

# Machine Intelligence 2

## 4.1 K-means Clustering

Prof. Dr. Klaus Obermayer

Fachgebiet Neuronale Informationsverarbeitung (NI)

SS 2018

# K-means Clustering

# Projection methods vs. clustering

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



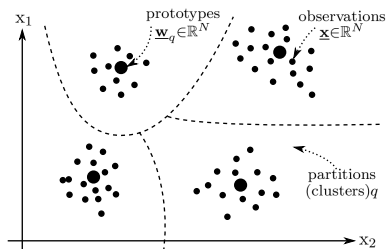
- ~ high-dimensional
- ~ groups, categories, hidden causes
- ~ interesting directions
- ~ "informative" manifolds

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

⇒ prototypes:  $\underline{\mathbf{w}}_q$ ,  $q = 1, \dots, M$  (M: number of clusters)

⇒ binary assignment variables  $m_q^{(\alpha)}$ :

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

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

# Cost function

→ average quadratic distance between observations and prototypes

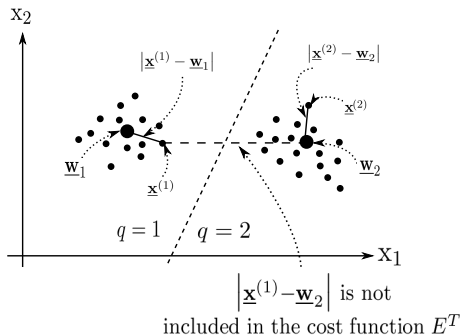
$$E^T[\{m_q^{(\alpha)}\}, \{\underline{\mathbf{w}}_q\}] = \frac{1}{p} \sum_{q, \alpha} m_q^{(\alpha)} (\underline{\mathbf{x}}^{(\alpha)} - \underline{\mathbf{w}}_q)^2$$

→ 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.  $\underline{\mathbf{w}}_q = \langle \underline{\mathbf{x}} \rangle + \underline{\eta}_q$ ,  $\underline{\eta}_q$  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**



# Batch K-means

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

$\Rightarrow$  "condition for extremum"

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

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

# Batch K-means

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

$\Rightarrow$  condition for minimum:

$$\begin{aligned} & \frac{\partial^2}{\partial \mathbf{w}_{qi} \partial \mathbf{w}_{q'j}} \left\{ \frac{1}{p} \sum_{q'',\alpha} m_{q''}^{(\alpha)} (\mathbf{x}^{(\alpha)} - \mathbf{w}_{q''})^2 \right\} \\ &= \frac{\partial}{\partial \mathbf{w}_{q'j}} \left\{ -\frac{2}{p} \sum_{\alpha} m_q^{(\alpha)} (x_i^{(\alpha)} - (\mathbf{w})_{qi}) \right\} = \left( \frac{2}{p} \sum_{\alpha} m_q^{(\alpha)} \right) \delta_{ij} \delta_{qq'} \end{aligned}$$

$\Rightarrow$  diagonal matrix with all positive entries  $\rightarrow$  condition for minimum is always satisfied.

$\Rightarrow$  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  $q$

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

change corresponding prototype according to

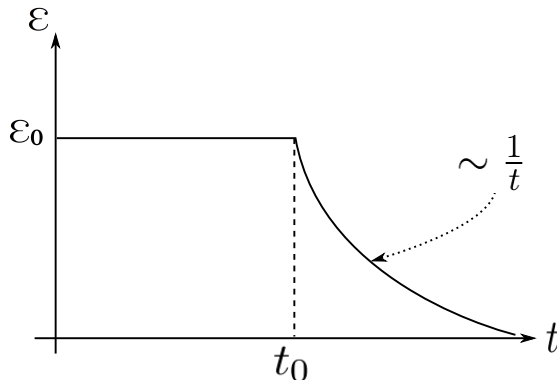
$$\Delta \underline{\mathbf{w}}_q = \varepsilon (\underline{\mathbf{x}}^{(\alpha)} - \underline{\mathbf{w}}_q)$$

change  $\varepsilon$

**end**

# On-line K-means

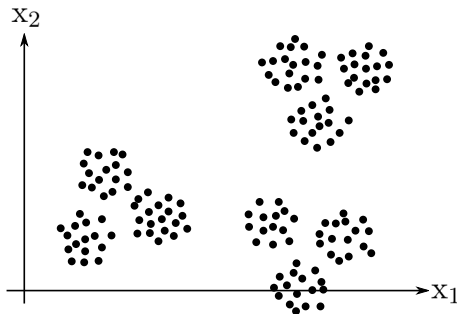
- more robust than batch-learning w.r.t convergence to local minima
- useful for streaming data
- quality of the found solution depends on choosing an appropriate "annealing" schedule for  $\varepsilon$ : Robbins-Monro conditions



