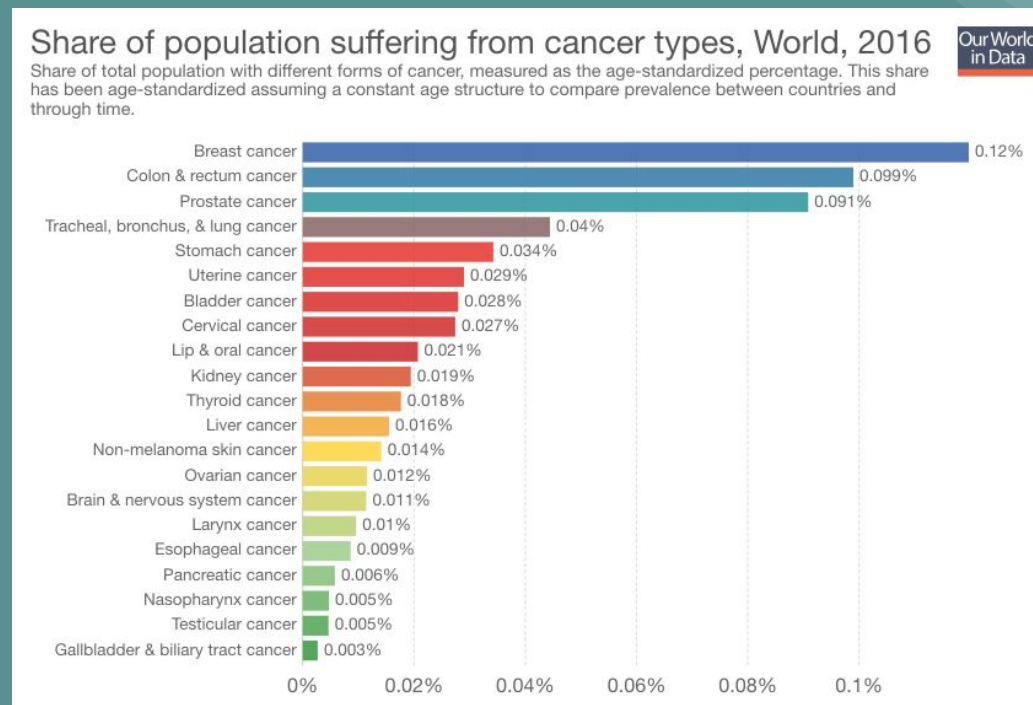# BIOPSY CLASSIFICATION OF CERVICAL CANCER

# Introduction

Third most common cancer in women worldwide.

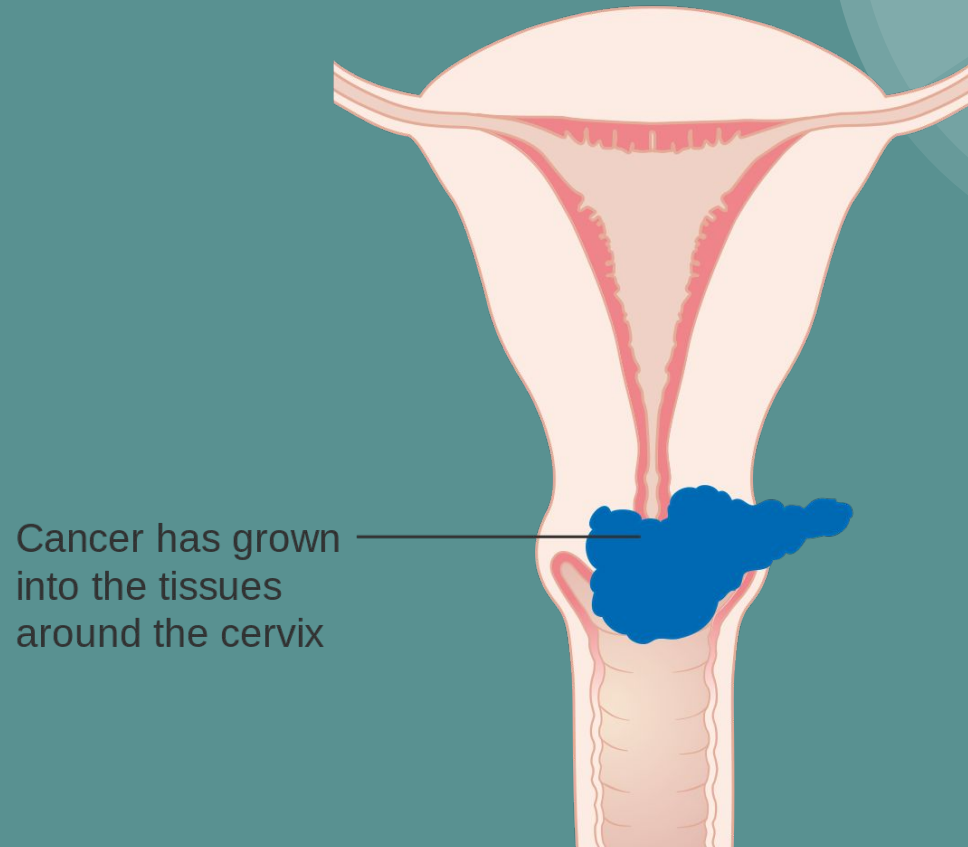500,000 women worldwide die of cervical cancer annually, of which 4,100 deaths are from the USA.

Most common cause of cancer death where proper tests are not available.

Easiest gynecologic cancer to prevent through screening.



Share of population suffering from cancer types, World, 2016
Share of total population with different forms of cancer, measured as the age-standardized percentage. This share has been age-standardized assuming a constant age structure to compare prevalence between countries and through time.

# Cervical Cancer

What causes cervical cancer: HPV (human papillomavirus). It's the most common sexually transmitted infection. **HPV** is usually harmless and goes away by itself, but some types can lead to cancer.

Cancer has grown into the tissues around the cervix

# Activities that increase the risk of HPV and for cervical cancer are:

- Numbers of sexual partners

- Use of birth control contraception pills for extended periods of time

- Numbers of Sexually Transmitted Diseases

- Giving birth to three or more children

# How it's diagnosed

Screening: The earlier cervical cancer is detected, the higher the chance of it being treated successfully. Pap tests are used to detect abnormal cells within the cervix.

