

# 特徴抽出によるスパムメールフィルタリング性能の向上

## Improvement of spam mail filtering performance by feature extraction

渡邊 隆志

Takashi Watanabe

### <要旨>

スパムフィルタリング機能が搭載されている電子メールクライアントにスパムメールを学習させたにもかかわらず、それらに類似したスパムメールが非スパムメールとして受信箱に入ってしまう現象が見られる。本研究は、スパムメールの特徴に着目し、それをスパムメールの判定基準にすることで、非スパムメールと誤判定されたスパムメール(通過スパムメール)を削減することを目的とする。従来の通過スパムメールを調査したところ、フィルタリングされたスパムメールとの件名と本文が類似していて、さらに本文中に URL が含まれており、かつ本文の文字数が少ないといった特徴が多いことを確認した。そこで、文字列の距離指標として使用されている Jaro-Winkler 距離で測定した件名と本文の類似度、本文の文字数、URL の有無という 4 つの特徴量をサポートベクターマシンで分類する新しい手法を提案する。実験を行った結果、既存のスパムフィルタリングに提案手法を追加する方法だけでなく、既存のスパムフィルタリングを提案手法に置き換える方法も有効であることが確認できた。

### <Abstract>

In these days most mail user agents include the spam filter function classifying spam mails and normal mails. However, some spam mails tend to be misdetected and pass this filtering easily and be stored in normal mail boxes. So the purpose of this study is set to reduce this decision error about spam mails (passed spam mails), based on features of spam mails. From the survey of the passed spam mails, we identified many passed spam mails have the following features; having similar subject and body, including URL in the body, and small number of characters in a body, in comparison with the normal mails. In this paper, regarding the features, we propose the new method of machine learning filtering using Jaro-Winkler distance, which focusses similarity in terms of subject and body, the existence of URL in a body, and the number of characters in a body. From the results of the experiments in use of the proposed method, it is found that the method is effective to decrease the passed spam mails. Furthermore we confirm that the proposed method could replace the conventional spam spam filter function.