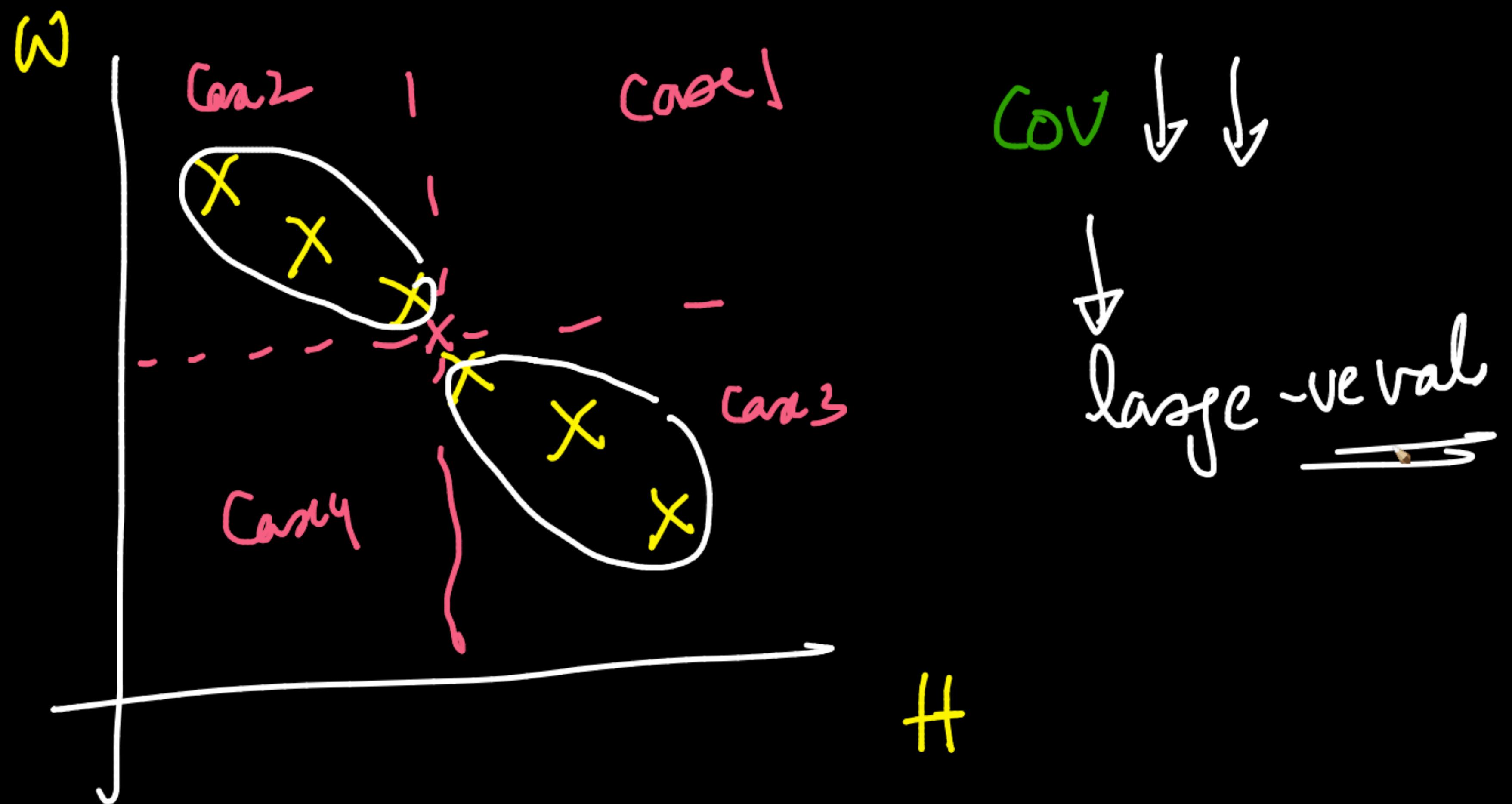
$$h_1 < \bar{h}$$

$$w_1 < \bar{w}$$

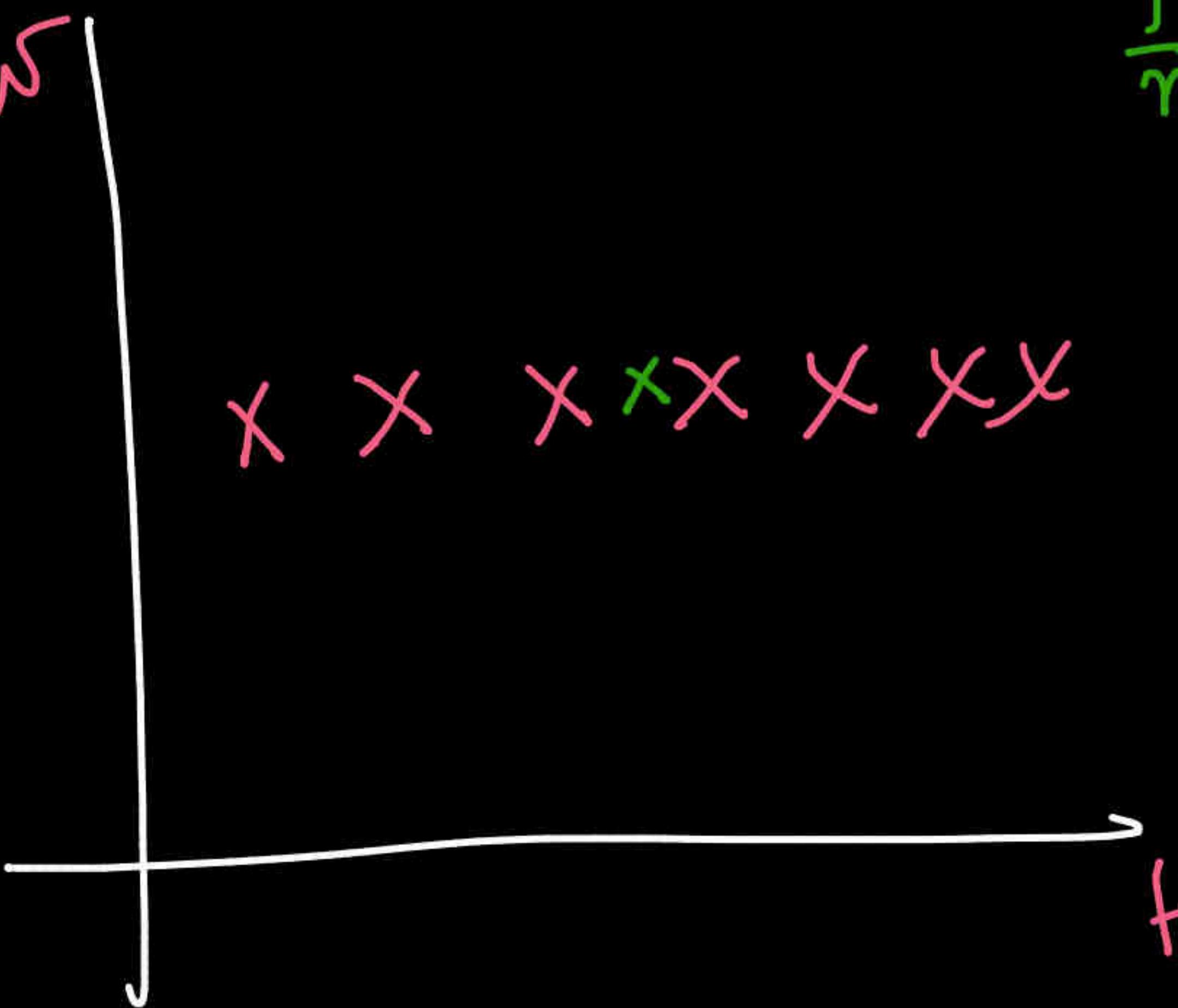
Cov ↑

Cov. high





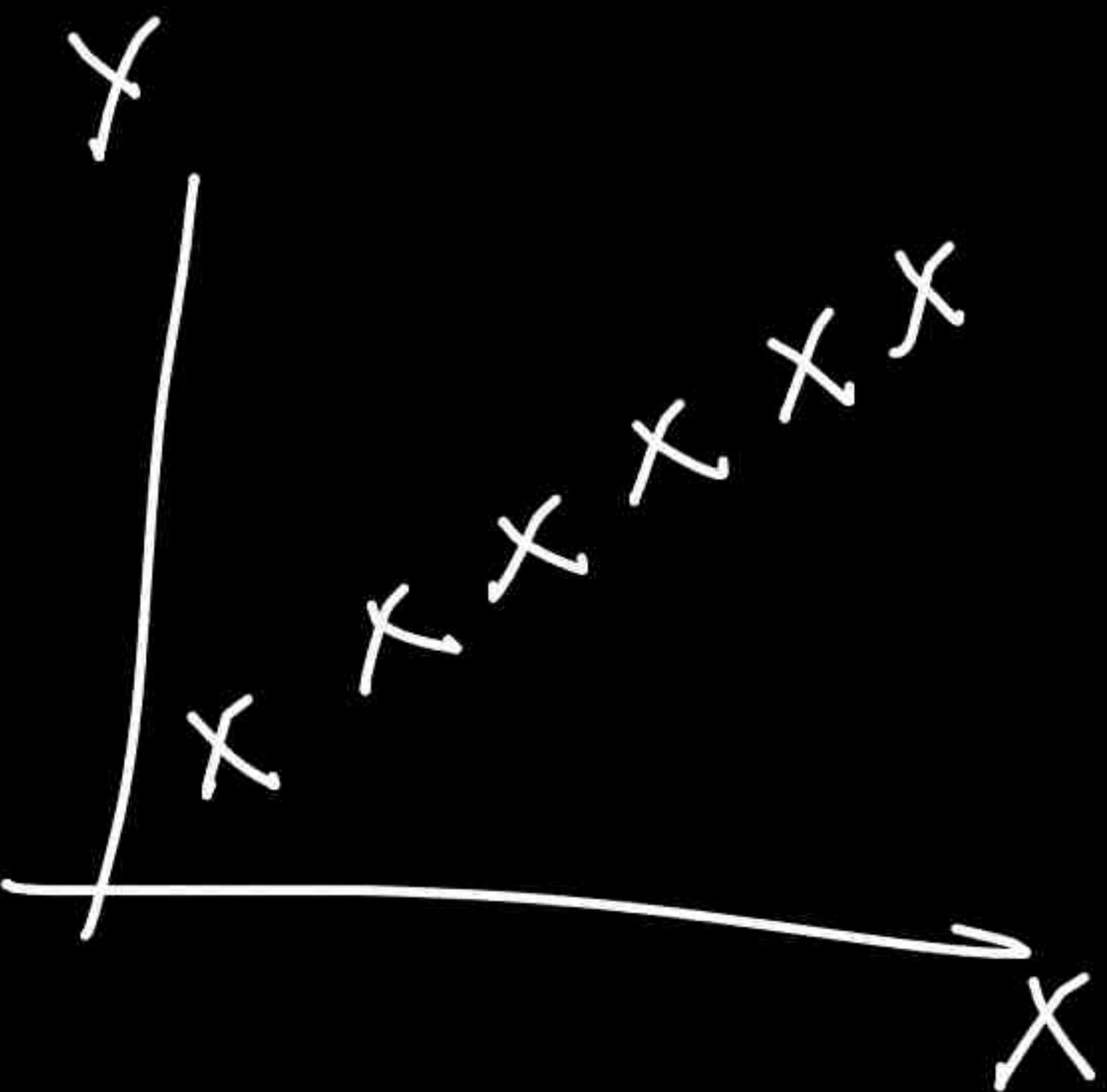
$y = \bar{w}$



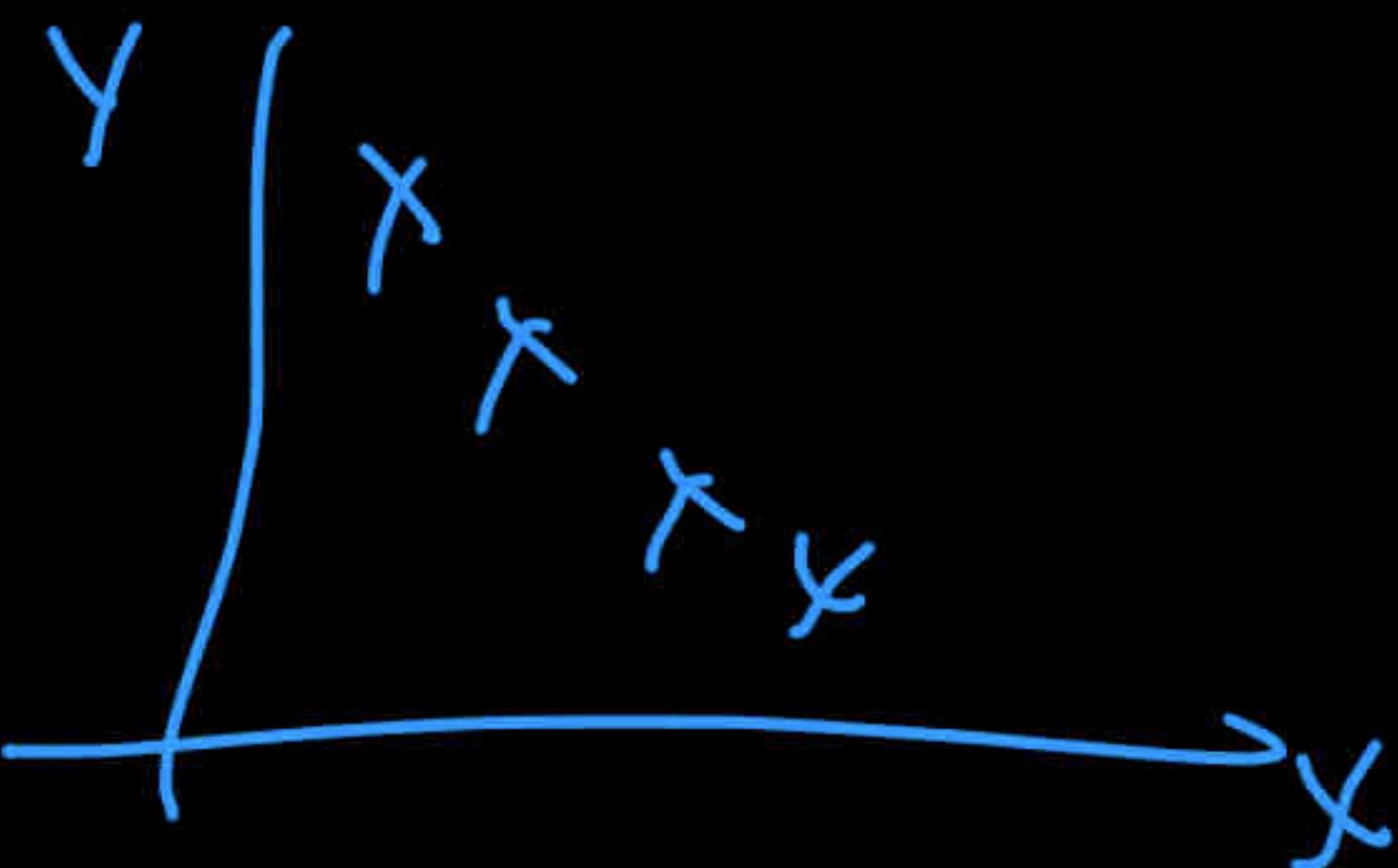
$$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

