

MACHINE LEARNING DATA - INTRODUCTION

WHO AM I

- Rémy Cazabet (remy.cazabet@univ-lyon1.fr)
- Class page: <http://cazabetremy.fr/Teaching/TIW/DAD.html>
- 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
- Interns application welcomed

WHO ARE YOU ?

DEFINITION

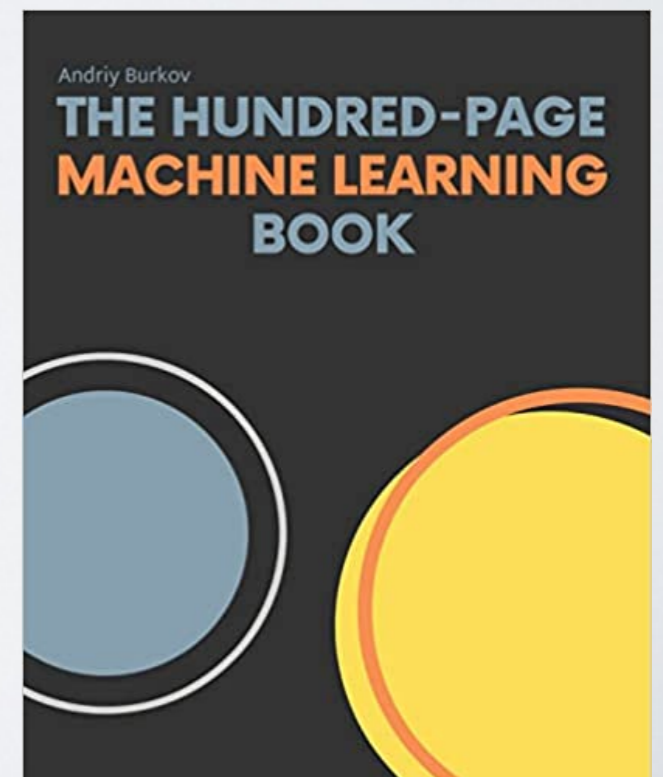
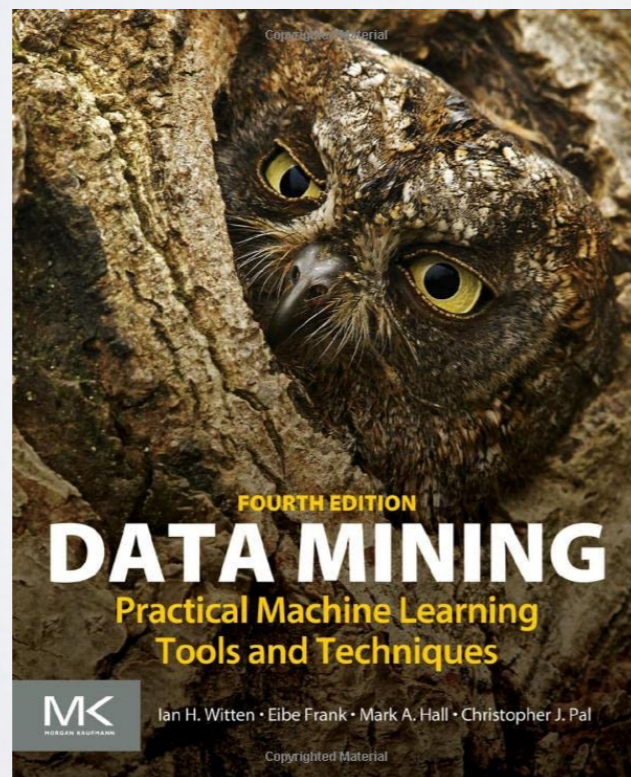
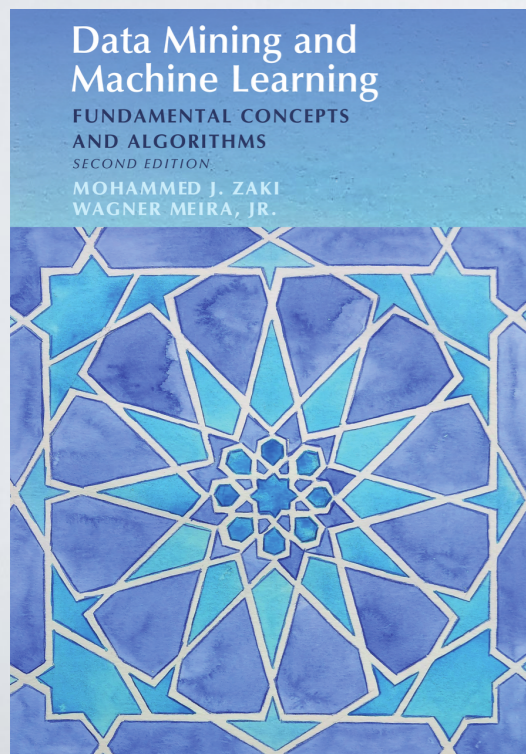
- Machine learning(ML) involves computers discovering how they can perform tasks without being explicitly programmed to do so. It involves computers learning from data provided so that they carry out certain tasks. It is a subset of Artificial Intelligence.
- [https://en.wikipedia.org/wiki/Machine_learning]

CLASS OVERVIEW

- Data description, preparation, etc.
- Unsupervised ML (beyond k-means)
- Supervised ML (beyond linear regression)
- Deep Neural Networks
- Network mining, GNN
- Large language models

THIS CLASS

- This class is based on:
 - Countless Wikipedia and blogs (use them too!)
- Some books
 - Borrow at my office



CLASS OVERVIEW

- Class with me: lecture + practical
- Two other lecturers
- Details on the lecture page:
 - <http://cazabetremy.fr/Teaching/DISS/ML.html>
- Exam:
 - Paper presenting 10%
 - Short project by other lecturers 10%
 - Final project 30% (small groups)
 - Final Exam 50%

TYPES OF DATA

DATA TYPES

- Data types : What kind of data (feature, variables) can we encounter?
 - ▶ People
 - Name, Age, Gender, Revenue, Birth Date, Address, etc.
 - ▶ House/Apartment
 - Surface area, Floor, Address, # of rooms, # of Windows, Elevator, etc.
- Types of features?

DATA TYPES

- Nominal
 - ▶ From “names”. No order between possible values
 - ▶ Color, Gender, Animal, Brand, etc. (Numbers: Participant ID, class...)
- Ordinal
 - ▶ Order between values, but not numeric
 - ▶ Size[small, medium, large], [Satisfied, ..., Unsatisfied]
- Interval
- Ratio

INTERVAL

- Numeric values, Difference is meaningful
 - ▶ $T^{\circ}: 30^{\circ} - 20^{\circ} = 15^{\circ} - 5^{\circ}$, But $30^{\circ} \neq 2 * 15^{\circ}$
 - ▶ $2022 - 2020 = 1789 - 1787$, but $1011 \neq 2022/2$
 - ▶ $=>0$ is not a meaningful value, is arbitrary
- ▶ No multiplicative relation, no ratio \Rightarrow You should not log-transform...
 - Log10: Increasing the value by 1 means multiplying by 10. But multiplying is wrong!

RATIO

- Numerical values, all operations are valid
 - Height, Duration, Revenue...
 - $=>0$ means “absence of value”.

OTHER TYPES

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

TRICKY CASES

- Real life is complex
- You will have to do modeling choices (feature engineering...)
- Possibles values: Blue, Cyan, White, Yellow, Orange, Red.
 - Nominal or Ordinal ?
- Survey: “rate X on a scale from 0 to 5”
 - What if labels are associated ? (“Bad”, “average”, ...)

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

MISSING VALUES

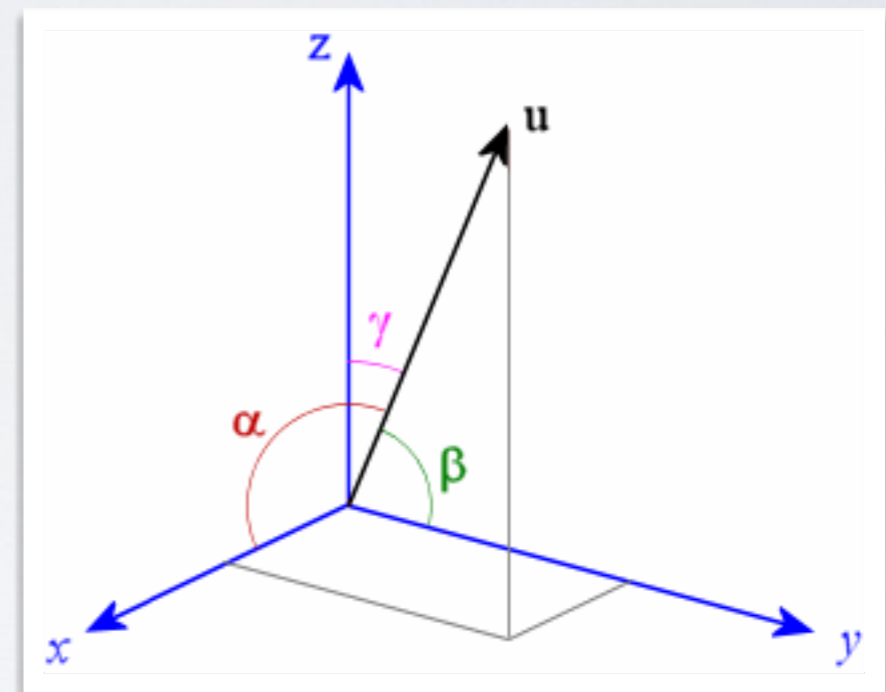
- Real life datasets are full of missing values
 - Impossible data: hair color for a bald person
 - More generally, failed to obtain them
- Few ML methods can deal with missing values
 - => Imputation
 - Naive: fill with average value
 - Use ML to fill missing values (other problems, introduce biases...)
 - 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...)
- If the data is plausible, no simple solutions
- Two common problems can be detected
 - 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...

