# On-line Spam Filter Fusion

Thomas Lynam & Gordon Cormack

*originally presented at SIGIR 2006*
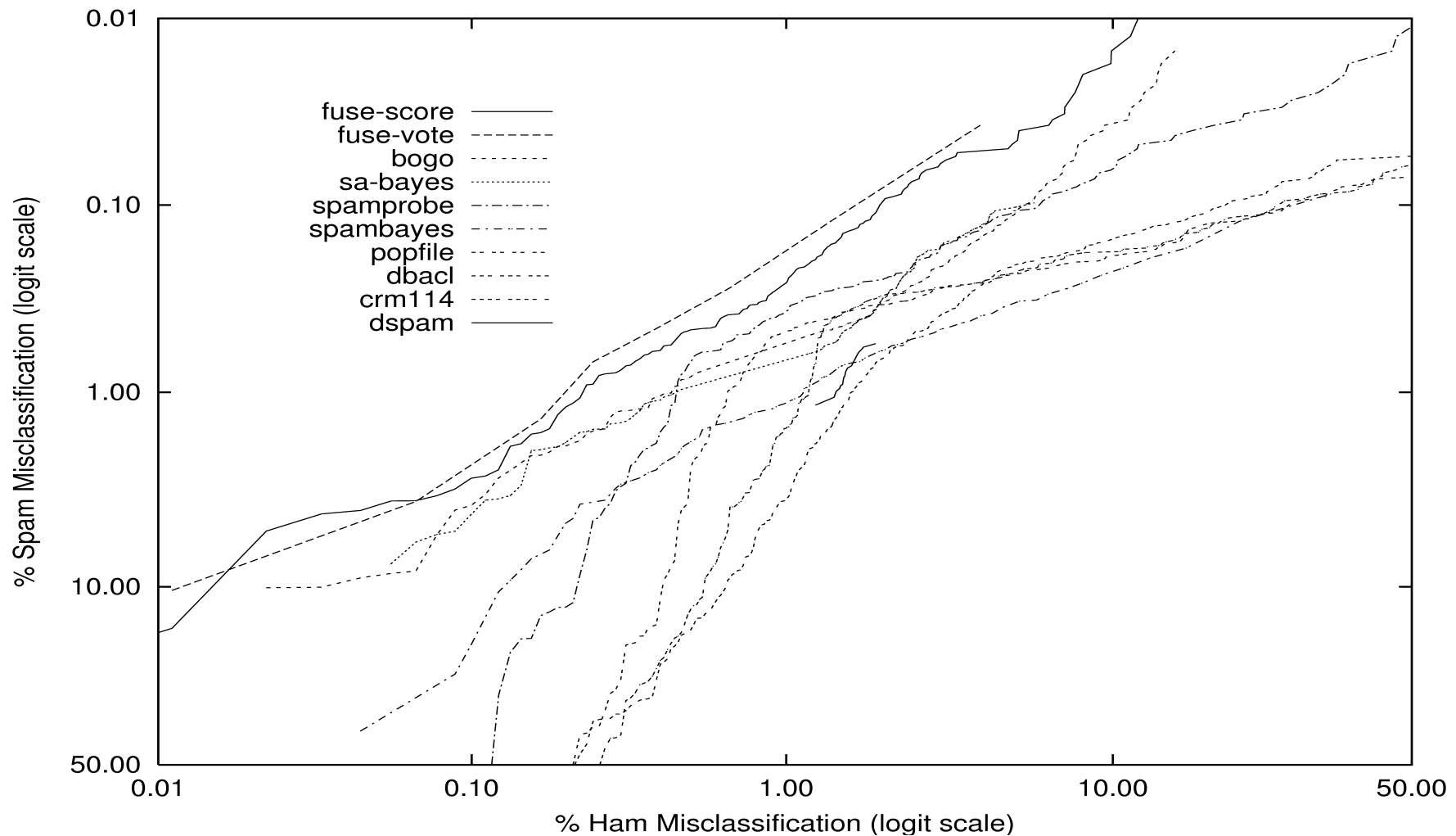
University of
Waterloo

