# Breast Cancer Classification

Group 10 Project
Jancey Liu
Weibing Wang
Gege Zhang
Junyi Liu
Zephyr Xiao

# Catalog:

- Question
- Methods & Model
- Visualization of Features
- Proportion
- ROC
- Coefficients of Features

# Question:

## Classifying whether a suspicious lump is benign or malignant based on 5 variables?

Mean Radius(mm): Mean of distances from center to points on the perimeter
Mean Texture: Standard deviation of gray-scale values
Mean Perimeter(mm): Mean size of the core tumor
Mean Area(mm^2): Mean slice area of the tumor
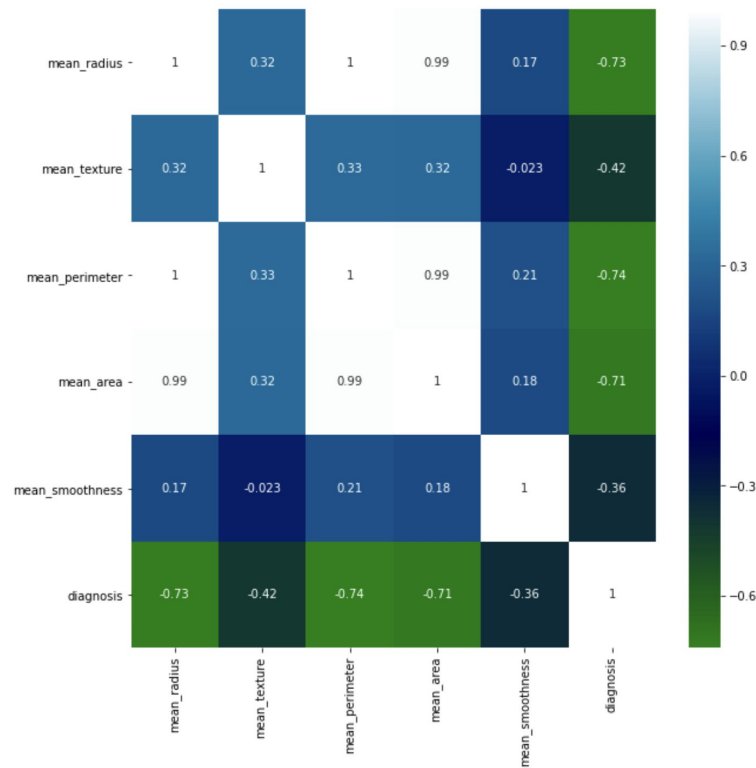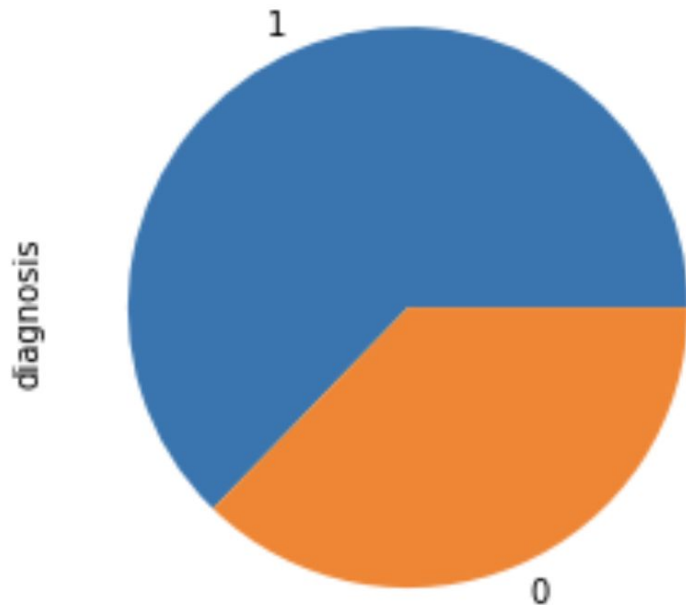Mean Smoothness: Mean of local variation in radius lengths

Categorical data:
Diagnosis : The diagnosis of breast tissues (0 = malignant, 1 = benign)

https://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+(original)

