

# DATA - INTRODUCTION

# WHO AM I

- Rémy Cazabet ([remy.cazabet@univ-lyon1.fr](mailto:remy.cazabet@univ-lyon1.fr))
- Class page: <http://cazabetremy.fr/Teaching/DISS/ML.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

# DEFINITION

- **Data mining** is the process of extracting and discovering patterns in large data sets involving methods at the intersection of machine learning, statistics, and database systems.[Wikipedia]
- Rather vague term

Data mining involves six common classes of tasks:<sup>[5]</sup>

- **Anomaly detection** (outlier/change/deviation detection) – The identification of unusual data records, that might be interesting or data errors that require further investigation.
- **Association rule learning** (dependency modeling) – Searches for relationships between variables. For example, a supermarket might gather data on customer purchasing habits. Using association rule learning, the supermarket can determine which products are frequently bought together and use this information for marketing purposes. This is sometimes referred to as market basket analysis.
- **Clustering** – is the task of discovering groups and structures in the data that are in some way or another "similar", without using known structures in the data.
- **Classification** – is the task of generalizing known structure to apply to new data. For example, an e-mail program might attempt to classify an e-mail as "legitimate" or as "spam".
- **Regression** – attempts to find a function that models the data with the least error that is, for estimating the relationships among data or datasets.
- **Summarization** – providing a more compact representation of the data set, including visualization and report generation.

# DEFINITION

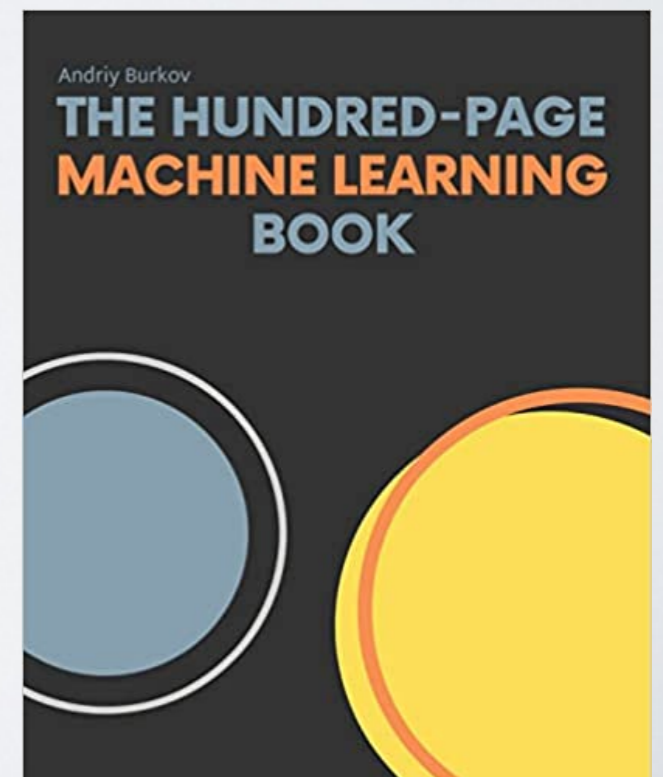
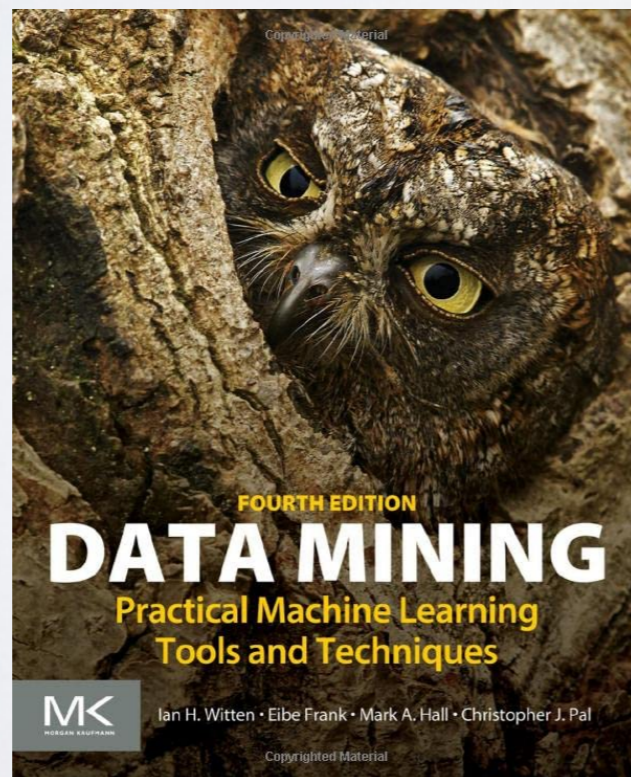
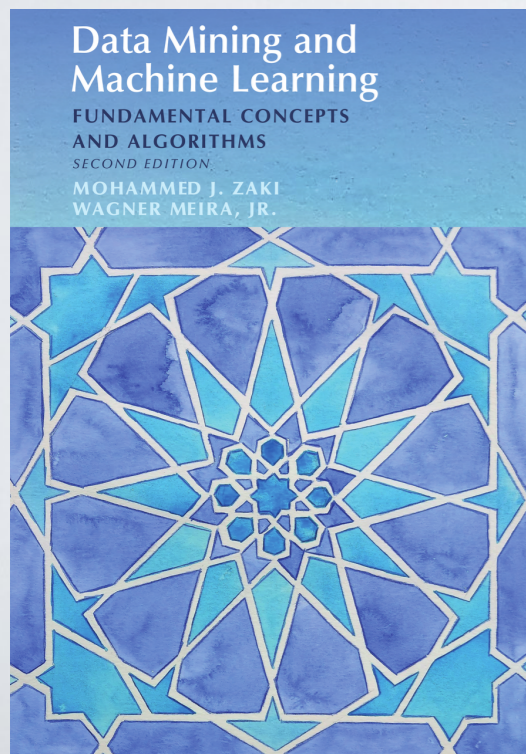
- In these Masters, my understanding
  - Making sense of data without prior knowledge
  - => Unsupervised
  - Discover what is hidden in the data
- Topics
  - Data
  - Clustering
  - Dimensionality reduction
  - Recommendation
  - Frequent Patterns
  - ?

# THIS CLASS

- Less math
  - Math are everywhere in ML. But most of it is applying simple math.
  - If you need to understand the hard one, it is simpler to take a book.
- More intuition
  - I want you to understand the large picture. You can focus on what you like.
- No learning by heart
  - Remember that the concept exist, so you can google it.
- And some practice
  - Huge amount of resources available for free

# THIS CLASS

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



# 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], Income [0-10k], [10k-15k], [15k-50k]...
- 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

# RATIO

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

# 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
- Directions expressed as angle (north+85°)
- 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

# RELATIVE AND ABSOLUTE

- The world values are divided in two types of things, and two types of interpretation
- Is there a larger difference between two persons:
  - Age 1 / 5? Revenue 1000€/1500€ ?
  - Age 91 / 95? Revenue 10 000 € / 10500 € ?
- Think about it:
  - In country 1, average salary is 100\$, p1 salary is 1000\$
  - In country 2, average salary is 1000\$, p2 salary is 2000\$
  - Should you consider that
    - p1 is well paid (10x average salary VS 2x for p2)
    - They are paid the same (1000\$ différence)



# RELATIVE AND ABSOLUTE

- If your values are expressed in absolute terms, but you think their interpretation should be in relative terms, you can transform them using the *log scale*: e.g.,
  - ▶ In  $\log_2$ , going from  $x$  to  $x+1$  means multiplying by 2.
  - ▶ In  $\log_{10}$ , going from  $x$  to  $x+1$  means multiplying by 10
- e.g, should you express earthquake strength in:
  - ▶ Energy released
  - ▶ Richter Magnitude (log scale)
  - ▶ Depends if you care about comparing small and large earthquakes (relative) or large from superlarge (raw scale => all those small ones look the same from there)

# CONTINUOUS OR DISCRETE ?

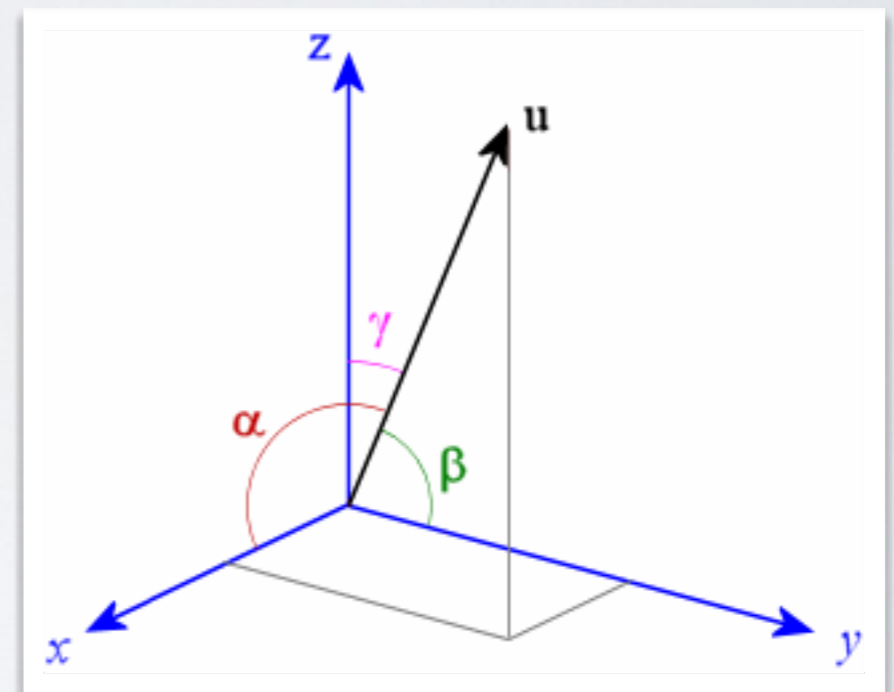
- Real data is never mathematically continuous (limit of precision in computers...)
- A dataset must be treated as continuous when the number of possible values is at least in the order of the number of observations.
- Think of plotting an histogram distribution...

# 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

# UNIVARIATE / MULTIVARIATE

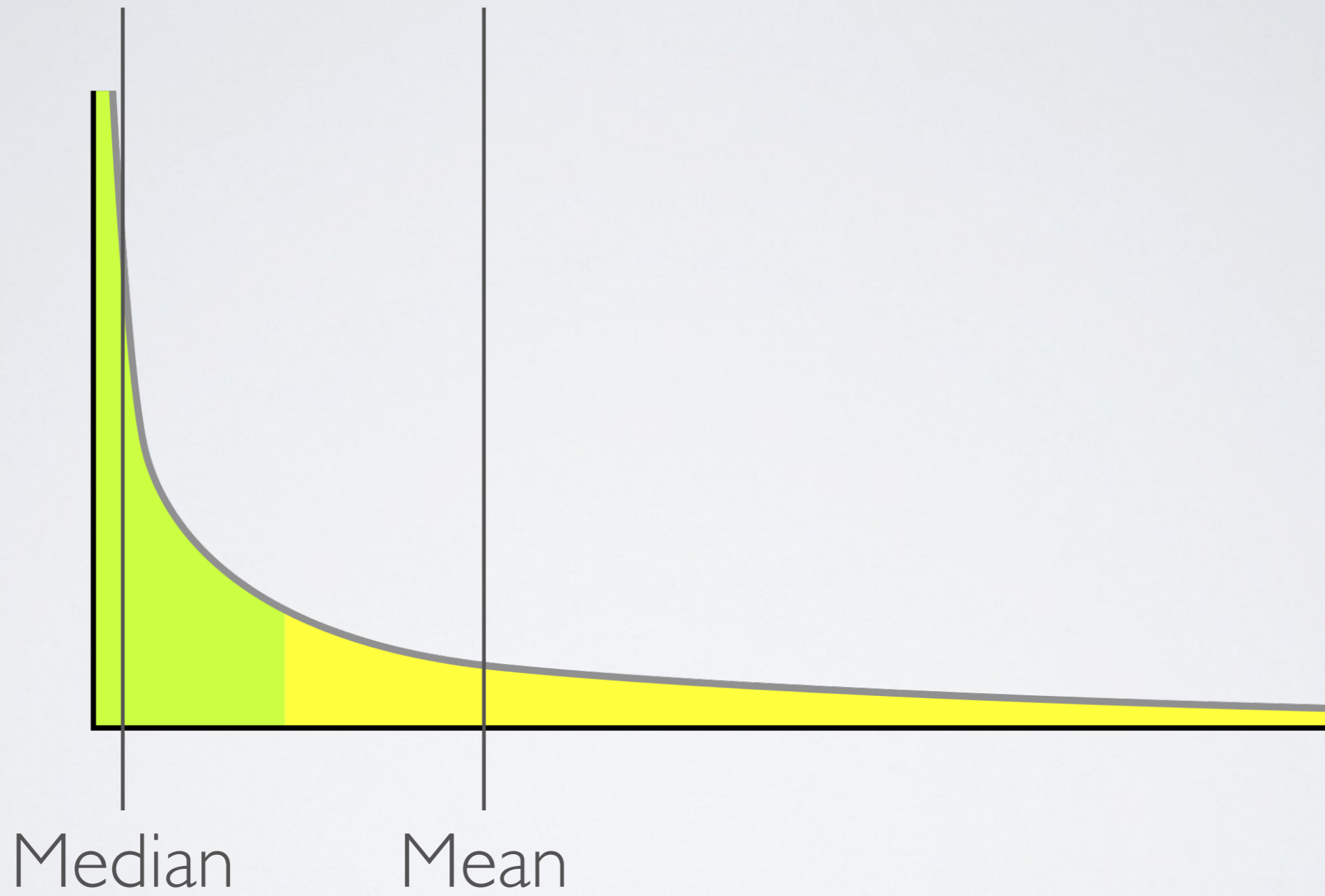
- Single *feature*: univariate
  - Age
- Real life: multivariate.
  - 2D (age, weight)
  - 3D (age, weight, height)
  - 4D (age, weight, height, genre)
  - ...



