

## SPR Day 12: KMeans

Read: PRML textbook 9.1 K-Means

Goals: (1) Understand problem setup  
cost function  
optimization problem

(2) Understand coord. descent approach

(3) Look ahead

- scaling to large data
- handling other distances
- towards probabilistic approach

# K-Means Clustering

## Problem Statement

Given:  $N$  observed data vectors, each  $D$ -dim.

$$\{x_n\}_{n=1}^N \quad \text{where } x_n \in \mathbb{R}^D \quad \forall n$$

Desired number of clusters  $K$  (integer,  $K > 0$ )

Goal: Determine locations of  $K$  cluster centers

$$\{\mu_k\}_{k=1}^K \quad \text{where } \mu_k \in \mathbb{R}^D \quad \forall k$$

such that across all datapoints, the distance between each point and its nearest center is minimized

We want a "low cost" clustering!

Cost function:

$$\begin{aligned} \text{cost}(X, \mu) &= \sum_{n=1}^N \min_{k \in \{1, 2, \dots, K\}} \text{dist}(x_n, \mu_k) \\ &= \sum_{n=1}^N \min_k \|x_n - \mu_k\|_2^2 \quad \text{Euclidean distance} \\ &= \sum_{n=1}^N \min_k \sum_d (x_{nd} - \mu_{kd})^2 \end{aligned}$$

Alternative: Can write in a way to make solving easier

Let  $r_n$  be a length  $K$  "one hot" vector

indicating the closest cluster to example  $x_n$

can write

$$\text{cost}(X, r, \mu) = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|_2^2$$

# KMeans Coordinate Descent

2 e.g. if  $K=4$ ,  $r_n$  could be:

	1	2	3	4
	0	0	0	1
	0	0	1	0
	0	1	0	0
	1	0	0	0

Goal: Find  $r = \{r_n\}_{n=1}^N$ ,  $r_n \in \text{onehot}(K)$   
and  $\mu = \{\mu_k\}_{k=1}^K$ ,  $\mu_k \in \mathbb{R}^D$

that minimize cost

$$\text{cost}(x, r, \mu) = J(x, r, \mu) = \sum_n \sum_k r_{nk} \|x_n - \mu_k\|_2^2$$

Optimization Problem:

$$\min J(x, r, \mu)$$

$$\sum_{n=1}^N r_n \in \text{onehot}(K)$$

$$\mu \in \mathbb{R}^{K \times D}$$

How to solve? What algorithm?

Can we just do gradient descent? No!  $r$  is not continuous

Look at coordinate descent.

Iteration  $t$  has two steps.

Step 1: Fix  $\mu = \mu^{t-1}$ . Find optimal  $r$  given  $\mu^{t-1}$ .

$$r^t = \min_{\sum_{n=1}^N r_n} J(x, r, \mu^{t-1})$$

$\mu^{t-1}$  fixed

Step 2: Fix  $r = r^t$ . Find optimal  $\mu$  given  $r^t$ .

$$\mu^t = \min_{\mu \in \mathbb{R}^{K \times D}} J(x, r^t, \mu)$$

# Solving Assignment Step

3

Step 1:

$$\min_{\{\tau_n\}_{n=1}^N} \sum_n \underbrace{\sum_k \tau_{nk} \|x_n - \mu_k\|_2^2}_{f^n(\tau_n)}$$

Because objective is a sum over  $N$  separate terms  
this looks like

$$\min_{\{\tau_n\}_{n=1}^N} f^1(\tau_1) + f^2(\tau_2) + \dots + f^N(\tau_N)$$

and we can solve each term separately?

$$\text{for each } n: \min_{\tau_n \in \text{onehot}(K)} \sum_k \tau_{nk} \|x_n - \mu_k\|_2^2$$

We can show that optimal update is:

$$\tau_{nk} = \begin{cases} 1 & \text{if } k = \operatorname{argmin}_{k \in \{1, 2, \dots, K\}} \|x_n - \mu_k\|_2^2 \\ 0 & \end{cases}$$

Proof by contradiction. Suppose we set  $\tau_{nj} = 1$ , but  
there exists  $k \neq j$  where  $\|x_n - \mu_k\| < \|x_n - \mu_j\|$ .

then we can improve cost, and thus  $\tau_{nj}$  not optimal  
(by setting  $\tau_{nk} = 1$  instead)

# Solving Mean Update Step

4

$$\text{Step 2: } \min_{\mu \in \mathbb{R}^{K \times D}} \sum_{k=1}^K \sum_{n=1}^N r_{nk} \|x_n - \mu_k\|_2^2$$

