# Comparison of anomaly detection-based spam email filters

Aleksej Milošević

Univerza v Mariboru, Fakulteta za elektrotehniko, računalništvo in informatiko
Smetanova 17, 2000 Maribor

aleksej.milosevic@student.um.si

**Napredek komunikacijske tehnologije vsakodnevno povezuje veliko število ljudi, vendar prav tako botruje čedalje večji količini odvečne elektronske pošte. Napadalci razvijajo napredne metode izogibanja filtrom odvečne pošte, na katere je potrebno odgovoriti z izpolnjevanjem in dopolnjevanjem obstoječih tehnik filtriranja. V članku predstavimo zaznavanje odvečne pošte z algoritmi zaznavanja anomalij osnovanim na strojnem učenju, katere glavna prednost je zmanjšana potreba po pred-označenih podatkih. V okviru raziskave smo primerjali tri algoritme zaznavanja anomalij: Simple Anomaly Detection Method, One-Class SVM in Isolation Forest. Ugotovili smo, da je optimalna metoda zaznavanja odvečne elektronske pošte odvisna od podatkovne množice in izbire bodisi legitimnih bodisi odvečnih dokumentov za zastopanje normalnosti.**


**With the widespread availability of broadband internet and affordable connected devices, spam email is a more prominent threat than ever before, incurring billions of dollars in costs yearly. Attackers use increasingly advanced methods of spam filter evasion techniques which forces the cybersecurity community to respond with innovative techniques to supplement existing countermeasures. One such technique is machine learning-based anomaly detection which reduces the need for pre-labelled data. In this paper we compare the effectiveness of three anomaly detection methods: Simple Anomaly Detection, One-Class SVM and Isolation Forest. We have found that the optimal method of anomaly detection depends on the dataset and the choosing legitimate or spam documents as representation of normality.**

## 1 Introduction

Electronic mail or e-mail is one of the most important and useful channels of communication that exist today because it enables cheap and efficient communication between interested parties with zero overhead to the participants. Spam, according to Geerthik is "an endless repetition of worthless text or image"(Geerthik, 2013). Email can be easily exploited because of the inherently unsafe nature of the Simple Mail Transfer Protocol or SMTP. The protocol has no internal mechanism to secure the contents and routing or to validate the sender (Geerthik, 2013)

According to the Federal Trade Commission, two thirds of all emails contain misleading information. The dangers of spam ranges from a minor annoyance such as having to delete the email landing unrestricted to an individual's inbox, to massive economic risks. According to calculations made by Rao and Reiley spam email costs between $14-18 billion per year while only generating a revenue of $160-360 million per year (Rao & Reiley, 2012).

 Mitigating spam mail is an evolving problem that requires innovative solutions. To this date social methods of fighting spam have proven cumbersome and ineffective and as such technical countermeasures were developed to ensure efficient protection. Before the rise of machine learning based spam filtering. filters followed a so called *'knowledge engineering'* approach which utilized a content-based heuristic filter. Such filters analyzed messages and classified them as either *'spam'* or *'ham'* based on the knowledge of regularities or certain observed patterns such as those defined by Guzella and Caminhas. (Guzella & Caminhas, 2009).

Heuristic filters, however, have a considerable drawback in the form of cumbersome and labor intensive updating of the rule-set in order to keep up with the newest trends. An example of such trends is the obfuscation of popular terms with special characters (e.g. "V!@g®@" instead of "Viagra"). Filters are supplemented by a plethora of other mitigation methods. Blacklisting involves keeping a list of e-mail or IP addresses of servers from which spam mail is likely to originate. A whitelist is the natural reverse of blacklists where an email server or a client keeps a list of trusted e-mail addresses or server IP's that are allowed communicate with the user. If an incoming message is from an IP or email address that is neither on the blacklist or the whitelist, the recipients email server may temporarily reject the message. If the email is in fact legitimate, the sender will attempt, after some time. to send the message again upon which the recipient will accept the email as legitimate and add the sender's credentials to a whitelist (Bhowmick & Hazarika, 2016).

Modern spam filters have evolved to use machine learning predictive models which are by their very nature adaptive and require less labor and manpower to maintain. Most machine learning and statistical approaches to filtering spam are supervised. This requires a prelabelled training set which is used to train an algorithm and build a predictive model. It is imperative that the number of instances in the training set is enough to reach satisfactory model accuracy. The biggest problem utilizing supervised predictive models is obtaining the large amount of labelled data necessary (Santos, Laorden, Ugarte-Pedrero, Sanz, & Bringas, 2011).

In this paper we compare three anomaly detection-based solutions. Section 2 describes in detail the structure of text-based spam filters. Section 3 describes the implementation of a vector space model anomaly detection method. Section 4 describes the execution of the experiment and displays the generated data. Section 5 compares and discusses classification results of three methods; implemented vector space model anomaly detection method, isolation forest and one-class SVM. The sixth and final section concludes the paper and proposes further research.

## 2    Spam filter structure

An email message is divided into a header, which contains the metadata of an email and the body which contains the actual message of the sender to the receiver. Before executing classification, the contents of the email need to be pre-processed in a form that is appropriate as input to the classifier. The pre-processing is conducted in the following steps (Bhowmick & Hazarika, 2016) :

1. Tokenization, which is the process of extracting individual words from the email
2. Lemmatization, which reduces the words to their root form
3. Stop-word removal, which removes the most common words in the English language (e.g is, and, the)
4. Representation conversion, which converts words into a machine-readable format

### 2.1    Vector space model representation

Emails are represented in the form of an Information Retrieval model. Specifically. we use a term frequency-inverse document frequency, or tf-idf model.

