

## 2. K-Nearest Neighbors

```
dataset = np.array([
    [1, 2, 3, 0],
    [2, 3, 1, 1],
    [3, 1, 2, 0],
    [4, 5, 1, 1],
    [3, 3, 4, 0]
])
```

### (a) Distance Calculation

The nearest 3 neighbors for K=3:

Neighbor Index: 1, Distance: 1.4142, Class Label: 1

Neighbor Index: 2, Distance: 2.0000, Class Label: 0

Neighbor Index: 4, Distance: 2.0000, Class Label: 0

Assigned class for K=3: 0

### (b) Impact of K

Assigned class for K=1: 1

Assigned class for K=5: 0

Small K Value:

- > Benefits:
  - (1) Very sensitive the local variation of the data set;
  - (2) Can deal with the situation where there are many small group of data.
- > Drawbacks:
  - (1) over-fitting, which means the model will be specific sensitive the training data and can not fit the new data very well;
  - (2) Will be disturbed by abnormal data, if the new observation is surrounded by some wrong data, the model will make the wrong decision;
  - (3) The model can not have a very good „understanding“ of the global data set.

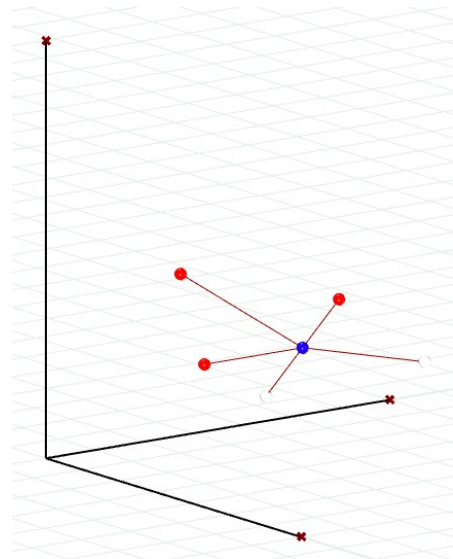
Big K Value:

- > Benefits:
  - (1) Will not be disturbed by wrong data, when there is a relative big data set, big K value will lead to a more average result;
  - (2) Have a better understanding of the whole data set.
- > Drawbacks:
  - (1) lack of precision, which means some small group will be ignored;
  - (3) Oversimplified.

### (c) Distance Weighting

Some type of distance-weighted voting:

**Inverse Distance Weighting** - In this type of voting, the closer the data is, the more impossible for the new observation to be the same. It will **reverse the result**, which means when K=3, the new observation(X1=3, X2=3, X3=2) will be classified as **Class Label: 1**.



**Gaussian Weighting** - In this type, the weight will be calculated based on the distance using a Gaussian function like  $w_i = e^{-\alpha \times d_i^2}$ . So when K=3, the nearest three neighbors are at indices 1, 2, and 4 in the data set, their new distance will be:

Neighbor 2 (index 1): New\_distance =  $1.4142 \times 0.1353 = 0.1913$

Neighbor 3 (index 2): New\_distance =  $2 \times 0.0183 = 0.1353$

Neighbor 5 (index 4): New\_distance =  $2 \times 0.0183 = 0.1353$

New class label will 1.

## Code

```
1 import numpy as np
2 from scipy.spatial import distance
3
4 dataset = np.array([
5     [1, 2, 3, 0],
6     [2, 3, 1, 1],
7     [3, 1, 2, 0],
8     [4, 5, 1, 1],
9     [3, 3, 4, 0]
10 ])
11
12 new_observation = np.array([3, 3, 2])
13
14 def knn_classification(K):
15     distances = []
16     for idx, obs in enumerate(dataset):
17         euclidean_distance = distance.euclidean(new_observation, obs[:3])
18         distances.append((euclidean_distance, obs[3], idx))
19
20     distances.sort(key=lambda x: x[0])
21     nearest_neighbors = distances[:K]
22
23     classes = [neighbor[1] for neighbor in nearest_neighbors]
24     majority_class = max(set(classes), key=classes.count)
25
26     return majority_class, nearest_neighbors
27
28 def knn_classification_weight(K):
29     distances = []
30     for idx, obs in enumerate(dataset):
31         euclidean_distance = distance.euclidean(new_observation, obs[:3])
32         distances.append((euclidean_distance, obs[3], idx))
33
34     distances.sort(key=lambda x: x[0])
35     nearest_neighbors = distances[:K]
36
37     classes = [neighbor[1] for neighbor in nearest_neighbors]
38     majority_class = max(set(classes), key=classes.count)
39
40     return majority_class, nearest_neighbors
41
42
43 def print_nearest_neighbors(nearest_neighbors, K):
44     print(f"\nThe nearest {K} neighbors for K={K}:")
45     for neighbor in nearest_neighbors:
46         distance, label, index = neighbor
47         print(f"Neighbor Index: {index}, Distance: {distance:.4f}, Class Label: {label}")
48
49 class_k1, nearest_neighbors_k1 = knn_classification(1)
50 class_k3, nearest_neighbors_k3 = knn_classification(3)
51 class_k5, nearest_neighbors_k5 = knn_classification(5)
52
53 print(f"Assigned class for K=1: {class_k1}")
54 print(f"Assigned class for K=3: {class_k3}")
55 print(f"Assigned class for K=5: {class_k5}")
56
57 print_nearest_neighbors(nearest_neighbors_k1, 1)
58 print_nearest_neighbors(nearest_neighbors_k3, 3)
59 print_nearest_neighbors(nearest_neighbors_k5, 5)
```

