**Introduction**

The program assignment 2 is a binary classification problem, it asks to correctly classify samples into different letters. To accomplish this goal, 5 different ML models should be applied, their hyperparameter should be tuned using 5-fold cross-validation, and some performance metrics (e.g. accuracy) should be used to help find the best model for this classification problem, then, dimension reduction should be applied, and the performance of different models before dimension reduction and after dimension reduction should be compared. Ultimately, some conclusions can be drawn from this.

I chose A and B for the third problem, before working on this problem, I tend to think M and Y will be the easiest to classify since the number of samples of these 2 letters is the largest among all 3 pairs of the 3 classification problems, which may could bring more information.

Dimension reduction should be regarded as useful for this problem, because before reducing the dimension, this dataset contains a sum of 16 features and is relatively too many, also, some of the features may be similar to each other hence should be eliminated.

A good dimension reduction method should be method that is able to eliminate the expected number of features while avoiding a significant accuracy loss/
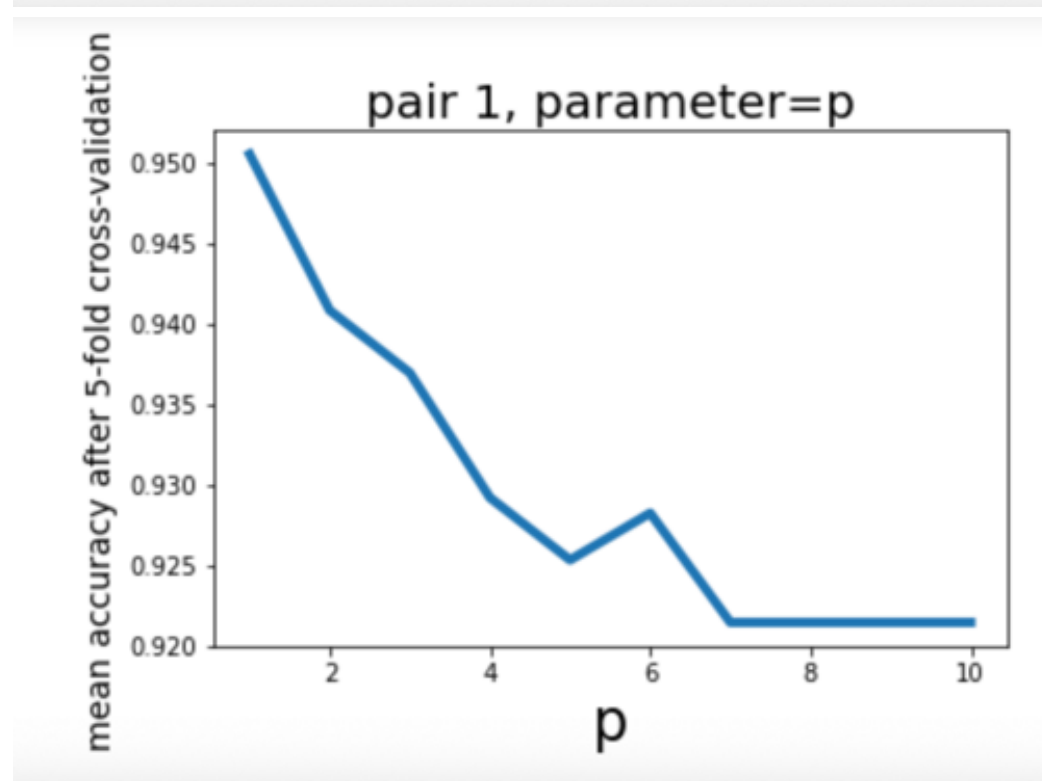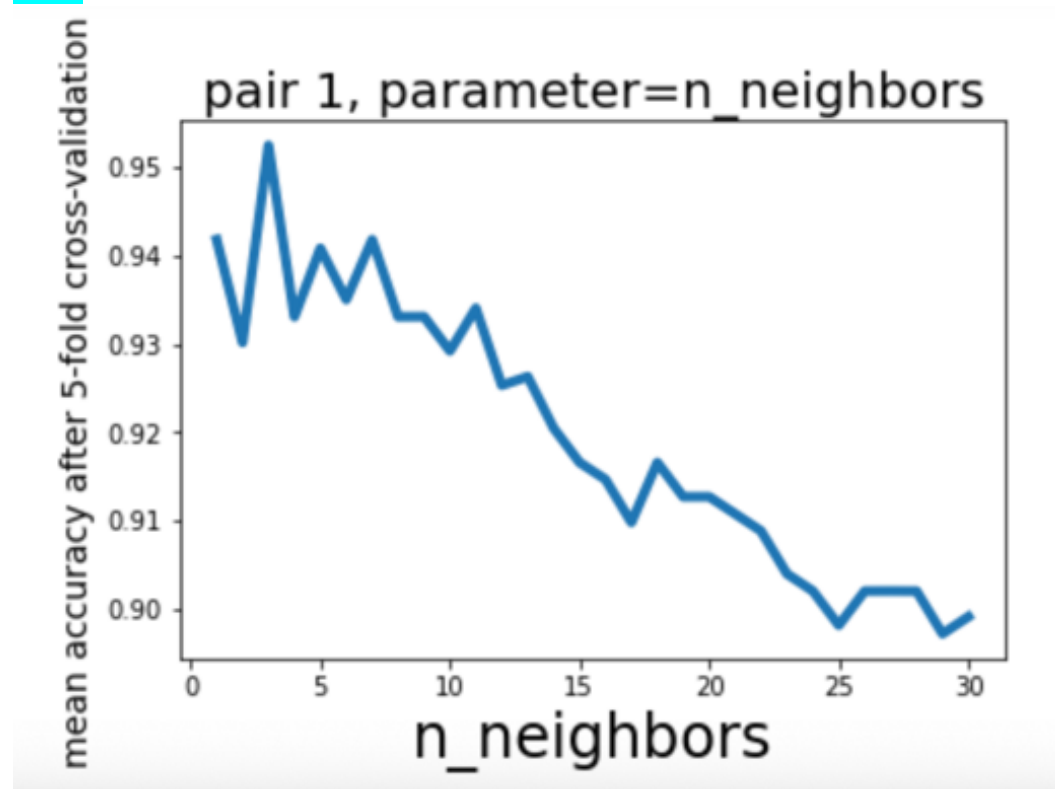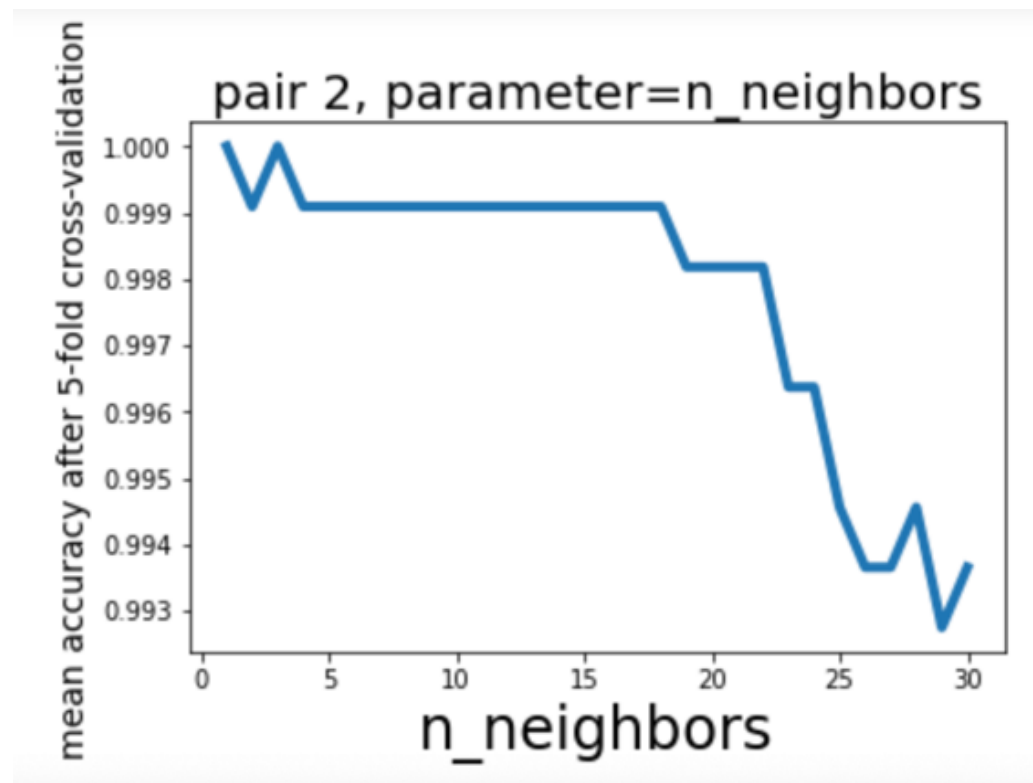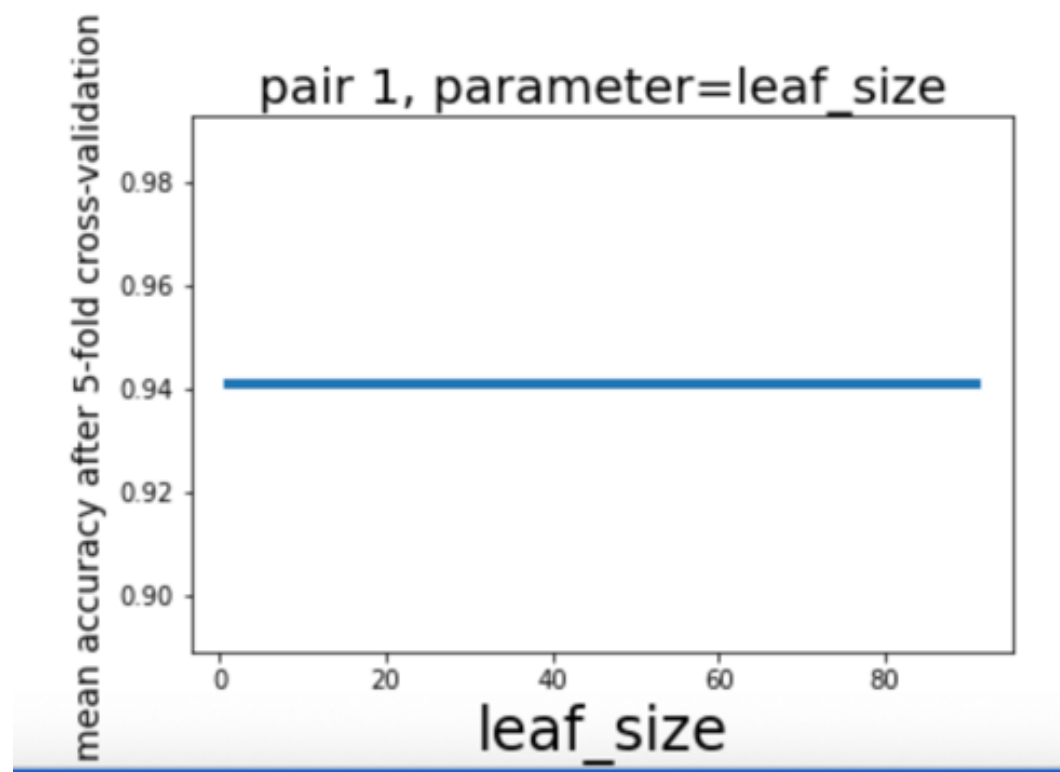
**Results**

