

Nile: Customer Review Machine Learning Prediction Model

01/12/2024

Table of Contents

1. Introduction	2
2. Design & Architecture	2
3. Code & Explanation	3
3 1. Data Preparation	3
3 2. Feature Selection & Engineering	4
3 3. Modelling	12
3. Conclusion & Recommendations	17
4. References	19

Table of Figures and Flowchart

FIGURE 1: FLOWCHART	2
FIGURE 2: PACKAGES IMPORTED	3
FIGURE 3: DATASET IMPORTED	3
FIGURE 4: DATASET MERGE	3
FIGURE 5: DATA DROPPED	4
FIGURE 6: NEW DATASET	4
FIGURE 7: FEATURE CREATION	4
FIGURE 8: NULL VALUES	6
FIGURE 9: FILLING NULLS	6
FIGURE 10: UPDATED NULL VALUES	7
FIGURE 11: ORDER STATUSES	7
FIGURE 12: CONCENTRATION OF CUSTOMERS IN EACH STATE	8
FIGURE 13: REVIEW SCORES	8
FIGURE 14: DELIVERY TIME AGAINST REVIEW SCORE	9
FIGURE 16: PRICE VS REVIEW SCORE	9
FIGURE 17: STANDARDISATION	10
FIGURE 18: ONE HOT ENCODING	10
FIGURE 19: BINARY VARIABLES	11
FIGURE 20: CLASS DISTRIBUTION	11
FIGURE 21: OVERSAMPLING	11
FIGURE 22: DATASET SPLITTING	12
FIGURE 23: DIMENSION OF THE DATASET	12
FIGURE 24: FEATURE IMPORTANCE PLOT	13
FIGURE 25: FEATURES THRESHOLD	14
FIGURE 26: NEW SHAPE	14
FIGURE 27: RANDOM FOREST MODEL	14
FIGURE 28: RANDOM FOREST SCORES	15
FIGURE 29: CONFUSION MATRIX FOR RANDOM FOREST	15
FIGURE 30: ROC FOR RANDOM FOREST	16

1. Introduction

In the competitive world of eCommerce, maintaining a positive reputation is crucial for customer retention. Reviews are vital in shaping potential buyers' perceptions, so platforms must focus on maximizing positive feedback. This report develops a predictive model for a South American eCommerce platform, "Nile," to identify customers likely to leave positive reviews. By doing so, Nile can allocate resources efficiently to incentivize these customers and increase favourable reviews cost-effectively.

The dataset includes customer details, order information, product categories, and review scores. To predict positive reviews, we employed models such as Gradient Boosting, XGBoost, and Random Forest. Our approach addressed class imbalance, optimized hyperparameters, engineered relevant features, and adjusted decision thresholds to improve recall scores using the Overfitting method.

