

**University of California, Riverside**

**CS205-  
Artificial Intelligence**

**Winter 2024**

**CS205 Assignment 2: Feature Selection.**

**06/12/2024**

<b>Name</b>	<b>SID</b>	<b>Email</b>	<b>Group Number</b>
<b>Parth Shinde</b>	<b>862466119</b>	<b>pshin022@ucr.edu</b>	<b>-</b>

I consulted the following in completing this assignment:

- CS 205 Lecture slides on **Project-2 breifing** by Dr. Eamonn Keogh.
- I have used the following non-trivial libraries for doing the assignment:-
  - **NumPy** : Used for loading the dataset and performing array manipulations and calculations.
  - **Sklearn**: Used for implementing K-Nearest Neighbor algorithm, finding the accuracy score and the also evaluating/validating the algorithm using LeaveOneOut algorithm.
  - **Pandas**: Used in the helper scatter-plot generation script.
  - **Matplotlib**: Used for generating the plots.
- I have used other trivial python-in built libraries that are present in the python3.9.
- I also consulted the following links for understanding sklearn, Leave-one-out validator & KNN classifier implementation.
  - <https://machinelearningmastery.com/loocv-for-evaluating-machine-learning-algorithms/>
  - <https://python-course.eu/machine-learning/k-nearest-neighbor-classifier-with-sklearn.php>

List of contents:

- Cover page: (this page)
- Description of assignment: Pages 3-4
- Example Traceback: Pages 5-25
- Summary of results: Page 26-29
- Code & Conclusion: Pages 29-35

## **Introduction**

In this assignment we must implement a feature selection algorithm to find the most optimal feature set for the dataset considering, we use the K-nearest neighbor classifier as our baseline model. The two algorithms given to us to implement are Forward selection and Backward elimination.

## **Description of Datasets**

The datasets given to us have the label in the first column of the dataset and all the other columns are the features given to us. We load this dataset using the numpy library and separate out the features and labels and pass that out to our classifier. The table below shows the shape of the dataset and some information about the dataset.

<b>Dataset</b>	<b>No of Features</b>	<b>No of instances or rows</b>
Small Dataset(2 labels)	12	500
Large Dataset(2 labels)	50	5000

## **K-Nearest-Neighbor Algorithm**

K-Nearest Neighbors (KNN) is a simple, non-parametric, instance-based learning algorithm that stores training instances and makes predictions based on them at prediction time, using distance metrics like Euclidean or Manhattan to find the nearest neighbors. The parameter  $k$  determines the number of neighbors to consider: a small  $k$  captures local structure, while a large  $k$  provides more stability. For classification, KNN uses majority voting among the  $k$ -nearest neighbors to determine the class, and for regression, it predicts values based on the average or weighted average of the  $k$ -nearest neighbors. The algorithm involves storing all training data, calculating distances from the test instance to all training instances, sorting to select the  $k$ -nearest neighbors, and making predictions based on these neighbors.

## **Forward Selection**

It is a greedy algorithm, which has a top to bottom approach to solving the search problem. It starts with no features (empty set of features) and iteratively adds the feature that improves the accuracy the most, until no more improvements are found.

## **Backward Elimination**

Similar to forward selection, it is a greedy algorithm, but it's a bottom to up approach. It starts with all features and iteratively removes the feature whose removal improves the accuracy the most, until no more improvements are found.

## Algorithm

- A general feature search algorithm (the greedy problems) follow the algorithm given below
- Initialize a feature set (empty or full in our case)
- Initialize a the initial best accuracy of the model to a baseline value(e.g baseline accuracy could be :-using all the features given.)
- While the stopping condition is not met:
  - For each feature in the feature set to consider:
    - Modify the current feature set by adding or removing the feature.
    - Evaluate the model's performance (accuracy) using the modified feature set.
    - If the accuracy improves:
      - Update the best accuracy.
      - Update the best feature set.
      - (add the best feature or remove the worst feature in feature set).
  - If the best feature set did not change during this iteration:
    - Terminate the search.
- Return the best feature set and its corresponding accuracy.

## Timing Analysis

According to the algorithm given above, the Worst-Case Time Complexity would approximately be  $O(n^2 * m * T)$ . Considering (m,T) to be the validation and time complexity of Classifier, respectively i.e in our case the KNN algorithm.

## Tracebacks

The tracebacks are given as follows, and the user inputs are highlighted in yellow, and the final answer is green.

### **1. Small Data and Forward Selection**

```
Type in the name of the file to test below:-CS205_small_Data_34.txt
Welcome to Parth's Feature Selection Algorithm.
Type the number of the algorithm you want to run.
1.Forward Selection
```

## 2.Backward Elimination

1

This dataset has 12 features (not including the class attribute), with 500 instances.

Running nearest neighbor with all 12 features, using "leaving-one-out" evaluation, I get an accuracy of 73.4%

Beginning search.

Using feature(s) {1} accuracy is 68.8%

Using feature(s) {2} accuracy is 85.6%

Using feature(s) {3} accuracy is 69.0%

Using feature(s) {4} accuracy is 70.2%

Using feature(s) {5} accuracy is 69.8%

Using feature(s) {6} accuracy is 69.0%

Using feature(s) {7} accuracy is 68.8%

Using feature(s) {8} accuracy is 67.6%

Using feature(s) {9} accuracy is 72.0%

Using feature(s) {10} accuracy is 70.4%

Using feature(s) {11} accuracy is 71.0%

Using feature(s) {12} accuracy is 74.2%

Feature 2 added to the set ,as the resultant feature set has best accuracy 85.6%. Current feature set: {2}

Using feature(s) {1, 2} accuracy is 81.2%

Using feature(s) {2, 3} accuracy is 84.4%

Using feature(s) {2, 4} accuracy is 83.8%

Using feature(s) {2, 5} accuracy is 83.8%

Using feature(s) {2, 6} accuracy is 80.8%

Using feature(s) {2, 7} accuracy is 83.6%

Using feature(s) {8, 2} accuracy is 84.0%

Using feature(s) {9, 2} accuracy is 84.4%

Using feature(s) {2, 10} accuracy is 84.4%

Using feature(s) {2, 11} accuracy is 83.6%

Using feature(s) {2, 12} accuracy is 96.8%

Feature 12 added to the set ,as the resultant feature set has best accuracy 96.8%. Current feature set: {2, 12}

Using feature(s) {1, 2, 12} accuracy is 92.4%

Using feature(s) {2, 3, 12} accuracy is 93.0%

Using feature(s) {2, 12, 4} accuracy is 94.2%

Using feature(s) {2, 12, 5} accuracy is 94.2%

Using feature(s) {2, 12, 6} accuracy is 92.6%

Using feature(s) {2, 12, 7} accuracy is 93.2%

Using feature(s) {8, 2, 12} accuracy is 94.6%

Using feature(s) {9, 2, 12} accuracy is 92.8%

Using feature(s) {2, 10, 12} accuracy is 93.2%

Using feature(s) {2, 11, 12} accuracy is 94.0%

Excuetion time using forward selection for CS205\_small\_Data\_\_34.txt is:- 35.39018177986145

Graph saved as Forward Selection\_CS205\_small\_Data\_\_34.txt.pdf

Finished search!! The best feature subset is {2, 12}, which has an

accuracy of 96.8%**Finished search!! The best feature subset is {2, 12}, which has an accuracy of 96.8%**

## 2. Small Data & Backward Elimination

Type in the name of the file to test below:-CS205\_small\_Data\_\_34.txt

Welcome to Parth's Feature Selection Algorithm.