The tf-idf is a product of term frequency and inverse document frequency. This is calculated by using the number of times a term *t* occurs in document *d*, where the term frequency of term *j* in document *i* is defined as (Laorden et al., 2014) :

$$tf_{i,j} = \frac{m_{i.j}}{\sum_k m_{k,j}}$$

- $m_{i,j}$ is the number of times the word $t_{i,j}$ appears in a document
- $\sum_k m_{k,j}$ is the number of words in a document

The inverse document frequency is defined as:

$$idf_i = \frac{|\varepsilon|}{|\varepsilon: t_i \in e|}$$

- $|\varepsilon|$ is the total number of documents containing the word $t_{i,j}$
- $|\varepsilon: t_i \in e|$ is the number of documents containing the word $t_{i,j}$

The term frequency – Inverse document frequency is then calculated as:

$$tfidf\ (i,j) = \ tf_{i,j} * idf_i$$

# 3   Anomaly detection

Anomaly detection is the process which enables finding of patterns in data that does not conform to expected behavior. As such, anomalies are patterns that do not conform to a notion of what is considered normal behavior (Chandola, Banerjee, & Kumar, 2007). The precise method of detecting anomalies differs from method to method.

An isolation forest explicitly identifies anomalous data points. Each tree conducts its search by isolating observations based on a randomly selected feature and randomly selected split value between the maximum and minimum of the feature. If individual trees collectively produce shorter path lengths, they are likely to be anomalies. The score anomaly score is defined as (Liu & Ting, 2018):

$$s(x,n) = \ 2^{-\frac{E(h(x))}{c(n)}}$$

- $h(x)$ is the path length of observation $x$
- $c(n)$ is the average path length of unsuccessful search in a Binary Search Tree
- $n$ is the number of external nodes

If the individual score of an observation is close to 1 it strongly indicates an anomaly. If the individual score is low, it indicates a normal observation. Uniform scores around 0,5 indicate that the observations are relatively uniformly distributed and that there are no clear anomalies.

A One-Class SVM is a variant of the Support Vector Machine algorithm which defines a hyperplane that separates normal data from anomalous data. The binary function captures a region in the area of a certain high enough probability density, using the chosen kernel function. For each point, the function returns +1 in a small region, which represents the training data points and -1 elsewhere. (Manevitz, 2001)

## 3.1   Simple anomaly detection method

Our implemented method as described in detail by Laorden et. al in the paper *Study on the effectiveness of anomaly detection for spam filtering*, analyses points in a feature space by extracting features from the document and measuring the distance from the point representing the email to the points that symbolize normality. Normality is represented by the model created

using the training data. The distance between two points is defined by Euclidean distance (Laorden et al., 2014) :

$$d(x,y) = \sum_{i=0}^{n} \sqrt{x_i^2 - y_i^2}$$

- $x$ is the first point
- $y$ is the second point
- $x_i$ is the $i$-th component of $x$
- $y_i$ is the $i$-th component of $y$

The result of the training method is a final distance value which considers every measure performed. The paper recommends three distinct metrics, also called combination rules:

- The mean value calculated from every distance in the training set
- The lowest distance value from every distance value in the training set
- The highest distance from every distance value in the training set

The algorithm establishes 10 different thresholds with which to determine whether the email is spam or not. The thresholds are chosen by first establishing the lowest and the highest one. The lowest threshold is determined as the lowest possible value at which no legitimate spam messages were misclassified. The highest threshold is determined as the lowest possible value at which no examples of legitimate mail were misclassified. For each of the thresholds three metrics are calculated; False Negative Ratio (FNR), False Positive Ratio (FPR) and the Weighted Accuracy (WA). False Negative Ratio is defined as (Demsar, 2006):

$$FNR(\beta) = \frac{FN}{FN + TP}$$

- FN is the number of documents misclassified as legitimate
- TP is the number of documents correctly classified as spam

False Positive Ratio is defined as:

$$FPR(\beta) = \frac{FP}{FP + TN}$$

- FP is the number of documents misclassified as spam
- TN is the number of correctly classified legitimate messages

Weighted Accuracy is defined as:

$$WA(\beta) = 1 - \frac{FNR + FPR}{2}$$

## 4    Experiment

We have utilized data from the LingSpam Corpus and EnronSpam public Corpus. The LingSpam Corpus contains messages originated from the *Linguistic list,* which is an email distribution about linguistics. The EnronSpam public corpus contains messages originating from the senior management of the now defunct company Enron. We have opted for the 'bare' datasets as opposed to the pre-processed ones. Due to computational limitations we have trimmed the dataset to a manageable size as seen in the following table:

Table 1 Number of documents in dataset per type

|  | LingSpam | EnronSpam |
|---|---|---|
| # of legitimate emails | 2412 | 8033 |
| # of spam emails | 481 | 2966 |

We have removed alpha-numeric and non-letter characters on both datasets. We have then performed stop word removal and stemming and transformed the datasets into a tf-idf matrix. For each dataset we have performed 10-fold Monte Carlo cross-validation in the ratio of 0.66:0.33 in favor of the training set (Tech & Gomez-gil, 2013). We have performed two sets of tests. The first set of models was trained on 'ham' emails with spam emails being regarded as anomalies. The second set of models was trained on spam emails with legitimate emails being regarded as anomalies.

## 4.1 Models trained on legitimate emails

Let us first consider the performance of algorithms on the LingSpam Corpus. Table 2 shows the classification results of IsolationForest on the LingSpam Corpus. The model has a disproportionally high FNR and a low FPR. This is an indicator that the trained model has a strong bias towards classifying documents as legitimate. Weighted accuracy is between roughly 0,678 and 0,693.

