

# Correlation analysis





Plus un pays mange de chocolat,  
plus il a de prix Nobel, révèle  
une étude

**Le Point**

Le nombre de prix Nobel dans un pays...  
manger du chocolat

Croquer du chocolat pour  
avoir le Nobel

**LE FIGARO**



Pour décrocher un Nobel

**INSOLITE**

Une étude scientifique affirme que la consommation



Likely?





Study the relationship between two parameters

## Correlation between chocolate consumption and the number of Nobel Prize winners in a country



Chocolate consumption /  
year



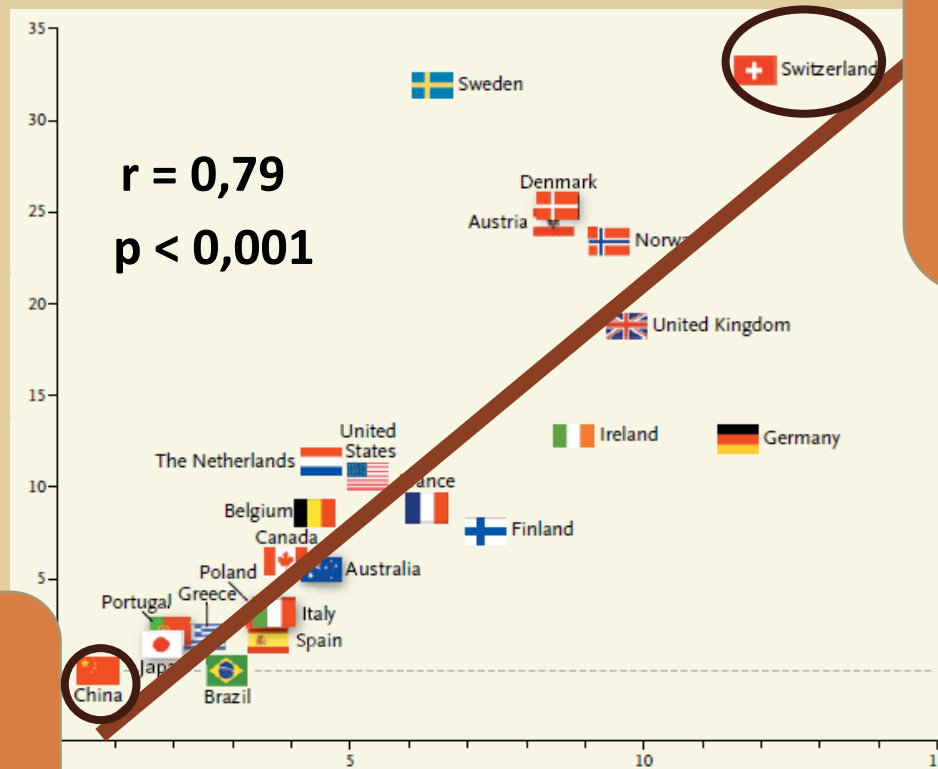
Number of Nobel Prizes  
awarded / year



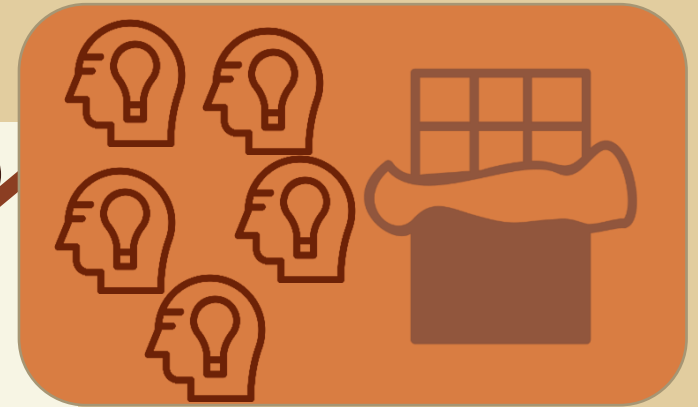
# Study Merlier, 2012



Nobel Prizes awarded /



Chocolate consumption / year



## Correlation



## Correlation: definition

### Correlation:

- measurement of association/relation between 2 variables
- indicates how the 2 variables vary together.

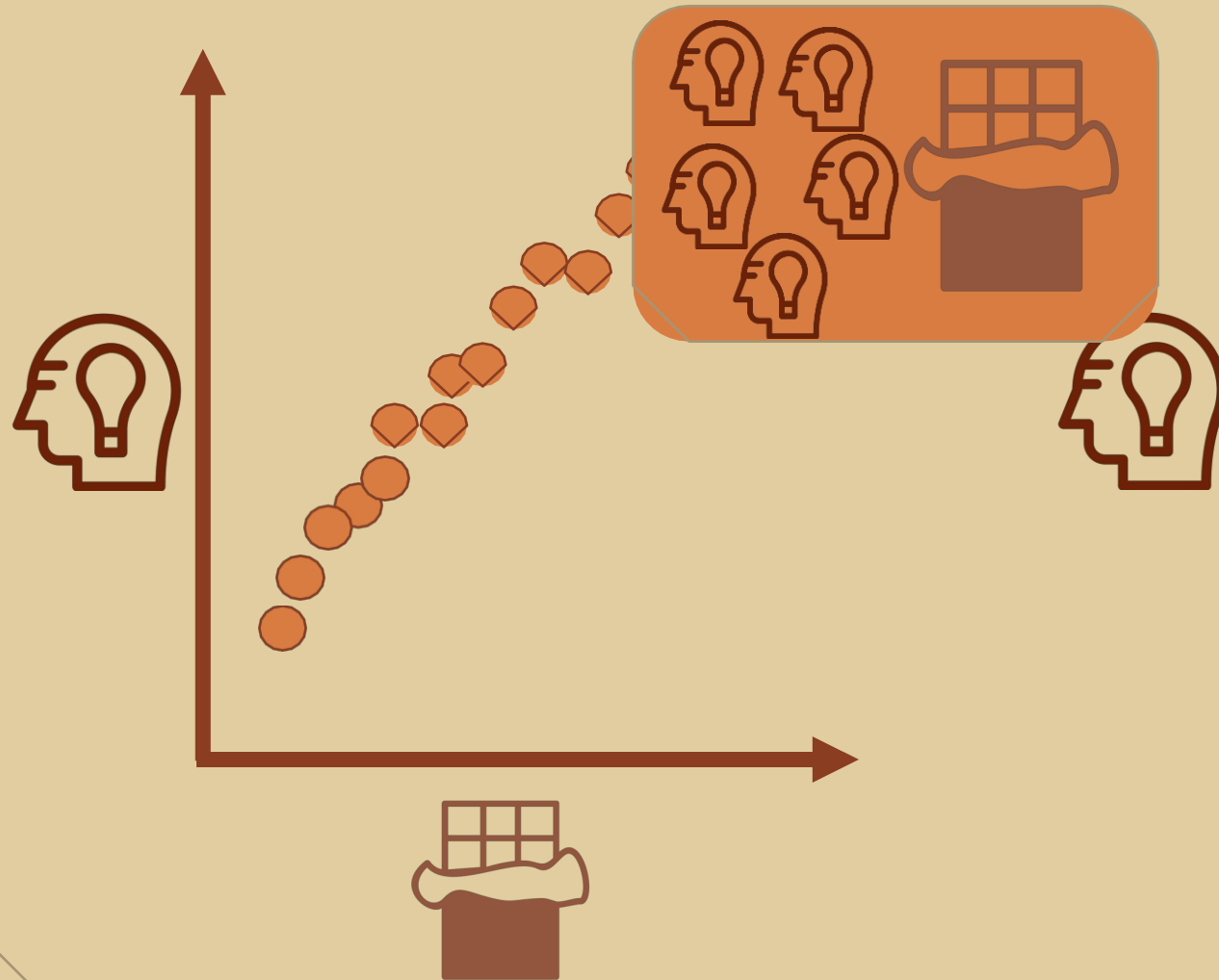


If the values of one variable increase or decrease at the same time as the values of the other variable, then there is a correlation between the two variables.

