

# Naive Bayes Spam Filter

Ruchira Patil<sup>1</sup>, Prathamesh Patil<sup>2</sup>, Mandar Patil<sup>3</sup>

<sup>1,2,3</sup>Student, Department of Computer Engineering, MGM CET, Navi Mumbai, India

**Abstract**—Spammers always find new ways to get spammy content to the public. Very commonly this is accomplished by using email, social media, or advertisements. Spam filters have been getting better at detecting spam and removing it, but no method is able to block 100% of it. Because of this, many different methods of text classification have been developed, including a group of classifiers that use a Bayesian approach. The Bayesian approach to spam filtering was one of the earliest methods used to filter spam, and it remains relevant to this day. Managing uncommon words, for the situation a word has never been met amid the learning stage, both the numerator and the denominator are equivalent to zero, both in the general equation and in the spamicity recipe. The product can choose to dispose of such words for which there is no data accessible. Words that typically show up in substantial amounts in spam may likewise be changed by spammers. For instance, «Viagra» would be supplanted with «Viaagra» or «V!agra» in the spam message.

**Index Terms**—Spam, Bayesian Filtering, Naive Bayes, Multinomial Bayes, Multivariate Bayes

## I. INTRODUCTION

Naive Bayes spam shifting is a gauge system for managing spam that can tailor itself to the email needs of individual clients and give low false positive spam location rates that are for the most part adequate to clients. It is one of the most established methods for doing spam separating, which was established in the 1990s

The Naive Bayesian classifier is based on Bayes' theorem with the independence assumptions between predictors. A Naive Bayesian model is easy to build, with no complicated iterative parameter estimation which makes it particularly useful for very large datasets. Despite its simplicity, the Naive Bayesian classifier often does surprisingly well and is widely used because it often outperforms more sophisticated classification methods.

## II. HISTORY

Although naive Bayesian filters did not become popular until later, multiple programs were released in 1998 to address the growing problem of unwanted email. The first scholarly publication on Bayesian spam filtering was by Sahami et al. in 1998. That work was soon thereafter deployed in commercial spam filters. However, in 2002 Paul Graham greatly decreased the false positive rate, so that it could be used on its own as a single spam filter.

## III. PROCESS

Specific words have specific probabilities of happening in spam email and in genuine email. For example, most email clients will oftentimes experience "Viagra" in spam email, yet will only sometimes observe it in other email. The filter doesn't know these probabilities ahead of time, and should first be prepared so it can develop them. To prepare the filter, the client should physically demonstrate whether another email is spam or not. For all words in each preparation email, the channel will alter the probabilities that each word will show up in spam or genuine email in its database. For example, Bayesian spam channels will normally have taken in a high spam likelihood for the words "Viagra" and "renegotiate", however a low spam likelihood for words seen just in genuine email, for example, the names of loved ones.

In the wake of preparing, the word probabilities (otherwise called probability capacities) are utilized to process the likelihood that an email with a specific arrangement of words in it has a place with either class. Each word in the email adds to the email's spam likelihood, or just the most intriguing words. This commitment is known as the back likelihood and is figured utilizing Bayes' hypothesis. At that point, the email's spam likelihood is processed over all words in the email, and if the aggregate surpasses a specific edge (say 95%), the channel will stamp the email as a spam.

As in some other spam sifting strategy, email set apart as spam would then be able to be naturally moved to a "Garbage" email envelope, or even erased inside and out. Some product execute isolate components that characterize a time period amid which the client is permitted to audit the product's choice.

The underlying preparing can generally be refined when wrong judgements from the product are distinguished (false positives or false negatives). That enables the product to powerfully adjust to the consistently advancing nature of spam.

Some spam channels join the consequences of both Bayesian spam separating and different heuristics (pre-characterized leads about the substance, taking a gander at the message's envelope, and so forth.), bringing about significantly higher sifting exactness, at times at the cost of adaptiveness.

## IV. MATHEMATICAL FOUNDATION

*Naive Bayes Theorem:*

The Naive Bayesian classifier is based on Bayes' theorem with the independence assumptions between predictors. A Naive Bayesian model is easy to build, with no complicated

iterative parameter estimation which makes it particularly useful for very large datasets.

$$P(c|x) = \frac{P(x|c)P(c)}{P(x)}$$

Likelihood
Class Prior Probability  
Posterior Probability
Predictor Prior Probability

$$P(c|X) = P(x_1|c) \times P(x_2|c) \times \dots \times P(x_n|c) \times P(c)$$

**Algorithm:**

Bayes theorem provides a way of calculating the posterior probability,  $P(c|x)$ , from  $P(c)$ ,  $P(x)$ , and  $P(x|c)$ . Naive Bayes classifier assume that the effect of the value of a predictor ( $x$ ) on a given class ( $c$ ) is independent of the values of other predictors. This assumption is called class conditional independence.

