

EEE-485 Final Project Report**Malicious URL Detection****1. Introduction**

As technology keeps progressing quickly and the use of the internet becomes more and more common in daily life, malicious URLs have become a real threat. These are links that are created artificially in order to promote scams, frauds, and phishing and are even used for cyberattacks. It is important for one to understand whether a URL is malicious or not before clicking the link; otherwise, one might be a victim of people with bad intentions. Hence, we aim to design and implement a model that can effectively predict whether a URL is safe or not with the help of various machine-learning algorithms.

2. Task

The main tasks of this project include data processing, model training, performance evaluation, and comparison. In this process, we first examine the data in our hands and add some new features to it, which will be explained in more detail in the Preprocessing section. Then, we use Logistic Regression, Multinomial Naive Bayes, and Neural Networks as machine-learning and deep-learning algorithms in order to test and train our model. Then, various complexities of these models are tested while measuring the performance and tuning the necessary hyperparameters. In the end, we will compare the performances and the efficiencies of these models according to some metrics, such as accuracy, and conclude which one fits our purpose the best.

3. Dataset

We used a dataset named “*Malicious URLs dataset*” that we found on *kaggle*, which has been formed by collecting URL data from different sources, including ISCX-URL2016 [1]. This dataset consists of 651191 URLs in total, a data frame that contains the names of the URLs together with their labels, out of which 428103 are considered to be benign, while the remaining are malicious. This dataset is great in terms of providing us with a huge amount of data along with their labels. However, it seems to lack information about some of the features that could be useful in training our model, and therefore, we performed some feature engineering on the existing data in order to extract some features, such as the length of the URL, length of the domain name; and the number of some special characters that appear in the URL, such as “-, @, ?, =”, or the number of letters and digits, and the counts of the special words like *http*, *https*, *www* and so on (explained in more detail in the Preprocessing section). For this dataset, we used 70-10-20 as our training, validation, and test splits, and we trained our final models on training + validation and reported metrics on the test dataset.

	url	type
0	br-icloud.com.br	phishing
1	mp3raids.com/music/krizz_kaliko.html	benign
2	bopsecrets.org/rexroth/cr1.htm	benign
3	http://www.garage-pirene.be/index.php?option=...	defacement
4	http://adventure-nicaragua.net/index.php?optio...	defacement

Fig.1: The First 5 Rows of The Dataset [1]

4. Programming Environment

During the implementations, the main programming language used is Python, and the version used is 3.9.7. Visual Studio Code is used as the IDE as it helps the organization and maintenance of the project together with the version control tools such as Git and GitHub. Throughout the project, we carried out the developments on a collaborative repository on GitHub.

We plan to benefit from NumPy library for efficient linear algebra operations, Pandas for data preparation and wrangling, Matplotlib for data visualization, Python’s tld library to extract useful information from URLs, and tqdm to create progress bars for loops.

5. Data Preprocessing

The initial dataset before preprocessing only consisted of long strings of URLs. Therefore, we had to do feature engineering and data preprocessing in order to come up with features that could be useful for our models [2][3].

First, we used a method to determine whether the URL is an IP address or not, and if it is, which type of URL it is. As a result, we came up with a feature named ‘*use_of_ip*’ that takes binary values [3][4]. Then, we benefited from the *get_tld()* function of the tld library and used it to extract

full-length domain (fld), top-level domain (tld), domain, and subdomain names if they exist [5][6]. We then found their lengths, which resulted in 4 new features: ‘url_length’, ‘subdomain_length’, ‘tld_length’, ‘fld_length’, and ‘path_length’. We thought length of these could be a useful feature as we observed that malicious URLs tend to be longer than the benign ones, and this is especially the case for the domain names. We also added another binary feature called, and ‘use_of_shortener’, which indicate the path of the URL (for example, the remaining part of the URL after ‘.com’ and whether the URL contains a shortening service (for example, Bitly or TinyURL) or not, respectively. We found the list of shortening services from a GitHub repository and searched for the existence of them [7]. For the final version of our model, we also added a method to see if the URL contains some suspicious words, such as: cash, free, lucky, bonus... and added a new feature ‘is_susp’ accordingly. Similarly, when we investigated our dataset, we noticed that many of the malicious URLs contained the word ‘php’, hence we added a new feature ‘is_php’, to see if the URL contains that word or not.

We then added more features to represent the number of occurrences of alphanumeric characters, digits, and punctuations, named ‘count_letters’, ‘count_digits’, and ‘count_puncs’. Then we counted the number of some special characters in the URL format in order to determine whether a URL could be malicious or not, and therefore, we added the following features as well: ‘count.’, ‘count@’, ‘count%’, ‘count?’, and ‘count=’ [2]. We used simple lambda functions and Python’s count() function in order to obtain these numbers. Again for our final model, we added features of counts of some other characters as well: ‘count+’, ‘count#’, ‘count/’, ‘count\$', ‘count!', ‘count*’, ‘count,’ ‘count_’, and ‘count:’. However, as these features are not directly correlated with the URL structure, incorporating all of them simultaneously resulted in a degradation of the models’ performance and thus, we employed techniques for feature selection to mitigate this issue. In addition to these features, we also attempted to incorporate the count of the ‘%20’ character (used to represent spaces in URLs) as an additional feature, however, through experimentation, we discovered that this particular feature adversely affects the performance of all machine learning models [8].

Additionally, we added ‘count_dirs’, and ‘first_dir_length’, which represent the number of the subdirectories in the URL (this is in fact equal to the number of ‘/’ characters) and the length of the first string, respectively. We did this extraction based on our previous observation that malicious URLs tend to be longer in length. We also added another two features named ‘url_length_q’ and ‘fld_length_q’, to distribute url_length and fld_length into 4 groups named ‘Short’, ‘Medium’, ‘Long’, and ‘Very Long’ by using percentiles based on the distribution of the data as. We obtained categorical features as a result of this operation, therefore, we later converted these values to integers using the fit_transform() method of OrdinalEncoder.

We also checked whether the URL starts with ‘http’ or ‘https’, and their number of occurrences as well as the occurrence of ‘www’. Hence, we added: ‘https’, ‘count-https’, ‘count-http’, and ‘count-www’ and all these features take binary values since a URL doesn’t contain more than one of these.

Then we added ‘letters_ratio’, ‘digit_ratio’, and ‘punc_ratio’, which we simply obtained by dividing count_letters, count_digits, and count_puncs by url_length since we thought they could be helpful for the Logistic Regression model, but we can’t use these features for the Naive Bayes models since they take non-integer values.

Finally, we added the label feature called ‘is_malicious’. In our initial dataset, the label feature consisted of 4 different values: Benign, Defacement, Phishing, and Malware. We decided to label URLs that are Benign as 0 and all the others that are malicious as 1. Therefore, we converted our problem into a binary classification in which we simply aim to determine whether a URL is benign or not. With respect to these changes, we added the new feature ‘is_malicious’. After doing this, the class distributions became around 65-35 and this skewness of the dataset didn’t create any problem for us when training our models and thus, we didn’t apply any oversampling or undersampling technique.

We saved this new data frame with the name ‘url_processed.csv’, and in both of our models, we used %70 of our data as the training set, %10 as the validation set, and %20 as the test set.

6. Detailed Description of the Proposed Methods

In this project, we will examine the effectiveness of different machine learning algorithms on our URL classification problem, which will include Logistic Regression, Multinomial Naive Bayes, and Neural Networks. All these three types of models are known to be efficient in classification

problems and hence, are suitable for our binary classification problem. Now, we will explain each of these methods in more detail.

6.1. Logistic Regression

Logistic Regression, first of all, is very efficient for classification problems and also efficient to train, and it's quite advantageous when the data is linearly separable. Besides that, the coefficients obtained in the Logistic Regression can be interpreted to get insights about the features, which can be helpful for the Final Report when picking the best subset of features for obtaining the most generalizable model. A Logistic Regression model tries to fit the logistic model to the data, which models the probability of an event and estimates the necessary parameters for that. For the Logistic Regression Model, we have the following maximum conditional likelihood estimators:

$$P(Y = 0 | X, w) = \frac{1}{1 + e^{(w_0 + \sum_i w_i X_i)}} = \frac{1}{1 + e^{w^T x}} \text{ and}$$

$$P(Y = 1 | X, w) = \frac{e^{(w_0 + \sum_i w_i X_i)}}{1 + e^{(w_0 + \sum_i w_i X_i)}} = \frac{e^{w^T x}}{1 + e^{w^T x}}$$

where $X = [1, X_1, X_2, \dots, X_p]^T$ is the feature vector and $w = [w_0, w_1, \dots, w_p]^T$ is the parameter (weight) vector. Using these, we can form the likelihood function as

$$L(w) = P(D|w) = \prod_{i: y_i=1} \frac{e^{w^T x_i}}{1 + e^{w^T x_i}} \prod_{i: y_i=0} \frac{1}{1 + e^{w^T x_i}}$$

where $D = \{(x_i, y_i)\}$ for $i = 1, 2, \dots, n$ and n denotes the total number of samples. For making the computations easier, let's find the log-likelihood function:

$$l(w) = \ln(P(D|w)) = \sum_{i=1}^n [y_i w^T x_i - \ln(1 + e^{w^T x_i})]$$

The derivative of this function is:

$$\nabla_w l(w) = \ln(P(D|w)) = \sum_{i=1}^n x_i [y_i - \frac{e^{x_i^T w^{(t)}}}{1 + e^{x_i^T w^{(t)}}}]$$

At the MLE estimate of w , the derivative of this function is 0. But this equation has no closed-form solution; therefore, we will use Gradient Update to converge to a local maximum. Hence, we can form the gradient update for the weight vector w can be derived as the following:

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \alpha \sum_j x_i^j [y^j - \hat{P}(Y^j = 1 | x_j, w^{(t)})] = w_i^{(t)} + \alpha \sum_j x_i^j [y^j - \frac{e^{x_j^T w^{(t)}}}{1 + e^{x_j^T w^{(t)}}}]$$

for t in range $[0, \# epochs]$

In the above notation, j denotes an index to a data sample and $x_0^j = 1$, y_j is the actual label for the data sample j , $w_0 = b$ is the bias term and $w_i^{(t)}$ denotes the i th element of the weight vector at iteration t and finally α denotes the learning rate for the update.

6.2. Naive Bayes Models with Bag-Of-Words Representation

Naive Bayes classifiers are based on the application of the Bayes Theorem, and they have ‘naive’ (conditional independence) assumptions for the features. Because of this, the Naive Bayes models handle the noise well. Also, the training using the maximum likelihood for these models is quite cheap, requiring linear time; the classification of a single data point is also quite fast by only picking the highest probability class. These models are useful in the classification for large datasets, as in the case of URL dataset since they’re highly scalable with the number of data samples and the number of features. We will use the Multinomial and Bernoulli models together with the bag-of-words representation. The bag-of-words document representation has the assumption that a feature’s (word)

existence is conditionally independent of its position given the class of the document, and we will use this assumption for both models.

6.2.1 Multinomial Naive Bayes Model

The Multinomial Naive Bayes Model takes multiple occurrences of a feature (word) into account and calculates probabilities with respect to occurrences of other features (words). Suppose we have a document D_i with n_i features in it, then the probability of D_i belonging to class y_k can be computed as

$$P(D_i|Y = y_k) = P(X_1 = x_1, X_2 = x_2, \dots, X_{n_i} = x_{n_i} | Y = y_k) = \prod_{j=1}^{n_i} P(X_j = x_j | Y = y_k)$$

In the above equation X_j represents the feature at j^{th} position in the document D_i and x_j represents the actual feature that's in the j^{th} position in the document, and lastly n_i denotes the number of positions in the document. In the above formulation, the length of the feature vector for the document, the length of the $X^{(i)}$, depends on the number of features in the document, which is n_i . But looking from the perspective of the set of features, we can also formulate this as

$$P(D_i|Y = y_k) = \prod_{j=1}^{|V|} P(X_j | Y = y_k)^{t_{w_j,i}}$$

In this notation, $|V|$ is the size of the feature set V , and X_j represents the appearance of the j^{th} feature and $t_{w_j,i}$ indicates how many times feature w_j appears in a document D_i . In classifying the documents, we are interested in the probability distributions over a set of possible classes, and hence, we need to calculate $P(Y = y_k | D_i)$. Using the Bayes Rule, this is equal to:

$$P(Y = y_k | D_i) = \frac{P(Y = y_k) \prod_{j=1}^{|V|} P(X_j | Y = y_k)^{t_{w_j,i}}}{\sum_k P(Y = y_k) \prod_{j=1}^{|V|} P(X_j | Y = y_k)^{t_{w_j,i}}}$$

We realize that the denominator is fixed for each class y_k and thus, we can say

$$P(Y = y_k | D_i) \propto P(Y = y_k) \prod_{j=1}^{|V|} P(X_j | Y = y_k)^{t_{w_j,i}}$$

The class for which this probability is maximized will be the prediction:

$$\hat{y}_i = \operatorname{argmax}_k P(Y = y_k | D_i) = \operatorname{argmax}_k P(Y = y_k) \prod_{j=1}^{|V|} P(X_j | Y = y_k)^{t_{w_j,i}}$$

Since each probability in the above multiplication is in the interval $[0, 1]$, we might face underflow issues in programming this method. Since the logarithm is a strictly increasing function in the interval $[0, 1]$, maximizing the above probability is the same as maximizing the logarithm of it. Therefore,

$$\hat{y}_i = \operatorname{argmax}_k P(Y = y_k | D_i) = \operatorname{argmax}_k [\ln(P(Y = y_k)) + \sum_{j=1}^{|V|} t_{w_j,i} * P(X_j | Y = y_k)]$$

where

$$P(X_j | Y = y_k) = \frac{T_{j,Y=y_k}}{\sum_{j=1}^{|V|} T_{j,Y=y_k}} \text{ and } P(Y = y_k) = \frac{N_k}{N}$$

and here, N indicates the total number of samples in the training set while N_k indicates the total number of samples of with true label y_k . Also, $T_{j,Y=y_k}$ is the number of occurrences of the j^{th} feature

in documents of class y_k in the training set, including the multiple occurrences of that feature. Here, we can also possibly apply smoothing to these estimated parameters, which is by adding imaginary occurrences (α times where α is the smoothing parameter) of every feature to every class, and parameters become:

$$P(X_j|Y = y_k) = \frac{T_{j,Y=y_k} + \alpha}{\alpha * |V| \sum_{j=1}^{|V|} T_{j,Y=y_k}}$$

This way, the model becomes more “fair” by hallucinating each feature in each class, not directly yielding to the elimination of a document because of the absence of a specific feature.

