



Artificial Intelligence

Dr. Mubashir Ahmad (Ph.D.)

K-Nearest Neighbor(KNN) Algorithm for Machine Learning

- K-Nearest Neighbor is one of the simplest Machine Learning algorithms based on Supervised Learning techniques.
- The K-NN algorithm assumes the similarity between the new case/data and available cases and puts the new case into the category that is most similar to the available categories.
- K-NN algorithm stores all the available data and classifies a new data point based on the similarity. This means when new data appears then it can be easily classified into a well-suited category by using the K- NN algorithm.

K-Nearest Neighbor(KNN) Algorithm for Machine Learning

- The K-NN algorithm can be used for Regression as well as for Classification but mostly it is used for Classification problems.
- K-NN is a non-parametric algorithm, which means it does not make any assumption on underlying data.
- It is also called a lazy learner algorithm because it does not learn from the training set immediately instead it stores the dataset and at the time of classification, it performs an action on the dataset.
- KNN algorithm at the training phase just stores the dataset and when it gets new data, then it classifies that data into a category that is much similar to the new data.

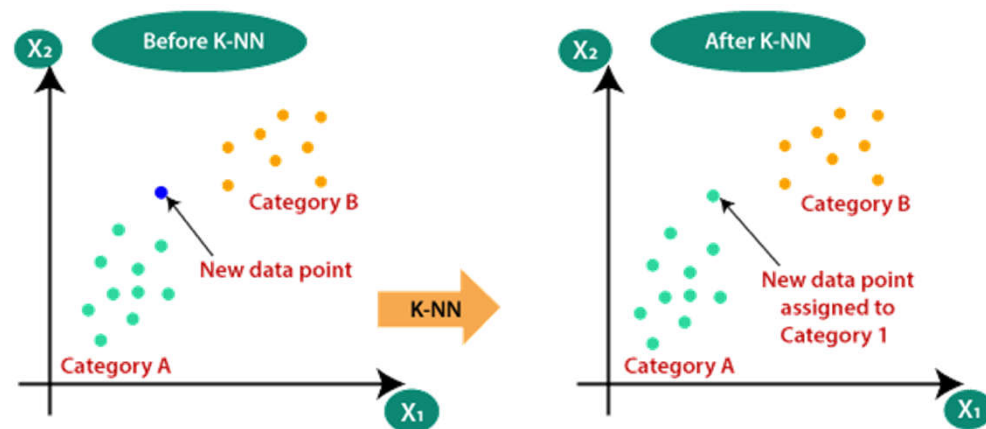
K-Nearest Neighbor(KNN) Algorithm for Machine Learning

- Suppose, we have an image of a creature that looks similar to cat and dog, but we want to know either it is a cat or dog. So for this identification, we can use the KNN algorithm, as it works on a similarity measure. Our KNN model will find the similar features of the new data set to the cats and dogs images and based on the most similar features it will put it in either cat or dog category.



Need of KNN Algorithm

- Suppose there are two categories, i.e., Category A and Category B, and we have a new data point x_1 , so this data point will lie in which of these categories? To solve this type of problem, we need a K-NN algorithm. With the help of K-NN, we can easily identify the category or class of a particular dataset. Consider the below diagram:

