

Learning with Cost Intervals

Xu-Ying Liu and Zhi-Hua Zhou
National Key Laboratory for Novel Software Technology
Nanjing University, Nanjing, 210093, China
{liuxy, zhouzh}@lamda.nju.edu.cn

ABSTRACT

Existing cost-sensitive learning methods require that the unequal misclassification costs should be given as precise values. In many real-world applications, however, it is generally difficult to have a precise cost value since the user maybe only knows that one type of mistake is much more severe than another type, yet it is infeasible to give a precise description. In such situations, it is more meaningful to work with a *cost interval* instead of a precise cost value. In this paper we report the first study along this direction. We propose the CISVM method, a support vector machine, to work with cost interval information. Experiments show that when there are only cost intervals available, CISVM is significantly superior to standard cost-sensitive SVMs using any of the minimal cost, mean cost and maximal cost to learn. Moreover, considering that in some cases other information about costs can be obtained in addition to cost intervals, such as the distribution of costs, we propose a general approach CODIS for using the distribution information to help improve performance. Experiments show that this approach can reduce 60% more risks than the standard cost-sensitive SVM which assumes the expected cost is the true value.

Categories and Subject Descriptors

I.2.6 [Artificial Intelligence]: Learning; H.2.8 [Database Management]: Database Applications—*Data mining*

General Terms

Algorithms; Design; Experimentation

Keywords

Cost-sensitive learning, misclassification costs, cost intervals

1. INTRODUCTION

In real-world applications, such as medical diagnosis, direct marketing and intrusion detection, different classification errors often lead to different losses. For example, in

medical diagnosis, the loss of misdiagnosing a patient as healthy is much more serious than misclassifying a healthy person as sick, because the former may lead to the loss of a life. Unfortunately, traditional research assumes that all the classification errors will result in the same loss. Thus, standard classification methods try to minimize the number of errors rather than the total cost. To deal with unequal costs, cost-sensitive learning has attracted much attention in recent years [14, 8, 19, 20, 17, 21, 13, 7, 12, 11].

Existing cost-sensitive learning methods work with fixed cost values. The cost information is provided by domain knowledge and is *taken as* a group of precise values. The classifiers will then be well tuned to reduce the total cost w.r.t. these particular cost values. However, in many real-world situations, although the user knows that one type of mistake is more severe than another type, it may be difficult to specify a precise cost value. The aspects that can cause imprecise costs include but not limited to:

- Inherent impreciseness. Some information is naturally stochastic, e.g., one may donate \$1 or \$5 randomly.
- Unknown information. Sometimes we can't know everything about a system. For example, customers' SSN and salary information can't be obtained because of privacy.
- Variations in modeling. In the process of modeling risk, the approach may sample data, transform input space, and use different methods or parameters. Due to these variations, the model may not provide precise assessment for costs.
- Expert opinions. Experts may have different opinions. And sometimes, experts can just give fair estimates of real risks.
- Dynamic environments. Environments always change. The risk may change with time, e.g., due to the appearance of a new competitor. While in some applications, the same model will be used in different scenarios.

In many situations, though precise information is not available, other useful information can be obtained. One of the most common and practical form to represent an imprecise value is to bound it with an interval. We call cost information represented by intervals *cost intervals*. There are at least three possible ways to obtain cost intervals:

- Natural upper and lower bounds. In some applications, costs naturally have upper and lower bounds. For example, the buying price of a stock is \$3, and someone determines to hold it when the price is between \$5 and \$8, then the profit is between \$2 and \$5 if the decision is right.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

KDD'10, July 25–28, 2010, Washington, DC, USA.

Copyright 2010 ACM 978-1-4503-0055-1/10/07 ...\$10.00.

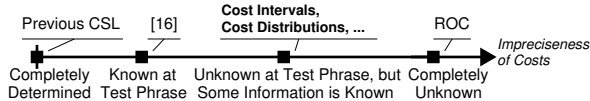


Figure 1: Related Work of Imprecise Costs.

- Transforming from confidence intervals. When there are no clear upper and lower bounds of cost, we can use a confidence interval of cost instead. For example, the 95% confidence interval indicates cost will appear in the interval with a probability of 0.95. Then it’s reasonable to use the lower and upper bounds of this interval to approximate the real situation. To extreme, the 100% confidence interval indicates the true cost value is definitely in it and it can be used as cost interval if its range is appropriate. Confidence intervals are easy to get in many applications. In medicine, the inclusion of confidence interval in analysis results has been a standard [9, 18]. For example, it was reported in 2007 that [18], compared with women having non-proliferative lesions, the risk for breast cancer is greater for women with atypical ductal hyperplasia of (IRR=5.0; 95% CI=2.26-11.0). It means that the incidence rate ratio is *about* 5.0 and the 95% confidence interval is [2.26, 11.0].
- Expert opinions. Sometimes, we can ask experts to provide cost intervals according to their experiences. It will be much easier for experts to give a cost interval than a “precise” cost value.