\Downarrow
0

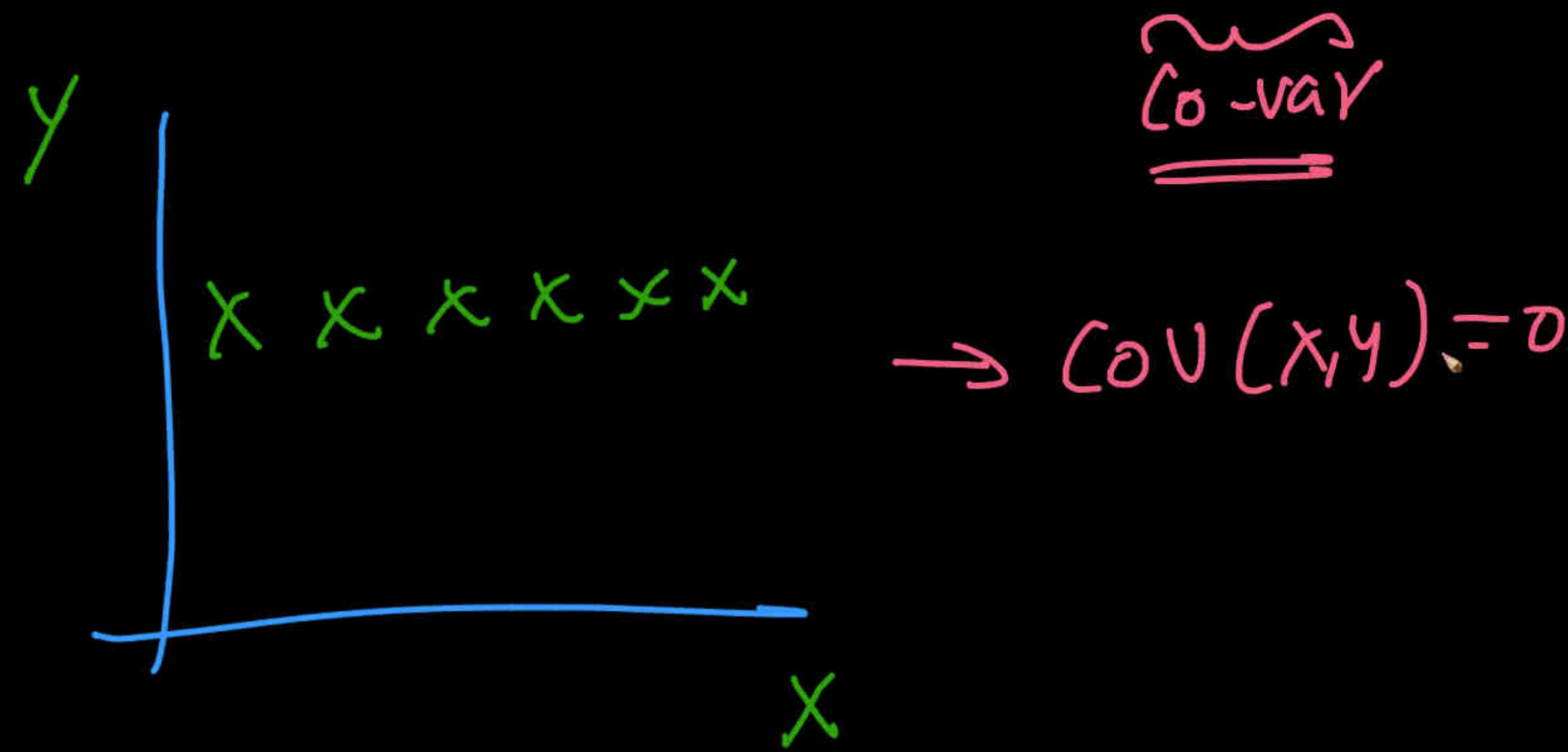
$\text{Cov} = 0$

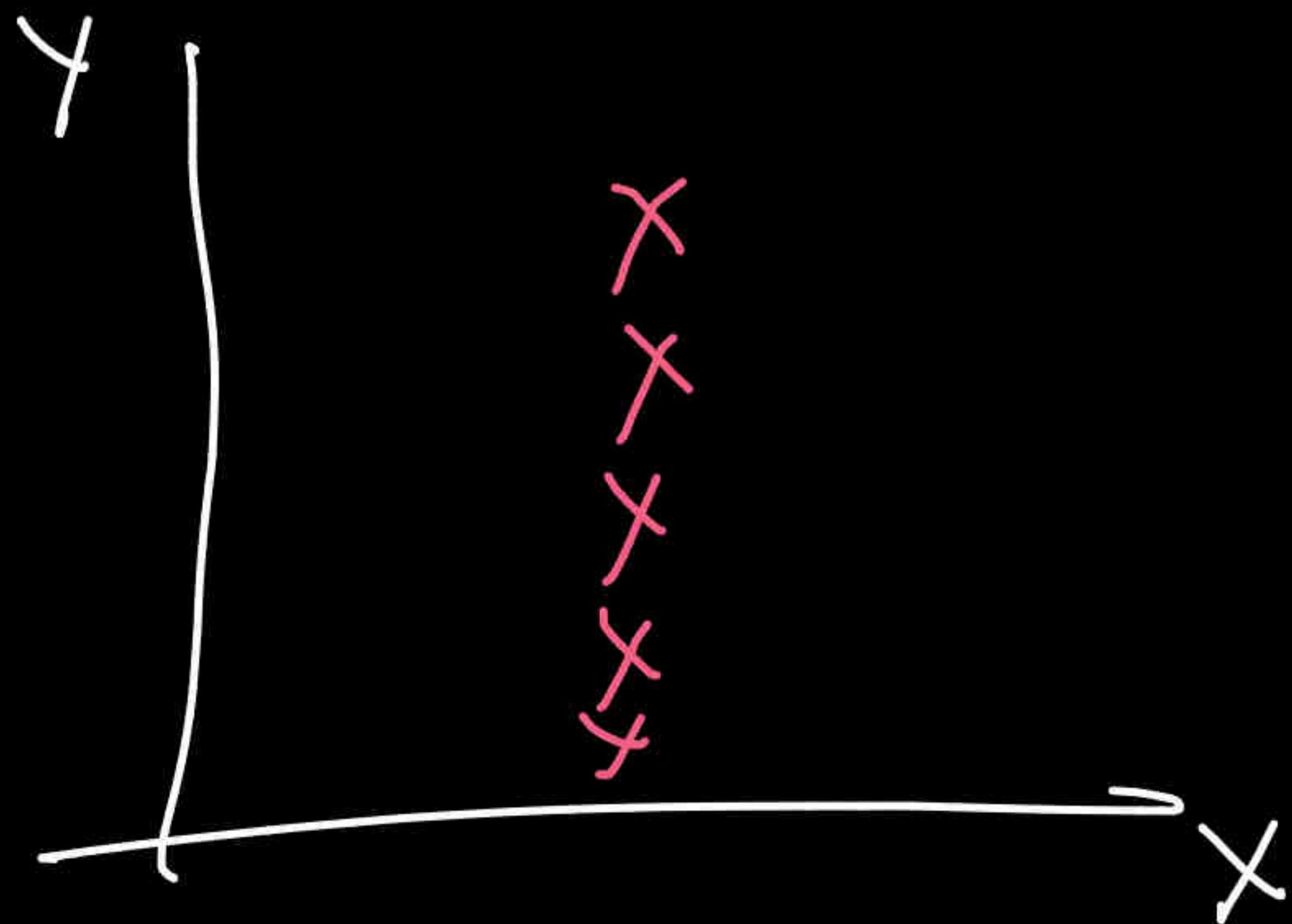


$\text{Cov}(x,y) \uparrow$: large +ve



$\text{Cov}(x,y) \downarrow$: large -ve





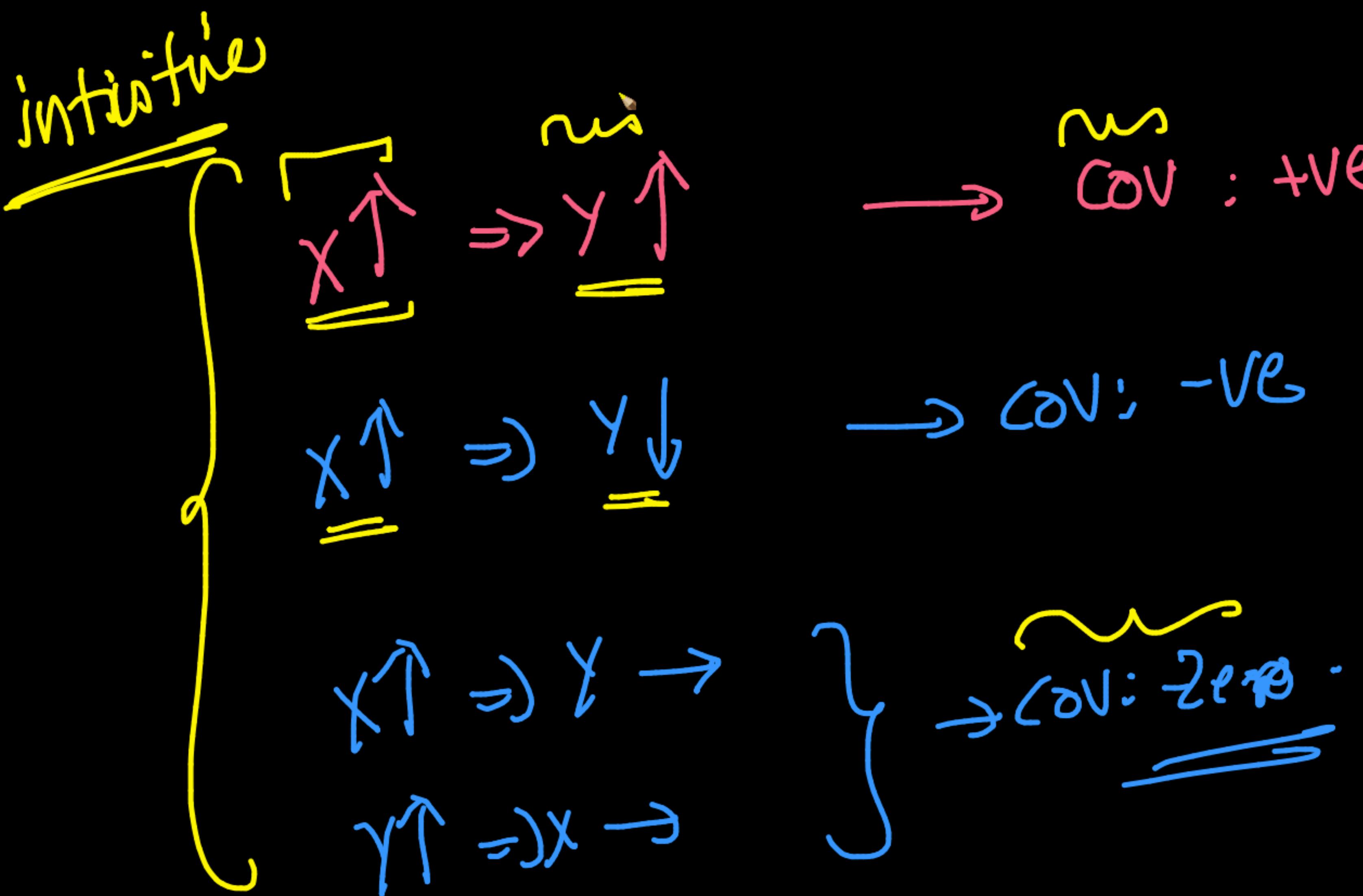
$$\text{Cov}(X, Y) = 0$$

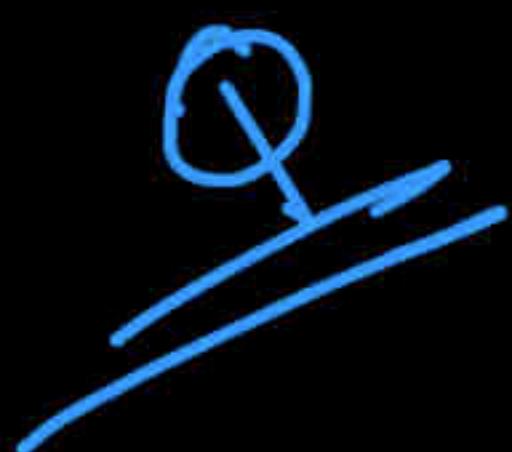
=

{ mean, Var; std-dev

$$\text{Covar} \leftarrow \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

\uparrow
 $\underbrace{x, y, z}$





$$\frac{1}{n} \sum_{j=1}^n (\tilde{x}_j - \bar{x})^2 \cdot (\tilde{y}_j - \bar{y})^2$$

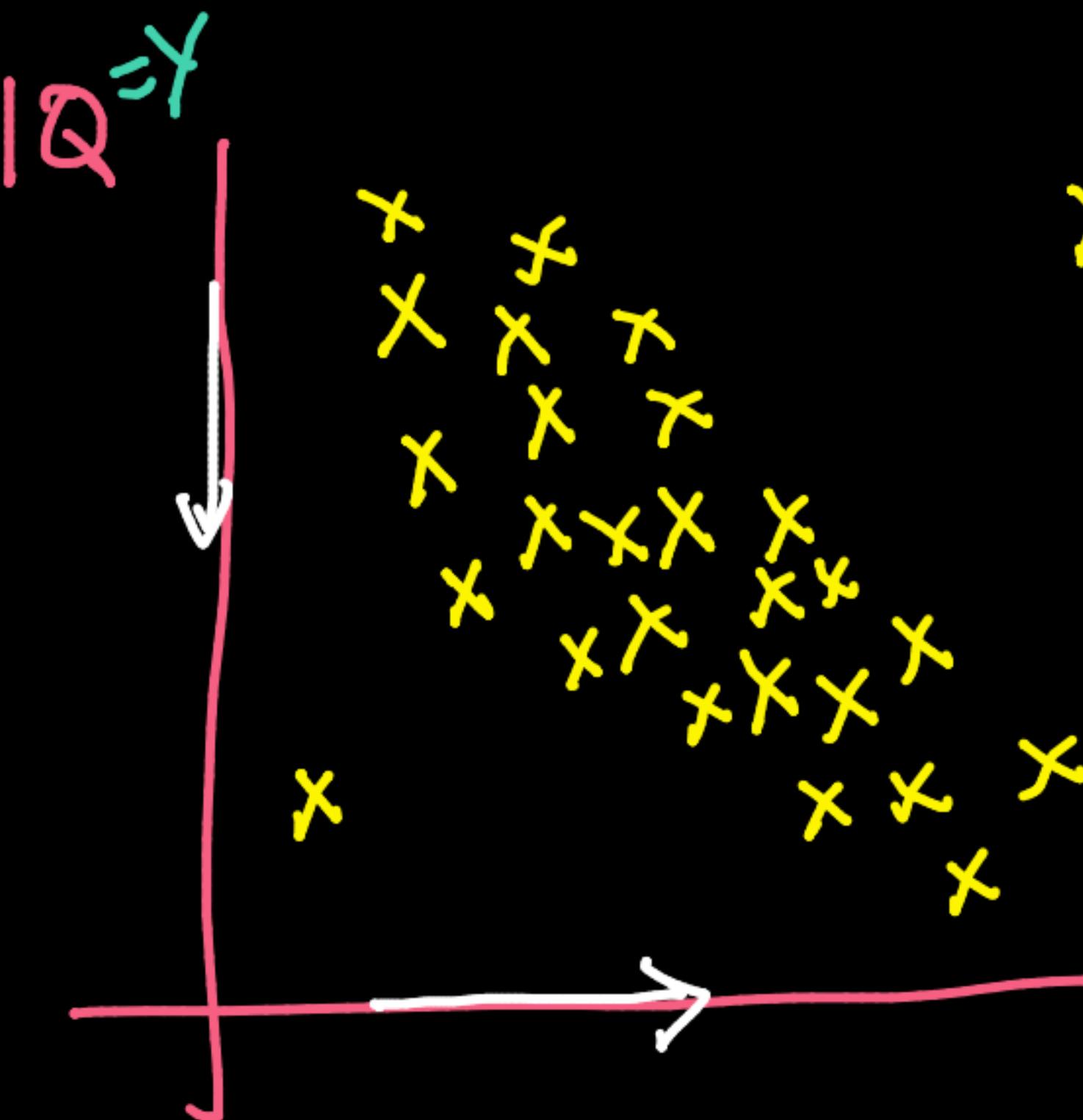
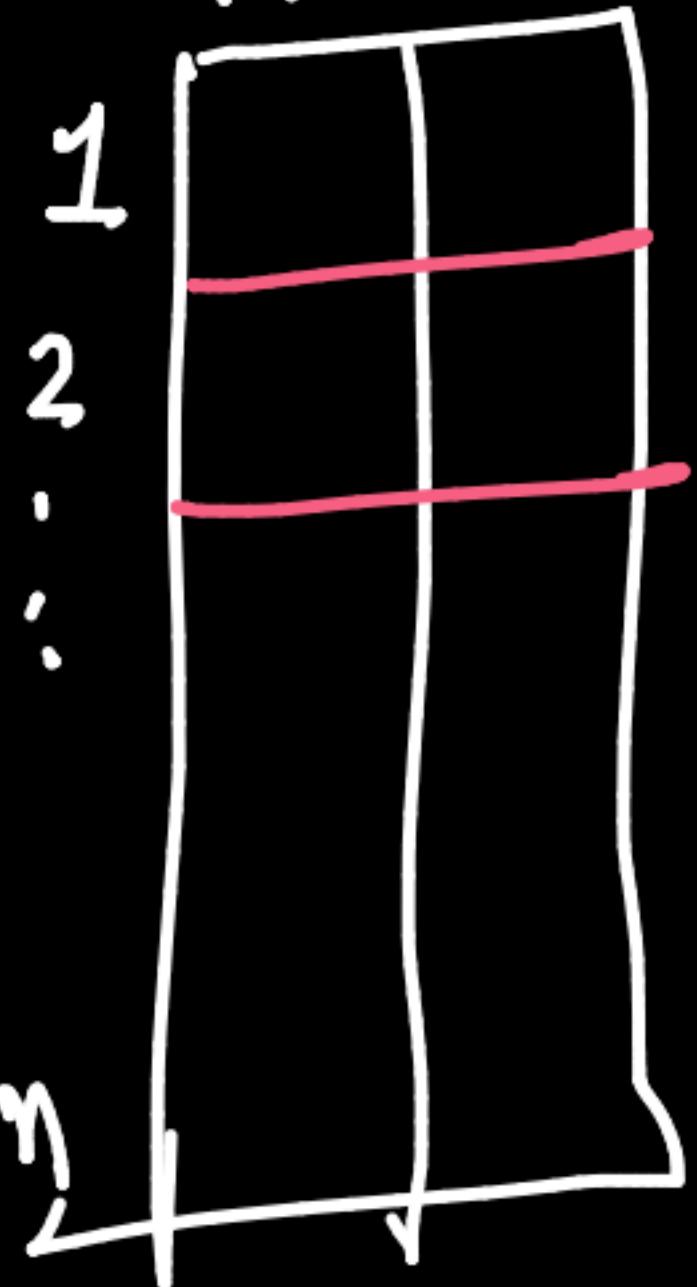
$\text{cov}(x,y)$ \rightarrow Correlation

a. $\text{Cov} > 0$

b. $\text{Cov} < 0$

c. $\text{Cov} \approx 0$

TV-has IQ



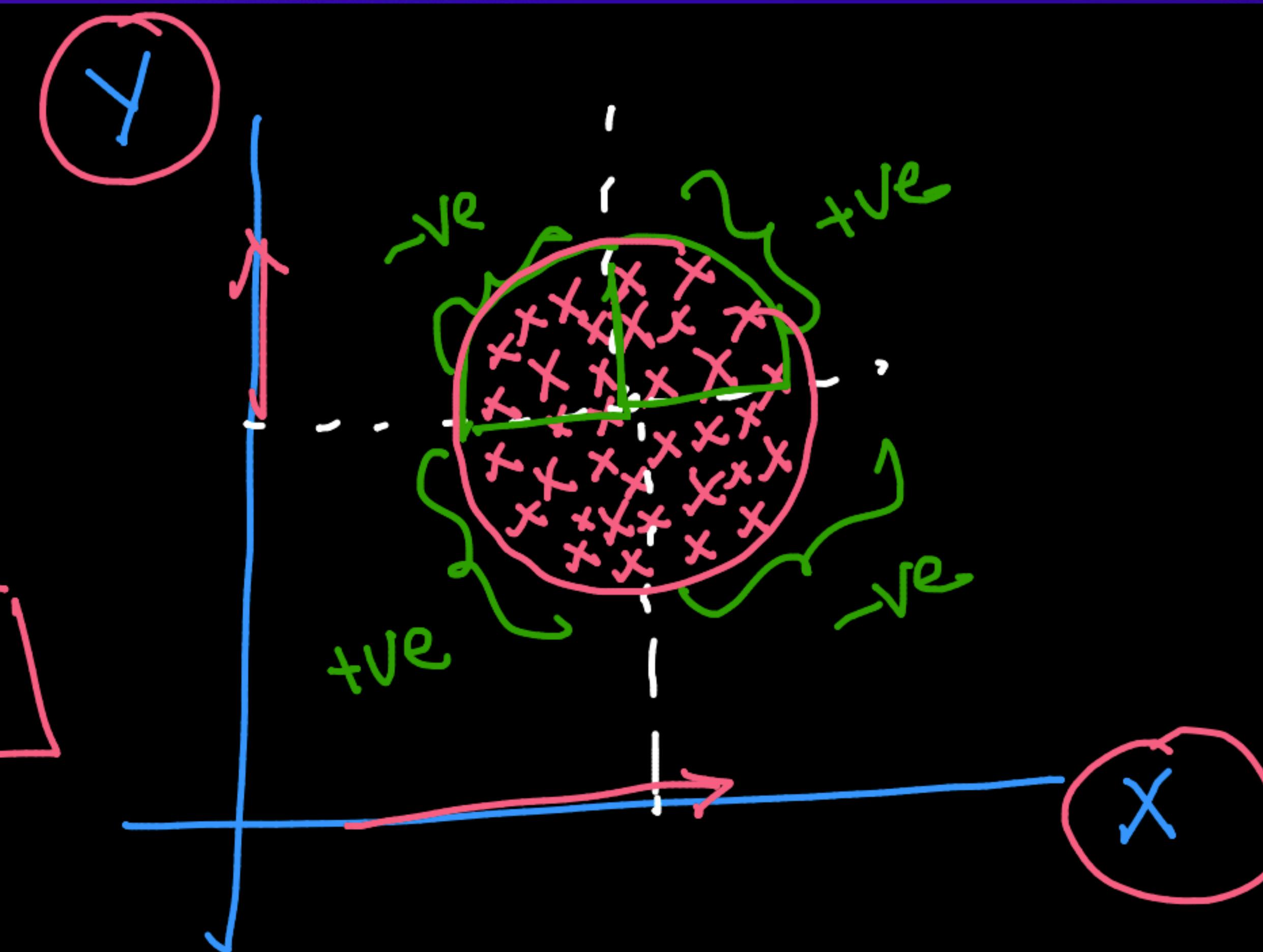
$$\text{Cov}(x, y) < 0$$

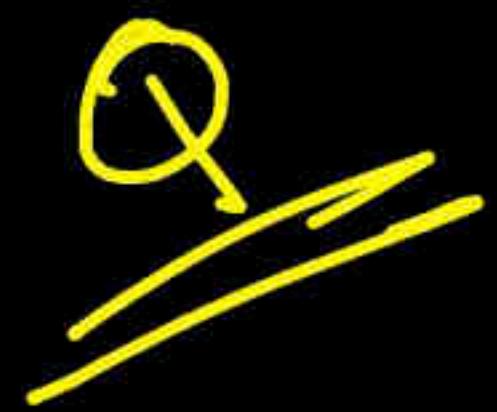
Task: as $\# \text{has TV} \uparrow$ does $|Q| \downarrow$

$$\text{a: } \text{Cov}(x_1 y) > 0$$

$$\text{b: } \text{Cov}(x_1 y) < 0$$

$$\text{c: } \text{Cov}(x_1 y) \approx 0$$





$$\text{Cov}(x, x) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})$$

↳ Var.(x)

Hoeffding's covariance identity [edit]

A useful identity to compute the covariance between two random variables X, Y is the Hoeffding's covariance identity:^[7]

$$\text{cov}(X, Y) = \int_{\mathbb{R}} \int_{\mathbb{R}} (F_{(X,Y)}(x, y) - F_X(x)F_Y(y)) dx dy$$

where $F_{(X,Y)}(x, y)$ is the joint cumulative distribution function of the random vector (X, Y) and $F_X(x), F_Y(y)$ are the [marginals](#).

Uncorrelatedness and independence [edit]

Main article: Correlation and dependence

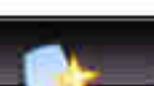
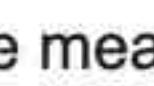
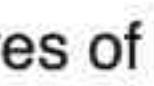
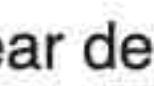
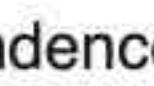
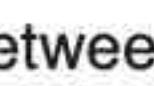
Random variables whose covariance is zero are called [uncorrelated](#).^{[4]:p. 121} Similarly, the components of random vectors whose [covariance matrix](#) is zero in every entry outside the main diagonal are also called uncorrelated.

If X and Y are [independent random variables](#), then their covariance is zero.^{[4]:p. 123[8]} This follows because under independence,

$$\text{E}[XY] = \text{E}[X] \cdot \text{E}[Y].$$

The converse, however, is not generally true. For example, let X be uniformly distributed in $[-1, 1]$ and let $Y = X^2$. Clearly, X and Y are not independent, but

$$\begin{aligned}\text{cov}(X, Y) &= \text{cov}(X, X^2) \\ &= \text{E}[X \cdot X^2] - \text{E}[X] \cdot \text{E}[X^2] \\ &= \text{E}[X^3] - \text{E}[X] \text{E}[X^2] \\ &= 0 - 0 \cdot \text{E}[X^2] \\ &= 0.\end{aligned}$$

In this case, the relations               <img alt="handwritten red circle with a cross through it" data-bbox="9580 940

Hoeffding's covariance identity [edit]

A useful identity to compute the covariance between two random variables X, Y is the Hoeffding's covariance identity:^[7]

$$\text{cov}(X, Y) = \int_{\mathbb{R}} \int_{\mathbb{R}} (F_{(X,Y)}(x, y) - F_X(x)F_Y(y)) dx dy$$

where $F_{(X,Y)}(x, y)$ is the joint cumulative distribution function of the random vector (X, Y) and $F_X(x), F_Y(y)$ are the [marginals](#).

Uncorrelatedness and independence [edit]

Main article: Correlation and dependence

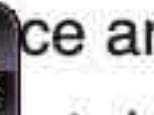
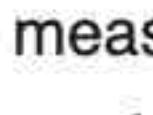
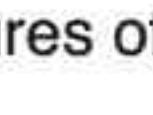
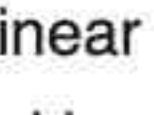
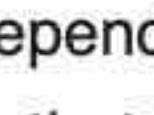
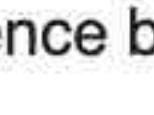
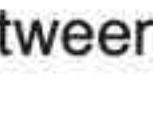
Random variables whose covariance is zero are called [uncorrelated](#).^{[4]:p. 121} Similarly, the components of random vectors whose [covariance matrix](#) is zero in every entry outside the main diagonal are also called uncorrelated.

If X and Y are [independent random variables](#), then their covariance is zero.^{[4]:p. 123[8]} This follows because under independence,

$$\mathbb{E}[XY] = \mathbb{E}[X] \cdot \mathbb{E}[Y].$$

The converse, however, is not generally true. For example, let X be uniformly distributed in $[-1, 1]$ and let $Y = X^2$. Clearly, X and Y are not independent, but

$$\begin{aligned} \text{cov}(X, Y) &= \text{cov}(X, X^2) \\ &= \mathbb{E}[X \cdot X^2] - \mathbb{E}[X] \cdot \mathbb{E}[X^2] \\ &= \mathbb{E}[X^3] - \mathbb{E}[X] \mathbb{E}[X^2] \\ &= 0 - 0 \cdot \mathbb{E}[X^2] \\ &= 0. \end{aligned}$$

In this case, the relations                 are measures of linear dependence between two random variables. This example shows that if two random variables are uncorrelated, that does not in general imply that they are independent.

$$\text{cov}(X, Y) = \int_{\mathbb{R}} \int_{\mathbb{R}} (F_{(X,Y)}(x, y) - F_X(x)F_Y(y)) dx dy$$

where $F_{(X,Y)}(x, y)$ is the joint cumulative distribution function of the random vector (X, Y) and $F_X(x), F_Y(y)$ are the [marginals](#).

Uncorrelatedness and independence [edit]

Main article: Correlation and dependence

Random variables whose covariance is zero are called [uncorrelated](#).^{[4]:p. 121} Similarly, the components of random vectors whose [covariance matrix](#) is zero in every entry outside the main diagonal are also called uncorrelated.

If X and Y are [independent random variables](#), then their covariance is zero.^{[4]:p. 123[8]} This follows because under independence,
 $E[XY] = E[X] \cdot E[Y]$.

The converse, however, is not generally true. For example, let X be uniformly distributed in $[-1, 1]$ and let $Y = X^2$. Clearly, X and Y are not independent, but

$$\begin{aligned}\text{cov}(X, Y) &= \text{cov}(X, X^2) \\ &= E[X \cdot X^2] - E[X] \cdot E[X^2] \\ &= E[X^3] - E[X] E[X^2] \\ &= 0 - 0 \cdot E[X^2] \\ &= 0.\end{aligned}$$

In this case, the relationship between Y and X is non-linear, while correlation and covariance are measures of linear dependence between two random variables. This example shows that if two random variables are uncorrelated, that does not in general imply that they are independent. However, if two variables are [jointly normally distributed](#) (but not if they are merely [individually normally distributed](#)), uncorrelatedness *does* imply independence.

Cov &
Indep

J_R J_R

where $F_{(X,Y)}(x,y)$ is the joint cumulative distribution function of the random vector (X, Y) and $F_X(x), F_Y(y)$ are the [marginals](#).

Uncorrelatedness and independence [edit]

Main article: [Correlation and dependence](#)

Random variables whose covariance is zero are called [uncorrelated](#).^{[4]:p. 121} Similarly, the components of random vectors whose [covariance matrix](#) is zero in every entry outside the main diagonal are also called uncorrelated.

If X and Y are [independent random variables](#), then their covariance is zero.^{[4]:p. 123[8]} This follows because under independence,
 $E[XY] = E[X] \cdot E[Y]$.

The converse, however, is not generally true. For example, let X be uniformly distributed in $[-1, 1]$ and let $Y = X^2$. Clearly, X and Y are not independent, but

$$\begin{aligned}\text{cov}(X, Y) &= \text{cov}(X, X^2) \\ &= E[X \cdot X^2] - E[X] \cdot E[X^2] \\ &= E[X^3] - E[X] E[X^2] \\ &= 0 - 0 \cdot E[X^2] \\ &= 0.\end{aligned}$$

In this case, the relationship between Y and X is non-linear, while correlation and covariance are measures of linear dependence between two random variables. This example shows that if two random variables are uncorrelated, that does not in general imply that they are independent. However, if two variables are [jointly normally distributed](#) (but not if they are merely [individually normally distributed](#)), uncorrelatedness *does* imply independence.

Relationship to inner products



10:47



wiki

if $\underline{X \& Y}$ are indep then $\underline{\text{cov}(X,Y) = 0}$ ✓



if $\underline{\text{cov}(X,Y) = 0}$ then $\underline{X \& Y}$

need not be
independent



wiki

$$\int_{\mathbb{R}} \int_{\mathbb{R}} F_{(X,Y)}(x,y) (F_X(x) - F_X(\bar{x})) (F_Y(y) - F_Y(\bar{y})) dxdy$$

where $F_{(X,Y)}(x,y)$ is the joint cumulative distribution function of the random vector (X, Y) and $F_X(x), F_Y(y)$ are the [marginals](#).

