

---

# Smart AI powered Spam Classification

# 1 • Introduction & Objectives

Most e-mail readers a nul•trWial amount of time regularly deleting junk e-mail (spam) messages, even as an expanding volume of such e-mail occupies server storage space and consumes network bandwidth. An ongoing challenge, therefore, rests within the development and refinement of automatic classifiers that can ffstirvuish egitimate e-mail frotn spam. Somepubshsed studies have examined swm detectors Lßng Naive Bayesian approaches and large feature sets of binary attributes that the existence of common and many commercial applications also use keywords in spam. Naive Bayesian t«hniques. Spammers rec%nize these attempts to their and have circumvent tt-æse filters, but these evasive tactics are developed tactics to themselves patterns that human can often identify quick". This work that human readers had the objectives of developing an alternative approach using a neural network (NN) classifier brained on a corpus of e-mail messages from several users. The features selection used in this work is one of the major improvements. beca use thefeature uses descriptive characteristics of words and messages similar to those that a human reader would use to identify spam, and the rno&l to select the best feature set, was based on forward feature selection. Another Ob\*ctWe in this work was to improve the spam detection near 95% of accuracy using Artificial Neural Networks; actually nomdv has reached rmore than 9% of accuracy using ANN.

## 1-1 - What is "spam"

Spam, terms, rmeans unwanted It has norrnany used to to unwanted email or userwt mess"es, and it is now also being used to to unwanted Instant (1M) and telephone Short Service (SMS) Spam email is unwanted, uninvited, and inevitably promotes something for sale. Often the terms junk emad. Unsolicited guk Email (UBE). or Unsolicited Commercial Email (LICE) are used to refer to spam email. Spam generally promotes Internet — based sales, but it also occasionally prormtes telephone• based or other methMs of Saks tm.

People wtw specialhze in sending spam are ca"ed spummers. Companies pay spamrners to send emads on their tkhatf, and the spammers have developed a range of computerized tools and techniques to send these messages. Spammers also run their own onune businesses and market them using spam email.

The term •spam emar genera"v precludes email from known sources, egardless of however unwanted the content is. One example of this woud be an endless "st of jokes sent frotn uquaintances. Email virus, Trojan Mrses, and other ma'ware (short for malicious software) are normal" categorized as spam either. althoe they share some common traits with sparm Emails that are not spam are often referred to as ham, particulartv in the anti.sparn community. Spam is subjective, and a mess%e considered spam by one recoent may bewelcomed by anottaer.

Anti-spam took be partianv effective in bl«king malware; however, they are best at bhcking spam. ial ant-virus software cm and should used to protect your inbox from other email

## 1.2 - Definitions

The ollowing definitions will be used th this w«k:

- Spam: Unsolicited commercial Email UCE. it is any email that has not been requested and contains an advertisement of some kind. ■ Ham: The opposite of spam— email that is wanted
- False Negative: A spam email message that was not detected successfully. ■ False Positive: A ham email message that was wrongly detected as spam.

### 1.3 - The History of spam

Here are important dates in the development of the internet:

- 1969: Two computers networked via a router
- 1971: email using a rudimentary system
- 1979: usenet (newsgroups) established
- IBO: The World Wide Web concept born
- 21<sup>st</sup> Å4: The Internet is a major global network responsible for billions of dollars of commerce.

There is one unmissable item from this timeline:

- 1978: The first email sent

Spam has part of the Internet from a relatively early stage in its development. The first spam email was sent on May 3<sup>rd</sup>, 1978, when the US Government funded Arpanet; as it was called then the first spammer was a DEC engineer called Gary Thuerk who invited recipients of his email to attend a product. This email was using the Arpanet, and caused an immediate response from the chief of the Arpanet. Major Raymond Czahor. at the violation of the non-commercial of the Arpanet.

Spam really took off 1994 when an Arizona attorney, Laurence Carter, automated the mailing of messages to many internet newsgroups (usenet) to advertise his firm's services. The resultant outcry from usenet users included the coining of the term 'spam'. when one respondent wrote —'Send cottons and cans of Spam to Cantor R Co.' - this sparked the rise of spam as is now experienced.