In this paper we report the first study along the direction of learning with cost intervals. Current cost-sensitive methods can be applied only when precise costs are given. To the best of our knowledge, there is no method learning with cost intervals. A related work is [16], which considers the situation that costs change over time. But it assumes that the true cost values are known at the time of classification. In our assumption, the true cost values are always unknown. ROC curve [2] can evaluate classifiers under imprecise class distributions or misclassification costs. Larger AUC (area under ROC curve) indicates better performance. This essentially assumes that nothing whatsoever is known about the relative severity of costs, which is a very rare situation in real-world problems. In our problem settings, cost intervals are known. The relationship to related work is shown in Figure 1.

An intuitive way to handle cost intervals is that, taking some value in a cost interval as the true cost, such as the minimal value, mean value, or maximal value, and then applying standard cost-sensitive learning methods. We will show theoretically that, they are not the best solutions. We propose CISVM (Cost-Interval-sensitive SVM), a support vector machine, to work with cost intervals. Experiments show that CISVM is significantly superior to all of them.

We also consider the situation in which more cost information is known in addition to cost intervals. For example, we can get distribution information of costs in some applications. We call them *cost distributions*. We propose a general approach CODIS to exploit the COSt DIStribution information to help improve performance. We will show theoretically that CODIS can minimize the expected risk, while standard cost-sensitive methods taking the expected cost value

as the true cost could not. Experiments show that CODIS can reduce 60% more risks than standard cost-sensitive SVM.

Section 2 formulates the problem of learning with cost intervals. Section 3 presents the CISVM method and Section 4 gives some analysis. Section 5 considers the situation in which cost distributions are known. Section 6 reports empirical results and Section 7 concludes the paper.

2. PROBLEM FORMULATION

In binary classification problems, suppose X is an input space and Y is a class label space with $y \in \{-1, +1\}$. Data are drawn i.i.d. from distribution $Pr(X, Y)$ and a training set is $S = \{(x_i, y_i)\}_{i=1}^n$. We consider *class-dependent costs*. That is, costs are only dependent on classes. The cost for any example of class y being misclassified to class y' is $cost(y, y')$. Thus, costs can be represented by a cost matrix. Since there are only 2 classes, assume that correct prediction costs 0, and let c_+ and c_- denote the cost of misclassifying a positive and negative example, respectively. Assume that positive class has higher cost, i.e., $c_+ \geq c_-$. Since the optimal decisions are unchanged when a cost matrix is multiplied by a positive constant [8], we can simplify the costs by fixing the cost of negative class so that we only need to consider the cost of positive class, i.e., $c_- = 1, c_+ = c (c \geq 1)$.

For a classifier $h : X \rightarrow Y$, the loss of misclassifying an example (x, y) (denoted by 0-c loss) is:

$$l(c, h(x), y) = cI(h(x) \neq y \wedge y = +) + I(h(x) \neq y \wedge y = -), \quad (1)$$

where $I(a) = 1$ if $a = true$ and 0 otherwise.

Previous cost-sensitive learning assumes that cost c is a fixed value. The goal is to find a classifier h^* minimizing the empirical risk

$$\tilde{R}(h, c) = \sum_{i=1}^n l(c, h(x_i), y_i). \quad (2)$$

When the true cost is unknown and a cost interval is available, we assume that the unique true cost C^* is a random value in the interval $[C_{min}, C_{max}]$. The goal is to learn a classifier H^* minimizing the *true risk* $\tilde{R}^* = \tilde{R}(h, C^*)$. Unfortunately, since the true cost is unknown, \tilde{R}^* can’t be obtained to guide the learning process. To overcome this difficulty, some risk \tilde{R}_s can be used instead to learn a classifier, which is in fact determined by some cost C_s , i.e., $\tilde{R}_s = \tilde{R}(h, C_s)$. Since in general \tilde{R}_s and C_s are different from the true risk \tilde{R}^* and the true cost C^* , respectively, we call them “surrogate risk” and “surrogate cost”. By minimizing surrogate risk \tilde{R}_s , the optimal classifier h_s^* is expected to minimize the true risk \tilde{R}^* . But this is infeasible since \tilde{R}^* is unknown. However, since the true cost can be any value in the cost interval, it is expected that any possible risk of h_s^* should be small enough. Obviously, not all surrogate risk will be good enough for this purpose, so an appropriate surrogate cost C_s must be carefully chosen. Thus, in order to learn a classifier making any possible risks small enough, we can formulate the problem of learning with cost intervals as Eq. 3, by considering C_s as a variable for learning.

$$\begin{aligned} \min_{h, C_s} \quad & \tilde{R}(h, C_s) \\ \text{s.t.} \quad & p(\tilde{R}(h, c) < \epsilon) > 1 - \delta, \forall c \in [C_{min}, C_{max}] \\ & C_{min} \leq C_s \leq C_{max}. \end{aligned} \quad (3)$$

Note that, it doesn’t matter for C_s to be different from C^* as long as $\tilde{R}(h, c)$ is small enough for all possible c values.