Type the number of the algorithm you want to run.

1. Forward Selection
2. Backward Elimination

2

This dataset has 12 features (not including the class attribute), with 500 instances.

Running nearest neighbor with all 12 features, using "leaving-one-out" evaluation, I get an accuracy of 73.4%

Beginning search.

Using feature(s) {2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12} accuracy is 73.6%  
Using feature(s) {1, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12} accuracy is 68.0%  
Using feature(s) {1, 2, 4, 5, 6, 7, 8, 9, 10, 11, 12} accuracy is 74.6%  
Using feature(s) {1, 2, 3, 5, 6, 7, 8, 9, 10, 11, 12} accuracy is 75.0%  
Using feature(s) {1, 2, 3, 4, 6, 7, 8, 9, 10, 11, 12} accuracy is 70.8%  
Using feature(s) {1, 2, 3, 4, 5, 7, 8, 9, 10, 11, 12} accuracy is 74.2%  
Using feature(s) {1, 2, 3, 4, 5, 6, 8, 9, 10, 11, 12} accuracy is 74.6%  
Using feature(s) {1, 2, 3, 4, 5, 6, 7, 9, 10, 11, 12} accuracy is 76.8%  
Using feature(s) {1, 2, 3, 4, 5, 6, 7, 8, 10, 11, 12} accuracy is 73.6%  
Using feature(s) {1, 2, 3, 4, 5, 6, 7, 8, 9, 11, 12} accuracy is 72.4%  
Using feature(s) {1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 12} accuracy is 72.6%  
Using feature(s) {1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11} accuracy is 70.6%

Feature 8 removed from the set, as the resultant feature set has best accuracy 76.8%. Current feature set: {1, 2,

5, 6, 7, 9, 10, 11, 12}

Using feature(s) {2, 3, 4, 5, 6, 7, 9, 10, 11, 12} accuracy is 77.2%  
Using feature(s) {1, 3, 4, 5, 6, 7, 9, 10, 11, 12} accuracy is 70.0%  
Using feature(s) {1, 2, 4, 5, 6, 7, 9, 10, 11, 12} accuracy is 76.0%  
Using feature(s) {1, 2, 3, 5, 6, 7, 9, 10, 11, 12} accuracy is 76.8%  
Using feature(s) {1, 2, 3, 4, 6, 7, 9, 10, 11, 12} accuracy is 74.4%  
Using feature(s) {1, 2, 3, 4, 5, 7, 9, 10, 11, 12} accuracy is 74.8%  
Using feature(s) {1, 2, 3, 4, 5, 6, 9, 10, 11, 12} accuracy is 75.8%  
Using feature(s) {1, 2, 3, 4, 5, 6, 7, 10, 11, 12} accuracy is 76.0%  
Using feature(s) {1, 2, 3, 4, 5, 6, 7, 9, 11, 12} accuracy is 75.8%  
Using feature(s) {1, 2, 3, 4, 5, 6, 7, 9, 10, 12} accuracy is 75.4%  
Using feature(s) {1, 2, 3, 4, 5, 6, 7, 9, 10, 11} accuracy is 71.8%

Feature 1 removed from the set, as the resultant feature set has best accuracy 77.2%. Current feature set: {2, 3, 4, 5, 6, 7, 9, 10, 11, 12}

Using feature(s) {3, 4, 5, 6, 7, 9, 10, 11, 12} accuracy is 71.0%  
Using feature(s) {2, 4, 5, 6, 7, 9, 10, 11, 12} accuracy is 79.0%  
Using feature(s) {2, 3, 5, 6, 7, 9, 10, 11, 12} accuracy is 79.4%  
Using feature(s) {2, 3, 4, 6, 7, 9, 10, 11, 12} accuracy is 78.0%  
Using feature(s) {2, 3, 4, 5, 7, 9, 10, 11, 12} accuracy is 77.4%  
Using feature(s) {2, 3, 4, 5, 6, 9, 10, 11, 12} accuracy is 79.0%  
Using feature(s) {2, 3, 4, 5, 6, 7, 10, 11, 12} accuracy is 76.2%  
Using feature(s) {2, 3, 4, 5, 6, 7, 9, 11, 12} accuracy is 76.0%  
Using feature(s) {2, 3, 4, 5, 6, 7, 9, 10, 12} accuracy is 76.2%  
Using feature(s) {2, 3, 4, 5, 6, 7, 9, 10, 11} accuracy is 72.8%

Feature 4 removed from the set, as the resultant feature set has best accuracy 79.4%. Current feature set: {2, 3, 5, 6, 7, 9, 10, 11, 12}

Using feature(s) {3, 5, 6, 7, 9, 10, 11, 12} accuracy is 69.2%  
Using feature(s) {2, 5, 6, 7, 9, 10, 11, 12} accuracy is 79.0%  
Using feature(s) {2, 3, 6, 7, 9, 10, 11, 12} accuracy is 80.4%  
Using feature(s) {2, 3, 5, 7, 9, 10, 11, 12} accuracy is 78.2%  
Using feature(s) {2, 3, 5, 6, 9, 10, 11, 12} accuracy is 78.0%  
Using feature(s) {2, 3, 5, 6, 7, 10, 11, 12} accuracy is 81.4%  
Using feature(s) {2, 3, 5, 6, 7, 9, 11, 12} accuracy is 80.2%

Using feature(s) {2, 3, 5, 6, 7, 9, 10, 12} accuracy is 78.0%

Using feature(s) {2, 3, 5, 6, 7, 9, 10, 11} accuracy is 74.8%

Feature 9 removed from the set,as the resultant feature set has best accuracy 81.4%. Current feature set: {2, 3, 5, 6, 7, 10, 11, 12}

Using feature(s) {3, 5, 6, 7, 10, 11, 12} accuracy is 72.2%

Using feature(s) {2, 5, 6, 7, 10, 11, 12} accuracy is 82.6%

Using feature(s) {2, 3, 6, 7, 10, 11, 12} accuracy is 81.4%

Using feature(s) {2, 3, 5, 7, 10, 11, 12} accuracy is 81.6%

Using feature(s) {2, 3, 5, 6, 10, 11, 12} accuracy is 82.4%

Using feature(s) {2, 3, 5, 6, 7, 11, 12} accuracy is 81.6%

Using feature(s) {2, 3, 5, 6, 7, 10, 12} accuracy is 80.4%

Using feature(s) {2, 3, 5, 6, 7, 10, 11} accuracy is 76.8%

Feature 3 removed from the set,as the resultant feature set has best accuracy 82.6%. Current feature set: {2, 5, 6, 7, 10, 11, 12}