# Question:

## Classifying whether a suspicious lump is benign or malignant based on 6 variables?
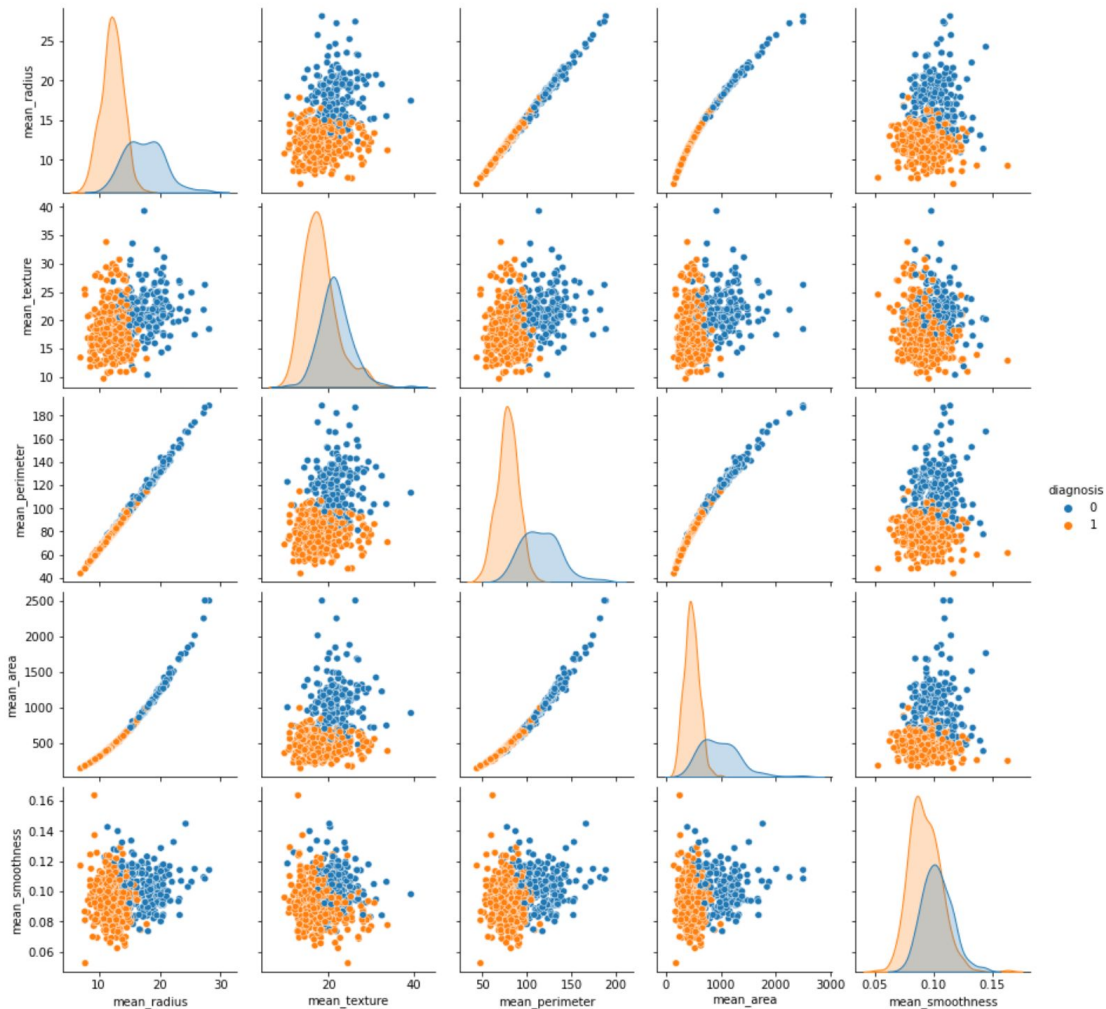
# Methods & Model

