

Machine Learning

Lecture 10 Clustering

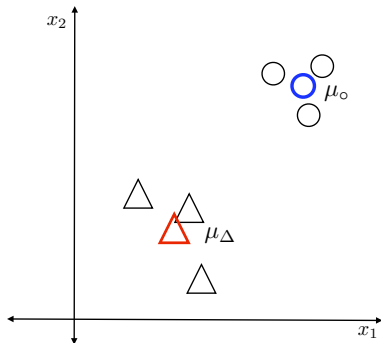
Felix Bießmann

Beuth University & Einstein Center for Digital Future



Clustering

Psychological Models of Categorization: Prototypes



Prototypes μ_{Δ} and μ_o :

$$\mu_{\Delta} = 1/N_{\Delta} \sum_n^{N_{\Delta}} \mathbf{x}_{\Delta,n}$$

$$\mu_o = 1/N_o \sum_n^{N_o} \mathbf{x}_{o,n}$$

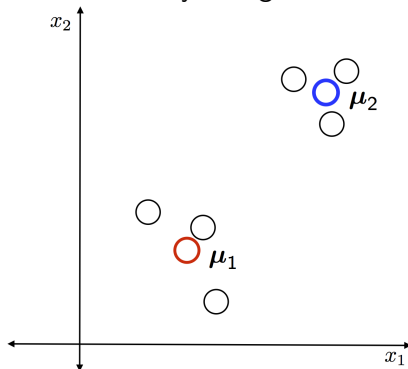
New data points \mathbf{x} are assigned to their closest cluster center μ^*

$$\mu^* = \underset{i}{\operatorname{argmin}} (\|\mu_i - \mathbf{x}\|_2) \quad (1)$$



Clustering

Psychological Models of Categorization: Prototypes



The only difference for clustering is:
We do not have labels.



K-means Clustering

Goal: Given data $\mathbf{x}_1, \dots, \mathbf{x}_N$ find cluster centers μ_1, \dots, μ_K such that the distances of data points to their respective cluster centre are minimized

$$\mathcal{J} = \sum_{n=1}^N \sum_{k=1}^K \mathbf{c}_{n,k} \|\mathbf{x}_n - \mu_{\mathbf{c}_k}\| \quad (2)$$

$$\text{where } \mathbf{c}_{n,k} \begin{cases} 1 & \text{if } \mathbf{x}_n \text{ belongs to cluster } k \\ 0 & \text{otherwise} \end{cases} \quad (3)$$



K-means Clustering

K-Means Algorithm

Re-iterating two steps:

1. Assign each data point \mathbf{x}_i to their closest cluster μ_k
2. Update μ_k to the mean of the members in that cluster



K-means Clustering Algorithm

Algorithm 1 K-means clustering

Require: data $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^D$, number of clusters k , iterations m .

- 1: Choose random data points as initial cluster centres $\boldsymbol{\mu}_1 \leftarrow \mathbf{x}_{i_1}, \dots, \boldsymbol{\mu}_k \leftarrow \mathbf{x}_{i_k}$ where $i_j \neq i_l$ for all $j \neq l$.
- 2: $\mathbf{c} \leftarrow \mathbf{0}_N$
- 3: $\mathbf{c}^{\text{old}} \leftarrow \mathbf{0}_N$
- 4: $i \leftarrow 0$
- 5: **while** $i < m$ **do**
- 6: **for** $j = 1$ to N **do**
- 7: Find nearest cluster centre $\mathbf{c}_j \leftarrow \operatorname{argmin}_{1 \leq l \leq k} \|\mathbf{x}_j - \boldsymbol{\mu}_l\|_2$
- 8: **end for**
- 9: **for** $j \leftarrow 1$ to k **do**
- 10: Compute new cluster centre $\boldsymbol{\mu}_j \leftarrow \frac{1}{|\{l: \mathbf{c}_l = j\}|} \sum_{l: \mathbf{c}_l = j} \mathbf{x}_l$
- 11: **end for**
- 12: **if** $\mathbf{c}^{\text{old}} = \mathbf{c}$ **then**
- 13: **break**
- 14: **end if**
- 15: $\mathbf{c}^{\text{old}} \leftarrow \mathbf{c}$
- 16: $i \leftarrow i + 1$
- 17: **end while**
- 18: **return** cluster centres $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k \in \mathbb{R}^D$, assignment vector $\mathbf{c}^{\text{old}} \in \mathbb{R}^n$



