# On Reoptimizing Multi-Class Classifiers[*]

Chris Bourke[†]        Kun Deng[†]        Stephen D. Scott[†]        Robert E. Schapire[‡]

N. V. Vinodchandran[†]

## Abstract

Significant changes in the instance distribution or associated cost function of a learning problem require one to reoptimize a previously-learned classifier to work under new conditions. We study the problem of reoptimizing a multi-class classifier based on its ROC hypersurface and a matrix describing the costs of each type of prediction error. For a binary classifier, it is straightforward to find an optimal operating point based on its ROC curve and the relative cost of true positive to false positive error. However, the corresponding multi-class problem (finding an optimal operating point based on a ROC hypersurface and cost matrix) is more challenging and until now, it was unknown whether an efficient algorithm existed that found an optimal solution. We answer this question by first proving that the decision version of this problem is NP-complete. As a complementary positive result, we give an algorithm that finds an optimal solution in polynomial time if the number of classes $n$ is a constant. We also present several heuristics for this problem, including linear, nonlinear, and quadratic programming formulations, genetic algorithms, and a customized algorithm. Empirical results suggest that under both uniform and non-uniform cost models, simple greedy methods outperform more sophisticated methods.

**Keywords:** Receiver Operator Characteristic (ROC), classifier reoptimization, multi-class classification.

## 1   Introduction

We study the problem of re-weighting classifiers to optimize them for new cost models. For example, given a classifier optimized to minimize classification error on its training set, one may attempt to tune it to improve performance in light of a new cost model. Equivalently, a change in the class distribution (the probability of seeing examples from each particular class) can be handled by modeling such a change as a change in cost model. More formally, we are concerned with finding a nonnegative weight vector $(w_1, \ldots, w_n)$ to minimize

$$\sum_{i=1}^{m} c\left( y_i, \operatorname*{argmax}_{1 \le j \le n}\{w_j\, f_j(x_i)\} \right) \ , \tag{1}$$

given labeled examples $\{(x_1, y_1), \ldots, (x_m, y_m)\} \subset \mathcal{X} \times \{1, \ldots, n\}$ for instance space $\mathcal{X}$, a family of confidence functions $f_j : \mathcal{X} \to \mathbb{R}^+$ for $1 \le j \le n$, and a cost function $c : \{1, \ldots, n\}^2 \to \mathbb{R}^+$.

This models the problem of reoptimizing a multi-class classifier in machine learning. A machine learning algorithm takes a set $S = \{(x_1, y_1), \ldots, (x_m, y_m)\} \subset \mathcal{X} \times \{1, \ldots, n\}$ of labeled training examples and selects a function $F : \mathcal{X} \to \{1, \ldots, n\}$ that minimizes misclassification cost on $S$, which is $\sum_{i=1}^{m} c(y_i, F(x_i))$, where cost $c(y_i, F(x_i))$ is a nonnegative function measuring the cost of predicting class $F(x_i)$ on example $x_i$ whose true label is $y_i$. For convenience, we will assume that the classifier $F$ is represented by a set of nonnegative functions $f_j : \mathcal{X} \to \mathbb{R}^+$ for $j \in \{1, \ldots, n\}$, where $f_j(x)$ is the classifier's confidence that $x$ belongs in class

---

$j$, and $F(x) = \text{argmax}_{1 \leq j \leq n}\{f_j(x)\}$. Each confidence function $f_j$ can be seen as a "base learner" (as in a one-versus-rest strategy).

An obvious solution would be to simply rerun the base learning algorithm to reoptimize each confidence function $f_j$ for the new cost function. However, this process may be very expensive or even impossible. Thus, the task before us is to reoptimize $F$ *without* discarding the family of base learners. As an example, consider the machine learning application of predicting where a company should drill for oil. In this example the set of instances $\mathcal{X}$ consists of candidate drilling locations, each described by a set of attributes (e.g. fossil history of the site, geographic features that quantify how well oil is trapped in an area, etc.). The set of classes could be a discrete scale from 1 to $n$, where 1 indicates no oil would be found, and $n$ indicates a highly abundant supply. To learn its classifier $F$, the learning algorithm was given a set $S \subset \mathcal{X} \times \{1, \ldots, n\}$ of instance-label pairs as well as a nonnegative, asymmetric cost function $c$, where $c(j, k)$ measures the cost of money and resources of thinking that an area of class $j$ really was class $k$. (This function not only indicates the cost of committing excessive resources to an area with too little oil, but also of committing too few resources to an area with abundant oil.) Once the function $F$ is learned and put into practice, it may become the case that the cost function changes from $c$ to $c'$, e.g. if new technologies in drilling and shipping of oil emerge. If this happens, then the function $F$ is no longer appropriate to use. One option to remedy this is to discard $F$ and train a new classifier $F'$ on $S$ under cost function $c'$. However this may not be an option if the original data $S$ is unavailable, say due to proprietary restrictions. In this case the best (or perhaps only) choice is to reoptimize $F$ based on a new (possibly smaller) data set. Such problems have been studied extensively (Fieldsend & Everson, 2005; Ferri et al., 2003; Hand & Till, 2001; Lachiche & Flach, 2003; Mossman, 1999; O'Brien & Gray, 2005; Srinivasan, 1999).

For learning tasks with only $n = 2$ classes, this problem is equivalent to that of finding the *optimal operating point* of a classifier given a ratio of true positive cost to false positive cost and has a straightforward solution via Receiver Operating Characteristic (ROC) analysis (Provost & Fawcett, 1997; Provost & Fawcett, 1998; Provost & Fawcett, 2001; Lachiche & Flach, 2003). ROC analysis takes a classifier $F$ that outputs confidences in its predictions (i.e. a ranking classifier), and precisely describes the tradeoffs between true positive and false positive errors. By ranking all examples $x \in S$ by their confidences $f(x)$ from largest to smallest (denoted $S = \{x_1, \ldots, x_m\}$), one achieves a set of $m + 1$ binary classifiers by setting thresholds $\{\theta_i = (h(x_i) + h(x_{i+1}))/2, 1 \leq i < m\} \cup \{h(x_1) - \epsilon, h(x_m) + \epsilon\}$ for some constant $\epsilon > 0$. Given a relative cost $c$ of true positive error to false positive error and a validation set $S$ of labeled examples, one can easily find the optimal threshold $\theta$ based on $S$ and $c$ (Lachiche & Flach, 2003). To do so, simply rank the examples in $S$, try every threshold $\theta_i$ as described above, and select the $\theta_i$ minimizing the total cost of all errors on $S$.

Though the binary case lends itself to straightforward optimization, working with multi-class problems makes things more difficult. A natural idea is to think of an $n$-class ROC space having dimension $n(n-1)$. A point in this space corresponds to a classifier, with each coordinate representing the misclassification rate of one class into some other class[1]. According to Srinivasan (1999), the optimal classifier lies on the convex hull of these points. Given this ROC polytope, a validation set, and an $n \times n$ cost matrix $M$ with entries $c(y, \hat{y})$ (the cost associated with misclassifying a class $y$ example as class $\hat{y}$), Lachiche and Flach (2003) define the optimization problem as finding a weight vector $\vec{w} \geq \vec{0}$ to minimize (1).

No efficient algorithm is known to optimally solve this problem for $n > 2$, and Lachiche and Flach (2003) speculate that the problem is computationally hard. We present a proof that the decision version of this problem is in fact NP-complete. As a complementary positive result, we give an algorithm that finds an optimal solution in polynomial time (w.r.t. the number of examples $m$) when the number of classes $n$ is constant. We also present several new heuristics for this problem, including an integer linear programming relaxation, a sum-of-linear fractional functions (SOLFF) formulation, and a quadratic programming formulation as well as a direct optimization of (1) with a genetic algorithm. Finally, we present a new custom algorithm based on partitioning the set of all classes into two *metaclasses*. This algorithm is similar to that of Lachiche and Flach (2003), but is more flexible in how a hypothesis is formed.

We compared all methods on a substantial number of data sets in four settings: optimization and generalization (performance on a training set and an independent testing set) in both uniform and non-uniform

---

[1]Assuming that cost is zero if the classification is correct, we need only $n(n-1)$ instead of $n^2$ dimensions.

cost settings. Though most algorithms were able to show an improvement over the base learner, in every setting, our MetaClass algorithm, an off-the-shelf genetic algorithm, and the algorithm of Lachiche and Flach (2003) consistently outperformed other methods. Among these three leaders, there are instances in which our MetaClass algorithm significantly outperforms the other two. However, overall, all three methods are equally good at improving the base classifier while generalizing well.

The rest of this paper is as follows. In Section 2 we discuss related work. In Section 3 we prove the decision version of this problem (which we call REWEIGHT) is NP-complete and in Section 4 we present an algorithm for producing an optimal solution that is efficient for a constant number of classes. Next, in Section 5 we discuss our heuristic approaches to this problem. We then experimentally evaluate our algorithms in Section 6 and conclude in Section 7.

## 2  Related Work

The success of binary ROC analysis gives hope that it may be possible to adapt similar ideas to multi-class scenarios. However, research efforts (Srinivasan, 1999; Hand & Till, 2001; Ferri et al., 2003; Lachiche & Flach, 2003; Fieldsend & Everson, 2005) have shown that extending current techniques to multi-class problems is not a trivial task. One key aspect to binary ROC analysis is that it is highly efficient to represent trade-offs of misclassifying one class into the other via binary ROC curves. In addition, the "area under the curve" (AUC) nicely characterizes the classifier's ability to produce correct rankings without committing to any particular operating point. Decisions can be postponed until a desired trade-off is required (e.g. finding the lowest expected cost).

Now consider the problem of classification in an $n$-class scenario. A natural extension from the binary case is to consider a multi-class ROC space as having dimension $n(n-1)$. A point in this space corresponds to a classifier with each coordinate representing the misclassification rate of one class into some other class. Following from Srinivasan (1999), the optimal classifier lies on the convex hull of these points.

Previous investigations have all shared this basic framework (Mossman, 1999; Srinivasan, 1999; Hand & Till, 2001; Ferri et al., 2003; Lachiche & Flach, 2003; Fieldsend & Everson, 2005; O'Brien & Gray, 2005). They differ, however, in the metrics they manipulate and in the approach they use to solve multi-class optimization problems. Mossman (1999) addressed the special case of three-class problems, focusing on the statistical properties of the volume under the ROC surface. This motivated the later work of Ferri et al. (2003), Lachiche and Flach (2003), and O'Brien and Gray (2005). Hand and Till (2001) extended the definition of two-class AUC by averaging pairwise comparisons. They used this new metric in simple, artificial data sets and achieved some success. Ferri et al. (2003) took a different approach in which they strictly followed the definition of two-class AUC by using "volume under surface" (VUS). They were able to compute the bounds of this measure in a three-class problem by using Monte Carlo methods. However, it is not known how well this measure performs in more complex problems.

