

Project 2

Jayden Guan

Luofeng Xu

Feature Selection Analysis Using Forward Selection and Backward Elimination

https://github.com/yanjun934523/CS205_Project2

Introduction

Feature selection plays a crucial role in machine learning and data analysis tasks as it helps identify the most relevant features that contribute to the predictive accuracy of a model. In this report, we analyze the performance of two popular feature selection algorithms: Forward Selection and Backward Elimination. These algorithms aim to find the optimal subset of features that maximizes the accuracy of a given model.

Methodology

Forward Selection starts with an empty set of features and iteratively adds the most predictive feature at each step until a stopping criterion is met. It evaluates the performance of the model after adding each feature and selects the one that improves the model the most.

Backward Elimination begins with the full set of features and iteratively removes the least informative feature at each step until a stopping criterion is satisfied. It evaluates the performance of the model after removing each feature and eliminates the one that has the least impact on the model's accuracy.

In the implementation of the code, we use the numpy library to calculate the distance between samples instead of using a for loop. By using numpy to calculate the distance between samples, we are taking advantage of its vectorized operations, which can significantly improve the efficiency of our code compared to using a for loop.

We conducted experiments using four datasets: CS170_small_Data__3.txt, CS170_large_Data__4.txt, CS170_XXXlarge_Data__7.txt, and the cancer dataset (a real world dataset on Kaggle).

<https://www.kaggle.com/datasets/erdemtaha/cancer-data>).

The cancer dataset contains 570 cancer cells and 30 features to determine

whether the cancer cells are benign or malignant. There are 2 types of cancers:
1. benign cancer (B) and 2. malignant cancer (M).

For each dataset, we applied both Forward Selection and Backward Elimination techniques to identify the most informative features.

Results

1. CS170_small_Data__3.txt:
 - Forward Selection:
Final Selected 2 Features: [7, 3]
Best Accuracy: 98.0%
Time Cost: 0.30 s
 - Backward Elimination:
Final Selected 2 Features: [3, 7]
Best Accuracy: 98.0%
Time Cost: 0.35 s
2. CS170_large_Data__4.txt:
 - Forward Selection:
Final Selected 2 Features: [9, 1]
Best Accuracy: 97.50%
Time Cost: 4.355 s
 - Backward Elimination:
Final Selected 14 Features: [0, 1, 2, 3, 4, 7, 8, 9, 10, 11, 12, 13, 14, 16]
Best Accuracy: 80.50%
Time Cost: 3.772 s
3. CS170_XXXlarge_Data__7.txt:
 - Forward Selection:
Final Selected 2 Features: [2, 7]
Best Accuracy: 97.75%
Time Cost: 1573 s
 - Backward Elimination:
Final Selected many Features: [1, 2, 3, 4, 5, 6, 8, 9, 10, 11, 12, 15, 16, 18, 19, 20, 22, 24, 25, 26, 27, 28, 29, 30, 31, 32, 34, 35, 37, 38, 39, 40, 41, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 55, 56, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 70, 72, 73, 74, 75, 76, 77, 78, 79]
Best Accuracy: 79.50%
Time Cost: 1296 s
4. Cancer Dataset:
 - Forward Selection:
Final Selected 9 Features: [23, 4, 26, 1, 3, 17, 7, 0, 19]
Best Accuracy: 99.12%

- Backward Elimination:
Final Selected **10** Features: [1, 17, 20, 21, 22, 23, 24, 25, 26, 29]
Best Accuracy: **99.12%**

Computational effort for search: We show the running time of forward selection and backward elimination on the three datasets in Table 1. We implemented the search in Python and numpy. We use a desk PC, with an Intel Core i5-8400 and 16 gigs of main memory.

Table 1

	CS170_small_Data__3 10 features 800 instances	CS170_large_Data__4 20 features 1600 instances	CS170_XXXlarge_Data__7 80 features 3200 instances
Forward selection	0.30 s	4.355 s	0.437 h
Backward elimination	0.35 s	3.772 s	0.360 h

