# kNN

Wang Jianfang

王建芳(in Chinese)

# Agenda

- Concept
- Example
- Code

# kNN

- **kNN**

- namely **k-NearestNeighbor** algorithm

-

# Classification of machine learning algorithms

- Supervised learning
- Unsupervised learning
- Semi-supervised learning

# Supervised learning

■ A function (model parameter) is learned from the given training data set. When new data comes, the result can be predicted according to this function. The training set of supervised learning requires input and output, which can also be said to be characteristics and objectives. The objectives of the training set are labeled manually.

❑ **Regression** :The prediction results are continuous, such as predicting tomorrow's temperature, 23,24,25 degrees.

❑ Classification:The results of the forecast are discrete, such as predicting the weather for tomorrow - cloudy, sunny, rainy.

# Unsupervised learning

- The input data is not marked, and there is no definitive result.

- The category of sample data is unknown, so it is necessary to classify the sample set (clustering) according to the similarity among samples to minimize the gap within the class and maximize the gap between the classes.

# Semi-supervised learning

- The training datasets contain both marked sample data and unmarked sample data
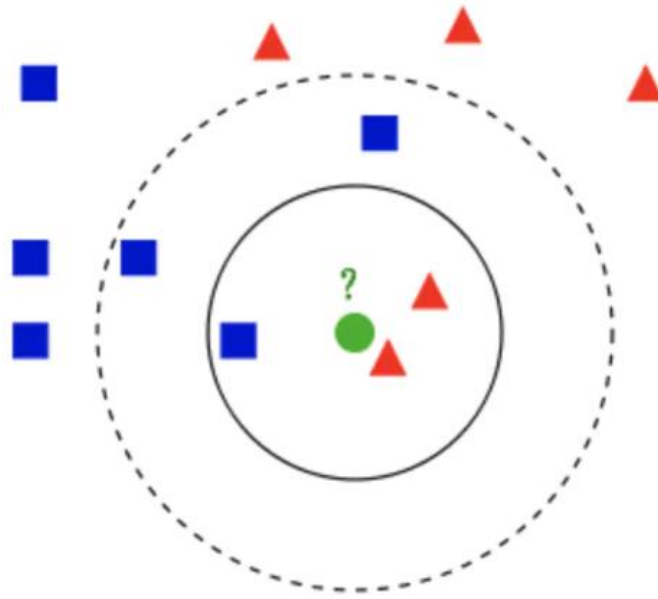
# Intro.

- K-Nearest Neighbor (KNN) is one of the simplest machine learning classification algorithms.

- In 1968, it was proposed by Cover and Hart, and its application scenarios include character recognition, text classification, image recognition and other fields.

- Its core idea is that a sample is most similar to the K samples in the data set, and if the majority of the K samples belong to a certain category, the sample also belongs to that category.

# Principle

- The core idea of the KNN algorithm is to represent the classification of the target data by using the nearest K samples.

- Specifically, there are training sample sets, and each sample contains data features and classification values.

- Input new data, comparing the data with each sample in the training sample sets, and find the nearest K. The classification with more occurrences can be used as the classification of new data in the k data,.

# Example



- We need to figure out what green looks like.
- When k is equal to 3, it's in the triangle.
- When k is equal to 5, it's a square.

# Example

- The principle of KNN is that when predicting a new value of $x$, it will judge which category $x$ belongs to according to the category of K point s closest to it. It sounds a little convoluted, but let's look at the picture.

# KNN Characteristics

- Therefore, this method has the following characteristics:
  - Supervised learning: the training sample set contains classification information
  - The algorithm is simple and easy to understand and implement
  - The results are influenced by K value, K generally does not exceed 20
  - The distance from each sample in the sample set needs to be calculated due to the large amount of calculation.
  - The imbalance of training sample set leads to inaccurate results

# KNN Method

- K-NN is used to divide each group of data into a certain class. Its pseudo code is as follows:

  ① Calculate <span style="color:red">the distance</span> between the point in the known class dataset and the current point;

  ② According to the increasing order of distance;

  ③ Select the K points with the minimum distance from the current;

  ④ The frequency of the first K points was determined;

  ⑤ The most frequent category of the first K points is determined as the prediction classification of the current point.

