

Introduction to Machine Learning

Lecture 2: Revisiting Clustering by Notion of Risk

Ali Bereyhi

`ali.bereyhi@utoronto.ca`

Department of Electrical and Computer Engineering
University of Toronto

Winter 2026

Recap: K -Means Clustering Algorithm

K -Means():

- 1: Initiate μ_1, \dots, μ_K
- 2: **while** μ_1, \dots, μ_K changing **do**
- 3: Set $\mathcal{C}_1, \dots, \mathcal{C}_K \leftarrow \text{Cluster_Assignment}(\mu_1, \dots, \mu_K)$
- 4: Update $\mu_1, \dots, \mu_K \leftarrow \text{Centroid_Update}(\mathcal{C}_1, \dots, \mathcal{C}_K)$
- 5: **end while**
- 6: Return μ_1, \dots, μ_K

$\text{Cluster_Assignment}(\mu_1, \dots, \mu_K)$:

$$\text{index of cluster for } x_n \leftarrow \underset{k \in \{1, \dots, K\}}{\operatorname{argmin}} \|x_n - \mu_k\|$$

$\text{Centroid_Update}(\mathcal{C}_1, \dots, \mathcal{C}_K)$:

$$\mu_k \leftarrow \text{average} \{x_n \in \mathcal{C}_k\} = \frac{1}{|\mathcal{C}_k|} \sum_{x_n \in \mathcal{C}_k} x_n$$

K -Means Clustering is **Good**!

When we derived the K -means clustering algorithm, we considered

K -centroid model for clustering, i.e.,

$$f(\mathbf{x}) = \operatorname{argmin}_{k \in \{1, \dots, K\}} \|\mathbf{x} - \boldsymbol{\mu}_k\|$$

for K centroids $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K \in \mathbb{R}^d$

K -means clustering is a **learning algorithm**

$$\mathcal{A} : \mathbb{D} \mapsto \boldsymbol{\mu}_1^*, \dots, \boldsymbol{\mu}_K^*$$

- ! We claim it finds a **good** model
- ? How can you define “**good**” set of centroids?
- ! Let's see

Clustering: *Alternative Formulation*

We have dataset \mathbb{D} : we want to

group samples into K clusters described by K centroids

For each $\mathbf{x}_n \in \mathbb{D}$, we define *K weights* $r_{n,1}, \dots, r_{n,K}$

$$r_{n,k} = \begin{cases} 1 & \text{if } \mathbf{x}_n \in \mathcal{C}_k \\ 0 & \text{otherwise} \end{cases}$$

Properties of $r_{n,k}$

$$\sum_{k=1}^K r_{n,k} = 1 \quad \text{and} \quad \sum_{n=1}^N r_{n,k} = |\mathcal{C}_k|$$

K-Means Clustering: *Alternative Formulation*

Cluster_Assignment(μ_1, \dots, μ_K):

1: **for** $n = 1 : N$ **do**

2: Assign K weights $r_{n,1}, \dots, r_{n,K}$ to sample x_n as

$$r_{n,k} = \begin{cases} 1 & \text{if } \mu_k = \text{closest centroid to } x_n \\ 0 & \text{otherwise} \end{cases}$$

3: **end for**

4: Return $r_{n,k}$ for $k = 1 : K$ and $n = 1 : N$

K-Means Clustering: Alternative Representation

Centroid_Update($\{r_{n,k}\}$):

- 1: **for** $k = 1 : K$ **do**
- 2: **if** $\sum_n r_{n,k} > 0$ **then**
- 3: Move μ_k to the center of cluster k , i.e.,

$$\mu_k = \frac{\sum_{n=1}^N r_{n,k} \mathbf{x}_n}{\sum_{n=1}^N r_{n,k}}$$

- 4: **else**
- 5: Leave μ_k unchanged
- 6: **end if**
- 7: **end for**
- 8: Return μ_1, \dots, μ_K

K-Means Clustering: Alternative Representation

We could iterate till we converge

K-Means() :

- 1: Initiate μ_1, \dots, μ_K
- 2: **while** μ_1, \dots, μ_K changing **do**
- 3: Set $\{r_{n,k}\} \leftarrow \text{Cluster_Assignment}(\mu_1, \dots, \mu_K)$
- 4: Update $\mu_1, \dots, \mu_K \leftarrow \text{Centroid_Update}(\{r_{n,k}\})$
- 5: **end while**
- 6: Return μ_1, \dots, μ_K

Defining Objective: Risk

? Where should be *good* centroids?

! Probably, close to the samples of that cluster!