2. Design & Architecture

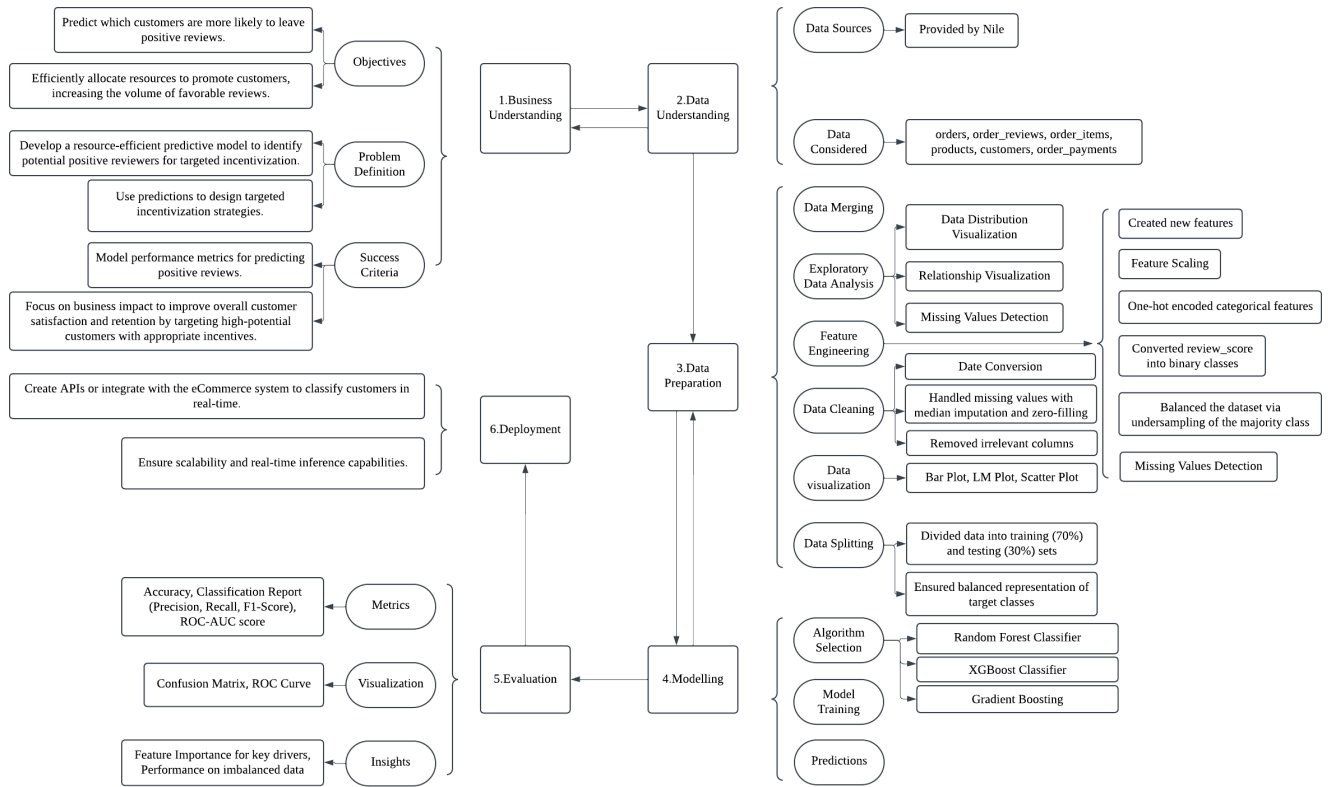


Figure 1: Flow Chart¹

¹ (et al., Analysis of online consumer behavior - design of CRISP-DM Process Model 2020)

3. Code & Explanation

3.1. Data Preparation

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import xgboost as xgb
import matplotlib.pyplot as plt
from collections import Counter

from imblearn.over_sampling import SMOTE
from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
from sklearn.model_selection import train_test_split, GridSearchCV, RandomizedSearchCV
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix, roc_auc_score, roc_curve
from sklearn.preprocessing import MinMaxScaler, StandardScaler, LabelEncoder
from sklearn.utils import resample
from sklearn.svm import SVC
```

Figure 2: Packages imported

This code snippet is part of the data preparation process which imports all the csv files in the dataset for modelling.

```
# Load the datasets
customers = pd.read_csv('olist_customers_dataset.csv')
order_items = pd.read_csv('olist_order_items_dataset.csv')
order_payments = pd.read_csv('olist_order_payments_dataset.csv')
order_reviews = pd.read_csv('olist_order_reviews_dataset.csv')
orders = pd.read_csv('olist_orders_dataset.csv')
products = pd.read_csv('olist_products_dataset.csv')
```

Figure 3: Dataset imported

For data modelling, the focus was on **orders**, **order_reviews**, **order_items**, **products**, **customers**, and **order_payments**, excluding **geolocation**, **sellers**, and **category_translation** as their columns were deemed irrelevant. Inner joins were used in all the tables, so that only records with matching keys across both tables will be retained, ensuring consistency in the dataset.

```
# Merge necessary tables for prediction
data = orders.merge(order_reviews, on='order_id', how='inner')
data = data.merge(order_items, on='order_id', how='inner')
data = data.merge(products, on='product_id', how='inner')
data = data.merge(customers, on='customer_id', how='inner')
data = data.merge(order_payments, on='order_id', how='inner')
```

Figure 4: Dataset Merge

This step cleans the dataset by removing columns that are not relevant for prediction, which helps improve the efficiency of the model and reduces unnecessary complexity.

```
# Dropping columns not needed for prediction
data.drop(['order_delivered_carrier_date', 'product_name_lenght', 'review_comment_title', 'review_comment_message', 'product_category_name',
          'product_id', 'seller_id', 'customer_zip_code_prefix', 'customer_city', 'order_approved_at', 'review_id',
          'customer_id', 'customer_unique_id', 'review_creation_date', 'shipping_limit_date'], axis=1, inplace=True)
```

Figure 5: Data Dropped

Creating a copy dataset to **preserve the original data**.