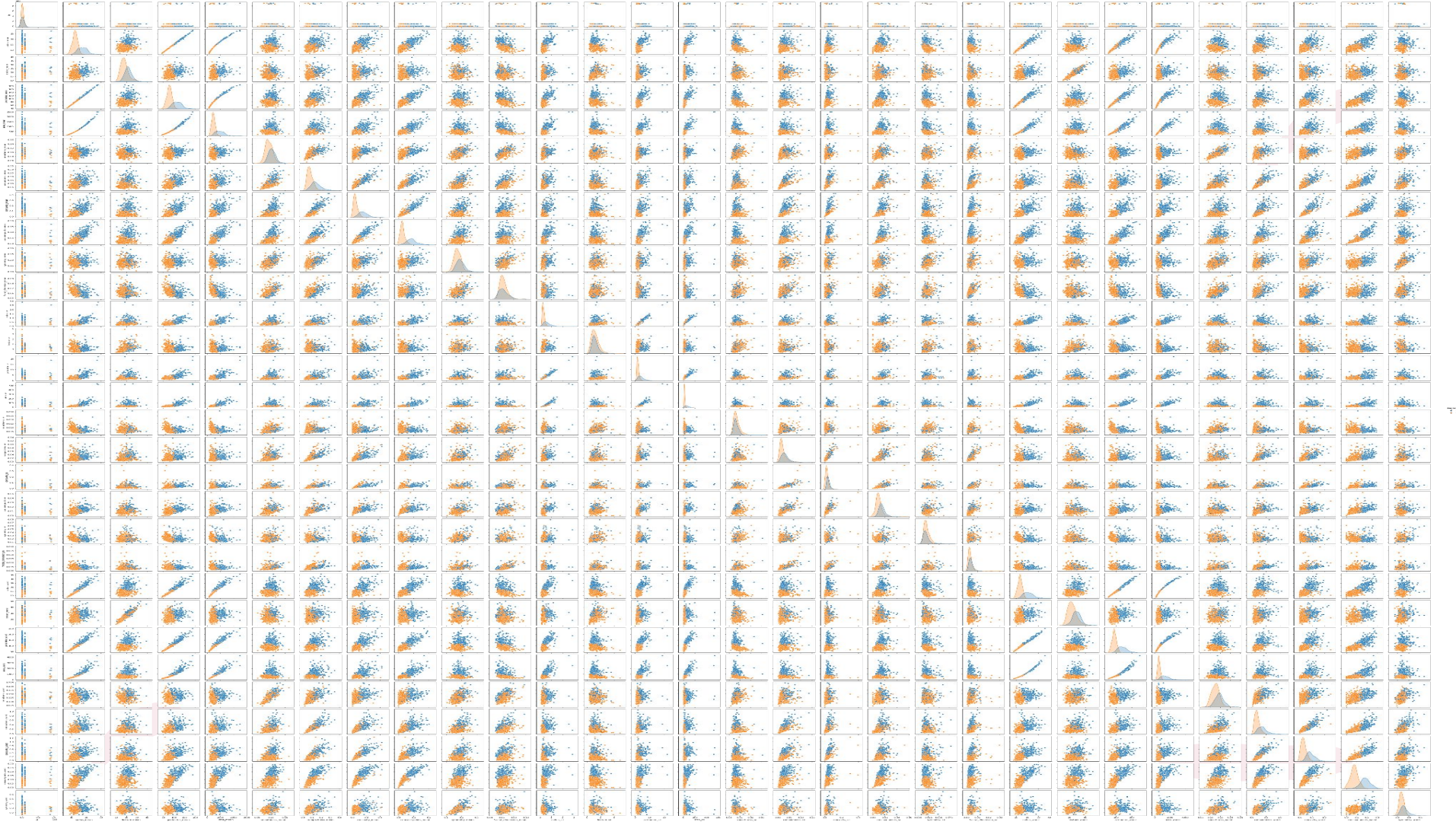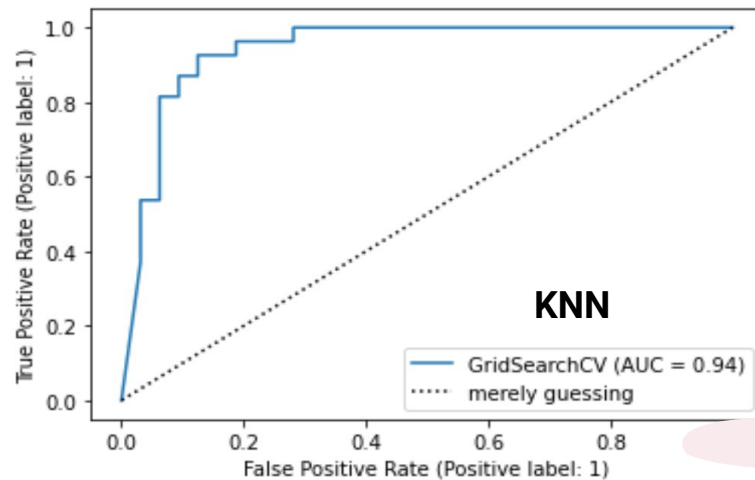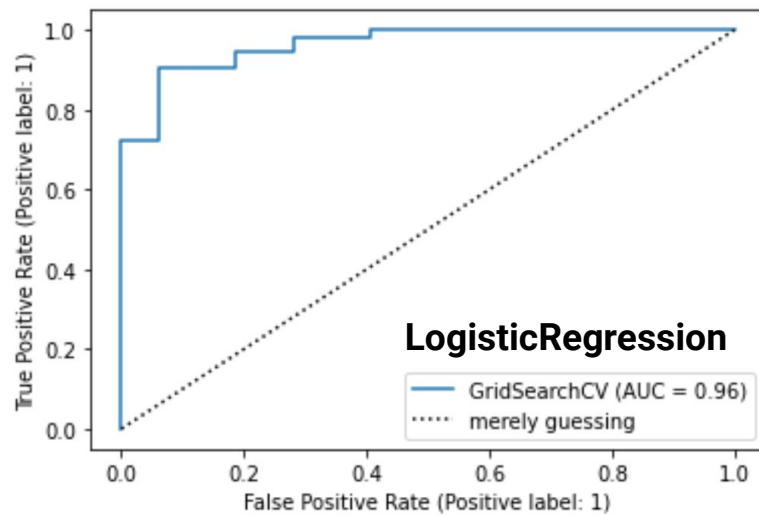