Let's aggregate how far samples are from their centroid

$$\mathcal{J}(\{r_{n,k}\}, \{\mu_k\}) = \frac{1}{N} \sum_{k=1}^K \sum_{n=1}^N r_{n,k} \|\mathbf{x}_n - \mu_k\|^2$$

- $r_{n,k} \|\mathbf{x}_n - \mu_k\|^2$ is non-zero only if $\mathbf{x}_n \in \mathcal{C}_k$
- Sum over n aggregates distances of samples in cluster k from μ_k
- Sum over k aggregates distances of all samples from *their clusters*

We compute how far each point is *in average* from its centroid

Notion of Optimality

? Where should be *good* centroids?

! Probably, close to the samples of that cluster!

Optimal Clustering with K -Centroid Models

Optimal assignments $\{r_{n,k}^\star\}$ and centroids $\{\mu_k^\star\}$ minimize the *risk*

$$\{r_{n,k}^\star\}, \{\mu_k^\star\} = \operatorname{argmin}_{\{r_{n,k}\}, \{\mu_k\}} \mathcal{J}(\{r_{n,k}\}, \{\mu_k\})$$

Clustering by Risk Minimization

Risk minimization for clustering is hard, so we can use alternating optimization

Risk_Minimization():

1: Initiate μ_1^*, \dots, μ_K^*

2: **while** μ_1^*, \dots, μ_K^* changing **do**

3: Minimize the risk for fixed centroids $\{\mu_k^*\}$

$$\{r_{n,k}^*\} \leftarrow \operatorname{argmin}_{\{r_{n,k}\}} \mathcal{J}(\{r_{n,k}\}, \{\mu_k^*\})$$

4: Minimize the risk for fixed assignments $\{r_{n,k}^*\}$

$$\{\mu_k^*\} \leftarrow \operatorname{argmin}_{\{\mu_k\}} \mathcal{J}(\{r_{n,k}^*\}, \{\mu_k\})$$

5: **end while**

6: Return $\mu_1^*, \dots, \mu_K^* \approx \mu_1^*, \dots, \mu_K^*$

Clustering via Risk Minimization

Minimize the risk for fixed centroids $\{\mu_k^*\}$

$$\begin{aligned}
 \{r_{n,k}^*\} &= \operatorname{argmin}_{\{r_{n,k}\}} \mathcal{J}(\{r_{n,k}\}, \{\mu_k^*\}) \\
 &= \operatorname{argmin}_{\{r_{n,k}\}} \frac{1}{N} \sum_{k=1}^K \sum_{n=1}^N r_{n,k} \|\mathbf{x}_n - \mu_k^*\|^2 \\
 &= \operatorname{argmin}_{\{r_{n,k}\}} \frac{1}{N} \sum_{n=1}^N (r_{n,1} \|\mathbf{x}_n - \mu_1^*\|^2 + \dots + r_{n,K} \|\mathbf{x}_n - \mu_K^*\|^2)
 \end{aligned}$$

So, the solution is given by setting for each n

$$r_{n,k}^* = \begin{cases} 1 & \text{if } k = \operatorname{argmin}_{j \in \{1, \dots, K\}} \|\mathbf{x}_n - \mu_j^*\| \\ 0 & \text{otherwise} \end{cases}$$

Clustering via Risk Minimization

Minimizing the risk for fixed centroids $\{\mu_k^*\}$ is accomplished as

```
Cluster_Assignment( $\mu_1, \dots, \mu_K$ ):
```

```
1: for  $n = 1 : N$  do
```

```
2:   Assign  $K$  weights  $r_{n,1}^*, \dots, r_{n,K}^*$  to sample  $x_n$  as
```

$$r_{n,k}^* = \begin{cases} 1 & \text{if } k = \operatorname{argmin}_{j \in \{1, \dots, K\}} \|x_n - \mu_j\| \\ 0 & \text{otherwise} \end{cases}$$

```
3: end for
```

```
4: Return  $r_{n,k}^*$  for  $k = 1 : K$  and  $n = 1 : N$ 
```

Clustering by Risk Minimization

Minimize the risk for fixed assignments $\{r_{n,k}^*\}$

$$\begin{aligned}\{\mu_k^*\} &= \operatorname{argmin}_{\{\mu_k\}} \mathcal{J}(\{r_{n,k}^*\}, \{\mu_k\}) \\ &= \operatorname{argmin}_{\{\mu_k\}} \frac{1}{N} \sum_{k=1}^K \sum_{n=1}^N r_{n,k}^* \|\mathbf{x}_n - \mu_k\|^2 \\ &= \operatorname{argmin}_{\{\mu_k\}} \frac{1}{N} \sum_{k=1}^K \left(\sum_{\mathbf{x}_n \in \mathcal{C}_k} \|\mathbf{x}_n - \mu_k\|^2 \right)\end{aligned}$$

