

## Project 2

Student Name: Yijian chai  
Linqing Li

Email: (Linqing LI) [lli260@ucr.edu](mailto:lli260@ucr.edu)

Date: June 4th, 2023

### Local Environment:

CPU: 11 th Gen Intel(R) Core i7-11700F

RAM: 32.0GB

System: windows 10

software platform: JupyterNotebook

### The resource I consulted:

(1) Forward Selection and Backward Elimination in KNN, also some approaches to optimize the scale of search space.

(2) Python 3.10

Code Link: <https://github.com/Lanceart/ai-proj/blob/main/proj2.ipynb>

### Index of this report:

- Cover page
- 1. Introduction
- 2. The result
- 3. Brief code introduction
- 4. Conclusion
- Appendix
- Citation

## 1. Introduction the algorithm

Backward elimination and forward selection are feature selection methods commonly used in supervised learning tasks such as nearest neighbor classification.

### (1) Backward Elimination:

Backward elimination is a feature selection method that starts from all available features. The idea is to start with the full feature set, and then iteratively remove the least important features (features that improve the model the least), one by one, until the desired number of features is reached, or the model does not improve significantly. The backward elimination algorithm is as follows:

- a) start.
- b) Fit the model and evaluate its performance.
- c) Remove the least significant feature
- d) Refit the model with the remaining features and re-evaluate its performance.
- e) Repeat steps c) and d)

### (2) Forward Selection:

Forward selection starts with no features and adds them one by one. At each iteration, the features that most improve the model are added to the feature set. This process continues until adding new features does not significantly improve the model. The forward selection algorithm is as follows:

- a) start.
- b) For each feature not already in the set, add it, fit the model, and evaluate its performance.
- c) Add the features that lead to the best performance boost to the set.
- d) Repeat steps b) and c) until adding new features does not significantly improve the model.

## 2. The result

Yijian (18/08/1999)

Linqing(31/01/2000)

So we should choose CS170\_small\_Data\_31.txt, CS170\_large\_Data\_18.txt, and CS170\_XXXlarge\_Data\_9.txt

(1)small : in the small part we use the unoptimized version code, you can see the result in Appendix.

The best feature subset is: [1, 6] which has an accuracy of 98.22 %

(2)Big: we use sample data, minimax like pruning, and stop in advanced.

The best feature subset is: [1, 7] which has an accuracy of 96.66 %

(3)XXXBig: we use the same approaches as (2)

Also for Backforward Elimination we use early abandon

The best feature subset is: [1, 74] which has an accuracy of 97.80 %

(4)Processing a real-world classification dataset

We did forward search on Diabetes Dataset, which is from UCI in 1998 which is the paper in Proceedings of the Symposium on Computer Applications and Medical Care (pp. 261--265).

In this data sets, it has 8 features and 691 instance.

Feature information:

Pregnancies: To express the Number of pregnancies

Glucose: To express the Glucose level in blood

BloodPressure: To express the Blood pressure measurement

SkinThickness: To express the thickness of the skin

Insulin: To express the Insulin level in blood

BMI: To express the Body mass index

DiabetesPedigreeFunction: To express the Diabetes percentage

Age: To express the age

The output is 0 or 1 means negative or positive.

So When want to analysis this code the first thing is clean the data, because the data form from Kaggle is kind difference from our data.

The final result analysis shows that diabetes has a greater correlation with Pregnancies, Glucose, and BMI indicators. It is easy to analyze, because gestational diabetes is a type of diabetes, and blood sugar and BMI indicators are closely related to diabetes symptoms.