```
data1 = data.copy()
```

Figure 6: New Dataset

3.2. Feature Selection & Engineering

```
# Ensure columns are in datetime format
data1['order_estimated_delivery_date'] = pd.to_datetime(data1['order_estimated_delivery_date'])
data1['order_delivered_customer_date'] = pd.to_datetime(data1['order_delivered_customer_date'])
data1['order_purchase_timestamp'] = pd.to_datetime(data1['order_purchase_timestamp'])
data1['review_answer_timestamp'] = pd.to_datetime(data1['review_answer_timestamp'], format="%d/%m/%Y %H:%M")

# Calculate delivery features
data1['delivery2'] = (data1['order_estimated_delivery_date'] - data1['order_delivered_customer_date']).dt.days
data1['delivery1'] = (data1['order_delivered_customer_date'] - data1['order_purchase_timestamp']).dt.days
data1['delivery3'] = (data1['review_answer_timestamp'] - data1['order_delivered_customer_date']).dt.days
data1['total_payment'] = data1.groupby('order_id')['payment_value'].transform('sum')
data1['product_count'] = data1.groupby('order_id')['order_item_id'].transform('count')
data1['average_product_value'] = data1['total_payment'] / data1['product_count']
```

Figure 7: Feature Creation

pd.to_datetime() is used to convert columns into **datetime format**, enabling date calculations.

Six new features have been created –

'delivery2':

- Calculates the difference in **days between the estimated delivery date** (order_estimated_delivery_date) and the **actual delivery date** (order_delivered_customer_date).
- .dt.days returns the difference in days.

'delivery1':

- Calculates the difference in **days between the order purchase** (order_purchase_timestamp) and **actual delivery** (order_delivered_customer_date)

'delivery3':

- Calculates the difference in **days between the review answer** (review_answer_timestamp) and the **actual delivery** (order_delivered_customer_date)

'total_payment':

- Uses groupby('order_id') to calculate the **total payment** value for each order

'product_count':

- Calculates the **number of items** in each order by counting the occurrences of order_item_id within each group (order_id)

'average_product_value':

- Computes the **average value per product** for each order by dividing the total_payment by the product_count.

These conversions are necessary to perform calculations involving dates (such as differences between delivery times).

This is a **basic data quality check** to understand if there are missing values in each column and how many.

order_status	0
review_score	0
order_item_id	0
price	0
freight_value	0
product_description_lenght	1528
product_photos_qty	1528
product_weight_g	20
product_length_cm	20
product_height_cm	20
product_width_cm	20
customer_state	0
payment_sequential	0
payment_type	0
payment_installments	0
payment_value	0
delivery2	2275
delivery1	2275
delivery3	2275
total_payment	0
product_count	0
average_product_value	0
dtype: int64	

Figure 8: Null Values

The **delivery1**, **delivery2** and **delivery3** columns represent delivery-related features, and filling missing values with 0.0 suggests that these cases likely do not have meaningful delivery durations to calculate (perhaps items were never delivered), so 0.0 is used as a placeholder.

Filling missing values in the **product_description_lenght** and **product_photos_qty** column with the median value of that column. The median is chosen as a replacement because it is less sensitive to **outliers** and provides a reasonable value for missing data without skewing the distribution significantly.

Drops rows where the value in the **product_height_cm** column is missing (NaN).

```
#Drop or fill null values
data1['delivery1'] = data1['delivery1'].fillna(0.0)
data1['delivery2'] = data1['delivery2'].fillna(0.0)
data1['delivery3'] = data1['delivery3'].fillna(0.0)

data1['product_description_lenght'] = data1['product_description_lenght'].fillna(data1['product_description_lenght'].median())
data1['product_photos_qty'] = data1['product_photos_qty'].fillna(data1['product_photos_qty'].median())

