

MACHINE LEARNING TECHNIQUES AND APPLICATIONS

M2 DISS

WHO AM I

- Rémy Cazabet (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

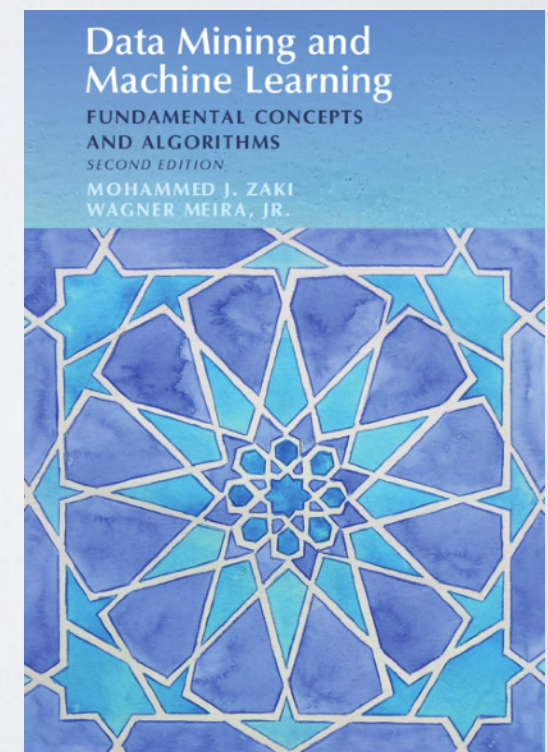
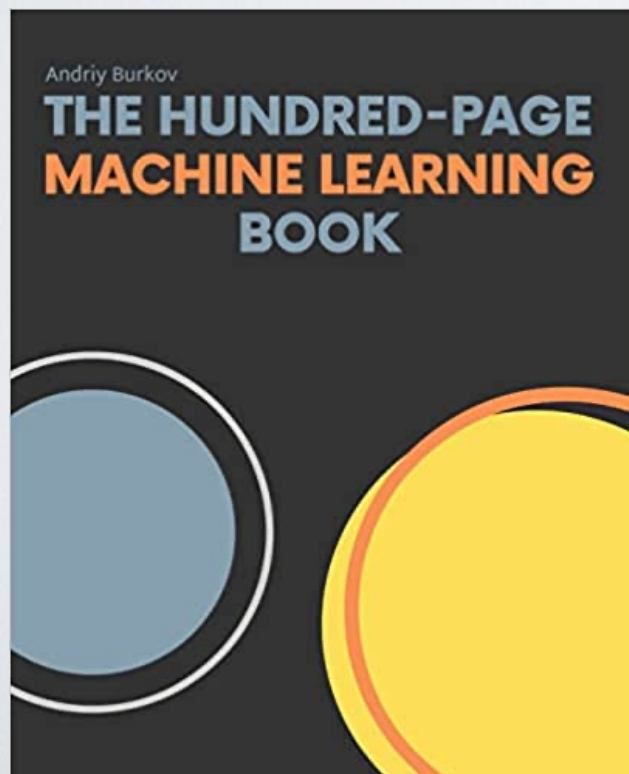
- Class with me: lecture + practical
- Two other lecturers
- Details on the lecture page:
 - Class page: <https://cazabetremy.fr/Teaching/DISS/ML.html>
 - All contents: slides, TP, data, corrections...
- Exam:
 - Paper presenting 10%
 - Short project by other lecturers 40%
 - Final Exam 50%

CLASS OVERVIEW

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

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]

DATA - INTRODUCTION

TYPES OF DATA

DATA TYPES

- Data types : What kind of data (feature, variables) can we encounter?

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...)
- Ordered:
 - Ordinal
 - Interval
 - Ratio

ORDERED

- Ordinal

- Order between values, but not numeric
- Size[small, medium, large], [Satisfied, ..., Unsatisfied], Income [0-10k],[10k-15k], [15k-50k]...

- Ratio

- Numerical values, all operations are valid
- Height, Duration, Revenue...

- Interval

- Numeric values, difference is meaningful
- T°: $30^{\circ} - 20^{\circ} = 10^{\circ}$, But $30^{\circ} \neq 2 * 15^{\circ}$
- $2022 - 2020 = 2$, but $1011 \neq 2022/2$
- $=>0$ is not a meaningful value, is arbitrary

OTHER TYPES

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

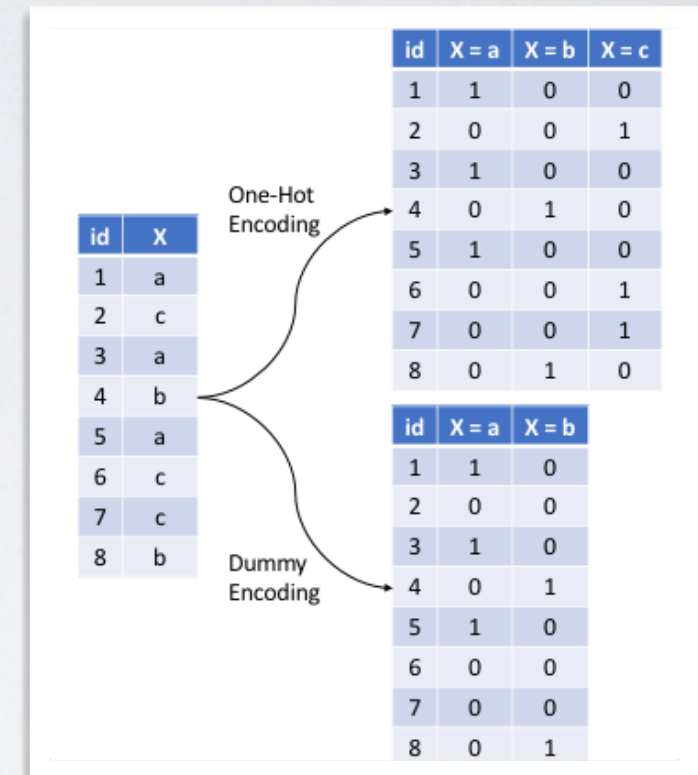
TRAPS

- 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 encoding



- Ordered:

- Ordinal => Transform to Interval/Ratio
- Interval/ratios => **usually forbidden to perform correlation, clustering, etc.**
- Ratio => :)

MISSING VALUES

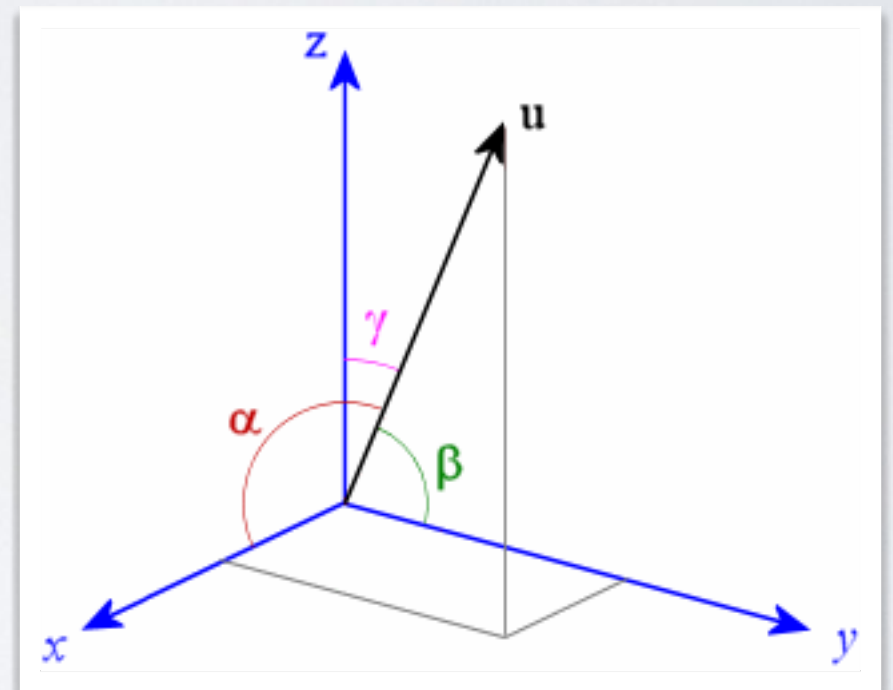
- Real-life datasets are full of missing values
 - Impossible data: *fur color* for a sphinx cat
 - More generally, failure to obtain them
- Few methods can deal with missing values
 - => Imputation
 - Naive: fill with average value
 - Use ML to fill-in missing values (other problems, introduce biases...)
 - 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...

UNIVARIATE / MULTIVARIATE

- Terminology:
 - Feature=variable="columns"
- 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
 - ...

