

Machine Intelligence 2

4.1 K-means Clustering

Prof. Dr. Klaus Obermayer

Fachgebiet Neuronale Informationsverarbeitung (NI)

SS 2017

K-means Clustering

Projection methods vs. clustering

observations: $\{\underline{\mathbf{x}}^{(\alpha)}\}, \alpha = 1, \dots, p; \quad \underline{\mathbf{x}} \in \mathbb{R}^N$



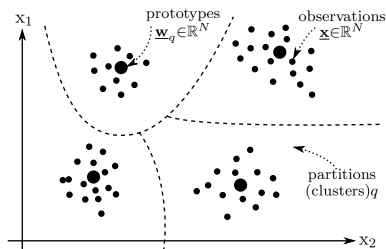
- ~ high-dimensional
- ~ groups, categories, hidden causes
- ~ interesting directions
- ~ "informative" manifolds

What is the relevant "structure"?

- ⇒ projection methods: search for "interesting" directions in feature space
- ⇒ clustering methods: grouping & categorization (and prototypes)

Central clustering

- ⇒ unsupervised formation of categories (partitions, clusters) according to predefined criteria
- ⇒ description of clusters by prototypes ← "central" clustering
- ⇒ goal: partitioning of observations $\underline{\mathbf{x}}^{(\alpha)}$, $\alpha = 1, \dots, p$; $\underline{\mathbf{x}}^{(\alpha)} \in \mathbb{R}^N$ according to similarity.



Cluster model



- \Rightarrow prototypes: $\underline{\mathbf{w}}_q$, $q = 1, \dots, M$ (M: number of clusters)
/Groups
- \Rightarrow binary assignment variables $m_q^{(\alpha)}$:



$$m_q^{(\alpha)} = \begin{cases} 1, & \text{if } \underline{\mathbf{x}}^{(\alpha)} \text{ belongs to cluster } q \\ 0, & \text{else} \end{cases}$$

- \Rightarrow normalization: $\sum_q m_q^{(\alpha)} = 1$

Cost function

Quality Criteria:

→ average quadratic distance between observations and prototypes

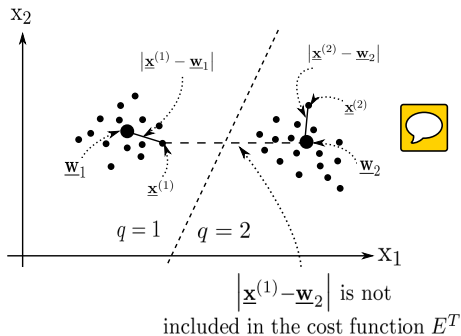
$$E^T[\{m_q^{(\alpha)}\}, \{\underline{\mathbf{w}}_q\}] = \frac{1}{p} \sum_{q, \alpha} m_q^{(\alpha)} (\underline{\mathbf{x}}^{(\alpha)} - \underline{\mathbf{w}}_q)^2$$

→ cost function implicitly quantifies our prior knowledge about the data

Model selection

Cost function:

- ⇒ cluster-centers: continuous variables
- ⇒ cluster-assignment: binary variables.
- ⇒ dissimilarity measure: squared Euclidean distance.



Batch K-means

Algorithm 1: batch K-means



random initialization of prototypes, e.g. $\underline{\mathbf{w}}_q = \underset{\text{center of mass}}{\langle \underline{\mathbf{x}} \rangle} + \underline{\eta}_q$, $\underline{\eta}_q$ small random vector

begin loop

(1) choose $m_q^{(\alpha)}$ such that E^T is minimal for the given prototypes



$$m_q^{(\alpha)} = \begin{cases} 1, & \text{if } q = \operatorname{argmin}_{\gamma} |\underline{\mathbf{x}}^{(\alpha)} - \underline{\mathbf{w}}_{\gamma}| \\ 0, & \text{else} \end{cases}$$

\Rightarrow assign every data point to its nearest prototype

(2) choose $\underline{\mathbf{w}}_q$ such that E^T is minimal for the -new- assignments



$$\underline{\mathbf{w}}_q = \frac{\sum_{\alpha} m_q^{(\alpha)} \underline{\mathbf{x}}^{(\alpha)}}{\sum_{\alpha} m_q^{(\alpha)}}$$

