

* K-Means Clustering Algorithm

- K-means is an unsupervised learning algorithm that partitions data into k distinct clusters by iteratively minimizing the distance between points and their cluster centers.
- Core idea :- K-means answers 'given k groups, how do i assign each point to a group such that points within a group are as similar as possible'.
- Mathematical objective :-

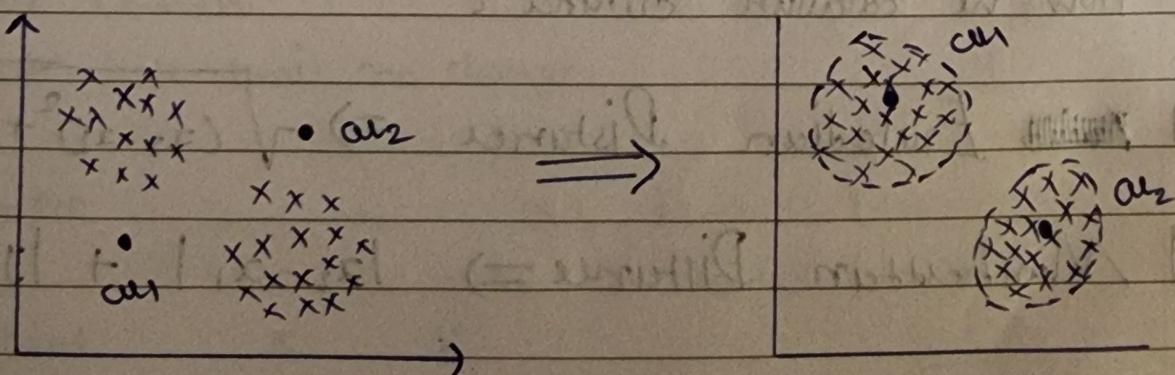
We minimize the within-cluster sum of squares (wss)

$$J = \sum_{k=1}^K \sum_{x_i \in C_k} \|x_i - c_k\|^2$$

→ Squared Euclidean distance

c_k = centroid

x_i = points

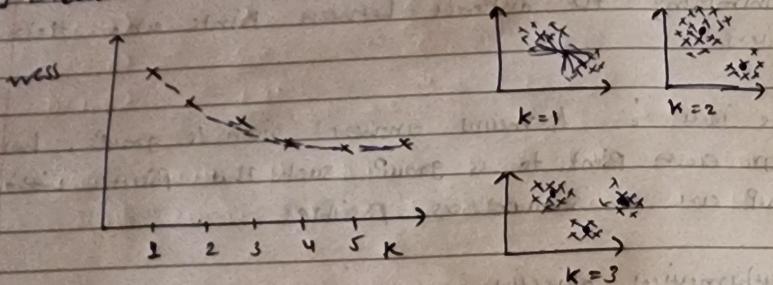


If $k = 2$, Centroids

- Now the question is how to find best k value?

- To find the optimal k value we have multiple methods

2) Elbow method (WCSS Analysis)



- When WCSS decrease sharply slow down we take that k value.

- when we only have k=2 the WCSS is too high as single centroid try to sit all points, k=2, now 2 centroid ~~will~~ try to sit its own space and due to similarity the WCSS decreases. ~~keep it same~~

* How we calculate distance?

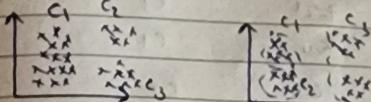
2) Euclidean Distance $\Rightarrow \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$

2) Manhattan Distance $\Rightarrow |x_2 - x_1| + |y_2 - y_1|$

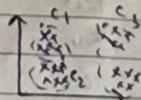
* There is one problem with random initialization of centroids which we call random initialization trap.

2) Duplicate Centroids

- 2) Poor Convergence
- 3) Inconsistent results
- 4) Suboptimal WCSS



It's random
is good



In worse case
It cluster like this.

- To address this issue we can use

- 1) K-means ++ :- use maximum distance between centroids
- 2) Multiple random starts :- not best

* Algorithm steps :-

- 1) Choose K (number of centroid)
- 2) Initialize K centroid randomly
- 3) Assign each point to nearest centroid
- 4) Update centroid.
- 5) Repeat 3-4 until no change.

* Example :-

Point	X	Y	
A	1	2	$c_1 = (1, 1)$
B	2	2	
C	6	5	
D	7	5	$c_2 = (7, 5)$

Step 2 :-

Step 2 :-

Point	c_{w1}	c_{w2}	cluster
A(1,1)	0	2.21	C1
B(2,1)	2	6.4	C1
C(6,5)	6.4	2	C2
D(7,5)	7.21	0	C3

- We calculate the distance from both centroid and assign the cluster to minimum distance to that cluster.