Fieldsend and Everson (2005), Lachiche and Flach (2003) and O'Brien and Gray (2005) developed algorithms to minimize the overall multi-class prediction accuracy and cost given some knowledge of a multi-class classifier. In particular, Fieldsend and Everson approximate the ROC Convex Hull (ROCCH) (Provost & Fawcett, 1997; Provost & Fawcett, 1998; Provost & Fawcett, 2001) using the idea of *Pareto front*. Consider the following formulation: let $R_{j,k}(\theta)$ be the misclassification rate of predicting examples from class $j$ as class $k$. This is a function of some generalized parameter $\theta$ that depends on the particular classifiers. For example, $\theta$ may be a combination of a weight vector $\vec{w}$ and hypothetical cost matrix $M$. The goal is to find $\theta$ that minimizes $R_{j,k}(\theta)$ for all $j, k$ with $j \neq k$. Consider two classifiers $\theta$ and $\phi$. Fieldsend and Everson say $\theta$ *strictly dominates* $\phi$ if all misclassification rates for $\theta$ are no worse than $\phi$ and at least one rate is strictly better. The set of all feasible classifiers such that no one is dominated by the other forms the Pareto front. Fieldsend and Everson present an evolutionary search algorithm to locate the Pareto front. This method is particularly useful when misclassification costs are not necessarily known.

More closely related to our work are the results of Lachiche and Flach (2003) and O'Brien and Gray (2005). Lachiche and Flach considered the case when the misclassification cost is known, and the goal is to find the optimal decision criterion that fits the training set. Recall that this can be solved optimally for the

binary case. In particular, only one threshold $\theta$ is needed to make the decision for two-class problems. Since there are only $m+1$ possible thresholds for $m$ examples, it is efficient enough to simply test all possibilities and select the one that gives the minimum average error (or cost). However, the situation is more complicated for multi-class problems.

Lachiche and Flach (2003) formulated the multi-class problem as follows. Suppose the multi-class learning algorithm will output a positive, real-valued function $f : \{x_1, \ldots, x_m\} \times \{C_1, \ldots, C_n\} \to \mathbb{R}^+$. Here, $f_j(x_i)$ gives the confidence that example $x_i$ belongs to class $j$. The decision criterion simply assigns example $x_i$ to the class with maximum score. Reweighting the classes involves defining a nonnegative weight vector $\vec{w} = (w_1, w_2, \ldots, w_n)$, and predicting the class for an example $x$ as

$$h(x) = \underset{1 \leq j \leq n}{\operatorname{argmax}} \left\{ w_j \, f_j(x) \right\} \ .$$

Since $\vec{w}$ has only $n-1$ degrees of freedom, we can fix $w_1 = 1$.

Lachiche and Flach (2003) used a hill-climbing heuristic to find a good weight vector $\vec{w}$. In particular, they took advantage of the fact that the optimal threshold for the two-class problem can be found efficiently. For each coordinate in the weight vector, they mapped the problem to a binary problem. The algorithm starts by assigning $w_1 = 1$ and all other weights 0. It then tries to decide the weight for one class at a time as follows. Let $S$ be the set of labeled examples and let $j$ be the current class for which we want to assign a "good" weight $w_j$. Then the set of possible weights for $w_j$ is

$$\left\{ \frac{\max_{i \in \{1, \ldots, j-1\}} f_i(x)}{f_j(x)} \ \middle| \ x \in S \right\} \ .$$

It is not difficult to see that at any stage there are at most $O(|S|)$ possible weights that can influence the prediction. Thus choosing the optimal weight in this setting can be easily achieved by checking all possibilities. Overall, their algorithm runs in time $\Theta(nm \log m)$. Though there is no guarantee that this approach can find an optimal solution, they gave empirical results suggesting that it works well for optimizing 1BC, a logic-based Bayes classifier (Lachiche & Flach, 1999).

Although only briefly mentioned by Lachiche and Flach (2003), this ROC thresholding technique is quite extensible to cost-sensitive scenarios. O'Brien and Gray (2005) investigated the role of a cost matrix in partitioning the estimated class probability space and as a replacement for the weights. Assuming that $M$ is a misclassification cost matrix, an optimal decision criterion would be

$$h(x) = \underset{1 \leq j \leq n}{\operatorname{argmax}} \left\{ \sum_{1 \leq k \leq n} c(j, k) \, \hat{p}_k(x) \right\} \ .$$

If $\hat{p}_k(x)$ is a good probability estimate of example $x$ belonging to class $k$, this prediction results in the lowest expected cost. However, if $\hat{p}_k(x)$ is not an accurate probability estimate, then to ensure optimality, the cost matrix $M$ has to be altered accordingly. Thus the cost matrix $M$ plays a similar role as the weight vector of Lachiche and Flach (2003) in defining the decision boundary in estimated probability space. O'Brien and Gray (2005) defined several standard operations to manipulate the cost matrix $M$ and proposed the use of a greedy algorithm to find the altered cost matrix (called a *boundary matrix*).

While the multi-class problem has been studied via heuristics, no one has yet answered the question as to whether this problem is hard, and no one has found efficient algorithms to solve restricted cases of the multi-class problem. Below we provide answers to both of these open questions as well as extend the current literature of heuristics.

## 3 Hardness

We now prove our hardness result of this problem. For convenience, in this section we use $f_{ij}$ to denote $f_i(x_j)$ and identify true/false with 1/0. We will show hardness for the uniform cost case, i.e. $c(j, k) = 1$ when $j \neq k$ and 0 otherwise. This of course implies hardness for the general case.

**Definition 1.** Problem **Reweight**

Given: nonnegative real numbers $f_{ij}$ ($i = 1, \ldots, m$, $j = 1, \ldots, n$), integers $y_i \in \{1, \ldots, n\}$, and an integer $K$.
Question: does there exist a vector of nonnegative real numbers $(w_1, \ldots, w_n)$ such that

$$\left| \left\{ i : \max_{j \neq y_i} \{w_j f_{ij}\} \geq w_{y_i} f_{iy_i} \right\} \right| \leq K \ ? \tag{2}$$

In other words, the problem is to find a vector $\vec{w} = (w_1, \ldots, w_n)$ that maximizes how often $w_j f_{ij}$ is maximized (over $j$) by the correct label $y_i$.

To prove the hardness of REWEIGHT, we will reduce from the minimum satisfiability problem MINSAT, shown to be NP-complete by Kohli et al. (1994).

**Definition 2.** Problem **MinSat** (Kohli et al., 1994)

Given: a set of disjunctions of pairs of literals

$$\begin{aligned}
\ell_{11} &\quad \vee \quad \ell_{12} \\
\ell_{21} &\quad \vee \quad \ell_{22} \\
&\quad \vdots \\
\ell_{m1} &\quad \vee \quad \ell_{m2} \ ,
\end{aligned}$$

where each $\ell_{ij}$ is a boolean variable $x_i$ or its negation $\neg x_i$. We are also given an integer $K$.
Question: does there exist a setting of $x_1, \ldots, x_n$ such that the number of satisfied clauses (disjuncts) is at most $K$?

**Theorem 1.** REWEIGHT *is* NP-*complete.*

*Proof.* First, it is easy to see that REWEIGHT is in NP. The certificate is simply the weight vector $\vec{w}$. This certificate is sized polynomially in the size of the input since its required precision is polynomially proportional to the precision of the input (the number of bits to represent each $f_{ij}$). We now reduce from MINSAT. Note that a special case of the constraint

$$\max_{j \neq y_i} \{w_j f_{ij}\} \geq w_{y_i} f_{iy_i}$$

used in (2) is an inequality of the form

$$w_{j_0} f_{ij_0} \geq w_{y_i} f_{iy_i} \tag{3}$$

for one particular $j_0 \neq y_i$. This can be seen simply by setting all of the other $f_{ij}$'s to zero, which gives

$$\max\{0, w_{j_0} f_{ij_0}\} \geq w_{y_i} f_{iy_i} \ . \tag{4}$$

Since in our construction the $w_j$'s and $f_{ij}$'s are nonnegative, these are equivalent. So, in what follows, we give constraints of the form of (3), but these really are of the form of (4). Thus while for the sake of clarity we map instances of MINSAT to inequalities, it is straightforward to convert these to a collection of $f_{ij}$ and $y_i$ values in REWEIGHT.

Given an instance of MINSAT as above, we create an instance of REWEIGHT. The new instance has $n' = 2n + 1$ weights: $v_0$; $w_1, \ldots, w_n$; and $w'_1, \ldots, w'_n$. The weight $v_0$ is forced to be strictly positive, and is used as a reference for all other weights. Roughly speaking, $w_i$ will correspond to boolean variable $x_i$ and $w'_i$ will correspond to its negation. More specifically, we will force $w_i$ to have a value close to $2v_0$ if $x_i$ is true, and a value close to $v_0$ otherwise; $w'_i$ will be forced to take just the opposite values (close to $v_0$ if $x_i$ is true, close to $2v_0$ if $x_i$ is false). We will also construct constraints corresponding to the MINSAT clauses that are satisfied if and only if the MINSAT clauses are satisfied.

To be more specific, we construct four classes of constraints. Each of these constraints actually gets repeated several times in the construction of the reduced instance, meaning that if the constraint holds, then it holds several times. In this way, the constraints can be assigned varying importance weights.

A. First, we force $v_0$ to be strictly positive. To do so, we include the constraint:

$$v_0 \leq 0 \ .$$

(Recall that the goal is to *minimize* how many of these constraints are satisfied, which effectively means that it will be forced to fail so that $v_0 > 0$.) This constraint gets repeated $r_A$ times, as specified below.

B. Next, we force each $w_i$ and $w_i'$ to have a value roughly between $v_0$ and $2v_0$. To do so, we simply include constraints:

$$
\begin{aligned}
w_i &\leq 0.99v_0 \\
w_i &\geq 2.01v_0 \\
w_i' &\leq 0.99v_0 \\
w_i' &\geq 2.01v_0
\end{aligned}
$$

for each $i$. Each of these is repeated $r_B$ times.

C. Next, we add constraints that will effectively force (for each $i$) exactly one of $w_i$ and $w_i'$ to be close to $v_0$, and the other to be close to $2v_0$. These are the constraints:

$$
\begin{aligned}
w_i &\leq 1.99w_i' \\
w_i' &\leq 1.99w_i \ .
\end{aligned}
$$

In the optimal solution, we will see that exactly one of these two constraints will hold. These constraints each get repeated $r_C$ times.

D. Finally, we encode the actual clauses of the MinSat instance. A MinSat clause of the form $x_i \vee x_j$ becomes the constraint

$$0.8w_i' \leq w_j \ .$$

A MinSat clause of the form $\neg x_i \vee x_j$ becomes the constraint

$$0.8w_i \leq w_j \ .$$

A MinSat clause of the form $x_i \vee \neg x_j$ becomes the constraint

$$0.8w_i' \leq w_j' \ .$$

Finally, a MinSat clause of the form $\neg x_i \vee \neg x_j$ becomes the constraint

$$0.8w_i \leq w_j' \ .$$

Each of these is repeated once.

The value $K$ for the instance of Reweight that we constructed is denoted $K'$ (reserving $K$ for the corresponding value of the original MinSat instance). We let:

$$
\begin{aligned}
r_C &= K + 1 \\
K' &= K + nr_C \\
r_B &= K' + 1 \\
r_A &= K' + 1 \ .
\end{aligned}
$$

This completes the construction, which is clearly polynomial in all respects. We now need to argue that the MinSat instance is "yes" if and only if the reduced Reweight instance is also "yes".

Suppose that $x_1, \ldots, x_n$ satisfies at most $K$ of the MINSAT clauses. In this case, we can easily construct settings of the weights so that at most $K'$ of the constructed constraints are satisfied.

