

# Classifier Loss Under Metric Uncertainty

David B. Skalak<sup>1</sup>, Alexandru Niculescu-Mizil<sup>2</sup>, and Rich Caruana<sup>2</sup>

<sup>1</sup> Highgate Predictions, LLC, Ithaca, NY 14850 USA

<sup>2</sup> Cornell University, Ithaca, NY 14853 USA

{skalak, alexn, caruana}@cs.cornell.edu

<http://www.cs.cornell.edu/>

**Abstract.** Classifiers that are deployed in the field can be used and evaluated in ways that were not anticipated when the model was trained. The final evaluation metric may not have been known at training time, additional performance criteria may have been added, the evaluation metric may have changed over time, or the real-world evaluation procedure may have been impossible to simulate. Unforeseen ways of measuring model utility can degrade performance. Our objective is to provide experimental support for modelers who face potential “cross-metric” performance deterioration. First, to identify model-selection metrics that lead to stronger cross-metric performance, we characterize the expected loss where the selection metric is held fixed and the evaluation metric is varied. Second, we show that the number of data points evaluated by a selection metric has substantial impact on the optimal evaluation. While addressing these issues, we consider the effect of calibrating the classifiers to output probabilities influences. Our experiments show that if models are well calibrated, cross-entropy is the highest-performing selection metric if little data is available for model selection. With these experiments, modelers may be in a better position to choose selection metrics that are robust where it is uncertain what evaluation metric will be applied.

**Keywords:** performance metric, evaluation, calibration, cross-metric.

## 1 Introduction

Most machine learning research on classification has assumed that it is best to train and select a classifier according to the metric upon which it ultimately will be evaluated. However, this characterization makes several assumptions that we question here. What if we don’t know the metric upon which the classifier will be judged? What if the classification objective is not optimal performance, but simply robust performance across several metrics? Does it make any difference how much data is available on which to base model performance estimates? What if we want at least to avoid the worst-performing selection metrics?

In this paper we give experimental results to begin to answer questions like the ones we have just posed. The results show that the choice of selection metric depends to a large degree on how much data is available to measure performance and depends also on whether the underlying models produce accurate probabilities.

It is not so far-fetched that we may not have as much knowledge of — and access to — the ultimate evaluation metric as is usually assumed. In some situations a modeler may have the discretion to build models that optimize one of several metrics but not have access to a classification algorithm that directly optimizes the evaluation metric. For example, the modeler may decide between optimizing cross-entropy or root-mean-squared error through the choice of model class and training algorithm. But if these models are evaluated with respect to the F-score metric, it would be important to compare expected performance losses in going from cross-entropy to F-score and from root-mean-squared error to F-score. These considerations arise in natural language processing (NLP) tasks, such as noun phrase coreference resolution, where classification models may be built to maximize accuracy, but where F-score or average precision provides the ultimate measure of success [1]. In fact, NLP tasks are often evaluated on multiple reporting metrics, compounding the cross-metric problem.

The complex data processing required for NLP systems often places NLP classifiers in a pipeline where they are judged according to the performance they enable in downstream modules that receive the class predictions. Embedded classifiers may be subjected to evaluation(s) that cannot easily be tested and that may change according to evolving criteria of the entire system.

A marketing group in a large organization may request a model that maximizes response lift at 10% of the universe of customers. After the model has been built, the marketing budget for the campaign is cut, but the marketing group has the campaign ready to roll out and so not have the time to commission another model. In that case the database marketing group may decide to contact only 5% of the customers. The model that optimized response at the 10% level will now be judged in the field according to a different criterion: response from 5% of the customers. (Alternatively, the marketing group may not even specify its performance criterion, but may request a model that “simply” yields optimal profits, accuracy, and lift.) What model should be selected to be robust to changes such as these?

The availability of multiple performance metrics also poses questions for machine learning research. For example, an author may want to use a test metric that would be most acceptable to a wide readership. An author might also want to apply a second test metric under which performance is most likely to vary meaningfully from the first, and therefore provide complementary guidance.