Biopsy: If cancer is suspected, the doctor will take a sample of cervical cells (biopsy) and conduct either a Punch or Cone Biopsy

**SEER Relative Survival Rates by Stage at Diagnosis**
**For Cervix Uteri Cancer, All Races, All Ages, Females**
**SEER 9 Registries for 1988-2003**

|            | Localized | Regional | Distant | Unstaged |
|------------|-----------|----------|---------|----------|
| Time zero  | 100.0%    | 100.0%   | 100.0%  | 100.0%   |
| 1-year     | 98.3%     | 84.3%    | 44.6%   | 82.4%    |
| 2-year     | 96.0%     | 69.3%    | 26.6%   | 75.1%    |
| 3-year     | 94.2%     | 61.1%    | 20.1%   | 71.0%    |
| 4-year     | 92.9%     | 56.2%    | 16.6%   | 67.7%    |
| 5-year     | 92.1%     | 53.3%    | 14.4%   | 66.4%    |
| 6-year     | 91.3%     | 51.1%    | 13.4%   | 64.7%    |
| 7-year     | 90.6%     | 49.6%    | 12.9%   | 64.3%    |
| 8-year     | 90.0%     | 48.0%    | 12.4%   | 63.9%    |
| 9-year     | 89.2%     | 46.8%    | 12.1%   | 62.6%    |
| 10-year    | 88.6%     | 46.0%    | 11.9%   | 61.4%    |

# Initial Data

```
     Age  Number of sexual partners  ...  Citology  Biopsy
0    18                         4.0  ...         0       0
1    15                         1.0  ...         0       0
2    34                         1.0  ...         0       0
3    52                         5.0  ...         0       0
4    46                         3.0  ...         0       0

[5 rows x 36 columns]
(858, 36)
```

```
     Age  Number of sexual partners  First sexual intercourse  ...  Dx:HPV  Dx  Biopsy
0    18                         4.0                      15.0  ...       0   0       0
1    15                         1.0                      14.0  ...       0   0       0
2    34                         1.0                         ?  ...       0   0       0
3    52                         5.0                      16.0  ...       1   0       0
4    46                         3.0                      21.0  ...       0   0       0

[5 rows x 33 columns]
(858, 33)
```

# Prior to the Preprocessing of the Data

| Model | Accuracy |
|---|---|
| Decision Tree | 94.44 |
| Random Forest | 83.33 |
| Support Vector Machine | 88.88 |
| *K-Nearest Neighbors | 83.33 |
| Naïve Bayes | 83.33 |

The table below demonstrates the comparison of accuracies of all the models run prior to any cleaning methodology

# Prior to the Preprocessing of the Data Continued

| Models | Precision | Recall | F1 Score |
|---|---|---|---|
| Decision Tree | 1 | 0.93 | 0.97 |
| Random Forest | 0.83 | 1.00 | 0.91 |
| Support Vector Machine | 0.88 | 1.00 | 0.94 |
| K-Nearest Neighbors | 0.83 | 1 | 0.91 |
| Naïve Bayes | 0.88 | 0.93 | 0.90 |

For women without cervical cancer

| Models | Precision | Recall | F1 Score |
|---|---|---|---|
| Decision Tree | 0.75 | 1 | 0.86 |
| Random Forest | 0 | 0 | 0 |
| Support Vector Machine | 1 | 0.33 | 0.50 |
| K-Nearest Neighbors | 0 | 0 | 0 |
| Naïve Bayes | 0.50 | 0.33 | 0.40 |