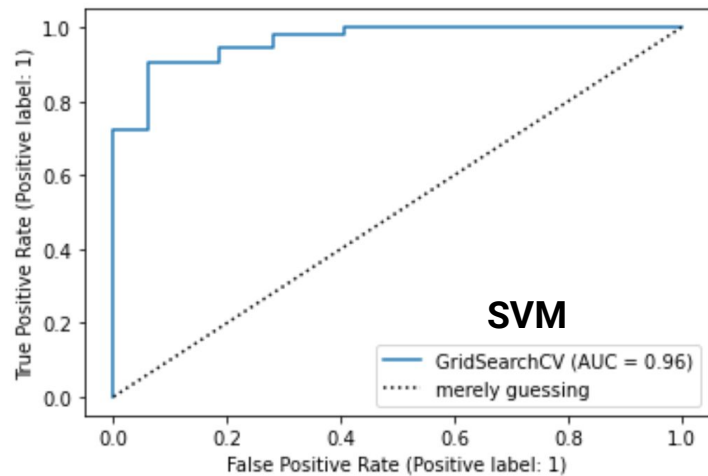| LogisticRegression | SVM | KNN | RandomForest |
|:---:|:---:|:---:|:---:|
| 91.8% | 91.8% | 88.2% | 94.1% |

