# Machine Learning
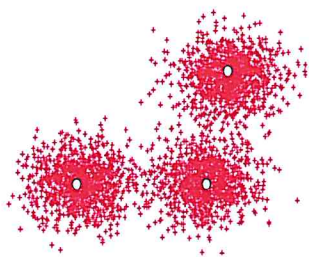## – winter term 2016/17 –

# Chapter 05: Clustering

Prof. Adrian Ulges
Masters "Computer Science"
DCSM Department
University of Applied Sciences RheinMain

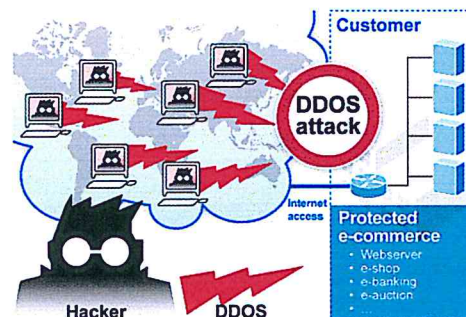# Unsupervised Learning = Learning without Labels <sub>images from [2], [1]</sub>

- ▶ **Clustering**: discover coherent groups of samples
- ▶ **Dimensionality reduction**: compressing samples
- ▶ **Itemset mining**: finding frequent substructures in the data
- ▶ **Anomaly detection**: detecting outliers in the data



Customers Who Bought This Item Also Bought

slide:ology: The Art and Science of Creating Grea... by Nancy Duarte
★★★★☆ (98)
$23.09

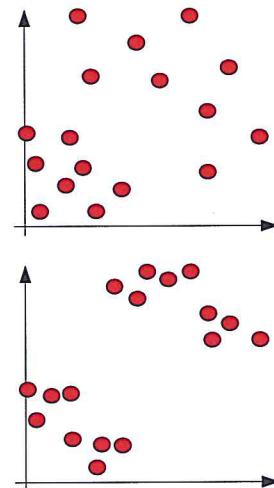The Naked Presenter: Delivering Powerful Present... by Garr Reynolds
$16.49

# Outline

3

# Clustering: Definition

- Clustering (or *cluster analysis*) is an unsupervised learning problem (*remember:* **samples only**, *no labels*)
- The challenge is to discover coherent subgroups (or *clusters*) of samples
- **Difference to classification**: In clustering, we try to *find* the classes <u>and</u> assign samples to them



## Challenges

1. Often, it is unclear by which criterion to cluster (*example: cluster users, but by which demographic attributes?*)
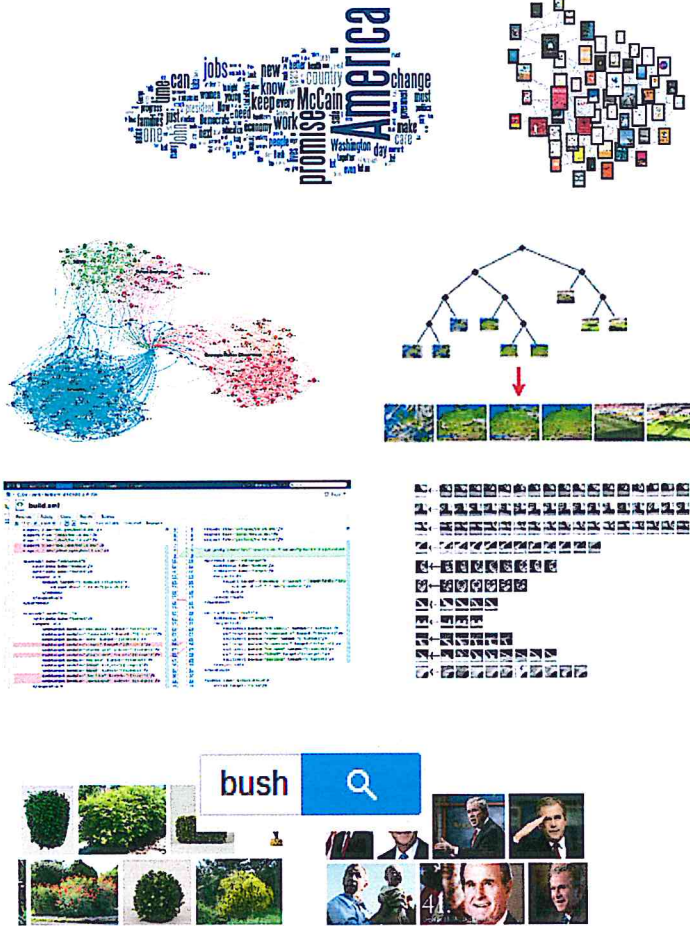2. Cluster granularity is unclear a priori
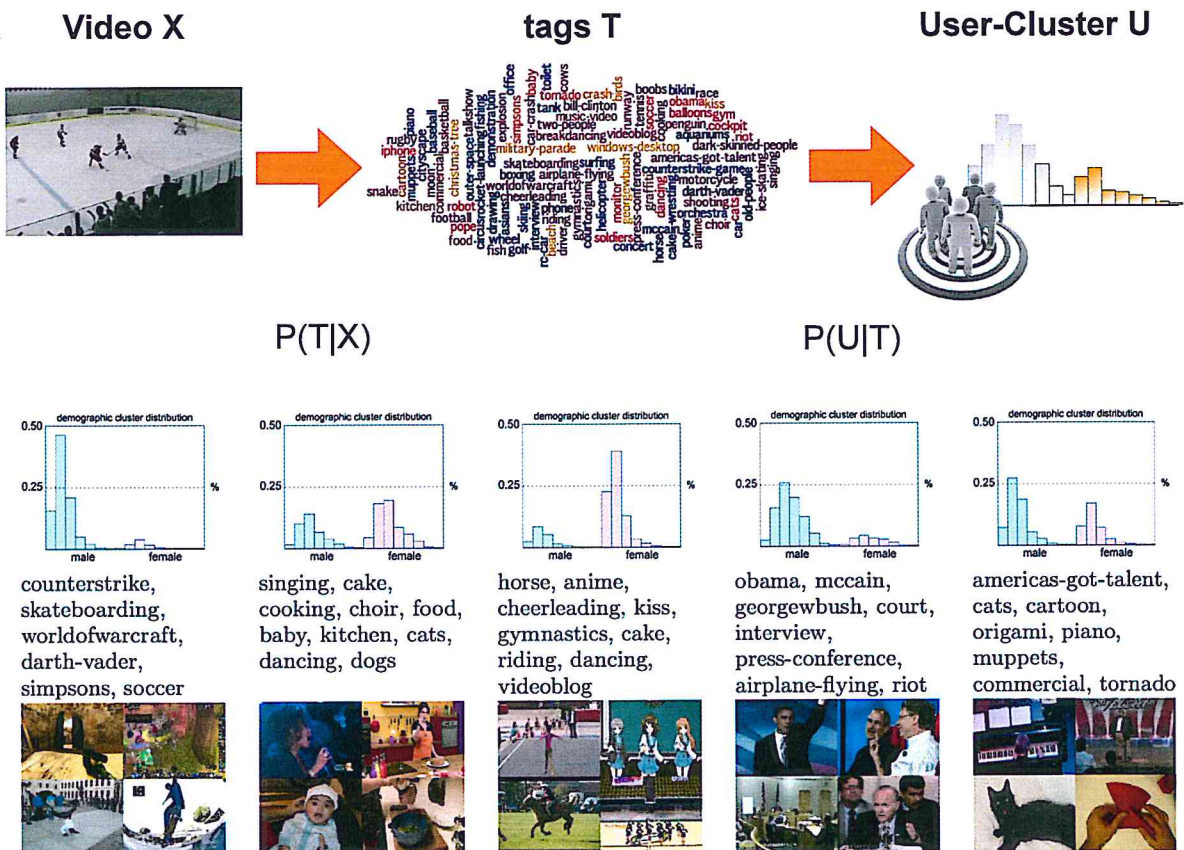
4

