

Semiconductors are the brains of modern electronics. They enable technologies critical to a country's economic growth, national security, and global competitiveness. From smartphones to planes, semiconductors evolved to improve technologies and do wonders for your entertainment and convenience. Industry sectors that rely a lot on semiconductors are Computing, Telecommunications, Household Appliances, Banking, Security, Automotive/Transportation, Healthcare and Manufacturing [1].

To build up the devices of a chip we need the basic substrate, the wafers. A wafer is a thin slice of semiconductor, such as a crystalline silicon (Si), used for the fabrication of integrated circuits. While the silicon is the basic material, there are wafers which include other elements.

Semiconductor manufacturers request wafers from suppliers based on the volume of orders and their experiments for process and technology enhancement and development. This dataset is a use-case from one of India's leading manufacturers of wafers. No wafer manufacturer want to create products with anomalies. There are two datasets, the Train.csv and Test.csv. The Train.csv includes features (without indicating the role during the production) and the class of Fail or Pass for the wafer. The role of this project is to find a model that forecasts the wafers with anomalies.

The dataset has been taken by kaggle:

<https://www.kaggle.com/datasets/arbazkhan971/anomaly-detection>

[1] <https://www.makeuseof.com/why-semiconductors-important/>

Import libraries

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

Read the data

```
In [2]: data = pd.read_csv("Train.csv")
data.head()
```

```
Out[2]:
```

	feature_1	feature_2	feature_3	feature_4	feature_5	feature_6	feature_7	feature_8	feature_9	feature_10	...	feature_1550	feature_1559	Class
0	100	160	16000	0	0	0	0	0	0	0	...	0	...	0
1	20	83	41500	1	0	0	0	0	0	0	1	0
2	99	150	15151	1	0	0	0	0	0	0	0	0
3	40	40	10000	0	0	0	0	0	0	0	0	0
4	12	234	195000	1	0	0	0	0	0	0	0	0

5 rows x 1559 columns

Data Cleaning

Are there Nan values in the dataset?

```
In [3]: data.isna().isna()
```

```
Out[3]:
```

	feature_1	feature_2	feature_3	feature_4	feature_5	...	feature_1555	feature_1556	feature_1557	feature_1558	feature_1559	Class
0	False	False	False	False	False	...	False	False	False	False	False	False
1	False	False	False	False	False	...	False	False	False	False	False	False
2	False	False	False	False	False	...	False	False	False	False	False	False
3	False	False	False	False	False	...	False	False	False	False	False	False
4	False	False	False	False	False	...	False	False	False	False	False	False

Length: 1559, dtype: bool

```
In [4]: colours = ['#000099', '#ffff00'] # specify the colours - yellow is missing. blue is not missing.
fig, ax = plt.subplots(1,2,figsize=(20,10))
fig.tight_layout(pad=0.0)
sns.heatmap(data.isna(), ax=ax[0], cmap=sns.color_palette(colours))
sns.heatmap(data.isna().T, ax=ax[1], cmap=sns.color_palette(colours));
```

There are no missing values in the dataset.

What type of data do we have?

```
In [5]: data.dtypes[data.columns]
```

```
Out[5]:
```

	feature_1	feature_2	feature_3	feature_4	feature_5	feature_1555	feature_1556	feature_1557	feature_1558	feature_1559	Class
0	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64
1	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64
2	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64
3	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64
4	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64	int64

Length: 1559, dtype: object

```
In [6]: data.dtypes[data.columns].nunique()
```

```
Out[6]:
```

	feature_1	feature_2	feature_3	feature_4	feature_5	feature_1555	feature_1556	feature_1557	feature_1558	feature_1559	Class
0	1	1	1	1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1	1	1	1	1
2	1	1	1	1	1	1	1	1	1	1	1
3	1	1	1	1	1	1	1	1	1	1	1
4	1	1	1	1	1	1	1	1	1	1	1

The data types of the variables are two, int64 and float64. How many features have int64 and how many float64?

- **int64**

```
In [7]: (data.dtypes[data.columns]=="int64").sum()
```

```
Out[7]:
```

	feature_1	feature_2	feature_3	feature_4	feature_5	feature_1555	feature_1556	feature_1557	feature_1558	feature_1559	Class
0	1	1	1	1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1	1	1	1	1
2	1	1	1	1	1	1	1	1	1	1	1
3	1	1	1	1	1	1	1	1	1	1	1
4	1	1	1	1	1	1	1	1	1	1	1

1558

- **float64**

```
In [8]: (data.dtypes[data.columns]=="float64").sum()
```

```
Out[8]:
```

	feature_1	feature_2	feature_3	feature_4	feature_5	feature_1555	feature_1556	feature_1557	feature_1558	feature_1559	Class
0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	0

So, there is only one float column. The rest columns have integer numbers.