Using feature(s) {5, 6, 7, 10, 11, 12} accuracy is 72.8%

Using feature(s) {2, 6, 7, 10, 11, 12} accuracy is 83.6%

Using feature(s) {2, 5, 7, 10, 11, 12} accuracy is 85.4%

Using feature(s) {2, 5, 6, 10, 11, 12} accuracy is 84.8%

Using feature(s) {2, 5, 6, 7, 11, 12} accuracy is 82.2%

Using feature(s) {2, 5, 6, 7, 10, 12} accuracy is 83.0%

Using feature(s) {2, 5, 6, 7, 10, 11} accuracy is 81.2%

Feature 6 removed from the set,as the resultant feature set has best accuracy 85.4%. Current feature set: {2, 5, 7, 10, 11, 12}

Using feature(s) {5, 7, 10, 11, 12} accuracy is 72.8%

Using feature(s) {2, 7, 10, 11, 12} accuracy is 84.8%

Using feature(s) {2, 5, 10, 11, 12} accuracy is 88.2%

Using feature(s) {2, 5, 7, 11, 12} accuracy is 88.6%

Using feature(s) {2, 5, 7, 10, 12} accuracy is 86.6%

Using feature(s) {2, 5, 7, 10, 11} accuracy is 80.6%

Feature 10 removed from the set,as the resultant feature set has best accuracy 88.6%. Current feature set: {2, 5, 7, 11, 12}

Using feature(s) {11, 12, 5, 7} accuracy is 74.8%

Using feature(s) {2, 11, 12, 7} accuracy is 88.8%

Using feature(s) {2, 11, 12, 5} accuracy is 90.6%

Using feature(s) {2, 12, 5, 7} accuracy is 91.8%

Using feature(s) {2, 11, 5, 7} accuracy is 82.8%

Feature 11 removed from the set,as the resultant feature set has best accuracy 91.8%. Current feature set: {2, 12, 5, 7}

Using feature(s) {12, 5, 7} accuracy is 75.6%

Using feature(s) {2, 12, 7} accuracy is 93.2%

Using feature(s) {2, 12, 5} accuracy is 94.2%

Using feature(s) {2, 5, 7} accuracy is 83.2%

Feature 7 removed from the set,as the resultant feature set has best accuracy 94.2%. Current feature set: {2, 12, 5}

Using feature(s) {12, 5} accuracy is 71.4%

Using feature(s) {2, 12} accuracy is 96.8%

Using feature(s) {2, 5} accuracy is 83.8%

Feature 5 removed from the set,as the resultant feature set has best accuracy 96.8%. Current feature set: {2, 12}

Using feature(s) {12} accuracy is 74.2%

Using feature(s) {2} accuracy is 85.6%

Execution time using backward elimination for CS205\_small\_Data\_\_34.txt is:131.12230706214905

Graph saved as Backward Elimination\_CS205\_small\_Data\_\_34.txt.pdf

Finished search!! The best feature subset is {2, 12}, which has an accuracy of 96.8%

### 3. Large Data & Forward Selection

Type in the name of the file to test below:-CS205\_large\_Data\_41.txt

Welcome to Parth's Feature Selection Algorithm.

Type the number of the algorithm you want to run.

- 1.Forward Selection
- 2.Backward Elimination

1

This dataset has 50 features (not including the class attribute), with 5000 instances.

Running nearest neighbor with all 50 features, using "leaving-one-out" evaluation, I get an accuracy of 70.0%

Beginning search.

Using feature(s) {1} accuracy is 70.1%  
Using feature(s) {2} accuracy is 71.3%  
Using feature(s) {3} accuracy is 70.4%  
Using feature(s) {4} accuracy is 71.3%  
Using feature(s) {5} accuracy is 72.5%  
Using feature(s) {6} accuracy is 71.4%  
Using feature(s) {7} accuracy is 70.3%  
Using feature(s) {8} accuracy is 70.2%  
Using feature(s) {9} accuracy is 71.7%  
Using feature(s) {10} accuracy is 70.6%  
Using feature(s) {11} accuracy is 70.8%  
Using feature(s) {12} accuracy is 71.5%  
Using feature(s) {13} accuracy is 70.7%  
Using feature(s) {14} accuracy is 69.9%  
Using feature(s) {15} accuracy is 70.4%  
Using feature(s) {16} accuracy is 70.9%  
Using feature(s) {17} accuracy is 70.1%  
Using feature(s) {18} accuracy is 70.5%  
Using feature(s) {19} accuracy is 70.9%  
Using feature(s) {20} accuracy is 70.0%  
Using feature(s) {21} accuracy is 85.0%  
Using feature(s) {22} accuracy is 69.4%  
Using feature(s) {23} accuracy is 71.2%  
Using feature(s) {24} accuracy is 71.5%  
Using feature(s) {25} accuracy is 71.1%  
Using feature(s) {26} accuracy is 70.1%  
Using feature(s) {27} accuracy is 71.6%  
Using feature(s) {28} accuracy is 71.3%  
Using feature(s) {29} accuracy is 70.6%  
Using feature(s) {30} accuracy is 71.6%  
Using feature(s) {31} accuracy is 70.4%  
Using feature(s) {32} accuracy is 71.6%  
Using feature(s) {33} accuracy is 70.2%  
Using feature(s) {34} accuracy is 70.7%  
Using feature(s) {35} accuracy is 69.9%  
Using feature(s) {36} accuracy is 70.3%  
Using feature(s) {37} accuracy is 71.2%  
Using feature(s) {38} accuracy is 70.0%  
Using feature(s) {39} accuracy is 70.7%  
Using feature(s) {40} accuracy is 71.7%