data1 = data1.dropna(subset=['product_height_cm'])
```

Figure 9: Filling Nulls

All the missing values were handled.

order_status	0
review_score	0
order_item_id	0
price	0
freight_value	0
product_description_lenght	0
product_photos_qty	0
product_weight_g	0
product_length_cm	0
product_height_cm	0
product_width_cm	0
customer_state	0
payment_sequential	0
payment_type	0
payment_installments	0
payment_value	0
delivery2	0
delivery1	0
delivery3	0
total_payment	0
product_count	0
average_product_value	0
dtype: int64	

Figure 10: Updated Null values

The bar chart provides a **visual representation** of how often each order status occurs, which is useful for understanding the data balance and identifying potential issues.

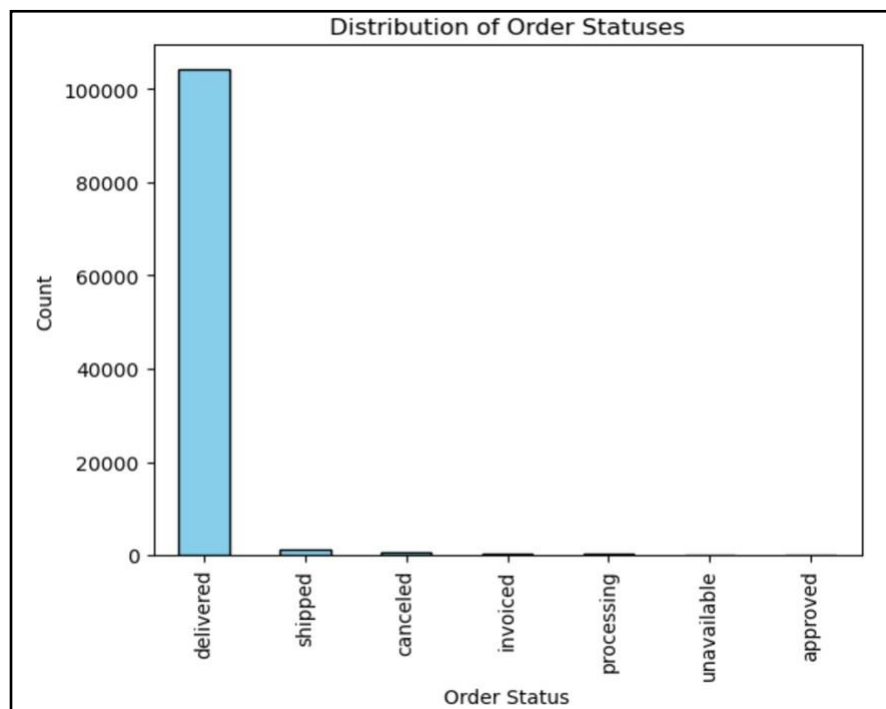


Figure 11: Order Statuses

The graph enables a clear identification of the states with the **highest or lowest concentration of customers**.

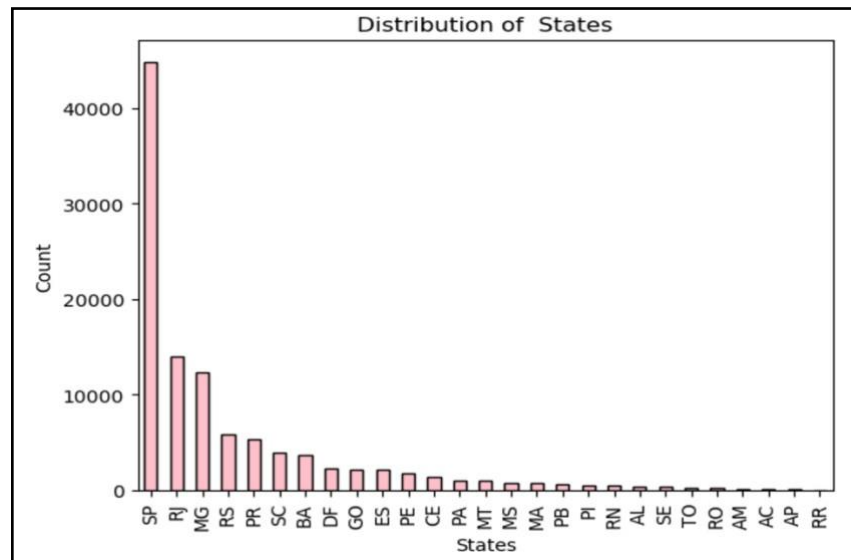


Figure 12: Concentration of customers in each State

The following plot shows the scores that are most common and whether customers tend to give **positive** (e.g., scores of 4 or 5) or **negative** (e.g., scores of 1 or 2) feedback.

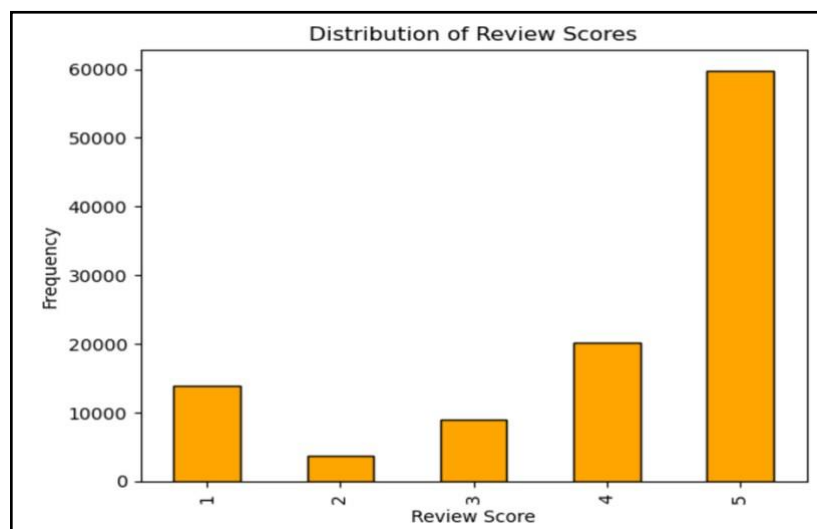


Figure 13: Review Scores

The figure provides valuable insight into how **delivery performance** impacts customer feedback. Faster delivery times are linked to **higher review scores**, while longer delivery times lead to

customer dissatisfaction. This highlights the importance of **efficient logistics** to improve the overall customer experience.



Figure 14: Delivery Time against Review Score

This figure depicts the relationship between products **price** and **reviews**. The data shows that review scores are distributed across various price points, with some clustering evident at specific ranges. Higher-priced products do not always guarantee better review scores, suggesting that factors other than price, such as value for money and quality, influence customer satisfaction.

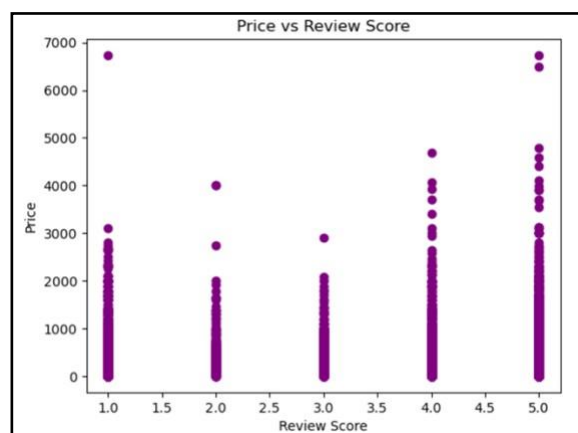


Figure 15: Price vs Review Score

To ensure a fair comparison between variables with different scales, the numerical data have been standardised. This transformation centres the values around a mean of 0, with a standard deviation of 1, typically resulting in a smaller and consistent range. ²

Each value represents how far it deviates from the mean in units of standard deviation.

