

Predicting the Future

Finite State Machines

Testing, Probability, Statistics

and other unpleasant things

Gordon V. Cormack

University of
Waterloo



Data Compression

50'000€ Prize for Compressing Human Knowledge

(widely known as the Hutter Prize)

Compress the 100MB file `enwik8` to less than the current record of about 16MB

- ♦ [The Task](#)
- ♦ [Motivation](#)
- ♦ [Detailed Rules for Participation](#)
- ♦ [Previous Records](#)
- ♦ [More Information](#)
- ♦ [Newsgroup on the contest and prize](#)
- ♦ [History](#)
- ♦ [Committee](#)
- ♦ [Donations](#)
- ♦ [Frequently Asked Questions](#)
- ♦ [Contestants](#)
- ♦ [Links](#)
- ♦ [Disclaimer](#)

*News: Alexander Ratushnyak is also the **second Winner!** Congratulations!*



... the contest continues ...



Being able to compress well is closely related to intelligence as explained below. While intelligence is a slippery concept, file sizes are hard numbers. Wikipedia is an extensive snapshot of Human Knowledge. If you can compress the first 100MB of Wikipedia better than your predecessors, your (de)compressor likely has to be smart(er). The intention of this prize is to encourage development of intelligent compressors/programs.

The Task

Create a compressed version (self-extracting archive) of the 100MB file `enwik8` of less than about 16MB. More precisely:

DMC

Information theory

Automata theory

Markov processes

Probability & stats

Machine learning

250 lines of code

Evaluation & measurement

Application (avionic telemetry)

theory + practice + evaluation + application



Predict Human Actions

NETFLIX

Netflix Prize

Home Rules Leaderboard Register Update Submit Download

NETFLIX

Browse Recommendations Friends Queue Buy DVDs

Home Genres New Releases Previews Netflix Top 100 Crit

Movies For You

Randy, the following movies were chosen based on your interest in:
[Bowling for Columbine](#)
[Carnivale: Season 1](#)
[Fahrenheit 9/11](#)

The Big One

★★★★★

More subversive
from
Michael

Carnivale: Season 2

★★★★★

Disc Series

Daniel Knaul
rivetingly cre
series conti
document t
entures of a motley cre

You really liked it...

Now own for just \$5.99

Shop as low

titles

Original artv

Lewis Black: Re
and Scre

Add

Welcome!

The Netflix Prize seeks to substantially improve the accuracy of predictions about how much someone is going to love a movie based on their movie preferences. Improve it enough and you win one (or more) Prizes. Winning the Netflix Prize improves our ability to connect people to the movies they love.

Read the [Rules](#) to see what is required to win the Prizes. If you are interested in joining the quest, you should [register a team](#).

You should also read the [frequently-asked questions](#) about the Prize. And check out how various teams are doing on the [Leaderboard](#).

Good luck and thanks for helping!

Guides:

Your longest run of agreement with the filters is 8 (for this session); your current run is 8 (keep playing!).

Scores for all your sessions will be used in determining the winner; just use the same email address each time you visit spamorham.org

Thanks gvcormac@uwaterloo.ca! On the last email you said **ham** and the spam filters said **ham**
So far you've found **0 possible filter errors** and verified **8 emails**

Click one of these buttons:

