

AI and Machine Learning, Homework6

Author: Ying Yiwen
Number: 12210159

Contents

1	Principle of KNN	2
2	Main Logic	2
2.1	Prediction	2
2.2	Normalization	2
3	Effect of Different Methods	3
4	Related Code Implementation	4
4.1	Data Preprocess	4
4.2	Evaluation	4
4.3	Main Function	5

1 Principle of KNN

KNN, or K-Nearest Neighbor algorithm, is a basic classification and regression method. Its main principle is to perform classification by measuring the distance between different feature values.

Training: In the KNN algorithm, we first need to have a training dataset that has been labeled with categories. For the training process, all we have to do is to add all the data points to the model.

Predicting: When we need to classify a new data point, we calculate its distance from all the points in the training dataset and select the K points with the closest distance. Then, we decide the category of the new data point based on the categories of these K points. If most of these K points belong to a certain category, then the new data point is categorized into that category.

How to define the distance and how to select K are the keys to the KNN algorithm. In practice, the most commonly used **distance metric** is the Euclidean distance. Euclidean distance calculates the straight line distance between two points in a multidimensional space. The formula is as follows:

$$\|\mathbf{x}^{(a)} - \mathbf{x}^{(b)}\|_2 = \sqrt{\sum_{j=1}^d \left(x_j^{(a)} - x_j^{(b)}\right)^2} \quad (1)$$

For the **choice of k**, if k is very small, it may be very sensitive to localized noise, leading to overfitting and the model does not take into account the overall distributional information. Whereas a larger k will give better results, according to the rule of thumb, k must be satisfied $k \leq \sqrt{n}$. However, too large k may cause the judgment to be based on samples that are very far away from the test point, and the decision boundary of the model becomes very vague, which can lead to underfitting.

KNN is able to generate complex decision boundaries naturally and it can perform well for problems with a large number of samples. However, the KNN algorithm is sensitive to sample noise and has high requirements for the number of samples. It may not be applicable to all scenarios.

2 Main Logic

2.1 Prediction

The main function of knn is the prediction of new point. With all samples added into the model, we can easily get the label of the new point.

```
1 # predict the label of a single point
2 dist, idx = self.kdtree.query(x, k=self.k, p=2) # Find k nearest neighbors
3 neighbors_labels = [self.y_train[i] for i in idx[0]] # Get labels of neighbors
4 prediction = max(set(neighbors_labels), key=neighbors_labels.count) # Majority votes
```

2.2 Normalization

At the same time, the normalization of samples may do well with knn. We can use mean normalization, which normalizes the data to a normal distribution with mean 0 and variance 1. This better balances the weights of the different features.

$$x' = \frac{x - \mu}{\sigma} \quad (2)$$

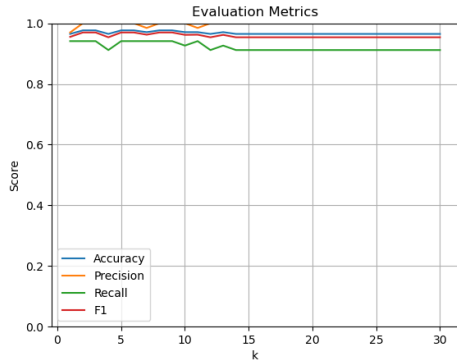
3 Effect of Different Methods

I implemented the model of knn using the principles mentioned before. For the core of the model, Euclidean distance is used. The size of k and normalization or not was discussed. Evaluating the model using accuracy, precision, recall, and F1 score, we found that knn achieved good results with this breast cancer dataset.

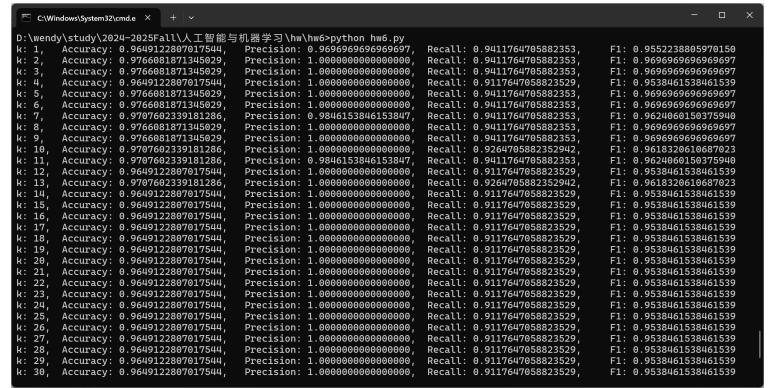
From the figure, we can see that:

When k is very small, it is true that it is not possible to take into account all the samples, and thus the model does not perform well. When k is increased to a certain level (roughly 5-10 here), the model performance is already excellent. When k is larger, the model performance does not improve significantly, but rather tends to decrease slightly.

Normalization is really useful for knn, as it balances the weights of the individual features so that they can all be taken into account, thus giving a better fit. Comparing the two graphs it can be seen that especially when the number of k is very small, normalization is able to improve the model very well; and for all cases of k, normalization is able to make the fit better.

