

INTRODUCTION

The Willy Wonka chocolate factory produces the **Golden Scrumpilicious Candy Bar**, which packaged in a golden wrapper. However, sometimes during production the bars end up being green which is a defect, but also very alluring and well selling.

Given data from the production line from 60 Oompa-Loompa's, the goal is to predict the production of the **green bars**.

In this presentation I will go over the overview of the solution, its design criteria, and its strengths.

METHODOLOGY

The training pipeline(right), depicts the approach to the solution. At first, an SVM model was applied to the imbalanced data, but seeing that data is hard to marginalise in Figure 1, it gave really poor results. Hence, Random forests (RF) were used since model architecture reduces overfitting. However, to improve results oversampling using SMOTE was performed to balance the data. Then using the best hyperparameters, the models were trained imbalanced and balanced. However, due to better recall, the balanced version was selected as the one to be used on the unknown test data set.

DATASET

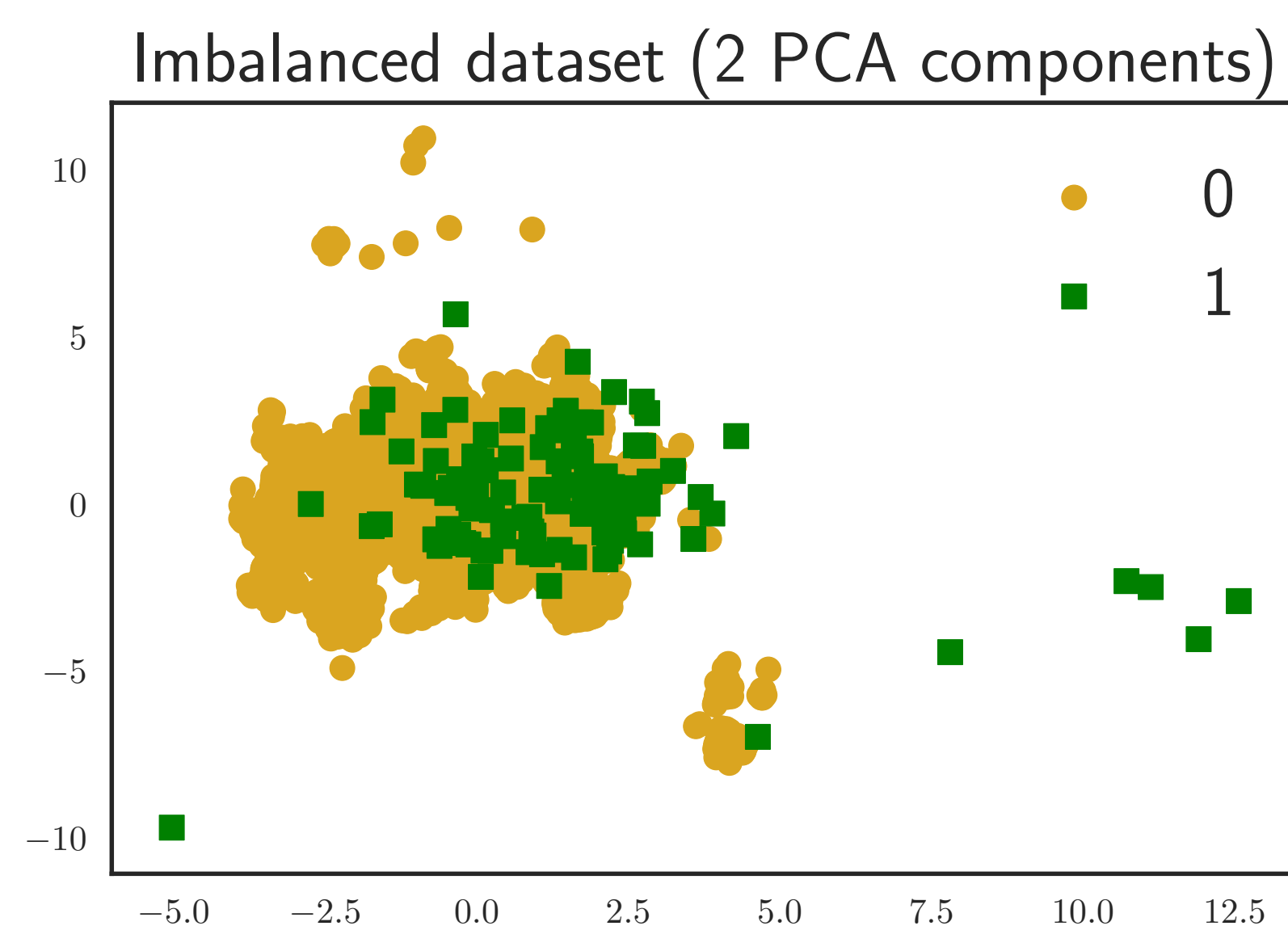


Figure 1: Normalised data with applied PCA dimensionality reduction. The plot shows how extremely difficult it is to distinguish between golden and green bars.