UNIVARIATE / MULTIVARIATE

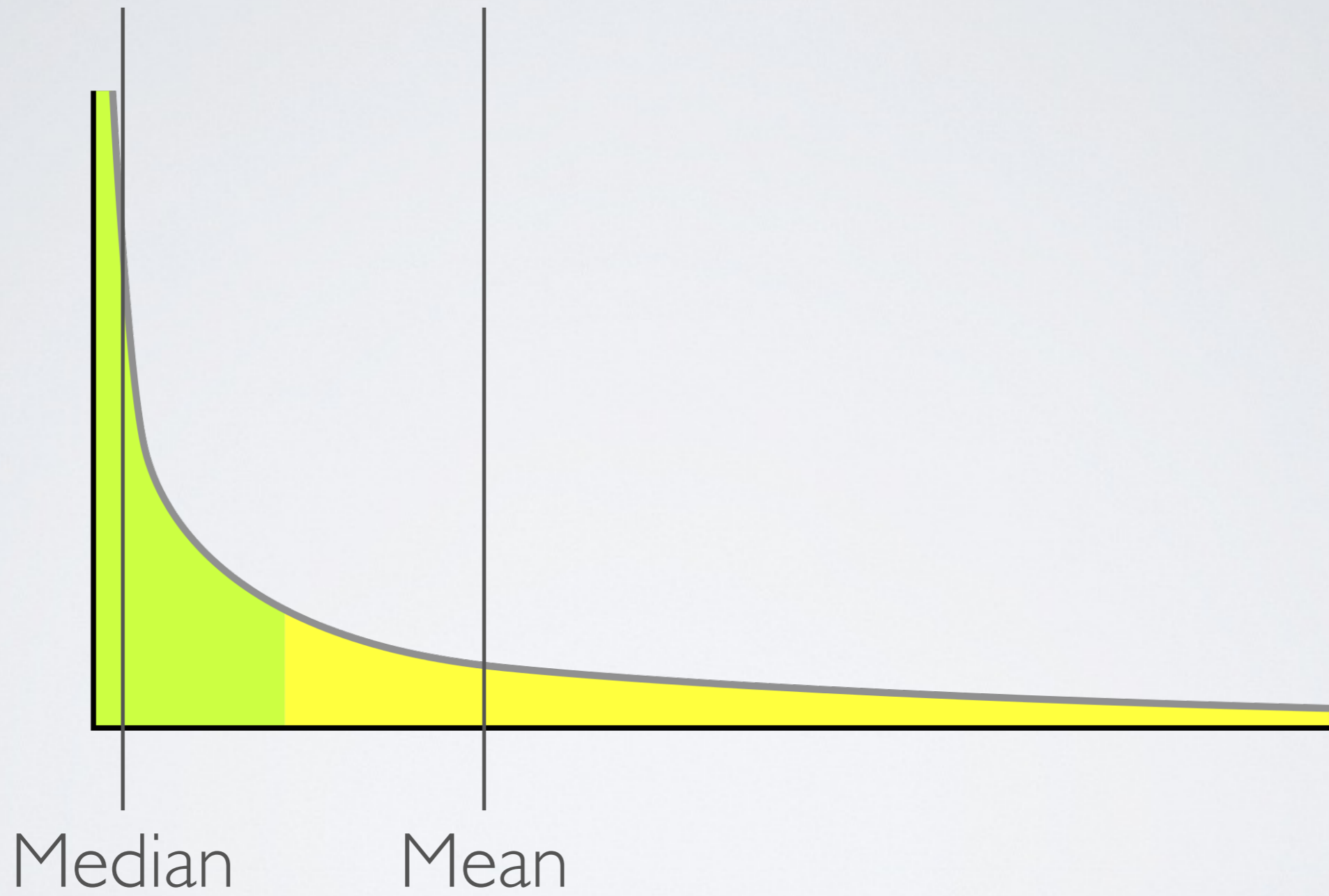
- Single *feature*: univariate
 - Age
- Real life: multivariate.
 - 2D (age, weight)
 - 3D (age, weight, height)
 - 4D (age, weight, height, genre)
 - ...



DESCRIBING A VARIABLE

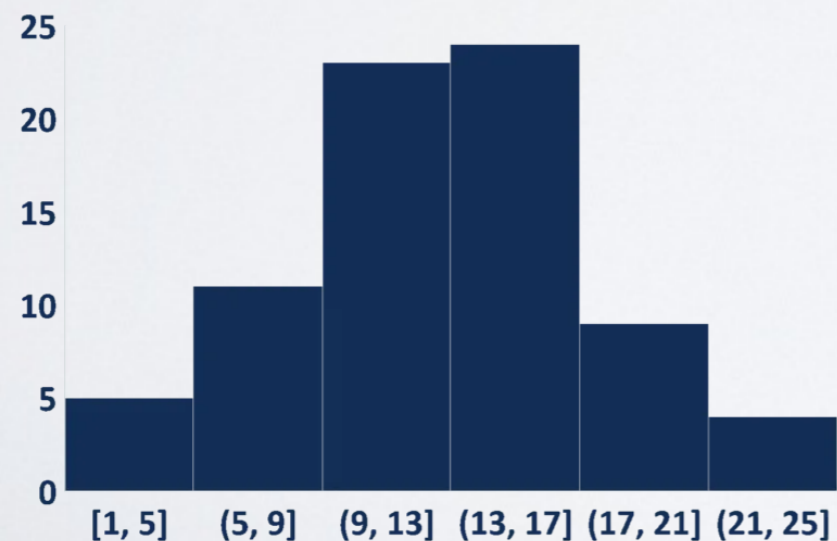
DESCRIBING VALUES

- Mean / Average
 - Be careful, not necessarily representative !
- Median
 - Be careful, not necessarily representative !
- Mode
 - Not necessarily representative
- Min/Max
 - ...

