# MACHNE LEARNING TECHNIQUES AND APPLICATIONS

## M2 DISS

# WHO AM I

- Rémy Cazabet ([remy.cazabet@univ-lyon1.fr](mailto:remy.cazabet@univ-lyon1.fr))

- Associate professor, LIRIS Laboratory, Lyon 1 University

- Team: Data Mining and Machine Learning (DM2L)

- Lyon's Institute of Complex Systems (IXXI)

# WHO AM I

- Research topics:
  - Large Network Analysis (Cryptocurrencies…)
  - Graph Clustering
  - Dynamic network
  - Graph Embedding
  - Graph Neural Networks

- Research internship.

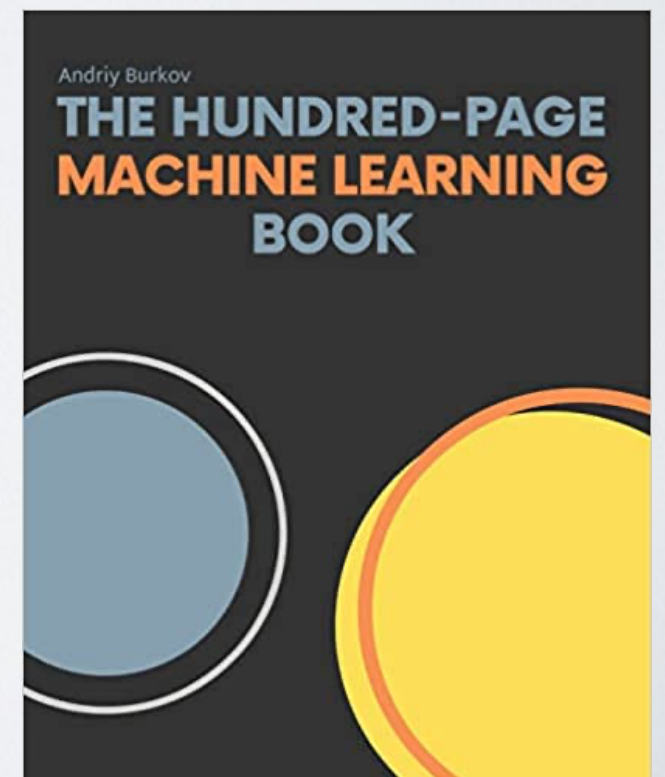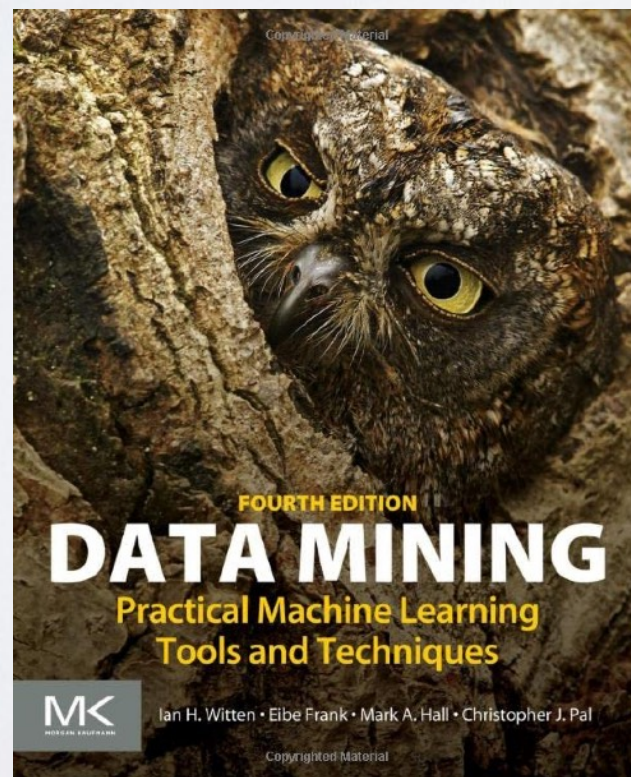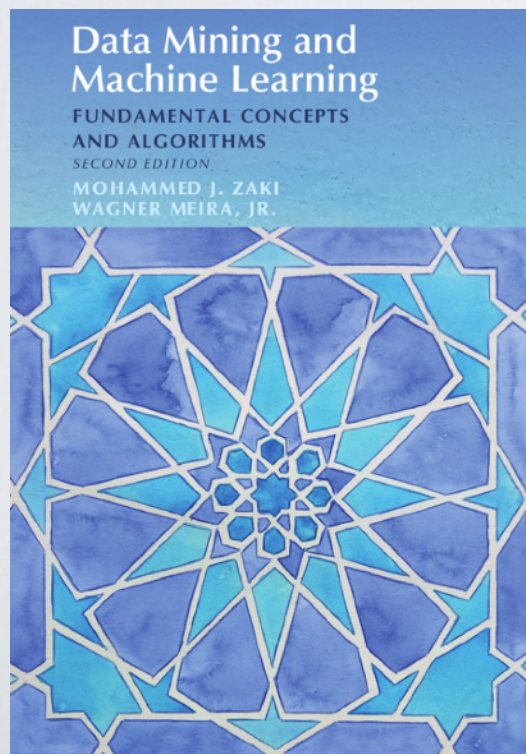- Maybe professional internships

# WHO ARE YOU?

# CLASS OVERVIEW

- Class with me: lecture + practical

- Three other lecturers

- Details on the lecture page:
  ‣ Class page: https://cazabetremy.fr/Teaching/DISS/ML.html
    - All contents: slides, TP, data, corrections…

- Exam:
  ‣ Short projects during semester: 50%
  ‣ Final Exam: 50%

# CLASS OVERVIEW

- Data description, preparation, etc.

- Unsupervised ML (beyond k-means)

- Supervised ML (beyond linear regression)

- Network Data Mining

- Deep Neural Networks

- Large Language Models

# THIS CLASS

- This class is based on:
  - ‣ Countless Wikipedia and blogs (use them too!)

- Some books
  - ‣ Borrow at my office

# DEFINITION

- Machine learning(ML) involves computers discovering how they can perform tasks without being explicitly programmed to do so. It involves computers learning from data provided so that they carry out certain tasks. It is a subset of Artificial Intelligence.

- [https://en.wikipedia.org/wiki/Machine_learning]

# DESCRIBING A DATASET

# TABULAR DATASET
## (For now)

**Tabular Data**

columns = attributes for those observations

| Player | Minutes | Points | Rebounds | Assists |
|--------|---------|--------|----------|---------|
| A | 41 | 20 | 6 | 5 |
| B | 30 | 29 | 7 | 6 |
| C | 22 | 7 | 7 | 2 |
| D | 26 | 3 | 3 | 9 |
| E | 20 | 19 | 8 | 0 |
| F | 9 | 6 | 14 | 14 |
| G | 14 | 22 | 8 | 3 |
| I | 22 | 36 | 0 | 9 |
| J | 34 | 8 | 1 | 3 |

Rows = observations

# OTHER TYPES

- Real Data can have many other forms
  ‣ Textual
  ‣ Relational (networks)
  ‣ Complex objects (picture, video, software…)

# TABULAR DATASET

- Size of the dataset
  - Number of observations
  - Number of variables

- Very large dataset?
  - =>Specific tools (Spark, Polars, etc.)

- Small dataset with many features?
  - Statistical tests, variable selection, etc.

# NATURE OF VARIABLES

# DATA TYPES

- Nominal/Categorical:
  ‣ From "names". No order between possible values
  ‣ Color, Gender, Animal, Brand, etc. (Numbers:Participant ID, class…)

- Numerical/Ordered:
  ‣ Interval
  ‣ Ratio

# NUMERICAL

- Ratio
  - ‣ Numerical values, all operations are valid
  - ‣ Height, Duration, Revenue…

- Interval
  - ‣ Numeric values, <u>difference</u> is meaningful
  - ‣ T°: 30°-20° = 15°-5°, But 30° 2*15°
  - ‣ 2022-2020=1789-1787, but 1011 2022/2
  - ‣ =>0 is not a meaningful value, is arbitrary
  - ‣ =>**Forbidden** to apply a log transformation
    - - Log convert sums into multiplications (e.g., +1 becomes twice as large)

# TRAPS

- Latitude and Longitude

- Hours expressed between 0 and 12/24, day of month, etc.
  - ‣ Convert in time since beginning of dataset ?

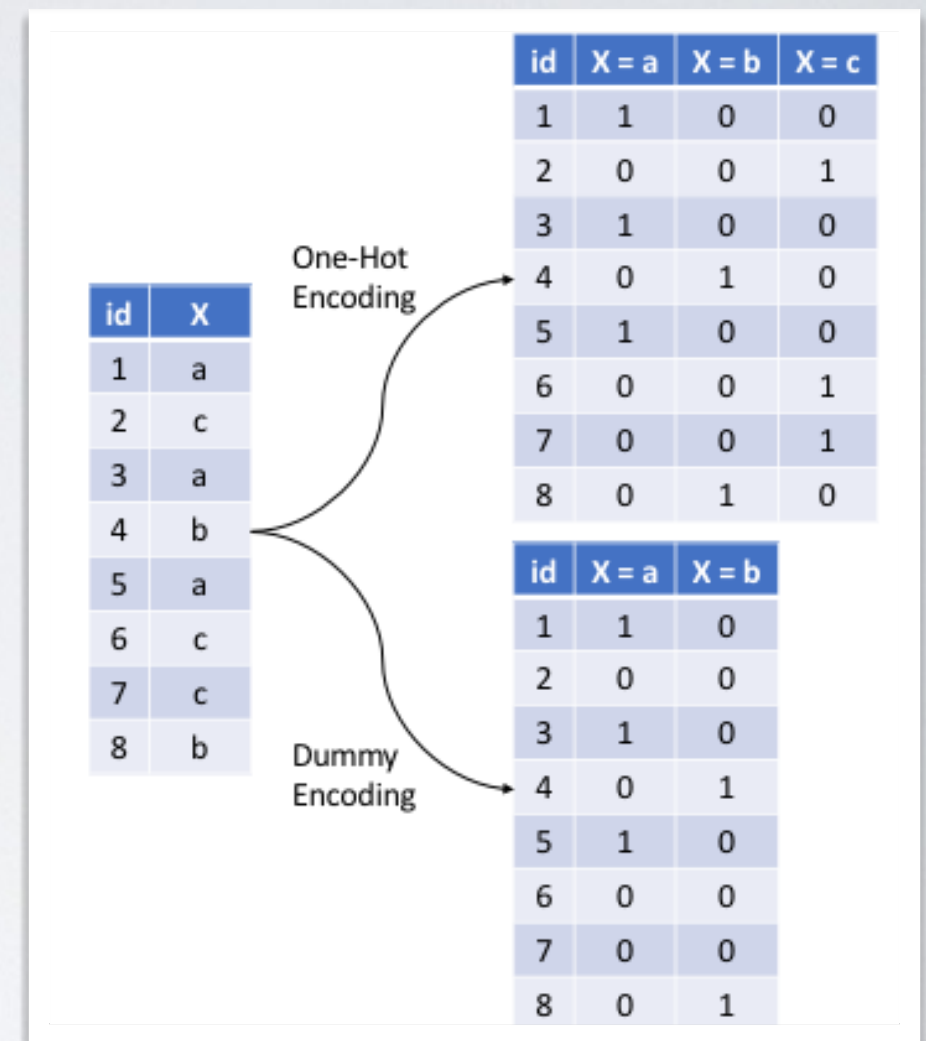- => Space and Time often handled with specific ML methods

# WHAT TO DO ?



- Nominal =>
  - ‣ One hot encoding
  - ‣ Also called
    - - Dummy encoding
    - - Indicator variables
    - - Binary vector encoring

- WARNING
  - ‣ Keeping numerical values for nominal variables is **WRONG**!!!