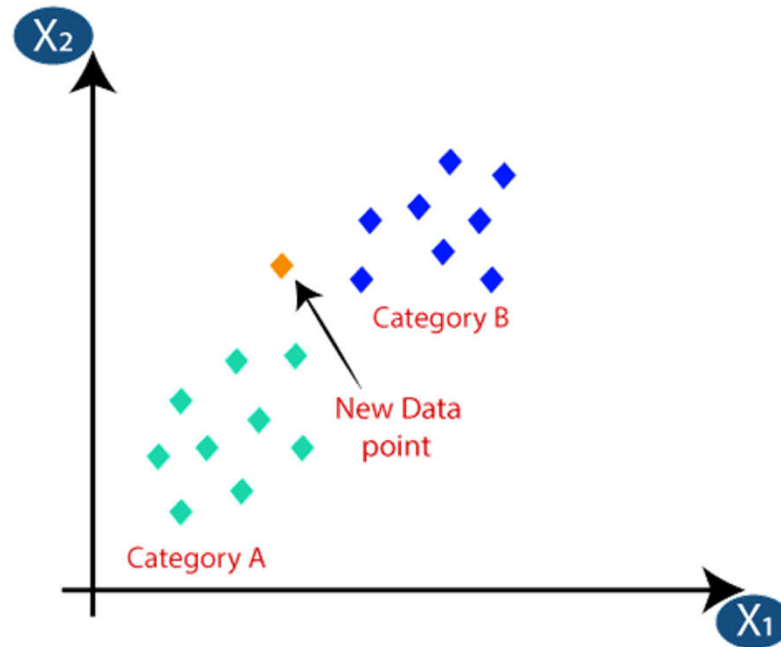


KNN Algorithm

- **Step-1:** Select the number K of the neighbors
- **Step 2:** Calculate the Euclidean distance of **K number of neighbors**
- **Step 3:** Take the K nearest neighbors as per the calculated Euclidean distance.
- **Step-4:** Among these k neighbors, count the number of the data points in each category.
- **Step-5:** Assign the new data points to that category for which the number of the neighbor is maximum.
- **Step-6:** Our model is ready.

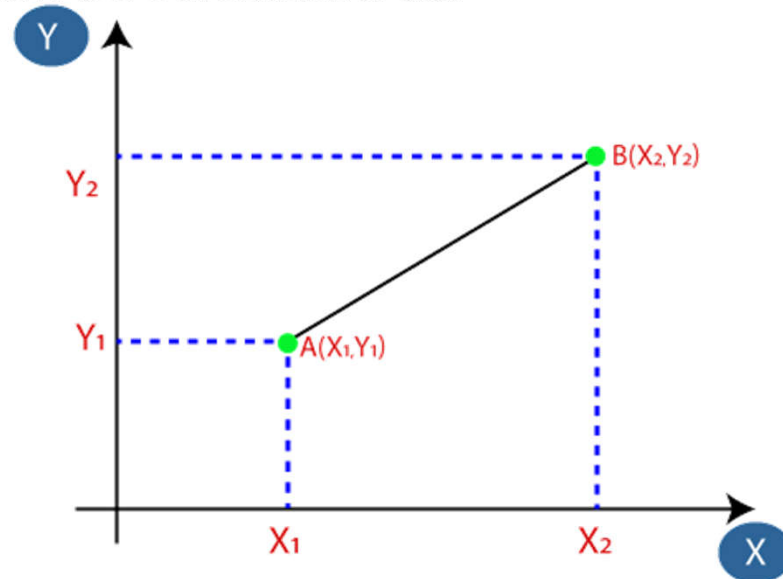
KNN Algorithm

- Suppose we have a new data point and we need to put it in the required category. Consider the below image:



KNN Algorithm

- Firstly, we will choose the number of neighbors, so we will choose the $k=5$.
- Next, we will calculate the Euclidean distance between the data points. The Euclidean distance is the distance between two points, which we have already studied in geometry. It can be calculated as:

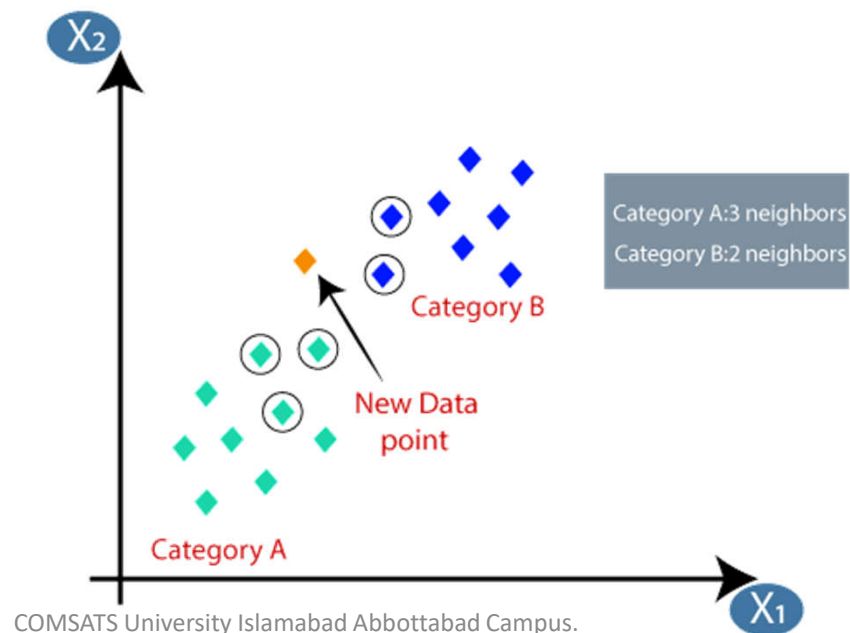


COMSATS

$$\text{Euclidean Distance between A and B} = \sqrt{(X_2 - X_1)^2 + (Y_2 - Y_1)^2}$$

KNN Algorithm

- By calculating the Euclidean distance we got the nearest neighbors, as three nearest neighbors in category A and two nearest neighbors in category B. Consider the below image:
- As we can see the 3 nearest neighbors are from category A, hence this new data point must belong to category A.