The pros and cons of each classifier are list in the following table.
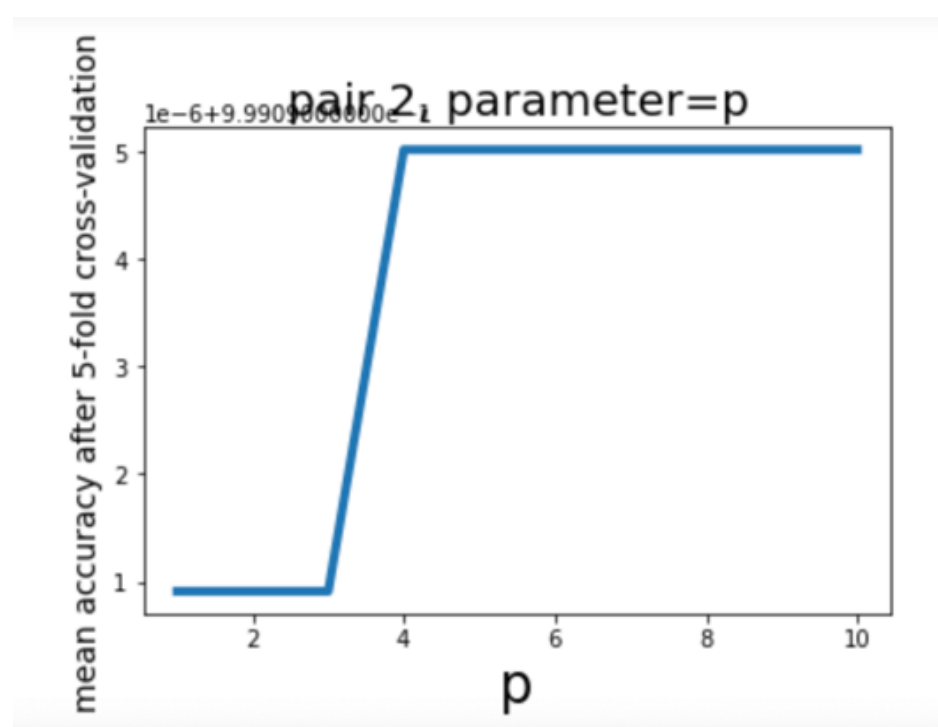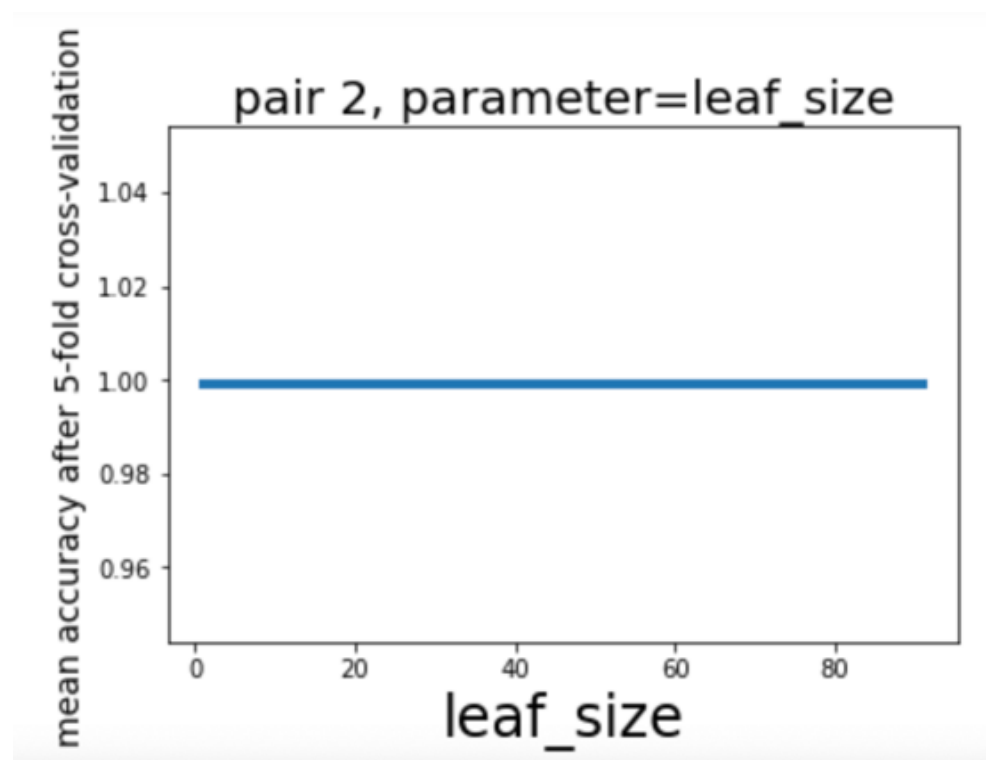
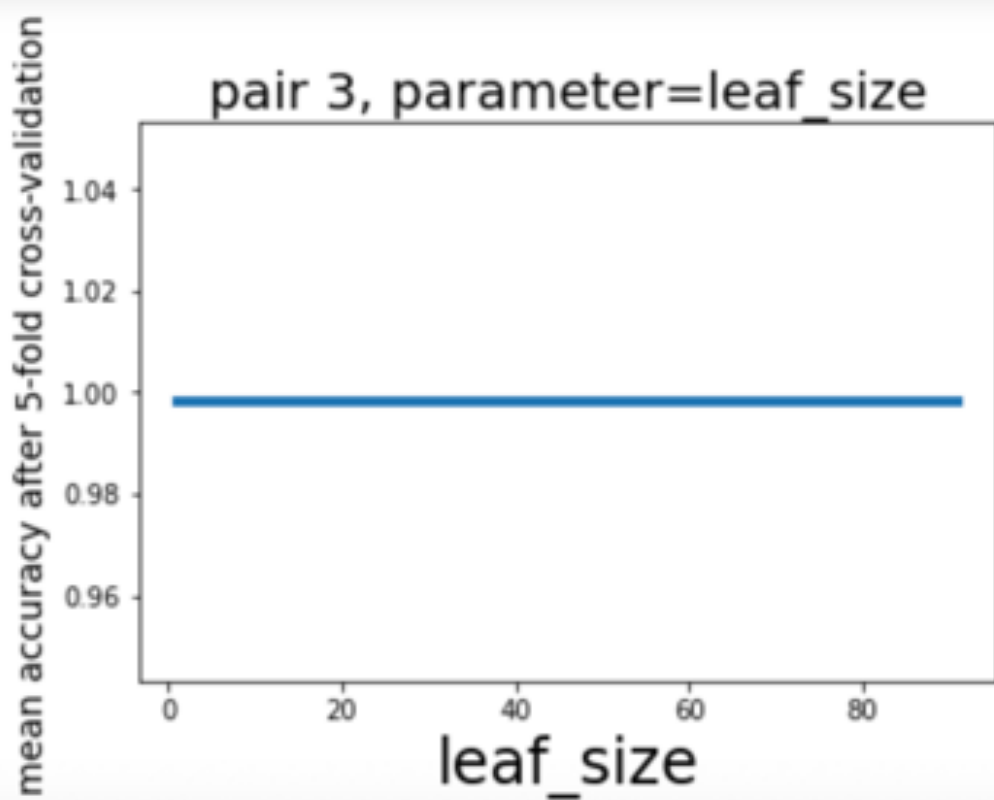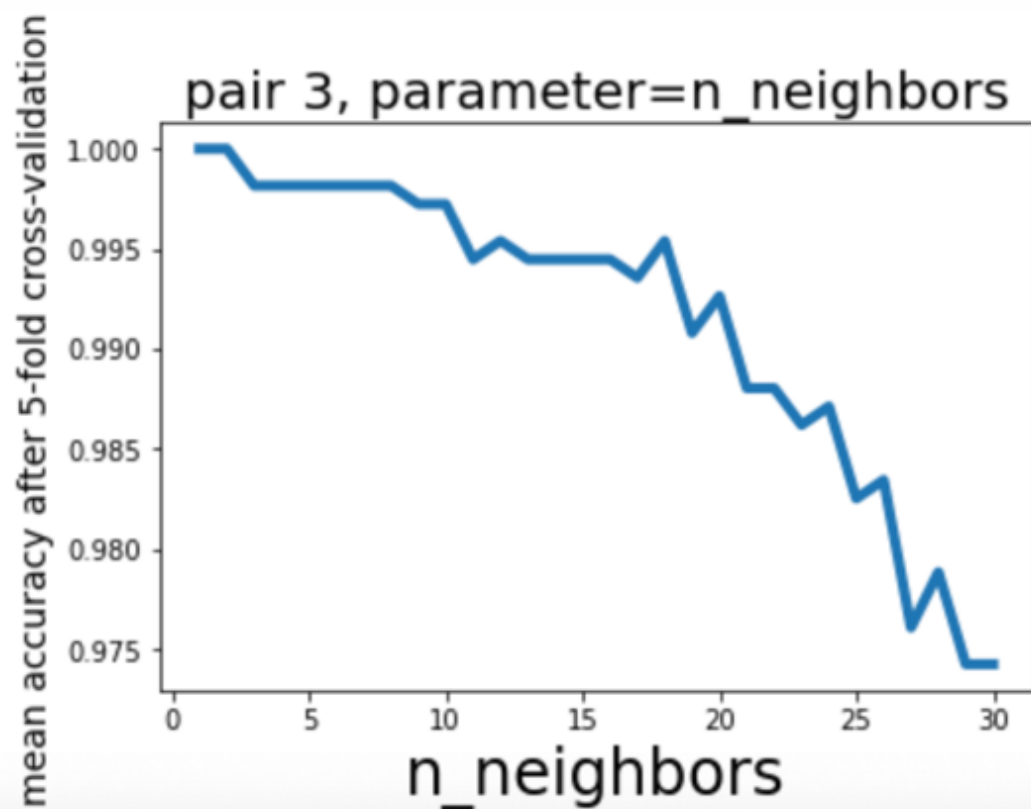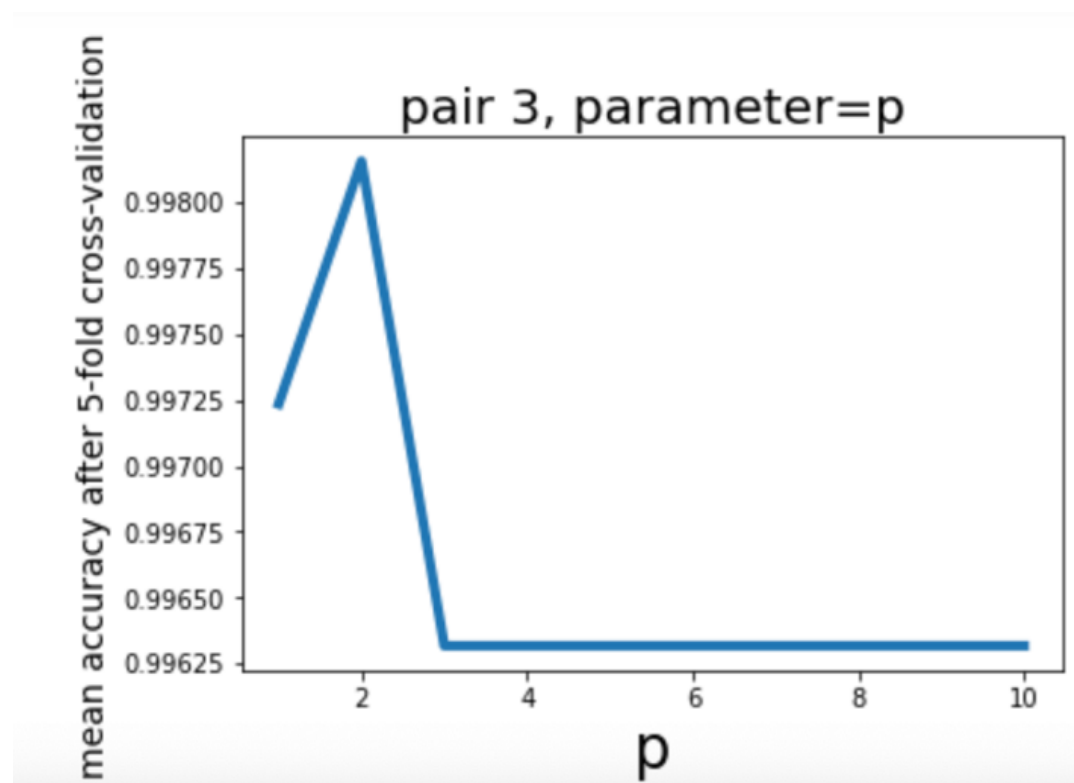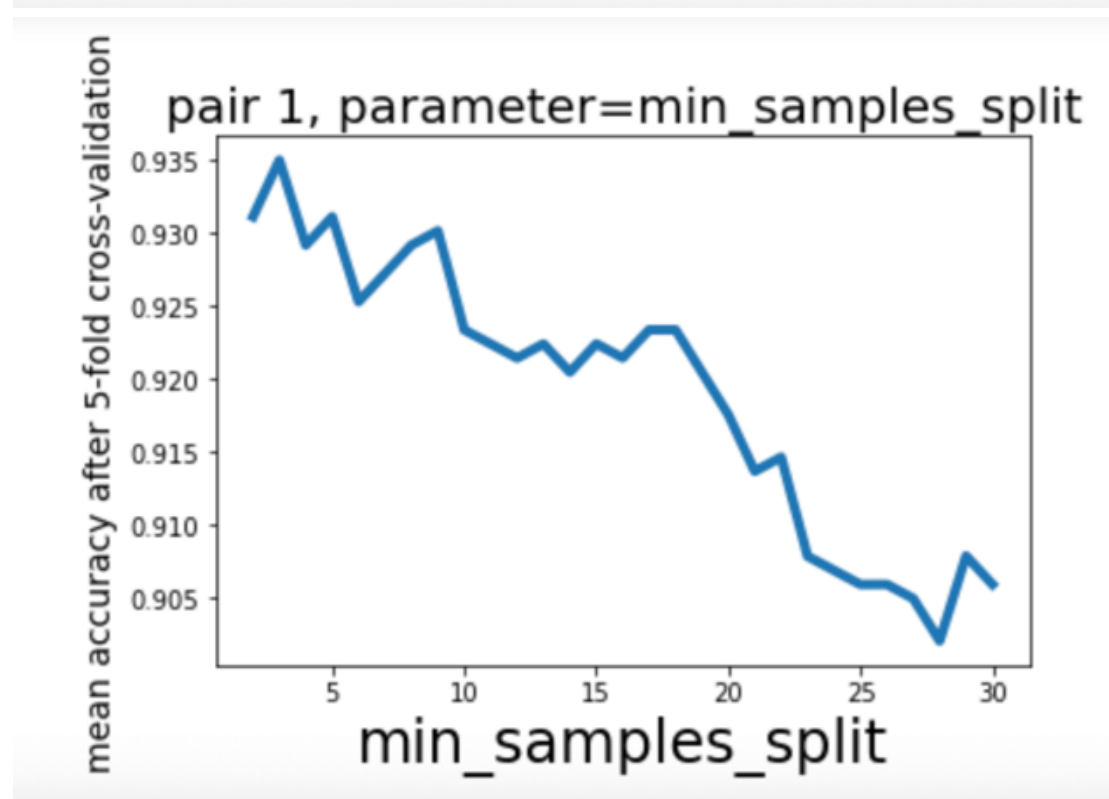|  | pros | cons |
|---|---|---|
| k-nearest neighbors | Requires no training time | Requires the storage of all the training data, and the testing time of calculating the distance between test samples and every training sample is relatively inefficient |
| Decision tree | Easy to apply | Easy to overfit |
| Random forest | Can have high accuracy while avoiding overfitting like a single decision tree | Takes a long-time training |
| SVM | Good robustness，able to avoid dimensional disasters | Difficult to implement for large training samples, Difficult in solving the problem of multiple classifications |
| Artificial neural network | Can fit just about any dataset with high accuracy | Difficult to interpret to people with no data science or machine learning background; long training time |

