

→ "UNSUPERVISED - LEARNING"

①

↳ "Clustering" \Rightarrow (Another set of problems in ML)

We have already learnt about ["Classification"] & ["Regression"] problems.

In classification & regression, we are given a dataset D comprising of $\{x_i, y_i\}$

In "Classification", y_i are the "class labels"

where as in case of "regression" y_i are real values.

In classification & regression, given a value of $[x_i]$ we want to predict the value of $[y_i]$

We basically trying to find a function

$$y = f(x) \text{ in both "Classification \& Regression"}$$

Note \Rightarrow The only difference b/w "Classification \& Regression" is, in case of classification $y_i \in \text{finite set of values (finite set)}$ whereas in case of regression $y_i \in \text{real values}$.

Eg: Classification $y_i \in \{0, 1\} \rightarrow$ 2 class classification
Regression $y_i \in \mathbb{R}$

Now let us move to clustering & see what it is?

In case of clustering, we are just given a bunch of datapoints without any class labels (y_i 's).

In clustering there are no " y_i 's"

Where as in case of regression & classification, we are given " x_i 's" & we want to predict the corresponding " y_i 's".
 $y = f(x)$

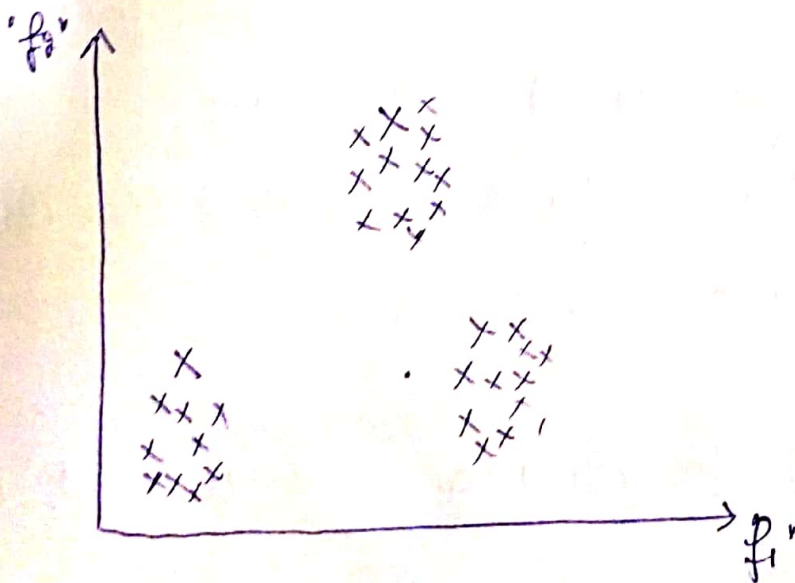
In all the techniques, like "logistic regression", "linear regression", "decision trees", "Naive Bayes", "random forest" etc, we are basically trying to find out this function $y = f(x)$ which maps " x_i 's" with " y_i 's".

In clustering, we do not have any such " y_i 's", we are just given " x_i 's" & our task is to group or cluster similar datapoints.

Task in clustering \Rightarrow To group or cluster "similar datapoints".

Let's try to understand it geometrically:-

Suppose we have two features f_1 & f_2 & we have a bunch of points as shown below -