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

begin loop

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

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

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

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

$M \leftarrow M + 1$

  do K-means clustering with these  $M$  prototypes

end



# 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.

$$1, 2, 3, \dots, M$$

$$9, 1, M, \dots, 7$$

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

# Validation measure

- Model free approaches: Stability based validation
- **Idea:** taking too many or too few clusters leads to unstable partitions
- **Data:**  $\mathbf{X}$  is a set containing all the data points *s.t.*  
 $\mathbf{X} = \{\underline{\mathbf{x}}^{(\alpha)}\}, \alpha = 1, \dots, p; \quad \underline{\mathbf{x}} \in \mathbb{R}^N$
- A solution of the clustering algorithm produces labellings set  $\mathbf{Y}$  *s.t.*  
 $\mathbf{Y} = \{y^{(\alpha)}\}$  where  $y^{(\alpha)} \in L := \{1, \dots, M\}$
- Dissimilarity between clustering solutions  $\mathbf{Y}_1$  and  $\mathbf{Y}_2$ :

$$d := \frac{1}{|\mathbf{Y}_1|} \sum_{\alpha} \mathbf{1} \left\{ y_1^{(\alpha)} \neq y_2^{(\alpha)} \right\}$$

Lange et. al., 2004, Neural Computation

# Validation measure

---

## Algorithm 4: Validation measure

---

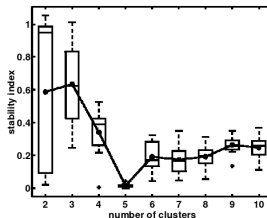
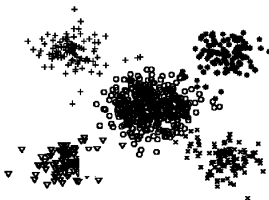
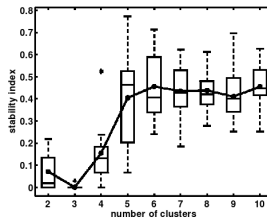
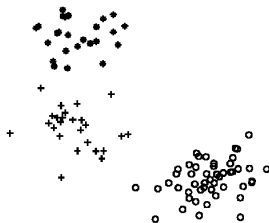
```

begin for each  $M \in \{M_{min}, \dots, M_{max}\}$ 
  begin loop for  $r$  splits of data
    Split  $\mathbf{X}$  into disjoint sets  $\mathbf{X}_1$  and  $\mathbf{X}_2$  randomly
    Apply clustering algorithm to both  $\mathbf{X}_1$  and  $\mathbf{X}_2$  to find  $\mathbf{Y}_1$  and  $\mathbf{Y}_2$ 
    Compute dissimilarity  $d_i := \frac{1}{|\mathbf{Y}_2|} \sum_{\alpha=1}^p \mathbf{1} \left\{ \mathbf{Y}_2^{(\alpha)} \neq \phi_1[\mathbf{X}_2^{(\alpha)}] \right\}$ 
    Where  $\phi_1$  denotes a classifier trained on  $(\mathbf{X}_1, \mathbf{Y}_1)$ 
  end
  Compute average dissimilarity:  $\hat{S}_A = \frac{1}{r} \sum_r d_i$ 
  Compute average dissimilarity  $\hat{S}_R$  of an clustering algorithm which assigns clusters
  randomly using the same data
  Calculate stability index:  $\bar{S}_M = \frac{\hat{S}_A}{\hat{S}_R}$ 
end
Return  $\hat{M} = \operatorname{argmin}_M(\bar{S}_M)$ 

```

---

# K-means: Gaussian data



## Further remarks

### Alternative clustering approaches

- density based models ("model-based"  $\Rightarrow$  Gaussian Mixture algorithm)
- 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 (improvements: context, smoothness ...)
- data: values  $\rightarrow$  label and residues, MDL principle for cluster number