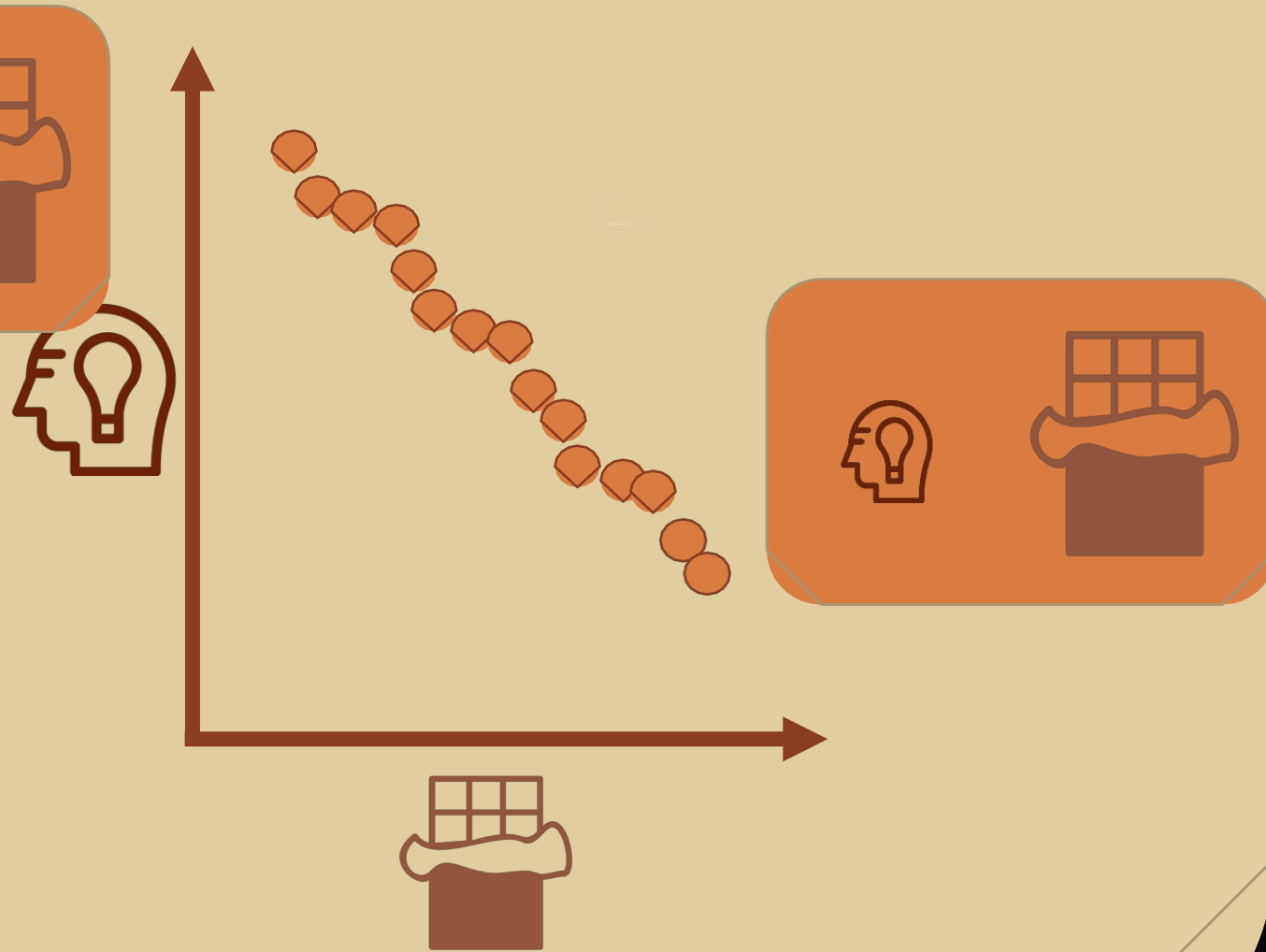
### Correlation:

- Strength of relationship
- Direction

Correlation positive  
linear



correlation negative  
linear





## Correlation: different tests

- Pearson
- Spearman
- Kendall



**Pearson correlation coefficient (r):** an indicator that measures the strength and direction of a linear relationship between two continuous variables.

**Negative correlation :**

When values from one variable increases, the values from the other variable decrease



$[-1, \dots$

$0,$

$\dots$

$1]$

**Positive correlation:**

When values from one variable increases, the values from the other variable increase.