\Rightarrow set $\underline{\mathbf{w}}_q$ to the center of mass of its assigned data

end

Model selection - batch K-means

Cost function:

$$E^T[\{m_q^{(\alpha)}\}, \{\mathbf{w}_q\}] = \frac{1}{p} \sum_{q, \alpha} m_q^{(\alpha)} (\mathbf{x}^{(\alpha)} - \mathbf{w}_q)^2$$

\Rightarrow "Center of mass is optimal" *why is it optimal?*

$$\frac{\partial}{\partial \mathbf{w}_q} \left\{ \frac{1}{2p} \sum_{q', \alpha} m_{q'}^{(\alpha)} (\mathbf{x}^{(\alpha)} - \mathbf{w}_{q'})^2 \right\} = -\frac{2}{p} \sum_{\alpha} m_q^{(\alpha)} (\mathbf{x}^{(\alpha)} - \mathbf{w}_q) \stackrel{!}{=} 0$$

$$\leadsto \mathbf{w}_q = \frac{\sum_{\alpha} m_q^{(\alpha)} \mathbf{x}^{(\alpha)}}{\sum_{\alpha} m_q^{(\alpha)}} \text{ \textit{number of data points in group}}$$

Batch K-means

Assumption: All Clusters have same size

Cost function:


$$E^T[\{m_q^{(\alpha)}\}, \{\mathbf{w}_q\}] = \frac{1}{p} \sum_{q,\alpha} m_q^{(\alpha)} (\mathbf{x}^{(\alpha)} - \mathbf{w}_q)^2$$

cost function enspricht variance

⇒ condition for minimum: taking second derivatives of cost function

second derivative

$$\frac{\partial^2}{\partial \mathbf{w}_{qi} \partial \mathbf{w}_{q'j}} \left\{ \frac{1}{p} \sum_{q'',\alpha} m_{q''}^{(\alpha)} (\mathbf{x}^{(\alpha)} - \mathbf{w}_{q''})^2 \right\}$$

$$= \frac{\partial}{\partial \mathbf{w}_{q'j}} \left\{ -\frac{2}{p} \sum_{\alpha} m_q^{(\alpha)} (x_i^{(\alpha)} - (\mathbf{w})_{qi}) \right\} = \left(\frac{2}{p} \sum_{\alpha} m_q^{(\alpha)} \right) \delta_{ij} \delta_{qq'}$$


⇒ Diagonal matrix with all positive entries → condition for minimum is always satisfied.

⇒ Note: Minimizing E^T is not convex optimization problem.

Batch K-means (continued)


$$E^T[\{m_q^{(\alpha)}\}, \{\underline{\mathbf{w}}_q\}] = \frac{1}{p} \sum_{q, \alpha} m_q^{(\alpha)} (\underline{\mathbf{x}}^{(\alpha)} - \underline{\mathbf{w}}_q)^2$$

- \Rightarrow If $\underline{\mathbf{w}}_q$ is center of mass $\implies E^T = \text{variance}$.
- $\Rightarrow E^T$ is non-increasing in every step and E^T is bounded from below \rightarrow K-means clustering converges to a (local) optimum of E^T .
- $\Rightarrow E^T$ at the solution can be interpreted as the "size" (variance) of the clusters.

On-line K-means

Algorithm 2: on-line k-Means

random initialization of prototypes, e.g.


$\underline{\mathbf{w}}_q = \langle \underline{\mathbf{x}} \rangle + \underline{\eta}_q$, $\underline{\eta}_q$ small random vector 

select learning step $0 < \varepsilon \ll 1$

begin loop

choose a data point $\underline{\mathbf{x}}^{(\alpha)}$

assign data point to its closest prototype q

$$q = \underset{\gamma}{\operatorname{argmin}} |\underline{\mathbf{x}}^{(\alpha)} - \underline{\mathbf{w}}_{\gamma}|$$