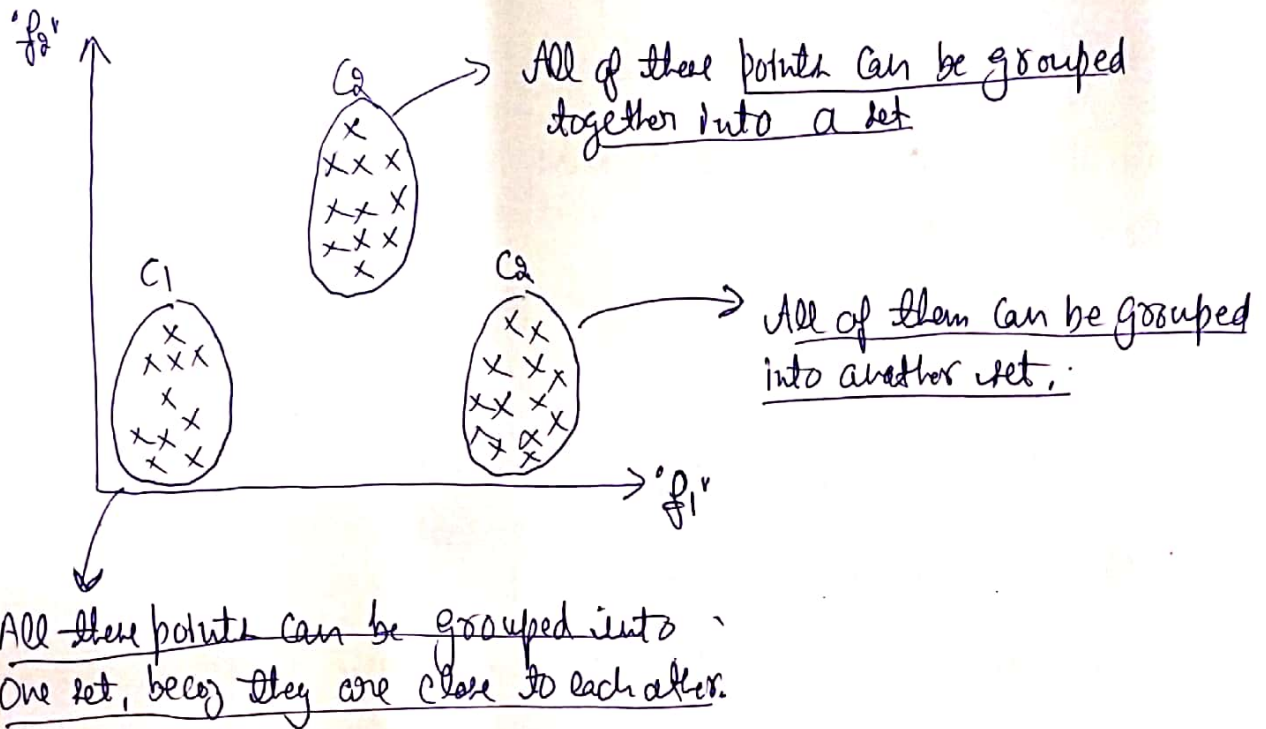


Here we do not have " y_i 's" anymore, we are simply given a bunch of points.

Q2) The task of clustering is to group "similar" points.

3

Geometrically speaking, if I say group similar points given the visual data, then the groups would look something like this.



In this case, we can say that these points can be grouped together into three clusters $\{C_1, C_2 \text{ \& } C_3\}$ & the definition of similarity here is, points are closed together (within a cluster) & they are farther away from other points (from other clusters).

The two things we are looking for are :-

a) Points in a cluster are closed together.

b) Points in different clusters are far away from each other.

Note:- The word "similar" is very very problem specific..

The task of clustering from a geometrical perspective is to group "similar points".

In clustering, from a mathematical stand point, all we are given is $D = \{x_i\}$ & no y_i 's are available.

Now, the big question here is, how do we measure, how well clustering algorithm performs.

In case of classification & regression, we have a bunch of metrics (performance metrics)

↳ AUC (Area under the curve)

↳ Precision

↳ Recall etc.

(All these metrics require y_i 's)

but in case of clustering, we don't have y_i 's, so how do we measure the performance of a clustering algorithm. we will see it later.

We will group the points in a cluster. on the basis of two things & this is the basis of clustering for most of the algo \Rightarrow

- a) Points in a cluster are close to each other.
b) Points in different clusters are far away from each other.

↳ This is the intuition behind clustering algorithms

There are three basic techniques that we will see :

1) K-Means & its variations

2) Hierarchical clustering

3) DBSCAN

} Different clustering algorithms

→ "Unsupervised Learning" :-

(5)

Clustering is often referred to as unsupervised learning.

Both classification & regression are called supervised learning algorithms or schemes. Because in both these cases, we are given $\{x_i\}$ & $\{y_i\}$ & using the $\{x_i\}$ & $\{y_i\}$ i.e. the training data, we are trying to find out a function $f(x) = y$. So, here we have y_i , which is helping or supervising us to find this function gracefully, whereas in case of clustering, we don't have y_i available, so we don't have any y_i to supervise our learning & hence it is often referred to as "Unsupervised Learning".