# MISSING VALUES

- Real-life datasets are full of missing values
  ‣ Impossible data: *fur color* for a sphinx cat
  ‣ More generally, failure to obtain them

- Few methods can deal with missing values
  ‣ =>Imputation
    - Naive: fill with average value
    - Use ML to fill-in missing values (other problems, introduce biases…)
    - Large literature, no good solution

# DATA QUALITY

- Data coming from the real world is often incorrect
  - ‣ Malfunctioning sensors (T°, speed…)
  - ‣ Human error or falsification (e.g., entered 100 instead of 1.00)
  - ‣ Undocumented change (e.g., Bicycle sharing station was moved…)

- Before applying a method blindly,
  - ‣ =>**check your data's quality!**
  - ‣ If the data is plausible, no simple solutions
  - ‣ Common
    - - Out-of-range values (e.g., a person's weight is negative or above 1000kg…)
    - - Zeros. (Weight of the person is 0. But in many cases, zero is possible too…)
    - - Variant: 01/01/1970…

# DESCRIBING A VARIABLE

# DESCRIBING VALUES

- Mean / Average
  ‣ Be careful, not necessarily representative !

- Median
  ‣ Be careful, not necessarily representative !

- Mode
  ‣ Not necessarily representative

- Min/Max
  ‣ …

# VARIANCE

- Variance:
  - ‣ Expectation of the <u>squared</u> deviation of a random variable from its mean

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

Also expressed as average squared distance
between all elements

$$\sigma^2 = \frac{1}{N^2} \sum_{i<j} \left(x_i - x_j\right)^2$$

# STANDARD DEVIATION

- Squared root of the Variance

$$\sigma = \sqrt{\sigma^2}$$

# ABSOLUTE DEVIATION

- MAD (Mean Absolute Deviation)
  - ‣ Deviation from mean or from median
  - ‣ (Variant: Median Absolute Deviation)

$$\text{MAD}(X) = \frac{1}{n} \sum_{i=1}^{n} \left| x_i - \bar{x} \right|$$

- So why are we using the Standard Deviation again?
  - ‣ The <u>mean</u> minimizes the expected squared distance
  - ‣ The <u>median</u> minimizes the MAD
  - ‣ Leads naturally to least square regression and PCA… see later.

# STATISTICAL DISTRIBUTIONS

# STDIV AND NORMAL DISTRIBUTION

Median     Mean

# DISTRIBUTION

- ## What is a distribution?
  - ‣ A description of the frequency of occurence of items
  - ‣ A generative function describing the probability to observe any of the possible events
  - ‣ Discrete or continuous
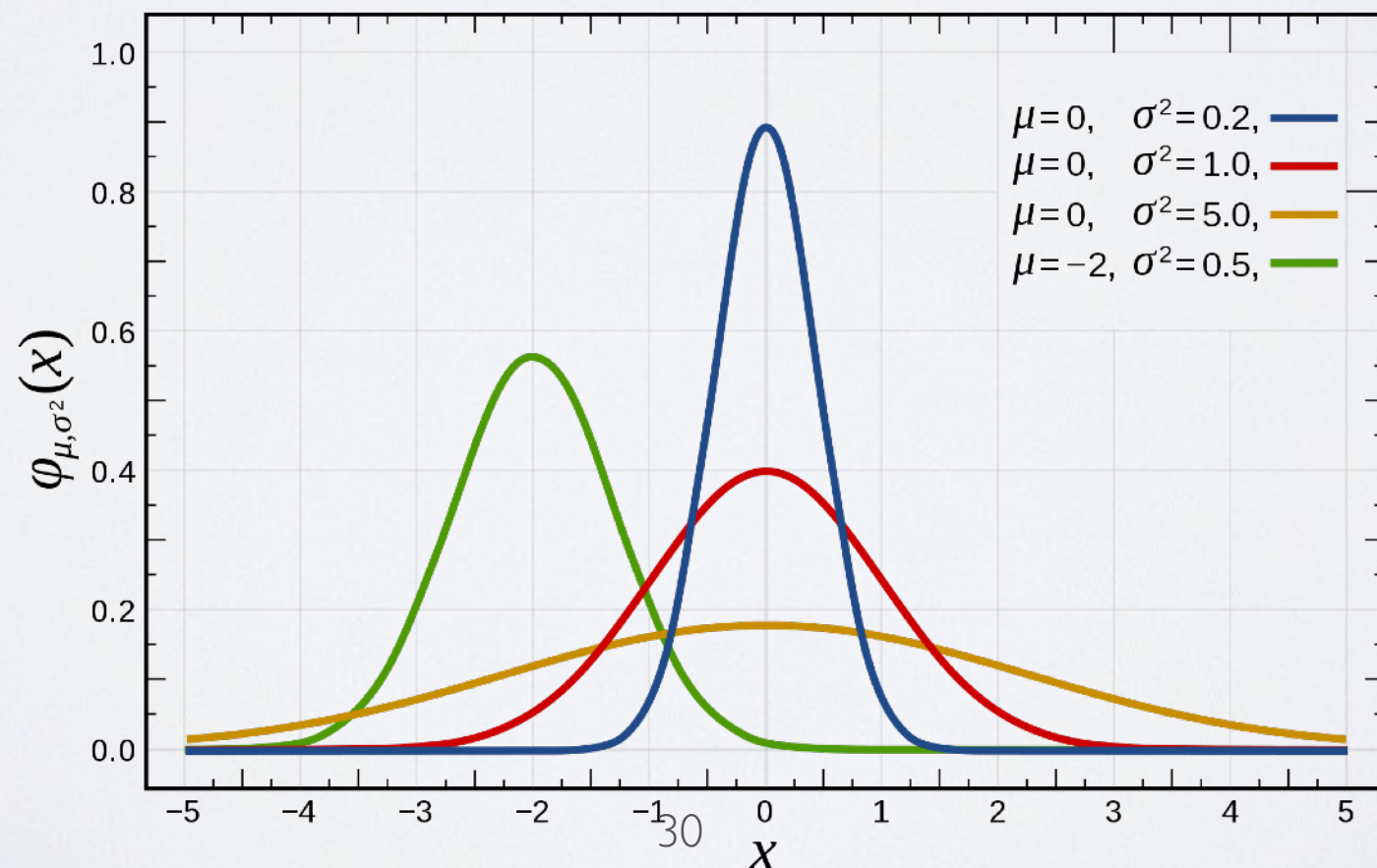


**Continuous Distribution**

# DISTRIBUTION (DISCRETE)



- =>25 observations in the interval (13,17]

- Raw values for a sample,

- or fraction
  ‣ 0.25
  ‣ 25%
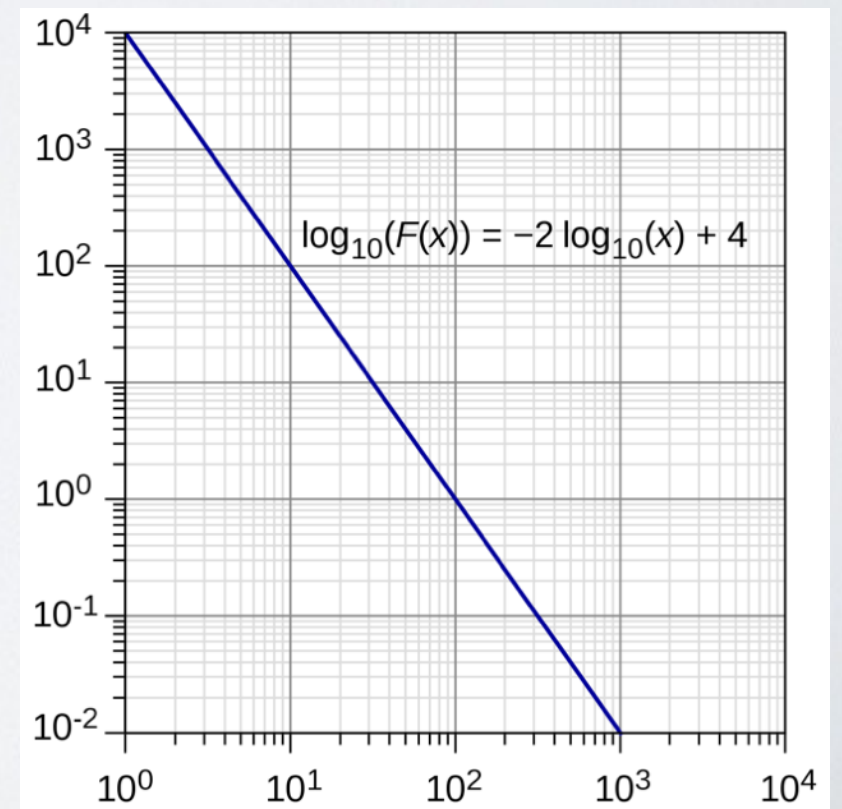  ‣ =>Sum to 1. Must be inferior to 1 for any value