Thus real-world considerations make evaluation more complicated than might be generally assumed. Performance metrics may change over time, may not be known, may be difficult to simulate, or may be numerous. In this paper we examine uncertain evaluation by providing experimental answers to two questions:

1. What selection metrics yield the highest performance across commonly applied evaluation metrics?
2. What is the effect of the number of data points available for making model selection judgments where the ultimate evaluation metric may be unknown?

In our experiments, we show one important factor is whether a classifier has been calibrated to output accurate probabilities. Context for all these results is

provided by a brief survey of closely related research (Section 2) and a discussion of the characteristic shape of distributions gleaned from plotting selection metric performance against evaluation metric performance (Section 6).

## 2 Related Research

As part of an extensive set of experiments, Huang and Ling defined a model selection ability measure called MSA to reflect the relative abilities of eight metrics to optimize one of three target metrics: accuracy, area under the ROC curve (“AUC”) and lift [2]. Given one of these three “goal” metrics, MSA measures the probability that one of the eight metrics correctly identifies which member of all pairs of models will be better on the goal metric. While this is an attractive summary approach, our experiments hew more closely to how we see model selection done in practice. Our experiments measure how one metric’s *best* performing models perform when measured by a second metric. Since practitioners tend to focus on superior models only, our methodology also reflects that bias. Our empirical study below also evaluates all our metrics as reporting methods rather than limiting the study to a proper subset of three goal metrics. The roles of probability calibration and classifier combination in reducing performance loss are also studied additionally here.

Several related efforts to develop algorithms to handle multiple performance criteria have also been made [3,4,5]. Additionally, Ting and Zheng [6] have provided an approach to deal with changes in costs over time.

In 2004 as part of a statistical study of AUC, Rosset showed empirically that, even where the goal is to maximize accuracy, optimizing AUC can be a superior strategy for Naive Bayes and k-nearest neighbor classifiers [7]. Joachims has extended support vector methodology to optimize directly non-linear performance measures that cannot be decomposed into measures over individual examples, and any measure derived from a contingency table [8]. Cortes and Mohri give a statistical analysis of accuracy and AUC and show that classifiers with the same accuracy can yield different AUC values when accuracy is low [9].

## 3 Experimental Design

### 3.1 Performance Metrics

The performance metrics we study are accuracy (ACC), lift at the 25th percentile (LFT), F-score (FSC), area under the ROC curve (ROC), average precision (APR), precision-recall break-even point (BEP), root-mean squared error (RMS), and mean cross-entropy (MXE). We also synthesize a hybrid metric that is defined as the equally-weighted mean performance under RMS, ROC and ACC (called “ALL”). We follow the definitions of these performance metrics found in Caruana and Niculescu [10], since they are implemented in the PERF code that was made available by Caruana in connection with the KDD Cup 2004.

We have also adopted the same conventions as to the normalization of classifier performance with respect to various metrics. Unfortunately, normalization is necessary in order to compare directly metrics with different measurement scales. Metrics have been normalized to values in  $[0, 1]$  where 0 represents the baseline performance of classifying all instances with the most frequent class in the data, and 1 corresponds to the best performance of any model developed in our lab on that data<sup>1</sup>.

### 3.2 Problems

Eleven binary classification problems were used in these experiments. ADULT, COV\_TYPE and LETTER are from the UCI Repository [11]. COV\_TYPE has been converted to a binary problem by treating the largest class as the positive and the rest as negative. We converted LETTER to boolean in two ways. LETTER.p1 treats “O” as positive and the remaining 25 letters as negative, yielding a an unbalanced problem. LETTER.p2 uses letters A-M as positives and the rest as negatives, yielding a well-balanced problem. HS is the IndianPine92 data set [12] where the difficult class Soybean-mintill is the positive class. SLAC is a problem from the Stanford Linear Accelerator. MEDIS and MG are medical data sets. COD, BACT, and CALHOUS are three of the datasets used in [13]. ADULT, COD, and BACT contain nominal attributes. For neural networks, SVMs, KNNs, and logistic regression we transform nominal attributes to boolean (one boolean per value). Each decision tree, bagged decision tree, boosted tree, boosted stump, random forest and naive Bayes model is trained twice, once with transformed attributes and once with the original ones.

