

# Bayesian spam filtering

**Keith Briggs**

Keith.Briggs@bt.com

`more.btexact.com/people/briggsk2`

Pervasive ICT seminar 2004 Apr 29 1500

bayes-2004apr29.tex TYPESET 2004 MAY 6 10:48 IN PDFL<sup>A</sup>T<sub>E</sub>X ON A LINUX SYSTEM

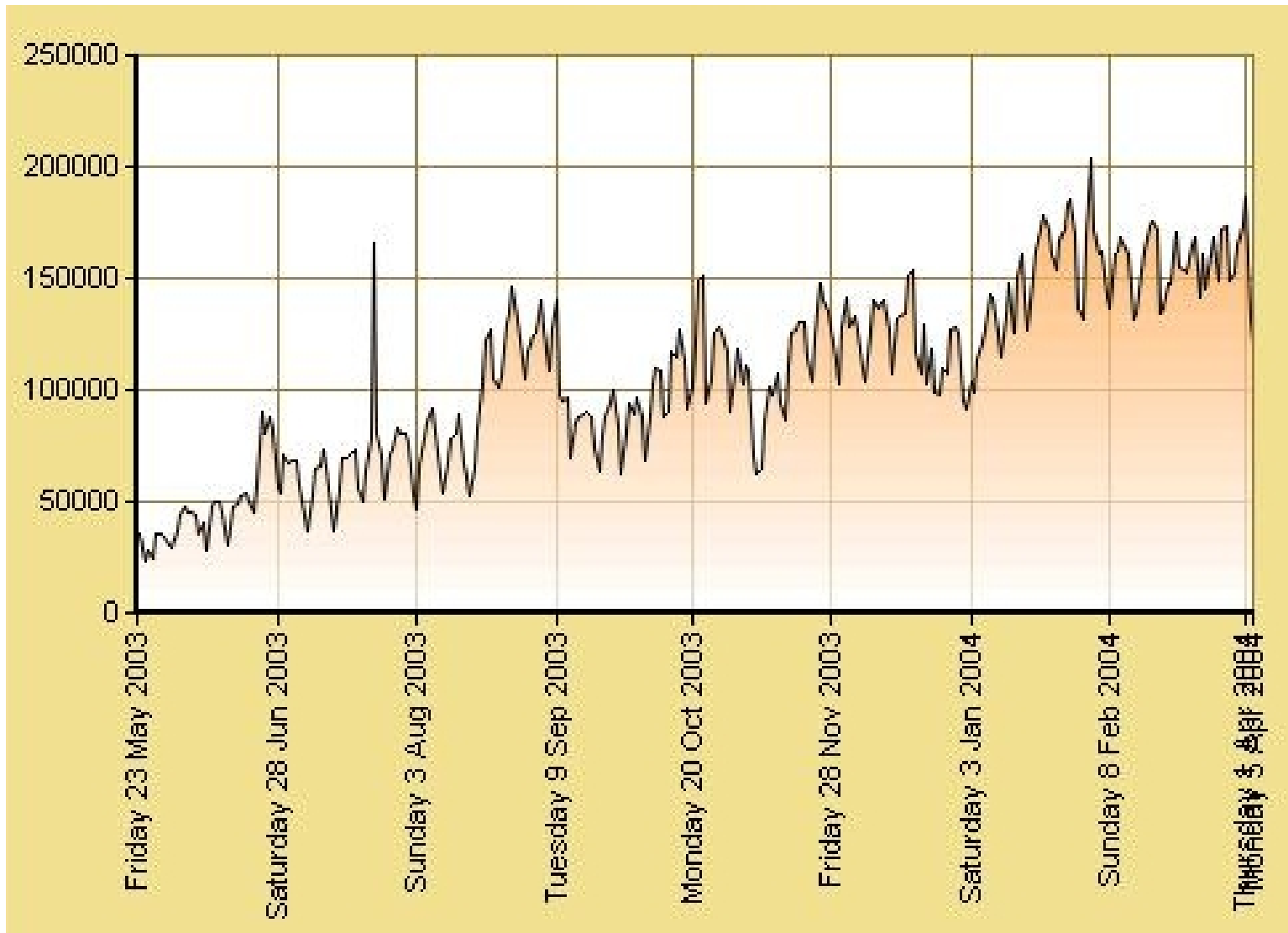
# Outline

- ★ the problem■
- ★ some 'solutions'■
- ★ probability theory■
- ★ Bayesian ideas■
- ★ . . . the real solution!■