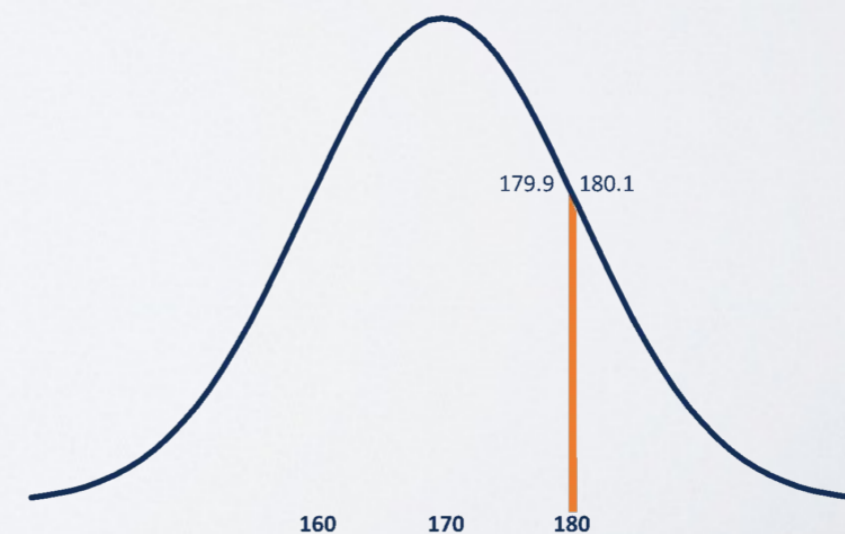


DISTRIBUTION

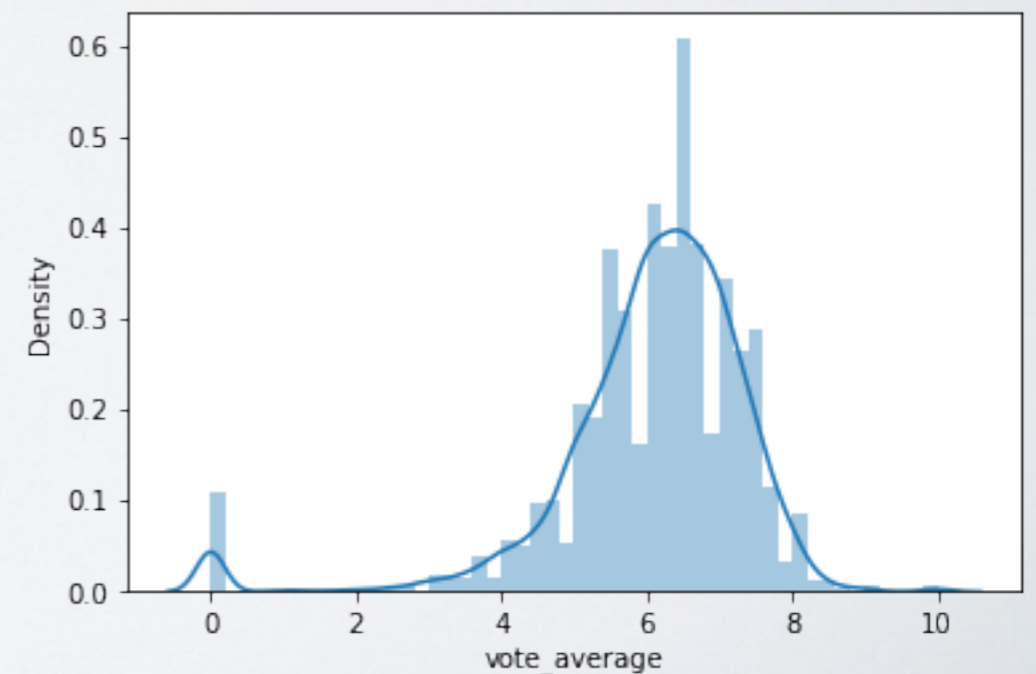
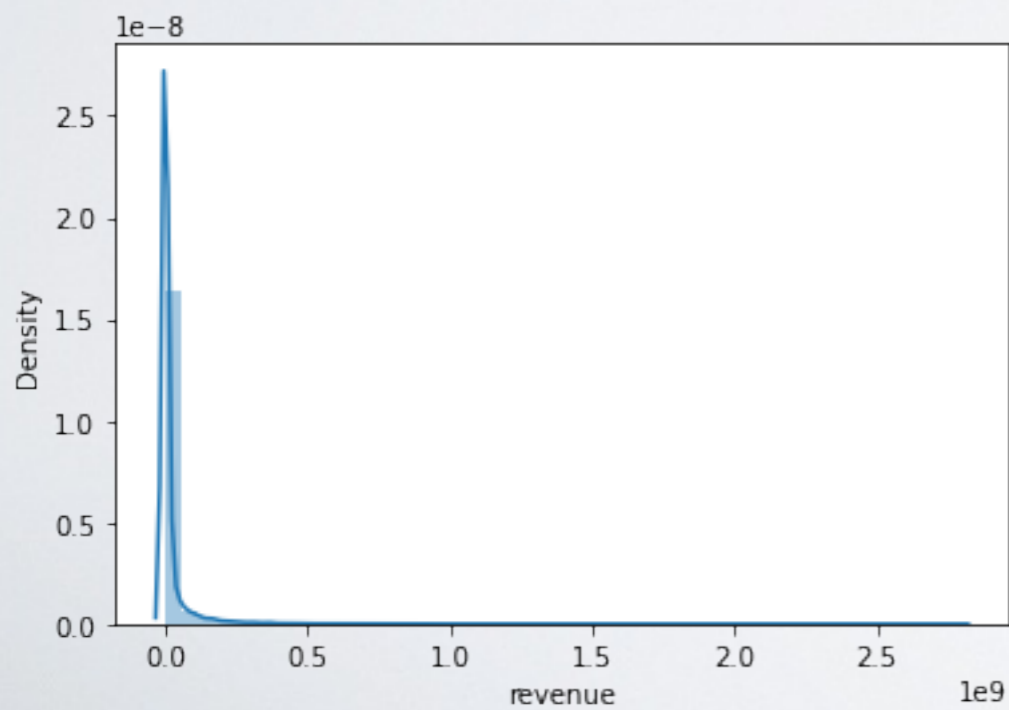
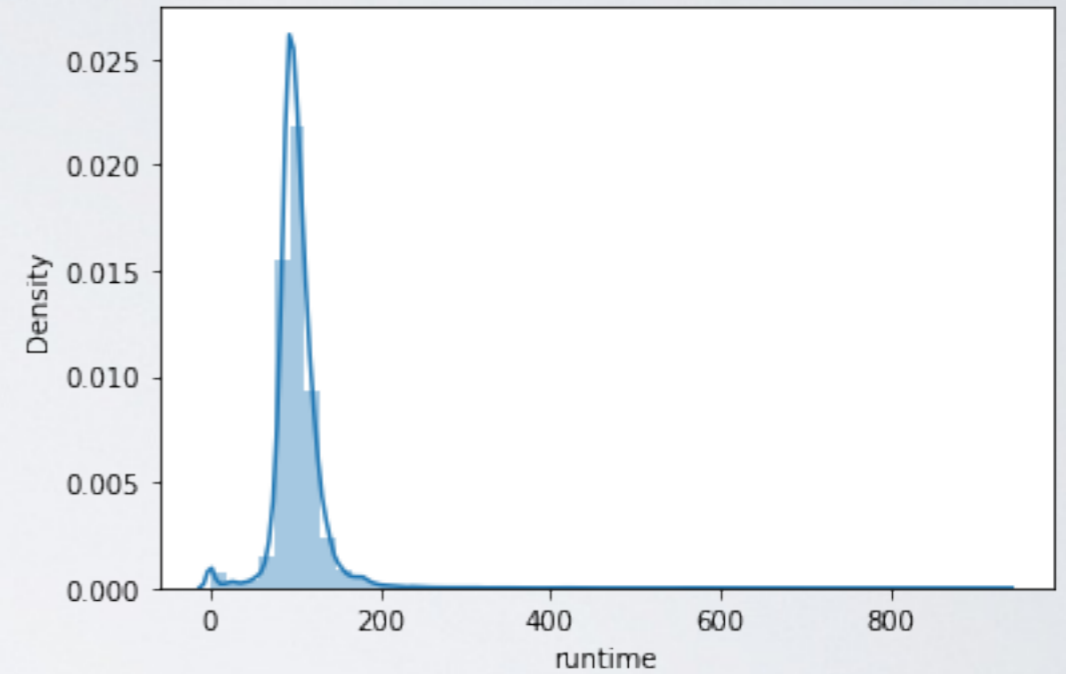
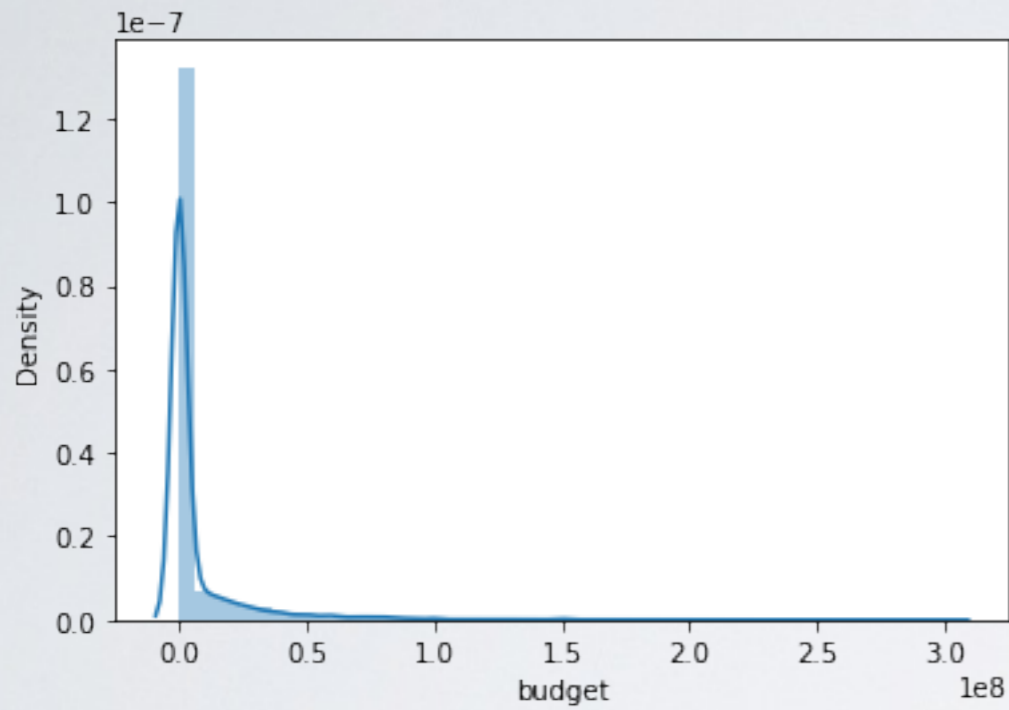
- What is a distribution?
 - ▶ A description of the frequency of occurrence of items
 - ▶ A generative function describing the probability to observe any of the possible events
 - ▶ Discrete or continuous