KNN Algorithm Example 1

Sepal Length	Sepal Width	Species
5.3	3.7	Setosa
5.1	3.8	Setosa
7.2	3.0	Virginica
5.4	3.4	Setosa
5.1	3.3	Setosa
5.4	3.9	Setosa
7.4	2.8	Virginica
6.1	2.8	Versicolor
7.3	2.9	Virginica
6.0	2.7	Versicolor
5.8	2.8	Virginica
6.3	2.3	Versicolor
5.1	2.5	Versicolor
6.3	2.5	Versicolor
5.5	2.4	Versicolor

KNN Algorithm Example 1

Sepal Length	Sepal Width	Species
5.2	3.1	?

Sepal Length	Sepal Width	Species
5.3	3.7	Setosa
5.1	3.8	Setosa
7.2	3.0	Virginica
5.4	3.4	Setosa
5.1	3.3	Setosa
5.4	3.9	Setosa
7.4	2.8	Virginica
6.1	2.8	Versicolor
7.3	2.9	Virginica
6.0	2.7	Versicolor
5.8	2.8	Virginica
6.3	2.3	Versicolor
5.1	2.5	Versicolor
6.3	2.5	Versicolor
5.5	2.4	Versicolor

KNN Algorithm Example 1

- Step 1: Euclidean distance formulae.

$$\text{Distance (Sepal Length, Sepal Width)} = \sqrt{(x - a)^2 + (y - b)^2}$$

$$\text{Distance (Sepal Length, Sepal Width)} = \sqrt{(5.2 - 5.3)^2 + (3.1 - 3.7)^2}$$

$$\text{Distance (Sepal Length, Sepal Width)} = 0.608$$

Sepal Length	Sepal Width	Species	Distance
5.3	3.7	Setosa	0.608

KNN Algorithm Example 1

- Step 1: Euclidean distance formulae.

Sepal Length	Sepal Width	Species	Distance
5.3	3.7	Setosa	0.608
5.1	3.8	Setosa	0.707
7.2	3.0	Virginica	2.002
5.4	3.4	Setosa	0.36
5.1	3.3	Setosa	0.22
5.4	3.9	Setosa	0.82
7.4	2.8	Virginica	2.22
6.1	2.8	Versicolor	0.94
7.3	2.9	Virginica	2.1
6.0	2.7	Versicolor	0.89
5.8	2.8	Virginica	0.67
6.3	2.3	Versicolor	1.36
5.1	2.5	Versicolor	0.60
6.3	2.5	Versicolor	1.25
5.5	2.4	Versicolor	0.75

KNN Algorithm Example 1 • Step 2: Find the Rank

Sepal Length	Sepal Width	Species	Distance	Rank
5.3	3.7	Setosa	0.608	3
5.1	3.8	Setosa	0.707	6
7.2	3.0	Virginica	2.002	13
5.4	3.4	Setosa	0.36	2
5.1	3.3	Setosa	0.22	1
5.4	3.9	Setosa	0.82	8
7.4	2.8	Virginica	2.22	15
6.1	2.8	Versicolor	0.94	10
7.3	2.9	Virginica	2.1	14
6.0	2.7	Versicolor	0.89	9
5.8	2.8	Virginica	0.67	5
6.3	2.3	Versicolor	1.36	12
5.1	2.5	Versicolor	0.60	4
6.3	2.5	Versicolor	1.25	11
5.5	2.4	Versicolor	0.75	7

COMSATS University Islamabad Abbottabad Campus.

