

K-Means Clustering

Dataset:

- * Point 1 : (2, 10)
- * Point 2 : (2, 5)
- * Point 3 : (8, 4)
- * Point 4 : (5, 8)
- * Point 5 : (7, 5)
- * Point 6 : (6, 4)

Initial Centroid:

- * Centroid 1: (2, 10)
- * Centroid 2: (5, 8)

Iteration 1

Point	Coordinates	Distance to centroid 1	Distance to centroid 2	Assigned cluster
1	(2, 10)	$\sqrt{(2-2)^2 + (10-10)^2} = 0$	$\sqrt{(2-5)^2 + (10-8)^2} = 3.61$	1
2	(2, 5)	$\sqrt{(2-2)^2 + (5-10)^2} = 5$	$\sqrt{(2-5)^2 + (5-8)^2} = 4.24$	2
3	(8, 4)	$\sqrt{(8-2)^2 + (4-10)^2} = 8.49$	$\sqrt{(8-5)^2 + (4-8)^2} = 5$	2
4	(5, 8)	$\sqrt{(5-2)^2 + (8-10)^2} = 3.61$	$\sqrt{(5-5)^2 + (8-8)^2} = 0$	2
5	(7, 5)	$\sqrt{(7-2)^2 + (5-10)^2} = 7.07$	$\sqrt{(7-5)^2 + (5-8)^2} = 3.61$	2
6	(6, 4)	$\sqrt{(6-2)^2 + (4-10)^2} = 7.21$	$\sqrt{(6-5)^2 + (4-8)^2} = 4.12$	2

Step 2:Update centroids

Box Cluster 1 (Point 1)

∴ New centroid 1 = (2, 10) [remain the same]

Box Cluster 2 (Point 2, 3, 4, 5, 6, 8)

$$\therefore \text{New centroid } 2 = \left(\frac{2+8+5+7+6}{5}, \frac{5+4+8+5+4}{5} \right)$$

$$= (5.6, 5.2)$$

Iteration 2

Point	Coordinates	Distance to centroid 1	Distance to centroid 2	Assigned cluster
1	(2, 10)	$\sqrt{(2-2)^2 + (10-10)^2} = 0$	$\sqrt{(2-5.6)^2 + (10-5.2)^2} = 6$	1
2	(2, 5)	$\sqrt{(2-2)^2 + (5-10)^2} = 5$	$\sqrt{(2-5.6)^2 + (5-5.2)^2} = 3.61$	2
3	(8, 4)	$\sqrt{(8-2)^2 + (4-10)^2} = 8.49$	$\sqrt{(8-5.6)^2 + (4-5.2)^2} = 2.68$	2
4	(5, 8)	$\sqrt{(5-2)^2 + (8-10)^2} = 3.61$	$\sqrt{(5-5.6)^2 + (8-5.2)^2} = 2.86$	2
5	(7, 5)	$\sqrt{(7-2)^2 + (5-10)^2} = 7.07$	$\sqrt{(7-5.6)^2 + (5-5.2)^2} = 1.414$	2
6	(6, 4)	$\sqrt{(6-2)^2 + (4-10)^2} = 7.21$	$\sqrt{(6-5.6)^2 + (4-5.2)^2} = 1.26$	2

Step 2: Update centroid:

Box Cluster 1 (Point 1)

∴ New centroid 1 = (2, 10) [unchanged]

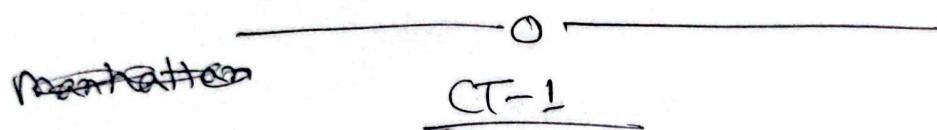
Box Cluster 2 (2, 3, 4, 5, 6)

$$\therefore \text{New centroid } 2 = \left(\frac{2+8+5+7+6}{5}, \frac{5+4+8+5+4}{5} \right)$$

$$= (5.6, 5.2) \text{ [Unchanged]}$$

After 2 iterations : Cluster 1 : Point 1 with centroid (2, 10)
 and cluster 2 : Point (2, 3, 4, 5, 6) with centroid (7, 4, 3)

The algorithm has converged as the centroid remain unchanged after the 2nd iteration.



