

CS 181

Clustering

March 2021

New task  $\underline{x} \rightarrow \text{summary}(\underline{x})$

Motivations understanding communication  
pre-training for supervised  
organizing

Data  $\{\underline{x}_1, \dots, \underline{x}_N\}$  often  $\underline{x}_n \in \mathbb{R}^D$   
could be  $\underline{x}_n \in \{0, 1\}^D$

Number of clusters  $K$  (may not be given)

Output: assignment of each example  
to a cluster

$\underline{z}_n$  1-hot  
"assignment"  $\begin{cases} z_{nk} = 0 & \text{if } \underline{x}_n \text{ is not in cluster } k \\ z_{nk} = 1 & \text{if } \underline{x}_n \text{ is in cluster } k \end{cases}$

What is a good clustering?

Idea: examples to be more similar to  
examples in same cluster than to  
examples in other clusters

Need measure of similarity:

eg.,  $d(x, x') = \|x - x'\|_2$  (L2)

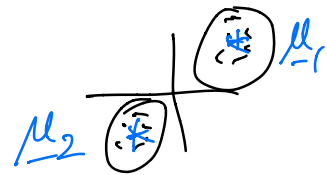
edit distance (strings)

Hamming distance (bit vectors)

⋮

**Alg 1**

K-means clustering



- Define "prototype"  $\mu_k \in \mathbb{R}^D$  per cluster

- Goal assign  $x_n$  + find  $\{ \mu_1, \dots, \mu_k \}$  s.t.

$$\min_{\{ \mu \}, \{ z \}} \sum_n \sum_k z_{nk} \|x_n - \mu_k\|_2^2$$

Note: Non-convex  
NP-hard

$$z_{nk} = \begin{cases} 1 & \text{if assigned to } k \\ 0 & \text{o.w.} \end{cases}$$

Addendum  
Objective criterion should be  $\|x_n - \mu_k\|_2^2$  not  $\|x_n - \mu_k\|_2$  as written in lecture

## Lloyd's algorithm

- 1) Randomly initialize prototypes  $\{\mu_k\}$
- 2) Repeat:

Step 1 Assign each example to its  
closest prototype

$$\arg \min_k \|x_n - \mu_k\|_2 \quad \left. \vphantom{\arg \min_k} \right\} \begin{array}{l} \text{for} \\ \text{each} \\ n \end{array}$$

Step 2 For each  $k$ , set  $\mu_k$  to  
the centroid (mean) of assigned  
examples

$$\mu_k := \frac{1}{N_k} \sum_n z_{nk} x_n$$

$$\text{where } N_k = \sum_n z_{nk}$$

# examples assigned  
to  $k$

Typical restart this multiple times  
& take the "best" solution

# Understanding Lloyd's algorithm

Recall objective

$$\min_{\{\mu\}, \{z\}} \sum_n \sum_k z_{nk} \|\underline{x}_n - \mu_k\|_2^2 \quad (*)$$

"Coordinate descent", alternating  $\{\mu\}$  +  $\{z\}$  updates

(Repeat)

Step 1: Fixing  $\{\mu\}$ , minimize the loss (\*)  
by assigning each  $\underline{x}_n$  to closest  
prototype

Step 2: Fixing  $\{z\}$ , minimize loss (\*)  
by choosing prototypes  $\{\mu\}$

$$\text{For } k: \quad L(\mu_k) = \sum_n z_{nk} (\underline{x}_n - \mu_k)^T (\underline{x}_n - \mu_k)$$

$$\frac{\partial L(\mu_k)}{\partial \mu_k} = -2 \sum_n z_{nk} (\underline{x}_n - \mu_k) = 0$$

$$\Leftrightarrow \mu_k = \frac{1}{N_k} \sum_n z_{nk} \underline{x}_n \quad \begin{array}{l} N_k \# \\ \text{examples} \\ \text{assigned to} \\ k \end{array}$$

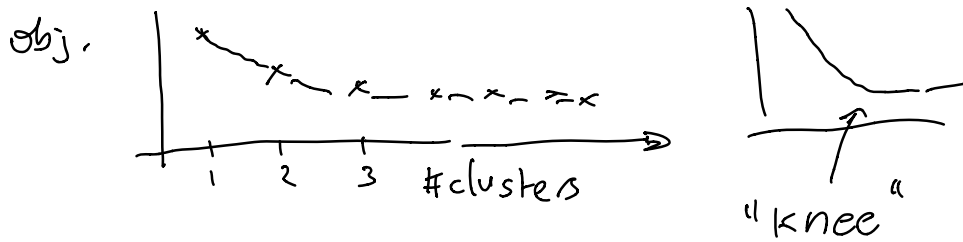
Again minimizing!