# Clustering: Applications

Clustering has nume-
rous **applications** in
various areas

- ▶ market research
- ▶ life sciences
- ▶ information
  retrieval
- ▶ computer vision
- ▶ social networks
- ▶ data mining

---

# Example: Demographic Clustering on YouTube [8]

| **Video X** | **tags T** | **User-Cluster U** |
|---|---|---|
| | $P(T\|X)$ | $P(U\|T)$ |



counterstrike,
skateboarding,
worldofwarcraft,
darth-vader,
simpsons, soccer

singing, cake,
cooking, choir, food,
baby, kitchen, cats,
dancing, dogs

horse, anime,
cheerleading, kiss,
gymnastics, cake,
riding, dancing,
videoblog

obama, mccain,
georgewbush, court,
interview,
press-conference,
airplane-flying, riot

americas-got-talent,
cats, cartoon,
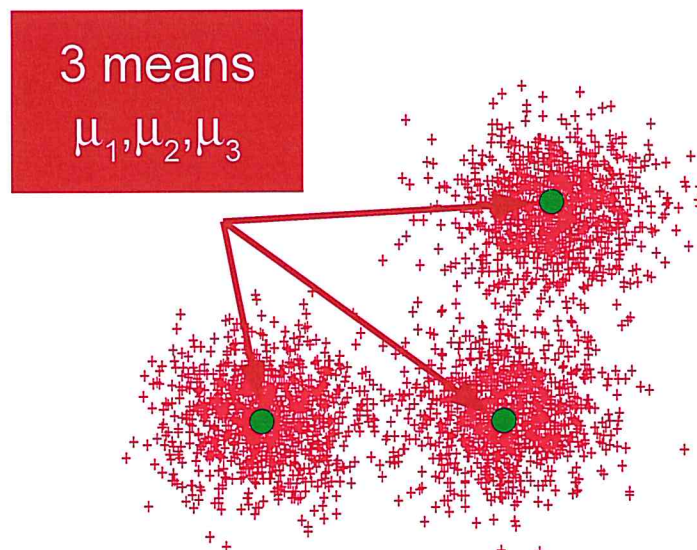origami, piano,
muppets,
commercial, tornado

# Outline

# Clustering: K-Means

We start with the "first choice" clustering algorithm: **K-Means**

- Given: samples $\mathbf{x}_1, ..., \mathbf{x}_n \in \mathbb{R}^d$
- We assume that samples are clustered around $K$ centers (the "K means") $\mu_1, ..., \mu_K \in \mathbb{R}^d$
- Each sample $\mathbf{x}_i$ belongs to a mean $k(i)$
- The clusters are **spheres** of **identical size**