### 3.3 Model Types

The 10 model types that we used in this experiment were: back-propagation neural networks, bagging of decision trees, boosting of decision trees, k-nearest neighbor, logistic regression, Naive Bayes, random forests, decision trees, boosting decision stumps and support vector machines. We create a library of approximately 2,000 models trained on training sets of size 4,000. We train each of these models on each of the 11 problems to yield approximately 22,000 models. The models are all as described in [14].

The output of such learning methods as boosted decision trees, boosted decision stumps, SVMs and Naive Bayes cannot be interpreted as well-calibrated posterior probabilities [15]. This has a negative impact on the metrics that interpret predictions as probabilities: RMS, MXE and ALL (which invokes RMS). To address this problem, we use post-training calibration to transform the predictions of all the methods into well-calibrated probabilities. In this paper calibration is done via Platt scaling [16].

---

<sup>1</sup> The performance upper bounds are available to interested researchers.

To fit the calibrated model we use a set of 1000 points reserved solely for calibration (i.e. they are not part of the training, validation or final test set).<sup>2</sup> While in practice one would use the same set of points both for calibration and for model selection, here we use separate sets in order to separate the effects of calibration from the effects of model selection on performance. The effect of calibration is further discussed in Section 5.

## 4 The Effect of Sample Size on Selection Metric Choice

In this section we discuss the effect of small data sample size on the decision of which selection metrics to use. Our primary objective in this section is to quantify the loss in selecting on one metric but reporting on another. To obtain the results in this section, we use the following methodology. For each problem, we train each of the approximately 2000 models on a 4000 points training set, and calibrate it using the extra 1000 points calibration set. All the trained models are then evaluated on a validation (selection) set, and the model with the best performance on the selection metric is found. Finally, we report the evaluation (reporting) metric performance of the best model on a final independent test set. To ensure that the results are not dependent on the particular train/validation/test set split, we repeat the experiment five times and report the average performance over the five trials.

To investigate how the size of the selection set affects the performance of model selection for different selection metrics, we consider selection sets of 100, 200, 500 and 1000 points. For comparison we also show results for “optimal” selection, where the final test set is used as the selection set.

We use the following experimental procedure. We are given a problem, a selection metric,  $s$ , and a reporting metric,  $r$ . We choose from our library the classifier  $C_s$  that has the highest normalized score under the selection metric  $s$ . We then measure the score of that classifier  $C_s$  under the reporting metric  $r$ . Call that score  $r(C_s)$ .

Next we identify the classifier  $C^*$  that has the highest performance on the reporting metric. Call that score  $r(C^*)$ . The difference  $r(C^*) - r(C_s)$  is the loss we report. The selection of  $C_s$  is done on a validation set and the reporting metric performance of both classifiers is computed on an independent test set.

Figure 1 shows the loss in performance due to model selection for nine selection metrics averaged across the nine reporting metrics. The tenth line, ORM (Optimize to the Right Metric), shows the loss of always selecting using the evaluation metric (i.e. select using ACC when the evaluation metric is ACC, ROC when the evaluation metric is ROC, etc.). On the X axis we vary the size of the selection set on a log scale. The right-most point on the graph, labeled

<sup>2</sup> This approach could give metrics affected by calibration (e.g., RMS and MXE) an advantage for model selection over metrics not affected by calibration (e.g., ACC, ROC, and LFT). To verify that this is not a problem we repeated the experiments with well-calibrated models that do not require post-training calibration (and thus do not use extra calibration data) and obtained similar results.

