

ROC graphs with instance-varying costs

Tom Fawcett

Institute for the Study of Learning and Expertise, 2164 Staunton Court, Palo Alto, CA 94306, USA

Available online 19 December 2005

Abstract

Receiver operating characteristics (ROC) graphs are useful for organizing classifiers and visualizing their performance. ROC graphs have been used in cost-sensitive learning because of the ease with which class skew and error cost information can be applied to them to yield cost-sensitive decisions. However, they have been criticized because of their inability to handle instance-varying costs; that is, domains in which error costs vary from one instance to another. This paper presents and investigates a technique for adapting ROC graphs for use with domains in which misclassification costs vary within the instance population.

© 2005 Elsevier B.V. All rights reserved.

Keywords: ROC analysis; Cost-sensitive learning; Classifier evaluation

1. Introduction

A receiver operating characteristics (ROC) graph is a simple technique for visualizing, organizing and selecting classifiers based on their performance. ROC graphs have long been used in signal detection theory to depict the trade-off between hit rates and false alarm rates of classifiers (Egan, 1975; Swets et al., 2000). ROC analysis has been extended for use in visualizing and analyzing the behavior of diagnostic systems (Swets, 1988; Fawcett, 2003). The medical decision-making community has an extensive literature on the use of ROC graphs for diagnostic testing (Zou, 2002). Swets et al. (2000) brought ROC curves to the attention of the wider public with their *Scientific American* article.

Recent years have seen an increase in the use of ROC graphs in the machine learning and pattern recognition communities. One advantage of ROC graphs is that they enable visualizing and organizing classifier performance without regard to class distributions or error costs (Provost and Fawcett, 1997). This ability becomes very important when dealing with skewed distributions or cost-sensitive learning. A researcher can graph the performance of a set

of classifiers, and that graph will remain invariant with respect to the operating conditions (class skew and error costs). As these conditions change, the region of interest of the graph may change, but the graph itself will not. In such cases, the researcher calculates, at classification time, the approximate operating conditions under which the set of classifiers will be used, overlays it onto the ROC graph, and uses the information to choose which classifier(s) to use. As conditions change, the ROC graph may be reconsulted, but the classifiers need not be re-evaluated. If one classifier is broadly superior, the researcher may decide to discard the others (Provost and Fawcett, 2001). Alternatively, cost curves (Drummond and Holte, 2000) may be used to depict classifier costs directly.

Unfortunately, such methods have an inherent limitation. ROC graphs plot true positive rate against false positive rate, treating all errors of a given type to be equivalent. In some domains, this assumption does not hold: the cost of a particular kind of error is not constant throughout the population but varies by example. A typical such domain is that of credit card fraud detection, in which transactions must be evaluated in real time and judged fraudulent or legitimate. A \$1000 transaction, if fraudulent, is much more costly than a fraudulent \$10 transaction, and classification error costs for the transactions should be correspondingly different. Such costs have been given a

E-mail addresses: tfawcett@acm.org, tom.fawcett@gmail.com

variety of names: example-specific costs, instance-varying costs, and case-conditional error costs (Turney, 2000).

ROC graphs have been criticized because of their inability to handle example-specific costs. In this paper, we present a straightforward transformation of ROC graphs, called ROCIV graphs, that accommodate example-specific costs. We show domains in which such graphs are useful; that is, in which standard ROC graphs may mislead as to classifier superiority. We prove that the area under the ROCIV curve has a natural interpretation related to the area under an ROC curve. Finally, we show several domains with instance-varying costs, illustrate the ROCIV graph on each, and argue for its superiority over alternative methods for evaluating classifiers on the domain.

2. A brief review of ROC graphs

This section will briefly review ROC graphs. A much more detailed introduction is provided in (Fawcett, 2006). Let $\{p, n\}$ be the positive and negative instance classes, and let $\{Y, N\}$ be the classifications produced by a classifier. Let $p(p|I)$ be the posterior probability that instance I is positive. The true positive rate of a classifier is

$$\text{TP rate} = p(Y|p) \approx \frac{\text{positives correctly classified}}{\text{total positives}}$$

The false positive rate of a classifier is

$$\text{FP rate} = p(Y|n) \approx \frac{\text{negatives incorrectly classified}}{\text{total negatives}}$$

We will use the term *ROC space* to denote the classifier performance space used for visualization in ROC analysis. On an ROC graph, TP rate is plotted on the Y-axis and FP rate is plotted on the X-axis. These statistics vary together as a threshold on a classifier's continuous output is varied between its extremes, and the resulting curve is called the

ROC curve. The ROC curve illustrates the error trade-offs available with a given classifier. Fig. 1 shows a typical ROC plot of three classifiers.

Algorithm 1 is a basic algorithm for generating an ROC graph from a test set. It exploits the monotonicity of thresholded classifications: any instance that is classified as positive with respect to a given threshold will be classified as positive for all lower thresholds. This algorithm assumes that the classifier assigns scores to instances, proportional to the probability that a given instance is positive. The function $f(i)$ is the score assigned to instance i by the classifier. In this algorithm, TP and FP start at zero. For each positive instance, we increment TP and for every negative instance we increment FP. We maintain a stack R of ROC points, pushing a new point onto R after each instance is processed. The final output is the stack R , which will contain points on the ROC curve.

Algorithm 1. Generating ROC points.

Inputs: L , the set of test examples; $f(i)$, the probabilistic classifier's estimate that example i is positive; P and N , the number of positive and negative examples.

Outputs: R , a list of ROC points increasing by *fp rate*.