The aim:

to understand why Bayesian methods work so well

# The problem



Total bt.com spam blocked per day

# RFC 706

Network Working Group  
Request for Comments: 706  
NIC #33861

Jon Postel (SRI-ARC)

## On the Junk Mail Problem

In the ARPA Network Host/IMP interface protocol there is no mechanism for the host to selectively refuse messages. This means that a Host which desires to receive some particular messages must read all messages addressed to it. Such a host could be sent many messages by a malfunctioning host. This would constitute a denial of service to the normal users of this host. Both the local users and the network communication could suffer. The services denied are the processor time consumed in examining the undesired messages and rejecting them, and the loss of network thruput or increased delay due to the unnecessary busyness of the network.

Request for Comments: 706

Nov 1975

# Solutions?

Method	version	pro	con
legal legal legal	CAN-SPAN opt-in opt-out		just politics you're kidding! "
blacklist whitelist	spamhaus etc.		centralized new emailers must register
financial penalty computational penalty	micropayment		who administers? new protocol
tripoli [5]		?	new protocol
challenge-response	e.g. about.mailblocks.com		lists attacks
filter filter	keyword probabilistic	it works!	it doesn't work
honeypot	beat them at their own game		admin load
commercial solutions	?	?	?



# Probability vs. degree of belief

★  $P(\text{event}) \equiv \lim_{n \rightarrow \infty} \frac{\#\text{events}}{n}$

- ▷ *objective*
- ▷ *must be able to repeat the experiment indefinitely*
- ▷ *rate of convergence of limit unspecified*
- ▷ *strictly speaking, this rules out using this definition in the real world*

■

★ 'degree of belief'  $B$  is more or less subjective

- ▷ *meaningful for a single, non-repeatable event*
- ▷ *your  $B$  might not be the same as my  $B$*
- ▷ *chance of rain tomorrow*
- ▷ *chance of horse winning a race*
- ▷ *spamminess of an email*

■

★ Cox's axioms are a set of reasonable assumptions under which there exists a function mapping beliefs to probabilities

# Probability theory

## ★ conditional probability

$$P(x = a|y = b) \equiv \frac{P(x = a, y = b)}{P(y = b)} \blacksquare$$

## ★ product rule

$$P(x, y|\mathcal{H}) = P(x|y, \mathcal{H})P(y|\mathcal{H}) = P(y|x, \mathcal{H})P(x|\mathcal{H}) \blacksquare$$

## ★ marginalization

$$\begin{aligned} P(x|\mathcal{H}) &= \sum_y P(x, y|\mathcal{H}) \\ &= \sum_y P(x|y, \mathcal{H})P(y|\mathcal{H}) \end{aligned}$$



# Bayes' theorem

★ Bayes' theorem - is just the product rule:

$$P(y|x, \mathcal{H}) = \frac{P(x|y, \mathcal{H})P(y|\mathcal{H})}{P(x|\mathcal{H})}$$



★ . . . with  $y$  interpreted as the data  $D$ , and  $x$  interpreted as parameters  $\theta$ :

$$P(\theta|D, \mathcal{H}) = \frac{P(D|\theta, \mathcal{H})P(\theta|\mathcal{H})}{P(D|\mathcal{H})}$$

# Prior, likelihood and posterior

- ★ we can think of Bayes' rule as:

$$\text{posterior} \propto \text{likelihood} * \text{prior}$$

- ★ for example, a single bit  $s$  sent twice over a noisy channel, received as  $r_1 r_2$ :