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

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

ABSOLUTE DEVIATION

- MAD (Mean Absolute Deviation)

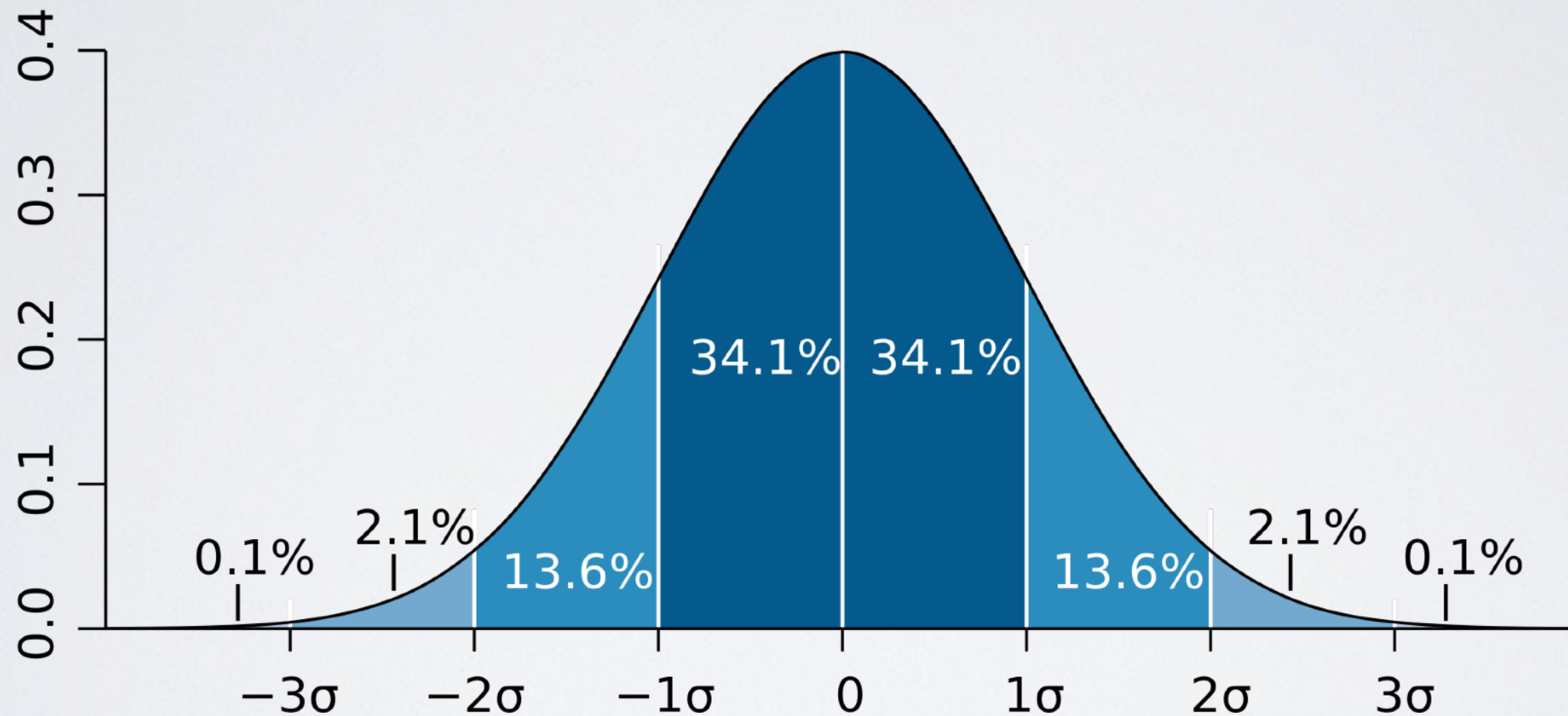
- Deviation from mean or from median
- (Variant: Median Absolute Deviation)

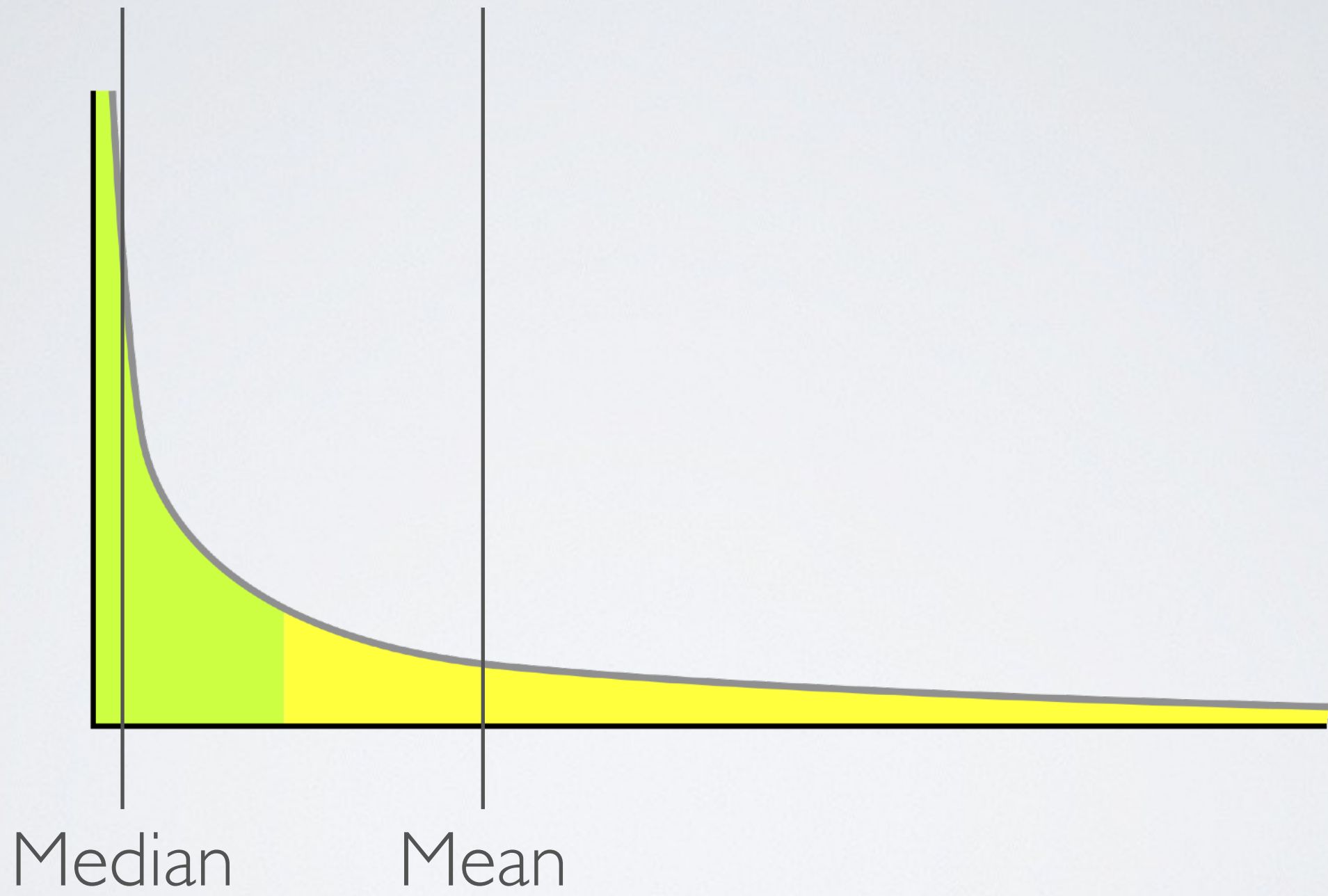
- $$\frac{1}{n} \sum_{i=1}^n |x_i - m(X)|$$

- So why are we using the Standard Deviation again?

- The mean minimizes the expected squared distance
- The median minimizes the MAD
- Leads naturally to least square regression and PCA... see later.