Require: $P > 0$ and $N > 0$

```

1:  $L_{\text{sorted}} \leftarrow L$  sorted decreasing by  $f$  scores
2:  $FP \leftarrow TP \leftarrow 0$ 
3:  $R \leftarrow \langle \rangle$ 
4:  $f_{\text{prev}} \leftarrow -\infty$ 
5:  $i \leftarrow 1$ 
6: While  $i \leq |L_{\text{sorted}}|$  do
7:   if  $f(i) \neq f_{\text{prev}}$  then
8:     push  $(\frac{FP}{N}, \frac{TP}{P})$  onto  $R$ 
9:      $f_{\text{prev}} \leftarrow f(i)$ 
10:  if  $L_{\text{sorted}}[i]$  is a positive example then
11:     $TP \leftarrow TP + 1$ 
12:  else /*  $i$  is a negative example */
13:     $FP \leftarrow FP + 1$ 
14:     $i \leftarrow i + 1$ 
15: push  $(\frac{FP}{N}, \frac{TP}{P})$  onto  $R$  /* This is (1,1) */
16: end
```

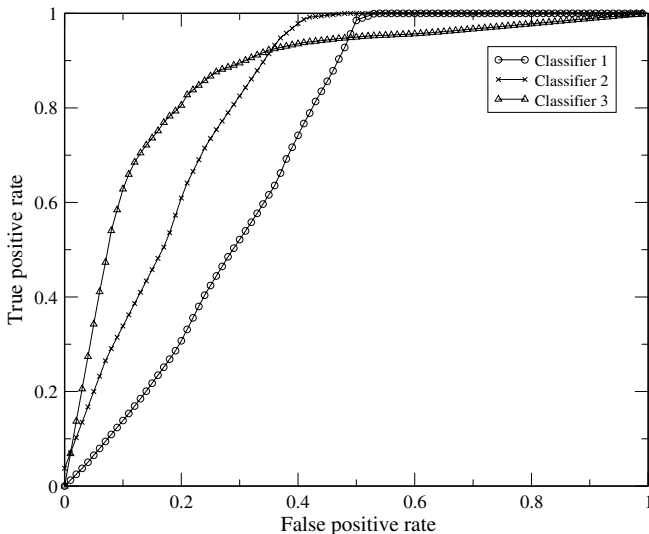


Fig. 1. An ROC graph of three scoring classifiers.

2.1. Error costs

Let $c(\text{hyp}, \text{class})$ be a two-place error cost function where hyp is the hypothesized class assigned to an instance by the classifier and class is the instance's real class. $c(Y, n)$ is the cost of a false positive error and $c(N, p)$ is the cost of a false negative error. In this paper, we shall treat costs and benefits equivalently: a cost is simply a negative benefit. Similarly, it is possible to “roll up” benefits into costs, defining error costs to include benefits not realized.

If a classifier produces posterior probabilities, decision analysis gives us a way to produce cost-sensitive classifications from the classifier (Weinstein and Fineberg, 1980). Classifier error frequencies can be used to approximate

probabilities (Pazzani et al., 1994). For an instance I , the decision to emit a positive classification is

$$[1 - p(\mathbf{p}|I)] \cdot c(\mathbf{Y}, \mathbf{n}) < p(\mathbf{p}|I) \cdot c(\mathbf{N}, \mathbf{p})$$

Regardless of whether a classifier produces probabilistic or binary classifications, its expected cost on a test set can be estimated as

$$\text{Cost} = \text{FP} \cdot c(\mathbf{Y}, \mathbf{n}) + \text{FN} \cdot c(\mathbf{N}, \mathbf{p})$$

Given a set of classifiers, a set of examples, and a precise cost function, most work on cost-sensitive classification uses an equation such as this to rank the classifiers according to cost and chooses the minimum. However, as discussed above, such analyses assume that the distributions are precise and static.

An advantage of ROC graphs is that they enable visualizing and organizing classifier performance without regard to class distributions or error costs. A researcher can graph the performance of a set of classifiers, and that graph will remain invariant with respect to the operating conditions (class skew and error costs). As these conditions change, the region of interest may change, but the graph itself will not. Identifying the region of interest is done using iso-performance lines.

2.2. Iso-performance lines

Provost and Fawcett (1998, 2001) show that a set of operating conditions may be transformed easily into a so-called *iso-performance line* in ROC space. Two points in ROC space, $(\text{FP}_1, \text{TP}_1)$ and $(\text{FP}_2, \text{TP}_2)$, have the same performance if

$$m = \frac{\text{TP}_2 - \text{TP}_1}{\text{FP}_2 - \text{FP}_1} = \frac{c(\mathbf{Y}, \mathbf{n})p(\mathbf{n})}{c(\mathbf{N}, \mathbf{p})p(\mathbf{p})} \quad (1)$$

This equation defines the slope of an iso-performance line. All classifiers corresponding to points on a line of slope m have the same expected cost. Each set of class and cost distributions defines a family of iso-performance lines. Lines “more northwest” (having a larger TP intercept) are better because they correspond to classifiers with lower expected cost.

The slope m is determined by external constraints: the class skew of the domain, represented by $p(\mathbf{n})/p(\mathbf{p})$, and the relative costs of false positive and false negative errors, represented by $c(\mathbf{Y}, \mathbf{n})/c(\mathbf{N}, \mathbf{p})$. All points (classifiers) along any such line will have equal expected cost. In essence, m represents the trade-off between the false positive error rate

and the true positive rate that is acceptable after taking into account the inherent class skew of the domain. For example, if $m = 9$, this represents the condition that a 1% increase in false positive rate is worth a 9% increase in true positive benefit (or, equivalently, a 9% decrease in false negative rate). However, this formulation assumes that all errors of a given type are the same. In other words, that false positive costs are constant within the population of negative examples and all false negative costs are constant within the population of positive examples.

3. Instance-varying costs

In some domains, the cost of a particular kind of error is not constant throughout the population but varies by example. Consider a simple credit card transaction domain used by Elkan (2001) in which the task is to decide whether to approve or refuse a given transaction. Elkan describes a benefit matrix for the task, shown in Fig. 2a. This cost matrix is justified with the following explanation. A refused fraudulent transaction has a benefit of \$20 because it may prevent future fraud. Refusing a legitimate transaction has a negative benefit because it annoys a customer. Approving a fraudulent transaction has a negative benefit proportional to the transaction amount (x). Approving a legitimate transaction generates a small amount of income proportional to the transaction amount ($0.02x$).

