



UNIVERSIDAD
DE MÁLAGA

Using GNUmail to Compare Data Stream Mining Methods for On-line Email Classification

José M. Carmona-Cejudo¹, Manuel Baena-García¹, José del
Campo-Ávila¹, João Gama², Albert Bifet³ and Rafael
Morales-Bueno¹

¹Universidad de Málaga, Spain

²University of Porto, Portugal

³University of Waikato, New Zealand

Manuel Baena-García
Castro Urdiales, October 2011

- 1 INTRODUCTION
- 2 GNU_SMAIL
- 3 EVALUATION
- 4 REPLICABLE EXPERIMENTATION
- 5 CONCLUSION

CONTEXT

EMAIL MINING

- Spam detection: a two-class problem usually solved with bayesian classifiers.
- **Email classification**: a multi-class problem to sort email into folders.

EMAIL CLASSIFICATION APPROACHES

- Batch learning. The whole dataset is available before the beginning of the learning process.
- **online learning**. Data are continually being received and processed over time.

HYPOTHESIS

There is a lack of systems to compare and evaluate different machine learning models for email classification

CONTEXT

EMAIL MINING

- Spam detection: a two-class problem usually solved with bayesian classifiers.
- **Email classification**: a multi-class problem to sort email into folders.

EMAIL CLASSIFICATION APPROACHES

- Batch learning. The whole dataset is available before the beginning of the learning process.
- **online learning**. Data are continually being received and processed over time.

HYPOTHESIS

There is a lack of systems to compare and evaluate different machine learning models for email classification

CONTEXT

EMAIL MINING

- Spam detection: a two-class problem usually solved with bayesian classifiers.
- **Email classification**: a multi-class problem to sort email into folders.

EMAIL CLASSIFICATION APPROACHES

- Batch learning. The whole dataset is available before the beginning of the learning process.
- **online learning**. Data are continually being received and processed over time.

HYPOTHESIS

There is a lack of systems to compare and evaluate different machine learning models for email classification.

CONTEXT

GNU_SMAIL

GNU_Smail is a framework that allows to compare different email classification algorithms.

CONTRIBUTIONS

We introduce next improvements to GNU_Smail:

- to carry out *replicable experimentation*.
- to evaluate data stream mining methods by using:
 - sliding windows.
 - fading factors.
- to use recently proposed statistical tests to compare the performance of online algorithms.

CONTENT

1 INTRODUCTION

2 GNUSMAIL

3 EVALUATION

4 REPLICABLE EXPERIMENTATION

5 CONCLUSION

GNUsmail: ARCHITECTURE AND CHARACTERISTICS

<http://code.google.com/p/gnusmail/>

CHARACTERISTICS

- Open source framework for online adaptive email classification.
- It contains modules for reading email, preprocessing text and learning.
- The email messages are read as the model is built.

ARCHITECTURE

- **Reading email module** can obtain email messages from local filesystem or remote IMAP server.
- **Text processing module** based on filters that extract attributes from emails.
- **Learning module** into which new algorithms, methods and libraries can be integrated.

TEXT PREPROCESSING MODULE

STRUCTURE

- A pipeline of (linguistic) operators which extract relevant features from every mail.
- Some ready-to-use filters are implemented as part of the GNUsmail core, and new ones can be incorporated.

FILTERS

- Relevant words based on the ranking provided by the tf-idf function.
- Sender, CC.
- Domain of sender.
- Capital letters proportion.
- Language.
- Number of receivers.

LEARNING MODULE

STRUCTURE

Based on WEKA and MOA frameworks:

- WEKA methods are used with small datasets in environments without time and memory restrictions.
- MOA methods are used in more demanding problems.

WEKA METHODS

- Multinomial Naïve Bayes
- IBk, k-nearest neighbours
- NN-ge (Nearest Neighbour with Generalised Exemplars)

MOA METHODS

LEARNING MODULE

STRUCTURE

Based on WEKA and MOA frameworks:

- WEKA methods are used with small datasets in environments without time and memory restrictions.
- MOA methods are used in more demanding problems.

WEKA METHODS

MOA METHODS

GNUsmall uses MOA by including its tools for evaluation, classification, and drift detection.

