

Low Level Design (LLD)

PHISHING DOMAIN DETECTION

Document Version: 0.1

Last Revised Date: 13-May-2023

Parvej alam Ansari

Document Version Control:

Version	Date	Author	Description
0.1	13-05-2023	Parvej alam	Introduction Architecure
0.1	13-05-2023	Parvej alam	Architecture Description Unit Test Cases

Contents

1	Introduction-----	3
1.1	About Low-Level Design Document-----	3
1.2	Scope-----	3
2	Architecture-----	3
3	Architecture Description-----	4
3.1	Data Gathering-----	4
3.2	Dataset Description-----	4
3.3	Data Ingestion from MongoDB-----	5
3.4	Data Validation-----	5
3.5	Data Transformation-----	5
3.6	Model Building using Model Trainer-----	5
3.7	Model Evaluation-----	6
3.8	Model Training using FAST API-----	6
3.9	Saving Model and Preprocessing File-----	6
3.10	Cloud Setup and Model Pushing-----	6
3.11	User Input and Application Output-----	6
4	Unit Test Cases-----	7

1. Introduction

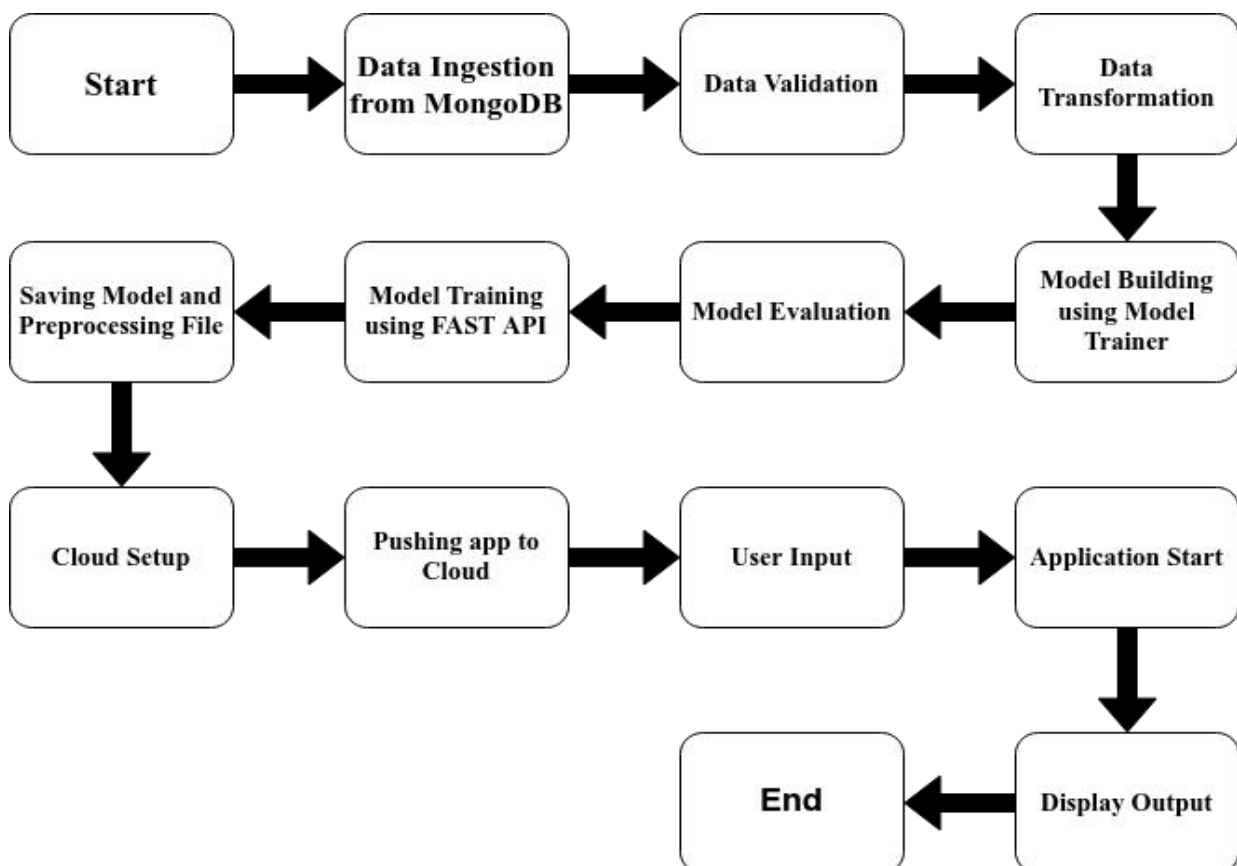
1.1 About Low-Level Design Document

The goal of LLD or a low-level design document (LLDD) is to give the internal logical design of the actual program code for Energy Efficiency System. LLD describes the class diagrams with the methods and relations between classes and program specs. It describes the modules so that the programmer can directly code the program from the document.

1.2 Scope

Low-level design (LLD) is a component-level design process that follows a step-by-step refinement process. This process can be used for designing data structures, required software architecture, source code and ultimately, performance algorithms. Overall, the data organization may be defined during requirement analysis and then refined during data design work.

2. Architecture



3. Architecture Description

3.1 Data Gathering

This data set we will be using is from the Mendeley Phishing Websites Dataset. The following is **Mendeley data set**:

Source:

Published: 24 September 2020 | Version 1 | DOI: 10.17632/72ptz43s9v.1

Contributor: Grega Vrbančič

Dataset link: <https://data.mendeley.com/datasets/72ptz43s9v/1> (csv format)

3.2 Dataset Description

Data Set Information:

These data consist of a collection of legitimate as well as phishing website instances. Each website is represented by the set of features which denote, whether website is legitimate or not. Data can serve as an input for machine learning process. In this repository the two variants of the Phishing Dataset are available But Full variant dataset is selected. Full variant - dataset_full.csv Short description of the full variant dataset: Total number of instances: 88,647 Number of legitimate website instances (labeled as 0): 58,000 Number of phishing website instances (labeled as 1): 30,647 Total number of features: 111