In this credit approval example, costs of mistakes are known beforehand because the credit amount is known at the time a classification must be made. In other domains, costs can only be known *after* an action has been taken. An example of this is when soliciting for charitable donations. If a charity decides not to solicit a prospective donor, the charity will likely never find out whether (or how much) the person would have donated.

To accommodate instance-varying costs, we extend the cost function. Costs may be expressed as three-place functions of the class, the hypothesized class and the instance x : $c(\text{hyp}, \text{class}, x)$.

There are several common ways in which researchers deal with example-specific costs. The most common approach is simply to calculate the total expected cost of a classifier (Elkan, 2001). Let X^+ and X^- be the set of positive and negative instances, respectively. If error costs are known exactly, the expected cost of a classifier may be calculated as

$$\sum_{x \in X^+} c(h(x), \mathbf{p}, x) + \sum_{x \in X^-} c(h(x), \mathbf{n}, x)$$

	fraudulent	legitimate
refuse	\$20	−\$20
approve	− x	$0.02x$

(a)

	fraudulent	legitimate
refuse	0	0
approve	$\$20 + x$	$0.02x + \$20$

(b)

Fig. 2. Cost matrices for the credit approval domain. (a) Original benefit matrix and (b) transformed cost-benefit matrix.

where $h(x)$ is the hypothesized class of instance x . This equation achieves an exact solution but it loses the advantage of ROC curves, which is to allow classifier performance to be visualized and compared over a range of performance conditions. However, this solution may be appropriate if costs and class distributions are known exactly and are known not to change.

Another method is to smooth out the costs: measure the average costs of false positives and false negatives and use these average values:

$$c(N, \mathbf{p}) \approx \sum_{x \in X^+} c(N, \mathbf{p}, x) / |X^+| \quad (2)$$

$$c(Y, \mathbf{n}) \approx \sum_{x \in X^-} c(Y, \mathbf{n}, x) / |X^-| \quad (3)$$

These two-place cost functions may then be used in Eq. (1) so that iso-performance lines can be used with ROC curves. Previous work on fraud detection (Fawcett and Provost, 1997) took this approach, and it may be acceptable if cost variance is small.

4. ROC Graphs with instance-varying costs

Another approach is to use a straightforward transformation of ROC graphs, explained in this section. Consider again the simple credit card transaction domain whose cost matrix is shown in Fig. 2a. For this domain, we assume that a Y decision corresponds to approving a transaction, and N means denying it. The default action will be to deny a transaction. To use a cost matrix for an ROC graph it must be transformed into a cost-benefit matrix where the costs are relative only to Y decisions.

First we subtract the first row from both rows in the matrix. Conceptually the resulting matrix corresponds to a baseline situation where all transactions are refused, so all fraud is denied and all legitimate customers are annoyed. We then negate the approve-fraudulent cell to turn it into a cost. This yields the cost-benefit matrix of Fig. 2b which forms the definition of the cost function $c(Y, \mathbf{p}, x)$ and $c(Y, \mathbf{n}, x)$.

In standard ROC graphs, the x -axis represents the fraction of total FP mistakes possible. In the instance-varying cost formulation, it will represent the fraction of *total FP cost* possible, so the denominator will now be

$$\sum_{x \in X^+} \$20 + x$$

Similarly the y -axis will be the fraction of *total TP benefits* so its denominator will be

$$\sum_{x \in X^+} 0.02x + \$20$$

Instead of incrementing TP and FP instance counts as in Algorithm 1, we increment $TP_benefit$ and FP_cost by the cost (benefit) of each negative (positive) instance as it is processed. The ROCIV points are the fractions of total

Algorithm 2. Generating ROCIV points and area from a dataset with example-specific costs.

Inputs: L , the set of test examples; $f(i)$, the probabilistic classifier's estimate that example i is positive; P and N , the number of positive and negative examples; $c(Y, class, i)$, the cost of judging instance i of class $class$ to be Y .

Outputs: R , a list of ROCIV points increasing by fp rate; A , the area under the ROCIV curve.

Require: $P > 0$ and $N > 0$

```

1: for  $x \in L$  do
2:   if  $x$  is a positive example then
3:      $P\_total \leftarrow P\_total + c(Y, \mathbf{p}, x)$ 
4:   else
5:      $N\_total \leftarrow N\_total + c(Y, \mathbf{n}, x)$ 
6:  $L\_sorted \leftarrow L$  sorted decreasing by  $f$  scores
7:  $FP\_cost \leftarrow TP\_benefit \leftarrow 0$ 
8:  $FP\_cost_{prev} \leftarrow TP\_benefit_{prev} \leftarrow 0$ 
9:  $A \leftarrow 0$  /* Area under the curve */
10:  $R \leftarrow \langle \rangle$  /* ROCIV points */
11:  $f_{prev} \leftarrow -\infty$ 
12:  $i \leftarrow 1$ 
13: while  $i \leq |L\_sorted|$  do
14:   if  $f(i) \neq f_{prev}$  then
15:     push  $(\frac{FP\_cost}{N\_total}, \frac{TP\_benefit}{P\_total})$  onto  $R$ 
16:      $A \leftarrow A + TRAP\_AREA(FP\_cost, FP\_cost_{prev}, TP\_benefit_{prev}, TP\_benefit)$ 
17:      $f_{prev} \leftarrow f(i)$ 
18:      $FP\_cost_{prev} \leftarrow FP\_cost$ 
19:      $TP\_benefit_{prev} \leftarrow TP\_benefit$ 
20:   if  $L\_sorted[i]$  is a positive example then
21:      $TP\_benefit \leftarrow TP\_benefit + c(Y, \mathbf{p}, L\_sorted[i])$ 
22:   else /*  $L\_sorted[i]$  is a negative example */
23:      $FP\_cost \leftarrow FP\_cost + c(Y, \mathbf{n}, L\_sorted[i])$ 
24:      $i \leftarrow i + 1$ 
25: push  $(\frac{FP\_cost}{N\_total}, \frac{TP\_benefit}{P\_total})$  onto  $R$  /* This is (1,1) */
26:  $A \leftarrow A + TRAP\_AREA(FP\_cost, FP\_cost_{prev}, TP\_benefit_{prev}, TP\_benefit)$ 
27:  $A \leftarrow A / (P\_total \times N\_total)$ 
28: end

```

benefits and costs, respectively. Conceptually this transformation corresponds to replicating instances in the instance set in proportion to their cost, though this transformation has the advantage that no actual replication is performed and non-integer costs are easily accommodated. We shall call these transformed ROC curves ROCIV curves.

