

Fighting Spam with Statistics

Spam - unsolicited commercial email - is a complex and growing problem, and threatens to derail the internet revolution. Joshua Goodman and David Heckerman of Microsoft Research describe some statistics-based methods for blocking spam, first by distinguishing it from wanted mail, and then by constructing puzzles they propose to use to challenge suspected spammers.

Some statistics about spam

Spam is a huge problem. A few statistics can tell us just how huge:

- Brightmail has reported that over 50% of mail on the internet is now spam.
- A recent Infoworld poll identified spam as the number 1 “information technology disaster of the past year.”
- From a report by the United States Federal Trade Commission [1]:
 - 66% of spam had false information somewhere in the message;
 - 18% of spam advertises “Adult” material.
- From a report by the Pew Internet and American Life Project [2]:
 - 25% of email users say that spam has reduced their use of email;
 - 12% of users spend a half hour or more per day dealing with spam.

Statistics also help explain why spam is such a problem.

- It costs only about one hundredth of a cent to send spam.
- 7% of email users say they have bought a product advertised in unsolicited email [2].

Given the tiny costs of spamming, even a tiny response rate makes spam economically viable.

Identifying spam—harder than it looks

Around a year ago, MSN’s Hotmail service proposed blocking all spam entering their system that was allegedly from Hotmail. Spammers often fake the address that their mail is from, including pretending that such mail is from Hotmail itself. Because there are no standards on the internet to say which IP addresses are allowed to send mail from which domain names, we usually cannot tell when a sender’s address is faked. (Recent proposals such as Microsoft’s Caller-ID standard will help fix this problem in the future.) The one exception today is that Hotmail knows its own valid internet addresses, and knows that mail from outside those addresses that says From Hotmail is fake. It sounded reasonable to block all such deceptive mail. There are, however, a few cases where this faking may be legitimate. For instance, when your friend sends you an on-line birthday card via email, the From address has your friend’s email address (perhaps at Hotmail), even though the email actually originates on the greeting card company’s servers, and comes from their internet address. Fortunately, we had acquired a dataset of spam and good mail. We were able to show that for every four spams allegedly from Hotmail that would be blocked, one good message would be lost – the exceptions were not as rare as had been thought. A 4-to-1 tradeoff is not close to tolerable. This simple analysis was

key to preventing an unacceptably large number of deleted good messages, including those birthday cards.

While there are no hard-and-fast rules for marking mail as spam – even a faked From address is only moderately bad – by combining different indicators together, we can often be almost certain that mail is spam, and safely delete it, or very confident that it is spam, and place it in a junk folder. The question is how to combine these indicators.

One of the most popular techniques for stopping spam is the Naïve Bayes method [3], often mistakenly simply called Bayesian Spam Filtering. In the Naïve Bayes method, we try to determine the probability that a given message is spam, or good. We start by using Bayes rule, which in this case tells us that

$$P(\text{spam}|\text{message}) = P(\text{message} | \text{spam}) \times P(\text{spam}) / P(\text{message}),$$

where

$$P(\text{message}) = P(\text{message} | \text{spam}) \times P(\text{spam}) + P(\text{message} | \text{not spam}) \times P(\text{not spam}).$$

$P(\text{spam})$ is simply the prior probability that any given message is spam. For instance, if half the messages you receive are spam, $P(\text{spam}) = 1/2$. Similarly for $P(\text{not spam})$. Also, we need to compute $P(\text{message}|\text{spam})$ and $P(\text{message}|\text{not spam})$ – these are respectively the probability of receiving any of the billions of possible spam messages and the probability of receiving any particular good piece of email, of the infinite possibilities. Again, this is difficult if not impossible. However, we can make some approximations. In particular, we will assume that the probability of every word in the message is independent of every other word, conditioned on knowing whether the message is spam. For instance, we assume that the probability of “click” is independent of the probability of “here,” given that the message is spam. Clearly, this approximation is not a great one, which is why this technique is called Naïve. However, once we have made this assumption, it is easy to compute the message probabilities. In particular,

$$P(\text{message} | \text{spam}) = P(\text{first word} | \text{spam}) \times P(\text{second word} | \text{spam}) \times \dots \times P(\text{last word} | \text{spam})$$

