

Feature Selection Analysis Using Forward Selection and Backward Elimination

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%

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_XXXlarge_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.