KNN Algorithm Example 1

- Step 2: Find the nearest neighbors, if $k=2$

Sepal Length	Sepal Width	Species	Distance	Rank
5.3	3.7	Setosa	0.608	3
5.1	3.8	Setosa	0.707	6
7.2	3.0	Virginica	2.002	13
5.4	3.4	Setosa	0.36	2
5.1	3.3	Setosa	0.22	1
5.4	3.9	Setosa	0.82	8
7.4	2.8	Virginica	2.22	15
6.1	2.8	Versicolor	0.94	10
7.3	2.9	Virginica	2.1	14
6.0	2.7	Versicolor	0.89	9
5.8	2.8	Virginica	0.67	5
6.3	2.3	Versicolor	1.36	12
5.1	2.5	Versicolor	0.60	4
6.3	2.5	Versicolor	1.25	11
5.5	2.4	Versicolor	0.75	7

If $k = 1$ – Setosa

If $k = 2$ – Setosa

KNN Algorithm Example 1

- Step 2: Find the nearest neighbors, if $k=5$

Sepal Length	Sepal Width	Species	Distance	Rank
5.3	3.7	Setosa	0.608	3
5.1	3.8	Setosa	0.707	6
7.2	3.0	Virginica	2.002	13
5.4	3.4	Setosa	0.36	2
5.1	3.3	Setosa	0.22	1
5.4	3.9	Setosa	0.82	8
7.4	2.8	Virginica	2.22	15
6.1	2.8	Versicolor	0.94	10
7.3	2.9	Virginica	2.1	14
6.0	2.7	Versicolor	0.89	9
5.8	2.8	Virginica	0.67	5
6.3	2.3	Versicolor	1.36	12
5.1	2.5	Versicolor	0.60	4
6.3	2.5	Versicolor	1.25	11
5.5	2.4	Versicolor	0.75	7

If $k = 1$ – Setosa

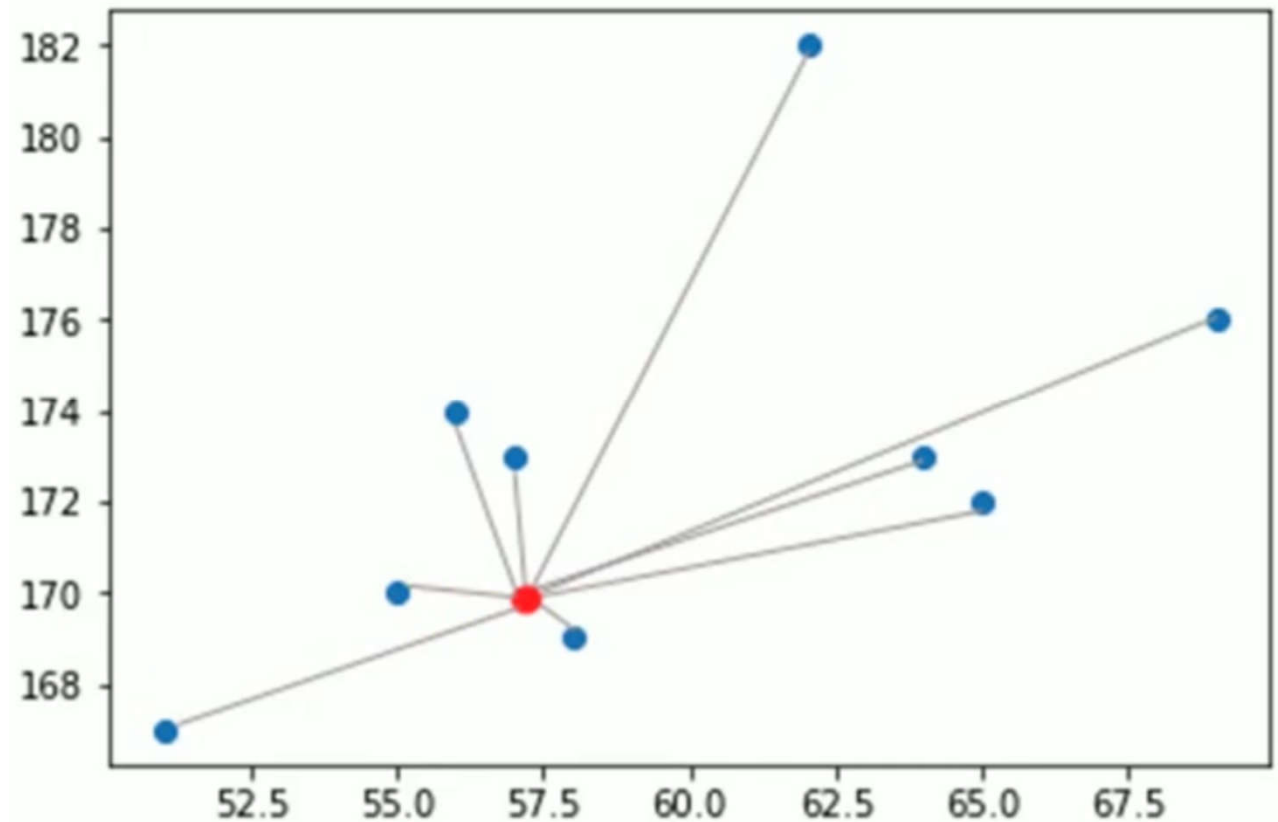
If $k = 2$ – Setosa

If $k = 5$ – Setosa

KNN Example 2:

Height (CM)	Weight (KG)	Class
167	51	Underweight
182	62	Normal
176	69	Normal
173	64	Normal
172	65	Normal
174	56	Underweight
169	58	Normal
173	57	Normal
170	55	Normal
170	57	?

KNN Example 2:



THE DISTANCE FORMULA

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

KNN Example 2:

Height (CM)	Weight (KG)	Class	Distance
167	51	Underweight	6.7
182	62	Normal	13
176	69	Normal	13.4
173	64	Normal	7.6
172	65	Normal	8.2
174	56	Underweight	4.1
169	58	Normal	1.4
173	57	Normal	3
170	55	Normal	2
170	57	?	

THE DISTANCE FORMULA

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

KNN Example 2:

Height (CM)	Weight (KG)	Class	Distance	Rank
169	58	Normal	1.4	1
170	55	Normal	2	2
173	57	Normal	3	3
174	56	Underweight	4.1	4
167	51	Underweight	6.7	5
173	64	Normal	7.6	6
172	65	Normal	8.2	7
182	62	Normal	13	8
176	69	Normal	13.4	9
170	57	?		

- If K=1, Normal
- If K=2, Normal
- If K=3, Normal
- If K=4, Normal

KNN Example 3

	Pepper	Ginger	Chilly	Liked
A	True	True	True	False
B	True	False	False	True
C	False	True	True	False
D	False	True	False	True
E	True	False	False	True

KNN Example 3

- The "Restaurant A" sells burger with optional flavors: Pepper, Ginger and Chilly.
- Every day this week you have tried a burger (A to E) and kept a record of which you liked.
- Using Hamming distance, show how the 3NN classifier with majority voting would classify

New Example - Q: pepper: false, ginger: true, chilly : true

- But How to calculate the distance for attributes with nominal or categorical values.
- Here we can use Hamming distance to find the distance between the categorical values.
- Let x_1 and x_2 are the attribute values of two instances.
- Then, in hamming distance, if the categorical values are same or matching that is x_1 is same as x_2 then distance is 0, otherwise 1.
- For example,
- If value of x_1 is **blue** and x_2 is **also blue** then the distance between x_1 and x_2 is **0**.
- If value of x_1 is **blue** and x_2 is **red** then the distance between x_1 and x_2 is **1**.

KNN Example 3

	Pepper	Ginger	Chilly	Liked	Distance
A	True	True	True	False	$1 + 0 + 0 = 1$
B	True	False	False	True	$1 + 1 + 1 = 3$
C	False	True	True	False	$0 + 0 + 0 = 0$
D	False	True	False	True	$0 + 0 + 1 = 1$
E	True	False	False	True	$1 + 1 + 1 = 3$

KNN Example 3

	Pepper	Ginger	Chilly	Liked	Distance	3NN
A	True	True	True	False	$1 + 0 + 0 = 1$	2
B	True	False	False	True	$1 + 1 + 1 = 3$	
C	False	True	True	False	$0 + 0 + 0 = 0$	1
D	False	True	False	True	$0 + 0 + 1 = 1$	2
E	True	False	False	True	$1 + 1 + 1 = 3$	

New Example - Q: pepper: false, ginger: true, chilly : true

Use Hamming Distance and

find the distance from Query Example (Q) to training examples (A-E)

KNN Example 3

New Example - Q: pepper: false, ginger: true, chilly : true

Use Hamming Distance and

find the distance from Query Example (Q) to training examples (A-E)

false

	Pepper	Ginger	Chilly	Liked	Distance	3NN
A	True	True	True	False	$1 + 0 + 0 = 1$	2
B	True	False	False	True	$1 + 1 + 1 = 3$	
C	False	True	True	False	$0 + 0 + 0 = 0$	1
D	False	True	False	True	$0 + 0 + 1 = 1$	2
E	True	False	False	True	$1 + 1 + 1 = 3$	