

Title: groh: profile1 (27.05.2015)

Date: Wed May 27 08:22:56 CEST 2015

Duration: 82:49 min

Pages: 72

PD Dr. Georg Groh



Finding Clusters in Profiles

Examples for profile elements that can be embedded in metric spaces:

- **Location & Velocity:** Metric space: $(\mathbb{R}^3, \|\cdot\|)$
- **Text** describing Interests: Metric space: $(\mathbb{R}^{|\text{Voc}|}, \|\cdot\|)$ where Voc denotes the Vocabulary of the text.

"I like to dance samba, bake pizza, watch tv and plant trees in the garden. I also like to bake cakes."



I	2
like	2
to	2
dance	1
samba	1
bake	2
pizza	1
watch	1
tv	1
and	1
plant	1
trees	1
in	1
the	1
garden	1
also	1
cakes	1

Often: Instead of term-frequency (tf) alone: use **term-frequency * inverse document frequency (idf)**; $idf = \log(\#of docs / \#of docs)$

Finding Clusters in Profiles

- How do we compute **clusters in metric spaces?**
- **Group models:** How do we compute **socially meaningful clusters** in metric spaces (and thus avoid quasi-groups)?
- First some **notations / basics:**
 - In graph clustering we had: A graph clustering $\mathbf{C}=\{C_1, C_2, \dots, C_K\}$ is a partition of V into non-empty subsets C_k
 - Now: **clustering** $\mathcal{C} : \mathcal{X} \rightarrow \mathcal{I}$: mapping of a metric value space \mathcal{X} to a set of cluster indices \mathcal{I}
 - Clusterings can be:
 - **exclusive or non-exclusive**
 - **crisp or fuzzy**
 - **hierarchical or non-hierarchical**

- How do we compute **clusters in metric spaces**?
- **Group models**: How do we compute **socially meaningful clusters** in metric spaces (and thus avoid quasi-groups)?
- First some **notations / basics**:
 - In graph clustering we had: A graph clustering $\mathbf{C}=\{C_1, C_2, \dots, C_K\}$ is a partition of V into non-empty subsets C_k
 - Now: **clustering** $\mathcal{C} : \mathcal{X} \rightarrow \mathcal{I}$: mapping of a metric value space X to a set of cluster indices I
 - Clusterings can be:
 - **exclusive or non-exclusive**
 - **crisp or fuzzy**
 - **hierarchical or non-hierarchical**



- How do we compute **clusters in metric spaces**?
- **Group models**: How do we compute **socially meaningful clusters** in metric spaces (and thus avoid quasi-groups)?
- First some **notations / basics**:
 - In graph clustering we had: A graph clustering $\mathbf{C}=\{C_1, C_2, \dots, C_K\}$ is a partition of V into non-empty subsets C_k
 - Now: **clustering** $\mathcal{C} : \mathcal{X} \rightarrow \mathcal{I}$: mapping of a metric value space X to a set of cluster indices I
 - Clusterings can be:
 - **exclusive or non-exclusive**
 - **crisp or fuzzy**
 - **hierarchical or non-hierarchical**



- **Exclusive** \rightarrow **non overlapping clusters**; **non-exclusive** \rightarrow overlapping clusters
- **Hierarchical** clustering \rightarrow imposes a tree structure (Dendrogram) on the C_k where an edge $C_i \rightarrow C'_j$ implies $C_i \subset C'_j$;
- **Crisp** clusterings: Conventional characteristic functions α_k for each Cluster C_k

$$\alpha_k : \mathcal{X} \rightarrow \{0, 1\} \text{ with } \alpha_k(x \in \mathcal{X}) = \begin{cases} 1 & x \in C_k \\ 0 & x \notin C_k \end{cases}$$

- **Fuzzy** clusterings: fuzzy membership function α_k for each Cluster C_k

$$\alpha_k : \mathcal{X} \rightarrow [0, 1]$$



- **Exclusive** \rightarrow **non overlapping clusters**; **non-exclusive** \rightarrow overlapping clusters
- **Hierarchical** clustering \rightarrow imposes a tree structure (Dendrogram) on the C_k where an edge $C_i \rightarrow C'_j$ implies $C_i \subset C'_j$;
- **Crisp** clusterings: Conventional characteristic functions α_k for each Cluster C_k

$$\alpha_k : \mathcal{X} \rightarrow \{0, 1\} \text{ with } \alpha_k(x \in \mathcal{X}) = \begin{cases} 1 & x \in C_k \\ 0 & x \notin C_k \end{cases}$$

- **Fuzzy** clusterings: fuzzy membership function α_k for each Cluster C_k

$$\alpha_k : \mathcal{X} \rightarrow [0, 1]$$



- **Exclusive** → non overlapping clusters; **non-exclusive** → overlapping clusters
- **Hierarchical** clustering → imposes a tree structure (Dendrogram) on the C_k where an edge $C_i \rightarrow C_j$ implies $C_i \subset C_j$;
- **Crisp** clusterings: Conventional characteristic functions α_k for each Cluster C_k

$$\alpha_k : \mathcal{X} \rightarrow \{0, 1\} \text{ with } \alpha_k(x \in \mathcal{X}) = \begin{cases} 1 & x \in C_k \\ 0 & x \notin C_k \end{cases}$$

- **Fuzzy** clusterings: fuzzy membership function α_k for each Cluster C_k

$$\alpha_k : \mathcal{X} \rightarrow [0, 1]$$



- Metric variant of **Single / Complete link clustering**: Hierarchical, crisp, non-overlapping
- **Completely analogous to graph clustering case**: Start with singletons and on each level of the dendrogram merge two clusters with minimal distance (cost)

- Single link:

$$d(C_{k_1}, C_{k_2}) = \min_{\{n_1, n_2 | x_{n_1} \in C_{k_1} \wedge x_{n_2} \in C_{k_2}\}} \|x_{n_1} - x_{n_2}\|$$

- Complete link:

$$d(C_{k_1}, C_{k_2}) = \max_{\{n_1, n_2 | x_{n_1} \in C_{k_1} \wedge x_{n_2} \in C_{k_2}\}} \|x_{n_1} - x_{n_2}\|$$



- Metric variant of **Single / Complete link clustering**: Hierarchical, crisp, non-overlapping
- **Completely analogous to graph clustering case**: Start with singletons and on each level of the dendrogram merge two clusters with minimal distance (cost)

- Single link:

$$d(C_{k_1}, C_{k_2}) = \min_{\{n_1, n_2 | x_{n_1} \in C_{k_1} \wedge x_{n_2} \in C_{k_2}\}} \|x_{n_1} - x_{n_2}\|$$

- Complete link:

$$d(C_{k_1}, C_{k_2}) = \max_{\{n_1, n_2 | x_{n_1} \in C_{k_1} \wedge x_{n_2} \in C_{k_2}\}} \|x_{n_1} - x_{n_2}\|$$



- Metric variant of **Single / Complete link clustering**: Hierarchical, crisp, non-overlapping
- **Completely analogous to graph clustering case**: Start with singletons and on each level of the dendrogram merge two clusters with minimal distance (cost)

- Single link:

$$d(C_{k_1}, C_{k_2}) = \min_{\{n_1, n_2 | x_{n_1} \in C_{k_1} \wedge x_{n_2} \in C_{k_2}\}} \|x_{n_1} - x_{n_2}\|$$

- Complete link:

$$d(C_{k_1}, C_{k_2}) = \max_{\{n_1, n_2 | x_{n_1} \in C_{k_1} \wedge x_{n_2} \in C_{k_2}\}} \|x_{n_1} - x_{n_2}\|$$



• General idea (also valid in graph clustering): **Optimize objective function** that formalizes clustering paradigm.