Cross validation results

## pair 1, parameter=leaf_size

mean accuracy after 5-fold cross-validation

leaf_size

## pair 2, parameter=n_neighbors

mean accuracy after 5-fold cross-validation

n_neighbors

pair 2, parameter=leaf_size

pair 2, parameter=p

pair 3, parameter=n_neighbors

pair 3, parameter=leaf_size

pair 3, parameter=p

pair 1, parameter=min_samples_leaf

pair 1, parameter=max_leaf_nodes

pair 2, parameter=max_depth

pair 2, parameter=min_samples_split

pair 2, parameter=min_samples_leaf

pair 2, parameter=max_leaf_nodes

pair 3, parameter=max_depth

pair 3, parameter=min_samples_split

pair 3, parameter=min_samples_leaf

pair 3, parameter=max_leaf_nodes

pair 1, parameter=max_depth



pair 2, parameter=max_depth

pair 3, parameter=max_depth

pair 1, parameter=coef0



pair 2, parameter=coef0

pair 3, parameter=coef0

pair 1, parameter=activation



pair 2, parameter=activation

pair 3, parameter=activation

Dimension reduction methods used

The three methods I used are filter method--f_score, Wrapper method--forward feature selection, and Embedded Method—random forest.
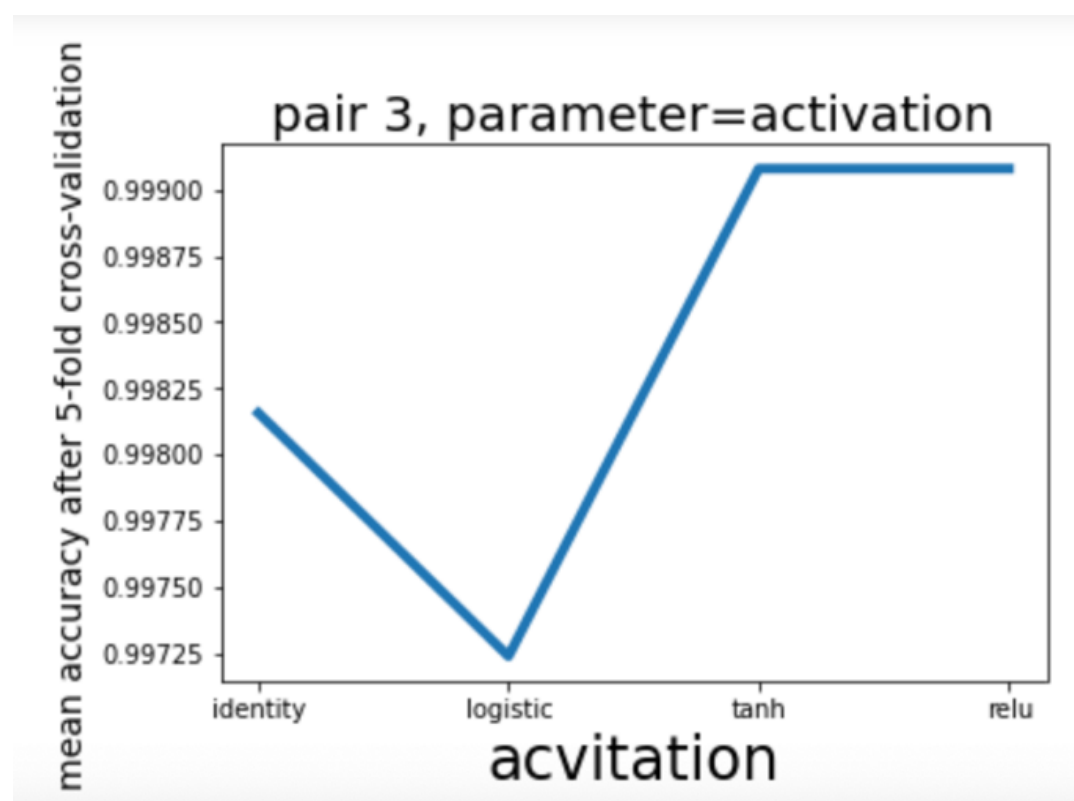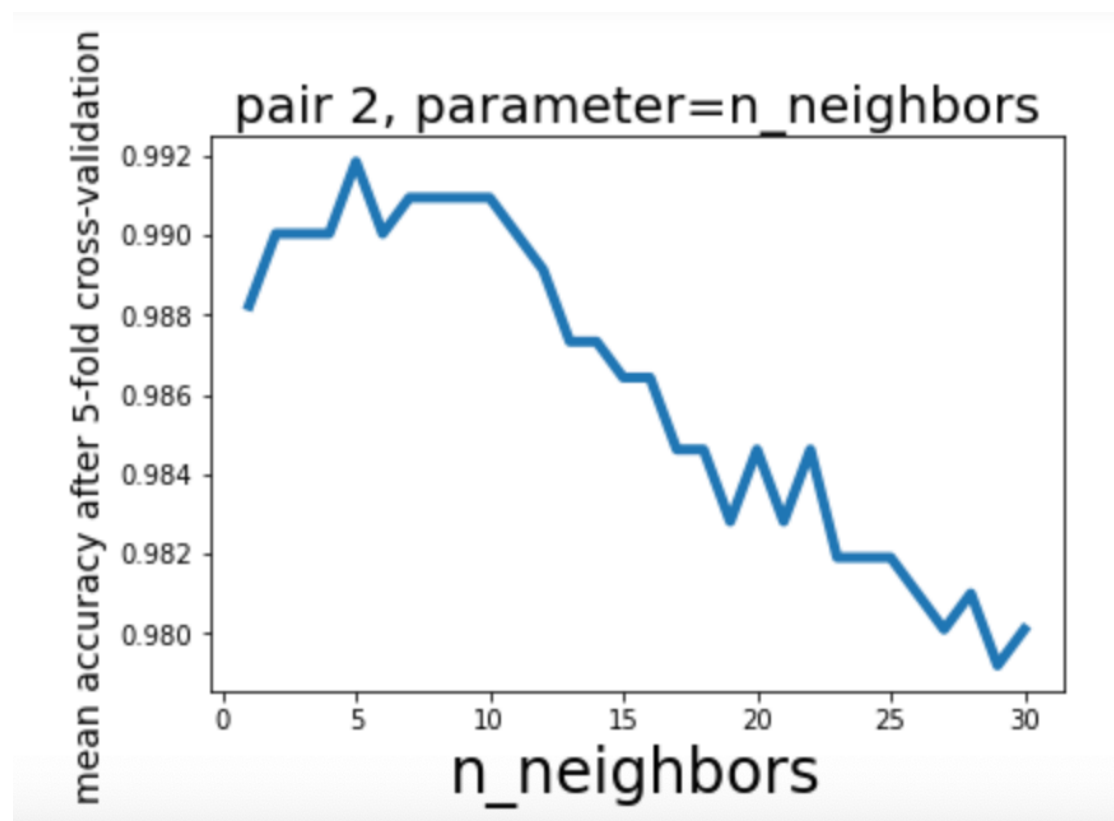