Using feature(s) {41} accuracy is 71.4%  
 Using feature(s) {42} accuracy is 71.0%  
 Using feature(s) {43} accuracy is 71.3%  
 Using feature(s) {44} accuracy is 71.1%  
 Using feature(s) {45} accuracy is 70.0%  
 Using feature(s) {46} accuracy is 70.0%  
 Using feature(s) {47} accuracy is 70.5%  
 Using feature(s) {48} accuracy is 71.2%  
 Using feature(s) {49} accuracy is 70.4%  
 Using feature(s) {50} accuracy is 73.7%  
 Feature 21 added to the set ,as the resultant feature set has best accuracy 85.0%. Current feature set: {21}  
 Using feature(s) {1, 21} accuracy is 83.5%  
 Using feature(s) {2, 21} accuracy is 83.4%  
 Using feature(s) {3, 21} accuracy is 84.3%  
 Using feature(s) {4, 21} accuracy is 84.7%  
 Using feature(s) {5, 21} accuracy is 84.4%  
 Using feature(s) {21, 6} accuracy is 83.7%  
 Using feature(s) {21, 7} accuracy is 85.1%  
 Using feature(s) {8, 21} accuracy is 83.4%  
 Using feature(s) {9, 21} accuracy is 84.9%  
 Using feature(s) {10, 21} accuracy is 84.0%  
 Using feature(s) {11, 21} accuracy is 84.9%  
 Using feature(s) {12, 21} accuracy is 83.7%  
 Using feature(s) {13, 21} accuracy is 84.9%  
 Using feature(s) {21, 14} accuracy is 84.2%  
 Using feature(s) {21, 15} accuracy is 84.5%  
 Using feature(s) {16, 21} accuracy is 83.9%  
 Using feature(s) {17, 21} accuracy is 83.9%  
 Using feature(s) {18, 21} accuracy is 83.6%  
 Using feature(s) {19, 21} accuracy is 84.6%  
 Using feature(s) {20, 21} accuracy is 83.5%  
 Using feature(s) {21, 22} accuracy is 83.5%  
 Using feature(s) {21, 23} accuracy is 85.1%  
 Using feature(s) {24, 21} accuracy is 84.4%  
 Using feature(s) {25, 21} accuracy is 83.0%  
 Using feature(s) {26, 21} accuracy is 84.4%  
 Using feature(s) {27, 21} accuracy is 84.9%  
 Using feature(s) {28, 21} accuracy is 83.6%  
 Using feature(s) {29, 21} accuracy is 84.4%  
 Using feature(s) {21, 30} accuracy is 83.6%  
 Using feature(s) {21, 31} accuracy is 82.5%  
 Using feature(s) {32, 21} accuracy is 85.2%  
 Using feature(s) {33, 21} accuracy is 85.3%  
 Using feature(s) {34, 21} accuracy is 84.6%  
 Using feature(s) {35, 21} accuracy is 84.2%  
 Using feature(s) {36, 21} accuracy is 84.1%  
 Using feature(s) {37, 21} accuracy is 83.4%  
 Using feature(s) {21, 38} accuracy is 83.9%  
 Using feature(s) {21, 39} accuracy is 83.6%  
 Using feature(s) {40, 21} accuracy is 83.4%  
 Using feature(s) {41, 21} accuracy is 84.3%  
 Using feature(s) {42, 21} accuracy is 85.7%  
 Using feature(s) {43, 21} accuracy is 83.1%  
 Using feature(s) {44, 21} accuracy is 84.0%  
 Using feature(s) {45, 21} accuracy is 84.5%

Using feature(s) {21, 46} accuracy is 84.2%  
 Using feature(s) {21, 47} accuracy is 84.7%  
 Using feature(s) {48, 21} accuracy is 84.1%  
 Using feature(s) {49, 21} accuracy is 84.3%  
 Using feature(s) {50, 21} accuracy is 97.3%  
 Feature 50 added to the set ,as the resultant feature set has best accuracy 97.3%. Current feature set: {50, 21}  
 Using feature(s) {1, 50, 21} accuracy is 95.3%  
 Using feature(s) {50, 2, 21} accuracy is 95.1%  
 Using feature(s) {50, 3, 21} accuracy is 95.1%  
 Using feature(s) {50, 4, 21} accuracy is 94.8%  
 Using feature(s) {50, 5, 21} accuracy is 95.6%  
 Using feature(s) {50, 21, 6} accuracy is 95.4%  
 Using feature(s) {50, 21, 7} accuracy is 95.2%  
 Using feature(s) {8, 50, 21} accuracy is 95.6%  
 Using feature(s) {9, 50, 21} accuracy is 95.9%  
 Using feature(s) {50, 10, 21} accuracy is 95.2%  
 Using feature(s) {50, 11, 21} accuracy is 95.1%  
 Using feature(s) {50, 12, 21} accuracy is 95.0%  
 Using feature(s) {50, 13, 21} accuracy is 94.8%  
 Using feature(s) {50, 21, 14} accuracy is 95.6%  
 Using feature(s) {50, 21, 15} accuracy is 95.4%  
 Using feature(s) {16, 50, 21} accuracy is 95.1%  
 Using feature(s) {17, 50, 21} accuracy is 95.2%  
 Using feature(s) {50, 18, 21} accuracy is 95.6%  
 Using feature(s) {50, 19, 21} accuracy is 95.3%  
 Using feature(s) {50, 20, 21} accuracy is 95.6%  
 Using feature(s) {50, 21, 22} accuracy is 95.1%  
 Using feature(s) {50, 21, 23} accuracy is 96.8%  
 Using feature(s) {24, 50, 21} accuracy is 94.9%  
 Using feature(s) {25, 50, 21} accuracy is 95.1%  
 Using feature(s) {50, 26, 21} accuracy is 95.7%  
 Using feature(s) {50, 27, 21} accuracy is 95.4%  
 Using feature(s) {50, 28, 21} accuracy is 94.9%  
 Using feature(s) {50, 29, 21} accuracy is 95.5%  
 Using feature(s) {50, 21, 30} accuracy is 95.7%  
 Using feature(s) {50, 21, 31} accuracy is 95.2%  
 Using feature(s) {32, 50, 21} accuracy is 95.2%  
 Using feature(s) {33, 50, 21} accuracy is 94.9%  
 Using feature(s) {50, 34, 21} accuracy is 95.3%  
 Using feature(s) {50, 35, 21} accuracy is 94.9%  
 Using feature(s) {50, 36, 21} accuracy is 95.3%  
 Using feature(s) {50, 37, 21} accuracy is 95.4%  
 Using feature(s) {50, 21, 38} accuracy is 95.3%  
 Using feature(s) {50, 21, 39} accuracy is 94.9%  
 Using feature(s) {40, 50, 21} accuracy is 95.5%  
 Using feature(s) {41, 50, 21} accuracy is 95.0%  
 Using feature(s) {50, 42, 21} accuracy is 95.3%  
 Using feature(s) {50, 43, 21} accuracy is 95.0%  
 Using feature(s) {50, 44, 21} accuracy is 95.1%  
 Using feature(s) {50, 45, 21} accuracy is 95.2%  
 Using feature(s) {50, 21, 46} accuracy is 94.8%  
 Using feature(s) {50, 21, 47} accuracy is 95.0%  
 Using feature(s) {48, 50, 21} accuracy is 95.9%  
 Using feature(s) {49, 50, 21} accuracy is 95.2%

Excution time using forward selection for CS205\_large\_Data\_\_41.txt is:-  
5115.121083021164  
Graph saved as Forward Selection CS205 large Data 41.txt.pdf  
Finished search!! The best feature subset is {50, 21}, which has an  
accuracy of 97.3%