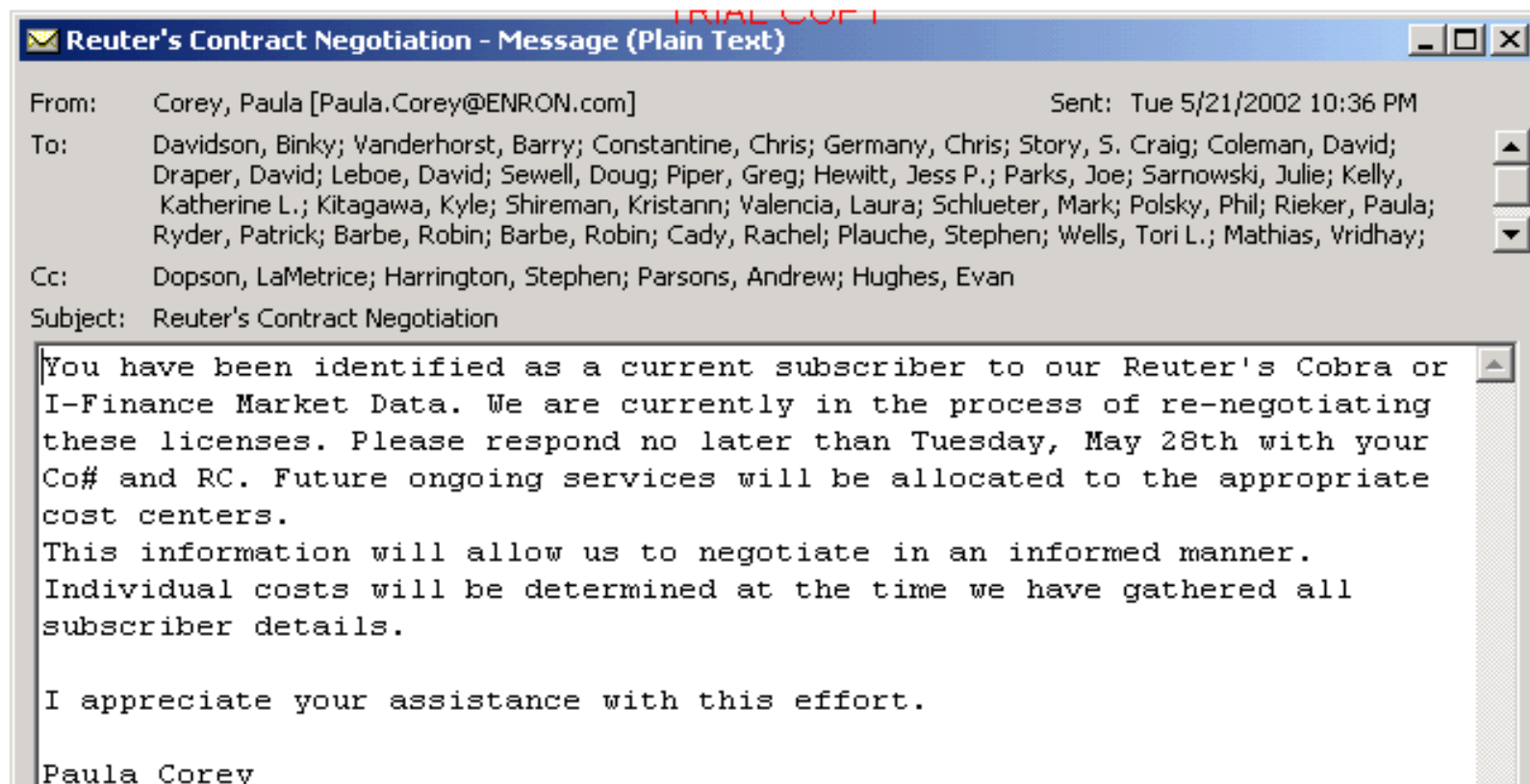
This is Spam

- I'm not sure

- This is Ham

☐ Flag message as funny

As displayed by Microsoft Outlook



Predictive models work for

data compression

spam detection

viruses, phishing, IM, SMS, blog, Web spam

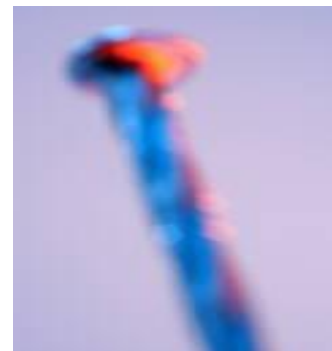
insensitive to language, alphabet, coding method

heterogeneous, multimedia, metadata

plagiarism detection, authorship attribution

intrusion detection

game playing



Need well defined tasks and evaluation!

Given a stream of bits

Represent the stream in fewer bits

Trick:

- predict each bit in turn (as a probability p)

- encode as $-\log_2(p)$ bits (on average)

 - arithmetic coding

 - optimal given p

Measure success! Compress some data!

What is Spam?

Unsolicited, unwanted email that was sent indiscriminately, directly or indirectly, by a sender having no current relationship with the recipient.

Depends on sender/receiver relationship

Not “whatever the user thinks is spam.”

Spam and non-spam examples

Spam

```
Hi,  
=20  
M e R / D / A  
V / a G R A  
P R O z & C  
A m o x / c i I l / n  
C i A L / S  
V A L / u M  
T r & m a d o I  
A m B / E N  
X & n a x  
L e V / T R A  
S O m &  
=20  
http://www.prosebutis.com <http://www.prosebutis.com>=20
```

Non-spam

Dear Gord:

Your C program has solved Ok the problem 11102 (Moonshine)
in 0.514 seconds using as much as 420 kbytes of virtual memory.
Congratulations!

