

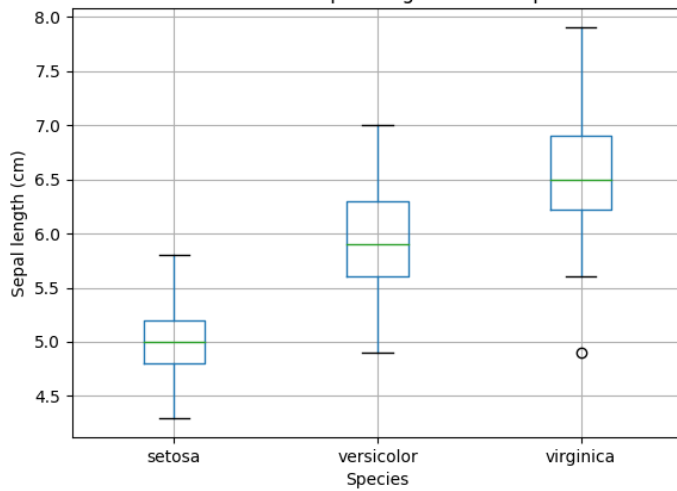
Q1:

```
Python Console
Time [sec] (for loop): 0.38894033432006836
Time [sec] (np loop): 0.009993553161621094
np loop is 38.919123962210136 times faster than for loop.
```

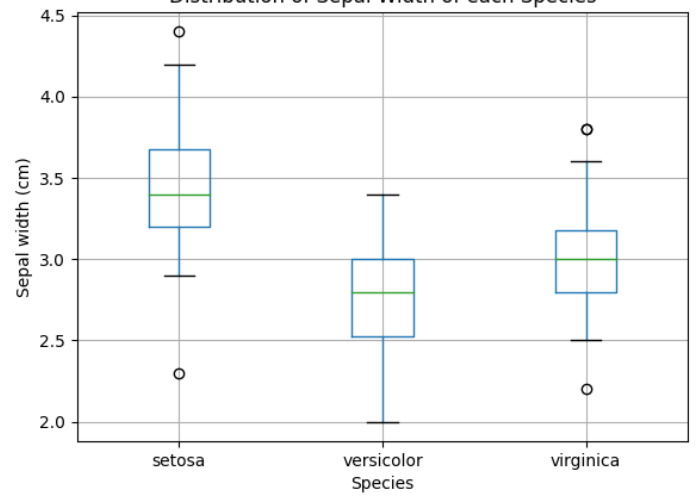
Q2:

b)

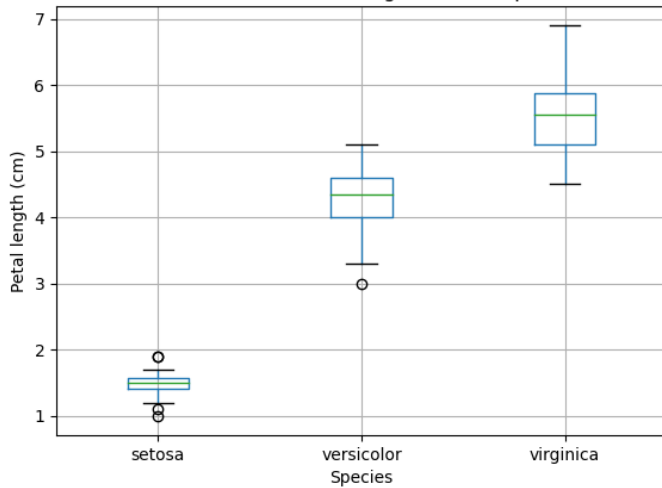
Distribution of Sepal Length of each Species