Table 2 Performance of the Isolation Forest algorithm on the LingSpam Corpus trained on legitimate emails

|  | False Negative ratio | False Positive ratio | Weighted Accuracy |
|---|---|---|---|
| Fold 1 | 0,768537074148 | 0,137774413323 | 0,689285714286 |
| Fold 2 | 0,792207792208 | 0,129594543388 | 0,688186813187 |
| Fold 3 | 0,780163599182 | 0,131855747558 | 0,693956043956 |
| Fold 4 | 0,789108910891 | 0,140684410646 | 0,679395604396 |
| Fold 5 | 0,830474268416 | 0,140430351076 | 0,671703296703 |
| Fold 6 | 0,791374122367 | 0,132803632236 | 0,686813186813 |
| Fold 7 | 0,792415169661 | 0,12433661865 | 0,691758241758 |
| Fold 8 | 0,804953560372 | 0,122051666043 | 0,696153846154 |
| Fold 9 | 0,816125860374 | 0,130003812429 | 0,678296703297 |
| Fold 10 | 0,819075712881 | 0,129241326725 | 0,678021978022 |

Table 3 shows the classification results of IsolationForest on the LingSpam Corpus. The algorithm has proven itself to be weak, considering the relatively high false positive ratio and the close to chance accuracy.

| | False Negative ratio | False Positive ratio | Weighted Accuracy |
|---|---|---|---|
| Fold 1 | 0,154362416107 | 0,573200992556 | 0,492146596859 |
| Fold 2 | 0,180722891566 | 0,513307984791 | 0,544502617801 |
| Fold 3 | 0,133333333333 | 0,536708860759 | 0,532984293194 |
| Fold 4 | 0,153333333333 | 0,511801242236 | 0,544502617801 |
| Fold 5 | 0,187134502924 | 0,508928571429 | 0,548691099476 |
| Fold 6 | 0,173652694611 | 0,536802030457 | 0,526701570681 |
| Fold 7 | 0,210191082803 | 0,546365914787 | 0,50890052356 |
| Fold 8 | 0,188311688312 | 0,534332084894 | 0,521465968586 |
| Fold 9 | 0,140127388535 | 0,53634085213 | 0,528795811518 |
| Fold 10 | 0,186666666667 | 0,526708074534 | 0,526701570681 |

Table 4 shows the classification results by Simple Anomaly Detection Method on the LingSpam Corpus using the 'Max' combination rule. Because the distance of the thresholds is so small the algorithm either classifies all documents as spam or all documents as legitimate. As such the 'Max' combination rule is of no use, at least in this instance.

*Table 4 Performance of Simple Anomaly Detection method using the Max combination rule, trained on legitimate emails*

| | Threshold | False Negative ratio | False Positive ratio | Weighted Accuracy |
|---|---|---|---|---|
| Max | 1,4142135623731 | 1,0 | 0,0 | 0,5 |
| | 1,4141024512620 | 0,0 | 1,0 | 0,5 |
| | 1,4139913401509 | 0,0 | 1,0 | 0,5 |
| | 1,4138802290398 | 0,0 | 1,0 | 0,5 |
| | 1,4137691179287 | 0,0 | 1,0 | 0,5 |
| | 1,4136580068175 | 0,0 | 1,0 | 0,5 |
| | 1,4135468957064 | 0,0 | 1,0 | 0,5 |
| | 1,4134357845953 | 0,0 | 1,0 | 0,5 |
| | 1,4133246734842 | 0,0 | 1,0 | 0,5 |
| | 1,4131024512620 | 0,0 | 1,0 | 0,5 |

Table 5 shows the classification results by Simple Anomaly Detection Method on the LingSpam Corpus using the 'Mean' combination rule, trained on legitimate emails. At the optimal threshold, the false negative ratio is at 0,241 and the false positive ratio at 0,28. The weighted accuracy at these values is 0,73.

*Table 5 Performance of Simple Anomaly Detection method using the Mean combination rule on the LingSpam Corpus, trained on legitimate emails*

| | Threshold | False Negative ratio | False Positive ratio | Weighted Accuracy |
|---|---|---|---|---|
| Mean | 1,411846939524870 | 0,9932885906040269 | 0,0 | 0,5033557046979866 |
| | 1,4081802728582100 | 0,912751677852349 | 0,018610421836228287 | 0,5343189501557113 |
| | 1,4045136061915400 | 0,6711409395973155 | 0,062034739454094295 | 0,633412160474295 |
| | 1,4008469395248700 | 0,44966442953020136 | 0,15384615384615385 | 0,6982447083118224 |
| | 1,397180272858210 | 0,24161073825503357 | 0,2816377171215881 | 0,7383757723116892 |
| | 1,3935136061915400 | 0,18791946308724833 | 0,413151364764268 | 0,6994645860742419 |
| | 1,3898469395248800 | 0,12751677852348994 | 0,5732009925558312 | 0,6496411144603393 |
| | 1,3861802728582100 | 0,10067114093959731 | 0,6997518610421837 | 0,5997884990091096 |
| | 1,382513606191540 | 0,040268456375838924 | 0,794044665012407 | 0,582843439305877 |

Table 6 shows the classification results by Simple Anomaly Detection Method on the LingSpam Corpus using the 'Min' combination rule. At the optimal threshold, the false negative ratio is at 0,04 and the false positive ratio at 0,4. The weighted accuracy at these values is 0,78.

*Table 6 Performance of Simple Anomaly Detection method using the Min combination rule on the LingSpam Corpus, trained on legitimate emails*

|  | Threshold | False Negative ratio | False Positive ratio | Weighted Accuracy |
|---|---|---|---|---|
| Min | 1,3719029835293500 | 0,9664429530201343 | 0,0 | 0,5167785234899329 |
|  | 1,2231252057515800 | 0,040268456375838924 | 0,39950372208436724 | 0,7801139107698969 |
|  | 1,074347427973800 | 0,026845637583892617 | 0,7233250620347395 | 0,6249146501906839 |
|  | 0,9255696501960357 | 0,006711409395973154 | 0,8374689826302729 | 0,577909803986877 |
|  | 0,7767918724182623 | 0,0 | 0,8734491315136477 | 0,5632754342431762 |
|  | 0,628014094640489 | 0,0 | 0,892059553349876 | 0,5539702233250621 |
|  | 0,4792363168627156 | 0,0 | 0,9168734491315137 | 0,5415632754342432 |
|  | 0,33045853908494227 | 0,0 | 0,9280397022332506 | 0,5359801488833746 |
|  | 0,18168076130716895 | 0,0 | 0,956575682382134 | 0,5217121588089331 |

The performance of the algorithms on the EnronSpam corpus is as follows. Table 7 shows the classification results of IsolationForest on the EnronSpam Corpus. The algorithm features relatively high false negative ratios, while the false positive ratio is at maximum of 0,1404. The weighted accuracy is at 0,683 on average.