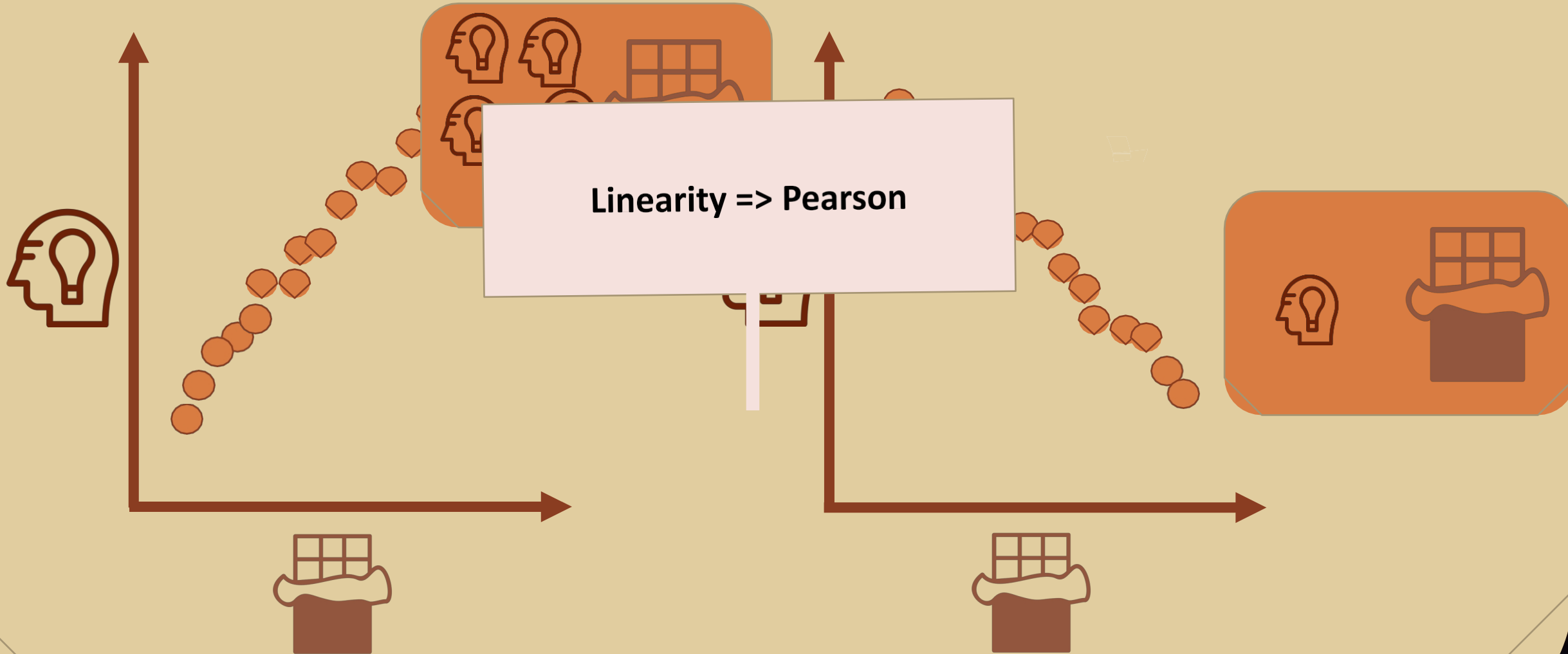
No linear correlation



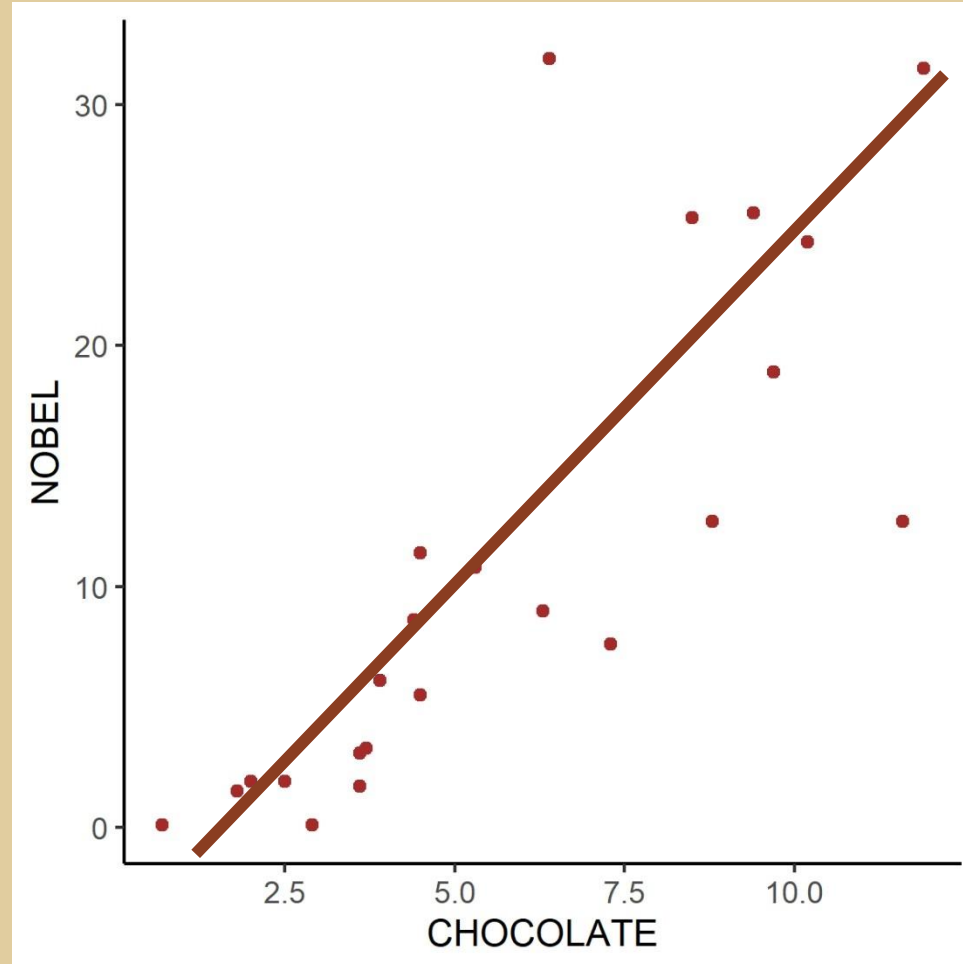
Correlation positive  
linear

correlation negative  
linear

Linearity => Pearson

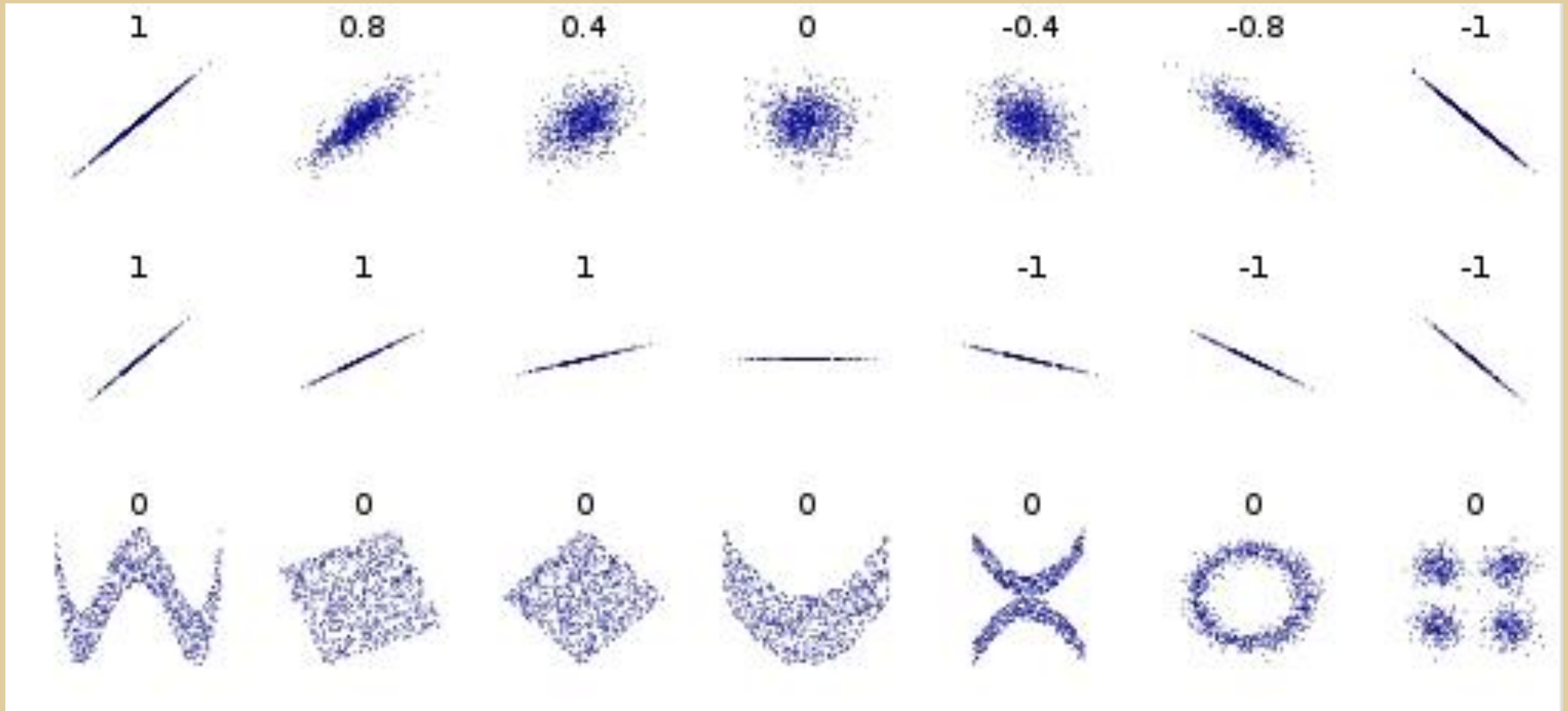


The purpose of the linear correlation coefficient is to quantify the more or less linear aspect of the scatterplot.





# Difference between sense of relationship and strength of relationship



## Hypothesis testing:

Rarely correlation of 0. When do we consider we have a significant correlation?

- $H_0$ : there is no correlation in the population ( $\rho = 0$ )
- $H_1$ : there is a correlation in the population ( $\rho \neq 0$ ).

The hypothesis test will then determine whether  $H_0$  can be rejected in favor of  $H_1$ , i.e. whether the observed correlation is significantly different from zero.



```
> cor.test(data$CHOCOLATE,data$NOBEL,method="pearson")
```

Pearson's product-moment correlation

```
data: data$CHOCOLATE and data$NOBEL  
t = 6.123, df = 21, p-value = 4.477e-06  
alternative hypothesis: true correlation is not equal to 0  
95 percent confidence interval:  
 0.5797205 0.9118788  
sample estimates:  
      cor  
0.8006078
```

## Application conditions for Pearson correlation (parametric test)

- 2 continuous variables
- paired data ( observations have values in both variables )



Can you place a point on a 2D plan if you only have one coordinate ?

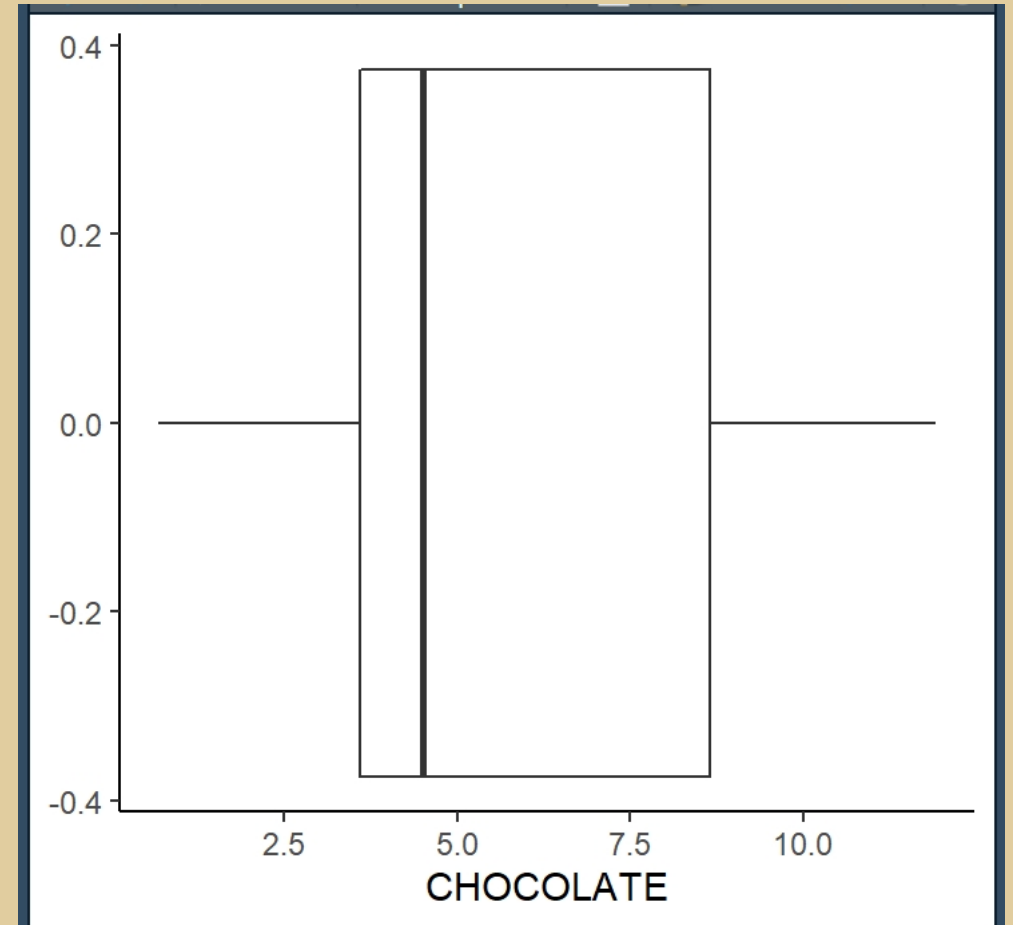
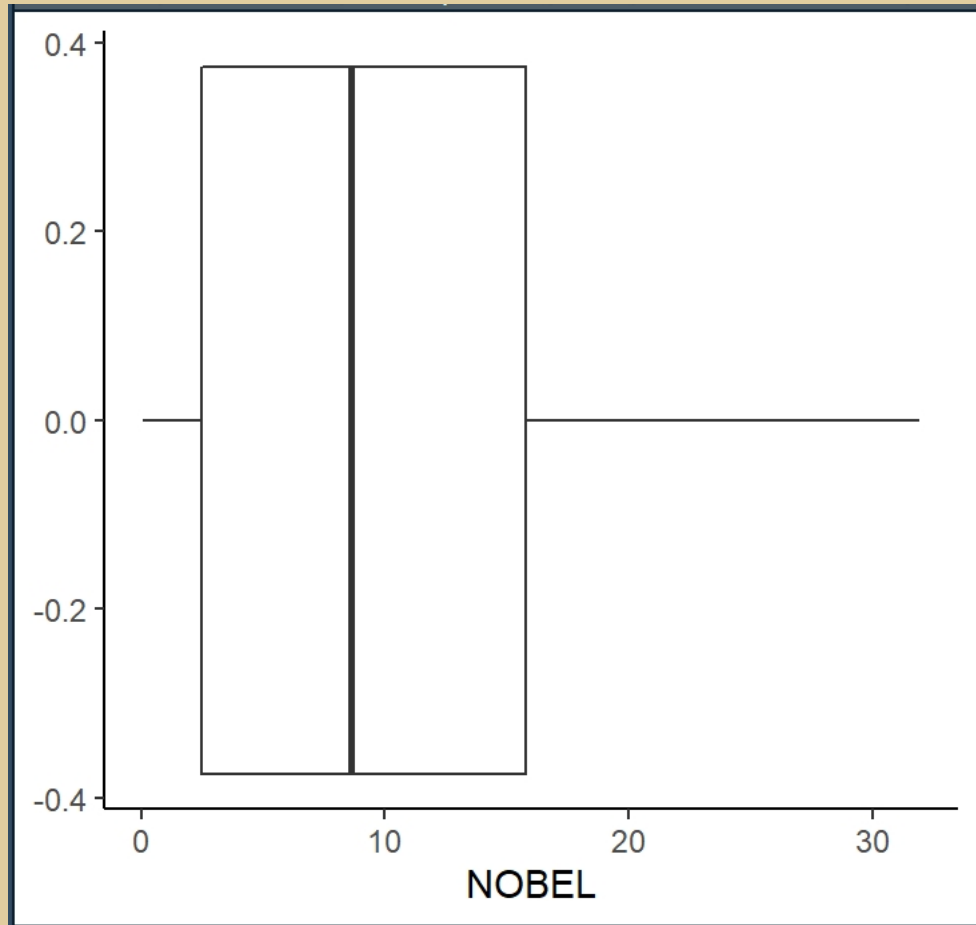


