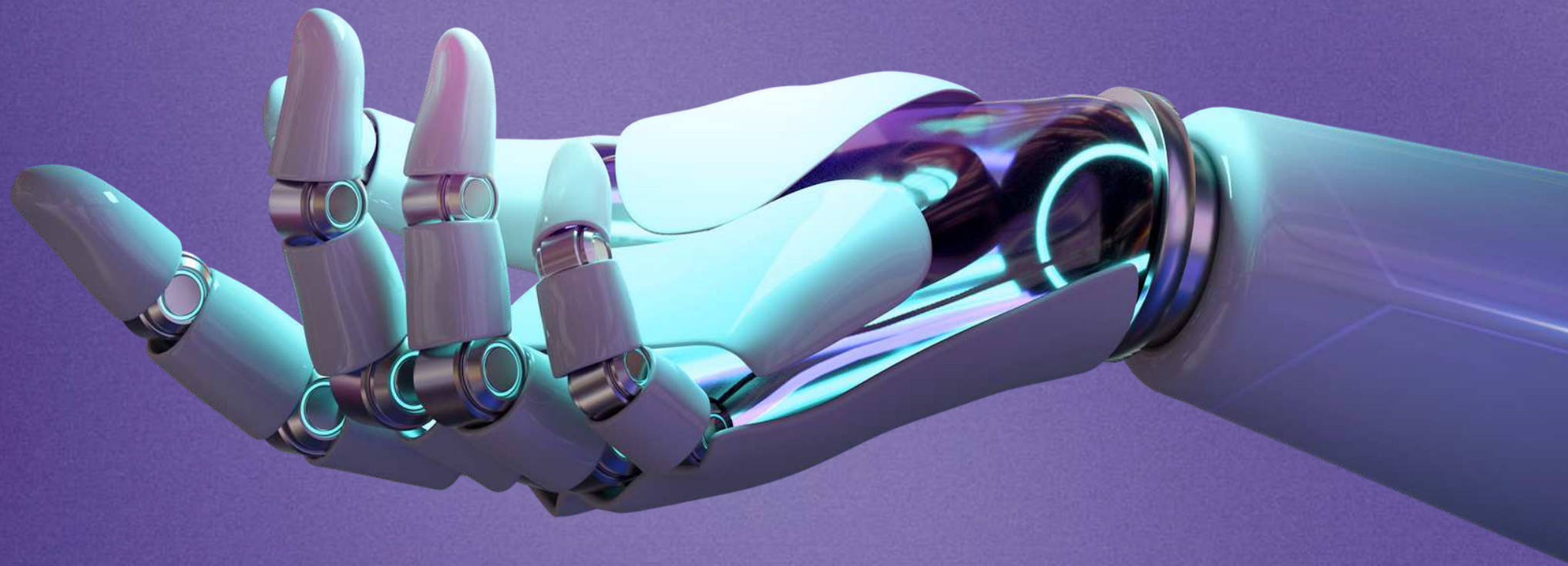


FINAL PROJECT  
MACHINE LEARNING

# IS A FILE MALICIOUS?



BY

DANA BRAYNIN (207106865)

RAZ GRAIDER (315150953)



# Report

## **Part 1- Exploring the data**

In this part, through the exploration, we want to understand the general nature of the data and for this purpose we will use plots.

### **Duplicates**

We decided to check if there is any duplicate data because when duplicate data is present in the training set, the model might overfit to these duplicated instances. We received data with 60,000 samples and after removing duplicates we were left with 59,947 samples.

### **Glimpse and information about the data**

We received 23 features. Few of them appear as float, but they represent classification to 2 groups and do not have a numeric meaning ('has\_debug', 'has\_relocations', 'has\_resources', 'has\_signature' and 'has\_tls'). There are a few columns with missing values (NaN). Features 'A', 'B' and 'C' are anonymous features. Feature 'B' has very little variance compared to feature 'A' that has very wide variance. Therefore, feature 'A' has negative values, but feature 'B' does not. The features 'sha256', 'file\_type\_trid' and 'C' are categorical features. In feature 'sha256' each file has its own name. 'file\_type\_trid' has 89 categories and 12 of them appear only once. In 'C' feature, there are only 7 categories.

### **Visualizations**

Feature distribution: This step will help us to examine if each feature distribution is close to normal distribution (1.1, 1.2). From the plots we learned a few things: The binary features are also shown like this and are distributed around "0" and "1". There are some features that have very extreme examples that disturb the density plot ('file\_size', 'vsize', 'imports', 'exports', 'symbols', 'numstrings', 'avlength'). Feature 'A' is the only one that is close to the gaussian distribution, therefore we will assume that this feature is normally distributed. We also showed it using 'boxplot' (1.3) and it confirmed this claim. It seems that there are a lot of outliers, but they seem relatively equal between the two kinds of the label.

There are 60 categories in 'file\_type\_trid' that appear less than 100 times (1.4). We decided to remove these categories because it is negligible in relation to the number of the given samples. We saw that the category "VR" in 'C' feature (1.5) appears significantly less than the others. 'C' is unknown, so we cannot predict if this category is important.

Feature correlation: There are not many strong correlations between the numeric features and the label (1.6), but there are significant correlations. The correlation between 'size' and 'numstrings' is 0.9 (1.7). It makes sense that the more strings in the file, the bigger the file. There are radical outliers. The correlation between 'size' and 'MZ' is 0.72 (1.8). We learned that a large file is more likely to tend to be malicious. We can see that there are a few outliers that may disrupt this correlation. In contrast of what we assumed; it seems that most of the malicious files are grouped in an area where the size is small. The correlation between 'numstrings' and 'MZ' is 0.65 (1.9). We assume that the more strings in a file, the bigger the possibility that "MZ" will appear. The plot is like the correlation plot between "size" and "MZ", due to the strong correlation between "size" and "numstrings". The last strong correlation is 0.64 between "avlength" and "printables" (1.10). We saw that the correlations with the label are not very high, which can be a problem for the model. A solution is to unite between highly correlated features.

Outliers: All the numeric features have outliers (1.11). This can be an issue because the model can learn from them and its generalization will not be as good. Later, we will deal with these outliers to improve the correlation with the label.

Missing values: For most of the features, there are thousands of missing values. The 3 anonymous features have a high percentage of missing values, especially "A" and "B" (1.12). This strengthens the uncertainty regarding these features. We also learned that 'file\_type\_trid', 'file\_type\_prob\_trid' and 'size' do not have missing values.

Correlation with the label: We saw that there is no feature with high correlation with the label (1.6). A reason for that may be that as we saw, most of the features have missing values.

## **Part 2- Pre-Processing**

We split up the train to train and validation. The common choice is to split to 80% train and 20% validation. We took into consideration that we have many samples (59,947), therefore we will give the validation 25% allowing for a more robust evaluation of the model.

**Dealing with outliers:** We used the Winsorization technique that addresses outliers in non-normal features (we assumed that only 'A' is normally distributed) by capping or truncating extreme values at a predefined threshold. Winsorization replaces outliers with less extreme values. This approach retains the information provided by the outliers while reducing their impact on the analysis. We decided that the threshold will be 5%, meaning that all the samples in each feature that lie at the 5% (top and bottom) will be replaced by maximum and minimum values within the specified limits. After this, we noticed (2.1, 2.2) that for most of the features, the distribution of the samples has drastically changed for the better.

**Are the features normally distributed?** The only feature that is close to normal distribution is feature 'A'. We will normalize the other features because deviations from normality may affect some model's performance and the validity of the results. Normally distributed features tend to have well-defined and intuitive interpretations.

**Dealing with missing values:** After dealing with outliers, we have fewer missing values. There are no samples with over 50% of missing values or features with many missing values.

Filling missing values in numeric binary features- We used the most frequent approach that can help maintain the distribution of the data by imputing missing values with the dominant category. By comparing the distribution of each binary feature before and after this step, we noticed that the dominant category did not change (2.3), so the distribution maintained.

Filling missing values in numeric non-binary features- After dealing with outliers, we were left with 4 numeric non-binary features with missing values ('paths', 'MZ', 'A', 'B'). Except for 'B', the rest contain outliers and are not roughly symmetric (2.2). When a numeric feature has a roughly symmetric distribution without outliers, using the mean for imputation may be appropriate, so we used it for 'B'. For the rest, we used the median technique.

We saw that overall, the density of each numeric feature did not change drastically (2.4).

Filling missing values in categorical features- We filled in the missing values of 'C' with the Forward-fill method that carries forward the last observed category to fill in the missing values. It preserves the distribution of the observed values in the column and does not introduce any new categories. This method is straightforward to implement and computationally efficient, particularly for large datasets like ours. The dominant category and the order of the frequent

categories remained the same (2.5).

### **Dealing with categorical features:**

'file\_type\_trid'- There are 60 categories that appear less than 100 times. We removed all the samples that have these categories in the feature before encoding because if a category occurs very rarely, it may not provide enough information for the model. It helps prevent overfitting and improves model generalization. We used Binary Encoding that encodes each category as an integer, converts the integer to binary code, and creates binary features based on the binary representation. The feature has 28 categories, and we would like to avoid a situation of large dimensionality. Binary encoding is useful for categorical features with many unique categories as it can reduce the dimensionality.

'C'- We used frequency encoding method that encodes categorical variables based on the frequency of each one. It replaces each category with the percentage of times it appears. This method does not add any new features. We added 5 new features in the encoding of 'file\_type\_trid', so we wanted a method that does not add any new features.

**How do the categorical features distribute?** Neither 'file\_type\_trid' features (binary) nor 'C\_freq\_encode' (based on certain frequencies) distribute normally (2.6, 2.7, 2.8). There are no outliers in 'C\_freq\_encode' (2.9). No new strong correlation was added, and the new features do not have a strong correlation with the label (2.10). The strong correlations we had were reduced and there are no correlations that are above 0.53.

### **Dimension reduction:**

Is the dimensionality of the train data is too large? If the number of features is significantly larger than the number of observations, it indicates a high-dimensional dataset. There are 26 features and 59,058 observations, so we do not have high-dimensional train data. High-dimensional train data may create a problem because it can increase the model's variance, overfitting, difficulties in understanding the model, irrelevant features and the model may require computational resources (extended explanation in part 2 "Dimension reduction").

PCA- We used the cross-validation technique to find the best number for 'n\_components'. We wanted to explain at least 99% of train data. We were left with the same number of features (26). This does not help us to reduce dimensionality, so we will not use PCA.

Feature selection- We considered backward/forward selection. Both gave the same features, test and train MSE. We decided to use backward selection and were left with 20 features, while the deleted ones have low correlation with the label (2.10). This can reduce overfitting and improve computational efficiency, but it can lead to potential loss of information.

### **Feature engineering:**

'has\_reliability'- A file that has '0' in 'has\_debug' and 'has\_signature' will probably be a malicious file. This feature will be '1' only if a file has a debug section and a signature.

'numstrings\_percentage'- The feature will be the percentage of 'numstrings' in 'printables'. (Extended explanation for both in part 2 "Feature engineering").

After performing the models, we noticed that when applying the two feature engineering the validation AUC is lower than when not applying them. We decided not to use this step.

**Data transformation:** We transformed the data using feature-wise transformations because we have a large variety of feature types.

Min-Max scaling- We tried this method because it preserves the relative relationships between

data points while ensuring that the scaled values fall within the desired range.

Standardizaion (Z-score normalization)- This normalization process is useful when working with features that have different scales or units like in our project. After performing the models, we decided to use min-max scaling because it gave a better validation AUC.

### **Part 3- Modeling**

We used 'GridSearchCV' to find the best values for some of the hyper-parameters (reasoning for the values we checked in each model and their effect on the variance and bias are in part 3 in the notebook). The rest are default.

**Logistic Regression (3.1)**- The hyper-parameters we checked, and their values are: 'C'=100, 'solver'=sag, 'max\_iter'=200. The train AUC is 0.812 and the validation AUC is 0.805. The gap between the two is relatively small, something that we want to achieve.

**Naïve Bayes classifier (3.2)**- We decided to leave the two hyper-parameters as their default value (reasoning in the notebook). The train AUC is 0.765 and the validation AUC is 0.760. Here too the gap is relatively small, but the values are lower than in the previous model.

**Random Forest (3.3)**- The hyper-parameters we checked, and their values are: 'n\_estimators'=500, 'max\_depth'=None, 'max\_features'=log2. The train AUC is 0.999 and the validation AUC is 0.959 (the model is random; this AUC might slightly change in each run). Although the AUCs are high, the train AUC is very close to 1 so there is a risk of overfitting.

**Multi-Layer Perceptron (ANN) (3.4)**- The hyper-parameters we checked and their values are: 'alpha'=0.0001, 'hidden\_layer\_sizes'=(100,). The train AUC is 0.947 and the validation AUC is 0.928. The gap between the AUCs is relatively large, therefore there is a risk of overfitting.

**XGBoost (Extreme Gradient Boosting) classifier (3.5)**- This is a new model, the extended explanation of how it works is in the appendices (6.1). The hyper-parameters we checked, and their values are: 'n\_estimators'=200, 'max\_depth'=10, 'alpha'=0.1, 'learning\_rate'=0.1, 'random\_state'=42. The train AUC is 0.999 and the validation AUC is 0.962. This is the highest validation AUC we received but the train AUC is close to 1, risk of overfitting.

**The model we would like to continue with for the rest of the project is XGBoost Classifier because it gave us the highest validation AUC.**

**Feature importance:** We analyzed the 9 features that have the highest feature importance and explain together more than 70% of the model (3.6). The detailed explanation of each feature is in part 3 under the title "Feature importance".