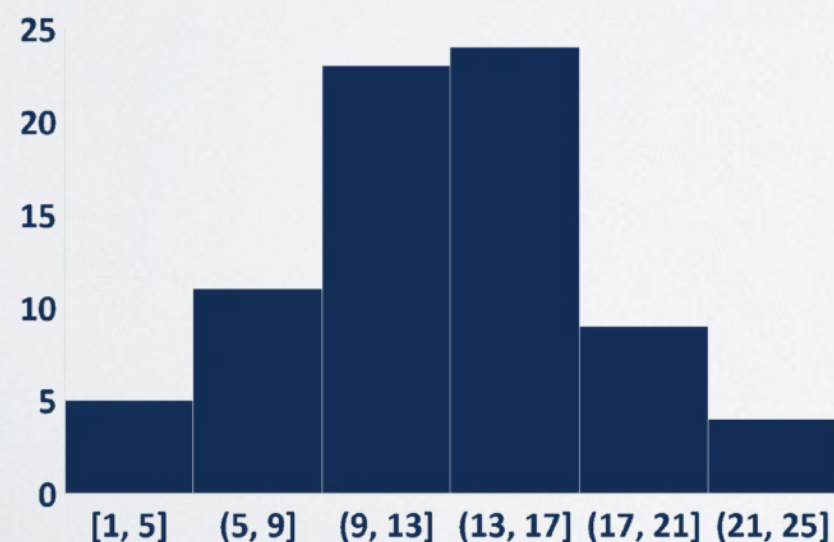
STDIV AND NORMAL DISTRIBUTION



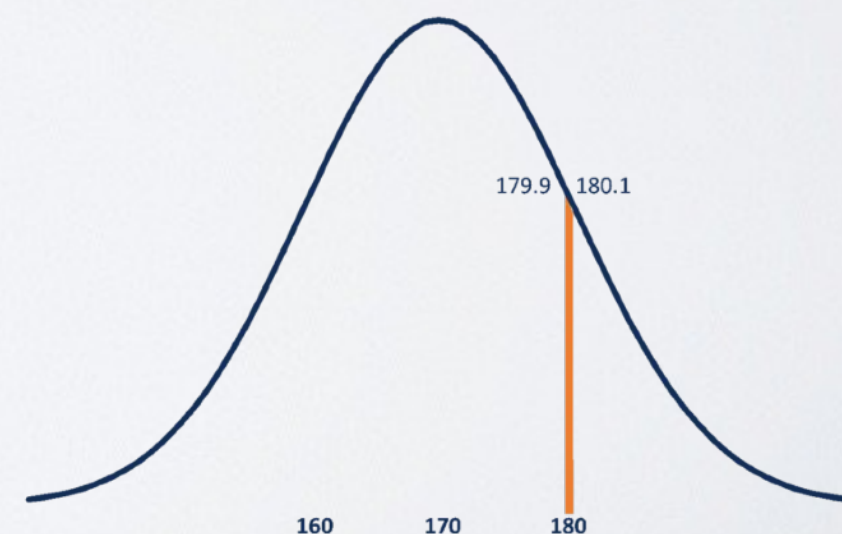


DISTRIBUTION

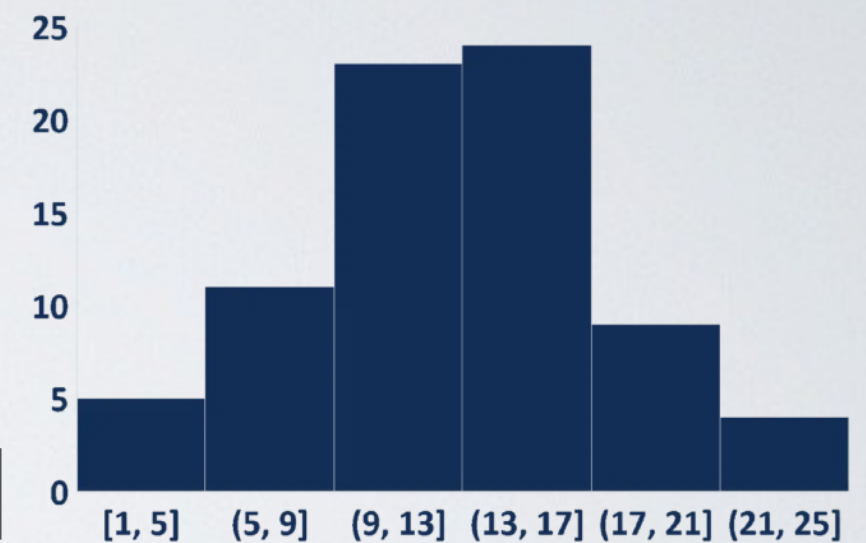
- What is a distribution?
 - A description of the frequency of occurrence of items
 - A generative function describing the probability to observe any of the possible events
 - Discrete or continuous



Continuous Distribution



DISTRIBUTION (DISCRETE)

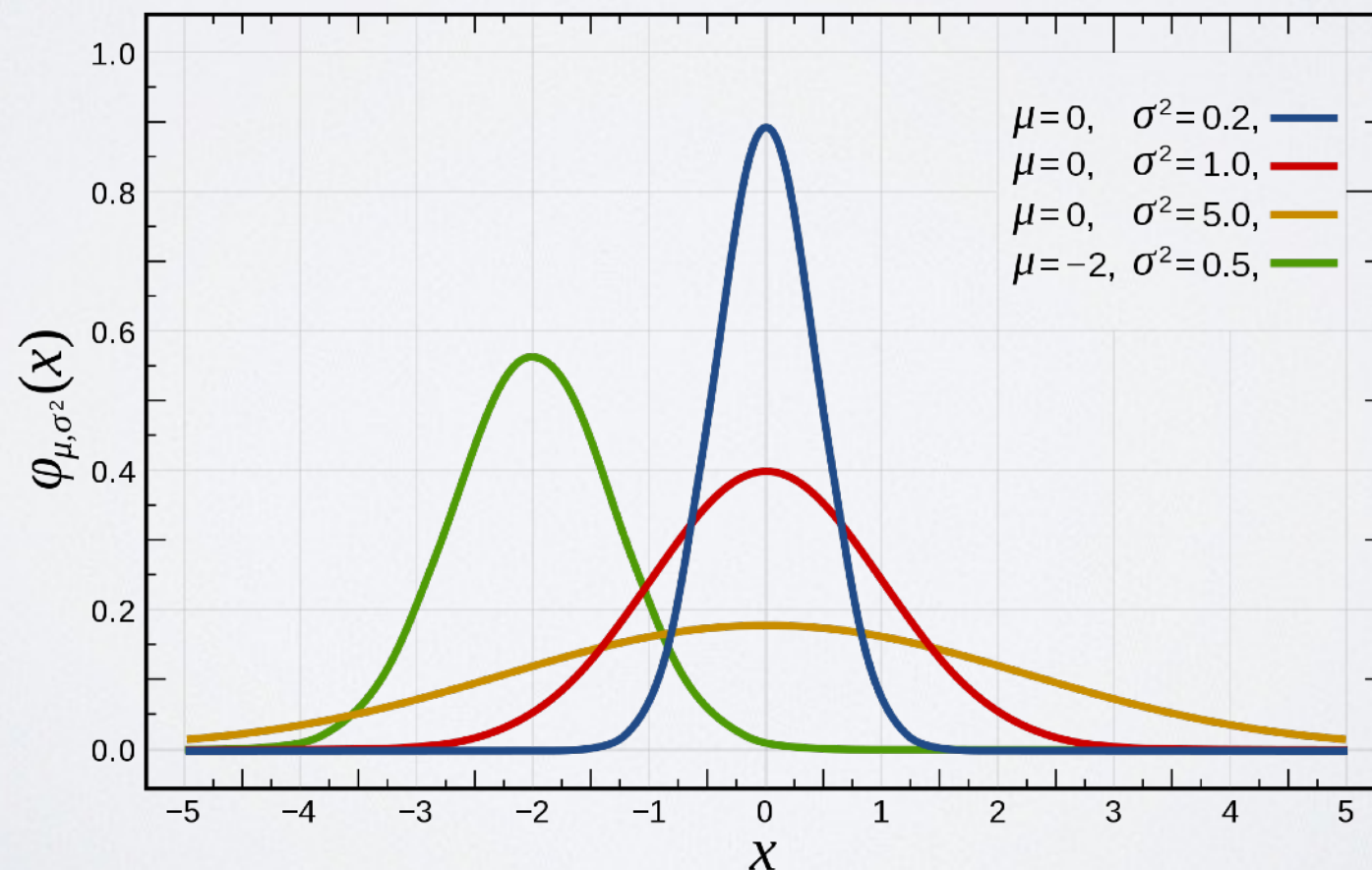


- \Rightarrow 25 observations in the interval $(13, 17]$

- Raw values for a sample,
- or fraction
 - 0.25
 - 25%
 - \Rightarrow Sum to 1. Must be inferior to 1 for any value

THEORETICAL DISTRIBUTIONS

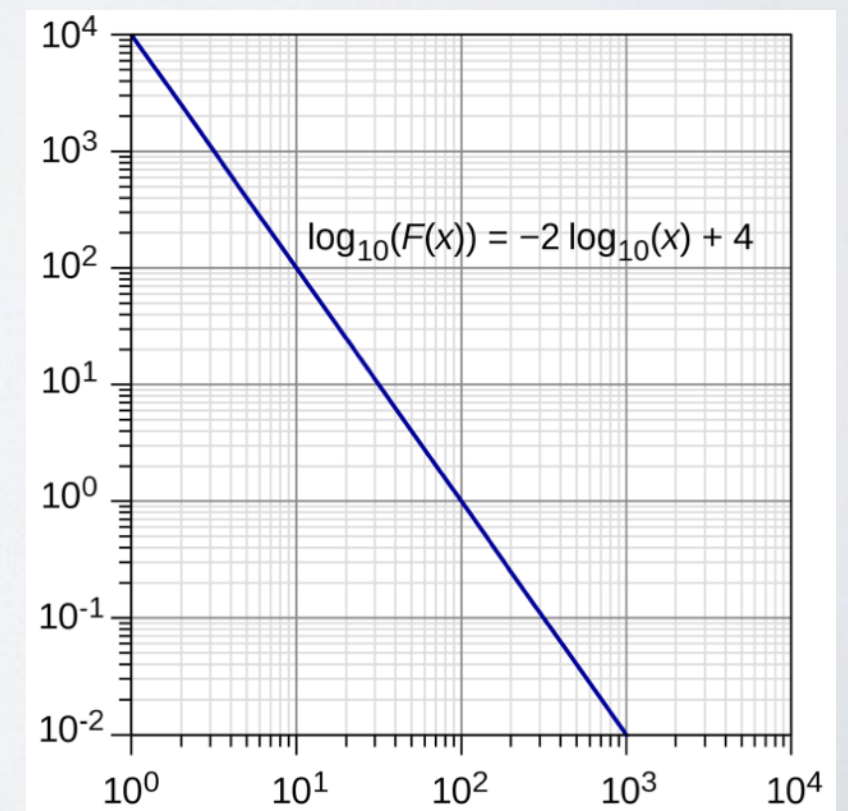
- Normal distribution
 - ▶ Many real variables follow it approximately (height, weight, price of a given product in various locations...
 - ▶ Random variations around a well-defined mean
 - ▶ Central limit theorem: average of many samples of a random variable converges to a normal distribution