3. THE CISVM APPROACH

Let $C_\mu = 0.5(C_{min} + C_{max})$ denote the mean cost, and $d = 0.5(C_{max} - C_{min})$. Then, a random value $c \in [C_{min}, C_{max}]$ can be represented as $c = C_\mu + \varepsilon$, where ε is a random value in interval $[-d, d]$ and $|\varepsilon| \leq d$.

Since there are infinite constraints in Eq. 3, it is intractable to get optimal solutions. To overcome this difficulty, CISVM tries to solve a relaxation with a small number of informative constraints. Considering the worst case risk is the upper bound of the risks w.r.t. any c in $[C_{min}, C_{max}]$, the optimal solution of minimizing it can make all the constraints in Eq. 3 hold. That is, if $h_s^* = \arg \min_h \sup_c \tilde{R}(h, c)$ and $\sup_c \tilde{R}(h_s^*, c) = \epsilon$, then $\tilde{R}(h_s^*, c) < \epsilon$ holds for any c . So, the worst case risk is appropriate to be used as surrogate risk $\tilde{R}(h, C_s)$ to guide the learning process. However, the worst case risk could be far away from the true risk. When $C^* = C_{min}$, it overestimates the risk the most. So the optimal solution of minimizing the worst case risk could not make the true risk small enough sometimes. That is, if $H^* = \arg \min_h \tilde{R}(h, C^*)$ and $\tilde{R}(H^*, C^*) = \epsilon^*$, then $\epsilon > \epsilon^*$ could happen. CISVM overcomes this difficulty by minimizing another risk $\tilde{R}(h, C')$ in the meanwhile to avoid overfitting to surrogate risk, where $C' \neq C_s$. If the distribution of cost were known, the expected risk could be suitable to evaluate a classifier in imprecise environments. Unfortunately, only cost intervals are known. However, in this situation, the risk w.r.t. the mean cost C_μ (called ‘‘the mean risk’’¹) has the smallest maximal distortion of the true risk (see Section 4), so it is the best choice to reflect how good a classifier performs on the entire interval. Therefore, CISVM works by minimizing the worst case risk and the mean risk in the meanwhile. Though similar to the worst case risk, the best case risk is often used for analysis in imprecise environments, it is too optimistic to guarantee other possible risks are small enough, so it is not a good candidate for either of $\tilde{R}(h, C_s)$ or $\tilde{R}(h, C')$.

Assuming that the prediction function is $f = w^T x + b$, CISVM utilizes a surrogate loss in the following form:

$$L(C_p, f(x), y) = I_{y=+}[C_p - yf(x)]_+ + I_{y=-}[1 - yf(x)]_+, \quad (4)$$

where $[a]_+ = \max(a, 0)$ and $I_a = I(a)$. It means that, the loss for a negative and positive example is $L^- = [1 - yf(x)]_+$ and $L^+ = [C_p - yf(x)]_+$, respectively. We use $L(C_p)$ for $L(C_p, f(x), y)$ when this will not cause confusion.

As we know, for fixed costs, existing methods gain cost sensitivity by introducing a parameter to bias toward expensive class. For example, sampling-based methods [8, 21] sample more examples of the expensive class. Instance-weighting-based methods [15] give higher weights to the expensive class examples. Thresholding-based methods [6, 8, 21] move decision threshold toward the inexpensive class. Other methods such as cost-sensitive SVMs [3] and cost-sensitive boosting methods [12] also have such parameters to control the tradeoff between different classes. We call this parameter *the cost parameter*. For fixed costs, it is reasonable to let the cost parameter equal to the true cost. Some work also indicated the best value of the cost parameter is usually unequal to the true cost [5, 13]. Chawla et al. [4] further proposed a method to automatically find the best value according to the results on a validation set. Similarly,

¹It is not exactly the mean risk since the distribution is unknown. We call it ‘‘the mean risk’’ for convenience.

a cost parameter is needed in the cost interval learning problem to achieve cost sensitivity, which is C_p in Eq. 4 to make tradeoff between different classes.

Let $\tilde{R}_L(C_p)$ denote the empirical loss on training set S :

$$\tilde{R}_L(C_p) = \sum_{i=1}^n L(C_p, f(x_i), y_i). \quad (5)$$

For the surrogate loss in form of Eq. 4, the worst case risk is $\tilde{R}_L(C_{max})$, which is guaranteed by Theorem 1. And the mean risk is $\tilde{R}_L(C_\mu)$. In addition, the best case risk is $\tilde{R}_L(C_{min})$.