For women with cervical cancer

# ROC Curves

# ROC Curves
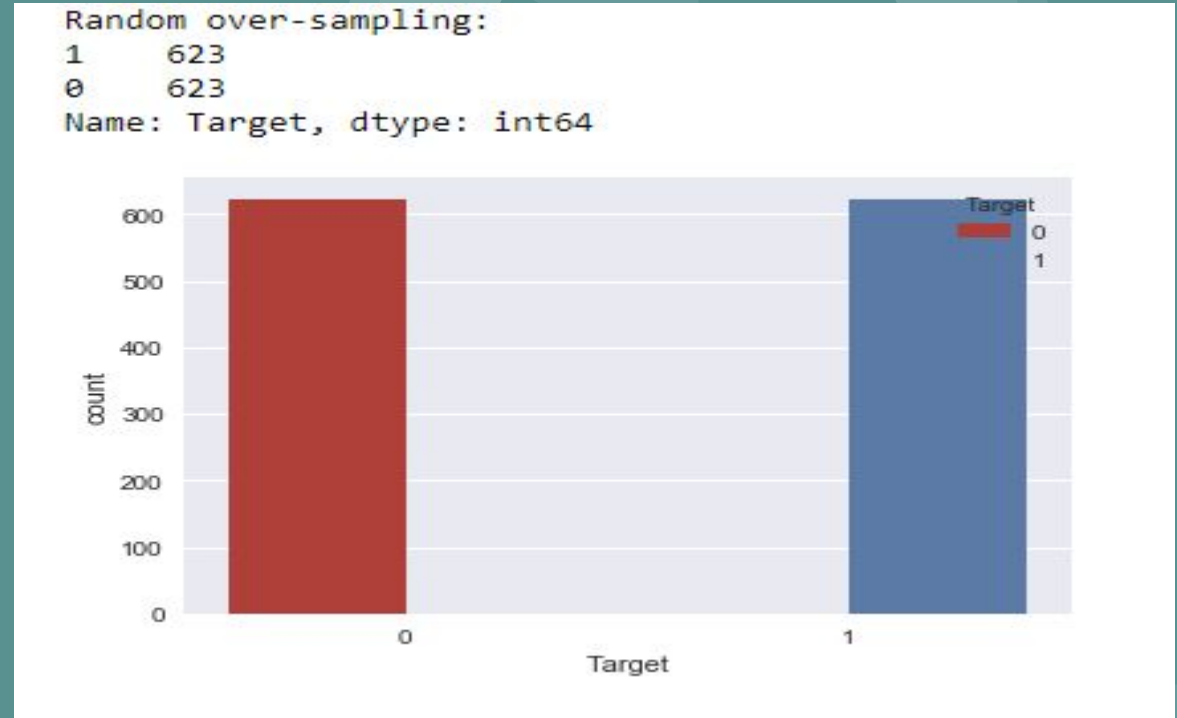
# ROC Curves

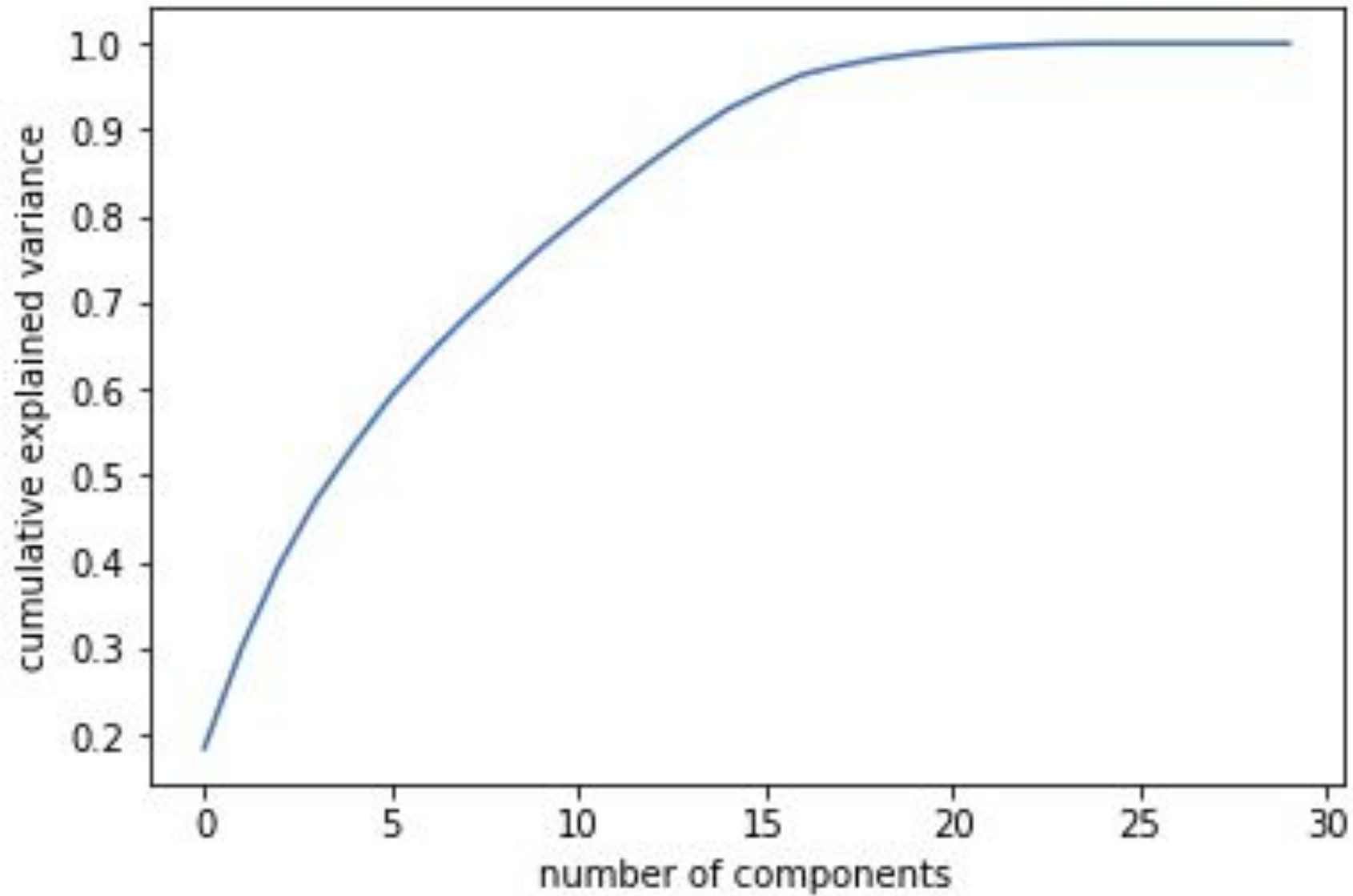# Confusion Matrices