*Table 7 Performance of the Isolation Forest algorithm on the EnronSpam Corpus trained on legitimate emails*

|  | False Negative ratio | False Positive ratio | Weighted Accuracy |
|---|---|---|---|
| Fold 1 | 0,768537074148 | 0,137774413323 | 0,689285714286 |
| Fold 2 | 0,792207792208 | 0,129594543388 | 0,688186813187 |
| Fold 3 | 0,780163599182 | 0,131855747558 | 0,693956043956 |
| Fold 4 | 0,789108910891 | 0,140684410646 | 0,679395604396 |
| Fold 5 | 0,830474268416 | 0,140430351076 | 0,671703296703 |
| Fold 6 | 0,791374122367 | 0,132803632236 | 0,686813186813 |
| Fold 7 | 0,792415169661 | 0,12433661865 | 0,691758241758 |
| Fold 8 | 0,804953560372 | 0,122051666043 | 0,696153846154 |
| Fold 9 | 0,816125860374 | 0,130003812429 | 0,678296703297 |
| Fold 10 | 0,819075712881 | 0,129241326725 | 0,678021978022 |

Table 8 shows the classification results of IsolationForest on the LingSpam Corpus. The algorithm features very low false negative ratios, while the false positive ratios are at roughly 0,5. The weighted accuracy is at 0,621 on average.

*Table 8 Performance of the One-Class SVM algorithm on the EnronSpam Corpus trained on legitimate emails*

|         | False Negative ratio | False Positive ratio | Weighted Accuracy |
|---------|---------------------|---------------------|-------------------|
| Fold 1  | 0,0280561122244     | 0,512490537472      | 0,62032967033     |
| Fold 2  | 0,033966033966      | 0,506631299735      | 0,623351648352    |
| Fold 3  | 0,0316973415133     | 0,497746055597      | 0,627472527473    |
| Fold 4  | 0,0346534653465     | 0,502281368821      | 0,627472527473    |
| Fold 5  | 0,0282542885974     | 0,501321253303      | 0,627472527473    |
| Fold 6  | 0,0280842527583     | 0,516836927734      | 0,617032967033    |
| Fold 7  | 0,0329341317365     | 0,491281273692      | 0,63489010989     |
| Fold 8  | 0,0299277605779     | 0,493822538375      | 0,62967032967     |
| Fold 9  | 0,0285152409046     | 0,495996950057      | 0,634615384615    |
| Fold 10 | 0,0275319567355     | 0,50552802135       | 0,628021978022    |

Table 9 shows the classification results by Simple Anomaly Detection Method on the EnronSpam Corpus using the 'Max' combination rule. As in the previous example of 'Max' combination rule, because the distance of the thresholds is so small the algorithm either classifies all documents as spam or all documents as legitimate. As such the 'Max' combination rule is of no use, at least in this instance.

*Table 9 Performance of Simple Anomaly Detection method using the Max combination rule on the EnronSpam Corpus, trained on legitimate emails*

|     | Threshold    | False Negative ratio | False Positive ratio | Weighted Accuracy   |
|-----|--------------|---------------------|---------------------|---------------------|
| Max | 1,414207849  | 0,996993987975952   | 0,0                 | 0,501503006012024   |
|     | 1,414096738  | 0,06613226452905811 | 0,6252838758516276  | 0,6542919298096572  |
|     | 1,413985627  | 0,01002004008016032 | 0,8849356548069645  | 0,5525221525564377  |
|     | 1,413874516  | 0,01002004008016032 | 0,9693414080242241  | 0,5103192759478078  |
|     | 1,413763405  | 0,01002004008016032 | 0,9924299772899319  | 0,49877499131495395 |
|     | 1,413652294  | 0,01002004008016032 | 0,9981074943224829  | 0,4959362327986784  |
|     | 1,413541183  | 0,01002004008016032 | 0,9992429977289932  | 0,4953684810954233  |
|     | 1,413430071  | 0,01002004008016032 | 0,9992429977289932  | 0,4953684810954233  |
|     | 1,413318960  | 0,0                 | 0,9992429977289932  | 0,5003785011355034  |
|     | 1,413096738  | 0,0                 | 1,0                 | 0,5                 |

Table 10 shows the classification results by Simple Anomaly Detection Method on the EnronSpam Corpus using the 'Mean' combination rule. At the optimal threshold, the false negative ratio is at 0,095 and the false positive ratio at 0,226. The weighted accuracy value is 0,839.

*Table 10 Performance of Simple Anomaly Detection method using the Mean combination rule on the EnronSpam Corpus, trained on legitimate emails*

|  | Threshold | False Negative ratio | False Positive ratio | Weighted Accuracy |
|---|---|---|---|---|
| Mean | 1,414408241388 | 1,0 | 0,0 | 0,5 |
|  | 1,410408241388 | 0,7875751503006012 | 0,0049205147615442 | 0,60375216746689272 |
|  | 1,406408241388 | 0,4759519038076152 | 0,0336866010598031 | 0,7451807475662908 |
|  | 1,402408241388 | 0,25250501002004005 | 0,0957607872823618 | 0,825867101348799 |
|  | 1,398408241388 | 0,09519038076152304 | 0,2263436790310370 | 0,83923297010372 |
|  | 1,394408241388 | 0,023046092184368736 | 0,3709311127933384 | 0,8030113975111464 |
|  | 1,390408241388 | 0,013026052104208416 | 0,5174110522331568 | 0,7347814478313175 |
|  | 1,386408241388 | 0,002004008016032064 | 0,6377744133232399 | 0,680110789330364 |
|  | 1,382408241388 | 0,001002004008016032 | 0,7229371688115064 | 0,6380304135902388 |
|  | 1,374408241388 | 0,0 | 0,8622255866767601 | 0,56888720666162 |

Table 11 shows the classification results by Simple Anomaly Detection Method on the EnronSpam Corpus using the 'Min combination rule. At the optimal threshold, the false negative ratio is at 0,04 and the false positive ratio at 0,1. The weighted accuracy value is 0,927.