- ▷  $P(s = 1 | r_1 r_2) = P(r_1 r_2 | s = 1) P(s = 1) / P(r_1 r_2)$

- ▷  $P(s = 0 | r_1 r_2) = P(r_1 r_2 | s = 0) P(s = 0) / P(r_1 r_2)$

- ★ that is, your **prior** (degree of belief before you observed that data  $r$ ), is **updated** by the information the data provides about the value of  $s$  (the likelihood), to provide your **posterior** degree of belief

## Thomas Bayes (1702-1761)

*I now send you an essay which I have found among the papers of our deceased friend Mr Bayes, and which, in my opinion, has great merit. . . In an introduction which he has writ to this Essay, he says, that his design at first in thinking on the subject of it was, to find out a method by which we might judge concerning the probability that an event has to happen, in given circumstances, upon supposition that we know nothing concerning it but that, under the same circumstances, it has happened a certain number of times, and failed a certain other number of times. — Richard Price.*



## € biased?

★ **Guardian** Friday 2002 January 04:

When spun on edge 250 times, a Belgian €1 coin came up heads 140 times and tails 110. *'It looks very suspicious to me'*, said Barry Blight, a statistics lecturer at the LSE *'If the coin were unbiased the chance of getting a result as extreme as that would be less than 7%'*

★ we compare the models  $\mathcal{H}_0$  (the coin is fair) and  $\mathcal{H}_1$  (the coin is biased), with uniform prior  $P(p|\mathcal{H}_1) = 1$  ■

★ the likelihood ratio is:

$$\frac{P(D|\mathcal{H}_1)}{P(D|\mathcal{H}_0)} = \frac{140!110!}{251!} \approx 0.48$$

★ thus the data give weak evidence (2.08/1) in favour of  $\mathcal{H}_0$ !

# Text classification theory

- ★ could be based on various choices of *features*: words, or  $n$ -grams ■
- ★ corpora  $C_1, C_2, \dots, C_k$  ■
- ★ priors  $\pi_1, \pi_2, \dots, \pi_k$  ■
- ★ models  $\mathbb{P}_{C_1}, \mathbb{P}_{C_2}, \dots, \mathbb{P}_{C_k}$  ■
- ★ if  $x$  is an unknown document, the posterior probability that  $x$  belongs to  $C_j$  is  $P(C_j|x) \propto \mathbb{P}_{C_j} \pi_j$  ■
- ★ decision rule: choose  $j$  to maximize  $P(C_j|x)$  ■
- ★ other uses
  - ▷ *sorting emails - work or personal*
  - ▷ *forwarding emails in French or German etc. to the right person*

# Digram measure

- ★ word  $w = w_1w_2 \dots w_k$  ■
- ★ reference measure  $R_C(w) \equiv p_C(\wedge, w_1)p_C(w_1, w_2) \dots p_C(w_k, \$)$  ■
  - ▷ *this is naïve - it assumes adjacent digrams are statistically independent*
- 
- ★ Dirichlet digram measure  $p_C(u, v) = \frac{\#(v|u)}{\sum_r \#(r|u)} \frac{+ \alpha \mu(v)}{+ \alpha}$  ■
- ★  $\alpha$  is a hyperparameter, and the optimum  $\alpha$  should be chosen from tests on various corpora
- ★ for spam we use two corpora, spam and ham
- ★ the latest softwares claim 99.9% accuracy
- ★ false positives are much worse than false negatives, so we can adjust our algorithm to account for this