6.2.2 Bernoulli Naive Bayes Model

The Bernoulli Model is interested in the absence or presence of the features and represents this with 0 and 1, where the latter indicates the presence of the feature. Unlike Multinomial, a probability for a parameter belonging to a specific feature is independent of the other features. Again, we will predict as the class for which the class probability is maximized, which is:

$$\begin{aligned} \hat{y}_i &= \underset{|V|}{\operatorname{argmax}}_k P(Y = y_k | D_i) \\ &= \underset{|V|}{\operatorname{argmax}}_k P(Y = y_k) \prod_{j=1}^{|V|} (t_j P(X_j|Y = y_k) + (1 - t_j)(1 - P(X_j|Y = y_k))) \end{aligned}$$

Again taking the logarithm, we have

$$\hat{y}_i = \underset{|V|}{\operatorname{argmax}}_k [\ln(P(Y = y_k)) + \sum_{j=1}^{|V|} \ln(t_j P(X_j|Y = y_k) + (1 - t_j)(1 - P(X_j|Y = y_k)))]$$

To calculate this, the needed parameters are

$$P(X_j|Y = y_k) = \frac{S_{j,Y=y_k}}{N_k} \text{ and } P(Y = y_k) = \frac{N_k}{N}$$

where N indicates the total number of samples in the training set while N_k indicates the total number of samples of with true label y_k . Also, $S_{j,Y=y_k}$ is the number of occurrences of the j^{th} feature in documents of class y_k in the training set, not including the multiple occurrences of that feature.

6.3. Neural Networks

Finally, we have chosen to implement Neural Networks as our preferred method due to their ability to learn non-linear and complex patterns in the data. Compared to previous methods, neural networks offer greater flexibility in fitting the data. A neural network is composed of multiple layers of interconnected neurons, and various parameters, such as the number of neurons in a layer, the connectivity between layers, and the number of hidden layers, can be adjusted or tuned when designing the model. After each layer, non-linear activation functions are applied to the outputs, enhancing the model's capacity to create more generalizable representations. By stacking multiple layers of neurons, the model can effectively capture intricate functions, albeit at the cost of optimization efforts. To update the model with new data, we have opted to test the batch gradient-descent (BGD) approach with batch sizes 32 and 64. The model weights are updated using back-propagation, which involves tracing from the deepest layer to the input layer and applying the chain rule to derive the necessary update equations for the parameters of the model. However, the process can sometimes get stuck in flat regions or oscillate around the minimum, leading to slower convergence. Momentum addresses this issue by introducing a "velocity" term that accumulates the past gradients and influences the current weight update. It adds a fraction of the previous velocity to the current gradient, making the updates more consistent and stable over time and therefore, we also introduced the momentum parameter in the trainings. In our URL classification task, we primarily experimented with the Rectified Linear Unit (ReLU) activation function for the hidden layers. This choice was motivated by its ability to alleviate the vanishing gradients problem and expedite the training process (essential because of the size of the dataset) compared to other activation functions.

Additionally, we intend to use the Binary Cross-Entropy loss function, sigmoid output activation, and a threshold to utilize our model as a classifier. To address potential overfitting issues, we might employ (if necessary) L2 and L1 regularization methods, which can enhance the model's generalization capabilities and prevent overfitting, thus improving its performance.

7. Preliminary Results

We implemented our Logistic Regression and Multinomial Naive Bayes models for this stage and tested and evaluated their final performances.

7.1. Logistic Regression

The new features that will be mentioned in 7.2.2 are added to the design matrix before training and testing the Multinomial Naive Bayes Model. We trained a total of 12 logistic regression models and tested their performances using our validation dataset. The first 4 of these 12 models have their weights initialized with normal distribution, the next 4 have their weights initialized with uniform distribution, and the remaining 4 with simply zeros. For each of the three different types of these initializers, we applied batch gradient descent by training with four different batch sizes, which are 32, 64, 128, and 256. We did this training using our train and validation data sets and fed each of our models with the training set a total of 20 times, so the epoch value picked was equal to 20. An example accuracy that we observed on our validation sets for each initialization method and different batch sizes can be seen in the plot below:

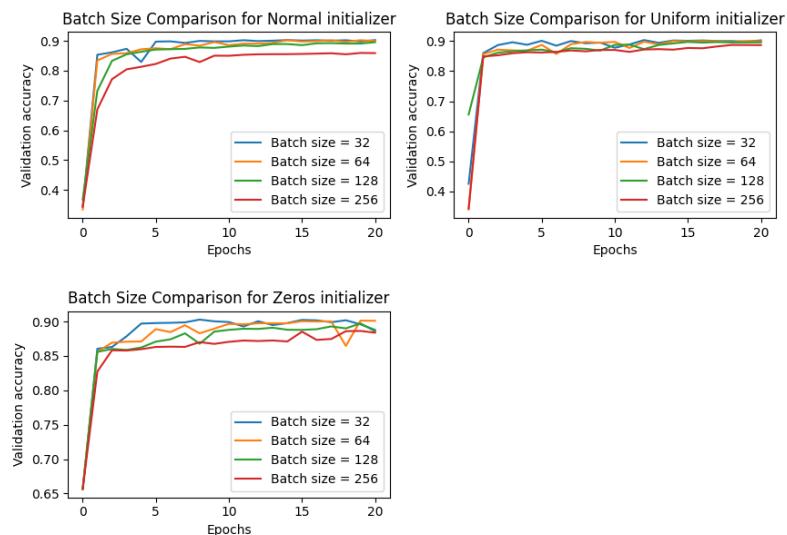


Fig.2: Validation accuracies for different initializations and batch sizes

The values for each accuracy points in the plot are:

```

Final Validation accuracy for Normal Initialization With Batch Size: 32 is : 0.899783
Final Validation accuracy for Normal Initialization With Batch Size: 64 is : 0.900475
Final Validation accuracy for Normal Initialization With Batch Size: 96 is : 0.885026
Final Validation accuracy for Normal Initialization With Batch Size: 128 is : 0.891721
Final Validation accuracy for Uniform Initialization With Batch Size: 32 is : 0.905373
Final Validation accuracy for Uniform Initialization With Batch Size: 64 is : 0.890800
Final Validation accuracy for Uniform Initialization With Batch Size: 96 is : 0.898309
Final Validation accuracy for Uniform Initialization With Batch Size: 128 is : 0.888758
Final Validation accuracy for Zeros Initialization With Batch Size: 32 is : 0.902670
Final Validation accuracy for Zeros Initialization With Batch Size: 64 is : 0.902241
Final Validation accuracy for Zeros Initialization With Batch Size: 96 is : 0.897419
Final Validation accuracy for Zeros Initialization With Batch Size: 128 is : 0.888143

```

On the machine equipped with Intel i7-10750h, the trainings took 17, 12 and 8 and 7 seconds on average, by using the tqdm library for 32, 64, 128, and 256 batch sizes, respectively.

The results of the above experiment were changing as each run used different seeds, but in several runs and as in the above example, the model in which Normal distribution was used with a batch size of 64 seemed to be one of the best performers in terms of accuracy, and it usually converged to a local minima faster than the other models. Therefore, we chose it as our final model,

trained it on our train and validation set, and then tested it on our test set. The final accuracy of our model can be seen in the plot below:

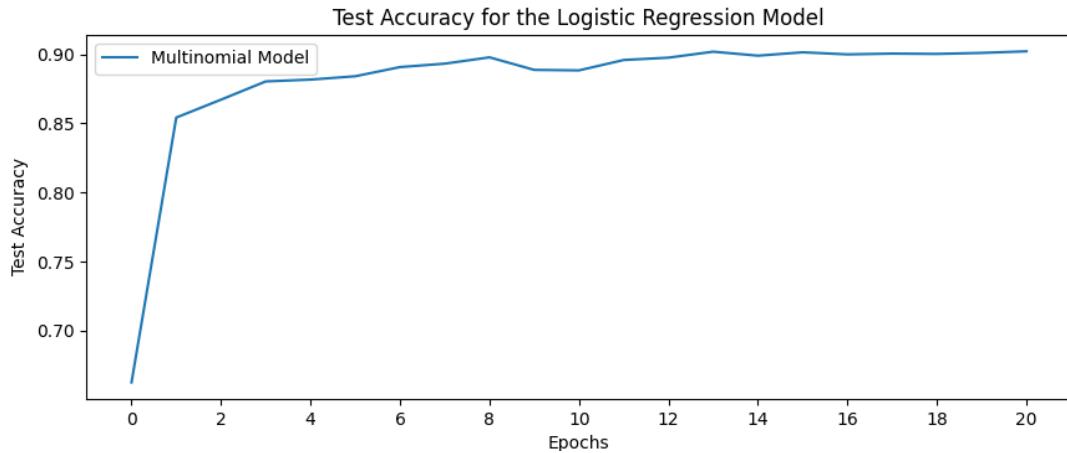


Fig.3: Final Test Accuracy for the Final Logistic Regression Model

For this test, the obtained confusion matrix was as follows:

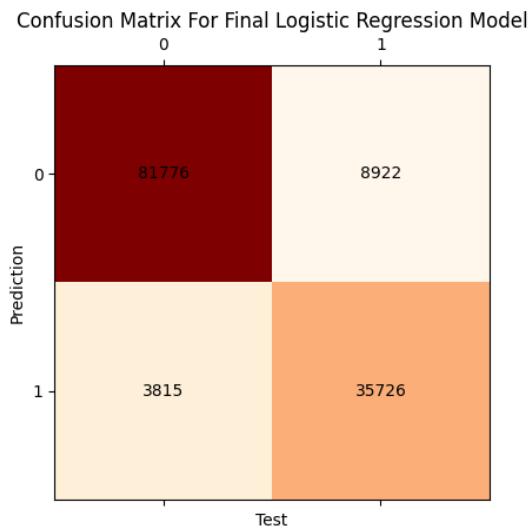


Fig.4: Confusion Matrix for the Final Logistic Regression Model

Final Test Accuracy: 0.90220

Final Test Precision: 0.9035

Final Test Recall: 0.8002

F1 Score For Final Logistic Regression Model: 0.84871

7.2.1. Multinomial Naive Bayes Model

The new features that will be mentioned in 7.2.2 are added to the design matrix before training and testing the Multinomial Naive Bayes Model. As mentioned earlier, we trained our Multinomial Naive Bayes model with 5 different smoothing parameters: 0, 2, 4, 6, and 8, to see the effects of this parameter on the model quality. The accuracy that we obtained on our validation set can be seen in the plot below:

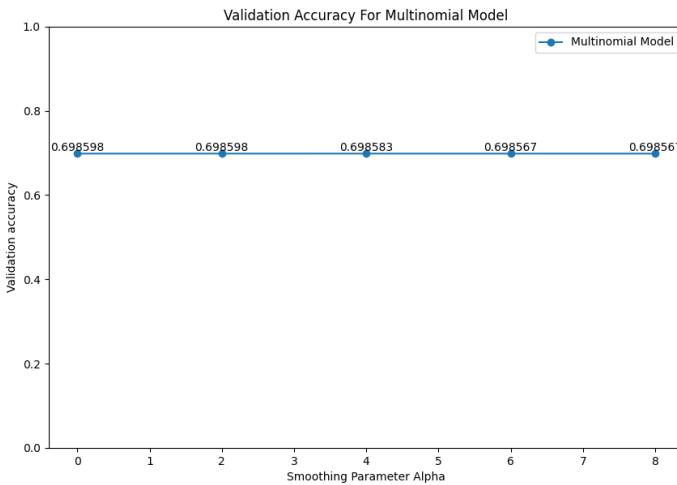


Fig.5: Validation Accuracy for the Multinomial Models With Different Smoothing Parameters

Again on the machine equipped with Intel i7-10750h, the trainings took around 0.04 seconds (Python's `time` library gave this result) for a single Multinomial Model since the calculations were quite fast as they're completed in linear runtime complexity and by avoiding for-loops with the help of numpy library.

Interestingly, as can be inferred from the results we obtained, increasing the smoothing parameter did not have much of an effect on our model after some point, and the model with no smoothing included seemed to have slightly better results on several runs. Therefore, we chose it as our final model and trained it from scratch on training + validation data, and tested it on our test set. According to this, the confusion matrix obtained was as follows.

