

# Classification

## K Nearest Neighbor Classifier

<https://towardsdatascience.com/machine-learning-basics-with-the-k-nearest-neighbors-algorithm-6a6e71d01761>

<https://zhuanlan.zhihu.com/p/25994179>

<https://zhuanlan.zhihu.com/p/26029567>

Choosing the value of K:

If k is too small → sensitive to noise points + overfitting (doesn't generalize well).

If k is too big → neighborhood may include points from other classes.

## The KNN Algorithm

1. Load the data
2. Initialize K to your chosen number of neighbors
3. For each example in the data
  - 3.1 Calculate the distance between the query example and the current example from the data.
  - 3.2 Add the distance and the index of the example to an ordered collection
4. Sort the ordered collection of distances and indices from smallest to largest (in ascending order) by the distances
5. Pick the first K entries from the sorted collection
6. Get the labels of the selected K entries
7. If regression, return the mean of the K labels
8. If classification, return the mode of the K labels

```
from collections import Counter
import math

def knn(data, query, k, distance_fn, choice_fn):
    neighbor_distances_and_indices = []

    # 3. For each example in the data
    for index, example in enumerate(data):
        # 3.1 Calculate the distance between the query example and the current
        # example from the data.
        distance = distance_fn(example[:-1], query)

        # 3.2 Add the distance and the index of the example to an ordered collection
        neighbor_distances_and_indices.append((distance, index))
```

```

# 4. Sort the ordered collection of distances and indices from
# smallest to largest (in ascending order) by the distances
sorted_neighbor_distances_and_indices = sorted(neighbor_distances_and_indices)

# 5. Pick the first K entries from the sorted collection
k_nearest_distances_and_indices = sorted_neighbor_distances_and_indices[:k]

# 6. Get the labels of the selected K entries
k_nearest_labels = [data[i][-1] for distance, i in k_nearest_distances_and_indices]

# 7. If regression (choice_fn = mean), return the average of the K labels
# 8. If classification (choice_fn = mode), return the mode of the K labels
return k_nearest_distances_and_indices, choice_fn(k_nearest_labels)

def mean(labels):
    return sum(labels) / len(labels)

def mode(labels):
    return Counter(labels).most_common(1)[0][0]

def euclidean_distance(point1, point2):
    sum_squared_distance = 0
    for i in range(len(point1)):
        sum_squared_distance += math.pow(point1[i] - point2[i], 2)
    return math.sqrt(sum_squared_distance)

def main():
    '''
    # Regression Data
    #
    # Column 0: height (inches)
    # Column 1: weight (pounds)
    '''
    reg_data = [
        [65.75, 112.99],
        [71.52, 136.49],
        [69.40, 153.03],
        [68.22, 142.34],
        [67.79, 144.30],
        [68.70, 123.30],
        [69.80, 141.49],
        [70.01, 136.46],
        [67.90, 112.37],
        [66.49, 127.45],
    ]

    # Question:
    # Given the data we have, what's the best-guess at someone's weight if they are 60 inches tall?
    reg_query = [60]
    reg_k_nearest_neighbors, reg_prediction = knn(
        reg_data, reg_query, k=3, distance_fn=euclidean_distance, choice_fn=mean
    )

    '''
    # Classification Data
    #
    # Column 0: age
    # Column 1: likes pineapple
    '''
    clf_data = [
        [22, 1],
        [23, 1],
        [21, 1],
        [18, 1],
    ]

```

```

    [19, 1],
    [25, 0],
    [27, 0],
    [29, 0],
    [31, 0],
    [45, 0],
]
# Question:
# Given the data we have, does a 33 year old like pineapples on their pizza?
clf_query = [33]
clf_k_nearest_neighbors, clf_prediction = knn(
    clf_data, clf_query, k=3, distance_fn=euclidean_distance, choice_fn=mode
)

if __name__ == '__main__':
    main()
view rawknn_from_scratch.py hosted with ❤ by GitHub

```

[Classification and Decision Trees](#)

[Naive Bayes and SVM](#)

[Support Vector Machines](#)

[Ensemble Methods](#)

[Questions](#)