

9.2 Non-linear Clustering

Machine Learning 1: Foundations

Marius Kloft (TUK)

- 1 Linear Clustering
- 2 Non-linear Clustering
- 3 Hierarchical Clustering

Kernel k -means clustering

The k -means clustering algorithm can be “kernelized”.

Kernel k -means clustering

The k -means clustering algorithm can be “kernelized”.

```
1: function KMEANS(parameter  $k$ , inputs  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ )
2:   initialize cluster centers  $\mathbf{c}_1, \dots, \mathbf{c}_k$ 
3:   repeat
4:     for  $i = 1 : n$  do
5:       label the input  $\mathbf{x}_i$  as belonging to the nearest cluster,
           
$$y_i := \arg \min_{j=1, \dots, k} \left\| \mathbf{x}_i - \mathbf{c}_j \right\|^2$$

6:     end for
7:     for  $j = 1 : k$  do
8:       compute cluster center  $\mathbf{c}_j$  as the mean of all inputs of the  $j$ th cluster,
           
$$\mathbf{c}_j := \text{mean}(\{\mathbf{x}_i : y_i = j\})$$

9:     end for
10:   until convergence criterion is met
11:   return cluster centers  $c_1, \dots, c_k$ 
12: end function
```

Kernel k -means clustering

The k -means clustering algorithm can be “kernelized”.

```
1: function KMEANS(parameter  $k$ , inputs  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ )
2:   initialize cluster centers  $\mathbf{c}_1, \dots, \mathbf{c}_k$ 
3:   repeat
4:     for  $i = 1 : n$  do
5:       label the input  $\mathbf{x}_i$  as belonging to the nearest cluster,
           
$$y_i := \arg \min_{j=1, \dots, k} \|\mathbf{x}_i - \mathbf{c}_j\|^2$$

6:     end for
7:     for  $j = 1 : k$  do
8:       compute cluster center  $\mathbf{c}_j$  as the mean of all inputs of the  $j$ th cluster,
           
$$\mathbf{c}_j := \text{mean}(\{\mathbf{x}_i : y_i = j\})$$

9:     end for
10:    until convergence criterion is met
11:  return cluster centers  $\mathbf{c}_1, \dots, \mathbf{c}_k$ 
12: end function
```

► Never compute mean explicitly

Kernel k -means clustering

The k -means clustering algorithm can be “kernelized”.

```
1: function KMEANS(parameter  $k$ , inputs  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ )
2:   initialize cluster centers  $\mathbf{c}_1, \dots, \mathbf{c}_k$ 
3:   repeat
4:     for  $i = 1 : n$  do
5:       label the input  $\mathbf{x}_i$  as belonging to the nearest cluster,
        $y_i := \arg \min_{j=1, \dots, k} \|\mathbf{x}_i - \mathbf{c}_j\|^2$ 
6:     end for
7:     for  $j = 1 : k$  do
8:       compute cluster center  $\mathbf{c}_j$  as the mean of all inputs of the  $j$ th cluster,
        $\mathbf{c}_j := \text{mean}(\{\mathbf{x}_i : y_i = j\})$ 
9:     end for
10:  until convergence criterion is met
11:  return cluster centers  $\mathbf{c}_1, \dots, \mathbf{c}_k$ 
12: end function
```

▶ Never compute mean explicitly ✓

▶ Compute $\|\mathbf{x}_i - \mathbf{c}_j\|^2$ via kernel function ✓

Kernel k -means clustering

- ▶ Let $\phi : \mathbb{R}^d \rightarrow \mathcal{H}$ be a kernel feature map, corresponding to a kernel $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$

Kernel k -means clustering

- ▶ Let $\phi : \mathbb{R}^d \rightarrow \mathcal{H}$ be a kernel feature map, corresponding to a kernel $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$
- ▶ Let $I_j := \{i \in \{1, \dots, n\} : y_i = j\}$ for all $j = 1, \dots, k$

Kernel k -means clustering

- ▶ Let $\phi : \mathbb{R}^d \rightarrow \mathcal{H}$ be a kernel feature map, corresponding to a kernel $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$
- ▶ Let $I_j := \{i \in \{1, \dots, n\} : y_i = j\}$ for all $j = 1, \dots, k$
- ▶ Then $\mathbf{c}_j := \frac{1}{|I_j|} \sum_{i \in I_j} \phi(\mathbf{x}_i)$

Kernel k -means clustering

- ▶ Let $\phi : \mathbb{R}^d \rightarrow \mathcal{H}$ be a kernel feature map, corresponding to a kernel $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$
- ▶ Let $I_j := \{i \in \{1, \dots, n\} : y_i = j\}$ for all $j = 1, \dots, k$
- ▶ Then $\mathbf{c}_j := \frac{1}{|I_j|} \sum_{i \in I_j} \phi(\mathbf{x}_i) \in \mathcal{H}$ (high dimensional!)