|  | Threshold | False Negative ratio | False Positive ratio | Weighted Accuracy |
|---|---|---|---|---|
| Min | 1,360757505 | 0,8887775551102205 | 0,0 | 0,5556112224448897 |
|  | 1,211090838 | 0,047094188376753505 | 0,0991672975018925 | 0,926869257060677 |
|  | 1,061424171 | 0,006012024048096192 | 0,336109008327025 | 0,8289394838124394 |
|  | 0,91175750475803 | 0,0 | 0,4651778955336866 | 0,7674110522331568 |
|  | 0,76209083809137 | 0,0 | 0,5582891748675246 | 0,7208554125662376 |
|  | 0,61242417142471 | 0,0 | 0,6404239212717638 | 0,679788039364118 |
|  | 0,46275750475805 | 0,0 | 0,7051476154428463 | 0,6474261922785769 |
|  | 0,31309083809139 | 0,0 | 0,7880393641180924 | 0,6059803179409537 |
|  | 0,16342417142472 | 0,0 | 0,858440575321726 | 0,5707797123391369 |
|  | 0,01375750475806 | 0,0 | 0,8834216502649508 | 0,5582891748675246 |

## 4.2 Models trained on spam emails

Again, let us consider the performance of the models trained on spam emails using the LingSpam Corpus. Table 12 shows the classification results of IsolationForest on the LingSpam Corpus, trained on spam emails. The model has a disproportionally high FNR and a low FPR. This is an indicator that the trained model has a strong bias towards classifying documents as legitimate. Weighted accuracy is poor as a result however, we can simply reverse the output of the classifier.

|  | False Negative ratio | False Positive ratio | Weighted Accuracy |
|---|---|---|---|
| Fold 1 | 0,961538461538 | 0,0872483221477 | 0,174869109948 |
| Fold 2 | 0,950570342205 | 0,066265060241 | 0,203141361257 |
| Fold 3 | 0,940506329114 | 0,0606060606061 | 0,211518324607 |
| Fold 4 | 0,95652173913 | 0,08 | 0,181151832461 |
| Fold 5 | 0,948979591837 | 0,105263157895 | 0,202094240838 |
| Fold 6 | 0,973350253807 | 0,059880239521 | 0,186387434555 |
| Fold 7 | 0,954887218045 | 0,0955414012739 | 0,186387434555 |
| Fold 8 | 0,940074906367 | 0,0974025974026 | 0,195811518325 |
| Fold 9 | 0,967418546366 | 0,0509554140127 | 0,183246073298 |
| Fold 10 | 0,919254658385 | 0,0466666666667 | 0,21780104712 |

Table 13 shows the classification results of One-Class SVM on the LingSpam Corpus, trained on spam emails. The model has an excellent false negative ratio of 0,0 with a weighed accuracy of roughly 0.9.

*Table 13 Performance of the One-Class SVM algorithm on the LingSpam Corpus trained on spam emails*

|  | **False Negative ratio** | **False Positive ratio** | **Weighted Accuracy** |
|---|---|---|---|
| Fold 1 | 0,0 | 0,697986577181 | 0,89109947644 |
| Fold 2 | 0,0 | 0,560240963855 | 0,902617801047 |
| Fold 3 | 0,0 | 0,684848484848 | 0,88167539267 |
| Fold 4 | 0,0 | 0,56 | 0,912041884817 |
| Fold 5 | 0,0 | 0,497076023392 | 0,910994764398 |
| Fold 6 | 0,0 | 0,574850299401 | 0,899476439791 |
| Fold 7 | 0,0 | 0,547770700637 | 0,909947643979 |
| Fold 8 | 0,0 | 0,564935064935 | 0,90890052356 |
| Fold 9 | 0,0 | 0,668789808917 | 0,890052356021 |
| Fold 10 | 0,0 | 0,486666666667 | 0,923560209424 |

Table 14 shows the classification results by Simple Anomaly Detection Method on the LingSpam Corpus using the 'Max' combination rule, trained on spam emails. Because the small distance between the lower and upper the algorithm either classifies all documents as spam or all documents as legitimate.

*Table 14 Performance of Simple Anomaly Detection method using the Max combination rule on the LingSpam Corpus, trained on spam emails*

|  | **Threshold** | **False Negative ratio** | **False Positive ratio** | **Weighted Accuracy** |
|---|---|---|---|---|
| Max | 1,41421356237309 | 0,9689826302729528 | 0,0 | 0,5155086848635235 |
|  | 1,41399134015087 | 0,008684863523573201 | 0,8791946308724832 | 0,5560602528019718 |
|  | 1,41376911792865 | 0,007444168734491315 | 0,8859060402684564 | 0,5533248954985261 |
|  | 1,41354689570642 | 0,00620347394540943 | 0,8993288590604027 | 0,5472338334970939 |
|  | 1,41332467348420 | 0,004962779156327543 | 0,9060402684563759 | 0,5444984761936482 |
|  | 1,41310245126198 | 0,002481389578163771 | 0,9395973154362416 | 0,5289606474927973 |
|  | 1,41288022903976 | 0,001240694789081885 | 0,9463087248322147 | 0,5262252901893517 |
|  | 1,41265800681754 | 0,001240694789081885 | 0,9731543624161074 | 0,5128024713974053 |
|  | 1,41243578459531 | 0,0 | 0,9798657718120806 | 0,5100671140939597 |
|  | 1,41199134015087 | 0,0 | 0,9865771812080537 | 0,5067114093959731 |

Table 15 shows the classification results by Simple Anomaly Detection Method on the LingSpam Corpus using the 'Mean' combination rule trained on spam emails. At the optimal threshold, the false negative ratio is at 0,08 and the false positive ratio at 0,194. The weighted accuracy value is 0,861.