THEOREM 1. *Suppose C_p is a random value in interval $[C_{min}, C_{max}]$, then the worst case risk w.r.t. loss function $L(C_p)$ is achieved when $C_p = C_{max}$:*

$$\tilde{R}_L(C_{max}) = \sup_{C_p} \tilde{R}_L(C_p). \quad (6)$$

PROOF.

$$\begin{aligned} \sup_{C_p} \tilde{R}_L(C_p) &= \sup_{C_p} \sum_{i=1}^n L(C_p, f(x_i), y_i) \\ &= \sup_{C_p} \sum_{y_i=+} [C_p - y_i f(x_i)]_+ + \sum_{y_i=-} [1 - y_i f(x_i)]_+ \\ &= \sup_{\varepsilon} \sum_{y_i=+} [C_\mu + \varepsilon - y_i f(x_i)]_+ + \sum_{y_i=-} [1 - y_i f(x_i)]_+ \\ &= \sum_{y_i=+} [C_\mu + \sup(\varepsilon) - y_i f(x_i)]_+ + \sum_{y_i=-} [1 - y_i f(x_i)]_+ \\ &= \sum_{y_i=+} [C_\mu + d - y_i f(x_i)]_+ + \sum_{y_i=-} [1 - y_i f(x_i)]_+ \\ &= \sum_{y_i=+} [C_{max} - y_i f(x_i)]_+ + \sum_{y_i=-} [1 - y_i f(x_i)]_+ \\ &= \sum_{i=1}^n L(C_{max}, x_i, y_i) = \tilde{R}_L(C_{max}). \quad \square \end{aligned}$$

COROLLARY 1. *The best case risk w.r.t. loss function $L(C_p)$ is achieved when $C_p = C_{min}$, which is $\tilde{R}_L(C_{min})$.*

If the optimal solutions of minimizing $\tilde{R}(h, C_s)$ and $\tilde{R}(h, C')$ were identical, the relaxed problem that CISVM tries to solve would be degenerated to minimizing $\tilde{R}(h, C_s)$ only, and thus its optimal solution would not be good enough for the original learning problem in Eq. 3. In fact, minimizing the worst case risk $\tilde{R}_L(C_{max})$ is not equivalent to minimizing the mean risk $\tilde{R}_L(C_\mu)$.

COROLLARY 2. *Suppose n_+ is the size of positive class. Then, $\tilde{R}_L(C_\mu) + d \cdot n_+ \geq \sup_{C_p} \tilde{R}_L(C_p)$. Minimizing $\tilde{R}_L(C_\mu)$ is equivalent to minimizing an upper bound of $\tilde{R}_L(C_{max})$, since $d \cdot n_+$ is a constant.*

PROOF.

$$\begin{aligned} \sup_{C_p} \tilde{R}_L(C_p) &= \sum_{y_i=+} [C_\mu + d - y_i f(x_i)]_+ + \sum_{y_i=-} [1 - y_i f(x_i)]_+ \\ &\leq \sum_{y_i=+} ([C_\mu - y_i f(x_i)]_+ + d) + \sum_{y_i=-} [1 - y_i f(x_i)]_+ \\ &= \tilde{R}_L(C_\mu) + d \cdot n_+. \quad \square \end{aligned}$$

Here, equality holds only when $yf(x) \leq C_\mu$ for all positive examples. Since this requires almost all positive examples be support vectors, equality does not hold in general.

To minimize the worst case risk $\tilde{R}_L(C_{max})$ and the mean risk $\tilde{R}_L(C_\mu)$ in the meanwhile, CISVM firstly learns a variation of support vector machine to minimize $\tilde{R}_L(C_{max})$, and

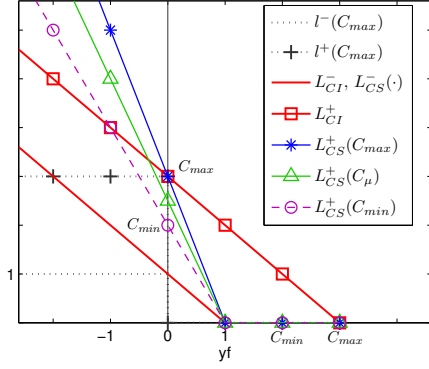


Figure 2: 5 loss functions l , L_{CI} , $L_{CS}(C_{max})$, $L_{CS}(C_{\mu})$, $L_{CS}(C_{min})$ are shown with each is plotted for positive example (marked by “+”) and negative example (marked by “-”) separately. $L_{CS}(\cdot)$ represents for any of $L_{CS}(C_{max})$, $L_{CS}(C_{\mu})$, and $L_{CS}(C_{min})$. L_{CI} is a convex upper bound of the true loss l .