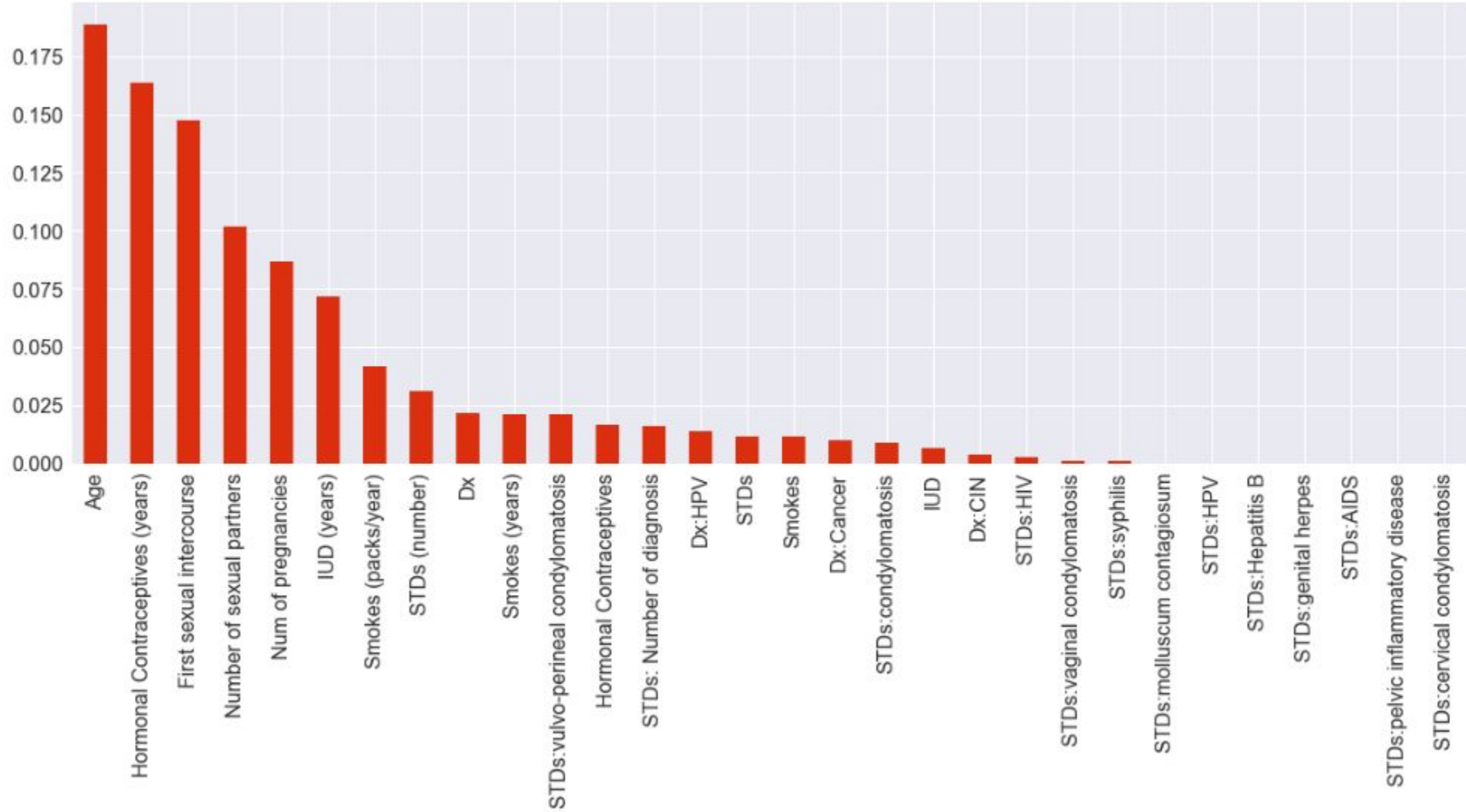# Imbalanced vs. Balanced Class Plots



Imbalanced Class Plot

Balanced Class Plot