# Visualization of Features
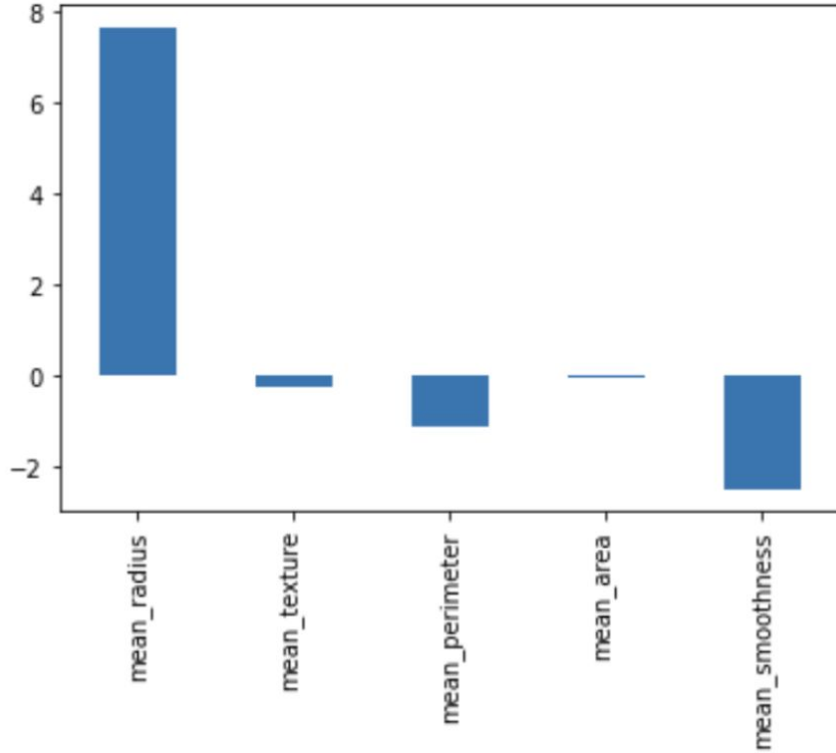
# Proportion

## Confusion Matrix

```
df:
     0   1
0   22  10
1    2  52
TN=22, FP=10, FN=2, TP=52
Classification Report:
              precision    recall  f1-score   support

           0       0.92      0.69      0.79        32
           1       0.84      0.96      0.90        54

    accuracy                           0.86        86
   macro avg       0.88      0.83      0.84        86
weighted avg       0.87      0.86      0.86        86
```
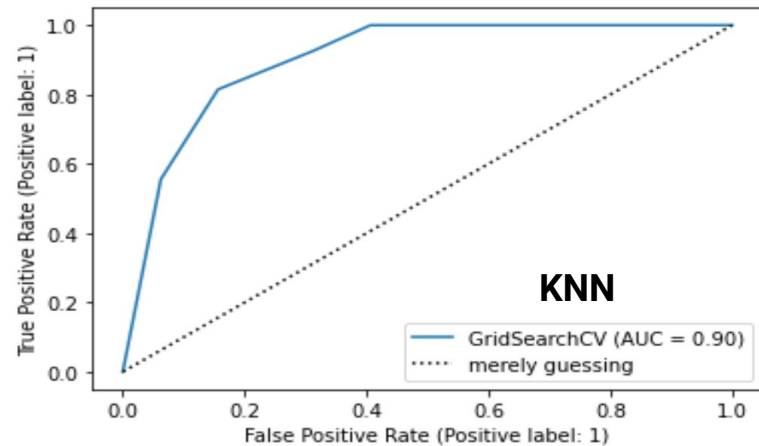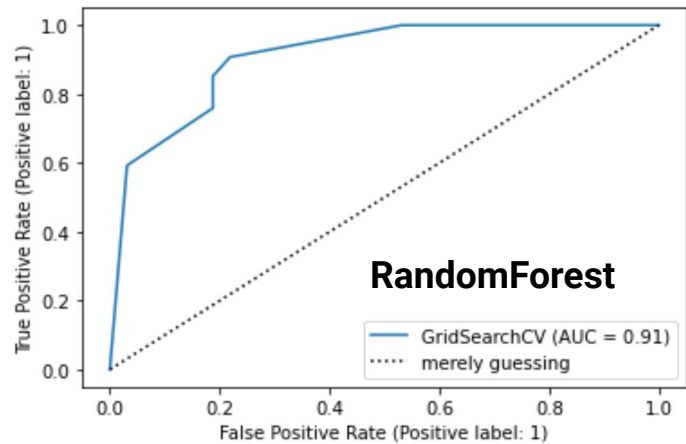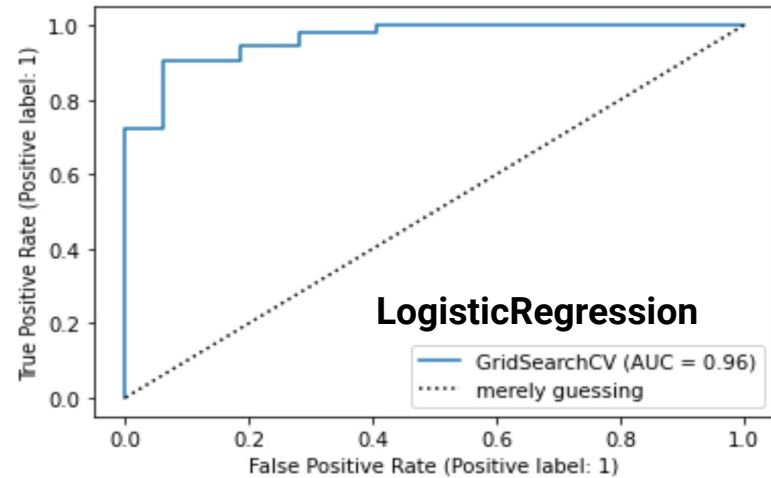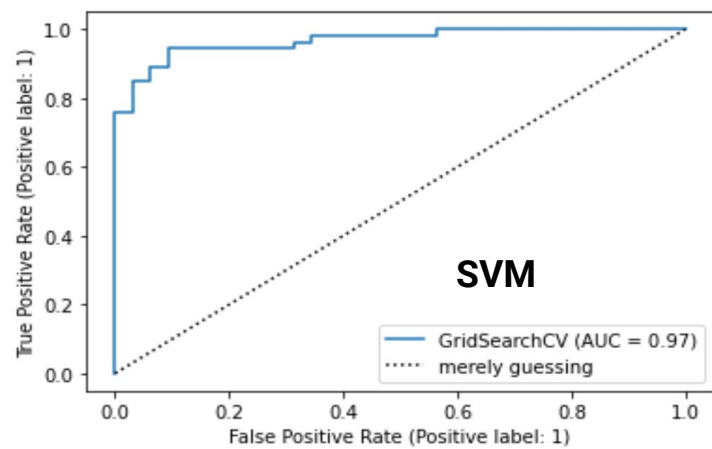
# Coefficients of Features (Importance of features) Before

