# Machine Learning and Data Mining WS21/22

# "4 Clustering I"

## Dr. Zeyd Boukhers

@ZBoukhers

## Institute for Web Science and Technologies

## University of Koblenz-Landau

## November 17, 2021

- Data dimension reduction
  - PCA
  - SVD

- What is clustering?
- What is unsupervised learning?
- How to evaluate clustering results?
- What are intrinsic and extrinsic evaluation measures?
- How does K-Means work?
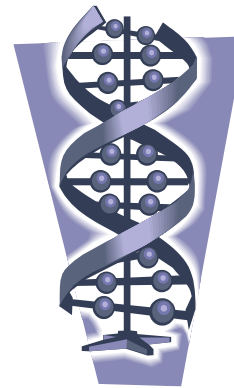- How to choose $k$ for K-Means?
- What is the EM algorithm?

User Profiles

Web document templates

Documents

Genetic sequences

Language dialects

- Identification of a finite set of *clusters* (= categories, "classes", groups) in the data.
- Objects in the same cluster should be as *similar* as possible.
    - High intra-similarity.
- Objects in different clusters should be as *dissimilar* as possible.
    - High inter-variance
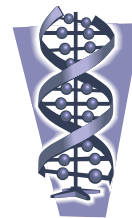- *"Unsupervised learning"* => no labels are given.

User profiles
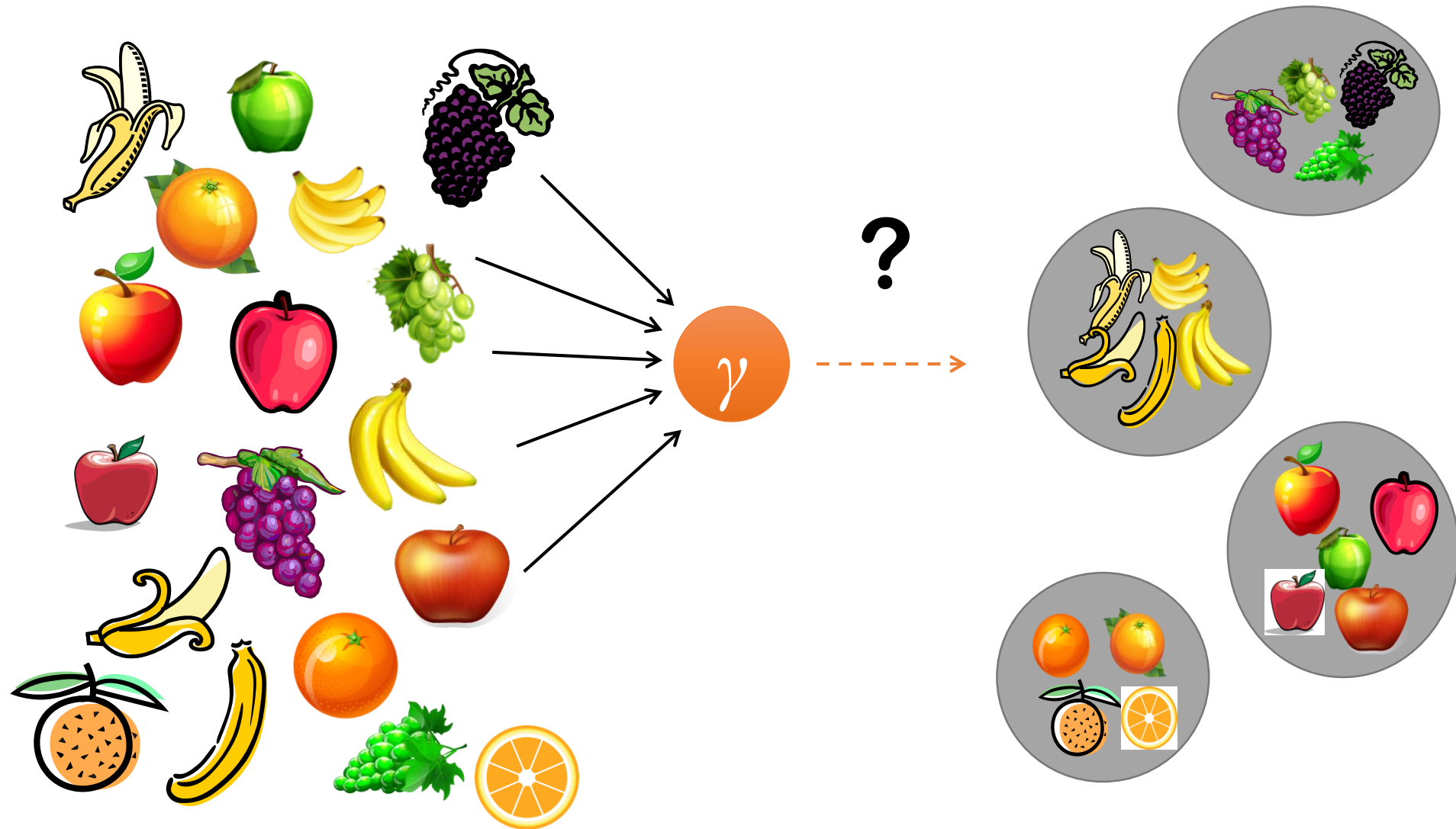
Web document templates

Documents

Genetic
sequences

Language dialects

- Objects (dataset):
  - $D = \{\boldsymbol{x}_1, \boldsymbol{x}_2, \dots, \boldsymbol{x}_N\}$

- An object is characterized by attributes
  - $\boldsymbol{x}_i = (x_{i,1}, x_{i,2}, \dots, x_{i,l})$

(green, round, even)

(orange, round, rough)

- Task:
  - Find groups
    - $\Psi = \{\psi_1, \psi_2, \dots, \psi_k\}$
  - Find function
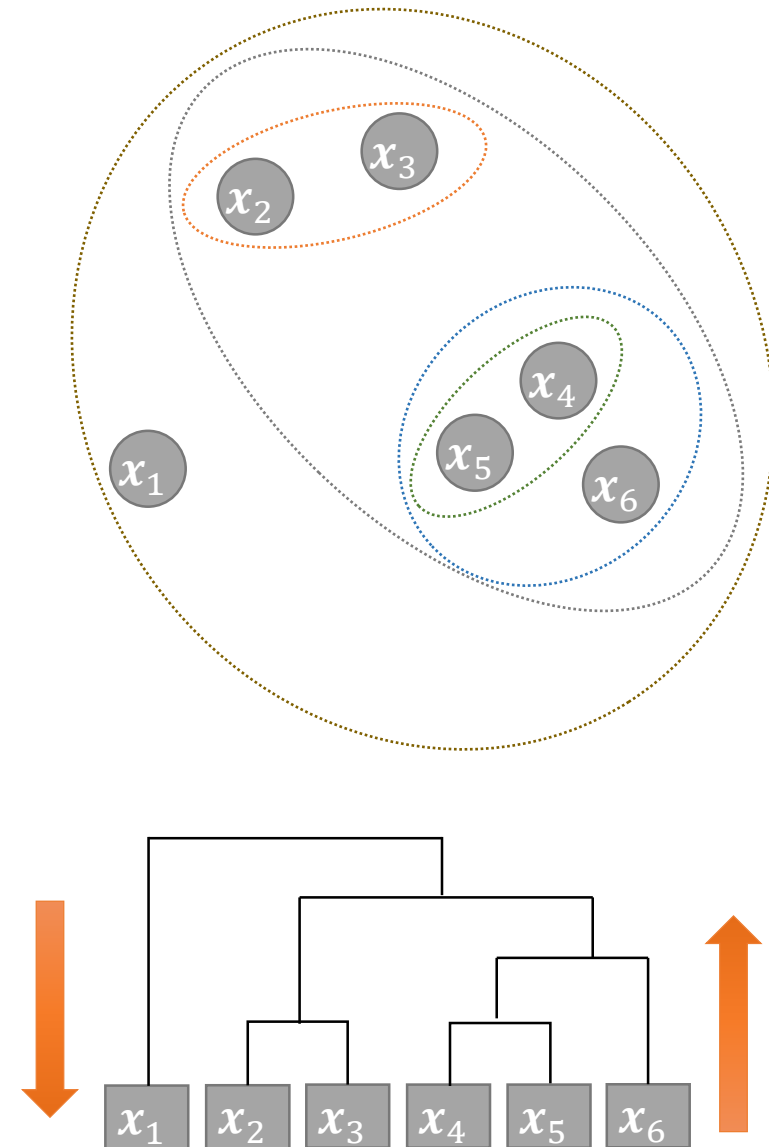    - $\gamma : D \to \Psi$