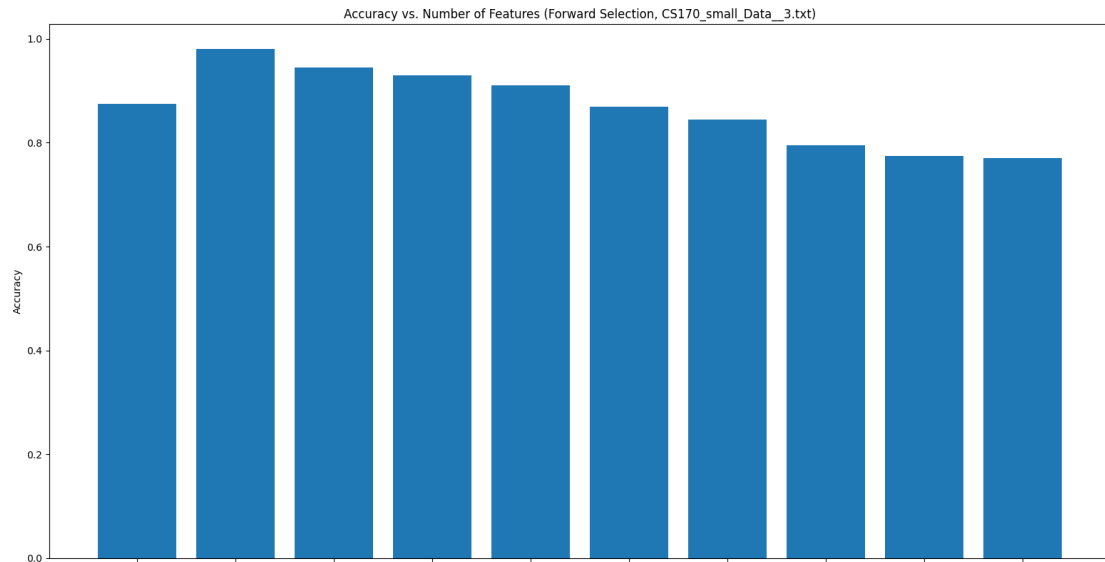


Figure 1

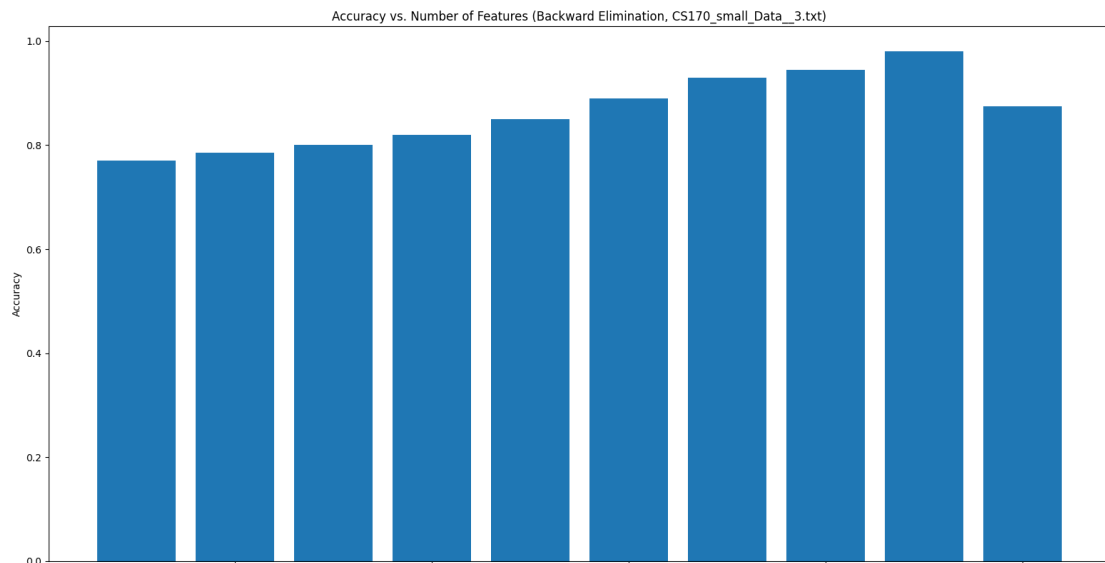


Figure 2

Figures 1 and 2 show the process diagrams of forward selection and back estimation. Figure 1 shows that during the forward selection process, after selecting two features (feature [7, 3]), the accuracy reached its maximum and then continuously decreased. Figure 2 shows the opposite performance of feedback estimation, with accuracy increasing as the number of discarded features increases, reaching its maximum when only two features remain.

Analyze

1. In CS170_small_Data__3.txt, forward selection and backward elimination select the same 2 features and get a high accuracy (98.0%).
2. In CS170_large_Data__4.txt and CS170_XXLarge_Data__7.txt, forward selection selects fewer features and gets higher accuracy (97.50%/ 97.75%) than backward elimination.
3. In cancer dataset, forward selection and backward elimination get the same high accuracy of 99.12%. However, forward selection still uses fewer features.

So, forward selection tends to select fewer features in our experiments and gets higher accuracy. Forward selection evaluates the performance of the model after adding each feature and selects the one that brings the most improvement in accuracy. This evaluation process allows forward selection to identify the features that contribute the most to accurate predictions. By focusing on the most impactful features, it can achieve high accuracy with a reduced number of features.

In contrast, backward elimination starts with a full set of features and eliminates one feature at a time based on its impact on the model's accuracy. This approach may result in a larger initial feature set and require more iterations to

arrive at the optimal subset of features. Overall, the iterative and evaluation-driven nature of forward selection enables it to identify a more concise and accurate subset of features compared to backward elimination. However, it's important to note that the performance of feature selection algorithms can vary depending on the dataset and the specific characteristics of the features themselves.