Uncorrelatedness and independence [edit]

Main article: Correlation and dependence

Random variables whose covariance is zero are called [uncorrelated](#).^{[4]:p. 121} Similarly, the components of random vectors whose [covariance matrix](#) is zero in every entry outside the main diagonal are also called uncorrelated.

If X and Y are [independent random variables](#), then their covariance is zero.^{[4]:p. 123[8]} This follows because under independence,

$$\text{E}[XY] = \text{E}[X] \cdot \text{E}[Y].$$

The converse, however, is not generally true. For example, let X be uniformly distributed in $[-1, 1]$ and let $Y = X^2$. Clearly, X and Y are not independent, but

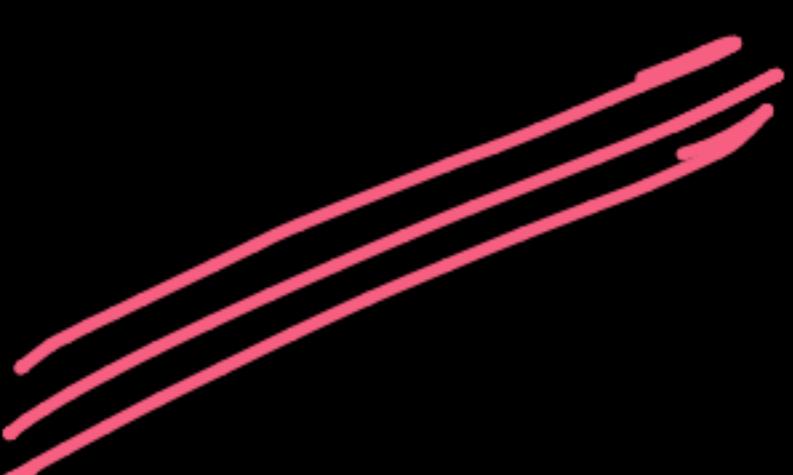
$$\begin{aligned} \text{cov}(X, Y) &= \text{cov}(X, X^2) \\ &= \text{E}[X \cdot X^2] - \text{E}[X] \cdot \text{E}[X^2] \\ &= \text{E}[X^3] - \text{E}[X] \text{E}[X^2] \\ &= 0 - 0 \cdot \text{E}[X^2] \\ &= 0. \end{aligned}$$

$$Y = X^2$$

In this case, the relationship between Y and X is non-linear, while correlation and covariance are measures of linear dependence between two random variables. This example shows that if two random variables are uncorrelated, that does not in general imply that they are independent. However, if two variables are [jointly normally distributed](#) (but not if they are merely [individually normally distributed](#)), uncorrelatedness *does* imply independence.

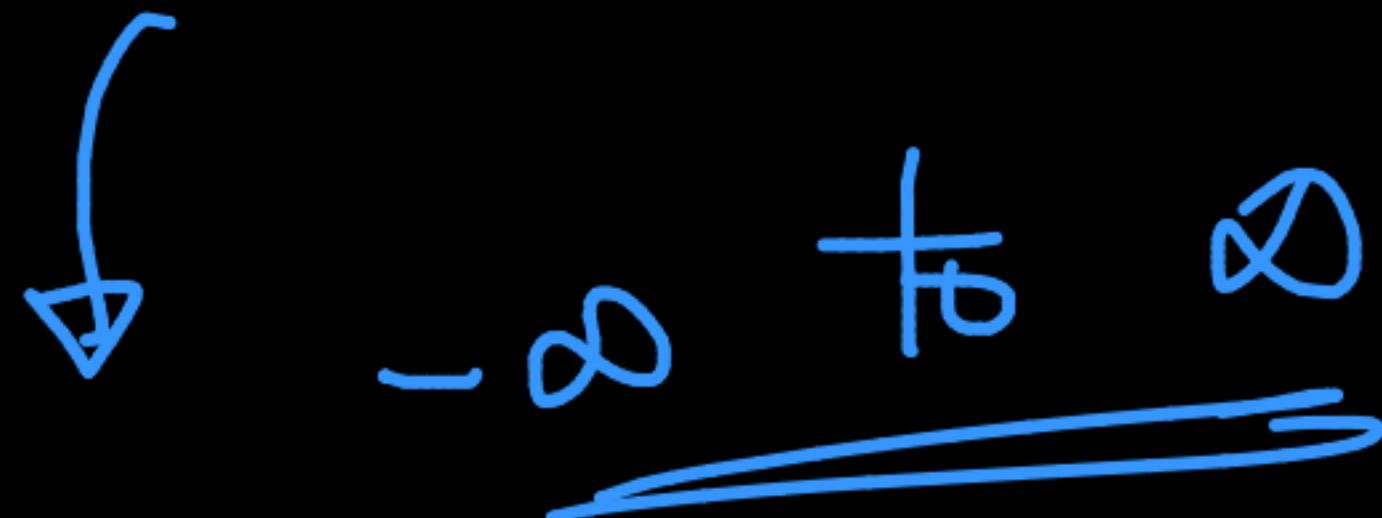
Relationship to inner product





$x \uparrow \quad y \uparrow \downarrow \Rightarrow$

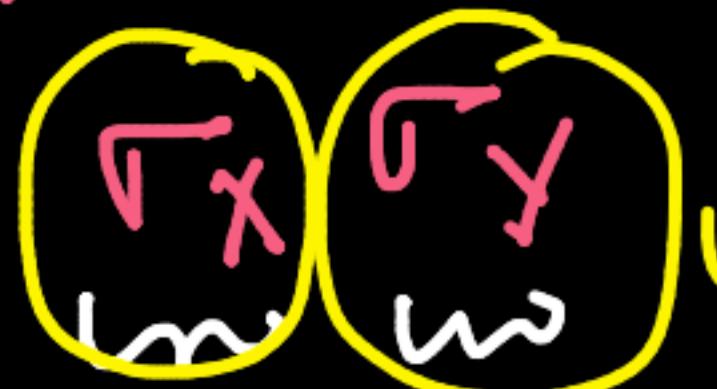
$$\left\{ \begin{array}{l} \text{Cov}(x,y) > 0 \\ \text{Cov}(x,y) < 0 \\ \text{Cov}(x,y) \approx 0 \end{array} \right.$$



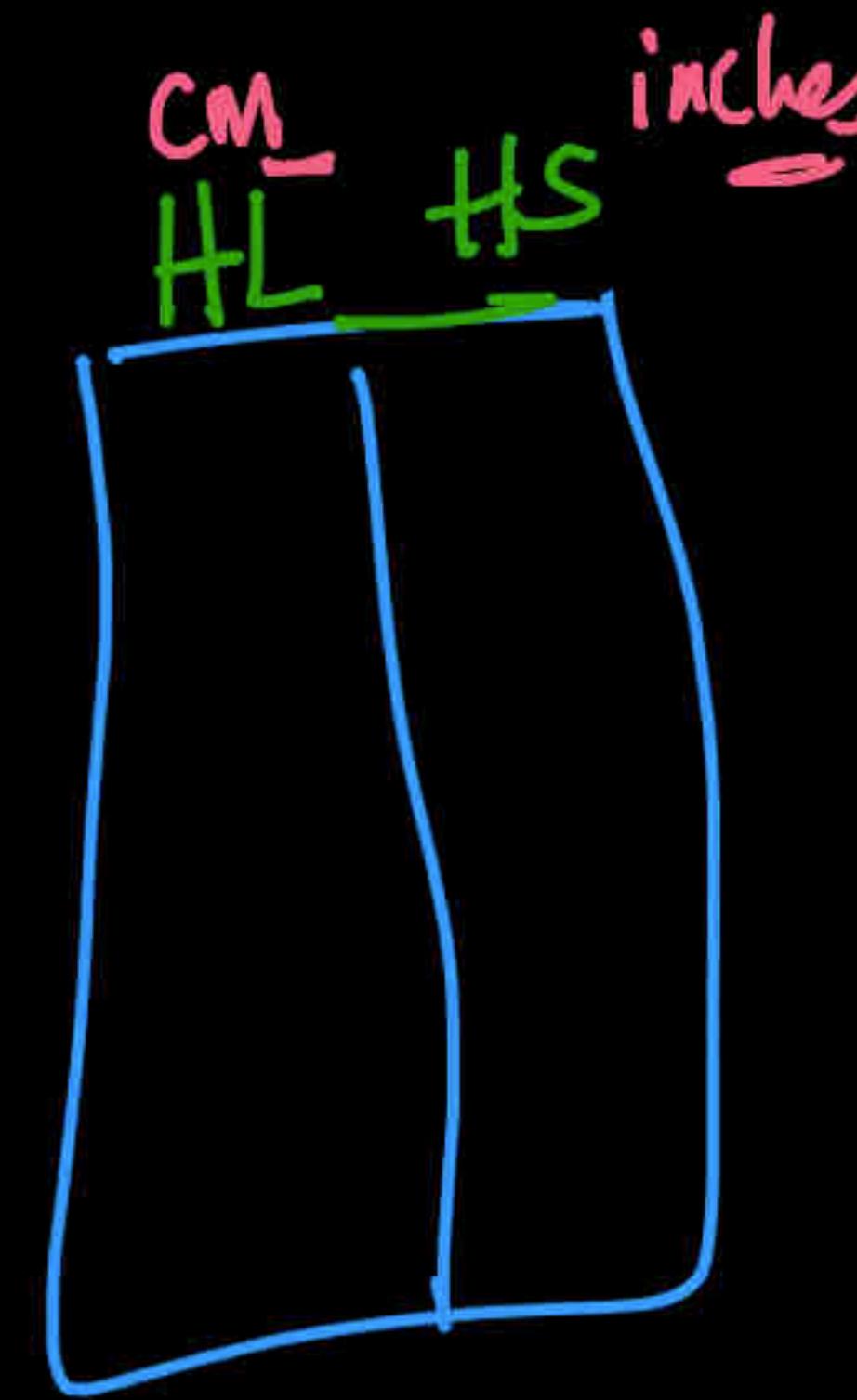
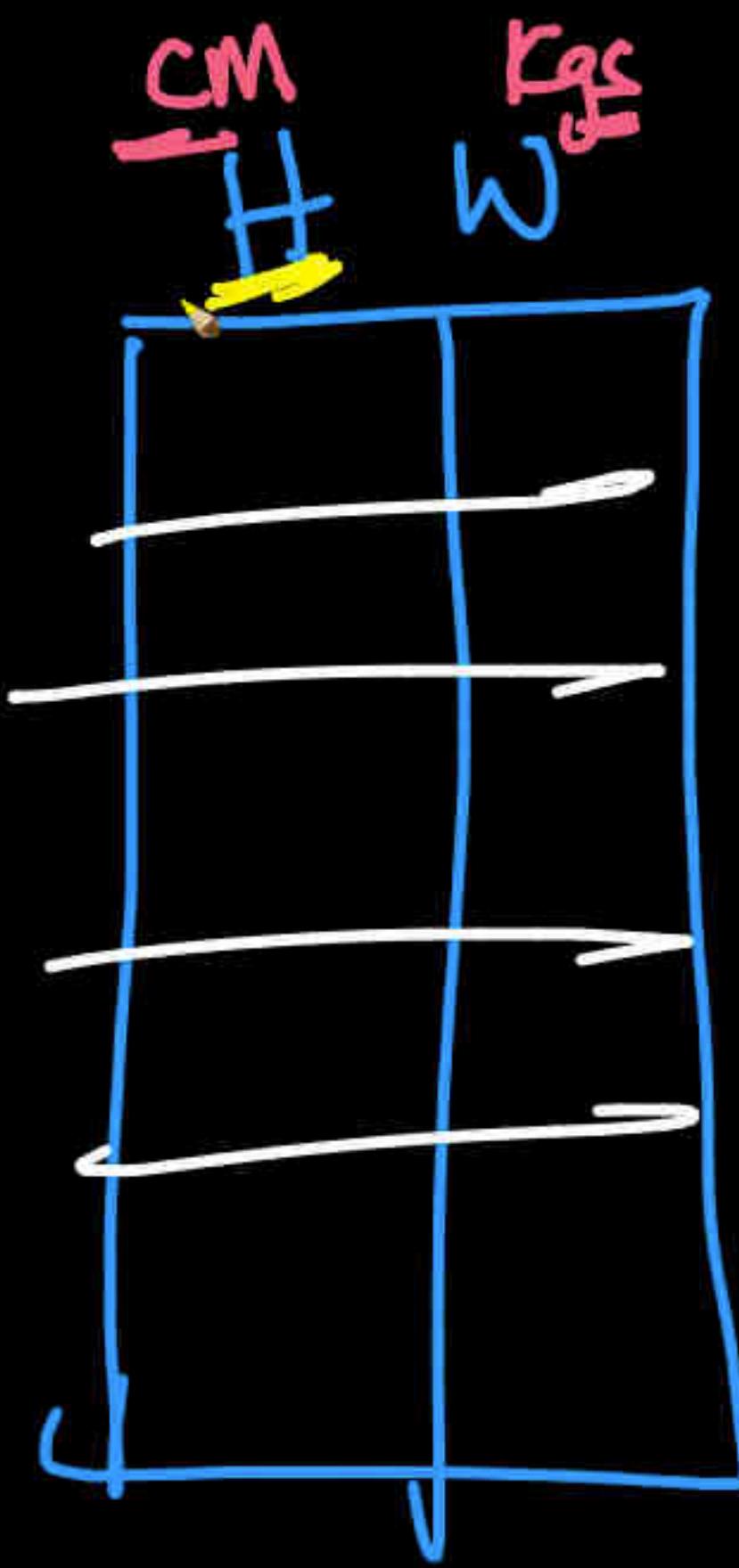
Correlation

Pearson-CC (x,y)

$$\rho_{x,y} = \frac{\text{Cov}(x,y)}{\sqrt{\text{Var}(x)} \sqrt{\text{Var}(y)}} \quad [-1 \quad +1]$$



independent of the units &
means



$$\text{Cov}(H, HS) > 0$$

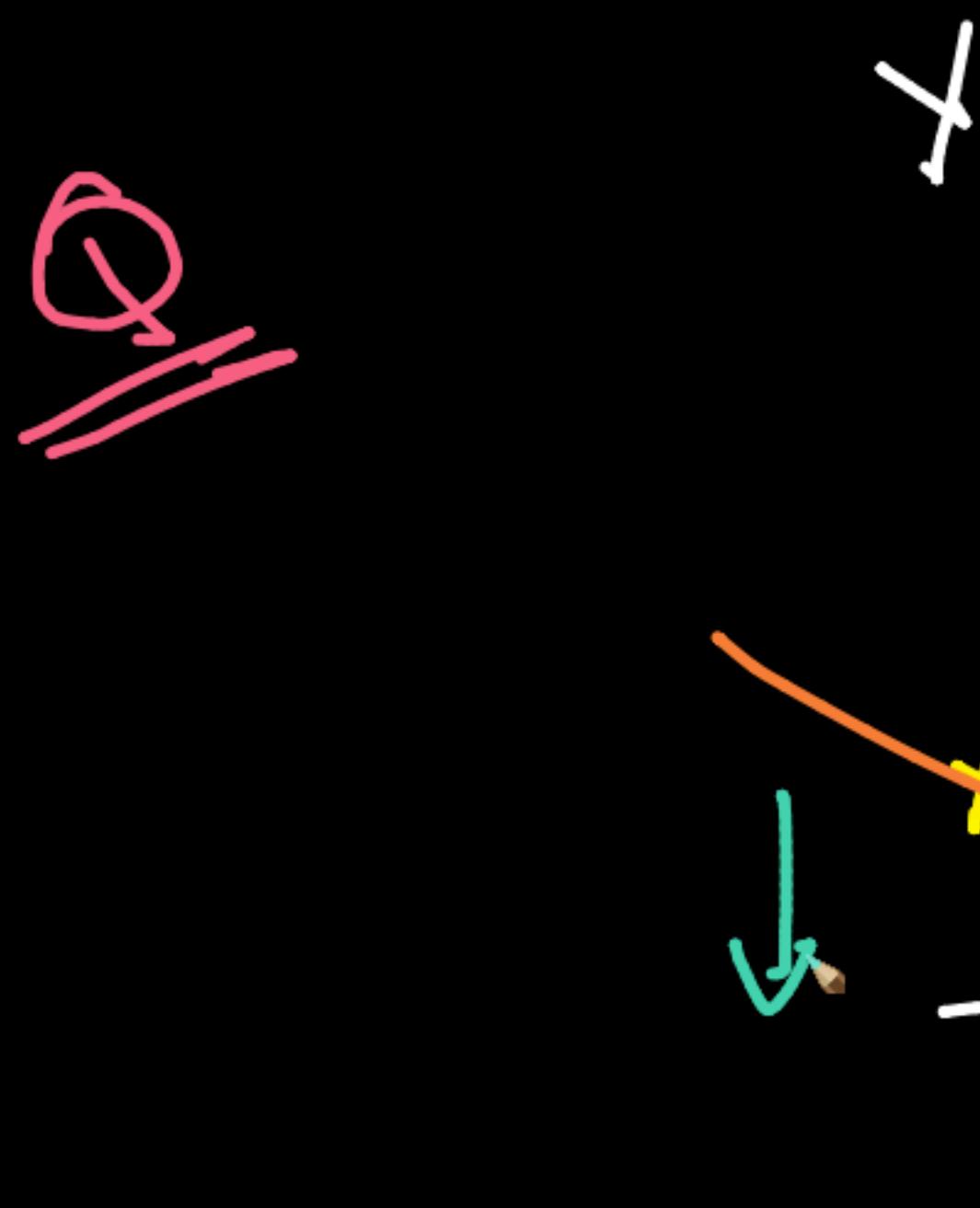
$$\text{Cov}(H, W) > 0$$

$$-1 \leq \rho_{X,Y} = \frac{\text{Cov}(X,Y)}{\sqrt{\sigma_X^2 \sigma_Y^2}} < 1$$

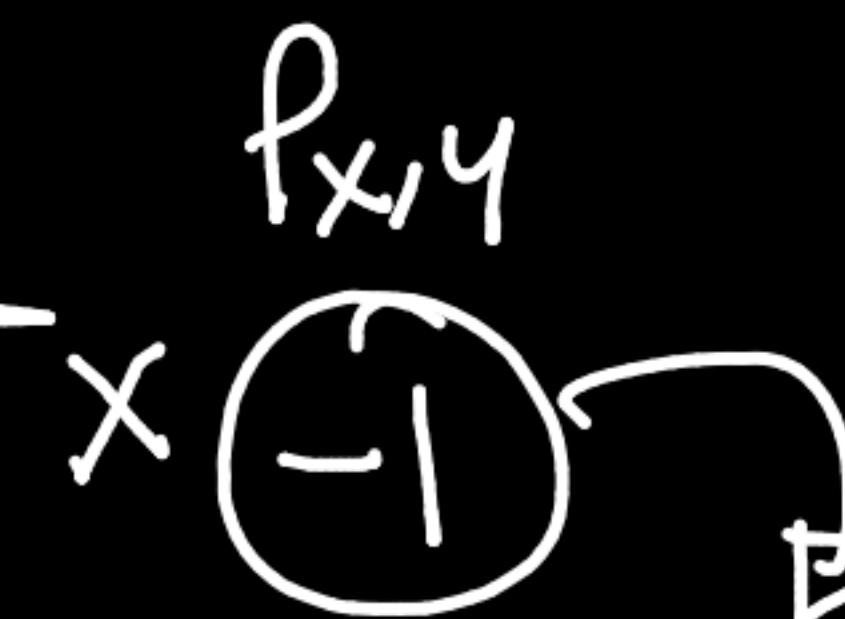


need not be

$$\frac{\text{Cov}(H_{cm}, \omega_{kg})}{\sqrt{H_{cm}} \sqrt{\omega_{kg}}} \not\equiv \frac{\text{Cov}(H_{ft}, \omega_{ls})}{\sqrt{H_{ft}} \sqrt{\omega_{ls}}}$$