### 3. Brief code introduction

#### (1) Data Cleaning

This part is a data cleaning for the last task in this project. Because most Kaggle dataset is CSV format, and typically, there is some annotation in first line. So we need to delete. And alignment the data, make the first column is the classification goal index.

```
from numpy import genfromtxt
# script for
files = input("Type in the name of the file to test: ")
my_data = genfromtxt(files, delimiter=',')
print(my_data[0][0])
my_data = np.delete(my_data,0,axis= 0)

print(my_data)

last_column = my_data[:, -1].reshape(-1, 1) # convert into vector
# delete the last
other_columns = my_data[:, :-1]
# put them into together
new_data = np.hstack((last_column, other_columns))
np.savetxt('new_data.txt',new_data)
```

#### (2) Forward

As the follow image show, there are two rectangle. I have to say they are the core of this algorithm, the red one is CrossValidation part, and you may notice there are two CrossValidation function in my code, because they will be used to stop early. In the CrossValidation\_ML function there is a tolerate errors counter. If programming detect there is much more errors than the tolerate errors we expect, it will brake the programm, and it really save much time for us.

And for the Blue rectangle, we can see a silde\_winows. Based on the early stop we design this part, because we want to stop when the accuracy gradually turns down. Stop instead of useless counting.

```
def ForwardSelection(data_nor,NF):
    print("Beginning search.\n")
    current_features = []
    final_acc = 0
    best_feature=[]
    my_windows = [0 for _ in range(NF)]
    counter = 0
    for i in range(1, NF+1):
        print("\n On level %d of the search tree" % (i),"contains", current_features)
        feature_select = 0
        cur_acc=0.0
        for j in range(1, NF+1):
            if j not in current_features:
                if j == 1:
                    acc = CrossValidation(current_features,data_nor,j,1)
                else:
                    acc = CrossValidation_ML(current_features,data_nor,j,1,(int) ((1-cur_acc) * len(data_nor) ))
                if acc> cur_acc:
                    cur_acc = acc
                    feature_select = j
        current_features.append(feature_select)
        print("\n On level %d of the search tree," % (i),"feature %d was added to the current set" % (feature_select))
        print("\n With ", len(current_features), " features, the accuracy is: ", cur_acc * 100, "%")

        my_windows[counter] = cur_acc
        if counter > 2 and Slide_window(my_windows, counter):
            break
        counter += 1
        if cur_acc >= final_acc:
            final_acc= cur_acc
            best_feature = list(current_features)

    print()
    print("Finish search!! The best feature subset is:", best_feature,"which has an accuracy of", final_acc * 100, "%")
```

```
#Performing Leave One Out Cross Validation and minimax like pruning
def CrossValidation_ML(features_in,data_nor,select_f, p, tolerate_error):
    error_counter = 0
    feature=list(features_in)
    if p==1:
        feature.append(select_f)
    if p ==2:
        feature.remove(select_f)
    count = 0
    dis = np.inf
    final=0
    for i in range(0,len(data_nor)):
        dis = np.inf
        for k in range(0,len(data_nor)):
            if not np.array_equal(k,i):
                ss = 0
                for j in range(0, len(feature)):#Calculating Euclidean Distance to determine the nearest neighbors
                    ss=ss+pow((data_nor[k][feature[j]] - data_nor[i][feature[j]]),2)
                d=np.sqrt(ss)
                if d < dis:
                    dis = d
                    final = k
            if (data_nor[final][0]==data_nor[i][0]):
                count += 1
            accuracy = (count / (len(data_nor)-1))#Calculating accuracy
        else:
            error_counter += 1

    if tolerate_error < error_counter:
        print("Using feature(s), the feature has been stop by minimax like pruning")
        return 0
    print("Using feature(s)",feature, "accuracy is", accuracy*100, "%")
    return accuracy
```

```
def Slide_window(my_window,layers):#Stop instead of useless counting
    return (my_window[layers-2] > my_window[layers-1]) and (my_window[layers-1] > my_window[layers])

#Normalizing data
```

Also in the main function we use sample idea to reduce the search scale. And we get 90% data to this program. Particularly, we won't use sample in the small\_data...\_.txt.