- $P(c|x)$  is the posterior probability of class (target) given predictor (attribute).
- $P(c)$  is the prior probability of class.
- $P(x|c)$  is the likelihood which is the probability of predictor given class.
- $P(x)$  is the prior probability of predictor.

*Computing the probability that a message containing a given word is spam*

How about we assume the speculated message contains "reproduction". A great many people who are accustomed to getting email realize that this message is probably going to be spam, all the more accurately a proposition to offer fake duplicates of surely understood brands of watches. The spam discovery programming, nonetheless, does not "know" such realities; everything it can do is process probabilities.

The equation utilized by the product to establish that, is gotten from Bayes' hypothesis

$$Pr(S|W) = \frac{Pr(W|S) \cdot Pr(S)}{Pr(W|S) \cdot Pr(S) + Pr(W|H) \cdot Pr(H)}$$

Where:

$Pr(S|W)$  is the probability that a message is a spam, knowing that the word "replica" is in it;

$Pr(S)$  is the overall probability that any given message is spam;

$Pr(W|S)$  is the probability that the word "replica" appears in spam messages;

$Pr(H)$  is the overall probability that any given message is not spam (is "ham");

$Pr(W|H)$  is the probability that the word "replica" appears in ham messages.

Most bayesian spam filtering calculations depend on recipes that are entirely legitimate (from a probabilistic angle) just if the words introduce in the message are autonomous occasions. This condition isn't by and large fulfilled (for instance, in

normal dialects like English the likelihood of finding a modifier is influenced by the likelihood of having a thing), however it is a helpful admiration, particularly since the measurable relationships between's individual words are generally not known. On this premise, one can get the accompanying equation from Bayes' hypothesis.

$$p = \frac{p_1 p_2 \dots p_N}{p_1 p_2 \dots p_N + (1 - p_1)(1 - p_2) \dots (1 - p_N)}$$

Where:

$Pr(S|W)$  is the probability that a message is a spam, knowing that the word "replica" is in it;

$Pr(S)$  is the overall probability that any given message is spam;

$Pr(W|S)$  is the probability that the word "replica" appears in spam messages;

$Pr(H)$  is the overall probability that any given message is not spam (is "ham");

$Pr(W|H)$  is the probwhere

**V. COMBINING INDIVIDUAL PROPERTIES**

$p$  is the probability that the suspect message is spam;

$p\{1\}$  is the probability  $p(S|W_{\{1\}})$  that it is a spam knowing it contains a first word (for example "replica");

$p\{2\}$  is the probability  $p(S|W_{\{2\}})$  that it is a spam knowing it contains a second word (for example "watches"); etc...

$p\{N\}$  is the probability  $p(S|W_{\{N\}})$  that it is a spam knowing it contains an Nth word (for example "home").probability that the word "replica" appears in ham messages.

**VI. OTHER HEURISTICS**

Impartial words like "the", "an", "a few", or "is" (in English), or their counterparts in different dialects, can be overlooked. All the more for the most part, some bayesian sifting channels basically disregard every one of the words which have a spamicity alongside 0.5, as they contribute little to a decent choice. The words mulled over are those whose spamicity is by 0.0 (unmistakable indications of honest to goodness messages), or beside 1.0 (particular indications of spam). A strategy can be for instance to keep just those ten words, in the inspected message, which have the best supreme esteem  $|0.5 - pI|$ .

Some product items consider the way that a given word seems a few times in the inspected message, [15] others don't.

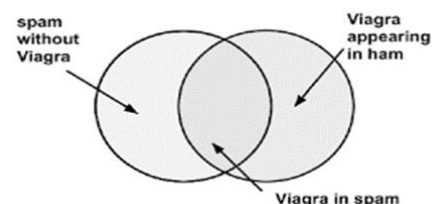


Fig. 1. Venn diagram

Some product items utilize designs (arrangements of words) rather than disengaged normal dialects words.[16] For instance, with a "setting window" of four words, they process the spamicity of "Viagra is useful for", rather than registering the spamicities of "Viagra", "is", "great", and "for". This technique gives greater affectability to setting and kills the Bayesian commotion better, to the detriment of a greater database.