- Unsupervised learning

Difference to classification:
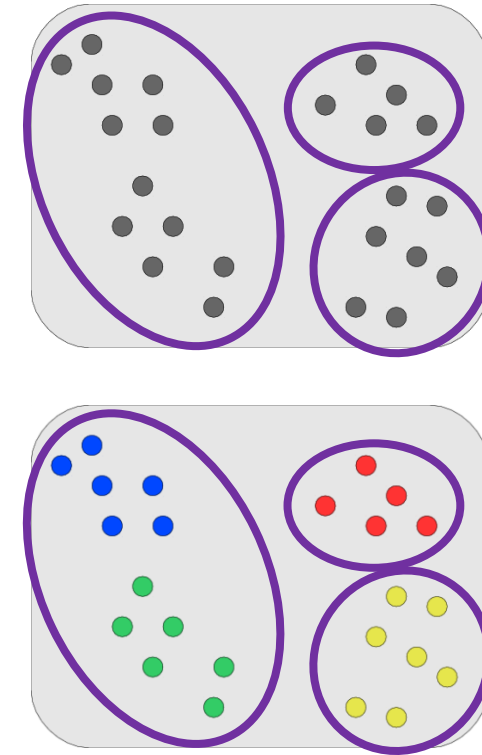Groups not given!

- ## Cluster types:
  - ### Flat vs. Hierarchical
  - ### Exclusive vs. Multiple clusters
    - $\gamma: D \to \wp(\Psi)$

- ## Function $\gamma$
  - ### Hard vs. Soft assignments
    - $\gamma: D \to \Psi \times \mathbb{R}$
  - ### Based on shape, density, estimates of distribution mixture

- ## Number of clusters
  - ### Provided (externally) or to be defined (given explicit hyperparameter) or found over the data (given other hyperparameters, e.g. density or density distribution).

- Intrinsic
  - Evaluate the quality of clusters directly.
    - E.g. compactness, separation of groups, etc.
- Extrinsic
  - Employ external knowledge
    - Ground truth from classification data
    - Assuming categories to be optimal clusters
  - Compare found with pre-defined clusters
    - Difficulty of finding a matching
- Indirect
  - User testing (satisfaction, task performance)
  - Application specific metrics

- Dunn Index
  - Notion of cluster separation

$$I_{\text{Dunn}}(\Psi) = \frac{\delta_{\min}}{\delta_{\max}}$$

- $\delta_{\min}$ smallest inter-cluster distance
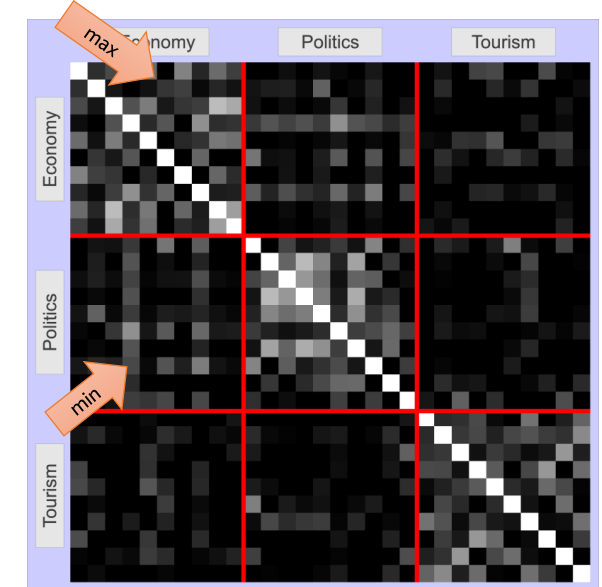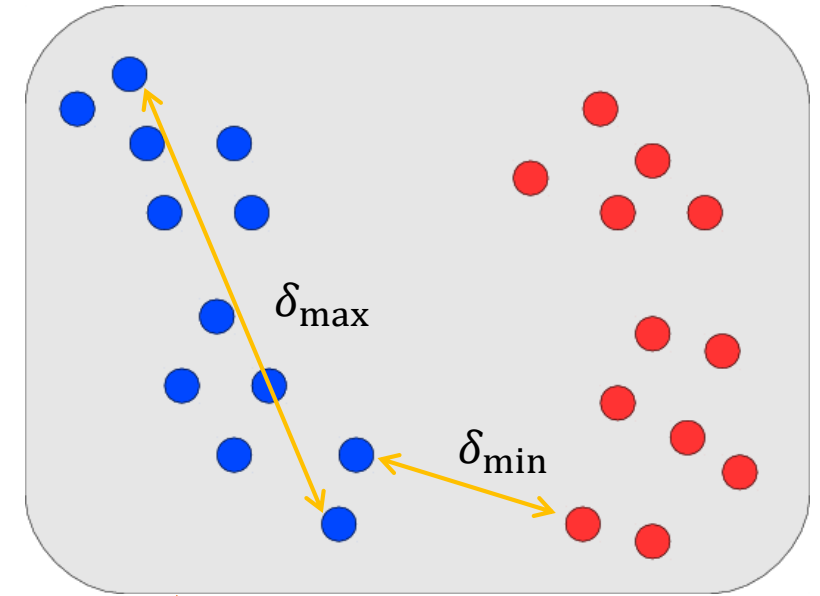- $\delta_{\max}$ largest intra-cluster distance

- Requires pair-wise distances
  - Distance matrix
  - Graphical representation
    - Minimal distance: white
    - Maximal distance: black

- Applicable also to ground truth
  - Notion of difficulty of cluster problem
  - Example

$$I_{\text{Dunn}}(\{c_E, c_P, c_T\}) = \frac{0.577}{1.414} = 0.435$$

14

- Silhouette coefficient $s(i)$
  for object $\boldsymbol{x}_i$
  - Average distance $a(i)$ to all other
    objects in the same cluster $\psi$

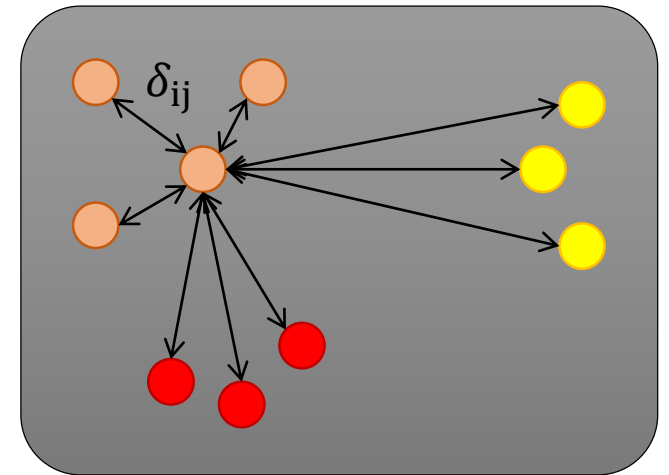$$a(i) = \sum_{\boldsymbol{x} \in \psi} \frac{1}{|\psi|} \delta(\boldsymbol{x}_i, \boldsymbol{x})$$

  - Average distance to any other cluster $\psi', \psi' \neq \psi$:

$$d(i, \psi') = \sum_{\boldsymbol{x} \in \psi'} \frac{1}{|\psi'|} \delta(\boldsymbol{x}_i, \boldsymbol{x})$$

  - Average distance $b(i)$ to the closest cluster $\psi', \psi' \neq \psi$:
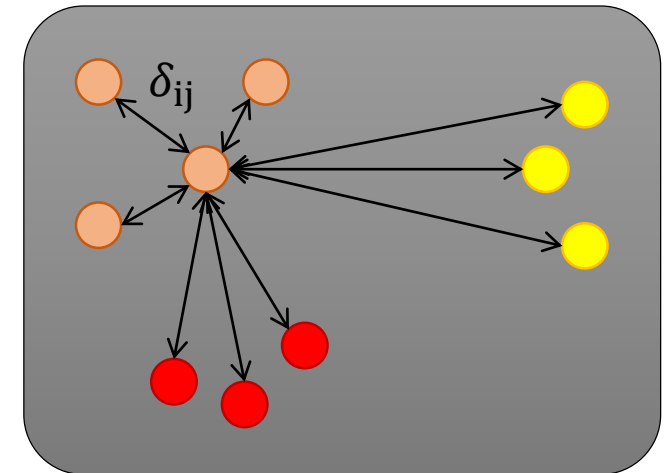
$$b(i) = \min_{\psi' \epsilon \Psi} d(i, \psi')$$