# THEORETICAL DISTRIBUTIONS

- Normal distribution
  - Many real variables follow it approximately (height, weight, price of a given product in various locations…
  - Random variations around a well-defined mean
  - Central limit theorem: <u>average</u> of many samples of a random variable converges to a normal distribution
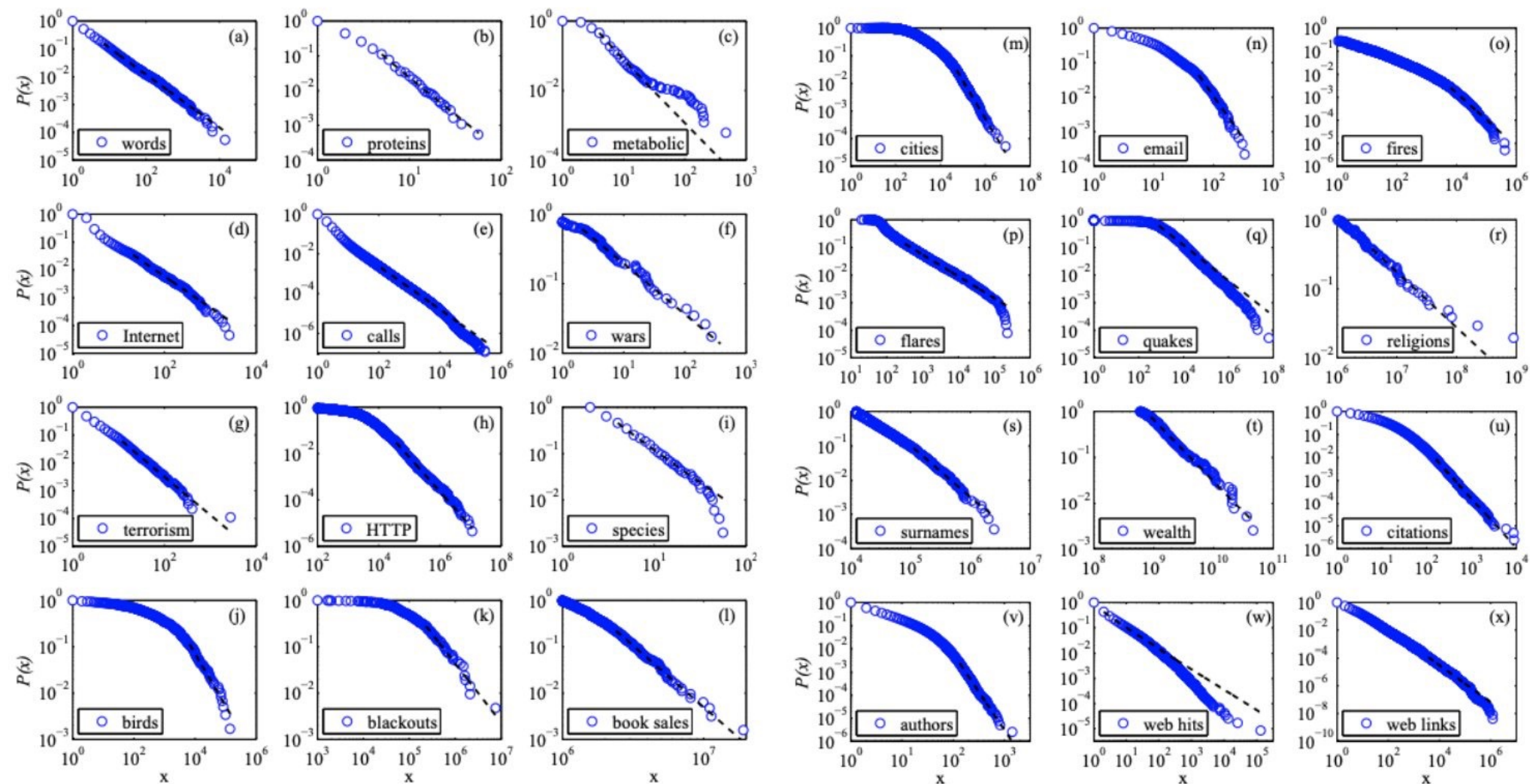
# THEORETICAL DISTRIBUTIONS

- Power Law distribution
  - ‣ A relative change in one quantity results in a proportional relative change in the other quantity, independent of the initial size of those quantities: one quantity varies as a power of another.
    - e.g., earthquakes 10 times more powerful are $x$ times less frequent.
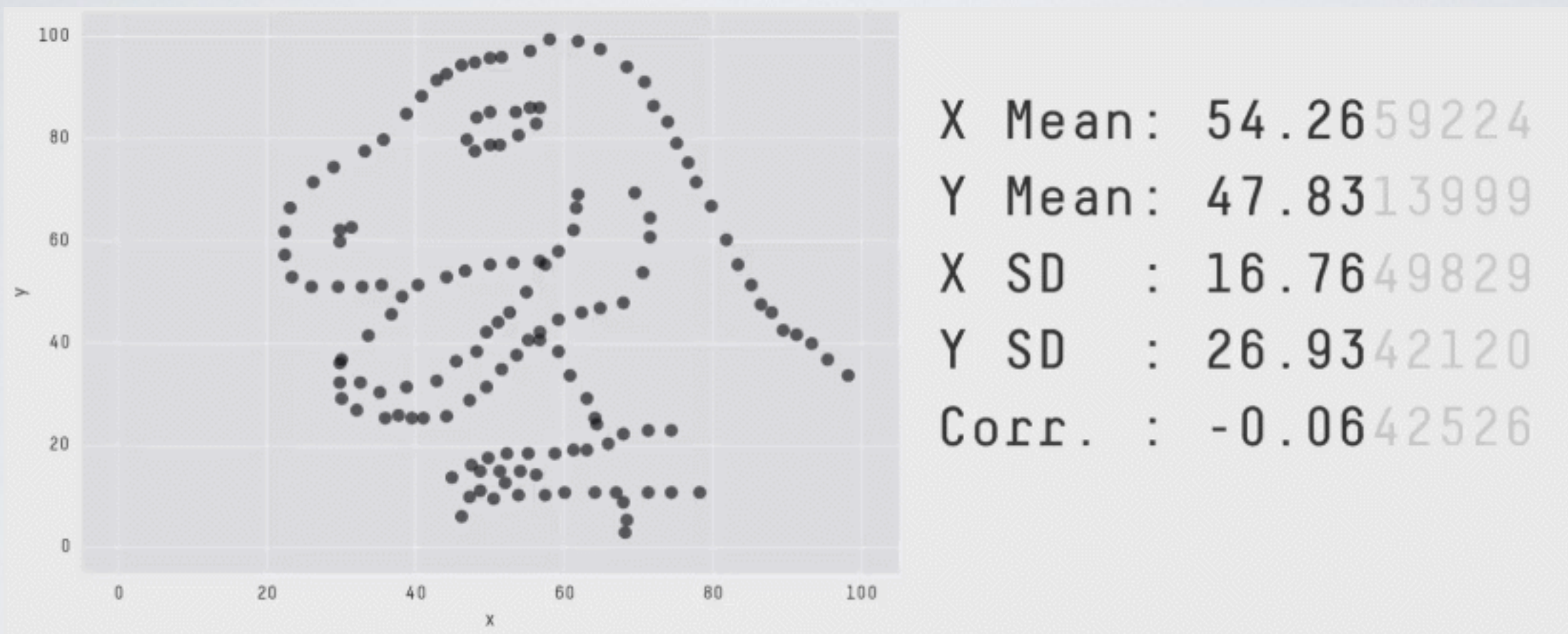    - e.g., cities 10 times bigger are $x$ time less frequent

$$\log_{10}(F(x)) = -2 \log_{10}(x) + 4$$

# THEORETICAL DISTRIBUTIONS

- Power Law distribution

# DESCRIPTIVE STATISTICS

X Mean: 54.2659224
Y Mean: 47.8313999
X SD  : 16.7649829
Y SD  : 26.9342120
Corr. : -0.0642526

The datasaurus

https://github.com/jumpingrivers/datasauRus

# DESCRIPTIVE STATISTICS

- My advice:
  - ‣ Plot the distribution.
  - ‣ Don't assume a theoretical distribution
  - ‣ Don't believe single-number statistics.

# VARIABLE INTERACTIONS

# COVARIANCE MATRIX

Covariance Matrix Formula

$$\begin{bmatrix} \text{Var}(x_1) & ...... & \text{Cov}(x_n, x_1) \\ \vdots & \ddots & \vdots \\ \text{Cov}(x_n, x_1) & ...... & \text{Var}(x_n) \end{bmatrix}$$