Spam email has increased in volume as the Internet has developed. In April PC Magazine reported that 10% of all email is spam.

### 1.4 • Spammers

Typically, spammers are paid to advertise particular websites, and companies, and are specialists in sending spam email. There are several well-known spammers who are responsible for a large proportion of spam and have evaded legal action-

Individuals of websites send their own but have extensive mailing lists and to bypass filters and avoid detection. They have a niche in the marketing industry, and their clients capitalize on this.

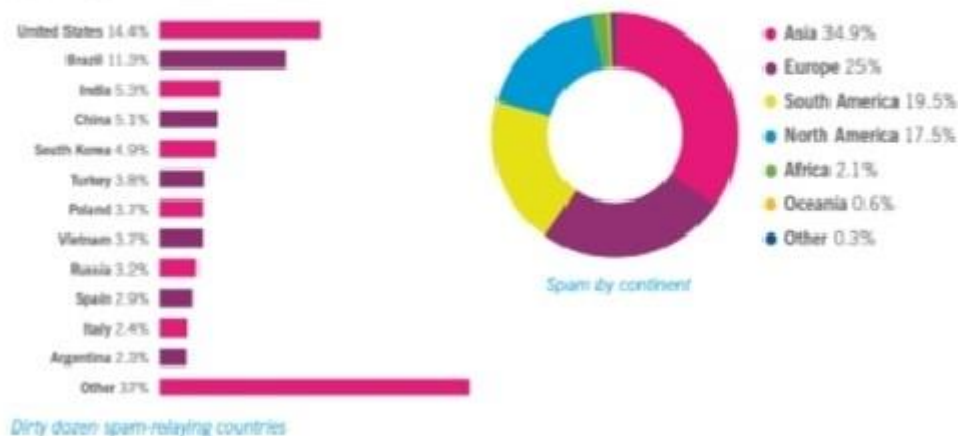
Most emails are now sent from Trojan loaded computers, as reported in a press release by broadband specialist. The owners of trojaned computers have been tricked into running software that allows a spammer to send spam email from the computer without knowledge of the Trojan software often exploits security holes in the system. browser,

or email of a user. When a malicious website is visited, the software is installed on the computer. unknown to users, their converter may source of thousands of spam email a day.

## 1.5 - Emerging spam for social networking attacks

With individuals and businesses on online social outlets, cyber-criminals have taken notice and started using them for their gain. Beyond the common nuisances, such as wasted company time and bandwidth, malware and malicious data theft issues have presented serious problems to social networks and their users. Spam is now common on "Real networking sites. and social engineering—trying to trick users to reveal vital data, or persuading people to visit dangerous web links—is on the rise.

Social network login credentials have become as valuable as email addresses, aiding the dissemination of social spam because these emails are more likely to be opened and trusted than standard messages. In many cases, spam and malware distribution are closely intertwined.



Source: SODHOS IR. — 2010 World Spam Review — Spain is part of the top 10

## 1.6 • The costs of spam

Spam is very cheap to send, the cost is insignificant as compared to conventional marketing techniques, so marketing by spam is very cost-effective. Despite very low rates of response, it translates into high costs for the victim.

## 1.7 - Costs to the Spammer

As of Tom Geller, Executive Director of SpamCon Foundation, estimated that the cost to send a single email was as little as one thousandth of a cent. Yet the cost to the recipient was around 10 cents.

The overheads in sending spam are low. The real costs are:

- **Internet connection:** there are lots of Internet Providers (ISPs) offering packages at around \$20/month. A spammer doesn't particularly need a Digital Subscriber Line (DSL) or cable modem service, a dial-up connection will also allow large quantities to be sent. In fact, accounts are preferable, as broadband accounts are routinely shut down when complaints about spam are received. Dial-up accounts are easy to set up and can be up and running within minutes, but DSL typically has a time of days.