Application Example: Geyser Eruptions



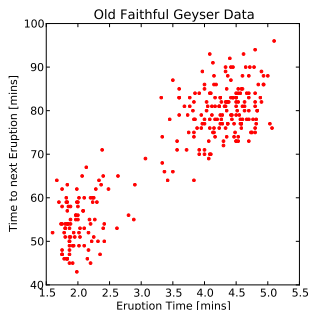
Old Faithful Geyser
Yellowstone National Park,
USA

Famous data set for clustering

- Old Faithful Eruptions
- Two dimensions
 1. Time of Eruption [mins]
 2. Time until next Eruption [mins]



Application Example: Geyser Eruptions

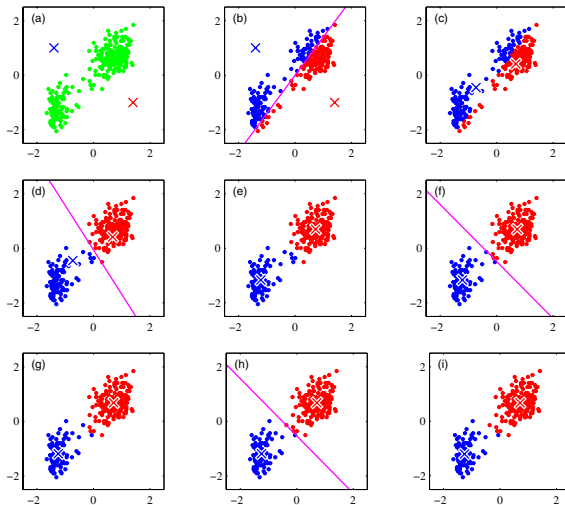


Famous data set for clustering

- Old Faithful Eruptions
- Two dimensions
 1. Time of Eruption [mins]
 2. Time until next Eruption [mins]

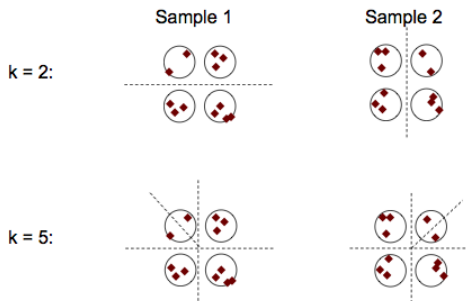


K-means Clustering Step-by-Step



Clustering Instability

Number of Clusters is a critical parameter



Clusterings are instable if number of clusters is too small or too large



Clustering Instability

- Number of clusters is critical hyper parameter
 - In supervised settings we use cross-validation to optimize hyper-parameters for accuracy on test data
 - How can we optimize the number of clusters?
- Choose that k that results in most **stable** clusterings
For a review see e.g. [von Luxburg, 2009]



Clustering Instability Algorithm

Algorithm 2 Clustering Instability

Require: data points $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$, clustering algorithm \mathcal{A} , maximal number of clusters K , iterations i .

Ensure: optimal number of clusters k^*

- 1: **for** $k = 2$ to K **do**
 - 2: Resample data set (e.g. random draws with replacement)
 - 3: **for** $it = 1$ to i **do**
 - 4: Cluster data using algorithm \mathcal{A} into k clusters
 - 5: **end for**
 - 6: Compute minimal (across all label permutations) distance between clusterings
 - 7: **end for**
 - 8: Chose that k that has the minimal instability over resamplings
-

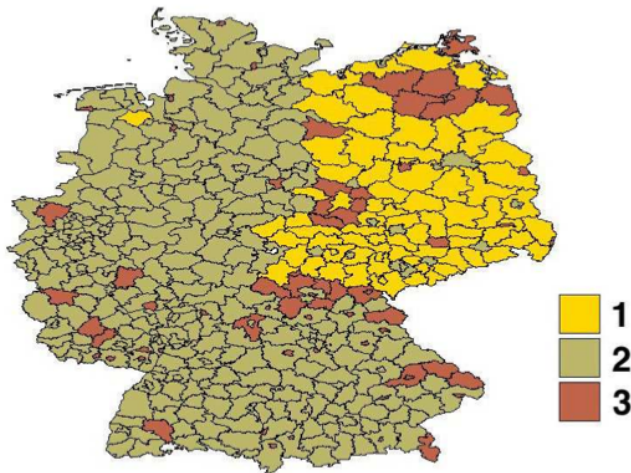


Application: Clustering 'Mentality' of German Population

- Until 1989 Germany was divided into
 - a free capitalistic western part
 - a communist eastern part
- Political systems influence mentality of people
- The survey "Perspektive Deutschland" investigated this
Survey asked questions about:
 - what people desire
 - what people are afraid of