Let's identify which column has float numbers.

```
In [9]: which_float = (data.dtypes[data.columns]=="float64")
which_float[which_float==True].index
```

```
Out[9]:
```

	feature_1	feature_2	feature_3	feature_4	feature_5	feature_1555	feature_1556	feature_1557	feature_1558	feature_1559	Class
0	False	False	False	False	False	False	False	False	False	False	False
1	False	False	False	False	False	False	False	False	False	False	False
2	False	False	False	False	False	False	False	False	False	False	False
3	False	False	False	False	False	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False	False	False	False

The data type of **feature_3** is float.

Let's try to understand the rest features, those with data type int64.

Let's see how many of them have a single value, 2 unique values and how many have 3 and more. I will not list the last variable which is the **Class**.

```
In [10]: binary_features = categorical_features
single_value_features = [], [], []
for x in (col for col in data.columns if col not in ['feature_3', 'Class']):
    number_of_categories = data[x].nunique()
    if number_of_categories==1:
        single_value_features.append(x)
    elif number_of_categories==2:
        binary_features.append(x)
    else:
        categorical_features.append(x)
```

Let's do a test, is the sum of the lengths of the lists plus 2 (for the columns **feature_3** and **Class**) equal to the total number of columns of the dataset?

```
In [11]: len(binary_features) + len(single_value_features) + len(categorical_features) + 2 == data.shape[1]
```

```
Out[11]:
```

	feature_1	feature_2	feature_3	feature_4	feature_5	feature_1555	feature_1556	feature_1557	feature_1558	feature_1559	Class
0	True	True	True	True	True	True	True	True	True	True	True
1	True	True	True	True	True	True	True	True	True	True	True
2	True	True	True	True	True	True	True	True	True	True	True
3	True	True	True	True	True	True	True	True	True	True	True
4	True	True	True	True	True	True	True	True	True	True	True

True

Binary features:

```
In [12]: len(binary_features)
```

```
Out[12]:
```

	feature_1	feature_2	feature_3	feature_4	feature_5	feature_1555	feature_1556	feature_1557	feature_1558	feature_1559	Class
0	1	1	1	1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1	1	1	1	1
2	1	1	1	1	1	1	1	1	1	1	1
3	1	1	1	1	1	1	1	1	1	1	1
4	1	1	1	1	1	1	1	1	1	1	1

1519

Categorical features:

```
In [13]: len(categorical_features)
```

```
Out[13]:
```

	feature_1	feature_2	feature_3	feature_4	feature_5	feature_1555	feature_1556	feature_1557	feature_1558	feature_1559	Class
0	1	1	1	1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1	1	1	1	1
2	1	1	1	1	1	1	1	1	1	1	1
3	1	1	1	1	1	1	1	1	1	1	1
4	1	1	1	1	1	1	1	1	1	1	1

2

single_value_features:

```
In [14]: len(single_value_features)
```

```
Out[14]:
```

	feature_1	feature_2	feature_3	feature_4	feature_5	feature_1555	feature_1556	feature_1557	feature_1558	feature_1559	Class
0	1	1	1	1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1	1	1	1	1
2	1	1	1	1	1	1	1	1	1	1	1
3	1	1	1	1	1	1	1	1	1	1	1
4	1	1	1	1	1	1	1	1	1	1	1

36

The categorical features are only 2, let's try to understand them:

```
In [15]: data[categorical_features].describe()
```

```
Out[15]:
```

	feature_1	feature_2
count	1763.000000	1763.000000
mean	53.094158	126.967067
std	55.842014	129.859641
min	1.000000	1.000000
25%	12.000000	33.500000
50%	39.000000	96.000000
75%	75.000000	159.000000
max	640.000000	640.000000

```
In [16]: sns.boxplot(data = data[categorical_features])
```

```
Out[16]:
```

These two variables have data types integer by definition, although they are continuous!! Thus, the variable which have a float data type and the categorical_features list will be merged into a single list:

```
In [17]: continuous_features = categorical_features + ['feature_3']
```

Let's see the single_value_features:

```
In [18]: data[single_value_features]
```

```
Out[18]:
```

	feature_57	feature_82	feature_106	feature_147	feature_262	feature_278	feature_284	feature_320	feature_362	feature_37
0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0