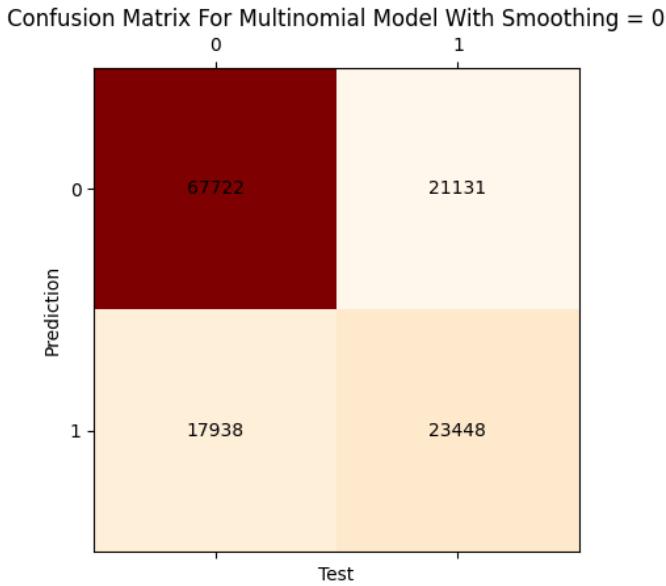


Fig.6: Confusion Matrix for the Final Multinomial Model

----- Multinomial Model -----
The number of correct predictions: 91170, wrong predictions: 39069,

accuracy: 0.7000207311174073

The number of true positives: 23448, true negatives: 67722

The number of false positives: 17938, false negatives: 21131

f1 score for multinomial model with smoothing = 0: 0.5455243413016926

7.2.2. Bernoulli Naive Bayes Model

In the second phase of the report, we added some new features related to the count of specific characters such as '+', '#', '/', '\$', '!', '*', ',', '_', ';' and so on, which were not directly related to the

structure of the URL. When we added all these features, for the first time, we observed a decrease in the quality of the Bernoulli Model and therefore decided to pick a best-performing subset of these features. In order to test the usefulness of these features, we decided to apply the forward selection method to extract the best set of features on the Bernoulli Model. The Bernoulli Model was picked since these new features mostly take binary values, 0 and 1, making them suitable with Bernoulli Model, and also, there was no parameter to tune in the Bernoulli Model. Besides that, the Bernoulli Model is quite accurate and computationally efficient; the same experiment would take much longer runtimes with Logistic Regression and Neural Network Models. With feature selection, the accuracies have been calculated on the validation dataset, and in each iteration, the character that improves the accuracy the most has been picked.

```
----- Bernoulli Model -----
Base Acc: 0.8632810700410019

----- Iteration 1 -----
New highest acc: 0.8642024601114882
New highest acc: 0.8691779664921144
Selected Char: /

----- Iteration 2 -----
New highest acc: 0.8699765045532026
Selected Char: +

----- Iteration 3 -----
New highest acc: 0.8699918610543774
New highest acc: 0.8701454260661251
Selected Char: ,

----- Iteration 4 -----
New highest acc: 0.8701607825672999
Selected Char: !

----- Iteration 5 -----
Selected Set Of Features:[ '/', '+', ',', '!' ]
```

The above selected features have been then added to the train, test, and validation datasets for all machine learning algorithms tested in the project. There is no parameter that we can tune for our Bernoulli Naive Bayes; therefore, we directly trained it on our train and validation and evaluated it on our test set. Below is the confusion matrix and the results of some other metrics:

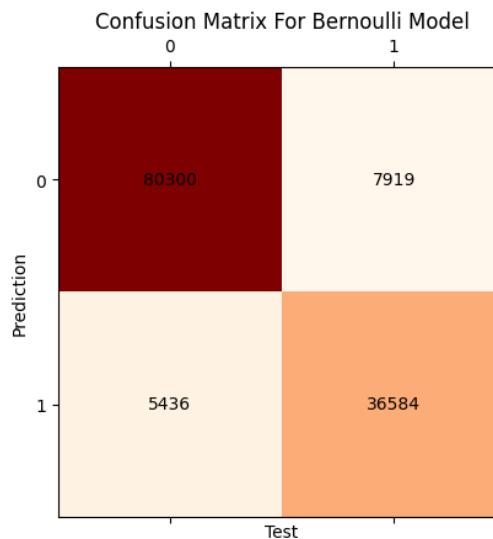


Fig.7: Confusion Matrix for the Final Bernoulli Model

```
----- Bernoulli Model -----
The number of correct predictions: 90257, wrong predictions: 39982,
accuracy: 0.8974577507505432
```

```
The number of true positives: 36584, true negatives: 80300
The number of false positives: 5436, false negatives: 7919
f1 score for bernoulli model: 0.8456479779942905
```

On the machine equipped with Intel i7-10750h, the training of the Bernoulli model took around 0.228995 seconds (Python's *time* library gave this result); still the calculations were quite fast but slower than the Multinomial case.

7.3 Neural Network Model

The neural networks that we tested in the project were fully connected ANN (Artificial Neural Network) models, and we tuned these models to obtain their best performance on our dataset. As our loss function in training, the Binary Cross Entropy loss function has been used as it's appropriate to the task of the project, which is binary classification. Besides, the sigmoid activation function has been used in the output layer in order to map the output of the model to a probability between 0 and 1, and batch gradient descent has been used to update the model. Together with this, a threshold of 0.5 has been chosen for the final classification. Also, in the training process, 500 epochs have been used together with a stopping criterion, which exits training if more than five epochs (patience=5) pass without an increase in the accuracy of less than 1e-7.

When adjusting the necessary hyperparameters to tune the model, there were four important ones, which are the number of neurons in the layers, batch size, momentum, and the learning rate. We also had the option to choose between ReLU and Sigmoid activation functions, but since our dataset is large, the performance of the training was also an important criterion for us. As a result of this, we chose ReLU as our activation function since it's fast to compute ReLU and its derivative. Also, usually, the Sigmoid function saturates and kills gradients, leading to a higher number of training epochs for the model to converge, which makes it an unfavorable option. The number of neurons that we tested were [64,1], [32, 32, 1], [64, 64, 1], where the numbers in the list indicate the number of neurons in the layers of our network. For the batch size, the values in the set {32, 64} have been tested, and for momentum and learning rate, the values from the sets {0.85, 0.95} and {0.01, 0.005, 0.001, 0.0005} have been tested, respectively. The best network has been picked based on the validation accuracy by training each configuration up to 200 epochs (maximum). Also, for the final best neural network model, both accuracy and F-1 score metrics have been reported in order to make a fair comparison with other selected machine learning algorithms. Based on the validation accuracy, the following histogram has been formed showing the distribution of the results for 48 different configurations of the network:

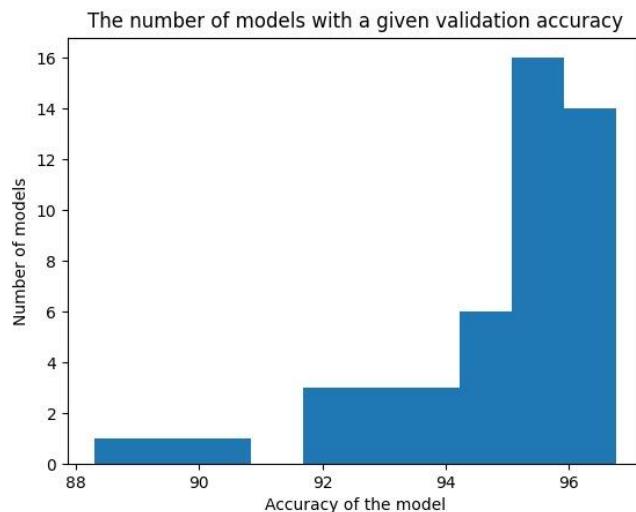


Fig.8: Histogram Showing The Distribution Of The Model Validation Accuracies For Different Configurations Of Hyperparameters

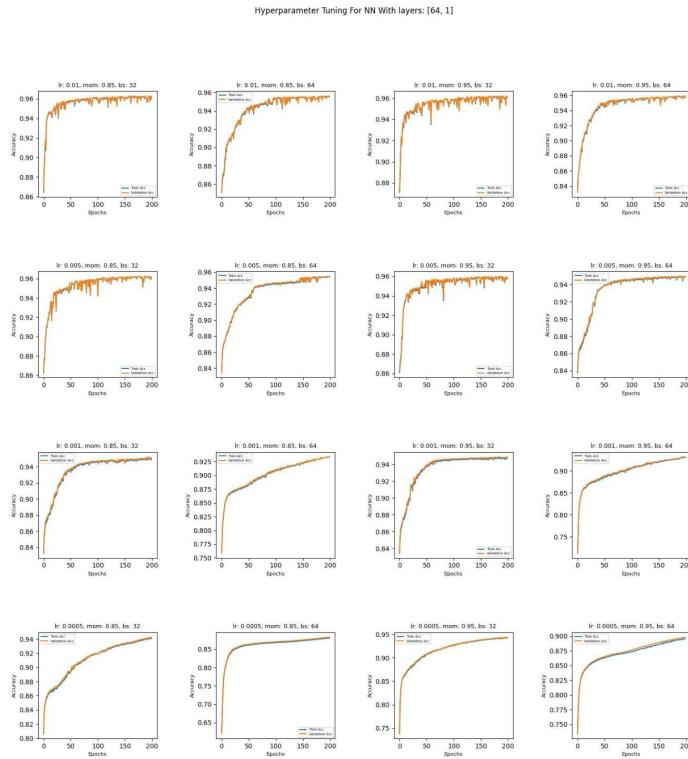


Fig.9: Plot Showing The Train and Validation Accuracies For Different Parameter Settings For The Model Size [64,1]

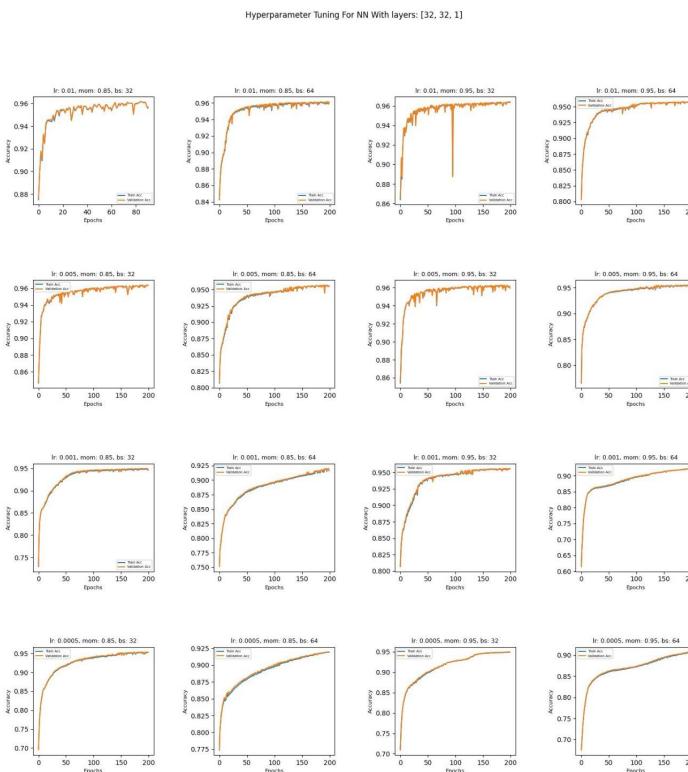


Fig.10: Plot Showing The Train and Validation Accuracies For Different Parameter Settings For The Model Size [32, 32, 1]

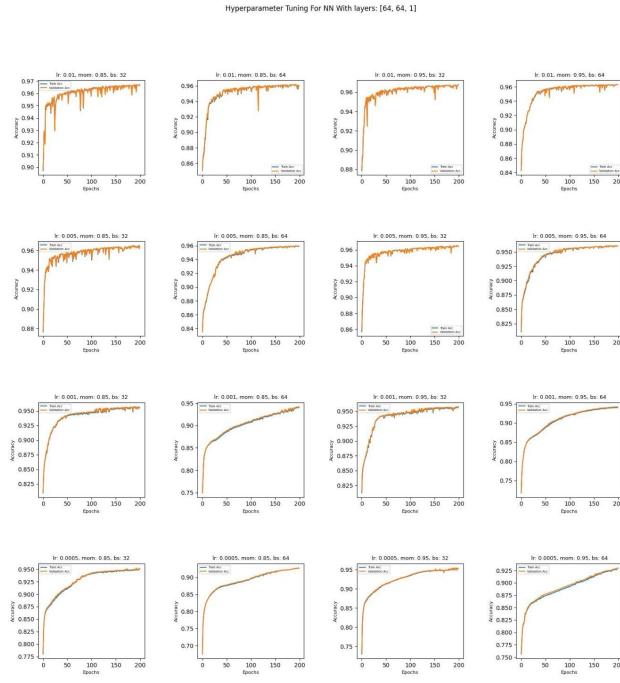


Fig.11: Plot Showing The Train and Validation Accuracies For Different Parameter Settings For The Model Size [64, 64, 1]

On the machine equipped with Intel i7-10750h, the training took around 6 and 3.5 seconds for each epoch when batch sizes are 32 and 64, respectively. The plots above illustrate the accuracies corresponding to various hyperparameter configurations for each model size. It can be inferred that the neural network's performance exhibits greater variability with different hyperparameter settings compared to the other machine learning models we investigated. Based on the validation accuracy, the best-performing network has been observed under the following configuration, which has been then trained again on the train + validation dataset and tested on the test dataset:

Best Neural Network Settings: Size=[64, 64, 1], Learning Rate=0.01, Batch Size:32, Momentum:0.95

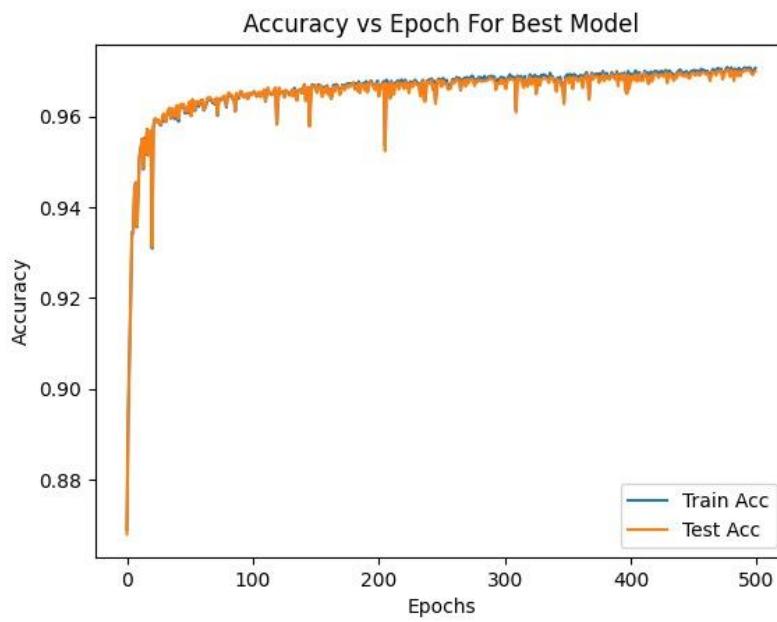


Fig.12: Plot Showing The Train and Validation Accuracies For Different Parameter Settings For The Model Size [64, 64, 1]

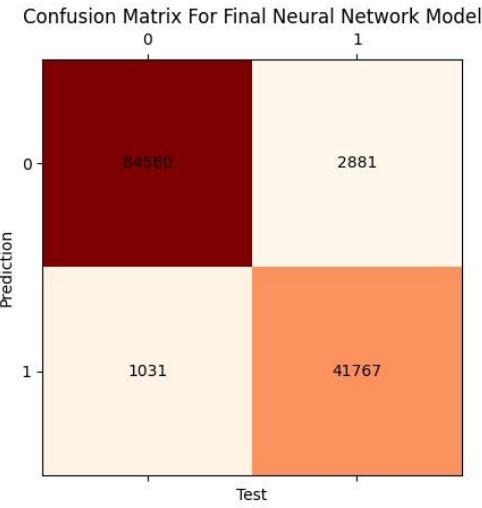


Fig.13: Plot Showing The Train and Validation Accuracies For Different Parameter Settings For The Model Size [64, 64, 1]

Final Test Accuracy: 0.969963

F1 Score For Final Neural Network Model: 0.9552638199574595

From the above plots, we observe that our best model does not overfit to the data after 500 epochs of training and also the train and test accuracies tend to go together though test accuracies are slightly lower. There are a few reasons behind this interesting fact that our model behaves similarly on test and train datasets, which is also the case for different models when tuning the parameters. With a large dataset like the URL dataset (around 600,000 data points), the model has access to a diverse range of examples. This ample data allows the model to learn better representations and generalize well to unseen data. Also, if the train and test datasets are sampled from the same distribution and are representative of the overall dataset, it is likely that the model learns patterns and generalizes well to both sets. Another thing that contributes is the dataset being balanced for two classes which enables the model to learn from a representative number of examples from each class, resulting in better generalization. Also, the complexity of the models tested was proper when the number of data points is considered. Consequently, the train and test accuracies tend to be similar as the model captures the underlying patterns effectively.