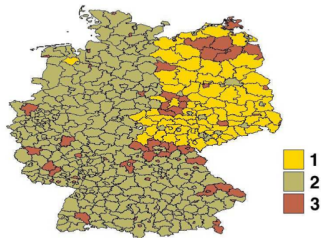
German Regions clustered by 'Mentality' of Population



...based on questionnaires, 'Perspektive Deutschland' poll 2005.



German Regions clustered by 'Mentality' of Population



People in cluster 1 say:

Helping others is important

I am afraid of loosing my job

People in cluster 2 say:

Reaching my own goals is important

I am afraid of loosing my health



Online K-Means

$$\begin{aligned}\mathcal{E}(\mu_k) &= \frac{1}{2}(\mathbf{x} - \mu_k)^2 \\ &= \frac{1}{2}\mathbf{x}^\top \mathbf{x} - \frac{1}{2}2\mathbf{x}^\top \mu_k - \frac{1}{2}\mu_k^\top \mu_k\end{aligned}\tag{4}$$

$$\frac{\partial \mathcal{E}(\mu_k)}{\partial \mu_k} = \mu_k - \mathbf{x}$$

Learning rate $\eta = \frac{1}{n_k}$ (n_k = number of samples in cluster k)

$$\text{Gradient Step } \mu_k \leftarrow \mu_k - \eta \frac{\partial \mathcal{E}(\mu_k)}{\partial \mu_k} = \mu_k - \frac{1}{n_k}(\mathbf{x} - \mu_k)$$



Online K-Means Algorithm

Algorithm 3 Online K-means clustering

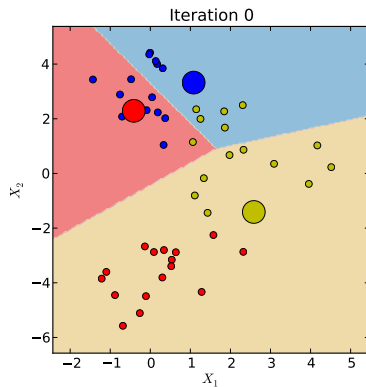
Require: data points $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$, number of clusters k , iterations m .

Ensure: cluster centres $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k \in \mathbb{R}^d$

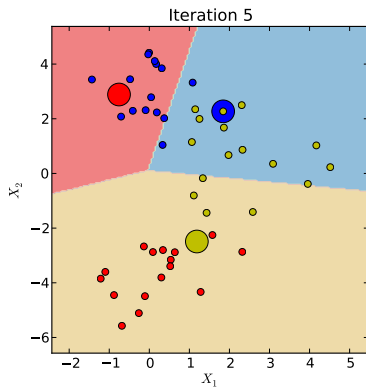
- 1: Choose random data points as initial cluster centres
 $\boldsymbol{\mu}_k \leftarrow \mathbf{x}_{i_j}, \dots, \boldsymbol{\mu}_k \leftarrow \mathbf{x}_{i_k}$ where $i_j \neq i_l$ for all $j \neq l$.
 - 2: Initialize cluster assignment counts $n_1, \dots, n_k \leftarrow 0$
 - 3: **for** $i = 1, \dots, m$ **do**
 - 4: Draw a new data point randomly \mathbf{x}_i
 - 5: Find nearest cluster centre $k^* \leftarrow \operatorname{argmin}_k \|\mathbf{x}_i - \boldsymbol{\mu}_k\|_2$
 - 6: Update cluster counts $n_{k^*} \leftarrow n_{k^*} + 1$
 - 7: Update cluster centers $\boldsymbol{\mu}_{k^*} \leftarrow \boldsymbol{\mu}_{k^*} + \frac{1}{n_{k^*}}(\mathbf{x}_i - \boldsymbol{\mu}_{k^*})$
 - 8: **end for**
-