# On-line vs Batch Classification

- Batch Hard Classifier
  - separate training and test data sets
  - Given ham/spam classification of training set
  - Compute ham/spam class for each message
- On-line Soft Classifier
  - Chronological sequence
  - Compute *spamminess* for each in sequence
    - ham/spam class by comparing to fixed threshold
  - Given ham/spam classifcation afterwards
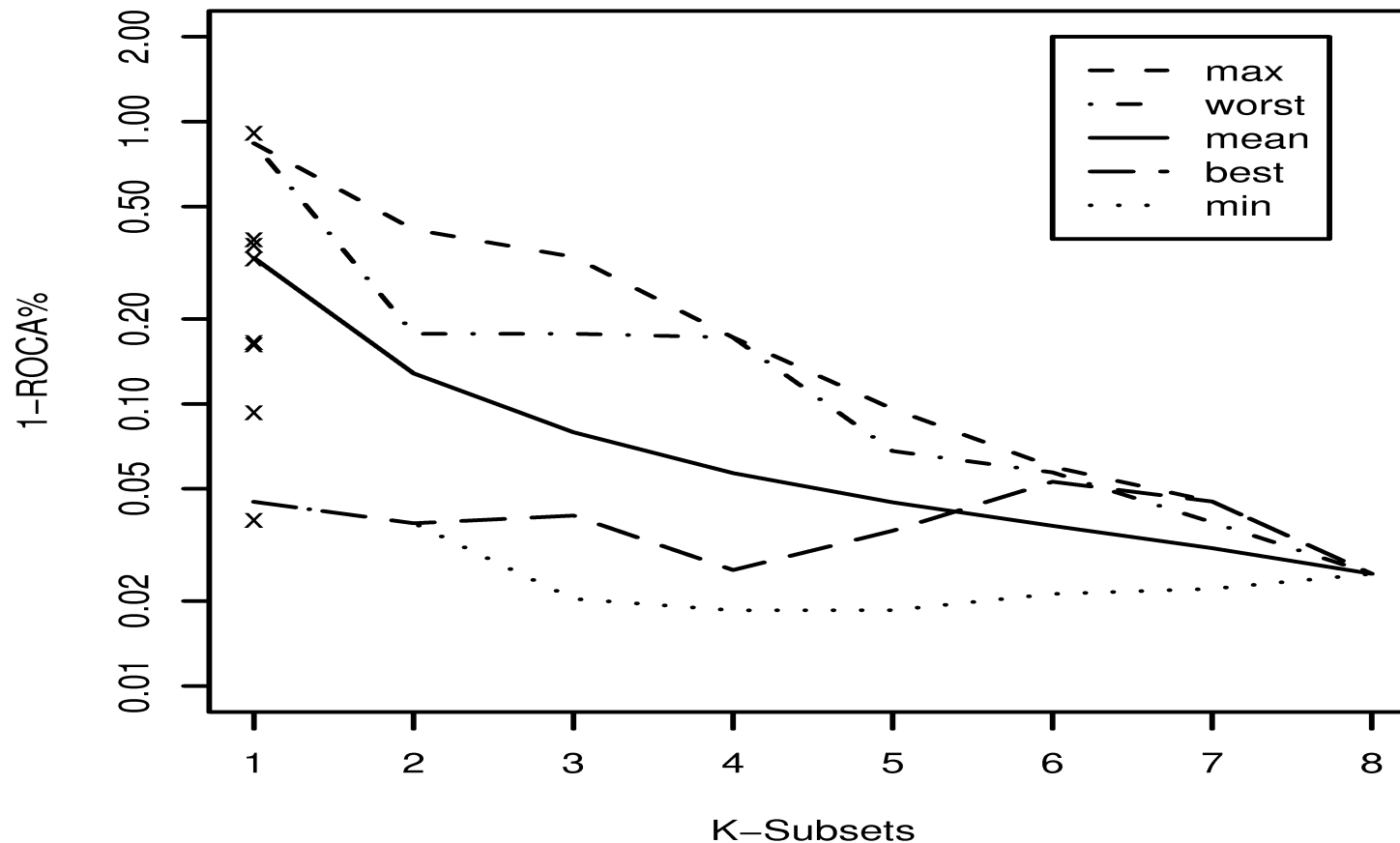    - Immediate, correct feedback (idealized user)

University of
Waterloo

# Measures of Success & Failure

- ROC Curve
- ROC Area *above* the curve (as percentage)
- Ham & spam misclassification rates
  - Sm(%) when threshold set for Hm(%) = .1
- 95% confidence intervals
  - For ROC area (logit transformed)
  - For difference between ROC areas (logit trans)
    - Significant result: difference interval excludes 0