3 means
$\mu_1, \mu_2, \mu_3$
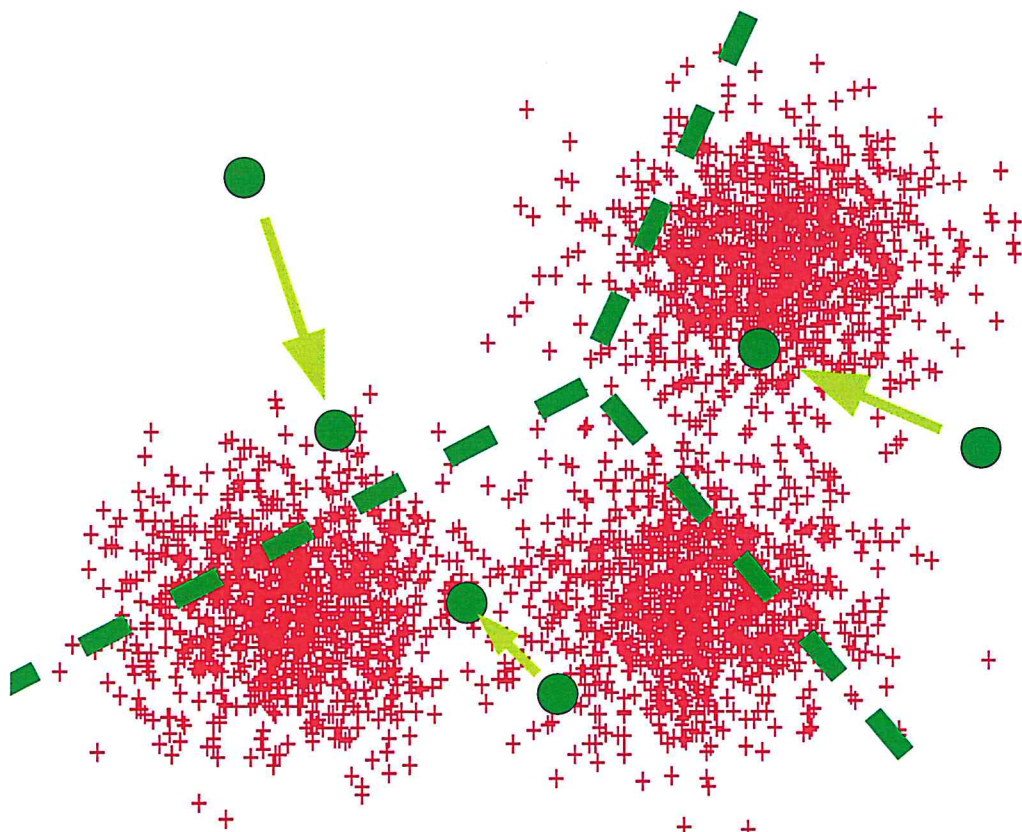
# K-Means: Approach

When trying to determine the clusters / the means, we face a **chicken-egg problem**

- ▶ If we knew the clusters, we could easily determine the means *(by averaging all samples of a cluster)*
- ▶ If we knew the means, we could determine the clusters *(by assigning each sample to its closest mean)*
- ▶ Approach (**interleaved optimization**): Alternately, fix the clusters/means and estimate the other

```
1   function KMEANS(x_1, ..., x_n, K)
2       initialize μ_1, ..., μ_K by random sampling from x_1, ..., x_n
3       repeat
4           for i = 1, .., n:        // assign each sample to its closest cluster
5               k(i) := arg min_{k=1,...,K} ‖x_i − μ_k‖
6           for k = 1, ...K:         // re-estimate each cluster's mean
7               X_k := {x_i | k(i) = k}
8               μ_k := 1/|X_k| ∑_{x∈X_k} x
9       until k(1), ..., k(n) do not change
10      return μ_1, ..., μ_K
```
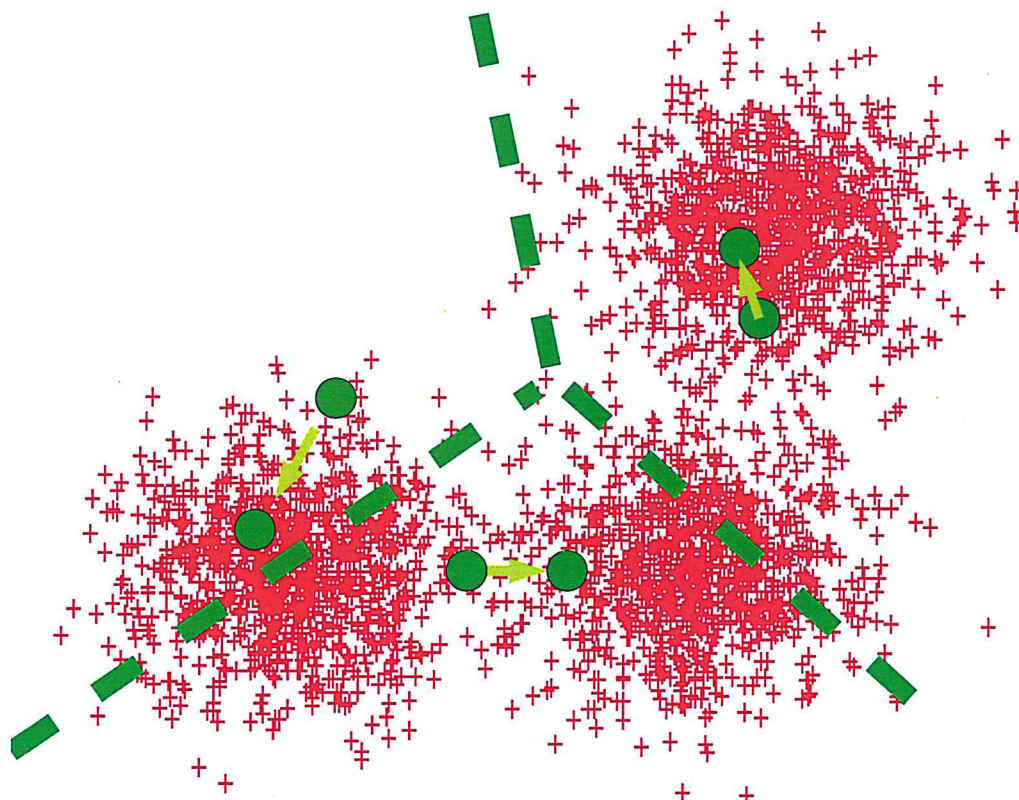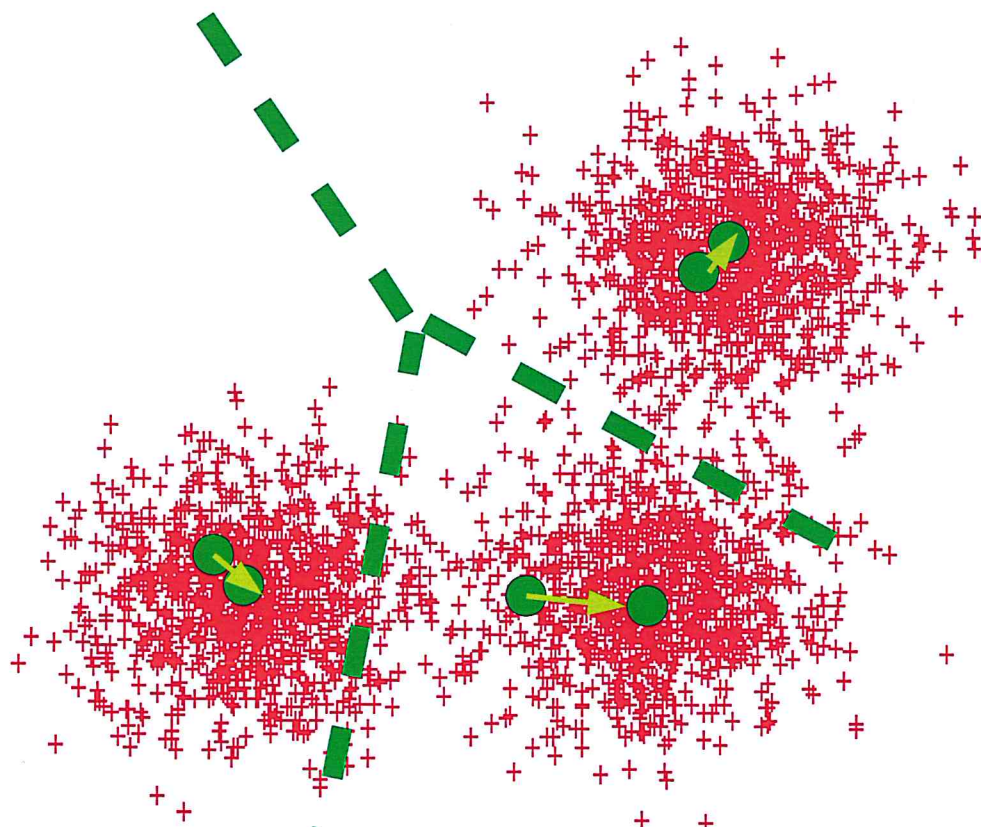
# K-Means: Example (Step 1)

