

# UNSUPERVISED LEARNING

35.1  
→ What is clustering?

Classification & Regression.

$\mathcal{D} = \{x_i, y_i\} \leftarrow$  Predict  $y_i$  given  $x_i$

$y_i \in \{0, 1\} \rightarrow$  2-class classifier.

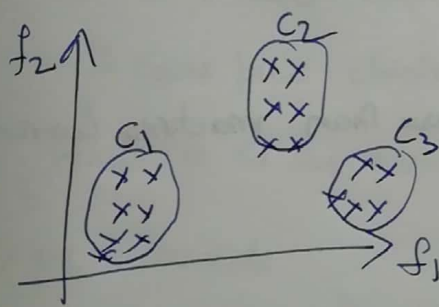
$y_i \in \mathbb{R} \rightarrow$  regression.

$$y = f(x)$$

In clustering

$\mathcal{D} = \{x_i\} \rightarrow$  no  $y_i$ 's

Task: group/cluster similar data points



→ points are close together & far away from other points

(a) Points in a cluster are close together

(b) Points in diff clusters are far away

Similar: is problem specific

$\mathcal{D} = \{x_i\}$  no  $y_i$ 's

Q How do you measure how well your algo. performed?  
↳ well will know later

Clustering: k-means, Hierarchical clustering, DBSCAN.

### 35.2 Unsupervised Learning

→ Clustering is always referred to as unsupervised learning

→ In case of classification & regression

$\{x_i, y_i\}$  we use a function  $f(x) = y$

Here  $y_i$  is supervising/helping to find function  $f(x)$

→ Semi-supervised learning:  $D = D_1 \cup D_2$

$\swarrow \quad \searrow$   
 $\{x_i, y_i\} \quad \{x_i\}$

Length of  $|D_1| \ll |D_2|$

### 35.3 Applications

→ Deeply studied in data mining course than machine learning

→ Tons of application

[https://en.wikipedia.org/wiki/Cluster\\_analysis#Applications](https://en.wikipedia.org/wiki/Cluster_analysis#Applications)

Ex: ① e-commerce: Amazon.com, Alibaba, eBay, Flipkart

Task: group "similar" customers based on their purchasing behaviour

$\swarrow$   
 $C_1, C_2, C_3, C_4, C_5$

$\rightarrow$  we can offer a discount (10%, 20%)

$\rightarrow$  we can offer a deal

$\downarrow$   
 $\left\{ \begin{array}{l} \text{money,} \\ \text{credit card} \end{array} \right\}$   $\left\{ \begin{array}{l} \text{estimate} \\ \text{income per} \\ \text{year} \end{array} \right\}$   
Geographical of Customer.

Ex ② Image segmentation (computer vision & image processing)

↳ grouping/clustering similar pixels into regions/segments

↳ Typical apply ML algo to perform obj detection

### Eg: 3 Amazon Food Review → class label.

review & text	Polarity (+ve, -ve)
r1 <input type="text"/>	+ve
r2 <input type="text"/>	-ve
r3 <input type="text"/>	+ve
⋮	⋮

(Manually reviewed & labelled)  
 ↳ Time Consuming  
 ↳ Costly / Expensive

Imagine we have 1M reviews (no, y's)

- ① cluster them into 10K groups using word-similarity (Bow, tfidf, w2v)  
 $\downarrow$   
 $c_1, c_2, c_3 \dots c_{10K}$
- ② review & label each cluster.

I will pick cluster  $c_i \rightarrow$  pick one pt  $\rightarrow$  +ve / -ve. All the points in  $c_i$  as +ve  
In great word (350 points)

- ③ ML-models

### 35.4 Metrics for clustering

$D = \{x_i\}$  ; no y's

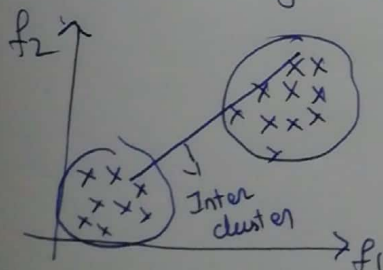
What are the performance measures?

clarity & separation  
 $\downarrow$   
 AUC  
 log-loss  
 Pre, Recall

$y_i$ 's ← class label  
 ← regression label

$\left\{ \text{ground truth} \right\}$

Geom: what is a clustering result?