--

PS: Check the board at <http://acm.uva.es/board/>

The Online Judge (Linux acm.uva.es 2.4.18-27.7.x i686)
Judge software version 2.8 [<http://acm.uva.es/problemset/>]
Wed May 24 23:19:30 UTC 2006



Objective: color **spam** red, **non-spam** green

Spam

```
Hi,  
=20  
M e R / D / A  
V / a G R A  
P R O z & C  
A m o x / c i I l / n  
C i A L / S  
V A L / u M  
T r & m a d o I  
A m B / E N  
X & n a x  
L e V / T R A  
S o m &  
=20  
http://www.prosebutis.com <http://www.prosebutis.com>=20
```

Dear Gord:

Your C program has solved Ok the problem 11102 (Moonshine)
in 0.514 seconds using as much as 420 kbytes of virtual memory.
Congratulations!

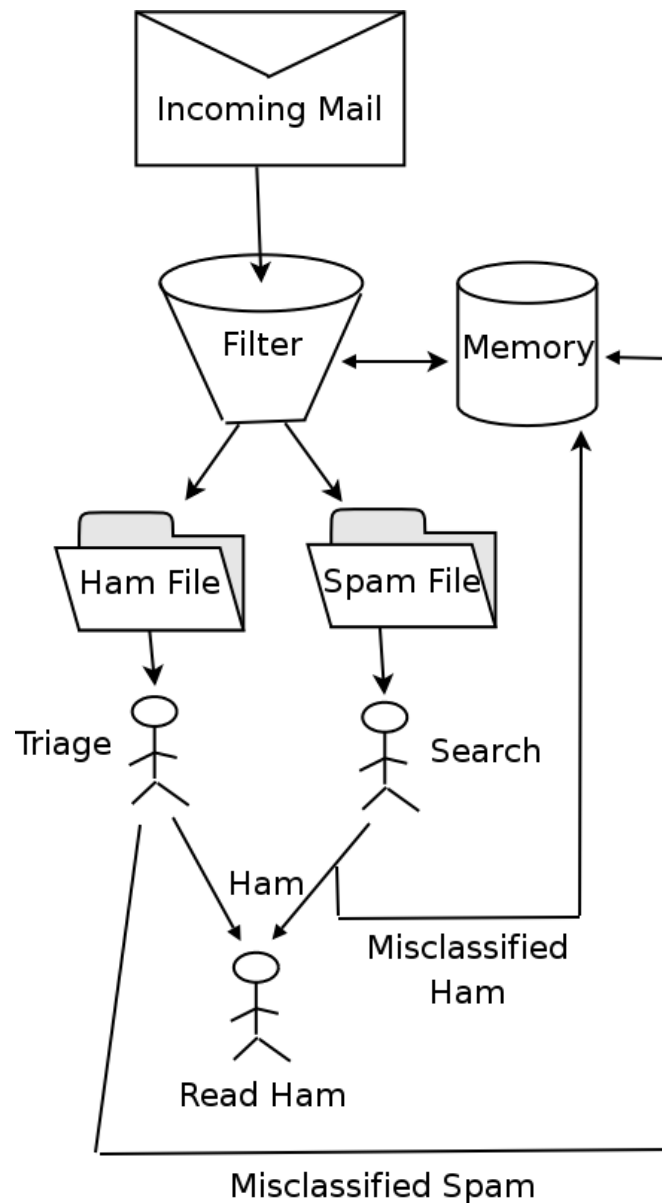
--

PS: Check the board at <http://acm.uva.es/board/>

The Online Judge (Linux acm.uva.es 2.4.18-27.7.x i686)
Judge software version 2.8 [<http://acm.uva.es/problemset/>]
Wed May 24 23:19:30 UTC 2006

Non-spam
(ham)

How is the coloring used?



Filter Classifies Email

Human addressee

Triage on ham File

Reads ham

Occasionally searches
for misclassified ham

Report misclassified
email to filter

Questions to answer

Method to color **spam** & **non-spam (ham)**?

How well does the method color?

How well is the overall purpose met?

Facilitating delivery of good email

“filtering spam” is just a means to the end

Models and prediction

Given a sequence of bits, predict the next one (x)

1011011011011011011 x

x is *probably* **0**

0101101110111101111 x

x is *probably* **1**

How '*probably*'?