The results obtained from the experimentation phase demonstrated that the Neural Network model yielded the highest accuracy of 0.97 and an F-1 score of 0.955, making it the most effective model for URL classification. This outcome highlights the power of Neural Networks in capturing intricate patterns and relationships within the data, enabling them to make accurate predictions.

7.4 Discussion Of The Obtained Results

As expected, we obtained the highest accuracy and highest F-1 scores from the Neural Network, followed by the Logistic Regression Model. Though the Bernoulli model was really close to Logistic Regression in terms of classification accuracy, the Multinomial Model was clearly worse than the Bernoulli Model. The reason behind this might be some of the features that play an important role in explaining the variance, which take values of only 0 and 1, such as *use_of_ip*, *count_http*, *count_https*, and *count_www*, and others. Since Bernoulli only focuses on the absence or presence of the features, these kinds of binary features that we have might have benefited the score of the Bernoulli Model, but at the same time, since the probabilities calculated are relative in the Multinomial Model, the effects of these features might have vanished. This usage of binary features is also not a problem in the Logistic Regression Model, as the model actually learns the ‘importance’ of each feature by estimating the corresponding coefficient for it. In the final report, the classification accuracies for both Bernoulli and Logistic Regression Models have been improved with the extraction of new features. Finally, because of the learning flexibility of neural networks for complex functions, we obtained a higher classification score from the neural networks, yet they required a great effort for

optimization in order to achieve this. Interestingly, our model also did not overfit the data since the large and diverse dataset provides the model with enough examples to learn from, and the models picked were not overly complex relative to the complexity of the problem, mitigating overfitting.

8. Encountered Challenges

We faced a number of challenges during the development of this project. The first and most obvious one was when we were in the early pre-processing stages. Since the original dataset we had did not have any features itself, we had to do plenty of research about URLs in order to come up with features that could actually be useful for our model and in order to achieve this. Once we were done with preprocessing, we started implementing our Logistic Regression model, which we had difficulties implementing from scratch without using predefined methods of some libraries such as scikit-learn. And for the future, we might have to do a more detailed analysis of the features that we use to train our Logistic Regression model, as some of these features may not be as relevant as they are for the Naive Bayes model and do not contribute much to the performance of our Logistic Regression model while increasing its complexity. For our Naive Bayes models, their implementations were comparably more straightforward than Logistic Regression, and we did not face major difficulties.

In our final version, when implementing the Neural Network method, we encountered more challenges as it is known to be more complex, and it's somewhat new for us to implement it from scratch. Besides that, the optimization of the network was compelling as there were many settings for us to tune, and also the training times were another hurdle. Hence, we needed to make some decisions and tune some of the parameters by hand, as the time wouldn't permit experimenting with each possible configuration. Another challenge was to come up with some additional features that could be useful for the training of our models. In order to find relevant features that could be used to determine whether a URL is malicious or not, we did some exhaustive research online and performed some experiments ourselves to see if the newly added features improved the performance of our models or not. As a result of this, we recognized that some of the features harmed the performance of the models, so we had to tackle this problem by selecting a subset of them with a forward selection process.

9. Gantt Chart

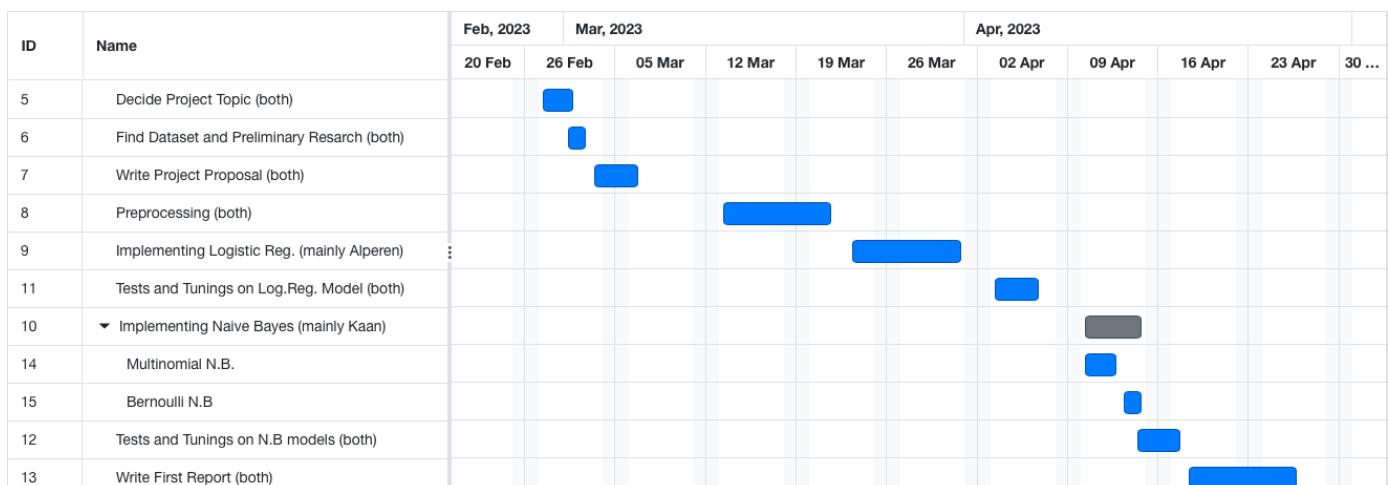


Fig.14: Gantt Chart Showing The Progress Plan Of The URL Classification Project Until The First Report

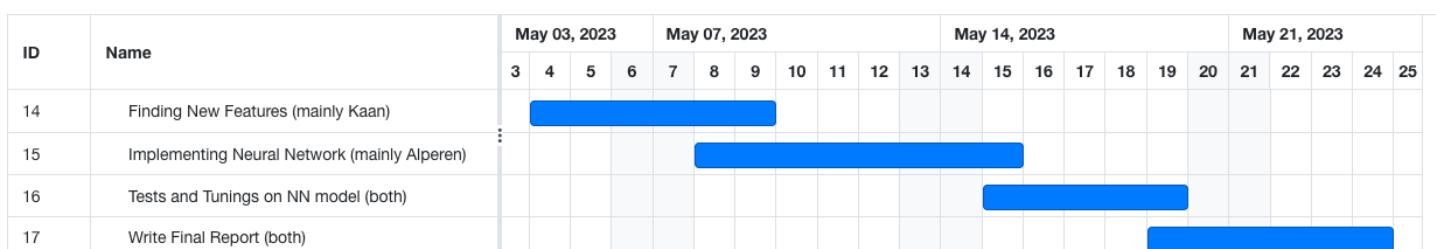


Fig.15: Gantt Chart Showing The Progress Plan Of The URL Classification Project Between First and Final Reports

10. Conclusion

The objective of this project was to classify URLs as either malicious or benign by extracting relevant features and utilizing various machine-learning models. This endeavor aimed to enhance the detection of URL-based threats. Extensive research and literature review were conducted to extract meaningful features from URL strings. The training and testing processes yielded significant results, achieving high accuracy and F-1 scores exceeding 0.9 in correctly classifying URLs. These outcomes underscore the efficacy of Neural Networks in capturing intricate patterns and relationships within the data. Additionally, the probabilistic models and Logistic Regression model demonstrated notable success. Nonetheless, further improvements are required to optimize these models for the URL dataset. Furthermore, employing comprehensive feature selection methods such as forward selection or backward elimination for each machine learning algorithm can help identify the most suitable feature sets. In cases where training times are considerable, such as with Neural Networks, employing statistical feature selection methods may prove advantageous. Advancements can be pursued through the exploration of ensemble methods that harness the strengths of different models. Additionally, incorporating more advanced feature extraction techniques, guided by ongoing research targeting emerging malicious URL techniques, holds promise for further advancements in this field.

11. References

- [1] Siddhartha, “Malicious urls dataset,” *Kaggle*, 23-Jul-2021. [Online]. Available: <https://www.kaggle.com/datasets/sid321axn/malicious-urls-dataset>. [Accessed: 26-Apr-2023].
- [2] “URL,” *Wikipedia*, 26-Apr-2023. [Online]. Available: <https://en.wikipedia.org/wiki/URL#Syntax>. [Accessed: 26-Apr-2023].
- [3] sid321axn, “Malicious URL detection using ML & Feat Engg,” *Kaggle*, 23-Jul-2021. [Online]. Available: <https://www.kaggle.com/code/sid321axn/malicious-url-detection-using-ml-feat-engg/notebook>. [Accessed: 26-Apr-2023].
- [4] *IPv4 and IPv6 address formats*. [Online]. Available: <https://www.ibm.com/docs/en/ts3500-tape-library?topic=functionality-ipv4-ipv6-address-formats>. [Accessed: 26-Apr-2023].
- [5] “TLD,” *PyPI*. [Online]. Available: <https://pypi.org/project/tld/>. [Accessed: 26-Apr-2023].
- [6] R. Ramakrishnan, “URL feature engineering and Classification,” *Medium*, 25-May-2021. [Online]. Available: <https://medium.com/nerd-for-tech/url-feature-engineering-and-classification-66c0512fb34d>. [Accessed: 26-Apr-2023].
- [7] 738, “738/awesome-URL-shortener: A curated list of awesome URL shortener,” *GitHub*. [Online]. Available: <https://github.com/738/awesome-url-shortener>. [Accessed: 26-Apr-2023].
- [8] R. E. Ikwu, “Extracting feature vectors from URL strings for malicious URL detection,” *Medium*. [Online]. Available: <https://towardsdatascience.com/extracting-feature-vectors-from-url-strings-for-malicious-url-detection-cbafc24737a> [Accessed: 23-May-2023].