Step 3 : Update the Centroid.

$$c_{w1} = \left(\frac{1+2}{2}, \frac{1+1}{2} \right) = (1.5, 1)$$

$$c_{w2} = \left(\frac{6+7}{2}, \frac{5+5}{2} \right) = (6.5, 5)$$

Step 4 : Reassign and do same until the below of any condition true

- 1) no assignment changed
- 2) centroid stop moving
- 3) WCSS stop decreasing

Result : C1 : {A,B}, w1(1.5, 1)
C2 : {C,D}, w2(6.5, 5)

$$\text{WCSS} : (1-1.5)^2 + (1-1)^2 + (2-1.5)^2 + (1-1)^2 = 10.5$$

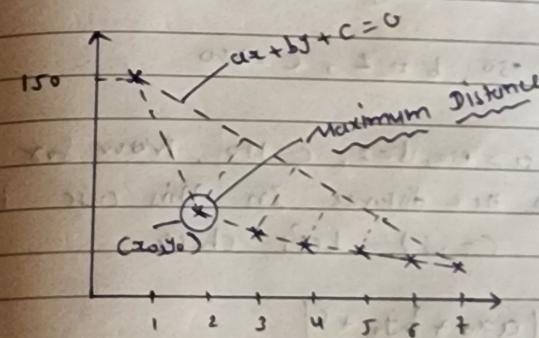
$$C2 : (6-6.5)^2 + (5-5)^2 + (7-6.5)^2 + (5-5)^2 = 0.5$$

* To validate the K value we have 2 methods

- 1) knee locator
- 2) silhouette

1) knee locator :-

- After plotting the graph between the WCSS and k value we try to find the maximum distance from the $k=1$ to $k=n$ line



- Now the question is how to find the maximum distance. Math help us.

- The distance $d = \frac{|ax_0 + by_0 + c|}{\sqrt{a^2 + b^2}}$, The perpendicular distance from (x_0, y_0) to line $ax + by + c = 0$

- In above case we have 2 points (1, 150) and (6, 0) and first we want line.

$$m = \frac{y_2 - y_1}{x_2 - x_1} = \frac{0 - 150}{6 - 1} = -30$$

$$y = mx + c$$

$$m = \frac{y - y_1}{x - x_1}$$

If I know any one point and slope I get bias and line equation

$$m(x - x_1) = y - y_1$$

$$y - 0 = -3(x - 1)$$

$$y = -3x + 1 \quad \text{Rewriting} \\ 3x + y - 1 = 0$$

$$\text{Standard form: } ax + by + c = 0$$

$$a = -3, b = 1, c = -1$$

- Now we know $a = -3, b = 1, c = -1$, Now for any point we can find the distance. In this case the nearest point is $(2, 3)$. Let's check.

$$d = \frac{|ax_0 + by_0 + c|}{\sqrt{a^2 + b^2}}$$

$$= \frac{|3(-2) + 1(3) + (-1)|}{\sqrt{3^2 + 1^2}} = \frac{|-6 + 3 - 1|}{\sqrt{10}} = \frac{4}{\sqrt{10}}$$

$$= \frac{4}{\sqrt{10}} = \frac{4}{\sqrt{10}} = 1.2649 = 1.26$$

which is maximum.

2) silhouette

- Silhouette answer: Is each point closer to its own cluster than to other clusters?

$$s_{ci} = \frac{b_{ci} - a_{ci}}{\max(a_{ci}, b_{ci})}$$

a_{ci} = Cluster c , i Point avg distance from (within cluster) all other point in same cluster c .

b_{ci} = Cluster b , i Point avg distance from (outer cluster) all other point in b cluster.

$a_{ci} > b_{ci}$ = Closer to other cluster \rightarrow wrong cluster

$a_{ci} < b_{ci}$ = Closer to own cluster \rightarrow good

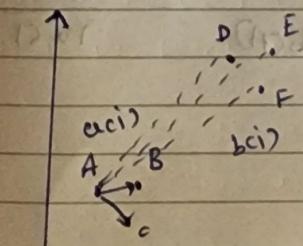
$a_{ci} = b_{ci}$ = Equally close to own and other

Denominator: $\max(a_{ci}, b_{ci})$ we to normalize the score to range $[-1, +1]$.

$s_{ci} = +1$ good

$s_{ci} = 0$ between

$s_{ci} = -1$ misclassified



$a(A)$ = avg distance within a cluster

points

$b(A)$ = avg distance other cluster point

- example :-

- we have 6 points, $k=2$

cluster 1 : $A(1,1)$, $B(2,1)$, $C(2,2)$

cluster 2 : $D(7,7)$, $E(9,8)$, $F(9,9)$

Step 1 : $a(A)$:

$$\text{distance}(A, B) = \sqrt{(1-2)^2 + (1-1)^2} = 1$$

$$d(A, C) = 2.21$$

$$a(A) = \frac{1 + 2.21}{2} = 1.61$$

Step 2 : $b(A)$: $\max(d(A, D), d(A, E), d(A, F))$

$$d(A, D) = 9.90$$

$$d(A, E) = 10.61$$

$$d(A, F) = 11.31$$

$$b(A) = 10.61$$

Step 3 : $s(A)$

$$s(A) = \frac{b(A) - a(A)}{\max(a(A), b(A))} = \frac{10.61 - 1.61}{10.61} = 0.89$$

A is now clustered

- calculate for all point

Point	$a(i)$	$b(i)$	$s(i)$
$A(1,1)$	1.21	10.61	0.89
$B(2,1)$	1.21	9.90	0.88
$C(2,2)$	1.21	9.23	0.87
$D(7,7)$	1.21	9.23	0.87
$E(9,8)$;	;	;
$F(9,9)$;	;	;

$\text{avg } s(i) = 0.88$ - good k value. else change and set by this