# K-Means: Example (Step 3...)

# K-Means: Properties

- ▶ K-Means corresponds to a local optimization of the sum of squared errors

$$E(\mu_1, ..., \mu_K) = \sum_{i=1}^{n} (\mathbf{x}_i - \mu_{k(i)})^2$$

- ▶ Computational effort: $O(K \cdot n \cdot d)$ per iteration. The number of iterations is often moderate.
- ▶ Convergence is guaranteed.

### Proof of Convergence

$$E_0 \underset{(1)}{\geqslant} E_0' \underset{(2)}{\geqslant} E_1 \geqslant E_1' \geqslant E_2 \geqslant E_2' \ldots \quad \geqslant 0$$

$\underbrace{\phantom{E_0 \geqslant E_0'}}$ re-assign samples to clusters

$\underbrace{\phantom{E_0' \geqslant E_1}}$ re-estimate cluster centers
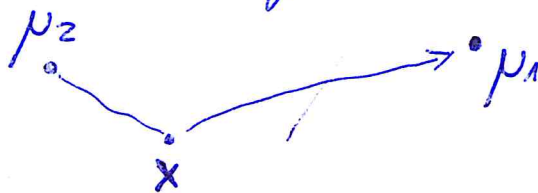
This sequence con-verges (monotonously decreasing + lower bound)

---

# K-Means: Properties

→ The KMeans Algorithm Converges!

### Proof of Convergence (cont'd)

(1) $E_k \geqslant E_k' \quad \forall k$

because for each sample $x_i$ the new center $\mu_{k'(i)}$ is at least as close as $\mu_{k(i)}$

(2) $E_k' \geqslant E_{k+1} \quad \forall k$

because $\bar{x} = \underset{y}{\arg\min} \sum_i \|x_i - y\|^2$

# K-Means: Properties

Proof of Convergence (cont'd)

---

# K-Means: Properties
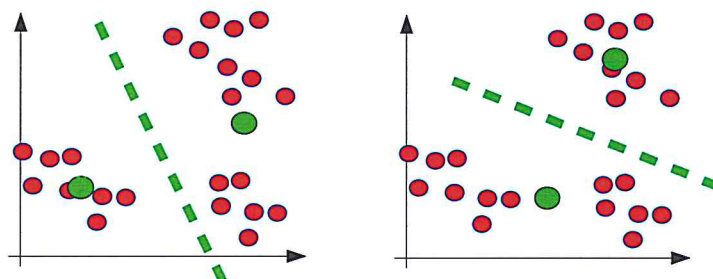
Does K-Means always lead to the same results?

*No: K-Means is a <u>local</u> search method!*

▶ **Problem 1**: The order of means can be permuted

$$\mu_1 = (0,0), \mu_2 = (1,1), \mu_3 = (5,3)$$
$$\mu_1 = (5,3), \mu_2 = (0,0), \mu_3 = (1,1)$$

▶ **Problem 2**: The resulting means can be completely different

▶ **Approach**: Restart multiple times, and keep the result with minimal error $E$.

▶ During the algorithm, **empty clusters** may occur. **Approach**: Reinitialize the corresponding center randomly and continue.
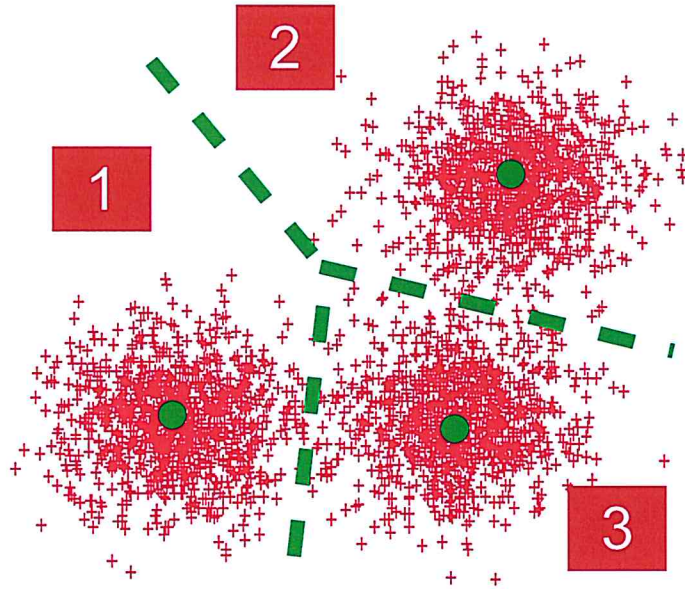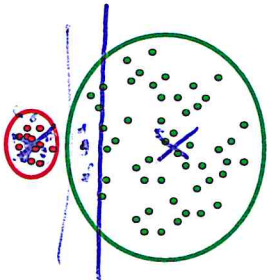
# K-Means: Properties (cont'd)