# Example digram measure

	ax	ey1	t	rih1	p	laa1	b	er	g	k	n	s	eh	th	mao1	aa	d	eh1	v	iy1	ae1	ow1	iy	ow	ey	uw	z	ih	ng	sh	ae	ah1	ao	f	hh	ay				
^	89	56	89	78	89	78	89	89	44	78	89	78	89	56	67	89	89	22	78	78	78	67	100	67	0	11	11	0	22	78	0	56	56	100	56	89	89	56		
ax	0	0	78	0	0	67	89	0	67	0	56	67	100	78	0	33	67	0	0	78	0	56	0	0	0	0	0	0	0	0	0	0	0	0	0	67	22	0		
ey1	0	0	67	0	0	44	56	0	44	0	0	56	67	56	0	0	67	0	0	44	0	22	0	0	0	0	0	0	56	0	0	67	0	0	0	33	0	0		
t	78	56	0	67	67	0	67	22	0	67	0	0	22	67	0	0	0	33	0	0	56	0	22	33	22	67	11	0	89	0	67	0	0	0	33	11	0	0	56	
r	78	67	56	0	56	11	56	56	22	0	56	56	33	33	11	22	56	44	0	44	67	33	67	56	33	78	33	33	67	22	56	0	22	33	78	0	33	11	56	
ih1	0	0	78	67	0	33	67	0	33	0	33	56	78	78	0	44	56	0	0	44	0	56	0	0	0	0	0	0	78	0	56	56	0	0	0	67	0	0		
p	67	22	56	67	44	0	56	56	0	56	0	0	0	56	11	22	44	33	33	0	67	0	78	44	44	11	0	0	0	0	33	0	33	0	22	0	0	0	44	
l	67	67	56	44	56	22	0	56	0	56	0	22	0	56	44	0	56	56	0	56	67	33	44	56	56	78	0	0	56	56	44	0	0	0	22	11	33	0	78	
aa1	0	0	67	78	0	44	44	0	44	0	33	44	78	56	0	11	56	0	0	56	0	0	0	0	0	0	0	0	56	0	0	0	0	0	0	0	33	0	0	
b	78	33	11	56	78	0	44	44	0	33	0	0	22	44	0	0	0	33	11	33	56	0	78	44	56	0	0	0	0	33	33	0	0	0	0	67	0	0	0	78
er	56	33	44	0	44	44	44	11	11	0	33	11	33	33	0	11	33	0	0	44	44	44	33	33	0	44	0	11	56	33	0	0	0	0	0	0	33	44	56	
g	56	33	0	67	56	0	33	11	0	33	0	0	22	0	0	0	11	22	0	0	44	0	22	56	44	0	0	22	22	67	0	0	0	0	0	0	0	0	22	
k	67	56	67	44	22	0	67	56	0	22	0	0	11	78	0	0	11	33	22	0	56	0	33	67	33	0	11	0	33	0	44	0	44	0	56	11	0	0	44	
n	67	56	78	11	56	11	67	67	0	44	11	56	33	67	0	11	22	44	0	89	56	22	33	44	67	78	22	0	56	56	56	0	33	0	56	22	56	33	44	
s	78	67	89	0	67	78	67	22	22	44	0	56	22	22	22	22	22	33	0	0	67	0	67	33	67	33	44	0	44	0	44	0	0	0	78	0	44	0	56	
eh	0	0	0	44	0	11	0	0	0	0	33	44	44	22	0	0	33	0	0	11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
th	22	0	0	56	44	0	0	22	0	0	0	0	0	22	0	0	0	22	0	0	11	0	11	0	0	0	0	11	0	0	0	0	0	0	0	0	0	0	0	
m	67	67	44	0	56	67	44	56	56	56	0	0	0	33	22	0	0	67	22	33	67	0	44	56	67	11	0	33	44	56	33	0	0	11	67	11	11	11	67	
ao1	0	0	44	89	0	0	78	0	0	0	22	0	33	44	0	0	11	0	0	11	0	0	0	0	0	0	0	0	44	0	56	0	0	0	0	0	56	0	0	
aa	0	0	0	33	0	0	11	0	0	0	0	11	22	11	0	0	11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
d	67	44	0	44	56	0	44	33	0	44	0	0	11	0	0	11	44	22	0	0	56	44	44	22	22	56	0	0	67	67	56	0	0	22	44	0	0	22	44	
eh1	0	0	56	78	0	56	67	0	22	0	56	67	78	67	0	11	56	0	0	56	0	67	0	0	0	0	0	0	0	33	0	11	44	0	0	0	44	0	0	
v	56	44	0	56	56	0	33	22	0	67	0	0	0	0	33	0	33	11	0	56	67	0	22	56	0	33	0	0	44	33	0	0	11	0	0	0	0	0	44	
iy1	44	0	44	0	0	44	44	0	11	0	0	44	67	44	0	33	56	0	0	56	0	56	0	56	0	0	0	0	0	67	56	0	67	0	0	0	44	0	0	
ae1	0	0	89	44	0	44	44	0	44	0	22	56	89	56	0	11	56	0	0	56	0	78	0	0	0	0	0	0	78	0	44	22	0	0	0	44	0	0		
ow1	33	0	33	0	0	33	56	0	11	0	0	22	67	67	0	56	33	0	0	22	0	44	0	0	0	0	0	0	67	33	0	33	0	0	0	0	0	0	0	
iy	67	33	11	0	0	11	44	33	0	33	0	33	0	22	11	22	33	22	0	33	33	33	0	33	44	0	0	0	78	33	0	0	0	0	0	0	0	11	0	
ow	0	0	22	0	11	0	11	0	11	0	0	11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