| | Threshold | False Negative ratio | False Positive ratio | Weighted Accuracy |
|---|---|---|---|---|
| Mean | 1,4036487469578 | 0,7109181141439206 | 0,0 | 0,6445409429280398 |
| | 1,3989820802911 | 0,42555831265508687 | 0,0268456375838926 | 0,7737980248805103 |
| | 1,3943154136245 | 0,21836228287841192 | 0,0738255033557047 | 0,8539061068829417 |
| | 1,3896487469578 | 0,08312655086848635 | 0,1946308724832214 | 0,8611212883241461 |
| | 1,3849820802911 | 0,019851116625310174 | 0,2953020134228188 | 0,8424234349759355 |
| | 1,3803154136245 | 0,007444168734491315 | 0,4093959731543624 | 0,7915799290555732 |
| | 1,3756487469578 | 0,003722084367245657 | 0,5570469798657718 | 0,7196154678834913 |
| | 1,3709820802911 | 0,0 | 0,5973154362416108 | 0,7013422818791946 |
| | 1,3663154136245 | 0,0 | 0,6174496644295302 | 0,6912751677852349 |
| | 1,3569820802911 | 0,0 | 0,7114093959731543 | 0,6442953020134228 |

Table 16 shows the classification results by Simple Anomaly Detection Method on the LingSpam Corpus using the 'Min' combination rule trained on spam emails. At the optimal threshold, the false negative ratio is at 0,28 and the false positive ratio at 0. The weighted accuracy value is 0,859.

*Table 16 Performance of Simple Anomaly Detection method using the Min combination rule on the LingSpam Corpus, trained on spam emails*

|  | Threshold | False Negative ratio | False Positive ratio | Weighted Accuracy |
|---|---|---|---|---|
| Min | 1,34420565491 | 0,2816377171215881 | 0,0 | 0,859181141439206 |
|  | 1,19987232158 | 0,003722084367245657 | 0,4161073825503356 | 0,7900852665412094 |
|  | 1,05553898824 | 0,001240694789081885 | 0,6040268456375839 | 0,6973662297866672 |
|  | 0,91120565491 | 0,001240694789081885 | 0,6510067114093959 | 0,6738762969007611 |
|  | 0,76687232157627 | 0,001240694789081885 | 0,7046979865771812 | 0,6470306593168684 |
|  | 0,62253898824294 | 0,0 | 0,7248322147651006 | 0,6375838926174497 |
|  | 0,47820565490961 | 0,0 | 0,7651006711409396 | 0,6174496644295302 |
|  | 0,33387232157628 | 0,0 | 0,7919463087248322 | 0,6040268456375839 |
|  | 0,18953898824295 | 0,0 | 0,8389261744966443 | 0,5805369127516778 |
|  | 0,0452056549096 | 0,0 | 0,8993288590604027 | 0,5503355704697986 |

The performance of the algorithms, trained on spam emails, on the EnronSpam corpus is as follows. Table 17 shows the classification results of IsolationForest trained on spam email. The model has a disproportionally high FNR and a low FPR. This is an indicator that the trained model has a strong bias towards classifying documents as illegitimate. Weighted accuracy is poor at first glance however, we can simply reverse the output of the classifier, which would represent the weighted accuracy at around 0,685.

*Table 17 Performance of the Isolation Forest algorithm on the EnronSpam Corpus trained on spam emails*

|  | False Negative ratio | False Positive ratio | Weighted Accuracy |
|---|---|---|---|
| Fold 1 | 0,886449659349 | 0,107214428858 | 0,327197802198 |
| Fold 2 | 0,882152330428 | 0,0959040959041 | 0,334065934066 |
| Fold 3 | 0,913974455297 | 0,0828220858896 | 0,309340659341 |
| Fold 4 | 0,913307984791 | 0,0881188118812 | 0,315659340659 |
| Fold 5 | 0,904492261231 | 0,0847628657921 | 0,318681318681 |
| Fold 6 | 0,901626939084 | 0,114343029087 | 0,314010989011 |
| Fold 7 | 0,932524639879 | 0,10878243513 | 0,294230769231 |
| Fold 8 | 0,931111943092 | 0,119711042312 | 0,28489010989 |
| Fold 9 | 0,870758673275 | 0,106194690265 | 0,342857142857 |
| Fold 10 | 0,893633244377 | 0,110127826942 | 0,325274725275 |

Table 18 shows the classification results of OneClass SVM on the EnronSpam Corpus, trained on spam emails. The model has a low FNR and an FPR of roughly 0,5. Weighted accuracy is on average 0,789.

|  | **False Negative ratio** | **False Positive ratio** | **Weighted Accuracy** |
|---|---|---|---|
| Fold 1 | 0,103709311128 | 0,5 | 0,787637362637 |
| Fold 2 | 0,0909435392194 | 0,529470529471 | 0,788461538462 |
| Fold 3 | 0,0961682945154 | 0,5081799591 | 0,793131868132 |
| Fold 4 | 0,093536121673 | 0,534653465347 | 0,784065934066 |
| Fold 5 | 0,0860702151755 | 0,556004036327 | 0,785989010989 |
| Fold 6 | 0,0813469542187 | 0,523570712136 | 0,797527472527 |
| Fold 7 | 0,0841546626232 | 0,545908183633 | 0,788736263736 |
| Fold 8 | 0,0846125046799 | 0,528379772962 | 0,797252747253 |
| Fold 9 | 0,0831109416698 | 0,525073746313 | 0,793406593407 |
| Fold 10 | 0,0880670987419 | 0,530973451327 | 0,788186813187 |

Table 19 shows the classification results by Simple Anomaly Detection Method on the EnronSpam Corpus using the 'Max' combination rule, trained on spam emails. Because the distance of the thresholds is small the algorithm either classifies all documents as spam or all documents as legitimate.

|  | **Threshold** | **False Negative Ratio** | **False Positive ratio** | **Weighted Accuracy** |
|---|---|---|---|---|
| Max | 1,41521142229 | 1,0 | 0,0 | 0,5 |
|  | 1,41498920007 | 1,0 | 0,0 | 0,5 |
|  | 1,41476697785 | 1,0 | 0,0 | 0,5 |
|  | 1,41454475563 | 1,0 | 0,0 | 0,5 |
|  | 1,41432253340 | 1,0 | 0,0 | 0,5 |
|  | 1,41410031118 | 0,0632096896290689 | 0,9749498997995992 | 0,4809202052856659 |
|  | 1,41387808896 | 0,0 | 0,9909819639278558 | 0,5045090180360721 |
|  | 1,41365586674 | 0,0 | 1,0 | 0,5 |
|  | 1,41343364452 | 0,0 | 1,0 | 0,5 |
|  | 1,41298920007 | 0,0 | 1,0 | 0,5 |

Table 20 shows the classification results by Simple Anomaly Detection Method on the EnronSpam Corpus using the 'Mean' combination rule trained on spam emails. At the optimal threshold, the false negative ratio is at 0,171 and the false positive ratio at 0,442. The weighted accuracy value is 0,694.