change corresponding prototype according to

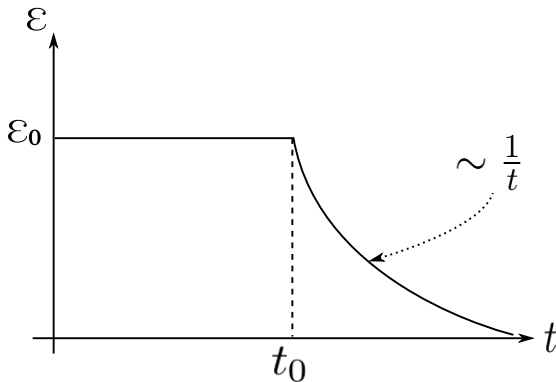

$$\Delta \underline{\mathbf{w}}_q = \varepsilon (\underline{\mathbf{x}}^{(\alpha)} - \underline{\mathbf{w}}_q)$$


change ε

end

On-line K-means

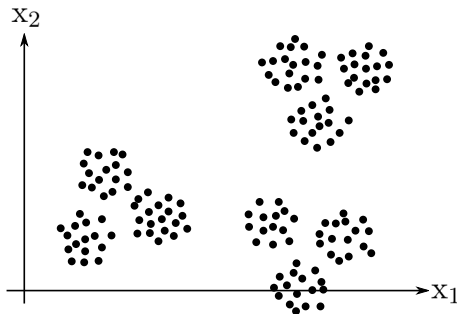
- more robust than batch-learning w.r.t convergence to local minima
- useful for streaming data
- quality of the found solution depends on choosing an appropriate "annealing" schedule for ε : Robbins-Monro conditions



Number of prototypes

- M : hyperparameter


Choice of resolution



1 cluster
3 cluster
9 cluster
many cluster?

⇒ additional assumptions
are needed!

Number of Prototypes: Choice of resolution


■ E_{\min}^T : average size of cluster (in terms of variance) 

→ large for few clusters – small for many clusters

→ zero, if number of cluster $\hat{=}$ number of data points

→ E_{\min}^T goes down if M increases.

■ choice of resolution

 → clusters “smaller” than the variance of noise probably do not capture meaningful structure

→ $E_{\min}^T \geq \sigma_{\text{noise}}^2$ which is a natural boundary on E_{\min}^T

Iterative refinement


Algorithm 3: iterative refinement

initialization: $\underline{\mathbf{w}}_1 = \frac{1}{p} \sum_{\alpha} \underline{\mathbf{x}}^{(\alpha)}$, $\underbrace{(E_{\min}^T)^*}_{\substack{\text{desired minimal} \\ \text{variance}}}$, $M = 1$ everything smaller is noise

begin loop

if $E_{\min}^T < (E_{\min}^T)^*$ then STOP

select partition $q \in \{1, \dots, M\}$ with largest variance

$$q = \operatorname{argmax}_{\gamma} \left(\frac{\sum_{\alpha} m_{\gamma}^{(\alpha)} (\underline{\mathbf{x}}^{(\alpha)} - \underline{\mathbf{w}}_{\gamma})^2}{\sum_{\alpha} m_{\gamma}^{(\alpha)}} \right)$$



add a new prototype: $\underline{\mathbf{w}}_{M+1} = \underline{\mathbf{w}}_q + \underbrace{\underline{\varepsilon}_q}_{\substack{\text{small random} \\ \text{vector}}}$

$M \leftarrow M + 1$

do K-means clustering with these M prototypes

end whole k means clustering for every iteration


Robustness of the clustering solution


- Solution should capture meaningful structure in the data
-  Multiple runs with different initializations should yield similar solutions.
- *Caveat:* Permutation of labels does neither change cost nor character of the solution.