## Considerations

1) How many clusters?

smaller, better interpretation

larger, better extraction of concepts



2) Parametric method

$D \times K$  parameters (prototypes)

3) Inflexible (linear decision boundary)

4) Fast! Assignment  $z_{uk}$  step

parallelizes by  $x_n$ ; prototype  
step parallelizes by  $\mu_k$

## Variations

1) k-means++ to initialize

2) K-medoids

replace  $\mu_k$  step with:

$\mu_k :=$  example assigned to  $k$  that minimizes total distance to the other examples in the cluster

$$\arg \min_{\substack{x_n \text{ s.t.} \\ z_{nk}=1}} \sum_{n'} z_{n'k} \|x_n - x_{n'}\|_2^2$$

3) L1 norm in place of L2 norm

$\mu_k :=$  median of the points assigned to cluster  $k$

# Alternate Hierarchical Agglomerative Clustering (HAC)

Data  $\{x_n\}$

Distance  $d(G, G')$

distance between groups

HAC:

- 1) Every example starts in its cluster
- 2) While the # clusters  $> 1$ ,  
Merge the two "closest" clusters

Notes

Forms a hierarchy of clusters

No need to specify  $K$  (# clusters)

Deterministic

Need two concepts

- 1)  $d(x, x')$  distance between points
- 2) "linkage" function, min, max, average, centroid  
 $\hookrightarrow$  Get  $d(G, G')$

## Comments

1) Average, centroid compromises between  $\mu_i$  +  $\mu_{j+1}$

2) Non-parametric (distance-based)

▣ Arbitrary cluster shapes

3) Scales as  $O(n^2)$  } k-means  
▣ Pairwise calculation }  $O(n)$

4) Can sometimes lead to overfitting

↙  
k-means  
compares  $x_n$  to  
prototype  $\mu_k$   
+ prototypes are averages

↘  
HAC  
compares pairs of  
examples  
(no averaging effect)



## Concept check

Data       $\underbrace{0 \dots 0}_l$        $\underbrace{\text{Random } 0/1}_{D-l}$       " $\underline{x}_0$ "  
                  $\underbrace{1 \dots 1}_l$        $\underbrace{\text{Random } 0/1}_{D-l}$       " $\underline{x}_1$ "

### Problem

Given  $\underline{x}, \underline{x}'$  in "0" cluster +  $\underline{z}$  in "1" cluster, the probability  $\underline{x}$  is closer to  $\underline{x}'$  than to  $\underline{z}$  goes to  $1/2$  as  $D \rightarrow \infty$   
(fixing  $l$ )

⊗ Noise comes to dominate } + KAC fails

⊗ K-means works OK here

$$\underline{\mu}_0 = (0, \dots, 0, \frac{1}{2}, \frac{1}{2}, \dots, \frac{1}{2})$$

$$\underline{\mu}_1 = (1, \dots, 1, \frac{1}{2}, \frac{1}{2}, \dots, \frac{1}{2})$$

With these prototypes we can correctly cluster data  
+ K-means converges

Why the problem with noise in KAC?

The "curse of dimensionality"

Suppose 1000 examples are uniform randomly distributed in a  $D$ -dimensional unit hypercube

Consider the squared-distance between random examples  $\underline{x}, \underline{z}$ :

$$\|\underline{x} - \underline{z}\|_2^2 = \sum_{j=1}^D (x_j - z_j)^2$$

} Sum of  $D$ , i.i.d. random variables

By central limit theorem, concentrates around  $D \times E[(x_j - z_j)^2]$  where  $x_j, z_j \sim U(0,1)$ . Distance is sqrt of this.

(compare to min distance 0, max distance  $\sqrt{D}$ )

Distr. of inter-example distances  $\rightarrow$  increasing concentration