Conclusion

Our analysis demonstrates the effectiveness of forward selection in selecting influential features while achieving high accuracy. Forward selection consistently outperformed backward elimination by selecting fewer features and maintaining the same or higher accuracy across different datasets. This can be attributed to its iterative approach of evaluating each added feature's impact on the model's performance. Overall, forward selection proves to be a reliable feature selection algorithm for predictive modeling tasks.

Code

cancer.py

```
1 import numpy as np
2 import pandas as pd
3
4 def euclidean_distance(x1, x2):
5     return np.sqrt(((x1 - x2) ** 2).sum(axis=1))
6
7 def nearest_neighbor(train_X, test_X):
8     distance = euclidean_distance(train_X, test_X)
9     idx = distance.argmin()
10    return idx
11
12 def forward_selection(train_X, test_X, train_y, test_y):
13     num_features = train_X.shape[1]
14     selected_features = []
15     best_accuracy = 0.0
16     best_features = None
17
18     for _ in range(num_features):
19         feature_accuracies = []
20
21         for feature in range(num_features):
22             if feature not in selected_features:
23                 selected_features.append(feature)
24                 accuracy = evaluate(train_X[:, selected_features], test_X[:, selected_features], train_y, test_y)
25                 feature_accuracies.append((feature, accuracy))
26                 selected_features.remove(feature)
27
28         best_feature, new_accuracy = max(feature_accuracies, key=lambda x: x[1])
29         selected_features.append(best_feature)
30
31         if new_accuracy > best_accuracy:
32             best_accuracy = new_accuracy
33             best_features = selected_features.copy()
34
35         print(f"Selected Feature: {best_feature}, Current Accuracy: {new_accuracy:.2%}, Current Features: {selected_features}")
36
37     return best_features, best_accuracy
38
39 def backward_elimination(train_X, test_X, train_y, test_y):
40     num_features = train_X.shape[1]
41     selected_features = list(range(num_features))
42
43     best_accuracy = evaluate(train_X[:, selected_features], test_X[:, selected_features], train_y, test_y)
44     best_features = selected_features.copy()
45
46     while len(selected_features) > 0:
47         feature_accuracies = []
48
49         for feature in selected_features:
50             current_features = selected_features.copy()
51             current_features.remove(feature)
52             accuracy = evaluate(train_X[:, current_features], test_X[:, current_features], train_y, test_y)
53             feature_accuracies.append((feature, accuracy))
54
55         worst_feature, new_accuracy = max(feature_accuracies, key=lambda x: x[1])
56
57         if new_accuracy >= best_accuracy:
58             selected_features.remove(worst_feature)
59             best_accuracy = new_accuracy
60             best_features = selected_features.copy()
61             print(f"Removed Feature: {worst_feature}, Current Accuracy: {new_accuracy:.2%}, Current Features: {selected_features}")
62         else:
63             break
64
65     return best_features, best_accuracy
66
67
```

```

67
68 def evaluate(train_X, test_X, train_y, test_y):
69     correct_predictions = 0
70
71     for test_X, test_y in zip(test_X, test_y):
72         predicted_label = train_y[nearest_neighbor(train_X, test_X)]
73         if predicted_label == test_y:
74             correct_predictions += 1
75
76     return correct_predictions / len(test_X)
77
78 # Function to load the data from a file
79 def load_data():
80     df = pd.read_csv('cancer.csv')
81     print(df.info())
82     print(df.shape)
83     return df
84
85 # User Interface
86 data = load_data()
87 data_arr = np.array(data)
88 y = data_arr[:, 1]
89
90 def normalize_features(data):
91     min_vals = np.min(data, axis=0)
92     max_vals = np.max(data, axis=0)
93     normalized_data = (data - min_vals) / (max_vals - min_vals)
94     return normalized_data
95
96 X = data_arr[:, 2:].astype(np.float64)
97 X = normalize_features(X)
98 train_X = X[:int(len(X)*0.8), :]
99 train_y = y[:int(len(y)*0.8)]
100 test_X = X[int(len(X)*0.8):, :]
101 test_y = y[int(len(y)*0.8):]
102
103 search_method = input("Select the search method (1: Forward Selection, 2: Backward Elimination): ")
104
105 if search_method == '1':
106     # Forward Selection feature search
107     selected_features, best_accuracy = forward_selection(train_X, test_X, train_y, test_y)
108     print("Final Selected Features:", selected_features)
109     print(f"Best Accuracy: {best_accuracy:.2%}")
110 elif search_method == '2':
111     # Backward Elimination feature search
112     selected_features, best_accuracy = backward_elimination(train_X, test_X, train_y, test_y)
113     print("Final Selected Features:", selected_features)
114     for feature in selected_features:
115         print(data.iloc[:, feature+2].name)
116     print(f"Best Accuracy: {best_accuracy:.2%}")
117 else:
118     print("Invalid search method selected. Please choose 1 or 2.")
119