The final algorithm for generating ROCIV points is Algorithm 2. This algorithm also calculates the area under the ROCIV curve at the same time that it generates the points. To do this, it sums the areas of successive trapezoids¹ as the points are being processed. The area under a ROCIV curve will be discussed further in Section 5.

¹ The function TRAP_AREA calculates these areas using the simple $base \times height/2$ formula for trapezoids.

Fig. 3 illustrates the difference between ROC curves and ROCIV curves via this transformation. Fig. 3a shows the ROC curves of two classifiers, A and B, with identical performance (their ROC curves are identical). This ROC curve shows raw classification performance with respect to positive and negative examples. Fig. 3 shows the ROCIV curves of A and B when example-specific costs are taken into account. Note that performance differs significantly, and each has regions of dominance. The ROC curve in Fig. 3a is a poor representation of the performance of either, and using it to select a low-cost classifier might be misleading.

Section 2.2 discussed how operating conditions could be transformed into an iso-performance line and used with ROC curves to choose the best performing classifier for those conditions. With ROCIV graphs, the iso-performance line slope is calculated the same way as with ROC curves, but the error costs $c(N, \mathbf{p})$ and $c(Y, \mathbf{n})$ are averages as calculated in Eqs. (2) and (3). Thus, when all costs of instances within a given class are equal, a ROCIV graph is identical to a ROC graph.

Algorithm 2 ranks the instances (in line 6) by their f scores generated by the classifier. As indicated in the *Inputs* section, these are assumed to be probabilities that an instance is positive. This is appropriate for evaluating classifiers on domains in which instance costs are not known at classification time. However, when the misclassification costs for instances are known, it is more natural that the classifier takes into account each cost and produce scores that reflect the misclassification cost—for example, the product of the instance's estimated posterior probability and its false negative cost: $p(\mathbf{p}|\mathbf{x}) \cdot c(N, \mathbf{p}, \mathbf{x})$.

5. Area under a ROCIV curve

The area under a conventional ROC curve (AUC) has an important statistical property: the AUC of a classifier is equivalent to the probability that the classifier will rank a randomly chosen positive instance higher than a randomly chosen negative instance. This is equivalent to the Wilcoxon test of ranks (Hanley and McNeil, 1982). In his Ph.D. thesis, Sing (2004) presents an elegant and intuitive proof of this. In this section, we show that the area under the ROCIV (which we shall call AUCIV) is very similar: it is equivalent to the same probability, with the stipulation that *instances are chosen in proportion to their costs*. The following proof borrows heavily from the notation and structure of Sing's (indeed, this proof may be seen as a variant of his).

Definition. Let n^+ and n^- denote the number of positive and negative examples, respectively. Denote the positive training examples by $x_1^+, \dots, x_{n^+}^+$ and the negative training examples by $x_1^-, \dots, x_{n^-}^-$. Let $\text{cost}(P) = \sum_{i=1}^{n^+} \text{cost}(x_i^+)$ and $\text{cost}(N) = \sum_{j=1}^{n^-} \text{cost}(x_j^-)$. We define two cost-weighted probability functions:

$$p_{\text{ivc}}(x_i^+) = \text{cost}(x_i^+) / \text{cost}(P)$$

$$p_{\text{ivc}}(x_j^-) = \text{cost}(x_j^-) / \text{cost}(N)$$

Theorem. The area under an ROCIV curve of a classifier h is equivalent to the probability that h will rank a randomly chosen positive instance higher than a randomly chosen negative instance, assuming that the probability of an instance being chosen is proportional to its cost. Formally:

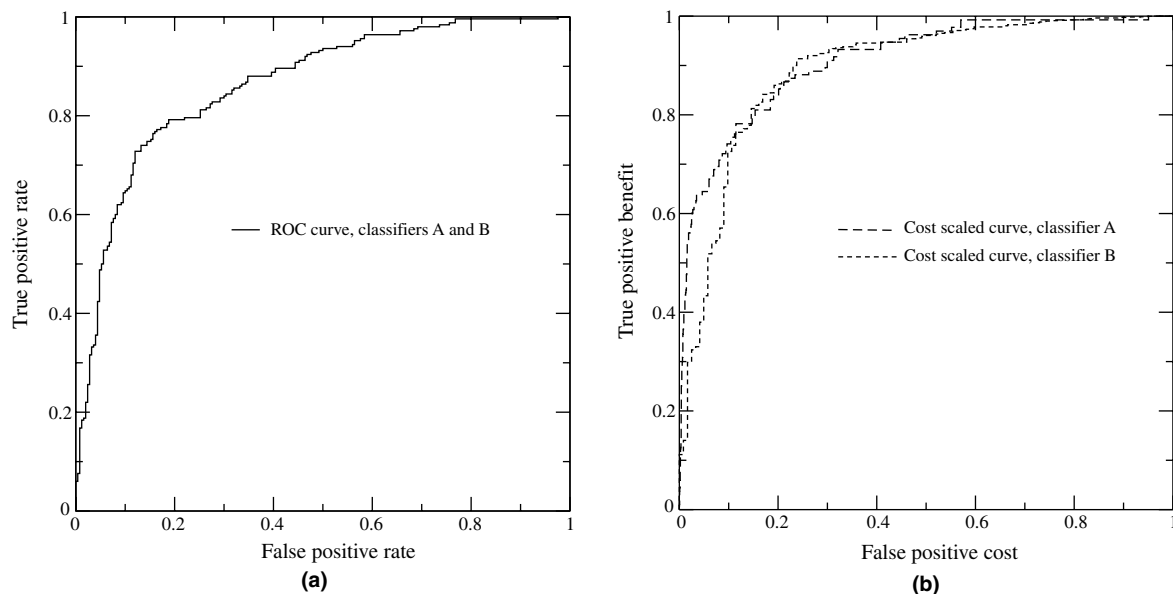


Fig. 3. Standard ROC curves vs ROCIV curves. (a) ROC curves of two classifiers and (b) ROCIV curves of the same classifiers.

$$\text{AUCIV} = \sum_{j=1}^{n^-} \sum_{i=1}^{n^+} p_{\text{ivc}}(x_i^+) \cdot p_{\text{ivc}}(x_j^-) \cdot \mathbf{I}_{h(x_i^+) > h(x_j^-)}$$

Proof. In order to simplify the proof, we assume without loss of generality that h assigns unique scores to each of the instances.