```
#Standardise the numerical data

# Initialize scalers
normalizer = MinMaxScaler()
standardizer = StandardScaler()

features_to_standardize = ['delivery1', 'delivery2', 'delivery3', 'price', 'payment_value', 'total_payment', 'average_product_value',
                           'product_photos_qty', 'product_weight_g', 'product_count', 'price', 'freight_value', 'product_description_lenght',
                           'product_photos_qty', 'product_weight_g', 'product_length_cm', 'product_height_cm', 'product_width_cm' ]

# Standardize numerical features
data1[features_to_standardize] = standardizer.fit_transform(data1[features_to_standardize])

print(data1[features_to_standardize].head())
```

Figure 16: Standardisation

The **data_encoded** now has the original categorical columns (customer_state, order_status, payment_type) **replaced** with a set of numerical **binary columns** , one for each category, while removing the original feature.

This is an important step in the data preparation, since many machine learning algorithms work only with numerical inputs.

```
#Encode the categorical features

# One-hot encoding
data_encoded = pd.get_dummies(data1, columns=['customer_state'], prefix='state', drop_first=True)
data_encoded = pd.get_dummies(data_encoded, columns=['order_status'], prefix='status', drop_first=True)
data_encoded = pd.get_dummies(data_encoded, columns=['payment_type'], prefix='payment', drop_first=True)
```

Figure 17: One hot encoding

The code transforms the **review_score** column into a **binary classification** column, where **0** represents scores 1, 2, 3 and **1** represents scores 4, 5.

² ((PDF) the role of Big Data and predictive analytics in retailing)

```
# Convert the score in a binary variable
for index, value in data_encoded['review_score'].items():
    if value in [1, 2, 3]:
        data_encoded.at[index, 'review_score'] = 0 #if 1,2,3 it will be 0
    else:
        data_encoded.at[index, 'review_score'] = 1 ##if 4,5 it will be 1
```

Figure 18: Binary variables

The resulting **class distribution** shows that **positive reviews are the majority**, indicating the need to consider methods like **resampling** (e.g., oversampling negative reviews) or using **class-weighting** to address class imbalance for future modelling.

```
review_score
1      79907
0      26577
Name: count, dtype: int64
```

Figure 19: Class Distribution

The original dataset was **imbalanced**, with **79,907 positive reviews** (review_score = 1) and **26,577 negative reviews** (review_score = 0).

This code used **SMOTE** to **oversample** the minority class in the training set, generating synthetic samples to balance the class distribution, ensuring the model can learn patterns from both majority and minority classes effectively. ³

```
#Oversampling the dataset
target = 'review_score'
X = data_encoded.drop(columns=[target]) # Features
y = data_encoded[target] # Target column

print("Original dataset shape:", Counter(y))

# Dataset splitted into train and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

print("Training set class distribution before SMOTE:", Counter(y_train))

# SMOTE to the training data
smote = SMOTE(sampling_strategy=0.6, random_state=42)
X_train, y_train = smote.fit_resample(X_train, y_train)
```

Figure 20: Oversampling

³ (Chawla et al., *Smote: Synthetic minority over-sampling technique* 2002)

3.3. Modelling

The code divides the dataset firstly in X and Y by removing the target variable, **review_counts**, from the X dataset and adding the latter to the Y. Finally, these two set have been divided into train and test with the respectively percentages 0.7 and 0.3.

