

マスメールデータベースと

それを用いたマスメール検出システム

松浦広明[†] 齋藤彰一[†] 上原哲太郎^{††} 泉裕[‡] 和田俊和[†]

[†]和歌山大学システム工学部 ^{††}京都大学大学院工学研究科

[‡]和歌山大学システム情報学センター

概要 マスメールによる被害は日に日に増大している。本研究では、マスメールを検出するために、マスメールデータベースと POP サーバからなるシステムを構築した。データベースでは、マスメールを自動的に収集し、各マスメールを正規化・解析したのち計算されたチェックサムを蓄積する。POP サーバでは、転送するメールのチェックサムを同様に計算し、データベースに問い合わせることで、マスメールかどうかを判断する。このシステムでは、チューニングや学習は不要である。処理速度が速く、データベースの内容が正しい限り誤検出は起きないのが大きな特徴である。このシステムを和歌山大学内の POP サーバで稼働させた結果、マスメールの認識率はおよそ 80-90%であった。

Mass-mail Database and Mass-mail Detection System.

Hiroaki Matsuura[†] Shoichi Saito[†] Tetsutaro Uehara^{††}

Yutaka Izumi[‡] Toshikazu Wada[†]

[†]Faculty of Systems Engineering, Wakayama University

^{††}Graduate School of Engineering, Kyoto University

[‡]Center for Information Science, Wakayama University

Abstract Problems caused by mass-mail have been more and more serious. We have constructed a system for detection of mass-mails, which is composed of a mass-mail database and a POP server. The database collects mass-mails automatically, and registers checksums calculated by analyzing the contents of the each mail. The POP server computes checksums of the handling mail and asks the database whether the checksums have been registered as mass-mail. This system requires no efforts of tuning and learning, runs fast, and does not mistake in mass-mail recognition as long as the database is valid. The server has run on a POP server in Wakayama University. Almost 80-90% of mass-mails have been detected successfully.

1. はじめに

不特定多数に大量に送信されるメール（マスメール）は、受信者へ大きな負担、公序良俗に反する内容、メールトラフィックの増大、ウイルスメールの拡散など、多くの問題を世界的な規模で引き起こしている。日に日に増え続けるマスメールを一刻も早く止めるために、様々な対応策が検討されている。

マスメールを判別する場合、コンテンツ同士を単純に比較する方法を用いることはできない。なぜなら、マスメールの送信者は、html の仕様を利用してコンテンツの表記方法を変化させたり、本文とは関係のない文字列をコンテンツに加えたりするなどの方法で、「意味は同じだが違う」コンテンツを大量に作成し、送信するためである[1]。

また、コンテンツのすべてではなく任意のキーワードを抽出・検索する方法がある。この方法は考え出された当初は効果的であったが、今では有用ではなくなった。マスメール送信者は、キーワードの綴りをわざと書き換えて、検出されないようにしているためである。

以上のような判別の困難さを解決するため、メールのコンテンツに含まれる各トークンの統計的性質に基づいてマスメールの識別を行う方法が盛んに研究されている。中でも多いのが Bayesian Filter を利用した種々のメールフィルタリングツールである[2]-[6]。これらのツールは、個人単位で利用することができ、チューニングできるという特長がある。その反面、学習させるまでに時間がかかること、学習のデータによって識別の精度が左右されること、処理に時間がかかること、日本語のような分かち書きされていないテキストは不得意であ

ること、などの不具合がある。また、最近のマスメール送信者は、一般的なビジネス文書やニュース記事を含むようなマスメールを送ることによって Bayesian Filter に誤判定・誤学習を起こさせるなど、マスメールが検出されないよう工夫しつつある。

マスメール検出システムは、処理が速く、誤検出がなく、ユーザが手軽に使えるようなシステムであることが望ましい。本研究では、マスメールを蓄積するデータベースを活用して、マスメールの判別を行う。

2. マスメール判別システム

2-1 概要

本研究のシステムは、マスメールを蓄積するサーバと、それを活用する POP サーバの二つのサーバからなる。

マスメールデータベース（図1）では、マスメールを自動的に収集する。集められたマスメールに対して、コンテンツのランダム要素や非表示文字を排除する「テキストの正規化」および「URL の抽出」を行う。得られたそれぞれのデータから MD5 法によりチェックサムを計算し、データベースに蓄える。データベースは、リモートからの利用が容易な DNS を用いる。