The dataset given consists of 62 columns, of which 60 are scalar features given by each Oompa-

Loompa, 1 is a time stamp, and the last is the target 'GREEN'. This task will consist of binary classification. However, the data is immensely imbalanced as green bar production is very rare.

Golden bars	Green bars
15704	102

Table 1: Gold vs Green bars in dataset

Feature engineering involved of decomposing the time stamp into 5 columns : "day", "dayWeek", "hour", "min", "month". Moreover, Features correlated with each other by a threshold $t = \pm 0.32$ were discarded. Hence, the training and testing was done on 22 features.

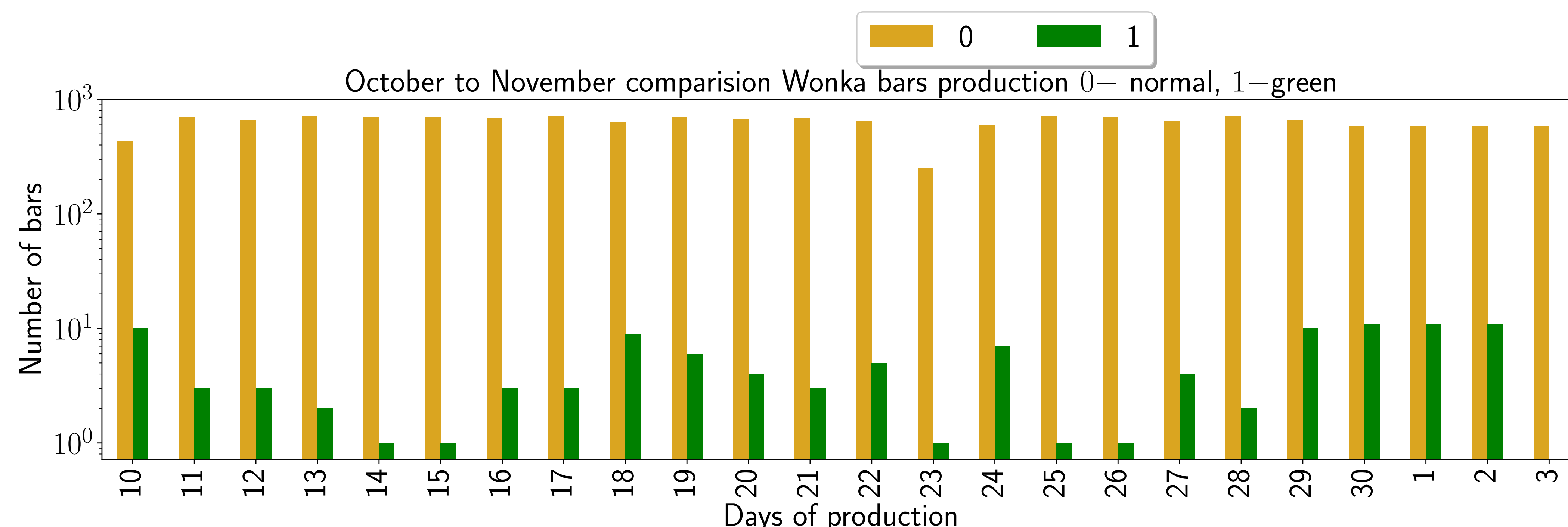
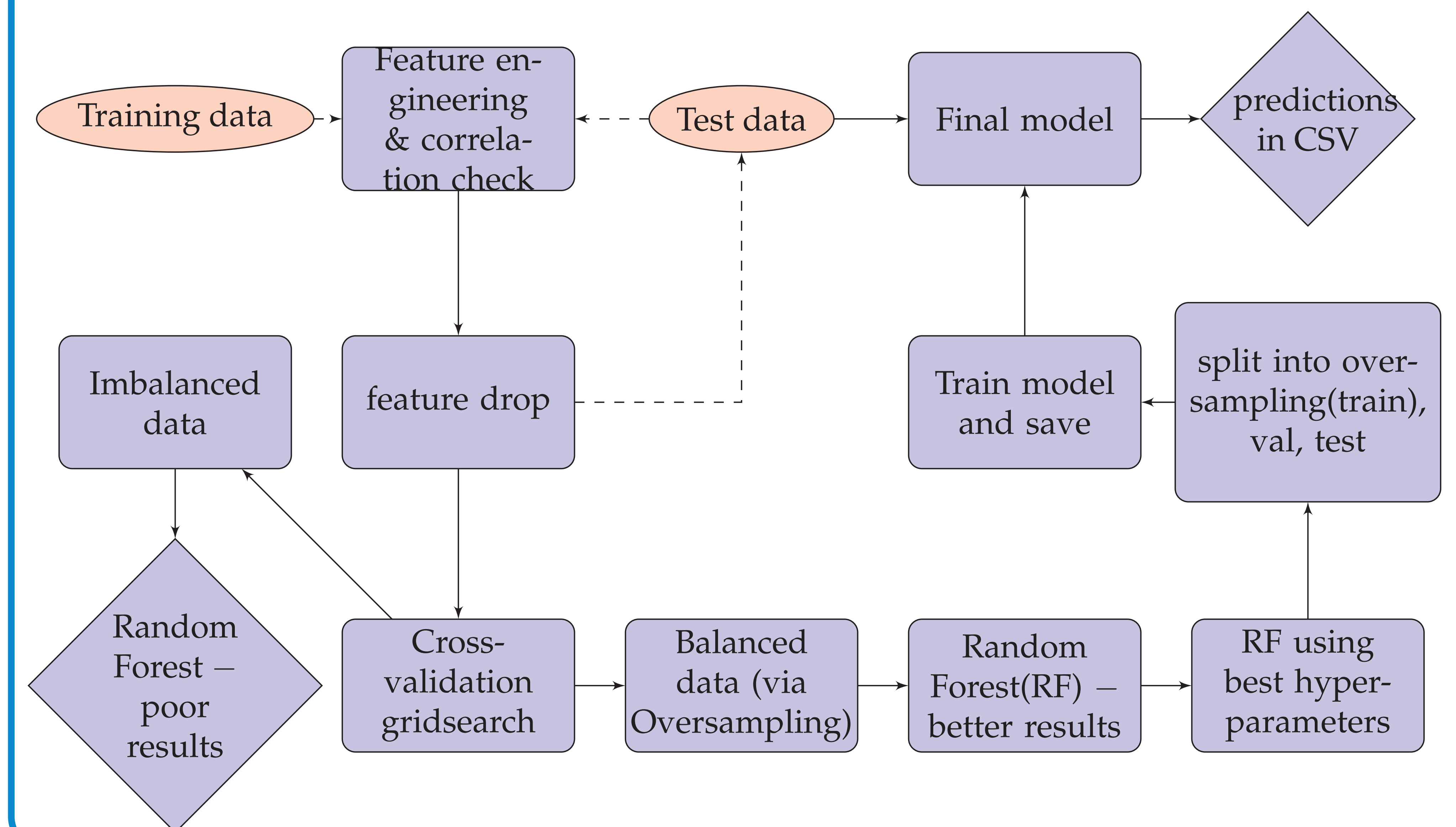