Continuous Distribution

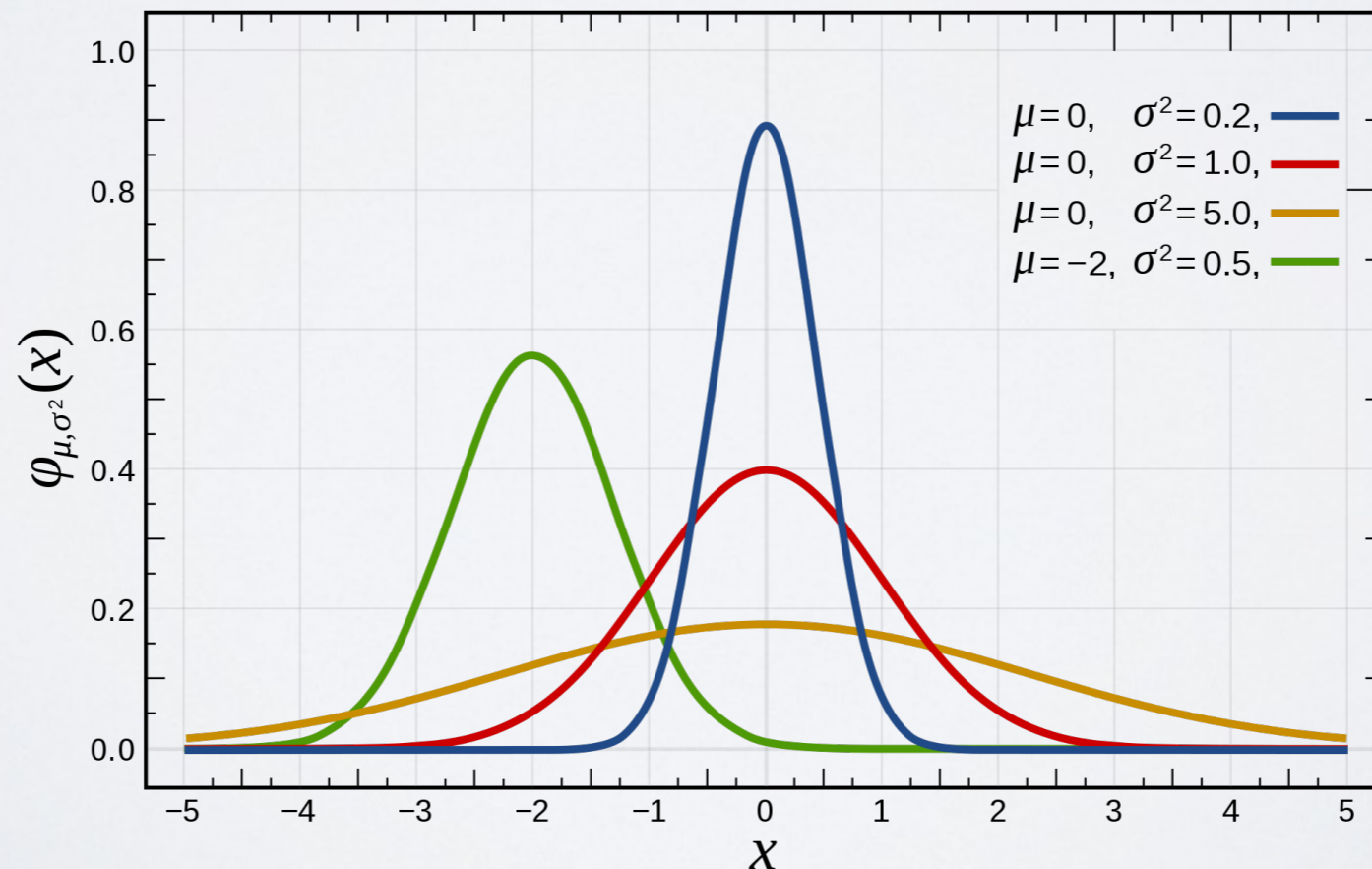


EMPIRICAL DISTRIBUTIONS



THEORETICAL DISTRIBUTIONS

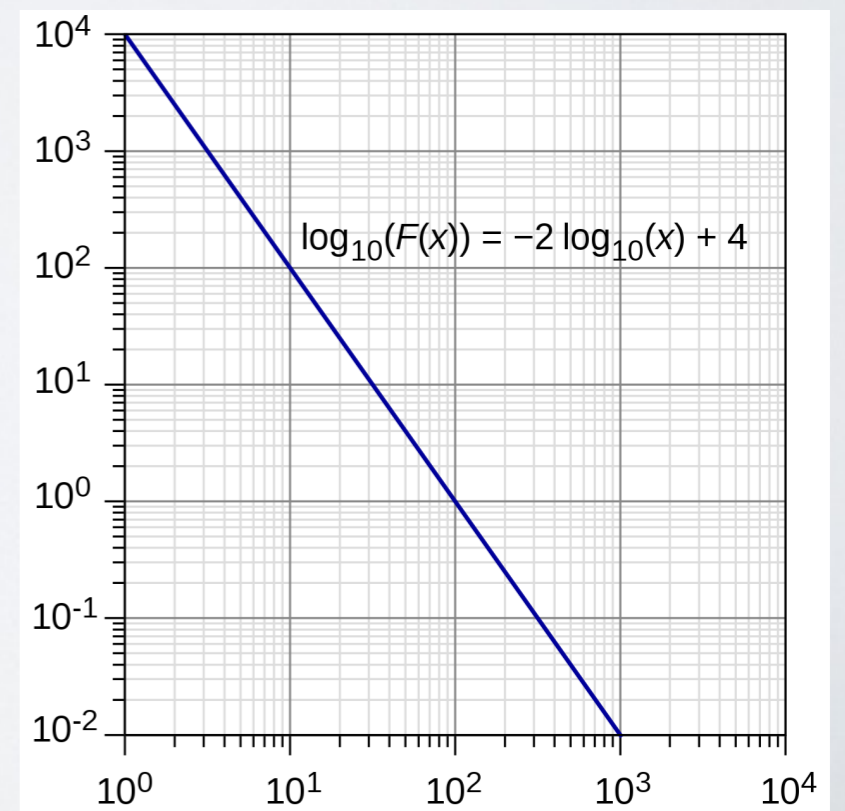
- Normal distribution
 - ▶ Many real variables follow it approximately (height, weight, price of a given product in various locations...)
 - ▶ Random variations around a well-defined mean
 - ▶ Central limit theorem: average of many samples of a random variable converges to a normal distribution



THEORETICAL DISTRIBUTIONS

- Power Law distribution

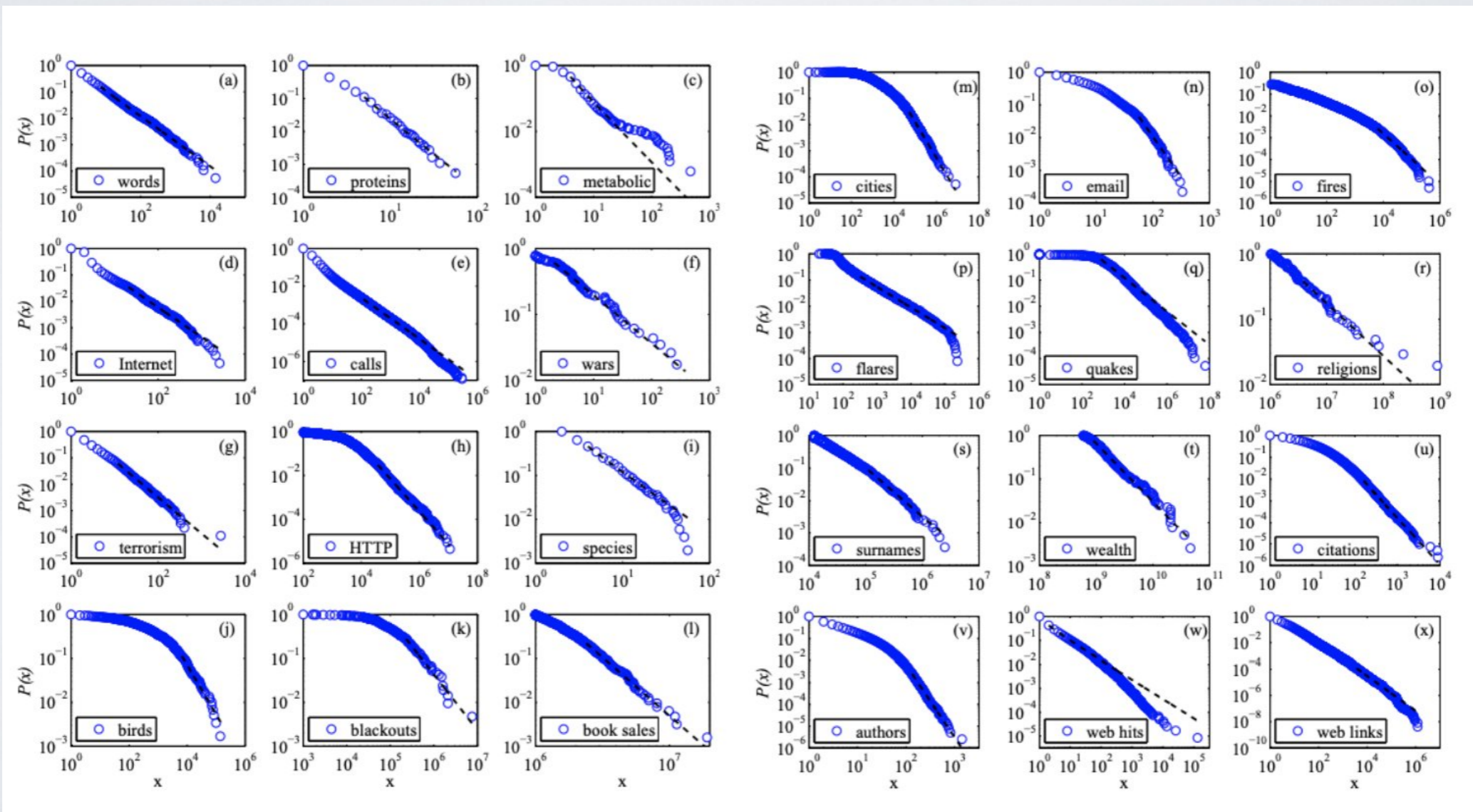
- ▶ A relative change in one quantity results in a proportional relative change in the other quantity, independent of the initial size of those quantities: one quantity varies as a power of another.
 - e.g., earthquakes 10 times more powerful are x times less frequent.
 - e.g., cities 10 times bigger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



THEORETICAL DISTRIBUTIONS

DISTRIBUTIONS

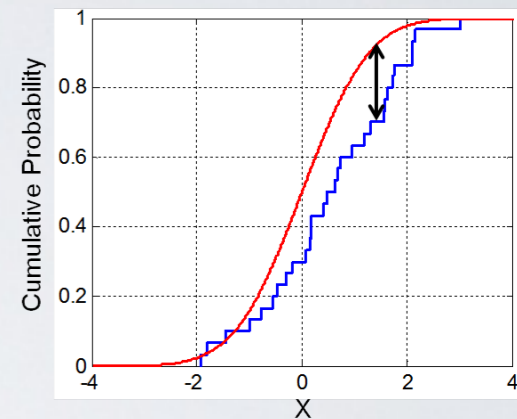
- Power Law distribution



DISTRIBUTION COMPARISON

- Statistical test

- ▶ P-value: **The probability** that my observed data could be observed if it were generated by the theoretical distribution XXX (null hypothesis)
 - Normality: Shapiro-Wilk, etc.
 - Categorical variables : Chi-squared χ^2
 - Etc. (search for the right test if you need it)
- ▶ High p-value: high probability to come from the null hypothesis
 - We usually set a p-value threshold, i.e., 0.05. (5% chance)
 - IF the p-value is below it, **I can conclude** that it is unlikely that my data has been generated by this exact null model (I can never be 100% sure)
 - IF the p-value is above, I can say that it is *possible* that it has been generated by it. However, it could also have been generated by another null hypothesis that I have not tried. **I cannot conclude.**



VARIANCE

- Variance:
 - Expectation of the squared deviation of a random variable from its mean