#### 4. Large Data & Backward Elimination

Type in the name of the file to test below:-CS205\_large\_Data\_\_41.txt  
Welcome to Parth's Feature Selection Algorithm.  
Type the number of the algorithm you want to run.  
1.Forward Selection  
2.Backward Elimination

2

This dataset has 50 features (not including the class attribute), with 5000 instances.  
Running nearest neighbor with all 50 features, using "leaving-one-out" evaluation, I get an accuracy of 70.0%  
Beginning search.  
Using feature(s) {2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50} accuracy is 70.8%  
Using feature(s) {1, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50} accuracy is 69.8%  
Using feature(s) {1, 2, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50} accuracy is 70.2%  
Using feature(s) {1, 2, 3, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50} accuracy is 69.9%  
Using feature(s) {1, 2, 3, 4, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50} accuracy is 70.5%  
Using feature(s) {1, 2, 3, 4, 5, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50} accuracy is 70.5%  
Using feature(s) {1, 2, 3, 4, 5, 6, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50} accuracy is 70.2%  
Using feature(s) {1, 2, 3, 4, 5, 6, 7, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50} accuracy is 71.0%  
Using feature(s) {1, 2, 3, 4, 5, 6, 7, 8, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35,

[illegible]

[illegible]

35, 36, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50} accuracy is 69.9%

Using feature(s) {1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50} accuracy is 70.9%

Using feature(s) {1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50} accuracy is 70.4%

Using feature(s) {1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50} accuracy is 70.2%

Using feature(s) {1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 42, 43, 44, 45, 46, 47, 48, 49, 50} accuracy is 70.5%

Using feature(s) {1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 43, 44, 45, 46, 47, 48, 49, 50} accuracy is 70.4%

Using feature(s) {1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 44, 45, 46, 47, 48, 49, 50} accuracy is 70.2%

Using feature(s) {1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 45, 46, 47, 48, 49, 50} accuracy is 70.5%

Using feature(s) {1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 46, 47, 48, 49, 50} accuracy is 70.3%

Using feature(s) {1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 47, 48, 49, 50} accuracy is 70.3%

Using feature(s) {1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 48, 49, 50} accuracy is 71.0%

Using feature(s) {1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 49, 50} accuracy is 70.8%

Using feature(s) {1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 50} accuracy is 69.9%

Using feature(s) {1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49} accuracy is 71.5%

Feature 50 removed from the set, as the resultant feature set has best accuracy 71.5%. Current feature set: {1, 2, 3, 4,

[illegible]

[illegible]





[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]



```

Using feature(s) {1, 2, 3, 4, 5, 6, 7, 9, 10, 11, 12, 13, 14, 15, 16, 17,
18, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 36, 37,
38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 49} accuracy is 72.1%
Using feature(s) {1, 2, 3, 4, 5, 6, 7, 9, 10, 11, 12, 13, 14, 15, 16, 17,
18, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 37,
38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 49} accuracy is 71.8%
Using feature(s) {1, 2, 3, 4, 5, 6, 7, 9, 10, 11, 12, 13, 14, 15, 16, 17,
18, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36,
38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 49} accuracy is 72.3%
Using feature(s) {1, 2, 3, 4, 5, 6, 7, 9, 10, 11, 12, 13, 14, 15, 16, 17,
18, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36,
37, 39, 40, 41, 42, 43, 44, 45, 46, 47, 49} accuracy is 71.9%
Using feature(s) {1, 2, 3, 4, 5, 6, 7, 9, 10, 11, 12, 13, 14, 15, 16, 17,
18, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36,
37, 38, 40, 41, 42, 43, 44, 45, 46, 47, 49} accuracy is 71.9%
Using feature(s) {1, 2, 3, 4, 5, 6, 7, 9, 10, 11, 12, 13, 14, 15, 16, 17,
18, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36,
37, 38, 39, 41, 42, 43, 44, 45, 46, 47, 49} accuracy is 72.3%
Using feature(s) {1, 2, 3, 4, 5, 6, 7, 9, 10, 11, 12, 13, 14, 15, 16, 17,
18, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36,
37, 38, 39, 40, 42, 43, 44, 45, 46, 47, 49} accuracy is 71.7%
Using feature(s) {1, 2, 3, 4, 5, 6, 7, 9, 10, 11, 12, 13, 14, 15, 16, 17,
18, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36,
37, 38, 39, 40, 41, 43, 44, 45, 46, 47, 49} accuracy is 72.2%
Using feature(s) {1, 2, 3, 4, 5, 6, 7, 9, 10, 11, 12, 13, 14, 15, 16, 17,
18, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30}Using feature(s) {1, 2, 3,
4, 5, 6, 7, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 20, 21, 22, 23, 24, 25,
26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43,
45, 46, 47, 49} accuracy is 72.3%
Using feature(s) {1, 2, 3, 4, 5, 6, 7, 9, 10, 11, 12, 13, 14, 15, 16, 17,
18, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36,
37, 38, 39, 40, 41, 42, 43, 44, 46, 47, 49} accuracy is 71.5%
Using feature(s) {1, 2, 3, 4, 5, 6, 7, 9, 10, 11, 12, 13, 14, 15, 16, 17,
18, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36,
37, 38, 39, 40, 41, 42, 43, 44, 45, 47, 49} accuracy is 72.0%
Using feature(s) {1, 2, 3, 4, 5, 6, 7, 9, 10, 11, 12, 13, 14, 15, 16, 17,
18, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36,
37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 49} accuracy is 71.9%
Using feature(s) {1, 2, 3, 4, 5, 6, 7, 9, 10, 11, 12, 13, 14, 15, 16, 17,
18, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36,
37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47} accuracy is 71.8%
Execution time using backward elimination for CS205_large_Data__41.txt
is:4230.797879457474
c:\All Data\AI\CS205-Proj2\new_feature.py:131: UserWarning: Tight layout
not applied. The bottom and top margins cannot be made
large enough to accommodate all axes decorations.
  fig.tight_layout()
Graph saved as Backward Elimination CS205_large_Data__41.txt.pdf
Finished search!! The best feature subset is {1, 2, 3, 4, 5, 6, 7, 9, 10,
11, 12, 13, 14, 15, 16, 17, 18, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29,
30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47,
49}, which has an accuracy of 72.6%

```

## **Summary of results:**

I ran all the above combinations for the given dataset and have found the following results. The timing results of the various dataset and algorithm combinations are presented in the following table.