Figure 2: Analysis of production of gold(0) and green(1) chocolate bars through October to early November 2018

TRAINING PIPELINE



RESULTS

First 5-fold cross-validation grid search, maximizing recall was performed to find the best hyperparameters (Best-hp). Despite the imbalanced *Imb* dataset showing better performance on cross-validation. When ran using the Best-hp on the entire dataset split into: train, validation, and test, the balanced *Bal* (oversampled) model performs much better. In fact, recall on the test set is 39% higher on the Best-hp balanced version than the imbalanced.

Model Type	Precision	Recall
CV <i>Bal</i> RF	0.86±0.12	0.39±0.13
CV <i>Imb</i> RF	0.67±0.06	0.49±0.13
Best-hp <i>Bal test</i>	0.67	0.93
Best-hp <i>Imb test</i>	0.82	0.54

Table 2: End results of training on balanced *Bal* and imbalanced *Imb* datasets + final training on the best hyperparameters.

CONCLUSIONS

Random forests are great estimators despite when collinerity is relatively high on an imbalanced dataset. However, reducing that collinerity and oversampling the training set, results in better overall performance. In this task however, we are largely concered with recall than precision. This is because, it is far more costly to wrongly predict a green bar as gold, rather than vice-versa and hence letting a green bar simply slip by. The percentage of green bars in the training dataset is 0.65% meanwhile, the percentage on the given test dataset predicted by the best model is just under 5 percent. Hence to conclude, I reckon that the model should yield good results once tested by the judges, as greens are rare as but the testing sample may slightly different. Another interesting fact was that, training without the use of a feature timestamp yielded awfully poor results with recall $< 10\%$. As Figure 2 potrays, some days just produce more green bars. Please note: logos in the top left and right belong to and were taken from www.loopqprize.ai