# DESCRIBING A VARIABLE

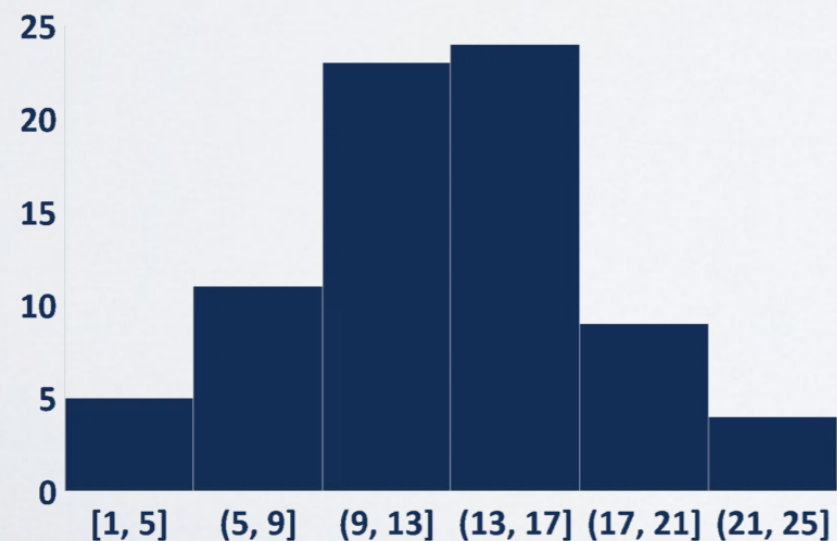
# DESCRIBING VALUES

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

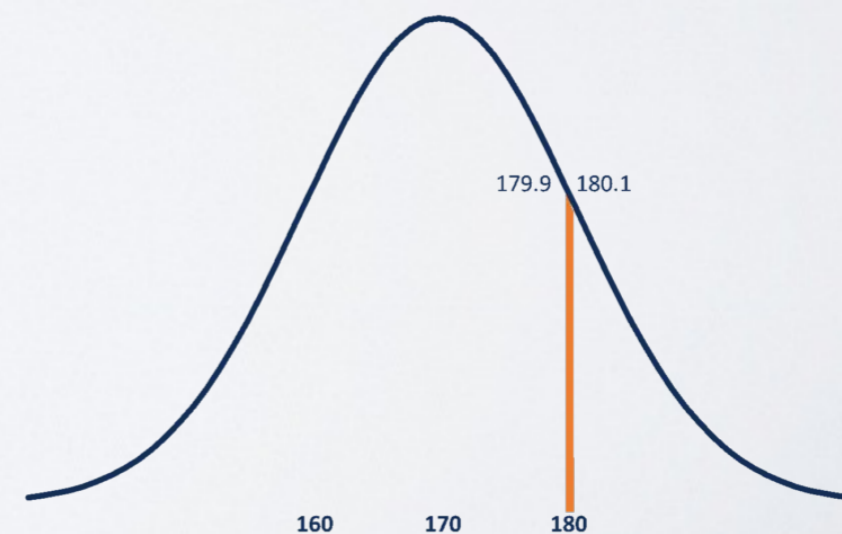


# DISTRIBUTION

- 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

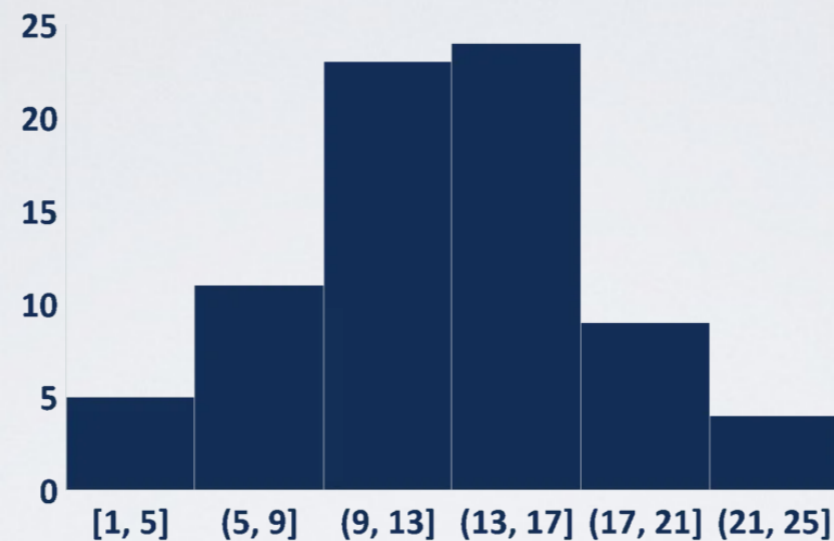


## Continuous Distribution





# DISTRIBUTION (DISCRETE)

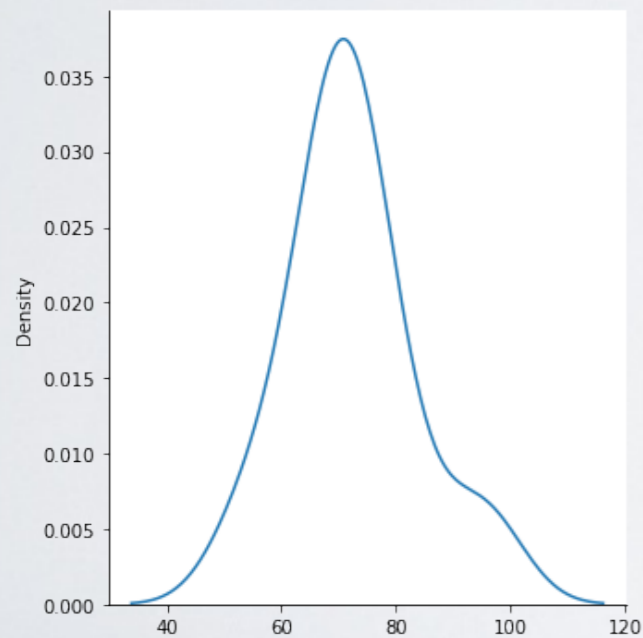


- $\Rightarrow$  25 observations in the interval (13, 17]
- Raw values for a sample,
- or fraction
  - 0.25
  - 25%
  - $\Rightarrow$  Sum to 1. Must be inferior to 1 for any value

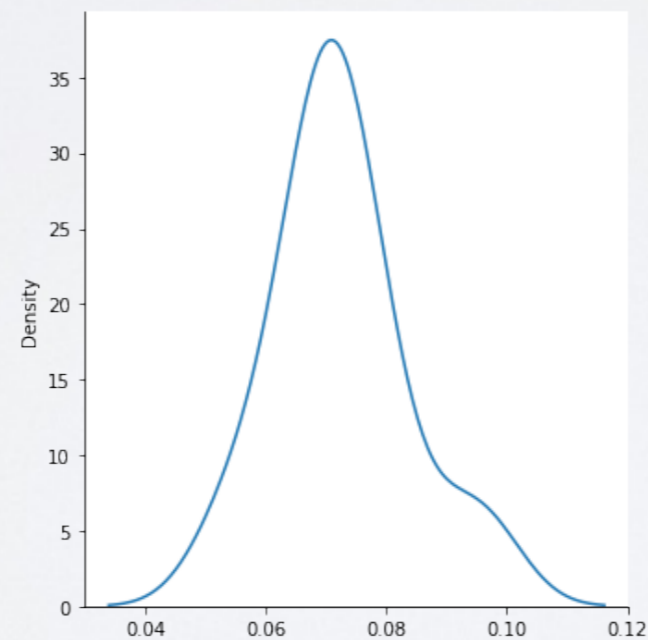
# DISTRIBUTION

- Pdf: Probability Density Function

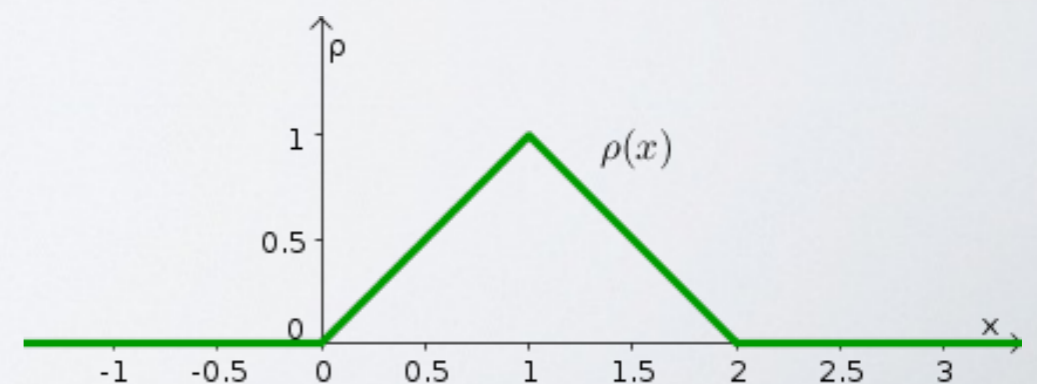
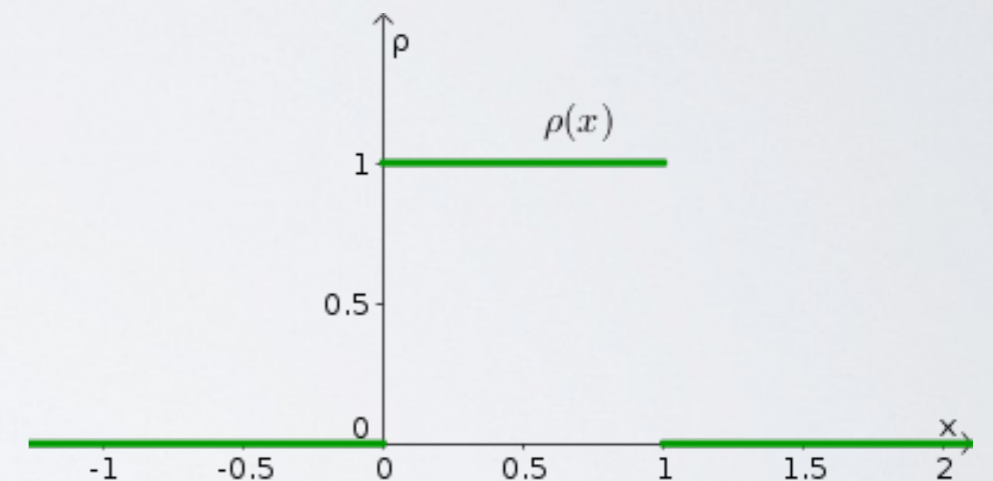
- ▶ Problem with continuous variables: If we draw at random a number  $\in [0,1]$ , the probability to be exactly any value is 0.
- ▶ Integral must be equal to 1
  - => Values can be above 1