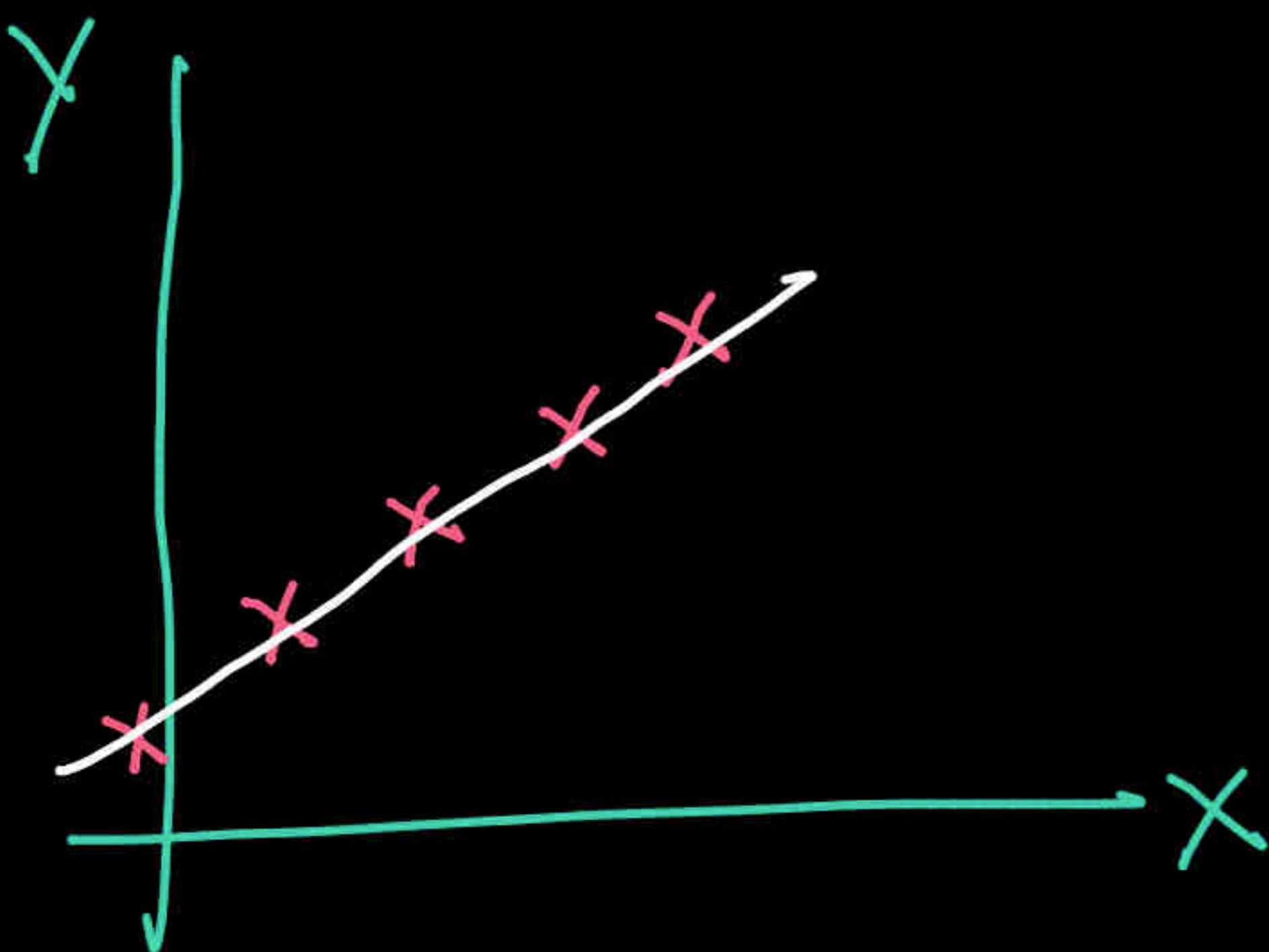


$$P_{CC} = \rho_{XY} = \frac{\text{Cov}(X,Y)}{\sigma_X \sigma_Y}$$



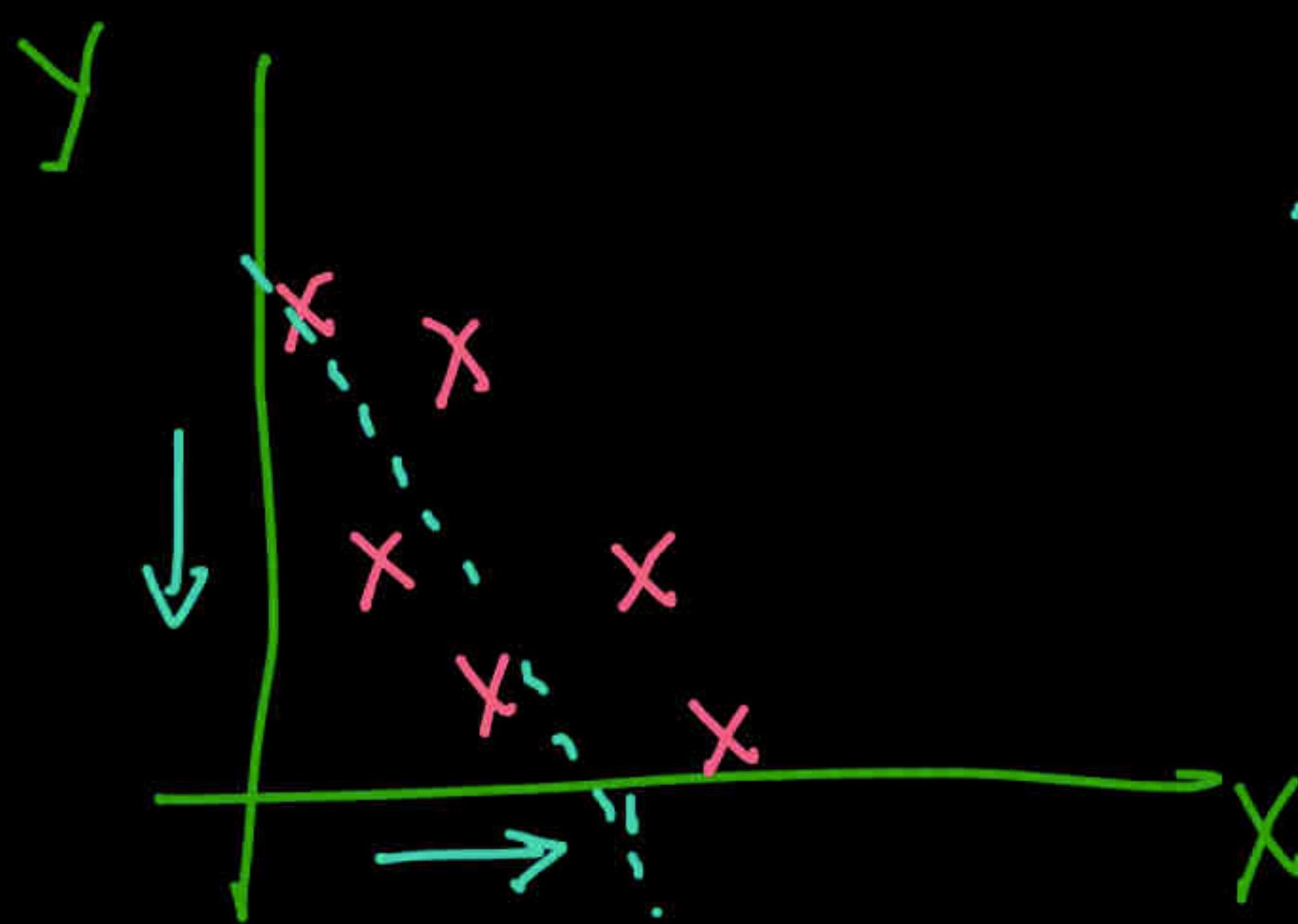
\times & Y are perfectly
in-linearly related
 \downarrow
there exists a line

$\rho_{XY} =$ \checkmark
~~close to -1~~
~~-ve not -1~~

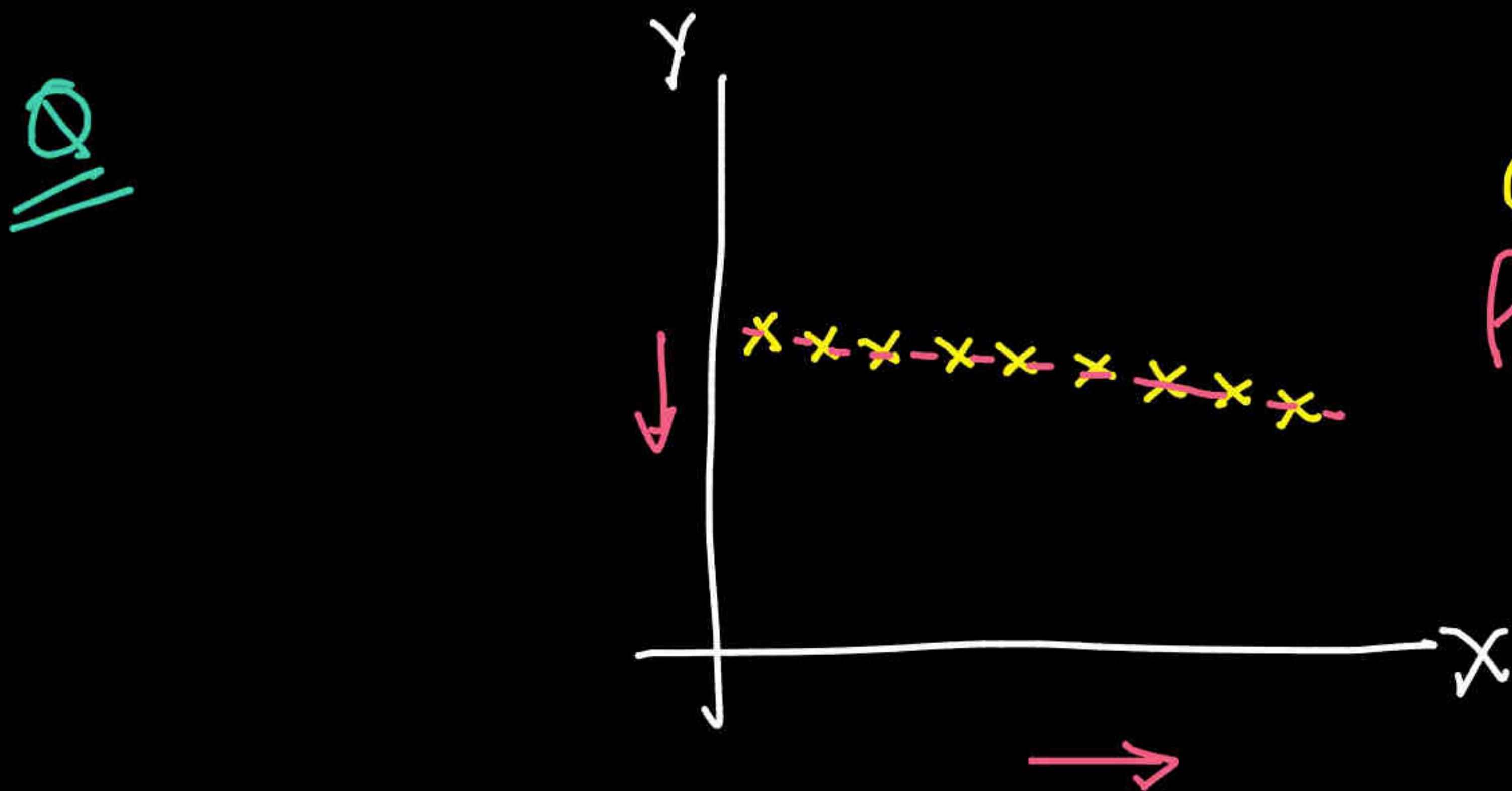


$$\rho_{x,y} = +1$$

X & Y are perfectly linearly related

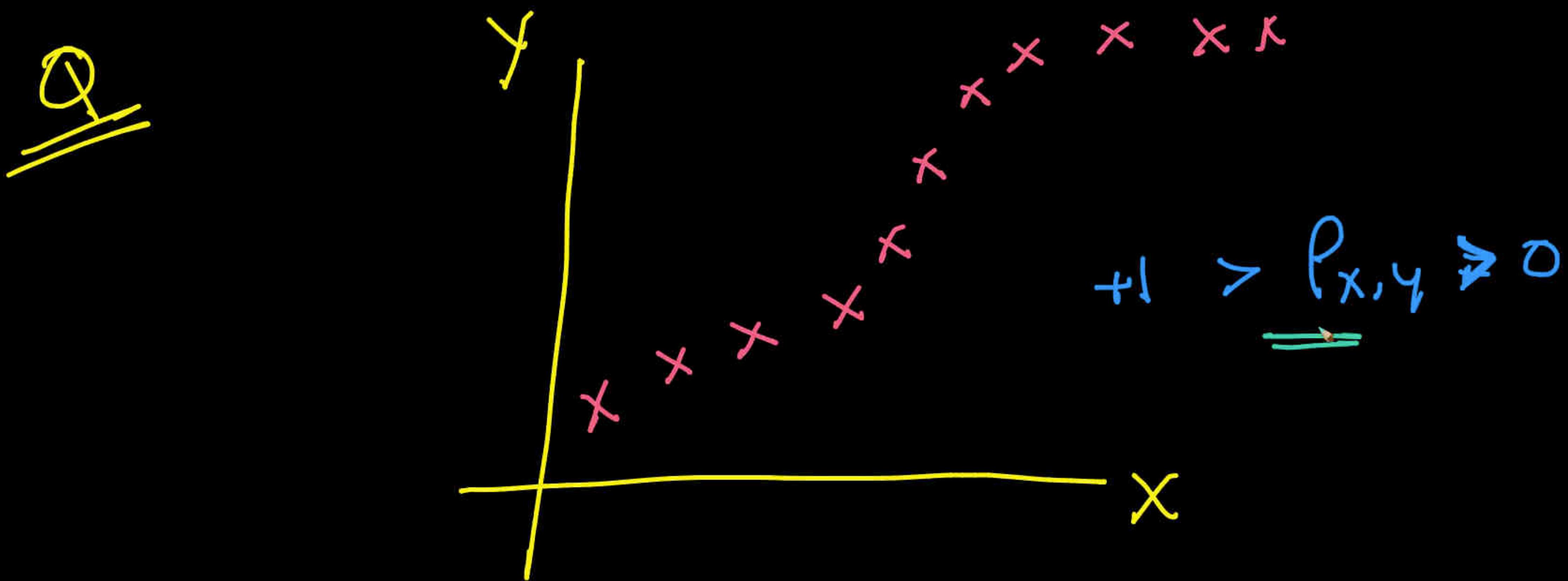


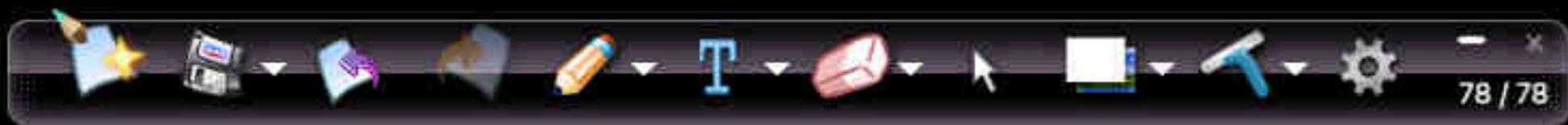
$$-1 < \rho_{x,y} < 0$$

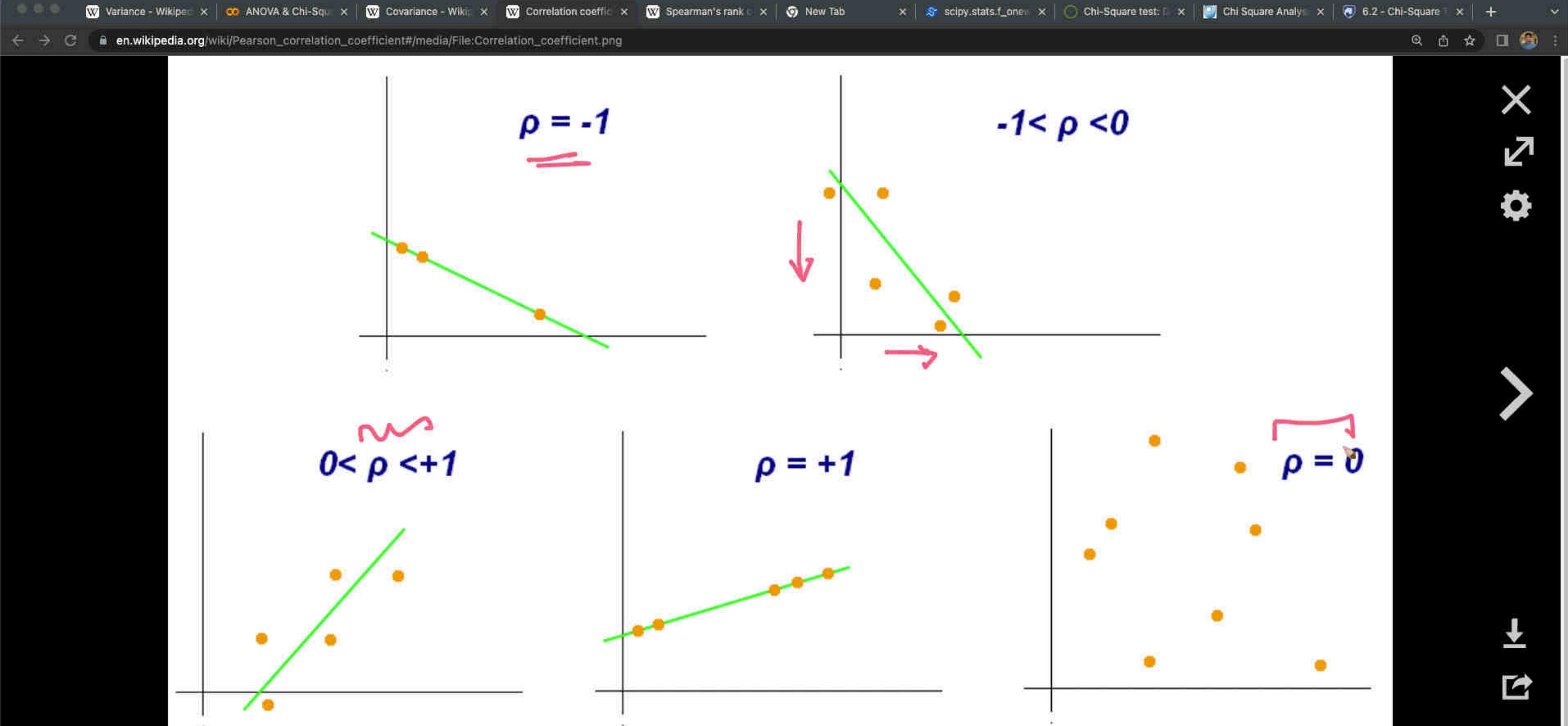


not measure the
rate of change

$$\rho_{xy} > -1$$

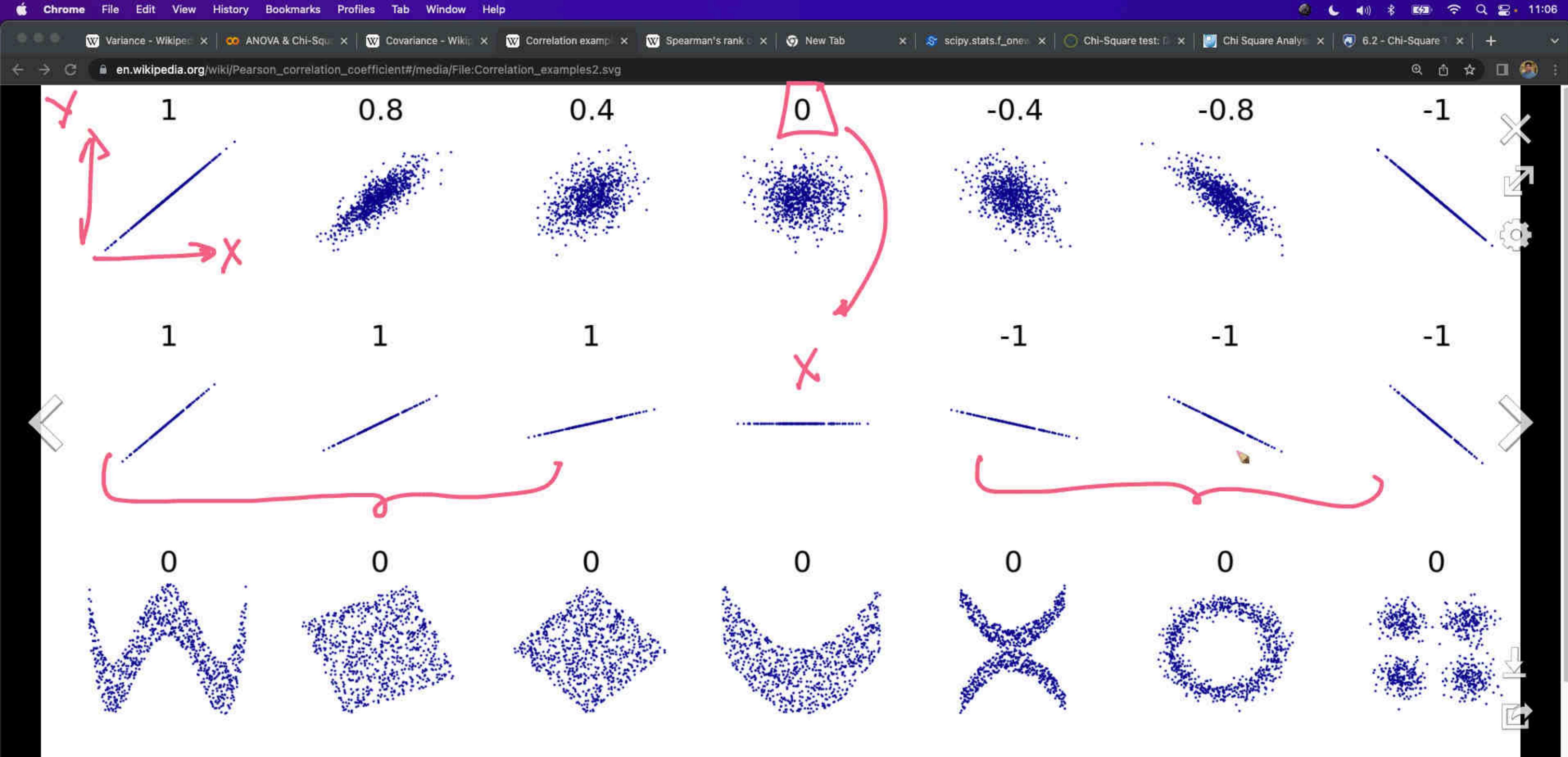






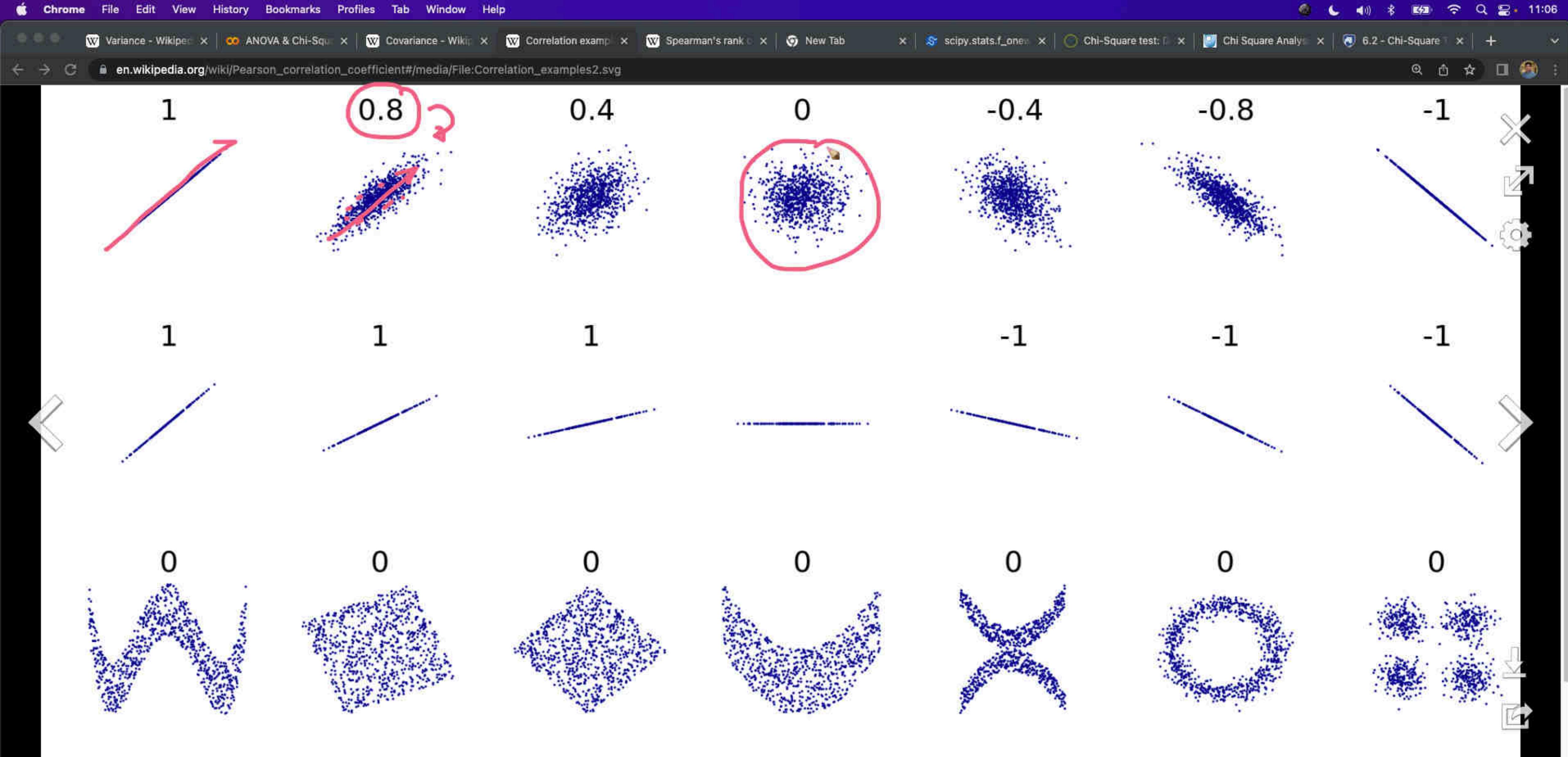
Examples of scatter diagrams with different values of correlation coefficient (ρ)