Attribute Information:

The dataset contains 112- attributes out of which 111 are inputs (or features) and 1 attribute is response in terms of legitimate or phishing (or outcome). The aim is to use the 111 features to classify whether the given URL is safe to search or not.

```
columns = ["qty_dot_url", "qty_hyphen_url", "qty_underline_url", "qty_slash_url",
"qty_questionmark_url", "qty_equal_url", "qty_at_url", "qty_and_url", "qty_exclamation_url",
"qty_space_url", "qty_tilde_url", "qty_comma_url", "qty_plus_url", "qty_asterisk_url",
"qty_hashtag_url", "qty_dollar_url", "qty_percent_url", "qty_tld_url", "length_url",
"qty_dot_domain", "qty_hyphen_domain", "qty_underline_domain", "qty_slash_domain",
"qty_questionmark_domain", "qty_equal_domain", "qty_at_domain", "qty_and_domain",
"qty_exclamation_domain", "qty_space_domain", "qty_tilde_domain", "qty_comma_domain",
"qty_plus_domain", "qty_asterisk_domain", "qty_hashtag_domain", "qty_dollar_domain",
"qty_percent_domain", "qty_vowels_domain", "domain_length", "domain_in_ip",
"server_client_domain", "qty_dot_directory", "qty_hyphen_directory", "qty_underline_directory",
"qty_slash_directory", "qty_questionmark_directory", "qty_equal_directory", "qty_at_directory",
"qty_and_directory", "qty_exclamation_directory", "qty_space_directory", "qty_tilde_directory",
"qty_comma_directory", "qty_plus_directory", "qty_asterisk_directory", "qty_hashtag_directory",
"qty_dollar_directory", "qty_percent_directory", "directory_length", "qty_dot_file",
"qty_hyphen_file", "qty_underline_file", "qty_slash_file", "qty_questionmark_file",
"qty_equal_file", "qty_at_file", "qty_and_file", "qty_exclamation_file", "qty_space_file",
"qty_tilde_file", "qty_comma_file", "qty_plus_file", "qty_asterisk_file", "qty_hashtag_file",
```

```
"qty_dollar_file", "qty_percent_file", "file_length", "qty_dot_params", "qty_hyphen_params",
"qty_underline_params", "qty_slash_params", "qty_questionmark_params", "qty_equal_params",
"qty_at_params", "qty_and_params", "qty_exclamation_params", "qty_space_params",
"qty_tilde_params", "qty_comma_params", "qty_plus_params", "qty_asterisk_params",
"qty_hashtag_params", "qty_dollar_params", "qty_percent_params", "params_length",
"tld_present_params", "qty_params", "email_in_url", "time_response", "domain_spf", "asn_ip",
"time_domain_activation", "time_domain_expiration", "qty_ip_resolved", "qty_nameservers",
"qty_mx_servers", "ttl_hostname", "tls_ssl_certificate", "qty_redirects", "url_google_index",
"domain_google_index", "url_shortened"]
```

