

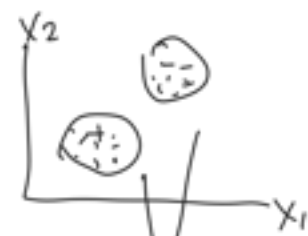
so far: $x \rightarrow y$ (prediction)
 New task $x \rightarrow \text{summary}(x)$

Why??

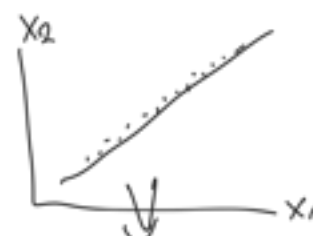
summarization: news/synthesis of info & organization
 visualization
 compression (image/video)
 latent causes [hypotheses]

What

unsupervised learning / representation learning



clustering
 \Leftrightarrow classification



manifold learning
 embedding
 \Leftrightarrow regression

discrete vs. continuous
supervised vs. unsupervised
~~non-prob~~ non-prob vs. probabilistic

How evaluate.

reconstruction error : $\mathbb{E} \mathcal{L}(x, \underbrace{\text{decode}(\underbrace{\text{encode}(x)}_{\text{summary}(x)})}_{\text{recovery of } x})$

but: we can also create very method-specific criteria
 (ex. clustering: stability, relative cluster sizes, etc.)

How to do it?
 (clustering)



how many?



what shape?

Set-up: Data: $x_1 \dots x_n$

Similarity/Distance $d(x, x')$ e.g. $d(x, x') = \|x - x'\|$

goal is to group similar points together

for now, suppose K groups, z_{nk} indicate group assignment

z_n to be a binary vec. of length K
 $z_{nk} = 1$ for k is the group of datum x_n $[0 1 0 0 0]$

Specific Approach #1: K-means

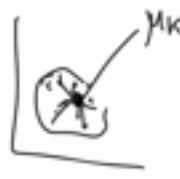
goal: define μ_k as the cluster "prototypes" $\{ \mu_1 \dots \mu_K \}$

assign z_{nk} , find μ_k s.t.

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

$\underbrace{\sum_k z_{nk}}_{\text{picking out which cluster } x_n \text{ was assigned to}}$
 $\underbrace{\|x_n - \mu_k\|_2^2}_{\text{loss}}$

encoding of x_n
 decoding of x_n



unfortunately, this is NP-hard!

can we find a local optima?

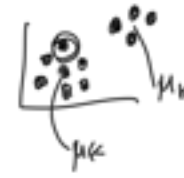
- iterative block update / coordinate ascent
- [K-means / Lloyd's Alg]

→ start by randomly assigning z_n 's

→ loop:

update $\{\mu\}$ given $\{z\}$:
 for each k , $\mu_k = \frac{1}{N_k} \sum_n z_{nk} x_n$ } set μ_k to the mean of x_n 's assigned to cluster k

update $\{z\}$ given $\{\mu\}$:
 for each n ,
 assign x_n 's to $\arg\min_k \|x_n - \mu_k\|$



Computation:

- nice, parallelized / distributed

Theory:

- nonconvex → need to do many random restarts / K-means++
 - but, procedure always is improving the obj.
- by assigning z_n to best $\|x_n - \mu_k\|$, loss has to go down (or stay constant)

$$\begin{aligned} \rightarrow \mathcal{L}(\mu_k) &= \sum_n z_{nk} (x_n - \mu_k)^T (x_n - \mu_k) \\ \frac{\partial \mathcal{L}}{\partial \mu_k} &= -2 \sum_n z_{nk} (x_n - \mu_k) = 0 \Rightarrow \mu_k = \frac{\sum_n z_{nk} x_n}{\sum_n z_{nk}} \end{aligned}$$

↳ implementation note: you can check this!

Concepts:

- what kind of boundary?

- what about different $d(x, x')$?

μ_k / $\mu_{k'}$

linear boundaries → may need to $x \rightarrow p(x)$, but K-means will only do linear boundaries

- ϵ assignment is still easy $\min_{k \in \{1, \dots, K\}} \|x - \mu_k\|$
 → μ_k may get hard
 - How many clusters?
 problem-specific
- K. medoids which forces μ_k to be a $\{x_n\}$

Specific Approach #2: deterministic, nonparametric (tree), non-linear boundaries
 Hierarchical Agglomerative Clustering



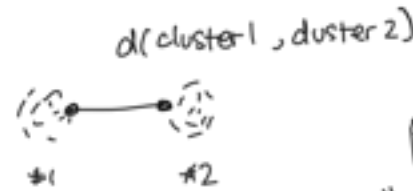
Alg: $\{x_n\}$, $d(x, x')$

- 1) everyone starts in their own cluster
- 2) while #clusters > 1 , merge closest clusters

Obvious question: what is the "closest" cluster? we have $d(x, x')$

"linkage"

- $\min_{\substack{n \in \text{cluster 1} \\ n' \in \text{cluster 2}}} d(x_n, x_{n'})$



"stringy" clusters

- $\max_{\substack{n \in C_1 \\ n' \in C_2}} d(x_n, x_{n'})$

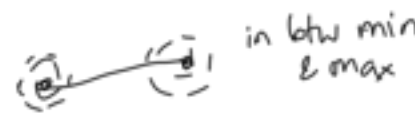


"round"



- average $\frac{1}{N_1 N_2} \sum_{n \in C_1} \sum_{n' \in C_2} d(x_n, x_{n'})$ in btw min & max

- centroid/ward $d(\mu_{\#1}, \mu_{\#2})$
 ↑ ↑
 mean mean
 cluster1 cluster2



Notes:

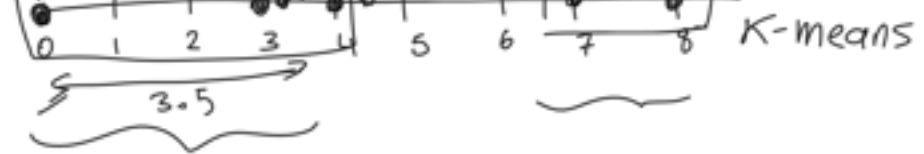
- very sensitive to choices of $d(x, x')$ & linkage
- high dims can get tricky w/ Euclidean $d(x, x')$ because all data will \approx equally far apart. [non-specific to HAC]
- obj is not as clear (2016, 2017 papers...)
 → cluster sizes
 → $\sum d(x, x')$ crossing a boundary



K-means, HAC: Concept Exercise



HAC: @ 2 clusters
 K-means: 2 clusters



one more example: μ large # of random [binary w.p. $\frac{1}{2}$].

$x_1:$

0	0	0	0
1	1	1	1

 $x_n:$

1	1	1	1
0	0	0	0

K-means: $0000 \frac{1}{2} \frac{1}{2} \frac{1}{2} \dots$
 $1111 \frac{1}{2} \frac{1}{2} \frac{1}{2} \dots$

HAC: for large μ , essentially random $d(x, x')$

Clustering