

# **Text Classification: Spam or Ham**

## **NLP Project (UML602)**

**Submitted By:**

Kashish Bansal – 101603155

Kashish Mittal – 101603156

**BE Third Year, COE**

**Submitted to:**

**Ms. Diksha Hooda**



**Computer Science and Engineering Department**

**TIET, Patiala**

**April 2019**

## Index

<b>S.No.</b>	<b>Topic</b>	<b>Page No.</b>
1	Introduction	1
2	Steps of working	3
3	Results	5
4	Applications	6
5	Future Scope	6
6	Conclusion	6
7	References	7

# 1. Introduction

Mobile spam is an increasing threat that may be addressed using filtering systems like those employed against email spam. We believe that email filtering techniques require some adaptation to reach good levels of performance on SMS spam, especially regarding message representation. In order to test this assumption, we have performed experiments on SMS filtering using email spam filters on mobile spam messages using a suitable feature representation, with results supporting our hypothesis.

A spam filter is a service feature designed to block spam from a user's inbox. Because a large amount of global messages are spam, effective spam filters are critical to maintaining clean and spam-free in-boxes. By keeping pesky spam messages away, spam filters increase user efficiency by sparing tedious manual sifting of legitimate messages and spam email deletion.

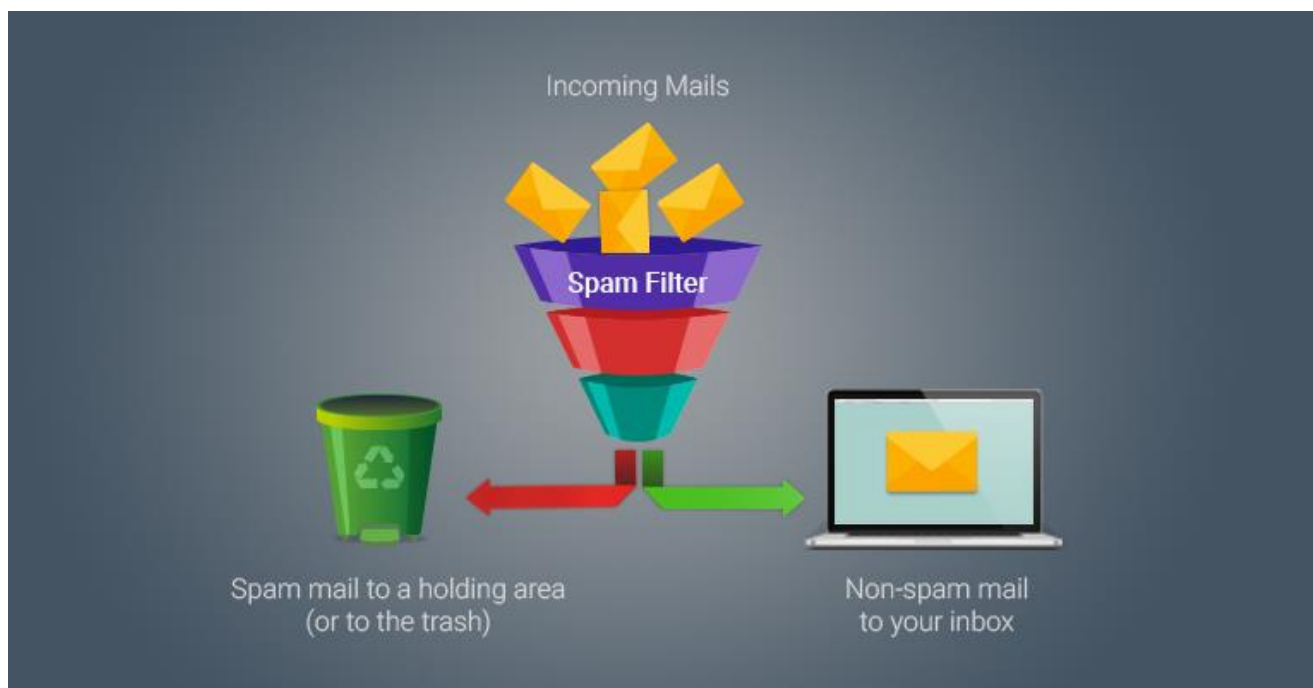


Fig:1 Spam Filter

In this, we will expand on this foundation and explore different ways to improve our text classification results. We will cover and use:

- Regular Expressions
- Feature Engineering
- Multiple scikit-learn Classifiers
- Ensemble Methods

## 1.1 Regular Expression:

Regular expressions are used to specify rules for matching strings of text. In other words, you can use regular expressions to match text in incoming messages. The settings in which you use the regular expressions will tell filter what to do with messages that have matching text.