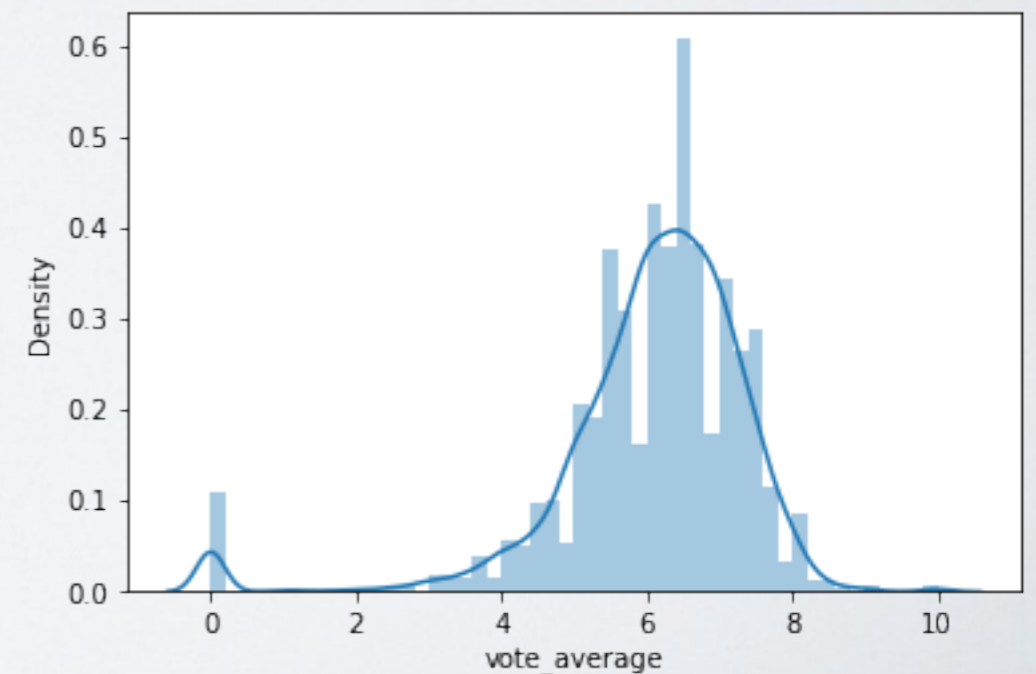
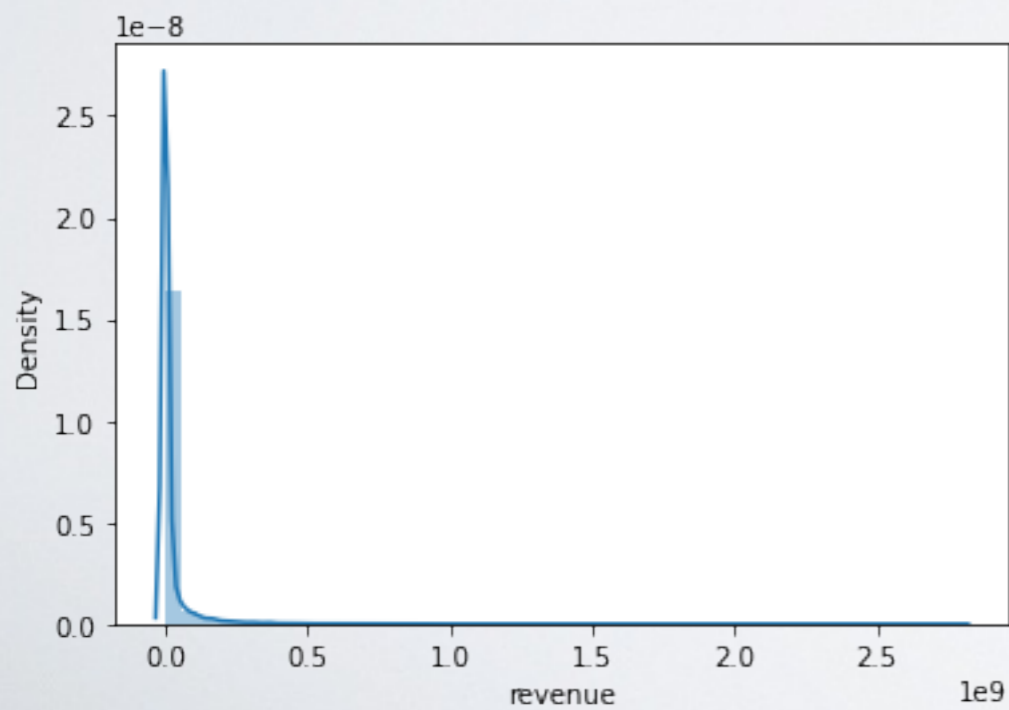
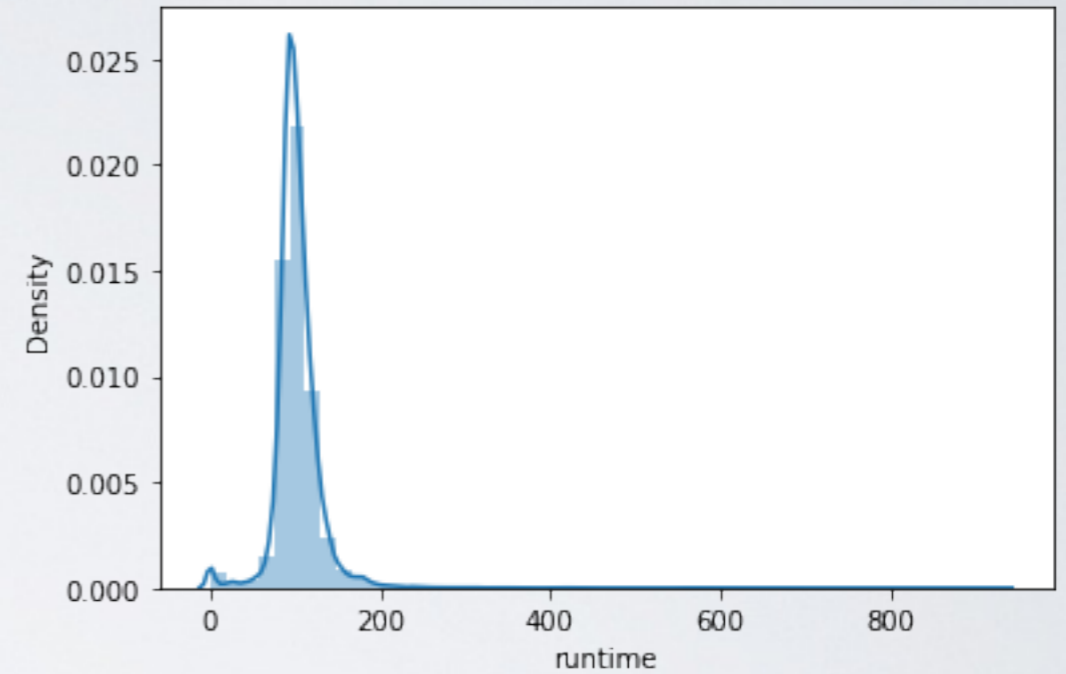
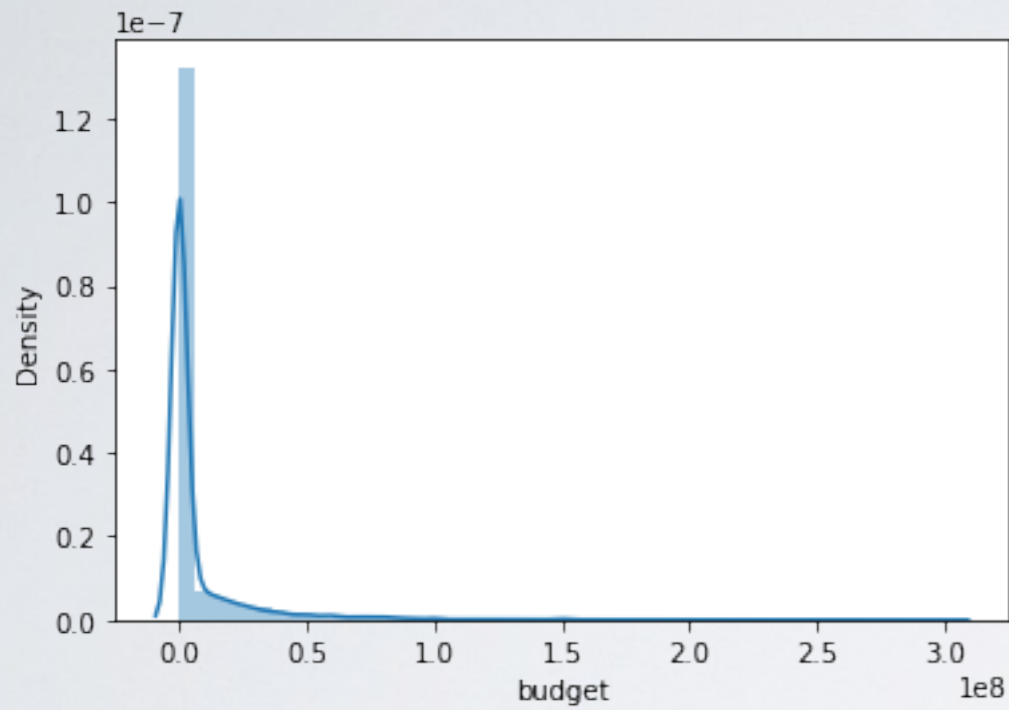
Weights in k



Weights in tons

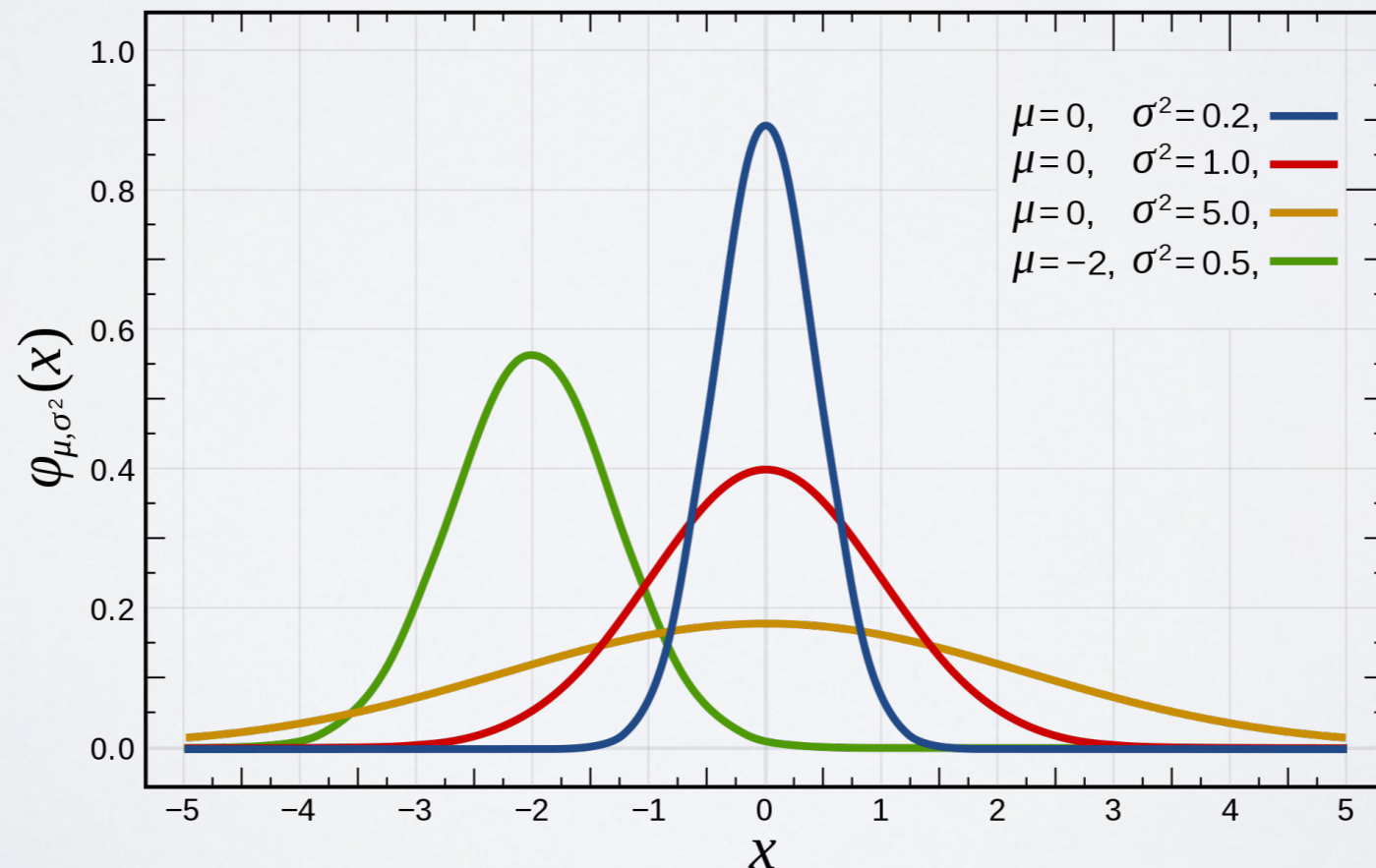


# EMPIRICAL DISTRIBUTIONS



# THEORETICAL DISTRIBUTIONS

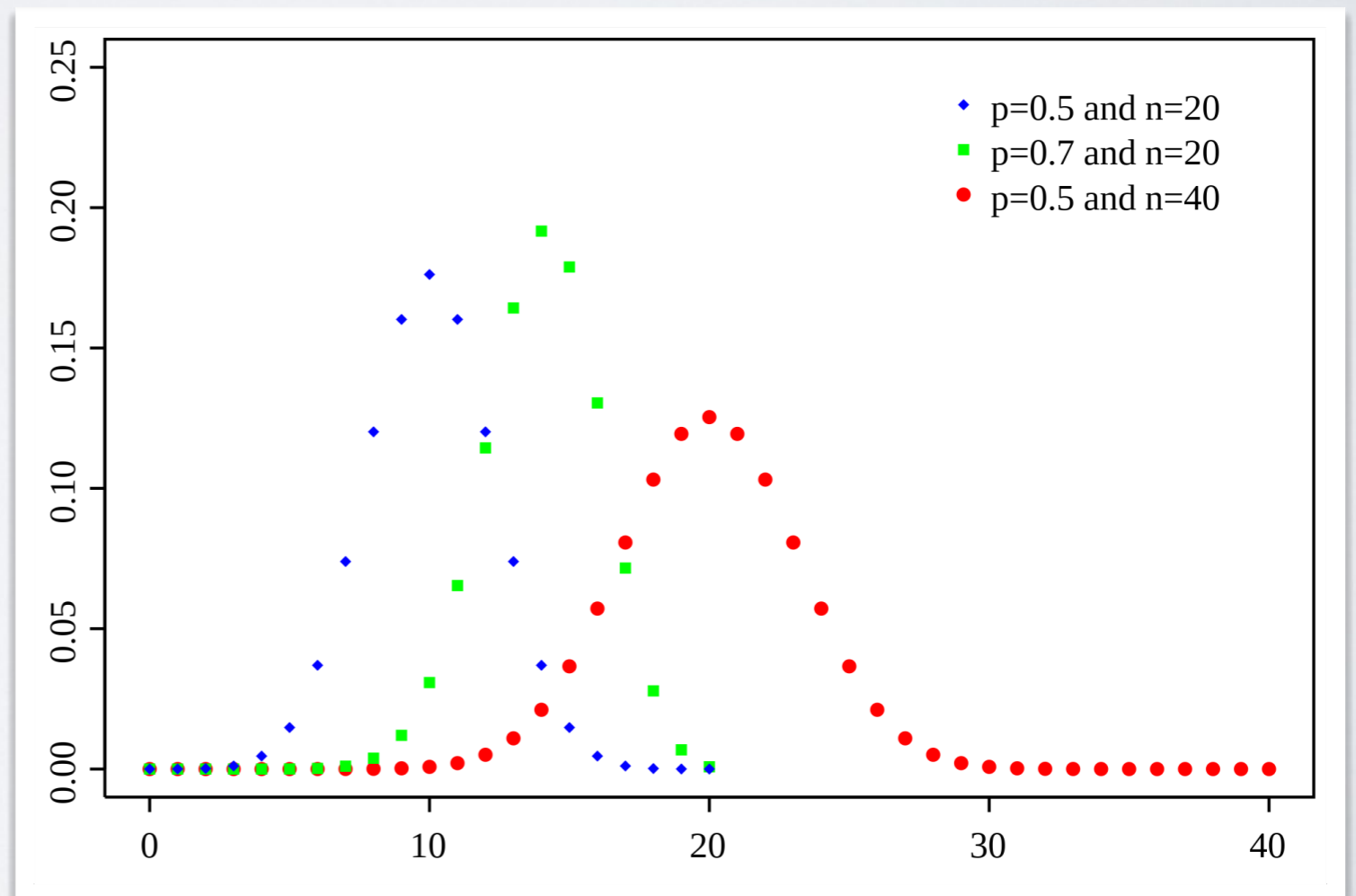
- Normal distribution
  - ▶ Many real variables follow it approximately (height, weight, price of a given product in various locations...)
  - ▶ Random variations around a well-defined mean
  - ▶ Central limit theorem: average of many samples of a random variable converges to a normal distribution