- Covariance matrix $\mathbf{K}$
  - ‣ Extension of Variance to multivariate data
  - ‣ $\text{Var}(X) = \text{E}\left[(X - \mu)^2\right]$
  - ‣ $\text{cov}(\mathbf{X}, \mathbf{Y}) = \text{K}_{\mathbf{XY}} = \text{E}\left[(\mathbf{X} - \text{E}[\mathbf{X}])(\mathbf{Y} - \text{E}[\mathbf{Y}])^{\text{T}}\right]$
    - How much observation X differs from the mean ? And Y ?
    - Multiply the respective divergences of X and of Y for each item
    - Take the average
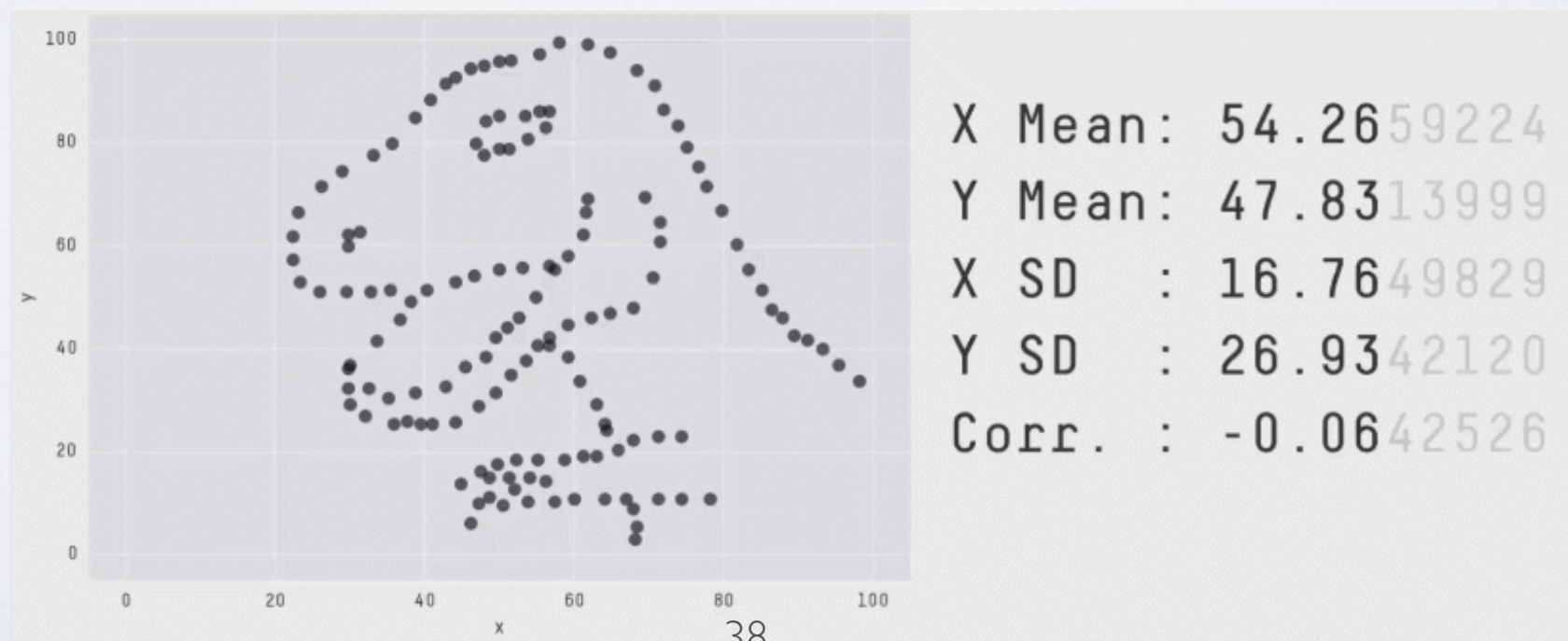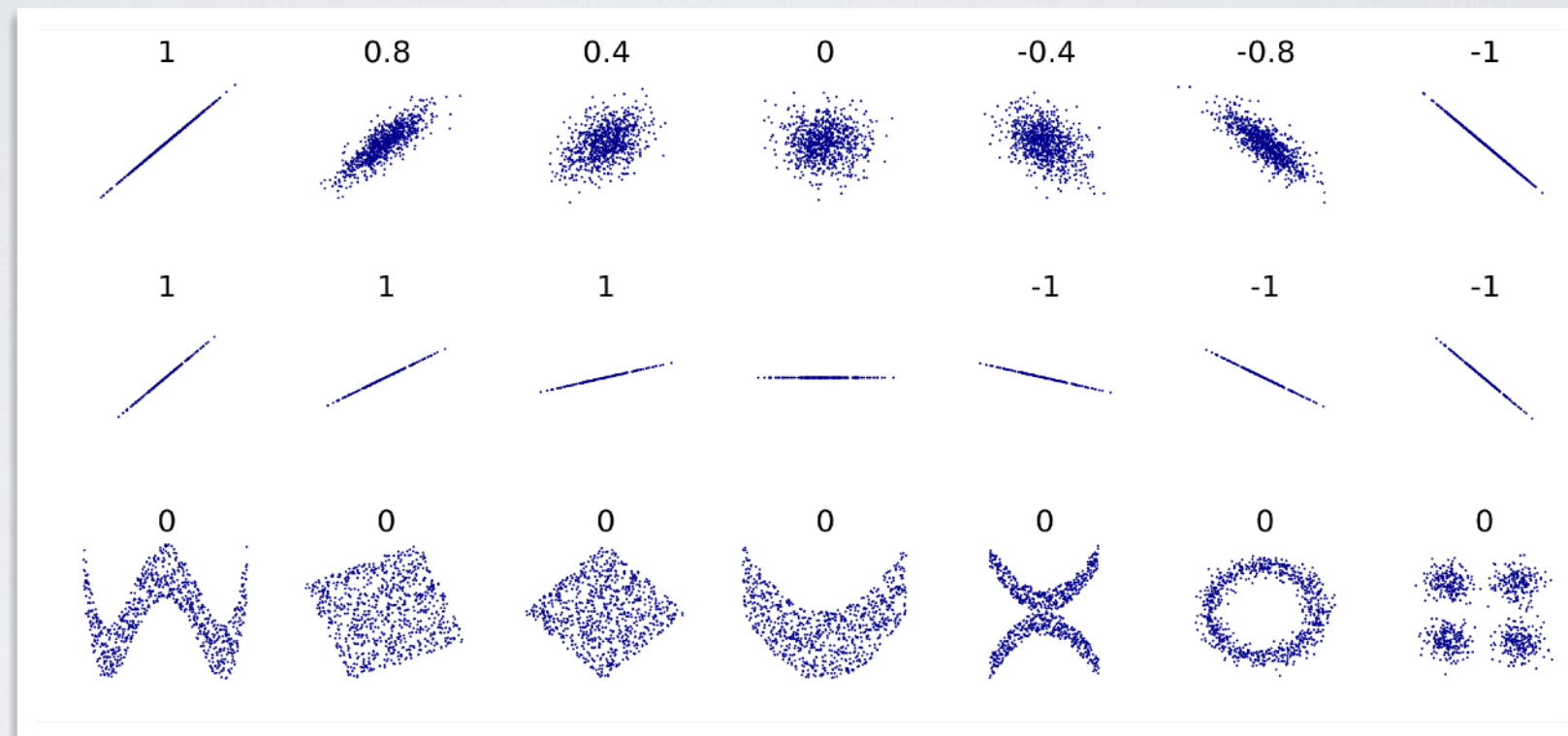  - ‣ $=> \text{cov}(\mathbf{X}, \mathbf{X}) = \text{Var}(\mathbf{X})$

- Covariance is hardly interpretable by itself.
  - ‣ If >0, divergences tend to be in the same direction
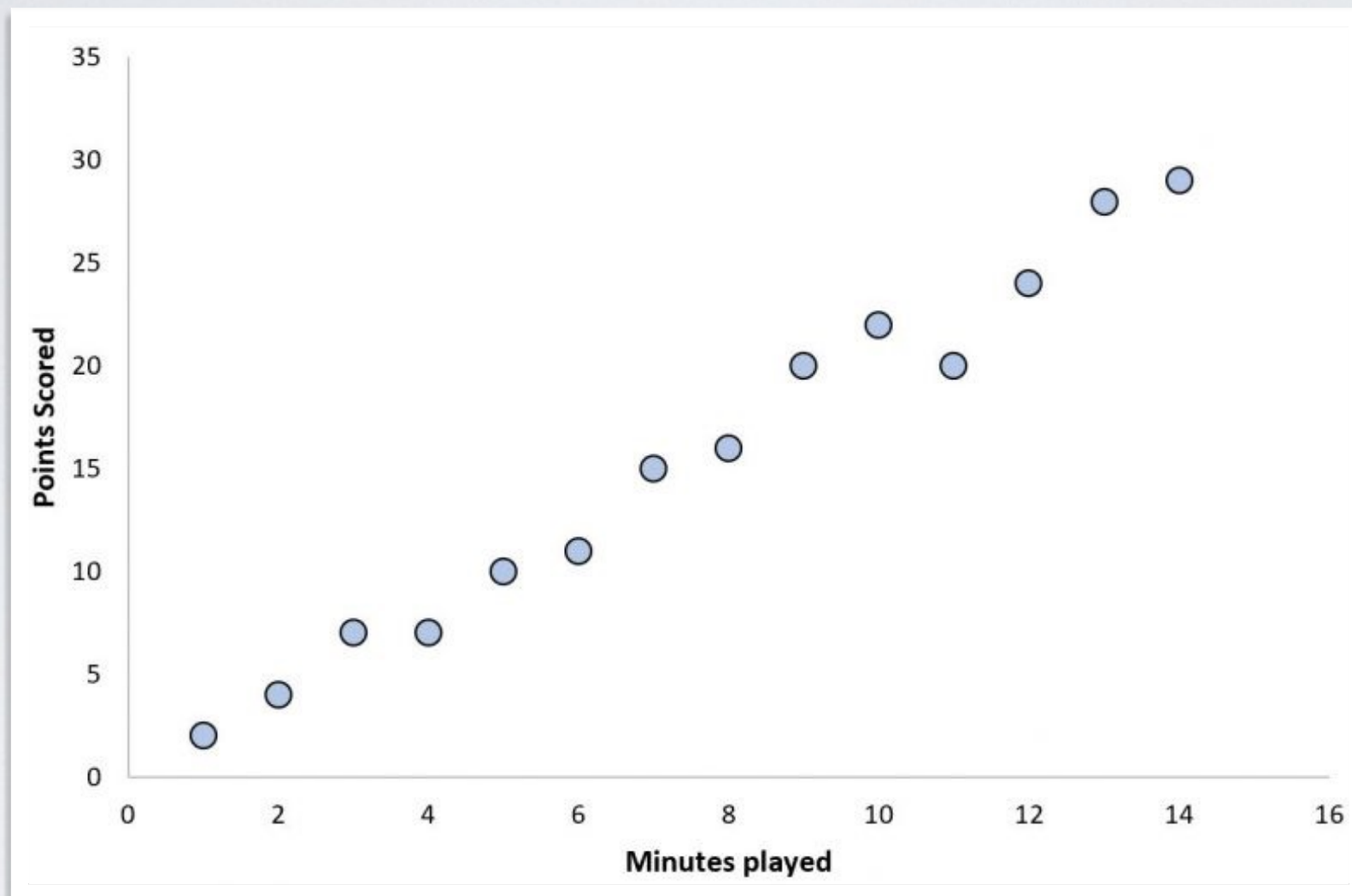  - ‣ Normalize it to obtain the "correlation coefficient"

# CORRELATION COEFFICIENT

- Pearson correlation coefficient : $\rho_{X,Y} = \dfrac{\mathrm{cov}(X, Y)}{\sigma_X \sigma_Y}$

  ‣ Normalize the Covariance by the Standard deviation.

  ‣ Independent from magnitude, i.e., no need to have normalized data

  ‣ Value in -1, +1.

