

Clustering

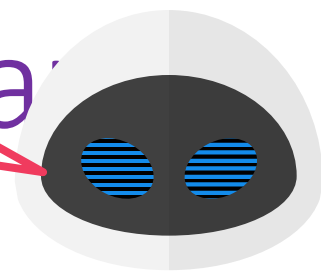
CS771: Introduction to Machine Learning

Purushottam Kar

How

Grouping similar objects together is the task of clustering

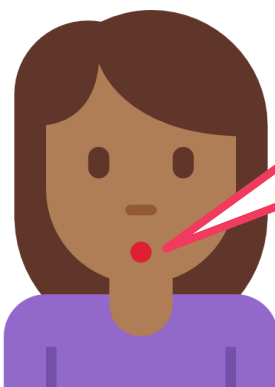
What



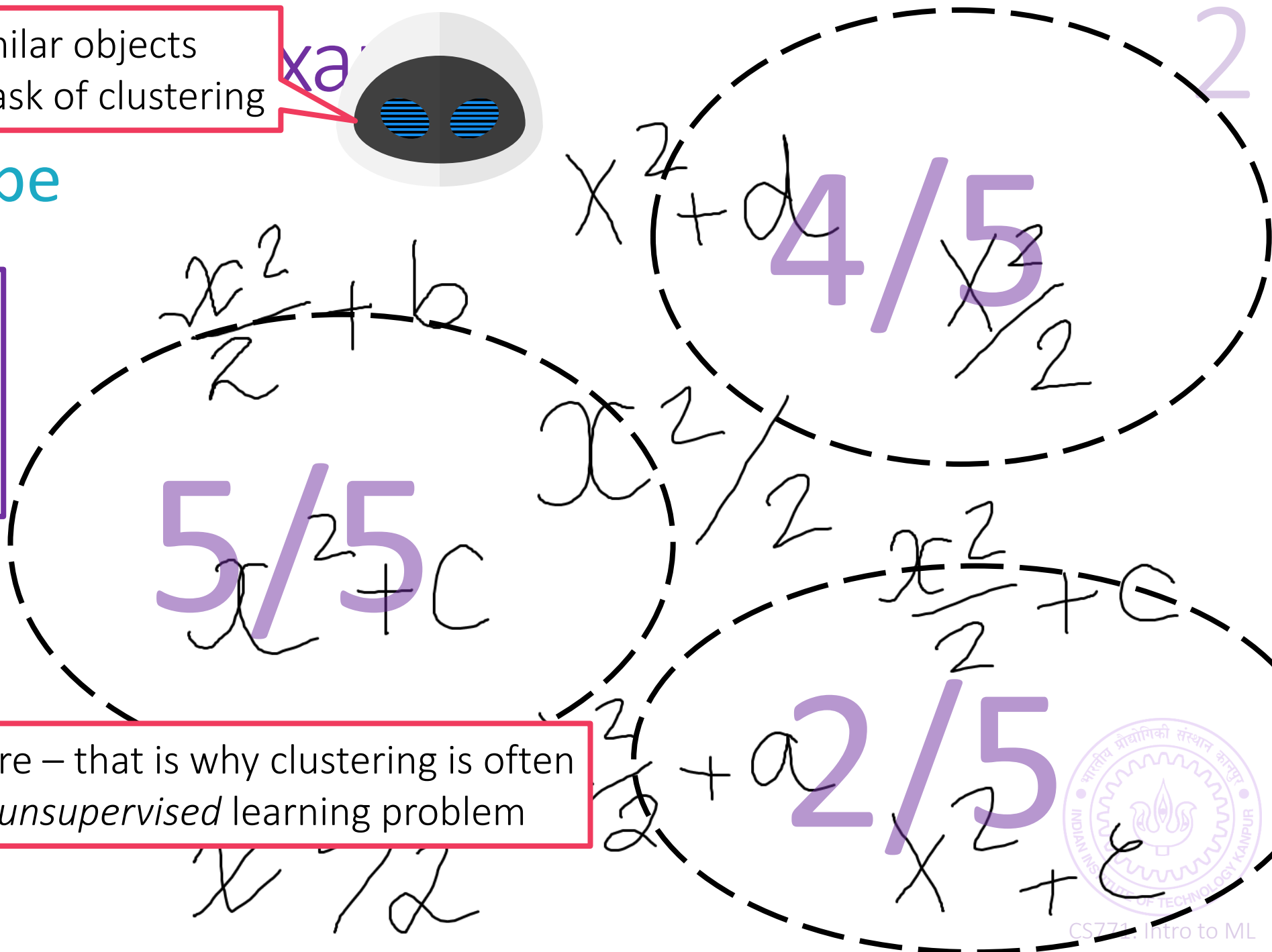
gradescope

Q1. $\int x = ?$

5 marks

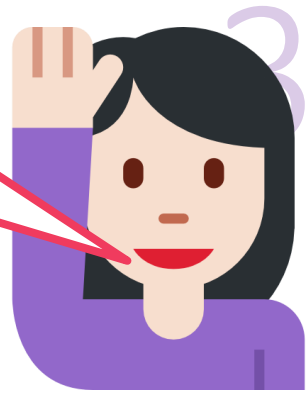


No labels here – that is why clustering is often called an *unsupervised* learning problem



Clustering

The technical term used in books/papers is centroid



Given a set S of n data points $\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^n \in \mathbb{R}^d$

Split this set into C disjoint clusters S_1, \dots, S_C i.e.

Assign every data point i to one of the subsets, say $z_i \in [C]$ (note that every data point is assigned to exactly one cluster) so that

Data points assigned to the same subset are “similar” to each other, e.g.

If $z_i = z_j = c$ for some $c \in [C]$ then $\|\mathbf{x}^i - \mathbf{x}^j\|_2$ is small

The K-means problem asks this problem a bit differently