$$\text{Var}(X) = \sigma^2 = \text{E} [(X - \mu)^2]$$

Also expressed as average squared distance
between all elements

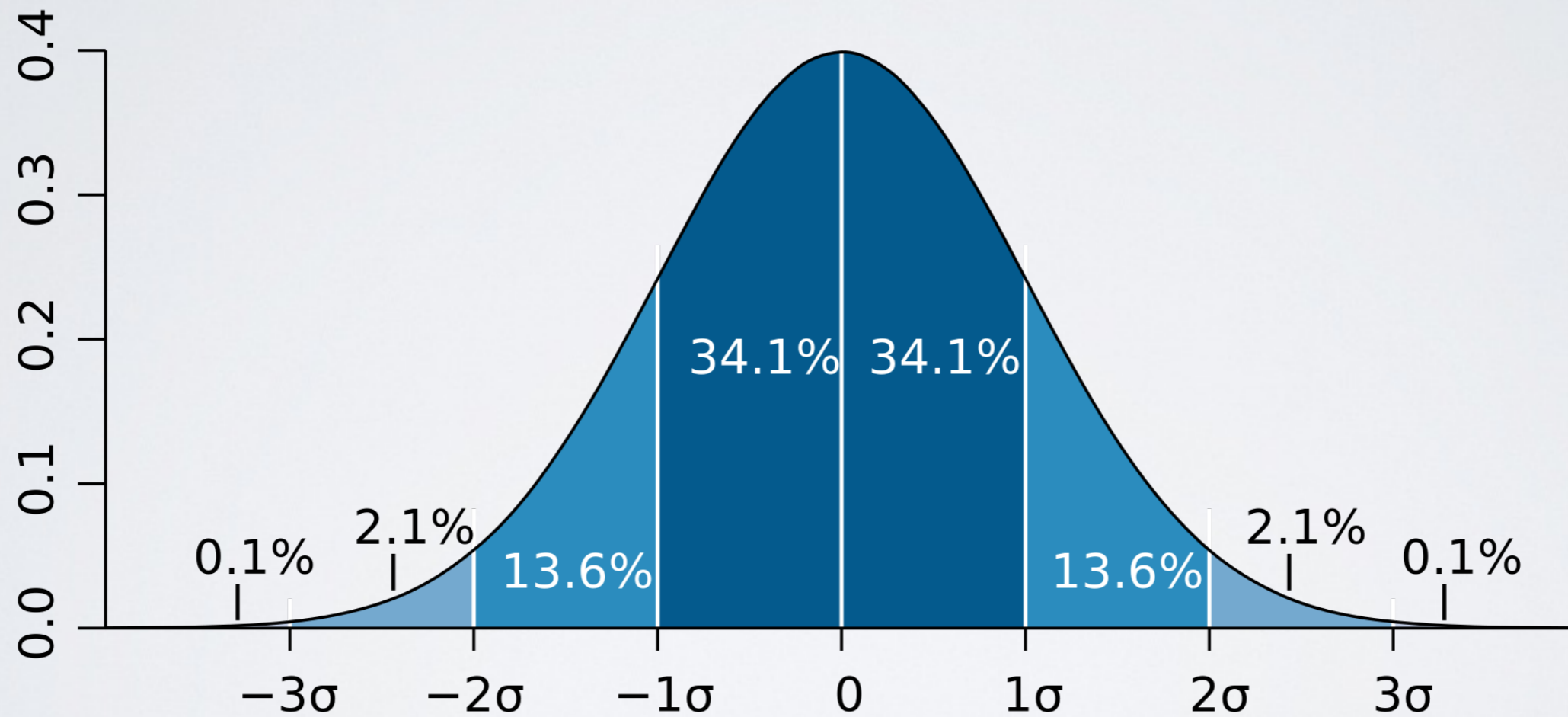
$$\sigma^2 = \frac{1}{N^2} \sum_{i < j} (x_i - x_j)^2$$

STANDARD DEVIATION

- Squared root of the Variance

$$\sigma = \sqrt{\sigma^2} = \sqrt{\mathbf{E} [(X - \mu)^2]}$$

RELATION WITH NORMAL DISTRIBUTION



VARIABLE INTERACTIONS

COVARIANCE MATRIX

Covariance Matrix Formula



- Covariance matrix \mathbf{K}

- ▶ Extension of Variance to multivariate data

- ▶ $\text{Var}(X) = \mathbf{E} [(X - \mu)^2]$

- ▶ $\text{cov}(\mathbf{X}, \mathbf{Y}) = \mathbf{K}_{\mathbf{XY}} = \mathbf{E} [(\mathbf{X} - \mathbf{E}[\mathbf{X}])(\mathbf{Y} - \mathbf{E}[\mathbf{Y}])^T]$

- 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

- ▶ $\Rightarrow \text{cov}(\mathbf{X}, \mathbf{X}) = \text{Var}(\mathbf{X})$

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

- Covariance is hardly interpretable by itself.

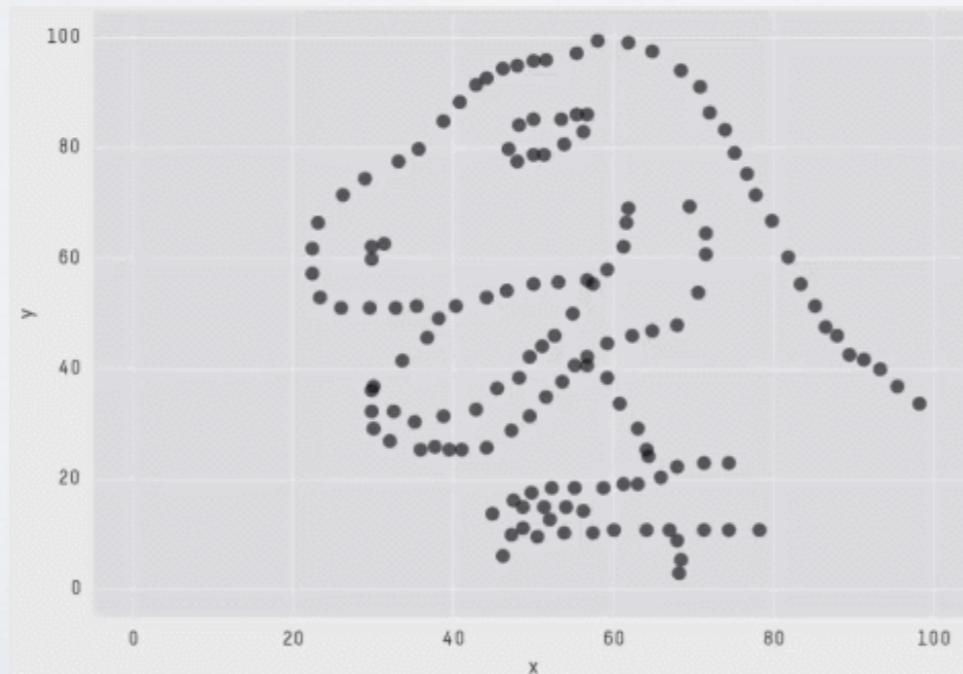
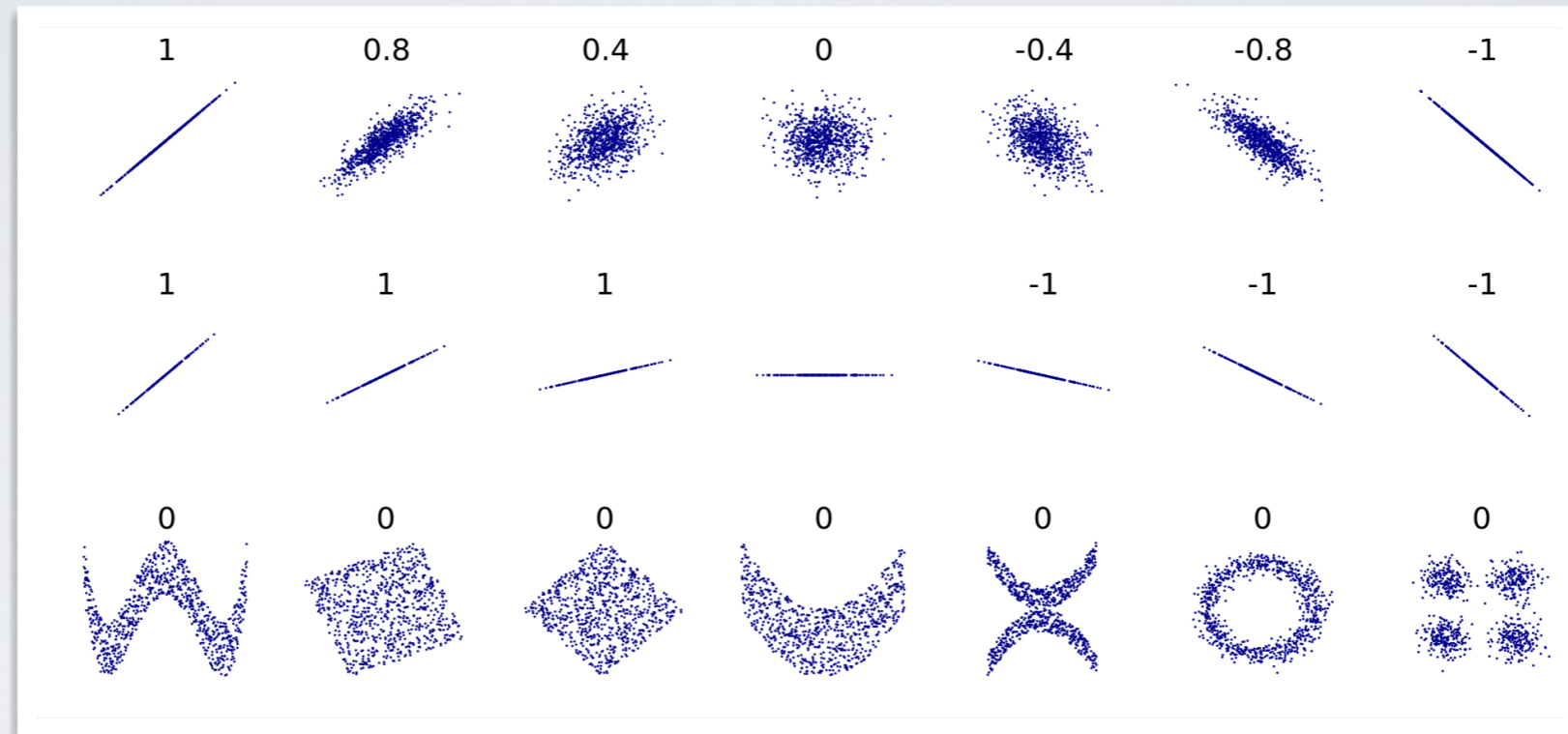
- ▶ If >0 , divergences tend to be in the same direction

- ▶ Normalize it to obtain the “correlation coefficient”

CORRELATION COEFFICIENT

- Pearson correlation coefficient : $\rho_{X,Y} = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}$
 - ▶ 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