THEORETICAL DISTRIBUTIONS

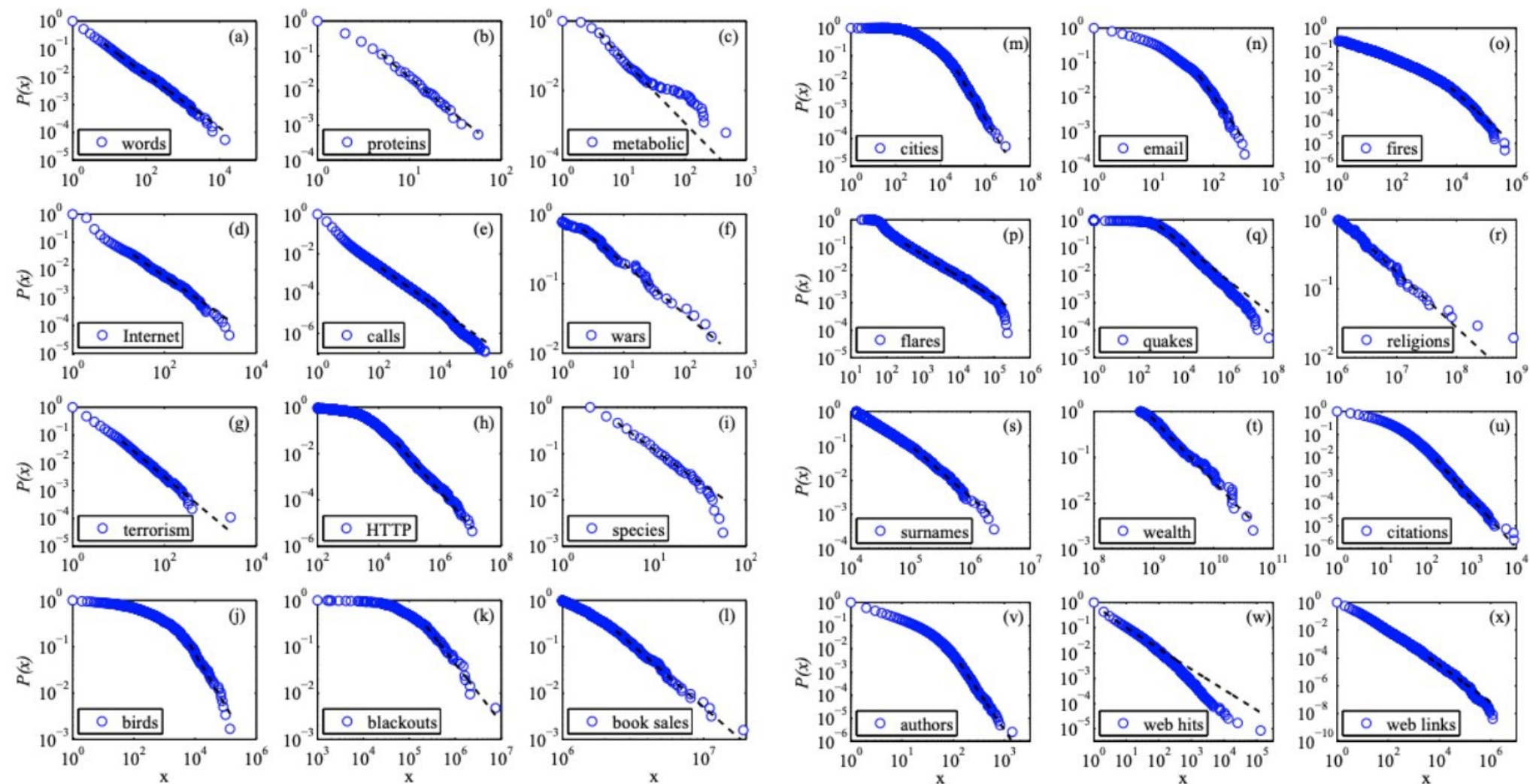
- Power Law distribution

- ▶ 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

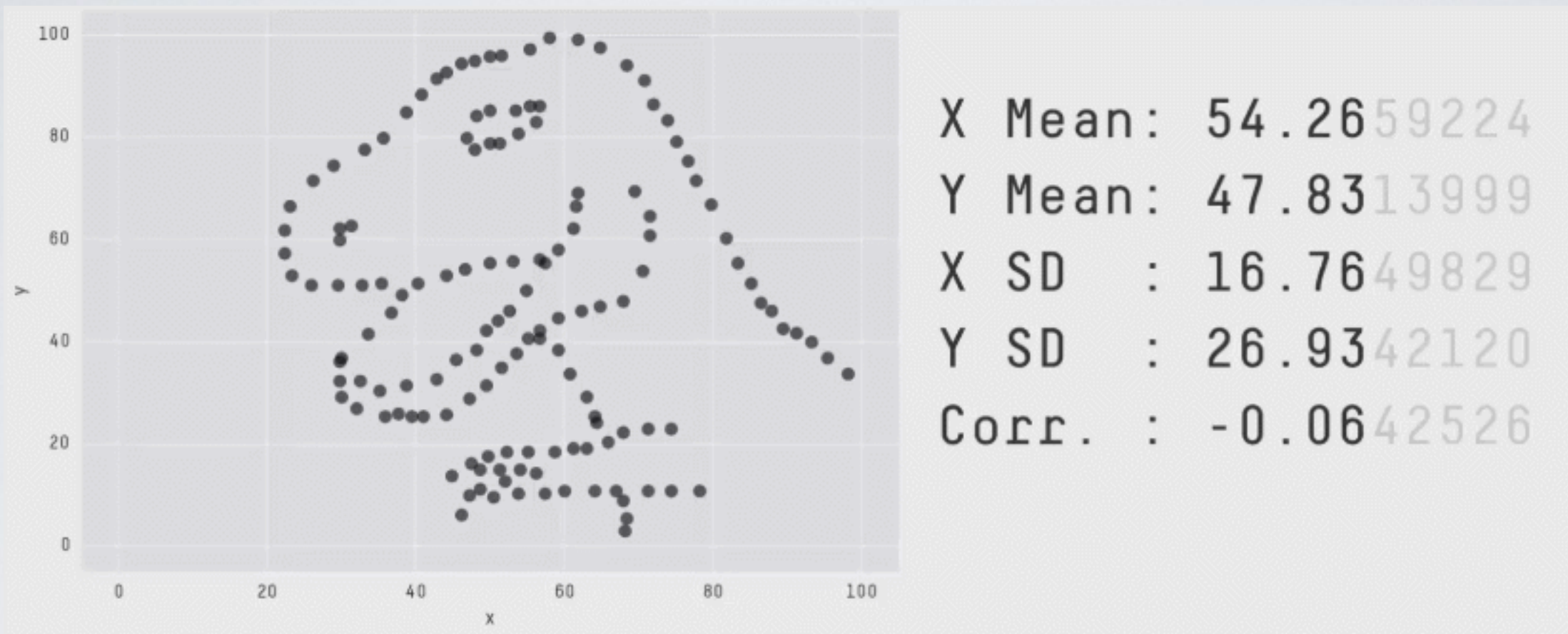


THEORETICAL DISTRIBUTIONS

- Power Law distribution



DESCRIPTIVE STATISTICS



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. Never ever.

STATISTICAL TESTS

WHAT IS IT?

- Questions such as:
 - Is my data following a normal distribution?
 - I could summarize it by mean and variance...
 - Are two variables coming from the same “population”
 - Is the probability of dying from COVID the same in two countries for 2 “identical” persons?
 - Are two variables independent?
 - “eating chocolate” and “having cancer”?
- You can use statistical tests:
 - Normality: Shapiro-Wilk, etc.
 - Categorical variables: Chi-squared χ^2 , etc.
 - Comparing distributions: Kolmogorov-Smirnov, t-test if assuming normality, etc.

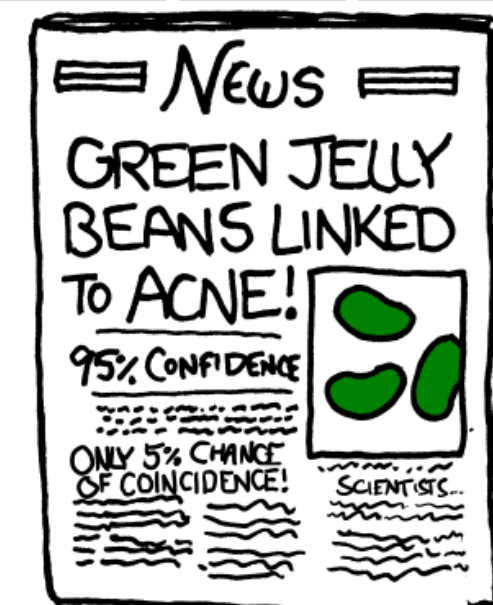
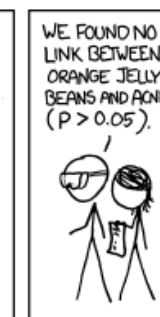
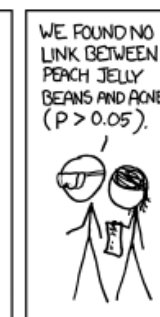
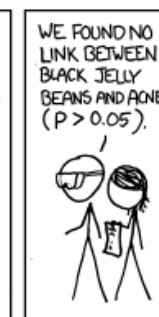
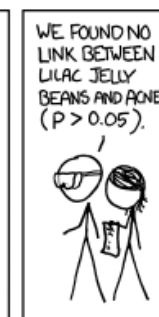
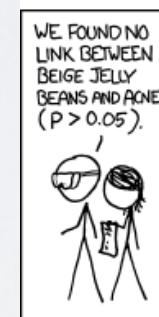
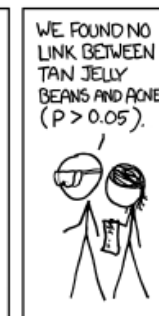
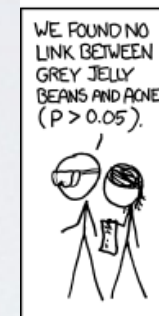
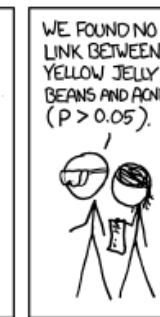
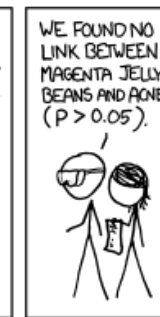
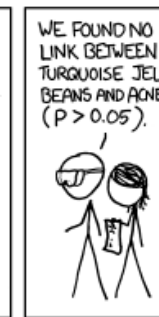
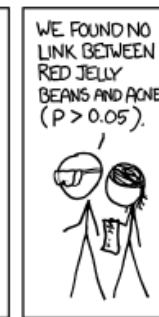
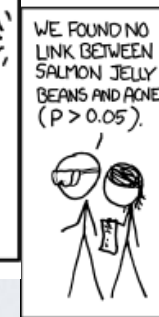
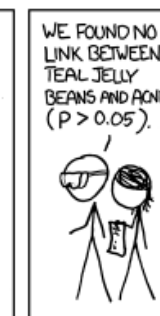
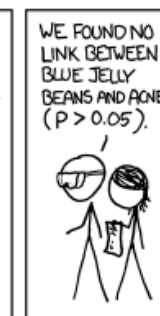
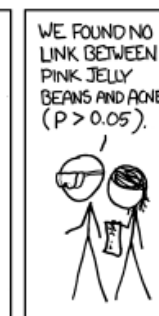
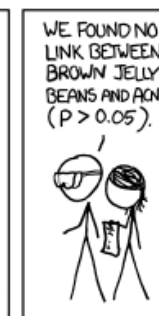
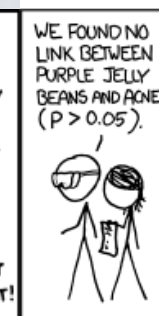
STATISTICAL TESTS

- “Can we reject the null hypothesis?”
 - p-value large \Rightarrow null hypothesis Likely True. (Probability obtain data if hypothesis True)
 - Normality test: Null hypothesis \Rightarrow distribution is normal.
 - Hypothesis testing: Null hypothesis \Rightarrow No relation between variables of interest

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
 - Nothing is exactly normal
 - No pairs of populations are exactly identical
 - No variables are independent
 - Having a cat and owning a SUV? Height of a person and their grades in high school? Etc.