More details



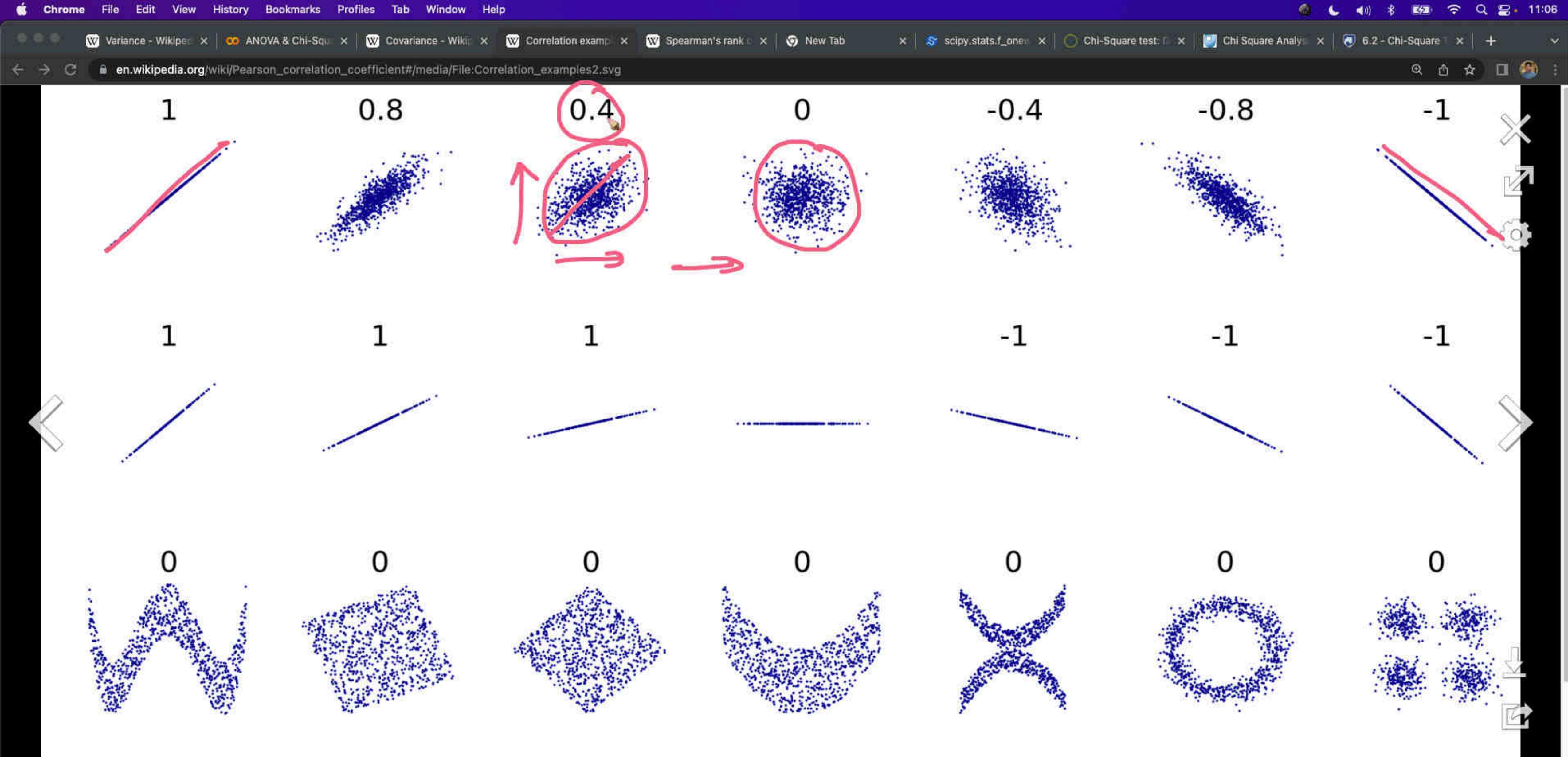
Several sets of (x, y) points, with the correlation coefficient of x and y for each set. Note that the correlation reflects the strength and direction of a linear relationship (top row), but not the slope of that relationship (middle), nor many ...

[More details](#)



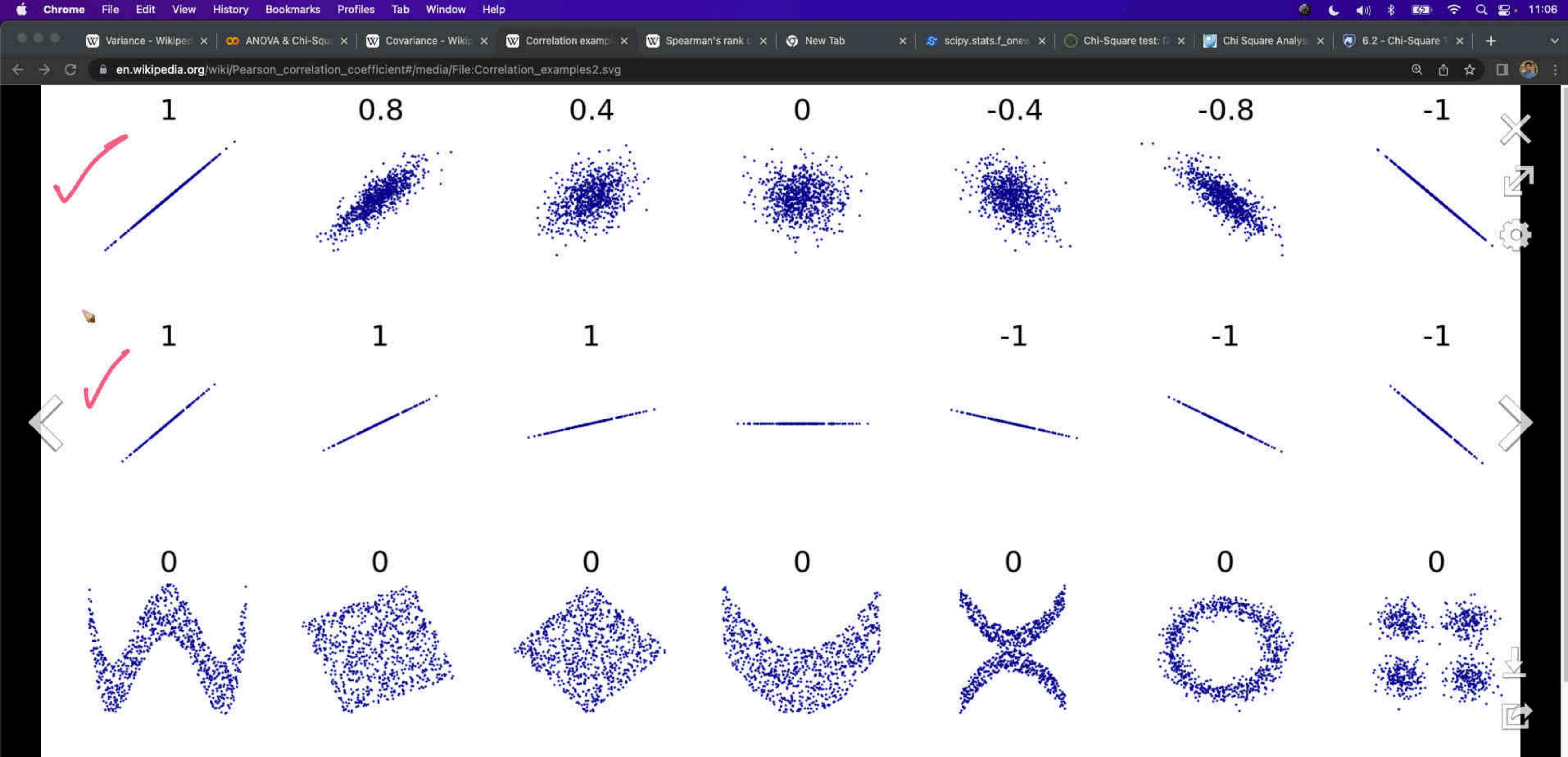
Several sets of (x, y) points, with the correlation coefficient of x and y for each set. Note that the correlation reflects the strength and direction of a linear relationship (top row), but not the slope of that relationship (middle), nor many ...

[More details](#)



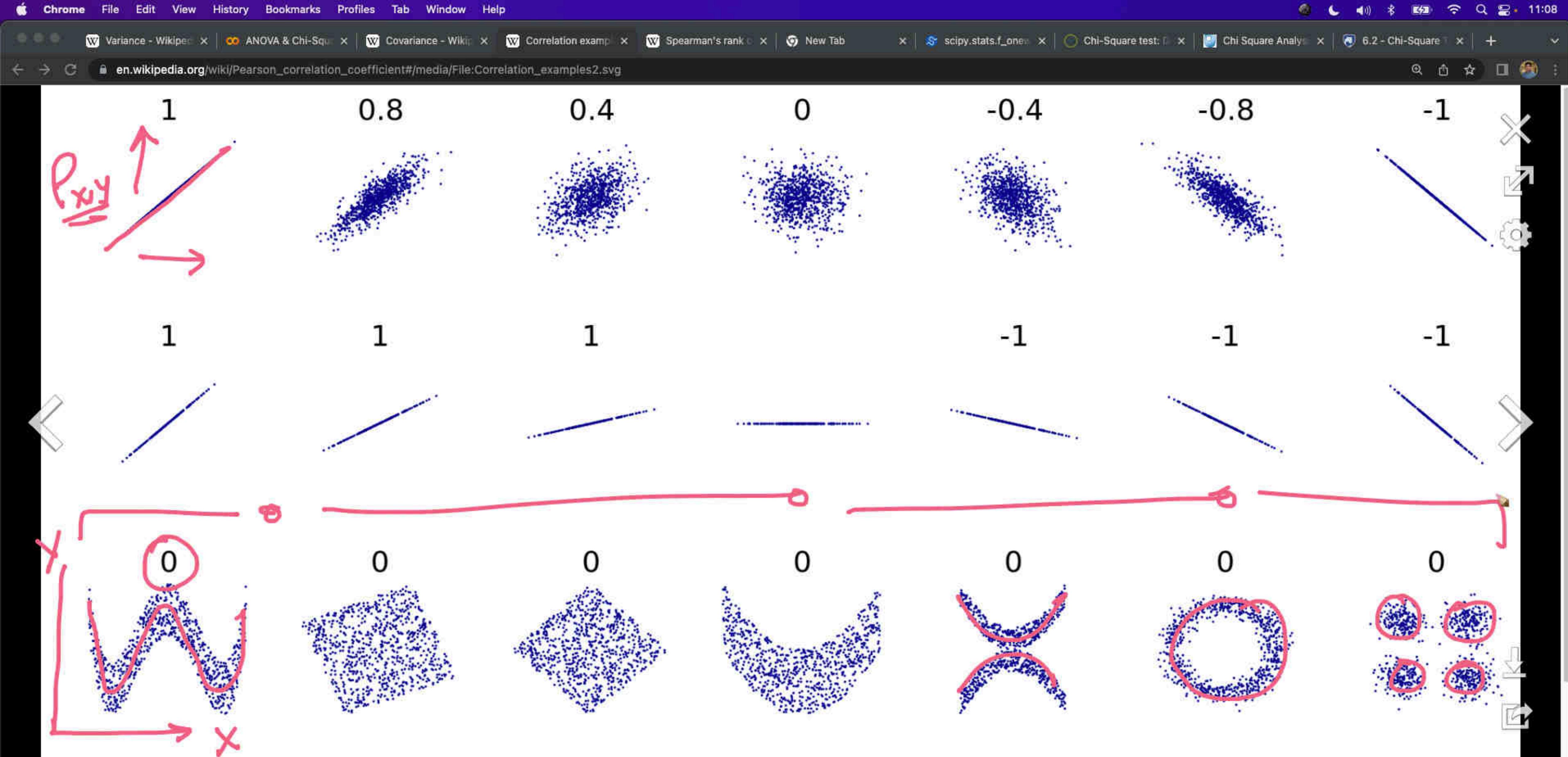
Several sets of (x, y) points, with the correlation coefficient of x and y for each set. Note that the correlation reflects the strength and direction of a linear relationship (top row), but not the slope of that relationship (middle), nor many ...

More details



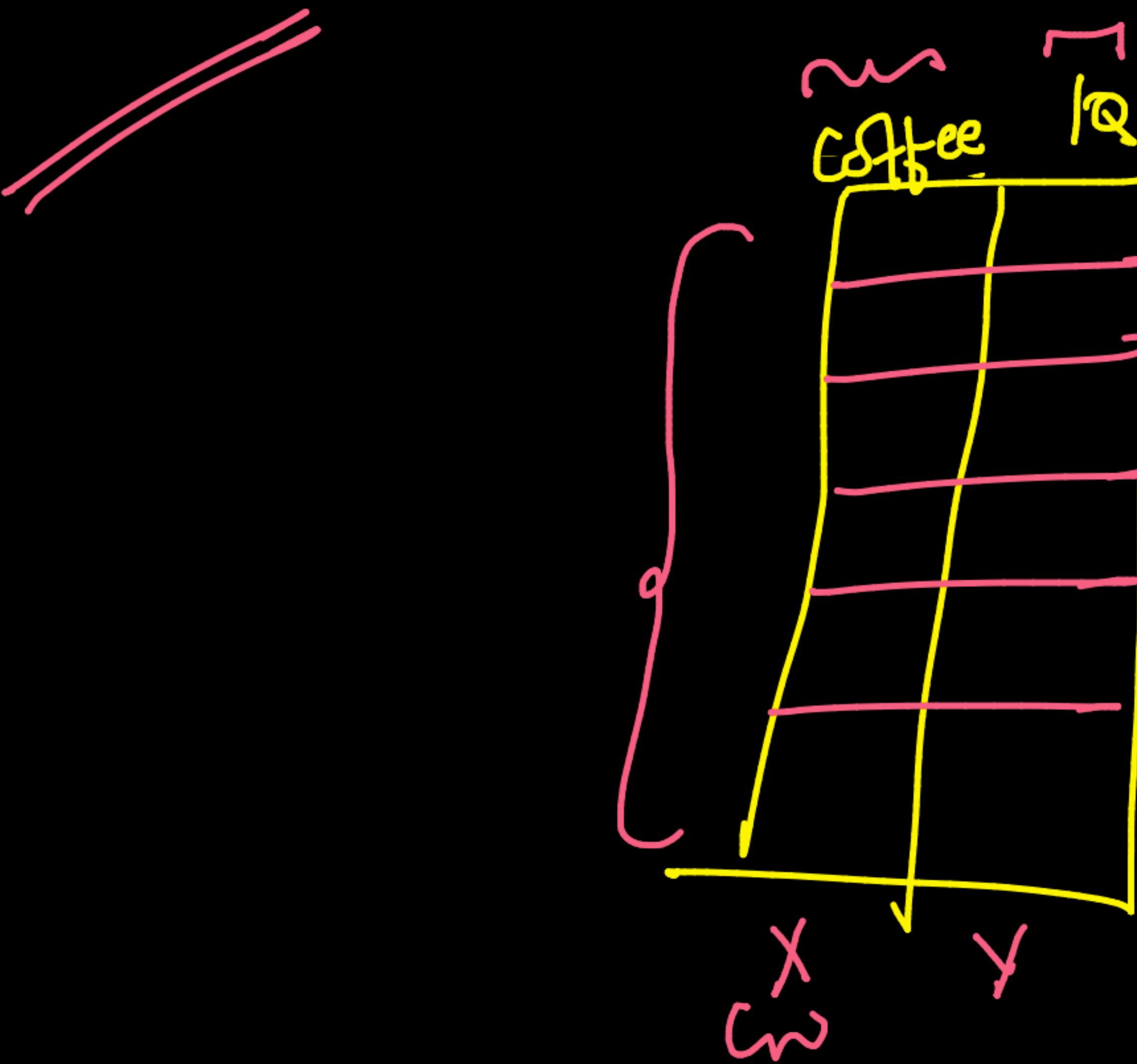
Several sets of (x, y) points, with the correlation coefficient of x and y for each set. Note that the correlation reflects the strength and direction of a linear relationship (top row), but not the slope of that relationship (middle), nor many ...

[More details](#)



Several sets of (x, y) points, with the correlation coefficient of x and y for each set. Note that the correlation reflects the strength and direction of a linear relationship (top row), but not the slope of that relationship (middle), nor many ...