# Proportion

## Confusion Matrix

```
df:
     0   1
0   24   8
1    3  51
TN=24, FP=8, FN=3, TP=51
Accuracy :0.872
Classification Report:
              precision    recall  f1-score   support

           0       0.89      0.75      0.81        32
           1       0.86      0.94      0.90        54

    accuracy                           0.87        86
   macro avg       0.88      0.85      0.86        86
weighted avg       0.87      0.87      0.87        86
```
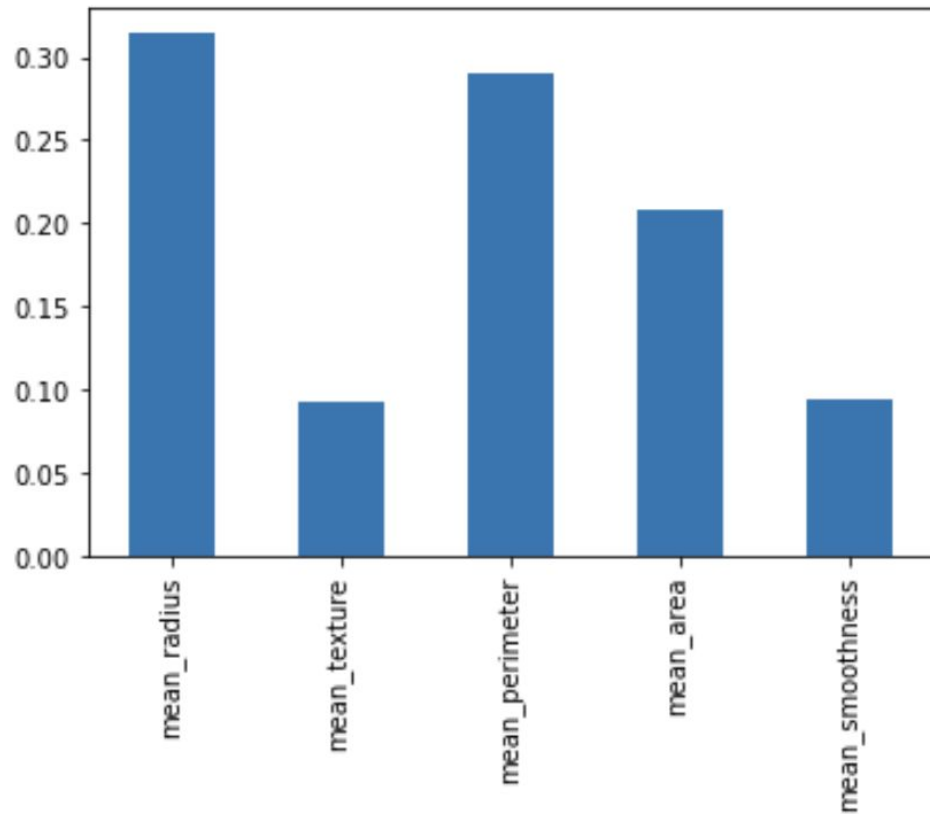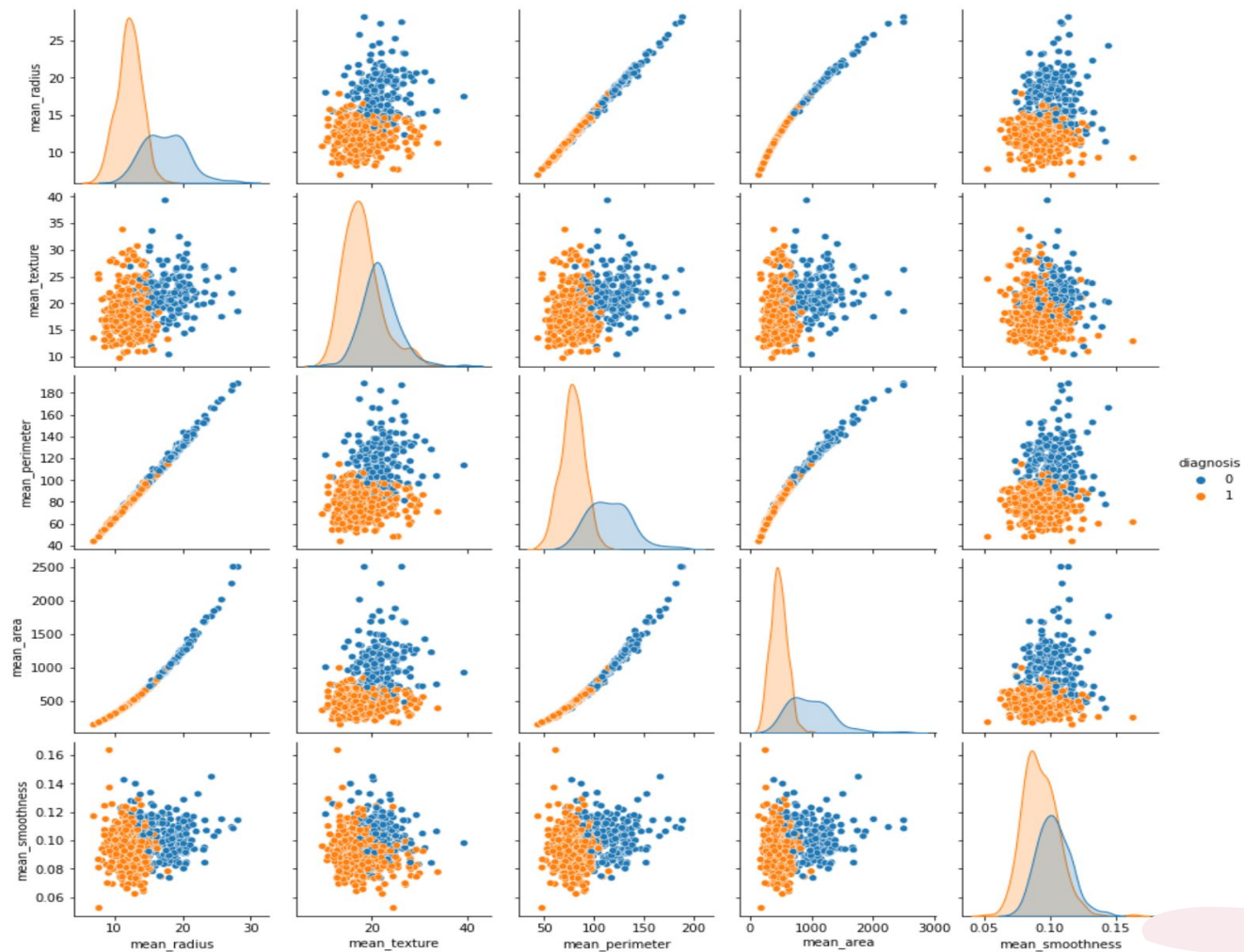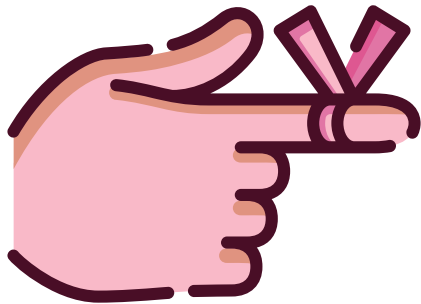
# Coefficients of Features (Importance of features) After

# Methods & Model

| LogisticRegression | SVM | KNN | RandomForest |
|:---:|:---:|:---:|:---:|
| 98.4% | 98.9% | 95.9% | 99.1% |

# Thank you for your listening！