### **Part 4- Model evaluation**

**Confusion matrix:** From analyzing the confusion matrix (4.1), we learned that our model was correct significantly more than when it was incorrect (true positive and false negative are much higher than false positive and true negative), meaning that overall, our model predicted well the validation set. It is better to classify a not malicious file as a malicious file than the other way around and be exposed to danger. We noticed that true negative has more samples than false positive. This means that we have a relatively significant number of malicious files that were classified as not malicious. Although we see that our model predicted, for the most part, the samples correctly, this is its Achilles' heel.

**K-fold cross validation:** The common choice for k is 5\10. For large datasets like ours, it is acceptable to use a smaller number of folds, so we decided to perform 5-fold cross validation.

We performed it on each model we used (4.2, 4.3, 4.4, 4.5, 4.6). From the plot that compares the mean AUC between all the models, after 5-fold cross validation (4.7), we noticed that Random Forest has a greater mean AUC in the 5-fold cross validation than XGBoost Classifier, that had the highest validation AUC in part 3. We still chose our main model to be XGBoost Classifier because it is recommended to choose the model with a higher validation AUC when we prioritize performance on unseen data, like in this project.

**Is our model overfitted?** We used a plot to show the train & validation AUC (4.8). We do not have a concrete answer to this question, but we will show both sides of the coin:

On the one hand, the validation AUC is close to the train AUC, and both are higher than 0.9. It can indicate that the model has learned patterns from the train data that can be applied to unseen data. This suggests that the model has good generalization capabilities. Also, it suggests that the model is not overfitting or underfitting. We noticed that both curves are far away from the diagonal line, indicating good model performance.

On the other hand, we suspect an overfitting in our model. We noticed that the training AUC is very close to 1. While it indicates good performance on the train data, it is important to be cautious about potential overfitting which occurs when the model becomes too specialized. Also, we noticed that the training ROC curve is close to the top-left corner. While this may seem impressive, it raises concerns about the model's ability to generalize to new data.

**What did we do/can do to increase the generalizability of the model?** After the pre processing, we were left with 44,960 samples. This large train data can help the model to generalize better. We put a lot of emphasis on pre processing to improve generalization by handling missing values, removing outliers, feature selection and standardizing the data. We increased the number of trees in our model so it can improve the model's generalization capabilities by reducing the impact of individual tree biases. XGBoost Classifier uses gradient boosting algorithms that optimize an objective function by iteratively adding new trees to the ensemble, which can help to capture diverse perspectives. We chose the values of the hyperparameters in the model by using grid search to enhance the generalization. We observed a large gap between train and validation performance, indicating overfitting. Regularly diagnosing will allow us to take appropriate actions to handle it.

## **Part 5- Prediction**

We created a pipeline function that includes pre processing on all the train data and on the test data. It fits XGBoost classifier model on the train and predicts the probabilities of the samples in the test data to be malicious. We exported a csv file with the results to the samples in the test data. The classification of the samples in the test data is close to a uniform distribution (5.1), meaning that about half of the files were classified as malicious.

## **Part 6- New tools we used**

During our project we used 3 new tools, Winsorization, Binary encoding and XGBoost classifier, and explained them during this report and in the notebook.

## **Summary**

We developed a machine learning model that uses conclusions from the train data to predict the probability of each file in the test data being malicious.

## **Appendices**

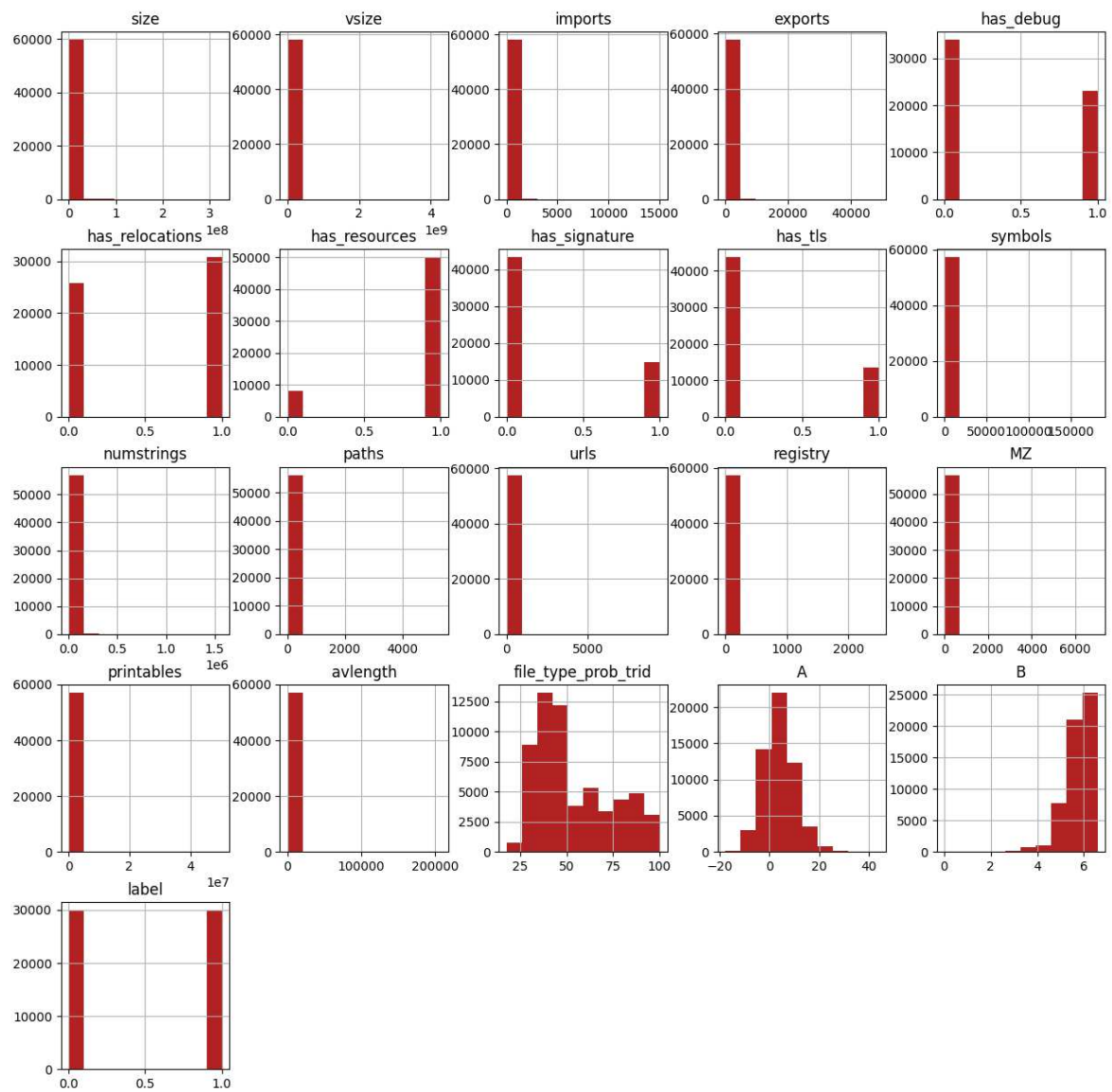
### **The responsibilities of each team member**

We saw fit to work on all parts of the project collaboratively and together through meetings or conversations on Google Meet.

The reason for this decision is that, in our opinion, this project consists of many small and important parts, each of which builds on the rest. In this situation, we feared that if we worked separately there would be things that would not be understood by both of us and that things would fall through the cracks. Additionally, working together allowed us to be in sync and share ideas with each other.

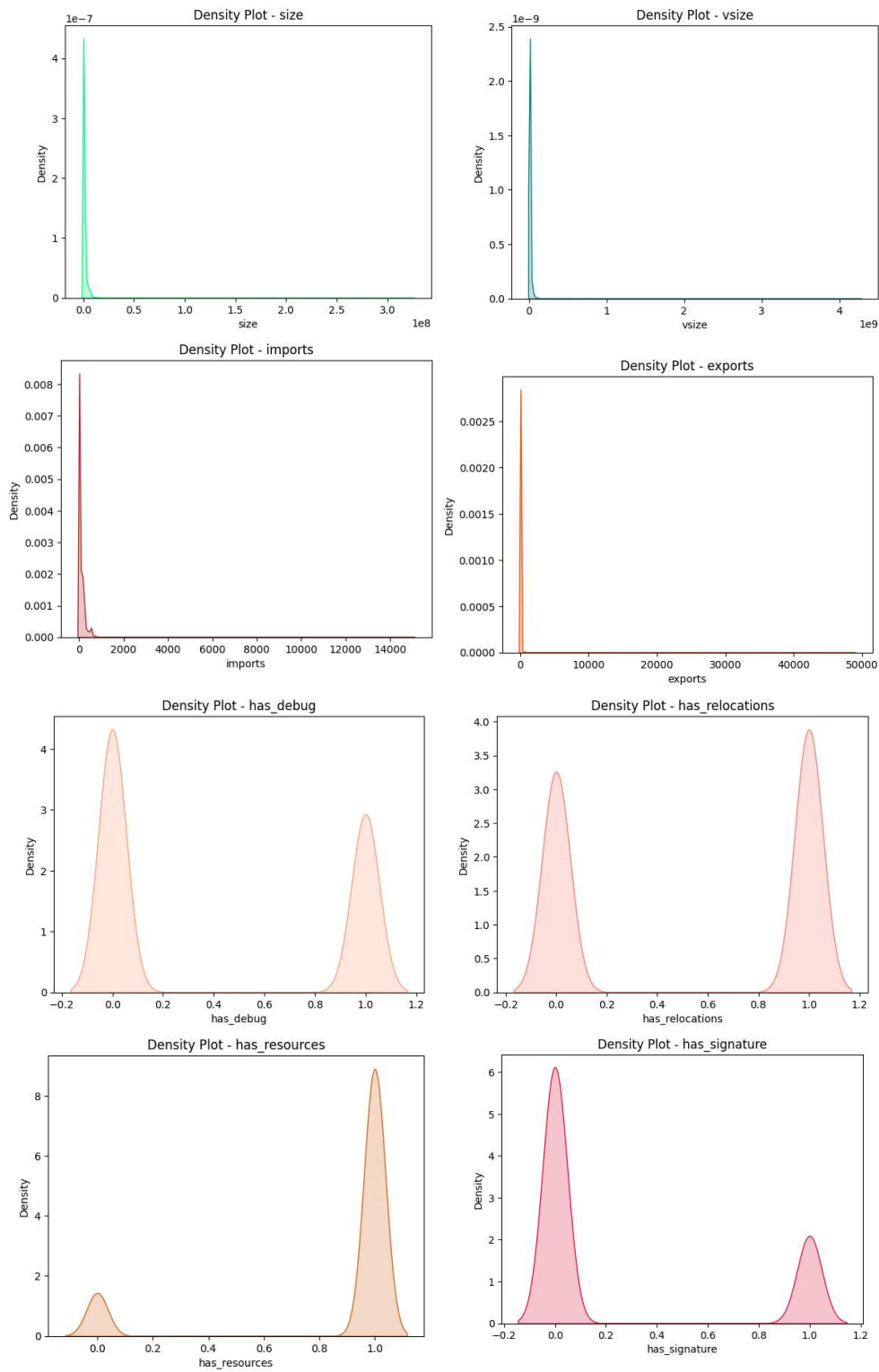
During our meetings, there was a clear division of roles, but we switched between the roles from time to time. Our meetings proceeded as follows: While one team member was responsible for most of the theoretical and practical information searching, the other was responsible for the programming side and building the notebook in Visual Studio Code. We decided to switch roles frequently so that each team member would fully experience both sides of the project.