  - Silhouette coefficient: $s(i) = \frac{b(i) - a(i)}{max(a(i), b(i))}$

- Silhouette coefficient
  - Values: $-1 \leq s(i) \leq 1$
  - Value close to 1:
    - $a(i)$ much smaller than $b(i)$
    - Distances within cluster very small in comparison to distances with other clusters
  - Value close to 0:
    - $a(i) \approx b(i)$
    - Same internal as external distance
  - Value close to -1:
    - $b(i)$ much smaller than $a(i)$
    - Other instances are (on average) closer than the instances in same cluster

- Aggregation: Average silhouette coefficient

$$s(i) = \frac{b(i) - a(i)}{max(a(i), b(i))}$$

$\delta_{ij}$

- External knowledge (ground truth of categories)

$$\Omega = \{\omega_1, \cdots, \omega_{k\prime}\}$$

- Compare the clusters $\Psi$ and categories $\Omega$

- Determine:

  - $n_j^{(u)}$ : number of objects from $\omega_u$ being clustered into $\psi_j$

- Purity

  - Ratio of strongest represented category

$$\text{Purity}(\psi_j) = \frac{1}{|\psi_j|} \cdot \max_{u=1,\cdots,k\prime} n_j^{(u)}$$

  - Aggregate over all clusters

$$\text{Purity}(\Psi) = \sum_{j=1}^{k} \frac{|\psi_j|}{N} \cdot \text{Purity}(\psi_j)$$

Purity of 1 can always be achieved!

- Mutual Information:
  - Mutual agreement between clusters and categories

$$\mathrm{MI}(\Psi) = \frac{1}{N} \sum_{j=1}^{k} \sum_{u=1}^{k'} n_j^{(u)} \cdot \log \frac{n_j^{(u)} \cdot N}{\sum_{m=1}^{k'} n_j^{(m)} \cdot \sum_{t=1}^{k} n_t^{(u)}}$$

  - Log base: 2 or $k \cdot k'$

- Rand Index:
  - Consider document pairs on categories and clusters
    - Agreements: same-same (ss), different-different (dd)
    - Disagreements: same-different (sd), different-same (ds)
  - Agreement-ratio:

$$I_{\mathrm{Rand}}(\Psi) = \frac{\mathrm{ss} + \mathrm{dd}}{\mathrm{ss} + \mathrm{dd} + \mathrm{sd} + \mathrm{ds}}$$

1. $\psi_1 = \{x_1, x_2, x_3, x_4, x, x_6, x_7, x_8, x_9, x_{10}, x_{11}, x_{20}\}$
2. $\psi_2 = \{o_{12}, o_{13}, o_{14}, o_{15}, o_{16}, o_{17}, o_{18}, o_{19}\}$
3. $\psi_3 = \{x_{21}, x_{22}, x_{23}, x_{24}, x_{25}, x_{27}, x_{28}, x_{30}\}$
4. $\psi_4 = \{x_{26}, x_{29}\}$

- Purity:

$$\text{Purity}(\psi_1) = \frac{10}{12} = 0.83$$

$$\text{Purity}(\Psi) = \frac{12}{30} \cdot 0.83 + \frac{8}{30} \cdot 1.0 + \frac{8}{30} \cdot 1.0 + \frac{2}{30} \cdot 1.0 = 0.93$$

- Mutual Information
  - Log base $k \cdot k'$
  - Several values are 0

$$\text{MI}(\Psi) = \frac{1}{30} \cdot \left( 10 \cdot \log\frac{10 \cdot 30}{12 \cdot 10} + 2 \cdot \log\frac{2 \cdot 30}{12 \cdot 10} + 8 \cdot \log\frac{8 \cdot 30}{8 \cdot 10} + 8 \cdot \log\frac{8 \cdot 30}{8 \cdot 10} + 2 \cdot \log\frac{2 \cdot 30}{2 \cdot 10} \right) = 0.370$$

# Example

1. $\psi_1 = \{x_1, x_2, x_3, x_4, x, x_6, x_7, x_8, x_9, x_{10}, x_{11}, x_{20}\}$
2. $\psi_2 = \{o_{12}, o_{13}, o_{14}, o_{15}, o_{16}, o_{17}, o_{18}, o_{19}\}$
3. $\psi_3 = \{x_{21}, x_{22}, x_{23}, x_{24}, x_{25}, x_{27}, x_{28}, x_{30}\}$
4. $\psi_4 = \{x_{26}, x_{29}\}$

- Agreements:
  - $ss = 103 = \binom{10}{2} + 1 + \binom{8}{2} + \binom{8}{2} + 1$
  - $dd = 280 = 10 \cdot 18 + 2 \cdot 10 + 8 \cdot 10$

- Disagreements
  - $sd = 32 = 2 \cdot 8 + 2 \cdot 8$
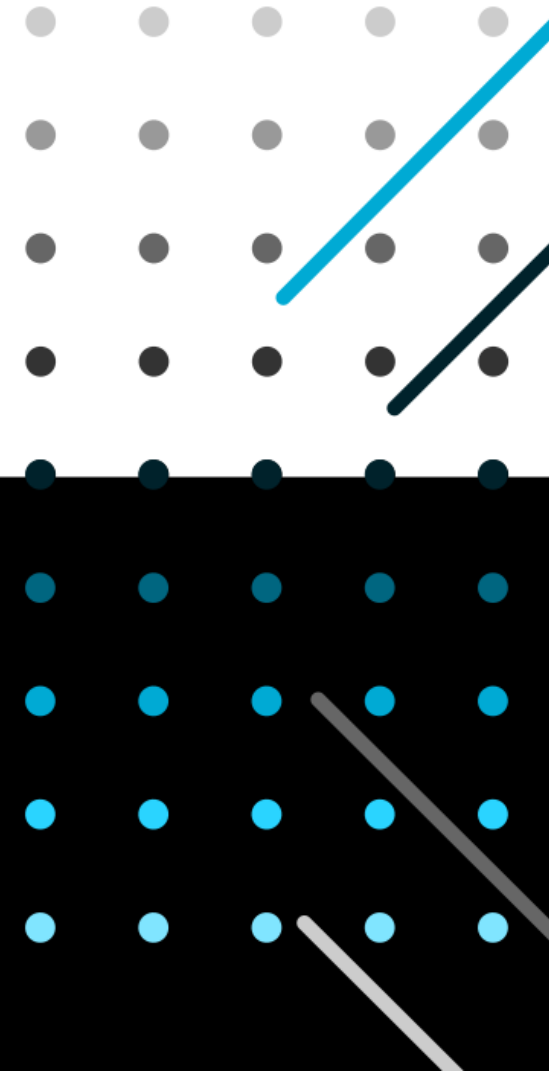  - $ds = 20 = 2 \cdot 10$

- Rand Index:

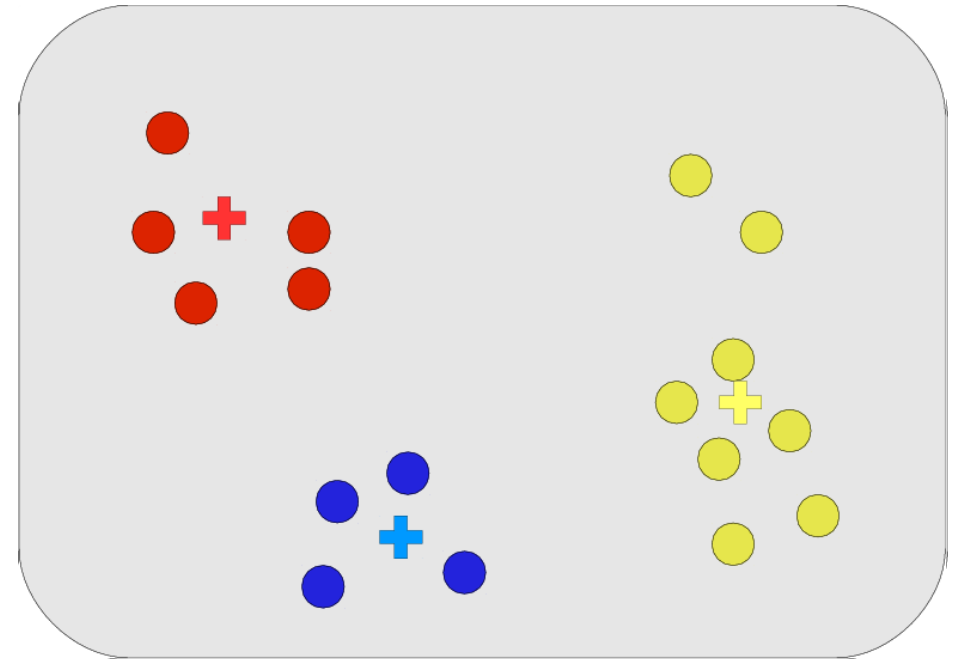$$I_{\text{Rand}}(\Psi) = \frac{383}{435} = 0.88$$

# K-Means

- General clustering algorithm

- Characteristics:
  - Flat clusters
  - No overlaps
  - Good runtime
  - Simple to implement

- Parameters
  - $k$ : number of clusters
  - Initial random seed

- Random seed $\Psi$ for $k$ clusters (e.g. single objects)

- Determine centroids $Z$ for clusters

- While centroids not stable
  - For all objects $\boldsymbol{x}_i$
    - Reassign $\boldsymbol{x}_i$ to cluster $\psi_j$ with minimal $\delta(\overrightarrow{\boldsymbol{x}_i}, \vec{z}_j)$
  - For all clusters $\psi_j$
    - Re-compute centroid $\vec{z}_j$

```
K-means(Set_of_points D, Integer k)
    Create an "initial" partitioning of D in k cluster;
    Compute a set Z'={Z'₁, ..., Z'ₖ} of centroids for the k Cluster;
    Z = {};
    repeat until Z = Z'
        Z = Z';
        Generate k clusters by assigning each data point to the nearest centroid
          Zⱼ;
        Compute the set Z'={Z'₁, ..., Z'ₖ} of centroids for the new clusters;
    return Z';
```
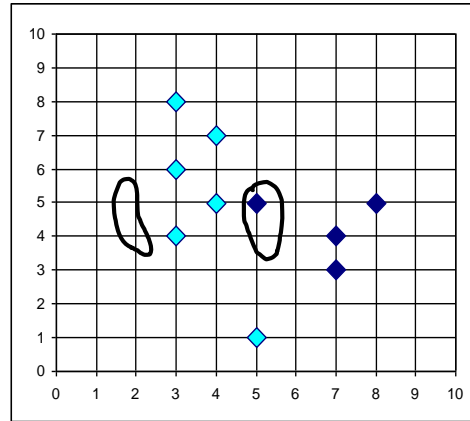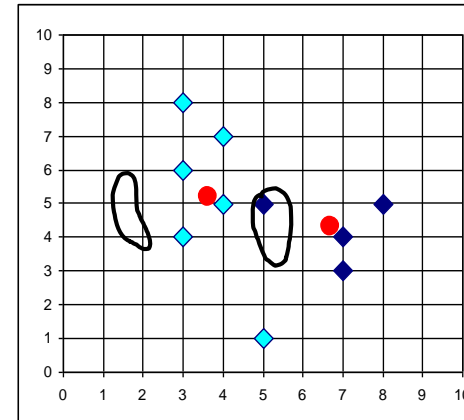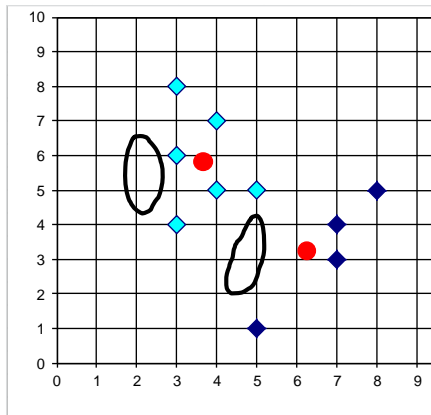
Compute
new centroids

Assign each point
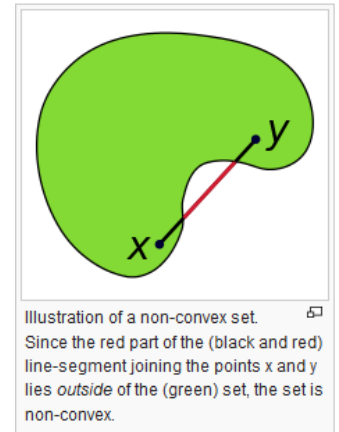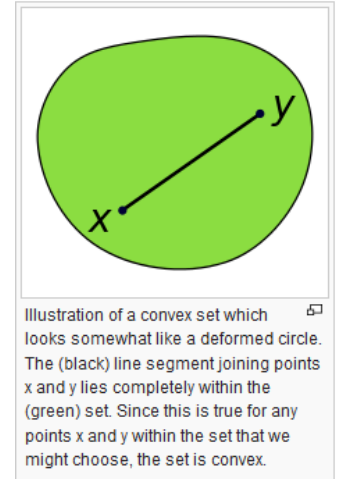to the nearest centroid

Compute
new centroids

- Advantages:
  - Efficiency:
    time complexity: $O(n)$ for each iteration,
    Number of iterations is usually very small (~ 5 - 10).
  - Simple implementation
  - Easy, good interpretability
  - $K$-means is the most popular clustering algorithm!

    $\Longrightarrow$

- Disadvantages:
  - Susceptible to noise and outliers since all objects influence the computation of centroids
  - Cluster have always convex form
  - The number of clusters $k$ is often difficult to determine
  - Strong dependency on initial partition (runtime + result!)

Illustration of a convex set which looks somewhat like a deformed circle. The (black) line segment joining points x and y lies completely within the (green) set. Since this is true for any points x and y within the set that we might choose, the set is convex.

Illustration of a non-convex set. Since the red part of the (black and red) line-segment joining the points x and y lies *outside* of the (green) set, the set is non-convex.
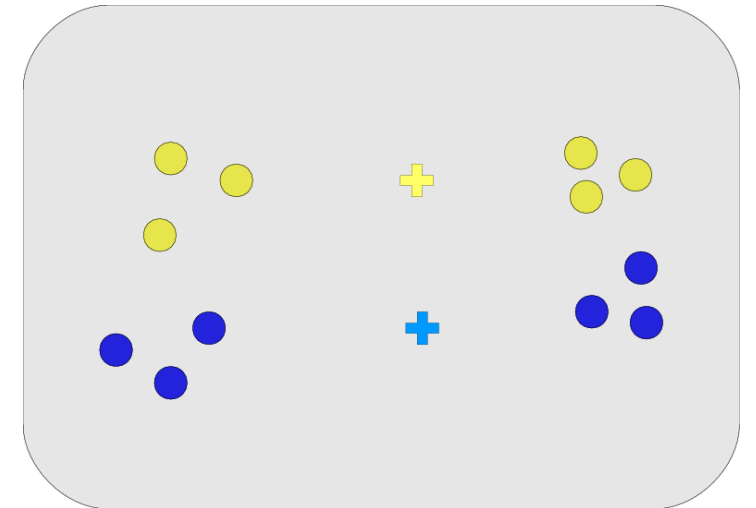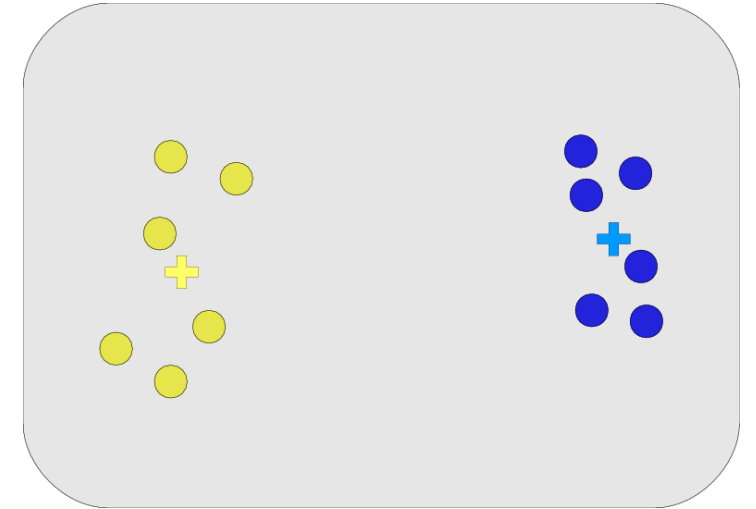
- Random seed
  - Assign all objects
  - Compute actual centroids