More details

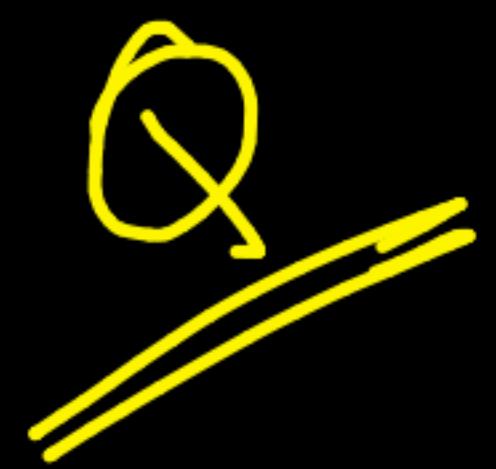


$P_{x,y} > 0$

① Correlation is +ve

② Coffee increases your IQ

Causation



measure Causalism



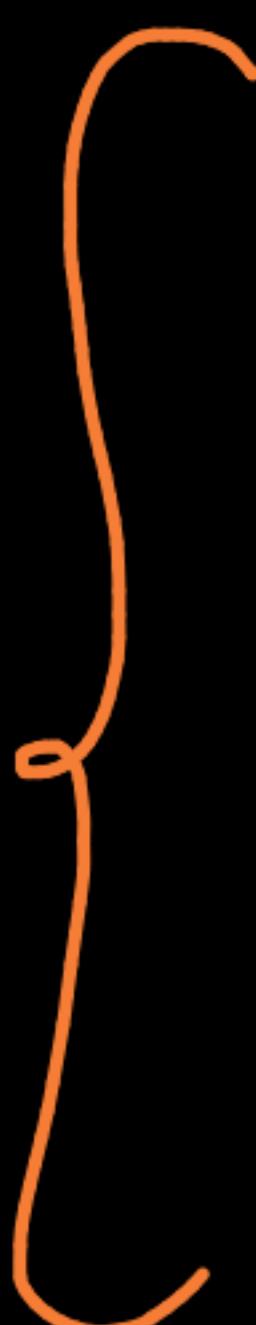
A | B

test

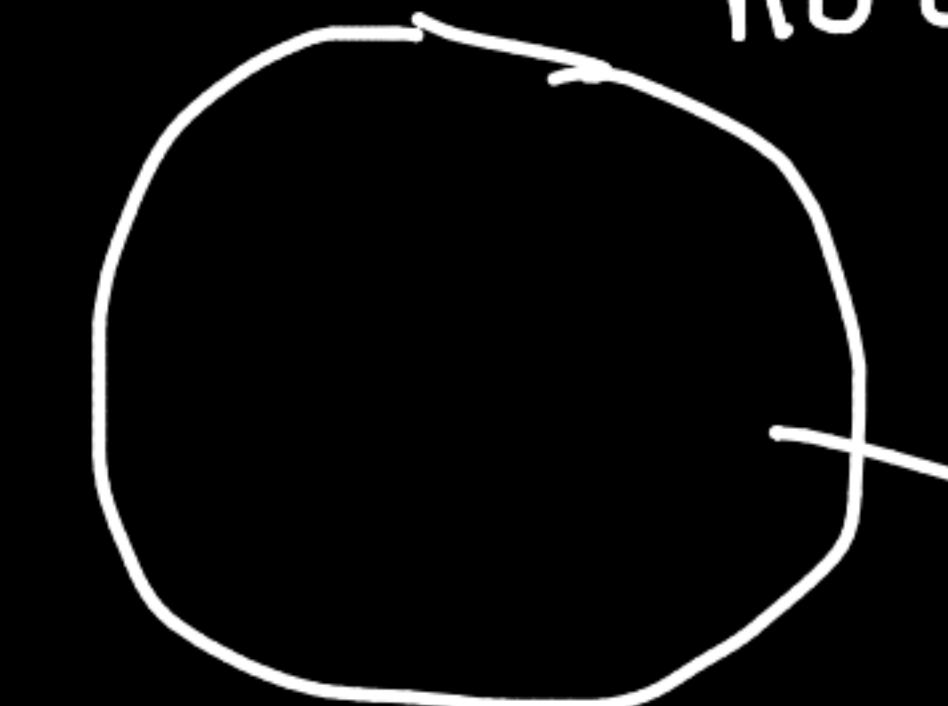
Simplest way

various
no coffee

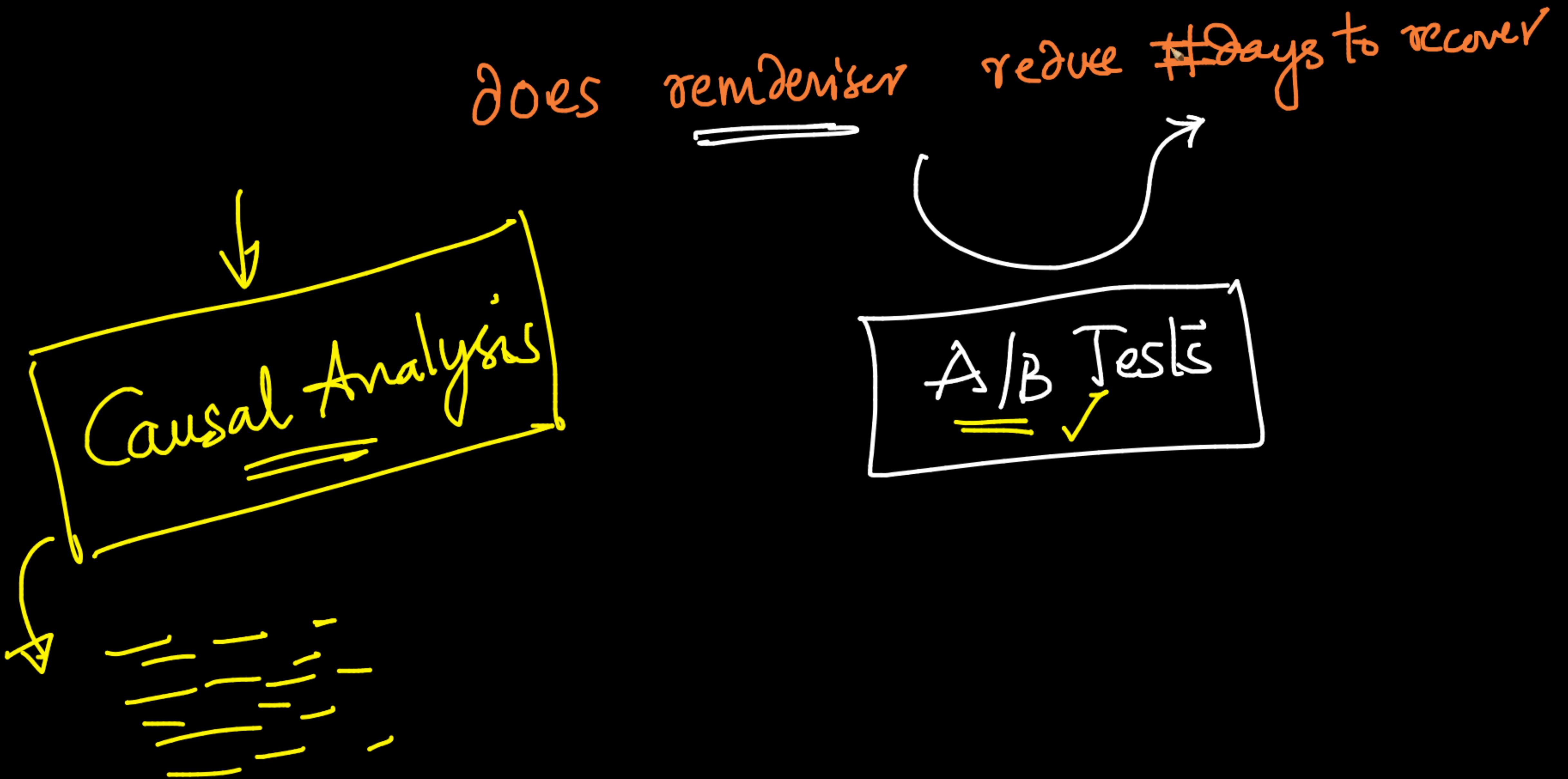
random
coffee



| Q

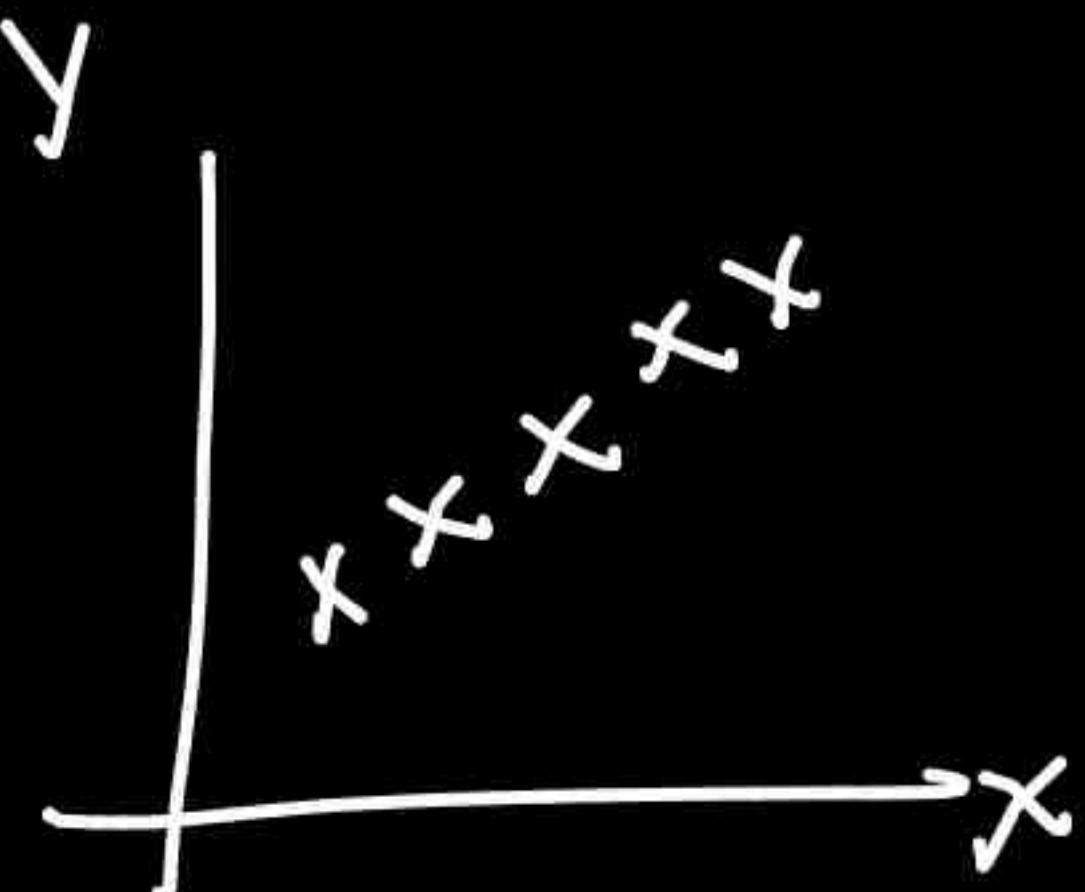


| Q

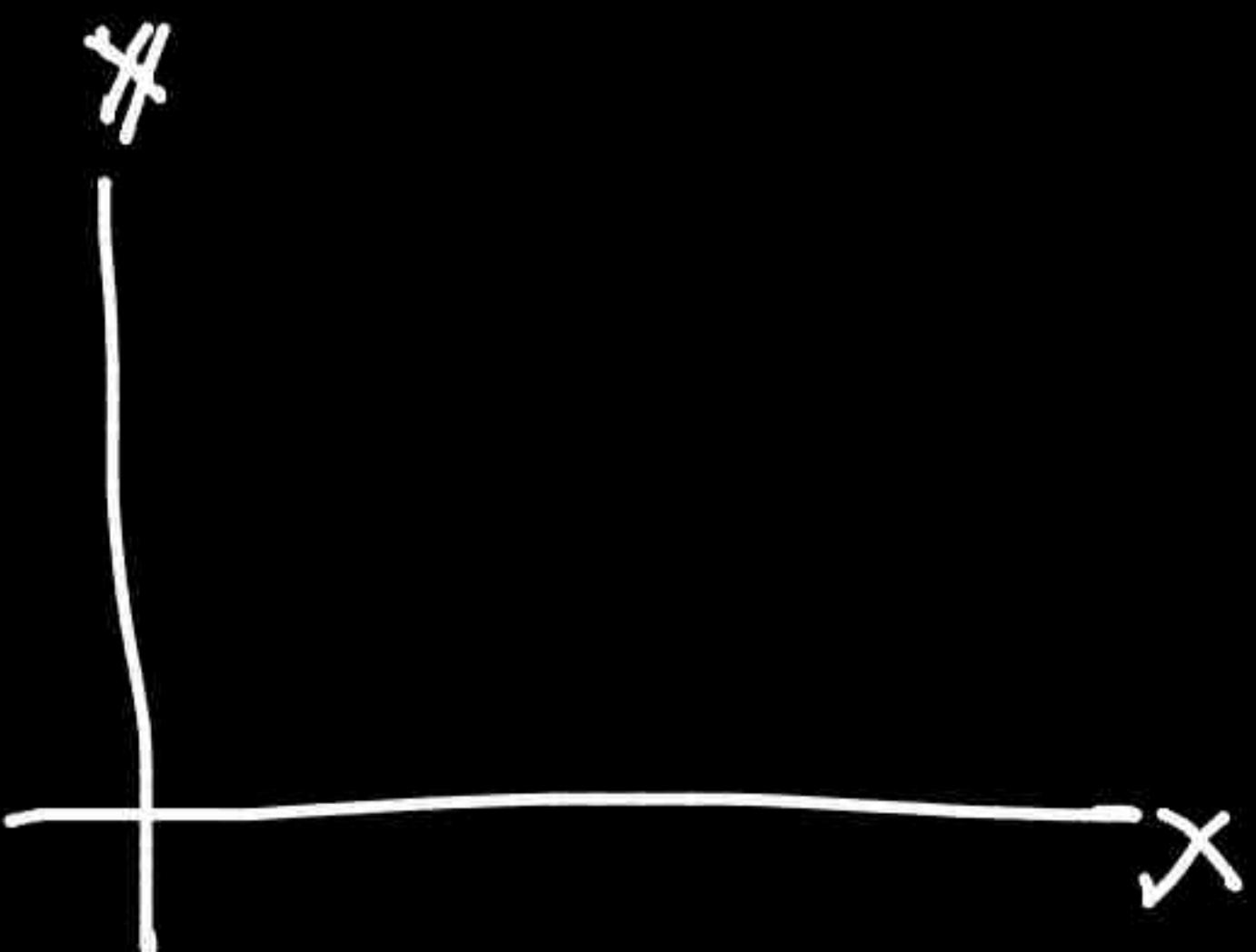


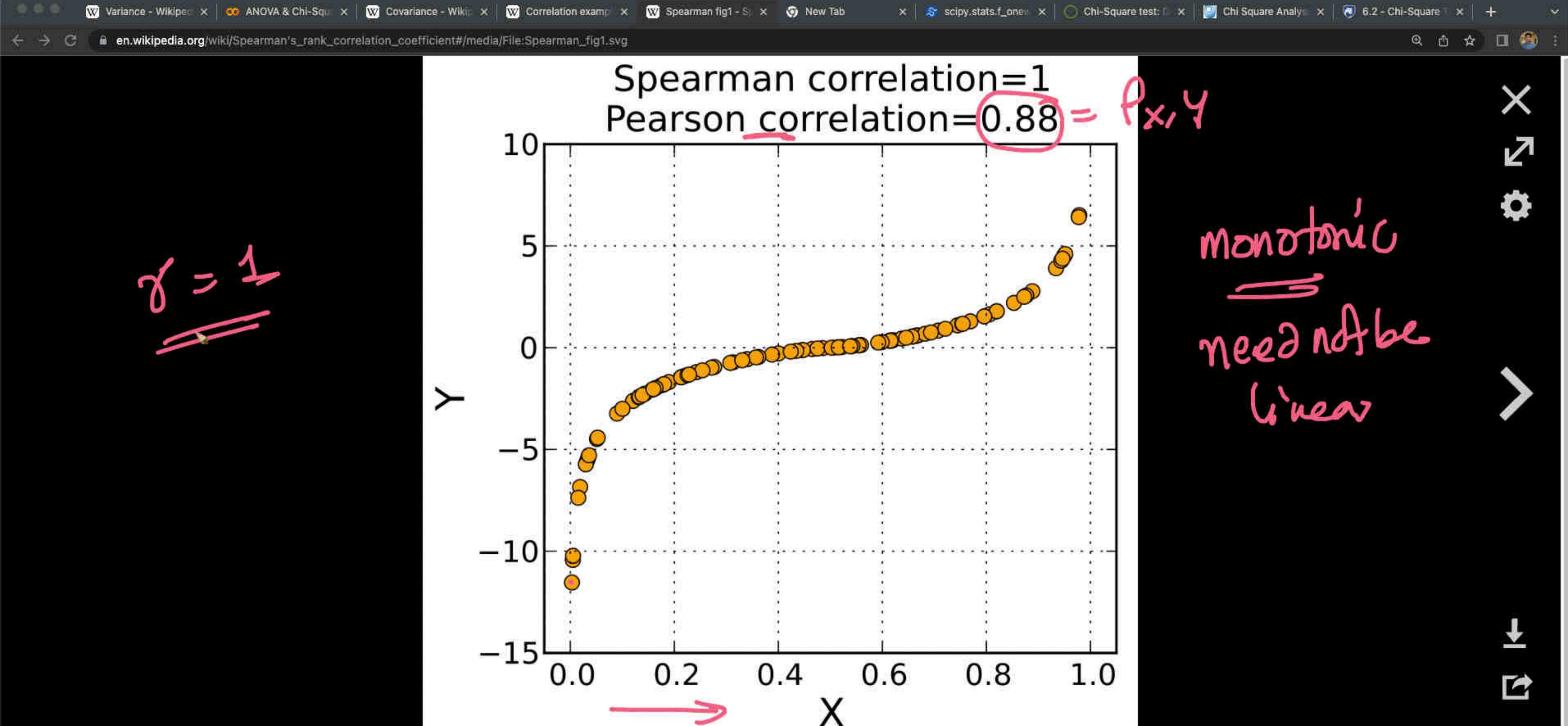


$$\rho_{x,y} = 1$$



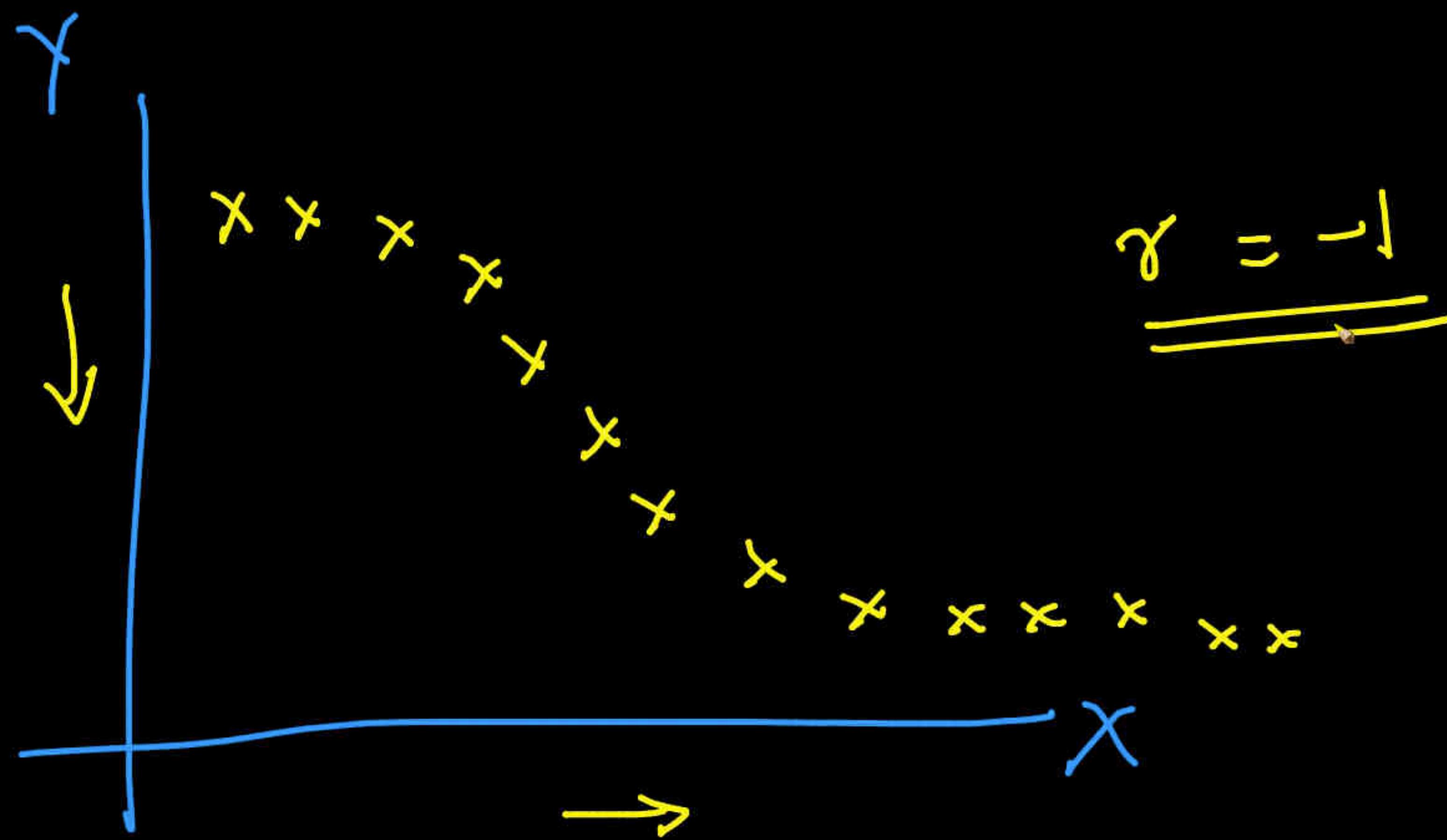
monotonic relationship
need not be linear





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 ha...

More details



Spearman
Rank CC.

spent in front of TV per week. [citation needed]

IQ, X_i	Hours of TV per week, Y_i
106	7
100	27
86	2
101	50
99	28
103	29
97	20
113	12
112	6
110	17

Firstly, evaluate d_i^2 . To do so use the following steps, reflected in the table below.

1. Sort the data by the first column (X_i). Create a new column x_i and assign it the ranked values 1, 2, 3, ..., n .
2. Next, sort the data by the second column (Y_i). Create a fourth column y_i and similarly assign it the ranked values 1, 2, 3, ..., n .
3. Create a fifth column d_i to hold the differences between the two rank columns (x_i and y_i).
4. Create one final column d_i^2 to hold the value of column d_i squared.

IQ, X_i	Hours of TV per week, Y_i	V	rank _x	rank _y	d	d^2
86	2					

IQ, X_i	Hours of TV per week, Y_i	rank x_i	rank y_i	d_i	d_i^2
86		2	1	1	0
97	20	2	6	-4	16
99	28	3	8	-5	25
100	27	4	7	-3	9
101	50	5	10	-5	25
103	29	6	9	-3	9
106	7	7	3	4	16
110	17	8	5	3	9
112	6	9	2	7	49
113	12	10	4	6	36

With d_i^2 found, add them to find $\sum d_i^2 = 194$. The value of n is 10. These values can now be substituted back into the equation

$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

to give

$$\rho = 1 - \frac{6 \times 194}{10(10^2 - 1)},$$

which evaluates to $\rho = -0.45$ (standard normal distribution).

With d_i^2 found, add them to find $\sum d_i^2 = 194$. The value of n is 10. These values can now be substituted back into the equation

$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

to give

$$\rho = 1 - \frac{6 \times 194}{10(10^2 - 1)},$$

which evaluates to $\rho = -29/165 = -0.175757575\dots$ with a [p-value = 0.627188](#) (using the [t-distribution](#)).

That the value is close to zero shows that the correlation between IQ and hours spent watching TV is very low, although the negative value suggests that the longer the time spent watching television the lower the IQ. In the case of ties in the original values, this formula should not be used; instead, the Pearson correlation coefficient should be calculated on the ranks (where ties are given ranks, as described above).

Determining significance [edit]

One approach to test whether an observed value of ρ is significantly different from zero (r will always maintain $-1 \leq r \leq 1$) is to calculate the probability that it would be greater than or equal to the observed r , given the [null hypothesis](#), by using a [permutation test](#). An advantage of this approach is that it automatically takes into account the number of tied data values in the sample and the way they are treated in computing the rank correlation.

Another approach parallels the use of the [Fisher transformation](#) in the case of the Pearson product-moment correlation coefficient. That is, [confidence intervals](#) and [hypothesis testing](#) can be carried out using the Fisher transformation:

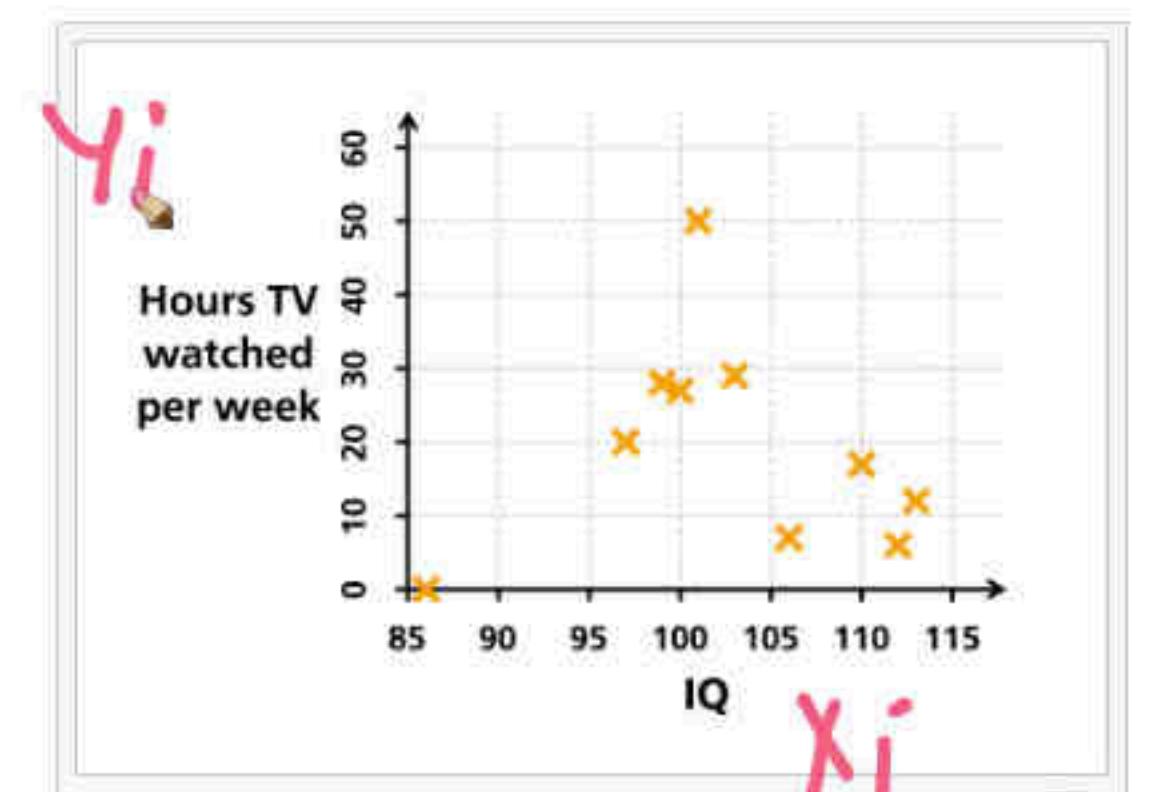


Chart of the data presented. It can be seen that there might be a negative correlation, but that the relationship does not appear definitive.