12. Appendix

utils.py

```
utils.py > ...
1 import re
2 import pandas as pd
3
4 from tld import get_tld
5 from typing import Tuple, Union
6 from urllib.parse import urlparse
7
8 def contains_ip_address(url: str) -> bool:
9     contains_ip = re.search(
10         ('([01]?\\d{1,2}\\d{1,2}[0-4]\\d{1,2}[0-5])\\.([01]?\\d{1,2}\\d{1,2}[0-4]\\d{1,2}[0-5])\\.([01]?\\d{1,2}\\d{1,2}[0-4]\\d{1,2}[0-5])\\.\\.' +
11         '([01]?\\d{1,2}\\d{1,2}[0-4]\\d{1,2}[0-5])\\/.+' # IPv4
12         '|([01]?\\d{1,2}\\d{1,2}[0-4]\\d{1,2}[0-5])\\.([01]?\\d{1,2}\\d{1,2}[0-4]\\d{1,2}[0-5])\\.([01]?\\d{1,2}\\d{1,2}[0-4]\\d{1,2}[0-5])\\.\\.' +
13         '([01]?\\d{1,2}\\d{1,2}[0-4]\\d{1,2}[0-5])\\/.+' # IPv4 with port
14         '|([01]?\\d{1,2}\\d{1,2}[0-4]\\d{1,2}[0-5])\\.([01]?\\d{1,2}\\d{1,2}[0-4]\\d{1,2}[0-5])\\.([01]?\\d{1,2}\\d{1,2}[0-4]\\d{1,2}[0-5])\\.\\.' +
15         '|([01]?\\d{1,2}\\d{1,2}[0-4]\\d{1,2}[0-5])\\.([0x[0-9a-fA-F]{1,2}]{1,2})\\.([0x[0-9a-fA-F]{1,2}]{1,2})\\/.+' # IPv4 in hexadecimal
16         '|([0-9a-fA-F]{1,4})\\{1,2}|([0-9a-fA-F]{1,4})\\{1,2}|' +
17         '|([0-9]+\\{1,2}|([0-9]+)\\{3}|([0-9]+)\\{1,4})|' +
18         '|((?:\\d|[01]?\\d{1,2}\\d{1,2}[0-4]\\d{1,2}[0-5])\\.\\{3}(?:\\d{1,2}\\d{1,2}[0-4]\\d{1,2}[0-5])\\{1,2}|\\d{1,2}\\d{1,2}[0-4]\\d{1,2}[0-5])\\{1,2}|', url) # IPv6
19     if contains_ip:
20         return 1
21     return 0
22
23 """ returns full length domain (fld), top level domain (tld), domain, and
24 subdomain names for a given URL. These can be useful as features """
25 def extract_tld(url: str, fix_protos: bool = False) -> Tuple[str, str, str, str]:
26     res = get_tld(url, as_object = True, fail_silently=True, fix_protocol=fix_protos)
27
28     fld = res.fld
29     tld = res.tld
30     domain = res.domain
31     subdomain = res.subdomain
32
33     return subdomain, domain, tld, fld
34
35 """ takes an instance and checks if the url is an IP address. If yes,
36 set four instances to None. If not, call extract_tld to process URL"""
37 def process_url(data: pd.Series) -> Tuple[str, str, str, str]:
38     try:
39         if data['use_of_ip'] == 0:
40             if str(data['url']).startswith('http:'):
41                 return extract_tld(data['url']), fix_protos=True)
```

```
utils.py > ...
1    def _get_tld(url: str) -> str:
2        parsed_url = parse.urlparse(url)
3        if parsed_url.netloc:
4            return parsed_url.netloc
5        else:
6            return parsed_url.path
7
8    def get_path(url: str) -> Union[str, None]:
9        try:
10            res = _get_tld(url, as_object=True, fail_silently=True, fix_protocol=True)
11            if res.parsed_url.query:
12                joined = res.parsed_url.path + res.parsed_url.query
13                return joined
14            else:
15                return res.parsed_url.path
16        except:
17            return None
18
19    def letter_count(url: str) -> int:
20        letter_count = 0
21        for c in url:
22            if c.isalpha():
23                letter_count += 1
24        return letter_count
25
26    def digit_count(url: str) -> int:
27        d_count = 0
28        for d in url:
29            if d.isnumeric():
30                d_count += 1
31        return d_count
32
33    def sub_directory_count(url: Union[str, None]) -> int:
34        if url:
35            dir_count = url.count('/')
36            return dir_count
37
38    def _extract_ccl_data(url: str, fix_protocol: bool):
39        try:
40            parsed_url = parse.urlparse(url)
41            if parsed_url.netloc:
42                return None, None, None, None
43            else:
44                return None, None, None, None
45
46    """ returns the path of the url
47    for example: www.bilkent.com/srs/cgpa -> srs/cgpa """
48    def get_path(url: str) -> Union[str, None]:
49
50        try:
51            res = _get_tld(url, as_object=True, fail_silently=True, fix_protocol=True)
52            if res.parsed_url.query:
53                joined = res.parsed_url.path + res.parsed_url.query
54                return joined
55            else:
56                return res.parsed_url.path
57        except:
58            return None
59
60    """ returns letter count """
61    def letter_count(url: str) -> int:
62
63        letter_count = 0
64        for c in url:
65            if c.isalpha():
66                letter_count += 1
67        return letter_count
68
69    """ return digits count """
70    def digit_count(url: str) -> int:
71
72        d_count = 0
73        for d in url:
74            if d.isnumeric():
75                d_count += 1
76        return d_count
77
78    """ the number of '/' characters gives the number of subdirectories """
79    def sub_directory_count(url: Union[str, None]) -> int:
80        if url:
81            dir_count = url.count('/')
82            return dir_count
```

```

utils.py > ...
82     |     return dir_count
83     | return 0
84
85     """ return length of the first directory """
86     def len_first_directory(url: Union[str, None]):
87         if url:
88             if len(url.split('/')) > 1:
89                 first_dir_len = len(url.split('/')[1])
90             return first_dir_len
91         return 0
92
93     """ check whether the URL has a shortening service or not
94     Shortening Service: a third-party website that converts that long URL to a short,
95     case-sensitive alphanumeric code."""
96     def check_shortening_service(url: str) -> int:
97
98         check = re.search('bit.ly|goo.gl|shorte.st|go2l.in|co.ow.ly|t.co|tinyurl|tr.im|is.gd|cli.gs|'
99                           '|yfrog.com|migre.me|ff.im|tiny.cc|ur4.eu|twit.ac|su.pr|twurl.nl|snipurl.com|'
100                          '|short.to|BudURL.com|ping.fm|post.ly|Just.as|bkite.com|snipr.com|flic.kr|loop.it.us|'
101                          '|doip.com|short.ie|k.am|wp.me|rubyurl.com|om.ly|to.ly|bit.co|lnkd.in|'
102                          '|db.tt|qr.ae|adf.ly|goo.gl|bitly.com|curl.tinyurl.com|ow.ly|bit.ly|ity.im|'
103                          '|q.gs|is.gd|po.st|bc.vc|twitthis.com|u.to|j.mp|bzurl.com|cutt.us|u.bb|yourls.org|'
104                          '|x.co|prettylinkpro.com|scrnch.me|filoops.info|vturl.com|qr.net|url.com|tweez.me|v.gd|'
105                          '|tr.im|link.zip.net., url)
106
107         return 1 if check else 0
108
109     """ check whether the URL starts with 'http' or 'https' """
110     def check_http(url):
111         htp = urlparse(url).scheme
112         check = str(htp)
113
114         return 1 if check == 'https' else 0
115
116     def suspicious_words(url):
117         match = re.search('client|admin|server|cash|PayPal|login|signin|bank|account|update|free|lucky|service|bonus|ebayisapi|sex|webscr', url)
118
119         return 1 if match else 0

```

preprocessing.py

```

preprocessing.py > ...
1  import pandas as pd
2  from utils import contains_ip_address, process_url, get_path, check_shortening_service, sub_directory_count, letter_count, digit_count, len_first_directory, check_http,
3  from sklearn.preprocessing import OrdinalEncoder
4
5  df = pd.read_csv("data/malicious_phish.csv")
6
7  print(df.head())
8
9  # create a column named is_ip by using the ip check function
10 df['use_of_ip'] = df['url'].apply(lambda i: contains_ip_address(i))
11
12 # get the number of URLs containing IP address
13 print(df['use_of_ip'].value_counts())
14
15 # an example url starting with the ip address
16 print(df[df['use_of_ip'] == 1].iloc[0]['url'])
17
18 print(df['url'][1])
19
20 # add new features
21 df[['subdomain', 'domain', 'tld', 'fld']] = df.apply(lambda x: process_url(x), axis=1, result_type="expand")
22
23 # General Features
24 df['url_path'] = df['url'].apply(lambda x: get_path(x))
25 df['use_of_shortener'] = df['url'].apply(lambda x: check_shortening_service(x))
26
27 # URL component length
28 df['url_length'] = df['url'].apply(lambda x: len(str(x)))
29 df['subdomain_length'] = df['subdomain'].apply(lambda x: len(str(x)))
30 df['tld_length'] = df['tld'].apply(lambda x: len(str(x)))
31 df['fld_length'] = df['fld'].apply(lambda x: len(str(x)))
32 df['path_length'] = df['url_path'].apply(lambda x: len(str(x)))
33
34 # count the number of alphanumeric characters, digits, and punctuations
35 df['count_letters'] = df['url'].apply(lambda i: letter_count(i))
36 df['count_digits'] = df['url'].apply(lambda i: digit_count(i))
37 df['count_puncs'] = (df['url_length'] - (df['count_letters'] + df['count_digits']))
38
39 df['is_php'] = df['url'].apply(lambda a: a.count('.php'))
40 df['count_%20'] = df['url'].apply(lambda a: a.count('%20')) # usually effects the model in a bad way
41

```

```

❷ preprocessing.py > ...
42     # check special character counts for the url
43     for c in '@%?=':
44         df['count' + c] = df['url'].apply(lambda a: a.count(c))
45         # print(df['count'+c][:5])
46
47     # extra features that are not directly related
48     # with the structure of a URL
49     featureList = ['+', '#', '/', '$', '!', '*', '.', '_', ':']
50     for c in featureList:
51         df['count' + c] = df['url'].apply(lambda a: a.count(c))
52
53     df['is_susp'] = df['url'].apply(lambda i: suspicious_words(i))
54
55     # Binary Label by converting benign to 0 and all other classes to 1
56     df['is_malicious'] = df['type'].apply(lambda x: 0 if x == 'benign' else 1)
57
58     # Inspect the malicious URLs
59     print([df[df['is_malicious']] == 1].head())
60
61     # the number of subdirectories (number of /) and the length can be helpful
62     df['count_dirs'] = df['url_path'].apply(lambda x: sub_directory_count(x))
63     df['first_dir_length'] = df['url_path'].apply(lambda x: len_first_directory(x))
64
65     # Binned Features
66     groups = ['Short', 'Medium', 'Long', 'Very Long']
67     # URL Lengths in 4 bins
68     df['url_length_q'] = pd.qcut(df['url_length'], q=4, labels=groups)
69     # FLD Lengths in 4 bins
70     df['fld_length_q'] = pd.qcut(df['fld_length'], q=4, labels=groups)
71
72     # Counts for https, http, www
73     df['https'] = df['url'].apply(lambda i: check_https(i))
74     df['count-https'] = df['url'].apply(lambda a: a.count('https'))
75     df['count-http'] = df['url'].apply(lambda a: a.count('http'))
76     df['count-www'] = df['url'].apply(lambda a: a.count('www'))
77
78     # Percentage Features
79     df['letters_ratio'] = df['count_letters'] / df['url_length']
80     df['digit_ratio'] = df['count_digits'] / df['url_length']
81     df['punc_ratio'] = df['count_puncs'] / df['url_length']
82

```

```

❷ preprocessing.py > ...
82
83     enc = OrdinalEncoder()
84     df[['url_length_q','fld_length_q']] = enc.fit_transform(df[['url_length_q','fld_length_q']])
85
86     print(df.head())
87
88     # write back the file
89     df.to_csv('data/url_processed.csv', index=False)

```

naive_bayes.py

```

❷ naive_bayes.py > ...
1  import numpy as np
2  import pandas as pd
3  import matplotlib.pyplot as plt
4  import time
5
6  df = pd.read_csv("data/url_processed.csv")
7
8  # shuffle the data frame
9  df = df.sample(frac = 1)
10
11 # pick only the numeric features
12 X = df[['use_of_ip', 'url_length', 'subdomain_length', 'tld_length', 'fld_length', 'path_length',
13          'count_letters', 'count_digits', 'count_puncs', 'count.', 'count@', 'count-',
14          'count%', 'count?', 'count=', 'count_dirs', 'use_of_shortener', 'first_dir_length',
15          'url_length_q', 'fld_length_q', 'https', 'count-https', 'count-http', 'count-www', 'is_susp', 'is_php']]
16
17 print(X.head())
18
19 y = df['is_malicious']
20
21 print(X.isna().sum())
22
23 # fill the nan values with 0
24 X = X.fillna(0)
25
26 # check again the number of nan values
27 print(X.isna().sum())
28
29 X.replace([np.inf, -np.inf], 0, inplace=True)
30
31 X['first_dir_length'] = X['first_dir_length'].astype('int64')
32 X['url_length_q'] = X['url_length_q'].astype('int64')
33 X['fld_length_q'] = X['fld_length_q'].astype('int64')
34
35 # divide the data into train, test and validation datasets
36 n = len(y)
37
38 train_size = int(n * 0.7)
39 valid_size = int(n * 0.1)
40 test_size = int(n * 0.2)
41

```

```

❷ naive_bayes.py > ...
42     X_train = X.iloc[:train_size]
43     y_train = y.iloc[:train_size]
44
45     X_valid = X.iloc[train_size: train_size + valid_size]
46     y_valid = y.iloc[train_size: train_size + valid_size]
47
48     X_test = X.iloc[train_size + valid_size:]
49     y_test = y.iloc[train_size + valid_size:]
50
51     print(X_train.shape, X_valid.shape, X_test.shape)
52
53     print(X_train.head())
54
55     # convert df into numpy
56     X_train = X_train.to_numpy()
57     X_valid = X_valid.to_numpy()
58     X_test = X_test.to_numpy()
59     y_train = y_train.to_numpy()
60     y_test = y_test.to_numpy()
61     y_valid = y_valid.to_numpy()
62
63     # calculate the spam and normal class probabilities
64     spam_count = np.count_nonzero(y_train)
65     normal_count = y_train.shape[0] - spam_count
66     p_spam = spam_count / y_train.shape[0]
67     p_normal = 1 - p_spam
68
69     # print the spam percentage
70     print("Spam percentage: " + str((p_spam)*100))
71
72     # add the number of words for each email
73     word_counts = np.sum(X_train, axis = 1)
74
75     # Get the total word counts for categories
76     spamWordCount = word_counts.T @ y_train
77     normalWordCount = word_counts.T @ (1 - y_train)
78
79     start_time = time.time()
80     # Get the frequencies for each word seperately
81     spamFrequencies = y_train.T @ X_train
82     normalFrequencies = (1 - y_train.T) @ X_train

```

```

❷ naive_bayes.py > ...
82     normalFrequencies = (1 - y_train.T) @ X_train
83
84     # Calculate the parameters by dividing with the total word count
85     spamProbabilities = spamFrequencies / spamWordCount
86     normalProbabilities = normalFrequencies / normalWordCount
87
88     # Take the log of the probabilities
89     with np.errstate(divide='ignore'):
90         logSpamProbabilities = np.log(spamProbabilities)
91         logNormalProbabilities = np.log(normalProbabilities)
92
93     logSpamProbabilities[np.isneginf(logSpamProbabilities)] = -1e+12
94     logNormalProbabilities[np.isneginf(logNormalProbabilities)] = -1e+12
95     print("---- % seconds ----" % (time.time() - start_time))
96
97     multinomial_acc_history = []
98     def fit_smoothened_multinomial(alpha, spamFrequencies, normalFrequencies, X_train, X_test, y_test):
99         regularizedSpamFrequencies = spamFrequencies + alpha
100        regularizedNormalFrequencies = normalFrequencies + alpha
101
102        # regularize the word counts for spam and normal categories, as well
103        regularizedSpamWordCount = spamWordCount + alpha * X_train.shape[1]
104        regularizedNormalWordCount = normalWordCount + alpha * X_train.shape[1]
105
106        # divide by the total word count to find the probabilities
107        regularizedSpamProbabilities = regularizedSpamFrequencies / regularizedSpamWordCount
108        regularizedNormalProbabilities = regularizedNormalFrequencies / regularizedNormalWordCount
109
110        # take the logarithm of the probabilities
111        regularizedLogSpamProbabilities = np.log(regularizedSpamProbabilities)
112        regularizedLogNormalProbabilities = np.log(regularizedNormalProbabilities)
113
114        num_correct = tp = tn = fp = fn = 0
115        for i in range(X_test.shape[0]):
116            row = X_test[i]
117            probSpam = np.log(p_spam)
118            probNormal = np.log(p_normal)
119
120            probSpam += (regularizedLogSpamProbabilities @ row)
121            probNormal += (regularizedLogNormalProbabilities @ row)