P-VALUES



P-VALUE	INTERPRETATION
0.001	HIGHLY SIGNIFICANT
0.01	
0.02	
0.03	
0.04	SIGNIFICANT
0.049	
0.050	
0.051	ON THE EDGE OF SIGNIFICANCE
0.06	
0.07	
0.08	HIGHLY SUGGESTIVE, SIGNIFICANT AT THE $p < 0.10$ LEVEL
0.09	
0.099	
≥ 0.1	HEY, LOOK AT THIS INTERESTING SUBGROUP ANALYSIS

DESCRIPTIVE STATISTICS

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

VARIABLE INTERACTIONS

COVARIANCE MATRIX

Covariance Matrix Formula



- Covariance matrix **K**

- Extension of Variance to multivariate data
- $\text{Var}(X) = E[(X - \mu)^2]$
- $\text{cov}(\mathbf{X}, \mathbf{Y}) = \mathbf{K}_{\mathbf{XY}} = E[(\mathbf{X} - E[\mathbf{X}])(\mathbf{Y} - E[\mathbf{Y}])^T]$
 - 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
- $\Rightarrow \text{cov}(\mathbf{X}, \mathbf{X}) = \text{Var}(\mathbf{X})$

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

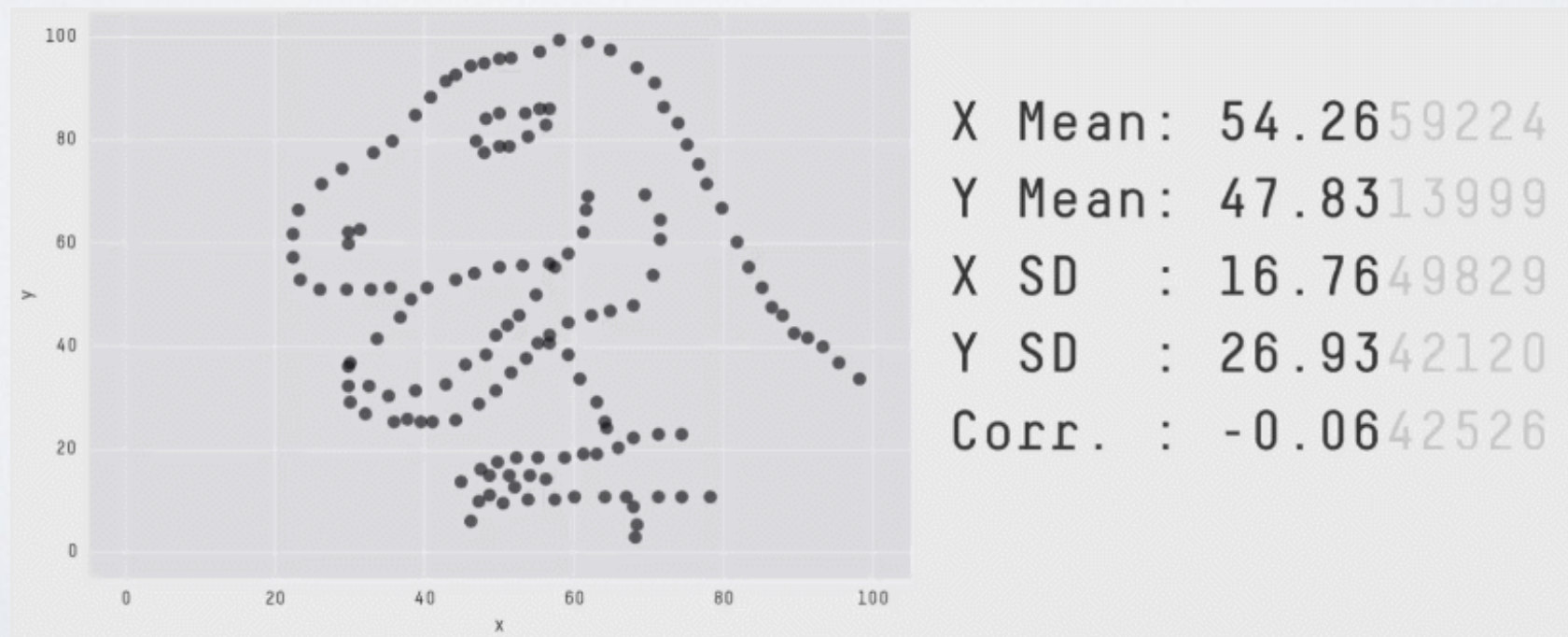
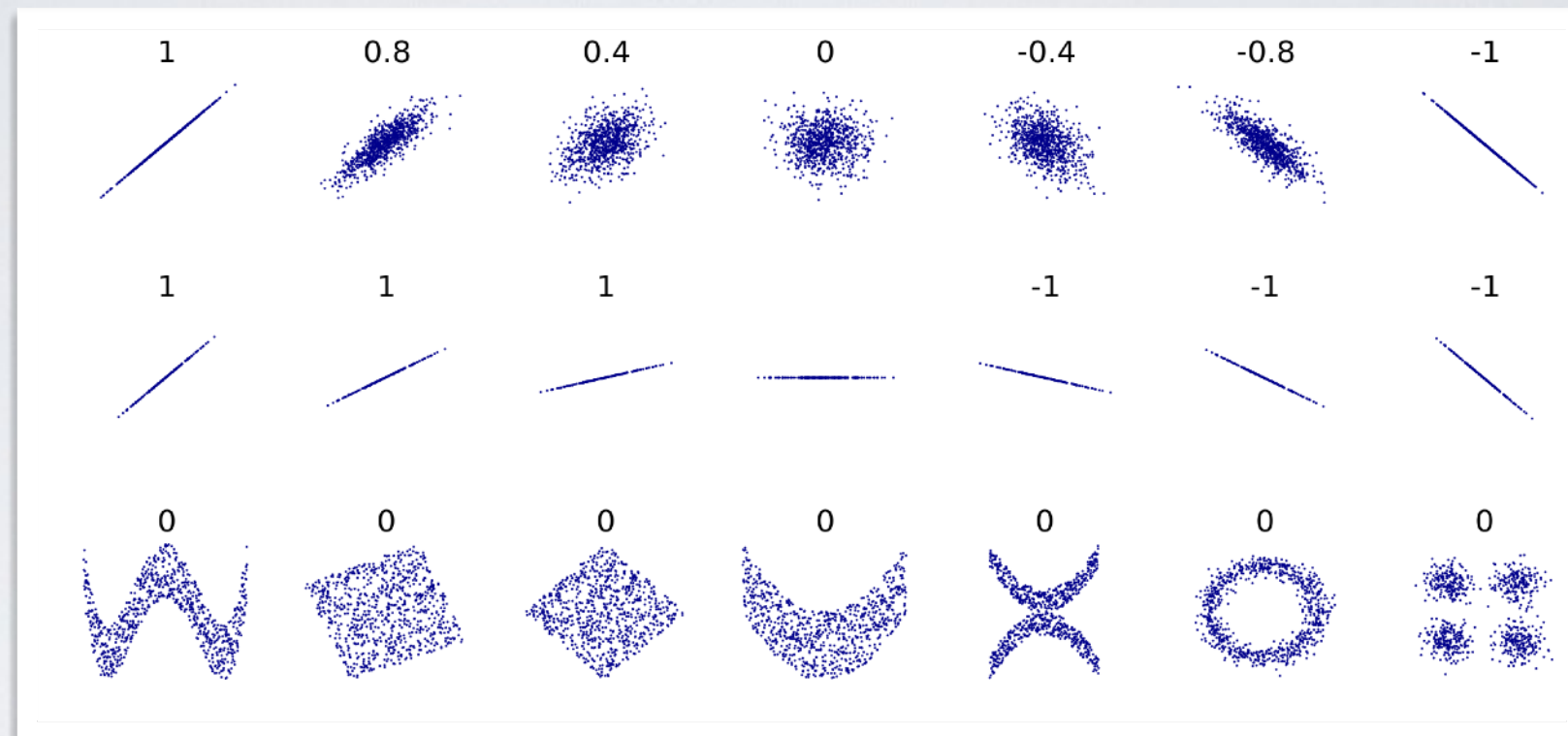
- Covariance is hardly interpretable by itself.