• K-Means: **Optimize intra cluster coherence:**

• Describe cluster C_k by **prototype** μ_k ; prototype need not be an actual pattern (If so, algorithm works with slight modifications as well)

• Determine cluster for each pattern x_n by **nearest neighbour rule:**

$$\mathcal{C}(x_n) = k_a \leftrightarrow \|x_n - \mu_{k_a}\| = \min_i \|x_n - \mu_k\|$$



• K-Means: **Optimize intra cluster coherence:**

• **Find prototypes** by optimizing objective function modeling intra cluster coherence as mean square error

$$J_{SQE} = \sum_{k=1}^K \sum_{\{n|x_n \in C_k\}} \|x_n - \mu_k\|^2$$

$$\frac{dJ_{SQE}}{d\mu_k} \stackrel{!}{=} 0 \implies \mu^k = \frac{1}{|C_k|} \sum_{\{n|x_n \in C_k\}} x_n$$

• \rightarrow cluster prototypes are barycenters („centers of gravity“) of their clusters.



• General idea (also valid in graph clustering): **Optimize objective function** that formalizes clustering paradigm.

• K-Means: **Optimize intra cluster coherence:**

• Describe cluster C_k by **prototype** μ_k ; prototype need not be an actual pattern (If so, algorithm works with slight modifications as well)

• Determine cluster for each pattern x_n by **nearest neighbour rule:**

$$\mathcal{C}(x_n) = k_a \leftrightarrow \|x_n - \mu_{k_a}\| = \min_i \|x_n - \mu_k\|$$



• K-Means: **Optimize intra cluster coherence:**

• **Find prototypes** by optimizing objective function modeling intra cluster coherence as mean square error

$$J_{SQE} = \sum_{k=1}^K \sum_{\{n|x_n \in C_k\}} \|x_n - \mu_k\|^2$$

$$\frac{dJ_{SQE}}{d\mu_k} \stackrel{!}{=} 0 \implies \mu^k = \frac{1}{|C_k|} \sum_{\{n|x_n \in C_k\}} x_n$$

• \rightarrow cluster prototypes are barycenters („centers of gravity“) of their clusters.



K-Means: Optimize intra cluster coherence:

- Find prototypes by optimizing objective function modeling intra cluster coherence as mean square error

$$J_{SQE} = \sum_{k=1}^K \sum_{\{n|x_n \in C_k\}} \|x_n - \mu_k\|^2$$

$$\frac{dJ_{SQE}}{d\mu_k} \stackrel{!}{=} 0 \implies \mu^k = \frac{1}{|C_k|} \sum_{\{n|x_n \in C_k\}} x_n$$

- cluster prototypes are barycenters („centers of gravity“) of their clusters.



K-Means: Optimize intra cluster coherence:

- Find prototypes by optimizing objective function modeling intra cluster coherence as mean square error

$$J_{SQE} = \sum_{k=1}^K \sum_{\{n|x_n \in C_k\}} \|x_n - \mu_k\|^2$$

$$\frac{dJ_{SQE}}{d\mu_k} \stackrel{!}{=} 0 \implies \mu^k = \frac{1}{|C_k|} \sum_{\{n|x_n \in C_k\}} x_n$$

- cluster prototypes are barycenters („centers of gravity“) of their clusters.



K-Means: Optimize intra cluster coherence:

- Find prototypes by optimizing objective function modeling intra cluster coherence as mean square error

$$J_{SQE} = \sum_{k=1}^K \sum_{\{n|x_n \in C_k\}} \|x_n - \mu_k\|^2$$

$$\frac{dJ_{SQE}}{d\mu_k} \stackrel{!}{=} 0 \implies \mu^k = \frac{1}{|C_k|} \sum_{\{n|x_n \in C_k\}} x_n$$

- cluster prototypes are barycenters („centers of gravity“) of their clusters.



K-Means: Optimize intra cluster coherence:

- Find prototypes by optimizing objective function modeling intra cluster coherence as mean square error

$$J_{SQE} = \sum_{k=1}^K \sum_{\{n|x_n \in C_k\}} \|x_n - \mu_k\|^2$$

$$\frac{dJ_{SQE}}{d\mu_k} \stackrel{!}{=} 0 \implies \mu^k = \frac{1}{|C_k|} \sum_{\{n|x_n \in C_k\}} x_n$$

- cluster prototypes are barycenters („centers of gravity“) of their clusters.



• K-Means: Optimize intra cluster coherence:

- Find prototypes by optimizing objective function modeling intra cluster coherence as mean square error

$$J_{SQE} = \sum_{k=1}^K \sum_{\{n|x_n \in C_k\}} \|x_n - \mu_k\|^2$$

$$\frac{dJ_{SQE}}{d\mu_k^k} \stackrel{!}{=} 0 \implies \mu^k = \frac{1}{|C_k|} \sum_{\{n|x_n \in C_k\}} x_n$$

- → cluster prototypes are barycenters („centers of gravity“) of their clusters.



• K-Means: Optimize intra cluster coherence:

- Find prototypes by optimizing objective function modeling intra cluster coherence as mean square error

$$J_{SQE} = \sum_{k=1}^K \sum_{\{n|x_n \in C_k\}} \|x_n - \mu_k\|^2$$

$$\frac{dJ_{SQE}}{d\mu_k^k} \stackrel{!}{=} 0 \implies \mu^k = \frac{1}{|C_k|} \sum_{\{n|x_n \in C_k\}} x_n$$

- → cluster prototypes are barycenters („centers of gravity“) of their clusters.