1763 rows x 36 columns

I will remove those variables

```
In [19]: data = data.drop(single_value_features, axis=1)
```

```
In [20]: data.head()
```

```
Out[20]:
```

	feature_1	feature_2	feature_3	feature_4	feature_5	feature_6	feature_7	feature_8	feature_9	feature_10	...	feature_1550	feature_1559	Class
0	100	160	16000	0	0	0	0	0	0	0	...	0	...	0
1	20	83	41500	1	0	0	0	0	0	0	1	0
2	99	150	15151	1	0	0	0	0	0	0	0	0
3	40	40	10000	0	0	0	0	0	0	0	0	0
4	12	234	195000	1	0	0	0	0	0	0	0	0

5 rows x 1523 columns

So, I have two lists:

- the **binary_features**, which include all variables with two categories
- the **continuous_features**, which include continuous variables

There are a lot of variables. Are there duplicated features?

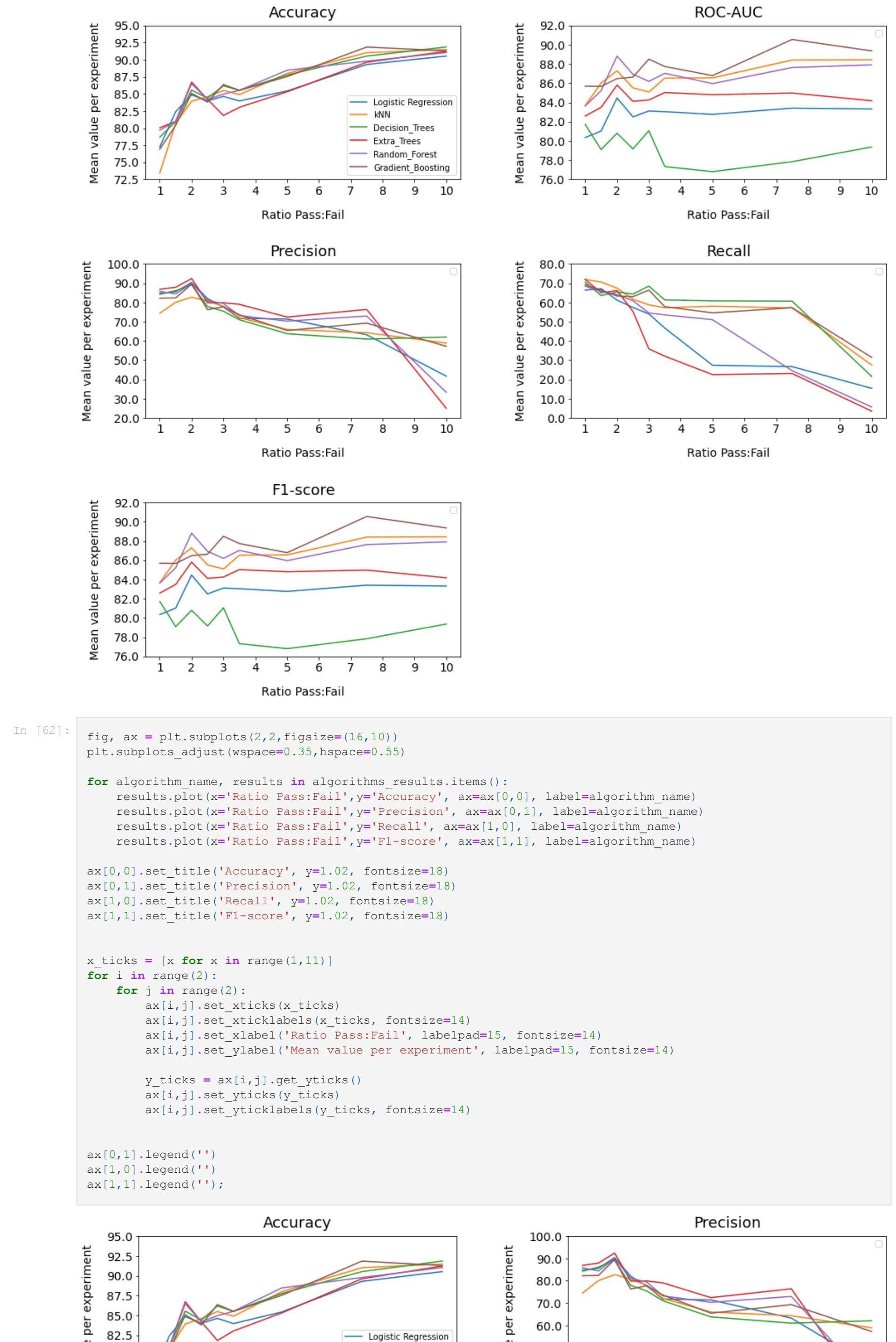
I can check the binary_variables. I will create a dataframe which will include two values, 0 & 1. Each feature will be tested if it is equal with another and if so the dataframe will include the 1, otherwise the 0.

```
In [21]: duplicated_columns = []
for x in binary_features:
    x_feature_check_with_y = []
    for y in binary_features:
        if data[x].equals(data[y])==True:
            x_feature_check_with_y.append(1)
        else:
            x_feature_check_with_y.append(0)
    duplicated_columns.append(x_feature_check_with_y)
```

```
In [22]: duplicated_columns = pd.DataFrame(duplicated_columns, columns=binary_features, index=binary_features)
duplicated_columns
```

```
Out[22]:
```