Given a clustering result $\mu_1, ..., \mu_K$, we can assign new samples **x** to clusters (this is called **vector quantization**):

$$k(\mathbf{x}) = \arg\min_k ||\mathbf{x} - \mu_k||$$

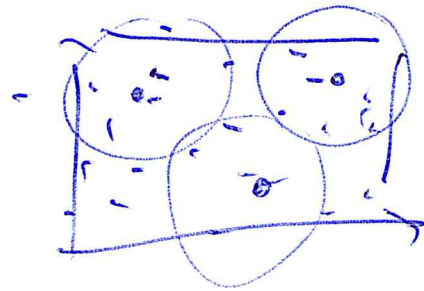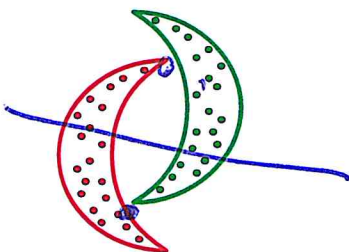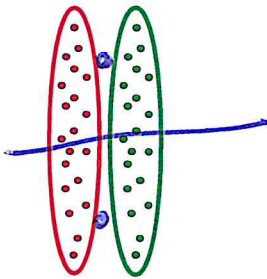# K-Means: Discussion



K ?

+ Simple, fast

− *local search*

# Outline

# Choosing $K$: Model Selection

*"Model selection is the task of selecting a statistical model from a set of candidate models, given data."*

(en.wikipedia.org)

Here: Model Selection = Choosing $K$

- $K$ too small (*undersegmentation*): clusters too diverse
- $K$ too high (*oversegmentation*): too many parameters, clusters too fine-grain
- Choosing the 'wrong' $K$ leads to **instable results**
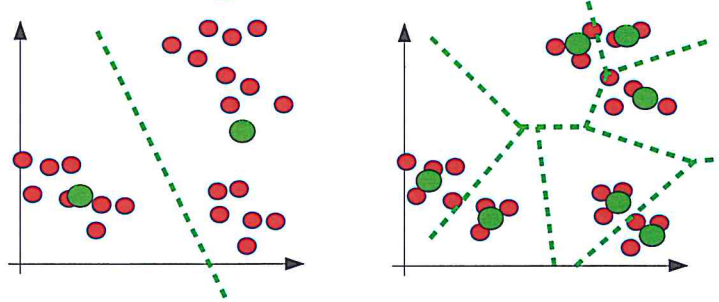
Approach 1: External Benchmark

- Sometimes, clustering is just one processing step of a **larger system**, and we can benchmark that larger system
- **Example**: User clustering for advertising ($\rightarrow$ *benchmark by click-through-rate*)

Goal: measure a model's **goodness-of-fit** without labels



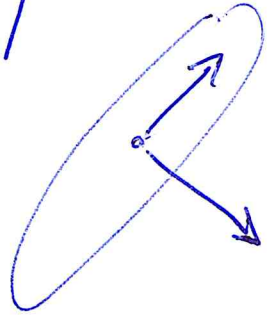Example: The **Bayes' Information Criterion (BIC)**

1. The clusters should be **compact** (small error E)
2. The model should be simple, i.e. have only **few parameters**

▸ Let $\theta$ be the model parameters to learn, and let $\#\theta$ be their number (e.g., in K-Means: $\#\theta = K \cdot d$)

▸ Test different values of $K$, and pick this one:

$$K^* = \arg\min_{K} \; -2ln\Big(p(\mathbf{x}_1, ..., \mathbf{x}_n|\theta)\Big) + \#\theta \cdot ln(n)$$

# BIC for K-Means: Derivation ✳

cluster centers

$$-2 \cdot ln\Big( p(x_1, ..., x_n | \theta) \Big)$$

independence

$$= -2 \cdot ln\Big( \prod_i p(x_i | \theta) \Big)$$

multivariate normal dists.
$$\Sigma = I$$

$$= -2 \cdot ln\Big( \prod_i \frac{1}{2\pi^{d/2}} e^{-\frac{1}{2} \frac{\|x_i - \mu_{K(i)}\|^2}{\text{......}}} \Big)$$

$$= -2\Big( \sum_i ln(\frac{1}{2\pi^{d/2}}) - \frac{1}{2}\|x_i - \mu_{K(i)}\|^2 \Big)$$