- Software: specialist SWn software is essential. A normal email client will restrict the number of messages that can be sent, and require the spammer to spend more time in front of the computer. Spammers usually write their own software, steal someone else's, or buy it. A spammer with some technical knowledge and starting from scratch can have software ready after a week. To pay someone to develop that software would cost the spammer €1000.
- A mailing list: most spammers will build up their own list of email addresses. For beginners, it is possible to buy a CD with 6 million email addresses on it for around €50. Ironically, these CDs are marketed via spam email. Email addresses that are guaranteed to currently be active sell for larger sums.
- A webserver; this is an optional cost. It allows a spammer to deliver web bug-images to validate their mailing list. Web bugs are discussed further in later chapters. Basic web hosting costs less than £10 a month.

less than plus monthly cost of less than 160. a spammer would have the software, internet and a of req- to operational.

A single computer can send off hundreds of spam. Spam varies. a typical message size might be around 100K. On a fast dial-up of around 56K it might take one second to send an email to one recipient. It would take only a minute to send it to 100 others. In other words, at least emails can be sent in an hour. For smaller emails, the sent hour would be greater. The spammer needs to invest 15 minutes of their time and software will continue to send spam for many months. With three phone lines, they could work for a total of an hour, and send approximately 10,000 emails an hour or 200,000 a day or more using a DQ line.

### i.e - Costs to the Recipient

The Union has a study into UCE in 2002. In the findings, it estimated the cost of spam to consumers and businesses to be around 8% of turnover. These costs are partly incurred through lost productivity or partly in direct costs. and 2% in direct costs incurred passed on.

The cost of a spam environment is estimated to be as high as 6m to 10m per year, per employee. For a 50-person company, this cost is £500,000 per year. Spam emails distract or take employees time and use disk space, processing power, and network bandwidth. Removing spam by hand is time consuming and laborious when there is a large amount of spam. In addition, there is a business risk, as genuine messages may be lost along with unwanted ones. can contaminate a company's reputation that some employees won't tolerate.

## L9 • Spam and the Law

In the legislative proceeding on spam has in progress since 1997. The latest legislation is the CAN SPAM act (15 U.S.C. 6801-6809) of 2003. This supersedes many State laws and is currently being used to prosecute spammers. However, it is not a deterrent. The Coalition Against Unsolicted Commercial Email (CAUCE) reported in June 2004 that despite several high-profile lawsuits by the Federal Trade Commission (FTC) and ISPs, spam volumes were still increasing. The CAN-SPAM act is seen as weak on two counts: that consumers have to explicitly opt-out from commercial and secondly, ISPs can take action against spammers,

In Europe, legislation exists that makes sending unsolicited email illegal. However, when Directive 2002/58/EC was passed there were several problems with it. Business-to-business emails were excluded — a business could email each and every account at any other business and stay within the law. Additionally, individual member States have to pass their own laws and penalties for offenders. The law requires spammers to use opt-in emailing where recipients have to explicitly request to receive commercial email rather than the opt-out

proposed in the USA. where anyone can receive spam and has to request to be removed from mailing lists.

In June 2004, a UK law was passed that gave spammers a choice of moving their operations to the UK due to the leniency of the laws there. The maximum penalty they face in the UK is £5,000 (pounds), while in Italy spammers face up to three years of imprisonment. In 2004, had been convicted under this in the UK.

In Australia, the Act came into effect in April. This makes it illegal, using the opt-in model. Additionally, there have been successful prosecutions for spam in Australia using previous laws.

The internet is a multinational network and domestic legislation cannot reach to another country. A US-based spammer would be at risk of prosecution if it spammed US citizens and advertised a product and sold in the US. But the spammer from the Far East would be at very little risk of prosecution. Domestic legislation will not affect the volume of but it may occasionally affect the types of products advertised via spam.

Spammers will often reroute spam via other nations, so spam is sent from the US to another country and relayed back to the US. This makes it more difficult to trace the source of the email and to prosecute them. Many countries have no anti-spam laws and there is little even risk to them. The difficulty by the Internet makes it difficult to trace spam email to its source. Anti-spam efforts are directed towards spammers through other means. In May 2008, the New York Times reported that the Direct Marketing Association is using trails in the real world to track spammers in the virtual world with success.