Let $A = \{a_1, \dots, a_n\}$ and $B = \{b_1, \dots, b_m\}$, $a_i, b_j \in R$. The *cost weighted rank* of a sample b_j with respect to A is the summed costs of the elements of A with higher value:

$$\text{rank}_c(b_j|A) = \sum_{i=1}^n \text{cost}(a_i) \cdot \mathbf{I}_{a_i > b_j}$$

Assume that positive and negative samples are labeled in descending order so that $h(x_1^+) > \dots > h(x_{n^+}^+)$ and $h(x_1^-) > \dots > h(x_{n^-}^-)$. Fig. 4 shows an unscaled ROCIV curve with the predictions arranged in descending h order. The curve steps upward with a positive example (marked by an empty circle) and to the right with a negative example (filled circle). The actual ROCIV curve would have the x -axis scaled by $1/\text{cost}(N)$ and the y -axis scaled by $1/\text{cost}(P)$. The area under the curve (shaded in the figure) is related to the AUCIV by

$$\text{AUCIV} = \frac{\text{Area}}{\text{cost}(P) \cdot \text{cost}(N)}$$

Each vertical column may be thought of as belonging to a negative example, marked by the filled circle at the top of its column. The area of an example's column is equal to the cost of that example times the sum of the costs of the positive examples that scored higher. This is the cost-weighted rank of the negative example with respect to the positive examples. Therefore,

$$\text{Area} = \sum_{j=1}^{n^-} \text{cost}(x_j^-) \cdot \text{rank}_c(x_j^- | \{x_1^+, \dots, x_{n^+}^+\})$$

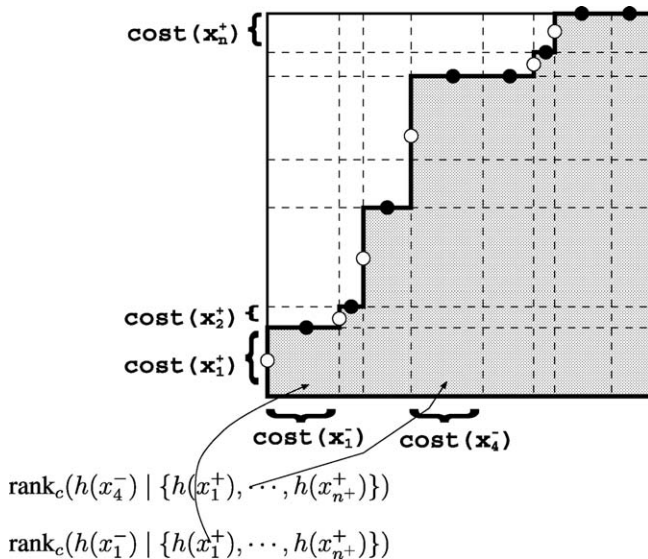


Fig. 4. A cost-scaled ROCIV plot.

Rewriting this in terms of the AUCIV we have

$$\begin{aligned} \text{AUCIV} &= \frac{\sum_{j=1}^{n^-} \text{cost}(x_j^-) \cdot \text{rank}_c(x_j^- | \{x_1^+, \dots, x_{n^+}^+\})}{\text{cost}(P) \cdot \text{cost}(N)} \\ &= \frac{\sum_{j=1}^{n^-} \sum_{i=1}^{n^+} \text{cost}(x_j^-) \cdot \text{cost}(x_i^+) \cdot \mathbf{I}_{h(x_i^+) > h(x_j^-)}}{\text{cost}(P) \cdot \text{cost}(N)} \\ &= \sum_{j=1}^{n^-} \sum_{i=1}^{n^+} \frac{\text{cost}(x_i^+)}{\text{cost}(P)} \cdot \frac{\text{cost}(x_j^-)}{\text{cost}(N)} \cdot \mathbf{I}_{h(x_i^+) > h(x_j^-)} \\ &= \sum_{j=1}^{n^-} \sum_{i=1}^{n^+} p_{\text{ivc}}(x_i^+) \cdot p_{\text{ivc}}(x_j^-) \cdot \mathbf{I}_{h(x_i^+) > h(x_j^-)} \quad \square \end{aligned}$$

6. Examples

To illustrate the effect of the ROCIV transformation, we present an empirical demonstration of its use for evaluation classifier performance on several domains with example-specific costs.

It should be emphasized that the point of the demonstrations here is to show classifier performance on realistic cost-varying domains. No claim is made that these are the best learning algorithms for these domains, or that the relative performance reported is typical for these algorithms.

6.1. Charitable donations

The first domain is a proprietary dataset used at Hewlett-Packard, provided by a third party. The data comprise about 60,000 records of solicitation response data for a charity. The independent variables measure, for each person, a basic history comprising recency and frequency of donations to the charity, and some features capturing periodicity of donations. The dependent (response) variable is the donation amount of a single mailing.

The classification problem is to determine whether a given person will donate. From the donation amount and the mailing costs, misclassification costs may be derived. The cost of a false positive is the cost of mailing a solicitation for which no donation is received. For this study, the mailing amount was assumed to be \$1. The benefit of a true positive is the donation amount x minus the \$1 mailing cost. For this mailing, when a donation is given the amounts varied greatly, from \$1 to \$1500 (mean $\$24 \pm \49).