University of
Waterloo

# Pilot Test ROC
## (Mr X corpus)



Thomas Lynam and Gordon Cormack - University of Waterloo - SIGIR 2006

# Pilot Tests K Subsets
## (Mr X corpus)



Thomas Lynam and Gordon Cormack - University of Waterloo - SIGIR 2006
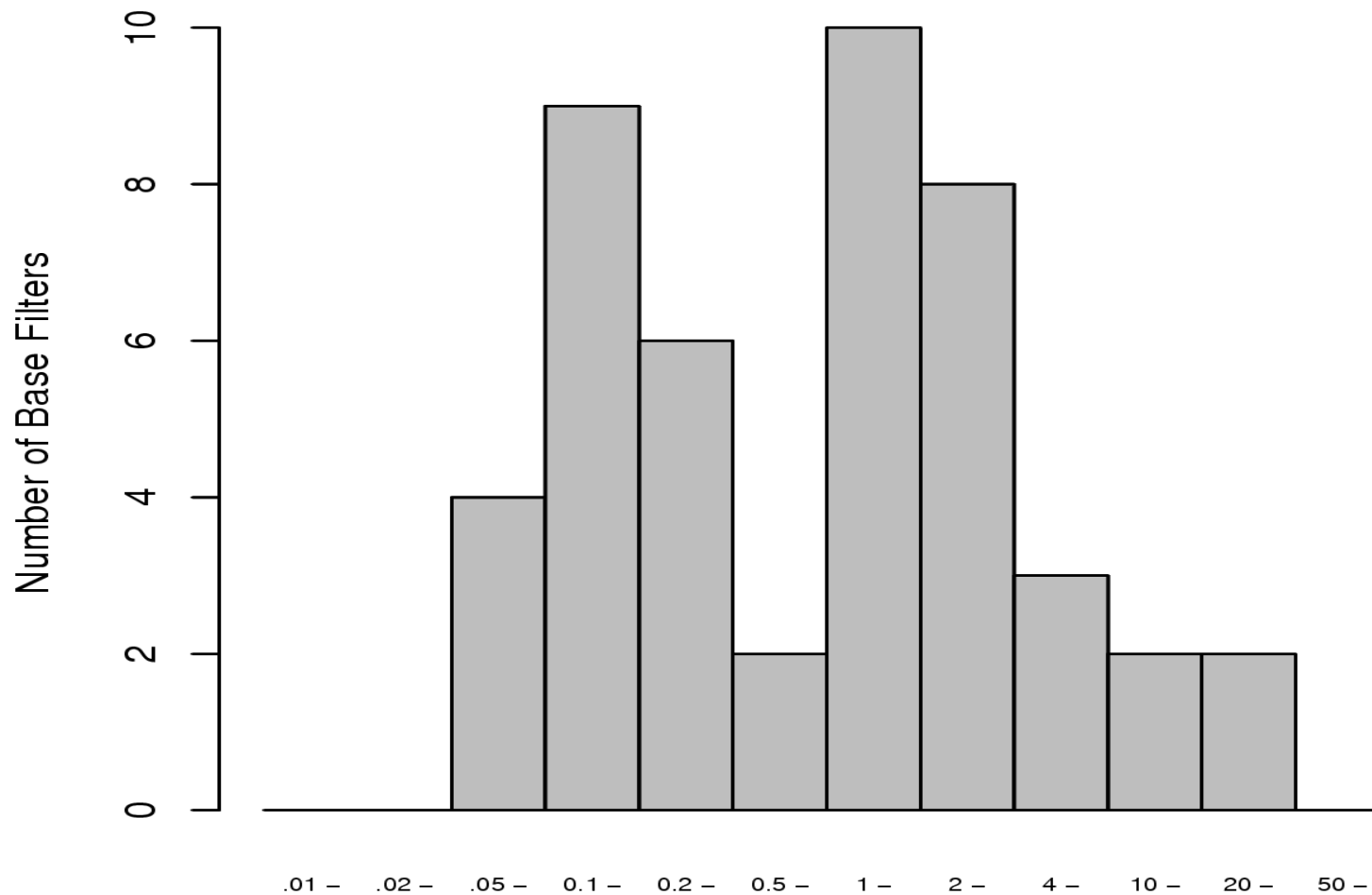
# TREC 2005 SPAM TRACK

- 4 corpora
  - 1 public, 3 private
- submit runs on public corpus
- submit filter to be run on private corpora
- 53 runs (different filters)
- 17 different organizations represented

University of
Waterloo

# TREC Spam Track Corpora

| | Ham | Spam | Total |
|---|---|---|---|
| Mr X | 9038 | 40048 | 49086 |
| S B | 6231 | 775 | 7006 |
| T M | 150685 | 19516 | 170201 |
| Full | 39399 | 52790 | 92189 |
| Aggregate | 205253 | 113129 | 318482 |

Thomas Lynam and Gordon Cormack - University of Waterloo - SIGIR 2006

# TREC Filter Performance Distribution



Number of Base Filters (y-axis: 0, 2, 4, 6, 8, 10)

(1–ROCA)% – Aggregate Pseudo–Corpus (x-axis: .01, .02, .05, 0.1, 0.2, 0.5, 1, 2, 4, 10, 20, 50)

Thomas Lynam and Gordon Cormack - University of Waterloo - SIGIR 2006

# Fusion Methods

- Best System (Baseline)
- Voting
- SumScore
- Log-odds Averaging
- SVM
- Logistic Regression

University of
Waterloo

# Log-odds Averaging

- 53 unknown systems
  - unknown min/max scores.
  - linear/nonlinear scoring
- How to normalize scores?