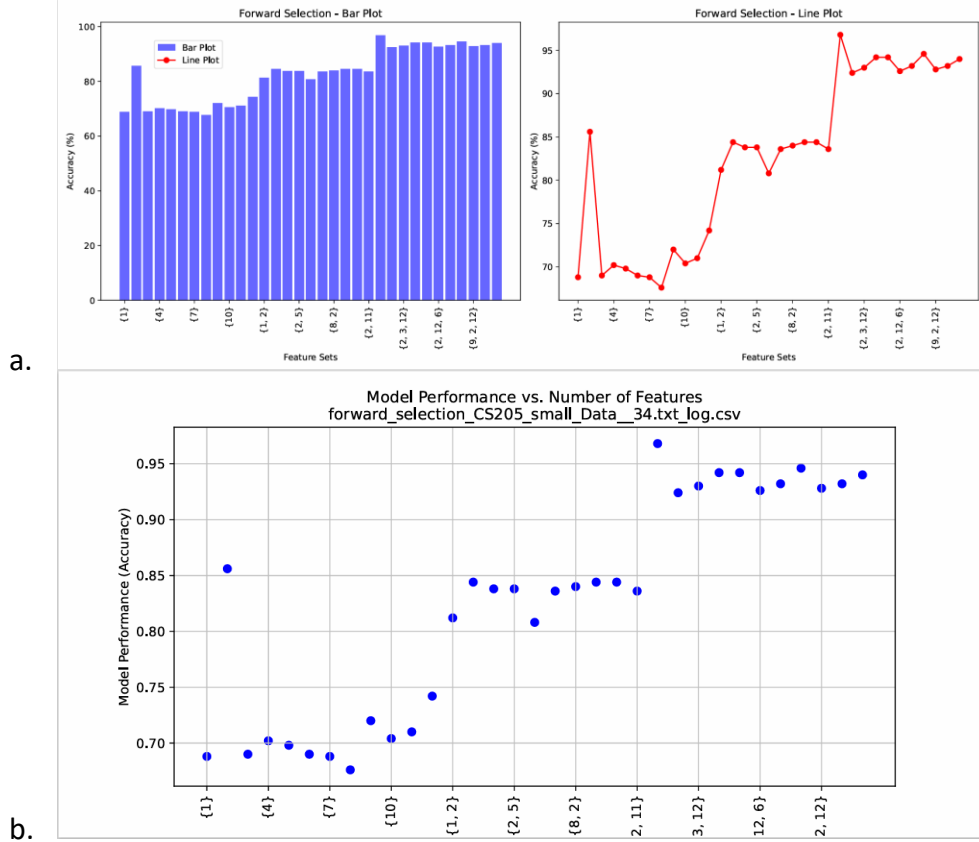
I have also generated graphs for each combination. The first graph for each combination is a mix of bar and line plot representations of the accuracy and feature set chosen. The second graph is a scatter plot showing the feature set versus accuracy.

The scripts are provided in this [github repo](#).

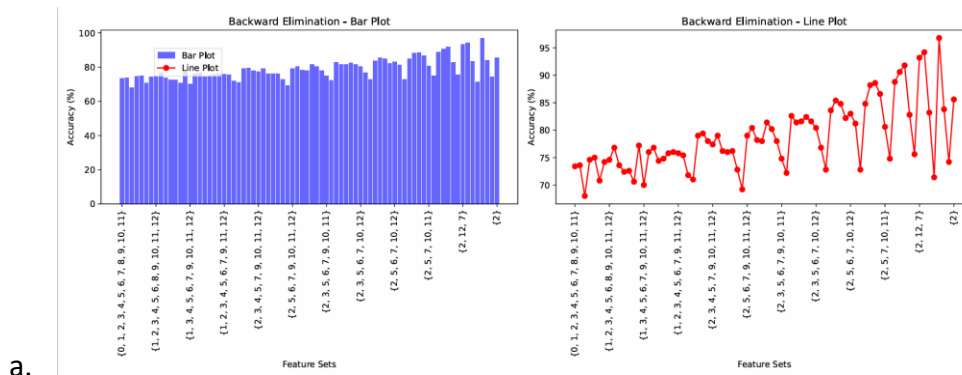
Timing Table:-

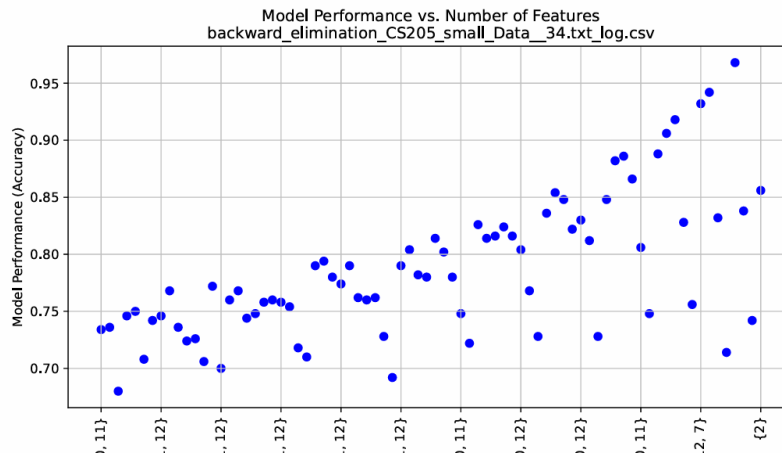
<b>Dataset</b>	<b>Feature selection technique</b>	<b>Excution time</b>
Small Dataset	Forward Selection	35 seconds
Small Dataset	Backward Elimination	131 seconds
Large Dataset	Forward Selection	5115 seconds(approx 1.42 hours)
Large Dataset	Backward Elimination	8023 seconds(approx 2.20 hours)