```

Main.py

```
1 import numpy as np
2 import time
3
4 def euclidean_distance(x1, x2):
5     return np.sqrt(((x1 - x2) ** 2).sum(axis=1))
6
7 def nearest_neighbor(train_data, test_instance):
8     distance = euclidean_distance(train_data, test_instance)
9     idx = distance.argmin()
10    return idx
11
12 def forward_selection(train_data, test_data):
13     num_features = train_data.shape[1] - 1
14     selected_features = []
15     best_accuracy = 0.0
16     best_features = None
17     train_label = train_data[:, 0]
18     train_data = train_data[:, 1:]
19     test_label = test_data[:, 0]
20     test_data = test_data[:, 1:]
21
22     for _ in range(num_features):
23         feature_accuracies = []
24
25         for feature in range(num_features):
26             if feature not in selected_features:
27                 selected_features.append(feature)
28                 accuracy = evaluate(train_data[:, selected_features], test_data[:, selected_features], train_label, test_label)
29                 feature_accuracies.append((feature, accuracy))
30                 selected_features.remove(feature)
31
32         best_feature, new_accuracy = max(feature_accuracies, key=lambda x: x[1])
33         selected_features.append(best_feature)
34
35         if new_accuracy > best_accuracy:
36             best_accuracy = new_accuracy
37             best_features = selected_features.copy()
38
39         print(f"Selected Feature: {best_feature}, Current Accuracy: {new_accuracy:.2%}, Current Features: {selected_features}")
40
41     return best_features, best_accuracy
42
43 def backward_elimination(train_data, test_data):
44     num_features = train_data.shape[1] - 1
45     selected_features = list(range(num_features))
46     train_label = train_data[:, 0]
47     train_data = train_data[:, 1:]
48     test_label = test_data[:, 0]
49     test_data = test_data[:, 1:]
50     best_accuracy = evaluate(train_data[:, selected_features], test_data[:, selected_features], train_label, test_label)
51     best_features = selected_features.copy()
52
53     while len(selected_features) > 0:
54         feature_accuracies = []
55
56         for feature in selected_features:
57             current_features = selected_features.copy()
58             current_features.remove(feature)
59             accuracy = evaluate(train_data[:, current_features], test_data[:, current_features], train_label, test_label)
60             feature_accuracies.append((feature, accuracy))
```



```

61
62     worst_feature, new_accuracy = max(feature_accuracies, key=lambda x: x[1])
63
64     if new_accuracy >= best_accuracy:
65         selected_features.remove(worst_feature)
66         best_accuracy = new_accuracy
67         best_features = selected_features.copy()
68         print(f'Removed Feature: {worst_feature}, Current Accuracy: {new_accuracy:.2%}, Current Features: {selected_features}')
69     else:
70         break
71
72     return best_features, best_accuracy
73
74
75 def evaluate(train_data, test_data, train_label, test_label):
76     correct_predictions = 0
77
78     for test_data_, test_label_ in zip(test_data, test_label):
79         predicted_label = train_label[nearest_neighbor(train_data, test_data_)]
80         if predicted_label == test_label_:
81             correct_predictions += 1
82
83     return correct_predictions / len(test_data)
84
85 # Function to load the data from a file
86 def load_data(filename):
87     data = np.loadtxt(filename)
88     return data
89
90 # User Interface
91 filename = input("Enter the file name: ")
92 data = load_data(filename)
93
94 train_data = data[:int(len(data)*0.8)]
95 test_data = data[int(len(data)*0.8):]
96
97 search_method = input("Select the search method (1: Forward Selection, 2: Backward Elimination): ")
98
99 if search_method == '1':
100     # Forward Selection feature search
101     start_time = time.time()
102     selected_features, best_accuracy = forward_selection(train_data, test_data)
103     end_time = time.time()
104     time_cost = end_time - start_time
105     print("Final Selected Features:", selected_features)
106     print(f"Best Accuracy: {best_accuracy:.2%}")
107     print(f"Time Cost: {time_cost:.4f} s")
108 elif search_method == '2':
109     # Backward Elimination feature search
110     start_time = time.time()
111     selected_features, best_accuracy = backward_elimination(train_data, test_data)
112     end_time = time.time()
113     time_cost = end_time - start_time
114     print("Final Selected Features:", selected_features)
115     print(f"Best Accuracy: {best_accuracy:.2%}")
116     print(f"Time Cost: {time_cost:.4f} s")
117 else:
118     print("Invalid search method selected. Please choose 1 or 2.")
119

```