~~Manhattan~~ Manhattan Distance \rightarrow

$$L_1 = |x_2 - x_1| + |y_2 - y_1|$$

Euclidean Distance \rightarrow

$$L_2 = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

Hence, suppose the value of K is 3.

Suppose the points A(2, 1), B(1, 3), C(7, 4), D(7, 3) are given. A and B from cluster I, where C and D from cluster II.

Apply the KNN algorithm using K=3 and Euclidean distance to determine which cluster (I or II) the query point E(8, 2) belongs to. Also illustrate a graph to visualize the clusters.

<u>Points</u>	<u>x_1</u>	<u>y_1</u>	<u>class</u>
A	2	1	cluster I
B	1	3	cluster II
C	7	4	cluster II
D	7	3	cluster II

we know, Euclidean distance

given, E(8,2) $\Rightarrow \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$

For A, $\sqrt{(8-2)^2 + (2-1)^2} \approx 6.08$

B, $\sqrt{(8-1)^2 + (2-3)^2} \approx 7.07$

C, $\sqrt{(8-7)^2 + (2-4)^2} \approx 2.24$

D, $\sqrt{(8-7)^2 + (2-3)^2} \approx 1.41$

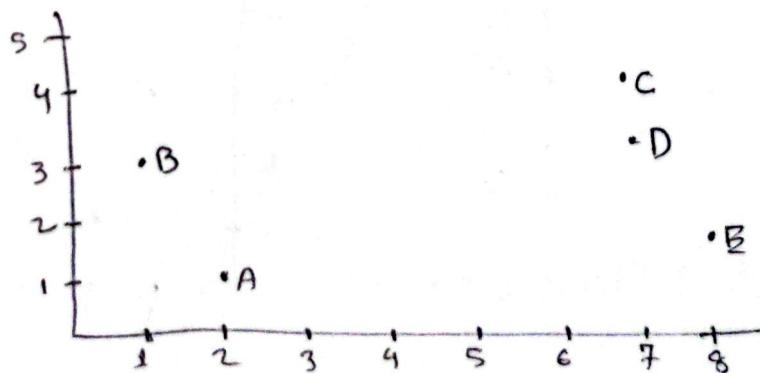
Sorting the distance, D(1.41)

C(2.24)

A(6.08)

B(7.07)

\therefore three nearest neighbors are D, C, A and 2 of them belongs to cluster B, we can classify E(8,2) into cluster II.



Date: / /

For Manhattan

	x_1	y_1	Class
A	2	1	cluster I
B	1	3	cluster I
C	7	4	cluster II
D	7	2	cluster II

\therefore We know Manhattan distance

$$E = (3, 2) \quad z |x_2 - x_1| + |y_2 - y_1|$$

\therefore We calculate the distance from $E(3, 2)$ to all given points.

$$\text{For } A: |3-2| + |2-1| = 2$$

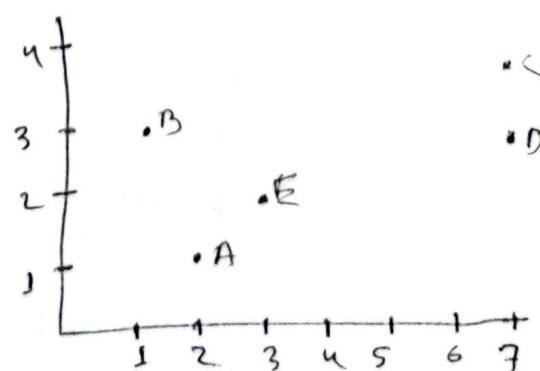
$$B = |3-1| + |2-3| = 3$$

$$C = |3-7| + |2-4| = 6$$

$$D = |3-7| + |2-3| = 5$$

\therefore Sorting the distance.