#### VII. DISADVANTAGES OF THE OLD SYSTEM

Contingent upon the execution, Bayesian spam sifting might be helpless to Bayesian harming, a system utilized by spammers trying to debase the viability of spam channels that depend on Bayesian separating. A spammer rehearsing Bayesian harming will convey messages with a lot of honest to goodness content (assembled from genuine news or abstract sources). Spammer strategies incorporate addition of irregular harmless words that are not ordinarily connected with spam, consequently diminishing the email's spam score, making it more inclined to slip past a Bayesian spam channel. Notwithstanding, with (for instance) Paul Graham's plan just the most huge probabilities are utilized, so cushioning the content out with non-spam-related words does not influence the discovery likelihood altogether.

Words that typically show up in substantial amounts in spam may likewise be changed by spammers. For instance, «Viagra» would be supplanted with «Viaagra» or «V!agra» in the spam message. The beneficiary of the message can in any case read the changed words, however every one of these words is met all the more seldom by the Bayesian channel, which frustrates its learning procedure. When in doubt, this spamming method does not work extremely well, on the grounds that the inferred words wind up perceived by the channel simply like the typical ones.

Another procedure used to attempt to crush Bayesian spam channels is to supplant content with pictures, either specifically included or connected. The entire content of the message, or some piece of it, is supplanted with a photo where a similar content is "drawn". The spam channel is typically unfit to break down this photo, which would contain the touchy words like «Viagra». In any case, since many mail customers debilitate the show of connected pictures for security reasons, the spammer sending connects to far off pictures may achieve less targets. Additionally, a photo's size in bytes is greater than the proportionate content's size, so the spammer needs more transfer speed to send messages straightforwardly including pictures. A few channels are more disposed to choose that a message is spam in the event that it has for the most part graphical substance. An answer utilized by Google in its Gmail email framework is to play out an OCR (Optical Character Recognition) on each mid to extensive size picture, examining the content inside.

#### VIII. IMPROVED SYSTEM

Managing uncommon words, for the situation a word has never been met amid the learning stage, both the numerator and

the denominator are equivalent to zero, both in the general equation and in the spamicity recipe. The product can choose to dispose of such words for which there is no data accessible.

All the more for the most part, the words that were experienced just a couple of times amid the learning stage cause an issue, since it would be a mistake to trust aimlessly the data they give. A basic arrangement is to just abstain from considering such problematic words also.

Applying again Bayes' hypothesis, and accepting the grouping amongst spam and ham of the messages containing a given word ("reproduction") is an irregular variable with beta dissemination, a few projects choose to utilize a remedied likelihood:

$$\Pr'(S|W) = \frac{s \cdot \Pr(S) + n \cdot \Pr(S|W)}{s + n}$$

Where:

$\Pr'(S|W)$  is the amended likelihood for the message to be spam, realizing that it contains a given word;

$s$  is the quality we provide for foundation data about approaching spam ;

$\Pr(S)$  is the likelihood of any approaching message to be spam;

$n$  is the quantity of events of this word amid the learning stage;

$\Pr(S|W)$  is the spamicity of this word.

This rectified likelihood is utilized rather than the spamicity in the joining recipe.

$\Pr(S)$  can again be taken equivalent to 0.5, to abstain from being excessively suspicious about approaching email. 3 is a decent incentive for  $s$ , implying that the educated corpus must contain in excess of 3 messages with that word to put more trust in the spamicity esteem than in the default value.

This equation can be reached out to the situation where  $n$  is equivalent to zero (and where the spamicity isn't characterized), and assesses for this situation to  $\Pr(S)$ .

#### IX. CONCLUSION

The detection of spam at a place close to the sending server is an important issue in the network security and machine learning techniques have a very important role in this topic. In this paper, applied again Bayes' hypothesis, and accepting the grouping amongst spam and ham of the messages containing a given word ("reproduction") is an irregular variable with beta dissemination. We also have managed uncommon words, for the situation a word has never been met amid the learning stage, both the numerator and the denominator are equivalent to zero, both in the general equation and in the spamicity recipe.

#### REFERENCES

- [1] Messaging anti-abuse working group report: First, second and third quarter 2011. November 2011.
- [2] Why bayesian filtering is the most effective anti-spam technology. 2011.
- [3] Messaging anti-abuse working group, "Email metrics report", <http://www.maawg.org/>, 2006.

- [4] Shlomo Hershkop, "Behaviour-based email analysis with applications to spam detection", Ph. D. Thesis, <http://www1.cs.columbia.edu/sh553/publicaions/>, 2006.
- [5] Wikipedia. Naive bayes spam \_ltering | wikipedia, the free encyclopedia, 2014.