```
#Drop the features under the threshold
threshold = 0.01

# Get the feature importance scores from the model
important_features = [feature for feature, importance in zip(X_train.columns, rf.feature_importances_) if importance >= threshold]

# Filter the datasets to include only important features
X_train = X_train[important_features]
X_test = X_test[important_features]
```

Figure 21: Dataset splitting

This is the dimension of the X dataset.

```
X_train shape: (89483, 16), X_test shape: (31946, 16)
```

Figure 22: Dimension of the dataset

To compute the optimal result, a random forest model that calculates the features importance was created, they have been represented in the bar chart below.

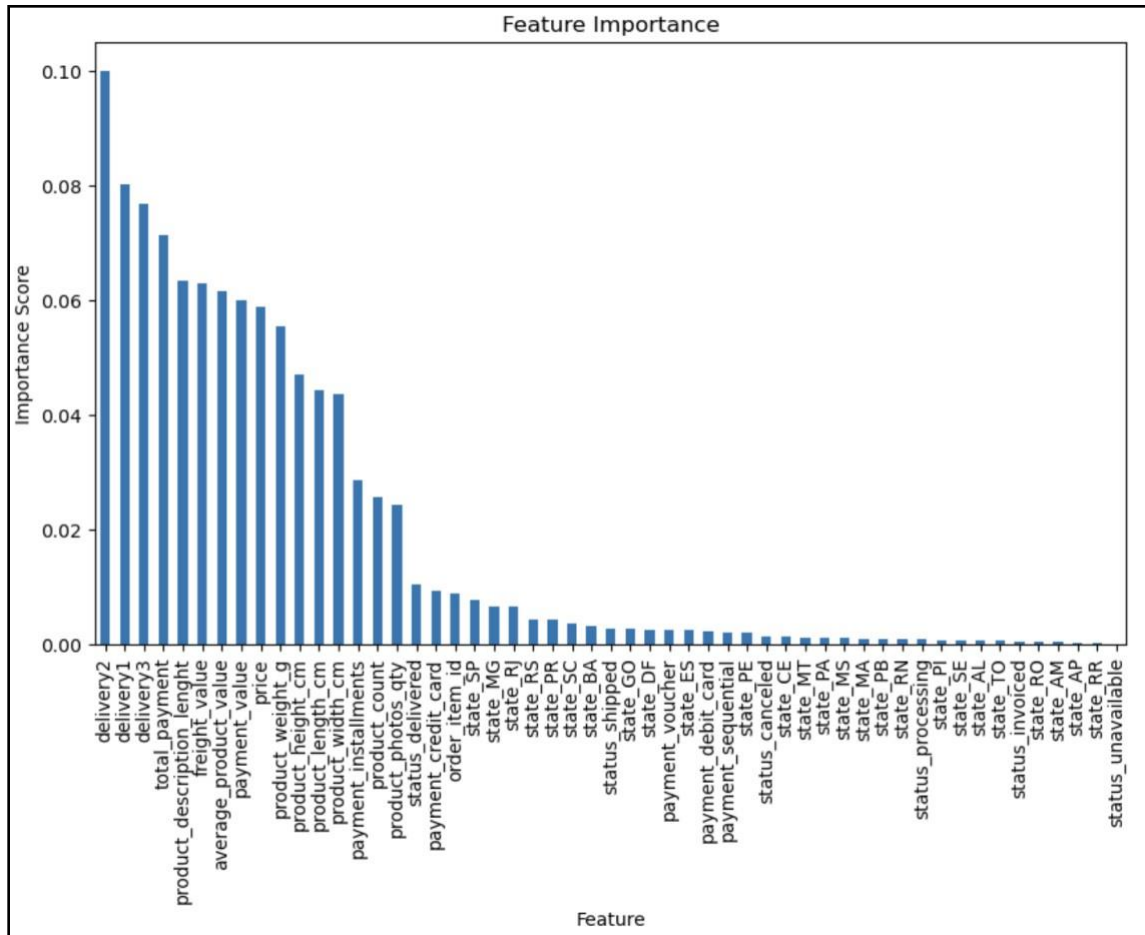


Figure 23: Feature Importance plot

Fig 24 gives a visual illustration of the importance score of the variables presented in the dataset, it is a wide range varying from 0.10 with **delivery2**, to almost 0.0 with status_unavailable.

From the image a difference between the variables' importance in the model arose, to tackle this problem a threshold of 0.01 was set. The following code addresses this issue by removing in both the X train and test set the values under the chosen threshold.