# THEORETICAL DISTRIBUTIONS

- Binomial distribution

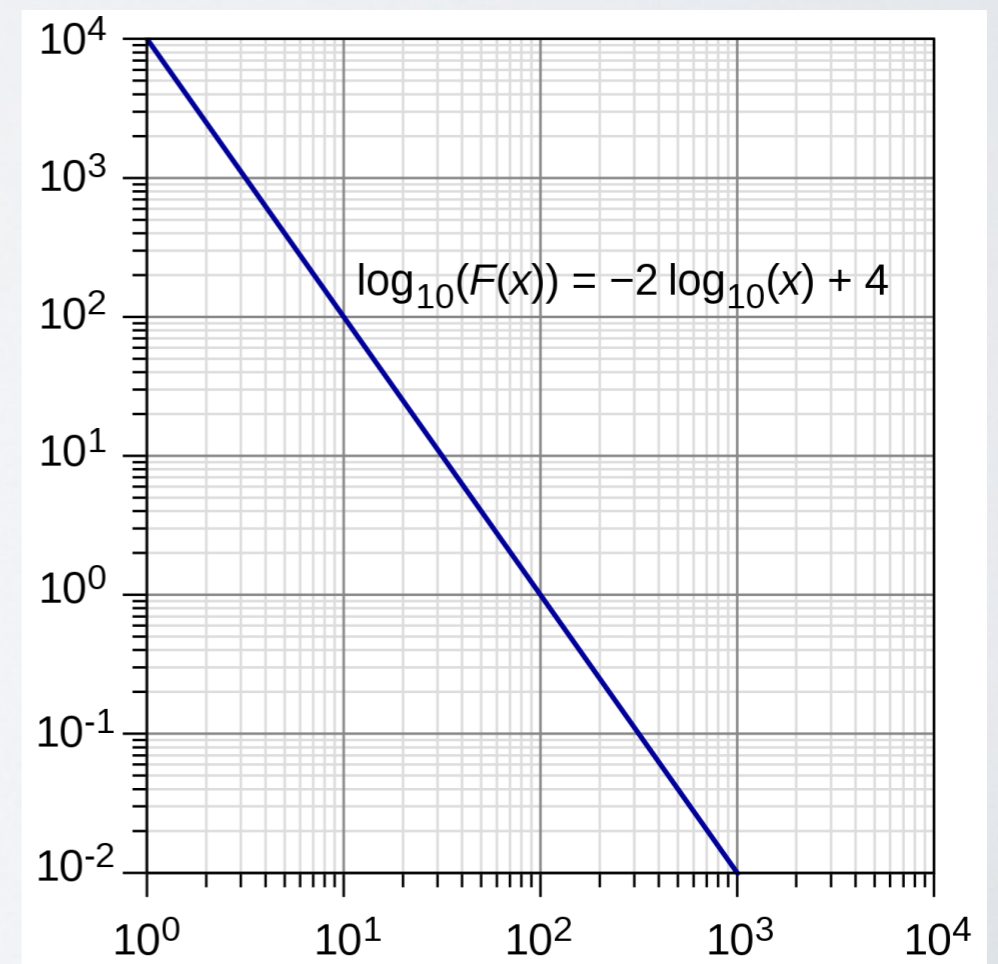
Number of successes in a sequence of  $n$  independent experiments, each asking a yes–no question, and each with its own Boolean-valued outcome: success (with probability  $p$ ) or failure (with probability  $q = 1 - p$ )



# THEORETICAL DISTRIBUTIONS

- Power Law distribution

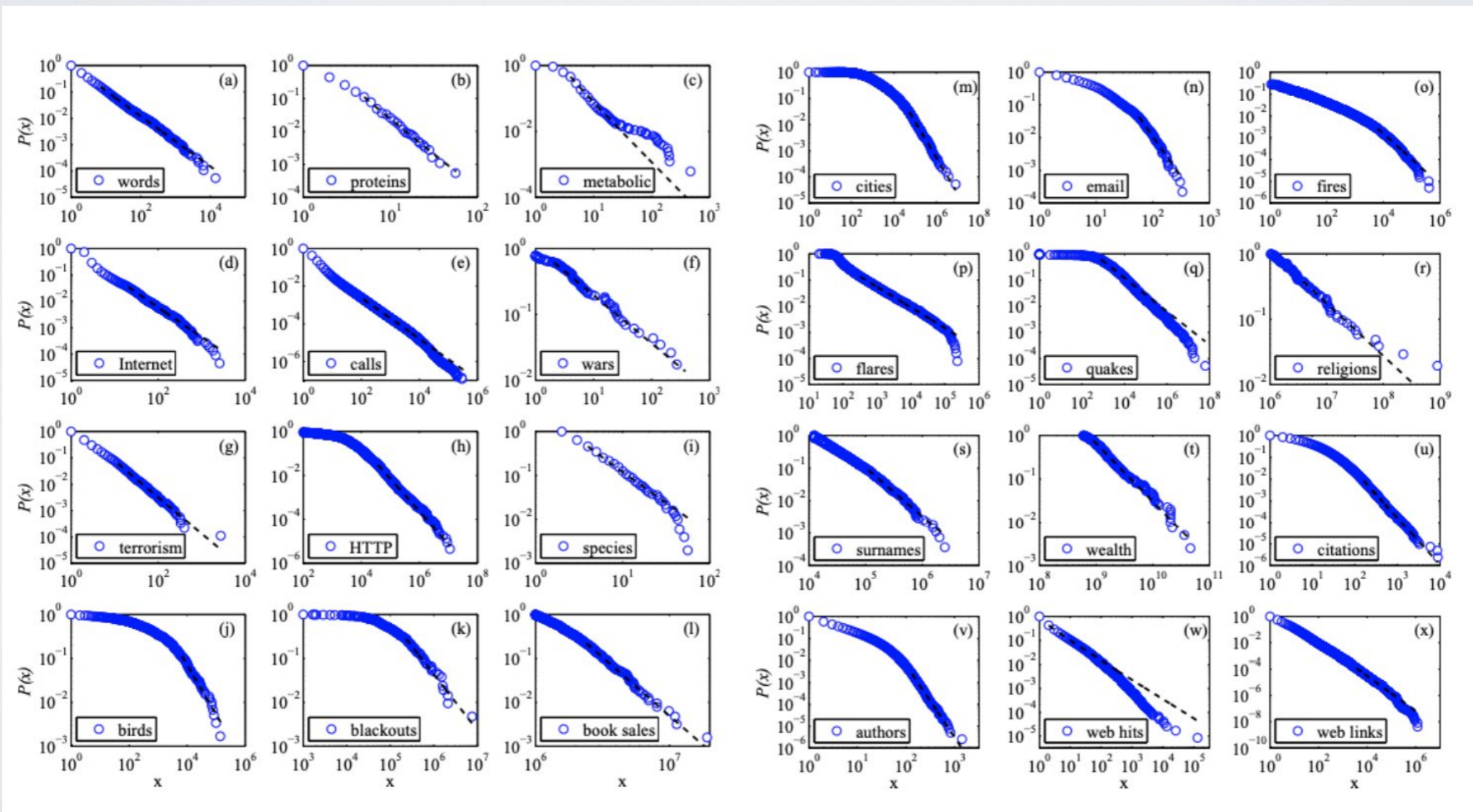
- ▶ A relative change in one quantity results in a proportional relative change in the other quantity, independent of the initial size of those quantities: one quantity varies as a power of another.



# THEORETICAL DISTRIBUTIONS

## DISTRIBUTIONS

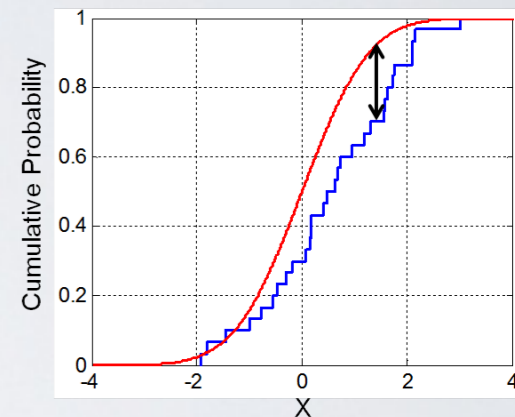
- Power Law distribution



# DISTRIBUTION COMPARISON

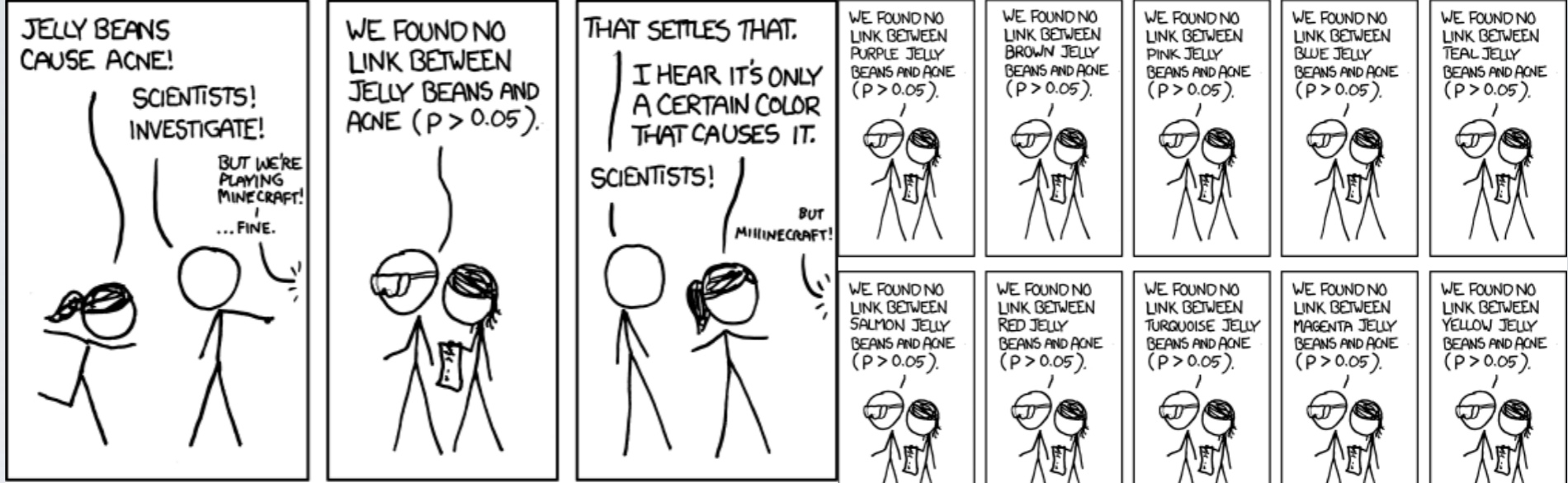
- Statistical test

- ▶ What is the probability that my observed data comes from the theoretical distribution XXX
  - Normality: Shapiro-Wilk, etc.
  - Categorical variables : Chi-squared  $\chi^2$
  - Etc.
- ▶ What is the probability that two distributions are identical ?
  - Kolmogorov-Smirnov test
  - Bootstrapping
- ▶ “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

