

High Level Design (HLD)

Spam Ham Classifier

Revision Number:

Last date of revision:

Document Version Control

Date Issued	Version	Description	Author
20-11-2021	1	Initial HLD-V1.0	Upendra Kumar

Contents

Document Version Control-----	2
Abstract-----	4
1 Introduction-----	5
1.1 Why this High-Level Design Document-----	5
1.2 Scope-----	5
1.3 Definitions-----	5
2 General Description-----	6
2.1 Product Perspective-----	6
2.2 Problem Statement-----	6
2.3 Proposed Solution-----	6
2.4 Further Improvements-----	6
2.5 Data Requirements-----	7
2.6 Tools Used-----	8
2.7 Constraints-----	9
2.8 Assumptions-----	9
3 Design Details-----	9
3.1 Process Flow-----	9
3.1.1 Model Training and Evaluation-----	10
3.1.2 Deployment Process-----	10
3.2 Event log-----	11
3.3 Error Handling-----	11
4 Performance-----	12
4.1 Reusability-----	12
4.2 Application Compatibility-----	12
4.3 Resource Utilization-----	12
4.4 Deployment-----	12
5 Conclusion-----	13
6 References-----	13

Abstract

Nowadays, most people have access to the Internet, and they cannot survive without Smartphone and computers. They not only use the Internet for fun and entertainment, but they also use it for business, stock marketing, searching, sending e-mails, and so on. Hence, the usage of the Internet is growing rapidly.

One of the threats for such technology is a spam. Spam has a lot of definitions, as it is considered one of the complex problems in e-mail services. **Spam is a junk mail/message, or an unsolicited mail/message. Spam e-mails are also those unwanted, unsolicited e-mails that are not intended for a specific receiver.** It is basically an online communication sent to the user without permission. It takes on various forms like adult content, selling products or services, job offers, and so forth. The spam has increased tremendously in the last few years.

The good, perfect, and official mails are known as *ham*. It is also defined as an e-mail that is generally desired. Today, more than 85% of mails or messages received by users are spam. It costs the sender very little time to send, but most of the costs are paid by the recipient or the service providers rather than by the sender.

The cost of spam can also be measured in lost human time, lost server time and loss of valuable mail/messages. The sent mail reserves a quota in the server, and the receiver may have a limited space, causing the server to reject another ham mail because it is out of space. Moreover, the reader may lose a lot of time in reading unuseful messages.

The Machine Learning field has a robust, ready-made and alternative way for solving this type of the problem. We can harness the power of Machine Learning to build a classifier type of thing, which help us to classify which message is/are spam or ham depending upon their context or content.

1 Introduction

1.1 Why this High-Level Design Document?

The purpose of this High Level Design (HLD) Document is to add the necessary details to the current project description to represent a suitable model for coding. This document is also intended to help detect contradictions prior to coding, and can be used as reference manual for how the modules interact at a high level.

The HLD will

- Present all of the design aspects and define them in detail
- Describe the user interface being implemented
- Describe the hardware and software interfaces
- Describe the performance requirements
- Include design feature and the architecture of the project
- List and describe the non-functional attribute like:
 - Security
 - Reliability
 - Maintainability
 - Portability
 - Reusability
 - Application compatibility
 - Resource utilization
 - Serviceability

1.2 Scope

The HLD document presents the structure of the system, such as the database architecture, application architecture (layers), application flow (Navigation), and technology architecture. The HLD uses non-technical to mildly-technical terms which should be understandable to the administrators of the system.

1.3 Definitions

- **Spam** — e-mails are those unwanted, unsolicited e-mails that are not intended for a specific receiver.
- **Ham** — The good, perfect, and official mails are known as *ham*.
- **SHC** — Spam Ham Classifier.

2 General Description

2.1 Product Perspective

The Spam Ham Classifier solution system is a data science-based machine learning model which help us to classify text messages into spam or ham depending upon their content.

2.2 Problem Statement

To create an AI solution for classifying spam and ham messages and to implement the following use cases.

