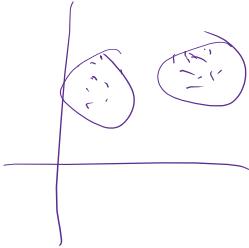


K-Means Clustering

- * Introduction to K-Means
- * K means core idea and Key Hyperparameters (k , init, max-iter)
- * choosing number of clusters and practical tips (Elbow method, scaling)
- * Implementing K-means on a Dataset
- * Model Evaluation, Limitations and Use Cases

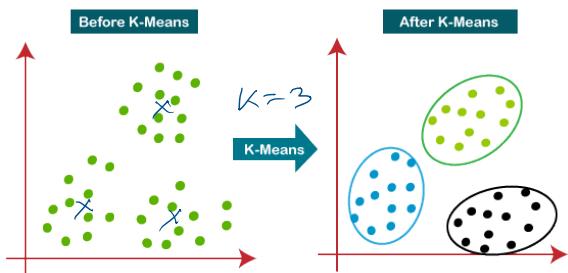


Introduction to Clustering and K-Means

* What is clustering?

Group data points such that -

- i) Intra-cluster similarity is high
 - ii) Inter-cluster similarity is low
- It's a unsupervised learning problem.

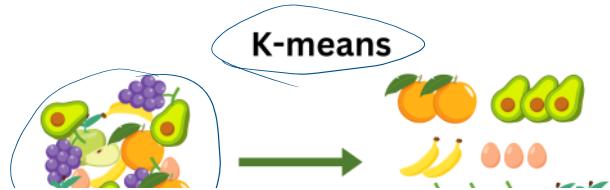
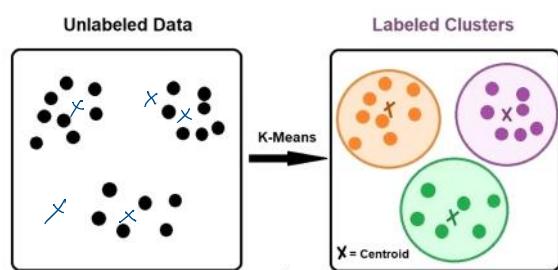


* What is K-Means?

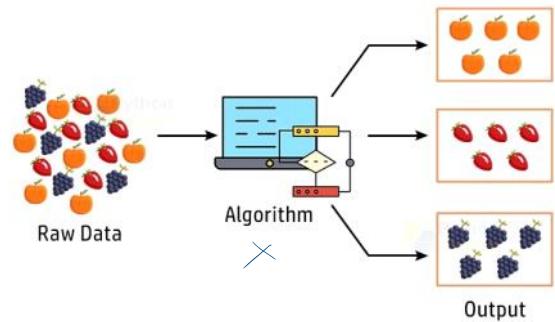
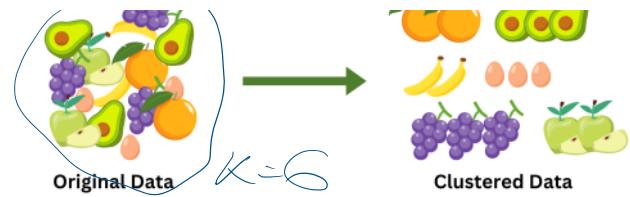
- K-Means is a partition based clustering algorithm.

It divides N data points into k disjoint clusters

= * assigning points to the nearest Centroid (mean)
, mean centroids as the mean



nearest centroid (mean)
 * updating centroids as the mean
 of assigned points



K-Means Core Idea and Key Hyperparameters

Core steps of K-means:

1. choose k (number of clusters)
2. Initialize k centroids (randomly or a smarter method)
3. Assign each points to the closest centroid
4. Update centroids (mean of points in each cluster)
5. Repeat steps 3 and 4 until assignments stop changing (convergence)

key hyperparameters:

- i) k : number of clusters
- ii) init: centroid initialization
 - ↳ random
 - ↳ k-means++
- iii) max_iter: maximum iteration