Prob($x = \mathbf{0}$ following 1011011011011011011)

Prob($x = \mathbf{1}$ following 0101101110111101111)

Model

abstracts the string of bits; used to predict *behavior*

0th order Markov model

Count the number of *zeros* & the number of *ones*:

1011011011011011011 \mathbf{x}

zeros: **6** *ones*: **13**

Use the proportion of *ones* to estimate

$$\text{Prob}(\mathbf{x} = \mathbf{1}) = 13/19 = 0.68$$

Doesn't seem like such a good estimate

how can we validate it?

intuition

testimonial

faith

experiment

1st order Markov model

Count the number of *ones* and *zeros* following a 0,
and the number following a 1

1011011011011011011 x

following 0: *zeros*: **0** *ones*: **6**

following 1: *zeros*: **6** *ones*: **6**

Use the proportion of *ones* following 1 to estimate

$$\text{Prob}(x = 1 \text{ following } 1) = 6/12 = 0.5$$

Still doesn't seem like such a good estimate

but better than 0th order

2nd order Markov model

Count the number of *ones* and *zeros* following 00,
and 01, and 10, and 11.

1011011011011011011 \mathbf{x}

following 00: *zeros*: **0** *ones*: **0**

following 01: *zeros*: **0** *ones*: **6**

following 10: *zeros*: **0** *ones*: **6**

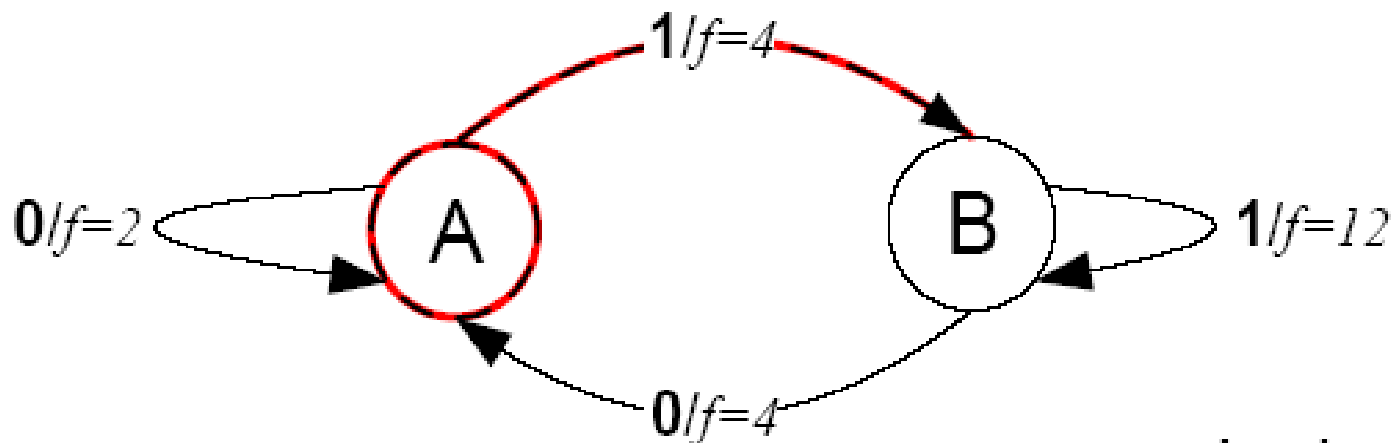
following 11: *zeros*: **5** *ones*: **0**

Use the proportion of *ones* following 11 to
estimate

$$\text{Prob}(\mathbf{x} = 1 \text{ following } 11) = 0/5 = 0$$

Overconfident! (Overfitted model)

Dynamic Markov model (DMC)