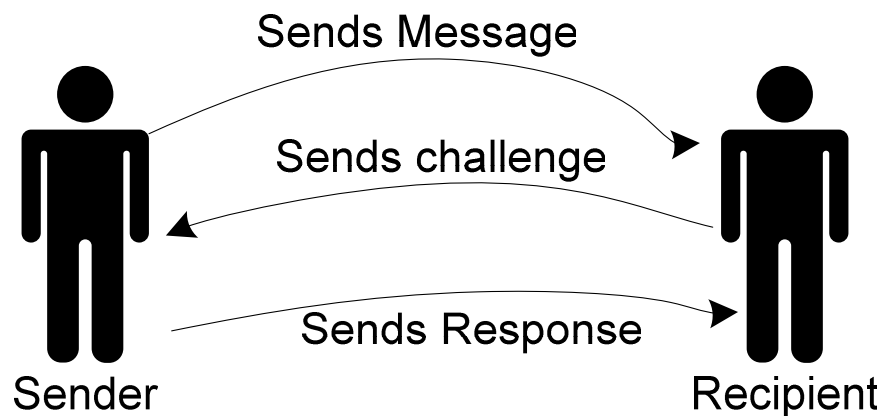
$$\dots \times P(\text{last word} | \text{not spam})$$

What is the probability of a particular word given spam? For each word, e.g. “click”, we simply count how many times the word occurred in all spam messages in a “training corpus” (a set of messages known to be spam or good) and how many spam messages there were overall. The ratio is an estimate of the probability that any particular spam word is the word “click.” We can do the same for words in non-spam messages.

Naïve Bayes is just one of many machine learning techniques, often used because it is the easiest to implement. It can be surprisingly effective, given the simplicity of the model. If an individual user is willing to carefully label thousands of his messages as spam or good to train the model, it can result in excellent performance for that user on his own mail. However, the independence assumption in Naïve Bayes limits its effectiveness. A Naïve Bayes model will tend to be confused by common spam phrases like “click here to unsubscribe” that sometimes also occur in good mail, because Naïve Bayes does not model the fact that these words occur together. If each word is 10 times as likely in spam as in good, a Naïve Bayes model will think any such message is 10,000 times more likely to be spam than good ($10 \times 10 \times 10 \times 10$)—a large overestimate of the true ratio. Other model types, such as Neural Networks, Graphical Models, Logistic Regression, and Support Vector Machines (SVMs) all make fewer assumptions, and can model these kinds of relationships between words implicitly or explicitly, at the expense of more complexity. At Microsoft, we use these more complex model types.

Challenge-Response systems

Probabilistic filters can do an excellent job of identifying spam. Nonetheless, like all filtering techniques, they are not perfect: they always occasionally mark some good mail as spam. At Microsoft Research, we have been exploring and adapting an idea due to Cynthia Dwork (now at Microsoft Research, Silicon Valley) and Moni Naor (at the Weizmann Institute of Science in Israel.) They proposed using “computational puzzles” to stop spam [4]. We are researching a variation on this idea, in which it is used in combination with a probabilistic filter [5]. If you send me mail, my filter examines it. If it appears good, it goes in my inbox. If it looks like junk, it goes in my junk folder. To make sure I eventually see any good mail sitting in my junk folder, my email program automatically sends you a “challenge”: a request to solve a computational puzzle. Your email program must spend perhaps fifteen seconds to solve this puzzle. Your email program solves this puzzle automatically, in the background, and you do not even notice it. It sends the puzzle solution back to my program, which verifies the solution, and moves your original message to my inbox.



If you are a legitimate sender, you are willing to spend the time to do this. But if you are a spammer, you cannot afford to solve these puzzles – you would be able to answer fewer than 6000 puzzles per day per computer, compared to the millions of messages you can normally send.

What do these puzzles look like? One solution is to use a hash function. A hash function is a function that takes a string of letters as input, and returns a number. The number is a deterministic, but nearly random, function of the input string. Slight changes in the input string result in large and unpredictable changes in the output of the function. The message recipient – the challenger – sends a challenge string based on the actual message to the message sender. The message sender must find some other string such that, when put in front of the first string, it results in a hash value of, say, zero. If the hash function returns a number between 1 and 15,000,000, the sender must try about 15,000,000 different strings to find one that results in a zero value. A typical hash function might require one microsecond – one millionth of a second. It would thus take, on average, fifteen seconds to find a puzzle solution. Meanwhile, the recipient only needs to verify that the solution is correct, meaning running the hash function once: he can do this in just one microsecond.