The solution for each k is given by

$$\mu_k^* = \frac{1}{|\mathcal{C}_k|} \sum_{\mathbf{x}_n \in \mathcal{C}_k} \mathbf{x}_n$$

Clustering by Risk Minimization

Minimizing the risk for fixed assignments $\{r_{n,k}^*\}$ is done by

Centroid_Update($\{r_{n,k}^*\}$):

- 1: **for** $k = 1 : K$ **do**
- 2: **if** $\sum_n r_{n,k}^* > 0$ **then**
- 3: Move μ_k to the center of cluster k specified by $r_{1,k}^*, \dots, r_{N,k}^*$
- 4: **else**
- 5: Leave μ_k unchanged
- 6: **end if**
- 7: **end for**
- 8: Return μ_1, \dots, μ_K

Clustering by *Risk Minimization*

Risk_Minimization():

1: Initiate μ_1^*, \dots, μ_K^*

2: **while** μ_1^*, \dots, μ_K^* changing **do**

3: Minimize the risk for fixed centroids $\{\mu_k^*\}$

$$\{r_{n,k}^*\} \leftarrow \text{Cluster_Assignment}(\mu_1^*, \dots, \mu_K^*)$$

4: Minimize the risk for fixed assignments $\{r_{n,k}^*\}$

$$\{\mu_k^*\} \leftarrow \text{Centroid_Update}(\{r_{n,k}^*\})$$

5: **end while**

6: Return $\mu_1^*, \dots, \mu_K^* \approx \mu_1^*, \dots, \mu_K^*$

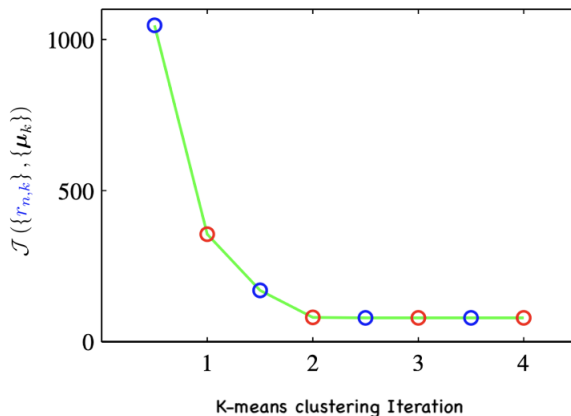
So we conclude that

$$\text{Risk_Minimization}() \equiv K\text{-Means}()$$

K -Means Clustering \equiv Risk Minimization

Risk also gives us means to **evaluate** the learned pattern

Back to our binary example



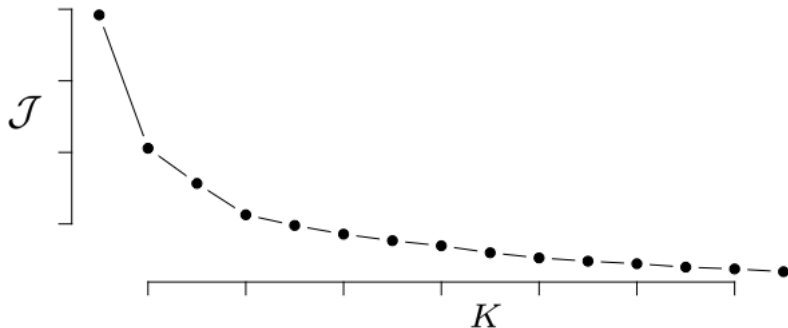
More Sophisticated Example: Segmentation¹

Each RGB pixel is a sample $x_n \in \mathbb{R}^3$: we cluster with $K = 10$ $K = 3$ $K = 2$



Choice of *Hyperparameter*

- ? How do we know K ?
- ! This is a *hyperparameter*

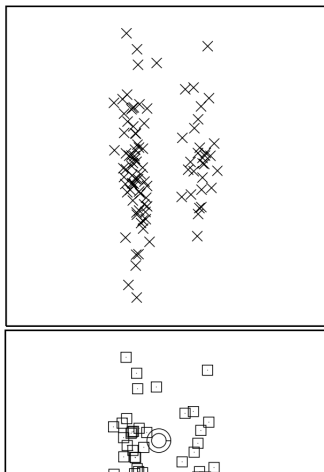


K -Means Clustering Always Converge

? Does K -means clustering always **converge** to a **stable state**?

! **Yes!** You can show it!

However, it does **not necessary** end with **what we want!**²



Soft Clustering

Recall that

$$\sum_{k=1}^K r_{n,k} = 1$$

? What if we think of any $r_{n,k} \in [0, 1]$?!