Let $v_0 = 1$ and let

$$w_i = \begin{cases} 1 & \text{if } x_i = 0 \\ 2 & \text{if } x_i = 1 \end{cases} ,$$

and let

$$w_i' = \begin{cases} 1 & \text{if } w_i = 2 \\ 2 & \text{if } w_i = 1 \end{cases} .$$

Then none of the constraints of types A and B is satisfied. Exactly half of the constraints of type C are satisfied, which means $n\, r_C$ constraints of type C are satisfied.

What about the constraints of type D? We claim that for each satisfied clause of the MINSAT solution, the corresponding type-D constraint is satisfied. If a clause of the form $x_i \vee x_j$ is satisfied, then $x_i = 1$ or $x_j = 1$, which means that $w_i = 2$ or $w_j = 2$, which means $w_i' = 1$ or $w_j = 2$, which means that $0.8w_i' \leq w_j$. Conversely, if $x_i \vee x_j$ is not satisfied, then $x_i = 0$ and $x_j = 0$, which means that $w_i = 1$ and $w_j = 1$, which means $w_i' = 2$ and $w_j = 1$, which means that $0.8w_i' \not\leq w_j$. (The arguments are the same when some of the variables appear negated in the clause.)

Thus, because at most $K$ of the clauses are satisfied, it follows that at most $K$ of the type-D constraints are satisfied. Therefore, the total number of satisfied constraints is at most $K + n\, r_C = K'$.

We now argue the other direction. Suppose that $v_0, w_1, \ldots, w_n, w_1', \ldots, w_n'$ satisfy at most $K'$ of the constructed constraints.

First of all, this means that

$$v_0 > 0$$

since $r_A > K'$. Also, since $r_B > K'$, this means that

$$0.99v_0 < w_i < 2.01v_0 \tag{5}$$

and

$$0.99v_0 < w_i' < 2.01v_0 . \tag{6}$$

Next, we claim that either

$$w_i \leq 1.99w_i' \tag{7}$$

or

$$w_i' \leq 1.99w_i . \tag{8}$$

Otherwise, if neither of these was true, then we would have

$$w_i > 1.99w_i' > (1.99)^2 w_i ,$$

which implies that $w_i < 0$. However, we have already established that $w_i > 0.99v_0 > 0$.

Further, we claim that at most one of (7) and (8) can be satisfied. We already have established that at least one constraint of each pair must be satisfied. If, in addition, both held for some pair, then the number of satisfied type-C constraints would be at least

$$(n-1)\, r_C + 2\, r_C = n\, r_C + r_C > K' .$$

So exactly one of each pair of type-C constraints is satisfied.

We next claim that for each $i$, exactly one of $w_i$ and $w_i'$ is in the interval $(0.99v_0, 1.1v_0)$ and the other is in $(1.9v_0, 2.01v_0)$. We know that either (7) and (8) is satisfied. Suppose $w_i' > 1.99w_i$. If $w_i' \leq 1.9v_0$ then

$$w_i < \frac{w_i'}{1.99} \leq \frac{1.9v_0}{1.99} < 0.99v_0 ,$$

a contradiction since we have already shown that $w_i > 0.99v_0$. Also, if $w_i \geq 1.1v_0$, then

$$w_i' > 1.99w_i \geq 1.99 \cdot 1.1v_0 > 2.01v_0 ,$$

7

again a contradiction since $w_i' < 2.01v_0$. So in this case, $w_i < 1.1v_0$ and $w_i' > 1.9v_0$. Moreover, because (5) and (6) hold, we have in this case that $w_i \in (0.99v_0, 1.1v_0)$ and $w_i' \in (1.9v_0, 2.01v_0)$. In this case, we assign the boolean variable $x_i$ the value 0. By a similar argument, if $w_i > 1.99w_i'$ then $w_i' \in (0.99v_0, 1.1v_0)$ and $w_i \in (1.9v_0, 2.01v_0)$. In this case, we assign the boolean variable $x_i$ the value 1.

We have established that exactly $n\, r_C$ type-C constraints are satisfied, and none of the type-A and type-B constraints is satisfied. Since at most $K'$ constraints are satisfied altogether, this means that at most $K$ type-D constraints are satisfied. We complete the reduction by showing that a type-D constraint is satisfied if and only if the corresponding MINSAT clause is satisfied (according to the assignment constructed above), which will mean that at most $K$ of these are satisfied.

Suppose $x_i \vee x_j$ is satisfied. Then $x_i = 1$ or $x_j = 1$, which means either $w_i > 1.9\, v_0$ or $w_j > 1.9\, v_0$, which means either $w_i' < 1.1\, v_0$ or $w_j > 1.9\, v_0$. We claim, in either case, that the constraint $0.8\, w_i' \le w_j$ is satisfied. For if $w_i' < 1.1\, v_0$, then because $w_j > 0.99\, v_0$, we have

$$0.8\, w_i' < 0.8 \cdot 1.1\, v_0 = 0.88\, v_0 < 0.99\, v_0 < w_j \ .$$

And if $w_j > 1.9\, v_0$ then because $w_i' < 2.01\, v_0$, we have

$$0.8\, w_i' < 0.8 \cdot 2.01\, v_0 = 1.608\, v_0 < 1.9\, v_0 < w_j \ .$$

Conversely, if $x_i \vee x_j$ is not satisfied then $x_i = 0$ and $x_j = 0$, which means that $w_i < 1.1\, v_0$ and $w_j < 1.1\, v_0$, which means that $w_i' > 1.9\, v_0$ and $w_j < 1.1\, v_0$, which means that

$$0.8\, w_i' > 0.8 \cdot 1.9\, v_0 = 1.52\, v_0 > 1.1\, v_0 > w_j \ ,$$

so the constraint $0.8\, w_i' \le w_j$ is not satisfied. $\qquad\square$

# 4 Constant-Class Algorithm

We now present an algorithm that finds an optimal solution for a non-uniform cost function in polynomial time when the number of classes $n$ is constant. Our algorithm takes as input a set of nonnegative real numbers $f_{ij}$ for $i = 1, \ldots, m$ and $j = 1, \ldots, n$, integers (labels) $y_i \in \{1, \ldots, n\}$, and a nonnegative cost function $c : \{1, \ldots, n\}^2 \to \mathbb{R}^+$. Assuming $n$ is a constant, in time polynomial in $m$ it will output a vector of weights $(w_1, \ldots, w_n)$ that minimizes (1).

Our algorithm is based on the observation that class $j$ will be predicted for instance $i$ if and only if $w_j/w_k > f_{ik}/f_{ij}$ for all $k \ne j$. Thus for a fixed $(j, k)$ pair, there are only $m$ values of $f_{ik}/f_{ij}$ that can affect the value of (1). We will call these values *breakpoints*. For each $(j, k)$ pair, one can easily compute all breakpoints, add in $-\infty$ and $+\infty$, sort them, and place them in an ordered set $B^{jk} = (-\infty, f_{1k}/f_{1j}, \ldots, f_{mk}/f_{mj}, +\infty)$. We use $B_\ell^{jk}$ to denote the $\ell$th element in $B^{jk}$. So a candidate range of values of $w_j/w_k$ is $\left(B_\ell^{jk}, B_{\ell+1}^{jk}\right]$.

We now define a *configuration* $C$ as a set of pairs of breakpoints across all $(j, k)$ pairs:

$$C = \bigcup_{j \in \{1, \ldots, n\}, k \ne j} \left\{ B_{\ell_{jk}}^{jk} < w_j/w_k \le B_{1+\ell_{jk}}^{jk} \right\} \ ,$$

for $\ell_{jk} \in \{1, \ldots, m+1\}$. We say that $C$ is *realizable* if there exists a set of nonnegative weights that satisfies all of its constraints. If no such set of weights exists, we say $C$ is *unrealizable*.

The idea of our algorithm is simple. It enumerates each configuration $C$ and then uses linear programming to test if $C$ is realizable[2]. If it is not, then we test the next configuration. If instead $C$ is realizable, then the weight vector returned by the linear programming algorithm is one of our candidate solutions to minimize (1).

---

[2] To handle the strict inequalities, we simply convert each constraint $a > b$ to $a \ge b + \epsilon$, constrain $1 \ge \epsilon \ge 0$, then maximize $\epsilon$ subject to the new constraints (note that we use the same $\epsilon$ for each strict constraint). If a solution is returned with $\epsilon > 0$ then we know the configuration is realizable. If no solution is found or if $\epsilon = 0$, then it is not realizable.

Our algorithm stores this weight vector with its cost and then checks the remaining configurations. Once all configurations have been checked, the algorithm returns the one with minimum cost.

Consider an optimal solution $\vec{w}^*$. Let $C^*$ be the configuration it satisfies. Any other positive weight vector $\vec{w}$ that also satisfies $C^*$ must also be optimal since it induces the same collection of predictions and hence has the same value of (1). Since our algorithm tries all configurations, it also tries $C^*$. Since $C^*$ is realizable, the algorithm finds a weight vector $\vec{w}$ satisfying it. Since $\vec{w}$ must be optimal, our algorithm returns an optimal weight vector.

Each $(j, k)$ pair has at most $m + 2$ breakpoints, which means that the number of $(B_\ell^{jk}, B_{\ell+1}^{jk})$ pairs for class pair $(j, k)$ is at most $m + 1$. So the number of configurations is at most

$$(m+1)^{2\binom{n}{2}} = (m+1)^{n^2 - n} \ ,$$

which is polynomial for constant $n$. (We can also combine $B^{jk}$ with $B^{kj}$ before sorting, which would reduce the number of configurations to $(2m+1)^{\binom{n}{2}}$.) Further, it takes polynomial time to test each configuration for realizability via linear programming and it takes polynomial time to compute the cost of each candidate solution. Therefore this algorithm takes polynomial time.

**Theorem 2.** *There exists an algorithm to produce an optimal solution to* (1) *that runs in polynomial time when the number of classes $n$ is constant.*

# 5 Heuristics

For even modest values of $n$ the time complexity of the constant-class algorithm in the previous section is not practical. For this reason, we present several alternative heuristics to solve the reoptimization problem, including several new mathematical programming formulations. First, we reformulate the objective function (1) as a relaxed integer linear program. We also give formulations as a sum of linear fractional functions (SOLFF) as well as a quadratic program. Besides these formulations, we describe a tree-based heuristic algorithm approach, MetaClass. Finally, (in Section 6) we present experimental results from these formulations. We give evidence that, under non-uniform costs, the objective function landscape for this problem is highly discontinuous and thus more amenable to global optimization methods such as genetic algorithms and margin maximization methods.