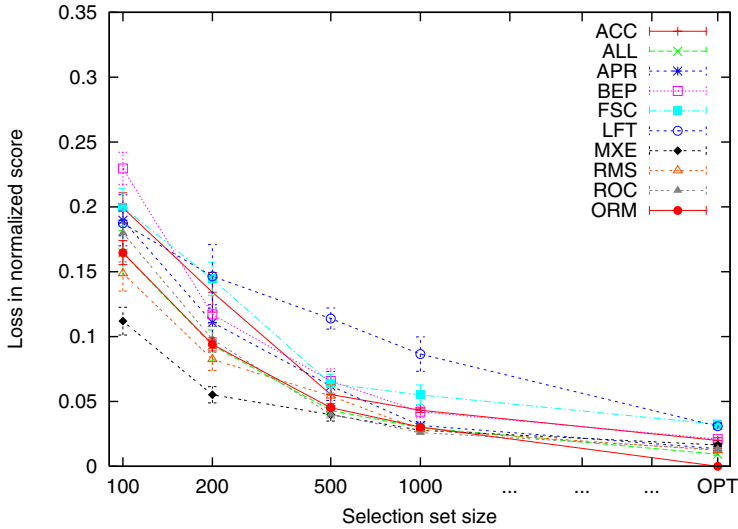


Fig. 1. Average across all nine reporting metrics

OPT, shows the loss when selection is done “optimally” (by cheating) using the final test set. This represents the best achievable performance for any selection metric, and can be viewed as a bias, or mismatch between the selection metric and the evaluation metric.<sup>3</sup>

The most striking result is the good performance of selecting on mean cross-entropy (MXE) for small sizes of the selection set. When the selection set has only 100 or 200 points, using cross-entropy as the selection metric incurs the lowest loss. In fact, at 100 and 200 points, selecting on MXE has the lowest loss for every individual reporting metric, not only on average! This may be a surprising result in that it undermines the common belief that it is always better to optimize to the metric on which the classifier will be evaluated.

We propose two hypotheses that would account for the superior performance of MXE for small data sets, but we do not yet have support for these possible explanations. MXE provides the maximum likelihood probability estimation of the binary targets. Under this hypothesis, MXE reflects the “correct” prior for target values as a binomial distribution [17]. Priors are particularly important where data are scarce. The second hypothesis recognizes that (of the metrics we consider) MXE assesses the largest penalty for large errors, which may be desirable where not much data is available.

For larger selection sets, MXE continues to be competitive, but ROC and ALL catch up when the selection set has 500 points. At 1000 points all metrics except BEP, ACC, FSC, and LFT have similar performance (on average across reporting metrics). This result suggests that, when the evaluation metric is uncertain, cross

<sup>3</sup> Of course, this bias/mismatch depends on the underlying set of classifiers to select among.

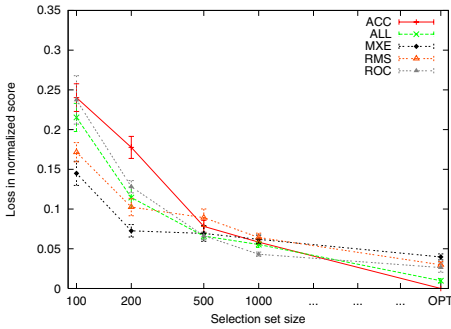


Fig. 2. Loss when reporting on ACC

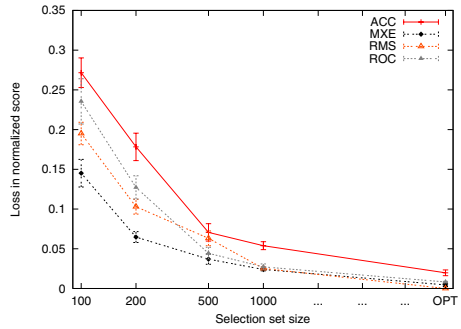


Fig. 3. Loss when reporting on RMS

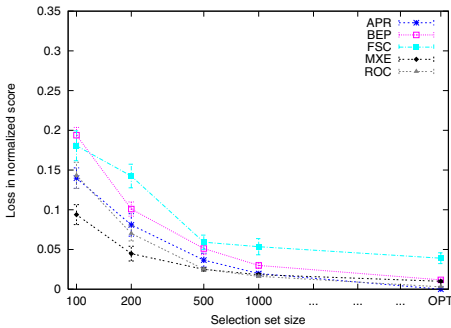


Fig. 4. Loss when reporting on APR

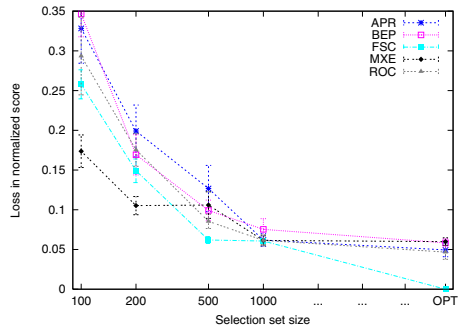


Fig. 5. Loss when reporting on FSC

entropy should be used as a selection metric, especially when validation data is scarce. When the validation set is larger, ROC, RMS and ALL also are robust selection metrics. LFT and FSC seem to be the least robust metrics, followed by BEP and ACC. Contrary to common belief, directly optimizing the evaluation metric (the ORM line) actually yields worse performance than both using MXE and RMS as a selection metric. Even for larger validation set sizes optimizing to the right metric does not yield a benefit on average.