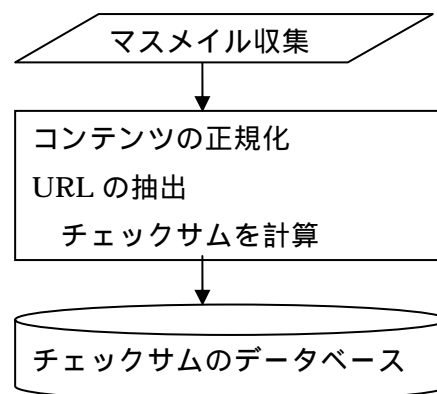


図1 マスメールデータベースの概要

POPサーバ(図2)では、処理中のメールに対してデータベースと同様の処理でチェックサムを計算し、データベースに問い合わせる。チェックサムがデータベースにあれば、処理中のメールをマスメールだとみなし、メールヘッダにキーワードとして"X-MMS-RESULT: 1"を挿入する。

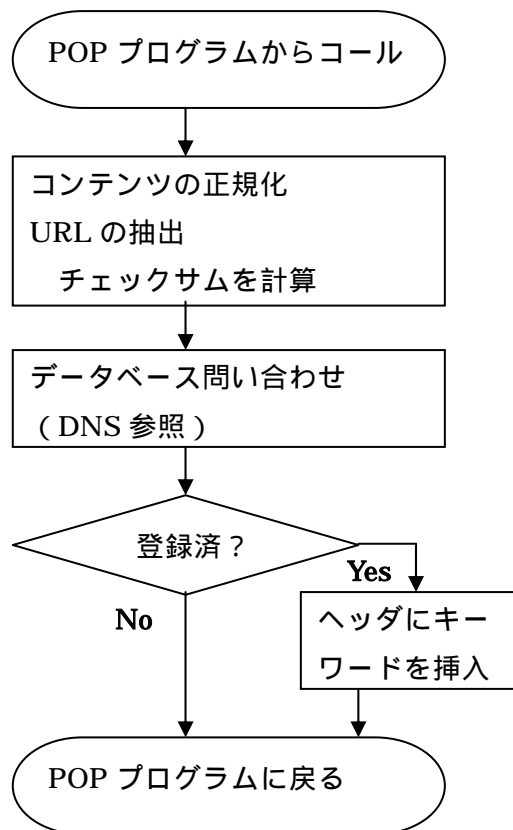


図2 POPサーバの概要

2-2 本システムの特徴

本システムでは、次のような利点が挙げられる。

- 1) データベースが正常である限り、誤認識は起こらない。
- 2) 認識器のような複雑な計算が不要なく、高速に動作する。
- 3) ユーザはメールリーダーで「メールヘッダ中のキーワードによる、フォルダ振

り分け」を設定するだけでよい。

- 4) チェックサム計算の部分をライブラリ化しているため、procmail のようなメールフィルタと組み合わせて実行したり、メール整理ツールを作成したりと、包括的なマスメールフィルタリングができる。

3 マスメールデータベース

3-1 マスメールの収集

現在行っているマスメールの収集方法を以下に述べる。

- ・マスメールを受信するためだけに用意したメールアドレス(おとりアドレス)から転送を行う。
- ・マスメールを受けたいメールアドレスをWEBページのコメント、もしくは、実質的にクリックできないmailtoアンカーに埋め込む。
- ・マスメールに「送信を止めるためにはメールを送るように」など指定されている場合は、マスメール送信者にメールを送る(これによって、マスメールの送信者はそのメールアドレスが有効であることを知る)。
- ・新規ドメインを取得し、そのドメイン宛てのすべてのメールをマスメールとして収集する。

3-2 データベースの状況

集められたマスメールの累積登録数(図3)、一日あたりの登録数(図4)、および時間帯あたりの登録数(図5)を示す。現在4万通弱のマスメールが収集されており、一日あたりの登録数は増加傾向にある。また、時間帯による変化はほとんどないことがわかる。