## 1. Small Data & Forward Selection.



## 2. Small Data & Backward Elimination

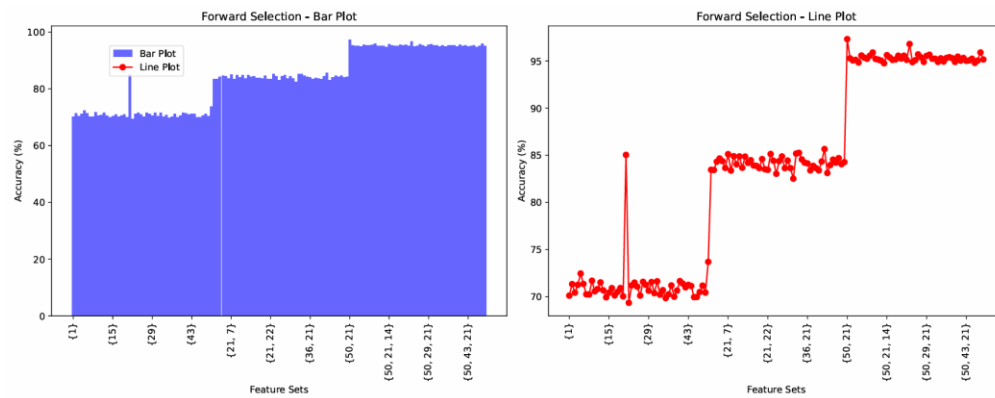




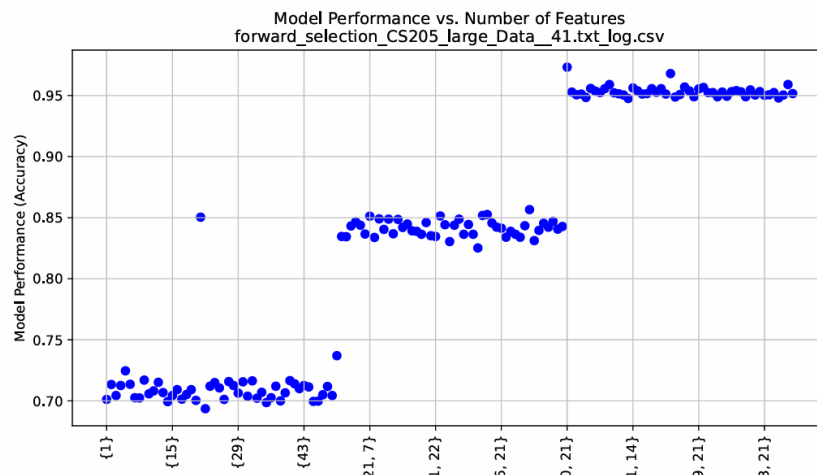
b.

### 3. Large Data & Forward Selection

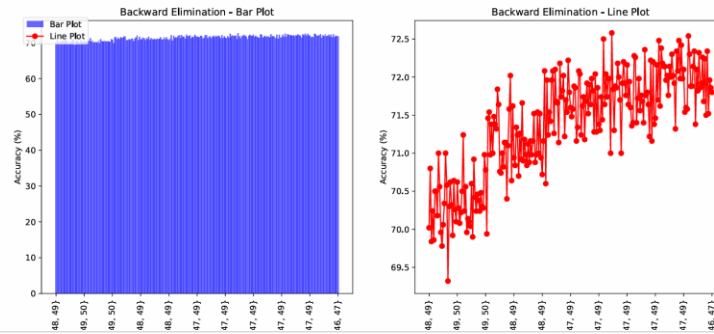
a.



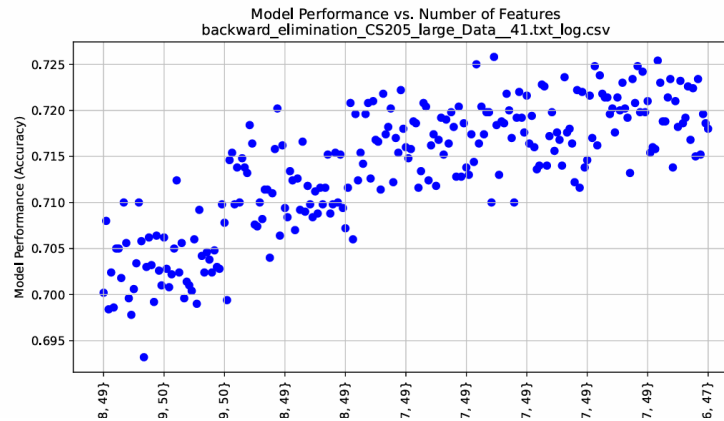
b.



### 4. Large Data & Backward Elimination



a.



b.

## Code Snippets

### 1.new\_feature.py

```
import numpy as np
import matplotlib.pyplot as plt
import csv
from sklearn.metrics import accuracy_score
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import LeaveOneOut
import time

# Load the dataset from the provided file path
def load_data(file_path):
    data = np.loadtxt(file_path)
    return data

# Separate the features and labels from the dataset
def separate_features_labels(data):
    # Features (excluding the first column which is the label)
    X = data[:, 1:]
    # Labels
    y = data[:, 0].astype(int)
```

```

    return X, y

# Perform k-nearest neighbor classification using Leave-One-Out cross-validation
def nearest_neighbor_loo(X, y, k=1):
    loo = LeaveOneOut()
    y_true, y_pred = [], []
    knn = KNeighborsClassifier(n_neighbors=k)

    for train_index, test_index in loo.split(X):
        X_train, X_test = X[train_index], X[test_index]
        y_train, y_test = y[train_index], y[test_index]
        knn.fit(X_train, y_train)
        y_pred.append(knn.predict(X_test)[0])
        y_true.append(y_test[0])

    return np.array(y_true), np.array(y_pred)

# Calculate the accuracy of the model
def calculate_accuracy(y_true, y_pred):
    return accuracy_score(y_true, y_pred)

# Forward Selection Algorithm with logging
def forward_selection(X, y, k=1):
    n_features = X.shape[1]
    selected_features = []
    remaining_features = list(range(n_features))
    best_accuracy = 0
    best_features = []
    log = []

    print("Beginning search.")

    while remaining_features:
        best_feature = None
        for feature in remaining_features:
            current_features = selected_features + [feature]
            X_subset = X[:, current_features]
            y_true, y_pred = nearest_neighbor_loo(X_subset, y, k)
            accuracy = calculate_accuracy(y_true, y_pred)
            print(f"Using feature(s) {set(f + 1 for f in current_features)}")
            accuracy is {accuracy*100:.1f}%")
            log.append((set(f + 1 for f in current_features), accuracy))
            if accuracy > best_accuracy:
                best_accuracy = accuracy
                best_feature = feature
                best_features = current_features
        if best_feature is not None:

```

```

        selected_features.append(best_feature)
        remaining_features.remove(best_feature)
        print(f"Feature {best_feature + 1} added to the set ,as the resultant
feature set has best accuracy {best_accuracy*100:.1f}%. Current feature set:
{set(f + 1 for f in selected_features)}")
    else:
        break

    return best_features, best_accuracy, log

# Backward Elimination Algorithm with logging
def backward_elimination(X, y, k=1):
    n_features = X.shape[1]
    selected_features = list(range(n_features))
    y_true, y_pred = nearest_neighbor_loo(X, y, k)
    best_accuracy = calculate_accuracy(y_true, y_pred)
    log = [(set(selected_features), best_accuracy)]

    print("Beginning search.")

    while len(selected_features) > 1:
        worst_feature = None
        for feature in selected_features:
            current_features = [f for f in selected_features if f != feature]
            X_subset = X[:, current_features]
            y_true, y_pred = nearest_neighbor_loo(X_subset, y, k)
            accuracy = calculate_accuracy(y_true, y_pred)
            print(f"Using feature(s) {set(f + 1 for f in current_features)}
accuracy is {accuracy*100:.1f}%")
            log.append((set(f + 1 for f in current_features), accuracy))
            if accuracy > best_accuracy:
                best_accuracy = accuracy
                worst_feature = feature
        if worst_feature is not None:
            selected_features.remove(worst_feature)
            print(f"Feature {worst_feature + 1} removed from the set,as the
resultant feature set has best accuracy {best_accuracy*100:.1f}%. Current feature
set: {set(f + 1 for f in selected_features)}")
        else:
            break

    return selected_features, best_accuracy, log

# plot accuracy
def plot_accuracies(log, title, file_path):
    feature_sets = [str(features) for features, acc in log]