Split S into C clusters S_1, \dots, S_C and find a ^{centroid} ~~prototype~~ for each cluster i.e. $\boldsymbol{\mu}^c \in \mathbb{R}^d$ s.t. if \mathbf{x}^i is assigned to cluster c i.e. $z_i = c$, then $\|\mathbf{x}^i - \boldsymbol{\mu}^c\|_2^2$ is small i.e. \mathbf{x}^i is close to ~~prototype~~ _{centroid} of its cluster



K-means clustering

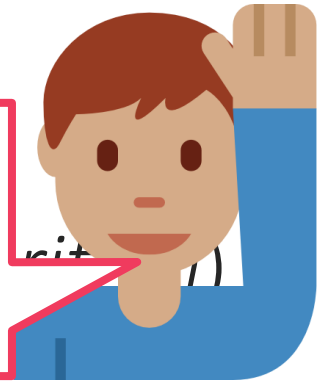
4

$$\min_{\{\mu^c \in \mathbb{R}^d\}, \{z_i \in [C]\}} \sum_{c=1}^C \sum_{i: z_i=c} \|\mathbf{x}^i - \mu^c\|_2^2$$

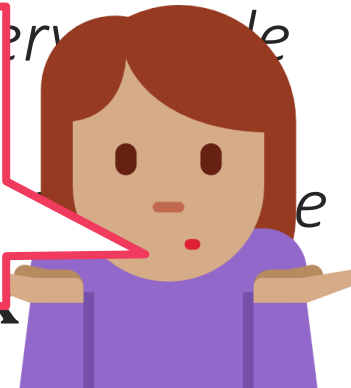
K-MEANS/LO

1. Initialize μ^c
2. For $i \in [n]$, update z_i using $\{\mu^c\}$
 1. Let $z_i = \arg \min_c \|\mathbf{x}^i - \mu^c\|_2^2$
3. Let $n_c = \#$ points assigned to c
4. Update $\mu^c = \frac{1}{n_c} \sum_{i: z_i=c} \mathbf{x}^i$
5. Repeat until convergence

The Lloyd's algorithm offers *monotonic progress*. It always lowers (or keeps the same) the value of $\sum_{c=1}^C \sum_{i: z_i=c} \|\mathbf{x}^i - \mu^c\|_2^2$ i.e. the objective function in each iteration. Can you show this?