# PCA Plot

# Feature Selection

# 7 Highest Predictive Powers

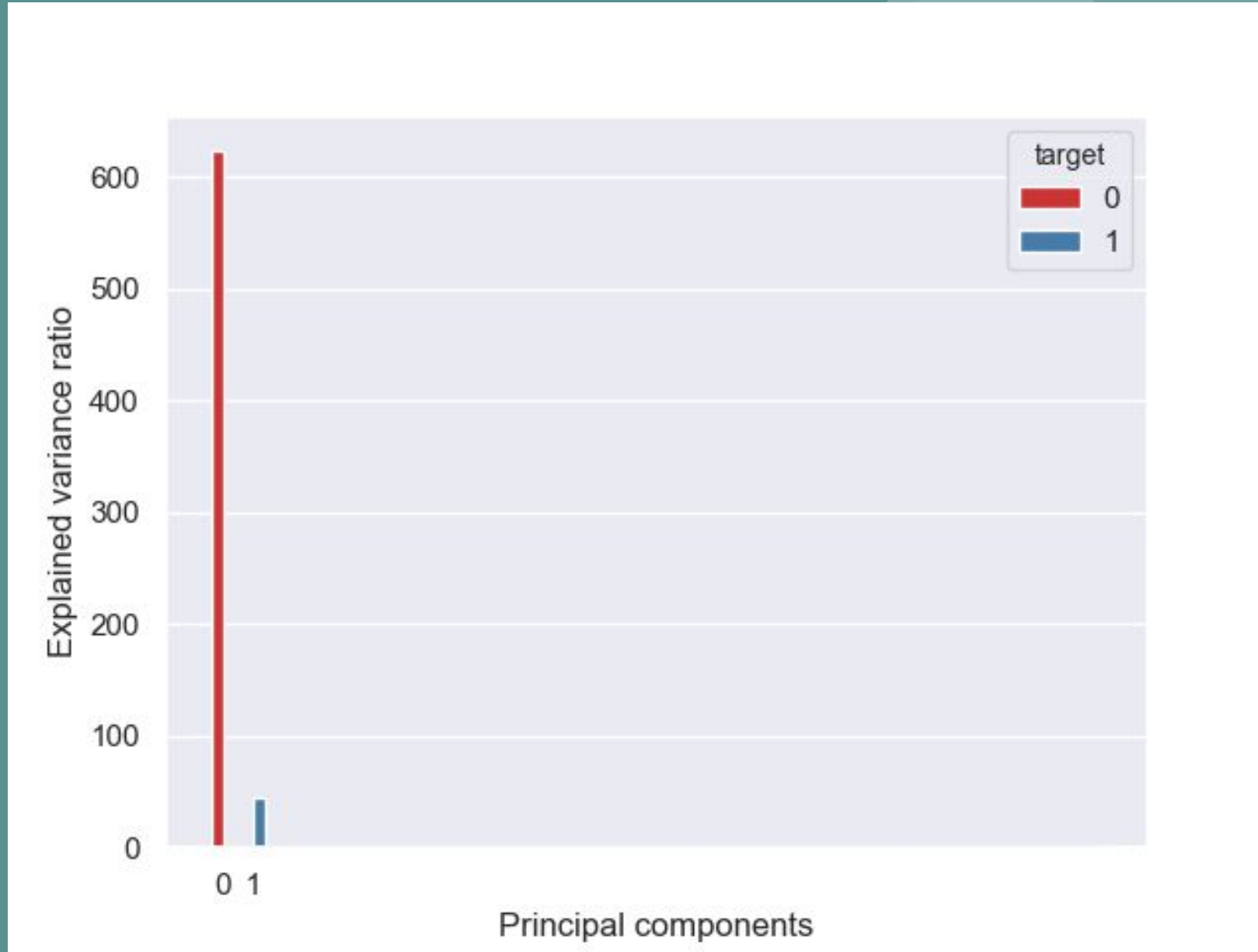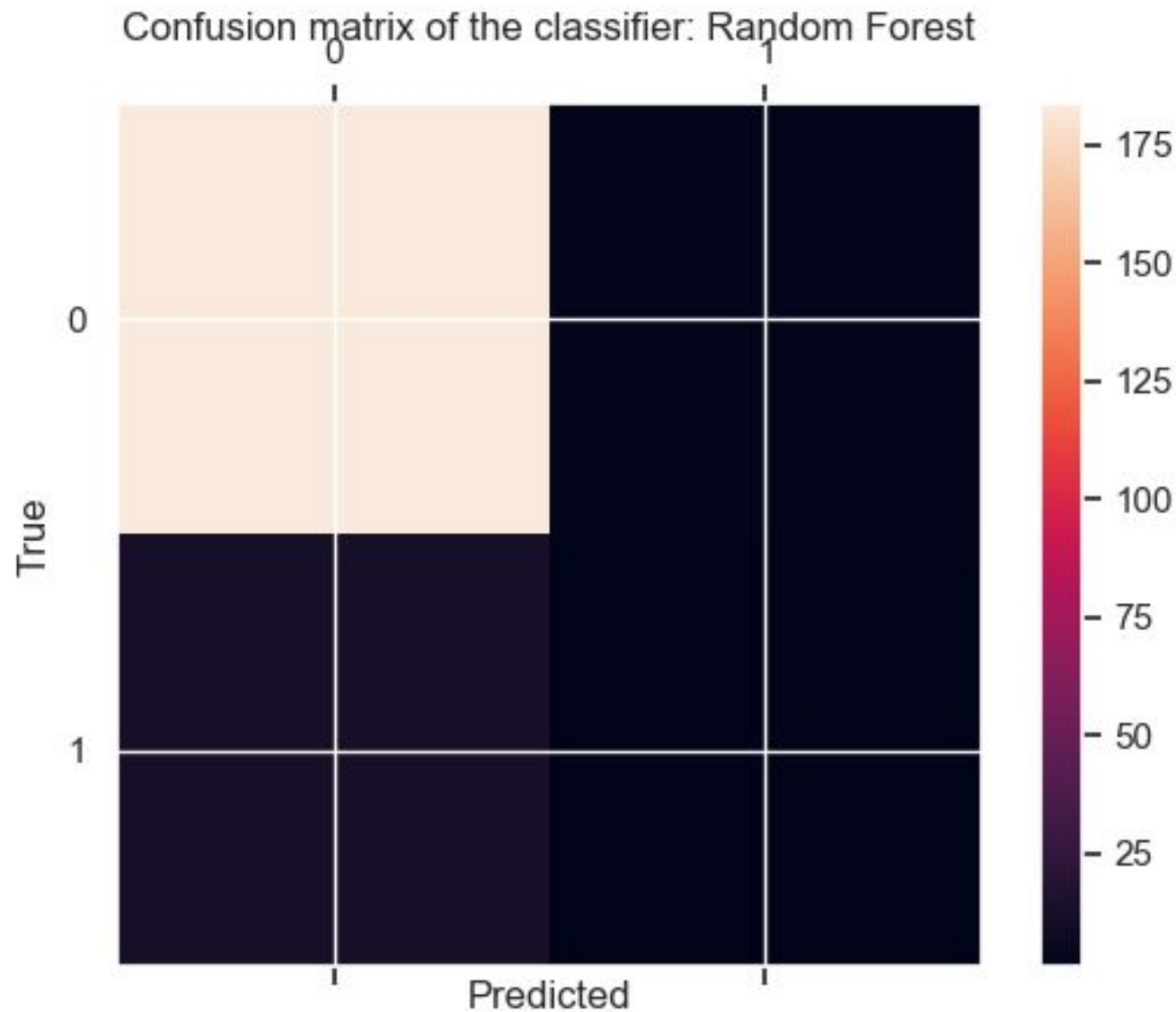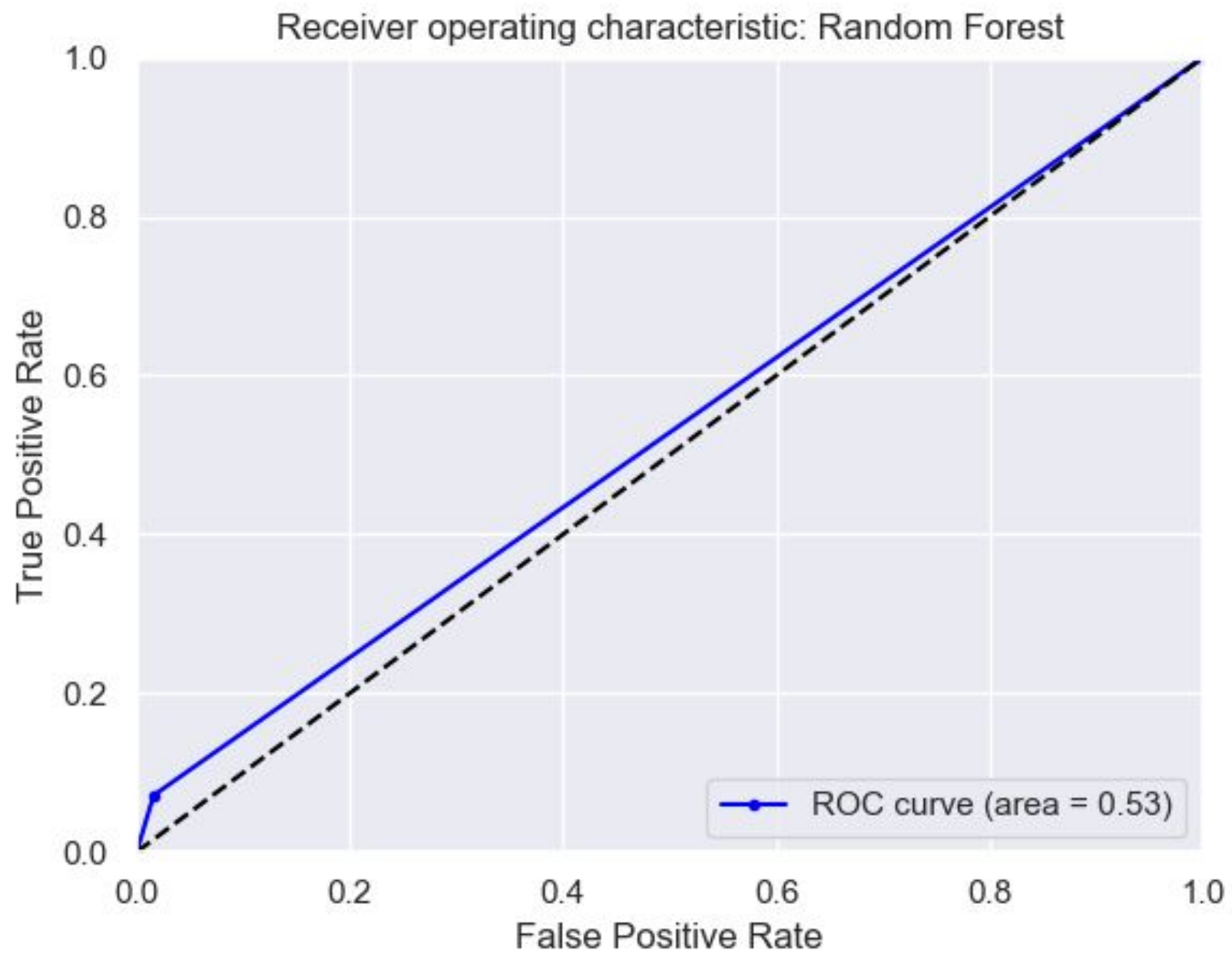