# Distance Functions

- Let X be an n-dimensional real vector space $R_n$ ,

$$x_i, x_j \in \mathcal{X}, \quad x_i = (x_i^{(1)}, x_i^{(2)}, \ldots, x_i^{(n)})^T, x_j = (x_j^{(1)}, x_j^{(2)}, \ldots, x_j^{(n)})^T,$$

- The $L_P$ distance of $x_i$, $x_j$ is defined as: $p > 1$

$$L_p(x_i, x_j) = \left(\sum_{l=1}^{n} |x_i^{(l)} - x_j^{(l)}|^p\right)^{\frac{1}{p}}$$

- When p = 1, it is called Manhattan distance:

$$L_1(x_i, x_j) = \sum_{l=1}^{n} |x_i^{(l)} - x_j^{(l)}|$$

- When p = 2, it is called Euclidean distance:

$$L_2(x_i, x_j) = \left(\sum_{l=1}^{n} |x_i^{(l)} - x_j^{(l)}|^2\right)^{\frac{1}{2}}$$

- When p = ∞ , it is the maximum value of each coordinate distance:

$$L_\infty(x_i, x_j) = max_l |x_i^{(l)} - x_j^{(l)}|$$

# Distance Functions Example

- Given the three points $x_1 = (1,1)^T$, $x_2 = (5,1)^T$, $x_3 = (4,4)^T$, in two-dimensional space, we try to find the nearest neighbor of $x_1$ under the $L_P$ distance when $p$ takes different values.

- Vectors

# Distance Functions Example

- The distance between $x_1$ and $x_2$:

- The Manhattan distance between $x_1$ and $x_3$ is:

- Then the European distance between $x_1$ and $x_3$ is:

- The $L_3$ distance between $x_1$ and $x_3$ is:

# Distance Functions Example

- Given the three points $x_1 = (1,1)^T$, $x_2 = (5,1)^T$, $x_3 = (4,4)^T$, in two-dimensional space, we try to find the nearest neighbor of $x_1$ under the $L_P$ distance when $p$ takes different values.

- solution:For $x_1$ and $x_2$, because the number of $x_1$ and $x_2$ in the first dimension is 1 and 5, and the number in the second dimension is 1. Therefore, when calculating the distance between $x_1$ and $x_2$, we only need to calculate $x_1^{(1)}$ and $x_2^{(1)}$, $L_P(x_1, x_2) = 4$.

# Distance Functions Example

- For $x_1$ and $x_3$, because the numbers of $x_1$ and $x_3$ in the first dimension are not the same, and the numbers in the second dimension are not the same, the Manhattan distance between $x_1$ and $x_3$ is:

$$L_1(x_1, x_3) = \sum_{l=1}^{n} |x_i^{(l)} - x_j^{(l)}| = \sum_{l=1}^{2} |x_i^{(l)} - x_j^{(l)}| = 3 + 3 = 6$$

- Then the European distance between $x_1$ and $x_3$ is:

$$L_2(x_i, x_j) = \left(\sum_{l=1}^{n} |x_i^{(l)} - x_j^{(l)}|^2\right)^{\frac{1}{2}} = \left(\sum_{l=1}^{2} |x_i^{(l)} - x_j^{(l)}|^2\right)^{\frac{1}{2}} = 3\sqrt{2} = 42.4$$

- The $L_3$ distance between $x_1$ and $x_3$ is:

$$L_3(x_i, x_j) = \left(\sum_{l=1}^{n} |x_i^{(l)} - x_j^{(l)}|^3\right)^{\frac{1}{3}} = 3.78$$

# Distance Functions

- From Step 3, we can see the distance between the position point we need to first ask for and each known point we choose, so there are three formulas for calculating the distance between two points:

Euclidean $\quad \sqrt{\displaystyle\sum_{i=1}^{k}(x_i - y_i)^2}$

Manhattan $\quad \displaystyle\sum_{i=1}^{k}|x_i - y_i|$

Minkowski $\quad \left(\displaystyle\sum_{i=1}^{k}\left(|x_i - y_i|\right)^q\right)^{1/q}$