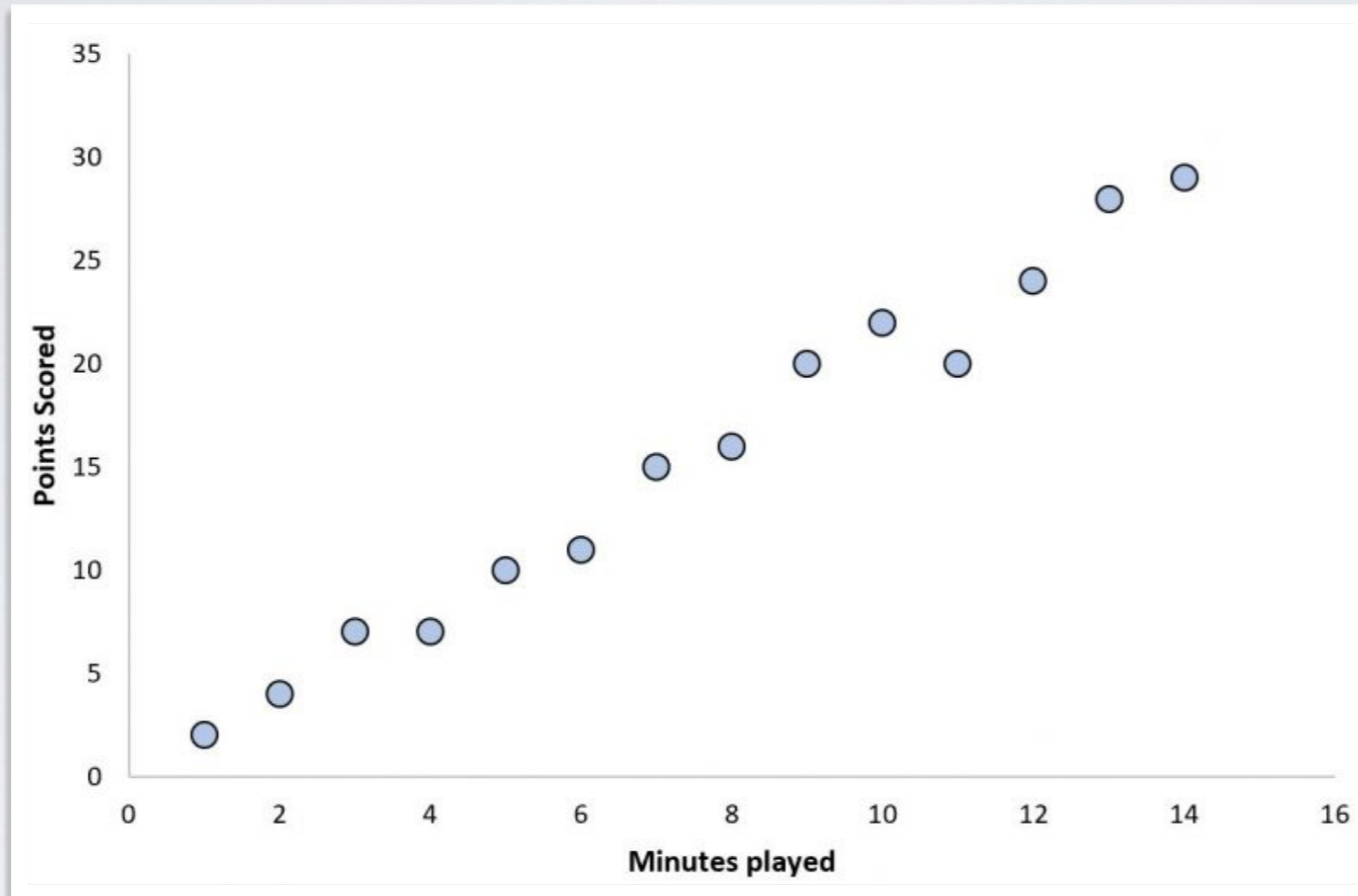


X Mean: 54.2659224
Y Mean: 47.8313999
X SD : 16.7649829
Y SD : 26.9342120
Corr. : -0.0642526

CORRELATION COEFFICIENT

- Other possible interpretation, e.g.
 - Cosine similarity of the vectors defined by the observations...
- 0.7 ? Is it a high or low value ?
 - It depends.

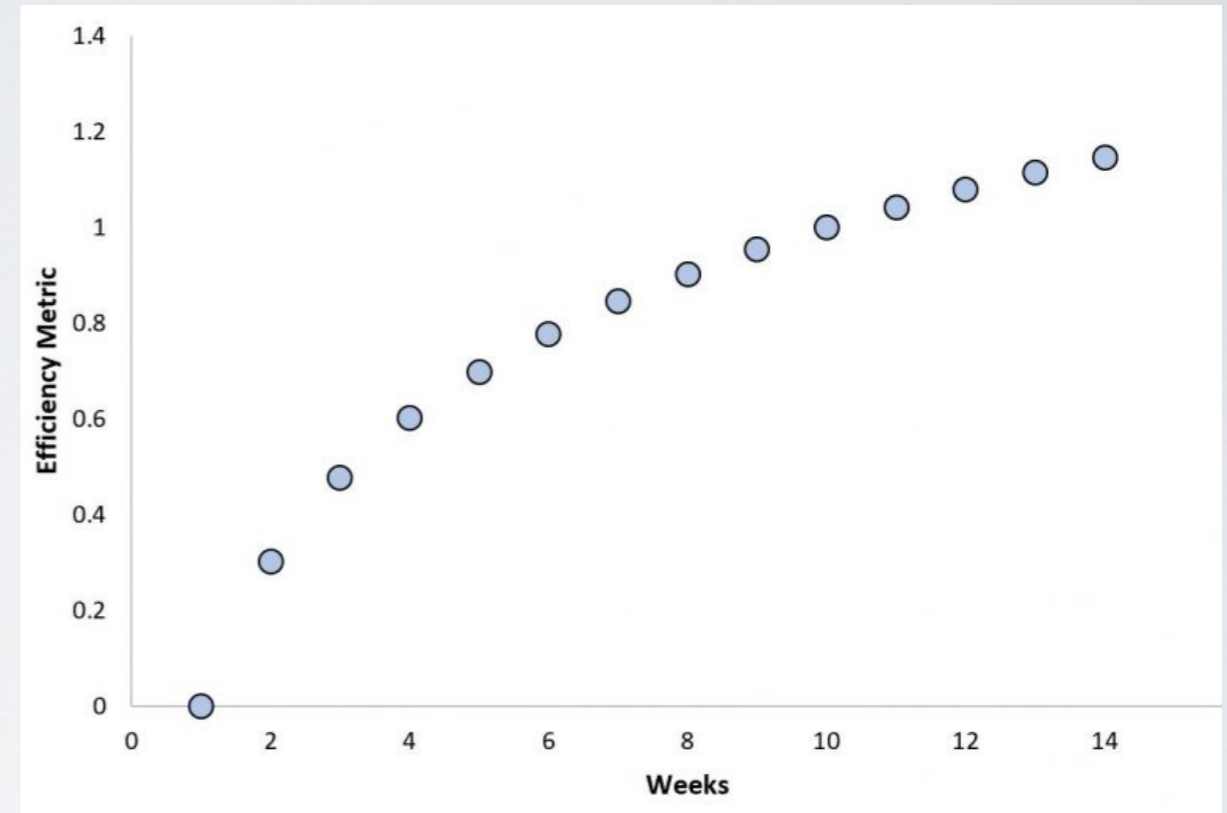
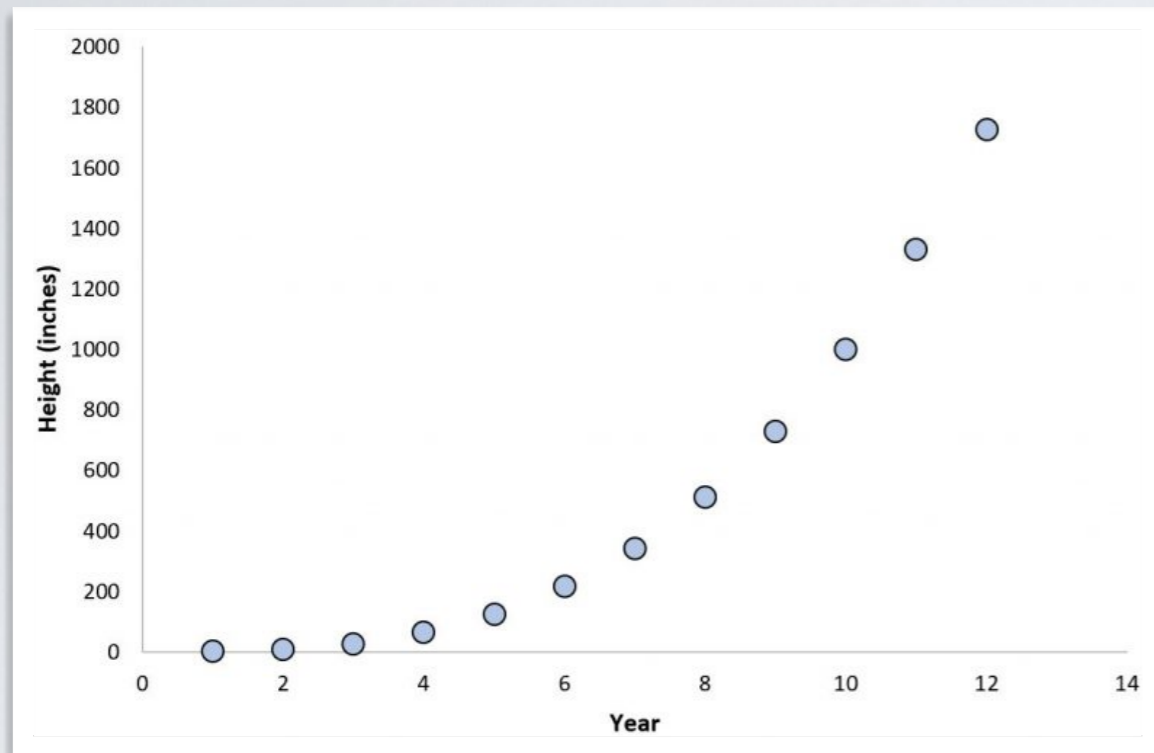
NONLINEAR RELATIONSHIPS



