

SPAMIA: Spam filtering by quantitative profiles

Marián Grendár, Jana Škutová, Vladimír Špitalský

Slovanet a.s., Záhradnícka 151, 821 08 Bratislava, Slovakia
marian.grendar, jana.skutova, vladimir.spitalsky@slovanet.net

Applied Statistics 2012, International conference
September 23 - 26, 2012, Ribno (Bled), Slovenia

This presentation was prepared as a part of the "SPAMIA" project, MŠ SR 3709/2010-11, supported by the Ministry of Education, Science, Research and Sport of the Slovak Republic, under the heading of the state budget support for research and development.

Content

Spam and its detection

- Spam

- Traditional approach to spam filtering

Quantitative profile approach

- Quantitative profiles

Results

- Test corpuses

- Performance of quantitative profiles

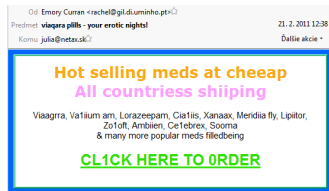
- Dimension of binary profiles

- Learning curves

Conclusions

Spam

- ▶ an unsolicited email message
- ▶ is usually send in a bulk to spread advert or viruses, or for phishing, scam, verification of email, . . .



Existing solutions for spam filtering

Methods

- ▶ heuristic rules
- ▶ naive Bayes filtering
- ▶ text-mining methods

Open-source products

SpamAssassin



Bogofilter



DSPAM



...

Comercial products

Disadvantages of existing solutions

- ▶ language dependence
- ▶ heuristic rules are fixed
- ▶ necessity to update these rules
- ▶ high vulnerability
- ▶ high computational costs

Quantitative profile approach

Spam and its detection

Spam

Traditional approach to spam filtering

Quantitative profile approach

Quantitative profiles

Results

Test corpuses

Performance of quantitative profiles

Dimension of binary profiles

Learning curves

Conclusions

Quantitative profile approach

- ▶ an email is represented by an m -dimensional vector of numbers with m fixed in advance
- ▶ QPs serve as an input to a classification algorithm

```
From: vba_test@postmaster.ru.uk [mailto:vba_test@postmaster.ru.uk]
Date: Thu, 17 Apr 2007 08:44:26 +0200
Return-Path: vba_test@postmaster.ru.uk
Received: from mxn1230 (vba-test.globeonline.ru [193.1.155.85])
  by openw1.mail.ru (Postfix) with SMTP; Thu, 17 Apr 2007 08:44:26 +0200
Received: from vba_test@postmaster.ru.uk [193.1.155.85]
  by openw1.mail.ru (Postfix) with SMTP; Thu, 17 Apr 2007 08:44:26 +0200
Received: from mail.globeonline.ru [193.1.155.85]
  by openw1.mail.ru (Postfix) with SMTP; Thu, 17 Apr 2007 08:44:26 +0200
From: "The Insider" <vba_test@postmaster.ru.uk>
Subject: "The Insider" - New Bulletin
Date: Thu, 17 Apr 2007 11:44:26 +0200
X-Mailbox: Archived by Mailman@vba-test.ru, 17 Apr 2007 11:44:26 +0200
Message-ID: <108012320402a9bde0102a2d6a4c01e@vba-test.ru>
X-Originating-IP: 17 Apr 2007 11:44:26 +0200
Status: 0
Content-Length: 134
Lines: 10

*** Bounce(s) removed ***

American games, resources, articles and stuff at American university
http://www.theinsider.org/news/article.asp?id=1476

-----
To be removed from this mailing list please use the form provided:
http://www.theinsider.org/news/emaila/theinsider/
```



$$QP = (qp_1, qp_2, \dots, qp_m)$$

Basic quantitative profiles

Binary profile: distances between occurrences of special character/characters (only first $k = 100$ occurrences for each email)

- ▶ **LP** line: lengths of lines
- ▶ **WP** word: lengths of words
- ▶ **BRP** brackets: distances between brackets
- ▶ ...

Histogram binary profile:

- ▶ **HWP**: histogram of lengths of words
- ▶ **HBRP**: histogram of distances between brackets
- ▶ ...