Figure 2 shows the performance for a few selection metrics when ACC is the evaluation metric. The figure shows ROC is superior as a selection metric to ACC even when the evaluation metric is ACC. ROC-based selection yields lower loss across all selection set sizes (except of course OPT, where ACC has zero loss by definition). This confirms the observation made by Rosset [7], which was discussed in Section 2. Although at low selection set sizes MXE has the best performance (followed by RMS), looking at the OPT point, we see that MXE has the largest bias (followed by RMS). Of all metrics ALL has the smallest bias.

In the Information Retrieval (IR) community, APR is often preferred to ROC as a ranking evaluation metric because it is more sensitive to the high end of the ranking and less sensitive to the low end. Figure 4 shows the loss in normalized

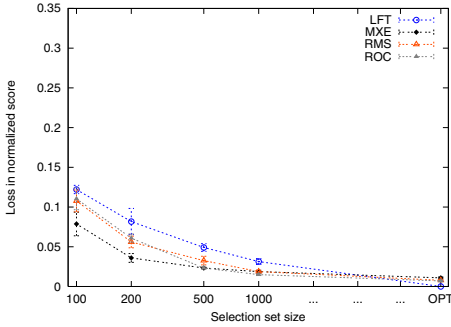


Fig. 6. Loss when reporting on LFT

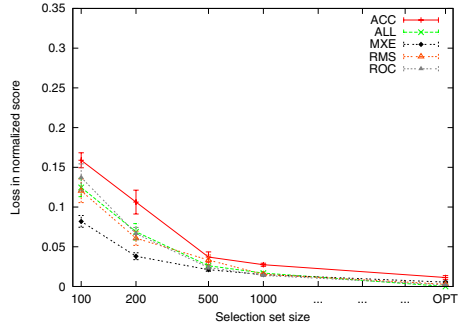


Fig. 7. Loss when reporting on ALL

score when the evaluation metric is APR. Besides APR and ROC, we also show the selection performance of MXE and two other IR metrics: BEP and FSC. The results suggest that selection based on ROC performs the same, or slightly better than selecting on APR directly. In fact ROC has a very low bias relative to APR, as shown by the OPT point in the graph. The other two IR metrics have lower performance, with FSC incurring a significantly higher loss.

Figure 5 depicts the loss in normalized score when using FSC as an evaluation metric. This figure may also be of interest to IR practitioners, since FSC is often relied upon in that field. The figure shows that, except for small validation set sizes, if FSC is the metric of interest, then FSC should also be used as a selection metric. For small validation sets, MXE again provides significantly lower loss. One other interesting observation is the large mismatch between FSC and the other metrics (the OPT point in the graph). This mismatch is one reason why, given enough validation data, FSC is the preferred selection metric when one is interested in optimizing FSC.

One other interesting case is shown in Figure 6 for LFT as the evaluation metric. The figure shows that even if one is interested in lift, one should not select based on it. MXE, RMS and ROC all lead to selecting better models.

Figure 7 shows the case when the performance is evaluated using a combination of multiple metrics. For the figure, the reporting metric is ALL which is an equally weighed average of ACC, RMS and ROC. Selecting on the more robust RMS or ROC metrics performs as well as selecting on the evaluation metric ALL. This is not the case for ACC, which is a less robust metric. For small validation sets, cross-entropy is again the best selection metric.