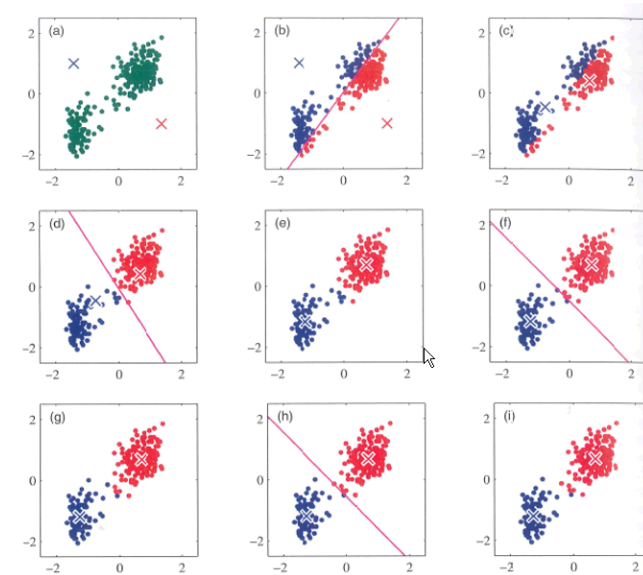
• K-Means: Optimize intra cluster coherence:

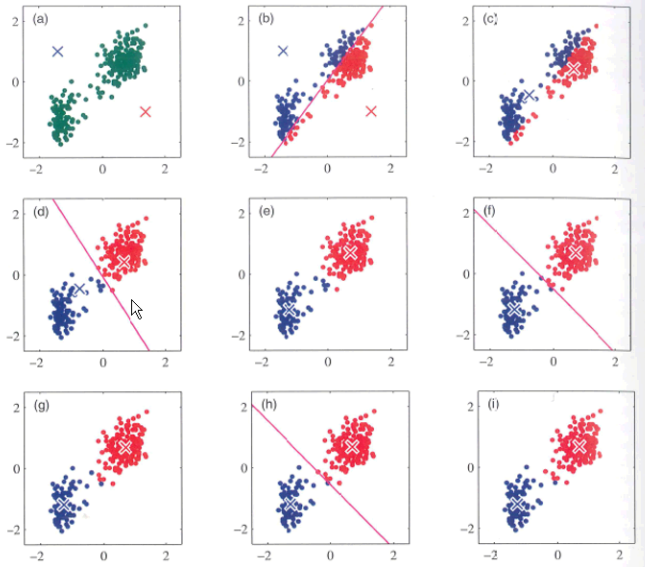
- Find prototypes by optimizing objective function modeling intra cluster coherence as mean square error

$$J_{SQE} = \sum_{k=1}^K \sum_{\{n|x_n \in C_k\}} \|x_n - \mu_k\|^2$$

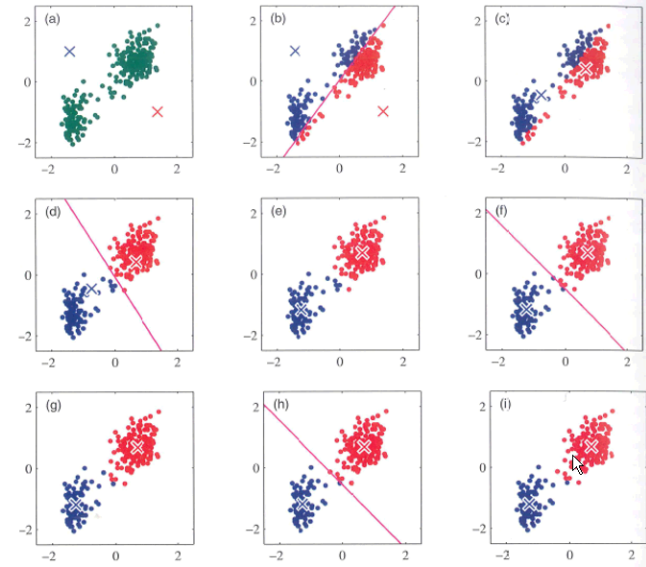
$$\frac{dJ_{SQE}}{d\mu_k^k} \stackrel{!}{=} 0 \implies \mu^k = \frac{1}{|C_k|} \sum_{\{n|x_n \in C_k\}} x_n$$

- → cluster prototypes are barycenters („centers of gravity“) of their clusters.





[3]



[3]

- Dunn Index:

$$D = \min_{k_1 \in [1, K]} \left(\min_{k_2 \in [1, K]} \left(\frac{d_1(\mathcal{C}_{k_1}, \mathcal{C}_{k_2})}{\max_{k_3 \in [1, K]} d_2(\mathcal{C}_{k_3})} \right) \right)$$

where $d_1(\mathcal{C}_{k_1}, \mathcal{C}_{k_2})$ is the distance function between two clusters defined by

$$d_1(\mathcal{C}_{k_1}, \mathcal{C}_{k_2}) = \min_{\{(n_1, n_2) | x_{n_1} \in \mathcal{C}_{k_1} \wedge x_{n_2} \in \mathcal{C}_{k_2}\}} \|x_{n_1} - x_{n_2}\|$$

(that is the single link distance from SAHN).

The “diameter” d_2 of the clusters is defined by

$$d_2(\mathcal{C}_i) = \max_{\{(n_1, n_2) | x_{n_1} \in \mathcal{C}_i \wedge x_{n_2} \in \mathcal{C}_i\}} \|x_{n_1} - x_{n_2}\|$$

- Dunn Index:

$$D = \min_{k_1 \in [1, K]} \left(\min_{k_2 \in [1, K]} \left(\frac{d_1(\mathcal{C}_{k_1}, \mathcal{C}_{k_2})}{\max_{k_3 \in [1, K]} d_2(\mathcal{C}_{k_3})} \right) \right)$$

where $d_1(\mathcal{C}_{k_1}, \mathcal{C}_{k_2})$ is the distance function between two clusters defined by

$$d_1(\mathcal{C}_{k_1}, \mathcal{C}_{k_2}) = \min_{\{(n_1, n_2) | x_{n_1} \in \mathcal{C}_{k_1} \wedge x_{n_2} \in \mathcal{C}_{k_2}\}} \|x_{n_1} - x_{n_2}\|$$

(that is the single link distance from SAHN).

The “diameter” d_2 of the clusters is defined by

$$d_2(\mathcal{C}_i) = \max_{\{(n_1, n_2) | x_{n_1} \in \mathcal{C}_i \wedge x_{n_2} \in \mathcal{C}_i\}} \|x_{n_1} - x_{n_2}\|$$

- Dunn Index:

$$D = \min_{k_1 \in [1, K]} \left(\min_{k_2 \in [1, K]} \left(\frac{d_1(\mathcal{C}_{k_1}, \mathcal{C}_{k_2})}{\max_{k_3 \in [1, K]} d_2(\mathcal{C}_{k_3})} \right) \right)$$

where $d_1(\mathcal{C}_{k_1}, \mathcal{C}_{k_2})$ is the distance function between two clusters defined by

$$d_1(\mathcal{C}_{k_1}, \mathcal{C}_{k_2}) = \min_{\{(n_1, n_2) | x_{n_1} \in \mathcal{C}_{k_1} \wedge x_{n_2} \in \mathcal{C}_{k_2}\}} \|x_{n_1} - x_{n_2}\|$$

(that is the single link distance from SAHN).