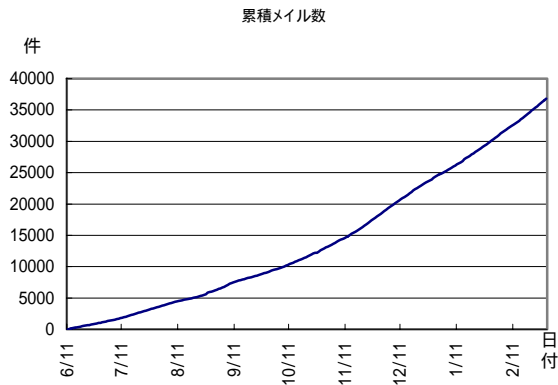


図3 累積メール数

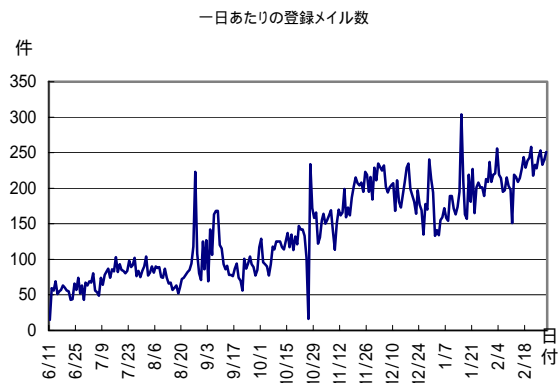


図4 一日あたりの登録メール数

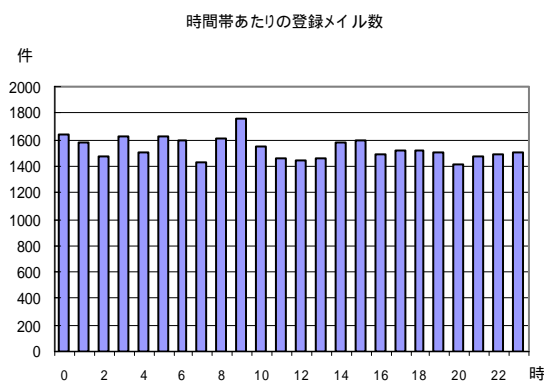


図5 時間帯あたりの登録メール数

4. チェックサムの計算

前述したように、データベースおよびPOPサーバでは、「テキストの正規化」および「URLの抽出」を行い、それぞれの結果に対して、チェックサムを計算する。

「URLの抽出」において、本来マスメイルとは関係ないURLがマスメイルに書かれている場合、誤ってそのURLをデータベースに登録してしまう危険性がある。そこで、「ホワイトリスト」を外部ファイルに定義して、マスメイルとは無縁のURLを保護している。

この一連の処理を行うルーチンはライブラリ化してあるため、容易にPOPサーバやメールフィルタリングツールに組み込むことができる。現在、ipop3d、qmail-pop3に組み込んで、動作を確認している。

チェックサム計算にかかる時間は3万通のメールでおよそ80秒であり、一通あたり0.00263秒であった（CPU: Intel Pentium 4 2.4GHz, OS: Red Hat Linux 7.3, Memory: 1GByte）。なお、アルゴリズムの詳細は、マスメイル送信者に本システムを回避する情報を与えかねないため公開できない。

5. POPサーバの運用

5-1 認識率

4.で述べたライブラリをipop3dに組み込み、和歌山大学システム工学部にて運用している。ユーザにより違いはあるが、マスメイルのおよそ80~90%を正しく認識している。しかしまだ10~20%のマスメイルを認識していないことから、データベースの拡充や正規化アルゴリズムの改善などが必要である。

一方、正常なメールをマスメイルと誤認識することはない。運用当初はデータベースに正常なメールが紛れ込んでいたことや、ホワイトリストが不十分であったため数件生じたが、現在では皆無である。

5-2 メイルトラフィックの状況

一日における全 POP 数と検出マスメール数 (図 6) , 全 POP 数に対する検出マスメール数の割合 (図 7) を示す . 平日・休日の周期がみられ 検出率はおよそ 6 ~ 12% を推移している .

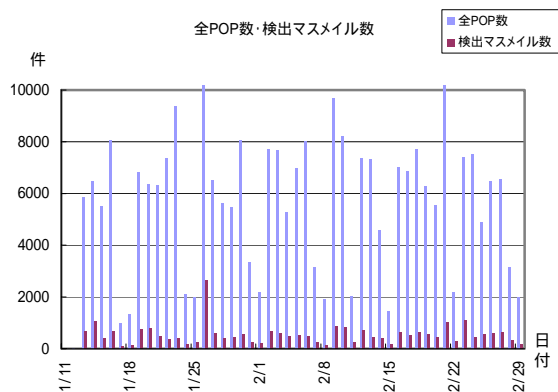


図 6 全 POP 数と検出マスメール数

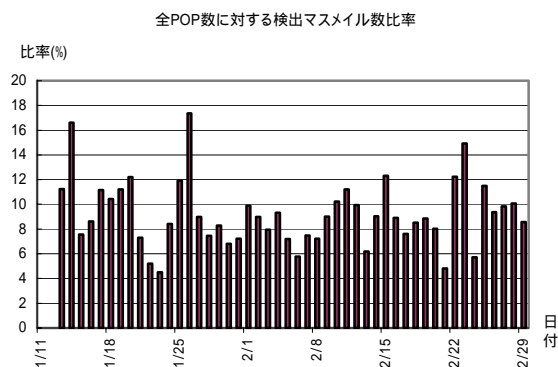


図 7 全 POP 数に対する
検出マスメール SPAM 数比率

5-3 マスメールの検出時に用いられたチェックサムの種類

一日におけるマスメールの検出が , 正規化したコンテンツのチェックサム (BODY) でなされたのか , 抽出した URL のチェックサム (URL) でなされたのかを検討した .