Several classification models were trained on this data, including Logistic regression, Naive bayes, a Neural network, J48 (a decision tree learner similar to C4.5) and PRIE. PRIE is a rule learning system designed to maximize ROC performance.²

Fig. 5 shows resulting ROC and ROCIV curves for this domain. Note that in the ROC curves, the performance of

² PRIE's algorithm has not been published previously but the details of its operation are not germane to the illustration here.

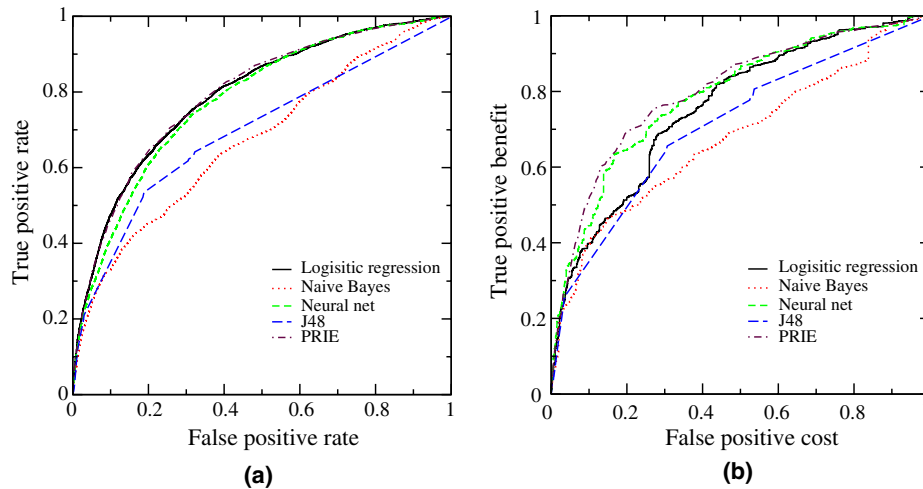


Fig. 5. Classifier performance, charitable donation data. (a) ROC curves and (b) ROCIV curves.

PRIE and logistic regression are nearly indistinguishable, and the performance of the neural net is very close. In the ROCIV curves, the relationships are somewhat different. The greatest difference is that logistic regression performs noticeably worse, dropping substantially below that of PRIE over most of the false positive range. The ROC curve of logistic regression is nearly monotonic, whereas its ROCIV curve has prominent concavities.

This differing performance may also be explained by looking at the costs of instances as they are ranked by each algorithm. These costs are shown in Fig. 6. The Y-axis shows the cost of each instance, with the instances ordered from highest assigned score (at the left) to lowest (at the right). There are many instances of cost -1 (the non-responders), which are difficult to see because of the scale. The ROCIV curves are essentially showing the simulta-

neous integrals of these instance costs: the true positive benefit is the integral of instance benefits above the line, and the false positive cost is the integral of instance costs below the line. An ideal classifier would rank all positive instances first, decreasing by benefit, followed by all negative instances, increasing by cost. One can see informally that the “cost mass” of the instances ranked by PRIE is slightly to the left of the cost mass of Logistic regression. This difference is reflected in the resulting ROCIV curves.

6.2. Credit scoring

The “German” domain is a dataset of German credit information provided by Dr. Hans Hofmann and donated to the UCI Machine Learning Repository (Hettich et al., 1998). Each record contains information about a person

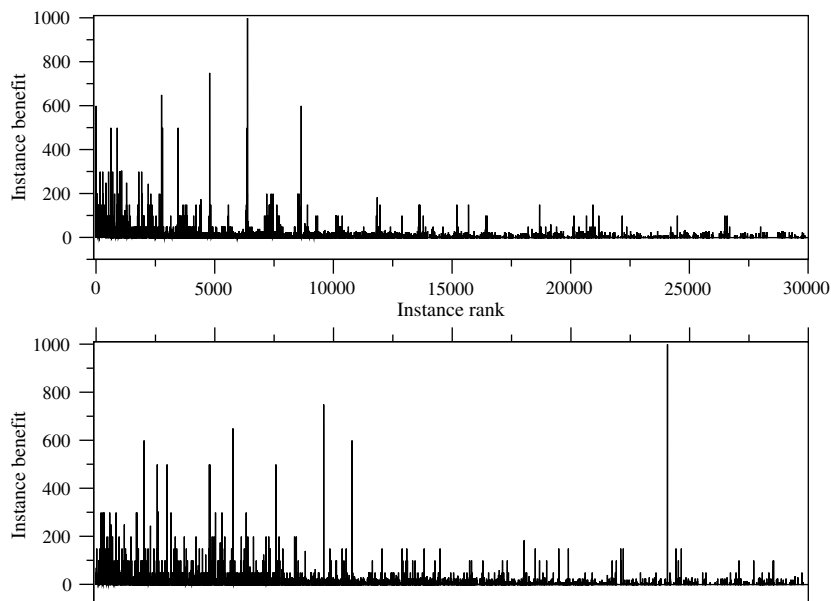


Fig. 6. Charitable donation data: Costs of instances as ordered by PRIE (top) and Logistic regression (bottom).

	fraudulent	legitimate
refuse	\$20	−\$20
approve	− x	$0.05x$

(a)

	fraudulent	legitimate
refuse	0	0
approve	$\$20 + x$	$0.05x + \$20$

(b)

Fig. 7. Matrices for the credit scoring domain. (a) Original benefit matrix and (b) transformed cost-benefit matrix.