then minimizes $\tilde{R}_L(C_{\mu})$ by parameter selection. In particular, to minimize the worst case risk, we set the surrogate loss of CISVM as:

$$\begin{aligned} L_{CI}(f(x), y) &= L(C_{max}, f(x), y) \\ &= I_{y=+}[C_{max} - yf(x)]_+ + I_{y=-}[1 - yf(x)]_+. \end{aligned} \quad (7)$$

L_{CI} is a convex upper bound of 0-c real loss function $l(c)$ in Eq. 1 for any cost c . See Figure 2 for an illustration. Then, CISVM minimizes the regularized loss of L_{CI} :

$$\begin{aligned} \min_{w, b, \xi \geq 0} \quad & \|w\|^2/2 + \lambda \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i(w^T \phi(x_i) + b) \geq C_{max} - \xi_i, \quad \forall i : y_i = +1 \\ & y_i(w^T \phi(x_i) + b) \geq 1 - \xi_i, \quad \forall i : y_i = -1 \end{aligned} \quad (8)$$

where $\phi(x)$ is a feature map induced by a kernel function. The dual form can be easily derived.

To minimize the mean risk $\tilde{R}_L(C_{\mu})$, CISVM selects λ and kernel parameters according to $\tilde{R}_L(C_{\mu})$. That is, among the classifiers achieving the optimal worst case risk, CISVM selects the one minimizing $\tilde{R}_L(C_{\mu})$. When a validation set is available, 0-c loss $l(C_{\mu})$ can be used instead of $L(C_{\mu})$ to assess the mean risk on the validation set.

Therefore, CISVM involves two parts: (1) minimizing the regularized worst case risk (Eq. 8) by learning a variation of SVM; (2) minimizing the mean risk by parameter selection on a validation set. The pseudo code is shown in Alg. 1. A cost interval is determined when C_{max} and C_{μ} are fixed, so CISVM provides a unique solution for each cost interval.

4. ANALYSIS

This section gives some analysis to compare existing cost-sensitive methods with CISVM on handling cost intervals.

Cost-sensitive SVM (CSSVM) [3] can only work with fixed costs. Its loss function for class-dependent cost is:

$$L_{CS}(C_p, f(x), y) = I_{y=+}C_p[1 - yf(x)]_+ + I_{y=-}[1 - yf(x)]_+, \quad (9)$$

where, C_p is the cost parameter. For a given cost c , $C_p = c$.

Intuitively, we can enable CSSVM to handle cost intervals by assuming that the true cost is a particular value in inter-

Algorithm 1 CISVM Algorithm

Input:

- $S = \{(x_i, y_i)\}_{i=1}^n$: training set, with $y_i \in \{-1, +1\}$
 - V : validation set
 - $[C_{min}, C_{max}]$: cost interval
 - $\{(\lambda, k)\}$: a set of parameters, where k is the kernel parameters
- 1: For each (λ_i, k_i) , train a SVM $f_i = w^T \phi(x) + b$ by solving the optimization problem in Eq. 8.
 - 2: Let $C_{\mu} = 0.5(C_{min} + C_{max})$ be the mean cost.
 - 3: Evaluate the mean risk on validation set V for each f_i :

$$\tilde{R}_V(f_i, C_{\mu}) = \sum_{(x, y) \in V} l(C_{\mu}, f_i(x), y), \quad l \text{ in Eq. 1}$$

Output: $H(x) = \arg \min_{f_i} \tilde{R}_V(f_i, C_{\mu})$

val $[C_{min}, C_{max}]$. In this paper, we consider three versions of CSSVM: CSMIN, CSMEAN and CSMAX, by setting C_p as C_{min} , C_{μ} and C_{max} , respectively. Their loss functions are shown in Figure 2, which are different from L_{CI} , the loss function of CISVM.

Let $\tilde{R}_{LCS}(C_p)$ be the empirical loss on training set S :

$$\tilde{R}_{LCS}(C_p) = \sum_{i=1}^n L_{CS}(C_p, f(x), y). \quad (10)$$

Then similar to the analysis in the previous section, for the loss function in form of Eq. 9, the worst case risk is $\tilde{R}_{LCS}(C_{max})$, the mean risk is $\tilde{R}_{LCS}(C_{\mu})$, and the best case risk is $\tilde{R}_{LCS}(C_{min})$. Hence, CSMAX minimizes the worst case risk $\tilde{R}_{LCS}(C_{max})$, CSMEAN minimizes the mean risk $\tilde{R}_{LCS}(C_{\mu})$, and CSMIN minimizes the best case risk $\tilde{R}_{LCS}(C_{min})$. The optimal solutions of minimizing the worst case risk and the mean risk are not equivalent.