```

```

accuracies = [acc * 100 for features, acc in log]
range_len = 10
x_ticks = range(len(feature_sets))
x_ticks_step = max(1, len(feature_sets) // range_len)
x_ticks_labels = [feature_sets[i] for i in range(0, len(feature_sets),
x_ticks_step)]
x_ticks_positions = [i for i in range(0, len(feature_sets), x_ticks_step)]

fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(15, 6))

ax1.bar(x_ticks, accuracies, color='b', alpha=0.6, label='Bar Plot')
ax1.set_xlabel('Feature Sets')
ax1.set_ylabel('Accuracy (%)')
ax1.set_xticks(x_ticks_positions)
ax1.set_xticklabels(x_ticks_labels, rotation=90)
ax1.set_title(f'{title} - Bar Plot')

ax2.plot(x_ticks, accuracies, color='r', marker='o', label='Line Plot')
ax2.set_xlabel('Feature Sets')
ax2.set_ylabel('Accuracy (%)')
ax2.set_xticks(x_ticks_positions)
ax2.set_xticklabels(x_ticks_labels, rotation=90)
ax2.set_title(f'{title} - Line Plot')

fig.tight_layout()
plt.suptitle(f'Accuracy of Different Feature Sets - {title}', y=1.05)
fig.legend(loc='upper left', bbox_to_anchor=(0.1,0.9))
# plt.show()
plt.savefig(f'{title}_{file_path}.pdf')
print(f'Graph saved as {title}_{file_path}.pdf')
plt.close()

# Function to save log data to CSV
def save_log_to_csv(log, filename):
    with open(filename, mode='w', newline='') as file:
        writer = csv.writer(file)
        writer.writerow(['Feature Set', 'Accuracy'])
        for features, acc in log:
            writer.writerow([str(features), acc])

# Main function to load data and perform feature selection
def main():
    file_path = input('Type in the name of the file to test below:-')
    # loading the data from the dataset using the load_data function
    data = load_data(file_path)
    # separating the features from the labels in the dataset.
    X, y = separate_features_labels(data)

```



```

n_instances, n_features = X.shape

print("Welcome to Parth's Feature Selection Algorithm.")
# print("Type in the name of the file to test :")
print("Type the number of the algorithm you want to run.")
print("\t 1.Forward Selection")
print("\t 2.Backward Elimination")

# creating a switch case like choice for the choice of the algorithm
algorithm_choice = int(input())

if algorithm_choice == 1:
    print(f"\nThis dataset has {n_features} features (not including the class
attribute), with {n_instances} instances.")
    # Running the intial nearest_neighbor using the Leave-One-Out evaluation
method
    y_true, y_pred = nearest_neighbor_loo(X, y)

    # calculating the accuracy of the intial predictions
    accuracy = calculate_accuracy(y_true, y_pred)

    print(f"Running nearest neighbor with all {n_features} features, using
“leaving-one-out” evaluation, I get an accuracy of {accuracy*100:.1f}%")
    start_time = time.time()

    # run the feature search algorithm
    best_features, best_accuracy, log = forward_selection(X, y)

    end_time = time.time()
    run_time = start_time - end_time
    print(f"Excution time using forward selection for {file_path}
is:{run_time}")

    # Save the output in the log file in a csv format
    save_log_to_csv(log, f'forward_selection_{file_path}_log.csv')
    # plot the accuracy to track the performance of the algorithm
    plot_accuracies(log, 'Forward Selection',file_path)

elif algorithm_choice == 2:
    print(f"\nThis dataset has {n_features} features (not including the class
attribute), with {n_instances} instances.")

    # Running the intial nearest_neighbor using the Leave-One-Out evaluation
method
    y_true, y_pred = nearest_neighbor_loo(X, y)

    # calculating the accuracy of the intial predictions

```

```

        accuracy = calculate_accuracy(y_true, y_pred)

        print(f"Running nearest neighbor with all {n_features} features, using
        "leaving-one-out" evaluation, I get an accuracy of {accuracy*100:.1f}%")

        start_time = time.time()

        # calculating the accuracy of the initial predictions
        best_features, best_accuracy, log = backward_elimination(X, y)

        end_time = time.time()
        run_time = end_time - start_time
        print(f"Execution time using backward elimination for {file_path}
        is:{run_time}")

        # Save the output in the log file in a csv format
        save_log_to_csv(log, f'backward_elimination_{file_path}_log.csv')
        # plot the accuracy to track the performance of the algorithm
        plot_accuracies(log, 'Backward Elimination',file_path)

    else:
        print("Invalid choice.")
        return

    print(f"Finished search!! The best feature subset is {set(f + 1 for f in
    best_features)}, which has an accuracy of {best_accuracy*100:.1f}%")

if __name__ == "__main__":
    main()

```

## 2.new\_scatter\_plot.py

```

import os
import matplotlib.pyplot as plt
import pandas as pd

# Function to generate and save plot from a CSV file
def generate_plot(file_name):
    df = pd.read_csv(file_name, encoding='utf-8-sig')

    # Extract the number of features and accuracies for plotting
    features = df['Feature Set']
    accuracies = df['Accuracy']

    # Set range length for x-ticks

```

```

range_len = 10
x_ticks = range(len(features))
x_ticks_step = max(1, len(features) // range_len)
x_ticks_labels = [features[i] for i in range(0, len(features), x_ticks_step)]
x_ticks_positions = [i for i in range(0, len(features), x_ticks_step)]

plt.figure(figsize=(10, 5))
plt.scatter(x_ticks, accuracies, color='blue', marker='o')
plt.xlabel('Number of Features')
plt.ylabel('Model Performance (Accuracy)')
plt.title(f'Model Performance vs. Number of Features\n{file_name}')
plt.xticks(x_ticks_positions, x_ticks_labels, rotation=90)
plt.grid(True)

# Save the plot as a PDF
output_file = f'{file_name}_scatter.pdf'
plt.savefig(output_file)
plt.close()

# Iterate over all CSV files in the current directory
for file_name in os.listdir('.'):
    if file_name.endswith('.csv'):
        generate_plot(file_name)

```

## Conclusion

- Both algorithms are complete but not optimal, as they gave similar results for the small dataset but different ones for the larger one. The backward elimination algorithm produced a worse feature set than forward selection. There may be more optimal solutions, but these were not found due to the greedy nature of the algorithms, which search for local maxima (in terms of accuracy) rather than global maxima.
- Forward selection converged faster in my case, but this may differ in other cases, considering both algorithms are greedy.
- A Bi-directional approach (using both forward selection and backward elimination) may give better results.
- The results overall seem intuitive based on the content discussed in the lectures.