	feature_4	feature_5	feature_6	feature_7	feature_8	feature_9	feature_10	feature_11	feature_12	feature_13	...	feature_1550	feature_1559
feature_4	1	0	0	0	0	0	0	0	0	0	...	0	...
feature_5	0	1	0	0	0	0	0	0	0	0	0
feature_6	0	0	1	0	0	0	0	0	0	0	0
feature_7	0	0	0	1	0	0	0	0	0	0	0



```
In [62]: fig, ax = plt.subplots(2,2,figsize=(16,10))
plt.subplots_adjust(wspace=0.35,hspace=0.55)

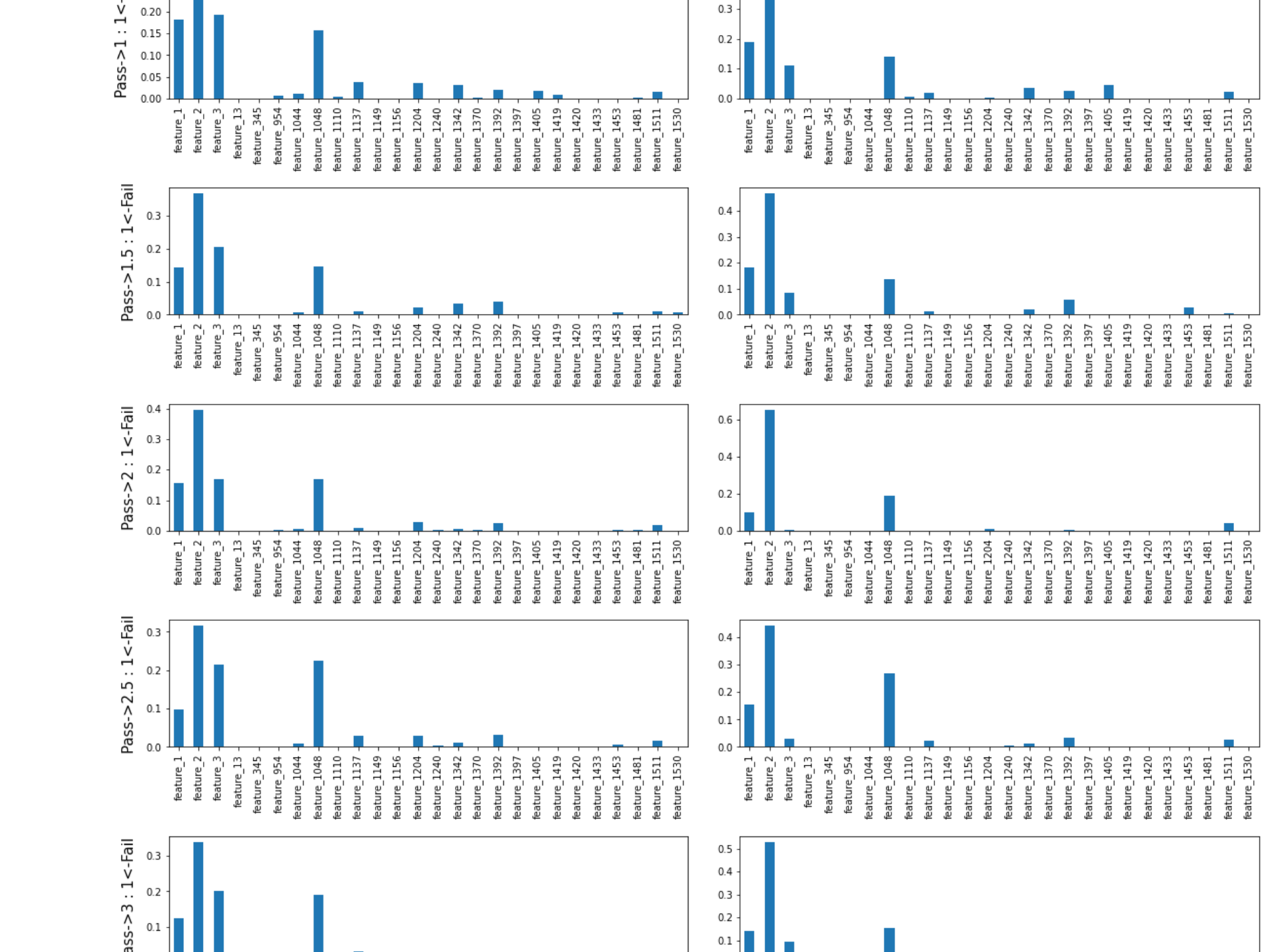
for algorithm_name, results in algorithms_results.items():
    results.plot(x='Ratio Pass-Fail',y='Accuracy',ax=ax(0,0),label=algorithm_name)
    results.plot(x='Ratio Pass-Fail',y='Precision',ax=ax(0,1),label=algorithm_name)
    results.plot(x='Ratio Pass-Fail',y='Recall',ax=ax(1,0),label=algorithm_name)
    results.plot(x='Ratio Pass-Fail',y='F1-score',ax=ax(1,1),label=algorithm_name)

ax[0,0].set_title('Accuracy',y=1.02,fontsize=18)
ax[0,1].set_title('Precision',y=1.02,fontsize=18)
ax[1,0].set_title('Recall',y=1.02,fontsize=18)
ax[1,1].set_title('F1-score',y=1.02,fontsize=18)

x_ticks = [x for x in range(1,11)]
for i in range(2):
    for j in range(2):
        ax[i,j].set_xticks(x_ticks)
        ax[i,j].set_xticklabels(x_ticks,fontsize=14)
        ax[i,j].set_xlabel('Ratio Pass-Fail',labelpad=5,fontsize=14)
        ax[i,j].set_ylabel('Mean value per experiment',labelpad=5,fontsize=14)

        y_ticks = ax[i,j].get_yticks()
        ax[i,j].set_yticklabels(y_ticks,fontsize=14)

ax[0,1].legend('')
ax[1,0].legend('')
ax[1,1].legend('')
```