## 5 The Effect of Model Probability Calibration on Selection Metric Choice

In this section we investigate how cross-metric optimization performance is affected by models with poor calibration such as boosted trees, boosted stumps,



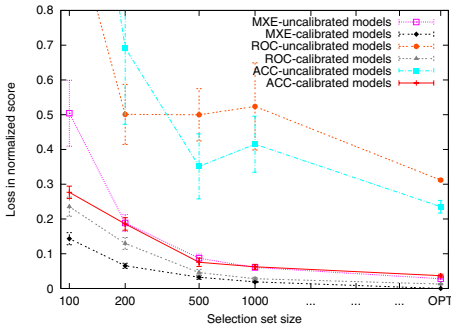


Fig. 8. Evaluation metric MXE

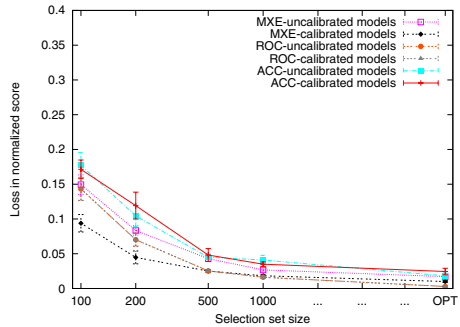


Fig. 9. Evaluation metric APR

SVMs and Naive Bayes. To this end, we repeat the experiments in Section 4, but use the original uncalibrated models instead of the Platt-calibrated ones.

As expected, having a mix of well calibrated and poorly calibrated models hurts cross-metric optimization. The effect of poorly calibrated models is twofold. On one hand, when selecting on a metric such as ROC, APR or ACC that does not interpret predictions as probabilities, and evaluating on a metric such as RMS, MXE or ALL that is sensitive to probability calibration, the selected model, while performing well on the “non-probability” measures, may be poorly calibrated, thus incurring a high loss on the “probability” measures.

This effect can be clearly seen in Figure 8. The figure shows the loss in normalized score when the reporting metric is MXE, and the selection metric is MXE, ROC or ACC. For each selection metric, two lines are shown: one for selecting from uncalibrated models, and the other for selecting from Platt-calibrated models. When selecting from uncalibrated models, using either ROC or ACC as selection metrics (the top two lines) incurs a very large loss in performance (note the scale). In fact, quite often, the MXE performance of the selected models is worse than that of the baseline model (the model that predicts, for each instance, the ratio of the positive examples in the training set). Using calibrated models eliminates this problem driving down the loss.

On the other hand, when selecting on one of the “probability” measures (RMS, MXE or ALL), the poorly calibrated models will not be selected because of their low performance on such metrics. Some of these models, however, do perform very well on “non-probability” measures such as ROC, APR or ACC. This leads to increased loss when selecting on probability measures and evaluating on non-probability ones because, in a sense, selection is denied access to some of the best models.

Figure 9 shows the loss in normalized score when the reporting metric is APR, and the selection metric is MXE, ROC or ACC. Looking at MXE as a selection metric we see that, as expected, the loss from model selection is higher when using uncalibrated models than when using calibrated ones. Since calibration does not affect ROC or APR, selecting on ROC and evaluating on APR yields