$$L_n = \log\left(\frac{\left|\{i < n \mid s_i \leq s_n \text{ and ith message is spam}\}\right| + \epsilon}{\left|\{i < n \mid s_i \geq s_n \text{ and ith message is ham }\}\right| + \epsilon}\right)$$
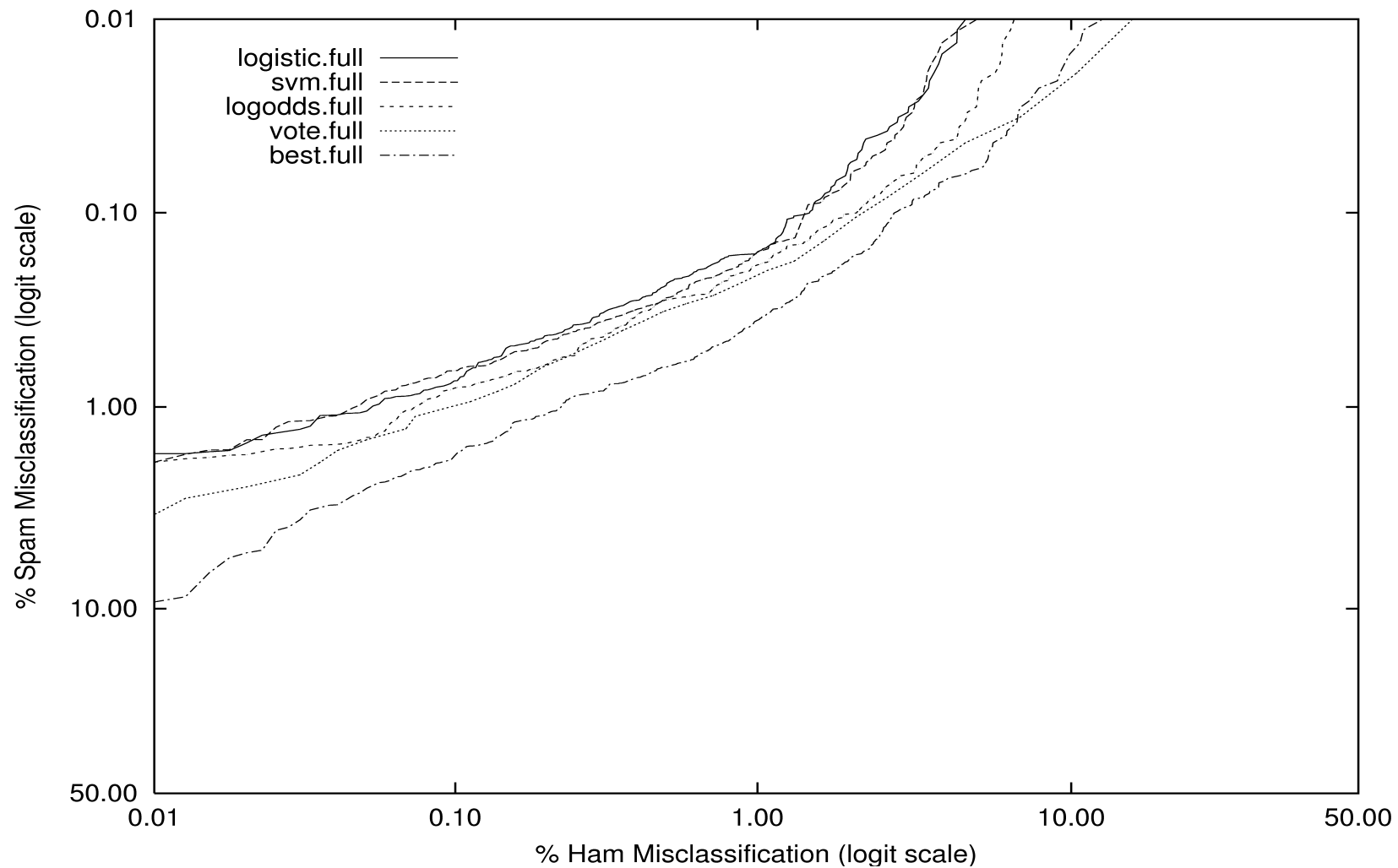
University of
Waterloo

# SVM Fusion

- SVM$^{light}$
  - default kernel and parameters
  - log-odds averaging used as features
- training set sizes of
  0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 20, 50, 100, 200, 500, 1000, 2000, 5000, 10000, 20000, 50000
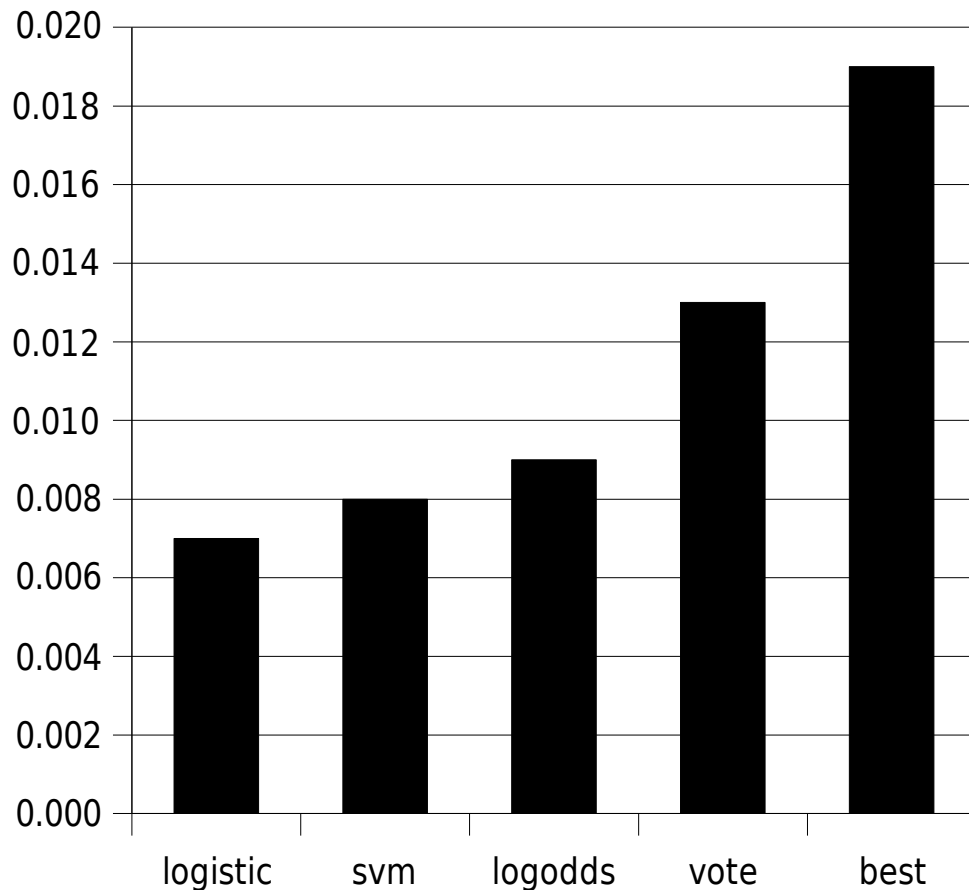- output used as spamminess score