Apart from supervised & Unsupervised learning,

we have an area in ML called "Semi-Supervised Learning"

In this we have a big dataset D which is basically a union of D_1 & D_2

$$D = D_1 \cup D_2$$

such that for D_1 , we have both $\{x_i\}$ & $\{y_i\}$

& for D_2 , we only have $\{x_i\}$

& typically $\|D_1\| \ll \|D_2\|$ { The size of D_1 is much

smaller than the size of D_2 }, which basically means that, we have a small portion of data with labels & a large portion of data without labels. This happens when the cost of obtaining labels, i.e. $\{y_i\}$ is very expensive.

PTO

It is called semi-supervised because on a small subset of data, we have y_i 's to supervise our learning, but there is also some data where we don't have it. So it is between "supervised" & "unsupervised" & hence it is referred to as "semi-supervised".

→ "Metrics of clustering" ⇒

The dataset given for clustering comprises of just " x_i 's" & no " y_i 's".

Now the question here is, how do we measure, how good is clustering, or, what are the performance metrics, which measure how good clustering is. It is a very very important problem, becoz for classification & regression, we have seen a bunch of metrics & all of them use " y_i 's".

y_i 's → "class-labels" in case of classification
→ "Regression values" in case of Regression.

↓
There are often referred to as ground truths becoz " y_i " is the truth, that we are already given & our job is, given an " x ", we need to determine " y "

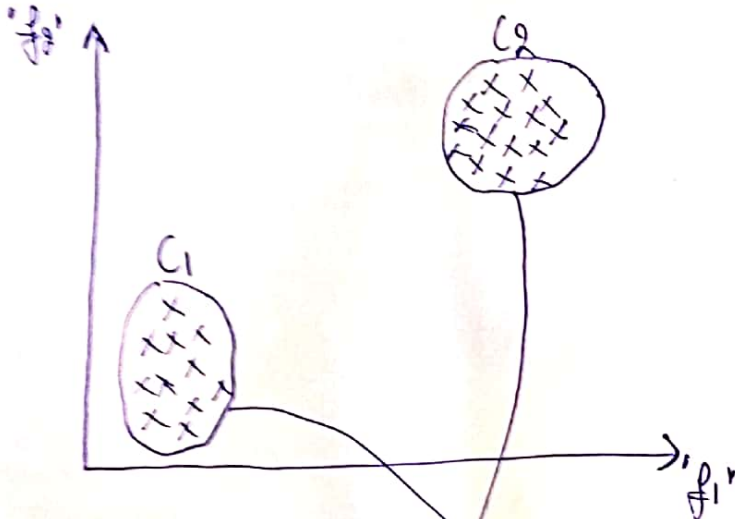
$$f(x) = y$$

↓
This is what whole "classification" & "Regression" is all about.

When we come to clustering, all we are given is " x_i 's" & no ground truths.

Now, before going into the measures, let's first understand the geometry. to get an intuition about, what is a good clustering result?

let's assume we have a cluster of points as shown, & remember we do not have any class labels here :-



Suppose, we have this data set $D = \{x_i\}$ where each x_i belongs to \mathbb{R}^2 . [$x_i \in \mathbb{R}^2$] & when we do a scatter plot this is what we get. as shown above.

Now if I say a 5 year old kid to group these points into two clusters. He will say that all these points are one cluster & all these points belong to another cluster.

let's call first group or cluster ' C_1 ' & second group or cluster ' C_2 '

There are two terms here :-

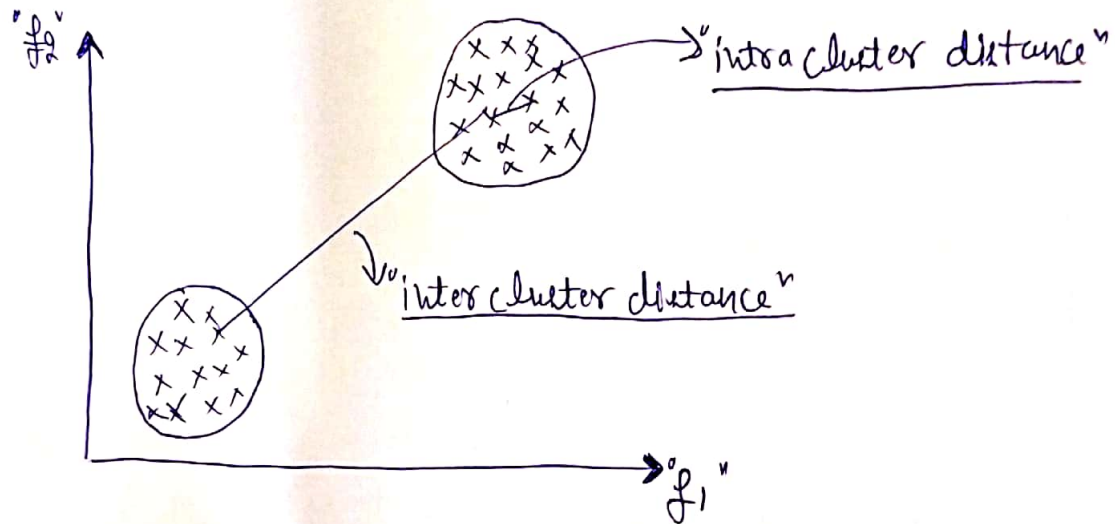
1) 'Intra-cluster'