LEMMA 1. *Suppose C_p is a random value in $[C_{min}, C_{max}]$, then the worst case risk w.r.t. loss function $L_{CS}(C_p)$ is achieved when $C_p = C_{max}$:*

$$\tilde{R}_{LCS}(C_{max}) = \sup_{C_p} \tilde{R}_{LCS}(C_p). \quad (11)$$

COROLLARY 3. *The best case risk w.r.t. loss function $L_{CS}(C_p)$ is achieved when $C_p = C_{min}$, which is $\tilde{R}_{LCS}(C_{min})$.*

COROLLARY 4. *$\tilde{R}_{LCS}(C_{\mu}) \leq \sup_{C_p} \tilde{R}_{LCS}(C_p)$. Equality holds only when $yf(x) \geq 1$ for all positive examples. Since this requires many positive examples being non-support-vectors, equality does not hold in general.*

The difference between CSMAX, CSMEAN CSMIN and CISVM is that they use different surrogate risks to guide the learning process. However, no matter what surrogate risk is used, it is generally unequal to the true risk, so there is always distortion of the true risk. In a good cost interval learning method, the maximal distortion of the true risk should be as small as possible. Based on this idea, we define the robustness of a cost-sensitive loss function on a cost interval as follows:

DEFINITION 1. *Suppose $\mathcal{L}(C_p, h(x), y)$ is a cost-sensitive surrogate loss, where C_p is the cost parameter. $l(C_p, h(x), y)$ is 0-c real loss in Eq. 1. Examples are drawn from $Pr(X, Y)$. Then, the distortion of $\mathcal{L}(C_p)$ on cost interval $[C_{min}, C_{max}]$*

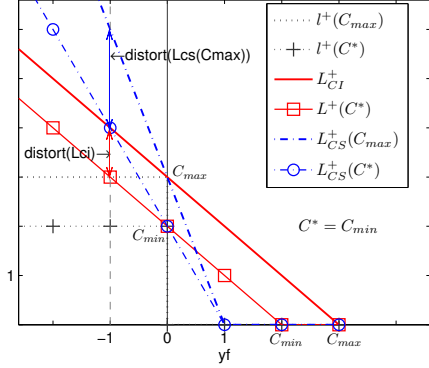


Figure 3: Comparison of robustness of CISVM and CSMAX. There are 3 loss functions: 0-c loss $l(C_{max})$, $L(C_{max})$ (i.e., L_{CI}), $L_{CS}(C_{max})$. The figure shows the loss for positive examples: $l^+(C_{max})$, L_{CI}^+ , $L_{CS}^+(C_{max})$. When the true cost $C^* = C_{min}$, the distortion of the true risk is largest for all of these loss functions, $\text{distort}(L_{CI}) = C_{max} - C_{min}$ and $\text{distort}(L_{CS}(C_{max})) = 2(C_{max} - C_{min})$.

is defined as the worst case distortion of $\mathcal{L}(C_p)$. Let $C^* \in [C_{min}, C_{max}]$ denote the true cost, then:

$$\text{distort}(\mathcal{L}(C_p)) = \sup_{C^*, h(x), y} |\mathcal{L}(C^*, h(x), y) - \mathcal{L}(C_p, h(x), y)|.$$

$\inf_{C_p} \text{distort}(l(C_p))$ can be used as a baseline when quantifying robustness, since l is the real loss and will not introduce additional distortion of risk.

The robustness of $\mathcal{L}(C_p)$ is then defined as

$$\text{robust}(\mathcal{L}(C_p)) = \frac{\inf_{C_p} \text{distort}(l(C_p))}{\text{distort}(\mathcal{L}(C_p))}. \quad (12)$$

The larger the robust value, the better the robustness. $\text{robust} = +\infty$ when there is no distortion.

THEOREM 2. The robustness of loss functions of CISVM and CSMAX is 0.5 and 0.25, respectively. Therefore, CISVM is more robust than CSMAX.

PROOF.

$$\begin{aligned} \inf_{C_p} \text{distort}(l(C_p)) &= \inf_{C_p} \sup_{C^*} |C_p - C^*| \\ &= 0.5(C_{max} - C_{min}) = d. \end{aligned}$$

$$\begin{aligned} \text{distort}(L_{CI}) &= \sup_{C^*, f(x), y} |[C^* - yf(x)]_+ - [C_{max} - yf(x)]_+| \\ &= \sup_{C^*} |C^* - C_{max}| = C_{max} - C_{min} = 2d. \end{aligned}$$

$$\begin{aligned} \text{distort}(L_{CS}(C_{max})) &= \sup_{C^*, f(x), y} |(C^* - C_{max})[1 - yf(x)]_+| \\ &= 2(C_{max} - C_{min}) = 4d. \quad \square \end{aligned}$$

Both CISVM and CSMAX try to minimize the worst case risk, but CISVM is more robust than CSMAX, which is shown in Theorem 2. The robustness of loss functions of CISVM and CSMAX is 0.5 and 0.25, respectively. See Figure 3 for an illustration. Similarly, for CSMEAN and CSMIN, the robust value is 0.5 and 0.25, respectively. CISVM also tries to minimize the mean risk whose loss function $L(C_\mu)$ has a robust value of 1.