# Logistic Regression

- LR-TRIRLS logistic regression package
- weights predict prior classification
- Negative weights considered over-fitting
- initial weight equal  1/number of filters
- training set sizes of
  0, 100, 200, 300, 400, 500, 600, 700, 800, 900, 1000, 1100, 2100, 4100, 9100, 19100, 39100, 69100, 99100, 129100, 159100.
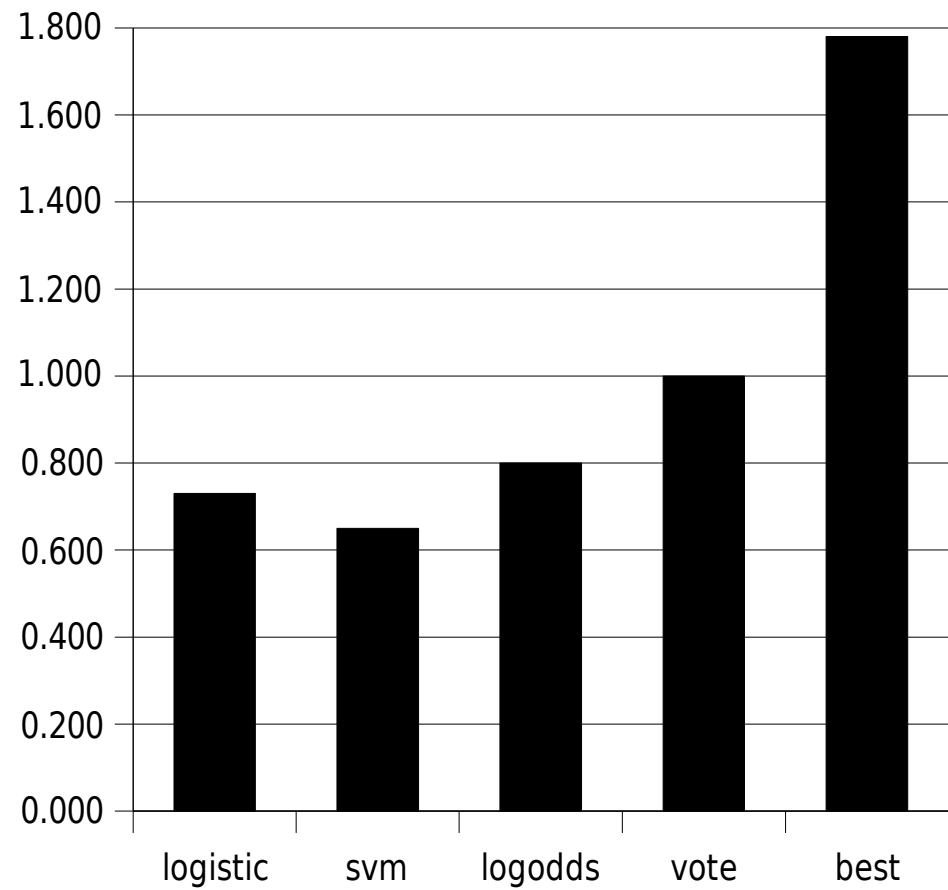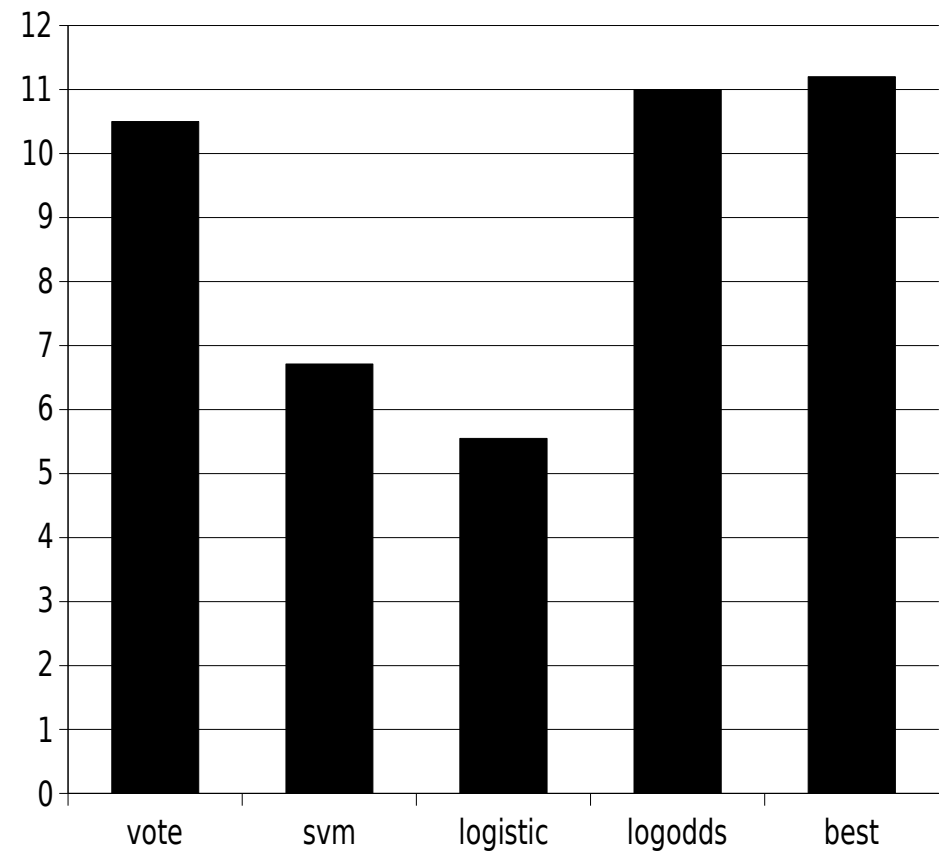- weighted average uses as spamminess score