The “diameter” d_2 of the clusters is defined by

$$d_2(\mathcal{C}_i) = \max_{\{(n_1, n_2) | x_{n_1} \in \mathcal{C}_i \wedge x_{n_2} \in \mathcal{C}_i\}} \|x_{n_1} - x_{n_2}\|$$

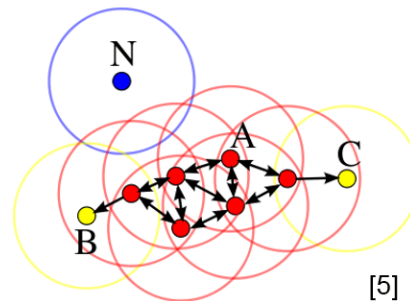
Example Application: Clustering locations

- Problem: How do we distinguish socially relevant clusters (candidates for groups) from quasi groups?
 - Compute clusterings over period of time: Good candidates: clusters that appear over and over again, clusters that appear periodically
 - Establish threshold for distance in clusters: Human “social distance”: A few meters (if groups are very small); few tens of meters (if groups are medium sized)
 - Include velocities: If divergent → no group

- K-Means is „OK“ as cluster algorithm, but has certain disadvantages:
 - favors spherical clusters
 - need to know K
 - no notion of noise

- Alternative → DBSCAN [4] (de facto state of the art):

- Idea: Two parameters: minPt, ϵ
- Rough idea: iterate: visit previously unseen pattern x:
 - if in ϵ -neighborhood $\{x'\}$ of x: $|\{x'\}| \geq \text{minPt}$ then start new cluster: include x and $\{x'\}$ and those of their ϵ -neighborhoods $\{x''\}$ that are dense enough ($|\{x''\}| \geq \text{minPt}$), etc.
 - else: x is noise

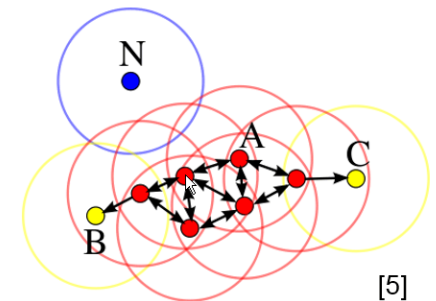


[5]

- K-Means is „OK“ as cluster algorithm, but has certain disadvantages:
 - favors spherical clusters
 - need to know K
 - no notion of noise

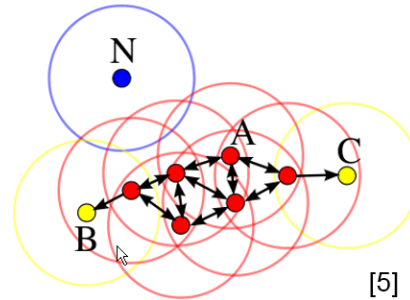
- Alternative → DBSCAN [4] (de facto state of the art):

- Idea: Two parameters: minPt, ϵ
- Rough idea: iterate: visit previously unseen pattern x:
 - if in ϵ -neighborhood $\{x'\}$ of x: $|\{x'\}| \geq \text{minPt}$ then start new cluster: include x and $\{x'\}$ and those of their ϵ -neighborhoods $\{x''\}$ that are dense enough ($|\{x''\}| \geq \text{minPt}$), etc.
 - else: x is noise



[5]

- K-Means is „OK“ as cluster algorithm, but has certain **disadvantages**:
 - favors **spherical clusters**
 - **need to know K**
 - **no notion of noise**



- **Alternative → DBSCAN [4]**
(de facto state of the art):
 - Idea: Two parameters: **minPt, ε**
 - Rough **idea**: **iterate**:
visit previously unseen pattern x:
if in ε-neighborhood {x'} of x: $|\{x'\}| \geq \text{minPt}$ then
start new cluster: include x and {x'} and those of their
ε-neighborhoods {x''} that are dense enough ($|\{x''}\}| \geq$
minPt), etc.
else: x is noise



- **Advantages** of DBSCAN:
 - We do not need to know K in advance
 - arbitrarily shaped clusters
 - notion of noise
- **Disadvantages**:
 - instead of having to know K, we need to „guess“ minPt and ε
instead (can be a problem for high dimensional pattern spaces (→
curse of dimensionality))
 - original DBSCAN has fixed (minPt, ε) → problems when cluster
density varies



K-Means Clustering

- Interesting aspect: How do we **determine correct number k of clusters**?
(Same problem with graph clustering: where to cut dendrogram?)
- Answer: Compute for every k clusterings; **choose the best clustering with
a cluster quality measure**
- **Cluster quality measures** for metric case: (countless variants exist in
literature; for an overview: e.g. [2]) (**Objective functions** modeling
clustering paradigm):
 - Dunn-Index
 - Entropy based indices
 -



Fuzzy C-Means Clustering

- K-Means was a crisp algorithm. Now: **fuzzy variant**
- Reformulate K-Means objective function with **membership matrix**
 r_{nk} : Membership of pattern x_n in class C_k

$$J_{SQE} = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$$

- optimization criterion
- together with non-overlapping constraint

$$dJ_{SQE}/d\mu_k = 0$$

$$\forall n (\exists k (r_{nk} = 1) \wedge ((k' \neq k) \rightarrow (r_{nk'} = 0)))$$

leads to well known K-Means

$$\mu_k = \sum_{n=1}^N r_{nk} x_n / \sum_{n=1}^N r_{nk} = (1/|C_k|) \sum_{n|x_n \in C_k} x_n$$



Fuzzy C-Means Clustering

- K-Means was a crisp algorithm. Now: **fuzzy variant**
- Reformulate K-Means objective function with **membership matrix**
 r_{nk} : Membership of pattern x_n in class C_k

$$J_{SQE} = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$$

- optimization criterion

$$dJ_{SQE}/d\mu_k = 0$$

- together with non-overlapping constraint

$$\forall n(\exists k(r_{nk} = 1) \wedge ((k' \neq k) \rightarrow (r_{nk'} = 0)))$$

leads to well known K-Means

$$\mu_k = \sum_{n=1}^N r_{nk} x_n / \sum_{n=1}^N r_{nk} = (1/|C_k|) \sum_{n|x_n \in C_k} x_n$$



Fuzzy C-Means Clustering

- K-Means was a crisp algorithm. Now: **fuzzy variant**
- Reformulate K-Means objective function with **membership matrix**
 r_{nk} : Membership of pattern x_n in class C_k

$$J_{SQE} = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$$

- optimization criterion

$$dJ_{SQE}/d\mu_k = 0$$

- together with non-overlapping constraint

$$\forall n(\exists k(r_{nk} = 1) \wedge ((k' \neq k) \rightarrow (r_{nk'} = 0)))$$

leads to well known K-Means

$$\mu_k = \sum_{n=1}^N r_{nk} x_n / \sum_{n=1}^N r_{nk} = (1/|C_k|) \sum_{n|x_n \in C_k} x_n$$