Probabilistic Assignment

In probabilistic assignment we assume $r_{n,k} \in [0, 1]$ such that

$$\sum_{k=1}^K r_{n,k} = 1$$

$r_{n,k}$ is hence a probability, i.e., probability of $x_n \in \mathbb{C}_k$

Risk as Expected Error

By this viewpoint, we can interpret the risk as an expected error

Risk \equiv Expected Error

Let $r_{n,k}$ be probability of $\mathbf{x}_n \in \mathbb{C}_k$

$$\begin{aligned}\mathcal{J}(\{r_{n,k}\}, \{\boldsymbol{\mu}_k\}) &= \frac{1}{N} \sum_{k=1}^K \sum_{n=1}^N r_{n,k} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2 \\ &= \mathbb{E} \{ \mathcal{E}(\mathbf{x}) \}\end{aligned}$$

with $\mathcal{E}(\mathbf{x})$ quantifying how *bad* we have classified, i.e.,

$$\mathcal{E}(\mathbf{x}) = \|\mathbf{x} - \text{cluster}(\mathbf{x})\|^2$$

Soft Clustering

We could revise our algorithm for *soft clustering*

Soft Clustering with K -Centroid Model

Optimal assignments $\{r_{n,k}^*\}$ and centroids $\{\mu_k^*\}$ minimize the *risk*

$$\{r_{n,k}^*\}, \{\mu_k^*\} = \underset{\{r_{n,k}\}, \{\mu_k\}}{\operatorname{argmin}} \mathcal{J}(\{r_{n,k}\}, \{\mu_k\})$$

for $\mu_k \in \mathbb{R}^d$ and $r_{n,k} \in [0, 1]$ such that

$$\sum_{k=1}^K r_{n,k} = 1.$$

Soft K -Means Clustering Algorithm

Soft_Cluster_Assignment(μ_1, \dots, μ_K):

1: **for** $n = 1 : N$ **do**

2: Assign K weights $r_{n,1}, \dots, r_{n,K}$ to sample x_n for some β as

$$r_{n,k} = \frac{e^{-\beta \|x_n - \mu_k\|^2}}{\sum_{k=1}^K e^{-\beta \|x_n - \mu_k\|^2}}$$

3: **end for**

4: Return $r_{n,k}$ for $k = 1 : K$ and $n = 1 : N$

? Does it remind you of any distribution?!

Soft K -Means Clustering Algorithm

Centroid_Update($\{r_{n,k}\}$):

1: **for** $k = 1 : K$ **do**

2: Move μ_k to the center of cluster k , i.e.,

$$\mu_k = \frac{\sum_{n=1}^N r_{n,k} \mathbf{x}_n}{\sum_{n=1}^N r_{n,k}}$$

3: **end for**

4: Return μ_1, \dots, μ_K

Soft K -Means Clustering Algorithm

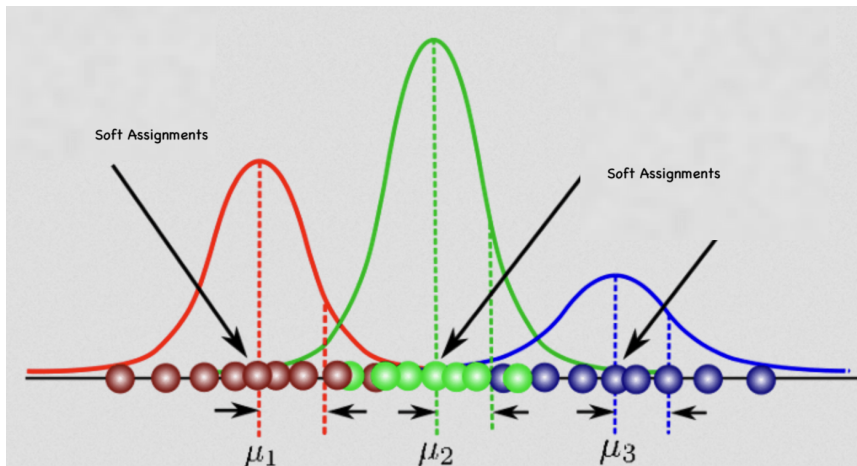
We could iterate till we converge

Soft_ K -Means():

- 1: Initiate μ_1, \dots, μ_K
- 2: **while** μ_1, \dots, μ_K changing **do**
- 3: Set $\{r_{n,k}\} \leftarrow \text{Soft_Cluster_Assignment}(\mu_1, \dots, \mu_K)$
- 4: Update $\mu_1, \dots, \mu_K \leftarrow \text{Centroid_Update}(\{r_{n,k}\})$
- 5: **end while**
- 6: Return μ_1, \dots, μ_K

Gaussian Prior

K-means assumes Gaussian distributed data!



Gaussian Prior

We understand it better if we discuss

Density Learning

This is what we do next!

Further Read

- MacKay
 - ↳ Chapter 20
- Bishop
 - ↳ Chapter 9

Soft K -means

K -means