

*** Applied Machine Learning Fundamentals ***

Clustering

M. Sc. Daniel Wehner

SAP SE

Winter term 2019/2020



Find all slides on [GitHub](#)

Lecture Overview

Unit I	Machine Learning Introduction
Unit II	Mathematical Foundations
Unit III	Bayesian Decision Theory
Unit IV	Probability Density Estimation
Unit V	Regression
Unit VI	Classification I
Unit VII	Evaluation
Unit VIII	Classification II
Unit IX	Clustering
Unit X	Dimensionality Reduction

Agenda for this Unit

① Introduction

Procedure

② Wrap-Up

Summary

Self-Test Questions

Lecture Outlook

Recommended Literature and further Reading

Section:
Introduction



k -Means: Introduction

- k -Means is a **clustering** algorithm and as such **unsupervised**
- A clustering algorithm tries to **find structure** in the data
- k defines the number of clusters to be found
- Once the clusters are found, they first have to be interpreted...
- ...and can then be used for prediction purposes

A cluster must be **internally homogeneous**, but simultaneously **externally heterogeneous**. (Elements of one cluster have to be very similar, but must differ significantly from elements in other clusters.)

k -Means: Example Use Cases

- **Behavioral segmentation**
 - Customer segmentation (e. g. [sinus milieus](#))
 - Creating profiles based on activity monitoring
- **Sorting sensor measurements**
 - Image grouping
 - Detection of activity types in motion sensors
- **Inventory categorization**
 - Grouping inventory by sales activity
 - Grouping inventory by manufacturing metrics
- Many, many more, ...

k -Means: Procedure

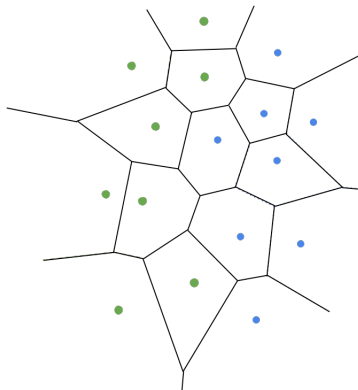
- **Vector quantization**
 - Represent data points by a single vector (here: called **centroid**) which is close to them
 - This is useful for **compression**!
- **How to:** Create k partitions ($\hat{=}$ clusters) of the data set \mathcal{D} , such that the sum of squared deviations from the cluster centroids is **minimal**:

$$\min_{\mu_j} \sum_{j=1}^k \sum_{\mathbf{x}^{(i)} \in \mathcal{D}_j} \|\mathbf{x}^{(i)} - \mu_j\|^2 \quad (1)$$

- With $\mathcal{D}_j \equiv j^{th}$ cluster, $\mu_j \equiv$ centroid of j^{th} cluster

Result: Voronoi Diagram

- The dots represent cluster centroids
- The lines visualize the **cluster boundaries**
- For a new point we can easily determine to which cluster it has to be assigned



k -Means Algorithm

- Input: $\mathcal{D} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)}\} \in \mathbb{R}^{n \times m}$, Number of clusters k
- Algorithm:

① $t \leftarrow 1$

② Randomly choose k means $\mu_1^{\langle 1 \rangle}, \mu_2^{\langle 1 \rangle}, \dots, \mu_k^{\langle 1 \rangle}$

③ While not converged:

3a Assign each $\mathbf{x}^{(i)} \in \mathcal{D}$ to the closest cluster:

$$\mathcal{D}_j^{\langle t \rangle} = \left\{ \mathbf{x}^{(i)} : \|\mathbf{x}^{(i)} - \mu_j^{\langle t \rangle}\|^2 \leq \|\mathbf{x}^{(i)} - \mu_{j^*}^{\langle t \rangle}\|^2; \forall j^* = 1, 2, \dots, k; \mathbf{x}^{(i)} \in \mathcal{D} \right\}$$

3b Update cluster centroids μ_j :

$$\mu_j^{\langle t+1 \rangle} = \frac{1}{|\mathcal{D}_j^{\langle t \rangle}|} \sum_{\mathbf{x}^{(i)} \in \mathcal{D}_j^{\langle t \rangle}} \mathbf{x}^{(i)}$$

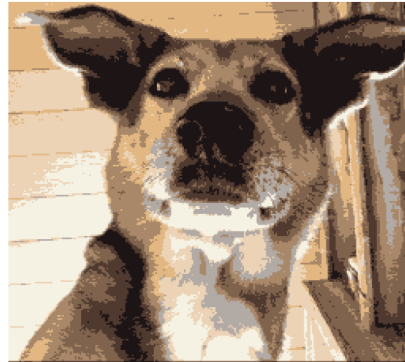
3c $t \leftarrow t + 1$

Image Compression

Original image



Compressed image



Section:
Wrap-Up



Summary





Self-Test Questions

1

What's next...?

Unit I	Machine Learning Introduction
Unit II	Mathematical Foundations
Unit III	Bayesian Decision Theory
Unit IV	Probability Density Estimation
Unit V	Regression
Unit VI	Classification I
Unit VII	Evaluation
Unit VIII	Classification II
Unit IX	Clustering
Unit X	Dimensionality Reduction

Recommended Literature and further Reading

Thank you very much for the attention!

Topic: *** Applied Machine Learning Fundamentals *** Clustering

Term: Winter term 2019/2020

Contact:

M.Sc. Daniel Wehner

SAP SE

daniel.wehner@sap.com

Do you have any questions?