This example implements a 1st order Markov model

A means *following 0*; **B** means *following 1*

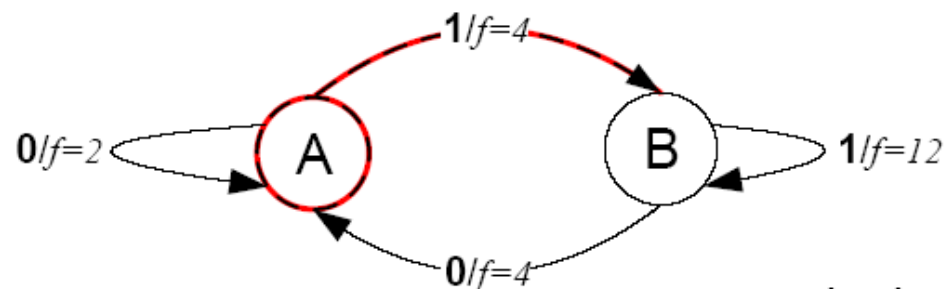
Outputs f on edges are frequencies

$\text{Prob}(1 \text{ following } \mathbf{A}) = 4 / (2 + 4) = 0.667$

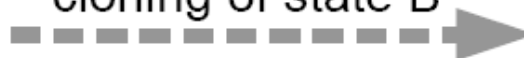
f incremented after each transition

DMC State Cloning

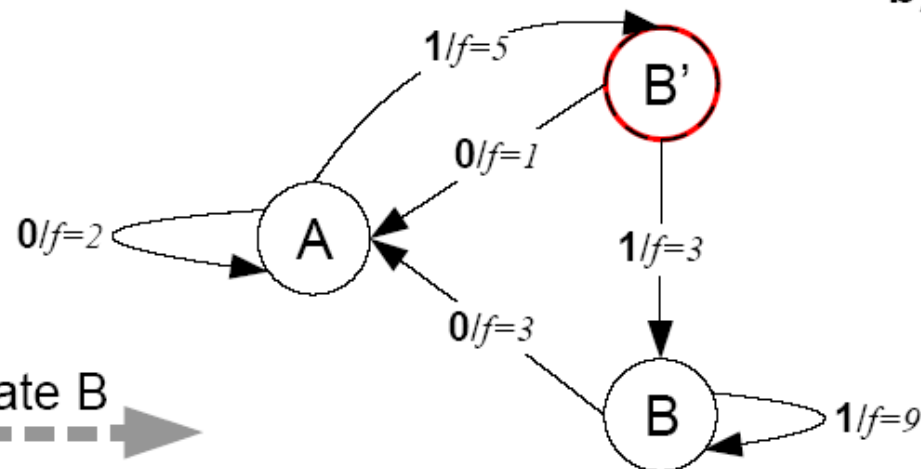
a)



cloning of state B



b)



State A, input 1, Prob 0.67

B cloned to create B'

B visited 16 times previously

f divided in 4:12 ratio in proportion to previous visits

4 from A; 12 from elsewhere

B should be cloned because it is visited from distinct contexts several times

f incremented as usual

Predict each bit in turn

DMC

Construct optimal code

arithmetic coding

the more probable the shorter the representation

$-\log_2 \text{prob}$

but how do you do a fraction of a bit?

many bits at a time

Google for dmc.c

Likelihood Ratio

Likelihood of a bit (say 0) in **spam**

1011011011011011011**0** **Prob**($\mathbf{x} = \mathbf{0}$)

Likelihood of same bit in **non-spam**

0101101110111101111**0** **Prob**($\mathbf{x} = \mathbf{0}$)

log-likelihood ratio

$$\textit{spamminess} = \log(\textbf{Prob}(\mathbf{x} = \mathbf{0}) / \textbf{Prob}(\mathbf{x} = \mathbf{0}))$$

coloring method

spam if $\textit{spamminess} > 0$; otherwise **non-spam**

more generally

spam if $\textit{spamminess} > \mathbf{t}$; otherwise **non-spam**

Combining Likelihoods

$$\begin{aligned} & \textit{spamminess}(\mathbf{x}_1 \mathbf{x}_2 \mathbf{x}_3 \dots \mathbf{x}_n) \\ &= \log(\text{Prob}(\mathbf{x}_1 \mathbf{x}_2 \mathbf{x}_3 \dots \mathbf{x}_n) / \text{Prob}(\mathbf{x}_1 \mathbf{x}_2 \mathbf{x}_3 \dots \mathbf{x}_n)) \\ &= \log(\text{Prob}(\mathbf{x}_1) / \text{Prob}(\mathbf{x}_1)) + \\ & \quad \log(\text{Prob}(\mathbf{x}_2) / \text{Prob}(\mathbf{x}_2)) + \\ & \quad \log(\text{Prob}(\mathbf{x}_3) / \text{Prob}(\mathbf{x}_3)) + \dots \\ & \quad \dots + \log(\text{Prob}(\mathbf{x}_n) / \text{Prob}(\mathbf{x}_n)) \end{aligned}$$

Email spamminess

Let

S be a string consisting of all known spam

N be a string consisting of all known non-spam

E be an email message

Define spamminess

$$\log (\text{Prob}(\text{E following } \text{S}) / \text{Prob}(\text{E following } \text{N}))$$