## 5.1 Mathematical Programming Formulations

### 5.1.1 Relaxed Integer Linear Program

We start by reformulating (1) as follows:

$$\underset{\vec{w}}{\text{minimize}} \left\{ \sum_{j=1}^n \sum_{k=1}^m c(j,k) \sum_{x_i \in C_j} I_{i,k} \right\} \ , \tag{9}$$

where $C_j \subseteq S$ is the set of instances of class $j$, $c(j,k)$ is the cost of misclassifying an example from class $j$ as $k$, and

$$I_{i,k} = \begin{cases} 1 & \text{if } w_k f_k(x_i) \geq w_\ell f_\ell(x_i), \ell \neq k \\ 0 & \text{otherwise} \end{cases} \ .$$

Recall that $f_k(x_i)$ is the base learner's confidence that example $x_i$ belongs to class $k$. Also, we assume $c(j, j) = 0$ for all classes $j$. Formalizing this as a constrained optimization problem, we make $\vec{I}$ a free

variable in (9) and minimize it subject to

$$I_{i,j} w_j f_j(x_i) = I_{i,j} \max_{1 \le k \le m}\{w_k f_k(x_i)\} \tag{10}$$

$$\sum_{j=1}^{m} I_{i,j} = 1 \tag{11}$$

$$I_{i,j} \in \{0,1\} \tag{12}$$

$$w_j \ge 0 \tag{13}$$

where each constraint holds for all $i \in \{1, \ldots, n\}$ and $j \in \{1, \ldots, m\}$. Equation (10) allows only the class that has the max value of $w_k f_k(x_i)$ to be indicated by $\vec{I}$ to be the predicted class of example $x_i$ and (11) forces exactly one class to be predicted per example $x_i$. We can change the optimization problem in two ways to get an equivalent problem. First, we change the "=" in (10) to "$\ge$". Second, we can relax (12) to be $I_{i,j} \in [0,1]$.

Note that (10) (even when amended with "$\ge$") will only be satisfied if $I_{i,j} = 0$ for all $C_j$ that don't maximize the RHS of (10). Thus, so long as we never have $w_k f_k(x_i) = w_{k'} f_{k'}(x_i)$ for some $k \ne k'$, the relaxation is equivalent to the original problem. Further, even if there is such a tie for classes $k$ and $k'$, it will not be an issue if the corresponding entries in the cost matrix are different, since an optimal solution will set $I_{i,k} = 1$ and $I_{i,k'} = 0$ if $c(j,k) < c(j,k')$. The potential problem of both $w_k f_k(x_i) = w_{k'} f_{k'}(x_i)$ and $c(j,k) = c(j,k')$ is fixed by (after optimizing) checking for any $I_{i,k} \notin \{0,1\}$ and arbitrarily choosing one to be 1 and the rest 0. Note that since there is a tie in this case, the prediction can go either way and the weight vector $\vec{w}$ returned is still valid.

Everything except (10) is linear. We now reformulate it. First, for each $i \in \{1, \ldots, n\}$, we substitute $\gamma_i$ for $\max_{1 \le k \le m}\{w_k f_k(x_i)\}$:

$$I_{i,j} w_j f_j(x_i) \ge \gamma_i I_{i,j} \tag{14}$$

$$w_k f_k(x_i) \le \gamma_i \ , \tag{15}$$

for all $i \in \{1, \ldots, n\}$ and $j, k \in \{1, \ldots, m\}$ where each $\gamma_i$ is a new variable. Obviously (15) is a linear constraint, but (14) is not even quasiconvex (Boyd & Vandenberghe, 2004). The complexity of this optimization problem motivates us to reformulate it a bit further.

Let us assume that $f_k(x_i) \in (0,1]$ (e.g. if $f_k(\cdot)$ are probability estimates from naïve Bayes or logistic regression). Now we can optimize (9) (again with $\vec{I}$ as a free variable) subject to:

$$\gamma_i - w_j f_j(x_i) + I_{i,j} \le 1 \tag{16}$$

$$\gamma_i \ge w_j f_j(x_i) \tag{17}$$

$$\sum_{j=1}^{m} I_{i,j} = 1 \tag{18}$$

$$I_{i,j} \in \{0,1\} \tag{19}$$

$$w_j \ge 0 \tag{20}$$

for all $i \in \{1, \ldots, n\}$ and $j \in \{1, \ldots, m\}$.

So long[3] as $w_j f_j(x_i) \in (0,1]$ and $I_{i,j} \in \{0,1\}$ for all $i \in \{1, \ldots, n\}$ and $j \in \{1, \ldots, m\}$, this is another equivalent optimization problem, this time a $\{0,1\}$ integer linear program. Unfortunately, we cannot relax (19) to $I_{i,j} \in [0,1]$ as we did before to get an equivalent problem. But we still use the relaxation as a linear programming heuristic. To help avoid overfitting, we also add a linear regularization term to (9):

$$\underset{\vec{w}, \vec{I}}{\text{minimize}} \left\{ \sum_{j=1}^{n} \sum_{k=1}^{m} c(j,k) \sum_{x_i \in C_j} I_{i,k} + \eta \|\vec{w} - \vec{1}\|_1 \right\} \tag{21}$$

where $\| \cdot \|_1$ is the 1-norm, $\vec{1}$ is the all-1s vector, and $\eta$ is a parameter. This regularizer penalizes large deviations from the original classifier.

---

[3]We can ensure this happens by bounding each $w_j$ appropriately.

### 5.1.2 Sum of Linear Fractional Functions Formulation

Another formulation comes from changing how predictions are made from deterministic to probabilistic. In this prediction model, given a new example $x$ to predict on, first compute $w_j f_j(x)$ for each $j \in \{1, \ldots, n\}$. Then predict class $j$ for example $x$ with probability

$$\frac{w_j f_j(x)}{\sum_{k=1}^n w_k f_k(x)} \ .$$

Assuming that an instance $x$ is uniformly drawn at random from the training set, the expected cost of this predictor is

$$\sum_{j=1}^n \sum_{k=1}^m c(j,k) \sum_{x_i \in C_j} \varphi(i,j) \ , \tag{22}$$

where

$$\varphi(i,j) = \frac{w_j f_j(x_i)}{\sum_{\ell=1}^n w_\ell f_\ell(x_i)}$$

subject to $w_j \geq 0$ for all $j \in \{1, \ldots, n\}$. We now have eliminated the variables $I_{i,j}$ and their integer constraints. However, we now have a nonlinear objective function in (22). Each individual term of the summation of (22) is a *linear fractional function*, which is quasiconvex and quasiconcave, and thus it is efficiently solvable optimally. However, the *sum of linear fractional functions* (SOLFF) problem is known to be hard (Matsui, 1996) and existing algorithms for this problem are inappropriate in solving (22) (they either restrict to few terms in the summation or to low-dimensional vectors). Instead, we apply a genetic algorithm to directly optimize (22).

### 5.1.3 Quadratic Programming Formulation

Convex programming is a special case of nonlinear programming in which the objective function and the inequality constraint functions are convex and the equality constraint functions are affine. The theory of convex programming is well-established (Rockafellar, 1970; Stoer & Witzgall, 1996). For a convex program, a local optimum is the global optimum and there are well-studied efficient algorithms to find such a global optimum.

We tried several quadratic programming methods based on the idea of margin maximization in support vector machines. We found from our experiments that the $\nu$-SVM-like formulation similar to that of Schölkopf and Smola (2001) gave the strongest result of the SVM-like formulations we tested.

$$\underset{\vec{w},\vec{b},\vec{\zeta},\rho}{\text{minimize}} \qquad \frac{1}{2}\|\vec{w}\|^2 + \frac{1}{m}\sum_{i=1}^m \zeta_i - \nu\rho \tag{23}$$

$$\text{s.t.} \quad w_j f_j(x_i) + b_j \leq w_{y_i} f_{y_i}(x_i) + b_{y_i} + \zeta_i - \rho \quad \forall i, \forall j \neq y_i \tag{24}$$

$$\vec{\ell}_w \leq \vec{w} \leq \vec{u}_w \tag{25}$$

$$\vec{\ell}_b \leq \vec{b} \leq \vec{u}_b \tag{26}$$

$$\vec{0} \leq \vec{\zeta} \tag{27}$$

$$0 \leq \rho \tag{28}$$

where $\vec{w} = (w_1, \ldots, w_n)$ is our weight vector and $\vec{b} = (b_1, \ldots, b_n)$ is a reweighting *offset*. The vector $\vec{\zeta} = (\zeta_1, \ldots, \zeta_m)$ serves as set of slack variables and $\rho$ as the margin with $\nu$ as a parameter. Furthermore, $\vec{\ell}_w, \vec{u}_w$ are the lower and upper bounds of $\vec{w}$ and $\vec{\ell}_b, \vec{u}_b$ are bounds for $\vec{b}$, all tunable parameters. In order to capture non-uniform costs we replace $\zeta_i$ with $c_i \zeta_i$ where $c_i = \max_{j=1,\ldots,m}\{c(y_i, j)\}$.

In our experiments, $\nu$ was fixed to be 0.1. The lower and upper bounds on $w$, $\vec{\ell}_w, \vec{u}_w$, were set to 0 and 1 respectively. We also tried several parameters for the offset, but little difference was observed. Thus, our experimental results use no offset (the lower and upper bounds on $\vec{b}$ were set to 0).

## 5.2 The MetaClass Heuristic Algorithm

We now present a new algorithm that we call MetaClass (Algorithm 1). This algorithm is similar to that of Lachiche and Flach (2003) in that we reduce the multi-class problem to a series of two-class problems. However, we take what can be considered a top-down approach while the algorithm of Lachiche and Flach (2003) can be considered bottom-up. Moreover, MetaClass has a lower time complexity. The output of the algorithm is a decision tree with each internal node labeled by two *metaclasses* and a threshold value. Each leaf node is labeled by one of the classes in the original problem. At the root, the set of all classes is divided into two metaclasses, $\mathcal{C}_1$ and $\mathcal{C}_2$. The criterion for this split may be based on any statistical measure. For simplicity, experiments were performed by splitting classes so that each metaclass would have roughly the same number of training examples by simply sorting classes according to the number of examples and alternately partitioning them into each metaclass. For each metaclass, our algorithm defines confidence functions $g_1(x_i)$ and $g_2(x_i)$ for each instance $x_i$, which are simply the sum of the confidences of the classes in $\mathcal{C}_1$ and $\mathcal{C}_2$, respectively. The ratio $G(x_i) = g_1(x_i)/g_2(x_i)$ is used to find a threshold $\theta$. We find $\theta$ by sorting the instances according to $G(x_i)$ and choose a threshold that minimizes the error rate (or cost). The threshold will be the average of $G(x_i)$ and $G(x_{i+1})$ for some instance $x_i$. Among equivalent thresholds (thresholds that induce the same error rate) we choose the median. We recursively perform this procedure on the two metaclasses until there is only a single class, at which point a leaf is formed.

The situation for non-uniform costs is slightly more complicated since it is not clear *which* class among those in metaclass an example is misclassified as. Recall that our cost function $c(y, \hat{y})$ represents the cost of misclassifying an instance $x$ of class $y$ as class $\hat{y}$. However, in this case we need a cost function to quantify the cost of misclassifying an example into a *set* of classes. Formally, we need a function $c' : \mathcal{C} \times 2^{\mathcal{C}} \to \mathbb{R}^+$. There are numerous natural extensions from $c$ to $c'$. For our experiments, $c'$ represents the average cost of misclassifying instances into metaclasses in $\mathcal{C}_1$ and $\mathcal{C}_2$. More formally, if $\mathcal{C}' \subseteq \mathcal{C}$, we define $c'(y, \mathcal{C}')$ to be 0 if $y \in \mathcal{C}'$ (that is, $x$'s true label class is in the metaclass) and

$$\frac{1}{|\mathcal{C}'|} \sum_{j \in \mathcal{C}'} c(y, j)$$

otherwise. The MetaClass algorithm is presented as Algorithm 1.