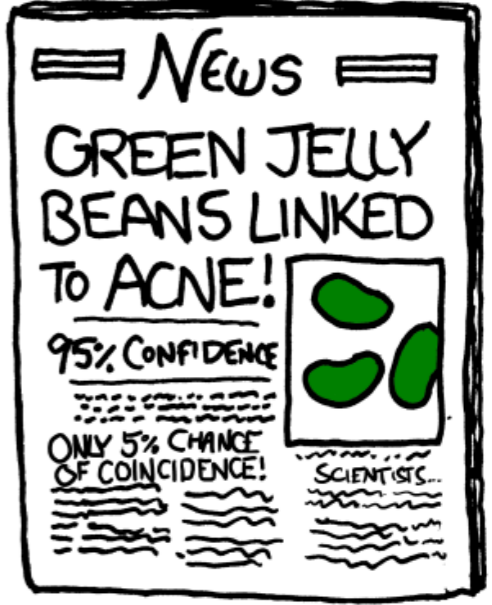
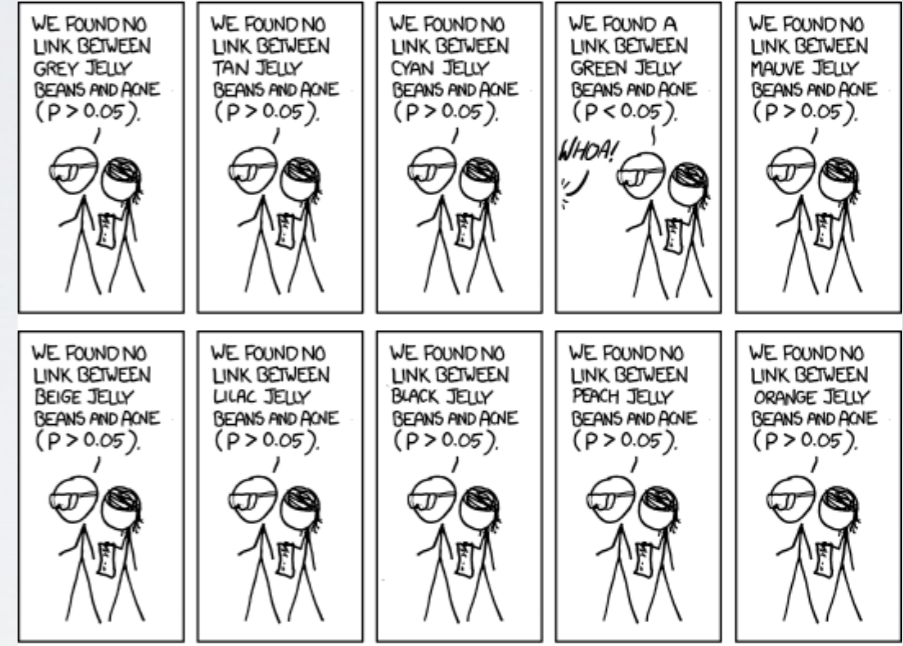




# P-VALUES



<u>P-VALUE</u>	<u>INTERPRETATION</u>
0.001	HIGHLY SIGNIFICANT
0.01	
0.02	
0.03	SIGNIFICANT
0.04	
0.049	OH CRAP. REDO CALCULATIONS.
0.050	
0.051	ON THE EDGE OF SIGNIFICANCE
0.06	
0.07	HIGHLY SUGGESTIVE, SIGNIFICANT AT THE $p < 0.10$ LEVEL
0.08	
0.09	
0.099	HEY, LOOK AT THIS INTERESTING SUBGROUP ANALYSIS
$\geq 0.1$	



# 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

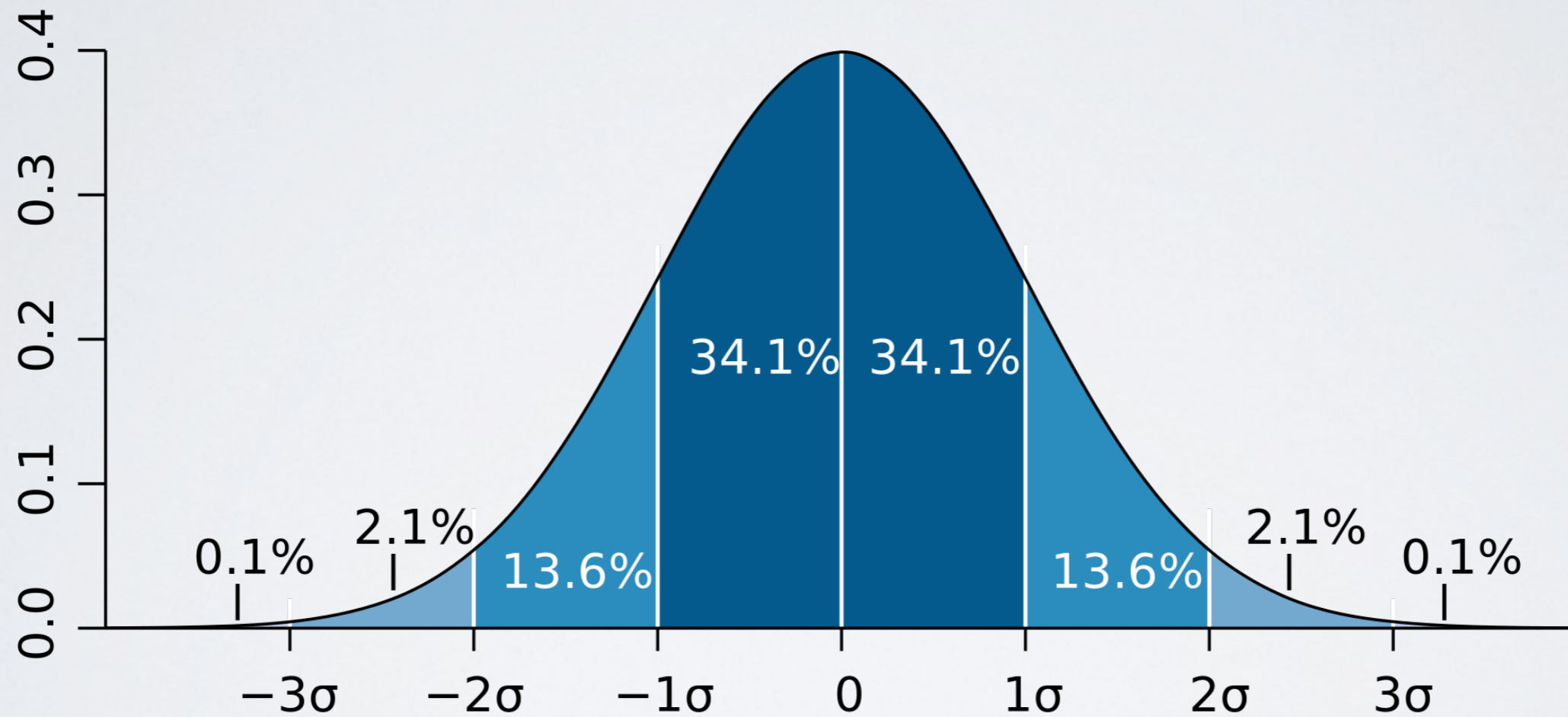
$$\sigma^2 = \frac{1}{N^2} \sum_{i < j} (x_i - x_j)^2$$

# STANDARD DEVIATION

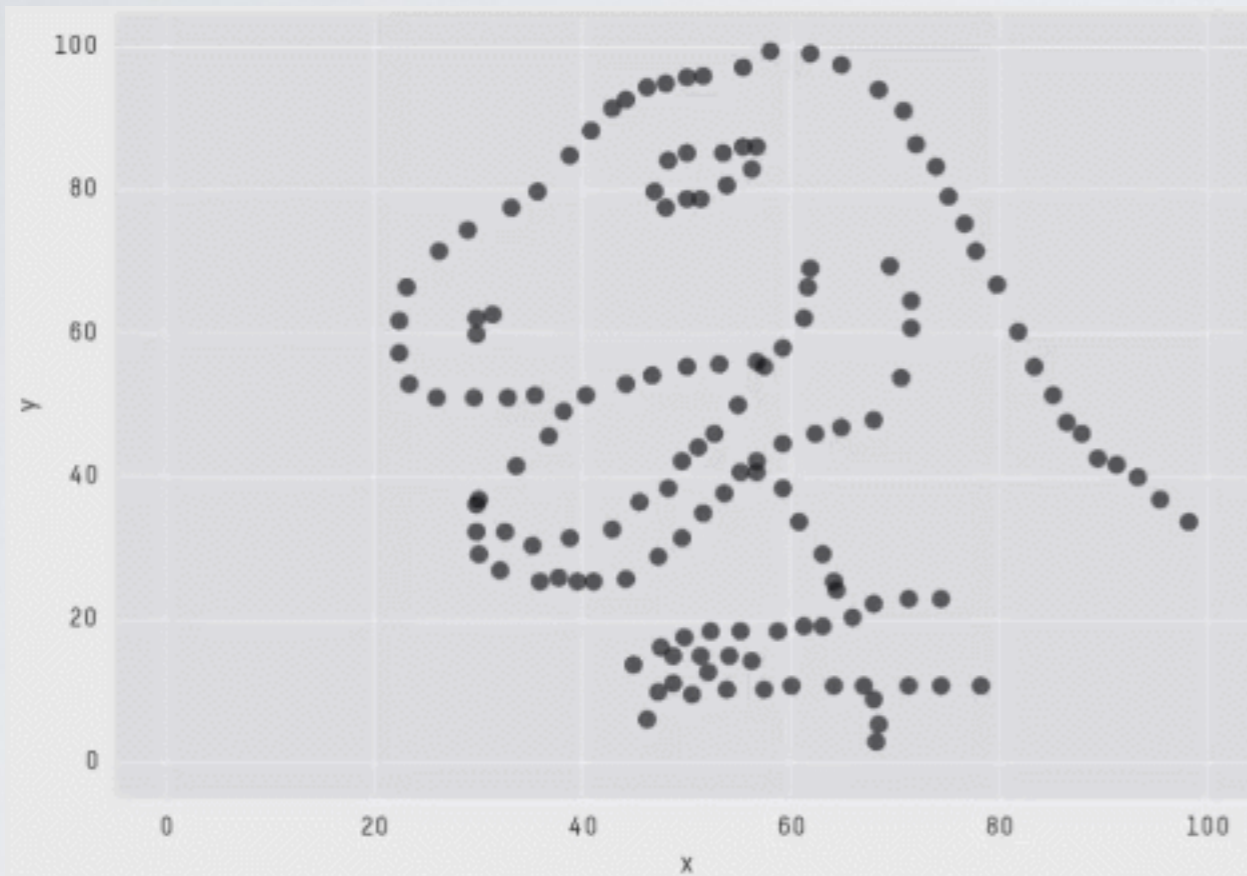
- Squared root of the Variance

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

# RELATION WITH NORMAL DISTRIBUTION



# 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>

# MEAN ABSOLUTE DEVIATION (MAD)

- MAD or AAD (Average Absolute Deviation)

- Deviation from mean or from median

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

- So why are we using the Standard Deviation again ?

- The mean minimizes the expected squared distance

- The median minimizes the MAD

- Nice relation with euclidean geometry (sum of variance is variance of the sum)

- Leads naturally to least square regression and PCA... see later.

- Nevertheless, not the unique true objective. Think of what you really want to measure... Sensibility to outliers... Undefined for power laws...