## Application conditions for Pearson correlation (parametric test)

- 2 continuous variables
- paired data ( observations have values in both variables )
- Independent observations
- linear relationship between the 2 variables
- Normal distribution of the two variables
- Outliers absence



## outliers





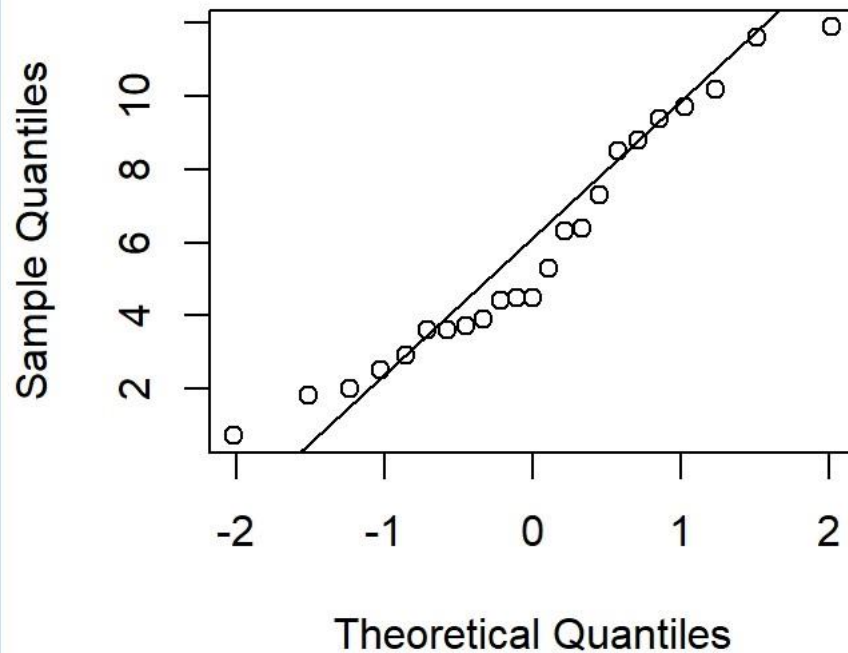
## Chocolate consumption

```
> shapiro.test(data$CHOCOLATE)
```

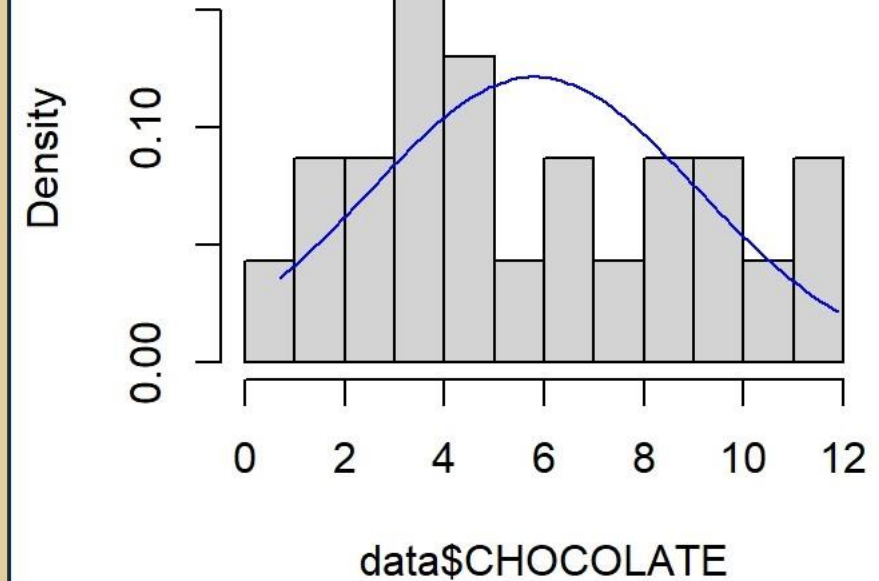
Shapiro-wilk normality test

```
data: data$CHOCOLATE  
W = 0.94223, p-value = 0.2006
```