University of
Waterloo

# ROC
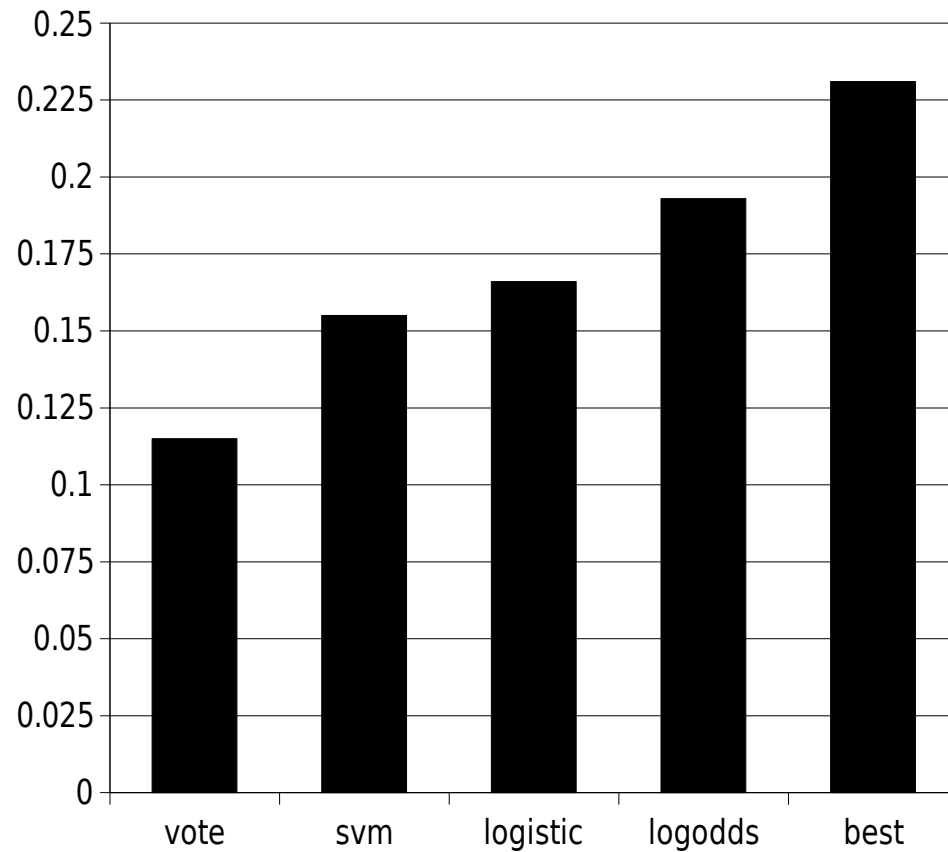## (Full Corpus)



Thomas Lynam and Gordon Cormack - University of Waterloo - SIGIR 2006

# Full Results



1-ROCA%

sm%@hm%=.1

Thomas Lynam and Gordon Cormack - University of Waterloo - SIGIR 2006