$1, 2, 3, \dots, M$ [Siehe Aufschrieb onenote](#)
 $9, 1, M, \dots, 7$

- $M!$ trivially equivalent optima \rightarrow robustness-criterion has to be adapted.
- Avoid "instability": many structurally different clustering solutions with equal cost

Validation measures

- Model free approaches: Stability based validation
- **Idea:** taking too many or too few clusters leads to unstable partitions
- $\mathbf{X} = \{\mathbf{x}^{(\alpha)}\}, \alpha = 1, \dots, p; \quad \mathbf{x} \in \mathbb{R}^N$; A solution of the clustering algorithm is $\mathbf{Y} = (y_1, \dots, y_p)$ where $y_i \in L := \{1, \dots, M\}$ 
- Comparing clustering solutions Y_1 and Y_2 :

$$d := \frac{1}{|\mathbf{Y}_1|} \sum_{\alpha} \mathbf{1}\{Y_{1,\alpha} \neq Y_{2,\alpha}\} $$

Caveats: Permutation, optimal classifier (e.g. KNN, see Lange et. al., 2004)

Validation measure

Algorithm 4: Validation measure

begin for each $M \in \{M_{min}, \dots, M_{max}\}$

 **begin** loop for r splits of data

split data \mathbf{X} into $\mathbf{X}_1^{(i)}$ and $\mathbf{X}_2^{(i)}$ and find corresponding clustering solution $\mathbf{Y}_1^{(i)}$ and $\mathbf{Y}_2^{(i)}$ by applying clustering algorithm

use $(\mathbf{X}_1^{(i)}, \mathbf{Y}_1^{(i)})$ to train a nearest neighbor classifier ϕ and compute $\phi[X_2]$

compute distance $d_i := \frac{1}{|\mathbf{X}_2|} \sum_{\alpha} \mathbf{1} \left\{ \mathbf{Y}_2^{(\alpha)} \neq \phi[\mathbf{X}_2^{(\alpha)}] \right\}$ (but considering the permutations)
 \mathbf{Y}_1

end

Average distance between partitions: $\hat{S}_{clustering} = \frac{1}{r} \sum_r d_i$



Sample s random clustering assignments and compute empirical average of the distances to estimate \hat{S}_{random}

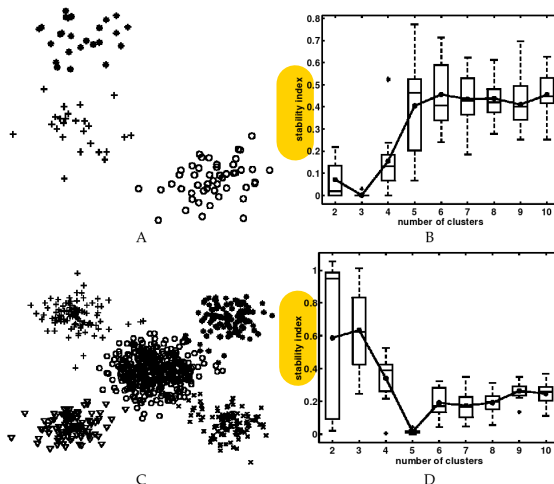
Calculate stability index: $\bar{S}_M = \frac{\hat{S}_{clustering}}{\hat{S}_{random}}$



end

Return $\hat{M} = \operatorname{argmin}_M(\bar{S}_M)$

K-means: Gaussian data



Caveats: Permutation, optimal classifier (e.g. KNN, see Lange et. al., 2004)

Further remarks

Alternative clustering approaches

- distribution based models ("model-based" \Rightarrow Gaussian Mixture algorithm)
- density based models
- hierarchical (connectivity based) clustering
 - single linkage (\sim nearest neighbor)
 - complete linkage
 - average linkage / within group ssq (Ward criterion)
 - agglomerative vs. divisive clustering

Current issues

- "big data": pre-processing (e.g. preselect spatial methods & KD-trees)
- Graph-based approaches & spectral clustering

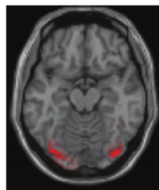
Applications

Image segmentation & compression

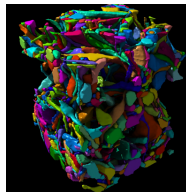
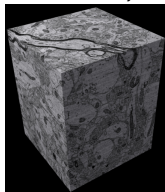


- k-means for pixels (e.g. RGB)
- segmentation via cluster-assignment (naive: context, smoothness ...)
- compression: $(N_{\text{pix}} \times S(R, G, B) \rightarrow N_{\text{pix}} \times k + \text{codebook})$

Neuroscience (fMRI, connectomics)



Litvak et.al(2001)



Seung lab