- HoeffdingTree
- OzaBag, OzaBoost
- DDM

CONTENT

- 1 INTRODUCTION
- 2 GNUsMAIL
- 3 EVALUATION**
- 4 REPLICABLE EXPERIMENTATION
- 5 CONCLUSION

EVALUATION OF DATA STREAM MINING METHODS

In data stream contexts, neither cross-validation nor other sampling procedures are suitable for evaluation.

PREQUENTIAL MEASURES

- A prediction is made for each new example.
- Once the real class is known we update a cumulative loss function.

FORGETTING MECHANISMS

- Sliding windows.
- Fading factors (preferred method).

COMPARING THE PERFORMANCE

ADAPTED McNEMAR STATISTIC (M)

$$M_i = \text{sign}(a_i - b_i) \times \frac{(a_i - b_i)^2}{a_i + b_i}$$

- $a_i = f_i + \alpha \cdot a_{i-1}$
- $b_i = g_i + \alpha \cdot b_{i-1}$
- f_i : 1 iff the example i is misclassified by the first classifier and not by the second one (0 otherwise).
- g_i : 1 iff the example i is misclassified by the second classifier and not by the first one (0 otherwise).

CONTENT

- 1 INTRODUCTION
- 2 GNUsMAIL
- 3 EVALUATION
- 4 REPLICABLE EXPERIMENTATION**
- 5 CONCLUSION

EXPERIMENTAL SETUP

INITIAL SETUP

- Based on ENRON email dataset.
- We have selected seven specific users.
- We have used only topic folders with more than two messages.
- GNUmail checks the availability of data, offering to download it.
- The messages are analyzed in chronological order.

ATTRIBUTES

- Sender username, sender domain.
- Number of recipients, body length, capital letters proportion, size of email, subject length.
- Most relevant words.

EXPERIMENTAL SETUP

ALGORITHMS

- OzaBag over NNge, using DDM for concept drift detection.
- NNge.
- Hoeffding Trees.
- Majority class.

OUTPUT

GNUsmall plots...

- the prequential-based metrics
- and the adapted McNemar test

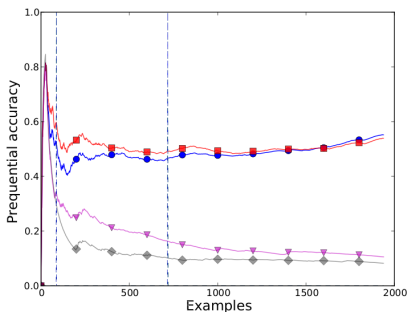
... to visually analyse the differences in performance.

RESULTS

FOLDER-WISE PREQUENTIAL ACCURACIES WITH BAGGING OF DDM AND NN-GE

Folder	Correct/Total	Percentage
beck-s (101 folders)	1071/1941	55.18%
europe	131/162	80.86%
calendar	104/123	84.55%
recruiting	89/114	78.07%
doorstep	49/86	56.97%
kaminsky-v (41 folders)	1798/2699	66.62%
universities	298/365	81.64%
resumes	420/545	77.06%
personal	154/278	55.4%
conferences	163/221	73.76%
lokay-m (11 folders)	1953/2479	78.78%
tw_commercial_group	1095/1156	94.72%
corporate	345/400	86.25%
articles	152/232	65.51%
enron_t_s	86/176	48.86%
williams-w3 (18 folders)	2653/2778	95.5%
schedule_crawler	1397/1398	99.91%
bill_williams_iii	1000/1021	97.94%
hr	74/86	86.05%
symsees	74/81	91.36%

RESULTS FOR BECK-S



(a) Prequential error

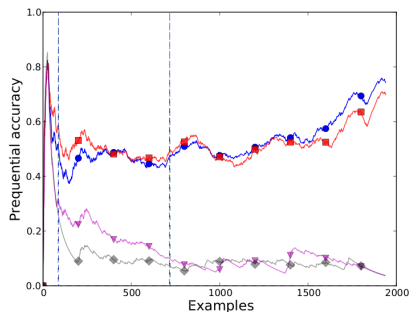
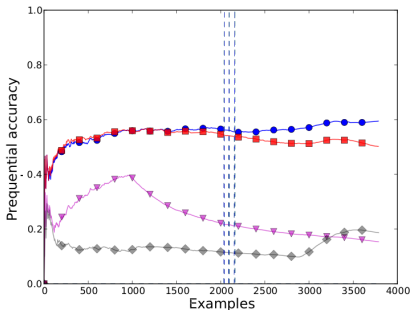
(b) Fading factors preq. ($\alpha = 0.995$)