### (3) The Backward

This part is mostly same as the (2). The different is we add early abandon in this part.

```
def BackwardElimination(data_nor,NF):
    param_stop= input("\nType the accuracy param to early abandon(0-100): \n You may input 100 if you don't want early abandon")
    param_stop= int(param_stop)

    print("Beginning search.\n")
    final_acc = 0
    best_feature=[]
    current_features = [i for i in range(1, NF+1)]

    for i in range(1, NF):
        print("\n On level %d of the search tree" % (i),"contains", current_features)
        feature_select = 0
        cur_acc = 0
        for j in range(1,NF):
            if (j in current_features):
                acc = CrossValidation(current_features,data_nor,j,2)
                if acc > cur_acc:
                    cur_acc = acc
                    feature_select = j
        if feature_select in current_features:
            current_features.remove(feature_select)
            print("\n On level ", i, " feature ", feature_select, " was removed from the current set")
            print("\n With ", len(current_features), " features, the accuracy is: ", cur_acc * 100, "%")
        if cur_acc >= final_acc:
            final_acc = cur_acc
            best_feature= list(current_features)
        if final_acc * 100 > param_stop:
            print()
            print("Early abandon!! The best feature subset is:", best_feature,"which has an accuracy of", final_acc * 100, "%")
            return

    print()
    print("Finish search!! The best feature subset is:", best_feature,"which has an accuracy of", final_acc * 100, "%")
```

#### 4. Conclusion

##### Backward Elimination:

Start with all available features and gradually remove the least important ones. After each deletion, the model is re-evaluated to check if the performance has improved or at least remained stable. This process continues until any further deletions significantly degrade the performance of the model.

Advantages: Initially consider all features, and gradually remove unimportant features, thus giving an initial model that considers the potential impact of all features.

Cons: Computationally expensive since it starts with a complex model.

Forward Selection: Start with no features and gradually add the most important features based on model performance. After each addition, the model is re-evaluated to ensure improved performance. This process continues until adding any new features does not significantly improve the performance of the model.

Pros: Initial computational cost is lower because it starts with a simpler model. It only adds features that provide significant improvements.

Cons: It may miss traits that are only effective in combination with others. Also, it can be more computationally expensive as the model complexity increases.

These methods are very intuitive, but computationally expensive, especially when dealing with a large number of features. They also do not take into account any possible interactions between features.

## Appendix:

### (1) Small data

Welcome to the Feature Selection Algorithm:

Type in the name of the file to test: CS170\_small\_Data\_\_31.txt

Type the algorithm you want to run:

- 1.Forward Selection
- 2.Backward Elimination

1

This dataset has 10 features (no including the class attribute), with 900 instances  
Beginning search.

On level 1 of the search tree contains []

Using feature(s) [1] accuracy is 85.76 %

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

On level 1 of the search tree, feature 1 was added to the current set

With 1 features, the accuracy is: 85.76 %

On level 2 of the search tree contains [1]

Using feature(s) [1, 2] accuracy is 83.2 %

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s) [1, 4] accuracy is 84.87 %

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s) [1, 6] accuracy is 98.22 %

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

On level 2 of the search tree, feature 6 was added to the current set

With 2 features, the accuracy is: 98.22 %

On level 3 of the search tree contains [1, 6]

Using feature(s) [1, 6, 2] accuracy is 94.77 %

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s) [1, 6, 4] accuracy is 94.99 %

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s) [1, 6, 8] accuracy is 95.11 %

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

On level 3 of the search tree, feature 8 was added to the current set

With 3 features, the accuracy is: 95.11 %

On level 4 of the search tree contains [1, 6, 8]

Using feature(s) [1, 6, 8, 2] accuracy is 90.1 %

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s) [1, 6, 8, 4] accuracy is 90.66 %

Using feature(s) [1, 6, 8, 5] accuracy is 91.21 %

Using feature(s) [1, 6, 8, 7] accuracy is 91.77 %

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

On level 4 of the search tree, feature 7 was added to the current set

With 4 features, the accuracy is: 91.77 %

Finish search!! The best feature subset is: [1, 6] which has an accuracy of 98.22 %

(2) Big data result

Welcome to the Feature Selection Algorithm:

Type in the name of the file to test: CS170\_large\_Data\_\_18.txt

Type the algorithm you want to run:

1.Forward Selection

2.Backward Elimination

1

This dataset has 20 features (no including the class attribute), with 1800 instances  
Beginning search.