* Elbow method

≈ railing

$$k=2$$

$$[(2, 3, 4), (10, 12, 14)]$$

* Elbow method

* Scaling

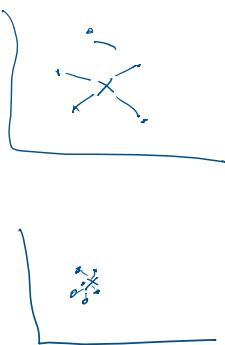
* K-means++

Choosing Number of Clusters and Practical Tips

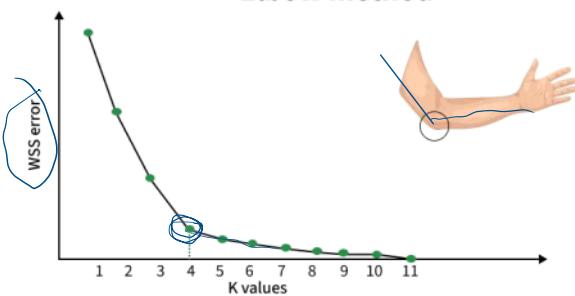
* Elbow Method:

Inertia

$$WCSS = \sum_{j=1}^k \sum_{x_i \in C_j} \|x_i - \mu_j\|^2$$



Elbow method



points:

(2,3), (3,4), (5,5)

$$\mu = \left(\frac{2+3+4}{3}, \frac{3+4+5}{3} \right)$$

$$= (3, 4)$$

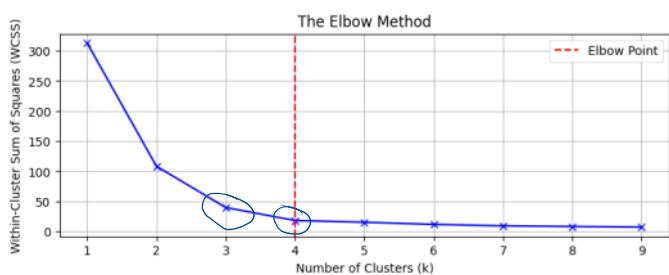
point	Distance to (3,4)	squared
(2,3)	$\sqrt{2}$	2
(3,4)	0	0
(5,5)	$\sqrt{2}$	2

$$\text{Inertia} = 2 + 0 + 2 = 4$$

Multiple clusters,

$$K = 2$$

+ Inertia = Inertia of cluster 1 + Inertia of cluster 2



$$K = 2$$

Intertia = Intertia of cluster 1 + Intertia of cluster 2

For Elbow method,

$$K = 1, 2, 3, \dots \text{None}$$

for each K , we calculate inertia

Then we plot it

K vs inertia

* Scaling %

Example:

Age range: 0 to 100

Income range: 0 to 10,000,000

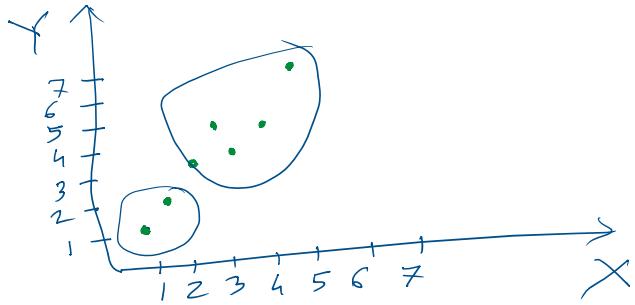
K-means++

Worked Example: K-Means with $k = 2$

A small 2D dataset of 7 points

point(ID)	x (variable 1)	y (variable 2)
→ 1	1.0	1.0
2	1.5	2.0
3	3.0	4.0
→ 4	5.0	-7.0
5	3.5	5.0
6	4.5	5.0
7	3.5	4.5