Thus, we contribute the success of CISVM to two aspects:

- (1) it minimizes the worst case risk and the mean risk in the meanwhile;
- (2) it is very robust.

5. UTILIZING COST DISTRIBUTION

In addition to cost intervals, sometimes we can get more information on the costs. In some applications, experts can provide the cost distribution information base on their experiences. The normal and uniform distributions are very popular and can be easily recognized from experience. For complex distributions, we can build models to assess costs. Then the values provided by different models can be regarded as samples from the underlying cost distribution. When cost distributions are known, many difficulties encountered in learning with cost intervals no more exist.

Let $R(h, c)$ be the expected risk of loss l for some cost c :

$$\begin{aligned} R(h, c) &= E_{P_{r(X,Y)}}[l(c, h(x), y)] \\ &= p_+ \cdot fn \cdot c + p_- \cdot fp, \end{aligned} \quad (13)$$

where $p_+ = p(y = +)$ and $p_- = p(y = -)$. fn is false negative rate, and fp is false positive rate.

Assume that cost c is independently drawn from distribution v with domain C^2 , which is independent of X . Then the goal is to find a classifier h minimizing the expected risk over v :

$$\begin{aligned} R_{CD}(h, v) &= E_{c \sim v}[R(h, c)] \\ &= E_{c \sim v} E_{P_{r(X,Y)}}[l(c, h(x), y)]. \end{aligned} \quad (14)$$

Note that, this is different from learning with example-dependent costs [19, 20, 3, 11], where examples $\{(x_i, y_i, c_i)\}$ are drawn from distribution D with domain $X \times Y \times C$, and costs are often assumed as $c_i = \text{cost}(x_i, y_i, y'_i)$, where y'_i is the prediction for x_i . The goal of example-dependent cost-sensitive learning is to minimize:

$$E_{(x_i, y_i, c_i) \sim D}[c_i I(h(x_i) \neq y_i)]. \quad (15)$$

In our setting, costs are independent of examples. That is, distribution v is not related to X . When c changes to c' , (e.g., due to time change), the cost of any misclassified positive example will be c' .

When a classifier h is learned, R_{CD} should be used for evaluation. It is equivalent to $R(h, E[c])$, see Eq. 16. That is, the expected cost $E[c]$ should be used to evaluate risk.

$$\begin{aligned} R_{CD}(h, v) &= E_{c \sim v}[R(h, c)] \\ &= E_{c \sim v}[p_+ \cdot fn \cdot c + p_- \cdot fp] \\ &= p_+ \cdot fn \cdot E[c] + p_- \cdot fp = R(h, E[c]). \end{aligned} \quad (16)$$

In the above we assume known cost distribution. When there are only samples of the underlying cost distribution available, let $S_c = \{c_i\}_{i=1}^{|c|}$ be the sample set of costs, then the goal is to minimize the empirical form of R_{CD} :

$$\tilde{R}_{CD}(h, S_c) = \frac{1}{|c|} \sum_{i=1}^{|c|} E_{P_{r(X,Y)}}[l(c_i, h(x), y)].$$

Eq. 16 shows for a fixed classifier h , the expected risk over distribution v is the risk w.r.t. the expected cost $E[c]$. An intuitive way to utilize cost distributions is to take $E[c]$ as the true cost and then exploit standard cost-sensitive methods. Unfortunately, there is no guarantee that this solution can minimize the expected risk. $E_{c \sim v}[R(h, c)] = R(h, E[c])$ holds only when h is independent of c . A learned classifier h is fixed and is therefore independent of c , so Eq. 16 holds.

²There may be other domains, e.g., time. Here, we only consider the marginal distribution of C .

Algorithm 2 CODIS Algorithm

Input:

- $S = \{(x_i, y_i)\}_{i=1}^n$: training set, with $y_i \in \{-1, +1\}$
- v : cost distribution
- A : learning algorithm for example-dependent costs
- T : learning rounds

```

1: for  $t = 1$  to  $T$  do
2:    $\hat{S} = \emptyset$ 
3:   for  $i = 1$  to  $n$  do
4:     Draw a cost  $c_i$  independently from  $v$ 
5:      $\hat{S} = \hat{S} \cup (x_i, y_i, c_i)$ .
6:   end for
7:    $h_t = A(\hat{S})$ : learn a classifier using  $A$  and  $\hat{S}$ .
8: end for

```

Output: $H(x) = \sum_{t=1}^T h_t(x)$

But in the learning process, h is a function in hypothesis space \mathcal{H} , so it is dependent on cost c .