- To classify spam text messages into spam.
- To classify ham text messages into ham.

2.3 Proposed Solution

The solution proposed here is a data science model based on machine learning can be implemented to perform above mention use cases. In first use case , we will take input from spam message and see whether proposed solution is going to detect it or not. And in second use case, we will take input from ham message, and check our solution whether it is performing or not in right way.

2.4 Further Improvements

The spam ham detection solution can be added with more use cases in IT domain and messaging services. SHC solution can also be synchronized with other IT domain and messaging solutions to give one step extra security to user from those people who try to cheat them by using spam messages.

2.5 Data Requirements

Data requirement completely depend on our problem statement.

we need text data of message which are already classified as Spam or Ham. Data Set should contain at least two attributes one class of message and other message body containing the complete message.

- **Class**: Message class either **Spam** or **Ham**.
- **Message**: Body of message containing complete message.

What is Spam Message?

A spam message is **any unrequested communication**, often sent via the internet or an electronic messaging service. On mobile, the most obvious types of unwanted text messages are unrecognized numbers and sent by auto-dialers, often promoting a product or service or trying to cheat receiver by unwanted messages reflecting offers and huge benefits.

"Ham" is e-mail that is not Spam!!



2.6 Tools used

Python programming language and frameworks such as NumPy, Pandas, Scikit-learn, Matplotlib, Plotly, Flask etc are used to build the whole model.



Visual Studio Code



- PyCharm is used as IDE
- Visual Studio Code is also used as IDE
- For visualization of the plots, Matplotlib, Seaborn and Plotly are used.
- AWS is used for deployment of the model.
- Heroku is also used for deployment of the model.
- Tableau/Power BI is used for dashboard creation.
- MongoDB is used for DB operations
- Python, Flask is used for backend development
- Github is used as Version Control System.
- Docker is used for containerization.
- CircleCI is used for CI-CD pipeline.

2.7 Constraints

The Spam Ham Detection solution system must be correct enough that it not mislead any message and as automated as possible and users should not be required to know any of the workings.

2.8 Assumptions

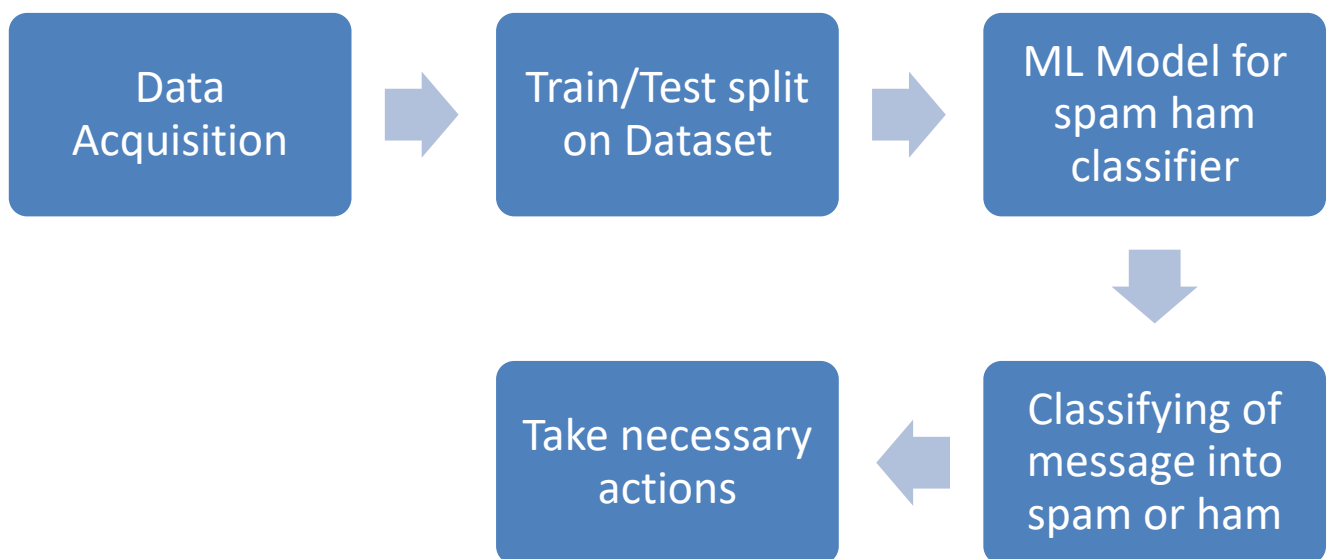
The main objective of the project is to implement the use cases as previously mentioned for new dataset that comes through platform which has this solution install in their device to capture text messages.

