

The data set contains 150 records in 3 categories, each with 50 pieces of data. Each record has 4 characteristics: sepal length, sepal width, petal length, and petal width. These 4 characteristics can be used to predict the iris flower belongs to which specific type such as iris-setosa, iris-versicolor and iris-virginica. If features with objects belonging to the same category such as all setosa have similar values, we can find well-defined classification in the data visualization. When the sepal length is fixed, the three types of Iris flowers are clearly distinguished in terms of petal width and petal length, although Versicolor and Virginica are slightly mixed. In particular, setosa's petal width and petal length are both significantly smaller than the other two. Secondly, Versicolor's petal width and petal length are smaller than those of Virginica. In the second scatter plot in the leftmost column, it can be seen that Versicolor and Virginiana are mixed with more points, and the distribution of Setosa is clear and obvious. When the sepal width is fixed, the three types of Iris flowers are well classified in terms of petal width and petal length. The distribution slightly worse on sepal length feature. Although Versicolor and Virginica have mixed separations, it is still very easy to distinguish different types of flowers. When the petal length is fixed, The three categories of Iris flowers are clearly distinguished in terms of sepal length, sepal width and petal width. Setosa's sepal width is slightly higher than the other two, and its sepal length is smaller than the other two. When the petal width is fixed, the three types of Iris flowers are clearly classified in three aspects: sepal length, sepal width, and petal length as well. Even Versicolor and Virginica still have few data points mixed together. The overall distribution is still obvious. These graphs show that a classifier trained using these functions may learn to classify various flower types reasonably.

## Question 2: KNN

```
In [29]: #Question 2: KNN

from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier

#1. First, divide the data into train, validation, and test sets (60%, 20%, 20%)

X_train, X_test, y_train, y_test= train_test_split(data.data,data.target, test_size=0.2
, random_state=42)

X_train, X_val, y_train, y_val= train_test_split(X_train,y_train, test_size=0.25, random_state=42)
```

```
In [30]: #2. Train the model with each classifier's default parameters. Use the train set and test the model on the test set. Store the accuracy of the model.
knn = KNeighborsClassifier(n_neighbors = 5) #5 is the default value

#Fit the model using X_train as training data and y_train as target values
knn.fit(X_train, y_train)
#Accuracy classification score
Accuracy_score_5=knn.score(X_test, y_test)
Accuracy_score_5
```

```
Out[30]: 0.9666666666666667
```

```
In [31]: #3. find k for KNN using the cross validation data
k_parameter= [1, 5, 10, 15, 20, 25, 30, 35]
cache=[] #cache to store the result
for par in k_parameter:
    knn=KNeighborsClassifier(n_neighbors = par)
    knn.fit(X_train, y_train)
    cache.append(knn.score(X_val, y_val))

    print(knn.score(X_val, y_val))
```

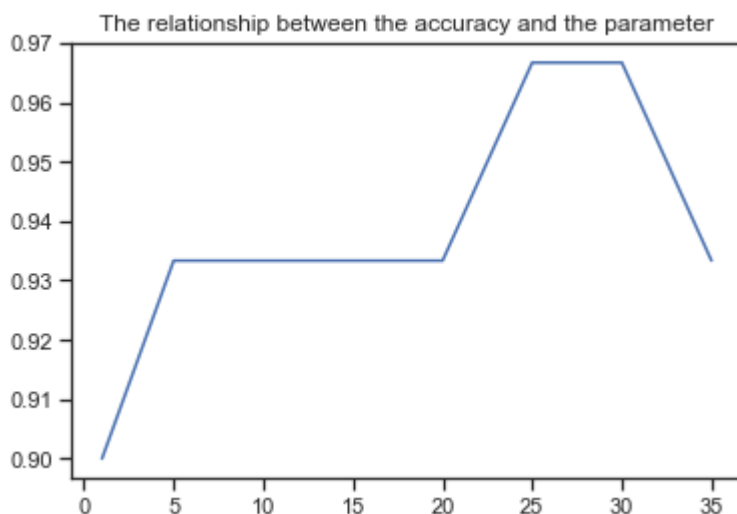
cache

```
0.9
0.9333333333333333
0.9333333333333333
0.9333333333333333
0.9333333333333333
0.9333333333333333
0.9666666666666667
0.9666666666666667
0.9333333333333333
```

