

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 (c-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 waters with anomalies.

The dataset has been taken by kaggle:

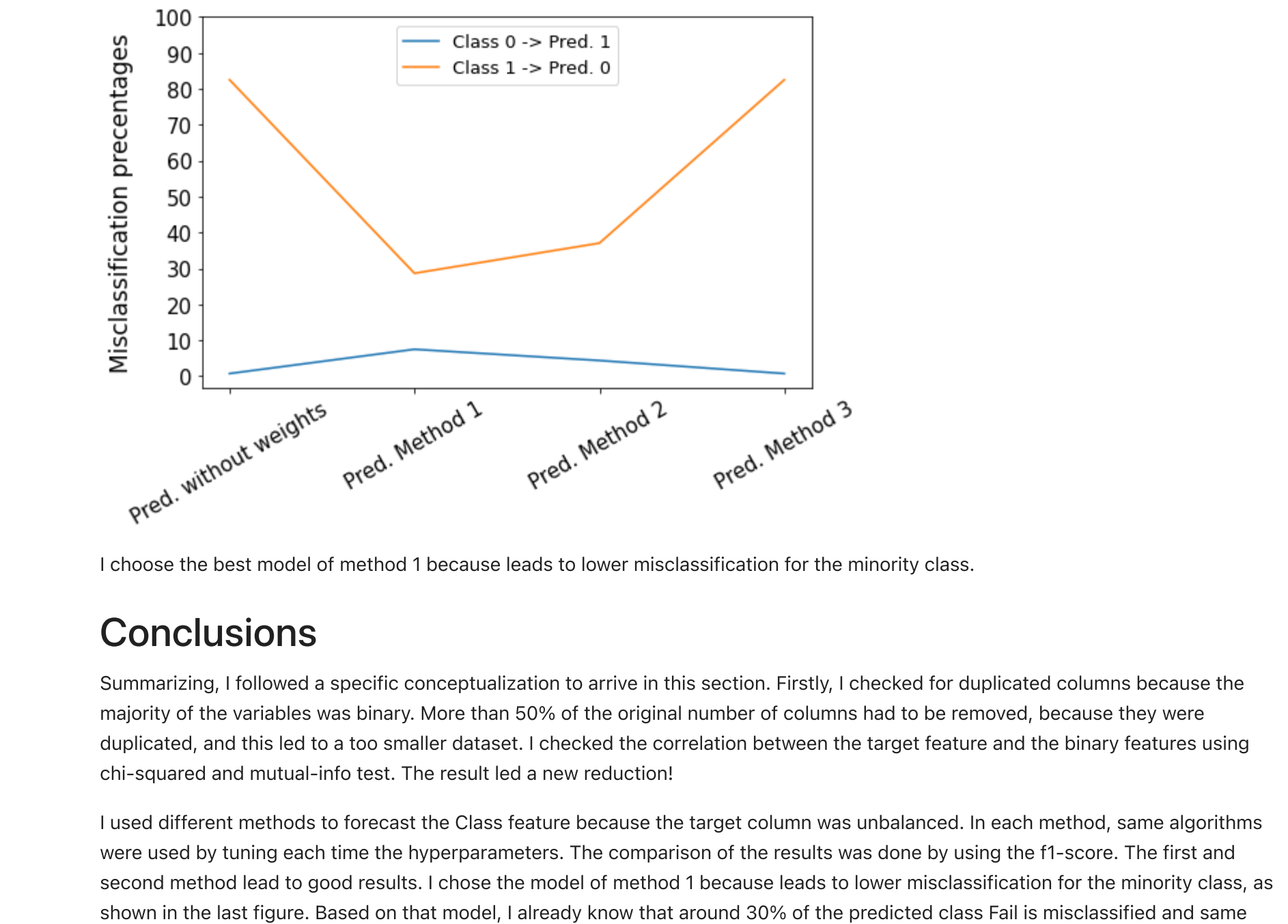
<https://www.kaggle.com/datasets/arbazkhan977/anomaly-detection>

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

Import libraries

```
import numpy as np
```

```
import matplotlib.pyplot as plt
import seaborn as sns
```

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

```
In [96]: test_set = pd.read_csv('Test.csv')
test_set.head()
```

	feature_1	feature_2	feature_3	feature_4	feature_5	feature_6	feature_7	feature_8	feature_9	feature_10	...	feature_1549	feature_1550
0	60.0	468.0	78000	1.0	0	0	0	0	0	0	...	0	0
1	108.0	179.0	1.6574	1.0	0	0	0	0	0	0	...	0	0
2	1.0	1.0	2.0000	0.0	0	0	0	0	0	0	...	0	0
3	60.0	468.0	78000	1.0	0	0	0	0	0	0	...	0	0
4	60.0	120.0	2.0000	1.0	0	0	0	0	0	0	...	0	0

5 rows × 1558 columns

Include the features after the chi_squared test

```
In [97]: test_set = test_set[continuous_features + data.columns[np.where(scores_chi2>=1)].tolist()]
test_set.head()
```

	feature_1	feature_2	feature_3	feature_13	feature_345	feature_954	feature_1044	feature_1048	feature_1110	feature_1137	...	feature_1549	feature_1550
0	60.0	468.0	78000	0	0	0	0	0	0	0	...	0	...
1	108.0	179.0	1.6574	0	0	0	0	0	0	0	...	1	...
2	1.0	1.0	2.0000	0	0	0	0	0	0	0	...	0	...
3	60.0	468.0	78000	0	0	0	0	0	0	0	...	0	...
4	60.0	120.0	2.0000	0	0	0	0	0	1	0	...	0	...

5 rows × 26 columns

Scale the data

```
In [98]: from sklearn.preprocessing import StandardScaler, MinMaxScaler
ss = StandardScaler()
mm = MinMaxScaler()

for col in test_set.columns[:3]:
    ss.fit(test_set[col].to_numpy().reshape(-1,1))
    test_set[col] = ss.transform(test_set[col].to_numpy().reshape(-1,1))
    test_set[col] = mm.fit_transform(test_set[col].to_numpy().reshape(-1,1)).reshape(1,-1)[0]
```

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.

```
In [99]: predictions_test_set = best_model_method_1.predict(test_set)
test_set['Predicted Class'] = predictions_test_set
```

```
In [100]: round((test_set['Predicted Class'].value_counts()/test_set.shape[0])*100,2).plot.pie(autopct='%1.2f%%',
figsize=(8,8), fontsize=20)
plt.title('Predicted Class', fontsize=20)
plt.ylabel('');
```

Predicted Class



Submission

```
In [94]: test_set['Predicted Class'].to_csv('Submission.csv', index=False)
```