Kernel k -means clustering

- ▶ Let $\phi : \mathbb{R}^d \rightarrow \mathcal{H}$ be a kernel feature map, corresponding to a kernel $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$
- ▶ Let $I_j := \{i \in \{1, \dots, n\} : y_i = j\}$ for all $j = 1, \dots, k$
- ▶ Then $\mathbf{c}_j := \frac{1}{|I_j|} \sum_{i \in I_j} \phi(\mathbf{x}_i) \in \mathcal{H}$ (high dimensional!)
- ▶ However, $\|\phi(\mathbf{x}_i) - \mathbf{c}_j\|^2 = \underbrace{\|\phi(\mathbf{x}_i)\|^2}_{=k(\mathbf{x}_i, \mathbf{x}_i)} - 2 \langle \phi(\mathbf{x}_i), \mathbf{c}_j \rangle + \|\mathbf{c}_j\|^2$

Kernel k -means clustering

- ▶ Let $\phi : \mathbb{R}^d \rightarrow \mathcal{H}$ be a kernel feature map, corresponding to a kernel $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$
- ▶ Let $I_j := \{i \in \{1, \dots, n\} : y_i = j\}$ for all $j = 1, \dots, k$

▶ Then $\mathbf{c}_j := \frac{1}{|I_j|} \sum_{i \in I_j} \phi(\mathbf{x}_i) \in \mathcal{H}$ (high dimensional!)

▶ However, $\|\phi(\mathbf{x}_i) - \mathbf{c}_j\|^2 = \underbrace{\|\phi(\mathbf{x}_i)\|^2}_{=k(\mathbf{x}_i, \mathbf{x}_i)} - 2 \langle \phi(\mathbf{x}_i), \mathbf{c}_j \rangle + \|\mathbf{c}_j\|^2$

can be completely expressed in terms of kernel functions:

▶ $\|\mathbf{c}_j\|^2 = \frac{1}{|I_j|^2} \sum_{i, i' \in I_j} k(\mathbf{x}_i, \mathbf{x}_{i'})$ and

$\langle \phi(\mathbf{x}_i), \mathbf{c}_j \rangle = \frac{1}{|I_j|} \sum_{i' \in I_j} k(\mathbf{x}_i, \mathbf{x}_{i'})$

□

Kernel k -means clustering

- ▶ Let $\phi : \mathbb{R}^d \rightarrow \mathcal{H}$ be a kernel feature map, corresponding to a kernel $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$
- ▶ Let $I_j := \{i \in \{1, \dots, n\} : y_i = j\}$ for all $j = 1, \dots, k$
- ▶ Then $\mathbf{c}_j := \frac{1}{|I_j|} \sum_{i \in I_j} \phi(\mathbf{x}_i) \in \mathcal{H}$ (high dimensional!)
- ▶ However, $\|\phi(\mathbf{x}_i) - \mathbf{c}_j\|^2 = \underbrace{\|\phi(\mathbf{x}_i)\|^2}_{=k(\mathbf{x}_i, \mathbf{x}_i)} - 2 \langle \phi(\mathbf{x}_i), \mathbf{c}_j \rangle + \|\mathbf{c}_j\|^2$

can be completely expressed in terms of kernel functions:

- ▶ $\|\mathbf{c}_j\|^2 = \frac{1}{|I_j|^2} \sum_{i, i' \in I_j} k(\mathbf{x}_i, \mathbf{x}_{i'})$ and
 $\langle \phi(\mathbf{x}_i), \mathbf{c}_j \rangle = \frac{1}{|I_j|} \sum_{i' \in I_j} k(\mathbf{x}_i, \mathbf{x}_{i'})$



But what if we want to cluster images?