## 2 - Spam Techniques

As spam increased in volume and became more of a problem, anti-spam techniques were developed to combat it. Tools to block spam were developed by a group of professionals. These tools were not always automated, but when used by system administrators of large sites, they successfully filter spam for a large number of users. In response, spammers evolved their techniques to increase the volume of spam delivered by working around and through the filters. As spam filters improved, spammers designed other methods of bypassing the filters and the cycle repeated. This resulted in the development of both spam and anti-spam techniques and over the years. This evolutionary process continues today.

Anti-spam tools use a wide range of techniques to reduce the volume of spam received by a user. A number of these techniques will be described in the following section. There are several anti-spam techniques based on Open Source tools that we will examine in the light of the various techniques it uses to filter spam.

### 2.1 • Spamming Techniques

Spammers have developed a complex arsenal of techniques for spamming. Important spamming techniques described in the following paragraphs.

## 2.2 • Open Relay Exploitation

An open relay is a mail server that allows anyone to send email through it. Spammers use open relays to send spam without the email being traced to its true origin.

## 2.3 • Collecting Email Addresses

Spammers had to find email addresses in order to send spam. Methods for collecting email addresses include harvesting from the Internet and Internet guessing email addresses. They use a variety of techniques, such as scanning newsgroups to simply find email addresses.

## 2.4 Hiding Content

Most people can detect spam from the email or server. It is often easy to discard warning emails without even looking at the body. One technique used by spammers is to hide the true content of their emails. Often, the subject of an email is a simple "Hi"; alternatively, an email might be a reply to a previous email, for example "Re: tonight". Other tricks that spammers use include using random names for example, by alluding to a credit card or loan missed payment or work-related subject.

As spam filters become more obvious, spam words, such as "Viagra", are spammers deliberately include misspelled words that are less likely to be filtered out; for example "Viagra" might become "iagra" or "V-iaggr@".

Although the human mind can easily translate the meaning of misspelled words, "being unconscious", a computer program will not associate these words with spam.

## 2.5 • Statistical Filter Poisoning

Statistical filter poisoning involves including many random words within an email to confuse a statistical filter. Statistical filters are described in the Anti-spam Techniques section.

## 2.6 • Unique Email Generation

To combat email content databases, which store content of known spam emails during the rounds, spammers generate unique email addresses to confuse the email content database; the spammer only needs to change one random character in the main body of the email. One technique is to use the recipient's name within the body of an email.

## 2.7 • Trojaned Machines

Spammers are limited by the speed of their Internet connection, be it DSL or dial-up. They are also directly traceable through ISP records. A recent trend among spammers is to use PC virus technology to infect innocent user computers with virus-like programs. These programs send

spam from the infected machines' PCs. Such infection is commonly known as a Trojan, after the surreptitious means.

Story of the Greeks invading the City of Troy by

### 3 - Anti-Spam Techniques

As the techniques to spam have become more sophisticated, so have the techniques to detect and filter spam from legitimate email. The main techniques are described in the following points. These techniques can be used on the email server by a system administrator.

or an anti-spam gateway can be purchased from an external vendor.

#### 3.1 • Keyword Filters

Filters are based upon common words or phrases in an email for example 'buy', 'last Chance', and 'Viagra'. Open source software includes a variety of keyword filters and allows easy configuration of rules.

3.2 • Open relay Blacklist relay (ORALS) identifies those relays that have been reported and added to these lists after being tested. Anti-spam tools can query open relay blacklists and filter out emails originating from these sources. Some open source software can integrate with several relay blacklist.

#### 3.3 • ISP Complaints

It has always been possible to complain to an ISP about a spammer. ISPs take complaints seriously, give a single warning, and after another complaint, they terminate the account of the offender. Other ISPs take a less active approach to spam that will rarely stop a spammer. Spammers naturally gravitate towards ISPs that are lenient with spammers.

ISP Complaints remain a manually managed technique, due to the effort that might be wasted if an automatic report is sent and email redirected to a not-spam. The website <http://www.spamcop.net> can examine an email; determine where spam reports should be directed, and send appropriate messages of complaint to the corresponding ISPs.

