

Filtering for Personal Web Information Agents

Gabriel L. Somlo

and

Adele E. Howe

Computer Science Dept.

Colorado State University



Objectives

- Incorporate text filtering into personal Web information agents
- Desired properties of agent-embedded filtering:
 - Avoid negative feedback
 - Learn quickly, with limited training
 - Incremental learning (avoid storing training instances)

Filtering Algorithms and Parameters

- TF-IDF representation + cosine similarity
 - 1- and 2-grams
 - stop-word pruning (y/n)
 - adaptive vs. min-max-ratio dissemination threshold
- Naïve Bayes Classifier
 - use terms from unlabeled documents?
 - how to avoid using labeled negatives?

Evaluation

- Data Set: TREC Disk #5
 - FBIS: 130,471 documents
 - LATimes: 127,742 documents

- Six topics:

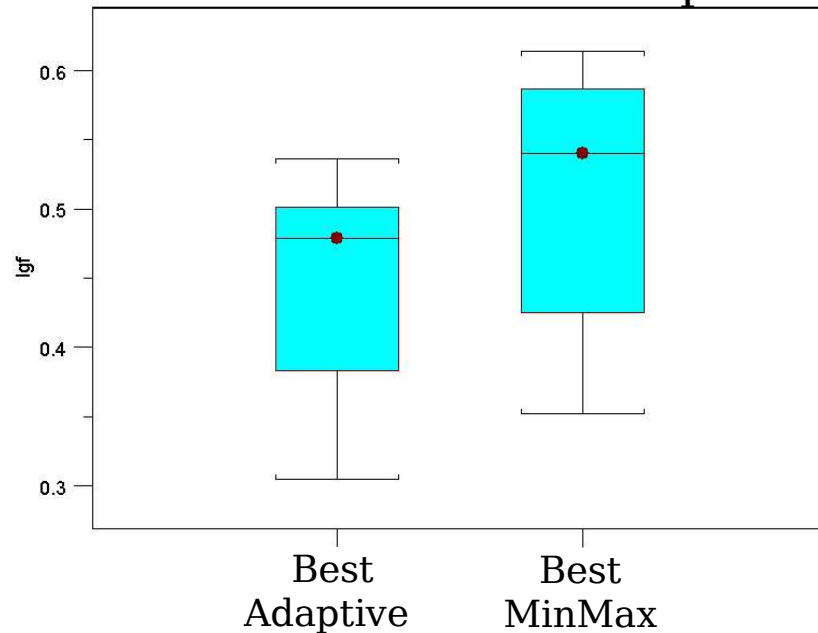
Topic	Corpus	# Relevant	# Non-Rel.	Med. Rel. Pos.
189	FBIS	584	695	62,595
301	FBIS	339	433	50,695
354	FBIS	175	715	64,424
374	LATimes	109	315	58,943
422	LATimes	98	840	70,875
426	LATimes	145	626	67,988

- Metric: harmonic mean $HM = \frac{2 \cdot \text{recall} \cdot \text{precision}}{\text{recall} + \text{precision}}$

TF-IDF Parameter Analysis

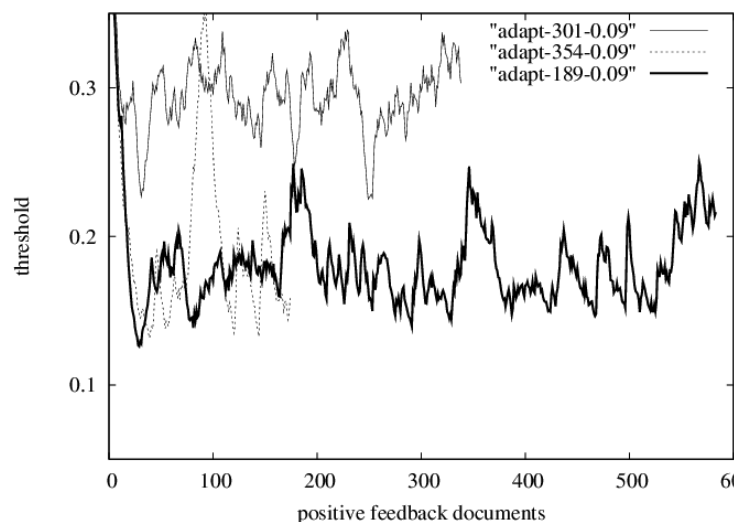
- 2-grams perform worse than single terms
- Stop-word removal does not improve HM
- Threshold learning: min-max outperforms adaptive learning