Measuring success

Collect email stream

adjudicate as **spam** or **ham**

gold standard

Filter email to

spam file if *spamminess* $> t$

ham file otherwise

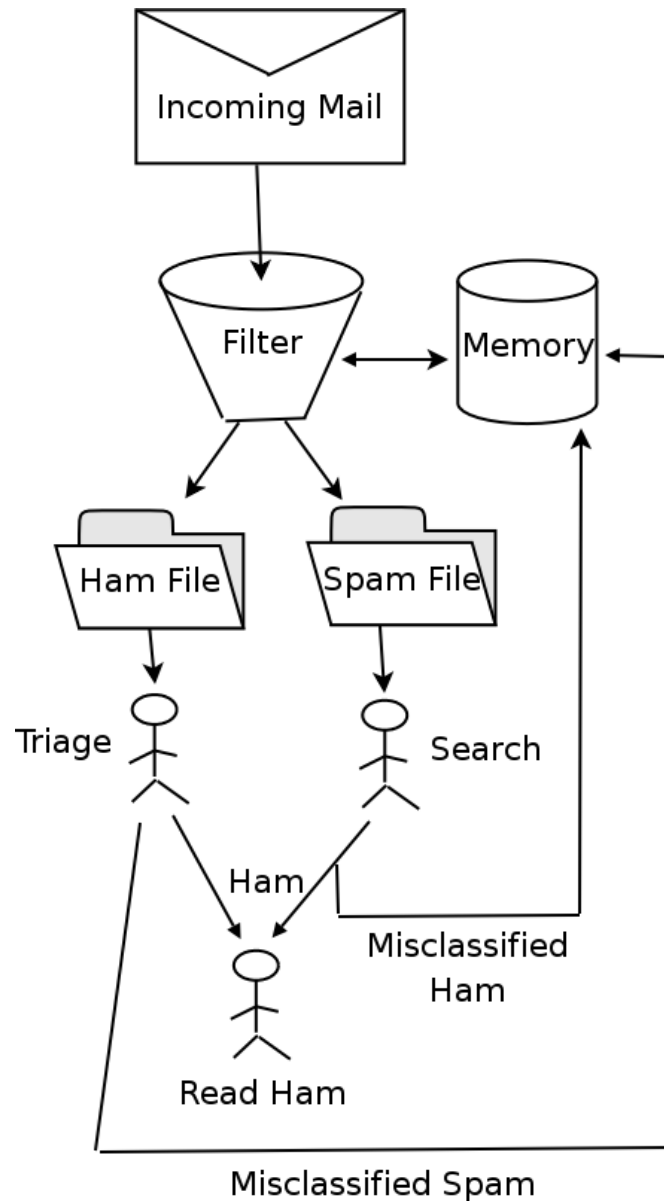
Idealized user

reports errors immediately

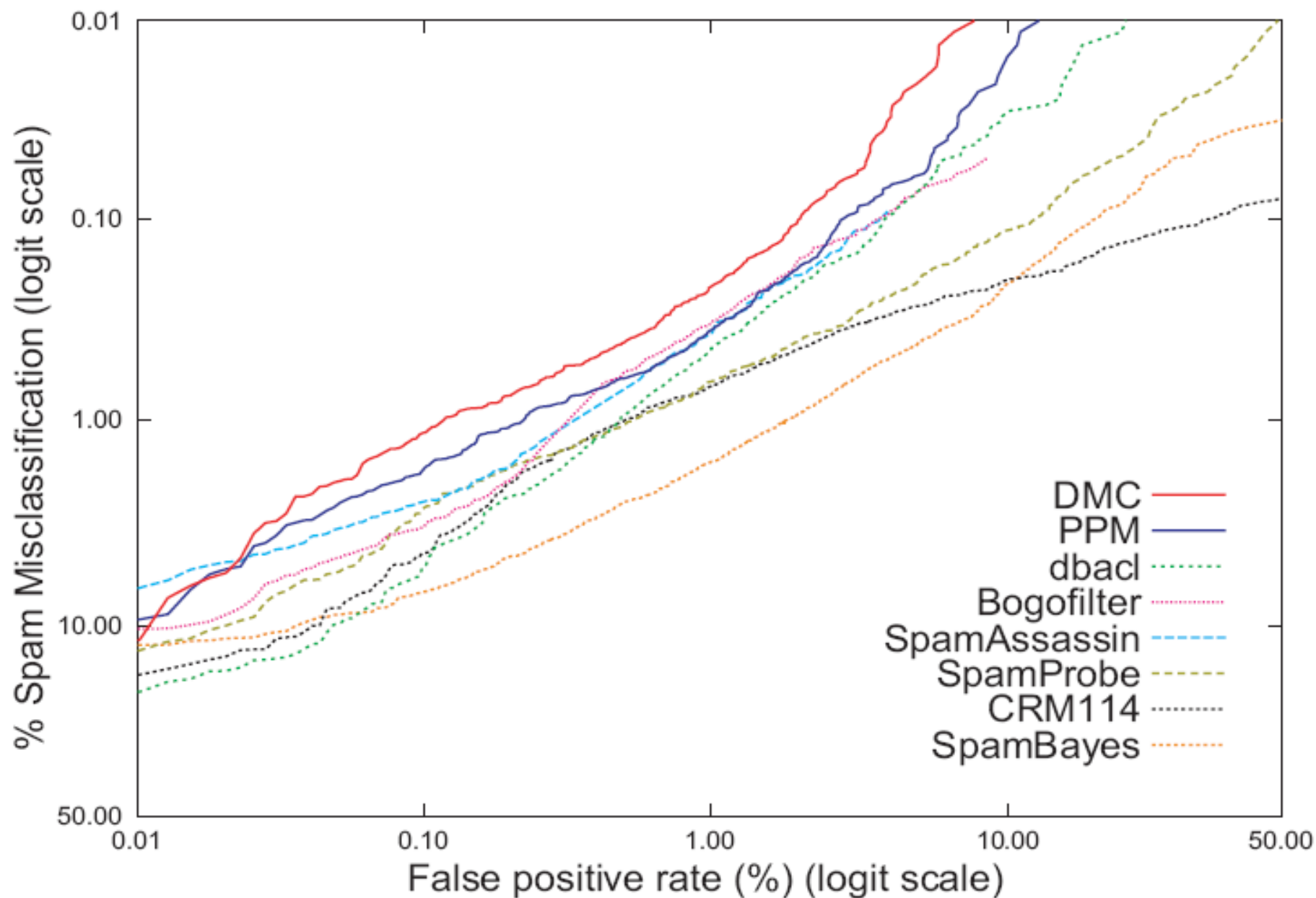
Measure

false positive rate

spam misclassification rate



Receiver Operating Characteristic Curve



Summary statistics

TREC public corpus

<i>Filter</i>	1-AUC (%)	SMR at 1% FP	SMR at 0.1% FP	SMR at 0.01% FP
DMC	[†] 0.013 (0.010 – 0.018)	0.22%	1.17%	14.47%
PPM	[†] 0.019 (0.015 – 0.023)	0.36%	1.78%	9.89%
dbacl ^b	0.037 (0.031 – 0.045)	0.45%	5.19%	19.77%
Bogofilter ^b	0.048 (0.038 – 0.062)	0.33%	3.41%	10.39%
SpamAssassin ^b	0.059 (0.044 – 0.081)	0.37%	2.56%	7.81%
SpamProbe	0.059 (0.049 – 0.071)	0.65%	2.77%	15.30%
CRM114 ^b	0.122 (0.102 – 0.145)	0.68%	4.52%	17.17%
SpamBayes ^b	0.164 (0.142 – 0.189)	1.63%	6.92%	12.55%

[†] improves on best TREC result ($p < .05$)