# VARIABLE INTERACTIONS

# COVARIANCE MATRIX

## Covariance Matrix Formula



- Covariance matrix  $\mathbf{K}$

- ▶ Extension of Variance to multivariate data

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

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

- How much variable  $X$  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

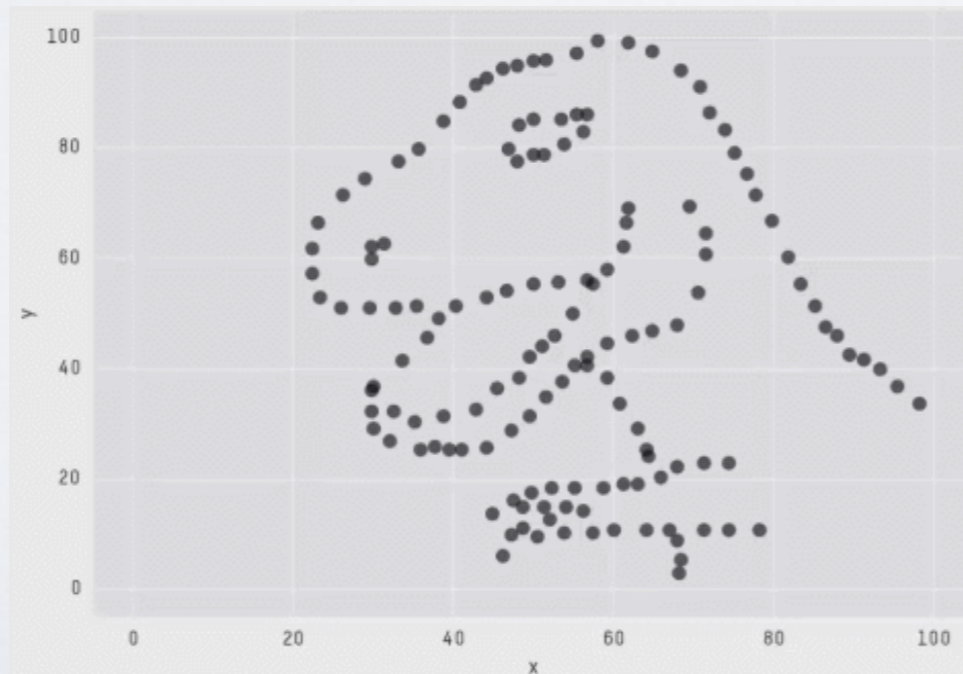
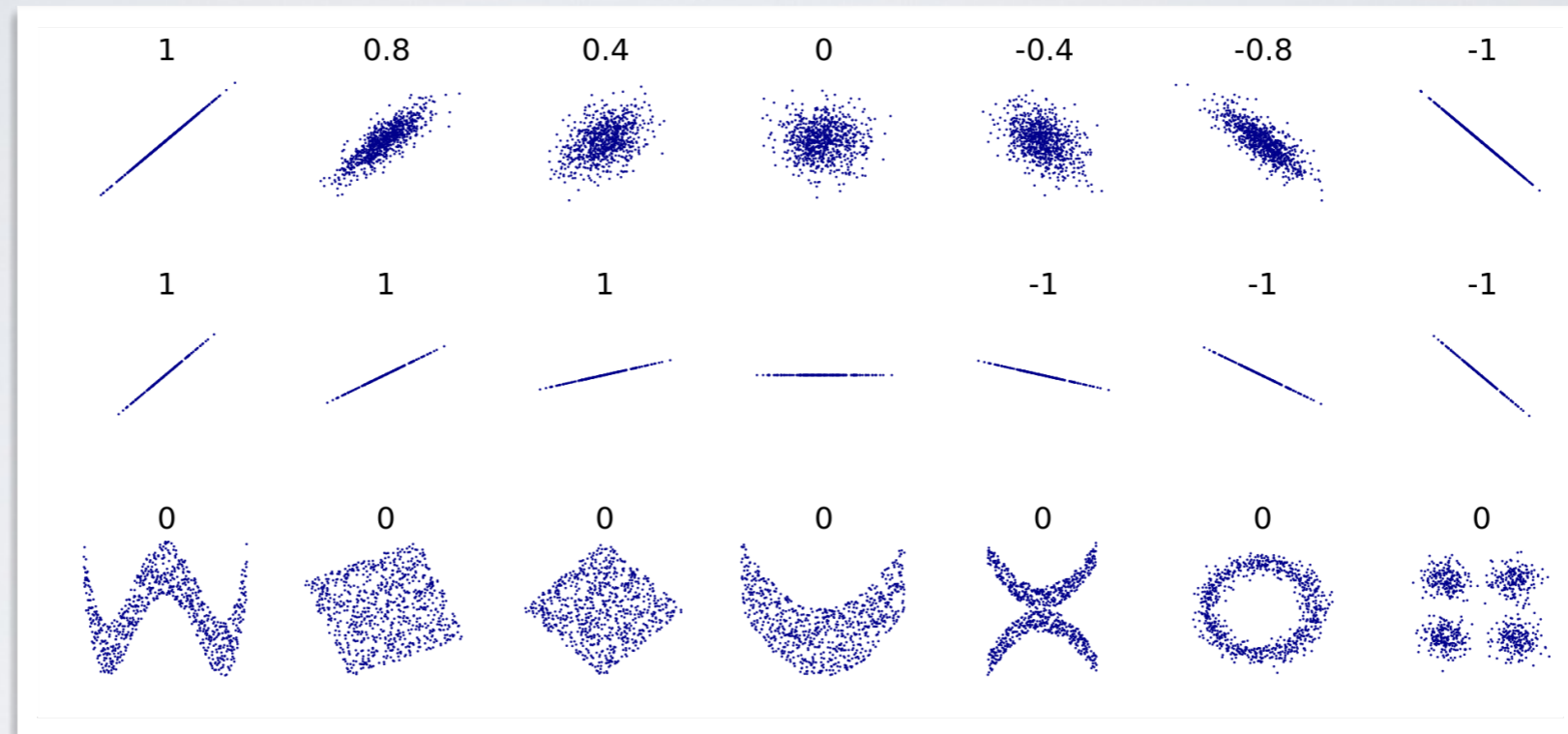
- ▶ Normalize it to obtain the “correlation coefficient”



# 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

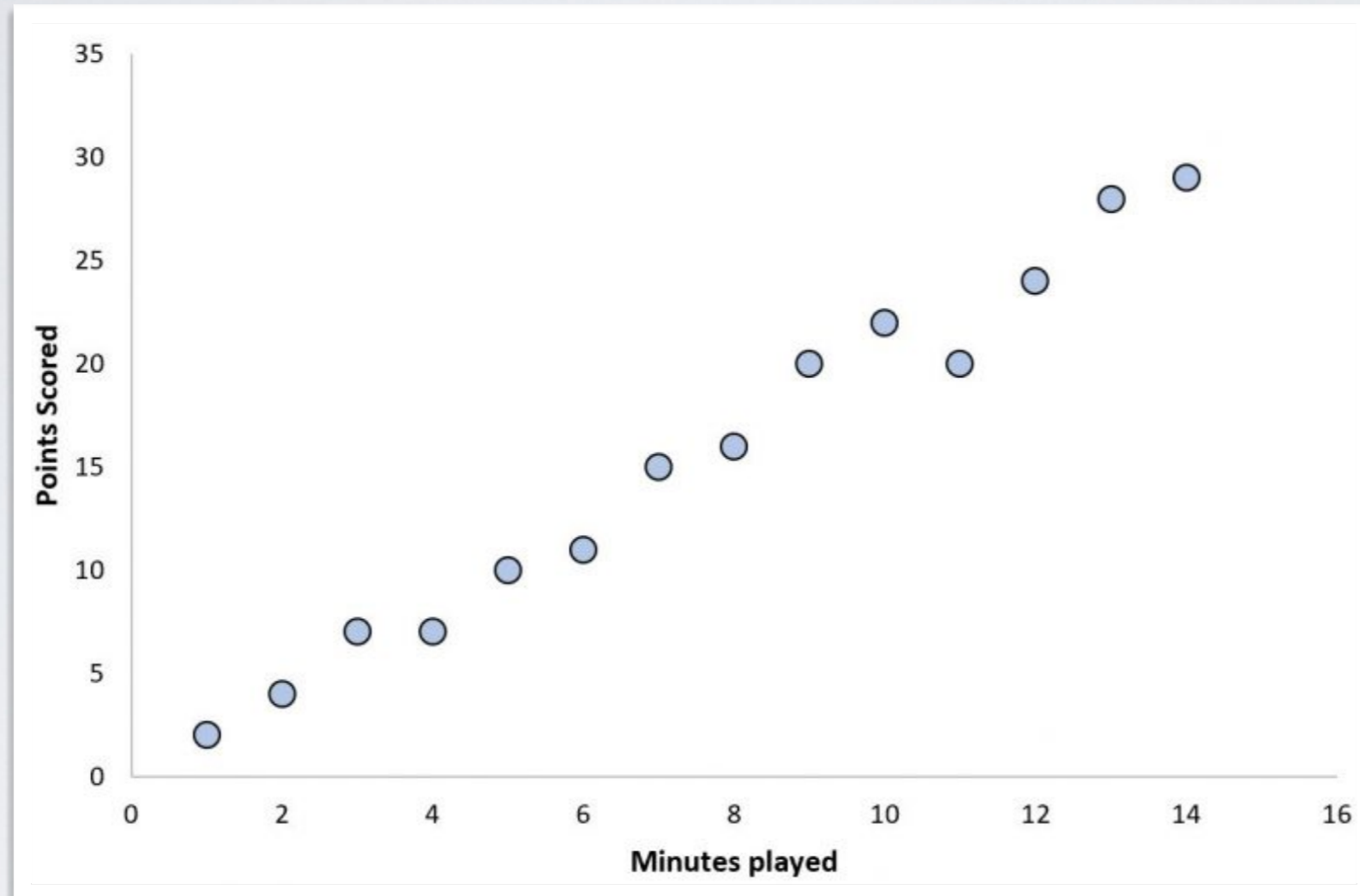


X Mean: 54.2659224  
Y Mean: 47.8313999  
X SD : 16.7649829  
Y SD : 26.9342120  
Corr. : -0.0642526

# 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.

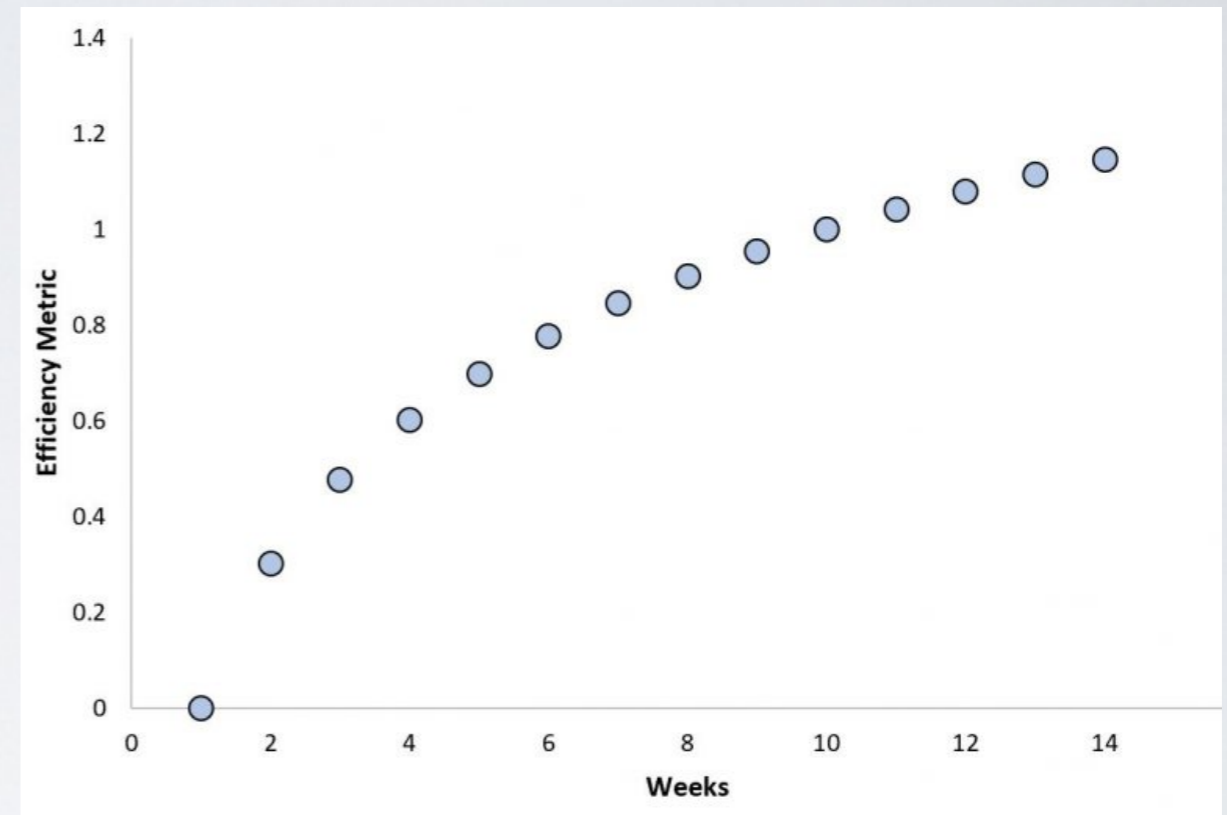
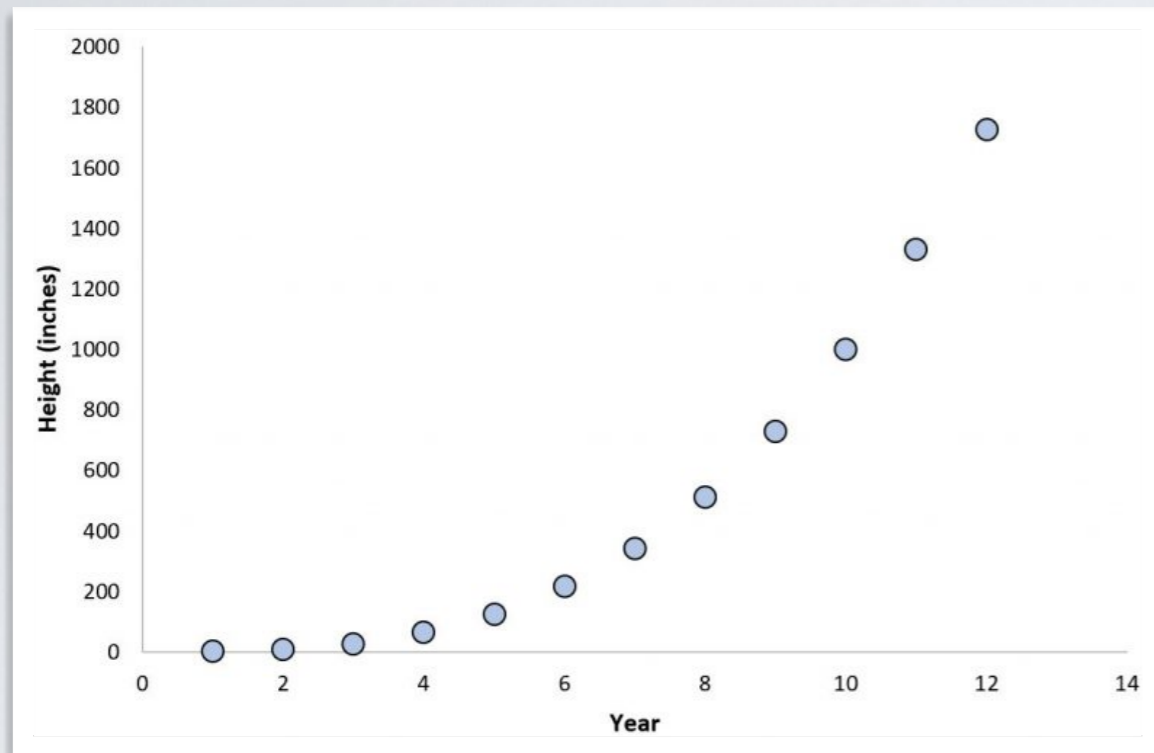
# NONLINEAR RELATIONSHIPS