## Example using dbacl [1]

```
# learn...
```

```
kbriggs:~/Bayes> dbacl -l ham ham/*
```

```
kbriggs:~/Bayes> dbacl -l spam spam/*
```

```
# test...
```

```
kbriggs:~/Bayes> echo "meet for lunch?" | dbacl -v -c ham -c spam  
ham
```

```
kbriggs:~/Bayes> echo "viagra" | dbacl -v -c ham -c spam  
spam
```

```
kbriggs:~/Bayes> echo "mortgage" | dbacl -v -c ham -c spam  
spam
```



## Can they beat us?

- ★ the old trick of `deliber@te` spel1ing errors and füñny diàçrítics now works in our favour!■
- ★ dictionary salad: lots of random words
  - fails because the spammer doesn't know our model■
- ★ it's a long story: some genuine text - a possible danger■
- ★ habeas haiku: copyright poem, attempt at legal protection - now a strong spam indicator■

## Honesty in signalling - some philosophy

- ★ spammers cannot violate the rules of any *machine* protocol like SMTP. . . ■
- ★ but they violate the rules of human netiquette all the time. . . ■
- ★ so perhaps we have to write some software which emulates human behaviour to detect the breach■
- ★ normal human communication has evolved mechanisms to detect sincerity: tone of voice, body language. . . ■
- ★ and in particular, complex rules of grammar which although they impose a cost on the speaker, the consequence is that the listener *knows* the speaker is making an effort, and is thus worth listening to■
  - ▷ *does this imply that any solution to the spam problem (or any other network security problem) must involve some similar mechanism?*
- ★ in any case, I now understand why correct spelling and grammar are so important in written documents!

## References

- [1] L Breyer *dbacl - a digramic Bayesian classifier*  
<http://dbacl.sourceforge.net/>
- [2] W Yerazunis *The spam-filtering accuracy plateau at 99.9% and how to get past it*  
[http://crm114.sourceforge.net/Plateau\\_Paper.pdf](http://crm114.sourceforge.net/Plateau_Paper.pdf)
- [3] D J C MacKay & L C Peto *A hierarchical Dirichlet language model*  
<http://www.inference.phy.cam.ac.uk/mackay/BayesLang.html>
- [4] P Gburzynski & J Maitan *Fighting the spam wars: a remailer approach with restrictive aliasing* ACM Trans. Internet Tech. 4 1-30 (2004)
- [5] Tripoli <http://www.pfir.org/tripoli-overview>
- [6] J D M Rennie & T Jaakkola *Automatic feature induction for text classification* MIT AI Lab, September 2002