A(2)
B(3)
D(5)
C(6)



$\therefore E(3, 2) \in \text{cluster I}.$

CT-2

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Sensitivity}/\text{Recall} = \frac{TP}{TP + FN}$$

$$\text{Specificity} = \frac{TN}{TN + FP}$$

$$\therefore F\text{-measure}/F\text{-score} = \frac{2 \times (\text{Precision} \times \text{Recall})}{(\text{Precision} + \text{Recall})}$$

Positive ←
Negative ↗

Instance	1	2	3	4	5	6	7	8	9	10	11	12
Actual classification	1	1	1	1	1	1	1	0	0	0	0	0
Predicted classification	0	0	1	1	1	1	1	1	0	0	0	0
Result	FN	FN	TP	TP	TP	TP	TP	FP	FP	TN	TN	TN

Date : 11/11/2023

CT-2

Instance	X	Actual Label ($X \cdot 2$)	Predicted Label ($(X+3) \cdot 2$)	Result
1	13	1	$16 \cdot 2 = 0$	FN
2	$13+1=14$	0	1	FP
3	$13+7=20$	0	1	FP
4	$13+3=16$	0	1	FP
5	$13+4=17$	1	0	FN
6	$13+2=15$	0	0	FN
7	$13+1=14$	0	1	FP
8	$13+7=20$	0	1	FP

$$\therefore \text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} = \frac{0+0}{0+0+5+3} = 0$$

$$\therefore \text{Precision} = \frac{TP}{TP + FP} = \frac{0}{0+5} = 0$$

$$\therefore \text{Sensitivity/ Recall} = \frac{TP}{TP + FN} = \frac{0}{0+3} = 0$$

$$\therefore \text{Specificity} = \frac{TN}{TN + FP} = \frac{0}{0+5} = 0$$

Hence,

~~RP~~ $TN = 0$

$RP = 5$

$FN = 3$

$TP = 0$

	Positive	Negative
Positive	TP	FN
Negative	FP	TN

confusion KNN = $\begin{bmatrix} 35 & 2 \\ 8 & 25 \end{bmatrix}$,

$$\begin{array}{ll} TP & FN \\ FP & TN \end{array} \text{ accuracy} = \frac{TP+TN}{TP+TN+FP+FN}$$

$$= \frac{35+25}{35+25+8+2} \approx 0.85$$

$$\therefore \text{Precision} = \frac{TP}{TP+FP} = \frac{35}{35+8} \approx 0.813$$

$$\therefore \text{Recall} = \frac{TP}{TP+FN} = \frac{35}{35+2} \approx 0.945$$

$$\therefore \text{Specificity} = \frac{TN}{TN+FP} = \frac{2}{2+8} \approx 0.2$$

Date : / /

Explain the concept of Dataset and label and feature.

① Dataset: A dataset in ML is a collection of data used to train and test a model. It contains multiple examples. (also called instances)

[A collection of samples]

② Features: Features are the input variable that describe each example in the dataset. The model learns pattern from these features.

[The input data]

③ Label: The Label (also called the target variable) is the value or category that the model is trying to predict.

[The expected Output]

Explain the weight and bias update process in linear regression.

Ans. In linear regression, we predict output(y) using:

$$y = w\alpha + b$$

w = weight (how much influence a feature)

b = Bias (shift the line up/down)

How do we update weight and bias?

We use Gradient Descent to adjust w and b to minimize error. The rules are.

$$w = w - \alpha \frac{\partial L}{\partial w}$$

$$b = b - \alpha \frac{\partial L}{\partial b}$$

Steps:

1. Compute the error (difference between predicted and actual values)
2. Find gradients to loss w.r.t. w & b.
3. Update w and b using the above formula.
4. Repeat until the error is minimized.

Explain the role of derivatives in the gradient descent algorithm.

1. Gradient = Derivative \rightarrow Shows how fast a function changes.

2. Direction of Change \rightarrow

- \times If the derivative is positive, the function is increasing - move left

- \times If the derivative is negative, the function is decreasing \rightarrow move right.

- \times If the derivative is zero, we may have reached a minimum.

3. Update Rule: Adjust parameters using:

$$\theta := \theta - \alpha \nabla J(\theta)$$

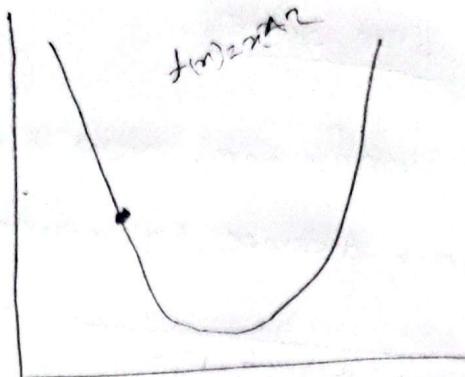
where α is the learning rate

4. Repeat Until Convergence \rightarrow keep updating until loss is minimized.

Derivatives tell us how to adjust model parameters \leftarrow to reduce error.