> **Input**  : A set of instances, $S = \{x_1, \ldots, x_m\}$; a set of classes, $\mathcal{C} = \{1, \ldots, n\}$;
> a learned confidence function $f : S \times \mathcal{C} \to \mathbb{R}^+$ and a tree node $T$
>
> **Output** : A decision tree with associated weights.
>
> **1** **if** $|\mathcal{C}| = 1$ **then**
> **2**    Stop, create a decision node and predict the class in $\mathcal{C}$
> **3** **end**
> **4** Split $\mathcal{C}$ into two metaclasses $\mathcal{C}_1, \mathcal{C}_2$ such that each metaclass has about an equal number of classes.
> **5** **foreach** *Instance* $x_i \in S$ **do**
> **6**    $g_1(x_i) = \sum_{j \in \mathcal{C}_1} f_j(x_i)$
> **7**    $g_2(x_i) = \sum_{j \in \mathcal{C}_2} f_j(x_i)$
> **8**    $G(x_i) = g_1(x_i)/g_2(x_i)$
> **9** **end**
> **10** Sort instances according to $G$
> **11** Select a threshold $\theta$ that minimizes
>
> $$\sum_{i=1}^{m} c'(y_i, M_\theta(x_i))$$
>
> where
> $$M_\theta(x_i) = \begin{cases} \mathcal{C}_1 & \text{if } \theta \geq G(x_i) \\ \mathcal{C}_2 & \text{otherwise} \end{cases}$$
>
> **12** Label $T$ with $\theta, \mathcal{C}_1, \mathcal{C}_2$
> **13** Create two children of $T$: $T_{\text{left}}, T_{\text{right}}$
> **14** Split $S$ into two sets,
>
> $$\begin{aligned} S_1 &= \{x_i \in S \mid M_\theta(x_i) = \mathcal{C}_1\} \\ S_2 &= \{x_i \in S \mid M_\theta(x_i) = \mathcal{C}_2\} \end{aligned}$$
>
> **15** Recursively perform this procedure on $S_1, \mathcal{C}_1, T_{\text{left}}$ and $S_2, \mathcal{C}_2, T_{\text{right}}$

*Algorithm 1: MetaClass*

Figure 1 gives an example of a tree built by the MetaClass algorithm on the UCI (Blake & Merz, 2005) data set Nursery, a 5-class data set. At the root, the classes are divided into two metaclasses, each with about the same number of training examples represented in their respective classes. In this case, the threshold $\theta = 0.8169$ favors the sum of confidences in metaclass $\mathcal{C}_1 = \{4, 3\}$ as an optimal weight.

Predictions for a new example $x$ are made as follows. Starting at the root node, we traverse the tree towards a leaf. At each node $T$ we compute the sum of confidences of $x$ with respect to each associated metaclass. We traverse left or right down the tree depending on whether $g_1(x)/g_2(x) \geq \theta$. When a leaf is reached, a final class prediction is made.

The number of nodes created by MetaClass is $\Theta(n)$, where $n$ is the number of classes. Since the split into two metaclasses ensures each has an equal number of classes, MetaClass results in a $\log(n)$-depth tree. At each level, the most complex step is sorting at most $m$ instances according to the confidence ratio. Thus, the overall time complexity is bounded by $\mathcal{O}(m \log(m) \log(n))$. This represents a speedup to the algorithm of Lachiche and Flach (2003), which requires $\Theta(nm \log(m))$ time. Classification is also efficient. At each node we compute a sum over an exponentially shrinking number of classes. The overall number of operations is thus

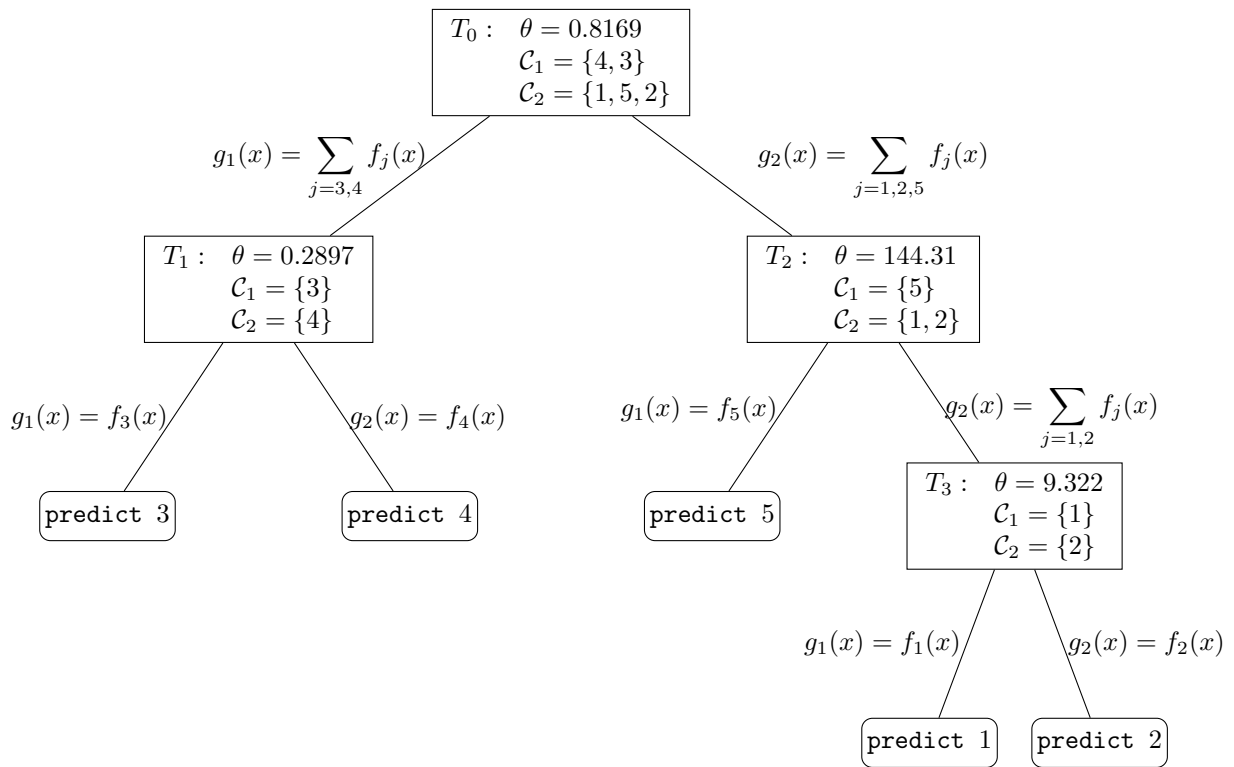$$\sum_{i=0}^{\log(n)-1} \frac{n}{2^i} \ ,$$

Figure 1: Example run of MetaClass on Nursery, a 5-class problem.

which is linear in the number of classes: $\Theta(n)$. This matches the time complexity of Lachiche and Flach's with respect to classification.

# 6   Experimental Results

The following experiments were performed on 24 UCI data sets (Blake & Merz, 2005), using Weka's naïve Bayes (Witten et al., 2005) as the baseline classifier and Matlab's optimization functions for reoptimization. We ran experiments evaluating improvements both in classification accuracy and under non-uniform costs. We used 10-fold cross validation for the error rate experiments. For the cost experiments, 10-fold cross validations were performed on 10 different cost matrices for each data set. Costs were integer values between 1 and 10, assigned uniformly at random. Costs on the diagonal were set to zero.

We evaluated the performance of each heuristic both in how well it was able to optimize the given cost function as well as how well the resultant classifier was able to generalize. To this end, we present performance on the training set as well as the testing set. Results for the optimization performance can be found in Table 1 while results for generalization can be found in Table 2. In each table, results for the uniform (error rate) and non-uniform cost models can be found in subtables (a) and (b) respectively. For uniform cost, the average over the 10 folds is reported while for non-uniform costs, the average over all folds and cost matrices is reported. Thus, the values for non-uniform costs represent the average cost over all 100 experiments per data set, per algorithm.

In all tables, for each data set, $n$ denotes the number of classes while $m$ denotes the total number of instances in each data set. The first column is the performance of our baseline classifier. For comparison, we have included results of our implementation for the algorithm of Lachiche and Flach (2003) on this baseline classifier (labeled "LF" in the tables). The results of the experiments on our heuristics can be found in the last five columns of each table. Here, "MC" is the MetaClass algorithm (Algorithm 1). "LP" is a linear programming solver (MOSEK ApS, 2005) on (21) with $\eta = 10^{-6}$. "GA1" is the sum of linear fractional functions formulation (22) using a genetic algorithm. "GA2" is a genetic algorithm optimization of (1). Both experiments used the GA implementation from Abramson (2005). Parameters for both used the default Matlab settings with a population size of 20, a maximum of 200 generations and a crossover fraction of 0.8. The algorithm terminates if no change is observed in 100 continuous rounds. In addition, the mutation function of Abramson (2005) is guaranteed to only generate feasible solutions (in our case, all weights must be nonnegative). Upon termination, a direct pattern search is performed using the best solution from the GA. The final column is the quadratic program (QP) from Section 5.1.3 using the CVX software package (Grant et al., 2006).

For all columns, values in *italics* indicate a worse performance than the baseline. Entries in **bold** indicate a significant difference from the baseline with at least a 95% confidence according to a Student's $t$ method. The overall best classifier(s) for each data set are underlined. Finally, the number of wins/losses (i.e. better/worse than the baseline) and those that are significant are summarized in each column. A win indicates that the heuristic improved upon the baseline while a loss indicates that the heuristic performed worse than the baseline.

To measure how well each algorithm is able to optimize, it would be ideal to compare it against an optimal solution. To do this, we could compute the optimal reweighting function either by solving the integer linear program of (9) or by running the constant-class algorithm presented in Section 4. Unfortunately, neither of these options is feasible. Even for $n = 3$ class data sets, the large number of instances forces an unmanageable number of variables in the ILP formulation (likewise for the constant-class algorithm). Even the best ILP packages cannot handle such large numbers of variables. In lieu of such comparisons, we instead report the performance of each heuristic on the training set (Table 1). That is, we report how well the heuristic was able to optimize the given problem. For data sets with $n = 2$, the algorithm of Lachiche and Flach (2003) and MetaClass are optimal since they simplify to the two-class ROC threshold problem. However, in some instances, their performances differ slightly. This is due to rounding errors and how each implementation deals with borderline examples, thus the differences are negligible (in particular, the largest discrepancies correspond to at most 11 instances over all 10 folds being classified differently by the two algorithms).

Table 1: Optimization Performance. The two tables present the performance rates of each heuristic. *Italicized* entries indicate the heuristic performed worse than the baseline (naïve Bayes) while non-italicized entries indicate an improvement or tie. A **bold** entry indicates a statistically significant difference from the baseline according to a student-$t$ test with a 95% confidence level. For each data set, the overall best performing heuristic(s) are underlined.