Basic quantitative profiles

Character profile: the number of occurrences of the characters

- ▶ **CP**: characters from \mathcal{A} (ASCII character set)

Grouped character profile: the number of occurrences of the groups of characters

- ▶ **CPG9**: numbers, spaces, brackets, operators, separators, upper/lower-case letters, forbidden characters, other
- ▶ **CPG11**: as CPG9, separately ! a \$

d -gram grouped character profile:

- ▶ **2CPG11**: pairs of groups of characters
- ▶ **3CPG11**: triples of groups of characters
- ▶ ...

Basic quantitative profiles

Moving window profile: CPGs for each parts of email

- ▶ **MWPCPG11**

Size profile:

- ▶ size of email
- ▶ sizes of selected headers
- ▶ sizes of parts of email according to content-type
- ▶ (optional) CPG of headers and parts
- ▶ **SP**
- ▶ **SPCPG11**

Graphical representation of line and character profile

```
From the_usaid@postmaster.us.us Thu Apr 17 16:44:24 2003
Return-Path: the_usaid@postmaster.us.us
Received: from ems1001 [192.1.153.47]
  by ap01b-postmaster.us.us [192.1.153.47] with Microsoft
  Exchange for the_usaid@postmaster.us.us; Thu, 17 Apr 2003 16:44:24 -0400
Received: from mail.yahoo.com by ems1001 with Microsoft Exchange;
  Thu, 17 Apr 2003 11:44:10 -0400
From: "The Insider" <the_usaid@postmaster.us.us>
To: "Obama-Bar" <obama@postmaster.us.us>
Subject: "The Insider" - News Bulletin
Date: Thu, 17 Apr 2003 11:44:10 -0400
X-Mailbox-Hashed: By Microsoft Word 10.0.2790.5509
Message-ID: <00000000000000000000000000000000>
X-OriginalArrivalTime: 17 Apr 2003 16:44:10.0714 (GMT)
Charset: 0
Content-Length: 534
Lines: 13

*** 00000000 0000 ***

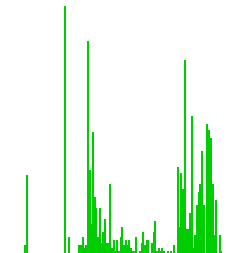
American public measures students and staff at American university
http://www.theinsider.org/news/bulletin.asp?id=2476

-----
To be removed from this mailing list please use the form provided:
http://www.theinsider.org/news/form/subscribe.asp/
```

(a) Email



(b) Line profile



(c) Character profile

Results

Spam and its detection

Spam

Traditional approach to spam filtering

Quantitative profile approach

Quantitative profiles

Results

Test corpuses

Performance of quantitative profiles

Dimension of binary profiles

Learning curves

Conclusions

Test corpuses

TREC 2007 corpus of 75 419 emails (spam 66.6%)

- ▶ train: 50 000 (68.3%)
- ▶ test: 25 419 (63.1%)

CEAS 2008 corpus of 137 705 emails (spam 80.3%)

- ▶ train: 90 000 (81.2%)
- ▶ test: 47 705 (77.9%)

Performance measures and classification algorithm

Performance measures

- ▶ **false negative rate fnr** (the ratio of misclassified spam) **at fixed low values of false positive rate fpr** (the ratio of misclassified ham)
- ▶ **the receiver operating characteristic (ROC) curve**, i.e. the graph of the true positive rate vs. the false positive rate, obtained as functions of the decision threshold