*Table 20 Performance of Simple Anomaly Detection method using the Mean combination rule on the EnronSpam Corpus, trained on spam emails*

|      | Threshold    | False Negative Ratio  | False Positive ratio  | Weighted Accuracy   |
|------|--------------|-----------------------|-----------------------|---------------------|
| Mean | 1,41399255   | 1,0                   | 0,0                   | 0,5                 |
|      | 1,41088144   | 0,8959121877365632    | 0,0651302605210420    | 0,5194787758711974  |
|      | 1,40777033   | 0,6298258894776685    | 0,1603206412825651    | 0,6049267346198832  |
|      | 1,40465922   | 0,35692657077971235   | 0,3046092184368738    | 0,669232105391707   |
|      | 1,40154811   | 0,17070401211203634   | 0,4418837675350701    | 0,6937061101764468  |
|      | 1,39843700   | 0,06510219530658592   | 0,5741482965931863    | 0,6803747540501139  |
|      | 1,39532589   | 0,02384557153671461   | 0,6993987975951904    | 0,6383778154340475  |
|      | 1,39221477   | 0,003406510219530658  | 0,8096192384769539    | 0,5934871256517578  |
|      | 1,38910366   | 0,000378501135503406  | 0,9008016032064128    | 0,549409947829042   |
|      | 1,38288144   | 0,0                   | 0,9889779559118237    | 0,5055110220440882  |

Table 21 shows the classification results by Simple Anomaly Detection Method on the EnronSpam Corpus using the 'Min combination rule trained on spam emails. At the optimal threshold, the false negative ratio is at 0,171 and the false positive ratio at 0,442. The weighted accuracy value is 0,694.

*Table 21 Performance of Simple Anomaly Detection method using the Min combination rule on the EnronSpam Corpus, trained on spam emails*

| | Threshold | False Negative Ratio | False Positive ratio | Weighted Accuracy |
|---|---|---|---|---|
| Min | 1,394503895 | 0,9992429977289932 | 0,0 | 0,5003785011355034 |
| | 1,243392784 | 0,04806964420893263 | 0,3016032064128257 | 0,8251635746891208 |
| | 1,092281673 | 0,004163512490537472 | 0,5440881763527055 | 0,7258741555783785 |
| | 0,94117056167920 | 0,002649507948523845 | 0,654308617234469 | 0,6715209374085036 |
| | 0,79005945056809 | 0,0 | 0,7344689378757515 | 0,6327655310621243 |
| | 0,63894833945699 | 0,0 | 0,7735470941883767 | 0,6132264529058116 |
| | 0,48783722834588 | 0,0 | 0,845691382765531 | 0,5771543086172345 |
| | 0,33672611723478 | 0,0 | 0,9328657314629258 | 0,5335671342685371 |
| | 0,18561500612367 | 0,0 | 0,9639278557114228 | 0,5180360721442886 |
| | 0,03450389501256 | 0,0 | 0,9769539078156313 | 0,5115230460921844 |

## 5 Results

In the following chapter the results of the described experiments are presented. The results are presented on a per-dataset basis, sorted by the training class (ham or spam) and the trained algorithm. For Isolation Forest and One-Class SVM we have used the average value aggregated from the 10 folds. For the Simple Anomaly Detection method, we have used the data from the threshold with the highest weighted accuracy.

Table 22 shows the results of the models trained on the LingSpam algorithm. Trained on legitimate documents, the worst performing algorithms are One-Class SVM and Simple Anomaly Method – MAX which have weighted accuracies near chance. The best algorithm is the Simple Anomaly Method – MIN, which has a weighted accuracy of 0,78. This method has a false negative ratio of 0,04 and a false positive ratio of 0,4. This means that the model is seldom wrong when it comes to identifying legitimate mail, but has some issues determining whether an email is spam. In the case of LingSpam, the results are markedly better when models are trained on spam email. As before, Simple Anomaly Detection – MAX fails at a weighted accuracy of near chance – 0,56. The best performing algorithm was the One-Class SVM with a weighted accuracy of 0,9. The false negative ratio is 0 and the false positive ratio is 0,55. This means that the method never misclassifies a spam document as a ham document.

*Table 22 Lingspam Corpus results*

| LingSpam | Algorithms | | | | |
|---|---|---|---|---|---|
| Weighted Accuracy | Isolation Forest | One-Class SVM | Simple Method- MAX | Simple Method - MEAN | Simple Method - MIN |
| Trained on 'ham' | 0,685 | 0,522 | 0,5 | 0,738 | 0,78 |
| Trained on 'spam' | 0,805 | 0,9 | 0,56 | 0,861 | 0,86 |
| | | | | | |
| False Negative Ratio | | | | | |
| Trained on 'ham' | 0,798 | 0,165 | 0 | 0,242 | 0,04 |
| Trained on 'spam' | 0,04 | 0 | 0,009 | 0,083 | 0,28 |
| | | | | | |
| False Positive Ratio | | | | | |
| Trained on 'ham' | 0,138 | 0,521 | 1 | 0,28 | 0,4 |
| Trained on 'spam' | 0,925 | 0,55 | 0,88 | 0,19 | 0 |

Table 23 shows the results of the models trained on the EnronSpam algorithm. As with LingSpam the worst performing methods trained on 'ham' are One-Class SVM and Simple Method – MAX. The best performing method is Simple Method – MIN with a weighted accuracy of 0,92. Both the false negative and false positive ratios are excellent at 0,05 and 0,01 respectively. Trained on 'ham' documents, the worst performing algorithm is Simple Anomaly Detection method – MAX with a weighted accuracy of 0,5. The best performing method is again Simple Method Min with a FNR of 0,5 and a FPR of 0,3. In general, the weighted accuracies are high despite the sometimes mediocre FNR and FPR because of the imbalanced dataset.