(a) Error Rates. The values in this table represent performance under a uniform cost model.

| Data Set | $n$ | $m$ | Naïve Bayes | LF | MC | LP Eq. (21) | GA1 Eq. (22) | GA2 Eq. (1) | QP |
|---|---|---|---|---|---|---|---|---|---|
| Audiology | 24 | 226 | 0.2143 | **0.0840** | **0.1425** | **0.1705** | **0.1145** | **0.1637** | **0.1754** |
| Bridges 2 (material) | 3 | 107 | 0.0757 | **0.0602** | **0.0405** | 0.0540 | *0.0882* | 0.0519 | 0.0570 |
| Bridges 2 (rel-l) | 3 | 107 | 0.1173 | **0.0778** | **0.0737** | 0.1069 | *0.1245* | 0.0747 | 0.1038 |
| Bridges 2 (span) | 3 | 107 | 0.0560 | 0.0508 | **0.0342** | *0.0570* | *0.0633* | **0.0321** | *0.0570* |
| Bridges 2 (type) | 6 | 107 | 0.1277 | **0.0904** | 0.1059 | 0.1277 | *0.1381* | 0.1090 | 0.1080 |
| Bridges 2 (t-or-d) | 2 | 107 | 0.0176 | **0.0124** | 0.0135 | *0.0218* | *0.0426* | **0.0124** | *0.0290* |
| Car | 4 | 1728 | 0.1291 | **0.0914** | 0.1024 | 0.1153 | *0.1639* | **0.0914** | 0.1158 |
| Post-Op | 3 | 1473 | 0.4830 | **0.4525** | 0.4666 | 0.4830 | 0.4805 | 0.4632 | 0.4693 |
| Horse-colic (code) | 3 | 368 | 0.2276 | **0.2173** | 0.2237 | 0.2276 | *0.2303* | **0.2149** | *0.2339* |
| Horse-colic (surgical) | 2 | 368 | 0.1657 | **0.1388** | **0.1370** | 0.1657 | 0.1497 | 0.1376 | 0.1515 |
| Horse-colic (site) | 63 | 368 | 0.2970 | **0.2161** | 0.2623 | **0.2844** | 0.2744 | 0.2412 | 0.2820 |
| Horse-colic (subtype) | 2 | 368 | 0.0000 | *0.0003* | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 |
| Horse-colic (type) | 8 | 368 | 0.0105 | **0.0048** | **0.0045** | 0.0090 | **0.0057** | 0.0072 | 0.0072 |
| Credit | 2 | 1000 | 0.2288 | **0.2245** | 0.2255 | 0.2288 | *0.2594* | **0.2245** | *0.2371* |
| Dermatology | 6 | 366 | 0.0094 | **0.0042** | 0.0048 | 0.0078 | **0.0072** | 0.0051 | 0.0069 |
| Ecoli | 8 | 336 | 0.1134 | **0.0886** | 0.0962 | *0.1150* | **0.1031** | 0.0988 | *0.1180* |
| Flags | 8 | 194 | 0.2147 | **0.1523** | 0.1597 | **0.2090** | 0.1626 | 0.1718 | 0.2044 |
| Glass | 7 | 214 | 0.4522 | **0.3073** | 0.3172 | *0.4532* | 0.3525 | 0.4076 | 0.4449 |
| Mushroom | 2 | 8124 | 0.0419 | **0.0176** | 0.0177 | 0.0367 | **0.0188** | 0.0176 | 0.0266 |
| Image Segmentation | 7 | 2310 | 0.1948 | 0.1642 | **0.1361** | 0.1946 | **0.1215** | 0.1574 | *0.1982* |
| Solar Flare (common) | 8 | 1066 | 0.2189 | **0.1695** | **0.1671** | *0.2198* | 0.1707 | 0.1810 | 0.2011 |
| Solar Flare (moderate) | 6 | 1066 | 0.0660 | **0.0337** | **0.0332** | 0.0660 | **0.0337** | 0.0439 | 0.0469 |
| Solar Flare (severe) | 3 | 1066 | 0.0277 | **0.0045** | **0.0044** | **0.0223** | **0.0046** | 0.0182 | 0.0124 |
| Vote | 2 | 435 | 0.0968 | **0.0878** | 0.0883 | 0.0968 | **0.0906** | 0.0888 | *0.1031* |
| Win/Loss | | | – | 23/1 | 23/0 | 11/5 | 15/8 | 23/0 | 16/7 |
| Significant Win/Loss | | | – | 21/0 | 22/0 | 7/1 | 14/5 | 22/0 | 13/5 |

(b) Non-uniform Costs. The values in this table represent performance under a non-uniform cost model.

| Data Set | $n$ | $m$ | Naïve Bayes | LF | MC | LP Eq. (21) | GA1 Eq. (22) | GA2 Eq. (1) | QP |
|---|---|---|---|---|---|---|---|---|---|
| Audiology | 24 | 226 | 1.2300 | **0.6019** | **0.7500** | **0.9663** | **0.7969** | **0.7173** | **0.9345** |
| Bridges 2 (material) | 3 | 107 | 0.3943 | **0.3408** | **0.1755** | 0.2916 | *0.4556* | 0.2565 | 0.3030 |
| Bridges 2 (rel-l) | 3 | 107 | 0.6772 | **0.4741** | 0.3759 | 0.6149 | *1.1963* | **0.3675** | 0.5873 |
| Bridges 2 (span) | 3 | 107 | 0.3290 | 0.3057 | **0.1999** | *0.3639* | *0.6190* | 0.2065 | *0.3740* |
| Bridges 2 (type) | 6 | 107 | 0.6728 | 0.6261 | **0.4416** | 0.6549 | *1.0157* | 0.4936 | 0.6077 |
| Bridges 2 (t-or-d) | 2 | 107 | 0.0966 | **0.0554** | **0.0554** | *0.1297* | *0.4614* | 0.0604 | *0.1334* |
| Car | 4 | 1728 | 0.7481 | **0.5314** | 0.5924 | 0.6863 | *1.6442* | **0.5103** | 0.6577 |
| Post-Op | 3 | 1473 | 2.8495 | **2.5125** | 2.6544 | *2.9183* | *3.6382* | 2.5371 | 2.7441 |
| Horse-colic (code) | 3 | 368 | 1.2740 | **1.1439** | 1.1797 | 1.2730 | *1.4168* | 1.1571 | *1.2972* |
| Horse-colic (surgical) | 2 | 368 | 1.0887 | **0.8043** | **0.7882** | *1.0926* | *1.2983* | 0.8309 | 0.9050 |
| Horse-colic (site) | 63 | 368 | 1.6023 | 1.4509 | 1.3801 | 1.5358 | *1.8110* | **1.1268** | 1.5908 |
| Horse-colic (subtype) | 2 | 368 | 0.0000 | *0.0017* | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 |
| Horse-colic (type) | 8 | 368 | 0.0597 | **0.0208** | **0.0197** | 0.0485 | **0.0325** | 0.0418 | 0.0397 |
| Credit | 2 | 1000 | 1.2160 | **0.9356** | 0.9372 | *1.2943* | *2.1381* | 0.9500 | **1.0035** |
| Dermatology | 6 | 366 | 0.0619 | *0.1426* | **0.0255** | 0.0520 | *0.0729* | 0.0325 | 0.0573 |
| Ecoli | 8 | 336 | 0.6024 | 0.5646 | 0.4768 | *0.6381* | *0.9202* | **0.4738** | *0.6535* |
| Flags | 8 | 194 | 1.2387 | **0.9360** | **0.8617** | 1.2092 | 1.0459 | 0.8970 | 1.1554 |
| Glass | 7 | 214 | 2.7112 | **1.6253** | **1.7130** | *2.7215* | *3.1616* | 2.2020 | 2.6379 |
| Mushroom | 2 | 8124 | 0.1926 | **0.0962** | **0.0962** | 0.1746 | 0.1247 | 0.0995 | 0.1348 |
| Image Segmentation | 7 | 2310 | 1.0715 | *1.2418* | **0.6963** | 1.0705 | 1.0546 | 0.8422 | *1.0791* |
| Solar Flare (common) | 8 | 1066 | 1.2229 | **0.8574** | 0.9157 | *1.2277* | **1.1494** | 0.8971 | **1.1500** |
| Solar Flare (moderate) | 6 | 1066 | 0.4174 | **0.2128** | **0.2076** | 0.4164 | **0.2709** | 0.2412 | 0.2935 |
| Solar Flare (severe) | 3 | 1066 | 0.1657 | **0.0260** | **0.0249** | 0.1468 | 0.0663 | 0.0948 | 0.0732 |
| Vote | 2 | 435 | 0.4527 | **0.3314** | 0.3335 | 0.4527 | *0.5320* | 0.4148 | *0.4539* |
| Win/Loss | | | – | 21/3 | 23/0 | 14/8 | 8/15 | 23/0 | 17/6 |
| Significant Win/Loss | | | – | 19/1 | 23/0 | 7/5 | 7/14 | 22/0 | 12/2 |

## 6.1 Generalization

Another measure of learning algorithms is how well they generalize beyond the training set. Unfortunately, none of the heuristics has any generalization guarantees other than what can be derived via standard PAC analysis. All algorithms could be adapted to use a validation set as yet another heuristic to improve generalization and prevent over-fitting. However, none of these techniques has any provable guarantees. In general, proving such guarantees is quite difficult. Results for error rates (uniform costs) and non-uniform costs are reported in Tables 2(a) and 2(b), respectively.

## 6.2 Comparison and Analysis

The significant wins and loses presented in Tables 1 and 2 provide a good reference to how well each algorithm performs individually with respect to the baseline as well as how much of an improvement each algorithm was able to achieve. As far as the optimization task (performance on the training set) is concerned, LF, MC and GA2 all showed statistically significant improvement over the baseline classifier in all or almost all data sets for both classification error and under non-uniform costs. These three methods also had superior performance in both settings with respect to the generalization task (performance on the test set).

For a more rigorous analysis, we performed a Friedman rank test as suggested by Demšar (2006) to compare the relative performances of multiple algorithms across multiple data sets. The Friedman test is a non-parametric statistical test similar to the parametric repeated measures ANOVA. It is especially applicable in this case since we cannot make any normality assumptions over dissimilar data sets from dissimilar domains. In each table (uniform/non-uniform cost on the train/test sets) for each data set, algorithms are ranked 1–7 based on their performance. A mean rank is taken for each algorithm and used to determine if some subset of algorithms is statistically significantly better than others. With a 95% confidence level, we can reject the hypothesis that all algorithms are equivalent for all four scenarios. That is, we can conclude that some set of algorithms outperform other methods with a high degree of certainty. Specific rankings can be found in Table 3.

We also performed a post-hoc analysis using a Tukey-Kramer pairwise test (Hochberg & Tamhane, 1987). The results of this analysis can be found in Figure 2. For each table, the pairwise Tukey-Kramer test induces a partial order among the heuristics such that a relation from algorithm $A$ to algorithm $B$ (indicated by a directed edge $A \rightarrow B$) implies that $A$ statistically significantly outperforms $B$. Transitive relations are implicit so that if $A$ outperforms $B$ and $B$ outperforms $C$, then $A$ also outperforms $C$. If no relation exists between two heuristics, then there is no statistically significant difference between them.

Following the post-hoc analysis, we can conclude that both greedy methods (MetaClass and the algorithm of Lachiche and Flach (2003)) and the second genetic algorithm formulation (GA2) are consistently the top three performers. Only in the case of uniform cost in the generalization setting does the first genetic algorithm outperform GA2. However they are not statistically significantly different. Moreover, with respect to the optimization task in both cost settings (as well as the non-uniform cost generalization), all three algorithms are statistically significantly better than all other methods.