Designing better puzzles

There is a problem with the simple puzzle described above: it has high variance. On average, the puzzle will require fifteen seconds to solve, but sometimes it will require only one second, sometimes thirty seconds, sometimes forty five seconds, and even occasionally more than a minute. This can be solved using a suggestion from Cynthia Dwork and Andrew Goldberg. We can use an easier puzzle, perhaps one that returns a random number between one and one million, but require the sender to find fifteen solutions. This lowers the variance by about a factor of fifteen. Variance is reduced linearly as the number of puzzles increases (assuming we keep the total expected time constant.) If we used a huge number of puzzles, say 100, the variance would become negligible. But the size of the solution would also be large, even larger than the original message in some cases.

An alternate solution, also due to Dwork and Goldberg is to use a hash which returns a number between, say, zero and 1,000,000,000 and then require the sender to find, say fifteen solutions all with the same value, and that value must be between zero and 1,000. For instance, he could find fifteen solutions with a value of zero, or fifteen solutions with a value of one, etc. We can think of this as balls dropping into bins at random, and waiting for one of the bins to get fifteen balls in it. This problem turns out to have much lower variance than simply finding fifteen solutions.

Conclusion

Spam is an extremely complex topic, and no single article can hope to describe all the ways to fight it. We've described here those techniques that rely most heavily on statistics, but we are exploring and pursuing many other ideas as well [6]. For example, we are currently exploring new industry standards that can help stop spam. In particular,

our challenge-response systems are expensive for very large, legitimate senders. We are exploring standard ways for these senders to become certified as non-spammers.

No single solution will stop spam. Laws will solve be one part of the solution. Improved filters using better statistical analyses will be another part. Finally, all legitimate senders will be able to bypass any filter mistakes: large senders through certification, smaller senders through low-variance computational puzzles. We anticipate a future in which the vast majority of spam is stopped, and all legitimate mail is delivered.

Joshua Goodman is a Researcher in the Machine Learning and Applied Statistics group at Microsoft Research. He has been on loan to Microsoft's Anti-Spam product team, since its inception. His previous work has been on language modeling (predicting word sequences) and on fast algorithms for logistic regression.

David Heckerman is founder and manager of the Machine Learning and Applied Statistics Group at Microsoft Research. Since 1992, he has been a Senior Researcher at Microsoft, where he has created applications including junk-mail filters, data-mining tools, handwriting recognition for the Tablet PC, troubleshooters in Windows, and the Answer Wizard in Office. His work includes Bayesian methods for learning probabilistic graphical models from data. David received his Ph.D. from Stanford University in 1990 and is a AAAI Fellow.

Bibliography:

- [1] United States Federal Trade Commission, False Claims in Spam, April 30, 2003.
- [2] D. Fallows. Spam: How it is hurting email and degrading life on the Internet, *Pew Internet and American Life Project*, October 2003.
- [3] M. Sahami, S. Dumais, D. Heckerman, and E. Horvitz. [A Bayesian Approach to Filtering Junk E-mail..](#) AAAI'98 Workshop on Learning for Text Categorization, July 27, 1998, Madison, Wisconsin.
- [4] C. Dwork and M. Naor, "[Pricing via Processing or Combatting Junk Mail](#)", *Lecture Notes in Computer Science* 740 (Proceedings of CRYPTO'92)}, 1993, pp. 137--147.
- [5] J. Goodman and R. Rounthwaite "SmartProof"
<http://www.research.microsoft.com/~joshuago/smartproof.pdf> , 2004
- [6] J. Goodman. "Spam Technologies and Policies",
<http://www.research.microsoft.com/~joshuago/spamtech.pdf> , 2004

Sidebar: How to Get Less Spam

- Choose an email address that is hard to guess. Many spammers use “dictionary attacks” where they try billions of common words separately and in combination, for example: Bob, Bob, Bob1970, StatisticsNut, etc. Choose a longer name with more words, like BobStatisticsNut1970, and a spammer is much less likely to guess it.
- Don’t give your email address to a company you don’t trust. Some companies sell the addresses they collect.
- Don’t put your email address on web pages. Many spammers “crawl” the web, looking for addresses that they then add to their lists. If you need to put your email address on the web, disguise it, or put it in an image. Spammers don’t use image recognition software (yet!)