Linear relationship

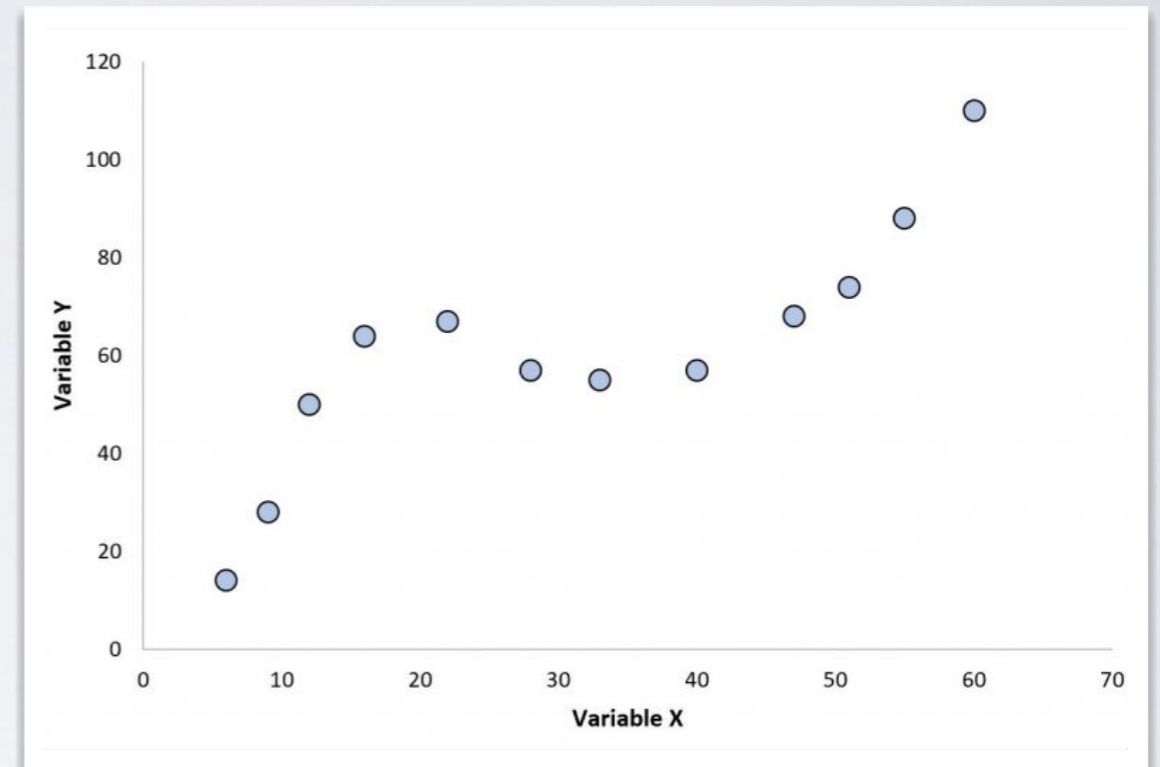
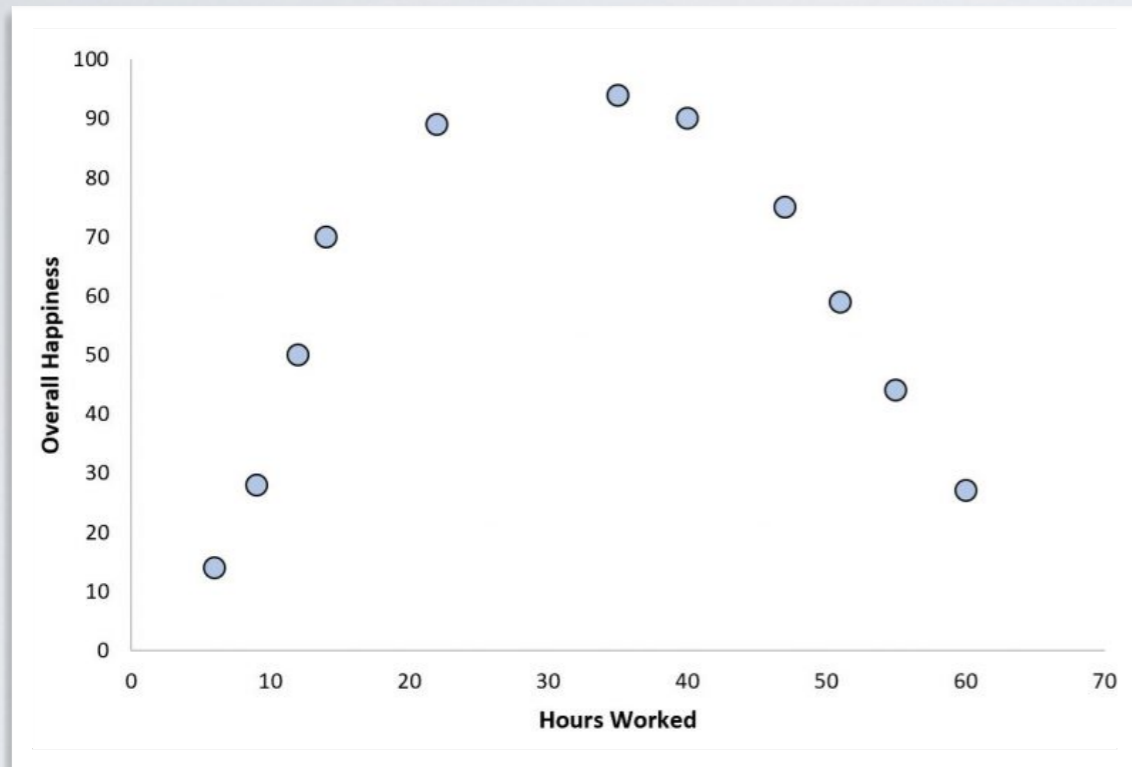
$$Y = a + bX + e$$

# NONLINEAR RELATIONSHIPS

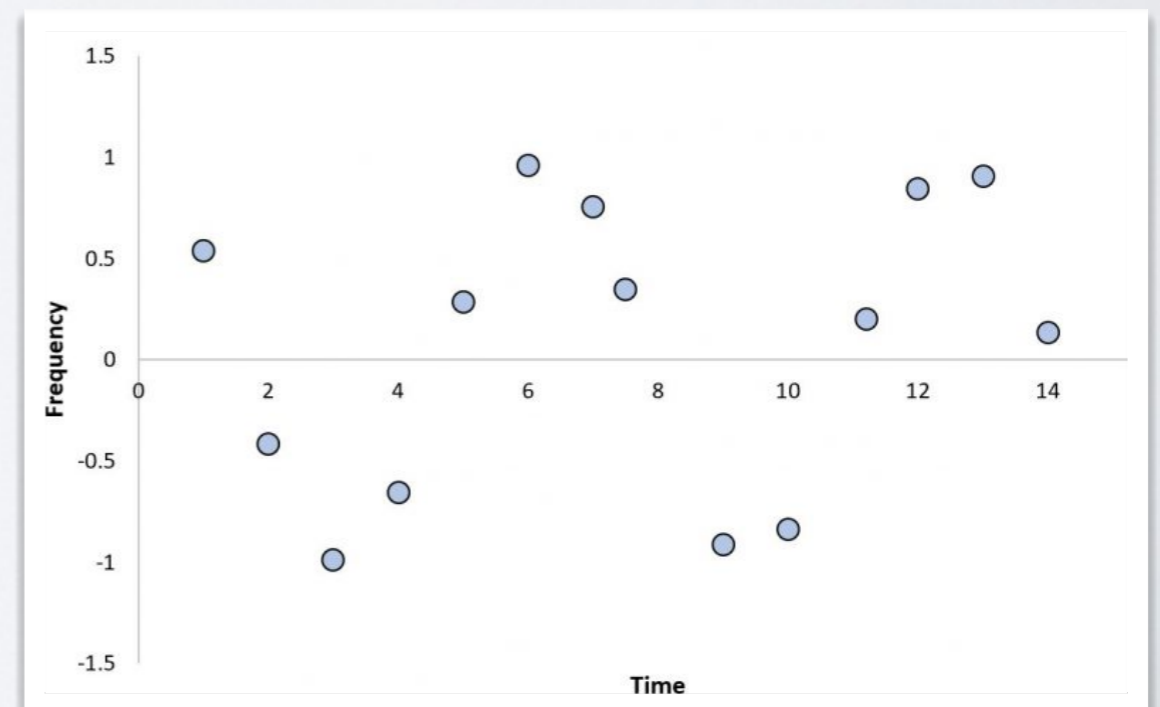


Monotonous, non-linear

# NONLINEAR RELATIONSHIPS



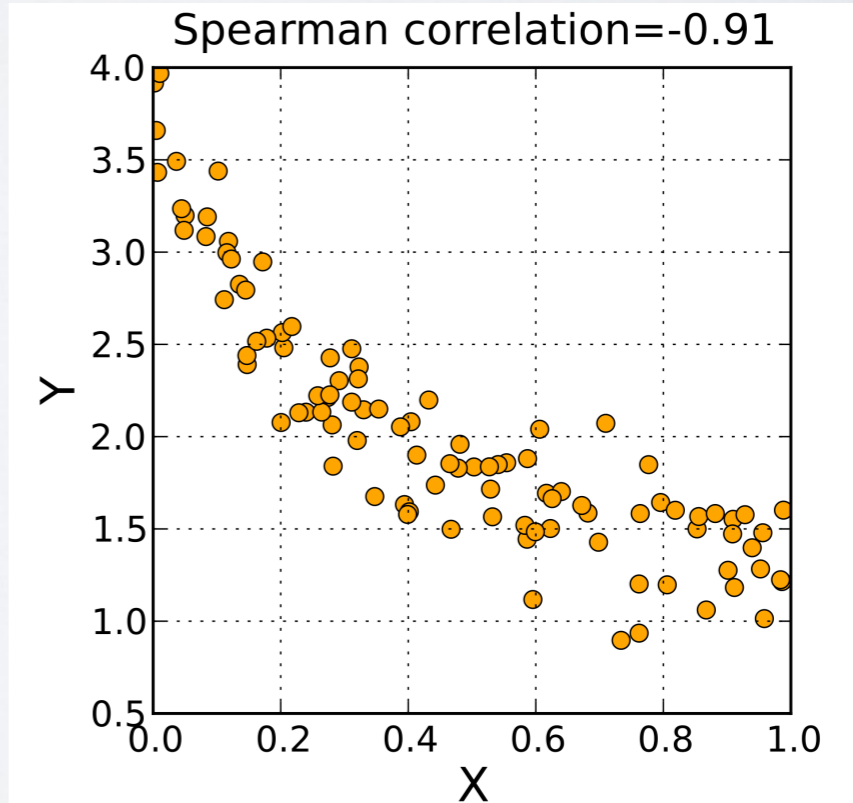
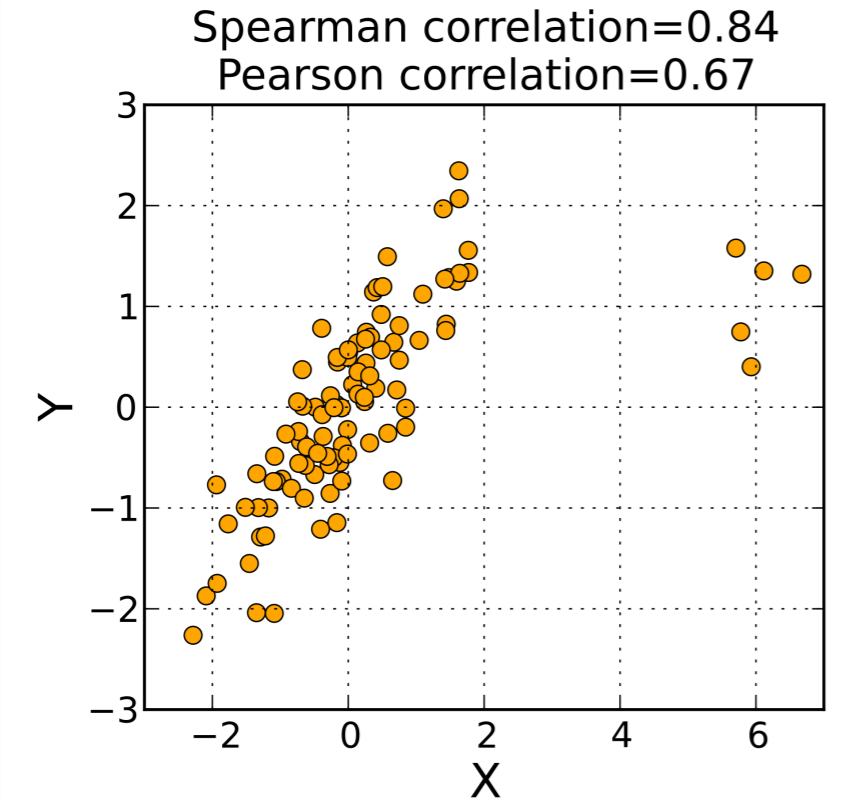
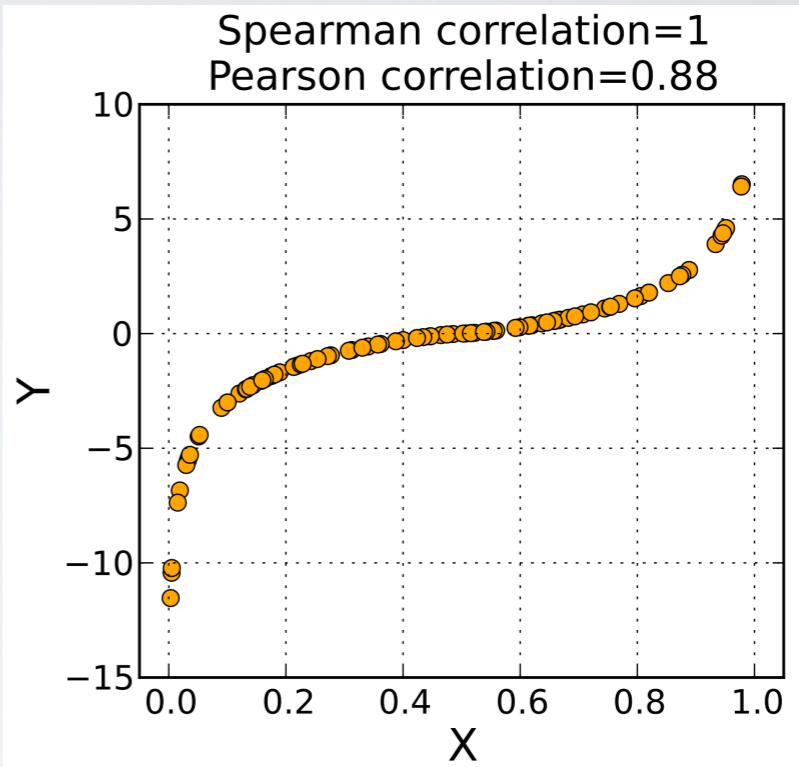
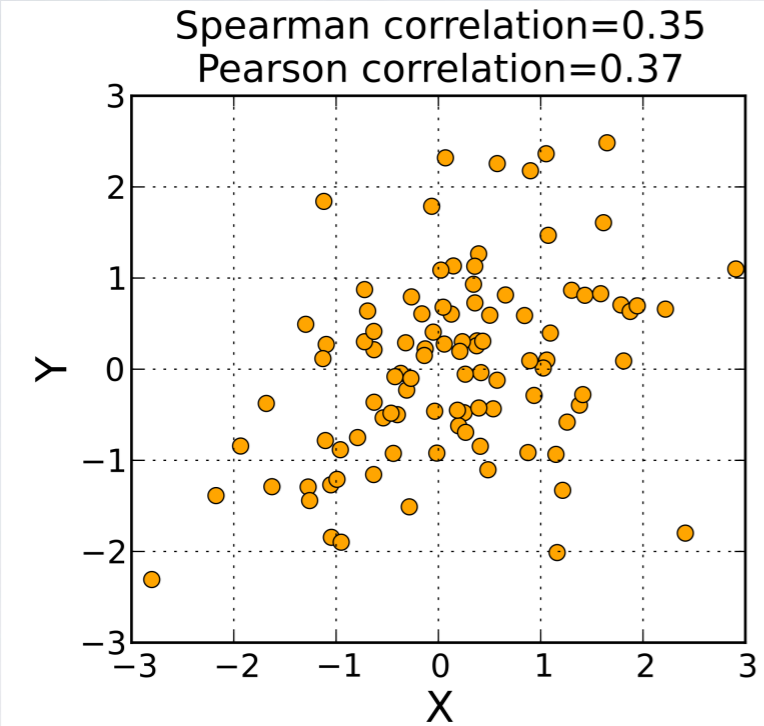
Non-monotonous,  
Non-linear