2) 'Inter-cluster'

↓
[It means within a cluster.]

↓
[It means across, or between clusters.]

So, one thing that has been done is, if we want to group the points into clusters. then the "intra-cluster" distance is kept small, which means if we take any two points within a cluster & if we look at the distance b/w those points, the intra-cluster distance is kept small. while the intercluster distance (the distance between points belonging to two diff. clusters) is kept large.



This whole thing about "intra-clusters" being small and "inter-cluster distance" being large is the basis of how we measure the clustering effectiveness..

In an ideal world, we want our "inter-cluster distance" to be very high & "intra-cluster distance" to be very low.

There is one metric called the "Dunn" index.

It is often referred to as:

$$D$$

Suppose, we have "K-clusters", $\{C_1, C_2, C_3, \dots, C_i, C_j, C_k\}$

Then Dunn index is defined as: \rightarrow

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

Note \Rightarrow \boxed{d} is different from $\boxed{d'}$.

$d(i,j)$ is the distance between cluster $[C_i]$ & $[C_j]$.

[Numerator is basically the maximum 'inter cluster distance']

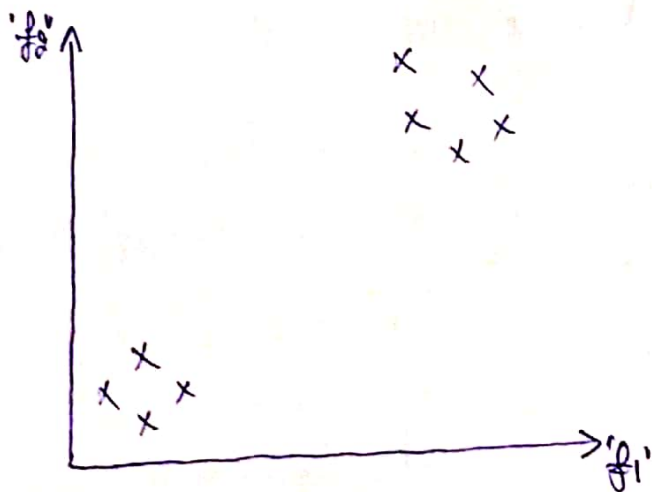
$d'(k)$ is the intra-cluster distance. It says loop at each of the clusters $\{C_1, C_2, \dots\}$ & find the cluster 'k' which has the maximum 'intra-cluster distance'

For Dunn index to be high, the 'inter-cluster distance' should be high & 'intra-cluster distance' should be low.

Note \Rightarrow If Dunn index is high, it implies very good clustering.

Let's now understand how to measure these distances.

Let's assume we have a '2d data'. Let's keep fewer points, so that it's easy for us to understand.



Here we have only two clusters.