## 1.2 Feature Engineering:

Feature engineering is the process of using domain knowledge of the data to create features that make machine learning algorithms work. If feature engineering is done correctly, it increases the predictive power of machine learning algorithms by creating features from raw data that help facilitate the machine learning process.

## 1.3 Scikit-learn

Scikit-learn provides a range of supervised and unsupervised learning algorithms via a consistent interface in Python. It is licensed under a permissive simplified BSD license and is distributed under many Linux distributions, encouraging academic and commercial use. The library is built upon the SciPy (Scientific Python) that must be installed before you can use Scikit-learn. This stack that includes:

- **NumPy**: Base n-dimensional array package
- **SciPy**: Fundamental library for scientific computing
- **Matplotlib**: Comprehensive 2D/3D plotting
- **IPython**: Enhanced interactive console
- **Sympy**: Symbolic mathematics
- **Pandas**: Data structures and analysis

Extensions or modules for SciPy are conventionally named SciKits. As such, the module provides learning algorithms and is named Scikit-learn.

## 1.4 Ensemble Method:

Ensemble learning helps improve machine learning results by combining several models. This approach allows the production of better predictive performance compared to a single model. That is why ensemble methods placed first in many prestigious machine learning competitions, such as the Netflix Competition, KDD 2009, and Kaggle.

## 2. Steps of working

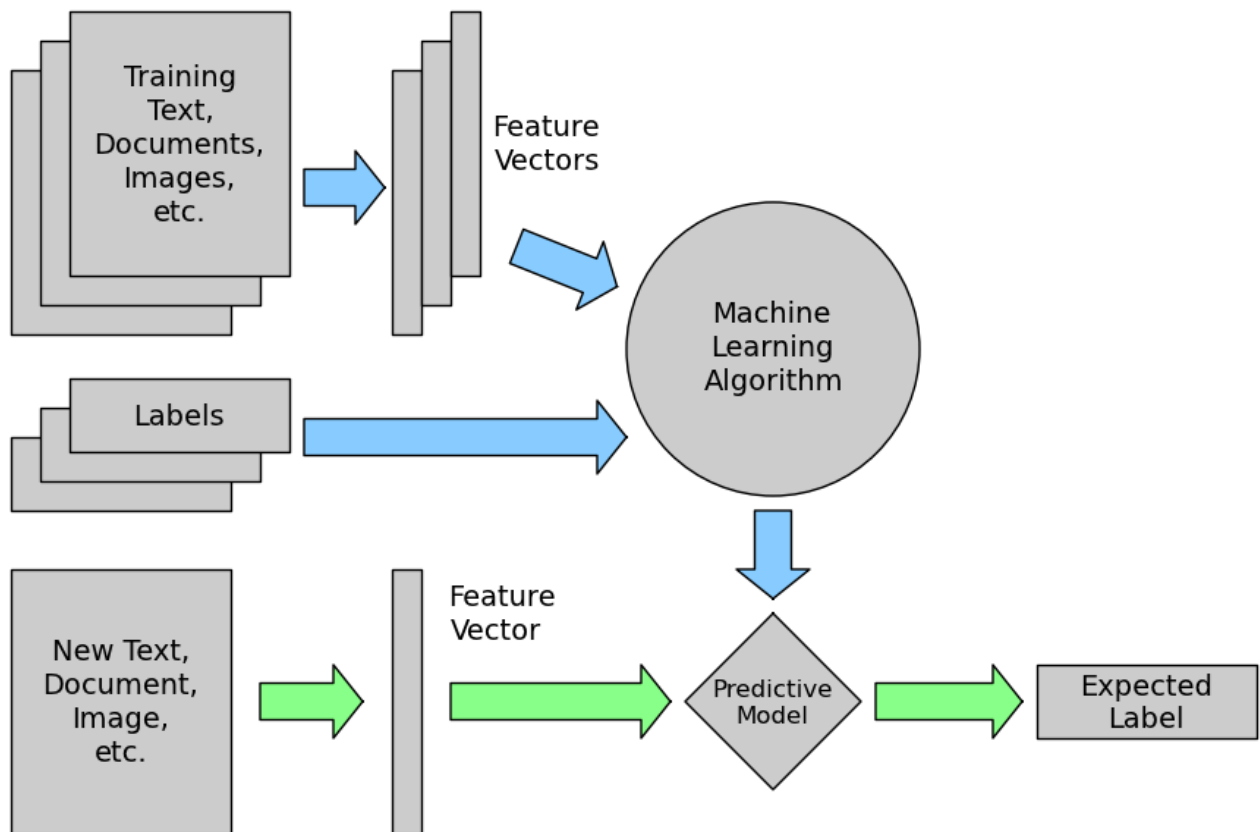


Fig 2: Basic Model

### 2.1 Import Necessary Libraries

To ensure the necessary libraries are installed correctly and up-to-date, print the version numbers for each library. This will also improve the reproducibility of our project.