Fuzzy C-Means Clustering

- K-Means was a crisp algorithm. Now: **fuzzy variant**
- Reformulate K-Means objective function with **membership matrix**
 r_{nk} : Membership of pattern x_n in class C_k

$$J_{SQE} = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$$

- optimization criterion

$$dJ_{SQE}/d\mu_k = 0$$

- together with non-overlapping constraint

$$\forall n(\exists k(r_{nk} = 1) \wedge ((k' \neq k) \rightarrow (r_{nk'} = 0)))$$

leads to well known K-Means

$$\mu_k = \sum_{n=1}^N r_{nk} x_n / \sum_{n=1}^N r_{nk} = (1/|C_k|) \sum_{n|x_n \in C_k} x_n$$



DBSCAN

- K-Means is „OK“ as cluster algorithm, but has certain **disadvantages**:
 - favors **spherical clusters**
 - **need to know K**
 - **no notion of noise**

- **Alternative → DBSCAN [4]**

(de facto state of the art):

- Idea: Two parameters: **minPt**, ϵ

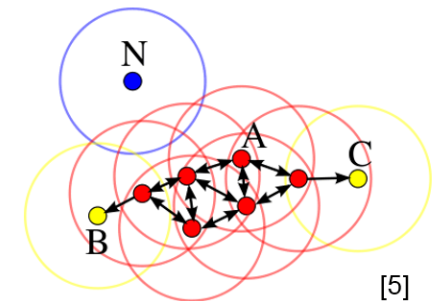
- Rough **idea**: **iterate**:

visit previously unseen pattern x :

if in ϵ -neighborhood $\{x'\}$ of x : $|\{x'\}| \geq \text{minPt}$ **then**

start new cluster: include x and $\{x'\}$ and those of their ϵ -neighborhoods $\{x''\}$ that are dense enough ($|\{x''\}| \geq \text{minPt}$), etc.

else: x is noise



[5]



- Dunn Index:

$$D = \min_{k_1 \in [1, K]} \left(\min_{k_2 \in [1, K]} \left(\frac{d_1(\mathcal{C}_{k_1}, \mathcal{C}_{k_2})}{\max_{k_3 \in [1, K]} d_2(\mathcal{C}_{k_3})} \right) \right)$$

where $d_1(\mathcal{C}_{k_1}, \mathcal{C}_{k_2})$ is the distance function between two clusters defined by

$$d_1(\mathcal{C}_{k_1}, \mathcal{C}_{k_2}) = \min_{\{(n_1, n_2) | x_{n_1} \in \mathcal{C}_{k_1} \wedge x_{n_2} \in \mathcal{C}_{k_2}\}} \|x_{n_1} - x_{n_2}\|$$

(that is the single link distance from SAHN).

The “diameter” d_2 of the clusters is defined by

$$d_2(\mathcal{C}_i) = \max_{\{(n_1, n_2) | x_{n_1} \in \mathcal{C}_i \wedge x_{n_2} \in \mathcal{C}_i\}} \|x_{n_1} - x_{n_2}\|$$



[7]

- K-Means was a crisp algorithm. Now: **fuzzy variant**
- Reformulate K-Means objective function with **membership matrix**
 r_{nk} : Membership of pattern x_n in class \mathcal{C}_k

$$J_{SQE} = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$$

- optimization criterion

$$dJ_{SQE}/d\mu_k = 0$$

- together with non-overlapping constraint

$$\forall n (\exists k (r_{nk} = 1) \wedge ((k' \neq k) \rightarrow (r_{nk'} = 0)))$$

leads to well known K-Means

$$\mu_k = \sum_{n=1}^N r_{nk} x_n / \sum_{n=1}^N r_{nk} = (1/|\mathcal{C}_k|) \sum_{n | x_n \in \mathcal{C}_k} x_n$$



- K-Means was a crisp algorithm. Now: **fuzzy variant**
- Reformulate K-Means objective function with **membership matrix**
 r_{nk} : Membership of pattern x_n in class \mathcal{C}_k

$$J_{SQE} = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$$

- optimization criterion

$$dJ_{SQE}/d\mu_k = 0$$

- together with non-overlapping constraint

$$\forall n (\exists k (r_{nk} = 1) \wedge ((k' \neq k) \rightarrow (r_{nk'} = 0)))$$

leads to well known K-Means

$$\mu_k = \sum_{n=1}^N r_{nk} x_n / \sum_{n=1}^N r_{nk} = (1/|\mathcal{C}_k|) \sum_{n | x_n \in \mathcal{C}_k} x_n$$



- Now **modify objective function** to:

$$J_{GSQE} = \sum_{n=1}^N \sum_{k=1}^K (r_{nk})^m \|x_n - \mu_k\|^2$$

- Exponent **m models degree of fuzzyness**:
 $m \rightarrow 1$: K-Means (crisp case);
 $m \rightarrow \infty$: $r_{nk} \rightarrow 1/K$ (where K is the number of clusters)

- **Optimize** the obj. fct. under the **conditions**:

$$\forall x_n : \sum_{k=1}^K \alpha_k(x_n) = \sum_{k=1}^K r_{nk} = 1$$

$$\forall \mathcal{C}_k : \sum_{n=1}^N \alpha_k(x_n) = \sum_{n=1}^N r_{nk} > 0$$



- Now **modify objective function** to:

$$J_{GSQE} = \sum_{n=1}^N \sum_{k=1}^K (r_{nk})^m \|x_n - \mu_k\|^2$$

- Exponent **m models degree of fuzziness**:
 $m \rightarrow 1$: K-Means (crisp case);
 $m \rightarrow \infty$: $r_{nk} \rightarrow 1/K$ (where K is the number of clusters)

- Optimize** the obj. fct. under the **conditions**:

$$\forall x_n : \sum_{k=1}^K \alpha_k(x_n) = \sum_{k=1}^K r_{nk} = 1$$

$$\forall \mathcal{C}_k : \sum_{n=1}^N \alpha_k(x_n) = \sum_{n=1}^N r_{nk} > 0$$



- Now **modify objective function** to:

$$J_{GSQE} = \sum_{n=1}^N \sum_{k=1}^K (r_{nk})^m \|x_n - \mu_k\|^2$$

- Exponent **m models degree of fuzziness**:
 $m \rightarrow 1$: K-Means (crisp case);
 $m \rightarrow \infty$: $r_{nk} \rightarrow 1/K$ (where K is the number of clusters)

- Optimize** the obj. fct. under the **conditions**:

$$\forall x_n : \sum_{k=1}^K \alpha_k(x_n) = \sum_{k=1}^K r_{nk} = 1$$

$$\forall \mathcal{C}_k : \sum_{n=1}^N \alpha_k(x_n) = \sum_{n=1}^N r_{nk} > 0$$



- Now **modify objective function** to:

$$J_{GSQE} = \sum_{n=1}^N \sum_{k=1}^K (r_{nk})^m \|x_n - \mu_k\|^2$$

- Exponent **m models degree of fuzziness**:
 $m \rightarrow 1$: K-Means (crisp case);
 $m \rightarrow \infty$: $r_{nk} \rightarrow 1/K$ (where K is the number of clusters)