^b TREC 2005 result

Sponsored by, held at

NIST – National Institute for Standards & Technology

<http://trec.nist.gov>

Goals

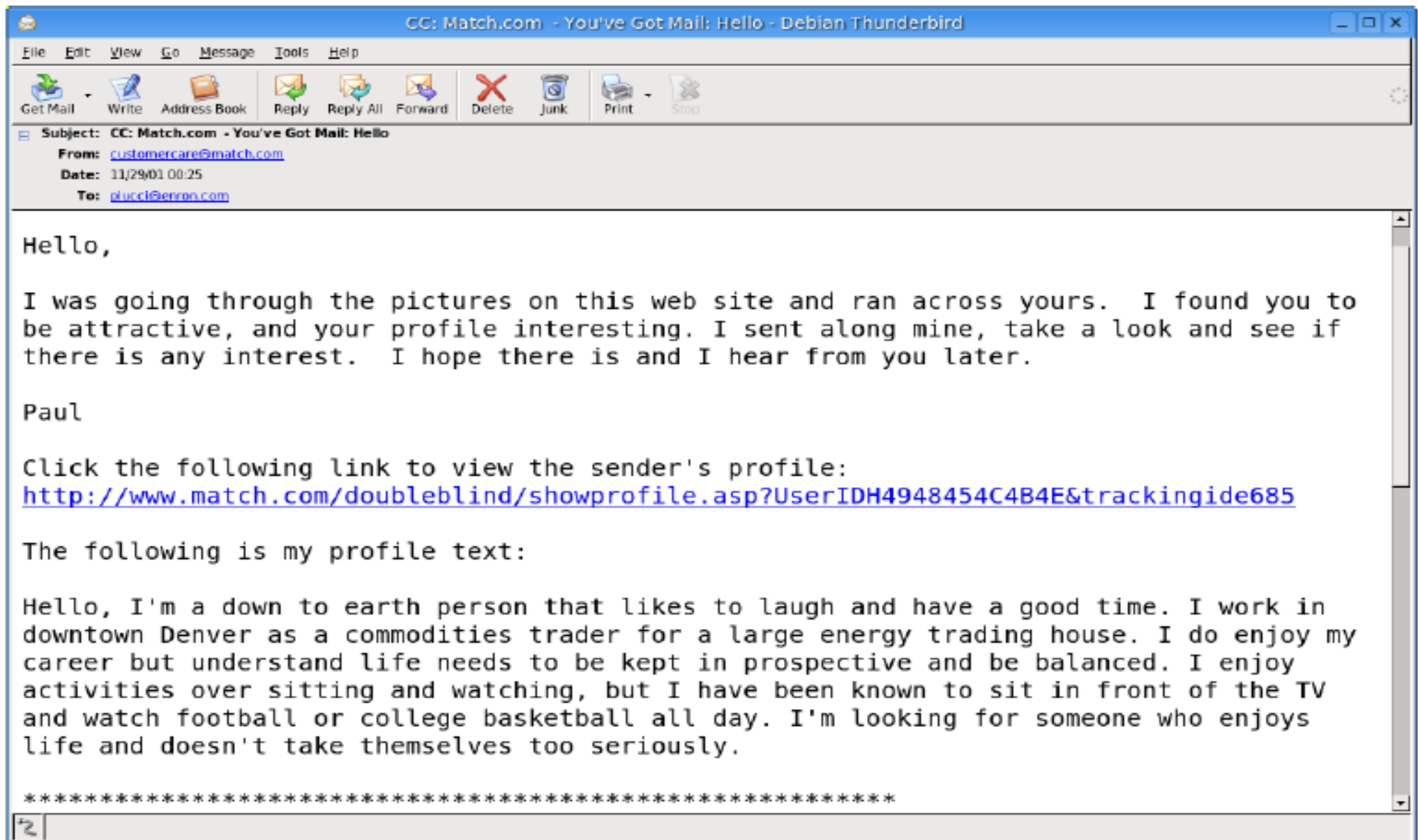
To increase the availability of appropriate evaluation techniques for use by industry and academia, including the deployment of new evaluation techniques more applicable to current systems.

Format

Participants do experiments in one or more *tracks*

Standardized evaluation of well-defined tasks

Spam or Ham?



Why Standardized Evaluation?

To answer questions!

Is spam filtering a viable approach?

What are the risks, costs, and benefits of filter use?

Which spam filter should I use?

How can I make a better spam filter?

What's the alternative?

Testimonials

Uncontrolled, unrepeatable, statistically bogus tests

Warm, fuzzy feelings

There's no Perfect Test

But a standardized test should

- Model real filter usage as closely as possible

- Evaluate the filter on criteria that reflect its effectiveness for its intended purpose

- Eliminate uncontrolled differences

- Be repeatable

- Yield statistically meaningful results

Future tests will

- Challenge assumptions in the current test

More information? *Google!*

cormack spam

TREC spam

DMC spam

DMC compression

ECML challenge

ROC curve

Markov model

PPM spam

OSBF Lua

Bogofilter

spamorham.org

spam conference

email anti-spam

likelihood ratio

machine learning

text classifier

Prediction by Partial Matching (PPM)

For each class:

left context occurrences

left context+prediction

log-likelihood estimate

compressed length

Smoothing/backoff:

zero occurrence problem

Adaptation:

increment counts

assuming in-class

ai.stanford.?



Context (509 spam, 1 ham)

ai.stanford.e



Prediction (0 spam, 1 ham)

ai.stanford.E



Prediction (509 spam, 0 ham)