- en.wikipedia.org/wiki/Spearman%27s_rank_correlation_coefficient
1. Sort the data by the first column (X_i). Create a new column x_i and assign it the ranked values $1, 2, 3, \dots, n$.
 2. Next, sort the data by the second column (Y_i). Create a fourth column y_i and similarly assign it the ranked values $1, 2, 3, \dots, n$.
 3. Create a fifth column d_i to hold the differences between the two rank columns (x_i and y_i).
 4. Create one final column d_i^2 to hold the value of column d_i squared.

l_{X/Y}

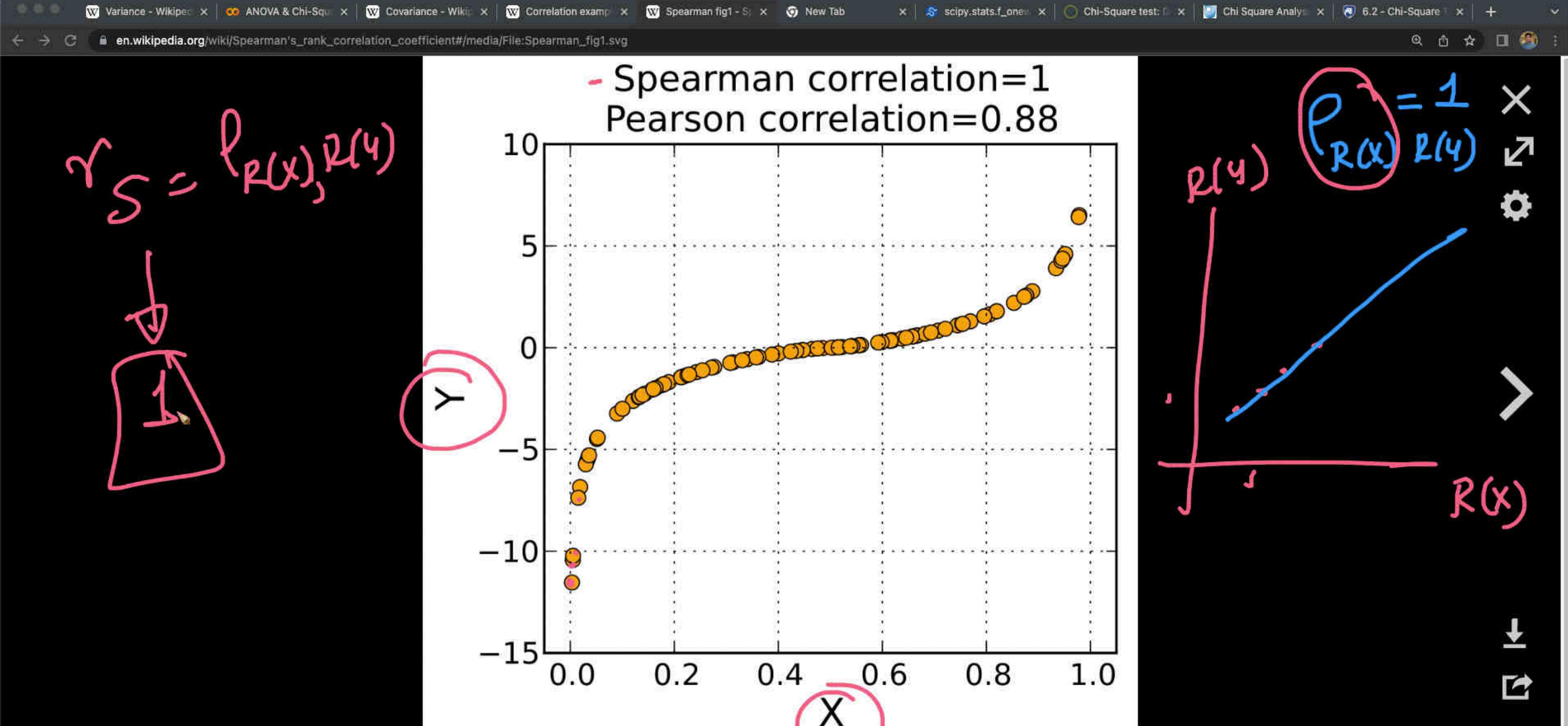
IQ, X_i	Hours of TV per week, Y_i	rank x_i	rank y_i	d_i	d_i^2
86	2	1	1	0	0
97	20	2	6	-4	16
99	28	3	8	-5	25
100	27	4	7	-3	9
101	50	5	10	-5	25
103	29	6	9	-3	9
106	7	7	3	4	16
110	17	8	5	3	9
112	6	9	2	7	49
113	12	10	4	6	36

With d_i^2 found, add them to find $\sum d_i^2 = 194$. The value of n is 10. These values can now be substituted back into the equation

$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

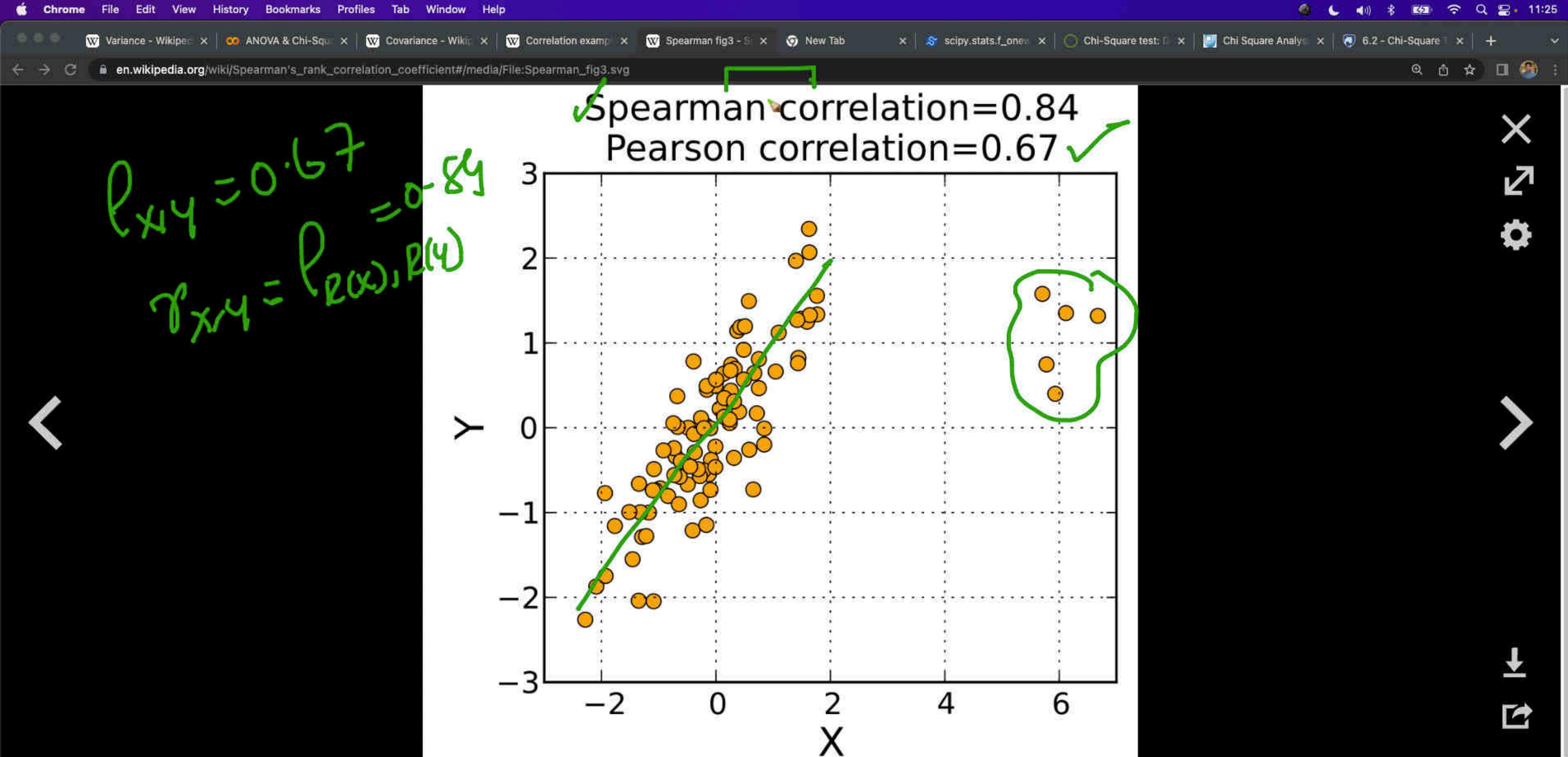
to give

$$\gamma_s = \overbrace{P_{R(x), R(y)}}^{\text{1}}$$



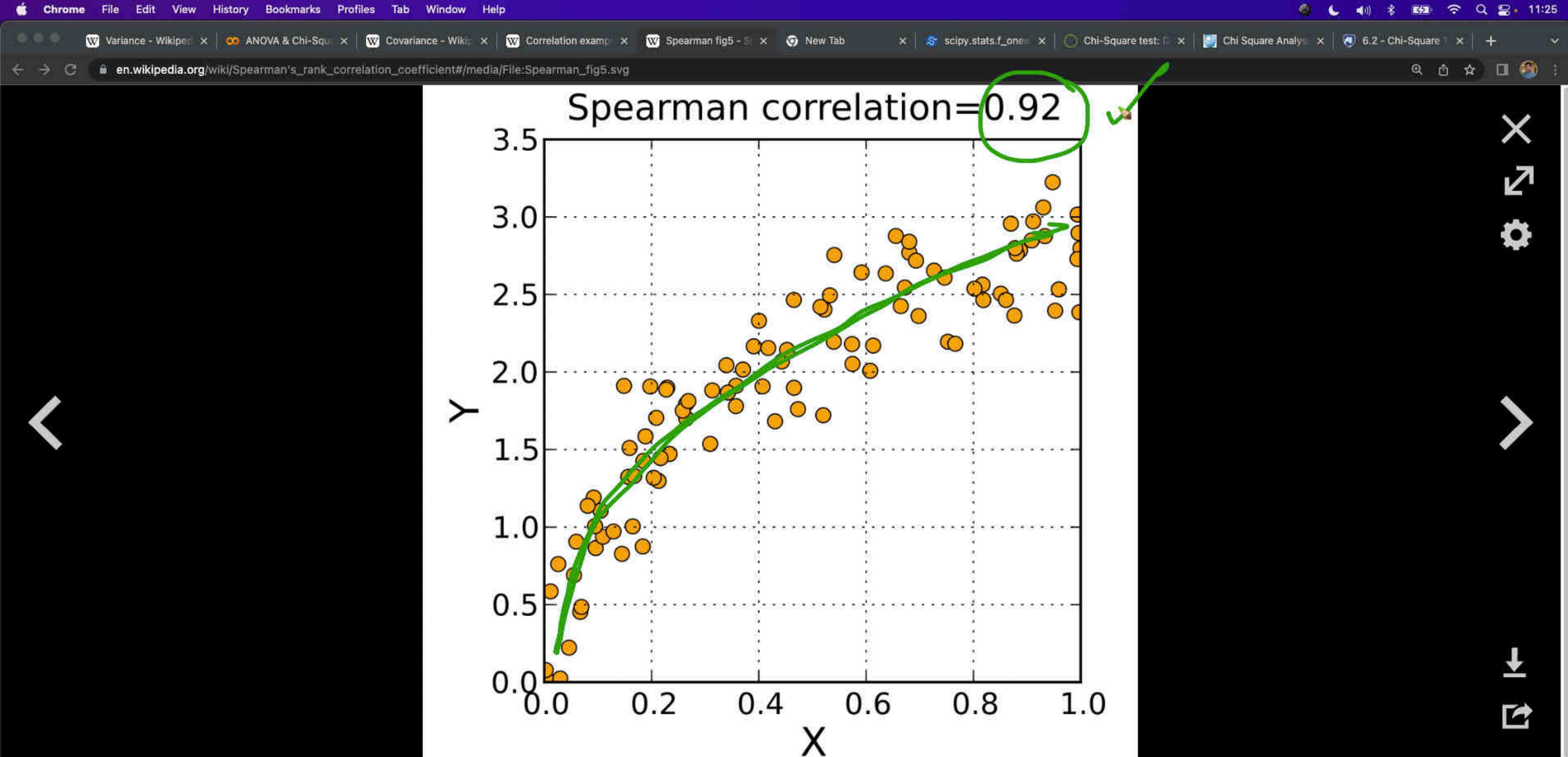
A Spearman correlation of 1 results when the two variables being compared are monotonically related, even if their relationship is not linear. This means that all data points with greater x values than that of a given data point will have greater y values.

More details

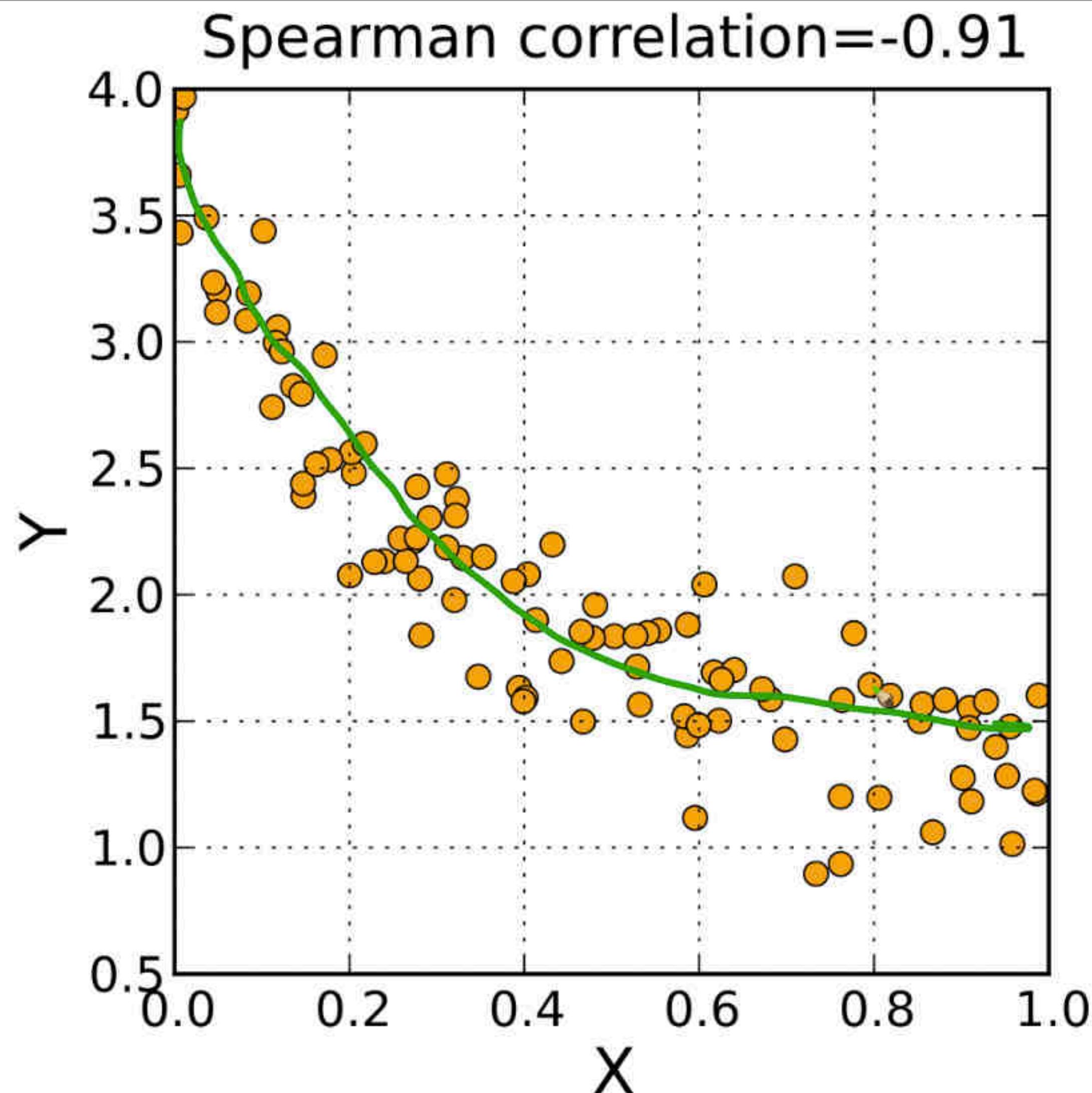


The Spearman correlation is less sensitive than the Pearson correlation to strong outliers that are in the tails of both samples. That is because Spearman's ρ limits the outlier to the value of its rank.

More details



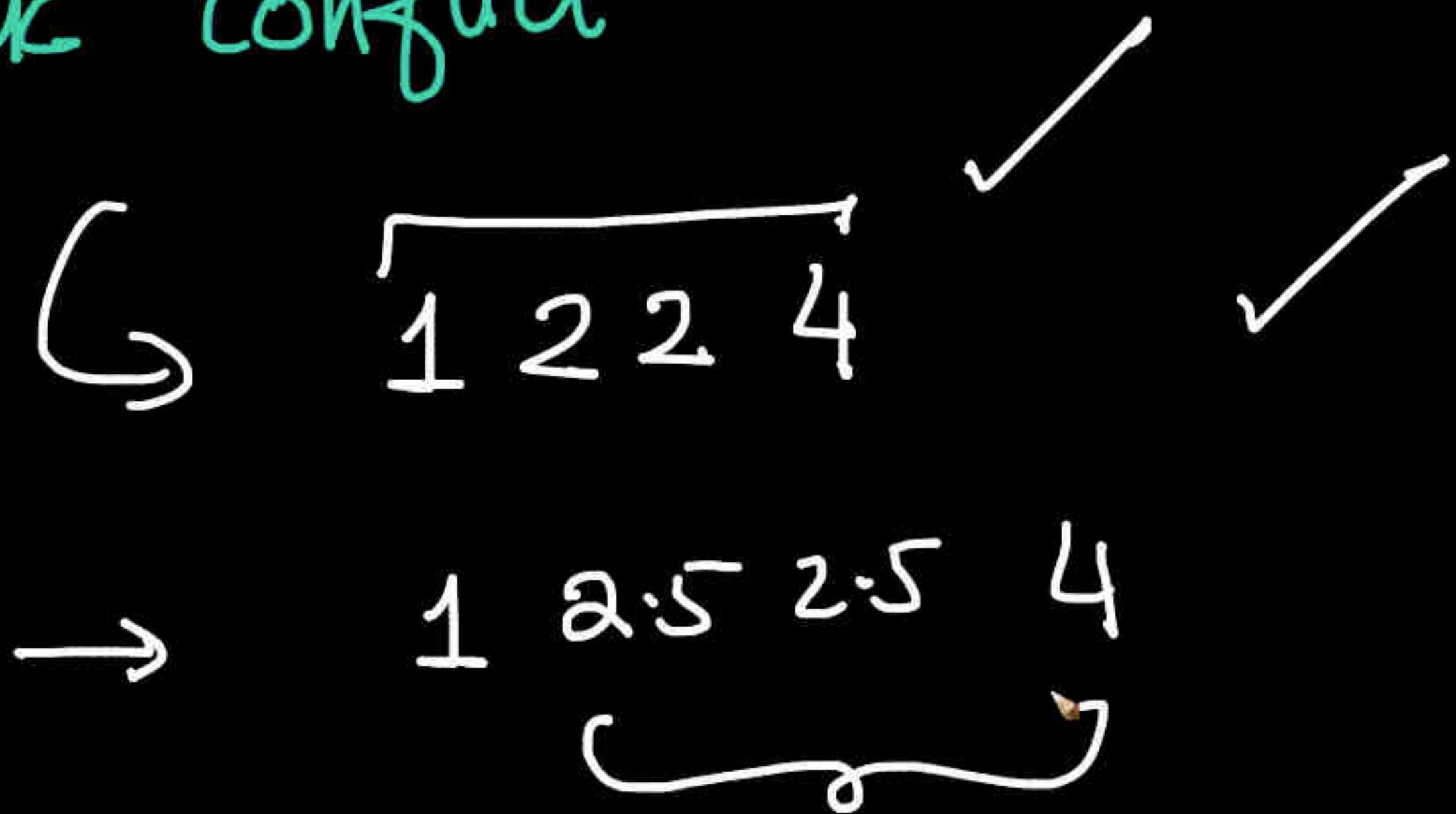
More details



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

More details

rank conflict





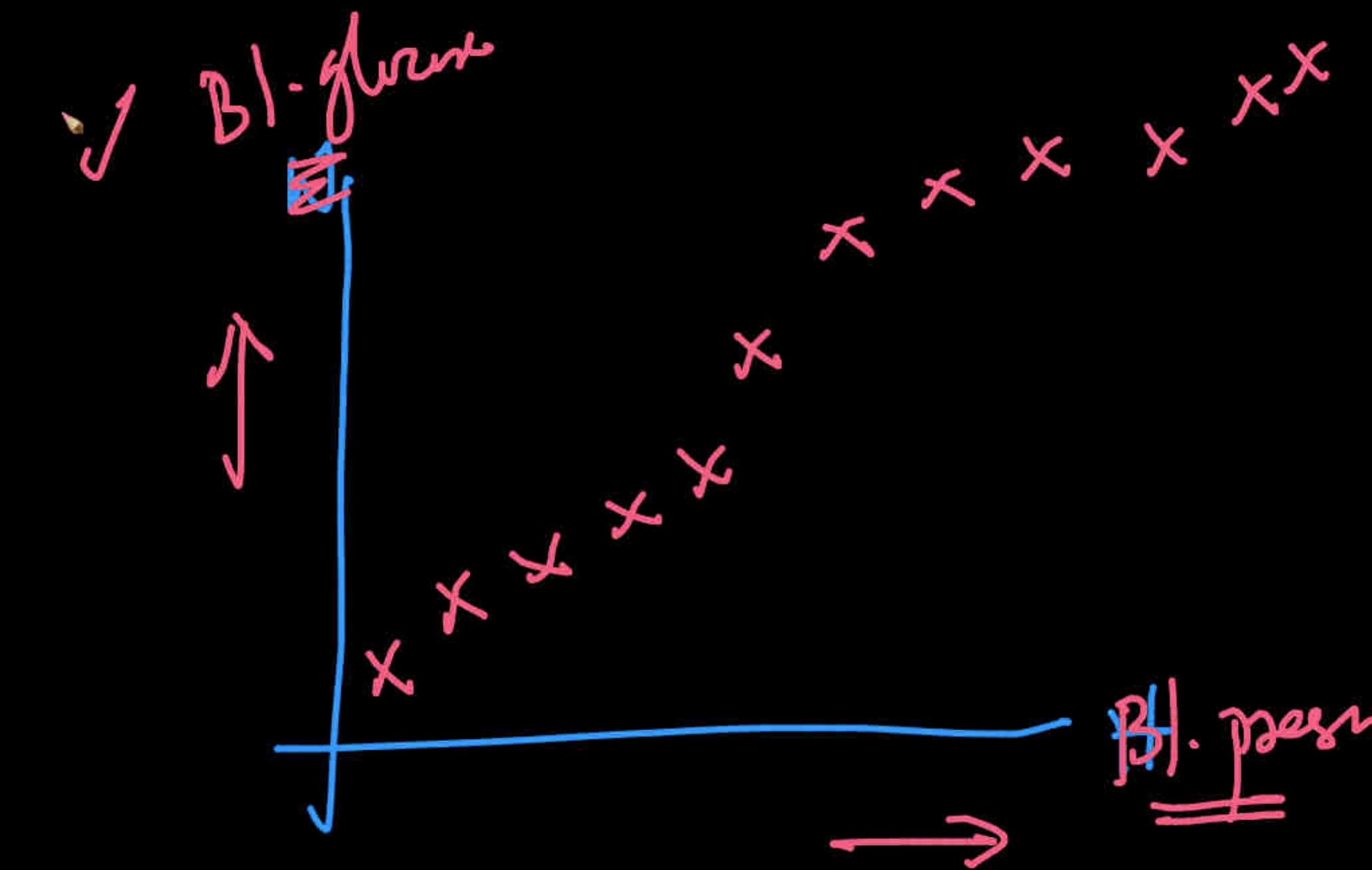


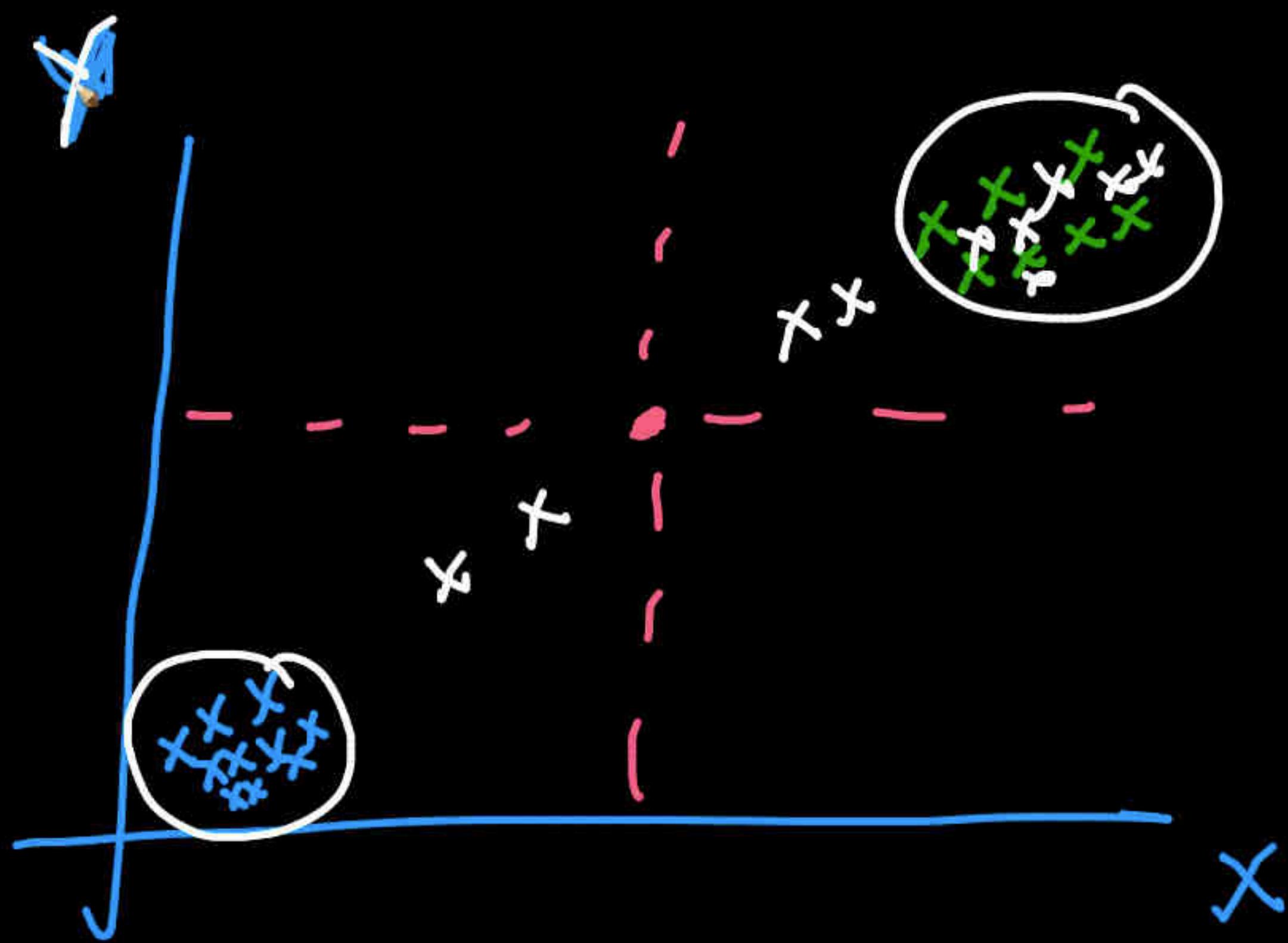
$-\infty \leq \text{Cov}(X, Y) \leq \infty \rightarrow$ As X changes
do Y change?