- 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

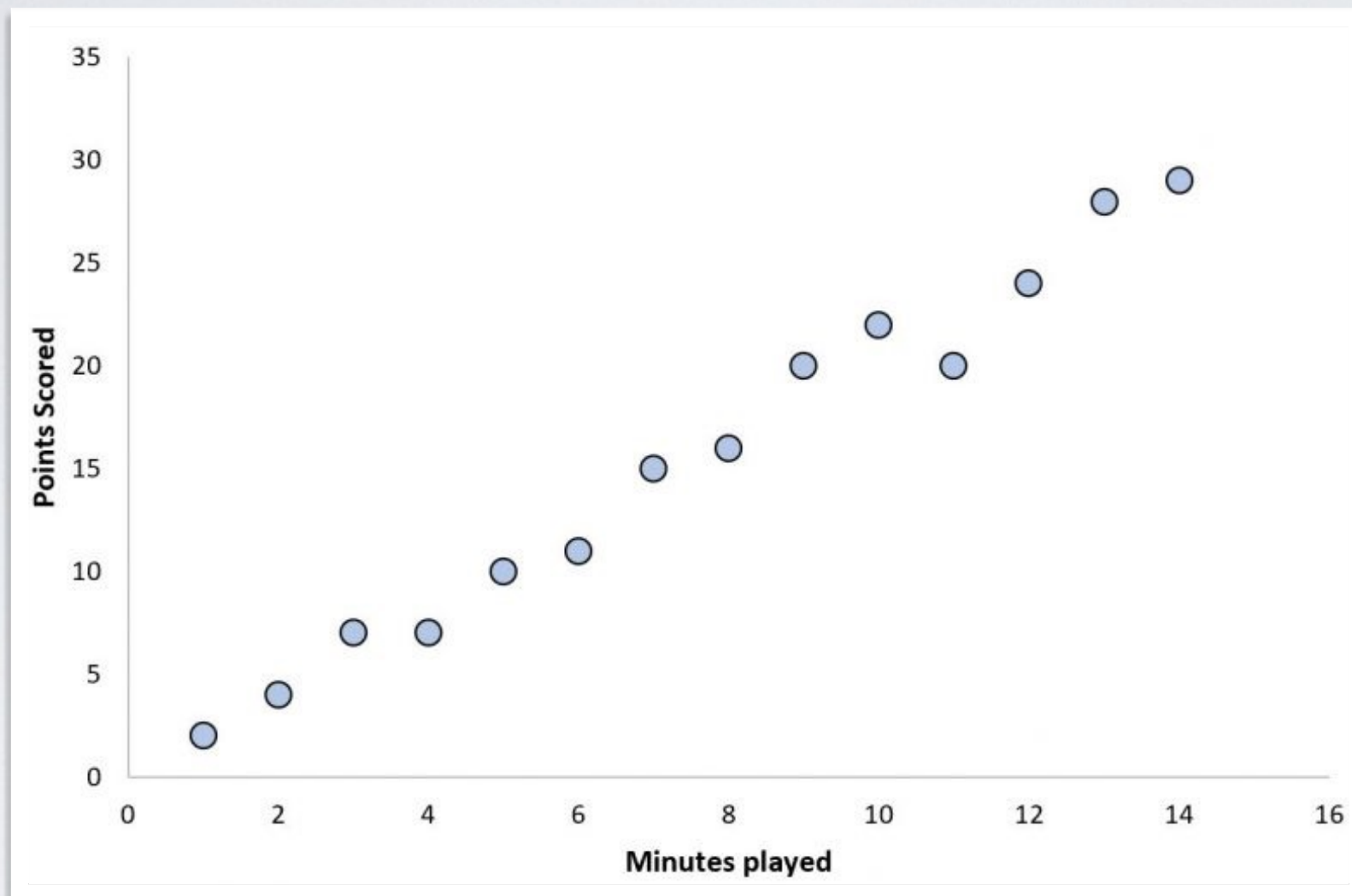
CORRELATION COEFFICIENT



CORRELATION COEFFICIENT

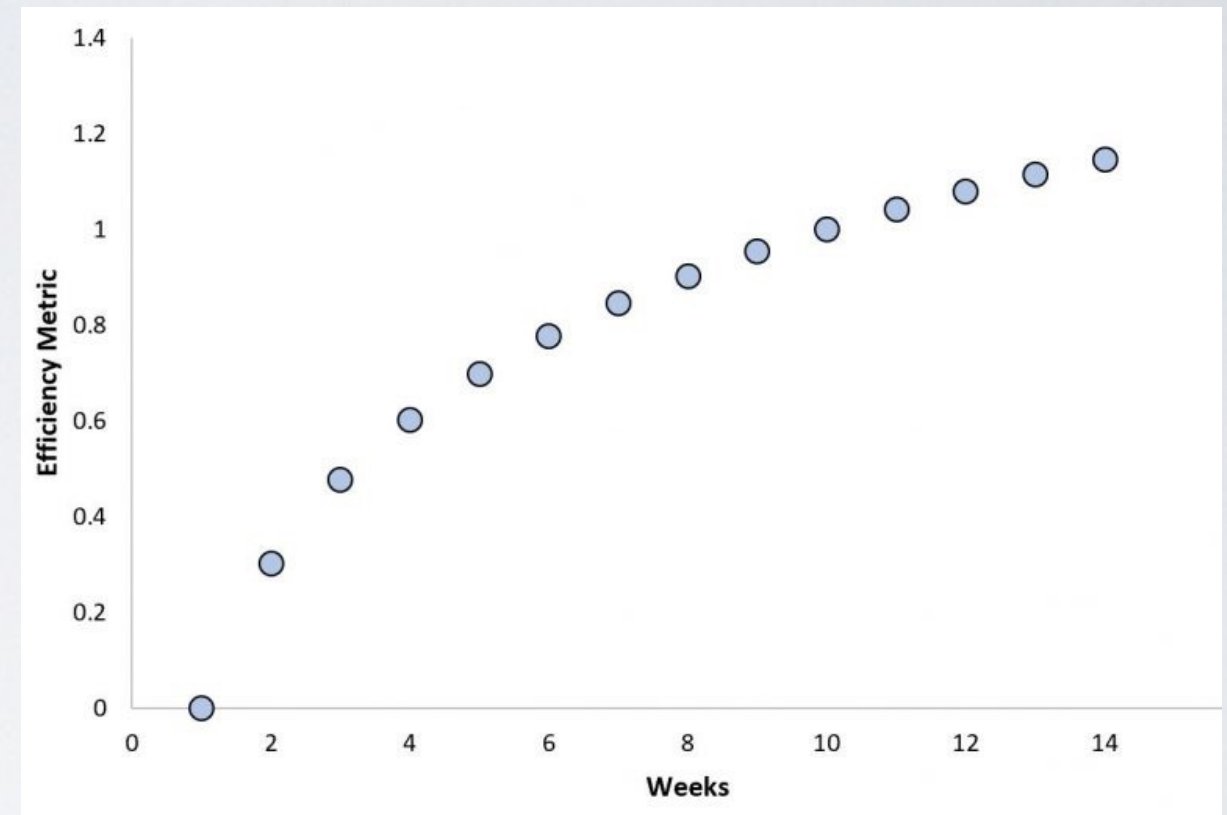
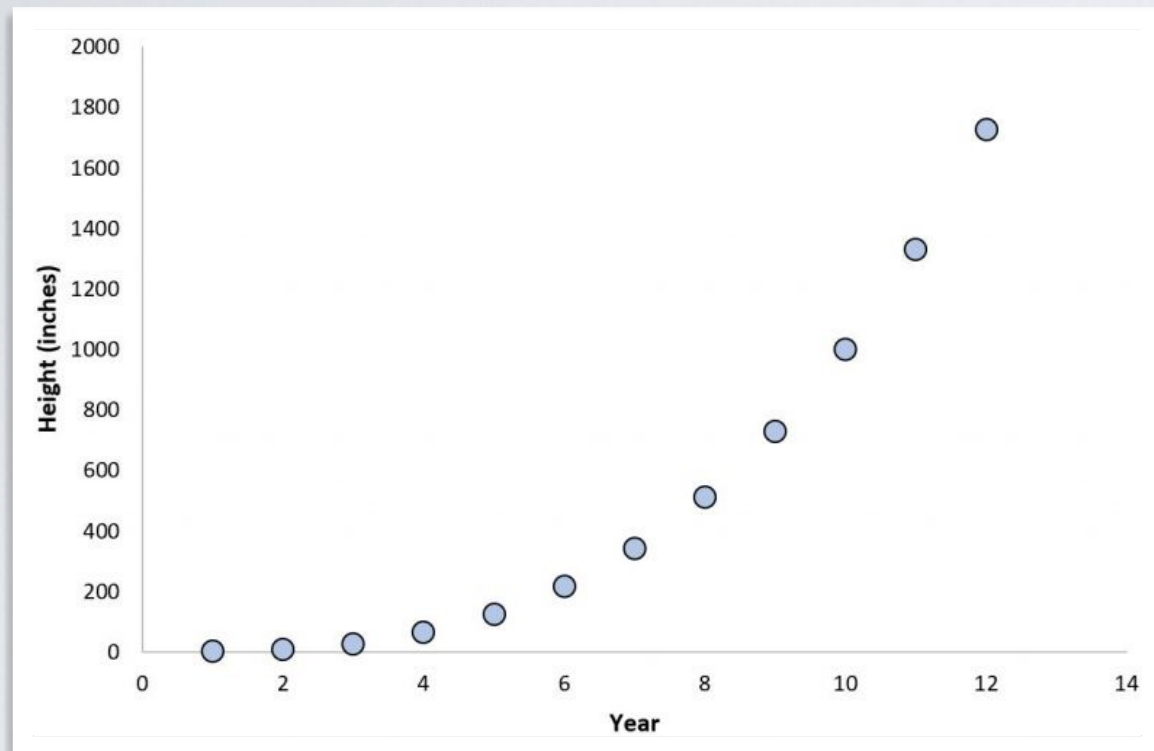
- Other possible interpretation, e.g.
 - Cosine similarity of the vectors defined by the observations...
- 0.7 ? Is it a high or low value ?
 - It depends.