- Stop criterion
  - No change of clusters (equivalent)
  - Small changes of the centroids
  - Fixed number of iterations

- Centroids
  - Non-standard metrics (e.g. string similarity)
    - Mean centroid cannot be computed
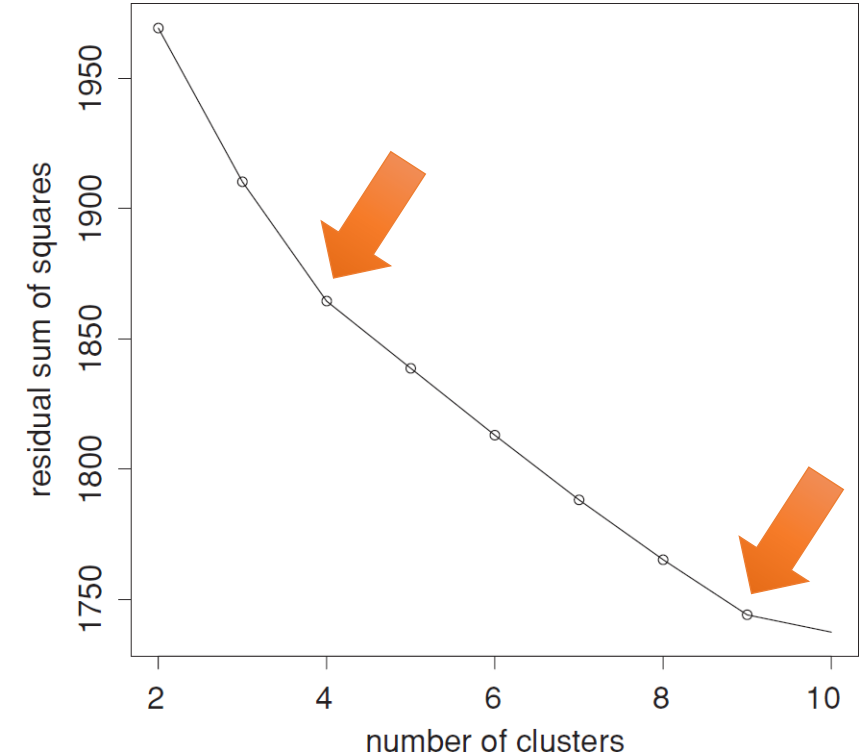  - K-Median
    - Use most central objects

- Choice of initial seed can cause different outcomes!
- Solution
  - Repeat with different seeds
  - Evaluate quality
    - Dunn index
    - Residual Sum of Squares

$$RSS(\Psi) = \sum_{j=1}^{k} \sum_{d_i \in \psi_j} \delta\left(\vec{d}_i, \vec{z}_j\right)^2$$
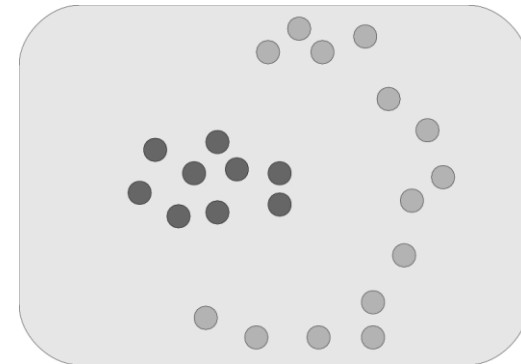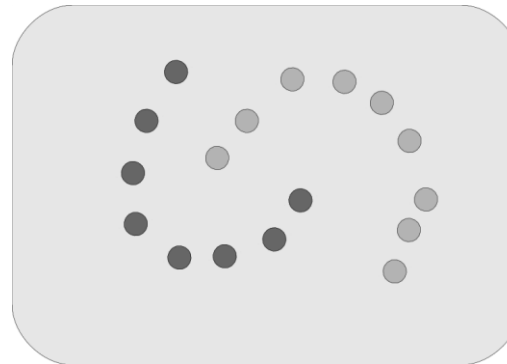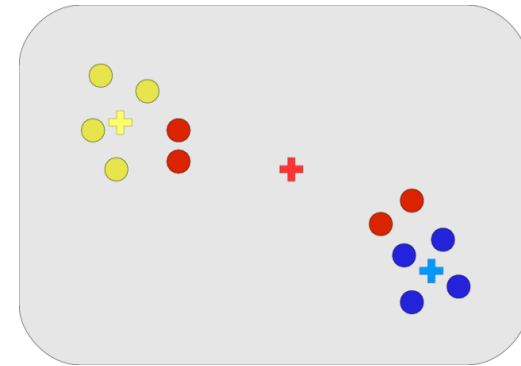
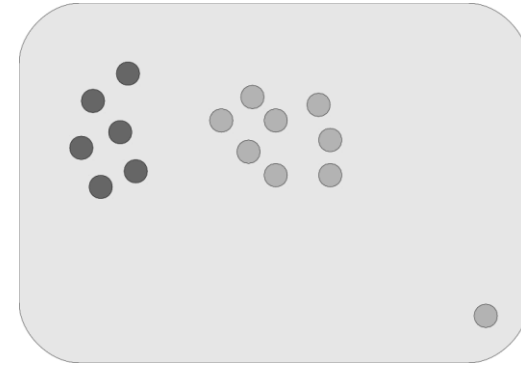- Choose best performing setting

- Important parameter!

- Knowledge about the data
  - Expert insights

- Development of RSS
  - Monotonous decline
  - Typically two points where decline slows down

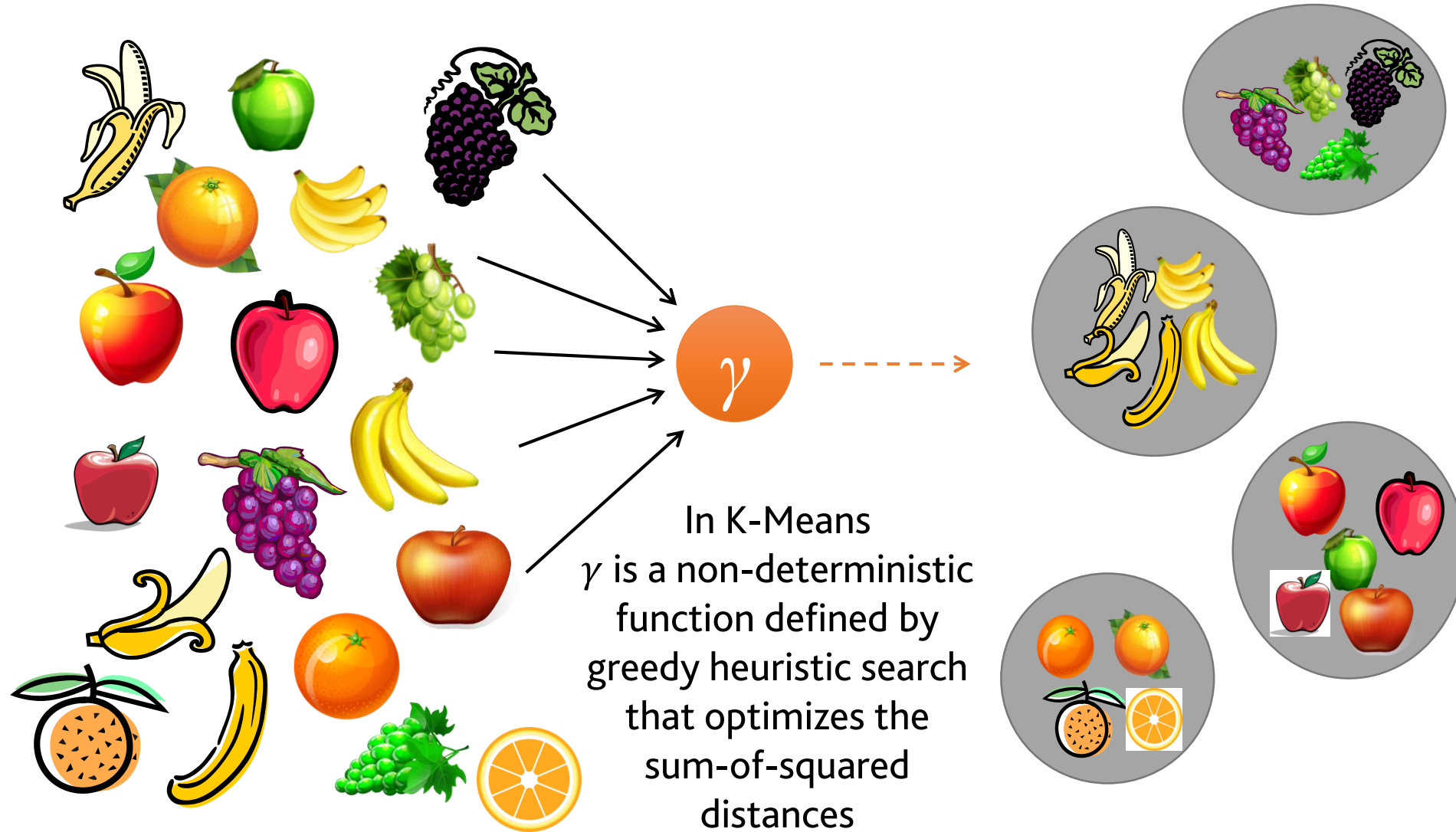- ## Outliers
  - ### Cause singleton clusters
  - ### Solution:
    - Remove and treat separately