- Optimize** the obj. fct. under the **conditions**:

$$\forall x_n : \sum_{k=1}^K \alpha_k(x_n) = \sum_{k=1}^K r_{nk} = 1$$

$$\forall \mathcal{C}_k : \sum_{n=1}^N \alpha_k(x_n) = \sum_{n=1}^N r_{nk} > 0$$



- Now **modify objective function** to:

$$J_{GSQE} = \sum_{n=1}^N \sum_{k=1}^K (r_{nk})^m \|x_n - \mu_k\|^2$$

- Exponent **m models degree of fuzziness**:
 $m \rightarrow 1$: K-Means (crisp case);
 $m \rightarrow \infty$: $r_{nk} \rightarrow 1/K$ (where K is the number of clusters)

- Optimize** the obj. fct. under the **conditions**:

$$\forall x_n : \sum_{k=1}^K \alpha_k(x_n) = \sum_{k=1}^K r_{nk} = 1$$

$$\forall \mathcal{C}_k : \sum_{n=1}^N \alpha_k(x_n) = \sum_{n=1}^N r_{nk} > 0$$



- K-Means was a crisp algorithm. Now: **fuzzy variant**
 - Reformulate K-Means objective function with **membership matrix**
- r_{nk} : Membership of pattern x_n in class C_k

$$J_{SQE} = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$$

- optimization criterion

$$dJ_{SQE}/d\mu_k = 0$$

- together with non-overlapping constraint

$$\forall n, k : (r_{nk} = 1) \wedge ((k' \neq k) \rightarrow (r_{nk'} = 0))$$

leads to well known K-Means

$$\mu_k = \sum_{n=1}^N r_{nk} x_n / \sum_{n=1}^N r_{nk} = (1/|C_k|) \sum_{n|x_n \in C_k} x_n$$



- Result:

$$r_{nk} = \left(\sum_{k'=1}^K \left(\frac{\|x_n - \mu_k\|}{\|x_n - \mu_{k'}\|} \right)^{\frac{2}{m-1}} \right)^{-1} \quad (\odot)$$

$$\mu_k = \frac{\sum_{n=1}^N r_{nk}^m x_n}{\sum_{n=1}^N r_{nk}^m} \quad (\odot \odot)$$

- the result assumes that no patterns and prototypes coincide

$$\forall n, k : \|x_n - \mu_k\| \neq 0$$

if they do coincide, set $r_{nk} = 1$ for $x_n = \mu_k$ and $r_{nk} = 0$ for $x_n \neq \mu_k$



- Limit $m \rightarrow \infty$ gives:

$$r_{nk} \xrightarrow{m \rightarrow \infty} \frac{1}{\sum_{k'=1}^K 1} = \frac{1}{K}$$

- Limit $m \rightarrow 1$ we get the nearest neighbor rule (K-Means) because:

$$r_{nk} = 1 / \left(\left(\sum_{k' \neq k} \left(\frac{\|x_n - \mu_k\|}{\|x_n - \mu_{k'}\|} \right)^{\frac{2}{m-1}} \right) + 1 \right)$$

in the limit $m \rightarrow 1$ the first sum in the denominator becomes ∞ if

$$\|x_n - \mu_k\| \neq \min_{1 \leq k' \leq K} \|x_n - \mu_{k'}\|$$

and it becomes 0 if

$$\|x_n - \mu_k\| = \min_{1 \leq k' \leq K} \|x_n - \mu_{k'}\|$$



- Limit $m \rightarrow \infty$ gives:

$$r_{nk} \xrightarrow{m \rightarrow \infty} \frac{1}{\sum_{k'=1}^K 1} = \frac{1}{K}$$

- Limit $m \rightarrow 1$ we get the nearest neighbor rule (K-Means) because:

$$r_{nk} = 1 / \left(\left(\sum_{k' \neq k} \left(\frac{\|x_n - \mu_k\|}{\|x_n - \mu_{k'}\|} \right)^{\frac{2}{m-1}} \right) + 1 \right)$$

in the limit $m \rightarrow 1$ the first sum in the denominator becomes ∞ if

$$\|x_n - \mu_k\| \neq \min_{1 \leq k' \leq K} \|x_n - \mu_{k'}\|$$

and it becomes 0 if

$$\|x_n - \mu_k\| = \min_{1 \leq k' \leq K} \|x_n - \mu_{k'}\|$$



- Limit $m \rightarrow \infty$ gives:

$$r_{nk} \xrightarrow{m \rightarrow \infty} \frac{1}{\sum_{k'=1}^K 1} = \frac{1}{K}$$

- Limit $m \rightarrow 1$ we get the nearest neighbor rule (K-Means) because:

$$r_{nk} = 1 / \left(\left(\sum_{k' \neq k} \left(\frac{\|x_n - \mu_k\|}{\|x_n - \mu_{k'}\|} \right)^{\frac{2}{m-1}} \right) + 1 \right)$$

in the limit $m \rightarrow 1$ the first sum in the denominator becomes ∞ if

$$\|x_n - \mu_k\| \neq \min_{1 \leq k' \leq K} \|x_n - \mu_{k'}\|$$

and it becomes 0 if

$$\|x_n - \mu_k\| = \min_{1 \leq k' \leq K} \|x_n - \mu_{k'}\|$$



- Limit $m \rightarrow \infty$ gives:

$$r_{nk} \xrightarrow{m \rightarrow \infty} \frac{1}{\sum_{k'=1}^K 1} = \frac{1}{K}$$

- Limit $m \rightarrow 1$ we get the nearest neighbor rule (K-Means) because:

$$r_{nk} = 1 / \left(\left(\sum_{k' \neq k} \left(\frac{\|x_n - \mu_k\|}{\|x_n - \mu_{k'}\|} \right)^{\frac{2}{m-1}} \right) + 1 \right)$$

in the limit $m \rightarrow 1$ the first sum in the denominator becomes ∞ if

$$\|x_n - \mu_k\| \neq \min_{1 \leq k' \leq K} \|x_n - \mu_{k'}\|$$

and it becomes 0 if

$$\|x_n - \mu_k\| = \min_{1 \leq k' \leq K} \|x_n - \mu_{k'}\|$$



- Limit $m \rightarrow \infty$ gives:

$$r_{nk} \xrightarrow{m \rightarrow \infty} \frac{1}{\sum_{k'=1}^K 1} = \frac{1}{K}$$