Date : _____

①



⑤



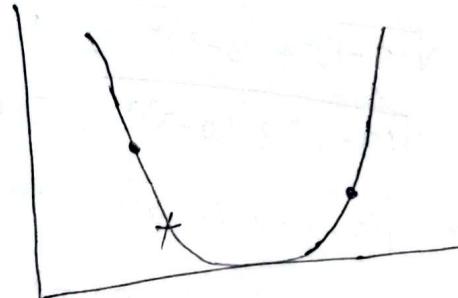
②



⑥

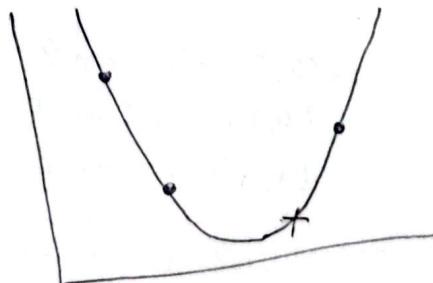


③



• previous

④



Date : / /

Mid Solve

②

We need to classify the new data point $(1, 3)$ using Euclidean distance, (3 -nearest Neighbors).

$$\text{For } (2, 3), \sqrt{(2-1)^2 + (3-1)^2}$$

$$\text{For } (-2, -1), \sqrt{(-2-1)^2 + (-1-3)^2} = 3.61$$

$$\text{For } (-1, 0), \sqrt{(-1-1)^2 + (0-3)^2} = 3.61$$

$$\text{For } (0, 2), \sqrt{(0-1)^2 + (2-3)^2} = 1.41$$

$$\text{For } (1, 1), \sqrt{(1-1)^2 + (1-3)^2} = 2.00$$

$$\text{For } (2, 3), \sqrt{(2-1)^2 + (3-3)^2} = 1.00$$

$$\text{For } (3, 0), \sqrt{(3-1)^2 + (0-3)^2} = 3.61$$

$$\text{For } (3, 2), \sqrt{(3-1)^2 + (2-3)^2} = 2.23$$

$K=3$. Hence, 3 closest points are,

$$(2, 3) \rightarrow \text{Class B}$$

$$(0, 2) \rightarrow \text{Class B}$$

$$(1, 1) \rightarrow \text{Class A}$$

\therefore the new $(1, 3)$ is classified as B

Impact of Increasing K :

neighbors.

1. If we increase ($K=5$ or 7), the model will consider more neighbors.
2. This may lead to smoother decision boundaries but can also increase bias if the dataset is small.
3. If K is too large compared to the dataset size, the model may favor the majority class, ignoring nearby points.

Explain the difference between linear regression and multivariate linear regression in modeling relationships between independent and dependent variables. Use two different dataset examples.

* Difference between Linear Regression and Multivariate Linear Regression:-

1. Linear Regression

- Involve one Independent variable (X) and one dependent variable (Y).
- The relationship: $y = \alpha + \beta X + b$

Example Dataset for Linear Regression:

<u>(X)</u>	<u>Y</u>
1	50
2	55
3	65
4	70

2. Multivariate Linear Regression:

- Involves multiple independent variables, predicting one dependent variable (Y).

$$y = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_p x_p + b$$

Example Dataset for Multivariate Linear Regression:

<u>Square Feet (X_1)</u>	<u>Bedrooms (X_2)</u>	<u>House Price</u>
1200	2	2000000
1500	3	3000000
1800	4	5000000

Date : / /

Given the following scenarios, identify whether the issue is likely due to high variance, high bias, or both;

- ① A model performs exceptionally well on the training set but poorly on the test set.
- ② A model performs poorly on both the training and test sets.
- ③ A model performs decently on the training set but its performance fluctuates significantly across different test sets.

Scenario	Issue	Explanation
I	High Variance	Overfits training data, performs poorly on test data.
II	High Bias	Underfits, performs poorly on both training and test data
III	High Variance	Performance fluctuates significantly across test sets.