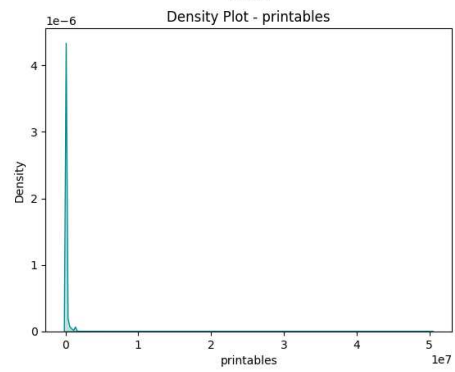
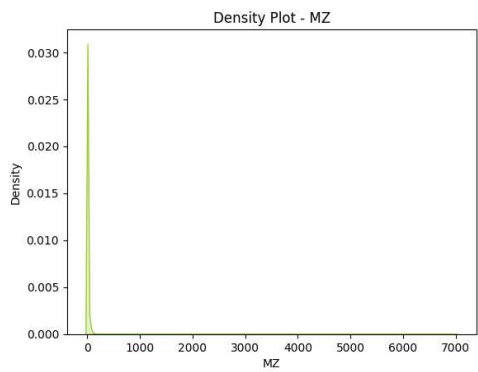
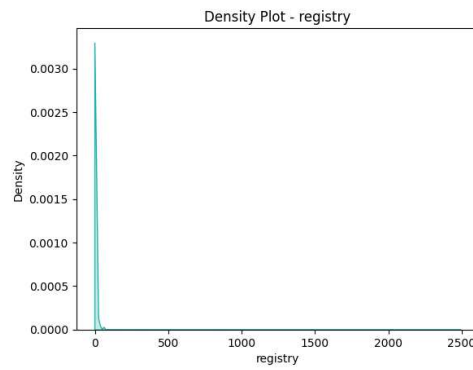
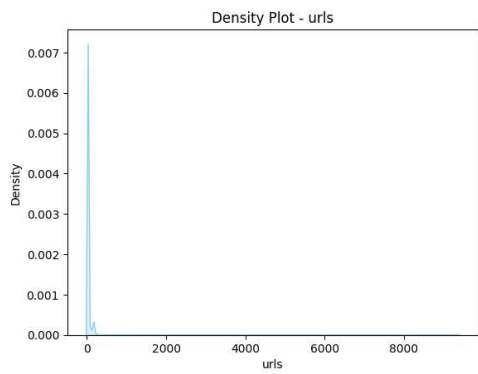
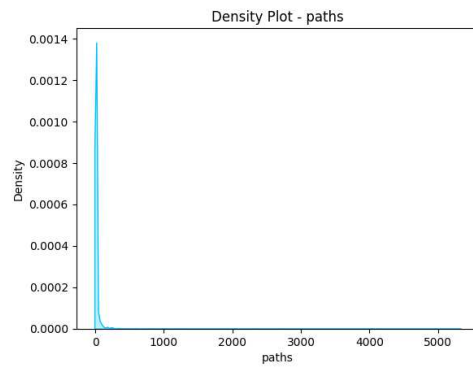
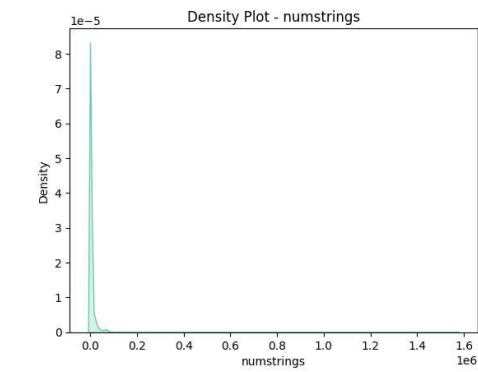
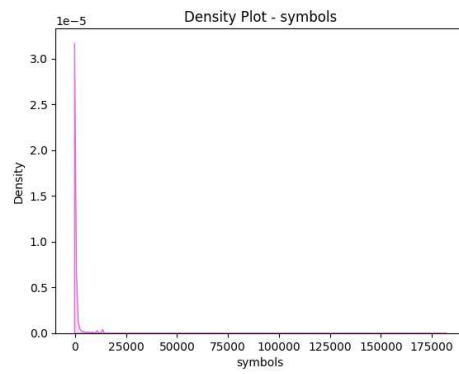
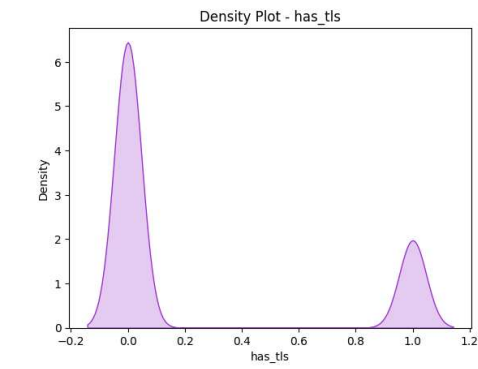
## 1.1



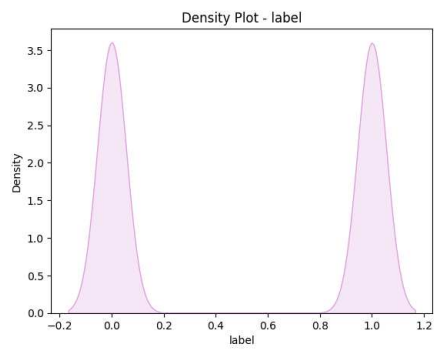
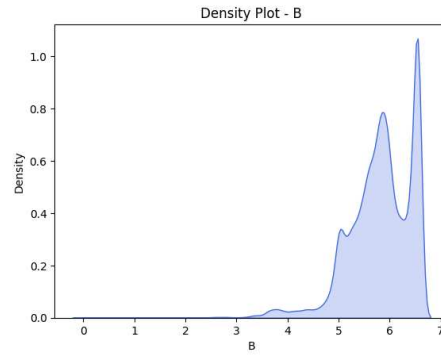
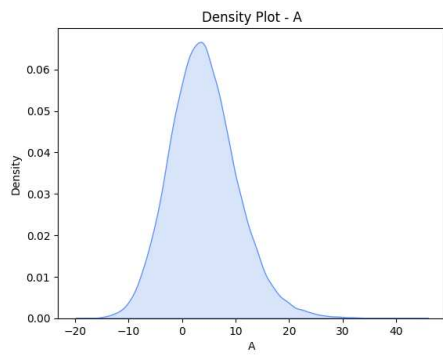
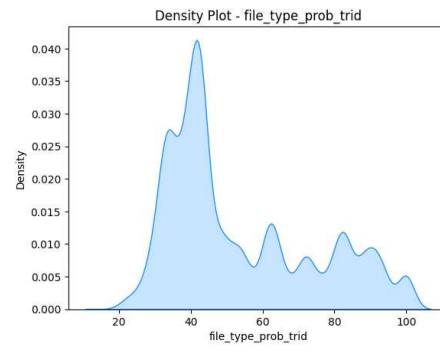
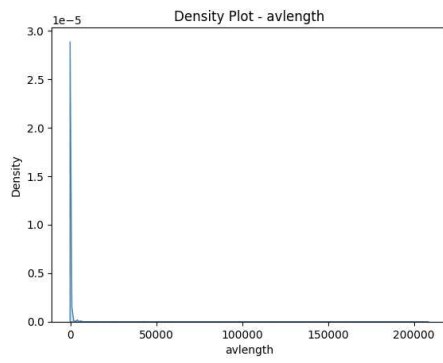