# S B Results



1-ROCA%

sm%@hm%=.1

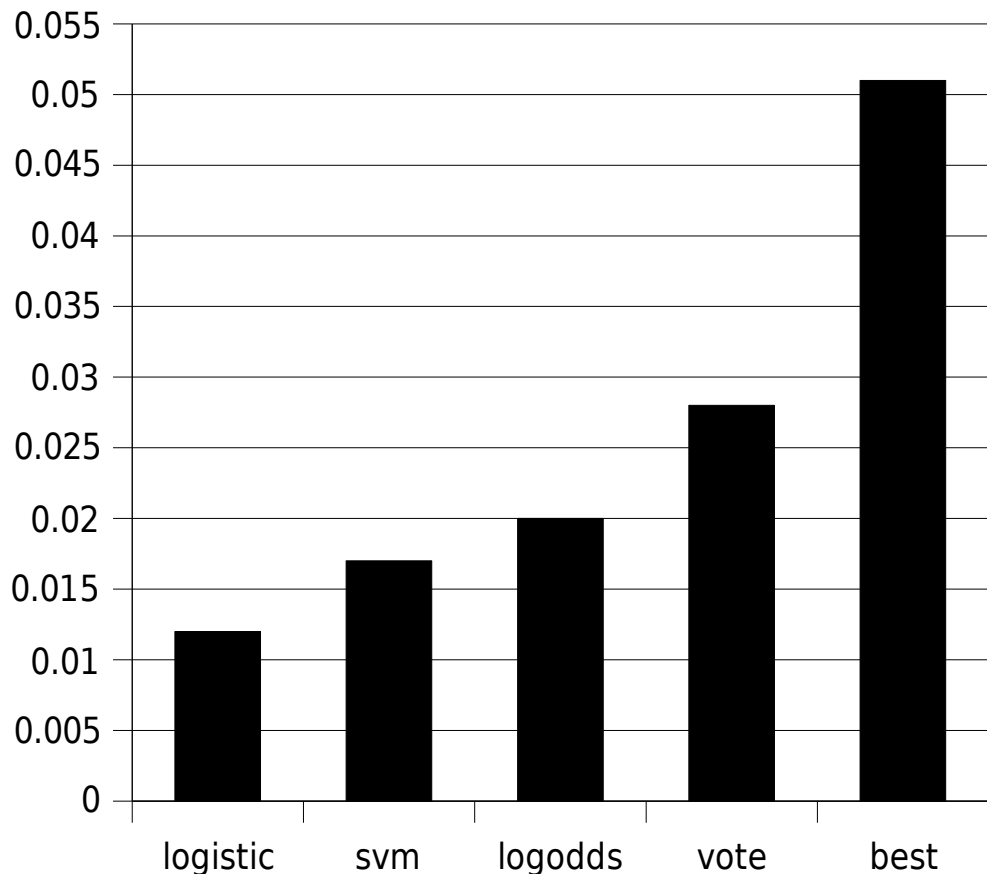Thomas Lynam and Gordon Cormack - University of Waterloo - SIGIR 2006
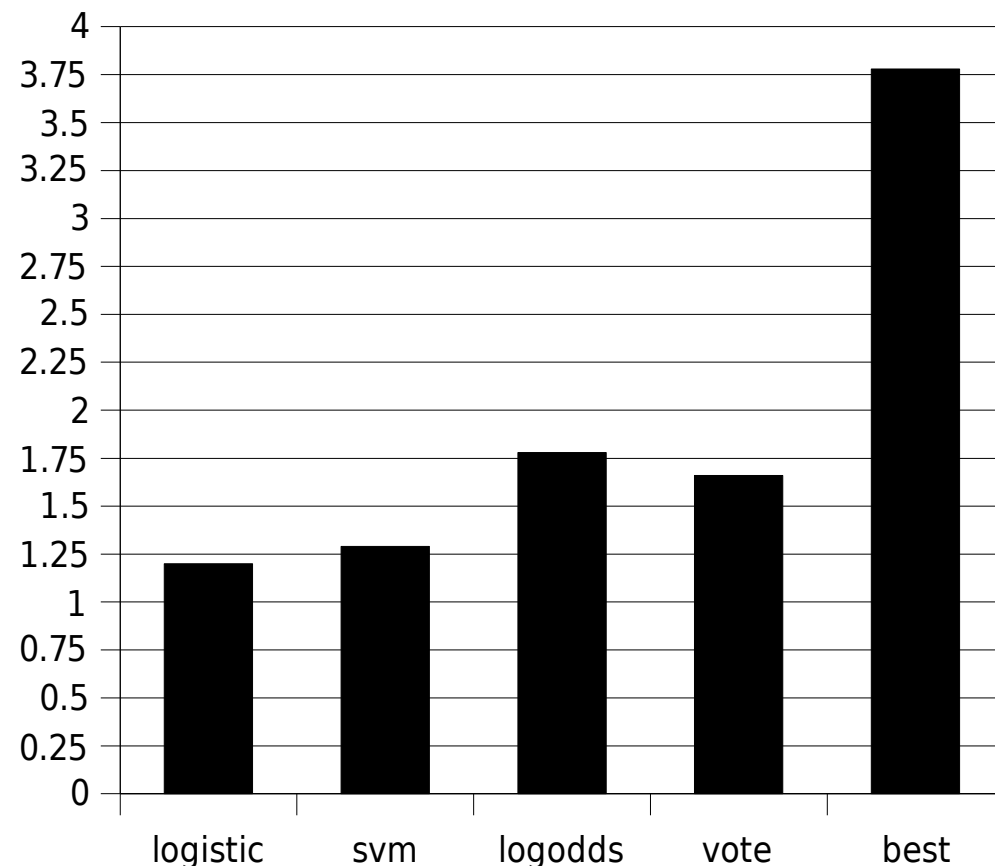
# Aggregate Results