Comparison of Adaptive and MinMax across topics

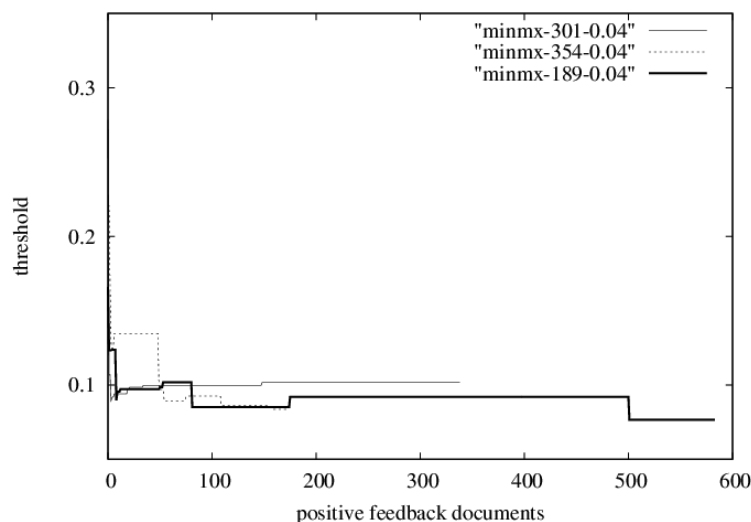


TF-IDF: Analysis of Learning the Dissemination Threshold

adaptive



min-max



- Static threshold = 0.1 comparable to best learned threshold – learning may be unnecessary!

Naïve Bayes Parameter Analysis

- Because corpus is biased toward non-relevant documents...
- Using terms from unlabeled documents is a terrible idea:

$$p(t|C) = \frac{1}{n_{\text{pos}}(C) + n_{\text{terms}}}$$

$n_{\text{pos}}(C) \equiv \#$ of term positions in class C
 $n_{\text{terms}} \equiv$ size of vocabulary

$$n_{\text{pos}}(+)\ll n_{\text{pos}}(-)$$

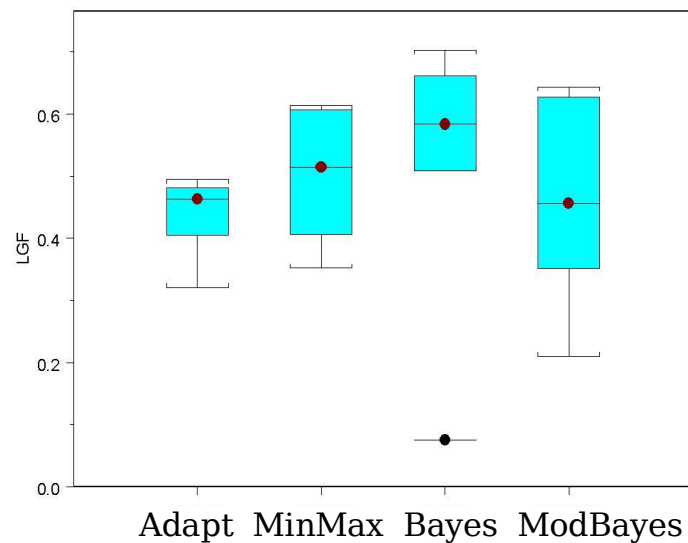
- Recall goes down as “positive” terms are discounted!

Modified Naïve Bayes without Negative Feedback

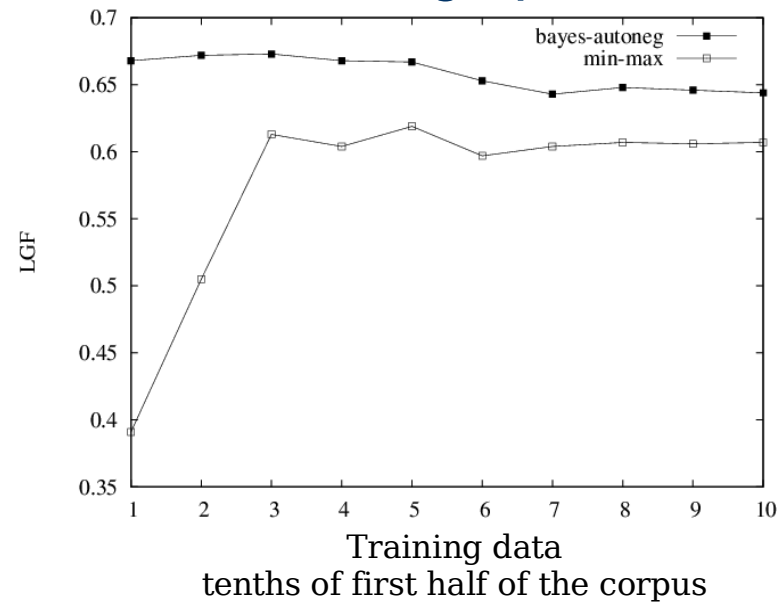
- Avoiding explicit negative feedback:
 1. Build initial classifier assuming all unlabeled docs are non-relevant
 2. Classify all unlabeled docs using initial classifier, and sort by $\Delta = [p(doc|-) - p(doc|+)]$
 3. From unlabeled docs with largest Δ , pick a number equal to that of the labeled relevant docs
 4. Build a new classifier from labeled relevant and picked non-relevant docs
- Equal numbers of relevant and non-relevant docs avoids problem shown on previous slide

Four Algorithms Compared

Harmonic Mean



learning speed



- Bayes performs better and learns faster !

Conclusions

Bayes: best performance, but requires negative feedback

ModBayes and MinMax not incremental

We may be able to bypass TF-IDF threshold learning and hard-code to 0.1

Bayes wins if we can convince users to supply negative feedback!