- Limit $m \rightarrow 1$ we get the nearest neighbor rule (K-Means) because:

$$r_{nk} = 1 / \left(\left(\sum_{k' \neq k} \left(\frac{\|x_n - \mu_k\|}{\|x_n - \mu_{k'}\|} \right)^{\frac{2}{m-1}} \right) + 1 \right)$$

in the limit $m \rightarrow 1$ the first sum in the denominator becomes ∞ if

$$\|x_n - \mu_k\| \neq \min_{1 \leq k' \leq K} \|x_n - \mu_{k'}\|$$

and it becomes 0 if

$$\|x_n - \mu_k\| = \min_{1 \leq k' \leq K} \|x_n - \mu_{k'}\|$$



- Now **modify objective function** to:

$$J_{GSQE} = \sum_{n=1}^N \sum_{k=1}^K (r_{nk})^m \|x_n - \mu_k\|^2$$

- Exponent **m models degree of fuzzyness**:
 $m \rightarrow 1$: K-Means (crisp case);
 $m \rightarrow \infty$: $r_{nk} \rightarrow 1/K$ (where K is the number of clusters)

- Optimize** the obj. fct. under the **conditions**:

$$\forall x_n : \sum_{k=1}^K \alpha_k(x_n) = \sum_{k=1}^K r_{nk} = 1$$

$$\forall \mathcal{C}_k : \sum_{n=1}^N \alpha_k(x_n) = \sum_{n=1}^N r_{nk} > 0$$

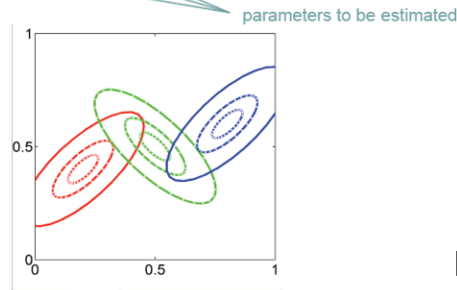


- Fuzzy C-Means is “OK” as a non-crisp clustering alg. but (as K-Means) favors spherical clusters → better approaches

- Example: **Gaussian Mixture Models (GMM)**

- Linear combination of Gaussians

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x|\mu_k, \Sigma_k) \text{ where } \sum_{k=1}^K \pi_k = 1, 0 \leq \pi_k \leq 1$$

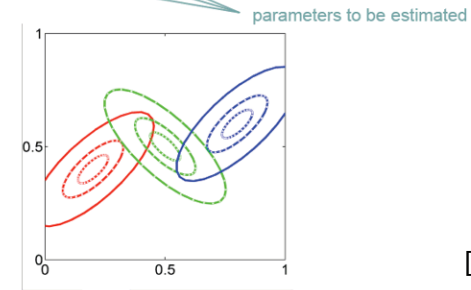


- Fuzzy C-Means is “OK” as a non-crisp clustering alg. but (as K-Means) favors spherical clusters → better approaches

- Example: **Gaussian Mixture Models (GMM)**

- Linear combination of Gaussians

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x|\mu_k, \Sigma_k) \text{ where } \sum_{k=1}^K \pi_k = 1, 0 \leq \pi_k \leq 1$$

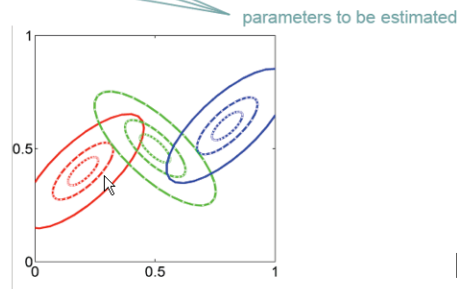


- Fuzzy C-Means is “OK” as a non-crisp clustering alg. but (as K-Means) favors spherical clusters → better approaches

- Example: **Gaussian Mixture Models (GMM)**

- Linear combination of Gaussians

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x|\mu_k, \Sigma_k) \text{ where } \sum_{k=1}^K \pi_k = 1, 0 \leq \pi_k \leq 1$$

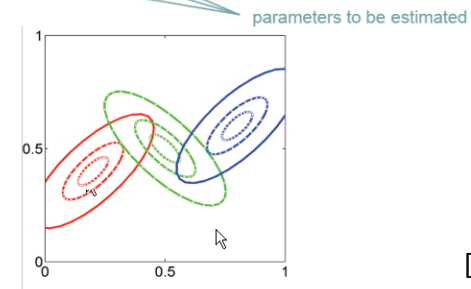


- Fuzzy C-Means is “OK” as a non-crisp clustering alg. but (as K-Means) favors spherical clusters → better approaches

- Example: **Gaussian Mixture Models (GMM)**

- Linear combination of Gaussians

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x|\mu_k, \Sigma_k) \text{ where } \sum_{k=1}^K \pi_k = 1, 0 \leq \pi_k \leq 1$$

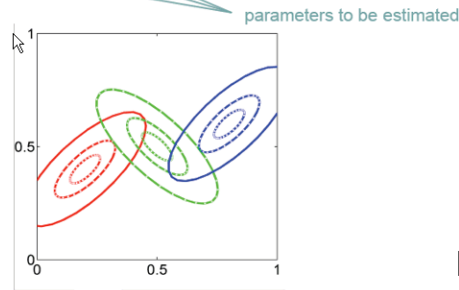


- Fuzzy C-Means is “OK” as a non-crisp clustering alg. but (as K-Means) favors spherical clusters → better approaches

Example: Gaussian Mixture Models (GMM)

- Linear combination of Gaussians

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x|\mu_k, \Sigma_k) \text{ where } \sum_{k=1}^K \pi_k = 1, \quad 0 \leq \pi_k \leq 1$$



- Fuzzy C-Means is “OK” as a non-crisp clustering alg. but (as K-Means) favors spherical clusters → better approaches

Example: Gaussian Mixture Models (GMM)

- Linear combination of Gaussians

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x|\mu_k, \Sigma_k) \text{ where } \sum_{k=1}^K \pi_k = 1, \quad 0 \leq \pi_k \leq 1$$

this is usually written as $p(x|\theta)$ denoting the dependency on the parameters $\theta = \{\pi_k, \mu_k, \Sigma_k\}_{k \in \{1,2,\dots,K\}}$

Writing this as a conditional probability makes sense in connection with Bayesian Machine Learning (see [8])



- Fuzzy C-Means is “OK” as a non-crisp clustering alg. but (as K-Means) favors spherical clusters → better approaches

Example: Gaussian Mixture Models (GMM)

- Linear combination of Gaussians

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x|\mu_k, \Sigma_k) \text{ where } \sum_{k=1}^K \pi_k = 1, \quad 0 \leq \pi_k \leq 1$$

this is usually written as $p(x|\theta)$ denoting the dependency on the parameters $\theta = \{\pi_k, \mu_k, \Sigma_k\}_{k \in \{1,2,\dots,K\}}$

Writing this as a conditional probability makes sense in connection with Bayesian Machine Learning (see [8])



- Fuzzy C-Means is “OK” as a non-crisp clustering alg. but (as K-Means) favors spherical clusters → better approaches

Example: Gaussian Mixture Models (GMM)

- Linear combination of Gaussians

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x|\mu_k, \Sigma_k) \text{ where } \sum_{k=1}^K \pi_k = 1, \quad 0 \leq \pi_k \leq 1$$

this is usually written as $p(x|\theta)$ denoting the dependency on the parameters $\theta = \{\pi_k, \mu_k, \Sigma_k\}_{k \in \{1,2,\dots,K\}}$

