

CS 181

Clustering

March 2021

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

What is a good clustering?

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

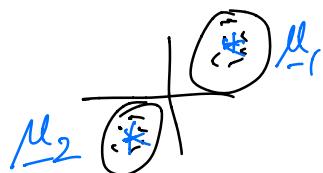
Need measure of similarity:

$$\text{e.g., } d(\underline{x}, \underline{x}') = \|\underline{x} - \underline{x}'\|_2 \quad (\ell^2)$$

edit distance (strings)

Hammng distance (bit vectors)

:



Alg 1

K-means clustering

- Define "prototype" $\underline{\mu}_k \in \mathbb{R}^D$ per cluster
- Goal: assign \underline{z}_n & find $\{\underline{\mu}_1, \dots, \underline{\mu}_K\}$
s.t.

$$\min_{\{\underline{\mu}\} \{\underline{z}\}} \sum_n \sum_k z_{nk} \|\underline{x}_n - \underline{\mu}_k\|_2$$

Note: Non-convex
NP-hard

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

Lloyd's algorithm

1) Randomly initialize prototypes $\{\mu_k\}$

2) Repeat:

Step 1 Assign each example to its
closest prototype } for
 $\arg \min_k \|x_n - \mu_k\|_2$ } each
n

Step 2 For each k, set μ_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

$$\min_{\{\mu\} \in \{\Xi\}} \sum_n \sum_k z_{nk} \|x_n - \mu_k\|_2 \quad (*)$$

"Coordinate descent", alternately $\{\mu\}$ + $\{\Xi\}$ update

(Repeat)

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

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

$$\text{for } k: L(\mu_k) = \sum_n z_{nk} (x_n - \mu_k)^T (x_n - \mu_k)$$

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

$$\Leftrightarrow \mu_k = \frac{1}{N_k} \sum_n z_{nk} x_n \quad N_k \# \text{ examples assigned to } k$$

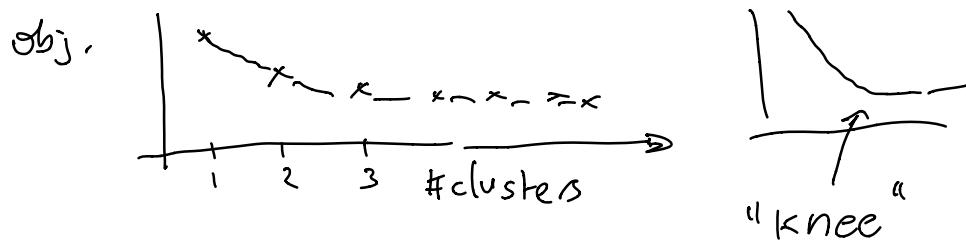
Again minimize!

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 μ_k

Variations

1) k-means ++ to initialize

2) K-medoids

replace μ_k step with:

$\underline{\mu}_k$:= example assigned to
k that minimizes total distance
to the other examples in
the cluster

$$\arg \min_{\underline{x}_n \text{ s.t.}} \sum_{n'} z_{n'k} \|\underline{x}_n - \underline{x}_{n'}\|_2^2$$
$$z_{nk} = 1$$

3) L1 norm in place of L2 norm

$\underline{\mu}_k$:= median of the
points assigned to cluster k

Alternate Hierarchical Agglomerative Clustering (HAC)

Data $\{x_i\}$ 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,
 \rightarrow Get $d(G, G')$ average, centroid

Comments

- 1) Average, centroid compromises between min + max
- 2) Non-parametric (instance-based)
 - ☒ Arbitrary cluster shapes
- 3) Scales as $O(n^2)$
 - ☒ Pairwise calculation
- 4) Can sometimes lead to overfitting

K-means

compares x_n to prototype μ_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)

④ Noise comes to dominate } + KAC fails

④ K-mean 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 HAC?

The "curse of dimensionality"

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

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

$$\|\mathbf{x}_j - \mathbf{z}_j\|_2^2 = \sum_{i=1}^D (x_{ji} - z_{ji})^2 \quad \left. \begin{array}{l} \text{sum of} \\ D, \text{ i.i.d.} \\ \text{random} \\ \text{variables} \end{array} \right.$$

By central limit theorem,

(concentrates around $D \times E[(x_{ji} - z_{ji})^2]$)

where $x_{ji}, z_{ji} \sim U(0, 1)$, Distance is sqrt of this.

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

④ Distr. of inter-example distances
→ increasing concentration