A case can now be made for preferring the simple greedy methods over the genetic algorithm. First, in all instances, one of the greedy methods either outperforms or is not statistically significantly worse than GA2. In fact, it is never the case that GA2 statistically significantly outperforms either of the two greedy methods. Second, MetaClass and the algorithm of Lachiche and Flach were both light on computational resources. Each fold required only a few seconds to execute for both algorithms. In contrast, the genetic algorithm (as well as all other methods) were extremely computation-intensive, requiring several minutes of execution time. In addition, the genetic algorithm required a large number of generations; terminating the GA early resulted in poor performance. Thus, the greedy methods that used local decisions were just as good, and in most cases performed better, than the more sophisticated methods.

We can also make a similar, though less strong, case for preferring MetaClass over the algorithm of Lachiche and Flach. Though both performed very well, in most cases their performance was statistically equivalent. However, for the non-uniform cost optimization task, MetaClass statistically significantly outperformed LF. Though the time complexity of MetaClass is faster than the algorithm of Lachiche and Flach,
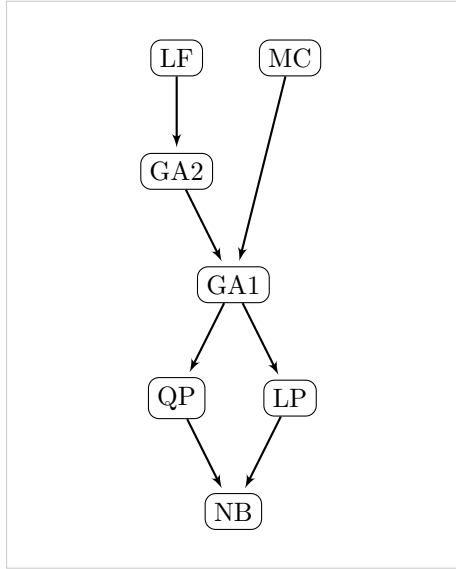
Table 2: Generalization Performance. The two tables present the test error rates of each heuristic. *Italicized* entries indicate the heuristic performed worse than the baseline (naïve Bayes) while non-italicized entries indicate an improvement or tie. A **bold** entry indicates a statistically significant difference from the baseline according to a student-*t* test with a 95% confidence level. For each data set, the overall best performing heuristic(s) are underlined.

(a) Error Rates. The values in this table represent performance under a uniform cost model.
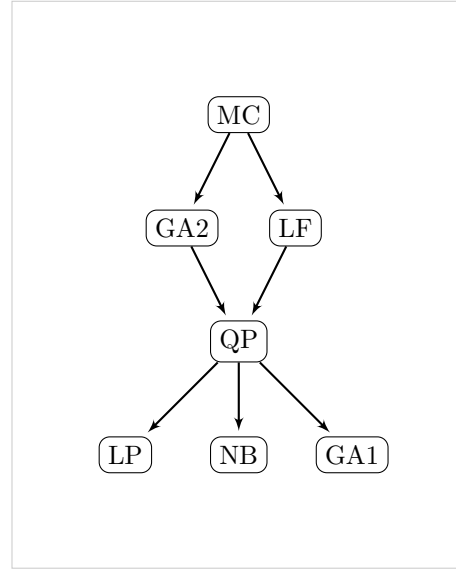
| Data Set | $n$ | $m$ | Naïve Bayes | LF | MC | LP Eq. (21) | GA1 Eq. (22) | GA2 Eq. (1) | QP |
|---|---|---|---|---|---|---|---|---|---|
| Audiology | 24 | 226 | 0.3094 | <u>0.2654</u> | *0.3227* | 0.2869 | 0.2826 | 0.2871 | 0.2737 |
| Bridges 2 (material) | 3 | 107 | <u>0.1581</u> | ***0.2427*** | ***0.2527*** | ***0.2709*** | *0.1763* | ***0.2136*** | ***0.2436*** |
| Bridges 2 (rel-l) | 3 | 107 | 0.3163 | *0.3527* | *0.3536* | 0.3081 | *0.3354* | *0.3454* | <u>0.3072</u> |
| Bridges 2 (span) | 3 | 107 | <u>0.2227</u> | *0.3290* | ***0.3290*** | *0.2427* | *0.4018* | *0.2799* | *0.2518* |
| Bridges 2 (type) | 6 | 107 | <u>0.4563</u> | ***0.5499*** | 0.4854 | 0.4654 | *0.4663* | 0.4654 | *0.4936* |
| Bridges 2 (t-or-d) | 2 | 107 | 0.1754 | ***0.2418*** | 0.2045 | 0.1754 | *0.1936* | *0.2499* | <u>0.1654</u> |
| Car | 4 | 1728 | 0.1464 | <u>**0.1191**</u> | **0.1249** | **0.1336** | *0.1723* | **0.1209** | **0.1331** |
| Post-Op | 3 | 1473 | 0.4948 | <u>0.4745</u> | 0.4833 | 0.4948 | *0.4989* | 0.4908 | 0.4853 |
| Horse-colic (code) | 3 | 368 | 0.3172 | 0.3148 | 0.3125 | 0.3172 | <u>0.2930</u> | 0.3120 | *0.3174* |
| Horse-colic (surgical) | 2 | 368 | 0.2089 | <u>**0.1737**</u> | **0.1764** | 0.2089 | **0.1791** | **0.1791** | 0.1980 |
| Horse-colic (site) | 63 | 368 | 0.7634 | *0.7770* | *0.7798* | 0.7634 | <u>0.7443</u> | *0.7660* | *0.7661* |
| Horse-colic (subtype) | 2 | 368 | <u>0.0027</u> | *0.0081* | *0.0081* | <u>0.0027</u> | <u>0.0027</u> | <u>0.0027</u> | <u>0.0027</u> |
| Horse-colic (type) | 8 | 368 | 0.0409 | 0.0407 | 0.0380 | 0.0409 | <u>0.0326</u> | 0.0353 | <u>0.0326</u> |
| Credit | 2 | 1000 | <u>0.2490</u> | *0.2569* | *0.2560* | <u>0.2490</u> | *0.2720* | *0.2579* | *0.2530* |
| Dermatology | 6 | 366 | 0.0272 | 0.0245 | *0.0274* | *0.0273* | 0.0272 | <u>0.0165</u> | 0.0219 |
| Ecoli | 8 | 336 | 0.1516 | 0.1427 | *0.1545* | *0.1545* | *0.1639* | <u>0.1368</u> | *0.1637* |
| Flags | 8 | 194 | 0.3760 | *0.4013* | *0.4318* | 0.3707 | <u>0.3705</u> | *0.3807* | *0.3860* |
| Glass | 7 | 214 | 0.5235 | <u>**0.3792**</u> | 0.4264 | 0.5235 | 0.4259 | 0.4958 | 0.5008 |
| Mushroom | 2 | 8124 | 0.0419 | **0.0179** | **0.0179** | **0.0374** | **0.0192** | <u>**0.0178**</u> | **0.0272** |
| Image Segmentation | 7 | 2310 | 0.1974 | 0.1809 | **0.1497** | 0.1974 | <u>**0.1259**</u> | 0.1727 | *0.1999* |
| Solar Flare (common) | 8 | 1066 | 0.2363 | **0.1726** | **0.1745** | 0.2363 | <u>**0.1707**</u> | **0.1979** | **0.2176** |
| Solar Flare (moderate) | 6 | 1066 | 0.0731 | <u>**0.0337**</u> | **0.0346** | 0.0731 | <u>**0.0337**</u> | **0.0506** | **0.0497** |
| Solar Flare (severe) | 3 | 1066 | 0.0281 | **0.0056** | <u>**0.0046**</u> | **0.0234** | <u>**0.0046**</u> | **0.0206** | **0.0141** |
| Vote | 2 | 435 | 0.0964 | *0.1081* | *0.1104* | 0.0964 | <u>0.0941</u> | *0.0965* | ***0.1056*** |
| Win/Loss | | | – | 14/10 | 11/13 | 6/5 | 13/9 | 14/9 | 13/10 |
| Significant Win/Loss | | | – | 7/3 | 7/2 | 3/1 | 6/0 | 7/1 | 5/2 |

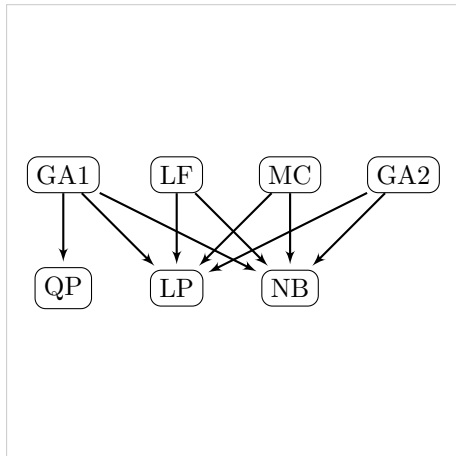(b) Non-uniform Costs. The values in this table represent performance under a non-uniform cost model.