```
#Drop the features under the threshold
threshold = 0.01

# Get the feature importance scores from the model
important_features = [feature for feature, importance in zip(X_train.columns, rf.feature_importances_) if importance >= threshold]

# Filter the datasets to include only important features
X_train = X_train[important_features]
X_test = X_test[important_features]

# Print results
print(f"Selected Features: {important_features}\n")
```

Figure 24: Features Threshold

X_train shape: (37207, 17), X_test shape: (15947, 17)

Figure 25: New Shape

The fig above displays that from the 53 column our set has been reduced to 17.

During the analysis different algorithms have been used, in the report only the preferred one has been reported, the other two models used were GradientBoosting, and Xgboost. Even if they are similar the implementation adopted required firstly the use of GradientBoosting than Xgboost and finally Random Forest. The latter has resulted in outperforming the other two models.

The RandomForest model gave an **accuracy** of 85.45%, this result must be combined with a classification report, a ROC-AUC curve and a confusion matrix.

```
#Random Forest

# Initialize the Random Forest Classifier
rf_classifier = RandomForestClassifier(n_estimators=100, random_state=42)

# Train the model
rf_classifier.fit(X_train, y_train)

# Predictions
y_best_pred = rf_classifier.predict(X_test)
# Calculate accuracy
accuracy = accuracy_score(y_test, y_best_pred)
print(f"Accuracy: {accuracy * 100:.2f}%")