Normal Q-Q Plot



Histogramme de données



## No. of Nobel Prizes

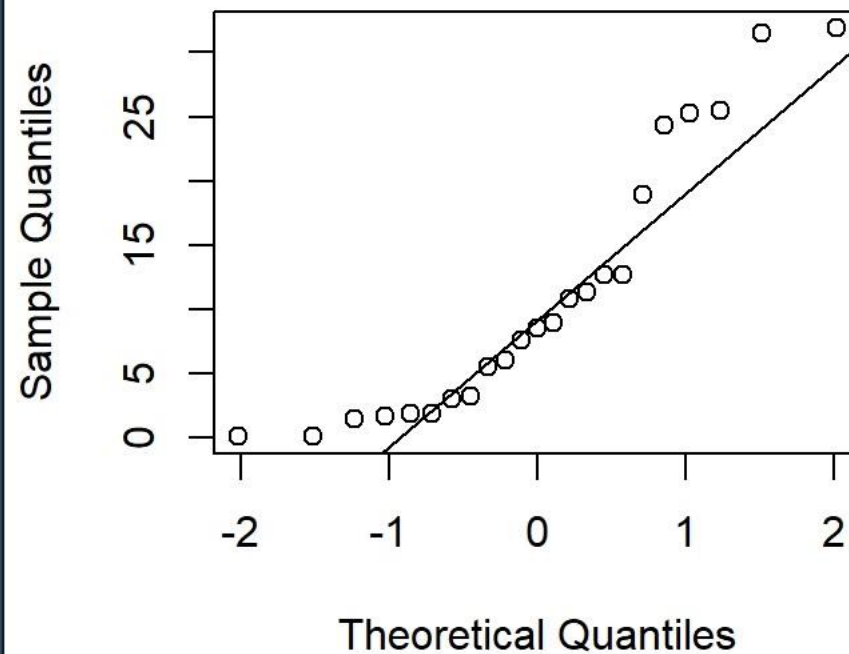
```
> shapiro.test(data$NOBEL)
```

Shapiro-Wilk normality test

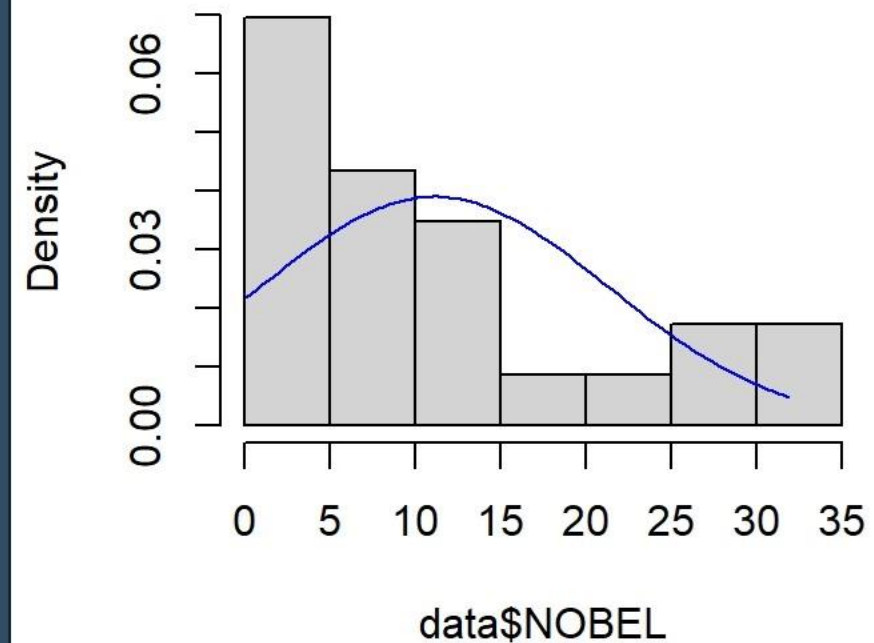
data: data\$NOBEL

W = 0.87014, p-value = 0.006449

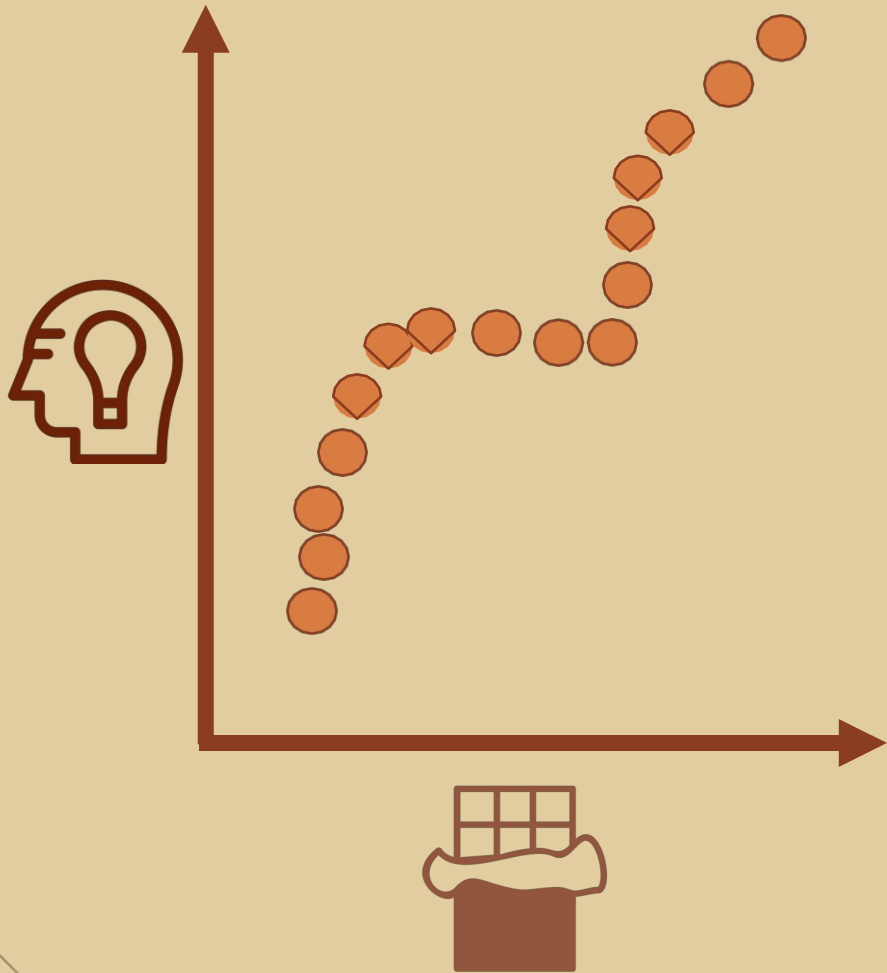
Normal Q-Q Plot



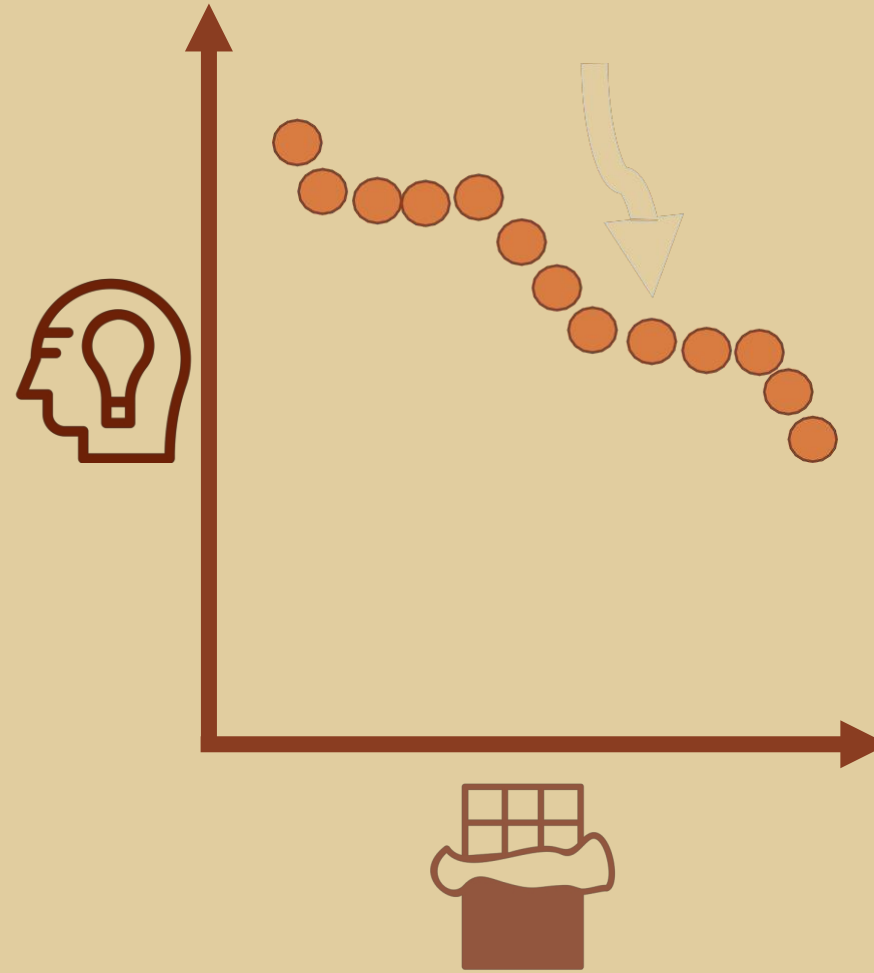
Histogramme de données




Monotonic evolution  
Non-linear positive  
correlation



Monotonic evolution  
Non-linear negative  
correlation








Linear  
path:  
Pearson

monotonous  
path:  
Sperman



Linear  
path:  
Pearson

Neither linear nor  
monotonic path: no  
correlation test

Monotonous  
path:  
Sperman



```
> cor.test(data$CHOCOLATE, data$NOBEL, method = "spearman")
```

Spearman's rank correlation rho

data: data\$CHOCOLATE and data\$NOBEL

S = 197.74, p-value = 4.003e-09

alternative hypothesis: true rho is not equal to 0

sample estimates:

rho

0.9023003



## Spearman correlation (non-parametric test)

Application conditions :

- The 2 variables must be numerical or ordinal.

*Ordinal variables are ordered categorical variables, for example, levels of education (primary, secondary, tertiary) or grades (A, B, C).*

- Paired data (each observation has a value for both variables)
- Observation independence
- Monotone relationship between the two variables

```
> cor.test(dat$CHOCOLATE, dat$NOBEL, method = "spearman")
```

Spearman's rank correlation rho

data: dat\$CHOCOLATE and dat\$NOBEL

S = 6.5912, p-value = 0.04986

alternative hypothesis: true rho is not equal to 0

sample estimates:

rho

0.8116794

## Correlation

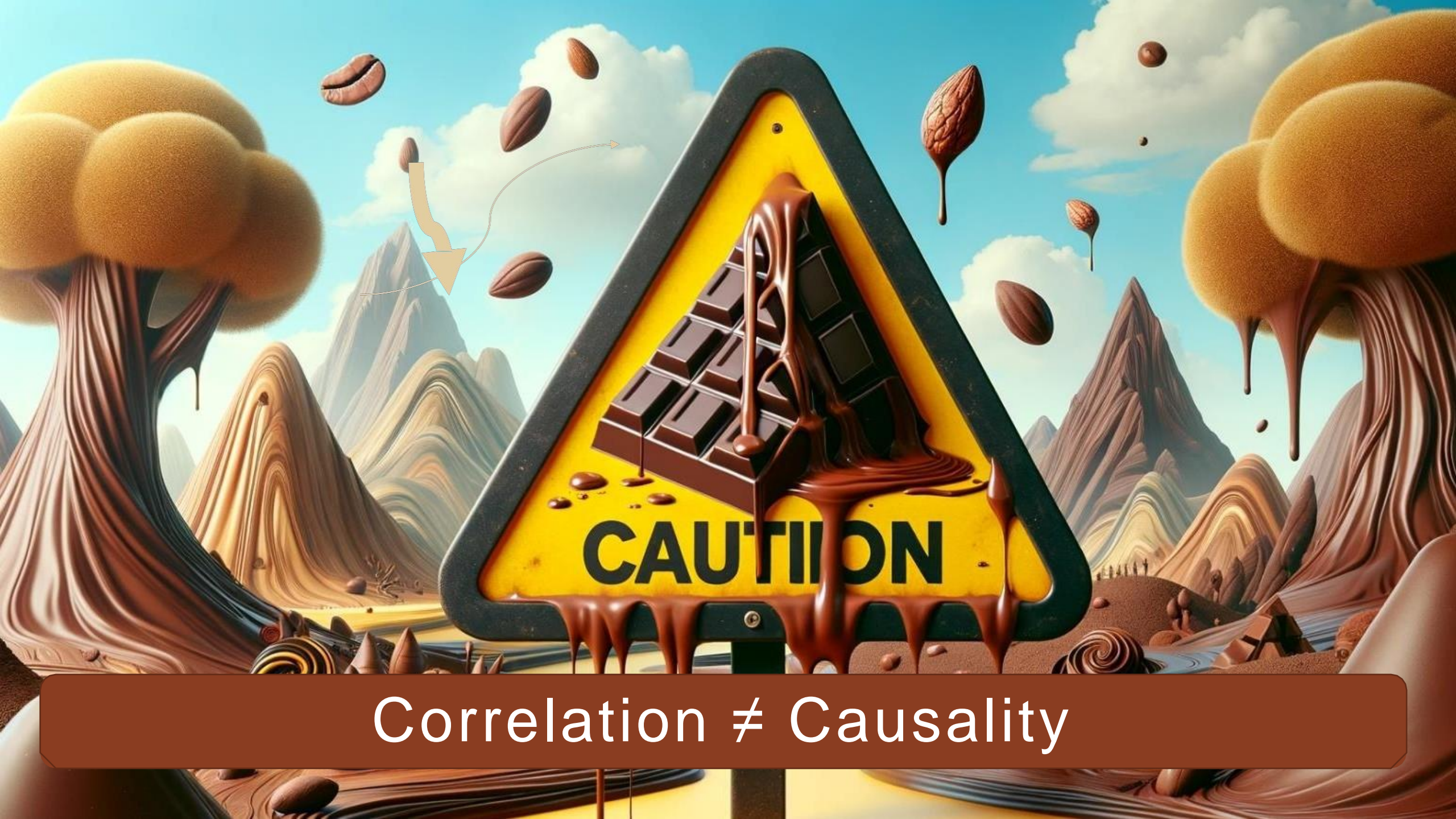
$r = 0,1$  (weak)

$r = 0,5$  (moderate)

$r = 0,7$  (strong)