$$\underbrace{\qquad\qquad}_{const.}$$

$$= \sum_i \|x_i - \mu_{K(i)}\|^2 = E$$

# The Bayes Information Criterion

$$K^* = \arg \min_K \; \underbrace{\sum_{i=1}^{n} \left( \mathbf{x}_i - \mu_{k(i)} \right)^2}_{\mathsf{E(K)}} + \underbrace{K \cdot d \cdot ln(n)}_{\text{model complexity}}$$
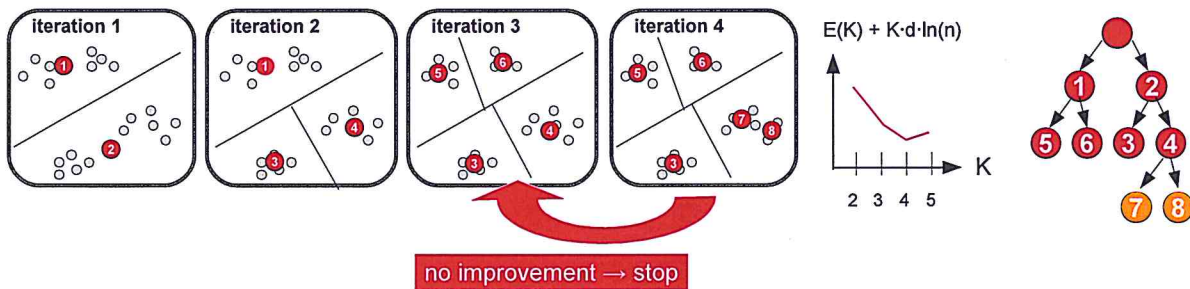
# Selecting $K$: Search Strategies

## Approach 1: **Naive**

- ▶ test values for $K$ in a reasonable range.
- ▶ For every $K$, re-run clustering and evaluate (**expensive!**)

## Approach 2: **Hierarchical Clustering** *(more efficient)*

- ▶ ... Iteratively, pick the largest cluster
- ▶ ... and apply $K$-Means to the samples in this cluster, obtaining $K$ new clusters
- ▶ ... stop once the overall quality (e.g., BIC) stops improving
- ▶ We obtain a **tree** of clusters

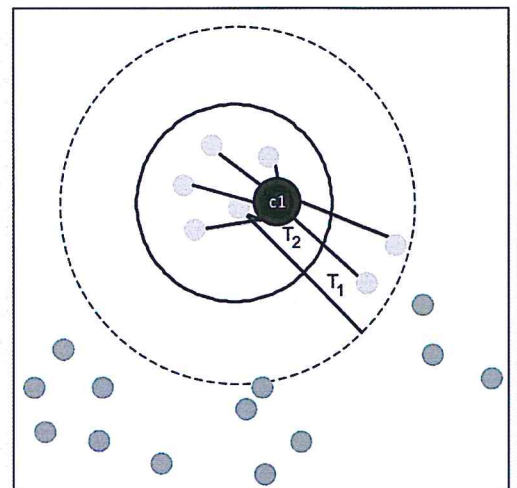# Selecting $K$: Canopy Clustering image from [7]

## Approach 3: Canopy Clustering

- ▶ A **greedy strategy** to find (potentially suboptimal) clusters on large datasets
- ▶ We use it to estimate $K$ and to initialize the means
- ▶ Canopy clusters can **overlap**!
- ▶ Canopy clustering uses **two thresholds**
  - ▶ $T_1$ (determines the number of clusters)
  - ▶ $T_2$ (determines the overlap of clusters) ($T_2 > T_1$)

```
1   function CLUSTER_CANOPY(X := {x_1, ..., x_n})
2       C := {}
3       while X <> {}:
4           choose a random sample x ∈ X
5           Y := {y ∈ X | ‖y − x‖ ≤ T_1}
6           Z := {y ∈ X | T_1 < ‖y − x‖ ≤ T_2}
7           C := C ∪ {x}
8           X := X\Y
9       return C
10
```
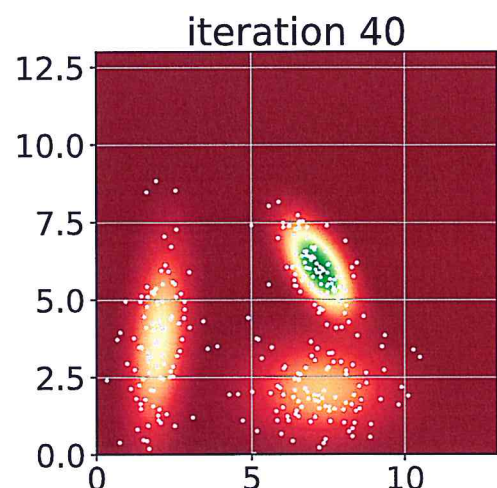
# Outline

# Expectation Maximization (EM)

▶ We can overcome some of the above limitations by **generalizing K-Means**, resulting in a famous approach called **Expectation Maximization (EM)**

## EM: Model

▶ We explain the data $\mathbf{x}_1, ..., \mathbf{x}_n$ by a **Gaussian mixture model**

$$\mathbf{x}_1, ..., \mathbf{x}_n \sim \sum_{k=1}^{K} P_k \cdot p(\mathbf{x}|\mu_k, \Sigma_k)$$

where $p$ is the **multivariate normal density** *(Chapter 3)*, $\mu_1, ..., \mu_K$ are $K$ centers, $\Sigma_1, ..., \Sigma_K$ are $K$ covariance matrices (the *shapes* of the clusters), and $P_1, ..., P_K$ are the cluster's proportions of the data (also called *priors*).