    - +1 means a perfect positive linear correlation, i.e., X=aY

    - -1 a negative one, i.e., X=-bY

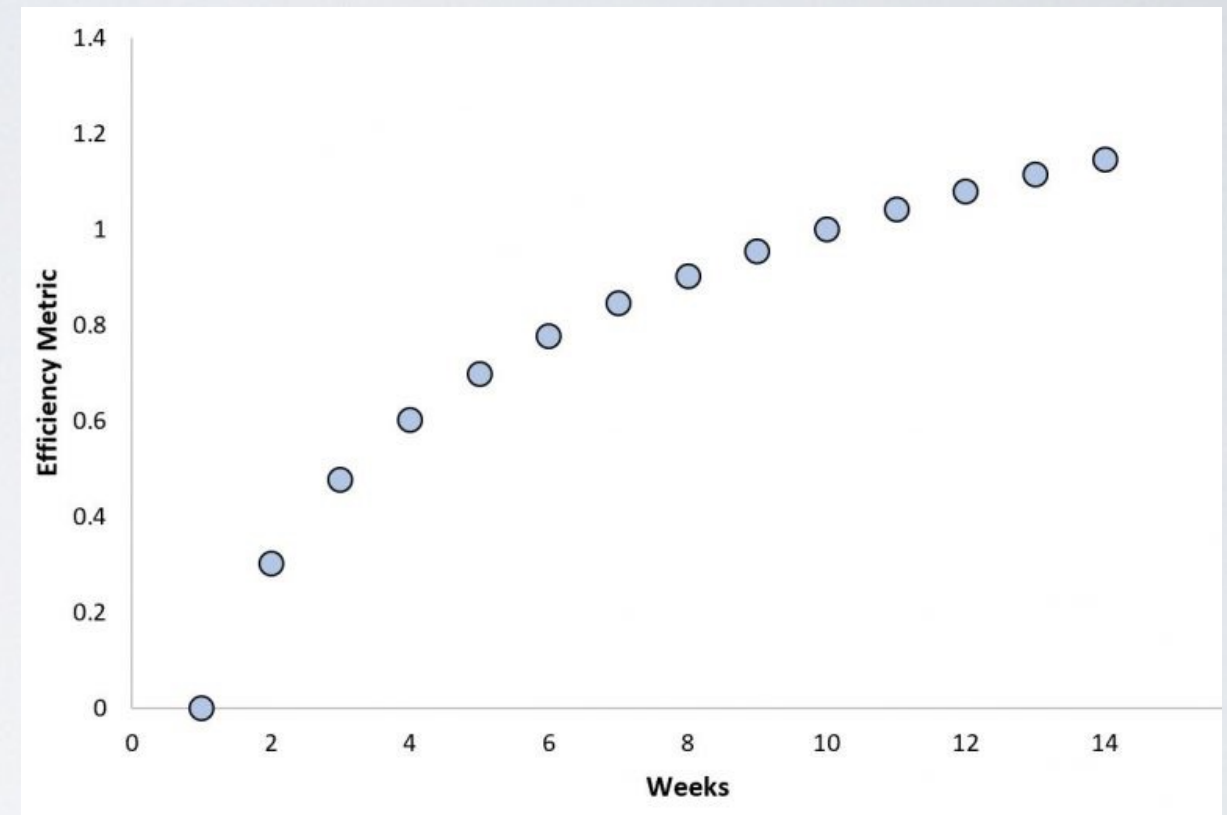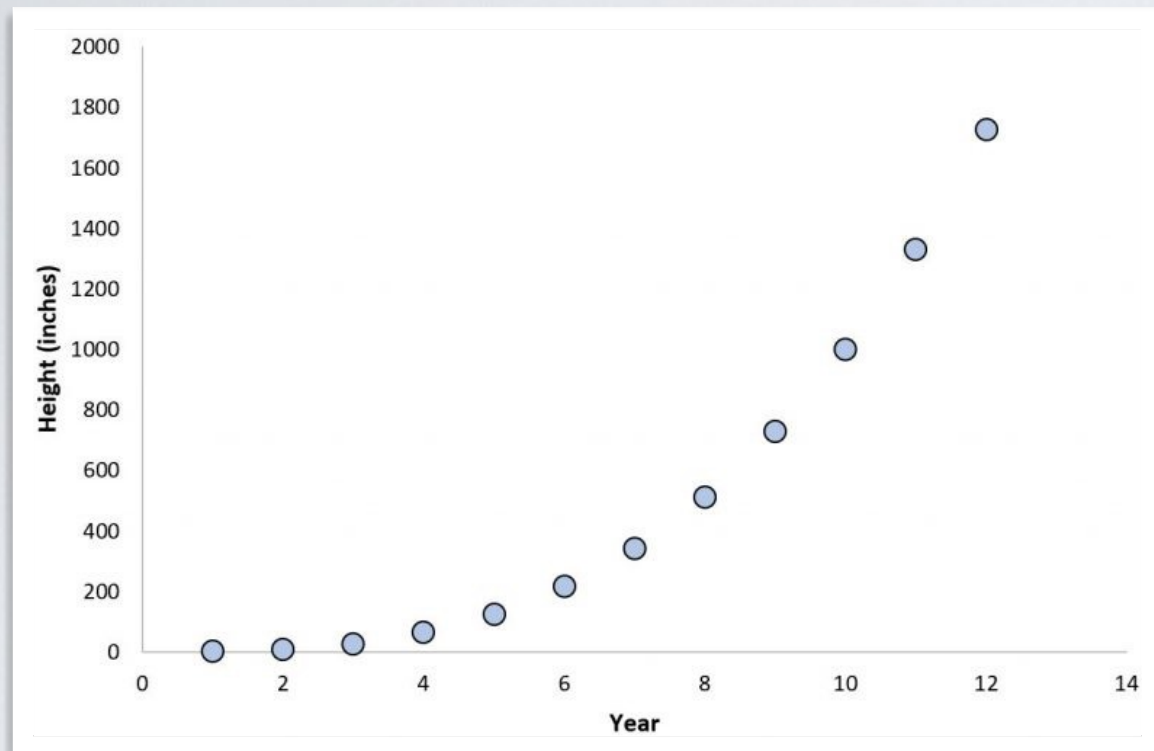  ‣ 0 can mean many different things
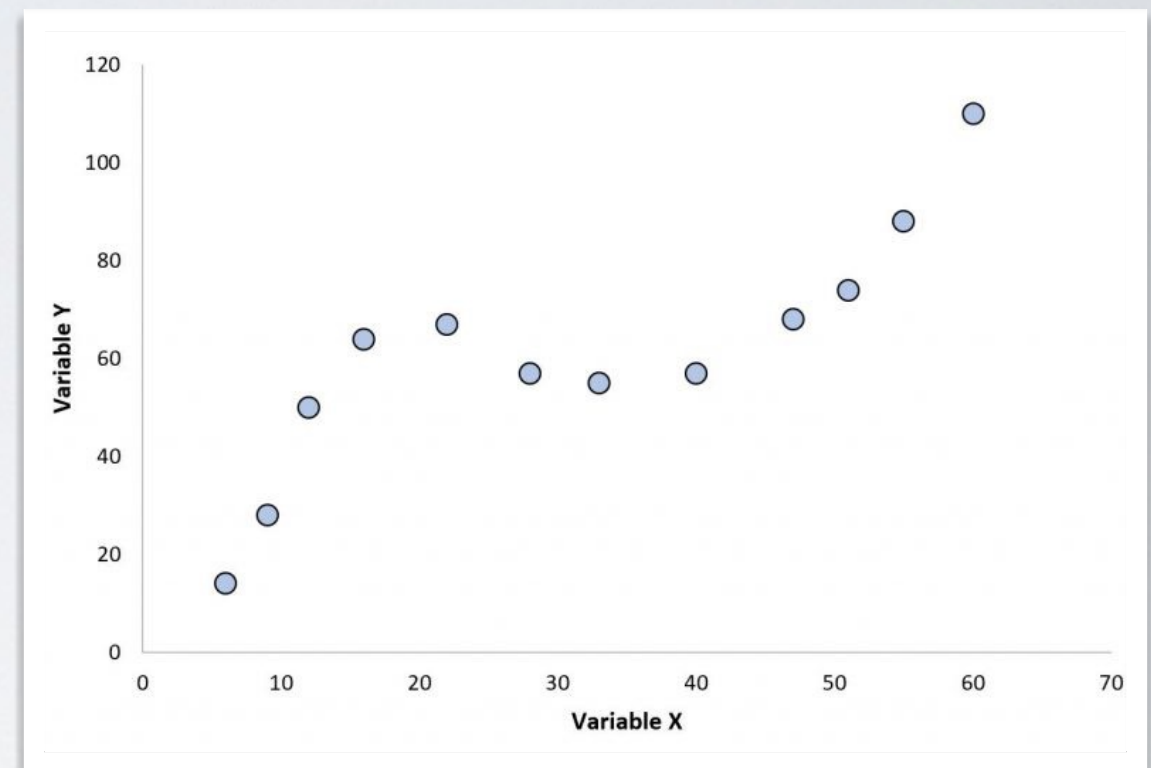
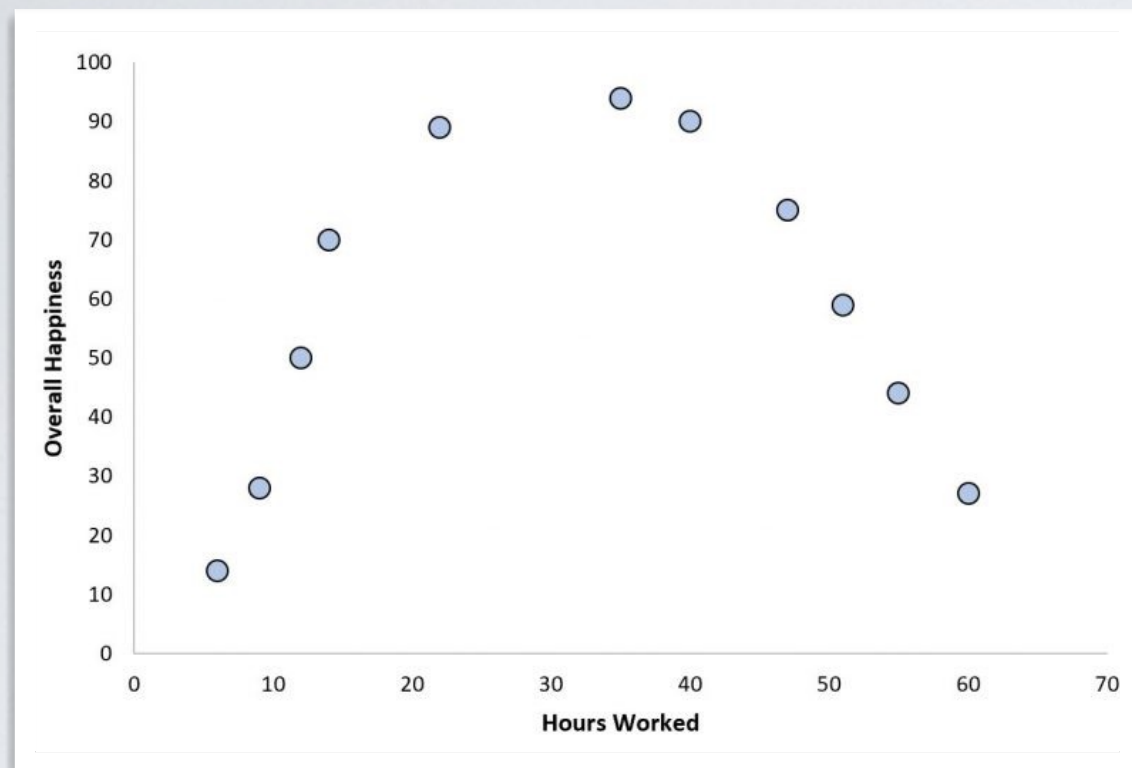# CORRELATION COEFFICIENT