Writing this as a conditional probability makes sense in connection with Bayesian Machine Learning (see [8])



- For a **distribution** $p(x|\theta)$ parametrized by a set of **parameters** θ and iid data $X = \{x_1, x_2, \dots, x_N\}$, simple machine learning corresponds to **finding the θ that best explains the data**
- **iid**: „identically independently drawn“ $\Rightarrow p(X|\theta) = \prod_i p(x_i|\theta)$
- $p(X|\theta)$ is called **likelihood**
- „finding the θ that best explains the data“:
Maximum Likelihood: $\theta_{ML} = \operatorname{argmax}_{\theta} p(X|\theta) \Rightarrow \nabla_{\theta} p(X|\theta) \stackrel{!}{=} 0$
- convenient: use **log** $p(X|\theta)$ instead of $p(X|\theta)$
 $\Rightarrow \log p(X|\theta) = \sum_i \log p(x_{-i}|\theta)$



- For a **distribution** $p(x|\theta)$ parametrized by a set of **parameters** θ and iid data $X = \{x_1, x_2, \dots, x_N\}$, simple machine learning corresponds to **finding the θ that best explains the data**
- **iid**: „identically independently drawn“ $\Rightarrow p(X|\theta) = \prod_i p(x_i|\theta)$
- $p(X|\theta)$ is called **likelihood**
- „finding the θ that best explains the data“:
Maximum Likelihood: $\theta_{ML} = \operatorname{argmax}_{\theta} p(X|\theta) \Rightarrow \nabla_{\theta} p(X|\theta) \stackrel{!}{=} 0$
- convenient: use **log** $p(X|\theta)$ instead of $p(X|\theta)$
 $\Rightarrow \log p(X|\theta) = \sum_i \log p(x_{-i}|\theta)$



from here we follow [3], so citations for images etc. are omitted



- For a **distribution** $p(x|\theta)$ parametrized by a set of **parameters** θ and iid data $X = \{x_1, x_2, \dots, x_N\}$, simple machine learning corresponds to **finding the θ that best explains the data**
- **iid**: „identically independently drawn“ $\Rightarrow p(X|\theta) = \prod_i p(x_i|\theta)$
- $p(X|\theta)$ is called **likelihood**
- „finding the θ that best explains the data“:
Maximum Likelihood: $\theta_{ML} = \operatorname{argmax}_{\theta} p(X|\theta) \Rightarrow \nabla_{\theta} p(X|\theta) \stackrel{!}{=} 0$
- convenient: use **log** $p(X|\theta)$ instead of $p(X|\theta)$
 $\Rightarrow \log p(X|\theta) = \sum_i \log p(x_{-i}|\theta)$



- For a **distribution** $p(x|\theta)$ parametrized by a set of **parameters** θ and iid data $X = \{x_1, x_2, \dots, x_N\}$, simple machine learning corresponds to **finding the θ that best explains the data**

- iid: „identically independently drawn“ $\Rightarrow p(X|\theta) = \prod_i p(x_i|\theta)$

- $p(X|\theta)$ is called **likelihood**

- „finding the θ that best explains the data“:
Maximum Likelihood: $\theta_{ML} = \operatorname{argmax}_{\theta} p(X|\theta) \Rightarrow \nabla_{\theta} p(X|\theta) \stackrel{!}{=} 0$

- convenient: use **log** $p(X|\theta)$ instead of $p(X|\theta)$
 $\Rightarrow \log p(X|\theta) = \sum_i \log p(x_{-i}|\theta)$



- For a **distribution** $p(x|\theta)$ parametrized by a set of **parameters** θ and iid data $X = \{x_1, x_2, \dots, x_N\}$, simple machine learning corresponds to **finding the θ that best explains the data**

- iid: „identically independently drawn“ $\Rightarrow p(X|\theta) = \prod_i p(x_i|\theta)$

- $p(X|\theta)$ is called **likelihood**

- „finding the θ that best explains the data“:
Maximum Likelihood: $\theta_{ML} = \operatorname{argmax}_{\theta} p(X|\theta) \Rightarrow \nabla_{\theta} p(X|\theta) \stackrel{!}{=} 0$

- convenient: use **log** $p(X|\theta)$ instead of $p(X|\theta)$
 $\Rightarrow \log p(X|\theta) = \sum_i \log p(x_{-i}|\theta)$



- For a **distribution** $p(x|\theta)$ parametrized by a set of **parameters** θ and iid data $X = \{x_1, x_2, \dots, x_N\}$, simple machine learning corresponds to **finding the θ that best explains the data**

- iid: „identically independently drawn“ $\Rightarrow p(X|\theta) = \prod_i p(x_i|\theta)$

- $p(X|\theta)$ is called **likelihood**

- „finding the θ that best explains the data“:
Maximum Likelihood: $\theta_{ML} = \operatorname{argmax}_{\theta} p(X|\theta) \Rightarrow \nabla_{\theta} p(X|\theta) \stackrel{!}{=} 0$

- convenient: use **log** $p(X|\theta)$ instead of $p(X|\theta)$
 $\Rightarrow \log p(X|\theta) = \sum_i \log p(x_{-i}|\theta)$



- For a **distribution** $p(x|\theta)$ parametrized by a set of **parameters** θ and iid data $X = \{x_1, x_2, \dots, x_N\}$, simple machine learning corresponds to **finding the θ that best explains the data**

- iid: „identically independently drawn“ $\Rightarrow p(X|\theta) = \prod_i p(x_i|\theta)$

- $p(X|\theta)$ is called **likelihood**

- „finding the θ that best explains the data“:
Maximum Likelihood: $\theta_{ML} = \operatorname{argmax}_{\theta} p(X|\theta) \Rightarrow \nabla_{\theta} p(X|\theta) \stackrel{!}{=} 0$

- convenient: use **log** $p(X|\theta)$ instead of $p(X|\theta)$
 $\Rightarrow \log p(X|\theta) = \sum_i \log p(x_{-i}|\theta)$



- For a **distribution** $p(x|\theta)$ parametrized by a set of **parameters** θ and iid data $X = \{x_1, x_2, \dots, x_N\}$, simple machine learning corresponds to **finding the θ that best explains the data**
- **iid**: „independently drawn“ $\Rightarrow p(X|\theta) = \prod_i p(x_i|\theta)$
- $p(X|\theta)$ is called **likelihood**
- „finding the θ that best explains the data“:
Maximum Likelihood: $\theta_{ML} = \underset{\theta}{\operatorname{argmax}} p(X|\theta) \Rightarrow \nabla_{\theta} p(X|\theta) \stackrel{!}{=} 0$
- convenient: use **log** $p(X|\theta)$ instead of $p(X|\theta)$
 $\Rightarrow \log p(X|\theta) = \sum_i \log p(x_i|\theta)$