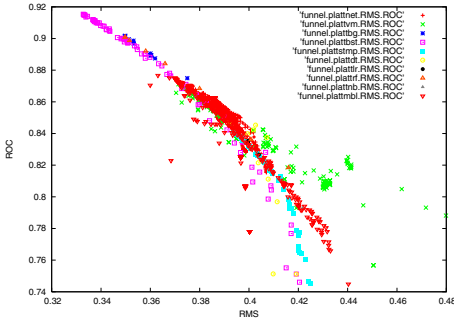


Fig. 10. Covertypes data funnel

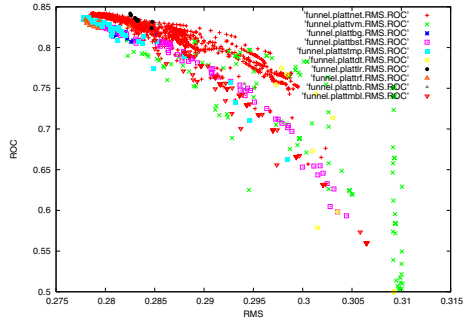


Fig. 11. Medis data funnel

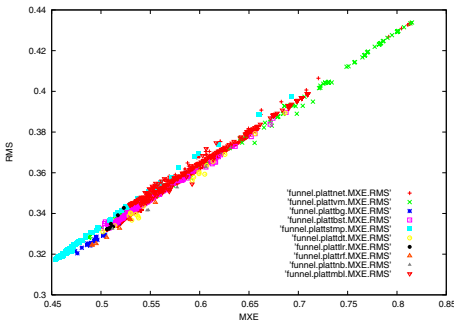


Fig. 12. Adult data funnel

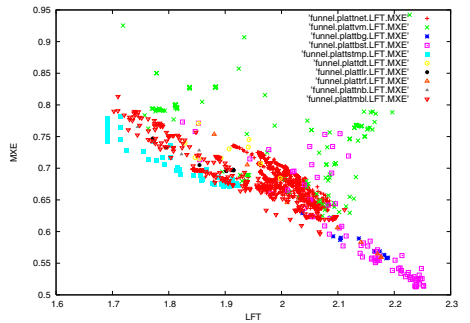


Fig. 13. Covertypes data funnel

the same results no matter if the models were calibrated or not. The same is not true when selecting using ACC because calibration can affect threshold metrics by effectively changing the threshold.

## 6 Visualizing the Joint Distribution of Selection and Evaluation Metric Performance

One way to gain further insight is to graph for each pair of metrics the distribution of performances for a large set of classifiers. For this experiment we rely on the pool of classifiers trained for the previous experiments. Recall that these classifiers came from 10 model classes. A variety of parameter settings for each model class yielded 8,910 classifiers, each of which may be evaluated according to its test-set performance for pairs of metrics. With 9 performance metrics, there are 36 plots of pairs of metrics to examine for each problem. Figures 10, 11, 12 and 13 show four of these that illustrate a variety of behaviors.

Figure 10 is a scatterplot of the ROC vs. RMS performance of the models on the Covertypes problem. In this figure boosted decision trees (purple boxes)

clearly dominate all other model types on both ROC and MXE. Bagged decision trees (blue stars) are the second best model. At the better-performing end of the spectrum (upper left of the plot) there is a strong correlation between performance on the two metrics. This correlation is reduced as performance worsens, leading to the broadening of the “funnel”.

Figure 11 shows a scatterplot for the same two metrics (ROC vs. RMS) but on the Medis problem. The shape of the funnel differs somewhat from that of Figure 10. On this problem, there is no one model type that dominates the other model types on both metrics. Calibrated boosted stumps have the best RMS, but neural nets and logistic regression yield somewhat better ROC. Also, there is less correlation between the two performance measures for different families of algorithms – the thread for each family is more distinguishable in this plot.

Figure 12 shows RMS vs. MXE performance for the Adult data set. As one might expect for two measures such as RMS and MXE that are so similar, this scatterplot shows a remarkably strong linear correlation between the two measures, with the best-performing models being neural nets and SVMs. Figure 13 shows MXE vs. LFT for Covertype. In this scatterplot there is a reasonably strong correlation between the two measures for most algorithms, but SVMs (green Xs) form a cloud of outliers with overall worse MXE.

One general feature of these “funnel” plots is that there is a narrowing at the high-performing end of the graph because it is difficult with most metrics to achieve near-optimal performance on one metric while achieving poorer performance on the other metric. When performance is poorer, however, often the funnel widens because when performance is poor on one metric it is possible to achieve a wide range of performances on other metrics. When performance is not optimal, it makes a larger difference what metric is used for selection.

The shape of the distribution of scores is seen many times for pairs of metrics. Often a wedge-shaped distribution can be seen reflecting the relatively wide variance in performance for classifiers that do not perform well along one or both of the two metrics. But we see a much tighter distribution at the vertex of the wedge for classifiers that do perform well under both metrics.

These distinctive distributions may provide a clue as to why calibrated classifiers suffer less cross-metric loss. Ensemble classifiers and calibrated classifiers both tend to yield higher-performing classifiers for a variety of metrics. In many graphs, they fall towards the narrow, extreme vertex of the wedge. At this thin edge of the plot, little variance in performance from metric to metric is seen. Consequently, cross-metric loss is lower in that region of the plot, which is inhabited by superior classifiers.