$D = \{x_i\}$   
 $x_i \in \mathbb{R}^2$

→ 2 clusters

→ Intra cluster: within a cluster

→ Inter cluster: b/w clusters

In a ideal world

↳ Inter-cluster dist  $\rightarrow$  V. high

↳ Intra-cluster dist  $\rightarrow$  V. low

Dunn Index:

$$D = \frac{\max_{i,j} d(c_i, c_j)}{\max_k d'(c_k)}$$

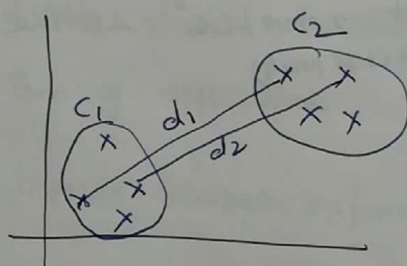
K - clusters.  
 $\{c_1, c_2, \dots, c_i, c_j, c_k\}$

maximal inter cluster distance

Intra cluster distance

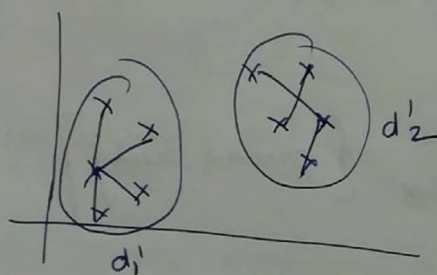
D is high  $\Rightarrow$  good clustering

$d(c_i, c_j) =$  dist b/w  $c_i$  &  $c_j$  farthest point.



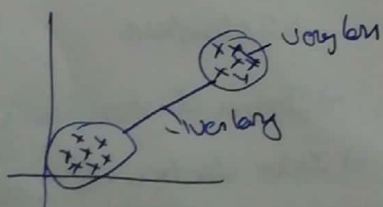
for Every pair:  $d_1, d_2, d_3, \dots, d_k$

Pick the max distance point



$\max(d'_1, d'_2)$

Ideal cluster

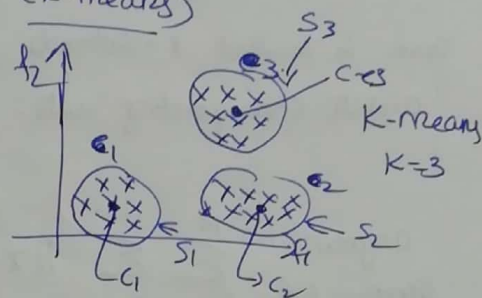




→ many other measures/metrics for measuring clustering

### 35.5 Centroids, Geometric Intuition (K-means)

- popular, simple
- Basic K-means.
- Variants of K-means.



$\left\{ \begin{array}{l} k: \# \text{ of clusters} \\ \uparrow \\ \text{hyper parameter} \end{array} \right. \rightarrow \text{CV} \rightarrow \text{ideas like}$

$c_1, c_2, c_3$  : centroids

$S_1, S_2, S_3 \in \text{sets}$

$S_1 \cap S_2 = \emptyset$  (new set)

$S_1 \cap S_3 = \emptyset$

$S_2 \cap S_3 = \emptyset$

→ when we say we want k-clusters

we will get k-centroids & k-sets

k-clusters : 1 2 3 ... k

k-centroids =  $c_1, c_2, c_3 \dots c_k$

k-sets :  $S_1, S_2, S_3 \dots S_k$

$$c_i = \frac{1}{n} \sum_{x_j \in S_i} x_j \leftarrow \text{mean point of } S_i$$

K-means : Centroid based clustering scheme

Big challenge : How to find k-centroids, once we get centroids we can easily find the sets by using nearest centroid idea.