- ## Empty clusters
  - ### Unlucky position of centroids
  - ### Solution:
    - Split large cluster

- ## Non-spheric shapes
  - ### Cannot be handled!

In K-Means
$\gamma$ is a non-deterministic
function defined by
greedy heuristic search
that optimizes the
sum-of-squared
distances

# Probabilistic Models

- Observations $D$ (our Data)
- Hidden (latent) parameters $\theta$

- Example: Throwing a coin: 10 x head, 2 x tail

- Observations $D$ (our Data)

- Hidden (latent) parameters $\theta$

- Example: Throwing a coin: 8 x head, 2 x tail
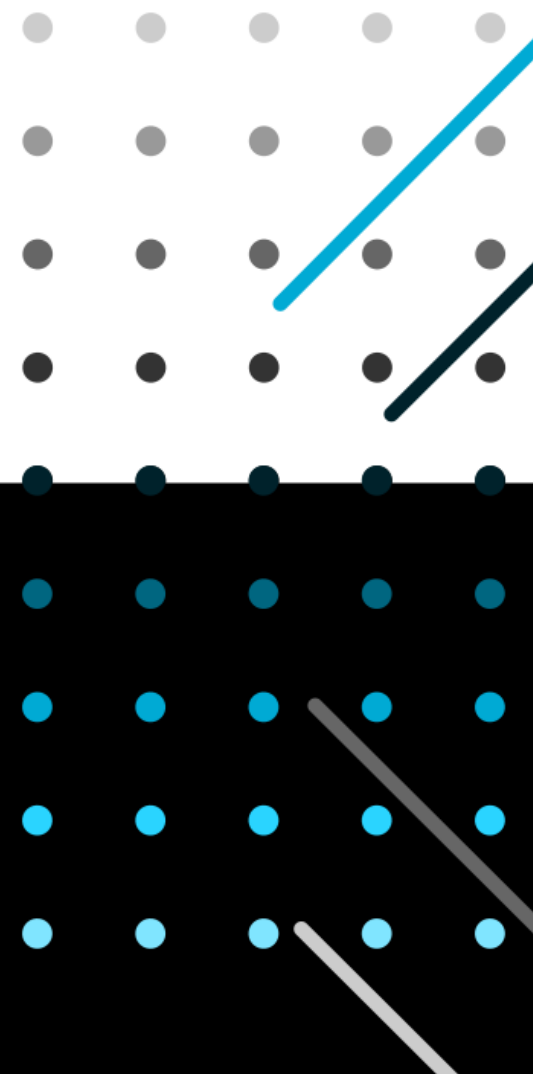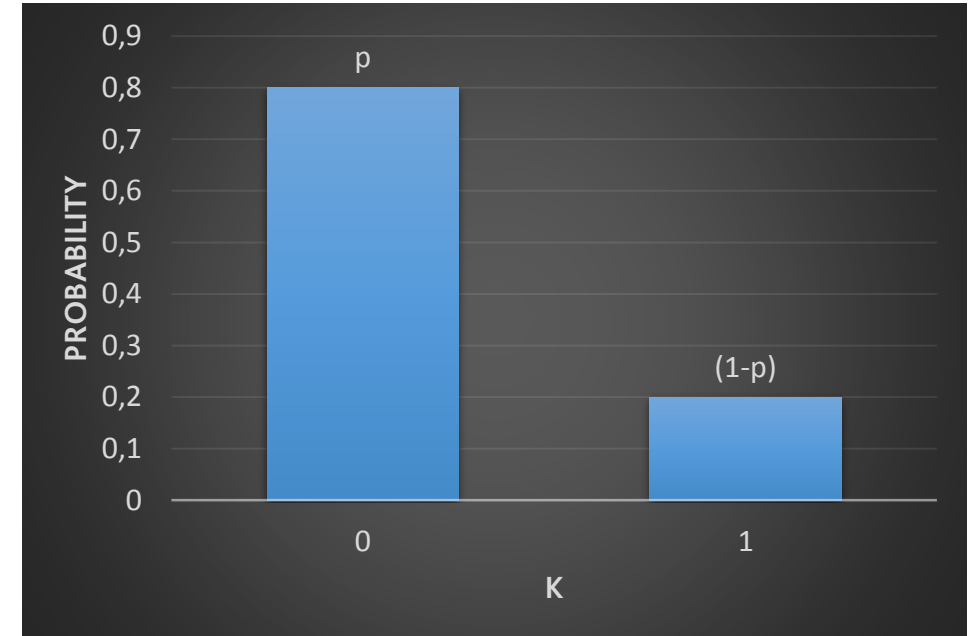
# Probabilistic models

- Observations $D$ (our Data)
- Hidden (latent) parameters $\theta$
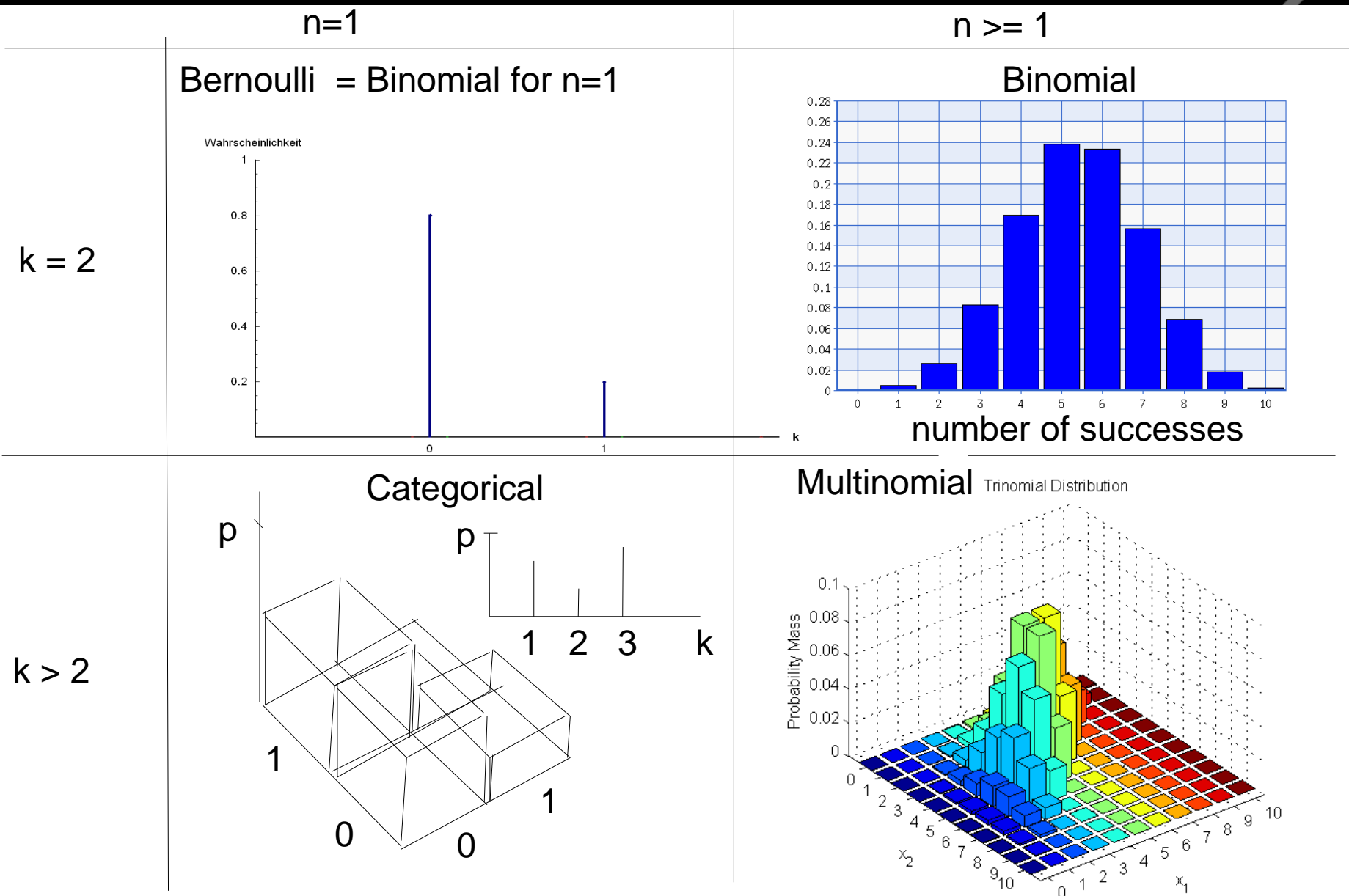
- Example: Throwing a coin: 8 x head, 2 x tail