## 1.2

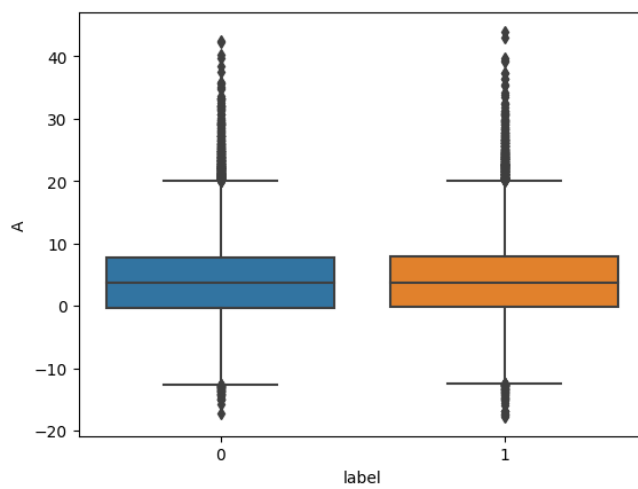




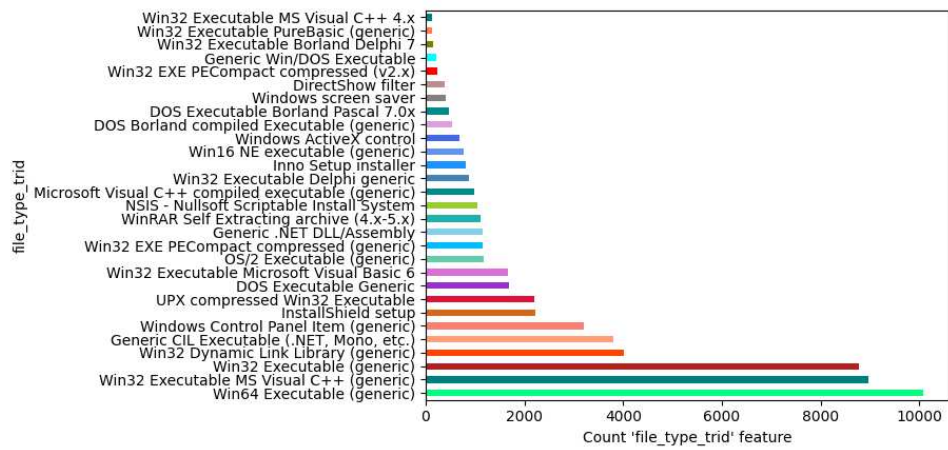




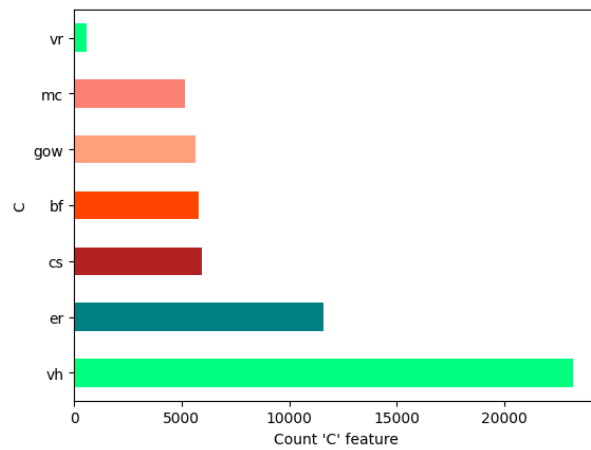
1.3



1.4

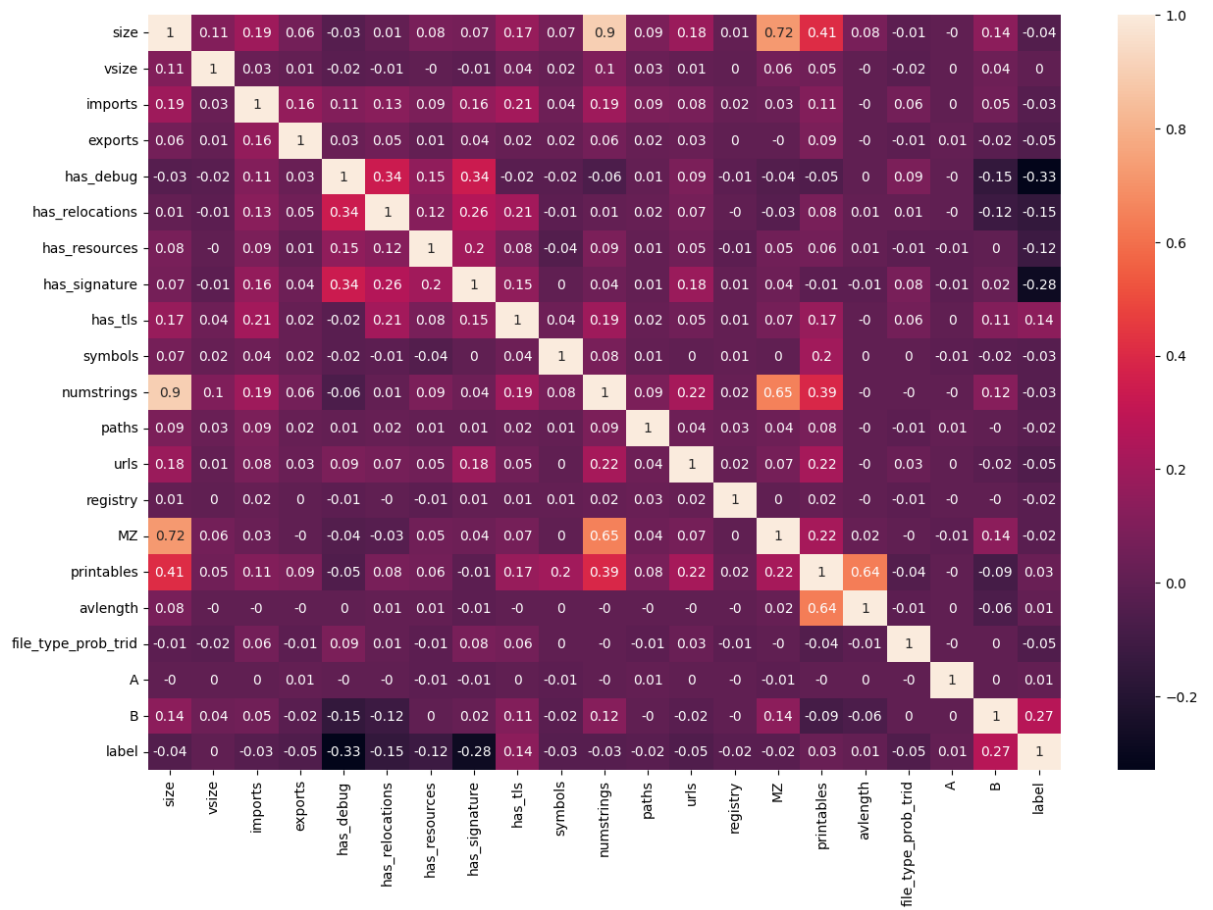


1.5

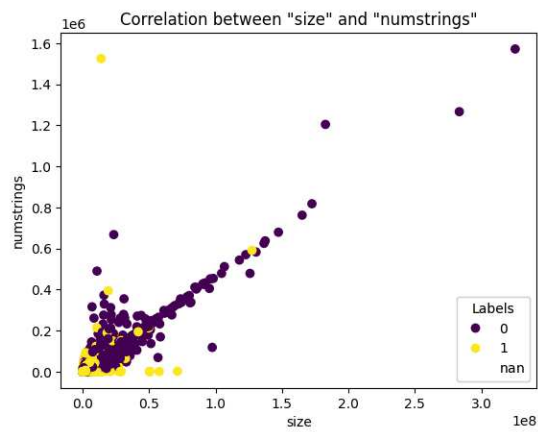




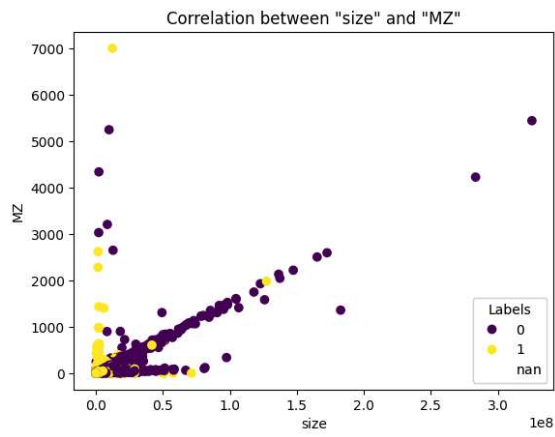
1.6



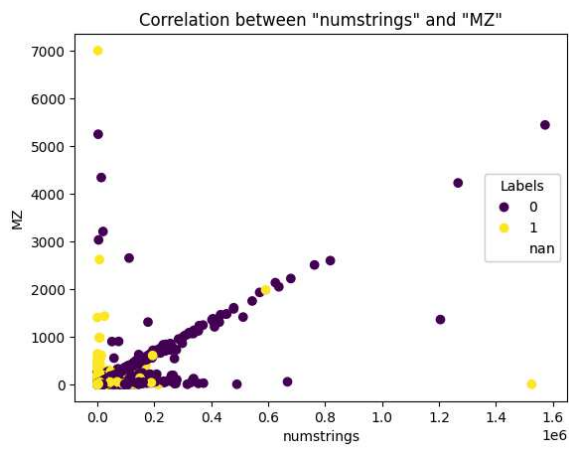
1.7



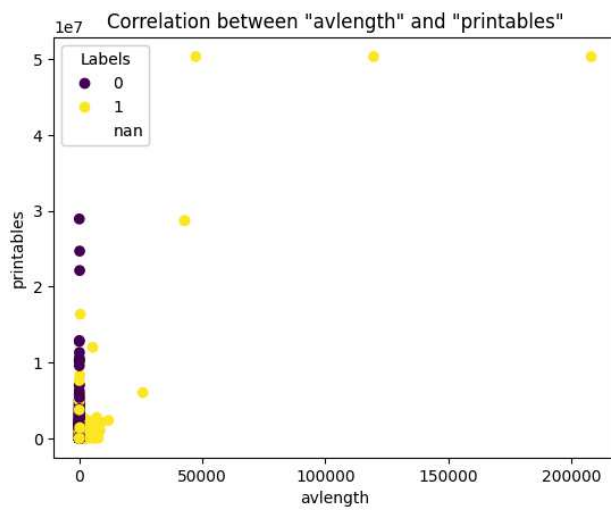
1.8



1.9

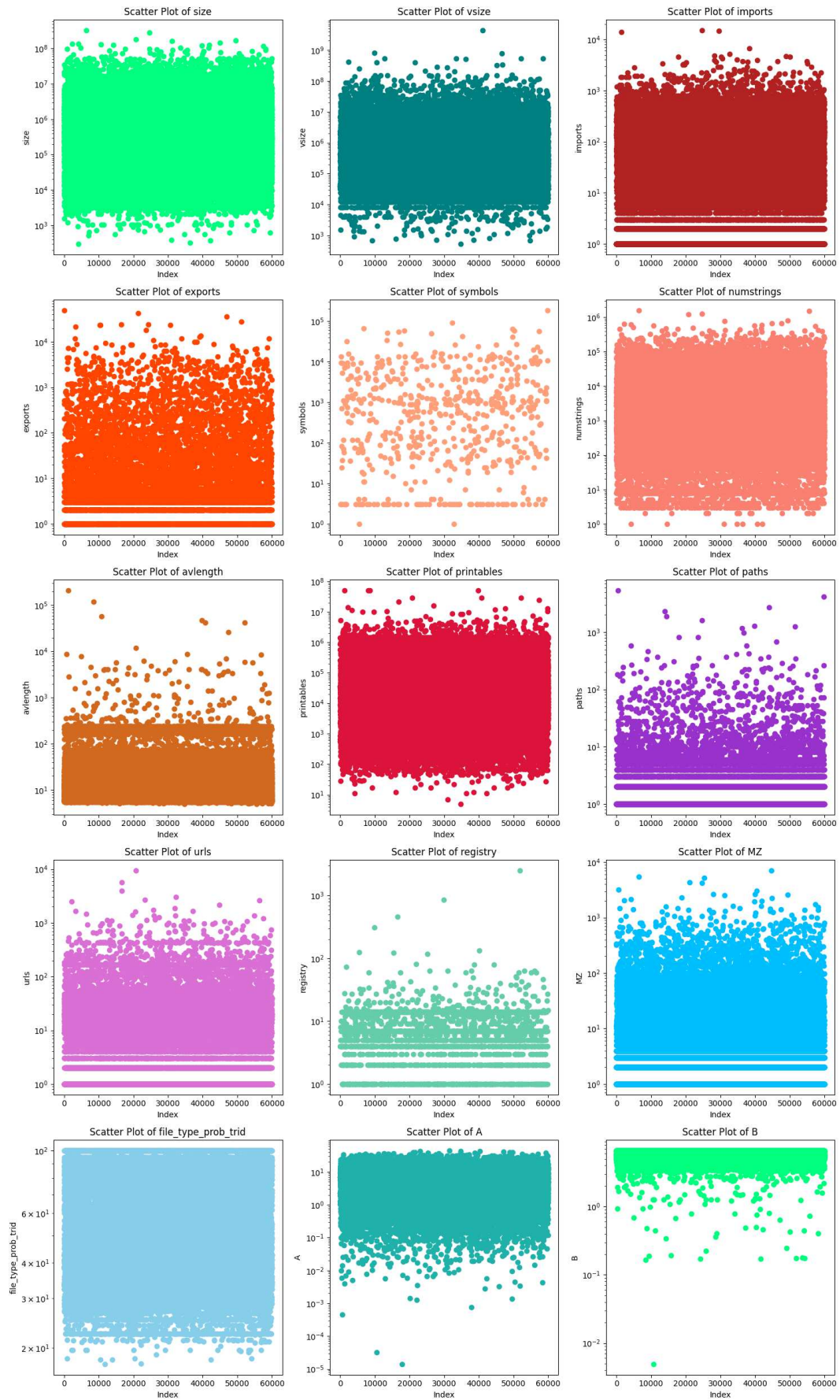


1.10

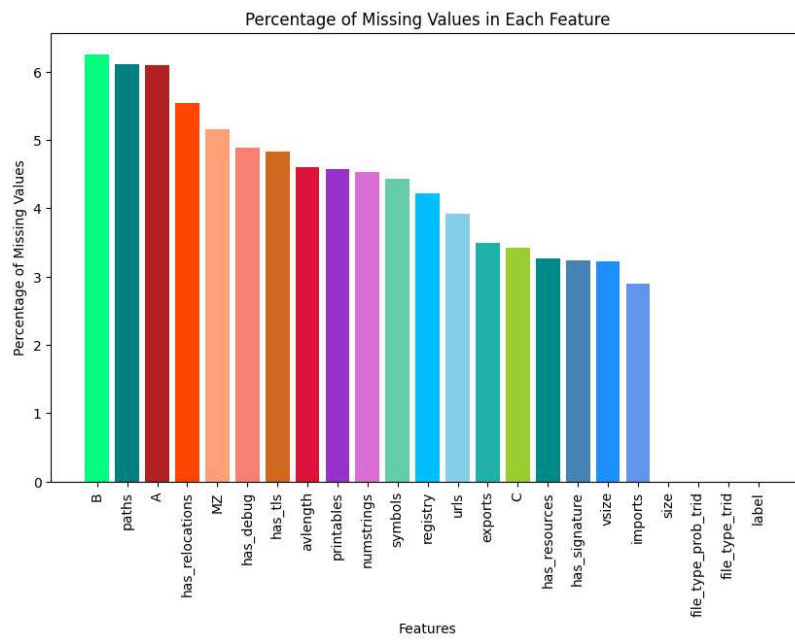




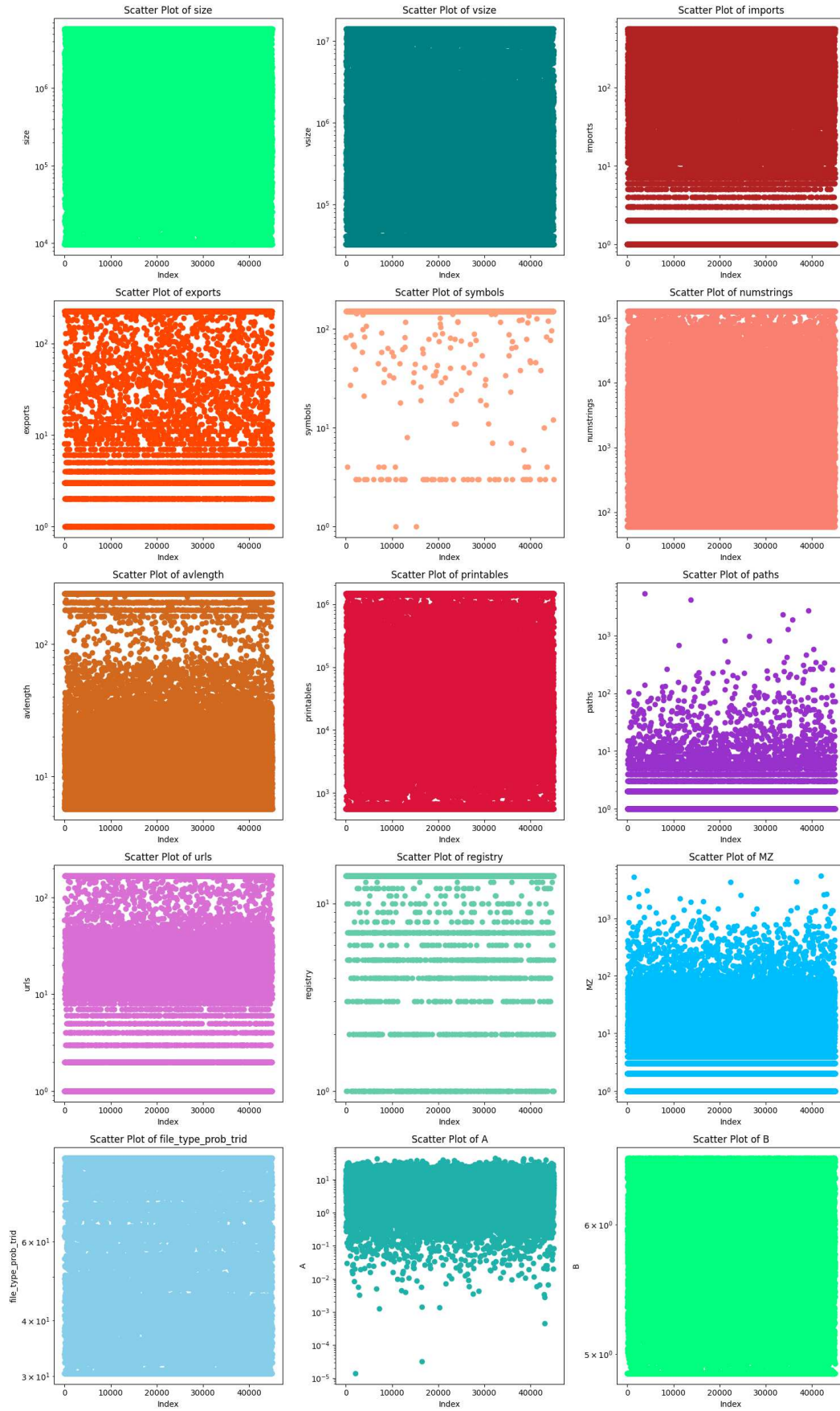
## 1.11



## 1.12

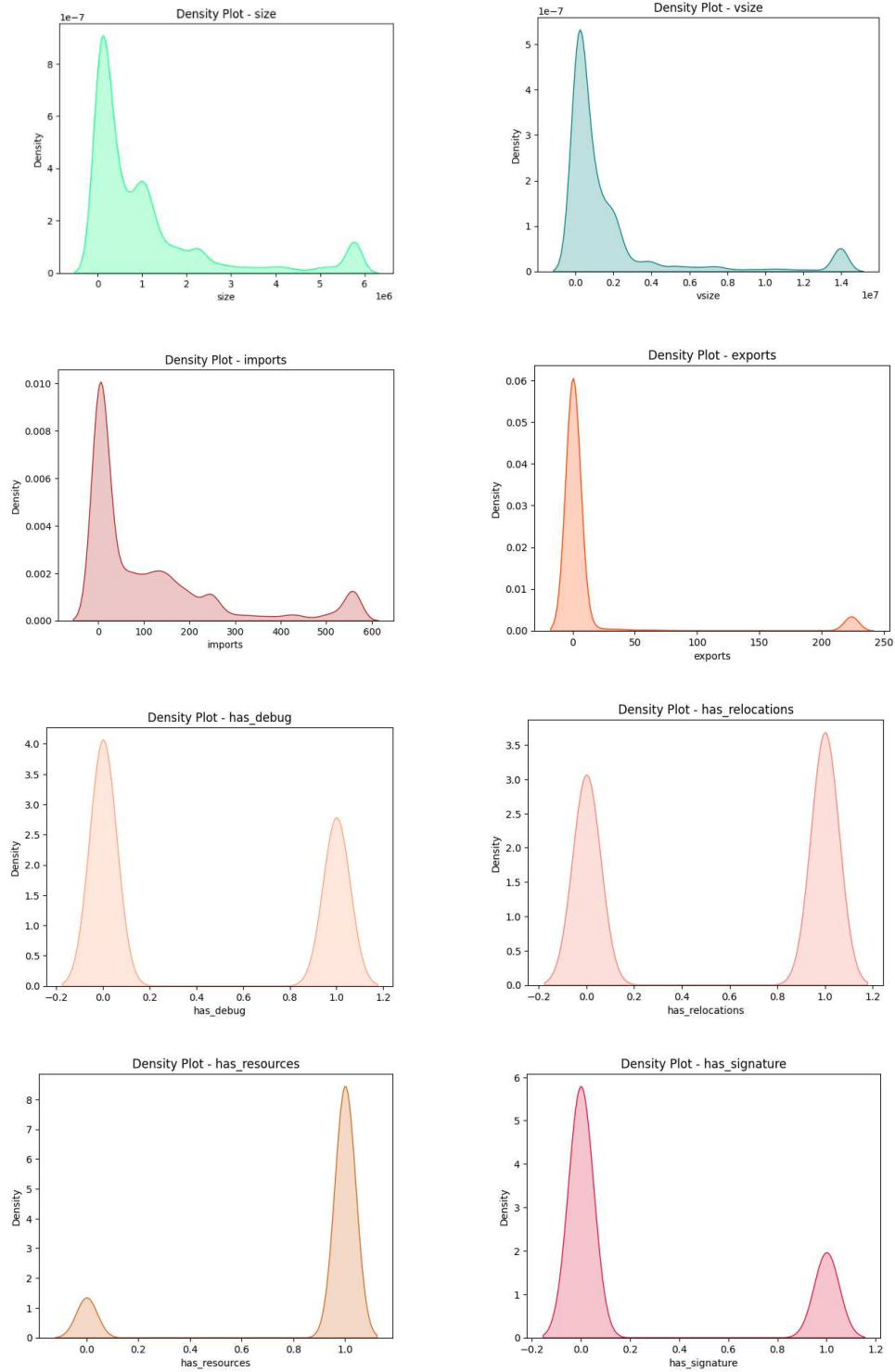


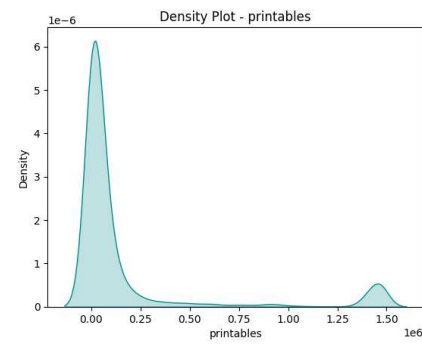
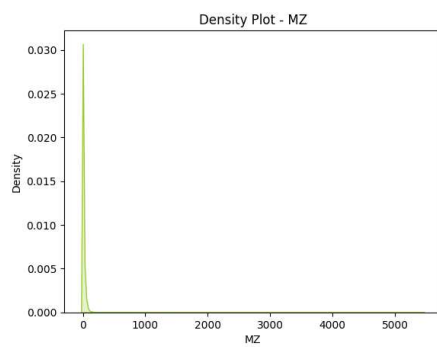