iteration 40

# Expectation Maximization (EM)

## Remarks

▶ In K-Means, we would have $P_1 = P_2 = \ldots = P_K = \frac{1}{K}$ and

$$\Sigma_1 = \Sigma_2 = \ldots = \Sigma_K = \begin{pmatrix} \sigma^2 & 0 & \ldots & 0 \\ 0 & \sigma^2 & \ldots & 0 \\ 0 & \ldots & \ldots & 0 \\ 0 & \ldots & 0 & \sigma^2 \end{pmatrix}$$
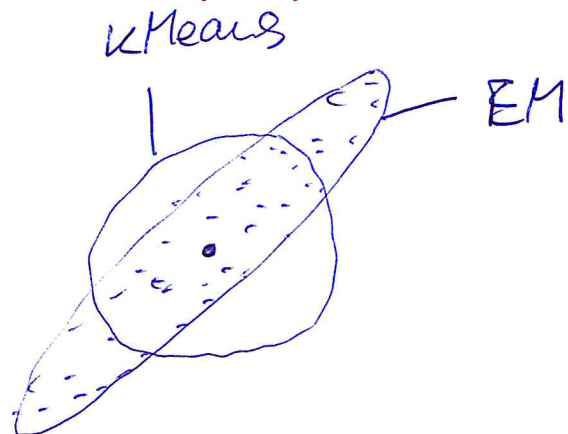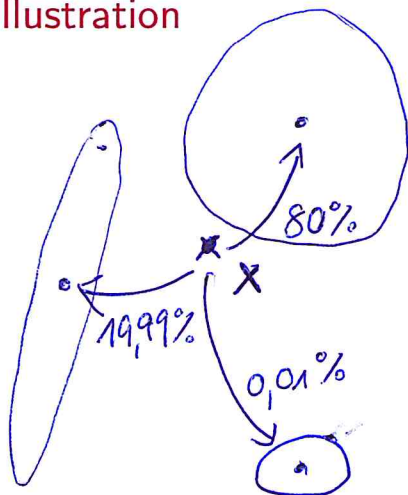
## Approach

▶ We **rename** the two alternating K-Means steps

E-Step  Re-assigning samples to clusters $\rightarrow$ "Expectation-Step"

M-Step  Re-estimating the cluster centers $\rightarrow$ "Maximization-Step"

▶ We **modify** these steps a bit

  ▶ **E-Step**: No hard assignment of samples to centers, but a **soft assignment** by computing the probability $P(k(i) = k \mid \mathbf{x}_i)$

  ▶ **M-Step**: Do not only estimate the cluster *centers*, but **parameters** in general (e.g., the clusters' shape+prior)

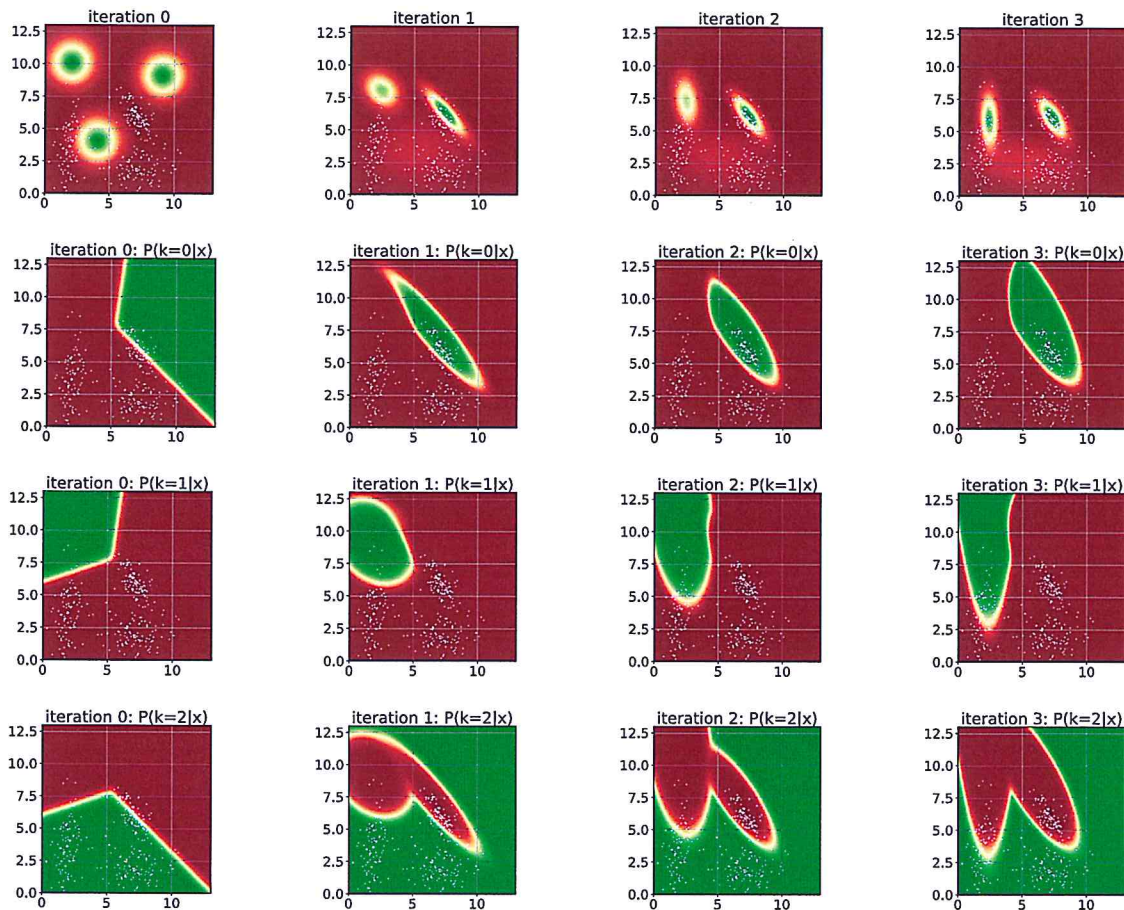# K-Means vs. Expectation Maximization (EM)