- $\theta$: Bernoulli ($p$)

|  | n=1 | n >= 1 |
|---|---|---|
| k = 2 | Bernoulli = Binomial for n=1 | Binomial |
| k > 2 | Categorical | Multinomial |



number of successes

|  | n=1 | n >= 1 | n → ∞ |
|---|---|---|---|
| **k = 2** | Bernoulli = Binomial for n=1 | Binomial | Gaussian |
| **k > 2** | Categorical | Multinomial | Multivariate Gaussian |

- Observations $D$ (our Data)
- Hidden (latent) parameters $\theta$

- Example: Throwing a coin: 8 x head, 2 x tail

- $\theta$: Bernoulli ($p$)



How to estimate $p$?

- Observations $D$ (our Data)
- Hidden (latent) parameters $\theta$

- Example: Throwing a coin: 8 x head, 2 x tail

- $\theta$: Bernoulli $(p)$

- Likelihood: $L = P(D|\theta) = \prod_{d_i \in D} P(d_i|\theta)$

- Observations $D$ (our Data)
- Hidden (latent) parameters $\theta$

- Example: Throwing a coin: 8 x head, 2 x tail
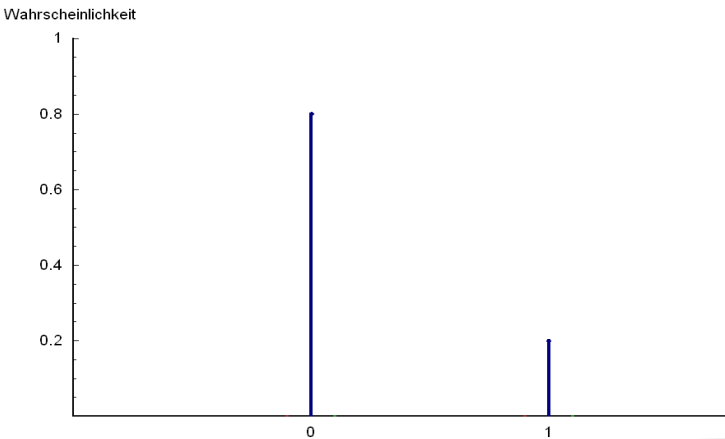
- $\theta$: Bernoulli ($p$)

- Likelihood: $L = P(D|\theta) = \prod_{d_i \in D} P(d_i|\theta)$

- Log-Likelihood: $\sum_{d_i \in D} \log(P(d_i|\theta))$

- Observations $D$ (our Data)
- Hidden (latent) parameters $\theta$

- Example: Throwing a coin: 8 x head, 2 x tail

- $\theta$: Bernoulli $(p)$

- Likelihood: $L = P(D|\theta) = \prod_{d_i \in D} P(d_i|\theta)$

- Log-Likelihood: $\sum_{d_i \in D} \log(P(d_i|p)) \quad |d_i \in \{H, T\}$

# Probabilistic models

- $\theta$: Bernoulli $(p)$

- Likelihood: $L = P(D|\theta) = \prod_{d_i \in D} P(d_i|\theta)$

- Log-Likelihood: $\sum_{d_i \in D} \log(P(d_i|p)) \quad |d_i \in \{H, T\}$

$$= n^T \cdot \log\big(P(T|p)\big) + n^H \cdot \log\big(P(H|p)\big)$$

$$= n^T \cdot \log(p) + n^H \cdot \log(1 - p)$$

$$\log L = n^T \cdot \log(p) + n^H \cdot \log(1 - p)$$

- Maximization:

$$\frac{\partial \log L}{\partial p} = \frac{n^T}{p} - \frac{n^H}{1 - p} = 0$$

$$\Leftrightarrow p = \frac{n^T}{n^T + n^H} = \frac{2}{2 + 8} = 0.2$$

# Expectation Maximization

- General clustering algorithm

- Characteristics:
  - Probabilistic approach
  - Soft assignments to clusters
  - Generalization of K-Means

- Parameters
  - $K$ : number of clusters
  - Initial random seed
  - Model for the distribution

- Data can be explained by a mixture of parametrized probability distributions – one per cluster.

- Data can be explained by a mixture of parametrized probability distributions – one per cluster.

- Problem: the true distributions are unknown, all we see is the data:

- Provide a model how the data is distributed, e.g.

  - Density of a normal distribution (1 dimensional)
  $$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$
    - Where $\mu$ is the mean and $\sigma$ the standard deviation

  - Density of a normal distribution (m-dimensional)
  $$f(\boldsymbol{x}) = \frac{1}{\sigma\sqrt{(2\pi)^m |S|}} e^{-\left(\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu})^T S(\boldsymbol{x}-\boldsymbol{\mu})\right)}$$
    - Where $\mu$ is an m-dimensional vector, $S$ an m×m covariance matrix and $|S|$ the determinant
  - Other distribution $P$ with parameters $\vartheta$

# Model estimation

- Estimate model parameters from data

- Maximum likelihood estimation:

  - 1 dimensional case:

$$\mu = \frac{1}{n}\sum_{x\in D} x_i \qquad \sigma = \sqrt{\frac{1}{n-1}\sum_{x\in D}(x_i - \mu)^2}$$

  - m-dimensional case:

$$\vec{\mu} = \frac{1}{n}\sum_{x\in D} \vec{x_i} \qquad S = \frac{1}{n-1}\sum_{x\in D}(\vec{x_i} - \vec{\mu})(\vec{x_i} - \vec{\mu})^T$$

  - Maximize log likelihood function:

$$\log(L(\vartheta|D)) = \log(\prod_{x\in D} P(\boldsymbol{x}|\vartheta)) = \sum_{x\in D}\log(P(\boldsymbol{x}|\vartheta))$$

But: we don't know which objects belongs to which cluster ...

- For each object $\boldsymbol{x}_i \in D$, we have a latent variables $z_{ij}$ modelling to which cluster $\omega_j$ it belong

- Two steps:
  - Expectation step
    - Calculate the expected values for $z_{ij}$ given the current model parameters $\vartheta$
    - So: how probable is it that object $\boldsymbol{x}_i$ belongs to cluster $\omega_j$ under the current model hypothesis $\vartheta$
  - Maximization step
    - Calculate the model parameters $\vartheta$ given the current estimates for the latent variables $z_{ij}$
    - So: how do the model parameters look like under the assumption that the assigment to clusters is correct.