# NONLINEAR RELATIONSHIPS
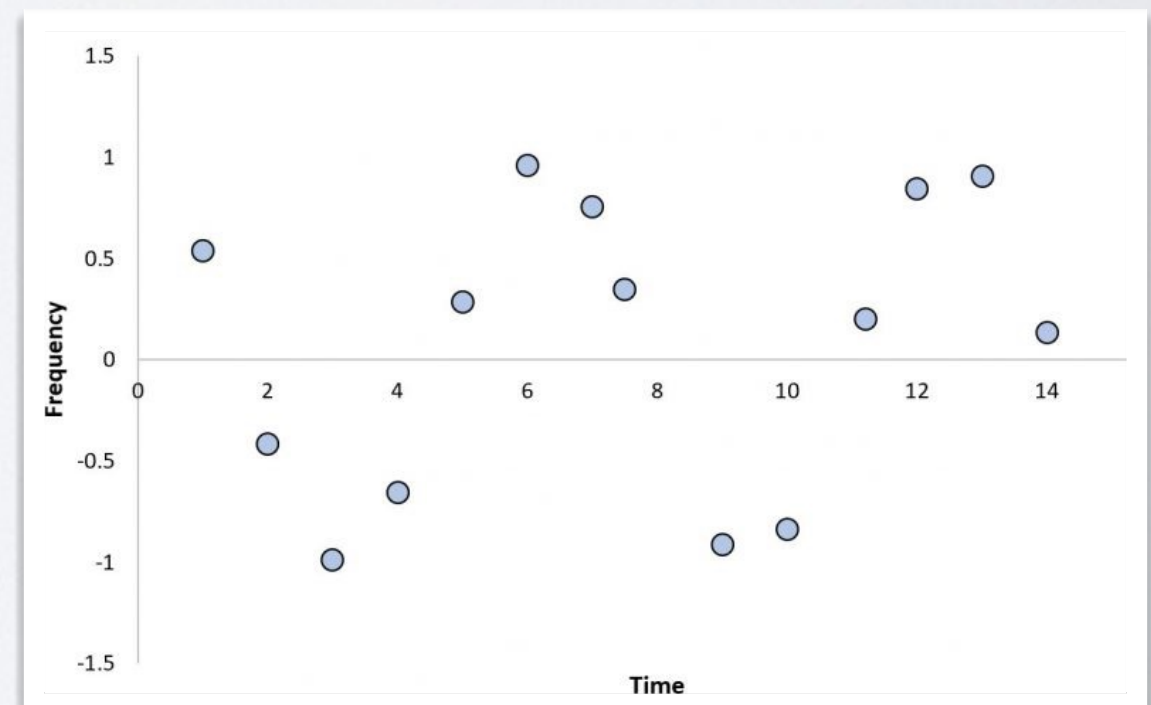


Linear relationship
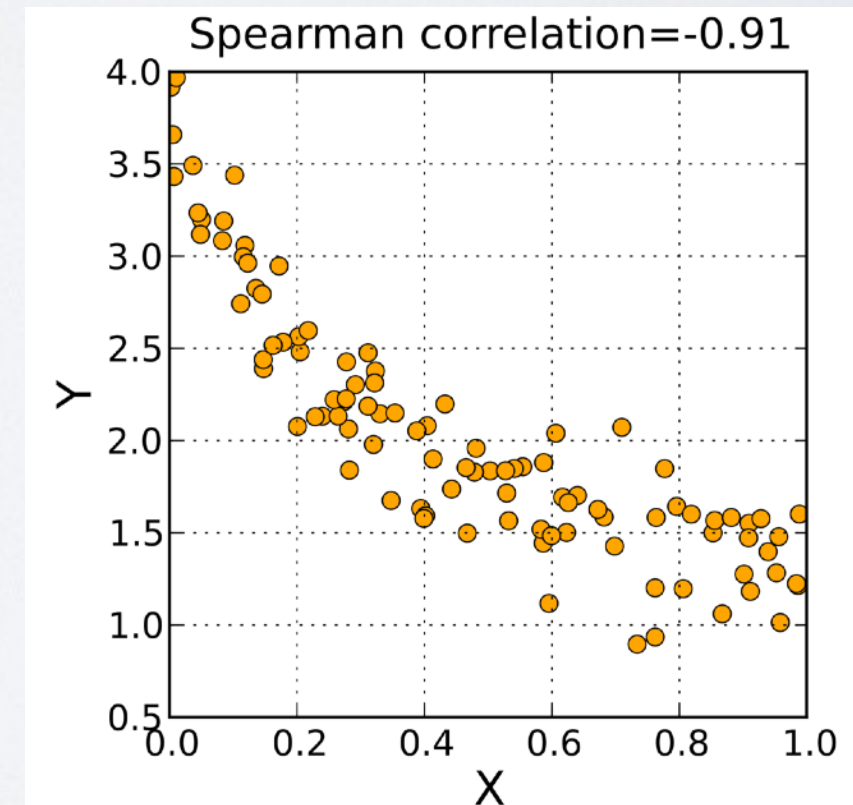Y=a+bX+e

# NONLINEAR RELATIONSHIPS



Monotonous, non-linear
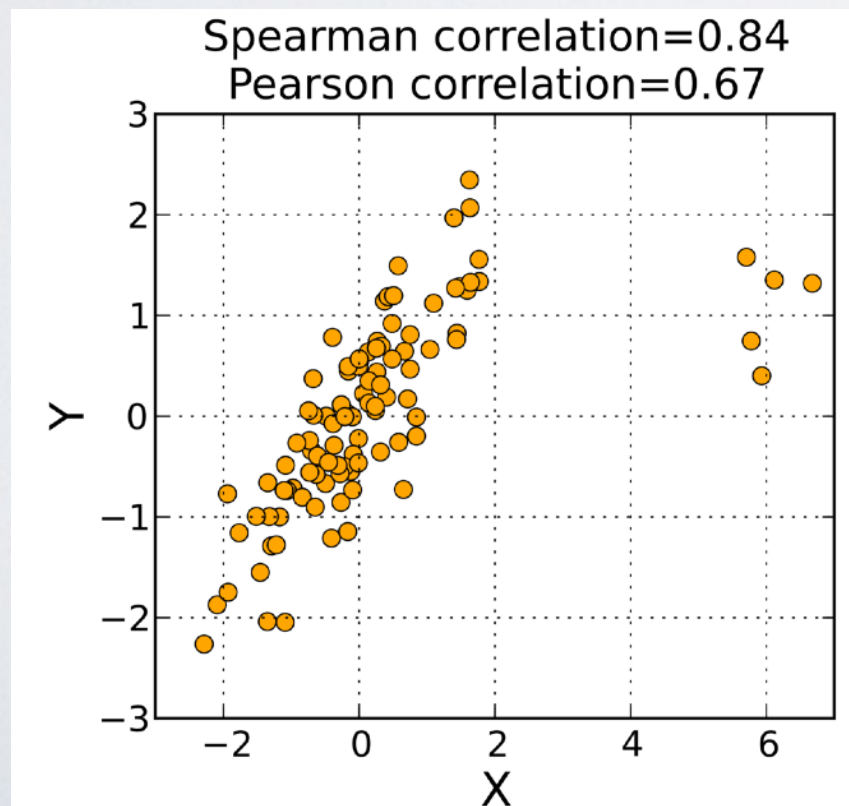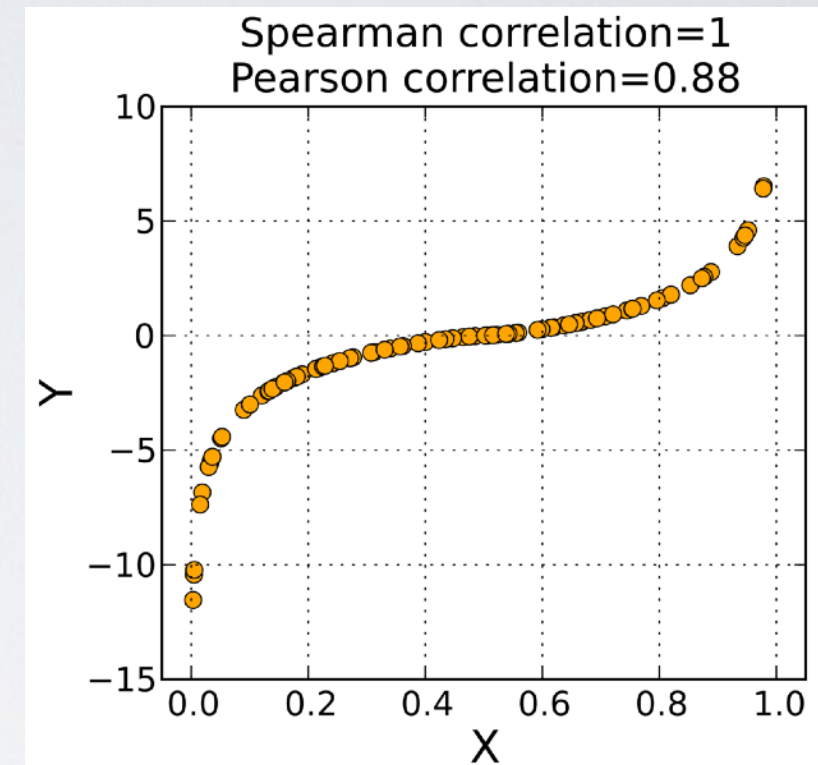
# NONLINEAR RELATIONSHIPS
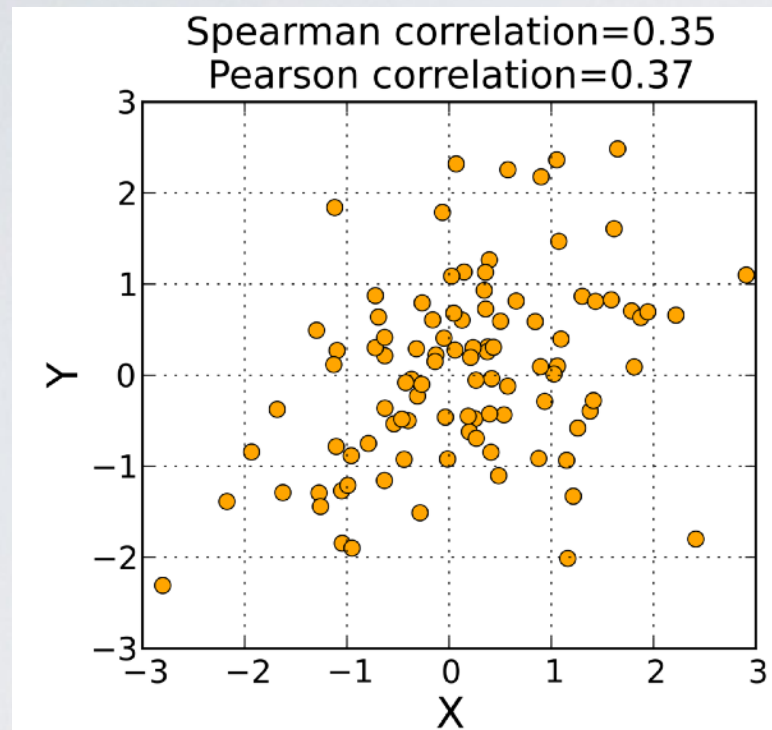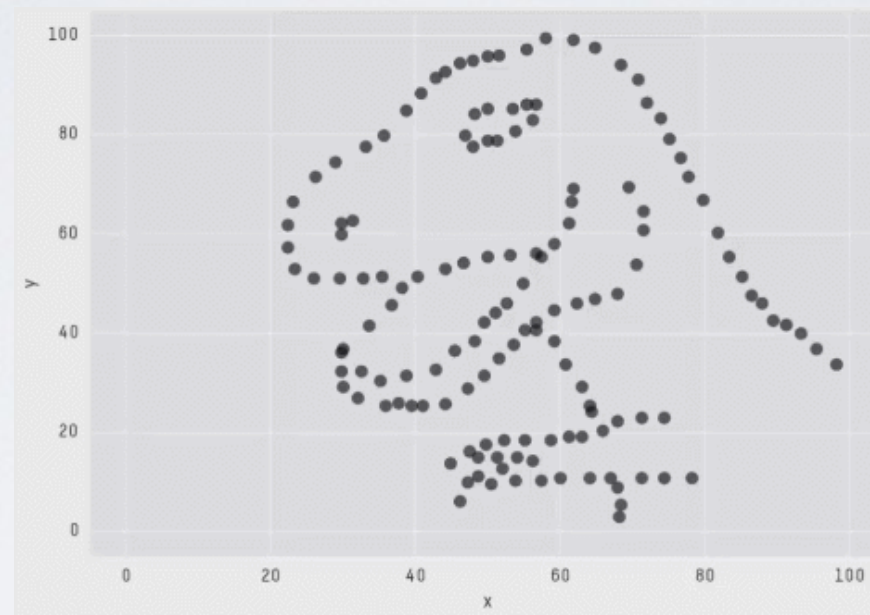