On level 1 of the search tree contains []

Using feature(s) [1] accuracy is 85.38 %

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

On level 1 of the search tree, feature 1 was added to the current set

With 1 features, the accuracy is: 85.34 %

On level 2 of the search tree contains [1]

Using feature(s) [1, 2] accuracy is 83.43 %

Using feature(s) [1, 3] accuracy is 84.27 %

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s) [1, 7] accuracy is 96.66481378543635 %

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

On level 2 of the search tree, feature 7 was added to the current set

On level 3 of the search tree contains [1, 7]

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

On level 3 of the search tree, feature 9 was added to the current set

On level 4 of the search tree contains [1, 7, 9]

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning



[illegible]

[illegible]

On level 1 of the search tree, feature 1 was added to the current set

With 1 features, the accuracy is: 83.58 %

On level 2 of the search tree contains [1]

Using feature(s) [1, 2] accuracy is 84.99 %

[illegible]





[illegible]







Using feature(s), the feature has been stop by minimax like pruning  
Using feature(s), the feature has been stop by minimax like pruning  
Using feature(s), the feature has been stop by minimax like pruning  
Using feature(s), the feature has been stop by minimax like pruning  
Using feature(s), the feature has been stop by minimax like pruning  
Using feature(s), the feature has been stop by minimax like pruning  
Using feature(s), the feature has been stop by minimax like pruning  
Using feature(s), the feature has been stop by minimax like pruning  
Using feature(s), the feature has been stop by minimax like pruning  
Using feature(s), the feature has been stop by minimax like pruning  
Using feature(s), the feature has been stop by minimax like pruning

On level 4 of the search tree, feature 66 was added to the current set

With 4 features, the accuracy is: 94.24 %

Finish search!! The best feature subset is: [1, 74] which has an accuracy of 97.80 %

#### (4) Diabetes data from kaggle data result

Welcome to the Feature Selection Algorithm:

Type in the name of the file to test: new\_data.txt

Type the algorithm you want to run:

- 1.Forward Selection
- 2.Backward Elimination

1

This dataset has 8 features (no including the class attribute), with 691 instances  
Beginning search.

On level 1 of the search tree contains []  
Using feature(s) [1] accuracy is 49.13 %  
Using feature(s) [2] accuracy is 63.33 %  
Using feature(s), the feature has been stop by minimax like pruning  
Using feature(s), the feature has been stop by minimax like pruning  
Using feature(s), the feature has been stop by minimax like pruning  
Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

On level 1 of the search tree, feature 2 was added to the current set

With 1 features, the accuracy is: 63.33 %

On level 2 of the search tree contains [2]

Using feature(s) [2, 1] accuracy is 64.35 %

Using feature(s) [2, 3] accuracy is 64.78 %

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s) [2, 6] accuracy is 69.42 %

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

On level 2 of the search tree, feature 6 was added to the current set

With 2 features, the accuracy is: 69.42 %

On level 3 of the search tree contains [2, 6]

Using feature(s) [2, 6, 1] accuracy is 70.29 %

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

On level 3 of the search tree, feature 1 was added to the current set

With 3 features, the accuracy is: 70.29 %

On level 4 of the search tree contains [2, 6, 1]

Using feature(s) [2, 6, 1, 3] accuracy is 66.38 %

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s) [2, 6, 1, 5] accuracy is 70.0 %

Using feature(s) [2, 6, 1, 7] accuracy is 70.14 %

Using feature(s), the feature has been stop by minimax like pruning

On level 4 of the search tree, feature 7 was added to the current set

With 4 features, the accuracy is: 70.14 %

On level 5 of the search tree contains [2, 6, 1, 7]

Using feature(s) [2, 6, 1, 7, 3] accuracy is 67.83 %

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s), the feature has been stop by minimax like pruning

Using feature(s) [2, 6, 1, 7, 8] accuracy is 69.28 %

On level 5 of the search tree, feature 8 was added to the current set

With 5 features, the accuracy is: 69.28 %

Finish search!! The best feature subset is: [2, 6, 1] which has an accuracy of 70.29 %

## Citation

[1] <https://www.kaggle.com/datasets/akshaydattatraykhare/diabetes-dataset>