#### 3.4 • Statistical Filters

Statistical filters are those that learn common words in both spam and ham. Subsequently, the data collected is used to examine emails and determine whether they are spam or ham. These filters are often based on the mathematical theory called Bayesian analysis. Statistical filters need to be trained by both ham and spam emails through the filter to learn the difference between the two. Ideally, a statistical filter should be trained regularly, and some anti-spam software allow statistical filters to be trained automatically.

#### 3.5 • Email Header Analysis

The software that spammers use often generates unusual headers in the emails produced. Anti-spam software detects these unusual headers and uses them to separate from ham. Some open source software includes many email tests.

#### 3.6 • Non-Spam Content Test



There are possibilities that ham emails could inadvertently trigger some anti-spam tests. For example, many emails are legitimately but unfortunately routed through a blacklisted open relay. Non-spam content test that an email is not spam. They are usually created for a specific individual or organization.

Spam content tests are rarely shared in public, as they are too industry-specific, and should not get into the hands of others who would use this information to their advantage. Software allows users to create rules that will subtract from the score of an email if certain content is received. An administrator might add negative rules for the names of products sold by one company or for a company-related person.

## 35 Whitelists

Whitelists are the opposite of blacklists — lists of email senders who are trusted to send ham and not spam. Emails from someone listed on a whitelist will normally not be marked as spam, no matter what the content of their email.

12

### 3.8 • Email Content Databases

Email content databases store the content of spam emails. These work because the same spam email will often be sent to hundreds or thousands of recipients. Email content databases store these emails and compare the content of emails to that contained in the database. A user reporting a spam email to such a database will assist all other users of the service. Some open source software can integrate with several email content databases.

### 3.9 • Sender Validation systems

A slightly different approach to spam is taken by sender validation systems. In these systems, when an email is received from an unknown source, the source is sent a challenge email. If a valid response is received to such an email, then the sender is added to a whitelist, the original email is delivered to the user, and the sender sent a challenge again.

This is effective as spammers forged sender and reply-to addresses and do not receive replies to the spam they send out. Consequently, the challenge is never received. In addition, spammers do not have the time to validation requests.

Some systems cleverly integrate with the user's outgoing email addresses to automatically add known contacts to a whitelist. Validation systems are proprietary and involve a small initial fee.

Sender validation systems are used when subscribing to a mailing list. Few email list administrators will respond to a challenge, so the user might end up not receiving emails from the list. With most systems, it is possible to manually add addresses to the whitelist to avoid a challenge.

or reverse required. but in the case of mailing lists, the addresses that emails are sent to may not be known until emails are received. Some open source software does sender validation

3.10 • Sender Policy Framework (SPF) the Sender Policy framework (SPF) can be used to ensure that an email is from a valid source. It validates that a sending email from a particular email address is permitted to send email from their current machine. SPF is a recent development and is being introduced relative, "quintessential". It uses additional Domain Name System (DNS) record to state which mail server can send email for a domain. Some open source software uses the current draft standards for

## 4 - Detecting Spam

### 4.1 • Contents Tests

Content tests analyze the message part of the email, and sometimes the headers. These tests typically look for key words or phrases within emails. Usually, when using content tests, a Scoring System is used. It is not just words associated with spam emails to appear in legitimate emails, so a score count of suspicious words is accumulated each email. Each word associated with spam increases the overall score of email. The final score is compared with a predefined threshold; this is used to decide whether an email is spam or ham.

Content tests need not focus on single words: phrases and sequences of punctuation are used. The symbols, phrases, and other symbols tested are normally generated by a developer, who analyzes spam and manually creates tests.

Sometimes the message headers are examined as part of a content test. The message headers include dates, time, and other attributes, such as the mail application used. Often, spam creation programs contain errors or misspellings in their headers that can be caught by spam.

Spammers attempt to avoid detection by deliberate misspelling and varying content slightly in each spam or spam run.

A simple example of a content test would be to locate the word "Viagra" within an email. A more complex content test is the sequence of characters "v?i?r?a?g" where the ? represents any character, and one or more instances of not present at all. For example, VIAGRA, V A G R A and would all match.