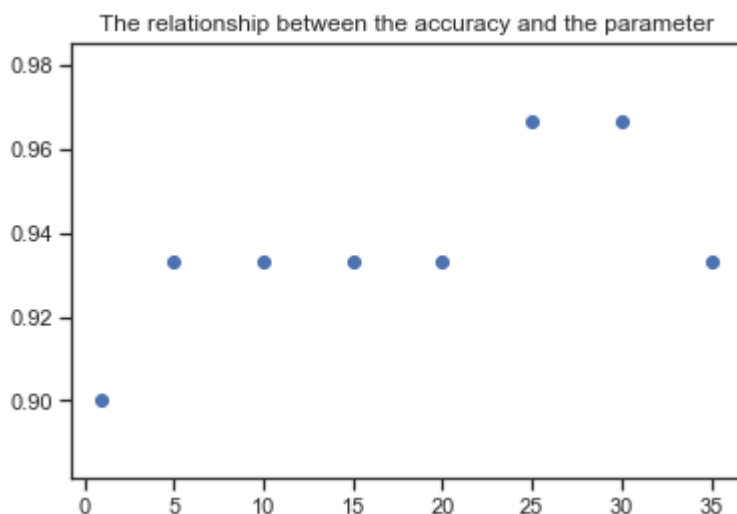
```
Out[31]: [0.9,
0.9333333333333333,
0.9333333333333333,
0.9333333333333333,
0.9333333333333333,
0.9666666666666667,
0.9666666666666667,
0.9333333333333333]
```

```
In [32]: #Plot a figure that shows the relationship between the accuracy and the parameter. Report the best k in terms of classification accuracy.
```

```
plt.title('The relationship between the accuracy and the parameter')
plt.plot(k_parameter, cache)
plt.show()
```



```
In [33]: plt.title('The relationship between the accuracy and the parameter')
plt.scatter(k_parameter, cache)
plt.show()
```



According to the two figures, when K is equal to 25 and 30, the accuracy is the highest. But based on speed of convergence and generalisation I pick 30 to be the best parameter. I set up the K to be 30 because of the accuracy. The lower value of k will predict a more locally model, and higher value of k will predict a more globally model. The smaller K can make the model more complex and consume too much to the computer, it is not good to the complexity(big O) and speed of convergence. It may also have overfitting problem.

```
In [34]: #4. Using the best found parameters, train the model using the train set and test the model on the test set. Report the accuracy of the model.
#Using the parameter is 25
knn25 = KNeighborsClassifier(n_neighbors = 25)
knn25.fit(X_train, y_train)
#Accuracy classification score when k=25
Accuracy_score_25=knn25.score(X_test, y_test)
Accuracy_score_25
```

Out[34]: 1.0

```
In [35]: #Using the parameter 30
knn35 = KNeighborsClassifier(n_neighbors = 30)
knn35.fit(X_train, y_train)
#Accuracy classification score when k=35
Accuracy_score_35=knn35.score(X_test, y_test)
Accuracy_score_35
```

Out[35]: 1.0

According to the above results, when  $k$  is equal to 25 and 30 to the KNN classifier, the accuracy is the highest. But based on speed of convergence and generalisation I pick 30 to be the best parameter. I set up the  $K$  to be 30 because of the accuracy. The lower value of  $k$  will predict a more locally model, and higher value of  $k$  will predict a more globally model. The smaller  $K$  can make the model more complex and consume too much to the computer, it is not good to the complexity(big O) and speed of convergence. It may also have overfitting problem. Its accuracy is 100% on the test set. The accuracy of the model is 100% on the test set when  $K=30$ .

## Question 3: SVM