NONLINEAR RELATIONSHIPS



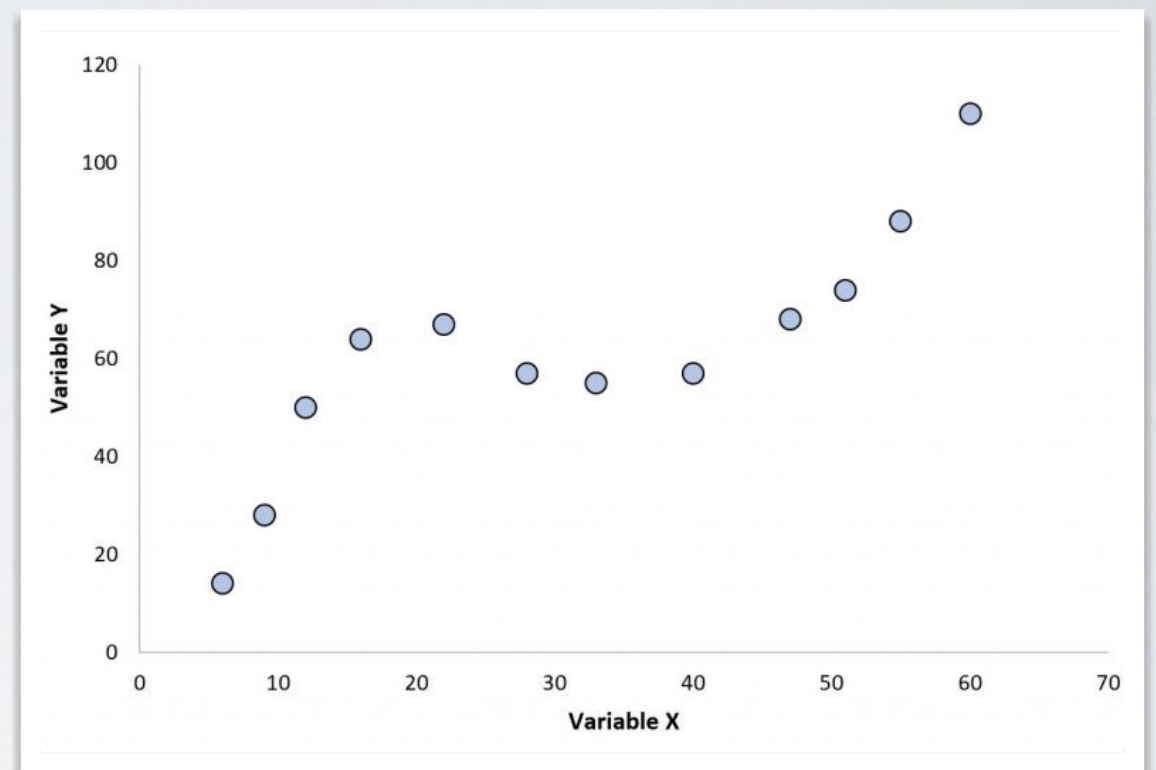
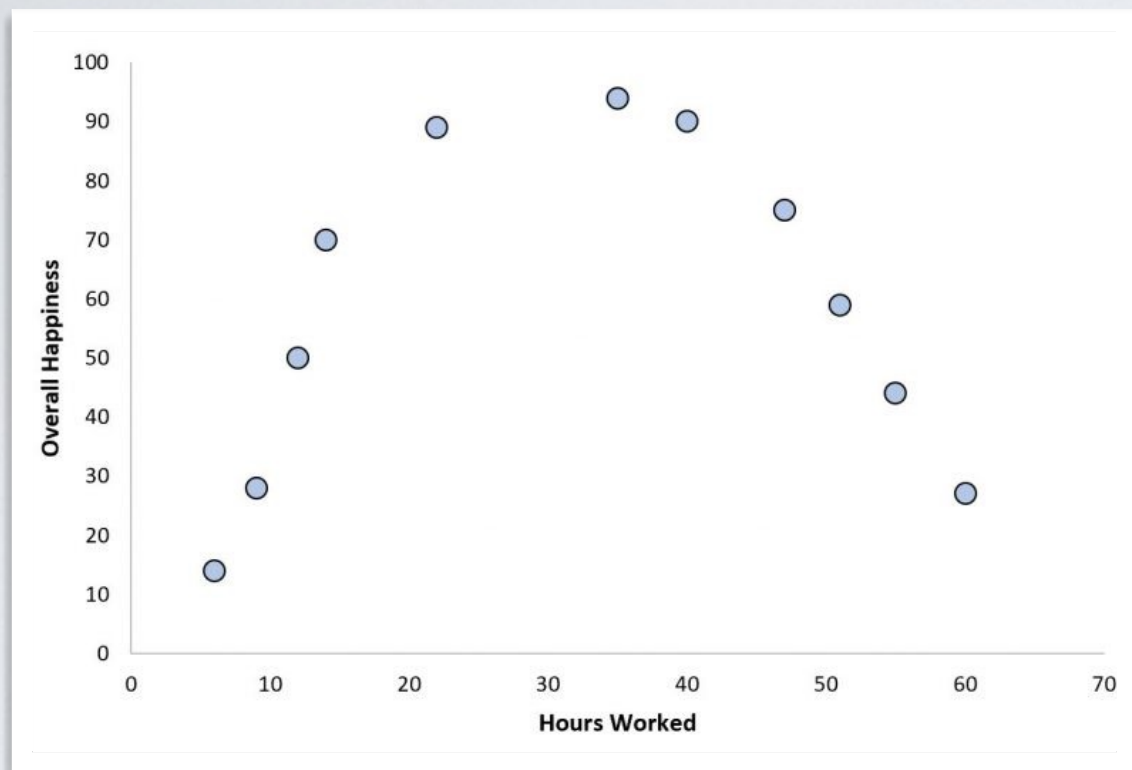
Linear relationship
 $Y = a + bX + e$