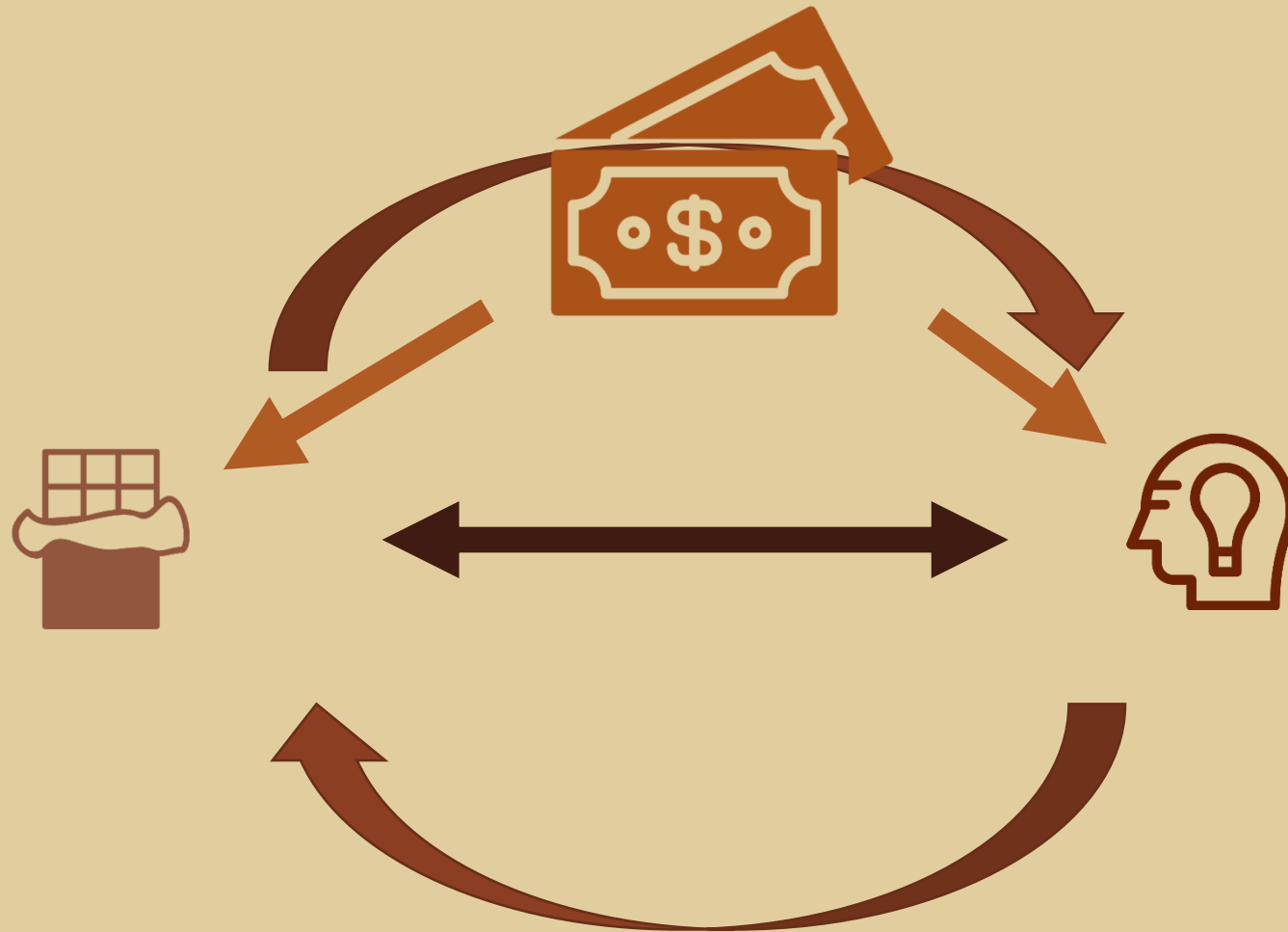
$r = 0,9$  (very strong)





Correlation  $\neq$  Causality

## Correlation vs. causation





# spurious correlations

*correlation is not causation*

[random](#) · [discover](#) · [next page](#) →

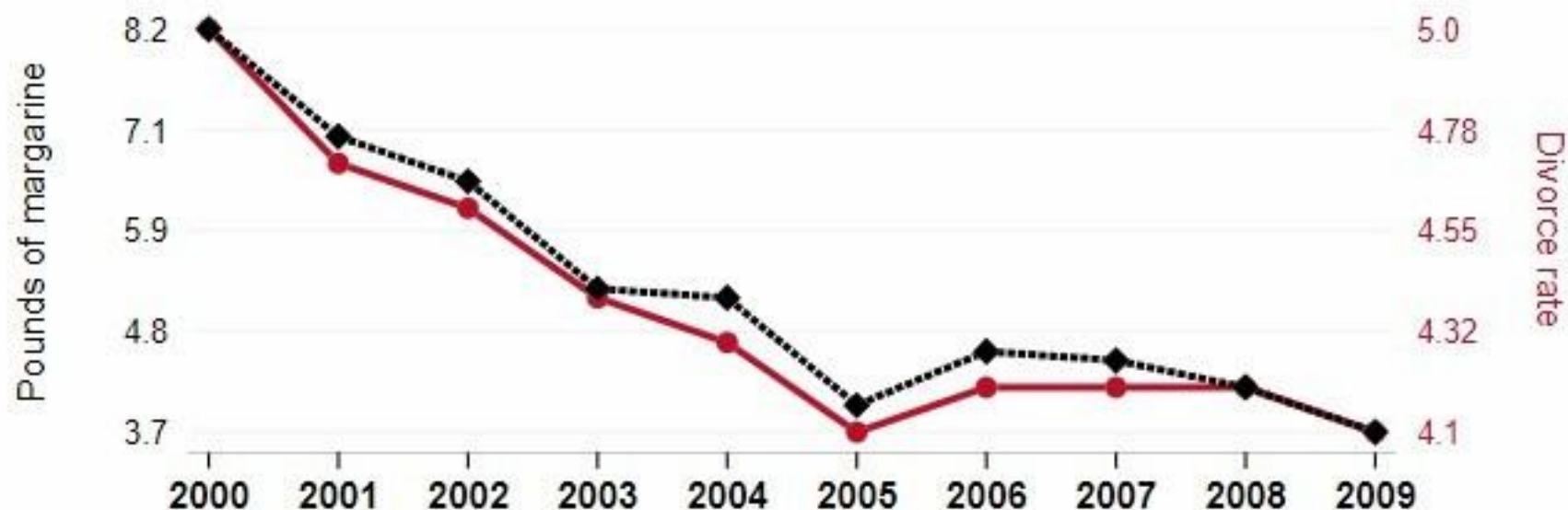
don't miss [spurious scholar](#),  
where each of these is an academic paper



## Per capita consumption of margarine

correlates with

## The divorce rate in Maine



◆ Per capita consumption of margarine in the United States - Source: US Department of Agriculture

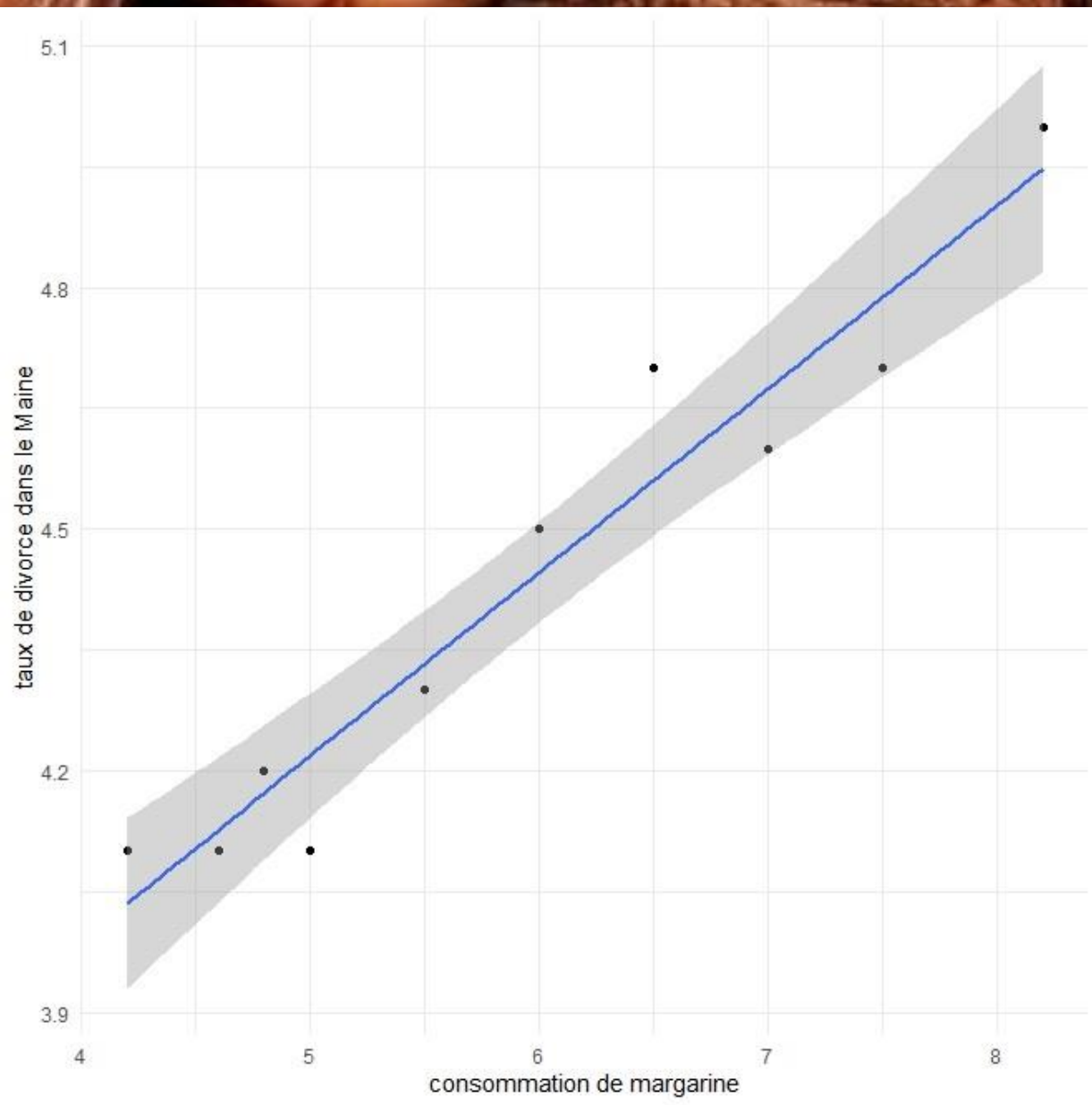
● The divorce rate in Maine - Source: CDC National Vital Statistics