3.3. Data Ingestion from MongoDB

Data Insertion into Database

- a. Database creation and connection - create a database with name passed. If the database is already created, open the connection to the database.
- b. Table creation in the database.
- c. Insertion of files in the table

Export Data from Database

Data export from database - The data in a stored database is exported as a .CSV file to be used for data pre-processing and model Training.

3.4. Data Validation

After data ingestion stage, the data should be validated in terms of number of columns, column names etc.

3.5. Data Transformation

As soon as data is successfully validated from data validation stage, The validated data is then sent to the data transformation stage where data is wrangled, cleaned by using some data imputation techniques. Furthermore if outliers are present then apply some outliers handling techniques. Next, if dataset is not having proper scaling then will apply proper scaling techniques to transform the data. Now data is ready to next stage i.e. Model Trainer.

3.6. Model Building using Model Trainer

In this stage, Number of classification type supervised machine learning models should be applied and will try to find out the best suitable model as per our available dataset by using some metric evaluation techniques. After getting the best suitable classification model, its training and testing accuracy will be determined again.

3.7. Model Evaluation

In this stage, selected and hypertuned model is evaluated using some criteria like training and testing accuracy with the help of proper metric evaluation techniques, overfitting and underfitting etc. once this selected model is properly evaluated, it will be passed to the pipeline stage.

3.8. Model Training using FAST API

After successful run of individual component, a pipeline will be created where each and every component discussed above will merge and that pipeline will be triggered using FAST API where the finalised model shall be trained with available training dataset.

3.9. Saving Model and Preprocessing File

In this stage, the trained model file plus the preprocessing file will be stored in pickled format.

3.10. Cloud Setup and Model Pushing

After the model/product/app is ready, the app is deployed to AWS by using some AWS services like Elastic Beanstalk (EB) and Docker hub. This cloud deployment will help user to access the application through any internet devices.

3.11. User Input and Application output

By using the app URL, user will get access to the application/model, where user will feed inputs to the model and that model will render the best outcome to its webpage in terms URL is Phishing or Non-phishing,

4. Unit Test Cases

Test Case Description	Pre-Requisite	Expected Result
Verify whether the Application URL is accessible to the user	1. Application URL should be defined	Application URL should be accessible to the user
Verify whether the Application loads completely for the user when the URL is accessed	1. Application URL is accessible 2. Application is deployed	The Application should load completely for the user when the URL is accessed
Verify whether the User is able to sign up in the application	1. Application is accessible	The User should be able to sign up in the application
Verify whether user is able to see input fields on logging in	1. Application is accessible 2. User is signed up to the application 3. User is logged in to the application	User should be able to see input fields on logging in
Verify whether user is able to edit all input fields	1. Application is accessible 2. User is signed up to the application 3. User is logged in to the application	User should be able to edit all input fields
Verify whether user gets Submit button to submit the inputs	1. Application is accessible 2. User is signed up to the application 3. User is logged in to the application	User should get Submit button to submit the inputs
Verify whether user is presented with recommended results on clicking submit	1. Application is accessible 2. User is signed up to the application 3. User is logged in to the application	User should be presented with recommended results on clicking submit