

Machine Learning Engineer Nanodegree

Capstone Project Proposal

Sms spam collection

Sake Revathi

Proposal:

Classifying sms spam collection

Domain Background

History:

Over recent years, as the popularity of mobile phone devices has increased, Short Message Service (SMS) has grown into a multi-billion dollars industry. At the same time, reduction in the cost of messaging services has resulted in growth in unsolicited commercial advertisements (spams) being sent to mobile phones. In parts of Asia, up to 30% of text messages were spam in 2012. Lack of real databases for SMS spams, short length of messages and limited features, and their informal language are the factors that may cause the established email filtering algorithms to underperform in their classification. In this project, a database of real SMS Spams from UCI Machine Learning repository is used, and after preprocessing and feature extraction, different machine learning techniques are applied to the database. Finally, the results are compared and the best algorithm for spam filtering for text messaging is introduced. Final simulation results using 10-fold cross validation shows the best classifier in this work reduces the overall error rate of best model in original paper citing this dataset by more than half.

The purpose of this project is to explore the results of applying machine learning techniques to Message spam detection. Short Message Service spam (sometimes called cell phone spam) is any junk message delivered to a mobile phone as text messaging through the SMS. The dataset for this project originates from the UCI Machine Learning Repository. More detail about dataset can be found on UCI dataset official website.

- This dataset has been collected from free or free for research sources at the Internet.
- The collection is composed of just one text file, where each line has the correct class followed by the raw message.

Problem Statement

In the project, we would try to analysis different methods to identify spam messages. We will use the different approach, based on word count and term-frequency inverse document-frequency transform to classify the messages. Following steps are required in order to achieve the objective:

- Download and pre-process the SMS Spam Collection v.1 dataset.
- Test and find best approach to classify the messages.
- Selection of approach and splitting the dataset into training and testing data.
- Initialize various classifier and train it.
- Evaluate the classifiers and finding best the model for a dataset.

Datasets and Inputs

The Reference Link: <https://www.kaggle.com/uciml/sms-spam-collection-dataset>

The dataset used for this project is SMS Spam Collection dataset originates from the UCI Machine Learning Repository. This dataset has been collected from free or free for research sources at the Internet. The collection is composed of just one text file, where each line has the correct class followed by the raw message. This dataset is tab-separated values (TSV) file. There are total 5572 entries in the dataset and has two column “Class” and “Text” where each row represent different message and Class contain two unique categories ham and spam. Dataset does not require any kind of cleaning, wrangling and there is no null value in any column.

Solution Statement

We are given labelled training data, so this makes it a supervised machine learning problem. For every message, it can be predicted whether it is ham or spam. The accuracy is quantifiable in terms of the f1_score. These performance scores can be compared against the public leader board scores available in the Kaggle website. A well-documented code with dataset will help anyone to replicate the work anywhere on any other machine. To begin with I would like to experiment with techniques which we are going to us are based on word count and term-frequency inverse document-frequency transform. After which I would like to test the approach using many different algorithms like Naive Bayes, Decision Tree, AdaBoost, K-Nearest Neighbours and Random Forest and test the accuracy score.

Benchmark Model

Benchmark models are available in Kaggle discussion forums which uses different Machine Learning algorithms. The available public and private leader board score in the Kaggle competition can be used to benchmark the performance of my algorithm. Also, it is possible to explore how the proposed model perform compared to existing models. The result shows that Naïve Bayes work better on the dataset with an accuracy_score of 0.70.

Evaluation Metrics

Accuracy is the first metric to be checked when the algorithms are evaluated, is the sum of true positives and the true negative outputs divided by the data size. Accuracy means how closer you are to the true value, whereas precision means that your data points are not widely spread. The Scikit-learn library provides a convenience report when working on classification problems to give you a quick idea of the accuracy of a model using a number of matrices, one of them is F1_score which work with all the model.

Project Design

The theoretical workflow of the project would look like:

- Download and pre-process the SMS Spam Collection v.1 dataset.
- Test and find best approach to classify the messages
- Selection of approach and splitting the dataset into training and testing data.
- Initialize various classifier and train it using training data.s
- Evaluate the classifiers and finding best the model for a dataset using testing data.