# Code

```python
import numpy
import operator


def createDataSet():
    group = numpy.array([[1.0, 1.1], [1.0, 1.0], [0, 0], [0, 0.1]])
    labels = ['A', 'A', 'B', 'B']
    return group, labels


def classify0(inX, dataSet, labels, k):
    dataSetSize = dataSet.shape[0]
    diffMat = numpy.tile(inX, (dataSetSize, 1)) - dataSet
    sqDiffMat = diffMat ** 2
    sqDistances = sqDiffMat.sum(axis=1)
    distances = sqDistances ** 0.5
    sortedDistIndicies = distances.argsort()
    classCount={}
    for i in range(k):
        voteIlabel = labels[sortedDistIndicies[i]]
        classCount[voteIlabel] = classCount.get(voteIlabel, 0) + 1
    sortedClassCount = sorted(classCount.items(), key=operator.itemgetter(1), reverse=True)
    return sortedClassCount[0][0]


data_set, labels = createDataSet()
result = classify0([0, 0], data_set, labels, 3)
print(result)
```
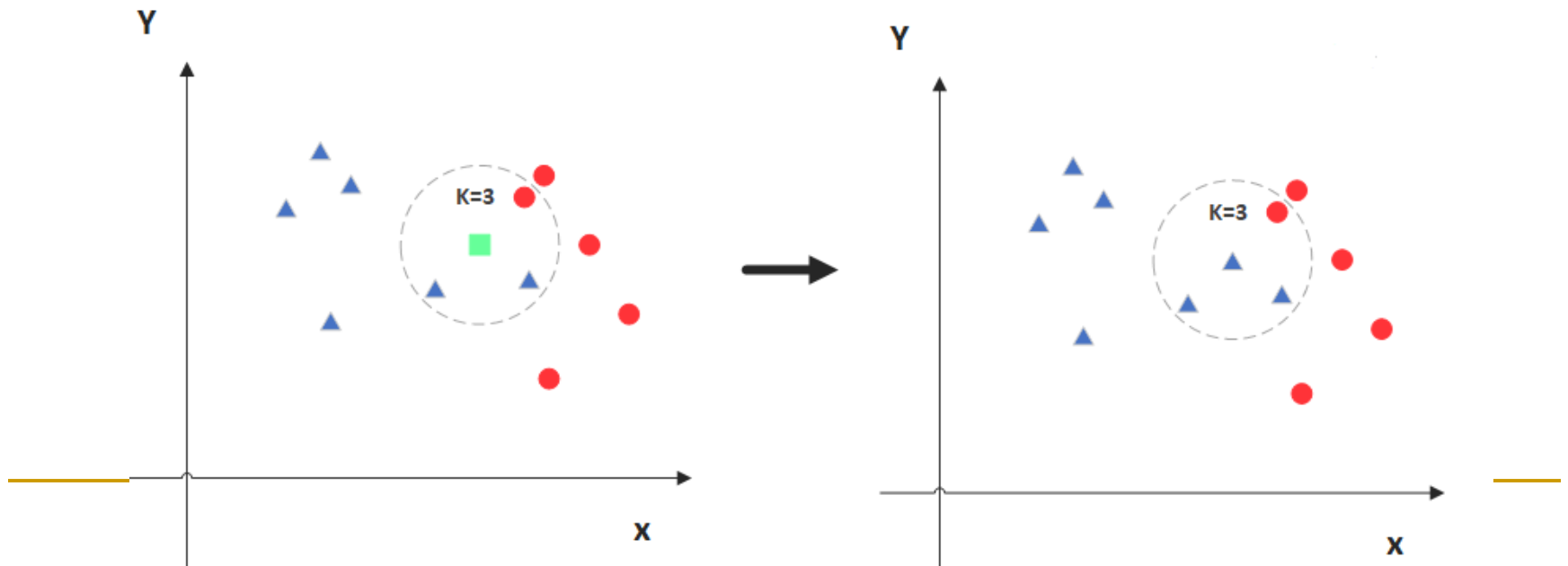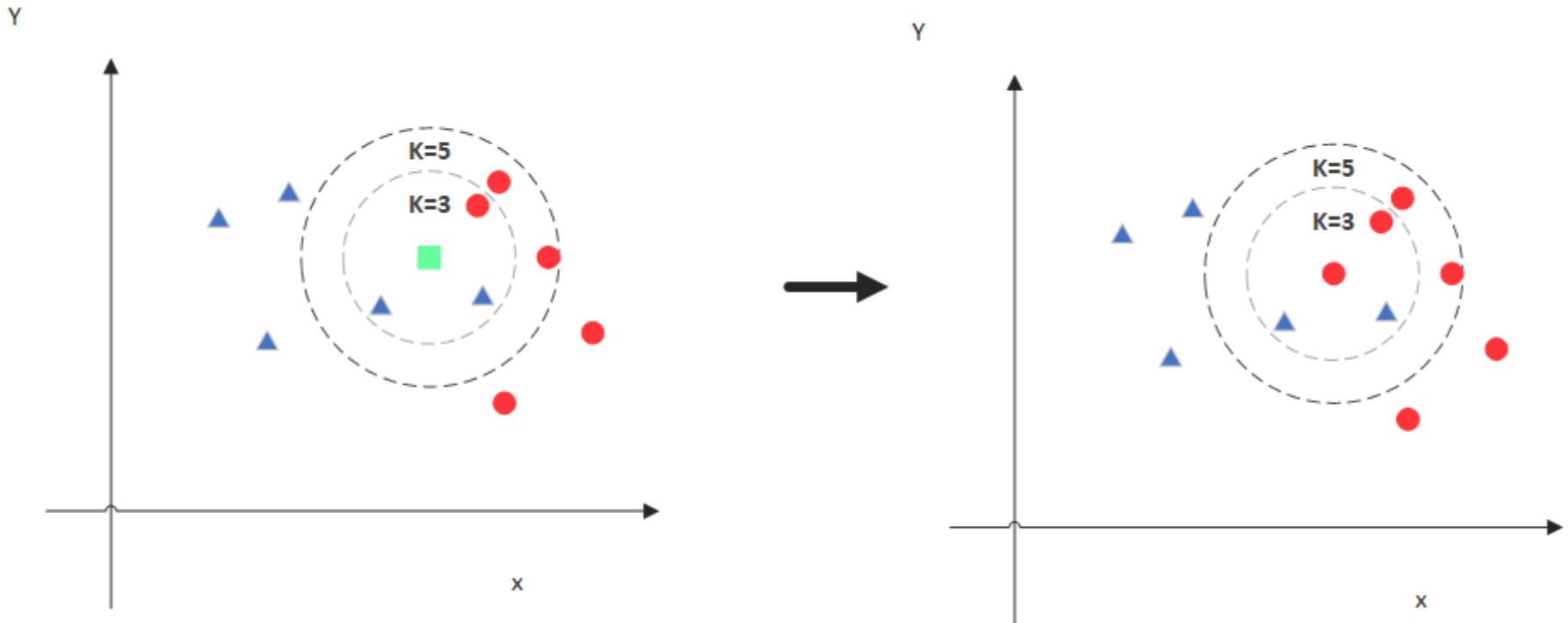
# Distance formula

- Where createdataset() is used to generate a simple dataset with tags.

- Classify0 () uses the Euclidean distance formula to calculate the distance between the new data and the data in the dataset.

- The green point on the graph is the point we're trying to predict, let's say K is equal to 3. The KNN algorithm then finds the three closest points (circled here) to see which have more categories. For example, in this case, there are more blue triangles, and the new green dots are classified as blue triangles.

■ But when K is equal to 5, the decision is different. This time there are more red circles, <span style="color:red">so the new green dots are classified as red circles</span>. From this example, we can see that K is important.

# Homework

- Implementation of KNN algorithm classification in Python (Iris Dataset)

# Questions and Comments?

**Thank you!!**