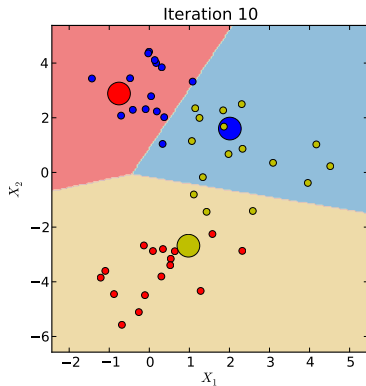
Online K-Means Algorithm - Example



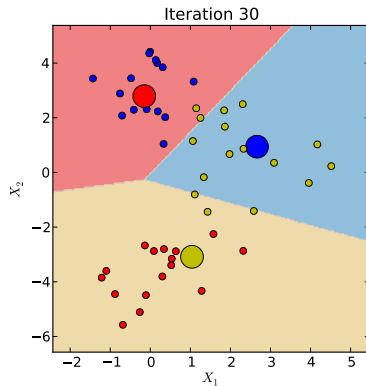
Online K-Means Algorithm - Example



Online K-Means Algorithm - Example



Online K-Means Algorithm - Example



Distance Measures for Real-Valued Data $\mathbf{x} \in \mathbb{R}^D$

Clustering Algorithms need a **distance function** $d(\mathbf{x}_i, \mathbf{x}_j)$

- For real valued data $\mathbf{x} \in \mathbb{R}^D$ we can use the Euclidean distance

$$d(\mathbf{x}_i, \mathbf{x}_j) = \|\mathbf{x}_i - \mathbf{x}_j\|_2^2 \quad (5)$$

- More robust (less sensitive to outliers) is the **city block distance** or \mathcal{L}_1 norm

$$d(\mathbf{x}_i, \mathbf{x}_j) = \|\mathbf{x}_i - \mathbf{x}_j\|_1 \quad (6)$$

- Another alternative is the correlation coefficient (also called cosine similarity)

$$d(\mathbf{x}_i, \mathbf{x}_j) = \frac{\mathbf{x}_1^\top \mathbf{x}_2}{\|\mathbf{x}_1\|_2 \|\mathbf{x}_2\|_2} \quad (7)$$

For standardized data $\sum_i \mathbf{x}_i = 0$, $\sum_i \mathbf{x}_i^2 = 1$ maximizing correlation is the same as minimizing euclidean distance.



Distance Measures for Non-Real-Valued Data

- For ordinal variables $\mathbf{x} \in \{\text{low, medium, high}\}^D$ we can transform the values into real-valued numbers (for three possible values e.g. $1/3, 2/3, 3/3$) and then apply distance functions for real-valued data
- For categorical variables $\mathbf{x} \in \{\text{red, green, blue}\}^D$ we can use a binary coding for the differences

$$d(\mathbf{x}_i, \mathbf{x}_j) = \sum_d^D \mathbf{x}_{id} \neq \mathbf{x}_{jd} \quad (8)$$

This metric is called **Hamming Distance**

For the sake of simplicity we only consider the euclidean distance



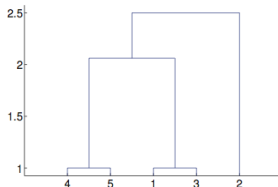
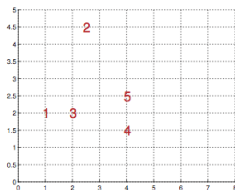
Hierarchical Clustering

- K-Means produces **flat** clusterings
- Often we are interested in a **hierarchy** of clusterings
- Examples:
 - Biological Species
 - Topics in Text Documents



Hierarchical Clustering

- K-Means produces **flat** clusterings
- Often we are interested in a **hierarchy** of clusterings
- Examples:
 - Biological Species
 - Topics in Text Documents



Hierarchical Clustering

- A popular approach to hierarchical clustering is
 1. Start with each data point as one cluster
 2. Successively merge (*agglomerate*) similar clusters

→ **Agglomerative Clustering**

As most clustering algorithms, these procedures are not defined via objective functions but via algorithms

→ Difficult to establish convergence criteria



Agglomerative Clustering

Merging requires distance function $d(C_i, C_j)$ for clusters C_i, C_j

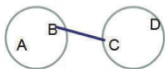


Agglomerative Clustering

Merging requires distance function $d(C_i, C_j)$ for clusters C_i, C_j

Single Linkage

Distance between two
closest points in C_i, C_j



$$d(C_i, C_j) = \min_{\mathbf{x}_i \in C_i, \mathbf{x}_j \in C_j} d(\mathbf{x}_i, \mathbf{x}_j)$$