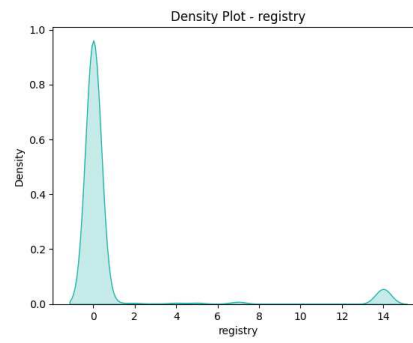
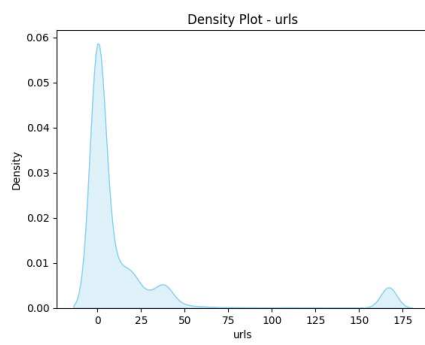
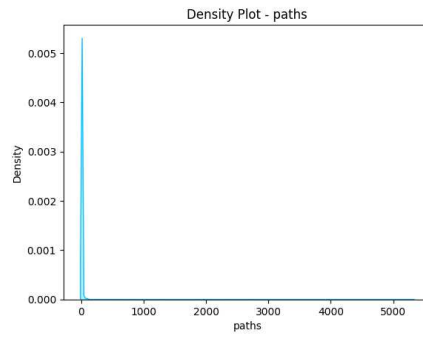
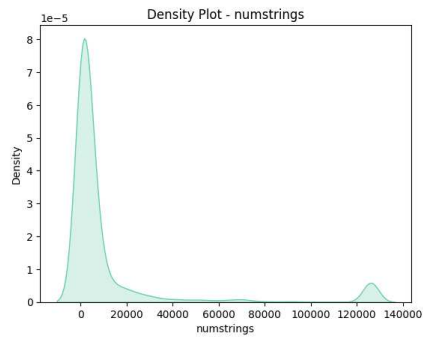
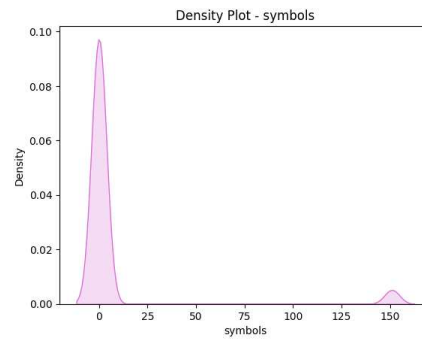
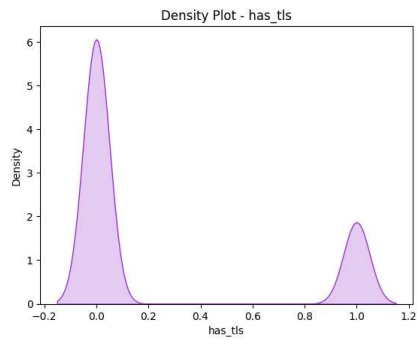
## 2.1

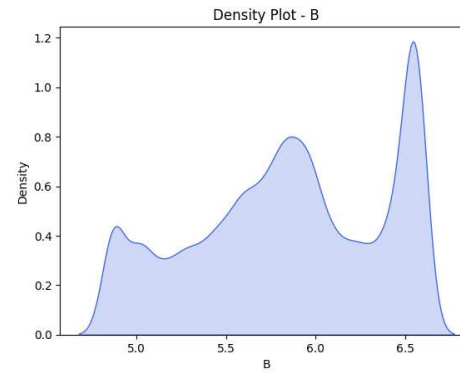
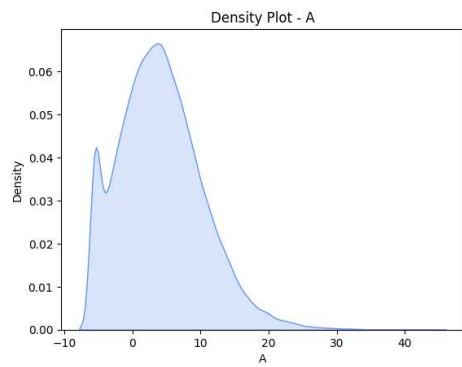
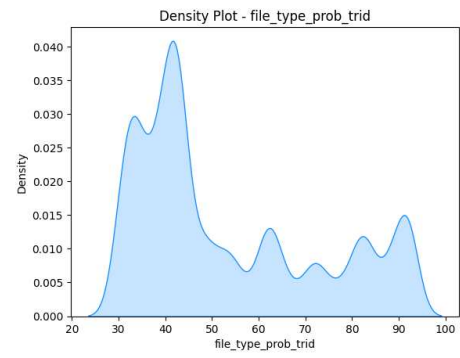
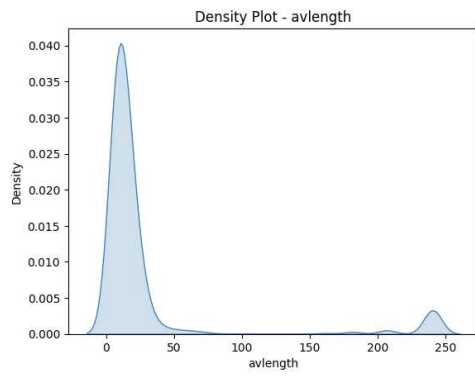




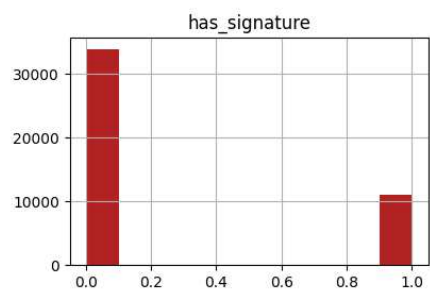
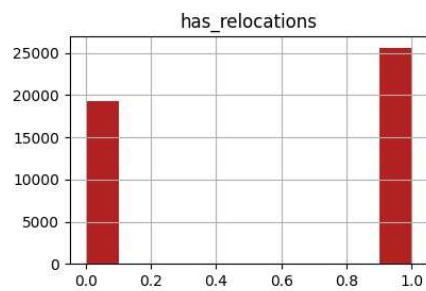
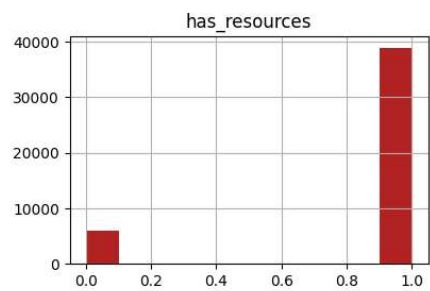
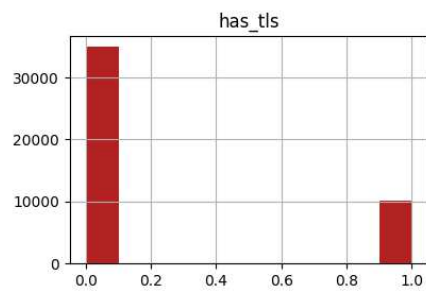
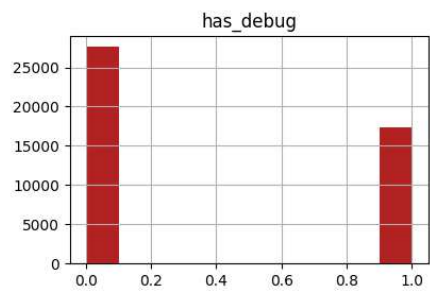
## 2.2