Illustration



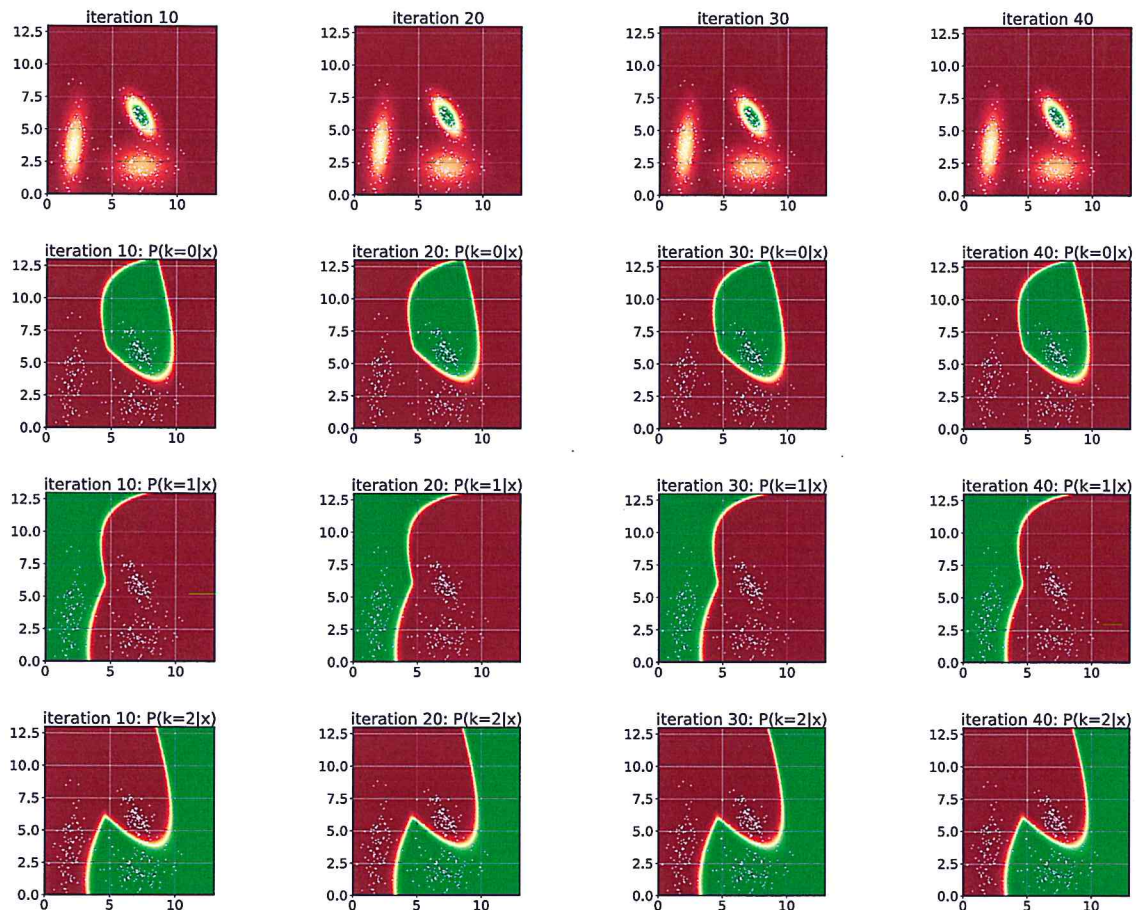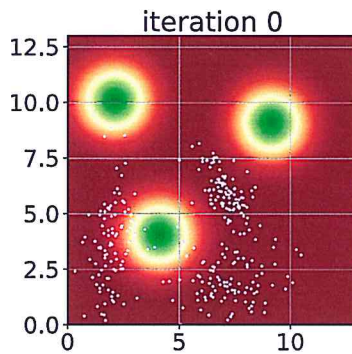| | K-Means | EM |
|---|---|---|
| E-Step | $k(i) := \arg\min_k \|\mathbf{x}_i - \mu_{k(i)}\|$ | $w_{ki} := P(k(i) = k \mid \mathbf{x}_i) = \frac{p(\mathbf{x}_i; \mu_k, \Sigma_k)}{\sum_{k'} p(\mathbf{x}_i; \mu_{k'}, \Sigma_{k'})}$ |
| M-Step | $\mu_k := \frac{\sum_{\mathbf{x} \in X_k} \mathbf{x}}{|X_k|}$ | $\mu_k := \frac{\sum_i w_{ki} \cdot \mathbf{x}_i}{\sum_i w_{ki}}$ |
| | — | $\Sigma_k := \frac{\sum_i w_{ki} \cdot (\mathbf{x}_i - \mu_k)(\mathbf{x}_i - \mu_k)^{\top}}{\sum_i w_{ki}}$ |
| | — | $P_k := \frac{\sum_i w_{ki}}{\sum_{k'} \sum_i w_{k'i}}$ |

# EM: Example

# EM: Example

# EM: Goodness-of-Fit
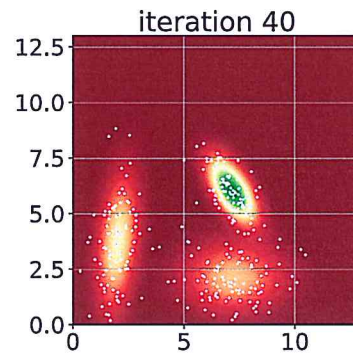
- ▶ Goal: restart EM many times, pick the 'best' model.
- ▶ Given an **EM model** $\Theta = (\mu_1, ..., \mu_K, \Sigma_1, ..., \Sigma_K, P_1, ..., P_K)$, we want to measure its **"goodness-of-fit"**.
- ▶ Approach: We measure the **likelihood** of the data

$$L(\mathbf{x}_1, ..., \mathbf{x}_n; \Theta) = \prod_i p(\mathbf{x}_i | \Theta)$$

$$= \prod_i \sum_k P_k \cdot p(\mathbf{x}_i; \mu_k, \Sigma_k)$$



low likelihood                                      high likelihood

# EM: Discussion

# EM as a general Learning Scheme

▶ EM for **Gaussian Mixture Models** is just a special case!

| symbol | general EM | **Gaussian Mixture Models** |
|--------|-----------|----------------------------|
| $X$ | (known) input data | the features $\mathbf{x}_1, ..., \mathbf{x}_n$ |
| $\Theta$ | parameters | means $\mu_1, ..., \mu_K$, shapes $\Sigma_1, ..., \Sigma_K$, priors $P_1, ..., P_K$ |
| $U$ | unknown data | the mapping from $\mathbf{x}_i$ to clusters $k$ |

## EM: General Learning Scheme

```
1  function EM(X)
2      initialize Θ randomly
3      repeat
4          compute P(U|X,Θ)                          // E-step
5          optimize parameters [6], obtaining a new Θ   // M-step
6      until convergence
7      return Θ
8
```

# Outline