Linear relationship

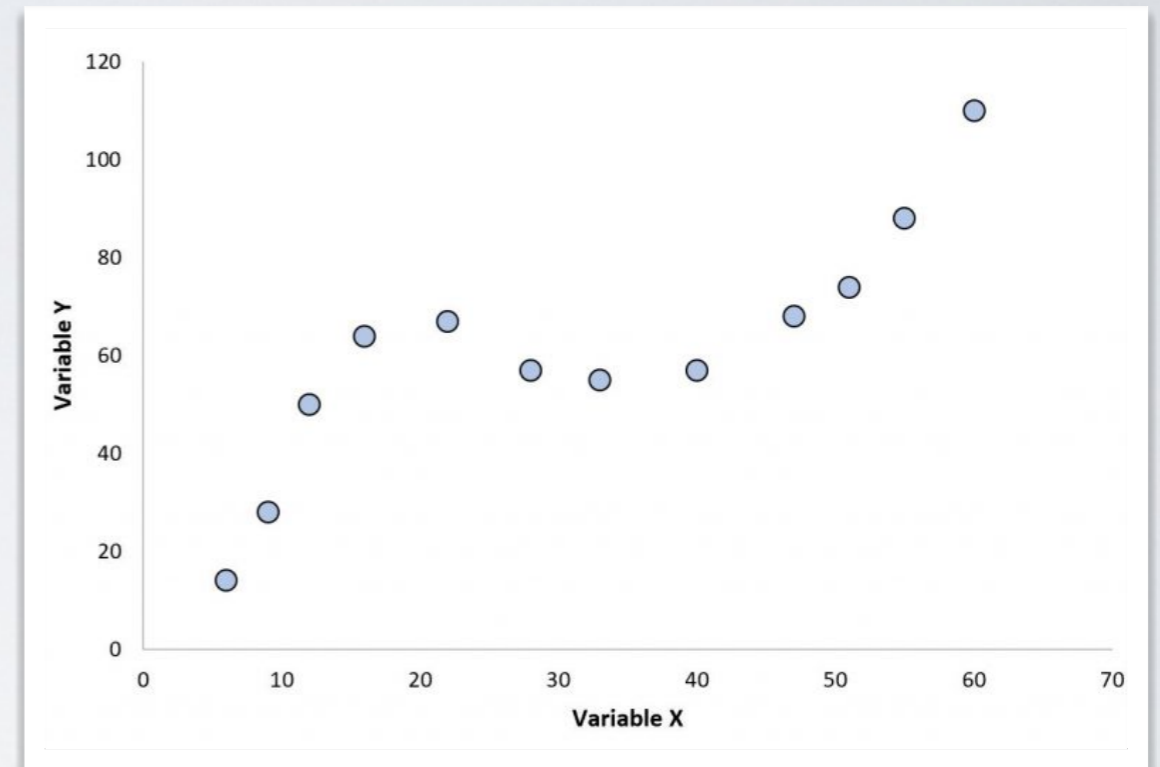
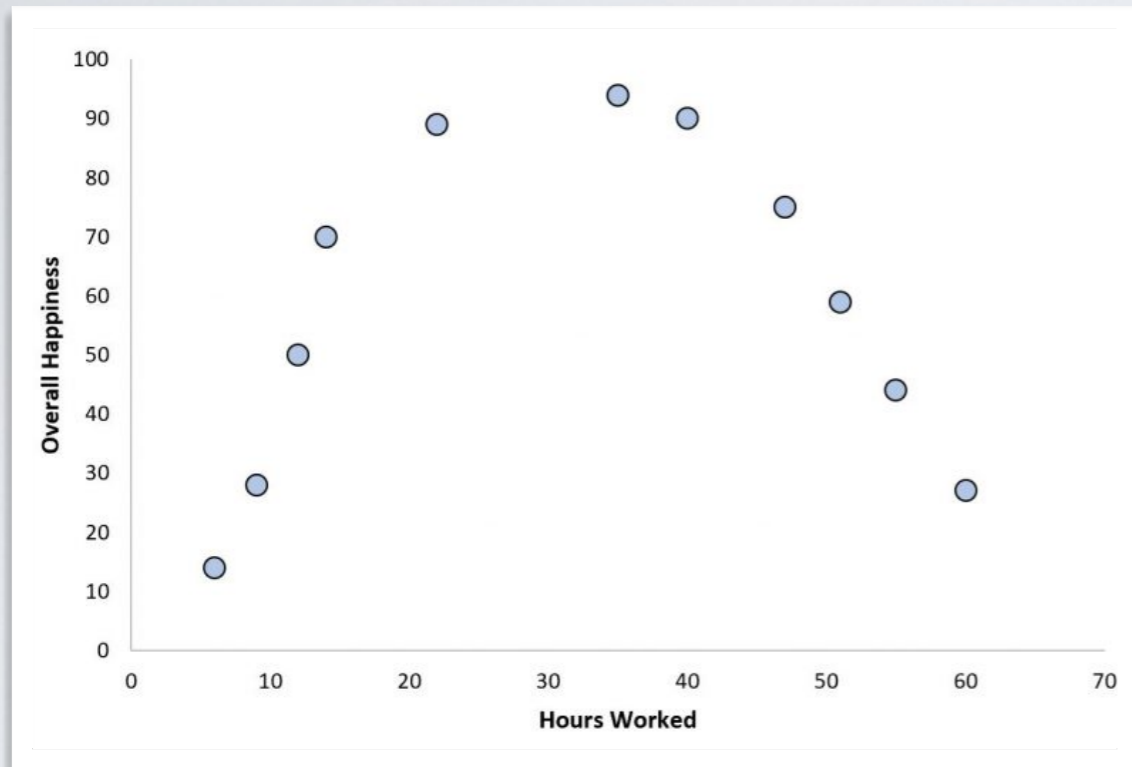
$$Y = a + bX + e$$