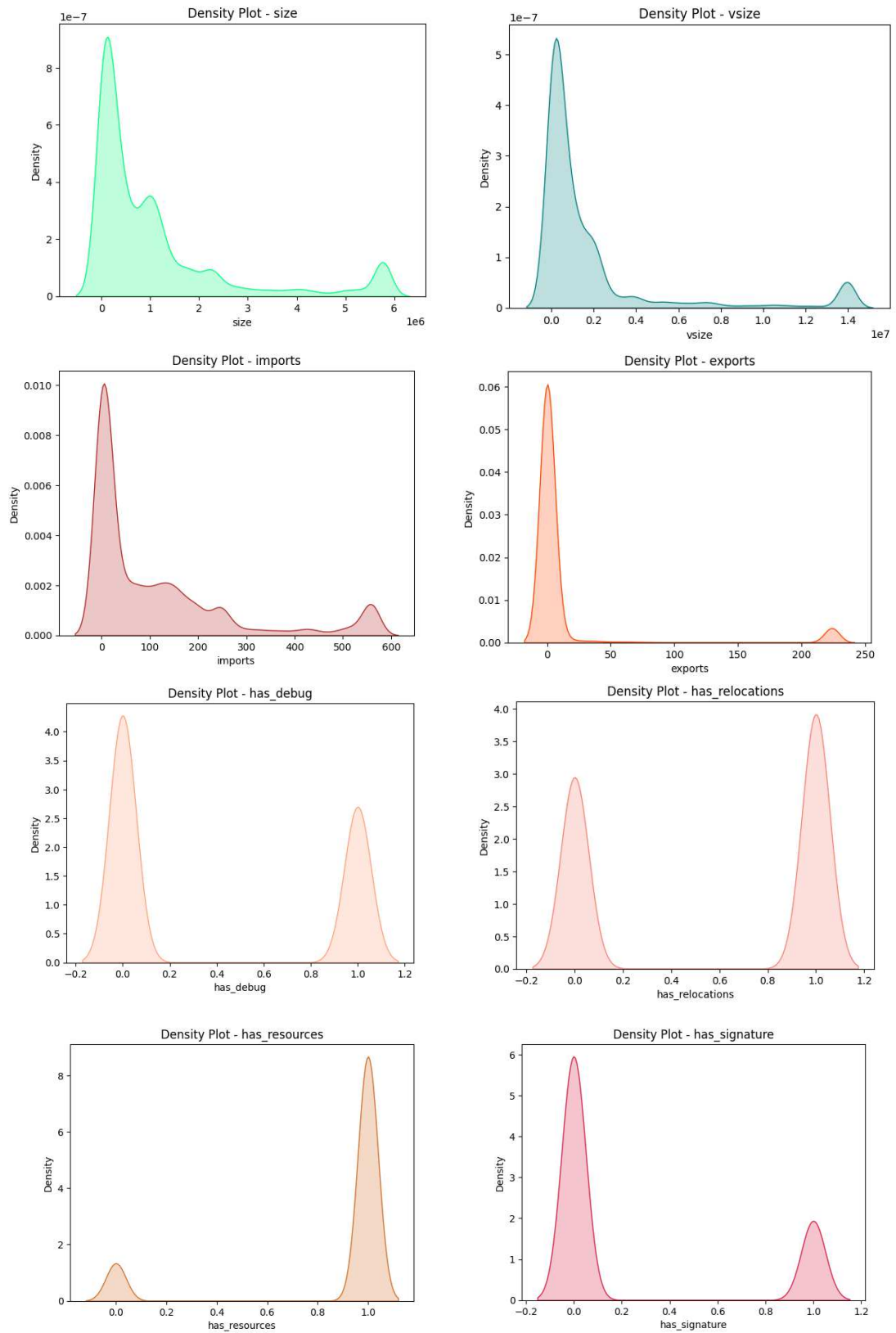




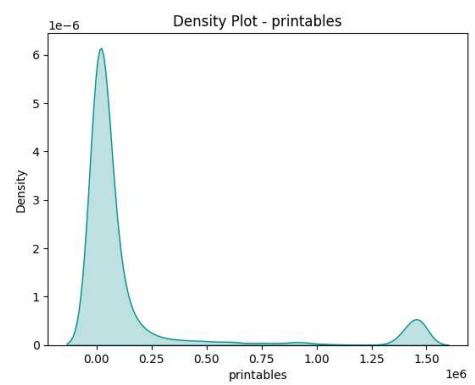
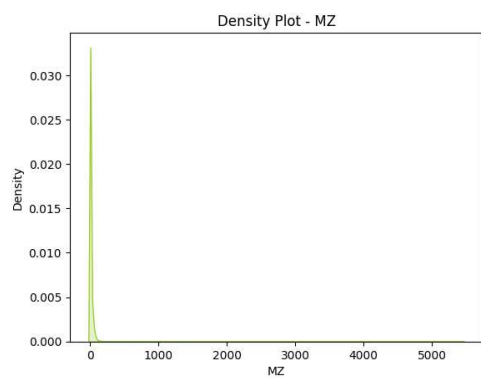
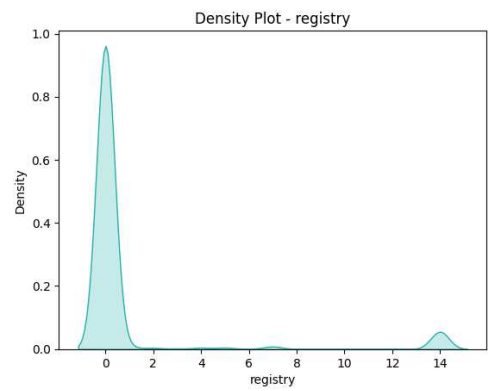
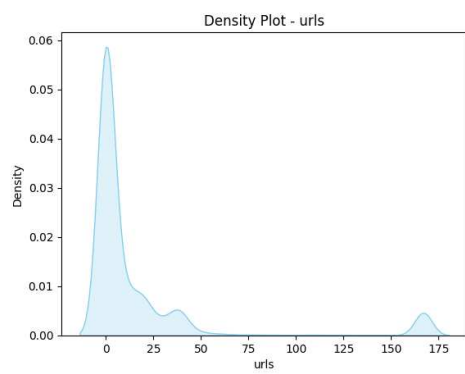
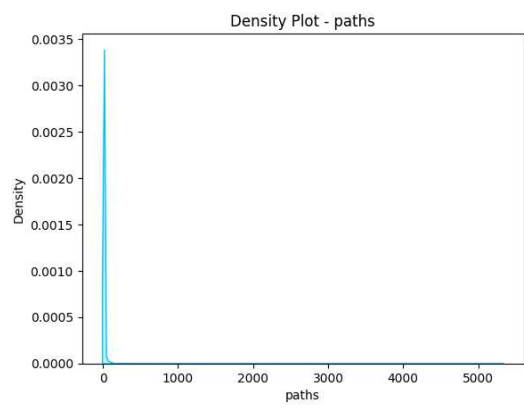
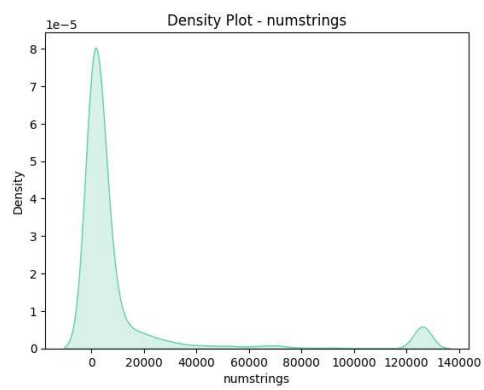
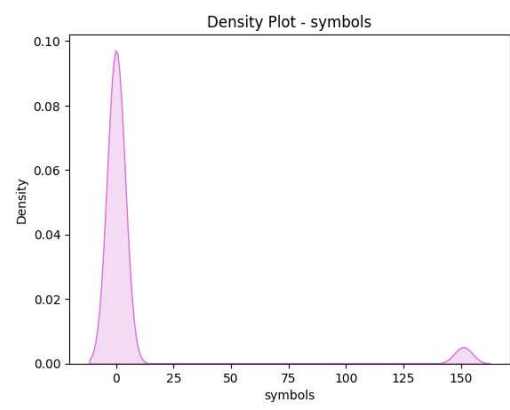
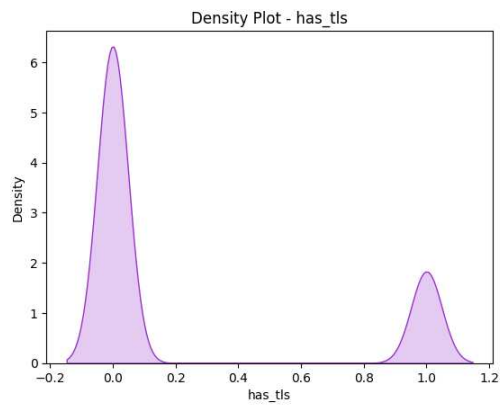
## 2.3

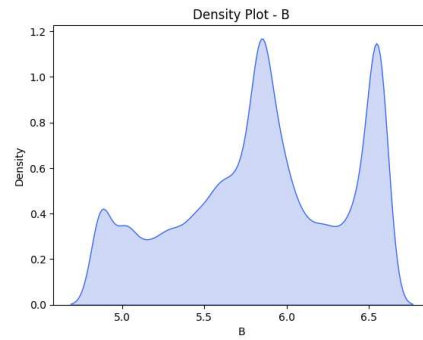
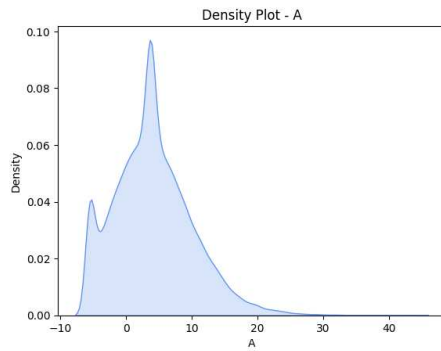
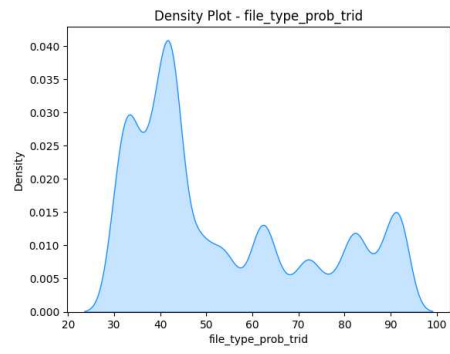
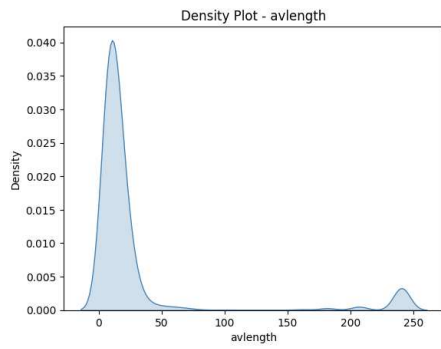