Python: 2.7.13 [Continuum Analytics, Inc.] (default, May 11 2017, 13:17:26) [MSC v.1500 64 bit (AMD64)]

NLTK: 3.2.5

Scikit-learn: 0.19.1

Pandas: 0.21.0

Numpy: 1.14.1

### 2.2 Load the Dataset

Now that we have ensured that our libraries are installed correctly, let's load the data set as a Pandas DataFrame. Furthermore, let's extract some useful information such as the column information and class distributions.

The data set we will be using comes from the UCI Machine Learning Repository. It contains over 5000 SMS labeled messages that have been collected for mobile phone spam research. It can be downloaded from the following URL:

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

## 2.3 Preprocess the Data

Preprocessing the data is an essential step in natural language process. In the following cells, we will convert our class labels to binary values using the LabelEncoder from sklearn, replace email addresses, URLs, phone numbers, and other symbols by using regular expressions, remove stop words, and extract word stems.

### 2.3.1 Regular Expressions

Some common regular expression metacharacters - copied from wikipedia

**^** Matches the starting position within the string. In line-based tools, it matches the starting position of any line.

**.** Matches any single character (many applications exclude newlines, and exactly which characters are considered newlines is flavor-, character-encoding-, and platform-specific, but it is safe to assume that the line feed character is included). Within POSIX bracket expressions, the dot character matches a literal dot. For example, `a.c` matches `"abc"`, etc., but `[a.c]` matches only `"a"`, `"."`, or `"c"`.

**[ ]** A bracket expression. Matches a single character that is contained within the brackets. For example, `[abc]` matches `"a"`, `"b"`, or `"c"`. `[a-z]` specifies a range which matches any lowercase letter from `"a"` to `"z"`. These forms can be mixed: `[abcx-z]` matches `"a"`, `"b"`, `"c"`, `"x"`, `"y"`, or `"z"`, as does `[a-cx-z]`. The `-` character is treated as a literal character if it is the last or the first (after the `^`, if present) character within the brackets: `[abc-]`, `[-abc]`. Note that backslash escapes are not allowed. The `]` character can be included in a bracket expression if it is the first (after the `^`) character: `[ ]abc]`.

**[^ ]** Matches a single character that is not contained within the brackets. For example, `[^abc]` matches any character other than `"a"`, `"b"`, or `"c"`. `[^a-z]` matches any single character that is not a lowercase letter from `"a"` to `"z"`. Likewise, literal characters and ranges can be mixed.

**\$** Matches the ending position of the string or the position just before a string-ending newline. In line-based tools, it matches the ending position of any line.

**( )** Defines a marked subexpression. The string matched within the parentheses can be recalled later (see the next entry, `\n`). A marked subexpression is also called a block or capturing group. BRE mode requires `( )`.

**\n** Matches what the `n`th marked subexpression matched, where `n` is a digit from 1 to 9. This construct is vaguely defined in the POSIX.2 standard. Some tools allow referencing more than nine capturing groups.

**\*** Matches the preceding element zero or more times. For example, `abc` matches `"ac"`, `"abc"`, `"abbbc"`, etc. `[xyz]` matches `""`, `"x"`, `"y"`, `"z"`, `"zx"`, `"zyx"`, `"xyzzy"`, and so on. `(ab)*` matches `""`, `"ab"`, `"abab"`, `"ababab"`, and so on.

**{m,n}** Matches the preceding element at least `m` and not more than `n` times. For example, `a{3,5}` matches only `"aaa"`, `"aaaa"`, and `"aaaaa"`. This is not found in a few older instances of regexes. BRE mode requires `{m,n}`.

## 2.4 Generating Features

Feature engineering is the process of using domain knowledge of the data to create features for machine learning algorithms. In this project, the words in each text message will be our features. For this purpose, it will be necessary to tokenize each word. We will use the 1500 most common words as features.