*Table 23 EnronSpam Corpus results*

| EnronSpam | Algorithms | | | | |
|---|---|---|---|---|---|
| Weighted Accuracy | Isolation Forest | One-Class SVM | Simple Method- MAX | Simple Method - MEAN | Simple Method - MIN |
| Trained on 'ham' | 0,683 | 0,621 | 0,65 | 0,83 | 0,92 |
| Trained on 'spam' | 0,685 | 0,791 | 0,5 | 0,69 | 0,83 |
| | | | | | |
| False Negative Ratio | | | | | |
| Trained on 'ham' | 0,801 | 0,028 | 0,06 | 0,09 | 0,05 |
| Trained on 'spam' | 0,088 | 0,089 | 0 | 0,17 | 0,05 |
| | | | | | |
| False Positive Ratio | | | | | |
| Trained on 'ham' | 0,131 | 0,501 | 0,63 | 0,23 | 0,01 |
| Trained on 'spam' | 0,898 | 0,535 | 1 | 0,44 | 0,3 |

To determine the superior method of training the models we have compared the average and maximum weighted accuracies for both corpora as shown in Table 24. On average, the LingSpam corpus did better trained on 'spam' with an average weighted accuracy of 0,792 versus the accuracy of 0,645 when trained on 'ham' documents. This is also reflected in the maximum weighted accuracies of 0,9 for 'spam' and 0,78 for 'ham'. The opposite is true for the EnronSpam corpus. Trained on 'ham' documents, the average weighted accuracy of the models is 0,74 compared to 0,69 when trained on 'spam' documents. This trend continues with the maximum weighted accuracies of 0,92 for the best performing model trained on 'ham' documents and 0,83 for the best model trained on 'ham'. We can recognize from this that the training method depends on the dataset.

*Table 24 Average and maximum weighted accuracies by type of training*

| Average Weighted Accuracies | LingSpam | EnronSpam |
|---|---:|---:|
| Trained on 'ham' | 0,645 | 0,7408 |
| Trained on 'spam' | 0,7972 | 0,6992 |
| | | |
| Maximum Weighted Accuracies | LingSpam | EnronSpam |
| Trained on 'ham' | 0,78 | 0,92 |
| Trained on 'spam' | 0,9 | 0,83 |

# 6 Conclusions

Continuously improving spam detection techniques is crucial to the effort of continuously improving safety on the internet. Increasingly sophisticated filter evasion methods demand an increasingly larger stack of technologies, each approaching the problem of spam detection from a different perspective.

In this paper, we have approached this problem from an anomaly detection perspective. We have implemented the Simple Anomaly Detection Method proposed by Laorden et. al and compared the resulting weighted average, false negative ratio and false positive ratio with the Isolation Forest and One-Class SVM algorithms. We have built separate models with legitimate and spam documents as representations of normality. We have found the results to be satisfactory, with the best method being the Simple Anomaly Detection using the minimum combination rule.

There are many avenues of research that are yet to be explored in this area. First and foremost, the research in this paper should be furthered by utilizing additional metrics like the ROC curve and AUC (Demsar, 2006). It would also be prudent to incorporate larger datasets and additional corpora to broaden the understanding of the cause and effect of dataset size on classification performance. Secondly, we should consider further research on anomaly detection ensembles where each member of the ensemble specializes in detecting types of anomalies. Thirdly, we should consider evaluating the performance of anomaly detection systems after utilizing with white-listing, grey-listing and black-listing and before using supervised methods.

# 7 References

Bhowmick, A., & Hazarika, S. M. (2016). E-Mail Spam Filtering: A Review of Techniques and Trends. In *Advances in Electronics , Communication and Computing* (pp. 583–590). https://doi.org/10.1007/978-981-10-4765-7

Chandola, V., Banerjee, A., & Kumar, V. (2007). *Anomaly Detection: A Survey. Information Sciences journal* (Vol. 32). https://doi.org/10.1038/nn.2116

Demsar, J. (2006). Statistical Comparisons of Classifiers over Multiple Data Sets. *Journal of Machine Learning Research 7*. https://doi.org/10.1016/j.jecp.2010.03.005

Geerthik, S. (2013). Survey on Internet Spam : Classification and Analysis. *International Journal of Computer Technology & Applications*, *4*(June), 384–391.

Guzella, T. S., & Caminhas, W. M. (2009). A review of machine learning approaches to Spam filtering. *Expert Systems with Applications*. https://doi.org/10.1016/j.eswa.2009.02.037

Laorden, C., Ugarte-Pedrero, X., Santos, I., Sanz, B., Nieves, J., & Bringas, P. G. (2014). Study on the effectiveness of anomaly detection for spam filtering. *Information Sciences*. https://doi.org/10.1016/j.ins.2014.02.114

Liu, F. T., & Ting, K. M. (2018). Isolation Forest. In *Eighth IEE International Confrence on Data mining*. https://doi.org/10.1109/ICDM.2008.17

Manevitz, L. M. (2001). One-Class SVMs for Document Classification. *Journal of Machine Learning*, *2*, 139–154.

Rao, J. M., & Reiley, D. H. (2012). The Economics of Spam. *Journal of Economic Perspectives*, *26*(3), 87–110. https://doi.org/10.1257/jep.26.3.87

Santos, I., Laorden, C., Ugarte-Pedrero, X., Sanz, B., & Bringas, P. G. (2011). Anomaly-based Spam Filtering. In *Proceedings of the 6th International Conference on Security and Cryptography SECRYPT* (p. In press). Retrieved from http://paginaspersonales.deusto.es/bosanz/publications/pdf/2011/Santos_2011_SECRYPT_Anomaly-based_Spam_Filtering.pdf

Tech, Y., & Gomez-gil, P. (2013). An assessment of ten-fold and Monte Carlo cross validations for time series forecasting, (April 2015). https://doi.org/10.1109/ICEEE.2013.6676075