Let's find out the optimal model:

```
In [64]: searching_optimal_fl, store_algorithm = [], []
for algorithm_name, results in algorithms_results.items():
    searching_optimal_fl.append(results['F1-score'].values)

algorithm, exp_sample = np.where(np.array(searching_optimal_fl)==np.array(searching_optimal_fl).max())
best_model_method_2 = best_estimator_per_size[exp_sample[0]][algorithm[0]]
best_model_method_2
```

```
Out[64]: GradientBoostingClassifier(max_leaf_nodes=8)
```

The best estimator is given by `best_model_method_2`

The ratio is Class Pass -> 75 : 1 <= Class Fail

Store the predictions into a new column.

```
In [66]: predictions_method_2 = best_model_method_2.predict(data_for_predictions.iloc[:,1:])
data_method_predictions['Pred. Method 2'] = predictions_method_2
```

Conclusions

This method was mainly an experiment in order to see how the size of the majority class play a role in the result. Precision and Recall are decreasing while the size of the majority class increases. On the one hand, the F1-score, the basic metric, indicates decrease for some models while the sample of the majority class increases, and the optimal algorithm, the GradientBoostingClassifier, shows a stability on the other hand. The features importances points out that the continuous features and the feature_1048 are important independent variables for the model.

- **## Method 3:** Reduce the dimensions using FactorAnalysis and employing classifiers with

Sort the correlated features to the target variable based on the kendall method. I sort them and give them to the FactorAnalysis.

```
In [70]: row_of_columns = data_for_predictions[data_for_predictions.columns[:-1]].corrwith(data_for_predictions.Class, axis=0)
components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [71]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [72]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [73]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [74]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [75]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [76]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [77]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [78]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [79]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [80]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [81]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [82]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [83]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [84]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [85]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [86]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [87]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [88]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [89]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [90]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [91]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [92]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [93]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [94]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [95]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [96]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [97]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [98]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [99]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [100]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [101]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

```
In [102]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```

Let's create a pandas series and sort the values to see the distribution.