### 4.2 • Header Tests

Header tests focus on the message headers. The tests are concerned with detecting fake headers and determining whether a message has been routed via an relay.

For example, a header test could flag an email that appears to have been sent over 72 hours ago, or sent at a future date. Most email servers have accurate clocks. However, spammers frequently use borrowed PCs, which may have inaccurate clocks and spam

Messages might have dates that are in the past or the future. Examining email headers is described in **more** detail later in the section.

These tests use up considerable amounts of CPU, Memory, and disk I/O resources.

### 4.3 • DNS - Based Blacklists

There are many DNS — based blacklists (DMBLs). These are also known simply as **blacklists** or

**blacklists**. They provide that is used MIA (Mail Transfer Agent Of Mail Transfer). The store and forward part of a messaging system (like email) and spam filters to indicate sites that are related to spammers. An MIA or spam filter may use one or more **blacklists**. Some open source software can use blacklists to filter spam. Blacklists can generate, be placed in one or more of these categories.

. A list: Of known open relay

- A list of known of spam
- A list of sites hosted by an that spammers in some way

Ever' has unique policies for adding and removing from the list. Some are very aggressive and block sources of spam, but also any address served by the same Internet Service Provider. The intention of this approach is to force ISPs to do business with spammers and thus force them out. This approach, called the Internet black hole of death has been used with success against major ISPs in the past.

Blacklists provide a spam filter or MTA with the ability to query the blacklist to see a particular IP address is listed. If the IP address is listed, then incoming email from that host is often rejected -

Generally, IP addresses are listed on a blacklist only if they have been reported. Reporting is either done by a human after examining the headers of an email message, or by an automatic system. Blacklists remove addresses from their list after a period of time, while **some** wait for proof that a spam problem has stopped,

Some blacklists will test a site to see if it is running an open relay. The tests are usually a task. **probing port25** in a way as the manual test.

It is the responsibility of system administrators and end users to not spammer's react to a listed blacklist. Some software submits suspected addresses to a blacklist automatically. • Which is a dangerous approach. Automatic systems can get confused if unforeseen circumstances and a relay blacklist (ORBL) could get flooded by false reports. It is better to provide the user with an option to report a relay to a blacklist than do it automatically. This relies on having software developers provide an option for users, and that • **users** having knowledge to determine whether to submit the request.

Blacklists generally use network I/O rather than CPU, and disk I/O. They use less **network** I/O than is used when receiving an email. These tests are suited to parallel systems (processing many emails at once) as the results take time to be retrieved and the machine is free to use other resources (CPU, memory, and disk I/O) for other tests.

#### 4.4 • Statistical Test

Various statistical techniques can be used to filter spam. These techniques involve a **training** phase, where a database of spam and ham emails is taught to the filter or passed through it to identify typical characteristics of spam and ham. This allows future emails to be identified based on the learning from emails. The various statistical techniques vary in their **choice** of tokens and the algorithms they use to predict whether an email is spam or ham. The tokens

used are words, but can include emails headers. HTML markup within emails. and others characters such as punctuation marks.

Statistical rely on regular training, They use the knowledge gained in training to estimate the probability that new emails are As spam change, the filter must adapt in order to continue to detect the spam

Statistical tests are resource intensive. using CPU, memory and disk I/O,

#### 4.5 Message Recognition

Often, a spammer will send the same message to many recipients. Although message headers may differ in each email, an email with the same body may be sent to many recipients. This has led to the creation of several anti-spam networks that contain a database of spam emails. Comparing incoming emails with the contents of this database it is possible to quickly filter out known spam messages.

To avoid sending the whole email across the network and comparing each character of line, a hash value is calculated and used. Hashing is a mathematical process that creates a unique signature from a larger message. It is very unlikely that two messages will have the same hash value, and so comparing hashes is statistically the same as comparing the whole message. As the hash values are much shorter than an email message, comparing hash values is significantly faster than comparing the whole message,

The calculation of a hash value is a CPU-intensive task, and there is some I/O and related latency involved while querying the database. This test is suited to parallel processing.