Again, can view as  $K$  separate functions each involving only  $\mu_k \in \mathbb{R}^D$

$$\text{For each } k: \min_{\mu \in \mathbb{R}^D} \sum_n r_{nk} \|x_n - \mu_k\|_2^2 = \underbrace{\sum_n r_{nk} (x_n - \mu_k)^T (x_n - \mu_k)}_{\mathcal{L}(\mu_k)}$$

Solve by: taking derivs, set to zero, find optimal  $\mu_k$   
const so  $\phi$

$$0 = \nabla_{\mu_k} \mathcal{L}(\mu_k) = \nabla_{\mu_k} \sum_n r_{nk} [x_n^T x_n - 2x_n^T \mu_k + \mu_k^T \mu_k]$$

$$0 = \text{"} = \sum_n r_{nk} [\phi - \nabla_{\mu_k} (2x_n^T \mu_k) + \nabla_{\mu_k} (\mu_k^T \mu_k)]$$

$$0 = \sum_n r_{nk} [-2x_n + 2\mu_k]$$

$$\sum_n r_{nk} x_n = \left( \sum_n r_{nk} \right) \mu_k$$

$$\mu_k = \frac{\sum_n r_{nk} x_n}{\sum_n r_{nk}}$$

Edge case:  
avoid divide by zero  
when  $\sum_n r_{nk} = \phi$

Set location of cluster  $k$  to the mean of all examples assigned to  $k$

# Iterative algorithm of K-Means

5

## Coordinate descent algo.

Input:  $\{x_n\}_{n=1}^N$  (examples)  
 $K$  (num clusters)  
 $\mu^0 = \{\mu_k\}_{k=1}^K$  (initial cluster locations)

Procedure:

Repeat for iters  $t=1, 2, \dots, T$

Step 1 [ for each  $n$  in  $1, 2, \dots, N$ :  
Update  $r_{nk}^t \leftarrow \begin{cases} 1 & \text{if } k = \text{argmin}_k \|x_n - \mu_k\|_2^2 \\ 0 & \text{o.w.} \end{cases}$

Step 2 [ for each  $k$  in  $1, 2, \dots, K$ :  
Update  $\mu_k^t \leftarrow \frac{\sum_n r_{nk}^t x_n}{\sum_n r_{nk}^t}$

Big Idea: Each step guaranteed to lower cost

$$J(x, r^0, \mu^0) \geq J(x, r^1, \mu^1) \geq J(x, r^2, \mu^2) \geq \dots \geq J(x, r^T, \mu^T)$$

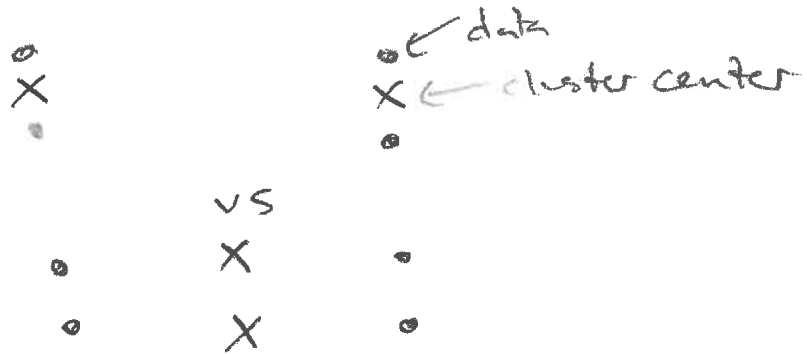
after iter 1                      after iter 1                      after iter 2                      after iter T  
step 1                                  step 2                                  step 1                      ...                      step T

Will eventually converge (stop changing:  $r^{t+1} = r^t$  for some  $t$ ).  
However, will it yield a global optimum? Not always

# Demos of KMeans

See slides

- illustration that local optima exist



# Scaling to Large Datasets

Can parallelize many steps

e.g. solve for  $r_1$   
 $r_2$  each on separate  
 $\vdots$   
 $r_N$  machines

Can process a few examples at a time

option 1: for each minibatch  $x^b$ :  
 E) find optimal assignments  $r^b$  given  $\mu^{t-1}$   
 SGD approach M) take gradient step  

$$\mu^t \leftarrow \mu^{t-1} - \epsilon \nabla J(x^b, r^b, \mu^t)$$

option 2:  
 incremental  
 or  
 memoized  
 approach

Store assignment statistics  
 for each batch  $b$

$$N_{bk} = \sum_{i \in b} r_{ik}$$

$$S_{bk} = \sum_{i \in b} r_{ik} x_i$$

for each minibatch  $x^b$ :

update local statistics

$$N_{bk}^{\text{new}}$$

$$S_{bk}^{\text{new}}$$

increment global stats

$$N_{ok} = N_{ok}^{\text{old}} + N_{ok}^{\text{new}}$$

update

$$\mu_k = S_{ok} / N_{ok}$$