NONLINEAR RELATIONSHIPS

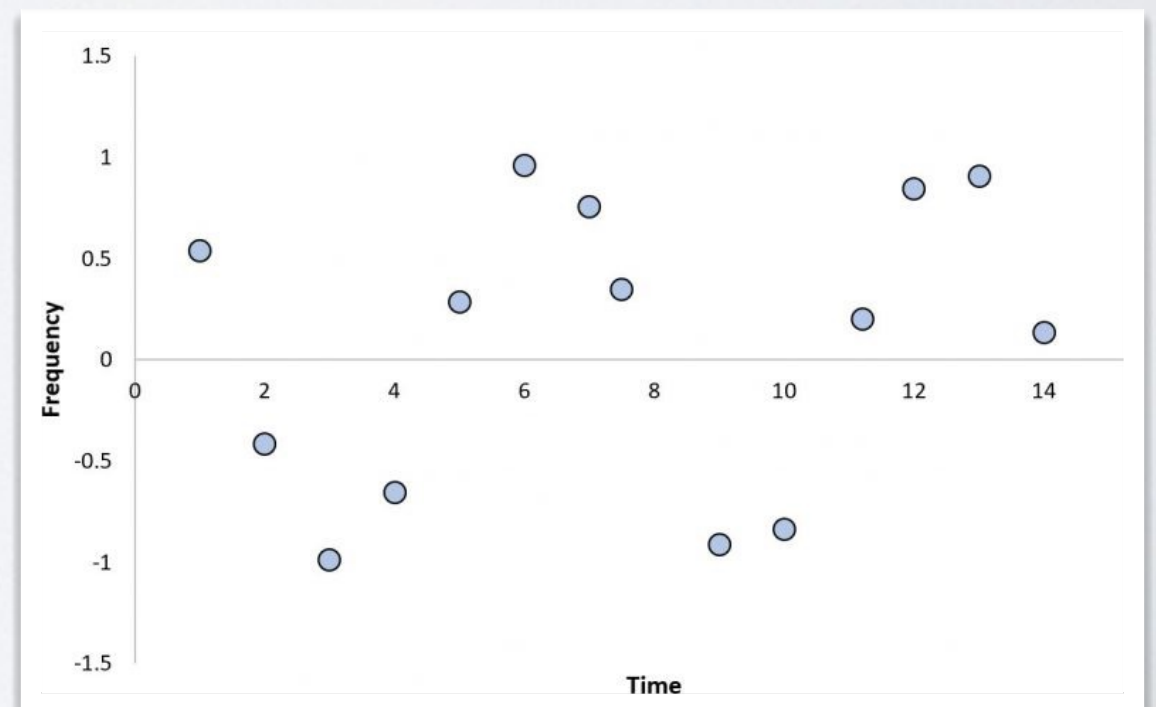


Monotonous, non-linear

NONLINEAR RELATIONSHIPS



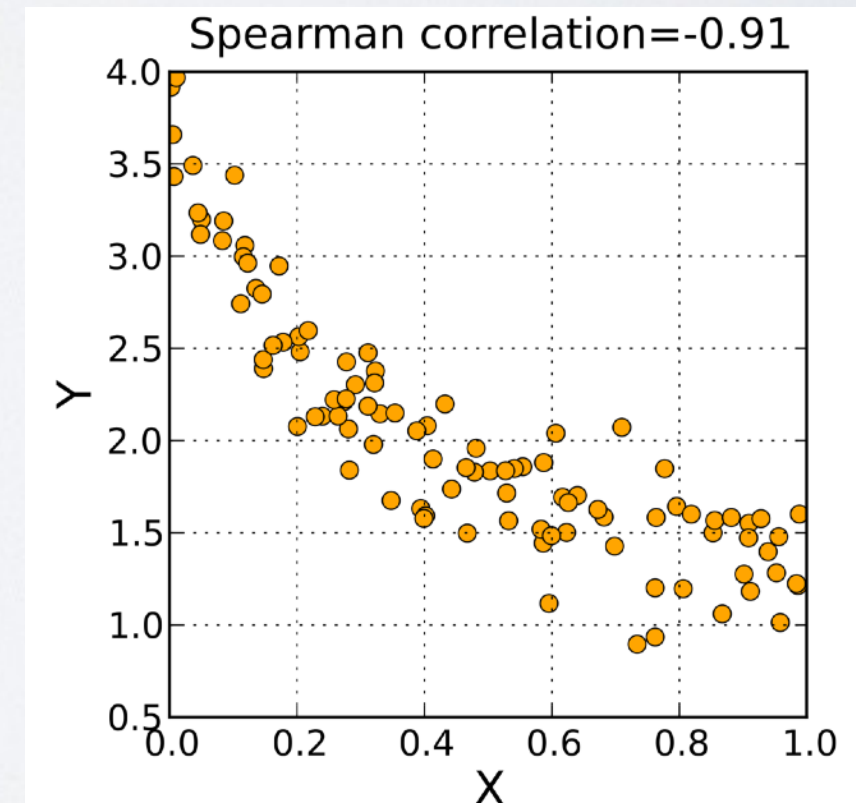
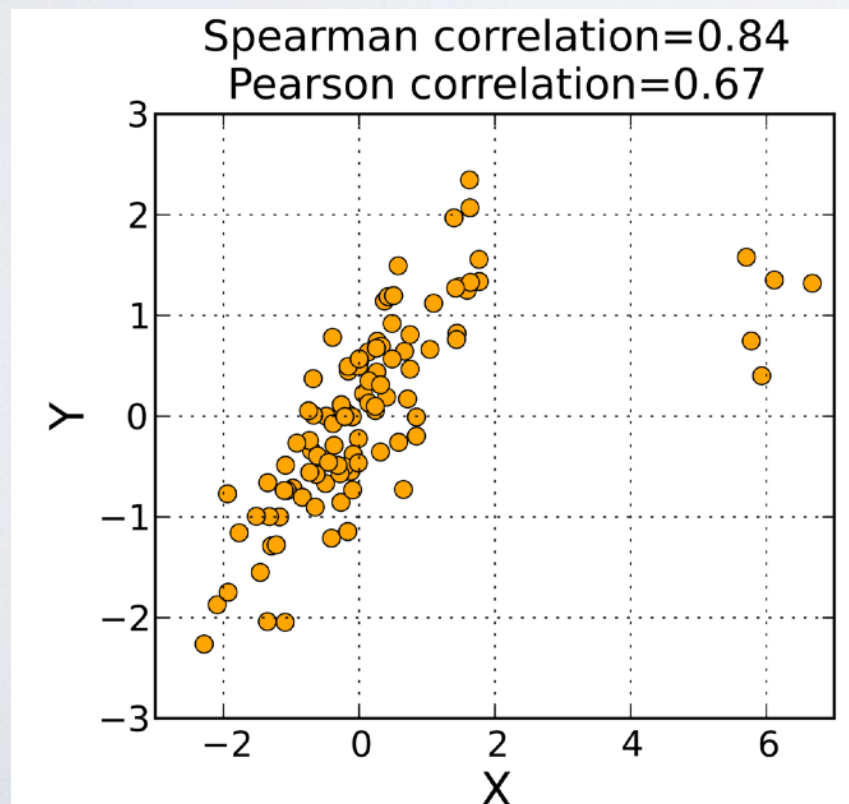
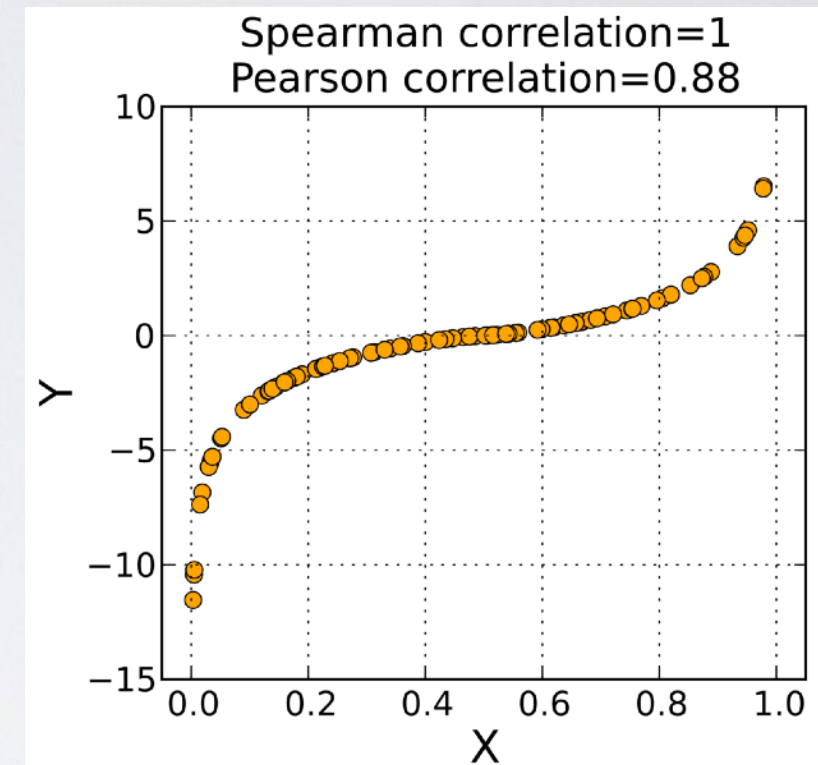
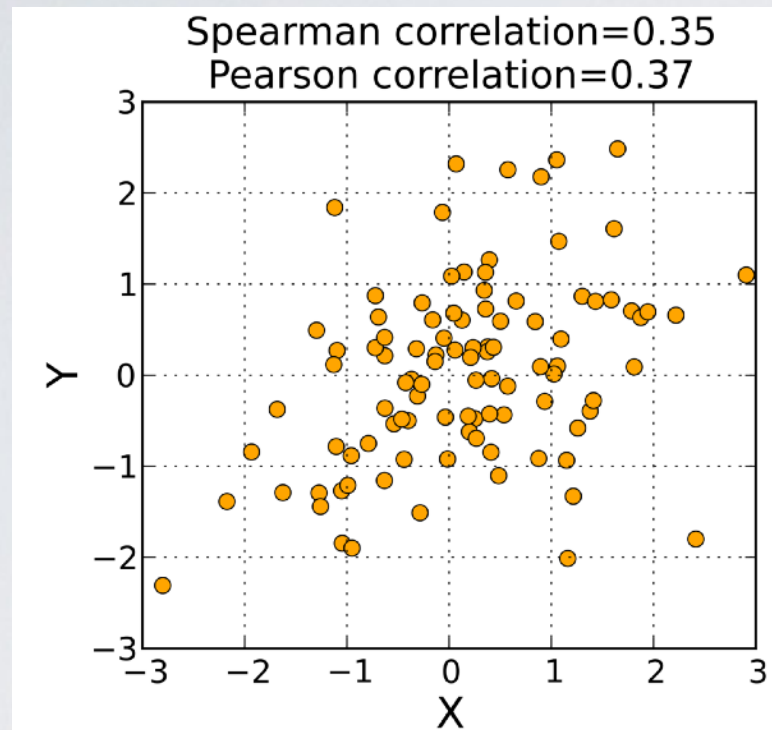
Non-monotonous,
Non-linear



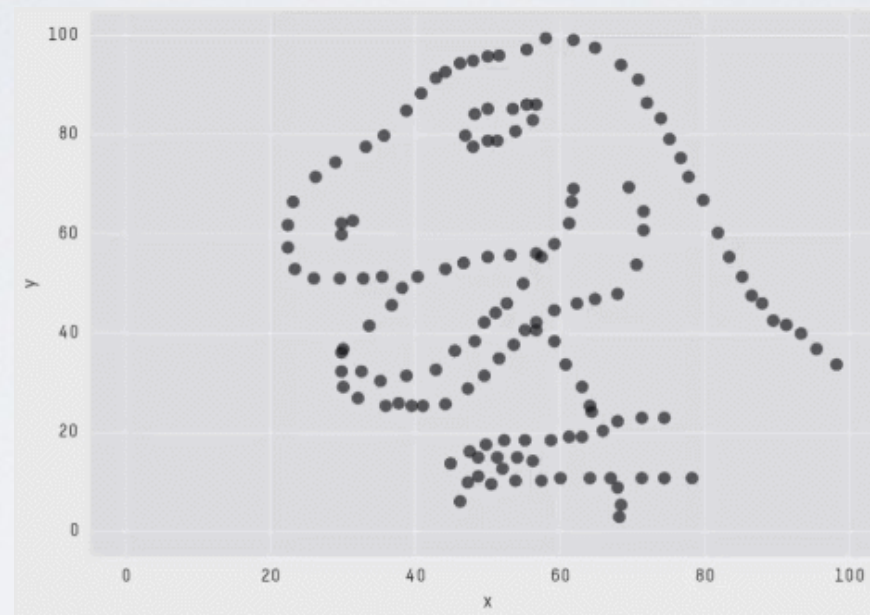
SPEARMAN'S CORRELATION

- 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



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

FEATURE SCALING

FEATURE SCALING: WHY



Y

Age: 20

Age: 20

Age: 90

m

Height: 1.82

Height: 1.82

Height: 1.50

g

Weight: 80 000

Weight: 81 000

Weight: 81 000

FEATURE SCALING: WHY

- We want to use euclidean distance to compute the “distance” between 2 people
 - $a = (y:20, m:1.82, g:80\ 000)$, $b = (y:20, m:1.82, g:81\ 000)$, $c = (y:90, m:1.50, g:80\ 020)$
 - **$d(a,b)=1000.0005$**
 - **$d(a,c)=72.8$**
 - 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

WARNING

- There is no magic recipe!
- Everything cannot be normalized
 - Percentage scores, grades, scores between 0 and 1 ...
 - You would make low values big.
 - Binary variables (one hot encoded or not)
 - Careful with variables having an “absolute meaning”
 - Number of observations, duration of time, positions, distances, etc.

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