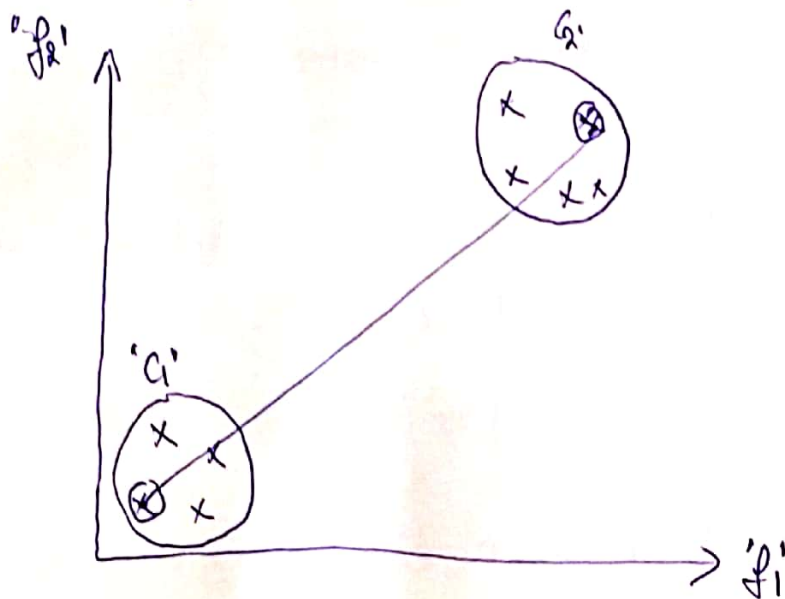
In Num index

$d(i, j)$ = distance between C_i 's & C_j 's farthest points.

Here we basically take every point in cluster $[C_i]$ & measure the distance of that from every other point in cluster $[C_j]$. We keep doing this for every pair of points. Where one point is from cluster $[C_i]$ & other point is from cluster $[C_j]$.

d_1, d_2, d_3, \dots

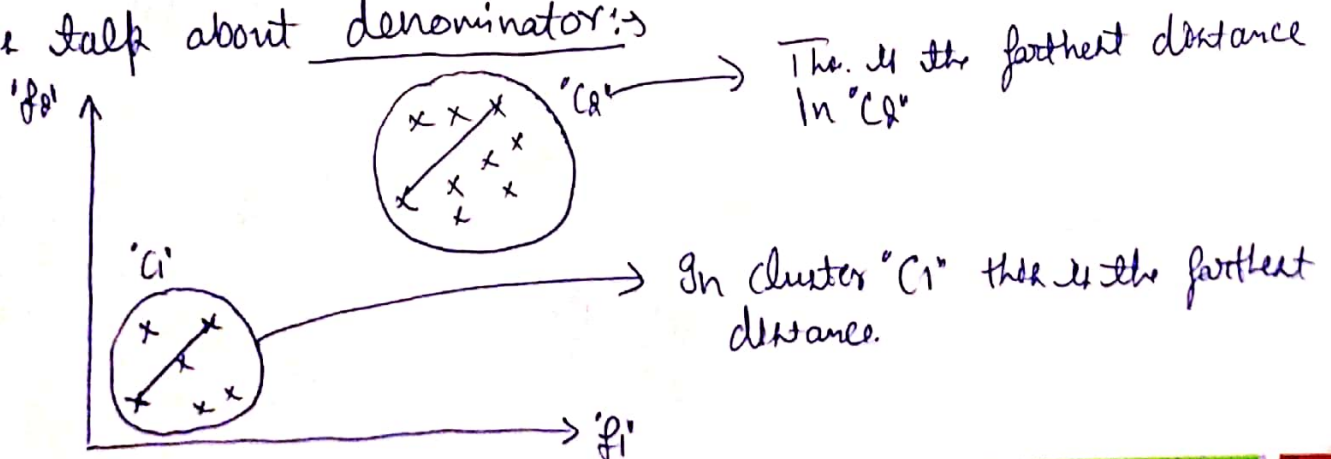
& from these distances take the maximum one.



So, $d(i, j)$ in the numerator is the distance b/w clusters C_1 & C_2 such that, it is the distance b/w two farthest points, such that the first point lies in C_1 & second point lies in C_2 .

This is about Numerator.

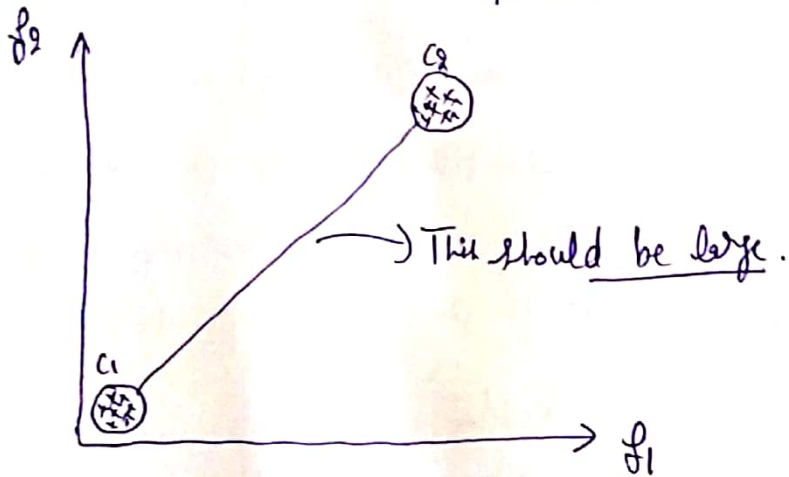
Let's talk about denominator:



So, denominator says, for every cluster, take the farthest distance & take the maximum among them. Let's call the farthest distance in cluster 'C1' as $[d_1']$ & farthest distance in cluster 'C2' as $[d_2']$ then take $\max \{d_1', d_2'\}$. That's what the denominator is. (11)

This is one metric of clustering. So ideally we want 'intra-cluster distance' to be as small as possible and 'inter-cluster distance' to be as large as possible, for Dunn Index to be large.

So, ideal clusters will look like this:-

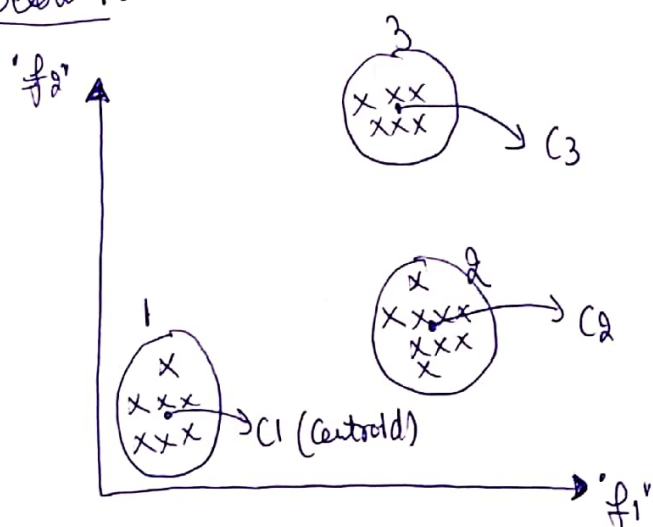


→ 'K-Means Clustering'

It is one of the most popular & also very very simple clustering algorithm. Before we understand it, let's look at the geometric intuition behind 'K-Means'.

There are various variants of K-Means like K-medoids, Kmeans++ etc. But first let's understand, what does K-means actually mean from a geometric perspective.

Let's say we have a bunch of points as shown below \Rightarrow



['K' in 'K-means' means # of clusters.]

Given this dataset, suppose we want to do 'K-means' clustering, with ' $K=3$ '. (K is # of clusters) it is a hyperparameter here, & it can be determined using an idea called 'cross-validation'.

So, what K-means with " $K=3$ " effectively does is, it groups these points into three clusters. For every cluster, this algorithm assigns something called a centroid, which is the central point.

C_1, C_2 & C_3 are centroids

& for each of these three clusters, we have a set

S_1, S_2 & S_3

Note \Rightarrow Set of points which are present in a cluster is called a set basically.

C_1, C_2, C_3 : centroids

S_1, S_2, S_3 : sets.

such that

$$S_1 \cup S_2 \cup S_3 = D$$

& no point should belong to two different sets.
ie $S_1 \cap S_2 = \emptyset, S_2 \cap S_3 = \emptyset$ & $S_1 \cap S_3 = \emptyset$.

③ This means, there should be no point which belongs to more than one set. & every point should belong to one of these clusters.

When you say, you want [K-clusters] what you get in K-means is basically [K-centroids] & they are referred to as [C_1, C_2, \dots, C_K] & [K-sets] of points referred to [S_1, S_2, \dots, S_K]

K-clusters mean

K-centroids: [$C_1, C_2, C_3, \dots, C_K$]

K-sets: [$S_1, S_2, S_3, \dots, S_K$]

① ② ③ \dots ⑫ \rightarrow clusters

Typically the centroid of any cluster, is the mean of all the points & it can be represented as:-

$$\left[C_i = \frac{1}{n} \sum_{x_j \in S_i} x_j \right] \rightarrow \text{Mean or Cental point in } S_i$$

Centroid is basically "geometric mean"