## 2.5 Scikit-Learn Classifiers with NLTK

Now that we have our dataset, we can start building algorithms! Let's start with a simple linear support vector classifier, then expand to other algorithms. We'll need to import each algorithm we plan on using from sklearn. We also need to import some performance metrics, such as `accuracy_score` and `classification_report`.

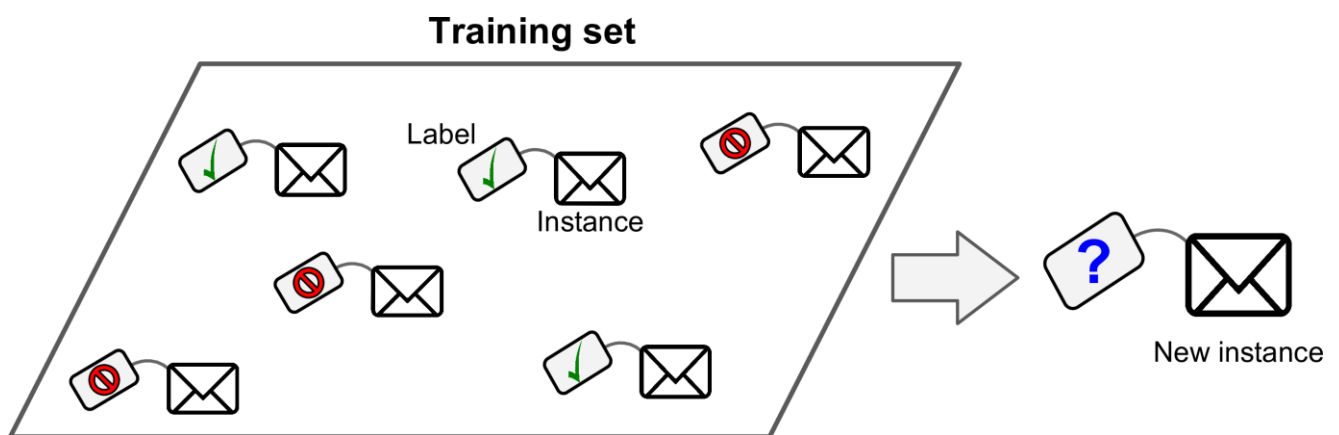


Fig 3: Training

## 3. Results

The dataset is trained with various models and the following results were obtained:

- K Nearest Neighbors Accuracy: 93.1083991385
- Decision Tree Accuracy: 97.4874371859
- SGD Classifier Accuracy: 97.9181622398
- Naive Bayes Accuracy: 97.5592246949
- SVM Linear Accuracy: 98.1335247667

Then the model was trained with ensemble classifier in which voting of all the above models were done and following result was obtained:

Voting Classifier: Accuracy: 98.1335247667

		predicted	
		ham	spam
actual	ham	1198	0
	spam	27	168

Fig: Confusion Matrix

## 4. Applications

- This model will help the company to apply for model for predicting different SPAM/HAM SMS in the text analytics market.
- This tools help understand the working knowledge of the different technique which are used in analytics.
- It also involves the design of a Text Classification Model to perform predictive analysis.

## 5. Future Scope

- Fake news detection
- Sarcasm detection
- Classify large data in short time which take more time if read manually
- Tagging content as categories

## 6. Conclusion

Automatic Text Classification is a machine learning technique. Document can be set to predefined categories best on textual content and extraction features. It has important applications in spam filtering and text mining. In the recent years the automatic categorization of texts into predefined categories is booming interest. Due to the increased availability of data & documents in the day to day basics in digital form and ensuing the need to organize them. In the research community the commanding, approach to this problem is based on machine learning techniques. The advantages of this approach efficiency, considerable savings in terms of expert manpower and straightforward portability to different domains. This report emphasizes on the text categorization that fall within the machine learning technique. How Automatic Text Classification can be used to classify SMS SPAM filtering classification model. An analysis of SMS SPAM filtering classification model has also been done using Automatic Text Classification.



## 7. References

- [https://github.com/eduonix/nlptextclassification/blob/master/NLP%20for%20Text%20Classification%20\(Jupyter%20Notebook\).ipynb](https://github.com/eduonix/nlptextclassification/blob/master/NLP%20for%20Text%20Classification%20(Jupyter%20Notebook).ipynb)
- <https://scikit-learn.org/stable/documentation.html>
- <https://medium.com/swlh/classify-emails-into-ham-and-spam-using-naive-bayes-classifier-ffddd7faalef>