### 35.6 Mathematical Formulation : Objective Function (K means)

$$\mathcal{D} = \{x_1, x_2, \dots, x_n\}$$

Task is to find K-centroids :  $c_1, c_2, \dots, c_k$   $\forall x_i \in S_j$

and its corresponding sets :  $S_1, S_2, \dots, S_k$   $\forall i, j \ S_i \cap S_j = \emptyset$

$$\underset{c_1, c_2, \dots, c_k}{\operatorname{argmin}} \sum_{i=1}^K \sum_{x \in S_i} \|x - c_i\|^2$$

Intra cluster distance minimized.  
↓  
S.t.  $x \in S_i$   
↑  
S<sub>i</sub> dist of  $x$  from  $c_i$   
 $\forall i, j \ S_i \cap S_j = \emptyset$

$$\begin{matrix} \{c_1, c_2, \dots, c_k\} \\ \downarrow \quad \downarrow \quad \downarrow \\ S_1 \quad S_2 \quad \dots \quad S_k \end{matrix} \rightarrow \text{using proximity idea}$$

→ In the above function there is no Intra cluster distance

→ V.V. hard to solve the above Equation.

Group hand

↳ Exponential time Complexity

In Computer Science : If we have a hard problem



approximation algo → using some hacks

→ Lloyd's algo

### 35.7 K-means Algorithm

Lloyd's algo

Step 1 : Initialization : (random initialization)

→ randomly pick  $k$  pts from  $\mathcal{D}$  and call them  $c_1, c_2, \dots, c_k$

Step 2 : Assignment

For each pt  $x_i$  in  $\mathcal{D}$

→ Select the nearest  $c_j$

→ Compute the dist  $(x_i, c_j)$   $\forall j = 1, 2, \dots, k$

→ add  $x_i$  to set  $S_j$  corresponding to centroid  $c_j$

→ we will have a point  $x_i$  assigned to a set  $S_j$   $j = 1, 2, \dots, k$   
 $i = 1$  to  $n$ .

Step 3: Recompute Centroids

→ recalculate/update  $c_j$ 's as follows

$$c_j = \frac{1}{|S_j|} \sum_{x_i \in S_j} x_i$$

→ # of points in set  $S_j$       mean-pt

Step 4: Repeat Step 2 & Step 3 until convergence

↙ assignment      ↘ Recalculate

What is Convergence?

Centroids don't change much

old Centroid =  $\{c_1, c_2, \dots, c_k\}$

new Centroid =  $\{c_1', c_2', \dots, c_k'\}$

$c_1 - c_1', c_2 - c_1', c_k - c_k'$

distance b/w new-centroid & old centroid is small

### 35.8. How to Initialize K-means++

Lloyd's algo + Initialization stage

↓  
Random Initialization ← random: pick  $k$ -pts randomly from  $D$   
↓  
 $c_1, c_2, c_3, \dots, c_k$

⇒ K-means has a problem called initialization sensitivity  
→ final clusters & centroids are very much dependent on initialization for

How to deal with this problem

① repeat K-means multiple times with different initializations

↳ Pick the best clustering based on  
    { smaller intra cluster distance  
    { large inter cluster distance

② K-means++

↳ Instead of random-init → it uses smart initialization

Initialization in K-means: (Lloyd) pick  $c_1, c_2, \dots, c_k$

① Pick the first centroid randomly →  $c_1$  from  $D$

②  $\forall x_i \in D$  create a distribution as follows  
dist

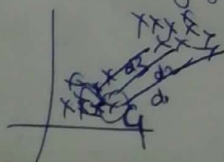
$$x_i \rightarrow \text{dist}^2(x_i, \text{nearest centroid}) \Rightarrow \|x_i - c_1\|^2$$

$x_1$	$d_1$
$x_2$	$d_2$
$x_3$	$d_3$
$\vdots$	$\vdots$
$\vdots$	$\vdots$
$x_n$	$d_n$

pick a pt from  $D - \{c_1\}$  with a prob. prop to  $d_i$

→ This is called a probabilistic approach.

→ Pick a point with large distance, because it is far from the centroid  $c_1$





(Q) Why do this probabilistically?

K-mean ÷ does get effected with outliers

### 35.9 Failure cases / limitations

→ K-means has problems when clusters are of different

\* Size

\* Densities

\* Non-globular shapes. (non-convex)



→ K-means has problems when data has outliers.