2000-2009,  $r=0.993$ ,  $r^2=0.985$ ,  $p<0.01$  - [tylervigen.com/spurious/correlation/5920](http://tylervigen.com/spurious/correlation/5920)

**[View details about correlation #5,920](#)**

*Spreading Love and Margarine: An Examination of the Butter-Splitter Correlation in Maine*





**People who drowned after falling out of a fishing boat**  
correlates with  
**Marriage rate in Kentucky**

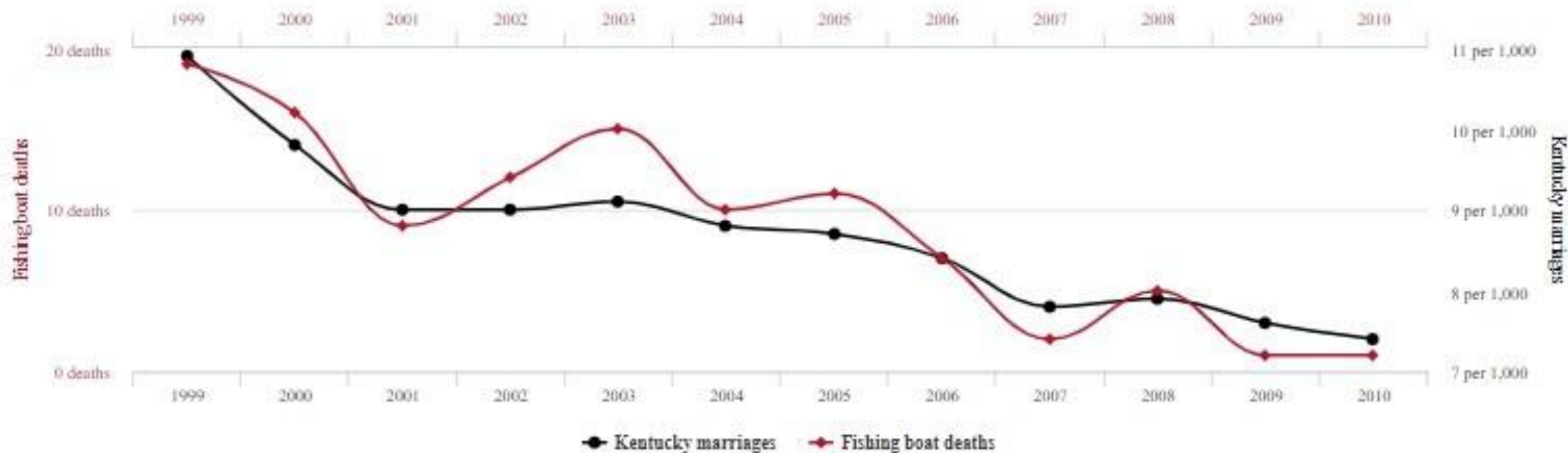
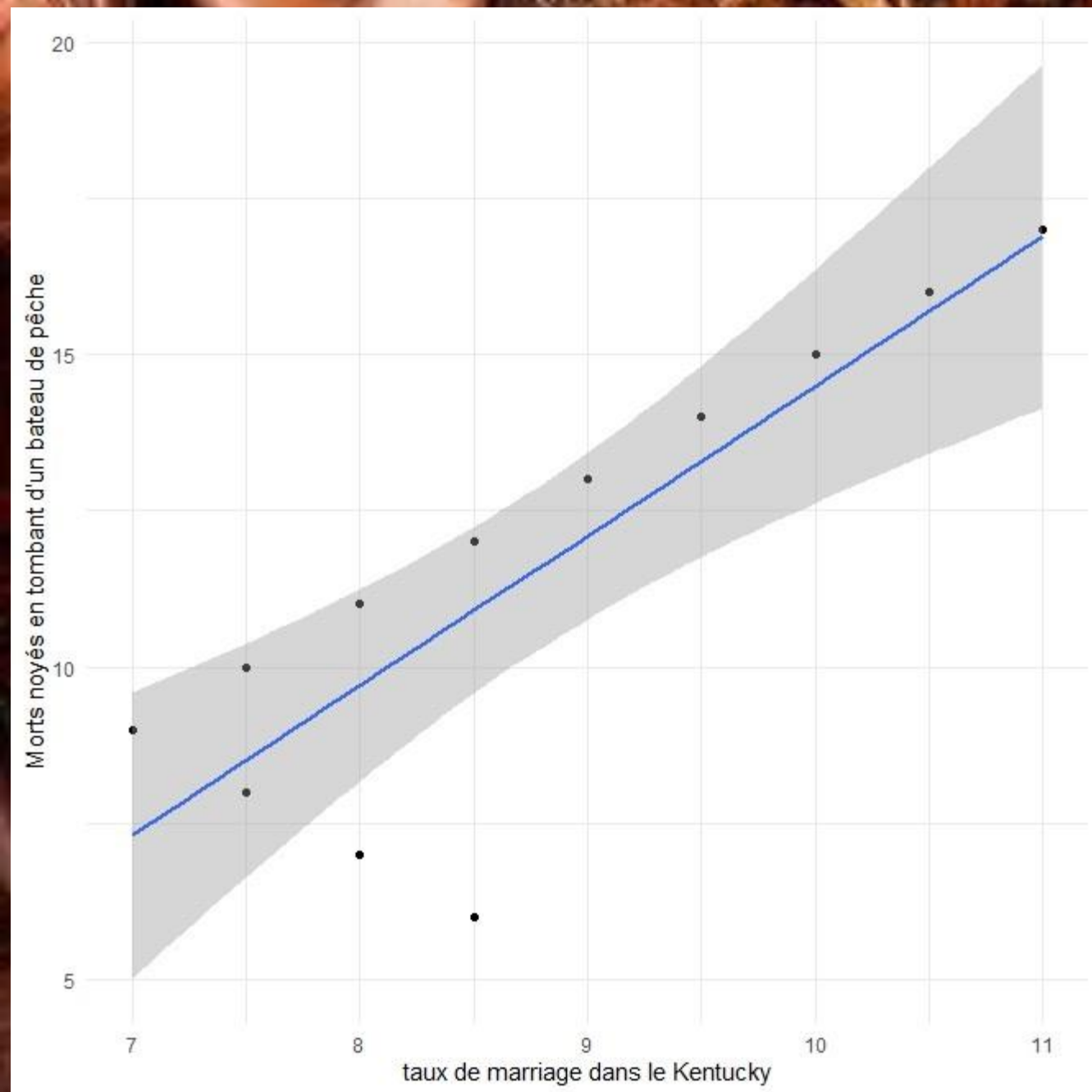


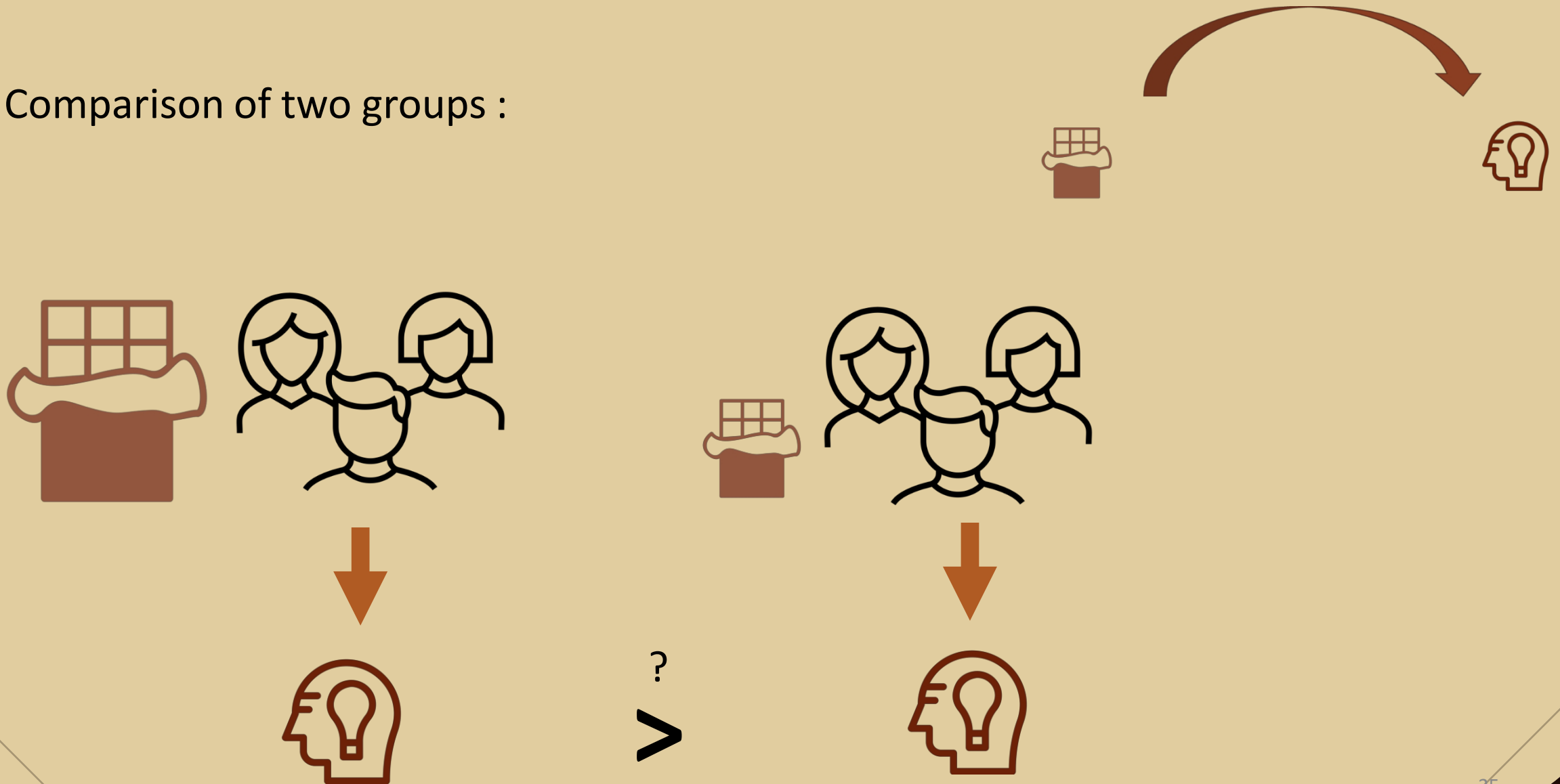
Figure 1.1: Spurious correlations, by Tyler Vigen, licenced under CC BY 4.0