NONLINEAR RELATIONSHIPS

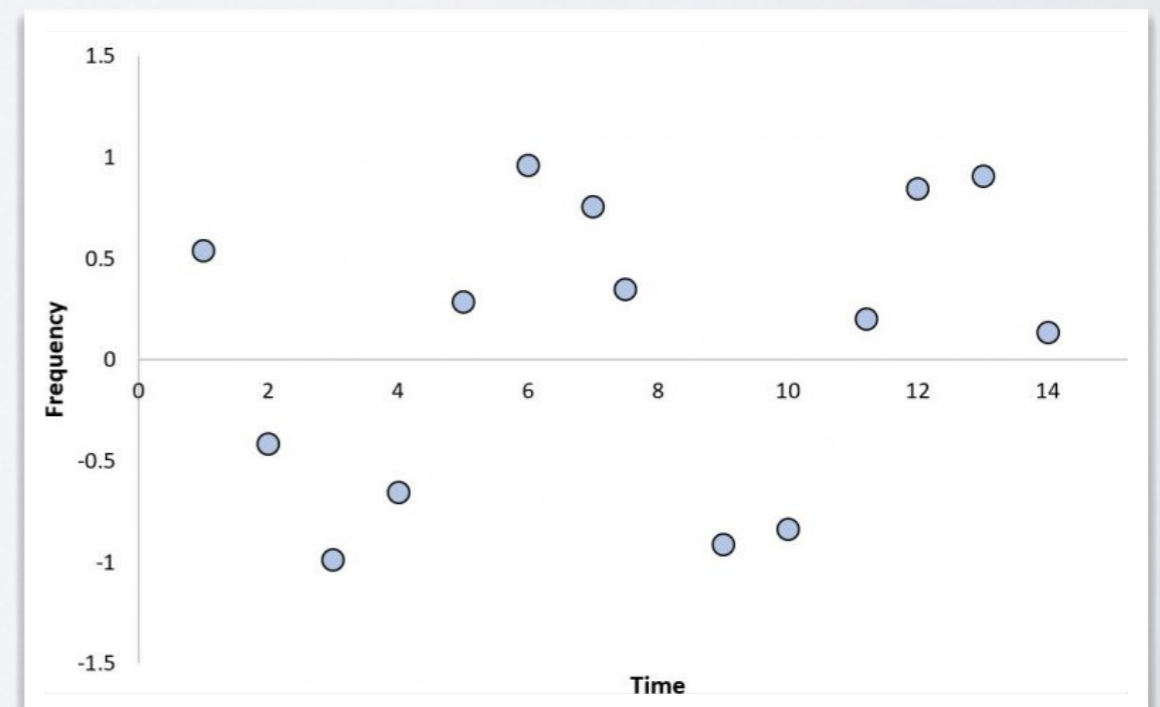


Monotonous, non-linear

NONLINEAR RELATIONSHIPS



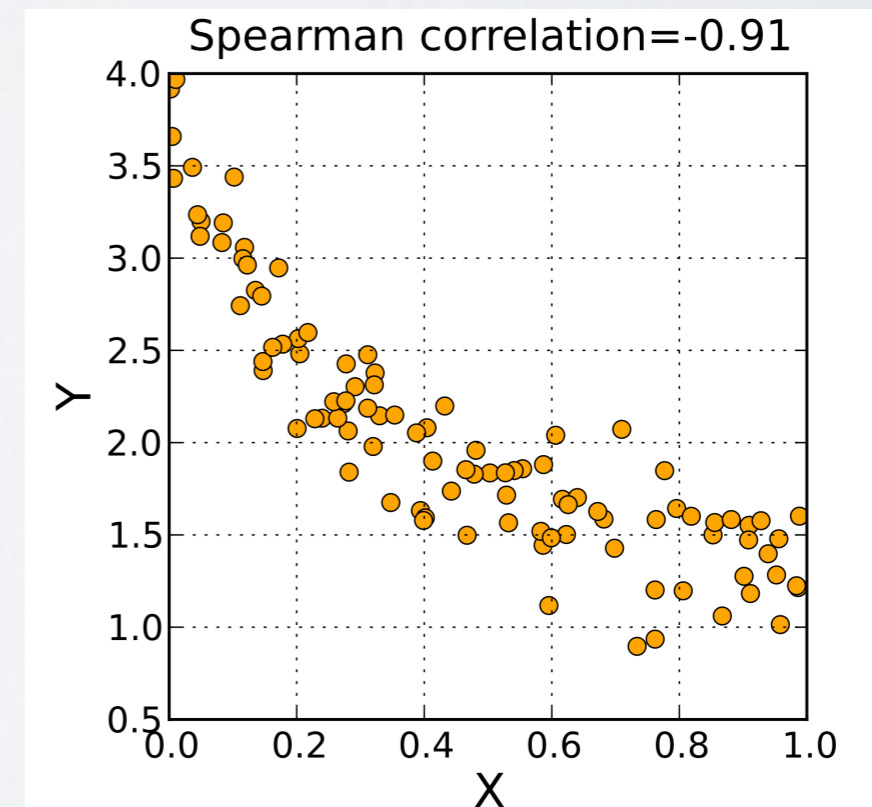
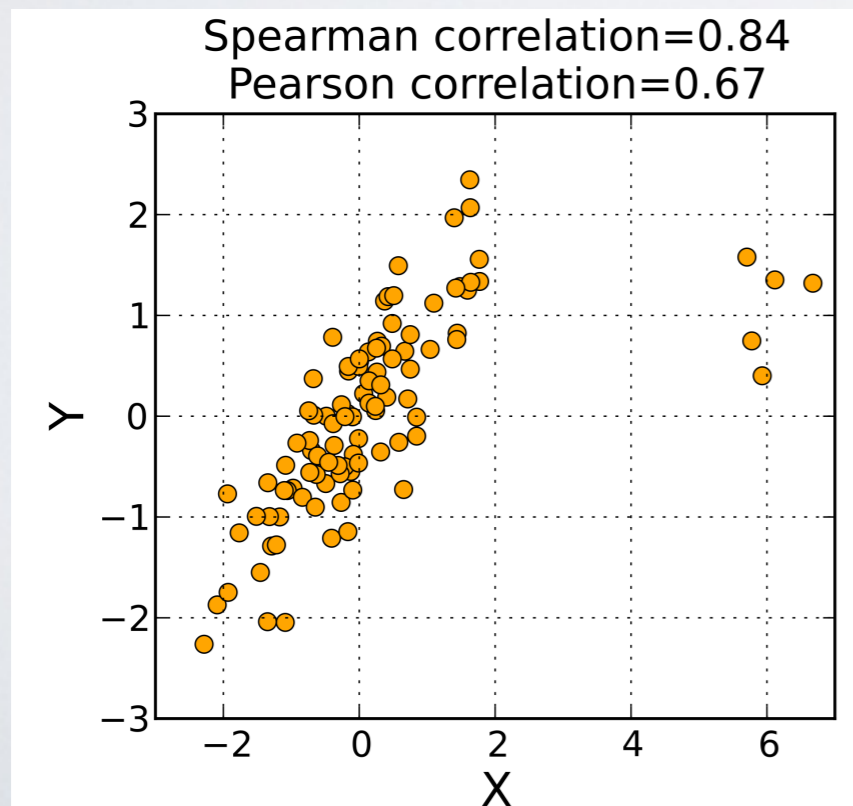
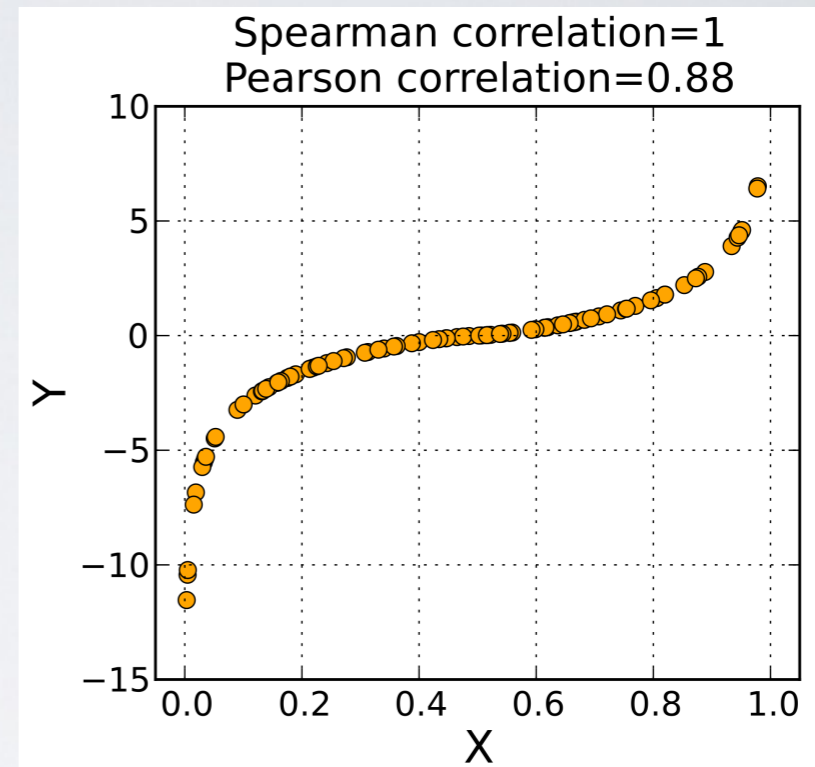
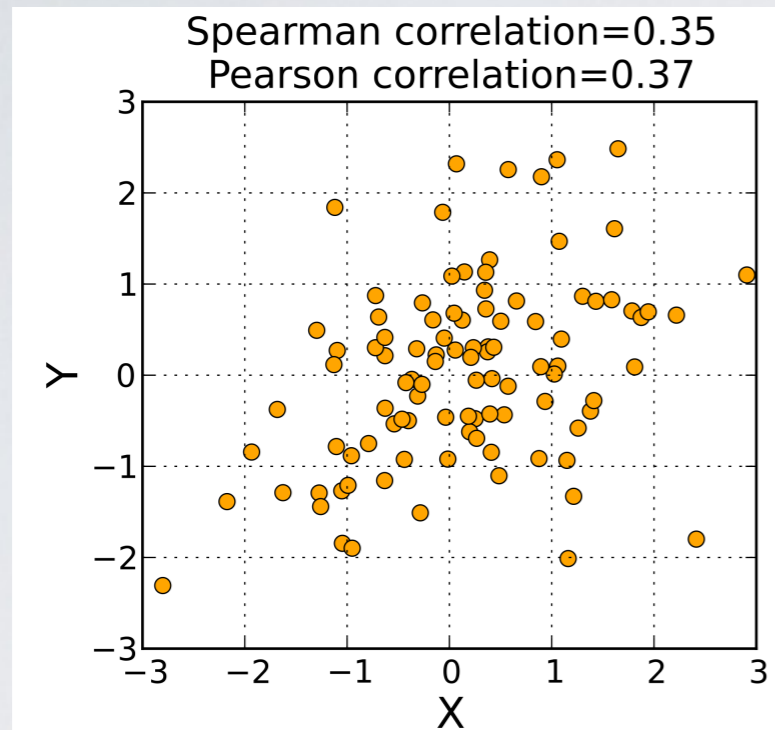
Non-monotonous,
Non-linear



SPEARMAN'S CORRELATION

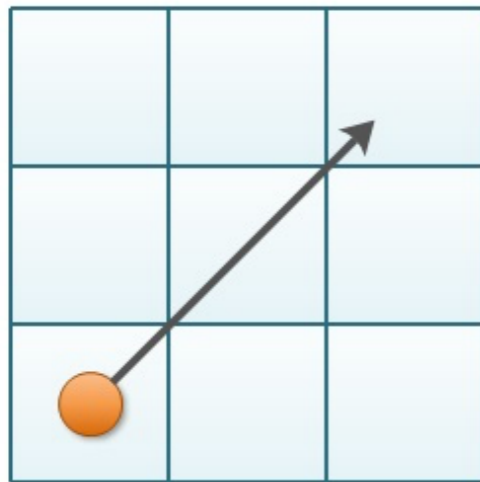
- Spearman's **rank** correlation coefficient
- Assesses how well the relationship between two variables can be described using a monotonic function
 - Not assuming a linear relation
- Pearson correlation coefficient between the rank variables
 - $r_s = \rho_{R(X), R(Y)} = \frac{\text{cov}(R(X), R(Y))}{\sigma_{R(X)}\sigma_{R(Y)}}$

SPEARMAN'S CORRELATION

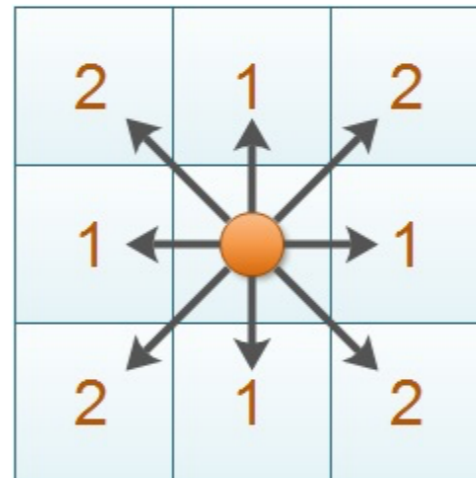


NOTIONS OF DISTANCE

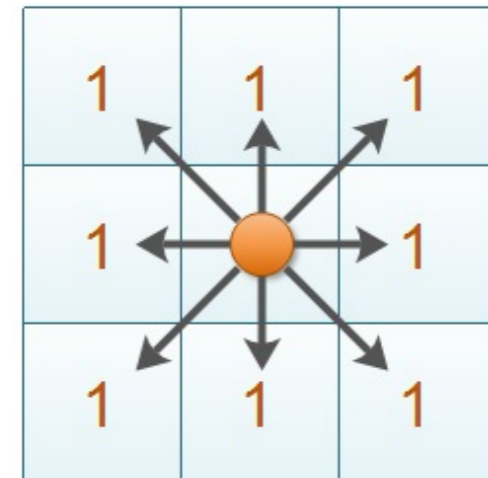
Euclidean Distance



Manhattan Distance



Chebyshev Distance



$$\sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2} \quad |x_1 - x_2| + |y_1 - y_2| \quad \max(|x_1 - x_2|, |y_1 - y_2|)$$

FEATURE SCALING

- We want to use euclidean distance to compute the “distance” between 2 people based on attributes age(y), height(m), weight(g).
 - ▶ $a = (y:20, m:1.82, g:80\ 000)$, $b = (y:20, m:1.82, g:81\ 000)$, $c = (y:90, m:1.50, g:80\ 020)$
 - $d(a,b) = 1000.0005$
 - $d(a,c) = 72.8$
 - ▶ That is not what we expected from our expert knowledge!
 - We should normalize/standardize data

FEATURE SCALING

- Rescaling (Normalization): $x' = \frac{x - \min(x)}{\max(x) - \min(x)} : [0, 1]$
- Mean normalization: $x' = \frac{x - \text{average}(x)}{\max(x) - \min(x)} : 0 = \text{mean}$
- Standardization (z-score normalization): $x' = \frac{x - \bar{x}}{\sigma}$
 - 0: mean, -1/+1: 1 standard deviation from the mean

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

EXPERIMENTS

- Go to the webpage of the class and do today's experiments
- The “Advanced” section is not mandatory, you can do it if you have time