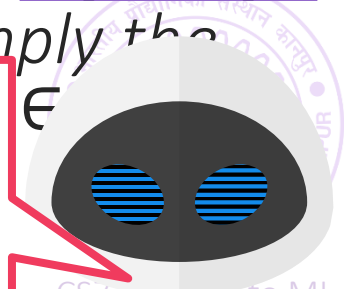


This looks a bit like coordinate minimization where we fix all but one coordinate and update that one coordinate to its optimal value



That optimal value of μ^c is simply the

True, coordinate minimization can be thought of as a special case of alternating optimization 😊



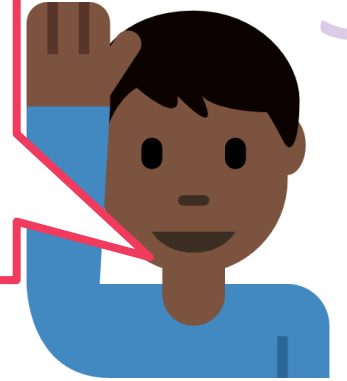
K-means++

Initializes k-means

Provable guarantee

Widely used in practice: especially beneficial if k is large

Note that a k-means++ always initializes centroids as actual data points. Also, no data point can be selected twice – if a data point \mathbf{x}^i gets selected once, then for all subsequent iterations, we will have $d_i = 0 = p_i$



5

K-MEANS++ INITIALIZER

1. Select first centroid randomly
 $\mu^1 = \mathbf{x}^i$, where $i \sim \text{UNIF}([n])$

2. For $j = 2, \dots, k$

1. For all $i \in [n]$, calculate $d_i = \min_{l \in 1, \dots, j-1} \|\mathbf{x}^i - \mu^l\|_2$

2. Set $\mu^j = \mathbf{x}^i$ where i is chosen with probability $p_i = \frac{d_i^2}{\sum_{s=1}^n d_s^2}$



Some applications of clustering

6

Can be used to make LwP a more powerful algorithm

Learn more than one prototype per class e.g. k prototypes by clustering data of each class into k clusters and using the centroids returned by the clustering algorithm as prototypes

A test point is assigned the class of its closest prototype

Note: this will increase training time, test time, and model size a bit

Seamlessly gives us the 1NN algorithm if we demand as many clusters (and hence as many centroids) as there are data points



Some applications of clustering

7

Identify *subpopulations* in data and improve ML performance

Example: have data for 1M customers but don't know age/gender

However, we suspect that age/gender significantly affects behaviour

Instead of running an ML algo (say SVM) on entire training data, first cluster training data and run ML algo separately on each cluster

If k clusters then k models will get learnt. For test data points, first find to which cluster they belong (using distance to centroid) and use that model

Increases model size and test time a bit but may increase accuracy too!

If we cluster these customers according to their onsite behaviour (which items did they view/like/buy), possible that we may accidentally discover gender/age groups within our data without knowing these details directly

Groups may not be perfectly clean but should improve ML performance



Some applications of clustering

Features for features!!

Reduce number of features (also called *dimensionality reduction*)

Example: have 1M features i.e. $\mathbf{x}^i \in \mathbb{R}^d$ for $d = 1M$ but we suspect many of these features are redundant

Example: synonyms in b

Can cluster features tog

Once we have \hat{d} feature clusters, say $C_1, \dots, C_{\hat{d}}$, we can create \hat{d} new features, e.g. by taking average of features within each cluster i.e. for each data point $\mathbf{x} \in \mathbb{R}^d$, create a new feature vector $\tilde{\mathbf{x}} \in \mathbb{R}^{\hat{d}}$

To do this, we first need

$$\text{where } \tilde{\mathbf{x}}_l = \frac{1}{|C_l|} \sum_{j \in C_l} \mathbf{x}_j \text{ for all } l \in [\hat{d}]$$