- Iterate until convergence (little change)

- Problem:
  - Expectation step needs model parameters to estimate latent variables
  - Maximization step needs latent variables to estimate model parameters

- Different options:
  - Hand selected initial model parameters
  - Random initialization
  - Choose random individuals
  - Perform k-means cluster to find initial clusters

# Example

- Cluster people by their height
  - Two classes
  - Assume initial model:
    - $\mu_1 = 110 \qquad \sigma_1 = 20$
    - $\mu_2 = 160 \qquad \sigma_2 = 20$

- Expectation:
  - $f_{\mu_1,\sigma_1}(124) = 0.01561$
  - $f_{\mu_2,\sigma_2}(124) = 0.00395$
  - Weights:
    - $z_{1,1} = 0.7982$
    - $z_{1,2} = 0.2018$
  - ...

| Gender | height |
|:------:|:------:|
| F | 124 |
| F | 115 |
| F | 121 |
| F | 139 |
| F | 98 |
| F | 135 |
| F | 131 |
| M | 170 |
| M | 166 |
| M | 155 |
| M | 167 |
| M | 158 |
| M | 175 |
| M | 143 |
| M | 163 |
| M | 160 |
| M | 145 |
| M | 176 |

- Given all weights:
- New model parameters
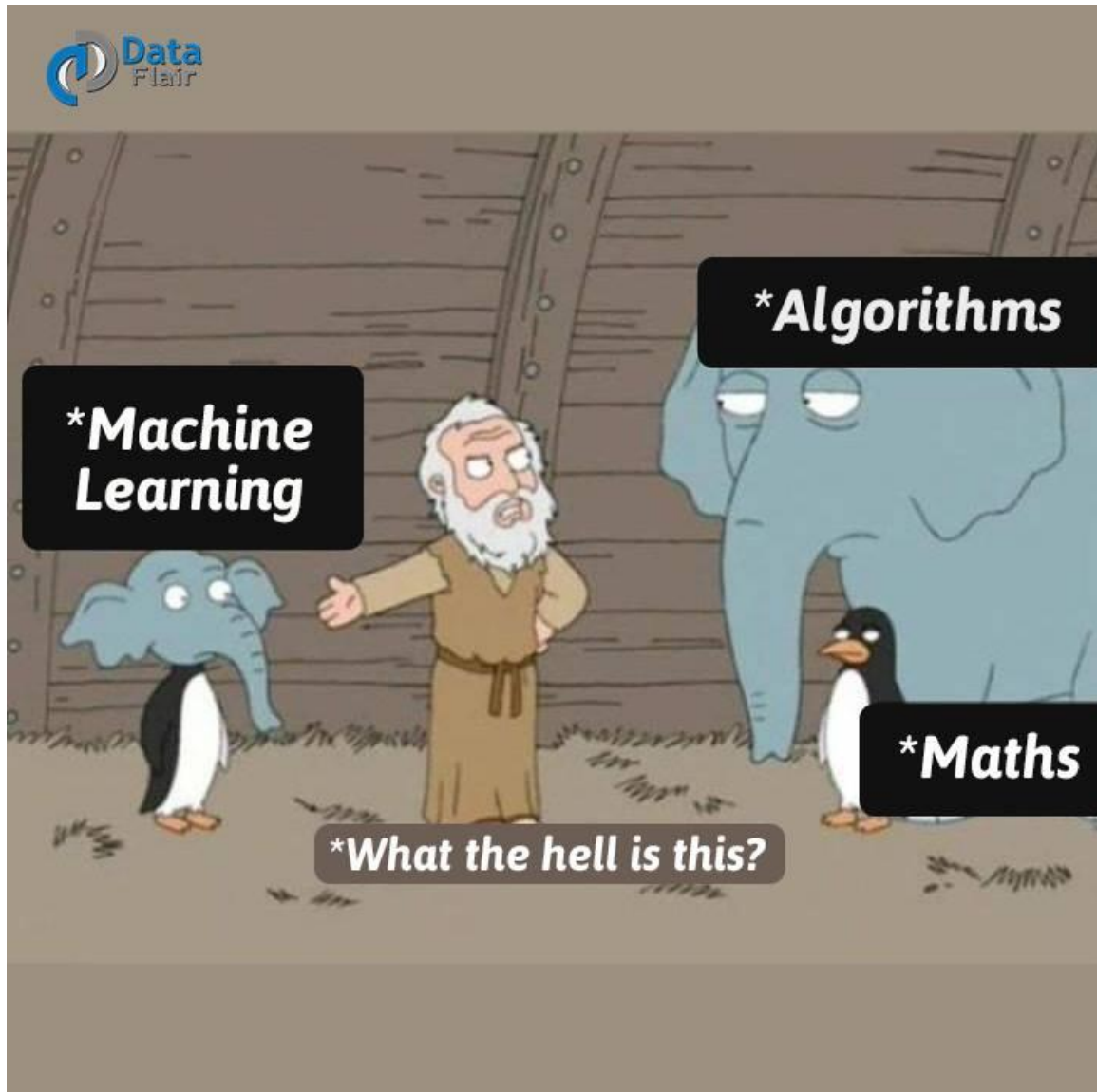
$$\mu_j = \frac{\sum_{x \in D} z_{i,j} \, x_i}{\sum_{x \in D} z_{i,j}}$$

$$\sigma_j = \sqrt{\frac{\sum_{x \in D} z_{i,j} (x_i - \mu)^2}{\sum_{x \in D} z_{i,j}}}$$

- Values:
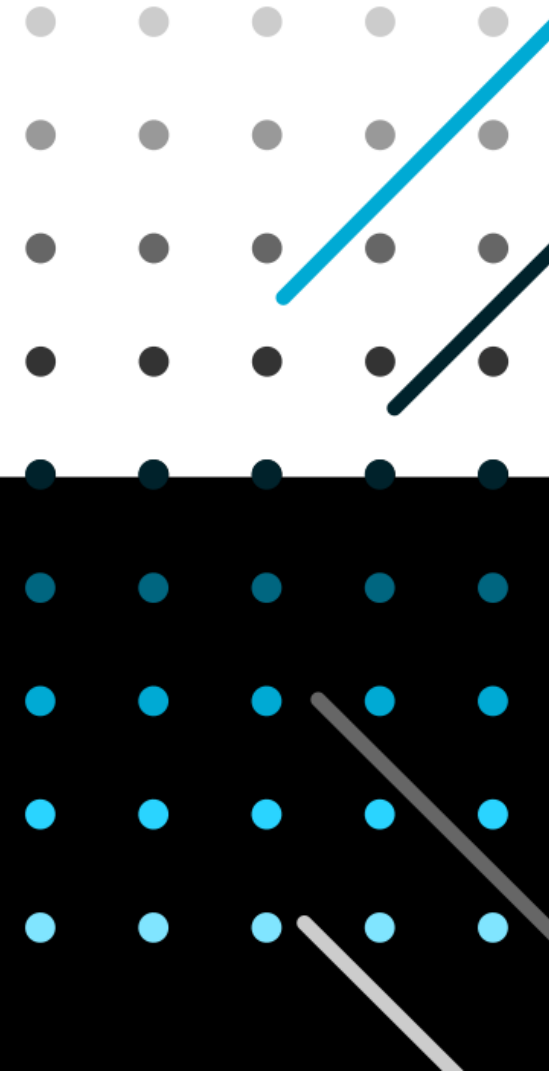  - $\mu_1 = 123.72$      $\sigma_1 = 15.98$
  - $\mu_2 = 157.72$      $\sigma_2 = 14.62$

| Gender | height |
|:------:|:------:|
| F | 124 |
| F | 115 |
| F | 121 |
| F | 139 |
| F | 98 |
| F | 135 |
| F | 131 |
| M | 170 |
| M | 166 |
| M | 155 |
| M | 167 |
| M | 158 |
| M | 175 |
| M | 143 |
| M | 163 |
| M | 160 |
| M | 145 |
| M | 176 |

# Summary

- Clustering
  - K-Means
  - Expectation-Maximization algorithm

# Thank you!

**Zeyd Boukhers**

E-mail:    Boukhers@uni-koblenz.de
Phone:    +49 (0) 261 287-2765
Web:      Zeyd.Boukhers.com

University of Koblenz-Landau
Universitätsstr. 1
56070 Koblenz