3 Design Details

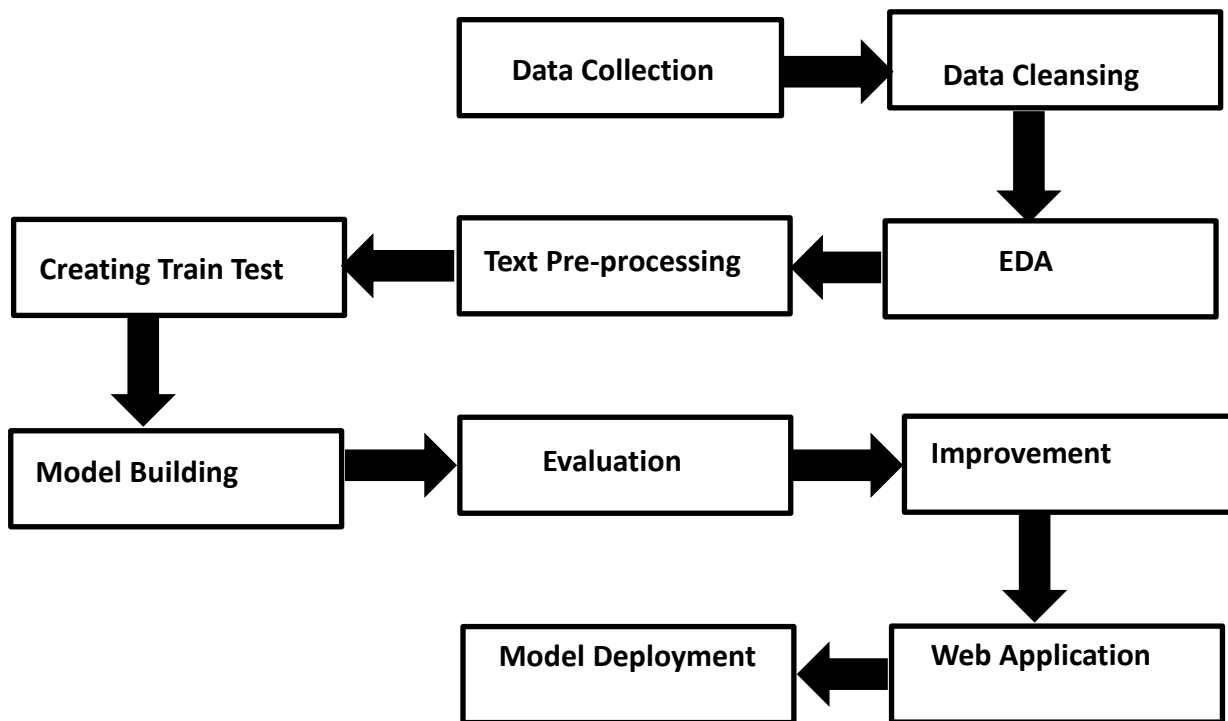
3.1 Process Flow

For detecting spam ham, we will use machine learning base model. Below is the process flow diagram is as shown below