Method 1: represent feature j using values it takes on the n train data

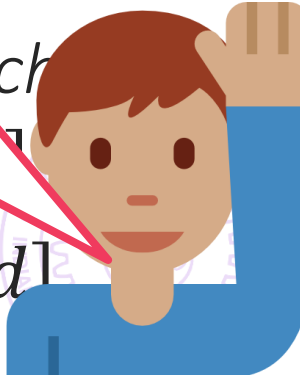
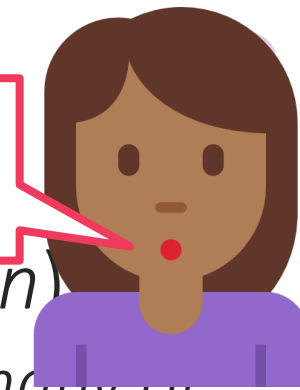
$$\mathbf{z}^j = [\mathbf{x}_j^1, \mathbf{x}_j^2, \dots, \mathbf{x}_j^n] \in \mathbb{R}^n \text{ for all } j \in [d]$$

Method 2: represent feature j u

Let n_c denote number of train c

$$\mathbf{z}^j = \left[\frac{1}{n_1} \sum_{i:y^i=1} \mathbf{x}_j^i, \frac{1}{n_2} \sum_{i:y^i=2} \mathbf{x}_j^i, \dots, \frac{1}{n_c} \sum_{i:y^i=c} \mathbf{x}_j^i \right] \in \mathbb{R}^c \text{ for all } j \in [d]$$

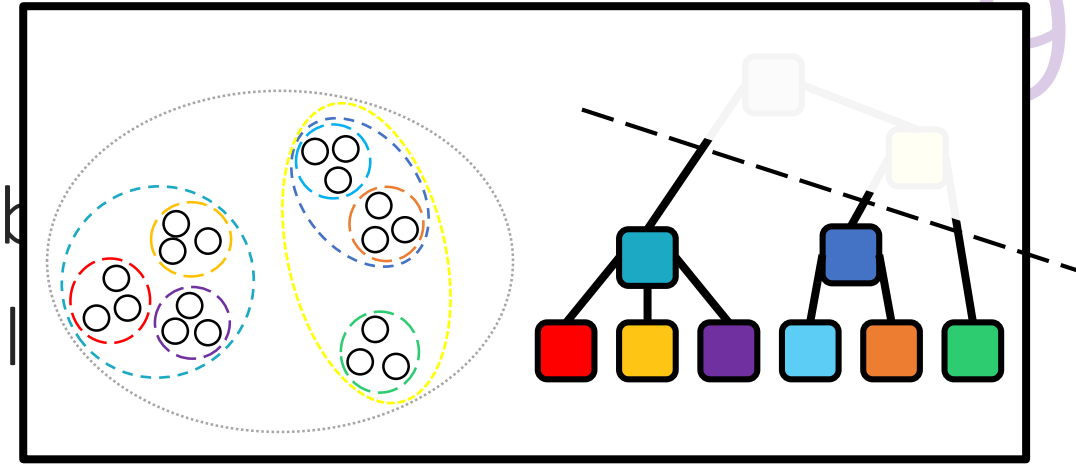
This trick is often called *feature clustering* or *feature agglomeration* and is a form of dimensionality reduction. Will see other dimensionality reduction techniques later



Variations in clustering

Might want to prevent empty clusters – k

Might want the algorithm to automatically
number of clusters \mathcal{C}



May Notice that imposing (resp. encouraging) balance among the clusters can
Aggl be seen as a form of a constraint (resp. regularization) on clustering. The
k-medoid problem also can be seen as a form of *constrained* clustering

Might not be happy with Euclidean distance as notion of similarity
– clustering with *Bregman divergences*

Several other problem variants known e.g. k medoids (uses general
 $d(\mathbf{x}^i, \boldsymbol{\mu}^c)$ instead of $\|\mathbf{x}^i - \boldsymbol{\mu}^c\|_2^2$ and $\boldsymbol{\mu}^c$ must be one of the data
points), *soft* k-means (a data point can belong to multiple clusters)

K-medoids preferable when centroids/prototypes must be real data points