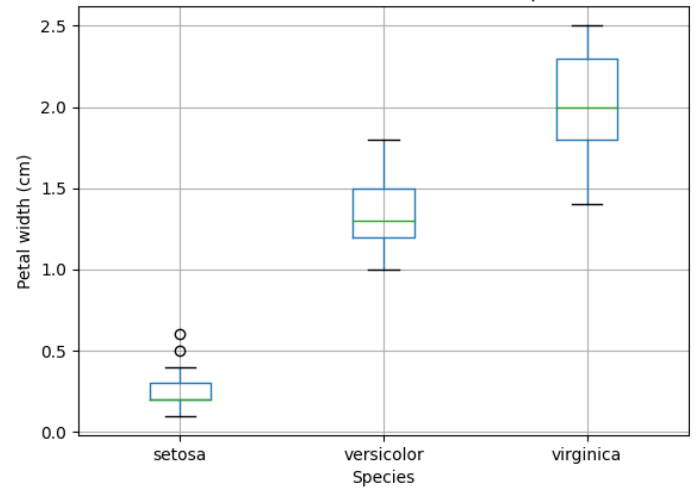
Distribution of Sepal Width of each Species



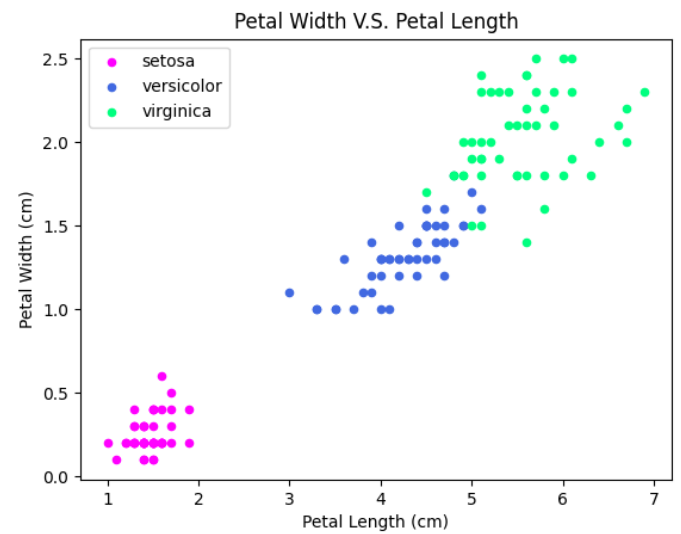
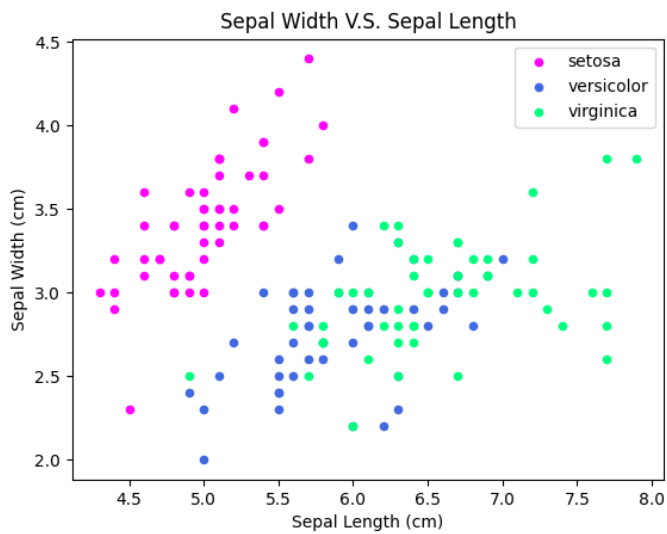
Distribution of Petal Length of each Species



Distribution of Petal Width of each Species



c)