FIGURE: Prequential based results for beck-s

RESULTS FOR KITCHEN-L



(a) Prequential error

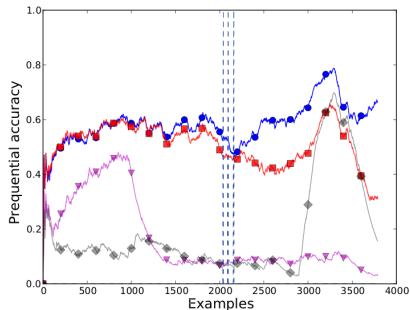
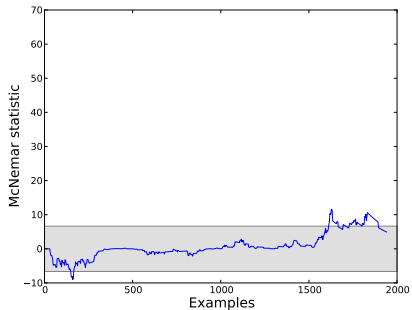
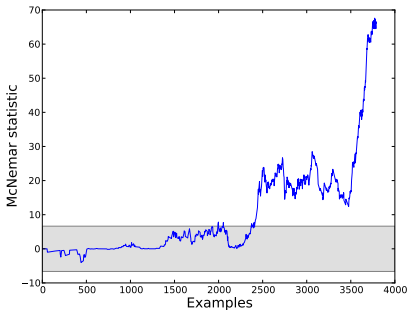
(b) Fading factors preq. ($\alpha = 0.995$)

FIGURE: Prequential based results for kitchen-l

ADAPTED McNEMAR TEST



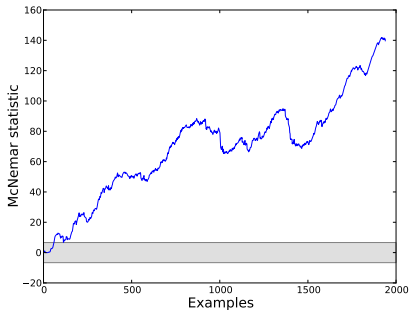
(a) beck-s



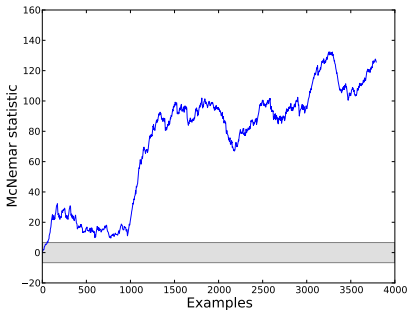
(b) kitchen-l

FIGURE: OzaBag vs. NN-ge using fading factors with $\alpha = 0.995$

ADAPTED McNEMAR TEST



(a) beck-s



(b) kitchen-l

FIGURE: OzaBag vs. Hoeffding tree using fading factors with $\alpha = 0.995$

CONTENT

- 1 INTRODUCTION
- 2 GNUMAIL
- 3 EVALUATION
- 4 REPLICABLE EXPERIMENTATION
- 5 CONCLUSION**

CONCLUSION AND FUTURE WORK

- We have presented different methods to evaluate data stream algorithms.
- We have incorporated to GNUmail recently proposed evaluation methods.
- Such evaluation methods improve prequential error measures.
- McNemar test is adequate as a tool to compare the online performance in the domain of email classification.
- Current online learning algorithm implementations needs to known all the attributes before the learning itself.
- Future methods should support online addition of new features.



UNIVERSIDAD
DE MÁLAGA

Using GNUmail to Compare Data Stream Mining Methods for On-line Email Classification

José M. Carmona-Cejudo¹, Manuel Baena-García¹, José del
Campo-Ávila¹, João Gama², Albert Bifet³ and Rafael
Morales-Bueno¹

¹Universidad de Málaga, Spain

²University of Porto, Portugal

³University of Waikato, New Zealand

Manuel Baena-García
Castro Urdiales, October 2011