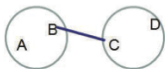


Agglomerative Clustering

Merging requires distance function $d(C_i, C_j)$ for clusters C_i, C_j

Single Linkage

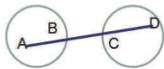
Distance between two
closest points in C_i, C_j



$$d(C_i, C_j) = \min_{\mathbf{x}_i \in C_i, \mathbf{x}_j \in C_j} d(\mathbf{x}_i, \mathbf{x}_j)$$

Complete Linkage

Distance between two
most distant points



$$d(C_i, C_j) = \max_{\mathbf{x}_i \in C_i, \mathbf{x}_j \in C_j} d(\mathbf{x}_i, \mathbf{x}_j)$$

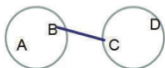


Agglomerative Clustering

Merging requires distance function $d(C_i, C_j)$ for clusters C_i, C_j

Single Linkage

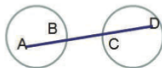
Distance between two
closest points in C_i, C_j



$$d(C_i, C_j) = \min_{\mathbf{x}_i \in C_i, \mathbf{x}_j \in C_j} d(\mathbf{x}_i, \mathbf{x}_j)$$

Complete Linkage

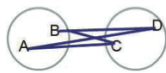
Distance between two
most distant points



$$d(C_i, C_j) = \max_{\mathbf{x}_i \in C_i, \mathbf{x}_j \in C_j} d(\mathbf{x}_i, \mathbf{x}_j)$$

Average Linkage

Average distance between
all $N_i N_j$ pairs



$$d(C_i, C_j) = \frac{1}{N_i N_j} \sum_{\mathbf{x}_i \in C_i, \mathbf{x}_j \in C_j} d(\mathbf{x}_i, \mathbf{x}_j)$$



Agglomerative Clustering Algorithm

Algorithm 4 Agglomerative Clustering

Require: Data points $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^D$, number of clusters k , distance function $d(.,.)$

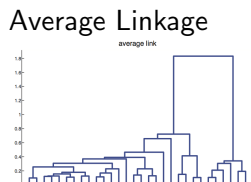
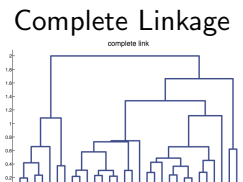
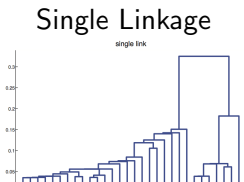
Ensure: Binary tree of clusters

- 1: Initialize each data point as cluster
 - 2: **for** $i = 1$ to N **do**
 - 3: $C_i \leftarrow i$
 - 4: **end for**
 - 5: Initialize each cluster as available for merging
 - 6: $\mathcal{S} \leftarrow \{1, \dots, N\}$
 - 7: **while** There are clusters to merge **do**
 - 8: Pick 2 most similar clusters to merge
 - 9: $j, k \leftarrow \operatorname{argmin}_{j,k} d(j, k)$
 - 10: Merge clusters to new cluster $C_l \leftarrow C_j \cup C_k$
 - 11: Mark j, k as unavailable for merging
 - 12: $\mathcal{S} \leftarrow \mathcal{S} \setminus \{j, k\}$
 - 13: **if** $C_l \notin \mathcal{S}$ **then**
 - 14: $\mathcal{S} \leftarrow \mathcal{S} \cup \{l\}$
 - 15: **end if**
 - 16: **end while**
-



Examples Hierarchical Clustering

Dendrograms (binary clustering trees) of yeast gene expression data



Taken from [Murphy, 2012]



Summary

- Clustering Algorithms find clusters in data
- Clustering Performance depends on distance function used
- K-Means is one of the most popular clustering algorithms
- For large data sets use Online K-Means
- Wrong number of clusters leads to unstable results
- Hierarchical clustering



References

- C. M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer US, 2007.
- S. Lloyd. Least squares quantization in pcm. *Information Theory, IEEE Transactions on*, 28(2):129–137, Mar. 1982. ISSN 0018-9448.
- K. P. Murphy. *Machine Learning: A Probabilistic Perspective*. Adaptive Computation and Machine Learning. The MIT Press, 1 edition, 2012. ISBN 0262018020,9780262018029.
- U. von Luxburg. Clustering stability: An overview. *Foundations and Trends in Machine Learning*, 2(3):235–274, 2009.