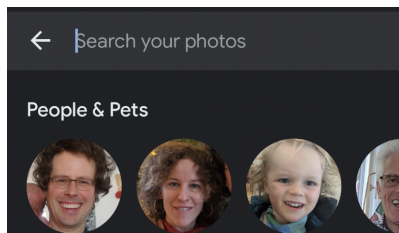
Example: Search Photo Collection For People

Given:

- ▶ Photo collection

Aim:

- ▶ Search for people



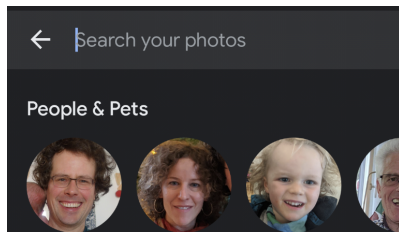
Example: Search Photo Collection For People

Given:

- ▶ Photo collection

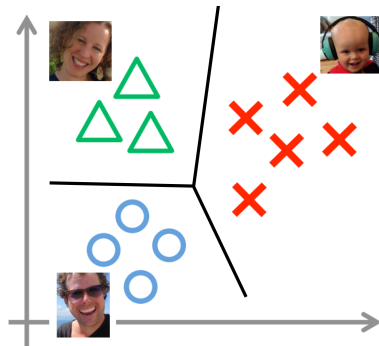
Aim:

- ▶ Search for people



Solution:

- ▶ Automatically crop out portraits
- ▶ **Cluster the portraits**



Deep Clustering

Can we make k -means deep?

Deep Clustering

Can we make k -means deep?

One can show:

- ▶ k -means minimizes the objective function

$$\min_{\mathbf{c}=(\mathbf{c}_1,\dots,\mathbf{c}_k)\in\mathbb{R}^{d\times k}} \sum_{i=1}^n \min_{j\in\{1,\dots,k\}} \|\mathbf{x}_i - \mathbf{c}_j\|^2$$

Deep Clustering

Can we make k -means deep?

One can show:

- ▶ k -means minimizes the objective function

$$\min_{c=(\mathbf{c}_1, \dots, \mathbf{c}_k) \in \mathbb{R}^{d \times k}} \sum_{i=1}^n \min_{j \in \{1, \dots, k\}} \|\mathbf{x}_i - \mathbf{c}_j\|^2$$

Analog to regression and classification, we could “deepify” k -means as follows:

Deep k -means

$$\min_{W, c} \frac{1}{2} \sum_{l=1}^L \|W^{(l)}\|_{\text{Fro}}^2 + \sum_{i=1}^n \min_{j \in \{1, \dots, k\}} \|\phi_W(\mathbf{x}_i) - \mathbf{c}_j\|^2$$

Deep Clustering

Can we make k -means deep?

One can show:

- ▶ k -means minimizes the objective function

$$\min_{c=(\mathbf{c}_1, \dots, \mathbf{c}_k) \in \mathbb{R}^{d \times k}} \sum_{i=1}^n \min_{j \in \{1, \dots, k\}} \|\mathbf{x}_i - \mathbf{c}_j\|^2$$

Analog to regression and classification, we could “deepify” k -means as follows:

Deep k -means

$$\min_{W, c} \frac{1}{2} \sum_{l=1}^L \|W^{(l)}\|_{\text{Fro}}^2 + \sum_{i=1}^n \min_{j \in \{1, \dots, k\}} \|\phi_W(\mathbf{x}_i) - \mathbf{c}_j\|^2$$

But this method does not always work so well

- ▶ ongoing research topic to get this working

Another Idea: Transfer Learning

We learned this trick already in the class on overfitting:

Another Idea: Transfer Learning

We learned this trick already in the class on overfitting:

- 1 Pre-training: Download an ANN with parameters W from the web that was pre-trained on huge amounts of images

Another Idea: Transfer Learning

We learned this trick already in the class on overfitting:

- 1 Pre-training: Download an ANN with parameters W from the web that was pre-trained on huge amounts of images
- 2 Deep representation of data: For each input \mathbf{x}_i , compute its vector of activations in the last hidden layer $\phi_W(\mathbf{x}_i)$

Another Idea: Transfer Learning

We learned this trick already in the class on overfitting:

- 1 Pre-training: Download an ANN with parameters W from the web that was pre-trained on huge amounts of images
- 2 Deep representation of data: For each input \mathbf{x}_i , compute its vector of activations in the last hidden layer $\phi_W(\mathbf{x}_i)$
- 3 Clustering: Run k -means on $\phi_W(\mathbf{x}_1), \dots, \phi_W(\mathbf{x}_n)$