## 2.4

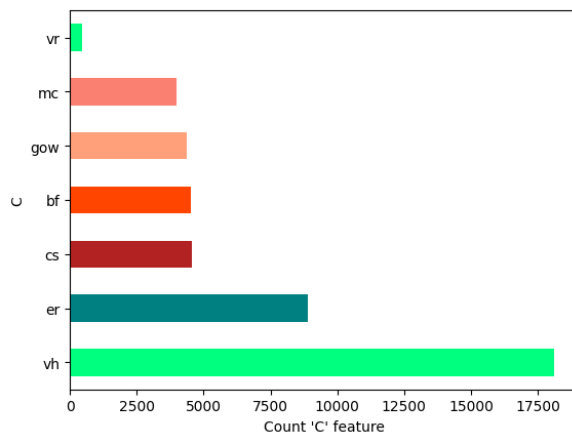




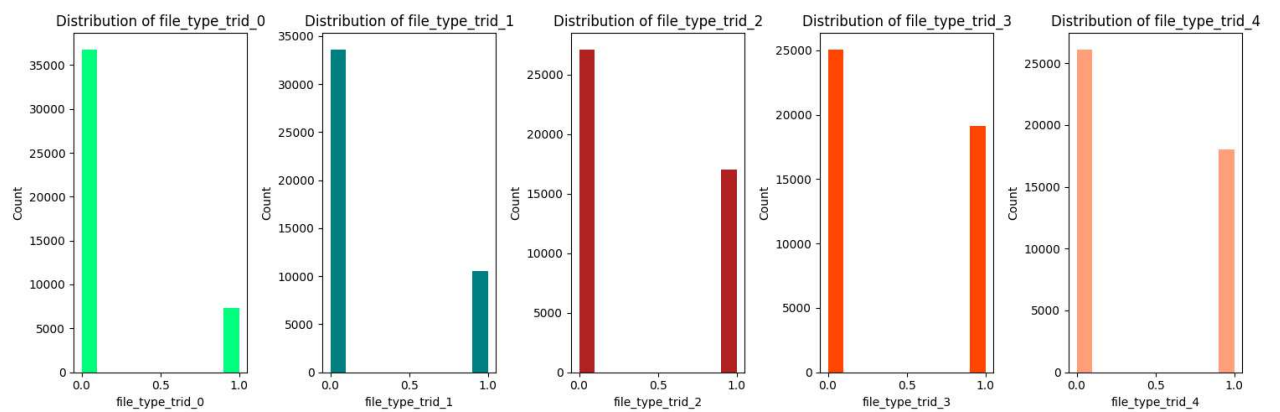




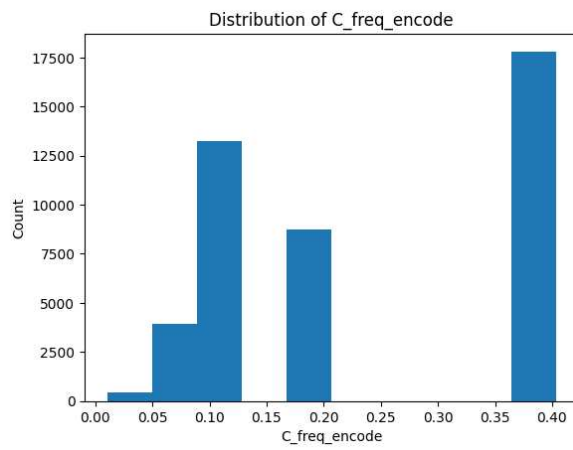
2.5



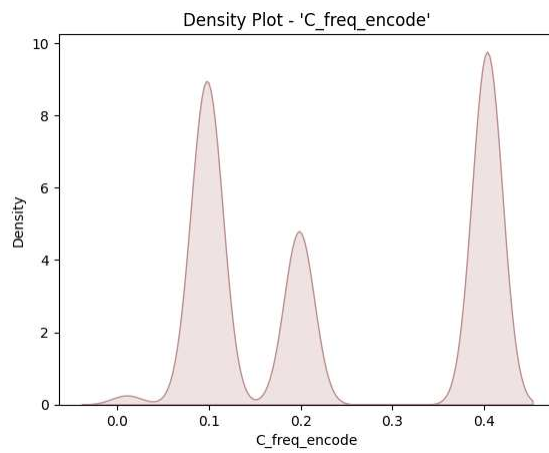
2.6



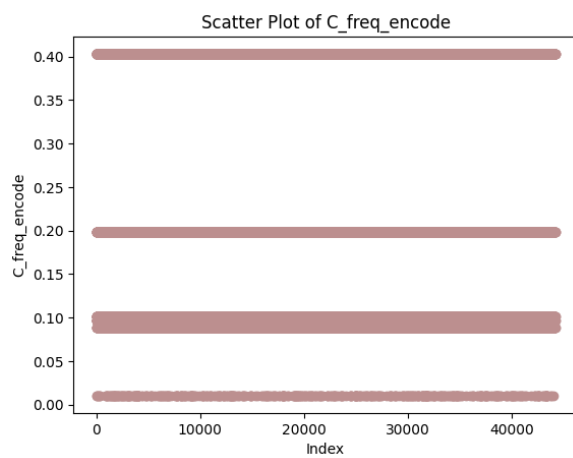
2.7



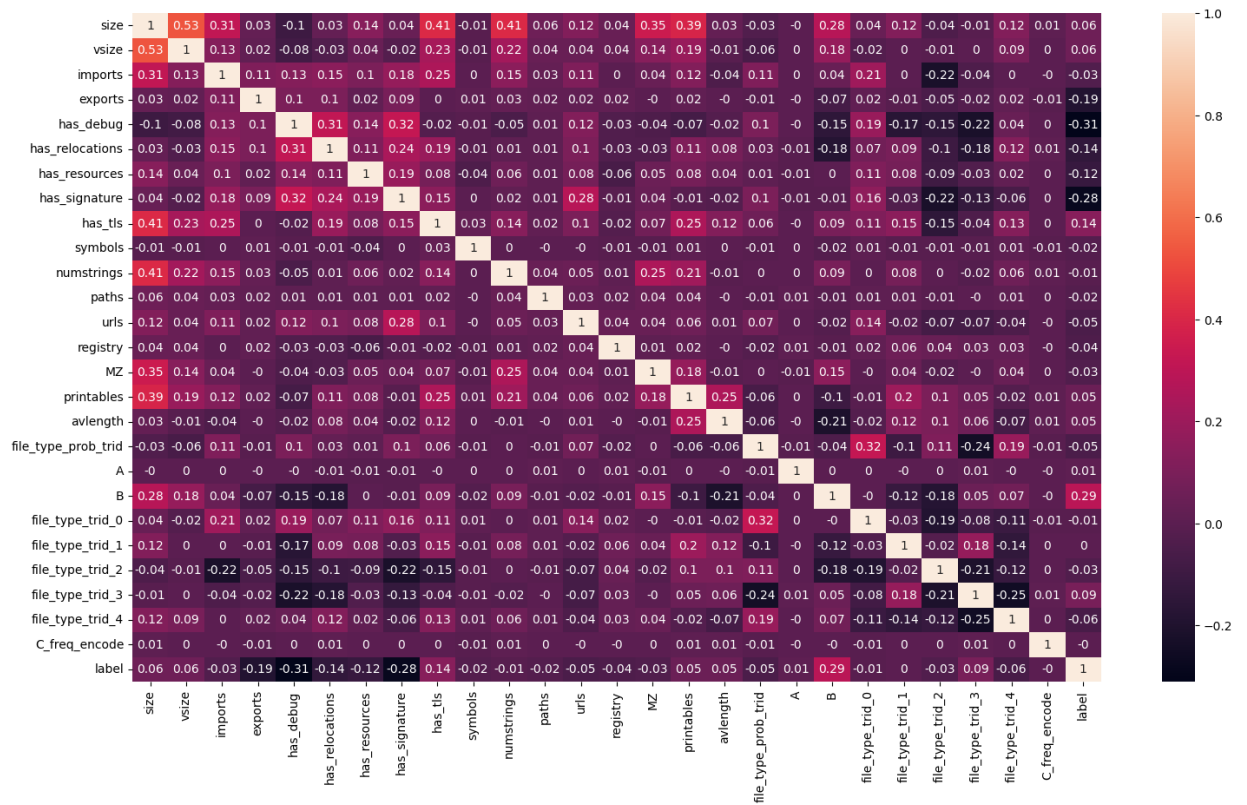
2.8



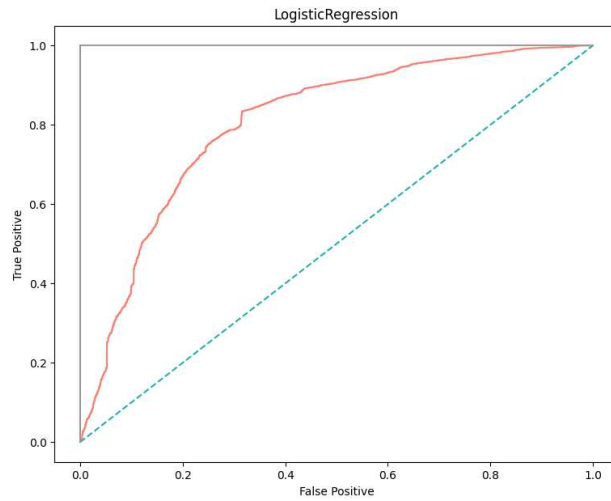
2.9



## 2.10

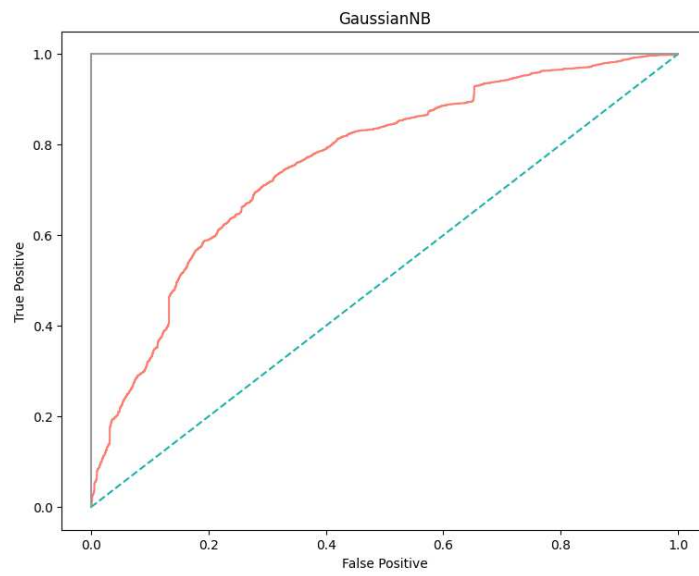


## 3.1

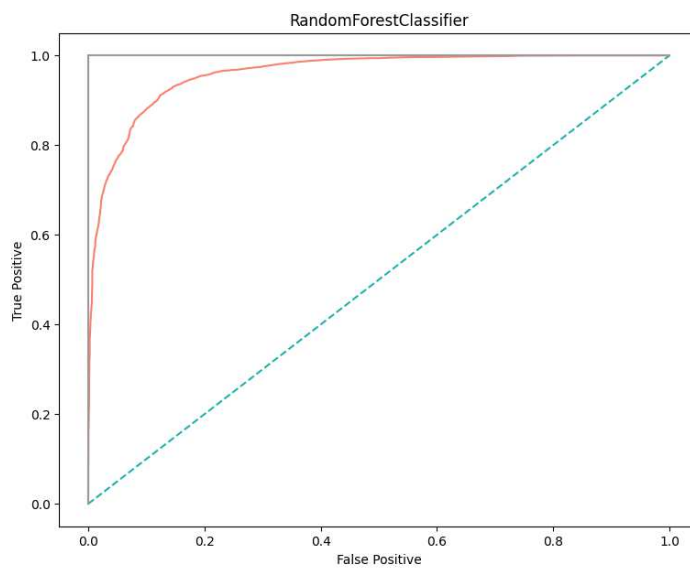




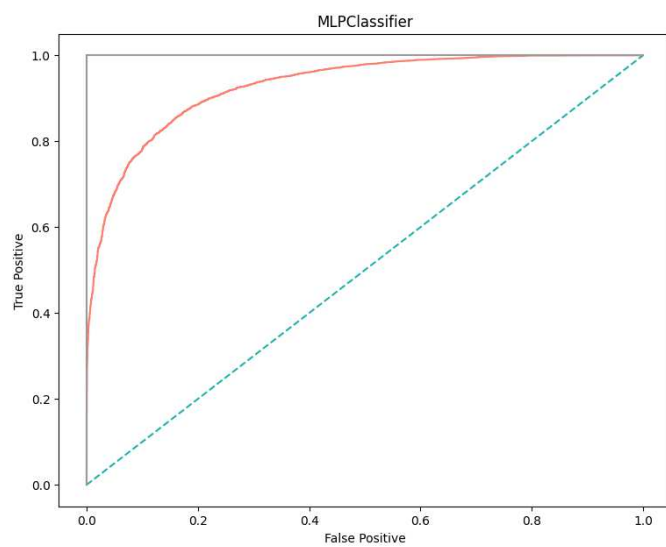
**3.2**



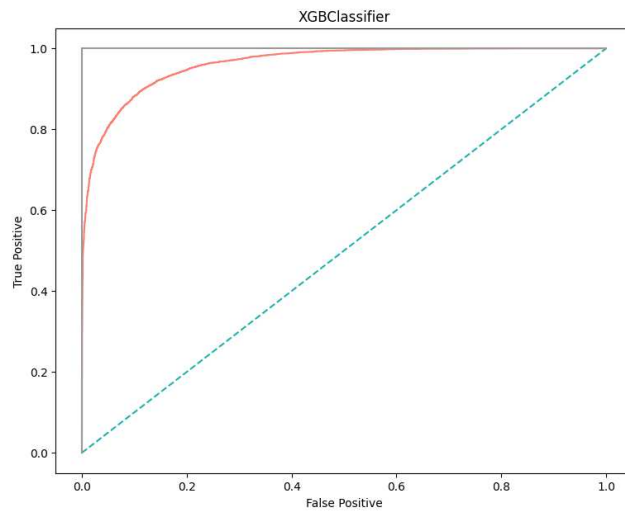
**3.3**



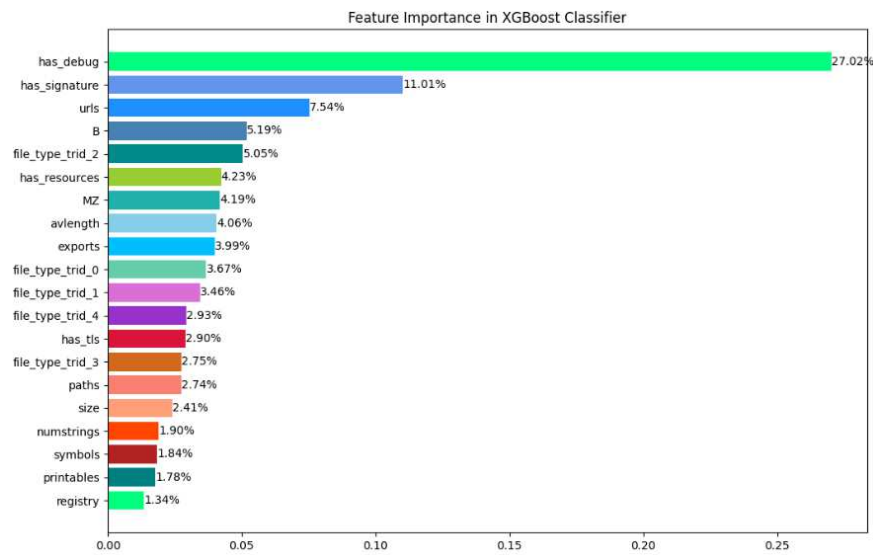
**3.4**



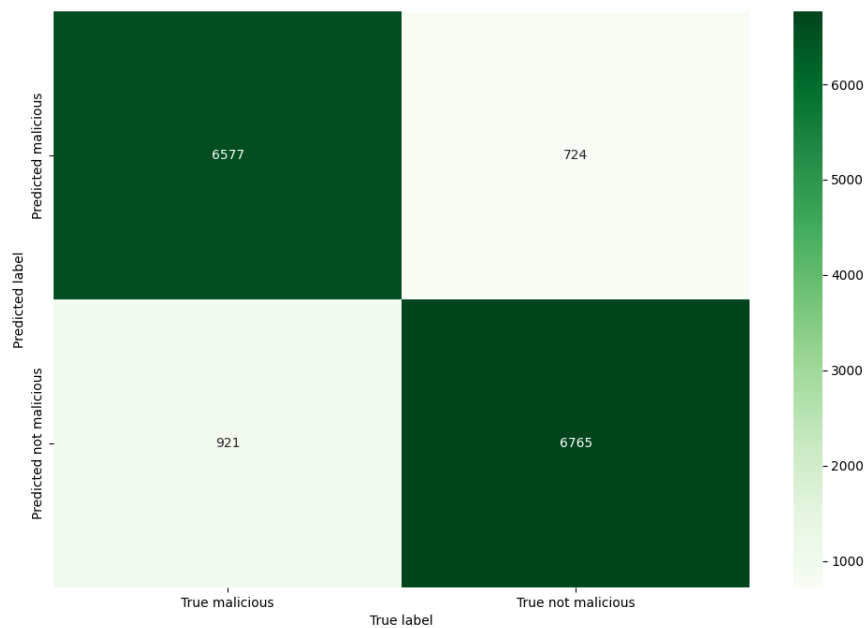
3.5



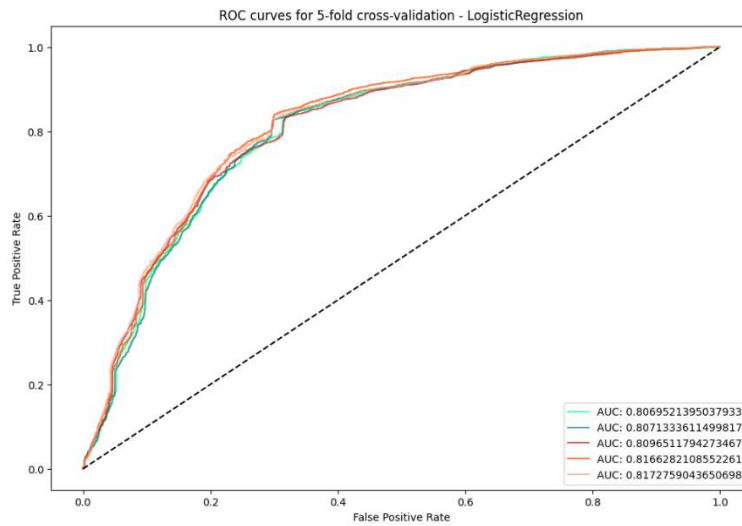
3.6



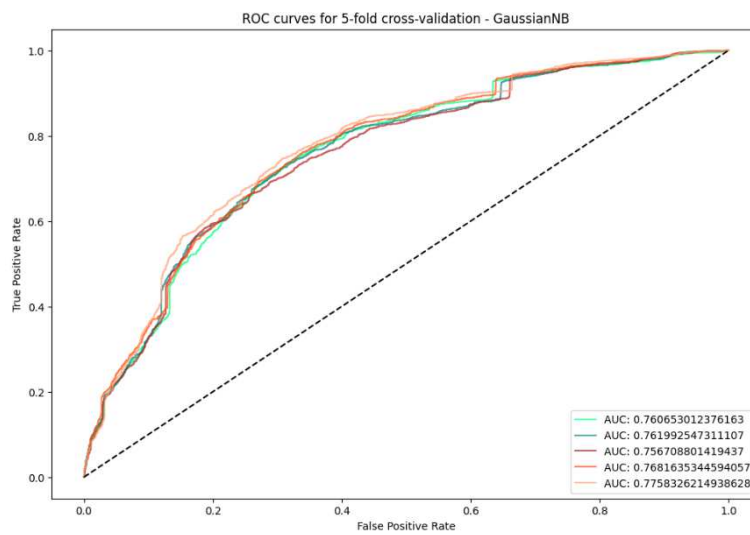
4.1



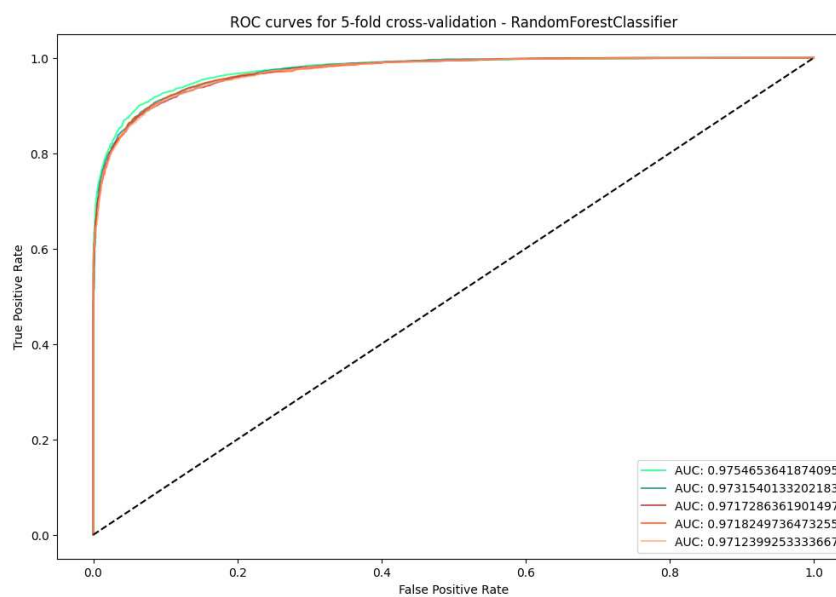
## 4.2



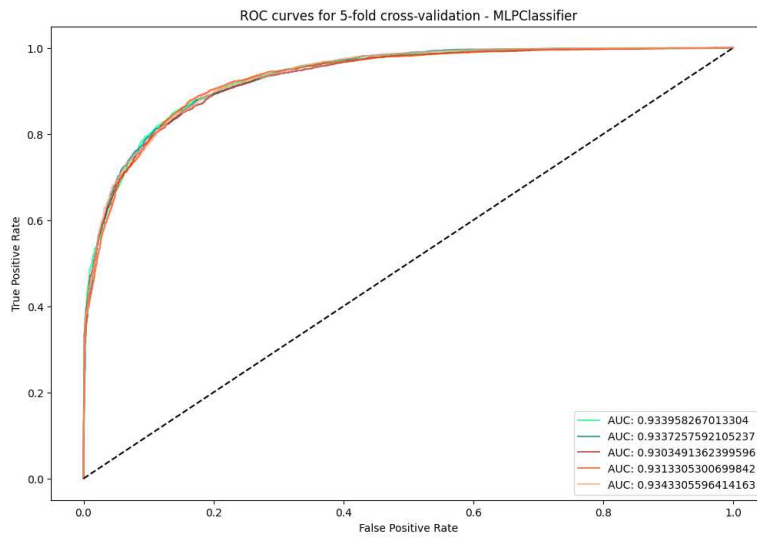
## 4.3



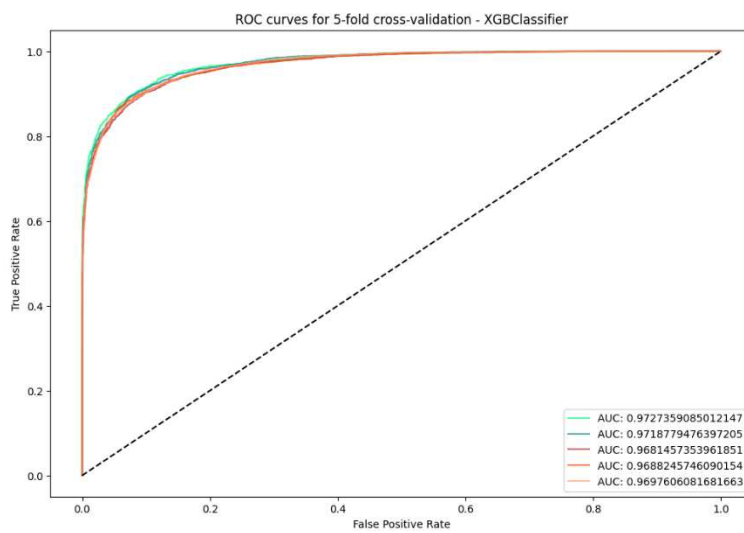
## 4.4



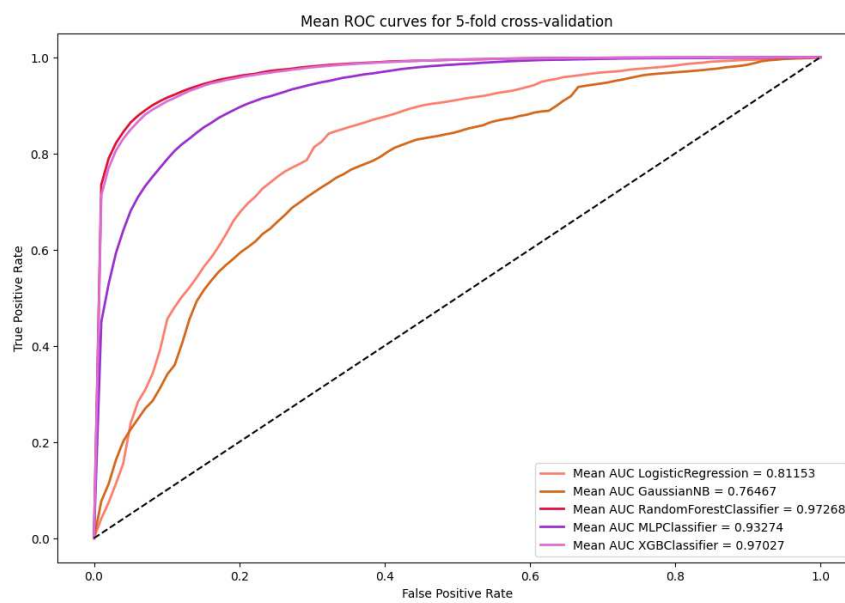
4.5



4.6

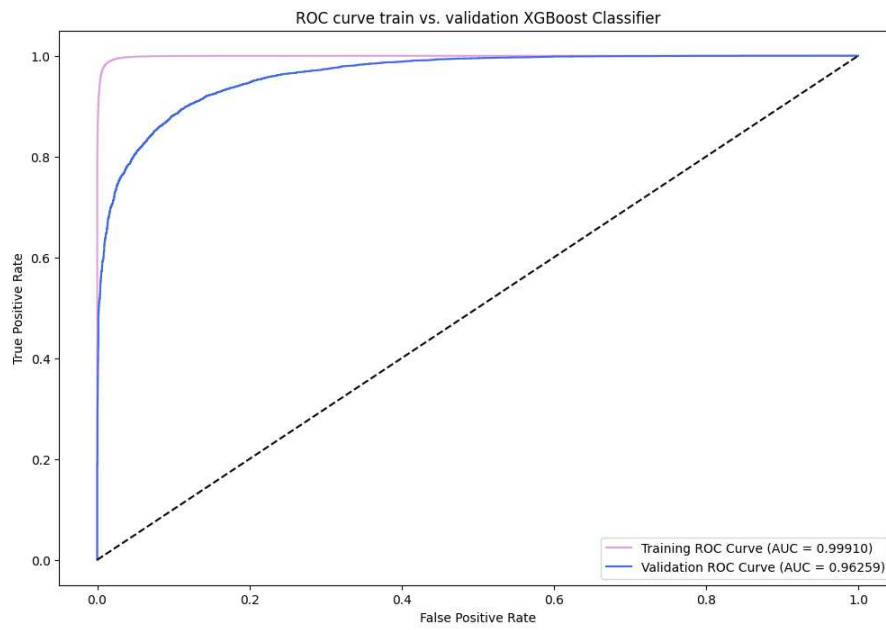


4.7



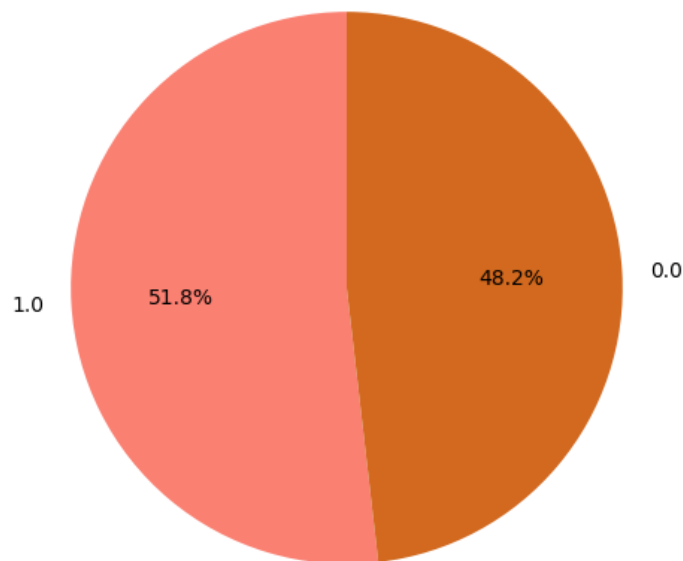


4.8



5.1

not malicious (0) vs. malicious (1)



## 6.1

### **XGBoost (Extreme Gradient Boosting) Classifier**

Boosting is an ensemble method, meaning it's a way of combining predictions from several models into one.

Gradient boosting is a specific type of boosting, called like that because it minimizes the loss function using a gradient descent algorithm.

XGBoost is a gradient boosting algorithm that uses decision trees as its "weak" predictors. Beyond that, its implementation was specifically engineered for optimal performance and speed.

To understand how XGBoost Classifier works, let's break down its key components and the underlying principles:

- **Loss Function:** XGBoost Classifier uses a loss function to quantify the difference between the predicted and actual labels. The goal is to minimize the loss function by iteratively improving the model.
- **Gradient Boosting:** XGBoost Classifier utilizes the gradient boosting framework. It builds an ensemble of weak models in a sequential manner, where each subsequent model corrects the mistakes made by the previous models. This iterative process improves the ensemble's predictive power.