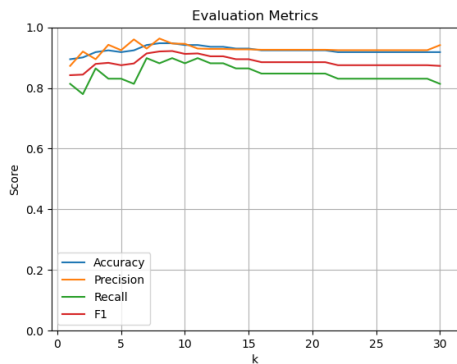


(a) Visualized Trend

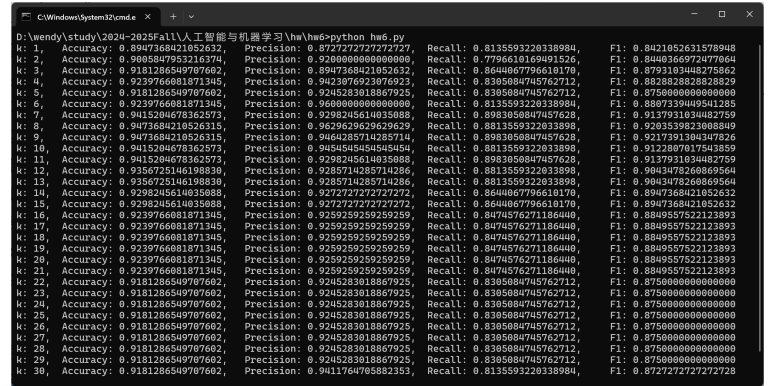


(b) Data in the Figure

Fig. 1: With Normalization



(a) Visualized Trend



(b) Data in the Figure

Fig. 2: Without Normalization

4 Related Code Implementation

4.1 Data Preprocess

In reading the data, we give the feature name for the features for our convinence.

```
1 def read_dataset():
2     # read dataset
3     column_names = [
4         'ID number', 'Diagnosis',
5         'radius_mean', 'texture_mean', 'perimeter_mean', 'area_mean', 'smoothness_mean', '
6         compactness_mean', 'concavity_mean', 'concave_points_mean', 'symmetry_mean', '
7         fractal_dimension_mean', 'radius_se', 'texture_se', 'perimeter_se', 'area_se', '
8         smoothness_se', 'compactness_se', 'concavity_se', 'concave_points_se', 'symmetry_se',
9         'fractal_dimension_se', 'radius_worst', 'texture_worst', 'perimeter_worst', '
10        area_worst', 'smoothness_worst', 'compactness_worst', 'concavity_worst', '
11        concave_points_worst', 'symmetry_worst', 'fractal_dimension_worst',
12    ]
13    data = pd.read_csv('wdbc.data', header=None, names=column_names)
14    return split_data(data)
```

Then, we can separate features and labels, so as to make it easy for the model to use the dataset.

```
1 def split_data(data):
2     # split data into features and labels
3     data = data.sample(frac=1)
4     # separate features and labels
5     x = data.iloc[:, 2:]
6     y = data.iloc[:, 1]
7     return split_test_and_train(x, y)
```

End up picking up the labels, we can split them into training set and testing set.

```
1 def split_test_and_train(x, y):
2     # split data into train and test
3     x = x.values
4     x = normalize_data(x)
5     x_train = x[:int(len(x)*0.7), :]
6     x_test = x[int(len(x)*0.7):, :]
7     y = y.values
8     y_train = y[:int(len(y)*0.7)]
9     y_test = y[int(len(y)*0.7):]
10    return x_train, x_test, y_train, y_test
```

Here is our choice for normalization, we can choose whether to use `x = normalize_data(x)` or not.

```
1 def normalize_data(x):
2     # normalize data
3     x = (x - np.mean(x, axis=0)) / np.std(x, axis=0)
4     return x
```

4.2 Evaluation

It's important for us to evaluate the model after training (though the training here is so simple). We still use accuracy, precision, recall and f1 score to evaluate it.

```
1 def evaluate(self, X_test, y_test):
2     # predict the label of multiple points
3     predictions = self.predict_multiple(X_test)
```

```

4  # calculate TP, TN, FP, FN
5  TP = np.sum((predictions == 'M') & (y_test == 'M'))
6  TN = np.sum((predictions == 'B') & (y_test == 'B'))
7  FP = np.sum((predictions == 'M') & (y_test == 'B'))
8  FN = np.sum((predictions == 'B') & (y_test == 'M'))
9  # calculate accuracy, precision, recall, f1
10 accuracy = (TP + TN) / (TP + TN + FP + FN)
11 precision = TP / (TP + FP)
12 recall = TP / (TP + FN)
13 f1 = 2 * precision * recall / (precision + recall)
14 return accuracy, precision, recall, f1

```

4.3 Main Function

All in all, the main function is like this:

```

1  # read dataset
2  x_train, x_test, y_train, y_test = read_dataset()
3
4  # test different k values
5  accuracy_list, precision_list, recall_list, f1_list = [], [], [], []
6  for k in range(1, 11):
7      # Create KNN instance and fit the model
8      knn = KNN(k=k)
9      knn.fit(x_train, y_train)
10     # Evaluate the model
11     accuracy, precision, recall, f1 = knn.evaluate(x_test, y_test)
12     accuracy_list.append(accuracy)
13     precision_list.append(precision)
14     recall_list.append(recall)
15     f1_list.append(f1)
16     print(f"k: {k},\tAccuracy: {accuracy:.16f},\tPrecision: {precision:.16f},\tRecall: {recall:.16f},\tF1: {f1:.16f}")
17
18 # plot the evaluation metrics
19 plot_evaluation_metrics(accuracy_list, precision_list, recall_list, f1_list)

```