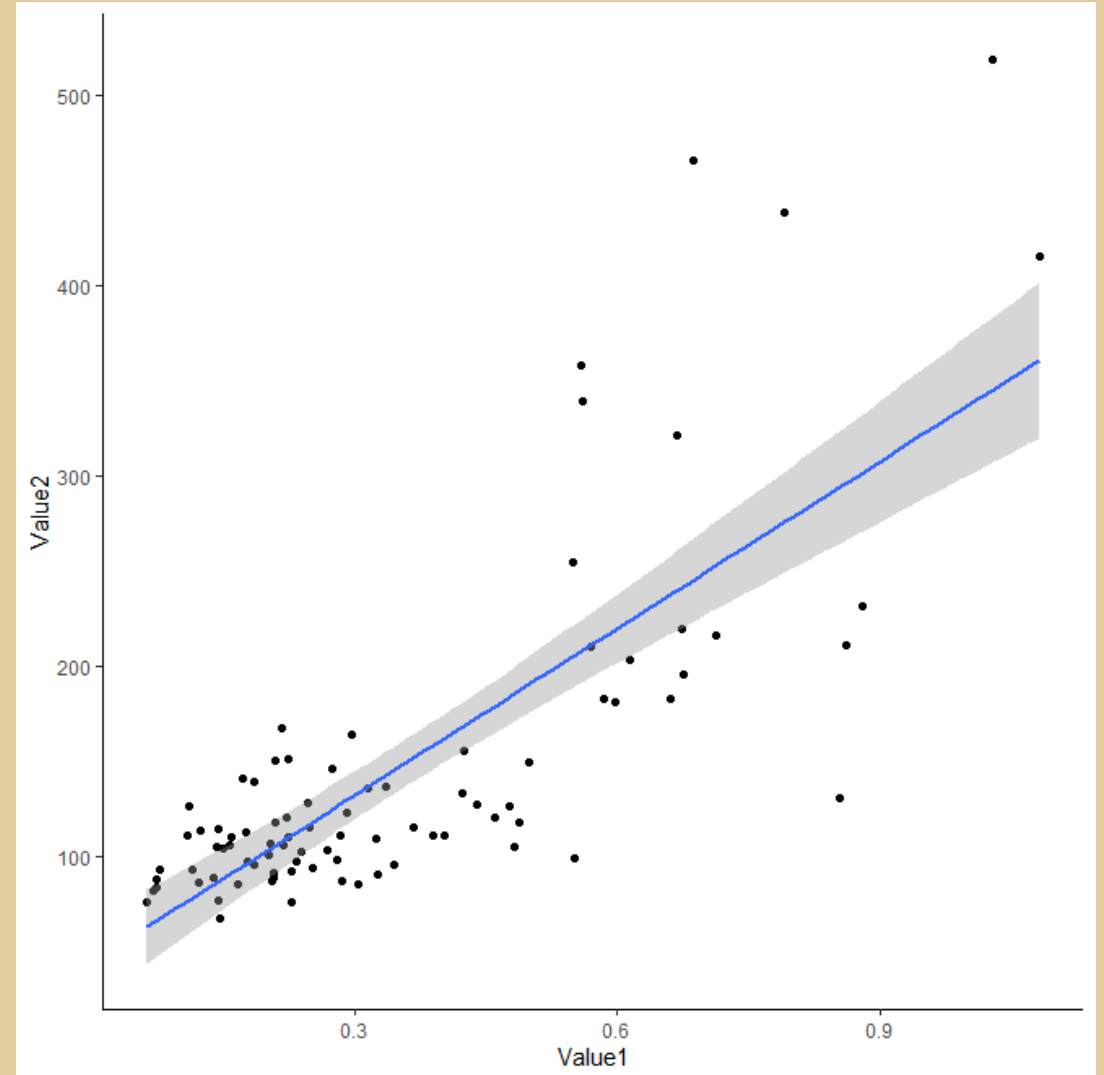
## Determining causality through experience

Comparison of two groups :





## Your field





# Your field

```
> print(cor_test)
```

Pearson's product-moment correlation

data: data\$value1 and data\$value2

t = 10.931, df = 85, p-value < 2.2e-16

alternative hypothesis: true correlation is not equal to 0

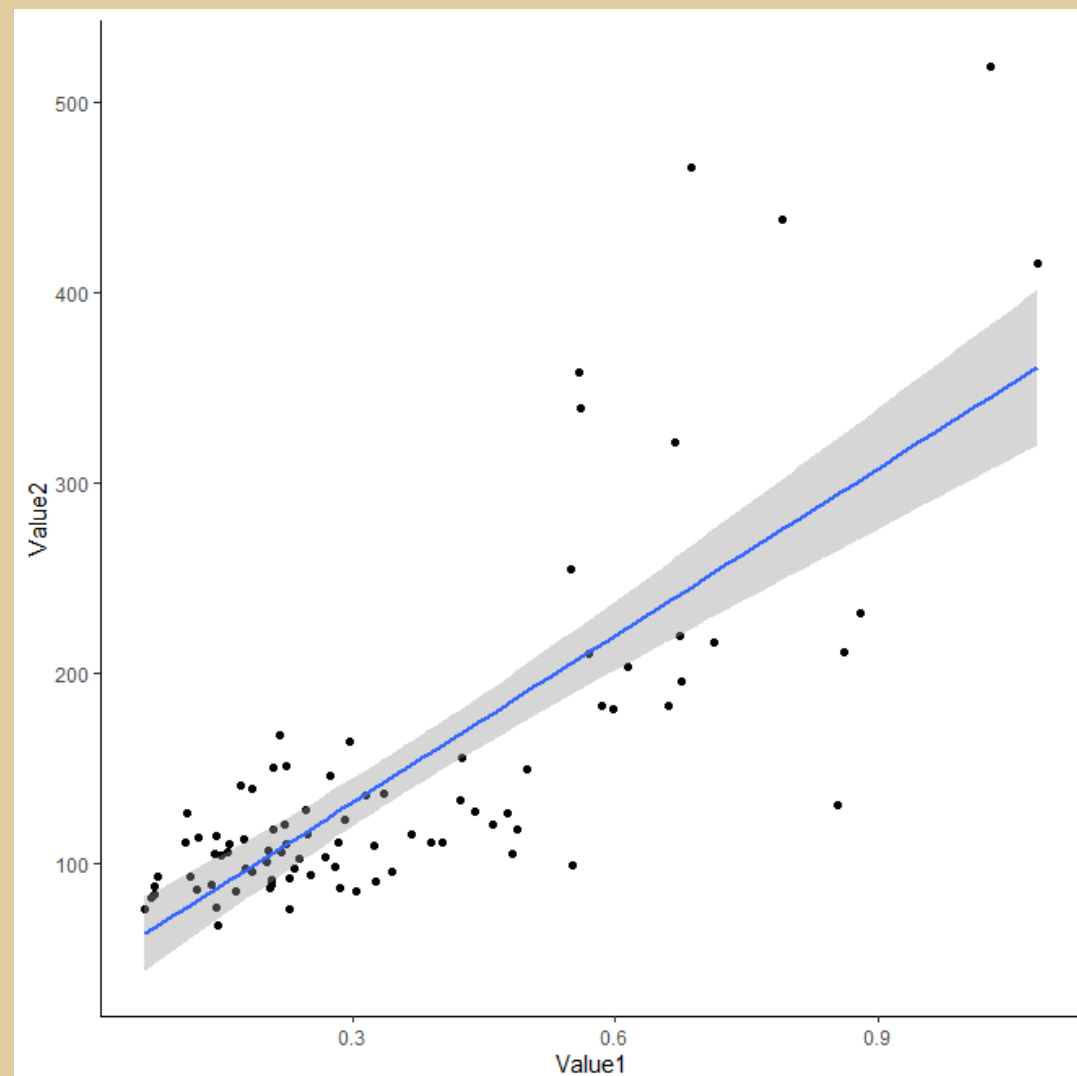
95 percent confidence interval:


0.6600598 0.8398366

sample estimates:

cor

0.7644224



A hand is holding a brown cardboard box, which is the central focus of the image. The box has a handle cutout at the top. Overlaid on the box is the text 'Take home message' and 'Test correlation'. Below this, a bullet point states '- Study the relationship between two quantitative variables'. The background is a blurred image of a person's arm and torso in a white shirt.

# Take home message

## Test correlation

- Study the relationship between two quantitative variables