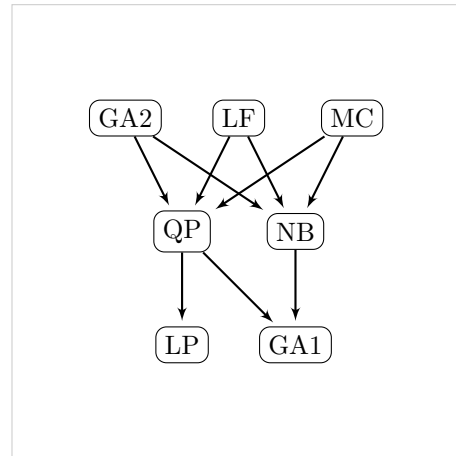| Data Set | $n$ | $m$ | Naïve Bayes | LF | MC | LP Eq. (21) | GA1 Eq. (22) | GA2 Eq. (1) | QP |
|---|---|---|---|---|---|---|---|---|---|
| Audiology | 24 | 226 | 1.7719 | 1.6826 | *1.8207* | 1.6385 | 1.7245 | <u>1.5194</u> | 1.6877 |
| Bridges 2 (material) | 3 | 107 | <u>0.8610</u> | ***1.2188*** | ***1.3630*** | ***1.4725*** | *1.0059* | ***1.2180*** | ***1.3585*** |
| Bridges 2 (rel-l) | 3 | 107 | 1.9355 | *2.2336* | *1.9764* | <u>1.9002</u> | ***2.5553*** | *1.9875* | 1.9163 |
| Bridges 2 (span) | 3 | 107 | <u>1.1871</u> | *1.7075* | *1.6680* | *1.2930* | ***1.9934*** | *1.4490* | *1.3508* |
| Bridges 2 (type) | 6 | 107 | 2.4584 | 2.3840 | <u>2.3612</u> | *2.6020* | ***2.7605*** | *2.5219* | *2.7270* |
| Bridges 2 (t-or-d) | 2 | 107 | 0.9654 | *1.0452* | 0.9729 | 0.9568 | *1.0794* | 0.9704 | <u>**0.8096**</u> |
| Car | 4 | 1728 | 0.8483 | **0.6653** | **0.7414** | 0.8073 | ***1.6665*** | <u>**0.6522**</u> | **0.7558** |
| Post-Op | 3 | 1473 | 2.9242 | **2.6791** | **2.7872** | *2.9993* | ***3.6611*** | **2.7056** | **2.8282** |
| Horse-colic (code) | 3 | 368 | 1.7914 | <u>1.6302</u> | 1.7247 | 1.7863 | *1.7950* | 1.6988 | 1.7726 |
| Horse-colic (surgical) | 2 | 368 | 1.3787 | <u>**1.0191**</u> | **1.0221** | *1.3890* | ***1.6364*** | **1.0629** | **1.1576** |
| Horse-colic (site) | 63 | 368 | 4.0891 | <u>3.9570</u> | *4.1976* | *4.1083* | *4.2647* | 4.0809 | *4.1560* |
| Horse-colic (subtype) | 2 | 368 | <u>0.0113</u> | *0.0432* | *0.0432* | <u>0.0113</u> | <u>0.0113</u> | <u>0.0113</u> | <u>0.0113</u> |
| Horse-colic (type) | 8 | 368 | 0.2225 | *0.2281* | 0.2086 | 0.2172 | <u>0.1704</u> | 0.1969 | 0.1894 |
| Credit | 2 | 1000 | 1.3202 | **1.0458** | **1.0477** | ***1.3950*** | ***2.1906*** | **1.0531** | **1.0805** |
| Dermatology | 6 | 366 | 0.1743 | *0.2649* | 0.1375 | 0.1564 | *0.1775* | <u>0.1146</u> | 0.1166 |
| Ecoli | 8 | 336 | 0.8104 | *0.8903* | 0.8168 | 0.8765 | ***1.2116*** | <u>0.7678</u> | *0.9218* |
| Flags | 8 | 194 | 2.1590 | *2.3070* | *2.3673* | 2.1407 | *2.3013* | <u>**1.9888**</u> | *2.1644* |
| Glass | 7 | 214 | 3.1308 | <u>**2.0044**</u> | **2.2886** | 3.1301 | *3.4910* | 2.6720 | 2.9600 |
| Mushroom | 2 | 8124 | 0.1929 | **0.0993** | **0.0994** | **0.1783** | **0.1262** | <u>**0.1031**</u> | **0.1382** |
| Image Segmentation | 7 | 2310 | 1.0854 | *1.2809* | <u>**0.7733**</u> | 1.0854 | *1.0989* | **0.8951** | *1.0883* |
| Solar Flare (common) | 8 | 1066 | 1.3079 | <u>**0.9174**</u> | **0.9662** | *1.3199* | **1.1622** | **0.9843** | 1.2331 |
| Solar Flare (moderate) | 6 | 1066 | 0.4644 | <u>**0.2210**</u> | **0.2216** | 0.4628 | **0.2749** | **0.2651** | **0.3030** |
| Solar Flare (severe) | 3 | 1066 | 0.1682 | **0.0361** | <u>**0.0340**</u> | 0.1556 | **0.0800** | **0.1069** | **0.0838** |
| Vote | 2 | 435 | 0.4510 | 0.3786 | <u>**0.3699**</u> | 0.4510 | ***0.5346*** | *0.4577* | *0.4605* |
| Win/Loss | | | – | 14/10 | 15/9 | 12/9 | 6/17 | 17/6 | 15/8 |
| Significant Win/Loss | | | – | 9/2 | 10/1 | 1/4 | 4/9 | 10/1 | 8/1 |

18

(a) Error - Optimization.

(b) Non-uniform Cost - Optimization.

(c) Error - Generalization.

(d) Non-uniform Cost - Generalization.

Figure 2: Relative Performances. An arrow from algorithm $A$ to $B$ indicates $A$ was statistically significantly better than $B$. Transitive relations are implicit. Statistics are according to a Tukey-Kramer pairwise comparison with a 95% confidence level.

Table 3: Mean Ranks. For each data set, a rank 1–7 is assigned to algorithms according to their performance. The tables below lists the mean rank for each algorithm over all data sets for the optimization/generalization problems in the uniform/non-uniform cost settings. The overall rank within each column is shown in parentheses.

| Algorithm | Training | | Testing | |
|-----------|----------|----------|---------|----------|
|           | Error    | Cost     | Error   | Cost     |
| NB  | 5.7604 (7) | 5.4167 (6) | 4.4875 (7) | 4.4250 (5) |
| LF  | 2.1771 (1) | 2.6146 (3) | 3.6542 (2) | 3.3188 (2) |
| MC  | 2.4604 (2) | 1.8708 (1) | 3.8354 (3) | 3.3229 (3) |
| LP  | 5.3125 (6) | 5.1896 (5) | 4.4583 (6) | 4.7542 (6) |
| GA1 | 4.4146 (4) | 5.7562 (7) | 3.6292 (1) | 5.0896 (7) |
| GA2 | 2.8750 (3) | 2.4979 (2) | 3.8458 (4) | 3.1292 (1) |
| QP  | 5.0000 (5) | 4.6542 (4) | 4.0896 (5) | 3.9604 (4) |

the data sets we used did not contain a large enough number of classes for the speed to become a significant advantage. However, in applications where the number of classes is large, this might become a significant issue. In addition, while both have similar resource requirements, the design of MetaClass is arguably more flexible. First, MetaClass learns what is essentially a decision tree, which is a provably more expressive hypothesis than a linear reweighting function. Any linear function can always be realized by a decision tree, but the converse is not true. Second, the choice of how we divide into separate metaclasses also gives us greater control over the learning process. In particular, since we can constrain the structure of the decision tree, we can customize to fit a specific learning domain using prior information. For instance, such an approach may fit well with hierarchical-class problems where we are given say, a taxonomy (i.e. an existing tree structure), and simply have to learn MetaClass's parameters.

# 7    Conclusions

Reoptimizing an already-learned classifier $f$ is an important problem in machine learning, particularly in applications where the cost model or class distribution of a learning problem deviates from the conditions under which a classifier $f$ was trained and the original training data for $f$ is unavailable. We answered the open problem concerning the hardness of this reoptimization problem by showing that the decision version is NP-complete. Though impractical, we also presented an algorithm that produces an optimal solution in polynomial time when the number of classes is constant. We also presented multiple algorithms for the multiclass version of this problem. Our experimental results showed that the greedy algorithms (our MetaClass algorithm and the algorithm of Lachiche and Flach (2003)) and one of our genetic algorithm approaches were roughly equally successful in reoptimizing under both uniform and non-uniform cost models. We found that, despite MetaClass's faster time complexity, in our experiments the number of classes was not large enough to give it a significant speed advantage. However, we also argued that MetaClass has several advantages over the algorithm of Lachiche & Flach in terms of its more expressive hypothesis class (trees) and its flexibility in partitioning the classes into metaclasses (e.g. in hierarchical classification settings). Future work is to precisely identify contexts in which these advantages can be exploited.

There are several other avenues for future research. In the non-uniform cost setting, our methodology omitted extreme skews in costs between classes. It may be of interest to examine such settings to see if the poor performing algorithms are able show improvement and to see if the better, greedy approaches continue to perform well. It would also be of interest to see if the greedy algorithms' relative performance holds when we change the base classifier. Do the greedy approaches still dominate if the base classifier is a more sophisticated learning algorithm like an SVM?

Finally, it is entirely possible that a better and more practical algorithm exists for the constant class setting. If so, how does its performance compare to the other methods?

# Acknowledgments

# References

Abramson, M. A. (2005). Genetic algorithm and direct search toolbox. `http://www.mathworks.com/`.

Blake, C., & Merz, C. (2005). UCI repository of machine learning databases. `http://www.ics.uci.edu/~mlearn/MLRepository.html`.

Boyd, S., & Vandenberghe, L. (2004). *Convex optimization.* Cambridge University Press.

Demšar, J. (2006). Statistical comparisons of classifiers over multiple data sets. *Journal of Machine Learning Research, 7*, 1–30.

Deng, K., Bourke, C., Scott, S. D., & Vinodchandran, N. V. (2006). New algorithms for optimizing multi-class classifiers via ROC surfaces. *Proceedings of the ICML Workshop on ROC Analysis in Machine Learning* (pp. 17–24).

Ferri, C., Hernández-Orallo, J., & Salido, M. (2003). Volume under the ROC surface for multi-class problems. *European Conference on Artificial Intelligence* (pp. 108–120).

Fieldsend, J., & Everson, R. (2005). Formulation and comparison of multi-class ROC surfaces. *Proceedings of the ICML Workshop on ROC Analysis in Machine Learning* (pp. 41–48).

Grant, M., Boyd, S., & Ye, Y. (2006). Disciplined convex programming. In L. Liberti and N. Maculan (Eds.), *Global optimization: From theory to implementation*, Nonconvex Optimization and its Applications, 155–210. Springer. Available at `http://www.stanford.edu/~boyd/cvx/`.

Hand, D., & Till, R. (2001). A simple generalisation of the area under the ROC curve for multiple class classification problems. *Machine Learning, 45*, 171–186.

Hochberg, Y., & Tamhane, A. C. (1987). *Multiple comparison procedures.* Wiley.

Kohli, R., Krishnamurti, R., & Mirchandani, P. (1994). The minimum satisfiability problem. *SIAM Journal of Discrete Mathematics, 7*, 275–283.

Lachiche, N., & Flach, P. (1999). 1BC: A first-order bayesian classifier. *Proceedings of the 9th International Workshop on Inductive Logic Programming* (pp. 92–103).

Lachiche, N., & Flach, P. (2003). Improving accuracy and cost of two-class and multi-class probabilistic classifiers using ROC curves. *Proceedings of the 20th International Conference on Machine Learning* (pp. 416–423).

Matsui, T. (1996). NP-hardness of linear multiplicative programming and related problems. *Journal of Global Optimization, 9*, 113–119.

MOSEK ApS (2005). The MOSEK optimization tools version 3.2. `http://www.mosek.com/`.

Mossman, D. (1999). Three-way ROCs. *Medical Decision Making*, 78–89.

O'Brien, D. B., & Gray, R. M. (2005). Improving classification performance by exploring the role of cost matrices in partitioning the estimated class probability space. *Proceedings of the ICML Workshop on ROC Analysis in Machine Learning* (pp. 79–86).

Provost, F. J., & Fawcett, T. (1997). Analysis and visualization of classifier performance: Comparison under imprecise class and cost distributions. *Proceedings of the Third International Conference on Knowledge Discovery and Data Mining (KDD-97)* (pp. 43–48).

Provost, F. J., & Fawcett, T. (1998). Robust classification systems for imprecise environments. *Proceedings of the 15th national conference on Artificial intelligence (AAAI)* (pp. 706–713).

Provost, F. J., & Fawcett, T. (2001). Robust classification for imprecise environments. *Machine Learning, 42*, 203–231.

Rockafellar, R. (1970). *Convex analysis (2nd ed.).* Princeton University Press.

Schölkopf, B., & Smola, A. J. (2001). *Learning with kernels.* Cambridge, MA, USA: MIT Press.

Srinivasan, A. (1999). *Note on the location of optimal classifiers in n-dimensional ROC space* (Technical Report PRG-TR-2-99). Oxford University Computing Laboratory, Oxford.

Stoer, I. J., & Witzgall, C. (1996). *Convexity and optimization in finite dimensions.* Springer-Verlag.

Witten, I. H., et al. (2005). Weka machine learning toolbox. `http://www.cs.waikato.ac.nz/ml/weka/`.