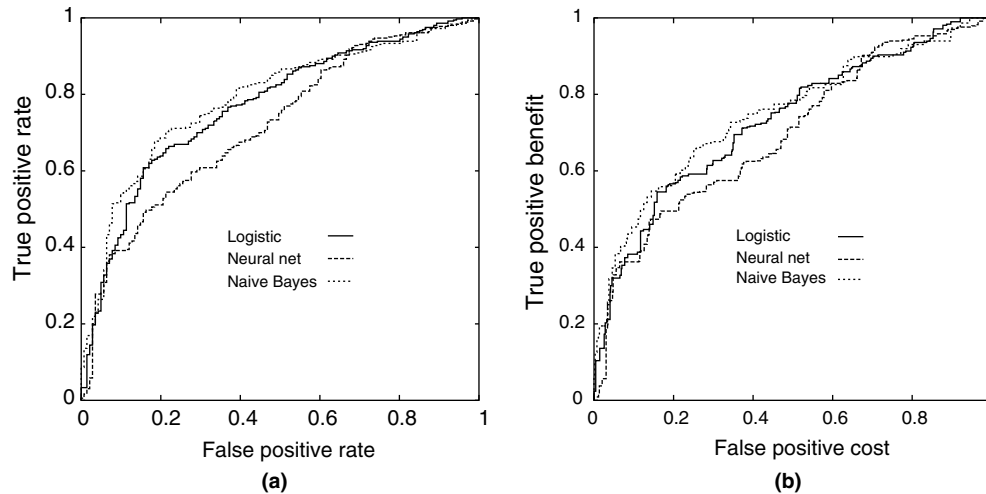


Fig. 8. Classifier performance, credit scoring domain. (a) ROC curves and (b) ROCIV curves.

requesting a loan, including various demographic information, a summary of the credit history, and the amount of the loan.

Classifiers were trained on this domain using the cost matrix in Fig. 7b, generated from 7a. Three classifier model types were induced using the Weka package (Witten and Frank, 2000): a multilayer perceptron, denoted “Neural net”; a Naive Bayes learner, and a simple logistic regression model. These model types were chosen because they were reasonably good at producing instance scores rather than simply assigning a class label to each instance.

Half of the instances were used for training, half for testing. The ROC curves on the test set are shown in Fig. 8a. Looking at the convex hull of the classifiers, each has a region of superiority: the neural net is most conservative ($0 \leq \text{FP} \leq 0.10$), followed by simple logistic regression ($0.10 < \text{FP} \leq 0.78$), followed by Naive Bayes ($0.78 < \text{FP} \leq 1.0$).

The corresponding ROCIV curves, shown in Fig. 8b, tell a different story. When considering individual instances’ costs and benefits, the regions of superiority do not correspond to those of Fig. 8a. Naive Bayes is superior over a much larger region, and the Neural net is virtually dominated completely.

6.3. Fraudulent phone call detection

As another demonstration of instance-varying costs we examined a simplified form of cell phone fraud detection. An historical dataset of 50,000 cell phone calls was used:

25,000 for training and 25,000 for testing. Fawcett and Provost (1997) previously described the domain in detail. The version used here is simplified in that the goal is to classify individual calls rather than to profile user behavior over time and to classify account days. The classifiers learned here are similar to the fraudulent call classifiers described by Fawcett and Provost (1997) in Section 6.3.

The independent variables in this domain are attributes of a given cell phone call, such as its originating and terminating locations, the length, the carrier used, time of day, and so on. The dependent variable is simply a binary flag indicating whether the call was fraudulent.

Fig. 9 shows the benefit and cost-benefit matrices for this domain. In this domain, the default action will be to do nothing (the account is not suspected of fraud). An alarm consists of flagging the account and temporarily disabling it. The top matrix is justified as follows.³ Alarming on a fraudulent call serves to inhibit some future fraud, a benefit valued at approximately \$1. Issuing an alarm on a legitimate customer temporarily shuts down the account and irritates the customer, a benefit of −\$2. Missing a fraudulent call costs a certain amount in overhead and toll costs (long distance and international charges are billed to the carrier) and eventual customer irritation. Finally,

³ While the costs and benefits used here are plausible, they are fabricated and should not be interpreted as actual values assigned by Bell Atlantic Mobile.

	fraudulent	legitimate		fraudulent	legitimate
no alarm	$-\$1 - \$0.20x$	$\$0.10x$	no alarm	0	0
alarm	$\$1$	$-\$2$	alarm	$\$2 + \$0.20x$	$-\$2 - \$0.10x$
(a)			(b)		

Fig. 9. Matrices for the cell phone fraud detection domain. (a) Original benefit matrix and (b) transformed cost-benefit matrix.

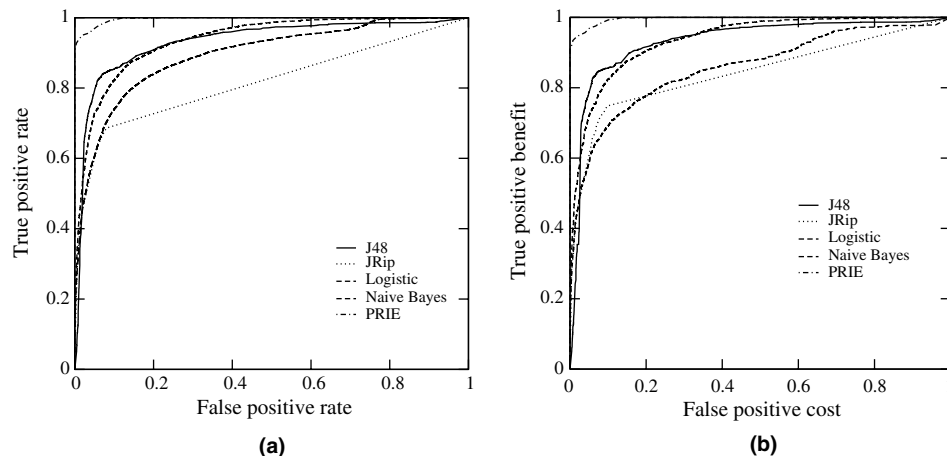


Fig. 10. Classifier performance, fraudulent call classification. (a) ROC curves and (b) ROCIV curves.

allowing a legitimate call generates a certain amount of revenue.

Following the transformation described in Section 4, the cost-benefit matrix in Fig. 9b is derived. The results are shown in Fig. 10 using a variety of classifiers. JRip is a rule learning classifier from the Weka package. The difference between the ROC and the ROCIV curves in this domain are very slight, perhaps because the relative performance of the classifiers are so distinct in this domain. The only substantial change is that Naive Bayes declines noticeably in the ROCIV curves, and no longer dominates JRip.