Step 1: Initialization

$$k = 2$$

$$m_1^{(0)} = (1.0, 1.0), m_2^{(0)} = (5.0, 7.0)$$

Step 2: Compute distances and assign clusters (Iteration 0)
we use Euclidean distance: $(x, y), (a, b)) = \sqrt{(x-a)^2 + (y-b)^2}$

point	Distance to m_1	Distance to m_2	Assign
1	$\sqrt{(1-1)^2 + (1-1)^2} = 0.0$	$\sqrt{(1-5)^2 + (1-7)^2} = 7.21$	C1
2	$\sqrt{(1.5-1)^2 + (2-1)^2} = 1.12$	$\sqrt{(1.5-5)^2 + (2-7)^2} = 6.10$	C1
3	3.61	3.61	C1
4	7.21	0.00	C2
5	4.72	2.50	C2
6	5.31	2.06	C2
7	4.30	2.92	C2

so, after iteration 0, assignments are:

Cluster C1: {1, 2, 3} \rightarrow

Cluster C2: {4, 5, 6, 7} \rightarrow

Step-3: Update centroids

for C1({1, 2, 3})

$$m_1^{(1)} = \left(\frac{1.0 + 1.5 + 3.0}{3}, \frac{1.0 + 2.0 + 4.0}{3} \right) = (1.83, 2.33)$$

for C2({4, 5, 6, 7})

$$m_2^{(1)} = \left(\frac{5.0 + 3.5 + 4.5 + 3.5}{4}, \frac{7.0 + 5.0 + 5.0 + 4.5}{4} \right) = (4.12, 5.38)$$

Step-3: Recompute distances and reassign (Iteration 1)

point	$d \text{ to } m_1^{(1)}$	$d \text{ to } m_2^{(1)}$	Assign
1	$\sqrt{(1-1.83)^2 + (1-2.33)^2} = 1.57$ 0.47 ↓	$\sqrt{(1-4.12)^2 + (1-5.38)^2} = 5.38$ 4.28	C1
2	2.04	1.78 ↓	C1
3	5.64	1.84 ↓	C2
4	3.15	0.73 ↓	C2
5	3.78	0.59 ↓	C2
6	2.74	1.08 ↓	C2
7			

So, after Iteration 1, assignments:

- Cluster C1 = {1, 2}
- Cluster C2 = {3, 4, 5, 6, 7}

Step-5: Update centroids again

$$\text{For } C1(\{1, 2\}) \quad m_1^{(2)} = \left(\frac{1.0 + 1.5}{2}, \frac{1.0 + 2.0}{2} \right) = (1.25, 1.50)$$

For $C2(\{3, 4, 5, 6, 7\})$

$$m_2^{(2)} = \left(\frac{3.0 + 5.0 + 3.5 + 4.5 + 3.5}{5}, \frac{4.0 + 7.0 + 5.0 + 5.0 + 4.5}{5} \right)$$

$$= (3.90, 5.10)$$

✓ Point	$d \text{ from } m_1^{(2)}$	$d \text{ from } m_2^{(2)}$	Assign

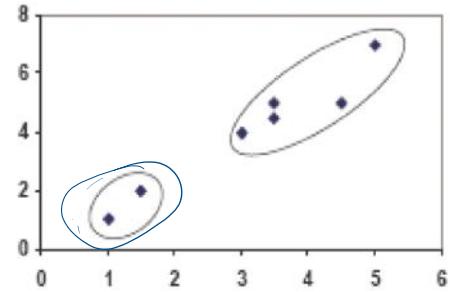
Final Assignments,

X

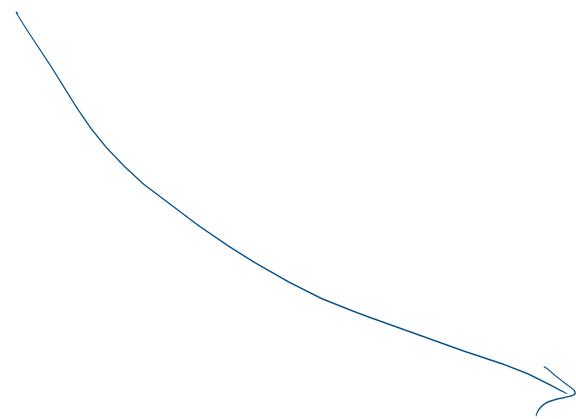
Cluster 1: $\{1, 2\}$

Cluster 2: $\{3, 4, 5, 6, 7\}$

PLOT



For, $k=3$



PLOT