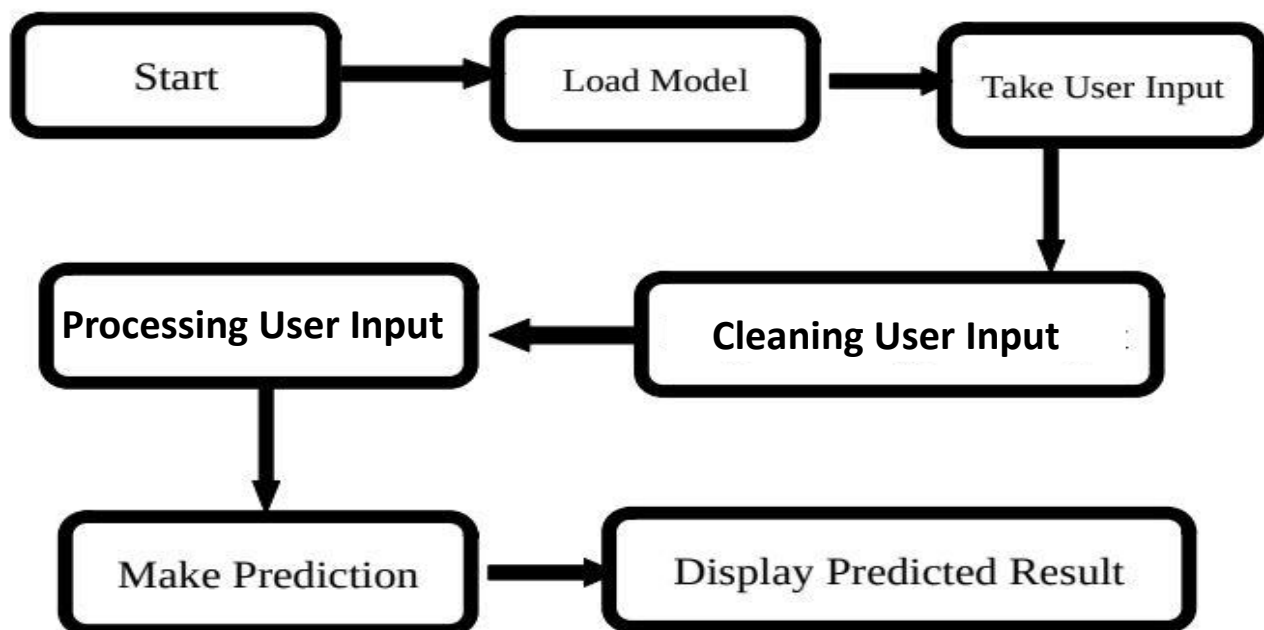
Proposed methodology



3.1.1 Model Training and Evaluation



3.1.2 Deployment Process



3.2 Event log

The system should log every event so that the user will know what process is running internally.

Initial Step-By-Step Description:

1. The System identifies at what step logging required.
2. The System should be able to log each and every system flow.
3. Developer can choose logging method. You can choose database logging/ File logging s well.
4. System should not hang even after using so many loggings. Logging just because we can easily debug issues so logging is mandatory to do.

3.3 Error Handling

Should errors be encountered, an explanation will be displayed as to what went wrong? An error will be defined as anything that falls outside the normal and intended usage.



4 Performance

The machine learning based Spam Ham Classifier solution will use to classify text messages into spam or ham. So that user will not come under false and cheat commitment of spammers. Also model retraining is very important to improve performance.

4.1 Reusability

The code written and the components used should have the ability to be reused with no problems.

4.2 Application Compatibility

The different components for this project will be using python as an interface between them. Each component will have its own task to perform, and it is the job of the Python to ensure proper transfer of information.

4.3 Resource utilization

When any task is performed, it will likely use all the processing power available until that function is finished.

4.4 Deployment



5 Conclusion

Spam Ham Classifier solution will take messaging service data to train our machine learning model and will evaluate its performance over usecases mentioned above. And then leverage its prediction to classify text messages into spam or ham and able to alert user from fraudulent message and spammers. This solution should be as accurate as possible, so that chances of misleading user will be taken good care of.

6 References

[UCI Machine Learning Repository For Data Set](https://archive.ics.uci.edu/ml/datasets/sms+spam+collection)

URL: <https://archive.ics.uci.edu/ml/datasets/sms+spam+collection>