```

```

☆ naive_bayes.py > ...
123     predicted = 0
124     if probSpam > probNormal:
125         predicted = 1
126
127     if predicted == y_test[i]:
128         num_correct += 1
129         if predicted == 1:
130             tp += 1
131         else:
132             tn += 1
133     else:
134         if predicted == 1:
135             fp += 1
136         else:
137             fn += 1
138
139     print("----- Multinomial Model With Additive Smoothing alpha = %d -----" % (alpha))
140     print("The number of correct predictions: " + str(num_correct) + ", wrong predictions: " + str(x_test.shape[0] - num_correct) + ", accuracy: " + str(num_correct / x
141     print("The number of true positives: " + str(tp) + ", true negatives: " + str(tn))
142     print("The number of false positives: " + str(fp) + ", false negatives: " + str(fn))
143     multinomial_acc_history.append(num_correct / x_test.shape[0]) # add the accuracy to the history
144
145 # apply smoothing and compare on the validation set to pick the best smoothing
146 for i in range(0, 10, 2):
147     fit_smoothened_multinomial(i, spamFrequencies, normalFrequencies, x_train, x_valid, y_valid)
148
149 # plotting the validation accuracy for different smoothing parameter values
150 x = list(range(0, 10, 2))
151 plt.figure(figsize=(18, 12))
152 plt.title('Smoothing Parameter Comparison')
153 plt.xlabel('Smoothing Parameter Alpha')
154 plt.ylabel('Validation accuracy')
155 plt.plot(x, multinomial_acc_history, '-o', label='Multinomial Model')
156 plt.yticks([x / 5 for x in range(0, 6)])
157 for (parameter, acc) in zip(x, multinomial_acc_history):
158     plt.text(parameter, acc, "{:.6f}".format(acc), va='bottom', ha='center')
159 plt.legend()
160 plt.title('Validation Accuracy For Multinomial Model')
161 plt.show()
162
163 # The best behaving multinomial model is one with smoothing = 0, report its accuracy again

```

```

☆ naive_bayes.py > ...
164 # by training on the train + validation set
165 x_train = np.concatenate((x_train, x_valid), axis=0)
166 y_train = np.concatenate((y_train, y_valid), axis=0)
167
168 # calculate the spam and normal class probabilities
169 spam_count = np.count_nonzero(y_train)
170 normal_count = y_train.shape[0] - spam_count
171 p_spam = spam_count / y_train.shape[0]
172 p_normal = 1 - p_spam
173
174 # print the spam percentage
175 print("Spam percentage: " + str((p_spam)*100))
176
177 # add the number of words for each email
178 word_counts = np.sum(x_train, axis = 1)
179
180 # Get the total word counts for categories
181 spamWordCount = word_counts.T @ y_train
182 normalWordCount = word_counts.T @ (1 - y_train)
183
184 # Get the frequencies for each word seperately
185 spamFrequencies = y_train.T @ x_train
186 normalFrequencies = (1 - y_train.T) @ x_train
187
188 # Calculate the parameters by dividing with the total word count
189 spamProbabilities = spamFrequencies / spamWordCount
190 normalProbabilities = normalFrequencies / normalWordCount
191
192 # Take the log of the probabilities
193 with np.errstate(divide='ignore'):
194     logSpamProbabilities = np.log(spamProbabilities)
195     logNormalProbabilities = np.log(normalProbabilities)
196
197 logSpamProbabilities[np.isneginf(logSpamProbabilities)] = -1e+12
198 logNormalProbabilities[np.isneginf(logNormalProbabilities)] = -1e+12
199
200 num_correct = tp = tn = fp = fn = 0
201 y_pred = []
202 for i in range(x_test.shape[0]):
203     row = x_test[i]
204     probSpam = np.log(p_spam)

```

```

❸ naive_bayes.py > ...
205     probNormal = np.log(p_normal)
206
207     probSpam += (logSpamProbabilities @ row)
208     probNormal += (logNormalProbabilities @ row)
209
210     predicted = 0
211     if probSpam > probNormal:
212         predicted = 1
213
214     y_pred.append(predicted)
215     if predicted == y_test[i]:
216         num_correct += 1
217         if predicted == 1:
218             tp += 1
219         else:
220             tn += 1
221     else:
222         if predicted == 1:
223             fp += 1
224         else:
225             fn += 1
226
227 print("----- Multinomial Model -----")
228 print("The number of correct predictions: " + str(num_correct) + ", wrong predictions: " + str(x_test.shape[0] - num_correct) + ", accuracy: " + str(num_correct / x_test.shape[0]))
229 print("The number of true positives: " + str(tp) + ", true negatives: " + str(tn))
230 print("The number of false positives: " + str(fp) + ", false negatives: " + str(fn))
231
232 # plot confusion matrix
233 confusion = pd.crosstab(y_test, y_pred)
234 fig, ax = plt.subplots()
235 ax.matshow(confusion, cmap='OrRd')
236 ax.set(xlabel='Test', ylabel='Prediction')
237
238 for i in range(2):
239     for j in range(2):
240         c = confusion[i][j]
241         ax.text(i, j, str(c), va='center', ha='center')
242 plt.title('Confusion Matrix For Multinomial Model With Smoothing = 0')
243 plt.show()
244
245 # print f1 score = 2 * precision * recall / (precision + recall)
246 precision = tn / (fp + tn)

```

```

❸ naive_bayes.py > ...
245     precision = tp / (fp + tp)
246     recall = tp / (tp + fn)
247     f1_score = (2 * precision * recall) / (precision + recall)
248     print("f1 score for multinomial model with smoothing = 0: ", f1_score)
249
250 """
251 Bernoulli Model is trained on train + validation and its accuracy is reported
252 on the test dataset since it doesn't have any parameters to tune
253 """
254 start_time = time.time()
255 def fit_bernoulli(x_train, x_test, y_train, y_test):
256     # create a copy of the train and test datasets
257     bernoulli_x_train = np.copy(x_train)
258     bernoulli_x_test = np.copy(x_test)
259
260     # make nonzero positions 1 for the bernoulli model
261     bernoulli_x_train[bernoulli_x_train != 0] = 1
262     bernoulli_x_test[bernoulli_x_test != 0] = 1
263
264     # get the frequencies of the words for spam and normal categories
265     bernoulliSpamFrequencies = y_train.T @ bernoulli_x_train
266     bernoulliNormalFrequencies = (1 - y_train.T) @ bernoulli_x_train
267
268     bernoulliSpamProbabilities = bernoulliSpamFrequencies / spam_count
269     bernoulliNormalProbabilities = bernoulliNormalFrequencies / normal_count
270
271     #print("---- % seconds ----" % (time.time() - start_time))
272     num_correct = tp = tn = fp = fn = 0
273     y_pred_berno = []
274     for i in range(bernoulli_x_test.shape[0]):
275         row = bernoulli_x_test[i]
276         probSpam = np.log(p_spam)
277         probNormal = np.log(p_normal)
278
279         # calculate the spam probabilities for each word separately
280         spamWordExists = (bernoulliSpamProbabilities * row)
281         spamWordDoesntExist = ((1 - bernoulliSpamProbabilities) * (1 - row))
282         spamWord = spamWordDoesntExist + spamWordExists
283         with np.errstate(divide='ignore'):
284             spamWord = np.log(spamWord)
285             spamWord[spamWord == -np.inf] = -1e+12 # replace -inf with extremely small values

```

```

☆ naive_bayes.py > ..
285     spamWord[np.isnan(spamWord)] = -1e+12 # replace 0 with extremely small values
286     probSpam += np.sum(spamWord, axis=0)
287
288     # calculate the normal probabilities for each word separately
289     normalWordExists = (bernoulliNormalProbabilities * row)
290     normalWordDoesntExist = ((1 - bernoulliNormalProbabilities) * (1 - row))
291     normalWord = normalWordDoesntExist + normalWordExists
292     with np.errstate(divide='ignore'):
293         normalWord = np.log(normalWord)
294     normalWord[np.isnan(normalWord)] = -1e+12 # replace 0 with extremely small values
295     probNormal += np.sum(normalWord, axis=0)
296
297     predicted = 0
298     if probSpam > probNormal:
299         predicted = 1
300
301     y_pred_berno.append(predicted)
302     if predicted == y_test[i]:
303         num_correct += 1
304         if predicted == 1:
305             tp += 1
306         else:
307             tn += 1
308     else:
309         if predicted == 1:
310             fp += 1
311         else:
312             fn += 1
313
314 return y_pred_berno, tp, tn, fp, fn
315
316 x_train = X.iloc[:train_size].to_numpy()
317 x_valid = X.iloc[train_size: train_size + valid_size].to_numpy()
318 y_train = y.iloc[:train_size].to_numpy()
319 y_valid = y.iloc[train_size: train_size + valid_size].to_numpy()
320
321 # get a base accuracy score
322 y_pred_berno, tp, tn, fp, fn = fit_bernoulli(x_train, x_valid, y_train, y_valid)
323
324 # list of possible features
325 featureList = ['+', '#', '/', '$', '!', '*', ',', '_', ':']

```

```

☆ naive_bayes.py > ..
325 featureList = ['+', '#', '/', '$', '!', '*', ',', '_', ':']
326 selectedFeatures = []
327 X_current = X.copy()
328 noIncrease = False
329 iterationNo = 1
330 maxAcc = (tp + tn) / (tp + tn + fp + fn) # initial accuracy
331 print('\n----- Bernoulli Model ----- \nBase Acc: ' + str(maxAcc))
332
333 # Apply forward selection to get the best features
334 while not noIncrease:
335     noIncrease = True
336     maxChar = None
337
338     print(f'\n----- Iteration {iterationNo} -----')
339     iterationNo += 1
340     for c in featureList:
341         if c not in selectedFeatures:
342             X_new = X_current.copy()
343             X_new['count'+c] = df['count'+c]
344             X_train_new = X_new.iloc[:train_size]
345             X_valid_new = X_new.iloc[train_size: train_size + valid_size]
346             y_pred_berno, tp, tn, fp, fn = fit_bernoulli(X_train_new, X_valid_new, y_train, y_valid)
347             acc = (tp + tn) / (tp + tn + fp + fn)
348
349             if acc > maxAcc:
350                 maxAcc = acc
351                 noIncrease = False
352                 maxChar = c
353                 print('New highest acc: ' + str(maxAcc))
354
355             if maxChar is not None:
356                 selectedFeatures.append(maxChar)
357                 X_current['count'+maxChar] = df['count'+maxChar]
358                 print('Selected Char: ' + str(maxChar))
359
360     print('Selected Set Of Features:' + str(selectedFeatures))
361     print(X_current.head())
362
363 # The best behaving multinomial model is one with smoothing = 0, report its accuracy again
364 # by training on the train + validation set
365 x_train = X_current.iloc[:train_size + valid_size].to_numpy()

```

```

❶ naive_bayes.py > ...
365     X_train = X_current.iloc[:train_size + valid_size].to_numpy()
366     X_test = X_current.iloc[train_size + valid_size: ].to_numpy()
367     y_train = y.iloc[:train_size + valid_size].to_numpy()
368
369     y_pred_berno, tp, tn, fp, fn = fit_bernoulli(X_train, X_test, y_train, y_test)
370     print("----- Bernoulli Model -----")
371     print("The number of correct predictions: " + str(num_correct) + ", wrong predictions: " + str(X_test.shape[0] - num_correct) + ", accuracy: " + str((tp + tn) / X_test.
372     print("The number of true positives: " + str(tp) + ", true negatives: " + str(tn))
373     print("The number of false positives: " + str(fp) + ", false negatives: " + str(fn))
374
375     # plot confusion matrix
376     confusion = pd.crosstab(y_test, y_pred_berno)
377     fig, ax = plt.subplots()
378     ax.matshow(confusion, cmap='OrRd')
379     ax.set(xlabel='Test', ylabel='Prediction')
380
381     for i in range(2):
382         for j in range(2):
383             c = confusion[j][i]
384             ax.text(j, i, str(c), va='center', ha='center')
385     plt.title('Confusion Matrix For Bernoulli Model')
386     plt.show()
387
388     # print f1 score = 2 * precision * recall / (precision + recall)
389     precision = tp / (fp + tp)
390     recall = tp / (tp + fn)
391     f1_score = (2 * precision * recall) / (precision + recall)
392     print("f1 score for bernoulli model: ", f1_score)

```

log_regression.py

```

❶ log_regression.py > ...
1  import numpy as np
2  import pandas as pd
3  import matplotlib.pyplot as plt
4
5  from typing import Literal, Tuple
6  from tqdm import tqdm
7
8  df = pd.read_csv("data/url_processed.csv")
9  np.random.seed(2023)
10
11 # shuffle the data frame
12 df = df.sample(frac = 1)
13
14 # our feature matrix
15 X = df[['use_of_ip', 'url_length', 'subdomain_length', 'tld_length', 'fld_length', 'path_length',
16          'count_letters', 'count_digits', 'count_puncs', 'count.', 'count@', 'count-',
17          'count%', 'count?', 'count+', 'count/', 'count!', 'count!', 'letters_ratio', 'digit_ratio', 'punc_ratio', 'count_dirs',
18          'use_of_shortener', 'first_dir_length',
19          'url_length_q', 'fld_length_q', 'https', 'count-https', 'count-http', 'count-www', 'is_susp', 'is_php']]
20
21 print(X.head())
22
23 # labels
24 y = df['is_malicious']
25
26 print(X.isna().sum())
27
28 # fill the nan values with 0
29 X = X.fillna(0)
30
31 # check again the number of nan values
32 print(X.isna().sum())
33
34 # divide the data into train, test and validation datasets
35 length = len(y)
36
37 # %70 train, %10 validation, %20 test
38 train_set_size = int(length * 0.7)
39 valid_set_size = int(length * 0.1)
40 test_set_size = int(length * 0.2)
41

```

```

❶ log_regression.py > ...
42
43     X_train = X.iloc[:train_set_size]
44     y_train = y.iloc[:train_set_size]
45
46     X_valid = X.iloc[train_set_size: train_set_size + valid_set_size]
47     y_valid = y.iloc[train_set_size: train_set_size + valid_set_size]
48
49     X_test = X.iloc[train_set_size + valid_set_size:]
50     y_test = y.iloc[train_set_size + valid_set_size:]
51
52     print(X_train.shape, X_valid.shape, X_test.shape)
53
54     print(X_train.head())
55
56     ## Implement the Logistic Regression model
57
58     def sigmoid(z):
59         return 1 / (1 + np.exp(-z))
60
61     def calc_accuracy(y_true, y_pred):
62         y_true = np.asarray(y_true)
63         y_pred = np.asarray(y_pred)
64         return np.mean(y_true == y_pred)
65
66     # logistic regression class
67     class LogisticRegression:
68
69         # param b : y interception
70         # param W : weights
71         # initializer : initial distribution of weights
72         def __init__(self, initializer: Literal['normal', 'uniform', 'zeros']):
73
74             self._b = None
75             self._W = None
76             self.initializer = initializer
77
78         @property
79         def b(self) -> np.ndarray:
80             return self._b
81
82         @property
83         def W(self) -> np.ndarray:

```

```

❶ log_regression.py > ...
84
85         def W(self) -> np.ndarray:
86             return self._W
87
88         def __repr__(self) -> str:
89             return f'LogisticRegression(initializer={self.initializer})'
90
91         def __str__(self) -> str:
92             # calls __repr__
93             return repr(self)
94
95         # initializes the parameters according to given features
96         def init_params(self, in_features: int,
97                         initializer: Literal['normal', 'uniform', 'zeros'] = None):
98             if initializer is None:
99                 initializer = self.initializer
100             if initializer == 'zeros':
101                 self._b = 0
102                 self._W = np.zeros(in_features)
103             else:
104                 random = np.random.default_rng()
105                 if initializer == 'uniform':
106                     self._b = random.uniform(-0.01, 0.01, size=1)
107                     self._W = random.uniform(-0.01, 0.01, size=in_features)
108                 elif initializer == 'normal':
109                     self._b = random.normal(0, 1, size=1)
110                     self._W = random.normal(0, 1, size=in_features)
111
112         # calculates the probability for the output of the sigmoid function
113         def __call__(self, X: np.ndarray) -> np.ndarray:
114             return sigmoid(self.b + X @ self.W)
115
116         # fits the model
117         def fit(self, X: np.ndarray, y: np.ndarray, X_valid: np.ndarray, y_valid: np.ndarray, epochs: int = 1,
118                batch_size: int = None, learning_rate: float = 0.01, shuffle: bool = True) -> np.ndarray:
119
120             # X : design matrix
121             # y : label vector
122             # X_valid : validation set design matrix
123             # y_valid : validation set label vector
124             y = np.asarray(y)

```

```

❶ log_regression.py > ...
123     X = np.asarray(X)
124     y = np.asarray(y)
125     X_valid = np.asarray(X_valid)
126     y_valid = np.asarray(y_valid)
127
128     if self.b is None or self.W is None:
129         self.init_params(X.shape[-1])
130
131     if batch_size is None:
132         batch_size = len(y)
133     n_batches = len(y) // batch_size
134
135     acc_log = calc_accuracy(y_valid, self.predict(X_valid))
136     history = [acc_log]
137     for _ in tqdm(range(epochs)):
138         idx = np.random.permutation(len(y)) if shuffle else np.arange(len(y))
139         for batch in range(n_batches):
140             batch_idx = idx[batch * batch_size: (batch + 1) * batch_size]
141             X_batch = X[batch_idx]
142             y_batch = y[batch_idx]
143             grad_b, grad_W = self._calc_gradients(X_batch, y_batch)
144             self._b = self._b - learning_rate * grad_b
145             self._W = self._W - learning_rate * grad_W
146             acc_log = calc_accuracy(y_valid, self.predict(X_valid))
147             history.append(acc_log)
148
149     return np.asarray(history)
150
151 # predict the output with respect to some threshold
152 def predict(self, X: np.ndarray, threshold: float = 0.5) -> np.ndarray:
153     return np.asarray(self(X) > threshold, dtype=np.int32)
154
155 # calculate gradients and update later
156 def _calc_gradients(self, X: np.ndarray, y: np.ndarray) -> Tuple[np.ndarray, np.ndarray]:
157     y_pred = self(X)
158     grad_b = np.mean(y_pred - y)
159     grad_W = X.T @ (y_pred - y) / len(y)
160     return grad_b, grad_W
161
162 ### Train the logistic regression model with changing batch sizes
163

```

```

❶ log_regression.py > ...
164     acc_batch_normal = []
165     acc_batch_uniform = []
166     acc_batch_zeros = []
167     def logistic_reg_trainer(initializer, batch_size):
168         model = LogisticRegression(initializer)
169         acc = model.fit(X_train, y_train,
170                         X_valid, y_valid,
171                         epochs=20,
172                         batch_size=batch_size,
173                         learning_rate=1e-3)
174         if initializer == 'normal': acc_batch_normal.append(acc)
175         if initializer == 'uniform': acc_batch_uniform.append(acc)
176         if initializer == 'zeros': acc_batch_zeros.append(acc)
177
178     fig, axs = plt.subplots(2, 2, figsize=(18, 12))
179     axs[1, 1].remove()
180
181     # weights initialized NORMAL dist with different batches
182     logistic_reg_trainer('normal', 32)
183     logistic_reg_trainer('normal', 64)
184     logistic_reg_trainer('normal', 128)
185     logistic_reg_trainer('normal', 256)
186     # plot normal
187     x = list(range(20 + 1))
188     #plt.figure(figsize=(18, 12))
189     axs[0, 0].set_title('Batch Size Comparison for Normal initializer')
190     axs[0, 0].set_xlabel('Epochs')
191     axs[0, 0].set_ylabel('Validation accuracy')
192     axs[0, 0].plot(x, acc_batch_normal[0], label='Batch size = 32')
193     axs[0, 0].plot(x, acc_batch_normal[1], label='Batch size = 64')
194     axs[0, 0].plot(x, acc_batch_normal[2], label='Batch size = 128')
195     axs[0, 0].plot(x, acc_batch_normal[3], label='Batch size = 256')
196     axs[0, 0].legend()
197
198     # weights initialized UNIFORM dist with different batches
199     logistic_reg_trainer('uniform', 32)
200     logistic_reg_trainer('uniform', 64)
201     logistic_reg_trainer('uniform', 128)
202     logistic_reg_trainer('uniform', 256)
203     # plot uniform
204     #plt.figure(figsize=(18, 12))

```

```

❶ log_regression.py > ...
205     axs[0,1].set_title('Batch Size Comparison for Uniform initializer')
206     axs[0,1].set_xlabel('Epochs')
207     axs[0,1].set_ylabel('Validation accuracy')
208     axs[0,1].plot(x, acc_batch_uniform[0], label='Batch size = 32')
209     axs[0,1].plot(x, acc_batch_uniform[1], label='Batch size = 64')
210     axs[0,1].plot(x, acc_batch_uniform[2], label='Batch size = 128')
211     axs[0,1].plot(x, acc_batch_uniform[3], label='Batch size = 256')
212     axs[0,1].legend()
213
214     # weights initialized ZEROS with different batches
215     logistic_reg_trainer('zeros', 32)
216     logistic_reg_trainer('zeros', 64)
217     logistic_reg_trainer('zeros', 128)
218     logistic_reg_trainer('zeros', 256)
219
220     # plot zeros
221     #plt.figure(figsize=(18, 12))
222     axs[1, 0].set_title('Batch Size Comparison for Zeros initializer')
223     axs[1, 0].set_xlabel('Epochs')
224     axs[1, 0].set_ylabel('Validation accuracy')
225     axs[1, 0].plot(x, acc_batch_zeros[0], label='Batch size = 32')
226     axs[1, 0].plot(x, acc_batch_zeros[1], label='Batch size = 64')
227     axs[1, 0].plot(x, acc_batch_zeros[2], label='Batch size = 128')
228     axs[1, 0].plot(x, acc_batch_zeros[3], label='Batch size = 256')
229     axs[1, 0].legend()
230
231     plt.subplots_adjust(wspace=0.2, hspace=0.5)
232     plt.show()
233
234     for i in range(1, 5):
235         print('Final Validation accuracy for Normal Initialization With Batch Size: %d is : %f' % (32 * i, acc_batch_normal[i - 1][-1]))
236     for i in range(1, 5):
237         print('Final Validation accuracy for Uniform Initialization With Batch Size: %d is : %f' % (32 * i, acc_batch_uniform[i - 1][-1]))
238     for i in range(1, 5):
239         print('Final Validation accuracy for Zeros Initialization With Batch Size: %d is : %f' % (32 * i, acc_batch_zeros[i - 1][-1]))
240
241     model = LogisticRegression("normal")
242
243     X_train_new = np.concatenate((X_train, X_valid), axis=0)
244     y_train_new = np.concatenate((y_train, y_valid), axis=0)
245     final_test_acc = model.fit(X_train_new, y_train_new,
246                                X_test, y_test,

```

```

❶ log_regression.py > ...
246     X_test, y_test,
247     epochs=20,
248     batch_size=64,
249     learning_rate=1e-3)
250
251     # Final Test Accuracy For The Best Logistic Regression Model
252     x = list(range(20 + 1))
253     plt.figure(figsize=(18, 12))
254     plt.xlabel('Epochs')
255     plt.ylabel('Test Accuracy')
256     plt.xticks(x for x in range(0, 22, 2)))
257     plt.plot(x, final_test_acc, label='Multinomial Model')
258     plt.legend()
259     plt.title('Test Accuracy for the Logistic Regression Model')
260     plt.show()
261
262     y_preds = model.predict(X_test)
263     confusion = pd.crosstab(y_test, y_preds)
264
265     # Report the final test accuracy
266     print('Final Test Accuracy: %f' % (final_test_acc[-1]))
267
268     # print f1 score = 2 * precision * recall / (precision + recall)
269     precision = confusion[1][1] / (confusion[0][1] + confusion[1][1])
270     recall = confusion[1][1] / (confusion[1][1] + confusion[1][0])
271     f1_score = (2 * precision * recall) / (precision + recall)
272     print("F1 Score For Final Logistic Regression Model: ", f1_score)
273
274     # plot confusion matrix
275     fig, ax = plt.subplots()
276     ax.matshow(confusion,cmap='OrRd')
277     ax.set(xlabel='Test', ylabel='Prediction')
278
279     for i in range(2):
280         for j in range(2):
281             c = confusion[j][i]
282             ax.text(i, j, str(c), va='center', ha='center')
283     plt.title('Confusion Matrix For Final Logistic Regression Model')
284     plt.show()

```

metrics.py

```
metrics.py > ...
1  import numpy as np
2
3  # mean squared error
4  def mse(y_true: np.ndarray, y_pred: np.ndarray) -> float:
5      y_true = np.asarray(y_true)
6      y_pred = np.asarray(y_pred)
7      y_true = y_true.reshape((len(y_true), 1))
8      return np.sum(np.square(y_true - y_pred)) / len(y_true)
9
10 # mean absolute error
11 def mae(y_true: np.ndarray, y_pred: np.ndarray) -> float:
12     y_true = np.asarray(y_true)
13     y_pred = np.asarray(y_pred)
14     y_true = y_true.reshape((len(y_true), 1))
15     return np.sum(np.abs(y_true - y_pred)) / len(y_true)
16
17
18 def r2(y_true: np.ndarray, y_pred: np.ndarray) -> float:
19     r = np.corrcoef(y_true.squeeze(), y_pred.squeeze())[0, 1]
20     return r ** 2
```

nn.py

```
nn.py > ⊕ calc_accuracy
1  import numpy as np
2  from activation_functions import ActivationFunction
3  from typing import Any, Type, Tuple, Iterable
4  from collections import defaultdict, namedtuple
5  from tqdm import tqdm
6  from metrics import mse, mae, r2
7
8  def relu(z):
9      return np.maximum(0, z)
10
11 def relu_backward(z):
12     return np.where(z > 0, 1, 0)
13
14 def sigmoid(z):
15     return 1 / (1 + np.exp(-z))
16
17 def sigmoid_backward(z):
18     return (1 - sigmoid(z)) * sigmoid(z)
19
20 def calc_accuracy(y_true, y_pred):
21     y_true = np.asarray(y_true).squeeze()
22     y_pred = np.asarray(y_pred).squeeze()
23     return np.mean(y_true == y_pred)
24
25 FullyConnectedLayerWeights = namedtuple('FullyConnectedLayerWeights', ['b', 'W'])
26 FullyConnectedLayerGradients = FullyConnectedLayerWeights
27
28 DEFAULT_METRICS = {
29     'MSE': mse,
30     'MAE': mae,
31     'R2': r2,
32 }
33
34 class NeuralNetwork:
35     def __init__(self, n_neurons):
36         self.input_units = None
37         # number of neurons in the layer
38         self.n_neurons = n_neurons
39         self.perceptron = None
40
41     def init_network(self):
42         pass
```

```

4 nn.py > calc_accuracy
41     def init_network(self):
42         rng = np.random.default_rng()
43
44         # init weights
45         self.perceptron = []
46         for i, n in enumerate(self.n_neurons):
47             n_prev = self.n_neurons[i - 1] if i != 0 else self.input_units
48             bound = np.sqrt(6 / (n + n_prev))
49             b = rng.uniform(-bound, bound, size=(n, 1))
50             W = rng.uniform(-bound, bound, size=(n, n_prev))
51
52             layer = FullyConnectedLayerWeights(b, W)
53             self.perceptron.append(layer)
54
55     def __call__(self, X):
56         Z = X.T
57         for layer in self.perceptron[:-1]:
58             V = layer.W @ Z + layer.b
59             Z = relu(V)
60         V = self.perceptron[-1].W @ Z + self.perceptron[-1].b
61         Z = sigmoid(V)
62         return Z.T
63
64         # predict the output with respect to some threshold
65     def predict(self, X: np.ndarray, threshold: float = 0.5) -> np.ndarray:
66         return np.asarray(self(X) > threshold, dtype=np.int32)
67
68         # forward pass through the layers by caching the outputs
69     def forward(self, X, y):
70         Z = X
71         V_cache, Z_cache = [], []
72         for layer in self.perceptron[:-1]:
73             V = layer.W @ Z + layer.b
74             Z = relu(V)
75
76             V_cache.append(V)
77             Z_cache.append(Z)
78
79         V = self.perceptron[-1].W @ Z + self.perceptron[-1].b
80         Z = sigmoid(V)
81
82         V_cache.append(V)
83
84         v
85
86         v
87
88         v
89         v
90
91         v
92
93         v
94         v
95         v
96         v
97         v
98         v
99         v
100        v
101        v
102        v
103        v
104        v
105        v
106        v
107        v
108        v
109        v
110        v
111        v
112        v
113        v
114        v
115        v
116        v
117        v
118        v
119        v
120        v
121        v
122        v
123        v