# 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
  - $r_s = \rho_{R(X), R(Y)} = \frac{\text{cov}(R(X), R(Y))}{\sigma_{R(X)}\sigma_{R(Y)}}$

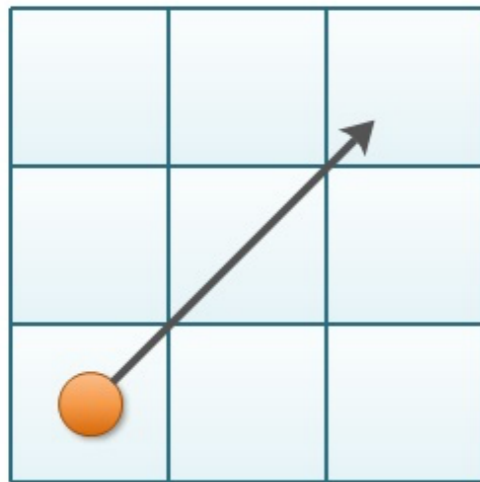
# SPEARMAN'S CORRELATION



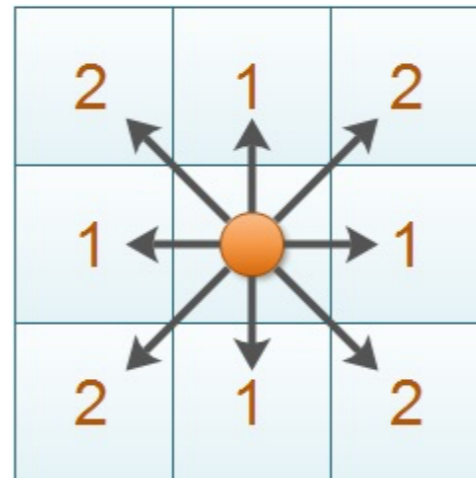


# NOTIONS OF DISTANCE

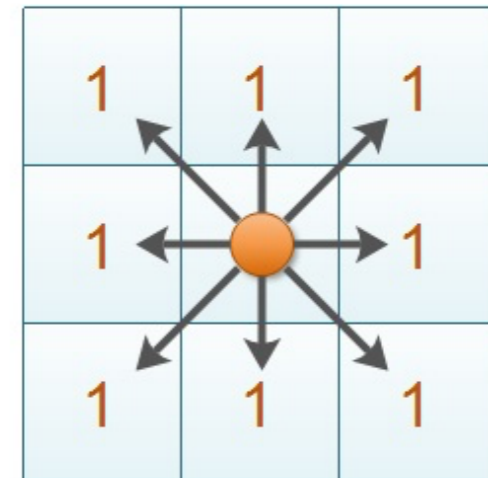
**Euclidean Distance**



**Manhattan Distance**

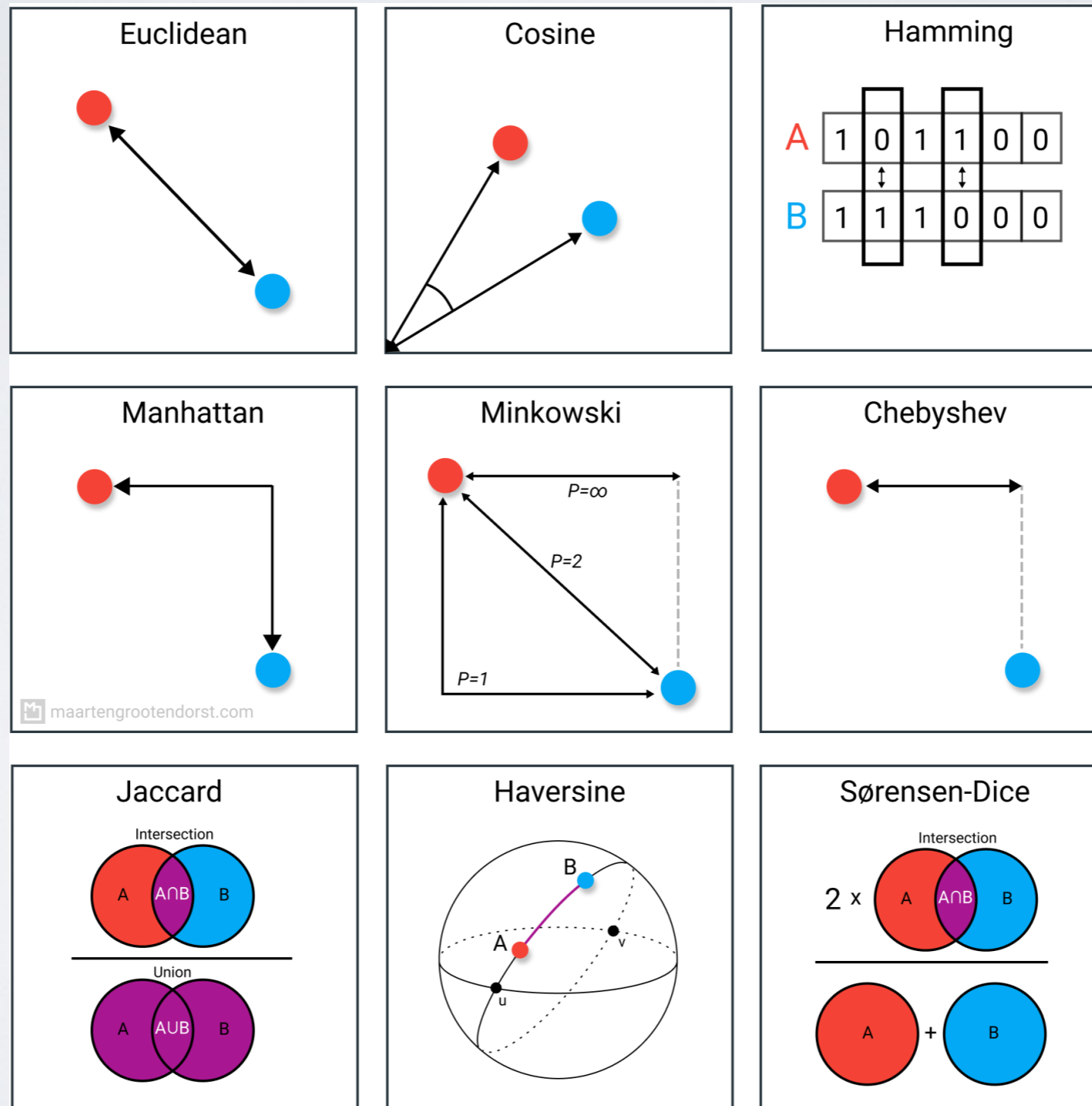


**Chebyshev Distance**



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

# NOTIONS OF DISTANCE



# FEATURE SCALING

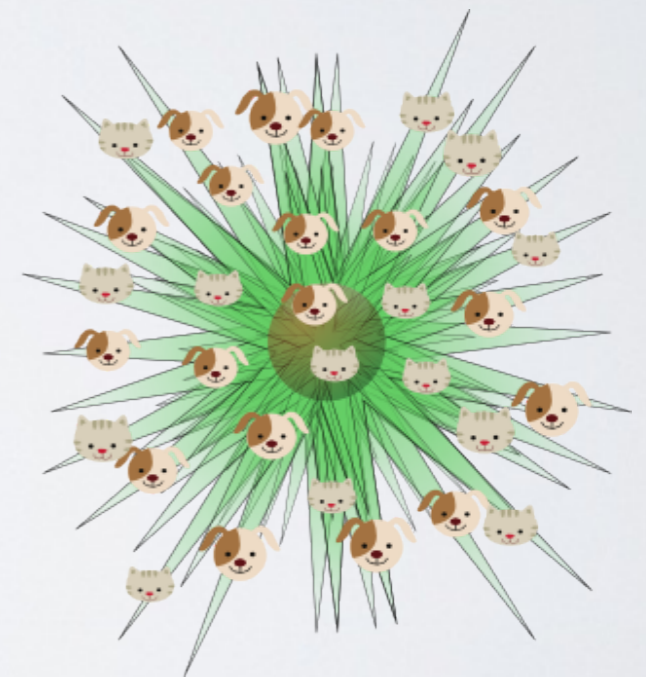
- We want to use euclidean distance to compute the “distance” between 2 people based on attributes age(y), height(m), weight(g).
  - ▶  $a = (y:20, m:1.82, g:80\ 000)$ ,  $b = (y:20, m:1.82, g:81\ 000)$ ,  $c = (y:90, m:1.50, g:80\ 020)$ 
    - $d(a,b) = 1000.0005$
    - $d(a,c) = 72.8$
  - ▶ That is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

# FEATURE SCALING

- Rescaling (Normalization):  $x' = \frac{x - \min(x)}{\max(x) - \min(x)}$
- Mean normalization:  $x' = \frac{x - \text{average}(x)}{\max(x) - \min(x)}$
- Standardization (z-score normalization):  $x' = \frac{x - \bar{x}}{\sigma}$

# CURSE OF DIMENSIONALITY

- Every observation is “far” from any other observation
- Imagine you focus on the 80% most frequent values for each variable:
  - ▶ 1 var: Covers 80% of the population
  - ▶ 2 var: 64% of the population
  - ▶ 3 var: 51%
  - ▶ 10 var: 10%
  - ▶ 100 var: 0.000000002%
- If you have many variables, you need huge datasets, else you cannot generalize if all observations are completely different



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 “Going further” section is not mandatory, you can do it if you have time and are interested