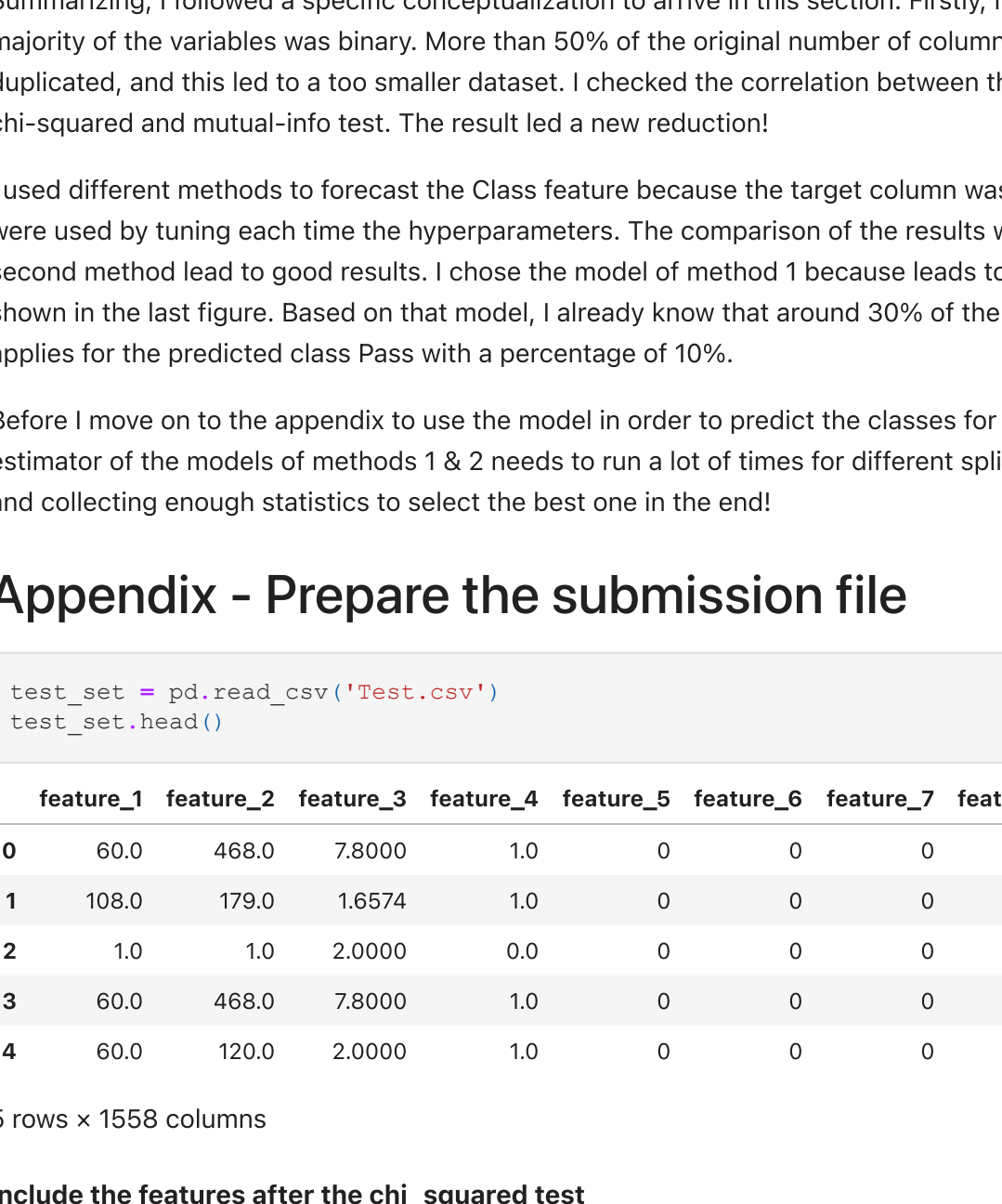
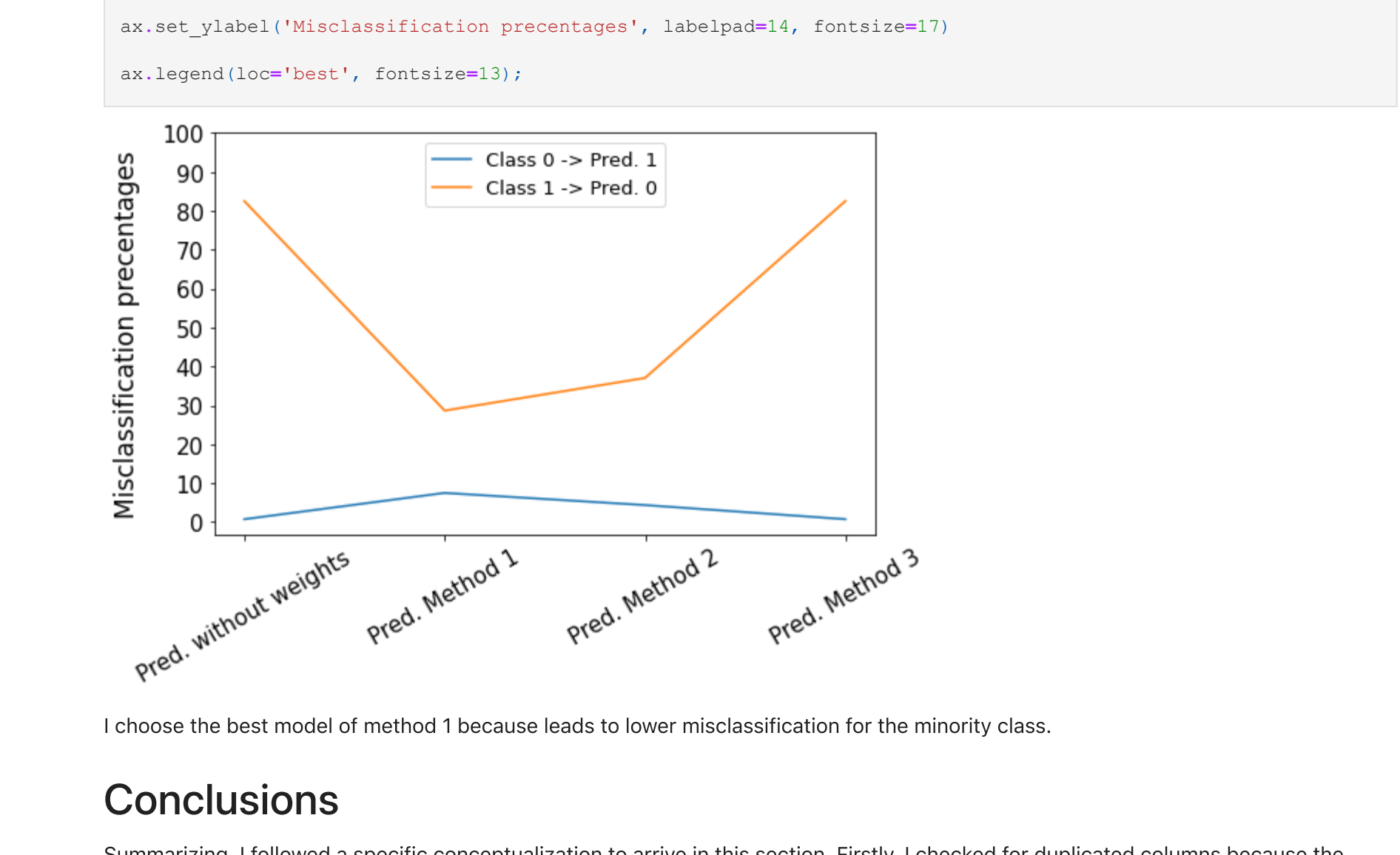
```
In [103]: components_dist_features = pd.Series(components_var, index=row_of_columns).sort_values(ascending=False)
```