Non-monotonous,
Non-linear

# SPEARMAN'S CORRELATION

- Spearman's **rank** correlation coefficient

- Assesses how well the relationship between two variables can be described using a monotonic function
  - ‣ Not assuming a linear relation

- Pearson correlation coefficient between the rank variables
  - ‣ $r_s = \rho_{R(X),R(Y)} = \dfrac{\text{cov}(R(X), R(Y))}{\sigma_{R(X)}\sigma_{R(Y)}}$

# SPEARMAN'S CORRELATION



Spearman correlation=0.35
Pearson correlation=0.37

Spearman correlation=1
Pearson correlation=0.88

Spearman correlation=0.84
Pearson correlation=0.67

Spearman correlation=-0.91

# DESCRIPTIVE STATISTICS



X Mean: 54.2659224
Y Mean: 47.8313999
X SD  : 16.7649829
Y SD  : 26.9342120
Corr. : -0.0642526

- My advice:
  ‣ Plot the relations
  ‣ Don't believe single-number statistics. Never ever.

# WARNING

- Correlation is not causation!!!
  ‣ "People having a Ferrari live longer in average"

- Confounding variable:
  ‣ an unobserved variable that affects both the cause being studied (Ferrari) and the effect observed (life expectation)
  ‣ =>The main problem of any study. It is impossible (apart from strictly controlled experiments) to avoid this problem.
  ‣ => **Be careful** when drawing conclusions from data
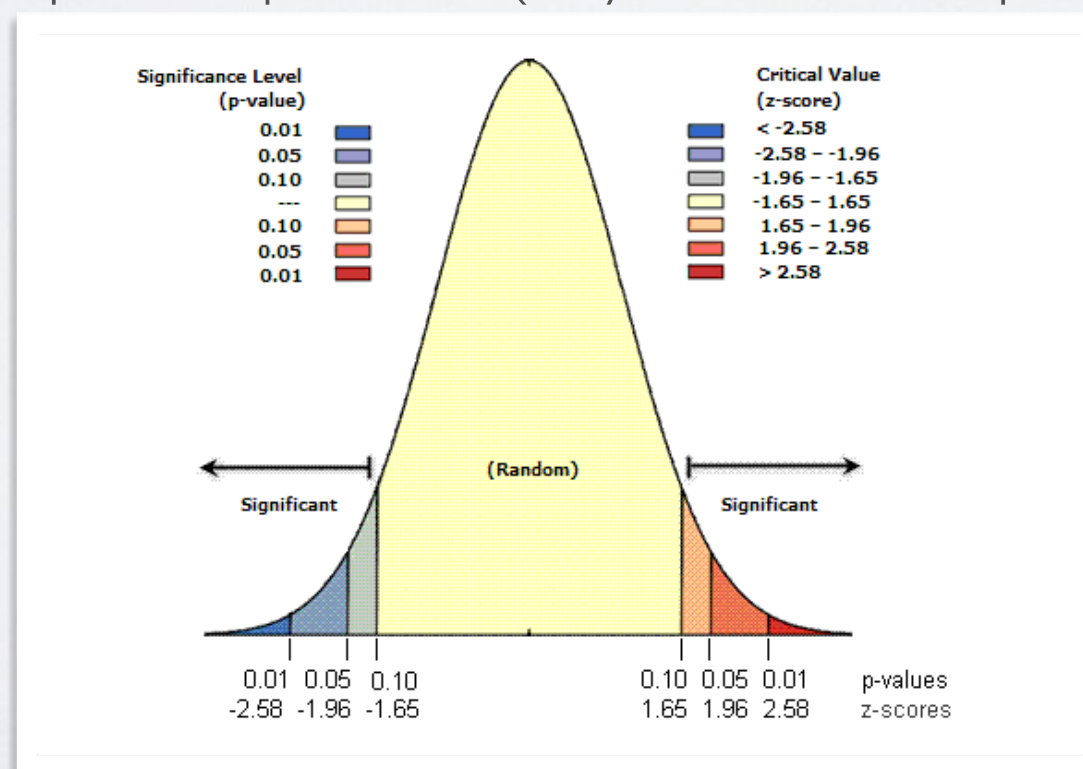
# STATISTICAL SIGNIFICANCE

# WHY?

- You observe a correlation between two variables

- How to be certain that this correlation is real, and not just due to random chance?
  - ‣ Imagine that you toss a coin 10 times and get 7 heads.
    - - Does it mean that the coin is biased, or is it expected by chance?
  - ‣ For correlations: In your dataset, tall people also tend to be wealthy people. Is it true, or just an effect of chance in your dataset?

# P-VALUE

- P-value: probability of observing a value as exceptional as the one you actually observed.
    ‣ [0,1]

- Can be computed analytically, or by simulations

# ANALYTICAL P-VALUE

- Assuming that a coin toss should follow a Binomial distribution, probability of observing 7 heads from 10 tosses?

- For correlations:
  - Assuming normal distributions of variables
  - Assuming bivariate normal relations between them
  - =>One can compute a p-value (beyond the scope of this class)
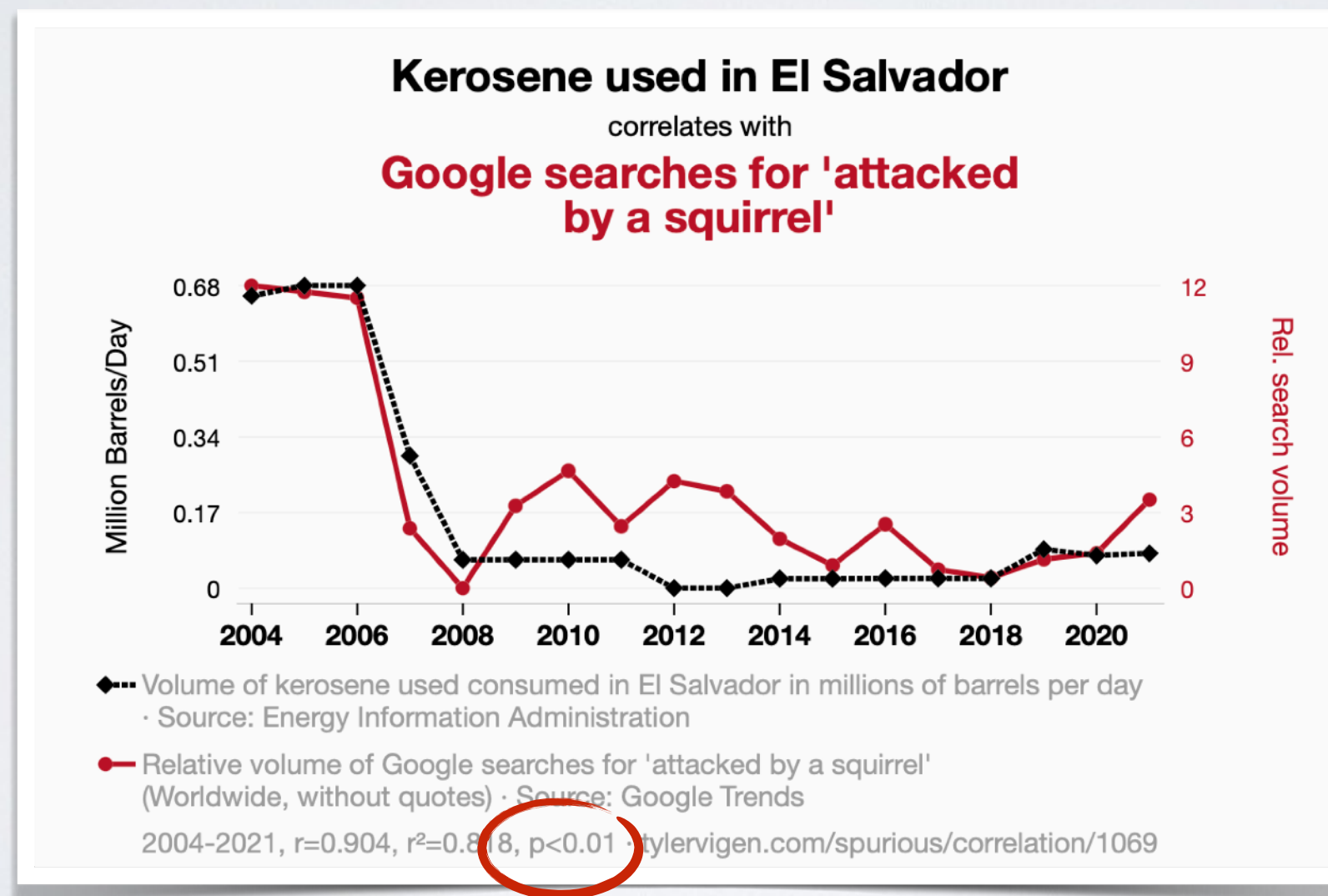
# SIMULATION-BASED P-VALUE