## 1-ROCA%


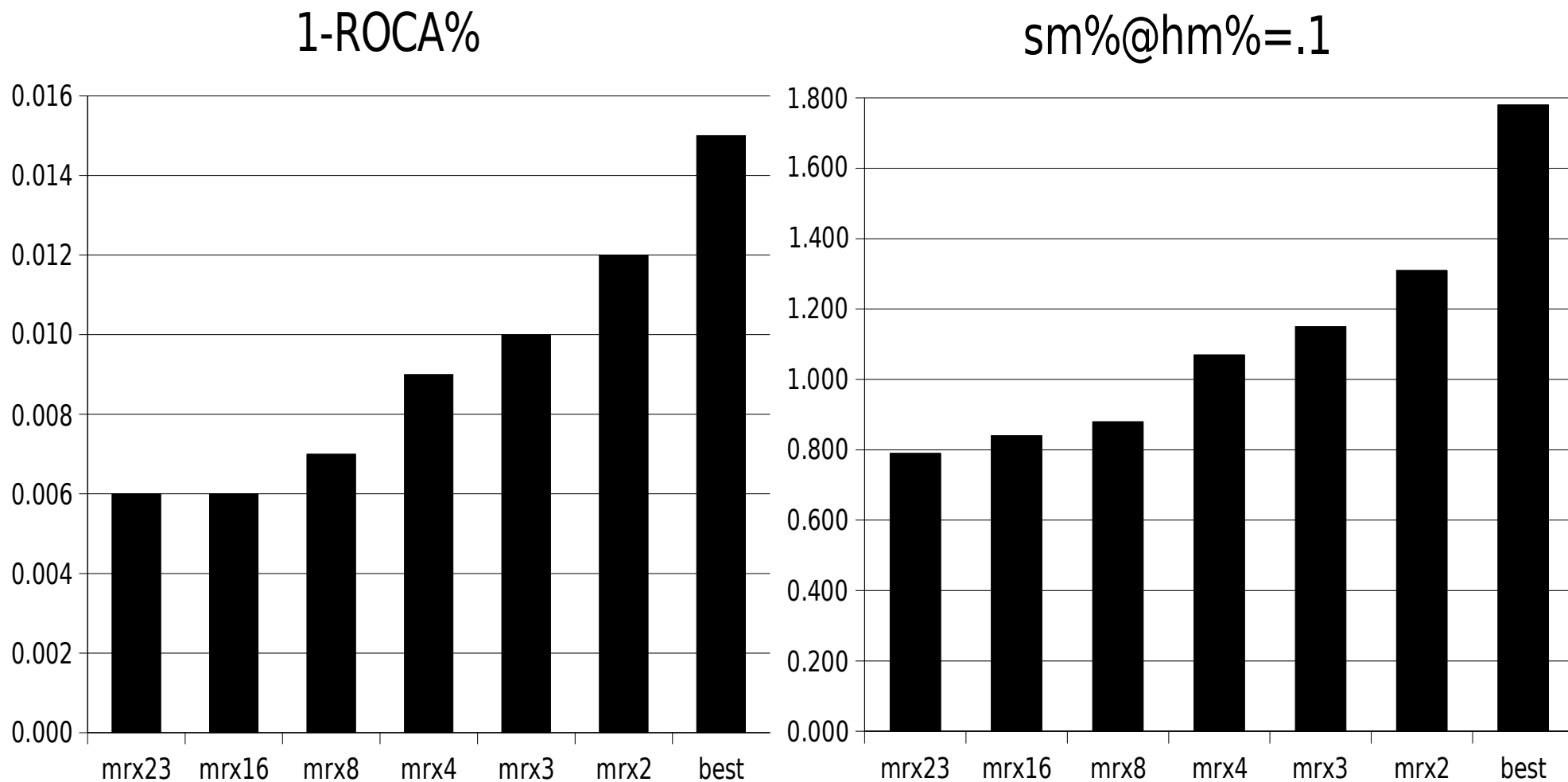
## sm%@hm%=.1

# Subset Experiment

- logistic regression subset selection
  - eliminate smallest filter weight
  - recompute logistic-regression weight
  - repeat
- train on Mr X and S B corpora
- subset size of
2, 3, 4, 8, 16 ..., largest subset with only postive weights

# Training on Mr X Corpus Results on Full Corpus



Thomas Lynam and Gordon Cormack - University of Waterloo - SIGIR 2006

# MrX-derived subsets on trec05p-1

| Subset | (1-ROCA)% | sm%@hm%=.1 |
|--------|-----------|------------|
| mrx23 | .007*** (.006-.009) | .79*** (.62-.99) |
| mrx16 | .007*** (.006-.009) | .84*** (.69-1.02) |
| mrx8 | .009*** (.007-.011) | .88*** (.71-1.08) |
| mrx4 | .012*** (.009-.015) | 1.07*** (.82-1.39) |
| mrx3 | .012*** (.010-.016) | 1.15*** (.92-1.44) |
| mrx2 | .016 (.012-.021) | 1.31** (1.01-1.68) |
| best | .019 (.015-.023) | 1.78 (1.42-2.22) |

University of Waterloo

# Base Filter Participation in Subsets
## (by Separate Performance)



Thomas Lynam and Gordon Cormack - University of Waterloo - SIGIR 2006

# TREC 06 MrX II Corpus



Thomas Lynam and Gordon Cormack - University of Waterloo - SIGIR 2006

# 1-ROCA(%) on Mrx II

- Logodds:  0.196 (.007 - .05)
- Vote:        0.224 (.009 - .05)
- Ofl:          0.363 (.02 - .06)
- Significance
  - Logodds – Ofl   p < .04   (96% confidence)
  - Vote – Ofl        p < .06   (94% confidence)

University of Waterloo

# Analysis

- All fusion methods substantially outperformed the best system
- On small corpus SVM and Logistic regression are less effective
- Voting seems more stable
- log-odds essential for other methods
- negative LR weights not always overfitting

University of Waterloo

# Conclusions

- Voting works surprisingly well
- Log-odds averaging works a little better
- Logistic Regression is slightly better
- SVM is the best for large corpus
- 53 filters not feasible
- predicting good small subsets possible

# Future Work

- explore meta analysis
- different methods of score normalization
- apply fusion to other areas

University of
Waterloo

# Questions?

University of
Waterloo

| Subset | (1-ROCA)% | | sm%@hm%=.1 | |
|---|---|---|---|---|
| mrx23 | .007*** | .006-.009 | .79*** | .62-.99 |
| mrx16 | .007*** | .006-.009 | .84*** | .69-1.02 |
| mrx8 | .009*** | .007-.011 | .88*** | .71-1.08 |
| mrx4 | .012*** | .009-.015 | 1.07*** | .82-1.39 |
| mrx3 | .012*** | .010-.016 | 1.15*** | .92-1.44 |
| mrx2 | 0.02 | .012-.021 | 1.31** | 1.01-1.68 |
| best | 0.02 | .015-.023 | 1.78 | 1.42-2.22 |

University of Waterloo

# SpamAssassin Corpus ROC curves

Thomas Lynam and Gordon Cormack - University of Waterloo - SIGIR 2006

# Mr X Corpus ROC Curves

Thomas Lynam and Gordon Cormack - University of Waterloo - SIGIR 2006