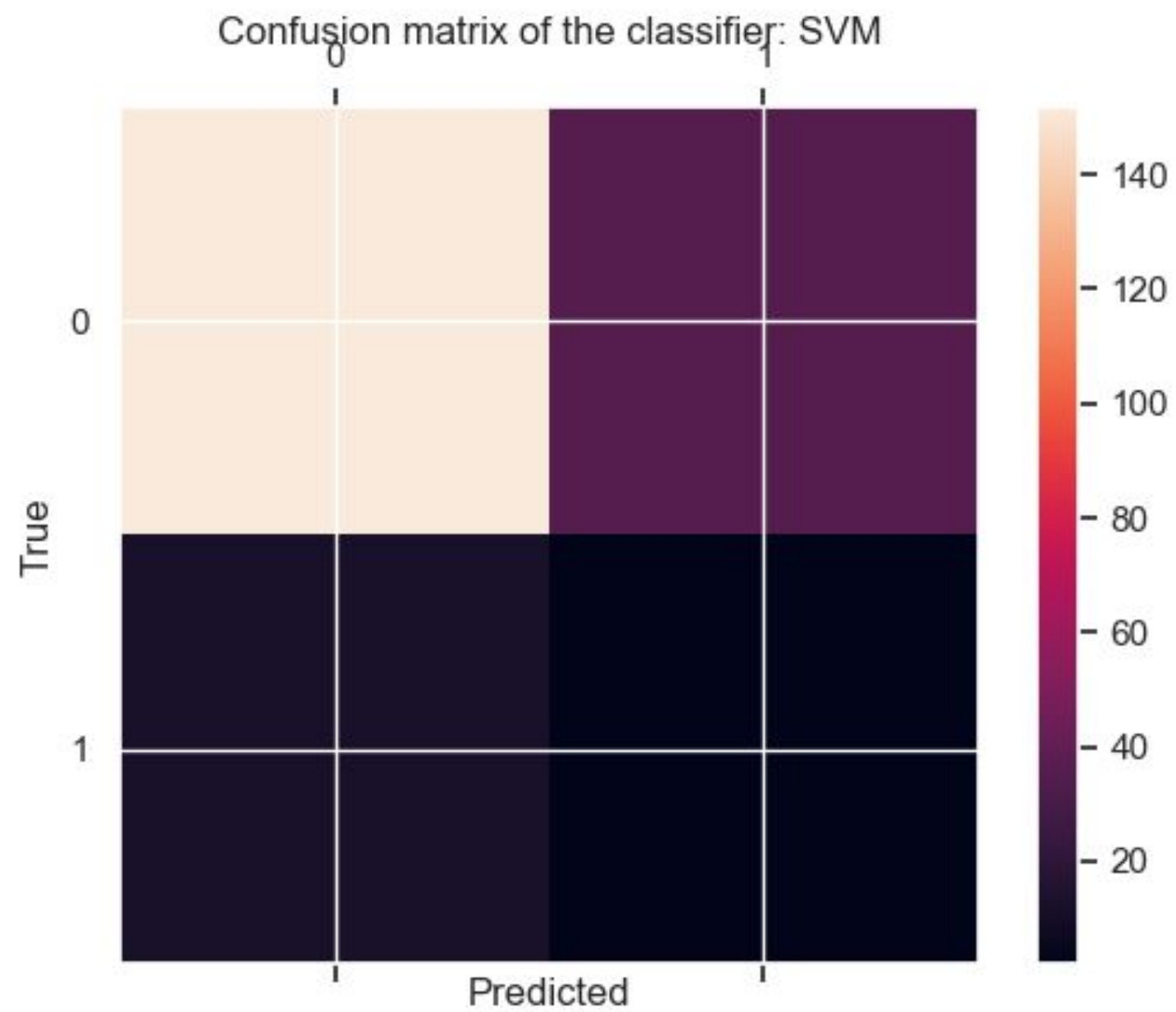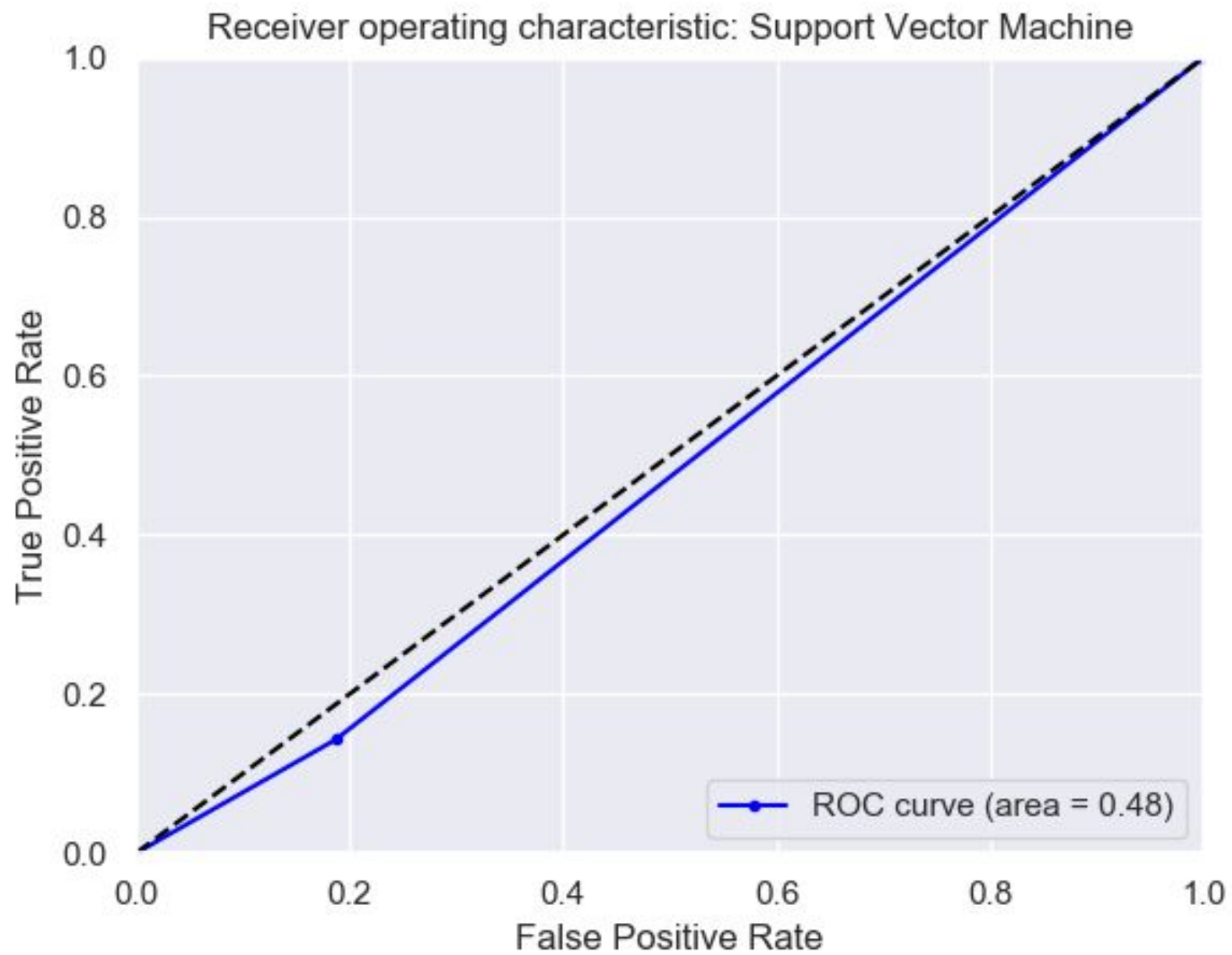| | |
|---|---|
| 0 | Age |
| 1 | Number of sexual partners |
| 2 | First sexual intercourse |
| 3 | Num of pregnancies |
| 4 | Smokes (packs/year) |
| 5 | Hormonal Contraceptives (years) |
| 6 | IUD (years) |
| 7 | STDs (number) |