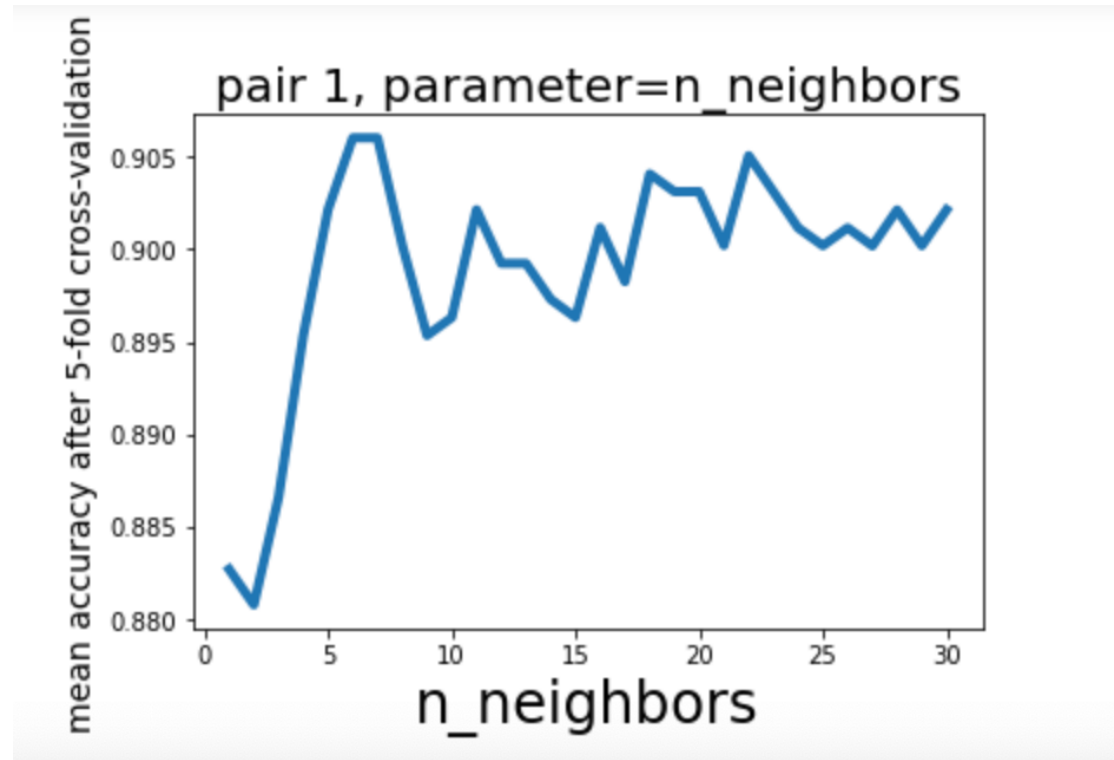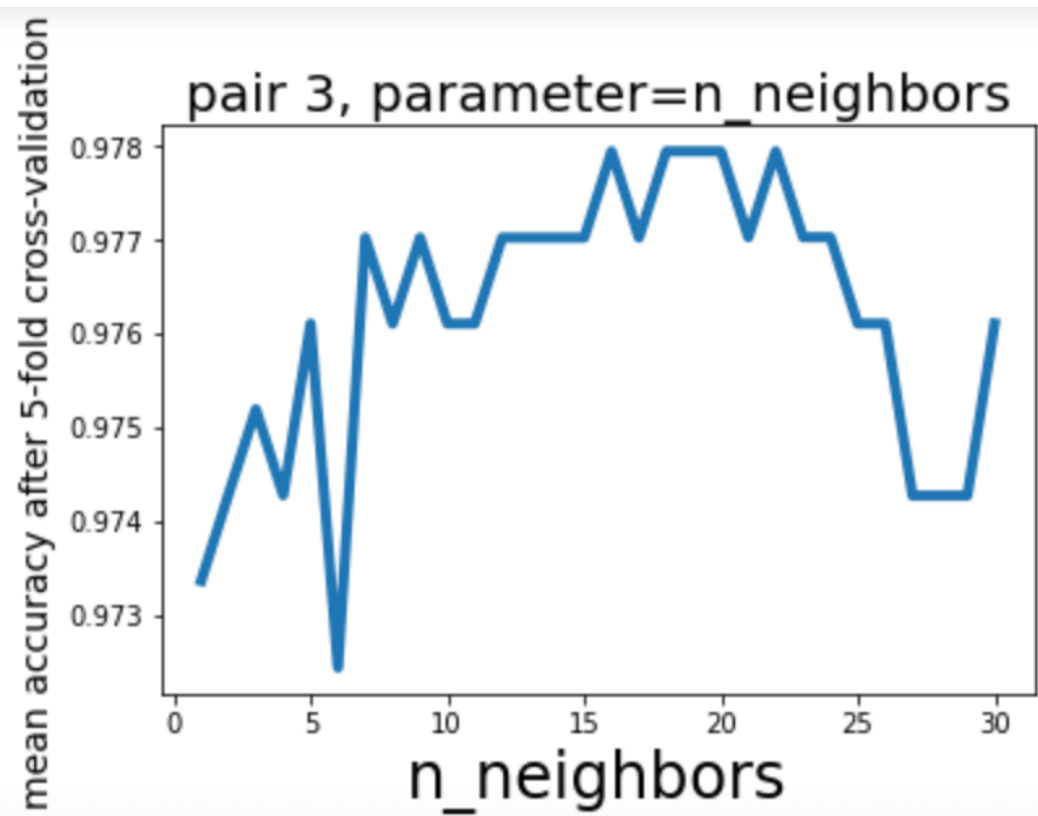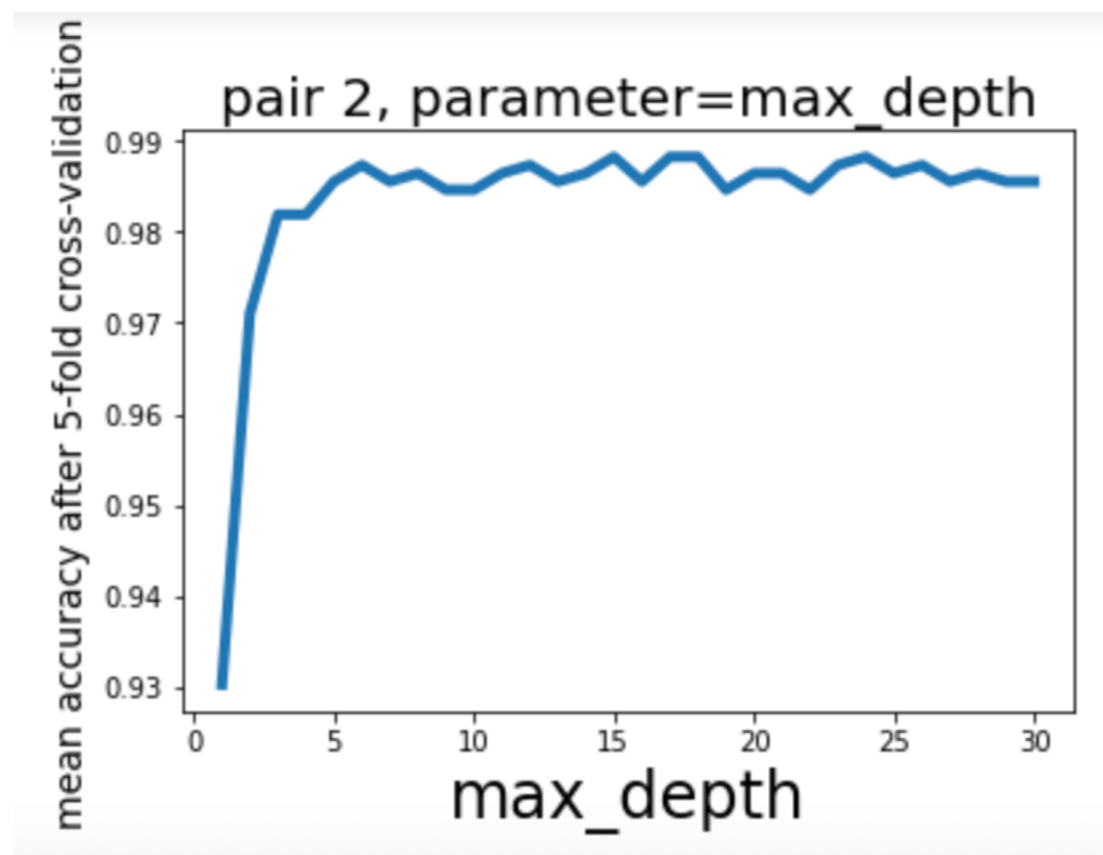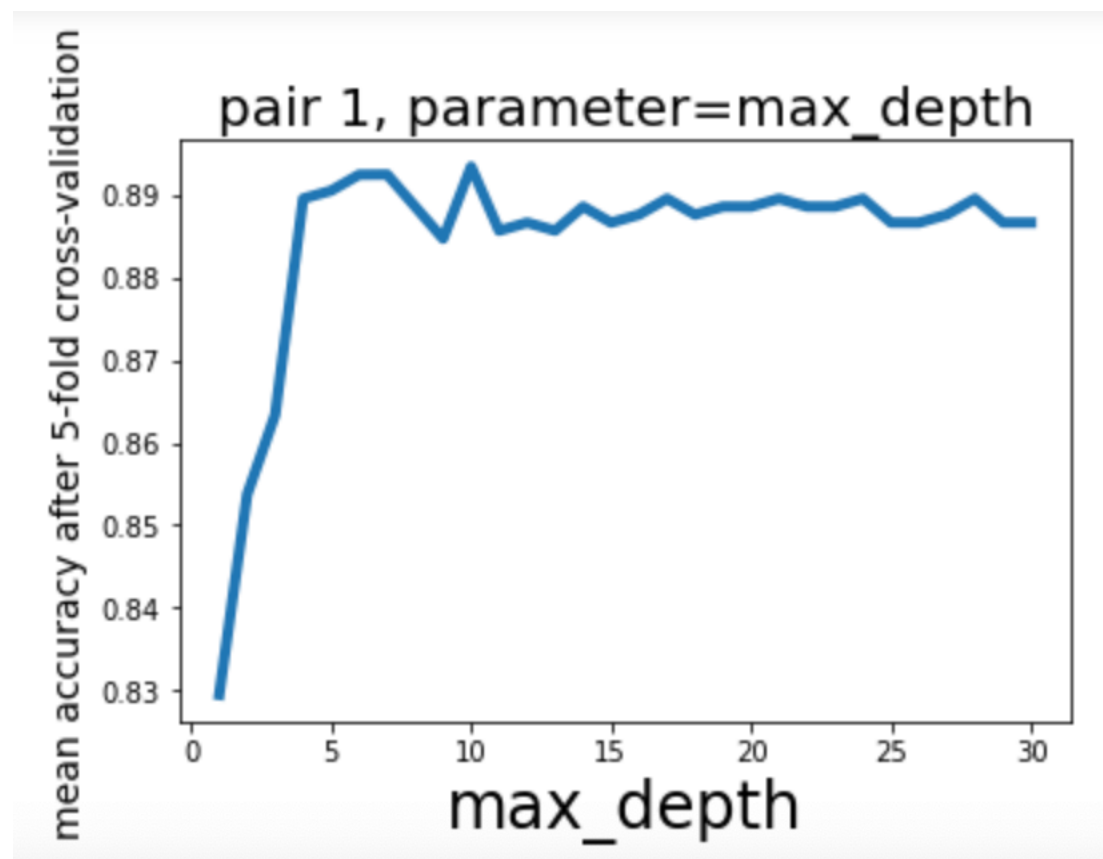
pair 3, parameter=n_neighbors
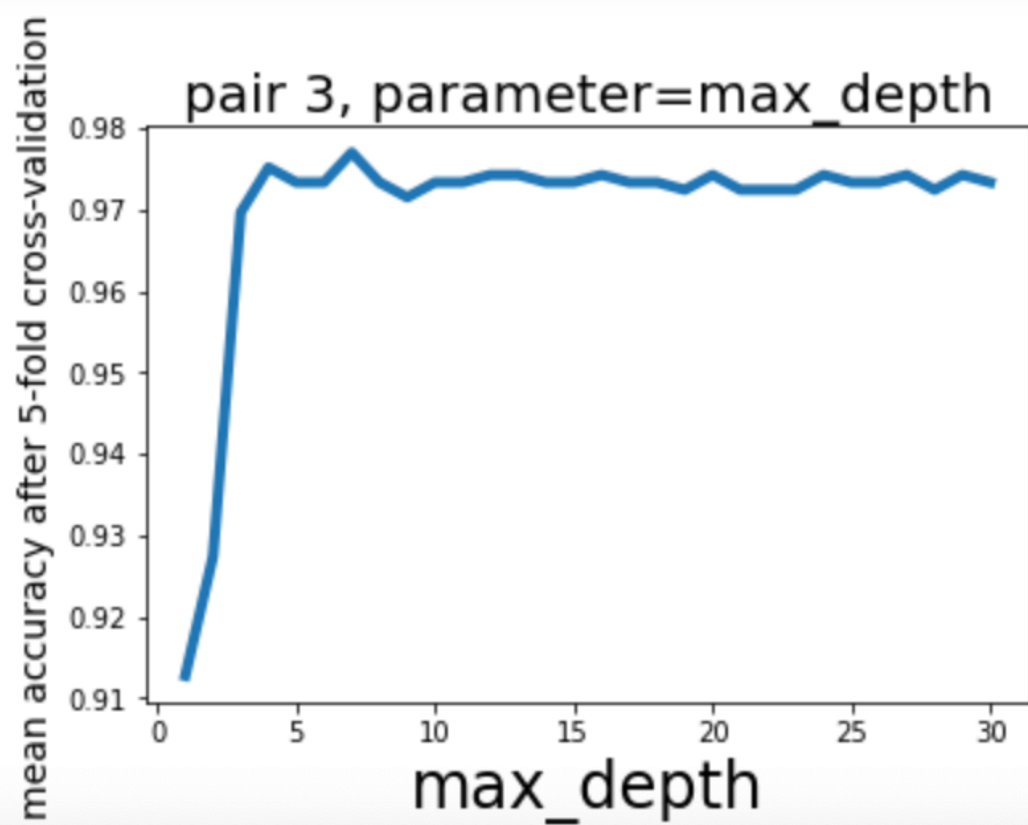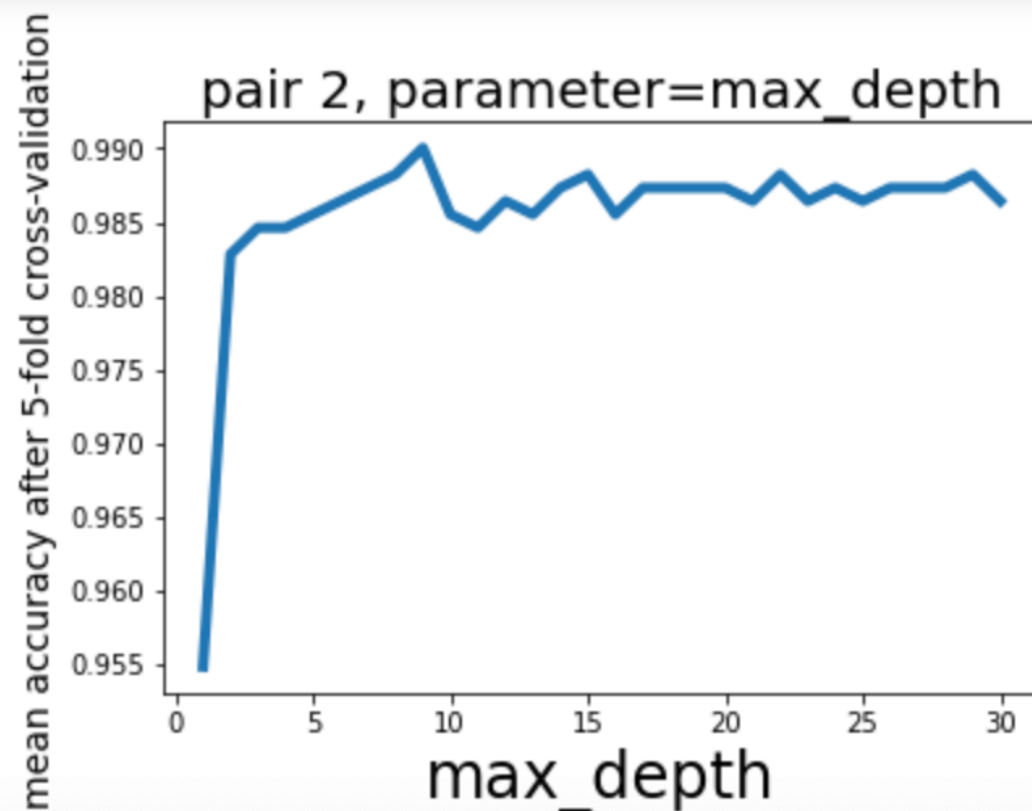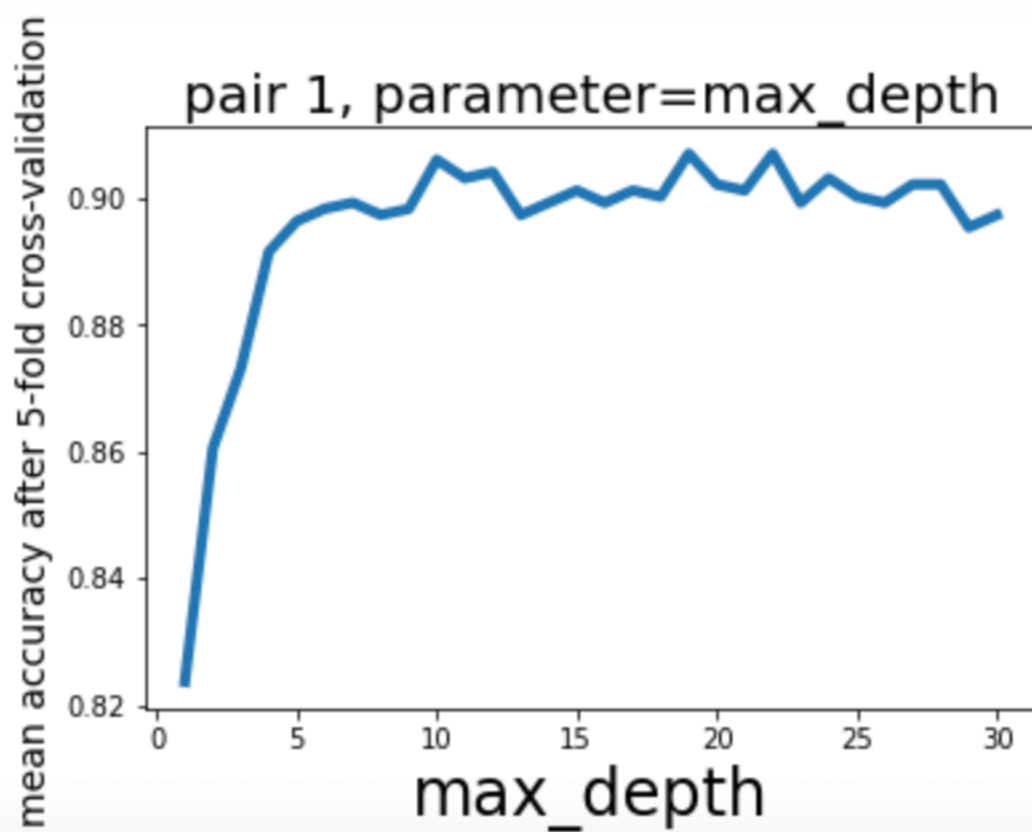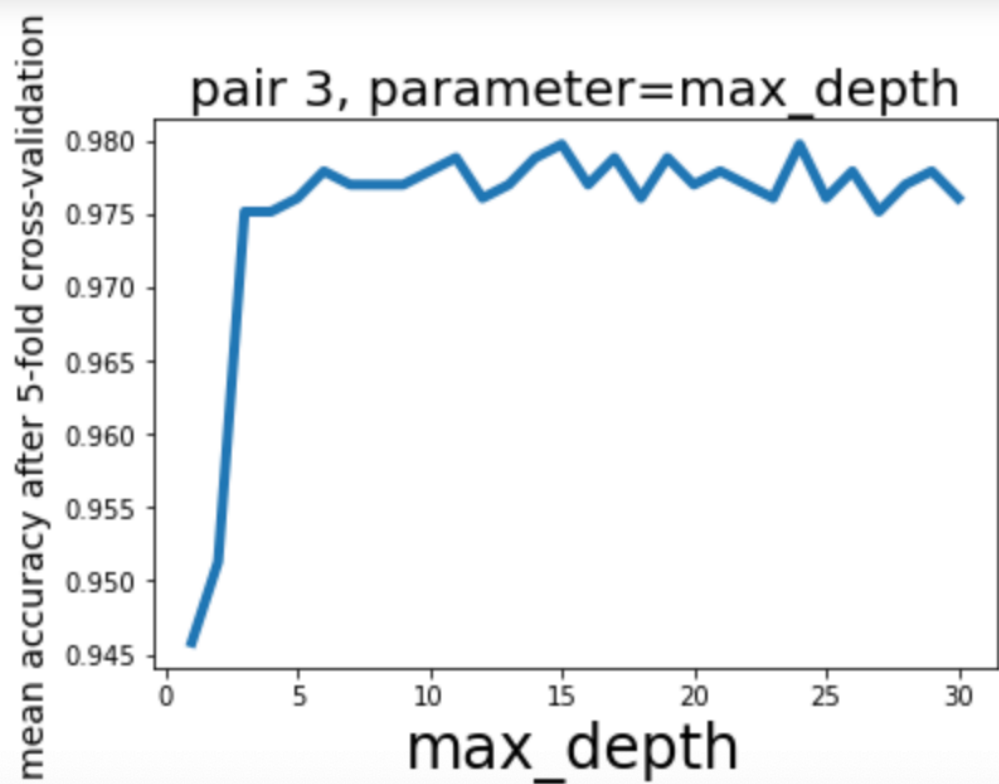
pair 3, parameter=max_depth