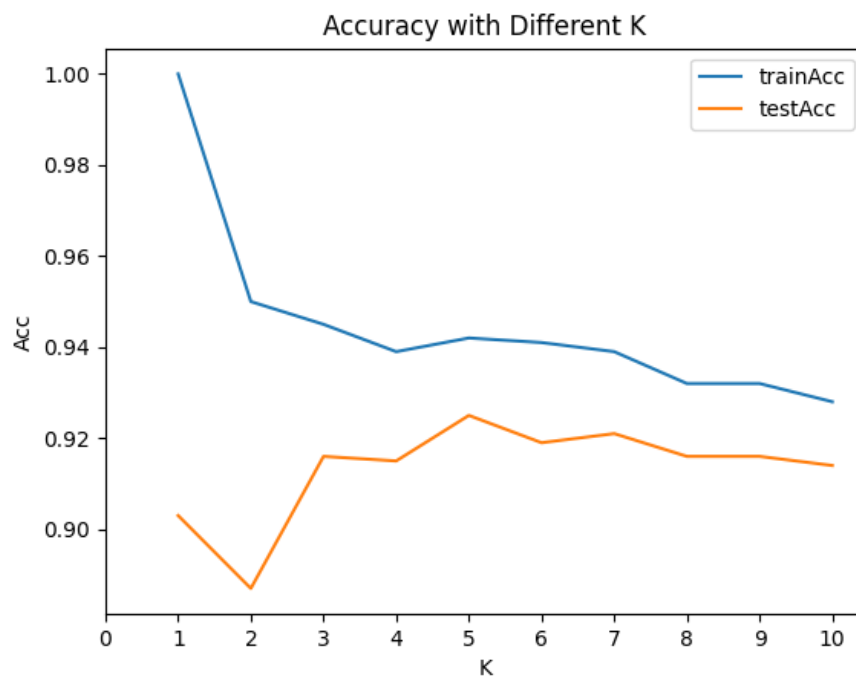


d) Rules to classify the species type:

1. If the petal length  $< 3\text{cm}$  and petal width  $< 1\text{cm}$ , the sample is setosa.
2. If  $3\text{cm} \leq \text{petal length} < 5\text{cm}$  and  $1\text{cm} \leq \text{petal width} < 1.75\text{cm}$ , the sample is versicolor
3. If petal length  $\geq 5\text{cm}$  and petal width  $\geq 1.75\text{cm}$ , the sample is virginica.

Q3:

d)



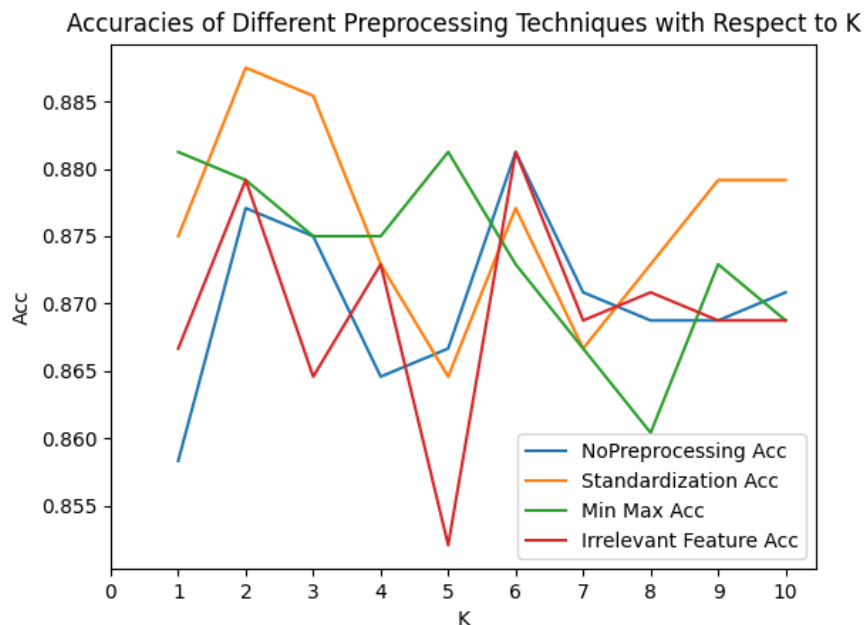
e)

The space complexity of the predict function is  $O(nd)$ , since we need to store all the training data, which is an  $n \times d$  array.

The time complexity of the predict function is  $O(n \times (nd + n \log n + k))$ . Assume the testing set is also size  $n$ . We iterate through every instance in testing set, and in each iteration, we compute the distance between the instance and every training sample, which uses  $O(nd)$ . Then, we run a quicksort on the distances which takes  $O(n \log n)$ , and find the first  $k$  distances using  $O(k)$ . Therefore, the time complexity is  $O(n \times (nd + n \log n + k))$ .

Q4:

d)



Both standard scale and min max range scale improve the performance of KNN compared to non-preprocessed data, especially when  $k$  is in ranges of 1 - 4 and 9 - 10.

This dataset is not very sensitive to irrelevant features. Since most of the time, dataset with irrelevant features (red line) shows a similar accuracy with non-preprocessed dataset (blue line).