# PC 1 vs. PC 2

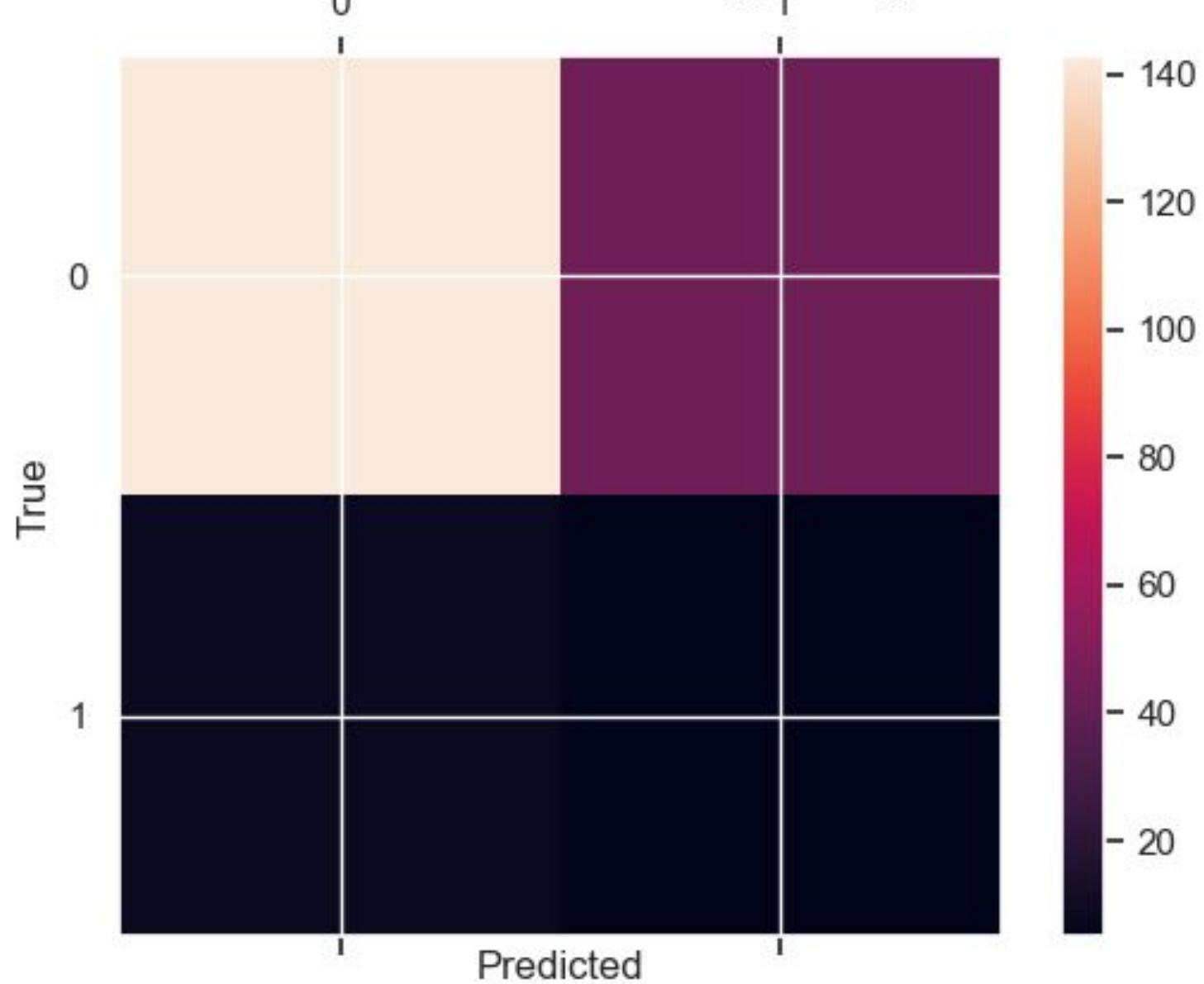# Variance Ratio vs. Principal Components
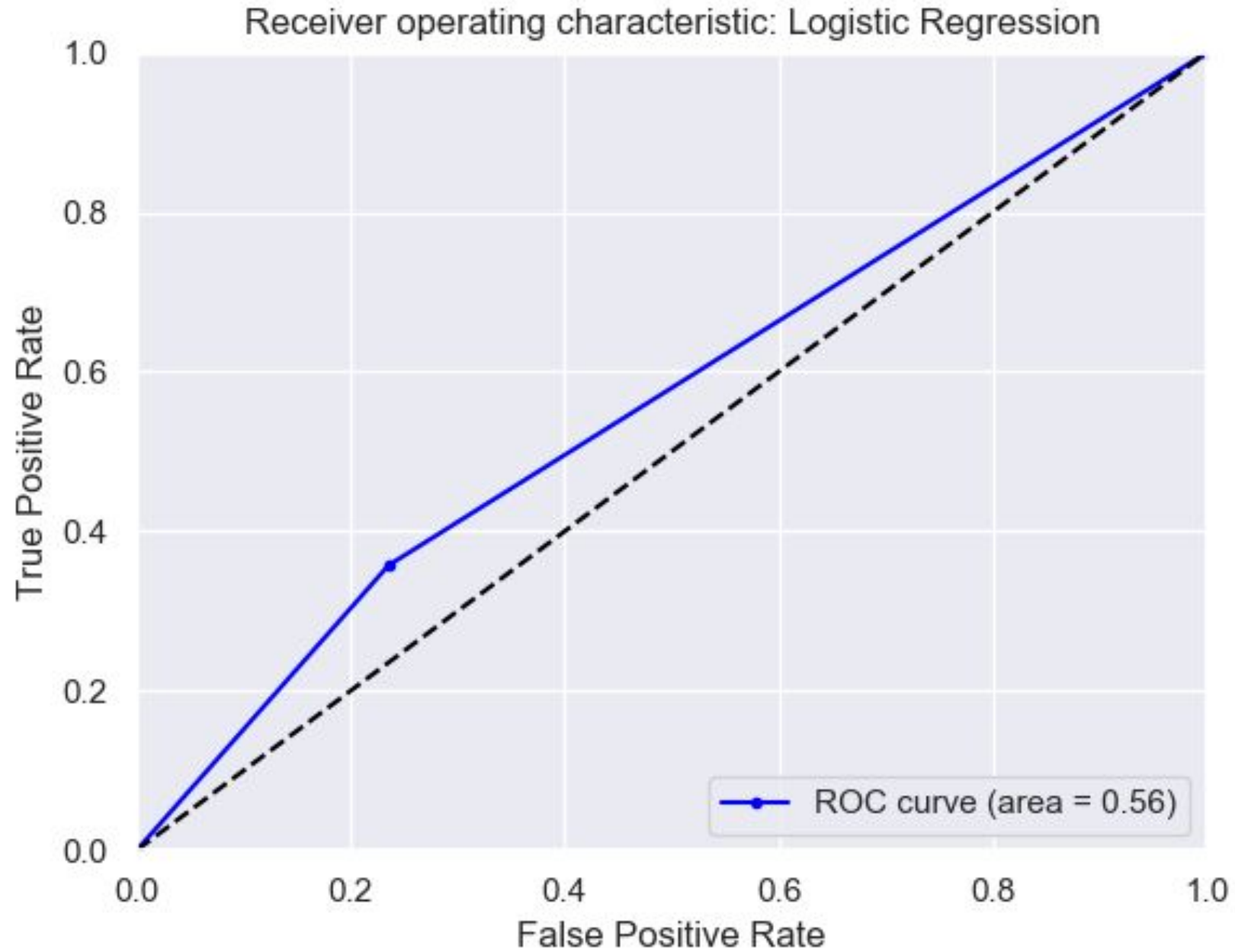
# After the Preprocessing of the Data



Confusion matrix of the classifier: Random Forest

Receiver operating characteristic: Random Forest

Receiver operating characteristic: Support Vector Machine

Confusion matrix of the classifier: Logistic Regression

Receiver operating characteristic: Logistic Regression

```
The accuracy of Logistic Regression  is: 0.736318407960199


[0 1]
Classification Report:
              precision      recall   f1-score     support

           0       0.94        0.76       0.84         187
           1       0.10        0.36       0.16          14
```

# Conclusion

- The 10 ten features with the highest predictive power were confirmed by the literature

- Preprocessing had a impact on the performance of our model

- The model could be run on other targets identified in the data

- Other classifiers could be run on this data