## 7 Conclusion

Our experiments have shown that when only a small amount of data is available, cross-entropy yields the strongest cross-metric performance. The experiments have also shown that calibration can affect the performance of selection metrics

in general, and of cross-entropy in particular. In general, MXE and ROC performed strongly as selection metrics and FSC, LFT, ACC, and BEP performed poorly. The next step in our research is to go beyond the empirical results presented in this paper and try to create a formal decomposition of cross-metric loss.

**Acknowledgments.** This work was supported by NSF Award 0412930.

## References

1. Munson, A., Cardie, C., Caruana, R.: Optimizing to arbitrary NLP metrics using ensemble selection. In: Proc. of Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing (HLT/EMNLP), pp. 539–546 (2005)
2. Huang, J., Ling, C.X.: Evaluating model selection abilities of performance measures. In: Evaluation Methods for Machine Learning, Papers from the AAAI workshop, Technical Report WS-06-06, AAAI, pp. 12–17 (2006)
3. Soares, C., Costa, J., Brazdil, P.: A simple and intuitive measure for multicriteria evaluation of classification algorithms. In: ECML 2000. Proceedings of the Workshop on Meta-Learning: Building Automatic Advice Strategies for Model Selection and Method Combination, Barcelona, Spain (2000)
4. Nakhaeizadeh, C., Schnabl, A.: Development of multi-criteria metrics for evaluation of data mining algorithms. In: Heckerman, D., Manilla, H., Pregibon, D. (eds.) Proceedings of the 3rd International Conference on Knowledge Discovery in Databases, Newport Beach, CA, AAAI Press, Menlo Park, CA (1997)
5. Spiliopoulou, M., Kalousis, A., Faulstich, L.C., Theoharis, T.: NOEMON: An intelligent assistant for classifier selection. In: FGML98. Number 11 in 98, Dept. of Computer Science, TU Berlin, pp. 90–97 (1998)
6. Ting, K.M., Zheng, Z.: Boosting trees for cost-sensitive classifications. In: Nédellec, C., Rouveirol, C. (eds.) ECML 1998. LNCS, vol. 1398, pp. 190–195. Springer, Heidelberg (1998)
7. Rosset, S.: Model selection via the auc. In: ICML '04: Proceedings of the Twenty-first International Conference on Machine Learning, p. 89. ACM Press, New York (2004)
8. Joachims, T.: A support vector method for multivariate performance measures
9. Cortes, C., Mohri, M.: Auc optimization vs. error rate minimization. In: Thrun, S., Saul, L., Scholkopf, B. (eds.) Advances in Neural Information Processing Systems 16, MIT Press, Cambridge (2004)
10. Caruana, R., Niculescu-Mizil, A.: Data mining in metric space: an empirical analysis of supervised learning performance criteria. In: KDD '04: Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining, pp. 69–78. ACM Press, New York (2004)
11. Blake, C., Merz, C.: UCI repository of machine learning databases (1998)
12. Gualtieri, A., Chettri, S.R., Crompton, R., Johnson, L.: Support vector machine classifiers as applied to aviris data. In: Proc. Eighth JPL Airborne Geoscience Workshop (1999)
13. Perlich, C., Provost, F., Simonoff, J.S.: Tree induction vs. logistic regression: a learning-curve analysis. *J. Mach. Learn. Res.* 4, 211–255 (2003)
14. Caruana, R., Niculescu-Mizil, A.: An empirical comparison of supervised learning algorithms. In: ICML '06: Proceedings of the 23rd International Conference On Machine Learning, pp. 161–168. ACM Press, New York (2006)

15. Niculescu-Mizil, A., Caruana, R.: Predicting good probabilities with supervised learning. In: Proc. 22nd International Conference on Machine Learning (ICML'05) (2005)
16. Platt, J.C.: Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. In: Smola, A., Bartlett, P., Schlkopf, B., Schuurmans, D. (eds.) *Advances in Large Margin Classifiers*, pp. 61–74. MIT Press, Cambridge (1999)
17. Mitchell, T.M.: *Machine Learning*. McGraw-Hill, New York (1997)