Classification algorithm

- ▶ **Random Forest** classifier

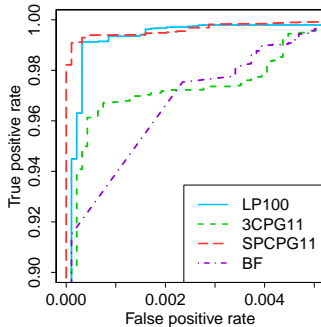
Performance of quantitative profiles

fnr (%) at fixed *fpr* = 0.1%

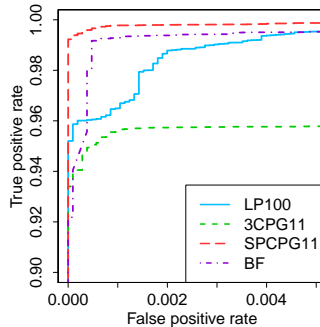
filter	TREC 2007	CEAS 2008
LP	0.65	3.46
WP	0.52	8.89
BRP	6.22	4.88
CP	14.61	4.98
3CPG11	3.26	4.42
MWPCPG11	17.26	5.45
SP	4.33	0.51
SPCPG11	0.60	0.22
SpamAssassin-RF	66.06	92.23
Bogofilter	7.98	0.71

Performance of quantitative profiles

ROC curves



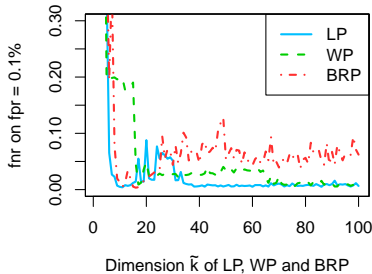
(a) TREC 2007



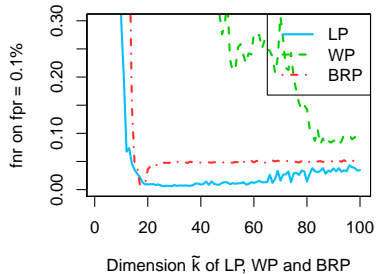
(b) CEAS 2008

Dimension of binary profiles

Dependence of BPs performance on its dimension



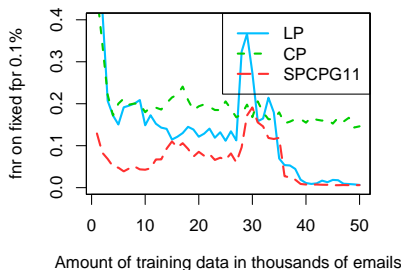
(a) TREC 2007



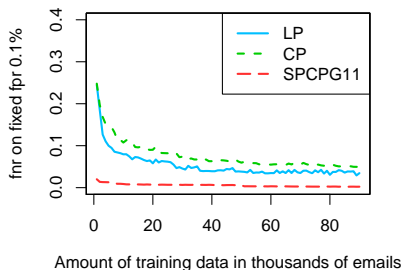
(b) CEAS 2008

Learning curves

Dependence of QPs performance on the size of training set



(a) TREC 2007



(b) CEAS 2008

Conclusions

Spam and its detection

Spam

Traditional approach to spam filtering

Quantitative profile approach

Quantitative profiles

Results

Test corpuses

Performance of quantitative profiles

Dimension of binary profiles






Learning curves

Conclusions

Conclusions

- ▶ quantitative profiles based Random Forest classifiers attain very good performance, at least comparable or better to that of Bogofilter and much better than optimized SpamAssassin
- ▶ the resulting filters are:
 - ▶ highly scalable
 - ▶ easy to parallelize (thanks to RF)
 - ▶ independent of language
 - ▶ easy to combine with other filters (thanks to QPs)

References

-  Breiman, L. (2001) *Random forests*. Machine Learning, 45, 5-32.
-  Grendár, M., Škutová, J. and Špitalský, V. (2011) *Spam filtering by quantitative profiles*. Appear in IJCSI Volume 9, Issue 5, September 2012. <http://arxiv.org/pdf/1201.0040v1.pdf>
-  Grendár, M., Škutová, J. and Špitalský, V. (2012) *Email categorization and spam filtering by random forest with new classes of quantitative profiles*. Proceedings of COMPSTAT 2012, 283–294.
-  R Development Core Team (2010) *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0. <http://www.R-project.org/>
-  Sroufe, P., Phithakkitnukoon, S., Dantu, R. and Cangussu, J. (2010) *Email shape analysis*. In *Distributed Computing and Networking*, Lecture Notes in Computer Science, K. Kant et al. (eds), 5935/2010, pp. 18-29.

Thank you for your attention.