Solution for all above problems is to increase the "k"

### 35.10 K-Medoids

Problem K-means ÷  $c_1, c_2, \dots, c_k$  → Centroids may not be interpretable

$$D = \{x_1, x_2, \dots, x_n\}$$

→ We will get data point instead of centroids

Partitioning around medoids (PAM) : K-medoids.

① Initialization ÷ K-means++ → Probabilistic method pick  $k$  pts from  $D$

② Assignment ÷ closest medoid → same as in K-mean.

$$\{x_i \in S_j \mid \text{if medoid is the closest medoid to } x_i\}$$

③ Update/recompute

$$K\text{-mean} \div c_j = \frac{1}{|S_j|} \sum_{x_i \in S_j} x_i$$

K-medoids (a) swap each medoid with a non medoid point  
(PAM) (b) If loss decrease keep the swap & undo the swaps

loss in k-means

$\downarrow$   
 $\min$

$$\sum_{i=1}^K \sum_{x \in S_j} \|x - m_j\|^2$$

$m_j$   
median  $j$

$m_1$   $x_1, x_2, x_3, x_4, x_5$   $m_2$   $x_6, x_7, x_8, x_9, x_{10}$

(a) loss value

$$x_1 = m_1; x_6 = m_2 \rightarrow \text{loss} \rightarrow l_1$$

(b) swap

$$m_1 = x_2$$

$$m_2 = x_6$$

compute

$$\rightarrow \text{loss value } \left. \begin{array}{l} m_1 = x_2 \\ m_2 = x_6 \end{array} \right\} \rightarrow \text{loss} = l_2$$

If  $l_2 < l_1$

$$m_1 = x_2$$

$$m_2 = x_6$$

else

$$m_1 = x_1$$

$$m_2 = x_6$$

k-means

$\hookrightarrow$  Interpretability

$\hookrightarrow$  Kernelization ; Sim; dist

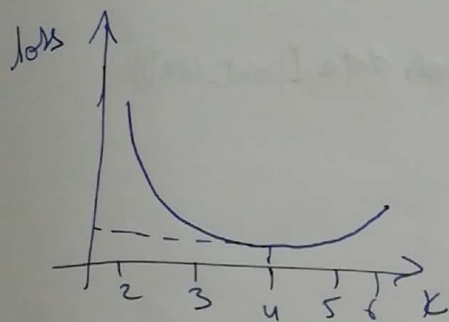
### 35.11 Determining the right 'k'

↑ hyper-parameter

- ① domain - knowledge : Food reviews  
↳ +ve & -ve  
(2 clusters)

### Elbow method (or) knee method

$$\text{loss} = \sum_{i=1}^k \sum_{x \in S_i} \|x - c_i\|^2 \rightarrow \text{minimize.}$$



best k is k=4

### 35.12 Code Samples

```
sklearn.cluster.KMeans(n_clusters=8, init='k-means++', n_init=10  
max_iter=300, tol=1e-4,)
```

#### k-medoids

```
from sklearn.metrics.pairwise import pairwise_distances
```

```
import numpy as np
```

```
import kmedoids
```

data =

```
D = pairwise_distances(data, metric='euclidean')
```

# Split into 2 clusters

```
M, c = kmedoids.KMedoids(D, 2)
```

```
print('medoids')
```

```
for point_idx in M:
```

```
    print(data[point_idx])
```

```
print(' ')
```

```
print('clustering result:')
```

```
for label in C:
```

```
    for point_idx in C[label]:
```

```
        print('label {0}: {1}'.format(label, data[point_idx]))
```

### 35.13 Time & Space Complexity

$$K\text{-means} \div O(n \times k \times d_i) \rightarrow \text{\#iterations}$$

$\uparrow \quad \uparrow \quad \uparrow$   
 $\text{\#pts} \quad \text{\#clusters} \quad \text{dim}$

Typically  $k \leq 10$   
 $i \leq 100$

$$\approx O(nd) \rightarrow \text{linear time complexity}$$

Space :  $O(nd + kd)$

$$\approx O(nd)$$

$\uparrow$   
Linear

★ K-means is quite fast