For the column Pred. without weights :
Classified 0 but predicted to 1: 12
Classified 1 but predicted to 0: 118

For the column Pred. Method 1 :
Classified 0 but predicted to 1: 118
Classified 1 but predicted to 0: 41

For the column Pred. Method 2 :
Classified 0 but predicted to 1: 69
Classified 1 but predicted to 0: 53

For the column Pred. Method 3 :
Classified 0 but predicted to 1: 12
Classified 1 but predicted to 0: 118



I choose the best model of method 1 because leads to lower misclassification for the minority class.

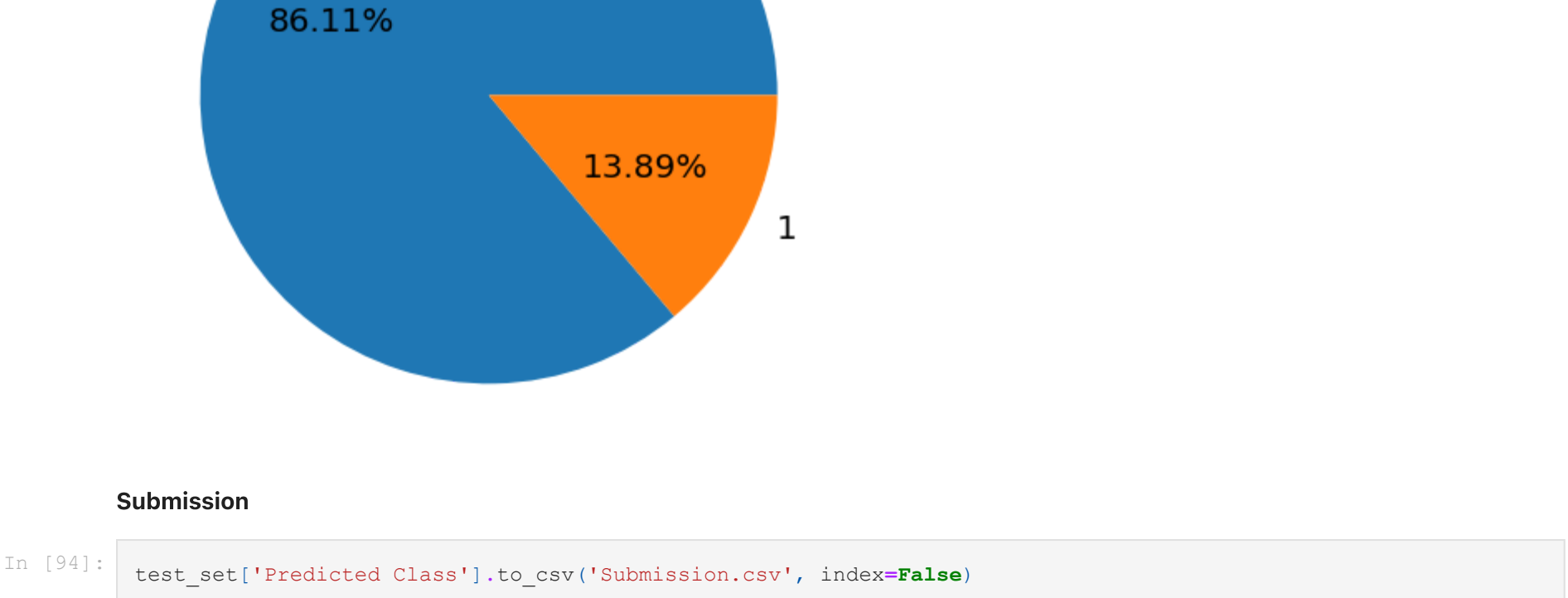
Conclusions

Summarizing, I followed a specific conceptualization to arrive in this section. Firstly, I checked for duplicated columns because the majority of the variables was binary. More than 50% of the original number of columns had to be removed, because they were duplicated, and this led to a too smaller dataset. I checked the correlation between the target feature and the binary features using chi-squared and mutual-info test. The result led a new reduction!

I used different methods to forecast the Class feature because the target column was unbalanced. In each method, same algorithms were used by tuning each time the hyperparameters. The comparison of the results was done by using the f1-score. The first and second method lead to good results. I chose the model of method 1 because leads to lower misclassification for the minority class, as shown in the last figure. Based on that model, I already know that around 30% of the predicted class fail is misclassified and same applies for the predicted class Pass with a percentage of 10%.

Before I move on to the appendix to use the model in order to predict the classes for the Test.csv file, I have to say that the best estimator of the models of methods 1 & 2 needs to run a lot of times for different splitted data by enhancing the hyperparameters and collecting enough statistics to select the best one in the end!

Appendix - Prepare the submission file



Include the features after the chi_squared test

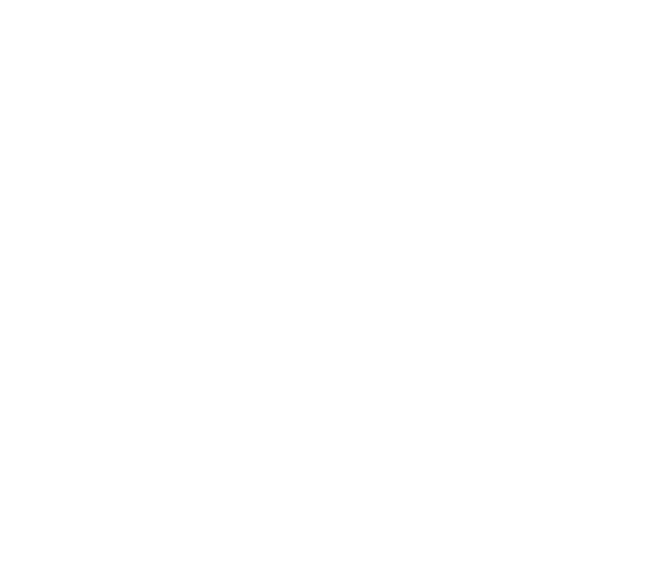


Scale the data



Apply the model of method 1

I already know that around 28% of Predicted class 1 must have been predicted as 0, and around 8% of class 0 to 1. Those are the misclassifications that I would wait for the test set, if I knew real Class.



Submission