7. Discussion

This paper has presented a straightforward transformation of ROC graphs, called ROCIV graphs, that accommodate instance-varying costs. The new curves have a intuitive interpretation, in that the axes are now scaled by instance costs within each class. The area under the ROCIV curves has a straightforward interpretation: it is equivalent to the probability that a randomly chosen positive instance will be ranked more highly than a randomly chosen negative instance, given that each is chosen in proportion to their costs. Finally, we have demonstrated the transformation on three domains and have seen cases in which the ROCIV curves show considerably different regions of classifier superiority than the corresponding ROC curves.

For clarity and simplicity, this paper has restricted discussion to problems with two classes. Given that the ROCIV graph uses a straightforward transformation of ROC space, it is likely that the results here would hold

for problems with more than two classes as well. In particular, Hand and Till (2001) show how the AUC of a multi-class problem may be calculated from pairwise two-class AUC measurements. Proving that this transformation holds also for the AUCIV remains as future work.

It is important to mention two caveats in adopting this transformation. First, while example costs may vary, ROC analysis requires that costs always be negative and benefits always be positive. For example, if a cost function were defined as $c(Y, p, x) = x - \$20$, with example x values ranging in $[0, 40]$, this would be violated for x in $[0, 20]$.

Second, incorporating error costs into the ROC graph in this way introduces an additional assumption into a researcher's testing environment. Traditional ROC graphs assume that the FP rate and TP rate metrics of the test population will be similar to those of the training population; in particular that a classifier's performance on random samples will be similar. This new formulation adds the assumption that the example costs will be similar as well. In other words, ROCIV curves assume that not only will the classifier continue to score instances similarly between the training and testing sets, but the costs and benefits of those instances will be similar between the sets too.

Adding this assumption partially violates the cost insensitivity of ROC curves. A standard ROC curve is insensitive to changes in both *intra*-class and *inter*-class error costs. A ROCIV curve remains insensitive to *intra*-class error cost variations but will be sensitive to *inter*-class cost variations. As such, a researcher using ROCIV curves should check inter-class error cost distributions in the training and testing environments to ensure that they are stable.

References

- Drummond, C., Holte, R.C., 2000. Explicitly representing expected cost: An alternative to ROC representation. In: Ramakrishnan, R., Stolfo, S. (Eds.), *Proc. KDD-2000*. ACM Press, pp. 198–207.
- Egan, J.P., 1975. *Signal Detection Theory and ROC Analysis*, Series in Cognition and Perception. Academic Press, New York.
- Elkan, C., 2001. The foundations of cost-sensitive learning, in: *IJCAI-01*, pp. 973–978. Available from: <<http://citeseer.ist.psu.edu/elkan01foundations.html>>.
- Fawcett, T., 2003. ROC graphs: Notes and practical considerations for researchers, Tech Report HPL-2003-4, HP Laboratories. Available from: <<http://www.purl.org/NET/tfawcett/papers/ROC101.pdf>>.
- Fawcett, T., 2006. An introduction to ROC analysis. *Pattern Recognition Lett.*, this issue, doi:10.1016/j.patrec.2005.10.010.
- Fawcett, T., Provost, F., 1997. Adaptive fraud detection. *Data Mining Knowledge Discovery* 1 (3), 291–316.
- Hand, D.J., Till, R.J., 2001. A simple generalization of the area under the ROC curve to multiple class classification problems. *Machine Learn.* 45 (2), 171–186.
- Hanley, J.A., McNeil, B.J., 1982. The meaning and use of the area under a receiver operating characteristic (ROC) curve. *Radiology* 143, 29–36.
- Hettich, S., Blake, C., Merz, C., 1998. UCI repository of machine learning databases, <http://www.ics.uci.edu/~mllearn/MLRepository.html>.
- Pazzani, M., Merz, C., Murphy, P., Ali, K., Hume, T., Brunk, C., 1994. Reducing misclassification costs. In: *Proc. 11th Internat. Conf. on Machine Learning*. Morgan Kaufmann, Los Altos, CA, pp. 217–225.
- Provost, F., Fawcett, T., 1997. Analysis and visualization of classifier performance: Comparison under imprecise class and cost distributions. In: *Proc. Third Internat. Conf. on Knowledge Discovery and Data Mining (KDD-97)*. AAAI Press, Menlo Park, CA, pp. 43–48.
- Provost, F., Fawcett, T., 1998. Robust classification systems for imprecise environments. In: *Proc. AAAI-98*. AAAI Press, Menlo Park, CA, pp. 706–713.
- Provost, F., Fawcett, T., 2001. Robust classification for imprecise environments. *Machine Learn.* 42 (3), 203–231.
- Sing, T., 2004. Learning localized rule mixtures by maximizing the area under the ROC curve, with an application to the prediction of HIV-1 coreceptor usage, Ph.D. thesis, Max-Planck-Institut für Informatik Saarbrücken, March.
- Swets, J., 1988. Measuring the accuracy of diagnostic systems. *Science* 240, 1285–1293.
- Swets, J.A., Dawes, R.M., Monahan, J., 2000. Better decisions through science. *Sci. Amer.* 283, 82–87.
- Turney, P.D., 2000. Types of cost in inductive concept learning, in: *Workshop on Cost-Sensitive Learning at the Seventeenth International Conference on Machine Learning*.
- Weinstein, M.C., Fineberg, H.V., 1980. *Clinical Decision Analysis*. W.B. Saunders Company, Philadelphia, PA.
- Witten, I., Frank, E., 2000. *Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations*. Morgan Kaufmann, San Francisco, Software. Available from: <<http://www.cs.waikato.ac.nz/~ml/weka/>>.
- Zou, K.H., 2002. Receiver operating characteristic (ROC) literature research. Available from: <<http://splweb.bwh.harvard.edu:8000/pages/ppl/zou/roc.html>>.