```

```

4 nn.py > calc_accuracy
42         V_cache.append(V)
43         Z_cache.append(Z)
44
45         # compute metrics
46         J = {metric: fn(y.T, Z.T) for metric, fn in DEFAULT_METRICS.items()}
47         return J, V_cache, Z_cache
48
49     def backward(self, X, y, V_cache, Z_cache):
50
51         # gradients for the model
52         gradients = []
53
54         delta = (1 / Z_cache[-1].shape[1]) * np.subtract(Z_cache[-1], y)
55         db = np.mean(delta, axis=1, keepdims=True)
56         dW = delta @ Z_cache[-2].T / Z_cache[-2].shape[-1]
57         gradients.append(FullyConnectedLayerGradients(db, dW))
58
59         # calculate gradients for the hidden layer
60         for i in reversed(range(1, len(self.perceptron) - 1)):
61             delta = (self.perceptron[i + 1].W.T @ delta) * relu_backward(V_cache[i])
62             db = np.mean(delta, axis=1, keepdims=True)
63             dW = delta @ Z_cache[i - 1].T / Z_cache[i - 1].shape[-1]
64             gradients.append(FullyConnectedLayerGradients(db, dW))
65
66         # treat the first hidden layer differently by using X (design matrix)
67         delta = (self.perceptron[1].W.T @ delta) * relu_backward(V_cache[0])
68         db = np.mean(delta, axis=1, keepdims=True)
69         dW = (delta @ X.T) / X.shape[-1]
70         gradients.append(FullyConnectedLayerGradients(db, dW))
71
72         # get the correct order
73         gradients.reverse()
74         return gradients
75
76
77     # forward and backward pass combined
78     def step(self, X, y):
79         J, V_cache, Z_cache = self.forward(X, y)
80         gradients = self.backward(X, y, V_cache, Z_cache)
81         return J, gradients
82
83     def fit(self,
84             X, # training features
85             y, # training labels
86             learning_rate=0.01,
87             epochs=1000,
88             batch_size=32,
89             validation_split=0.2,
90             early_stopping=False,
91             patience=10,
92             verbose=1,
93             shuffle=True,
94             validation_data=None,
95             **kwargs):
96
97         # validate inputs
98         if not isinstance(X, np.ndarray):
99             raise ValueError("X must be a NumPy array")
100            ...
101            ...
102            ...
103            ...
104            ...
105            ...
106            ...
107            ...
108            ...
109            ...
110            ...
111            ...
112            ...
113            ...
114            ...
115            ...
116            ...
117            ...
118            ...
119            ...
120            ...
121            ...
122            ...
123            ...

```

```

❷ nn.py > ⚡ calc_accuracy
123     X, # training features
124     y, # training labels
125     X_valid, # validation features
126     y_valid, # validation labels
127     alpha=0.1, # learning rate
128     momentum=0.85, # momentum
129     epochs=50, # number of training iterations
130     batch_size=32, # batch size for training
131     patience=5, # patience
132     min_delta=1e-7, # patience criteria
133     shuffle=True,
134     cold_start=False
135     ):
136
137     X = np.asarray(X)
138     y = np.asarray(y)
139     X_valid = np.asarray(X_valid)
140     y_valid = np.asarray(y_valid)
141
142     if cold_start or self.perceptron is None:
143         self.input_units = X.shape[-1]
144         self.init_network()
145
146     n_batches = len(X) // batch_size
147     history = defaultdict(list)
148
149     delta_weights_prev = None
150     n_no_improvement = 0
151     for epoch in (progress_bar := tqdm(range(epochs))):
152         # get the shuffled data
153         train_indices = np.random.permutation(len(X)) if shuffle else np.arange(len(X))
154         batch_avg_losses = []
155         for batch in range(n_batches):
156             # get the batch data from the shuffle
157             batch_indices = train_indices[batch * batch_size: (batch + 1) * batch_size]
158             X_batch = X[batch_indices].T
159             y_batch = y[batch_indices].T
160
161             # take one step: forward and backward combined
162             J, gradients = self.step(X_batch, y_batch)
163
164             batch_avg_losses.append(J)
165
❸ nn.py > ⚡ calc_accuracy
165     batch_avg_losses.append(J)
166
167     # update weights
168     delta_weights = self.add_velocity(delta_weights_prev, gradients, momentum)
169     self.update_params(delta_weights, alpha)
170
171     # updates saved for momentum
172     delta_weights_prev = delta_weights
173
174     # evaluate the model on validation dataset
175     valid_avg_losses = {metric: np.mean(fn(y_valid, self(X_valid)))
176                         for metric, fn in DEFAULT_METRICS.items()}
177     train_avg_losses = {metric: np.mean([loss[metric] for loss in batch_avg_losses])
178                         for metric in DEFAULT_METRICS.keys()}
179
180     valid_acc_log = calc_accuracy(y_valid, self.predict(X_valid))
181     train_acc_log = calc_accuracy(y, self.predict(X))
182
183     history['train_acc'].append(train_acc_log)
184     history['valid_acc'].append(valid_acc_log)
185
186     for metric in DEFAULT_METRICS.keys():
187         history[f'train_{metric}'].append(train_avg_losses[metric])
188         history[f'valid_{metric}'].append(valid_avg_losses[metric])
189
190     # for output display
191     progress_bar.set_description_str(f'n_neurons={"-".join([str(i) for i in self.n_neurons])}, '
192                                         f'alpha={alpha}, momentum={momentum}, batch_size={batch_size}')
193     progress_bar.set_postfix_str(f'train_acc=[{train_acc_log:.7f}], '
194                                 f'valid_acc=[{valid_acc_log:.7f}]')
195
196     # exit training if there is no significant improvement in the model
197     if epoch > 2 and history['valid_MSE'][-2] - history['valid_MSE'][-1] < min_delta:
198         n_no_improvement += 1
199         if n_no_improvement > patience:
200             break
201         else:
202             n_no_improvement = 0
203     return history
204
205     # calculate the updates by introducing momentum
206     @staticmethod
❹ nn.py > ⚡ calc_accuracy
206
207     def add_velocity(delta_weights_prev, gradients, momentum):
208         delta_weights = [
209             FullyConnectedLayerWeights(*[momentum * delta_w_prev + (1 - momentum) * grads
210                                         for delta_w_prev, grads in zip(delta_weights_prev_, gradients_)])
211             for delta_weights_prev_, gradients_ in zip(delta_weights_prev, gradients)]
212         if delta_weights_prev else gradients
213         return delta_weights
214
215     # apply updates on to the weights
216     def update_params(self, delta_weights, alpha):
217         self.perceptron = [
218             FullyConnectedLayerWeights(*[w - alpha * delta for w, delta in zip(layer, delta_weights_)])
219             for layer, delta_weights_ in zip(self.perceptron, delta_weights)]

```

nn_test.py

```
❷ nn_test.py > ...
1  import pandas as pd
2  import numpy as np
3  import matplotlib.pyplot as plt
4  from nn import NeuralNetwork
5  from itertools import product
6  from nn import calc_accuracy
7  from os import path
8
9  df = pd.read_csv("data/url_processed.csv")
10 np.random.seed(2023)
11
12 # shuffle the data frame
13 df = df.sample(frac = 1)
14
15 # our feature matrix
16 X = df[['use_of_ip', 'url_length', 'subdomain_length', 'tld_length', 'path_length',
17         'count_letters', 'count_digits', 'count_puncs', 'count.', 'count@', 'count-',
18         'count%', 'count?', 'count=', 'count+', 'count/', 'count,', 'count!',
19         'letters_ratio', 'digit_ratio', 'punc_ratio', 'count_dirs',
20         'use_of_shortener', 'first_dir_length',
21         'url_length_q', 'fld_length_q', 'https', 'count_https', 'count_http',
22         'count-www', 'is_susp', 'is_php']]
22
23 print(X.head())
24
25 # labels
26 y = df['is_malicious']
27
28 print(X.isna().sum())
29
30 # fill the nan values with 0
31 X = X.fillna(0)
32
33 # check again the number of nan values
34 print(X.isna().sum())
35
36 # divide the data into train, test and validation datasets
37 length = len(y)
38
39 # %70 train, %10 validation, %20 test
40 train_set_size = int(length * 0.7)
41 valid_set_size = int(length * 0.1)
```

```
❷ nn_test.py > ...
41  valid_set_size = int(length * 0.1)
42  test_set_size = int(length * 0.2)
43
44  X_train = X.iloc[:train_set_size]
45  y_train = y.iloc[:train_set_size]
46
47  X_valid = X.iloc[train_set_size: train_set_size + valid_set_size]
48  y_valid = y.iloc[train_set_size: train_set_size + valid_set_size]
49
50  X_test = X.iloc[train_set_size + valid_set_size:]
51  y_test = y.iloc[train_set_size + valid_set_size:]
52
53  nn_layers_list = [
54      [64, 1],
55      [32, 32, 1],
56      [64, 64, 1],
57  ]
58  nn_alphas = [1e-2, 5e-3, 1e-3, 5e-4]
59  nn_momentums = [0.85, 0.95]
60  nn_batch_sizes = [32, 64]
61
62  nn_hyperparams = list(product(nn_alphas, nn_momentums, nn_batch_sizes))
63
64  def plot_from_history(history, title):
65      h = pd.DataFrame.from_dict(history)[['train_acc', 'valid_acc']]
66      h.set_axis(['Train Acc', 'Test Acc'], axis=1)
67      h.plot()
68      plt.xlabel('Epochs')
69      plt.ylabel('Accuracy')
70      plt.title(title)
71      plt.show()
72
73  def subplot_mse(hist_list, nn_hyperparams, layers):
74
75      fig, axs = plt.subplots(4, 4, figsize=(18, 18))
76      for i in range(4):
77          for j in range(4):
78              axs[i,j].set_title('lr: ' + str(nn_hyperparams[4*i+j][0]) + ', mom: ' + str(nn_hyperparams[4*i+j][1]) + ', bs: ' + str(nn_hyperparams[4*i+j][2]), fontdict={
79                  'fontsize': 8})
80              axs[i,j].set_xlabel('Epochs', fontdict={'fontsize': 8})
81              axs[i,j].set_ylabel('Accuracy', fontdict={'fontsize': 8})
82              axs[i,j].plot(hist_list[4 * i + j]['train acc'], label='Train Acc')
```

```

❷ nn_test.py > ...
51     axs[i,j].plot(hist_list[4 * i + j][train_acc], label='Train Acc')
52     axs[i,j].plot(hist_list[4 * i + j][valid_acc], label='Validation Acc')
53     axs[i,j].legend(fontsize=5)
54 plt.suptitle('Hyperparameter Tuning For NN With layers: ' + str(layers))
55 plt.subplots_adjust(wspace=0.5, hspace=0.7)
56 plt.savefig(path.join('plots', f'{layers}-validation-plot.jpg'))
57 plt.show()
58
59 accList = []
60 layer_hist_list = []
61 for layers in nn_layers_list:
62     hist_list = []
63     for alpha, momentum, batch_size in nn_hyperparams:
64         nn = NeuralNetwork(n_neurons=layers)
65         history = nn.fit(X_train, y_train, X_valid, y_valid, alpha=alpha, batch_size=batch_size, momentum=momentum, epochs=200, patience=5)
66         hist_list.append(history)
67         accList.append(history['valid_acc'][-1]) # push the last accuracy obtained
68     layer_hist_list.append(hist_list)
69
70 for i, layers in enumerate(nn_layers_list):
71     subplot_mse(layer_hist_list[i], nn_hyperparams, layers)
72
73 # Get the best parameters by using the validation accuracy
74 # as the metric
75 accList = np.asarray(accList)
76 nn_best_index = np.argmax(accList)
77 nn_best_layer_index = nn_best_index // len(nn_hyperparams)
78 nn_best_layer = nn_layers_list[nn_best_layer_index]
79 nn_best_params = nn_hyperparams[nn_best_index % len(nn_hyperparams)]
80 nn_best_alpha = nn_best_params[0]
81 nn_best_momentum = nn_best_params[1]
82 nn_best_batchsize = nn_best_params[2]
83
84 # Plot the distribution of the accuracies
85 plt.figure()
86 plt.title('The number of models with a given validation accuracy')
87 plt.xlabel('Accuracy of the model')
88 plt.ylabel('Number of models')
89 plt.hist(accList * 100)
90 plt.show()
91
92 print(f'\nBest Neural Network Settings: Size:{nn_best_layer}, Learning Rate:{nn_best_alpha}, Batch Size:{nn_best_batchsize}, Momentum:{nn_best_momentum}')
❸ nn_test.py > ...
122 print(f'\nBest Neural Network Settings: Size:{nn_best_layer}, Learning Rate:{nn_best_alpha}, Batch Size:{nn_best_batchsize}, Momentum:{nn_best_momentum}')
123
124 # train the best model on train + validation dataset, report metrics
125 # on the test dataset
126 X_train = X.iloc[:train_set_size + valid_set_size]
127 y_train = y.iloc[:train_set_size + valid_set_size]
128 best_model = NeuralNetwork(n_neurons=nn_best_layer)
129 history = best_model.fit(X_train, y_train, X_test, y_test, alpha=nn_best_alpha, batch_size=nn_best_batchsize, momentum=nn_best_momentum, epochs=500, patience=5)
130
131 # plot the test and train accuracies of the best model
132 plot_from_history(history, 'Accuracy vs Epoch For Best Model')
133
134 y_test = y_test.to_numpy()
135 y_preds = best_model.predict(X_test)
136 confusion = pd.crosstab(y_test.reshape((-1,)), y_preds.reshape((-1,)))
137
138 # Report the final test accuracy
139 print('Final Test Accuracy: %f' % (calc_accuracy(y_test.reshape((-1,)), y_preds.reshape((-1,))))))
140
141 # print f1 score = 2 * precision * recall / (precision + recall)
142 precision = confusion[1][1] / (confusion[0][1] + confusion[1][1])
143 recall = confusion[1][1] / (confusion[1][1] + confusion[1][0])
144 f1_score = (2 * precision * recall) / (precision + recall)
145 print("F1 Score For Final Neural Network Model: ", f1_score)
146
147 # plot confusion matrix
148 fig, ax = plt.subplots()
149 ax.matshow(confusion, cmap='OrRd')
150 ax.set(xlabel='Test', ylabel='Prediction')
151
152 for i in range(2):
153     for j in range(2):
154         c = confusion[j][i]
155         ax.text(i, j, str(c), va='center', ha='center')
156 plt.title('Confusion Matrix For Final Neural Network Model')
157 plt.show()

```