全体の検出マスメール数に対する , BODY での検出率と URL での検出率を , 図 8 に示す (一つのマスメールが BODY と

URL の両方で検出される場合もあるため , BODY の比率と URL の比率との和が 100% 以上になることがある) . この結果から , URL による検出が BODY による検出のおよそ 2 ~ 3 倍ほど高い比率を示していることが分かった . URL はマスメールの判別に大きな手がかりとなっていると言える . BODY の検出率が低いのは , コンテンツは URL よりもランダム要素が大きく , 正規化の段階でランダム要素を完全に除去できなかった理由もあると考えられる .

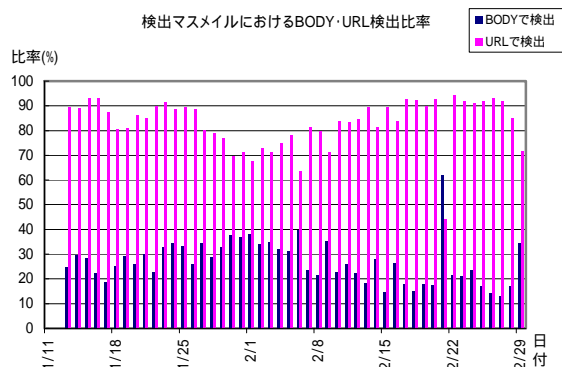


図 8 検出マスメールにおける
BODY ・ URL 検出比率

5-4 チェックサム登録日時とマスメール検出日時との時差

POP サーバが任意のメールをマスメールと判断した日時と , 判断の理由となったチェックサムがマスメールデータベースに登録された日時を比較した . これによってチェックサム登録からマスメール検出までの時間を知ることができる . ただし , データベースには同じチェックサムを持つマスメールが一つ以上登録されているため , 求める解は複数ある . 解のうちの最短の時差を求めた . その結果を図 9 に示す .

これによると , POP サーバでマスメールだと判断されたチェックサムのうち , 50%

超が 24 時間以内に届いたマスメールから登録されたチェックサムであり、また、ほとんどすべてが 10 日以内に届いたマスメールから登録されたチェックサムである。このことから、POP サーバで検出されるマスメールとマスメールデータベースに届くマスメールは、ほぼ同時期に届けられていることがわかる。より多くのマスメールをより早くデータベースに登録すれば検出率の向上が見込まれる。

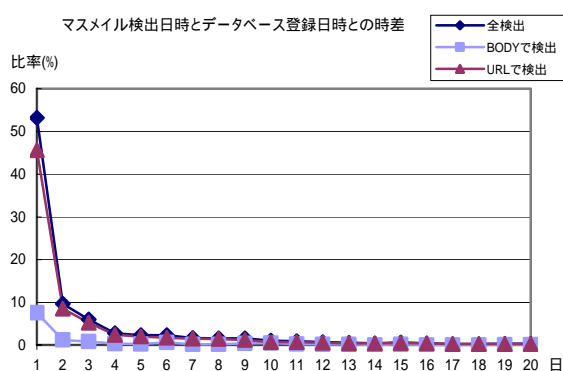


図9 マスメール検出日時とデータベース登録日時との時差

6. まとめと今後の展開

自動的に収集したマスメールを正規化・解析してチェックサムを計算し、そのチェックサムデータを配信するマスメールデータベースと、それを参照してマスメールを検出する POP サーバとを用いて、マスメールを検出するシステムを構築した。高速性、正確性、簡便性などが秀でている。

現在よりマスメール検出率を上げるために次の対策を行っている。

- 1) ランダム要素を取り除くアルゴリズムの性能を上げる。
- 2) 検出率を上げるためには数多くのマスメールを集める必要がある。spam archive[7]のマスメールデータベー

スを利用することを検討中である。

- 3) 現アルゴリズムでは、「正規化したコンテンツ」と「抽出した URL」でチェックサムを計算している。メールアドレスや電話番号をも抽出する新アルゴリズムをテスト中である。現行版よりも高い検出率を得ている。

これらを反映させた新システムは、今年度中に運用を開始する。

本システムは和歌山大学システム工学部で運用中であるが、来年度は、和歌山大学・京都大学・大阪大学など各研究機関のメール管理者の協力を仰ぎ、運用を拡大してゆく予定である。

参考文献

- [1] John Graham-Cumming, "The Spammer's Compendium", 2003 SPAM CONFERENCE (http://sourceforge.net/docman/display_doc.php?docid=14869&group_id=63137).
- [2] "A Plan for Spam", <http://www.paulgraham.com/spam.html>.
- [3] "Getting Started with Mozilla Mail Spam Filtering", <http://www.mozilla.org/mailnews/spam-howto.html>.
- [4] "McAfee Security", <http://www.mcafee.com/myapps/msk/default.asp>.
- [5] "Spam Assassin", <http://spamassassin.org/>.
- [6] "POP File", <http://popfile.sourceforge.net/>.
- [7] "SpamArchive.org", <http://www.spamarchive.org/>.