$-1 \leq \rho_{X,Y} \leq +1 \rightarrow$ linear relationship

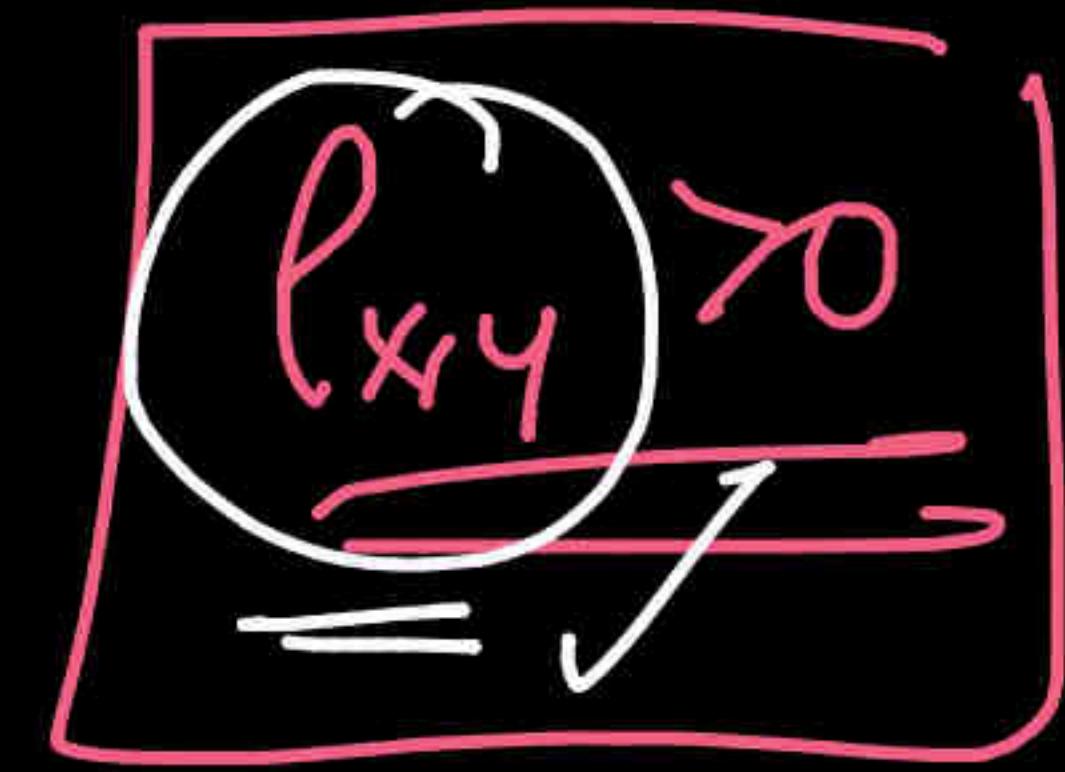
✓ $\left\{ -1 \leq \gamma_{X,Y} \leq +1 \right\} \rightarrow$ monotonicity

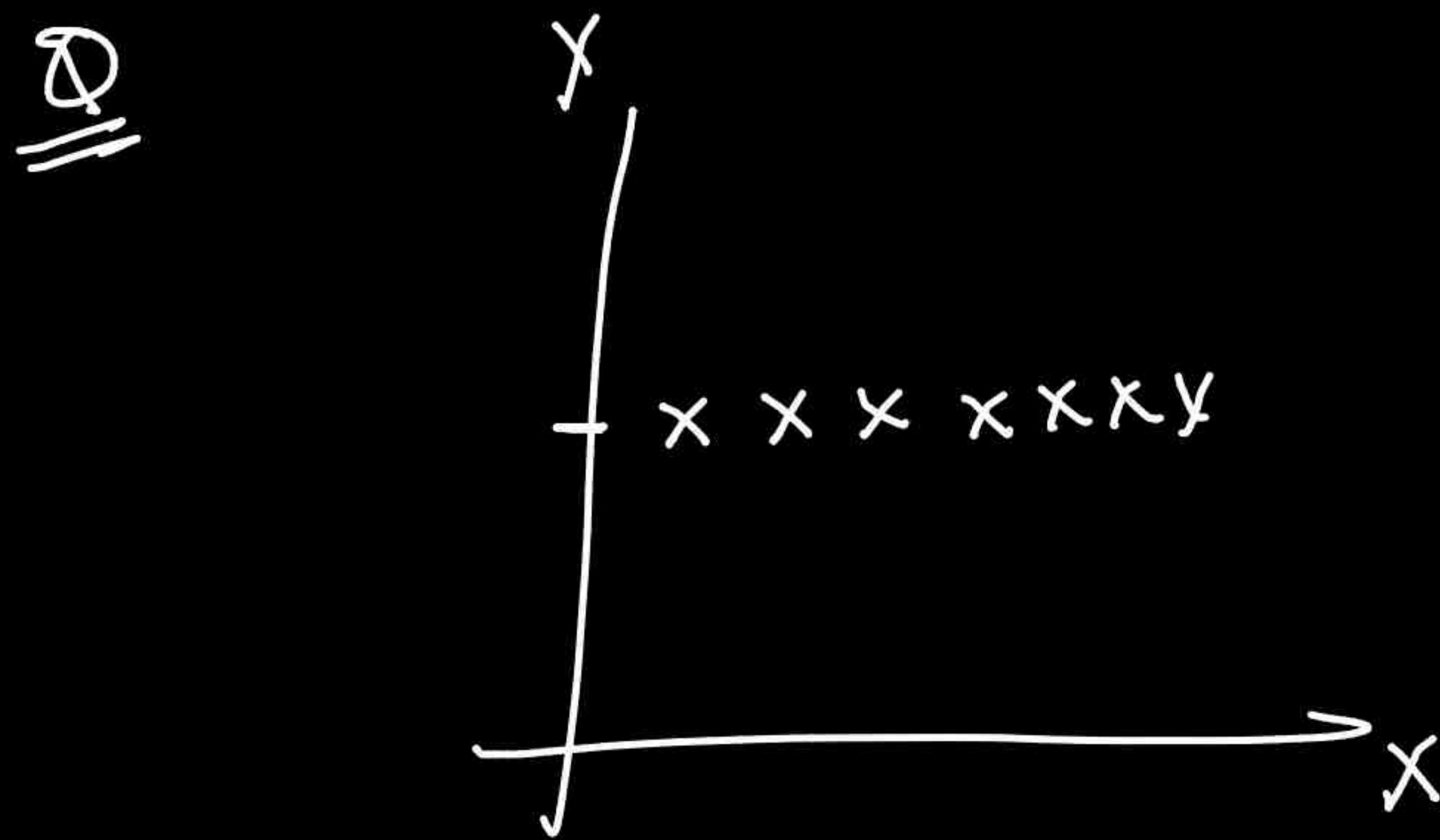
less prone to outliers





Cov \uparrow





$$\frac{Cov}{\sigma_x \sigma_y} = \frac{0}{D} = \text{Undefined}$$

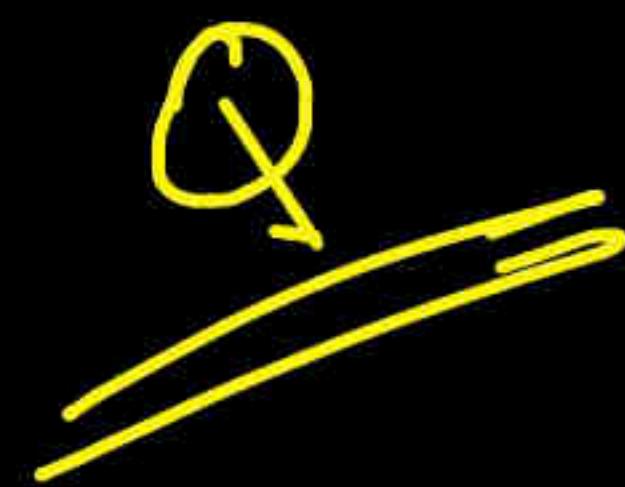
$$\gamma_s = +1$$

$$\left\{ \begin{array}{l} y = x^3 \\ \hline \end{array} \right.$$

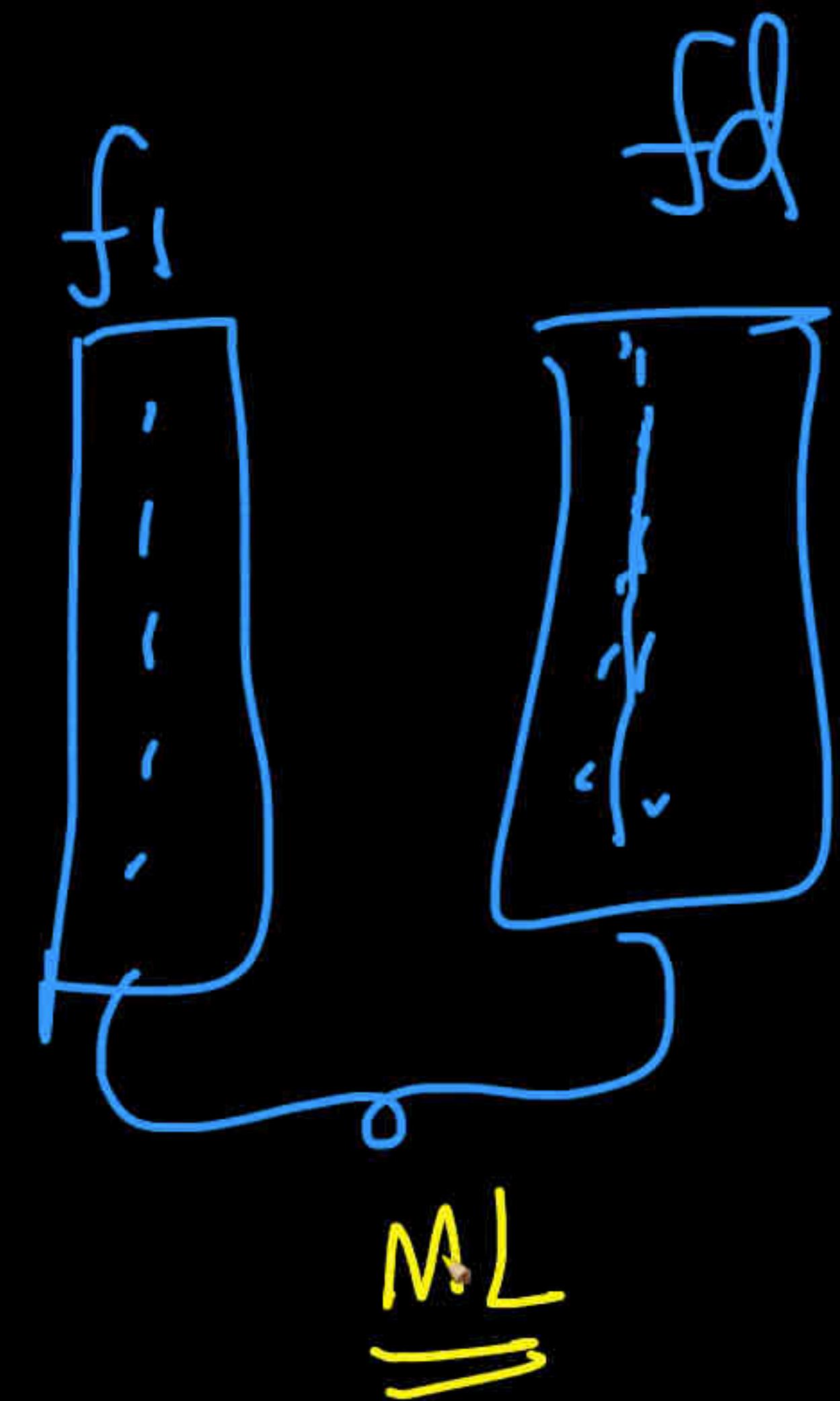
$$y = x^2 + x + 2$$

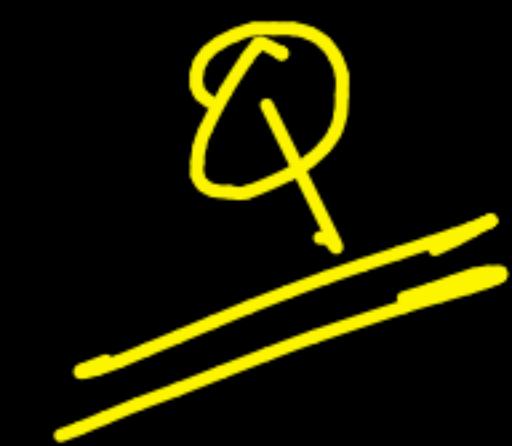
$$\left\{ \begin{array}{l} y = \sqrt{x} \\ \hline \end{array} \right.$$

linear-seg (ML)



Data Analysis





$X \sim \text{disl(params)}$

discrete

$$E(X) = \sum_{i=0}^{\infty} p(x=i) \cdot i$$

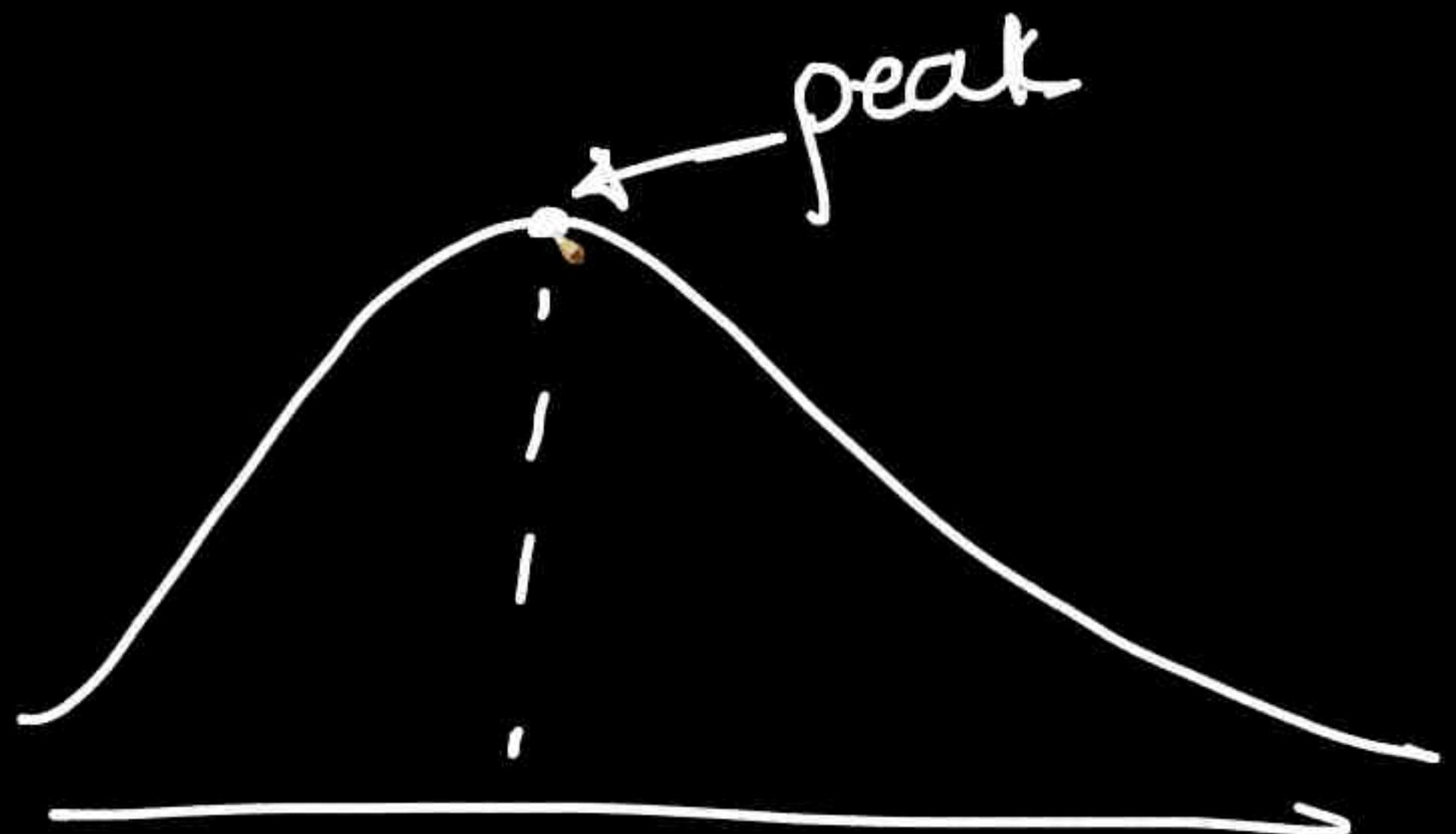
PMF

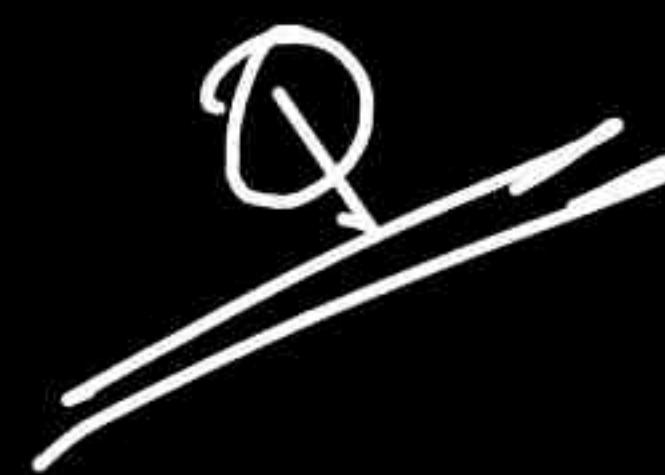
cont:

$$E(X) = \int x f(x) \cdot dx$$

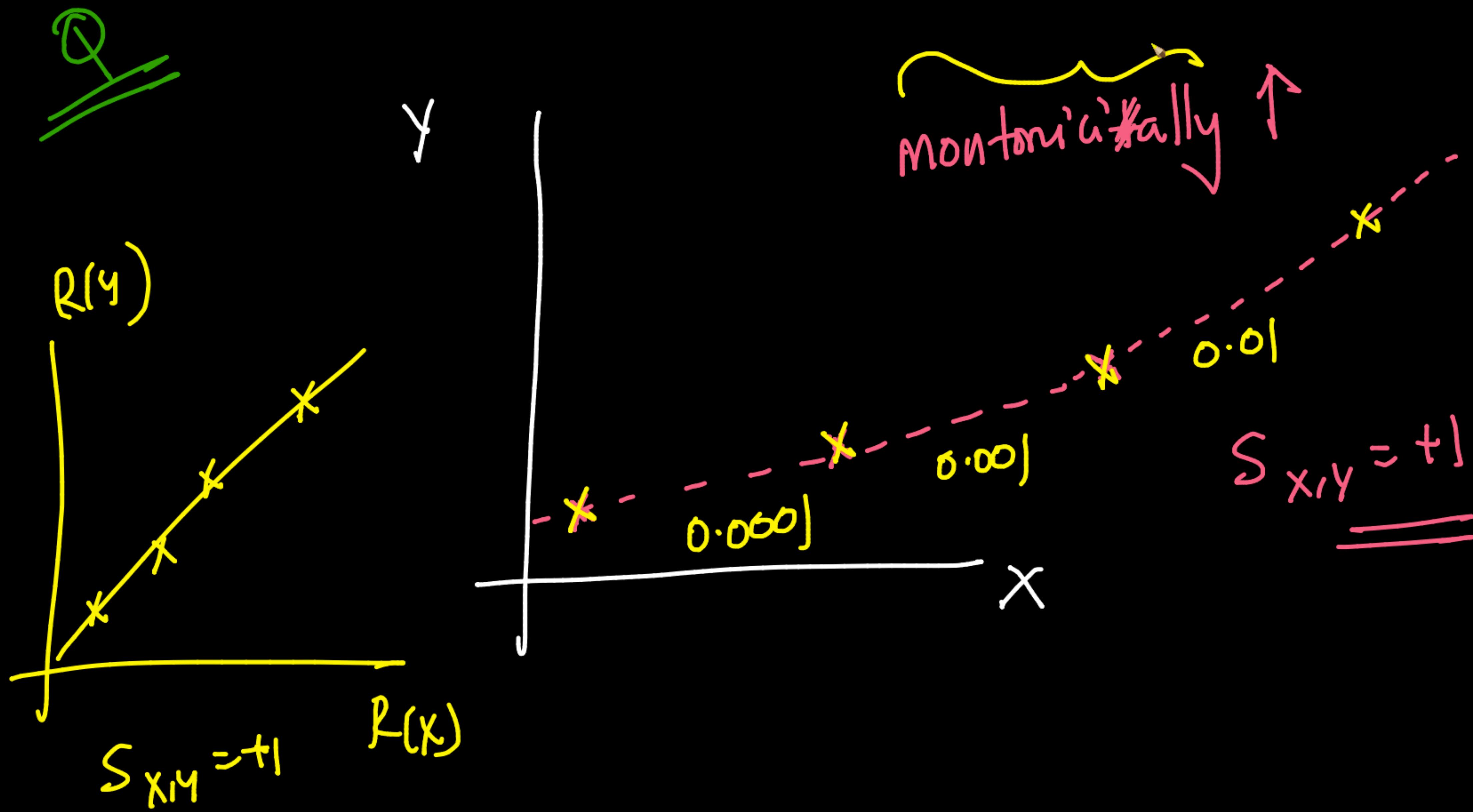
PDF

mode →





Technique → 



ρ denotes the usual Pearson correlation coefficient, but applied to the rank variables,
 $\text{cov}(\text{R}(X), \text{R}(Y))$ is the covariance of the rank variables,
 $\sigma_{\text{R}(X)}$ and $\sigma_{\text{R}(Y)}$ are the standard deviations of the rank variables.

Only if all n ranks are distinct integers, it can be computed using the popular formula

$$r_s = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)},$$

where

$d_i = \text{R}(X_i) - \text{R}(Y_i)$ is the difference between the two ranks of each observation,
 n is the number of observations.

[Proof]

[show]

Identical values are usually^[4] each assigned fractional ranks equal to the average of their positions in the ascending order of the values, which is equivalent to averaging over all possible permutations.

If ties are present in the data set, the simplified formula above yields incorrect results: Only if in both variables all ranks are distinct, then $\sigma_{\text{R}(X)} \sigma_{\text{R}(Y)} = \text{Var}(\text{R}(X)) = \text{Var}(\text{R}(Y)) = (n^2 - 1)/12$ (calculated according to biased variance). The first equation — normalizing by the standard deviation — may be used even when ranks are normalized to $[0, 1]$ ("relative ranks") because it is insensitive both to translation and linear scaling.



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.