- **Decision trees as base learners:** XGBoost Classifier employs decision trees as the base learners. These trees are shallow, usually with a small depth, to prevent overfitting. Each decision tree predicts the class labels based on a subset of features. The construction of decision trees is done using gradient boosting, which involves fitting the tree to the negative gradient of the loss function.
- **Gradient Calculation:** During the training process, XGBoost Classifier calculates the gradient (partial derivative) of the loss function with respect to the predictions made by the ensemble so far. The gradient reflects the direction and magnitude of the improvement needed to minimize the loss function. This calculation is performed for each training instance.
- **Weighted Residuals:** After calculating the gradient, XGBoost Classifier converts it into weighted residuals. The weights are determined based on the loss function and the predicted probabilities or scores from the previous iteration. These weighted residuals represent the "errors" made by the ensemble, and the subsequent weak model (decision tree) aims to correct these errors.
- **Tree Construction:** XGBoost Classifier builds decision trees sequentially, with each tree constructed to minimize the weighted residuals. It uses a technique called "regularized learning" to control the complexity of the trees and prevent overfitting. The regularization terms, such as L1 and L2 regularization, are added to the loss function during tree construction.
- **Ensemble Combination:** The individual decision trees are added to the ensemble one by one, with each tree contributing to the final prediction. The predictions from all the trees

are combined through weighted averaging, where the weights are determined based on the performance of the trees on the training data.

The values are typically calculated using the log of odds and probabilities. The output of the tree becomes the new residual for the dataset, which is used to construct another tree. This process is repeated until the residuals stop reducing or for a specified number of times. Each subsequent tree learns from the previous trees and is not assigned equal weight, unlike how Random Forest works.

Historically, XGBoost has performed quite well for structured tabular data.