pair 1, parameter=max_depth



pair 2, parameter=max_depth

pair 3, parameter=max_depth

pair 3, parameter=coef0

pair 1, parameter=activation



pair 2, parameter=activation

pair 3, parameter=activation

Pair 1

|  | Accuracy on validation set before dimension reduction |
|---|---|
| KNN | 0.8245977011494252 |
| Decision tree | 0.8372413793103449 |
| Random forest | 0.8510344827586207 |
| SVM | 0.8917241379310346 |
| ANN | 0.8917241379310343 |

Pair 2

|  | Accuracy on validation set before dimension reduction |
|---|---|
| KNN | 0.9875 |
| Decision tree | 0.98125 |
| Random forest | 0.98125 |
| SVM | 0.98125 |
| ANN | 0.9875 |

Pair 3

|  | Accuracy on validation set before dimension reduction |
|---|---|
| KNN | 0.955241935483871 |
| Decision tree | 0.9358870967741936 |
| Random forest | 1.0 |
| SVM | 0.955241935483871 |
| ANN | 0.9679435483870968 |

Pair 1

|  | Accuracy on validation set after dimension reduction |
|---|---|
| KNN | 0.8581609195402299 |
| Decision tree | 0.8301149425287356 |
| Random forest | 0.871264367816092 |
| SVM | 0.864367816091954 |
| ANN | 0.83816091954023 |

Pair 2

|  | Accuracy on validation set after dimension reduction |
|---|---|
| KNN | 0.99375 |
| Decision tree | 0.975 |
| Random forest | 0.98125 |
| SVM | 0.98125 |
| ANN | 0.9872983870967742 |

Pair 3

|  | Accuracy on validation set after dimension reduction |
|---|---|
| KNN | 0.9423387096774194 |
| Decision tree | 0.955241935483871 |
| Random forest | 0.9487903225806452 |
| SVM | 0.9233870967741936 |
| ANN | 0.9233870967741936 |

According to accuracy, I will choose random forest as my final model for this problem, because its accuracy is relatively robust among 3 pairs and is also higher than the average level of 5 ML models.

On the one hand, dimension reduction lowers the accuracy to a certain extent, but it is acceptable since those are not a significant loss of accuracy; on the other hand, dimension reduction decrease the running time of all the 5 models which increases efficiency.

If I was given this same task for a new dataset, I will first normalize the data, because in most cases, after normalization, the performance of the models will be better.