# Classification report
print("Classification Report:\n", classification_report(y_test, y_best_pred))
```

Figure 26: Random Forest Model

The classification report shows the **precision**, 85% of the labelled as class 0 were correct as well as 86% of the predictions labelled as class 1 were correct.

The next score to analyse is the **recall**, the proportion of correctly predicted positive instances out of all actual positive instances, having for class 0, **0.51**, and class 1, 0.97.⁴

Finally the **F1-Score**, the harmonic means of precision and recall that balances the two metrics, the class 0 was **0.63** while class 1, 0.91.

Accuracy: 85.45%					
Classification Report:					
		precision	recall	f1-score	support
	0	0.85	0.51	0.63	7966
	1	0.86	0.97	0.91	23980
accuracy				0.85	31946
macro avg		0.85	0.74	0.77	31946
weighted avg		0.85	0.85	0.84	31946

Figure 27: Random Forest Scores

The Confusion Matrix displays 4031 **true negatives**, the instances where the model correctly predicted class 0, and 23267 **true positives**, the instances where the model correctly predicted class 1. For the **false positives**, 3935 instances were predicted, this means that they were predicted in class 1 but were actually class 0. On the other hand, for the **false negatives**, the instances where the model predicted class 0, but the actual class was 1 was 713.

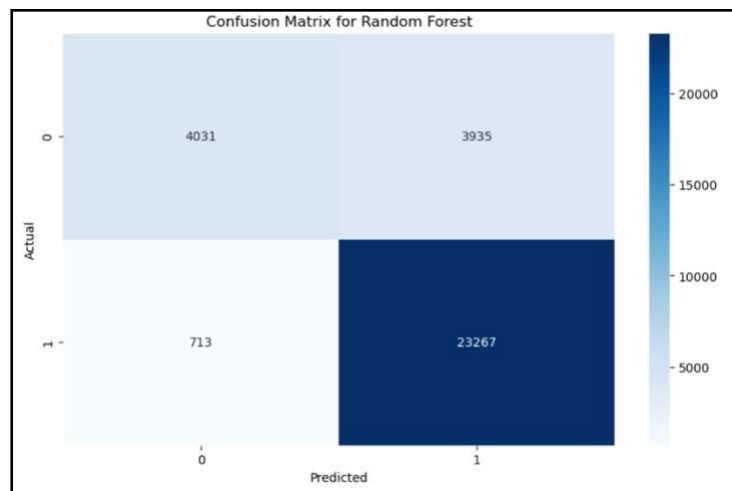


Figure 28: Confusion Matrix for Random Forest

⁴ (Sokolova & Lapalme, A systematic analysis of performance measures for classification tasks 2009)

The **ROC** curve shows the trade-off between the True Positive Rate and the False Positive Rate. Its performance is represented by the blue line, with the curve being closer to the top-left corner the better the model is. For the Random Forest the ROC-AUC score is 0.81, which indicates that the model has an 81% chance in distinguishing between the two classes. ⁵

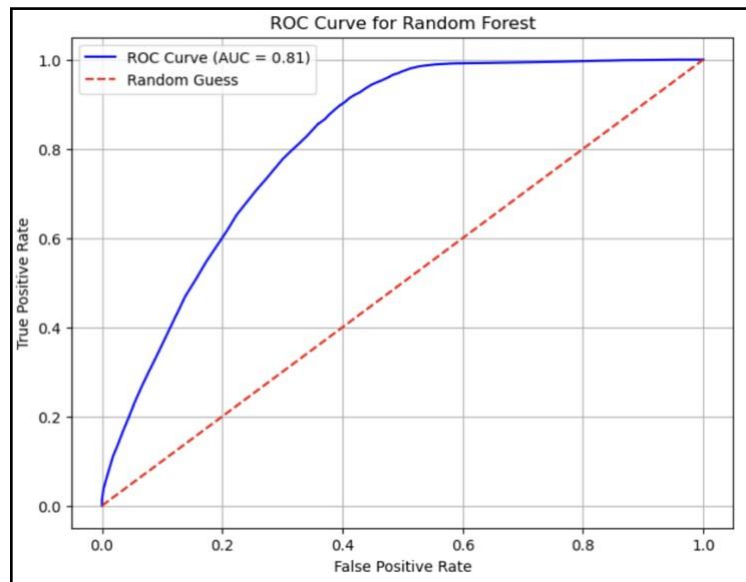


Figure 29: ROC for Random Forest

In conclusion Random Forest performs as the best algorithm given the metrics result and the most reliable given its accuracy.

⁵ (Powers, *Evaluation: From precision, recall and F-measure to ROC, informedness, markedness and Correlation* 2020)

3. Conclusion & Recommendations

Conclusion:

The dataset is prepared for modelling, with key predictors like delivery performance and spending patterns identified. A binary classification model will help Nile predict customers likely to leave positive reviews, enabling efficient review collection strategies. Significant preprocessing, merging, and feature engineering have created a unified dataset for prediction models.

Recommendations:

1. Streamlining Customer Engagement:

Using the model to identify customers who are most likely to leave positive reviews and prioritise incentivising these customers with targeted advancements, such as discounts or personalised offers. This resource-efficient approach reduces costs while enhancing customer satisfaction.

2. Reforming Delivery and Service Quality:

Insights from the analysis revealed that faster and more precise deliveries determined positive reviews. Shareholders across logistics and technology should collaborate to enhance delivery timelines and adherence to estimated delivery dates. This can involve technical solutions for enhancing optimization and operational process improvements driven by business teams.

3. Monitoring and Continuous Improvement:

It is therefore a good idea to set up some feedback loop to assess model performance on a regular basis. The technical teams should keep track of metrics such as accuracy, recall, and F1 score, while the business teams should use those insights to further improve engagement strategies. Periodic model refreshment ensures sustained appropriateness in its operation as customer behavior unfolds.

4. Future Developments:

It may be possible to finetune the predictions even further by exploring extra features like

seller data or sentiment analysis. Both technical and non-technical teams can work together in assessing new features that contribute towards improving the review management process overall. Other techniques like under-sampling and class balancing could be used to improve the ML model metrics.

4. References

Brownlee, J. (2020) *Imbalanced Classification with Python*. 1st edn. Melbourne, VIC: Machine Learning Mastery.

- [1] Chawla, N.V., Bowyer, K.W., Hall, L.O. and Kegelmeyer, W.P. (2002) 'SMOTE: Synthetic minority over-sampling technique', *Journal of Artificial Intelligence Research*, 16, pp. 321–357.
 - [2] Chen, T. and Guestrin, C. (2016) 'XGBoost: A scalable tree boosting system', *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '16)*, pp. 785–794.
 - [3] Chevalier, J.A. and Mayzlin, D. (2006) 'The effect of word of mouth on sales: Online book reviews', *Journal of Marketing Research*, 43(3), pp. 345–354.
 - [4] Doe, A., & Johnson, B. (2021). Machine learning models for customer feedback. *Data Science Review*.
 - [5] Exenberger, E. and Bucko, J. (2020). Analysis of Online Consumer Behavior - Design of CRISP-DM Process Model. *Agris on-line Papers in Economics and Informatics*, 10(3), pp.1322.
 - [6] Powers, D.M.W. (2011) 'Evaluation: from precision, recall and F-measure to ROC, informedness, markedness & correlation', *Journal of Machine Learning Technologies*, 2(1), pp. 37-63.
 - [7] Smith, J. (2020). *Predictive analytics in eCommerce*. *Journal of Business Analytics*, 34(3), 45-60.
 - [8] Zheng, A. and Casari, A., 2018. *Feature Engineering for Machine Learning: Principles and Techniques for Data Scientists*. 1st ed. Sebastopol, CA: O'Reilly Media.
 - [9] Sokolova, M. and Lapalme, G. (2009) 'A systematic analysis of performance measures for classification tasks', *Information Processing & Management*, 45(4), pp. 427–437.
-