- You observed a correlation $c$ between variables $X$ and $Y$

- Model-based p-value
  - ‣ 1)Compute the distributions of $X$ and $Y$
  - ‣ 2)Repeat $n$ times
    - Simulate values for $X$ and $Y$.
    - Compute the correlation for each simulation
  - ‣ 4)Count how many times you observed a value of $c$ as exceptional as the true one

- Permutation based p-value
  - ‣ 1) Repeat $n$ times:
    - Shuffle the values of $Y$
    - Compute the correlation for each shuffling
  - ‣ 2)Count how many times you observed a value of $c$ as exceptional as the true one

# STATISTICAL TESTS

- Useful when you have **very little data** and that you **cannot obtain more**

- If you have large datasets, in general, these tests are useless
  ‣ No distribution is exactly normal
  ‣ No variables are exactlly independent
    - Having a cat and owning a SUV? Height of a person and their grades in high school? Etc
    - =>Any relation is "significant"
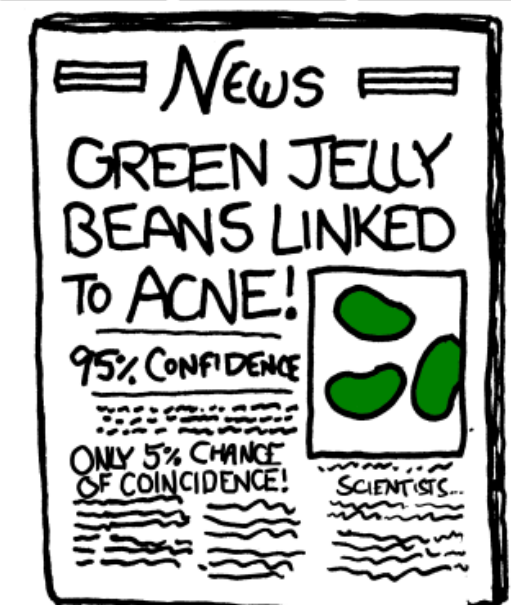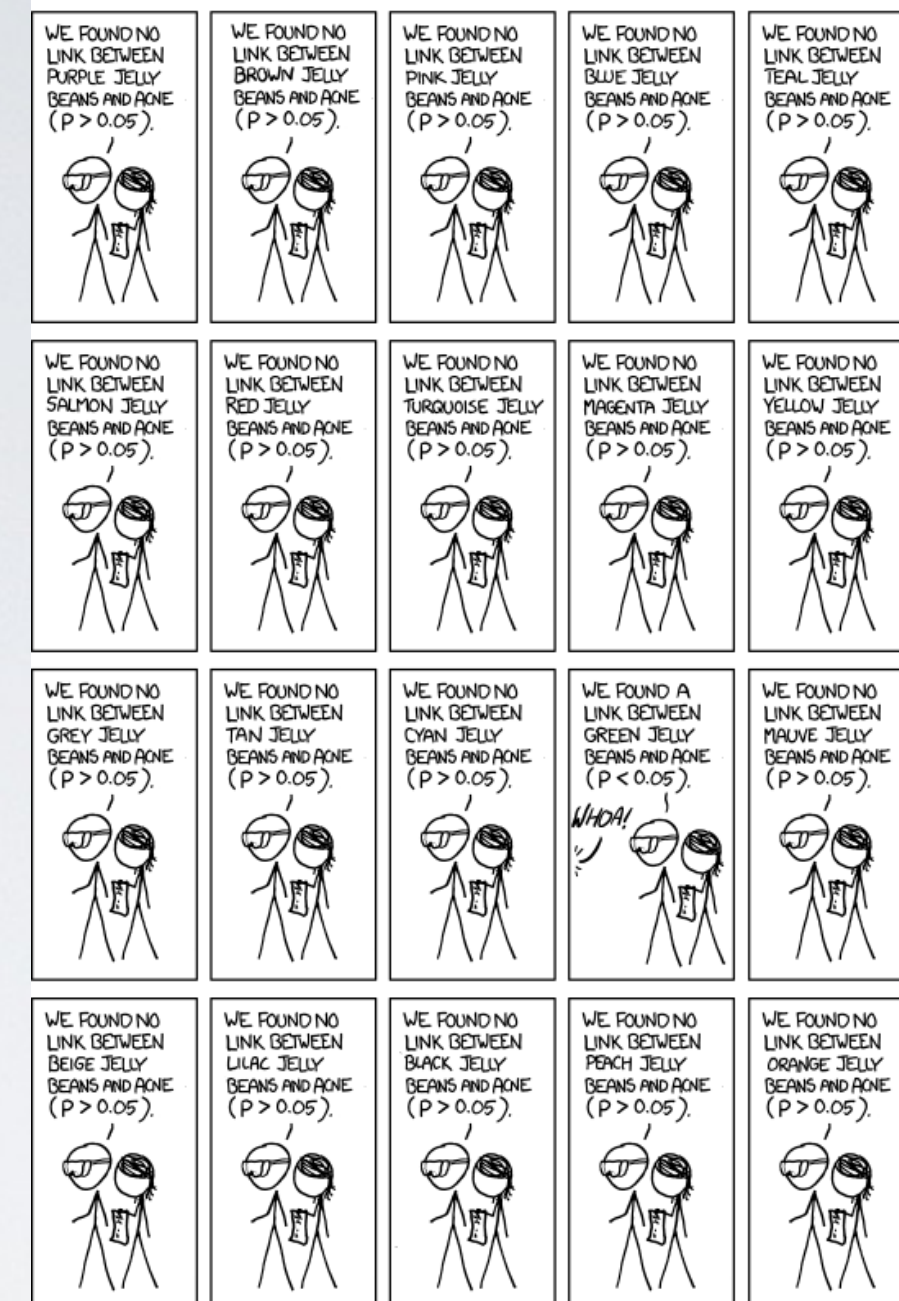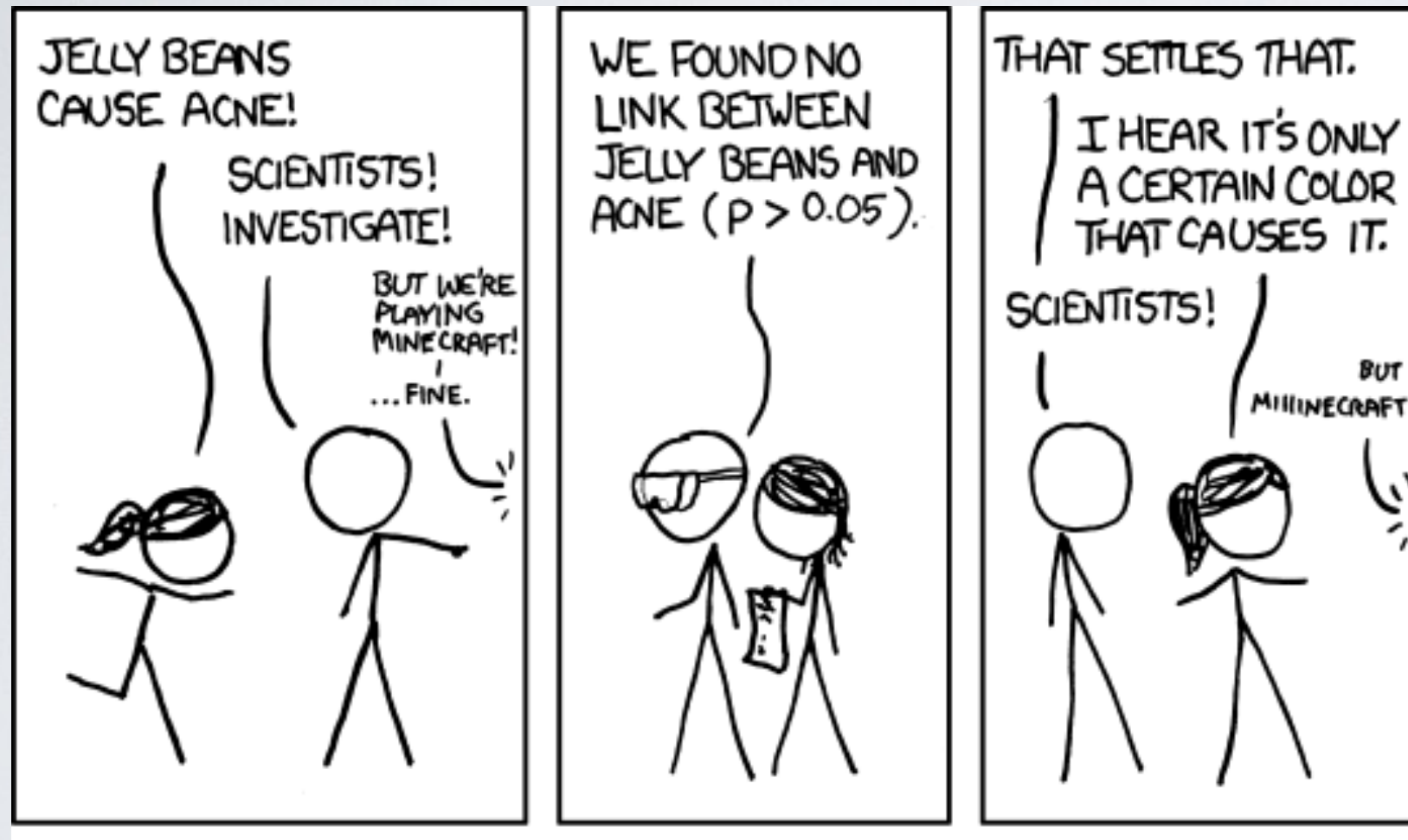
# SPURIOUS CORRELATIONS

https://www.tylervigen.com/spurious-correlations



**Kerosene used in El Salvador**
correlates with
**Google searches for 'attacked by a squirrel'**

Volume of kerosene used consumed in El Salvador in millions of barrels per day · Source: Energy Information Administration

Relative volume of Google searches for 'attacked by a squirrel' (Worldwide, without quotes) · Source: Google Trends

2004-2021, r=0.904, r²=0.818, p<0.01 · tylervigen.com/spurious/correlation/1069

# P-VALUES

# P-VALUES

# STATISTICAL SIGNIFICANCE

- My advice:
  ‣ Plot the data
  ‣ If the relation is not so obvious that you have no doubts, don't believe it
  ‣ Get more data :)

# SOME "GOLDEN RULES"

# SOME "GOLDEN RULES"

- In real life:
  - ‣ Your data does not follow a normal distribution. Nor a power law, nor any other theoretical distribution
  - ‣ Your features are always correlated
  - ‣ You always have non-linear relationships

# SOME "GOLDEN RULES"

- GIGO: Garbage in, Garbage out

# SOME "GOLDEN RULES"

- Real data is always garbage

# SOME "GOLDEN RULES"

- Get to know your data
  - ‣ Exploratory Analysis