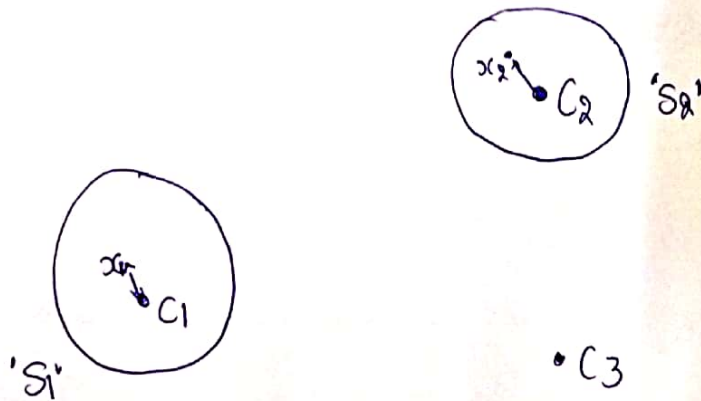
"K-means" is basically centroid based clustering scheme., because here we define each cluster using its centroid.

Every point is assigned to a cluster corresponding to the nearest centroid.

This is the idea behind clustering & the big challenge in K-means clustering is, how to find the "K-centroids"; becoz when we find the "K-centroids", we can easily find the sets.

So, once we get the ['K-centroids'], we can easily get ['K-sets'], by justly using nearest centroid idea.

Suppose, we have a centroid ' C_1 ', ' C_2 ' & ' C_3 ' as shown below & we have a datapoint ' x '. Now since ' x ' is closest to ' C_1 ', we will assign it to set ' S_1 '.



Similarly if we have a point ' x_2 ', which is closer to ' C_2 ', so we will assign it to set ' S_2 ' & so-on.

→ "K-Means Mathematical formulation & Objective function"

let's now look at the mathematical formulation of K-Means clustering :-

We are given a dataset D comprising of n points

$$D = \{x_1, x_2, x_3, \dots, x_n\}$$

our task is to find ' K ' centroids $\{C_1, C_2, \dots, C_K\}$ &

their corresponding sets of points $\{S_1, S_2, \dots, S_K\}$

such that each point in a set like S_1 has the nearest centroid C_1 .

P.T.O

such that each point x_i belongs to at least one cluster set S_j (15)

$$\forall i \quad x_i \in S_j$$

& given two clusters, their intersection is a null set

$$\forall i, j \quad S_i \cap S_j = \phi$$

Mathematically speaking these are the constraints

$$\left\{ \begin{array}{l} \forall i \quad x_i \in S_j \\ \forall i, j \quad S_i \cap S_j = \phi \end{array} \right\}$$

It says every point should belong to at least one cluster & between two clusters i & j there should be no common points

Let's write the objective function

Squared distance of x from centroid

$$\arg \min_{\substack{C_1, C_2, \dots, C_K \\ (S_1, S_2, \dots, S_K)}} \sum_{i=1}^K \sum_{x \in S_i} \|x - C_i\|^2$$

↓
assign all the cluster.

$$\text{s.t. } \left\{ \begin{array}{l} x \in S_i \\ S_i \cap S_j = \phi \end{array} \right\} \quad \text{Constraints}$$

We need to find C_1, C_2, \dots, C_K . & once we find the centroids, it is easy to find the sets using proximity idea. i.e. (Given a point, assign it to the nearest centroid.)

So, we want to find the 'K-centroids' which minimize the sum over each cluster. We want to minimize the distance of points from the centroid.

$\|x - C_i\|^2 \rightarrow$ squared distance of ' x ' from centroid.

So, what it is telling us is, find the centroids such that, each of the points is assigned to the nearest centroid so that the intra-cluster distance is minimized. (15)

So, intuitively, it says to minimize the ["intra-cluster distance"]
it is not saying anything about ["inter-cluster distance"]

$$\left\{ \begin{array}{l} \text{arg min} \\ c_1, c_2, \dots, c_k \end{array} \right. \underbrace{\sum_{i=1}^k \sum_{x \in S_i} \|x - c_i\|^2}_{\substack{\text{Sum of squared distance from centroid} \\ \text{in cluster } i}} \quad \text{s.t.} \quad \left\{ \begin{array}{l} x_i \in S_i \\ S_i \cap S_j = \emptyset \end{array} \right.$$

↓
Summate over all clusters
(2)

(1)

→ This mathematical problem is very hard to solve. It is an NP hard problem. Its time complexity is exponential $O(2^n)$.

So, we will solve this using approximation algorithms.