Recall the cost parameter C_p described in Section 3. Cost-sensitive methods use C_p to make tradeoff between different classes according to their costs. Therefore, C_p can be regarded as a function dependent on the true cost c : $C_p = g(c)$. Since h is dependent on C_p in the learning process, it is a function of c : $h = h(C_p) = h(g(c))$. We then have:

PROPOSITION 1. v is a cost distribution. h is a function in hypothesis space \mathcal{H} to be learned. It has cost parameter C_p : $h(x) = h(x, C_p)$. C_p is dependent on the true cost c : $C_p = g(c)$. Then $h(x) = h(x, g(c))$, and the following holds when $g(c) \neq c$ or h is a nonlinear function of $g(c)$.

$$E_{c \sim v}[R(h(g(c)), c)] \neq R(h(E[c]), E[c]). \quad (17)$$

Inequality holds in general, since h is generally a nonlinear function of $g(c)$, and there is no guarantee that $g(c)$ should be c even for fixed costs.

Therefore, standard methods taking the expected cost $E[c]$ as the true cost and then minimizing the corresponding risk can not minimize $E_{c \sim v}[R(h(g(c)), c)]$. However, it can be easily derived that cost distribution learning problem can be reduced to a special case of example-dependent cost-sensitive learning problem, which is guaranteed by Theorem 3.

THEOREM 3. Suppose example (x, y) is i.i.d. drawn from $Pr(X, Y)$, and cost c is i.i.d. drawn from $v(C)$. Suppose Pr and v are independent. Then new example (x, y, c) can be considered to be i.i.d. drawn from a joint distribution D with domain $X \times Y \times C$, with $p_D(x, y, c) = p_{Pr}(x, y)p_v(c)$, where $p_A(a)$ denotes the pdf of sample a drawn from distribution A . We have the following holds for any function h with cost parameter $C_p = g(c)$:

$$E_{c \sim v}[R(h(g(c)), c)] = E_{(x, y, c) \sim D}[l(c, h(x, g(c)), y)]. \quad (18)$$

PROOF.

$$\begin{aligned}
E_{c \sim v}[R(h(g(c)), c)] &= E_{c \sim v} E_{(x, y) \sim Pr}[l(c, h(x, g(c)), y)] \\
&= E_{(x, y) \sim Pr, c \sim v}[l(c, h(x, g(c)), y)] \\
&= E_{(x, y, c) \sim D}[l(c, h(x, g(c)), y)]. \quad \square
\end{aligned}$$

Cost distribution v is independent of X by definition, so Theorem. 3 holds. Based on this theorem, CODIS, a general method exploiting example-dependent cost-sensitive methods is proposed to utilize cost distributions. Firstly, a cost sample c_i is drawn from v (or provided by a risk model) for

Table 1: Cost Intervals

Width	Counts	Cost Intervals
3	10	[1, 3], [4, 6], [7, 9], [10, 12], [13, 15], [16, 18], [19, 21], [22, 24], [25, 27], [28, 30]
5	6	[1, 5], [6, 10], [11, 15], [16, 20], [21, 25], [26, 30]
11	5	[1, 11], [5, 15], [10, 20], [15, 25], [20, 30]
15	4	[1, 16], [5, 20], [10, 25], [15, 30]

each example (x_i, y_i) in training set S to form a new example set $\hat{S} = \{(x_i, y_i, c_i)\}_{i=1}^n$. Secondly, a standard example-dependent cost-sensitive method A is called to learn a classifier minimizing the risk of \hat{S} . Furthermore, to reduce the variance caused by sampling from v , we can draw cost multiple times from v for a single example (x_i, y_i) since v and Pr are independent. Therefore, we can repeat the first two steps several times and ensemble all the classifiers. The pseudo code is shown in Alg. 2. We use 0-1 voting for ensemble. It is possible to use other kinds of ensemble to get even better performance and which will be studied in the future. The learning algorithm A can be any standard cost-sensitive method which works with example-dependent costs, such as Costing [20] and example-dependent cost-sensitive SVM [3], which both work by assigning each example a weight proportional to its own cost c_i .

6. EXPERIMENTS

6.1 Settings

We compare CISVM and CSMIN, CSMEAN, CSMAX, with standard SVM as baseline on 15 binary UCI data sets [1]: *breast-cancer*, *breast-w*, *credit-a*, *credit-g*, *diabetes*, *german*, *haberman*, *heart-statlog*, *ionosphere*, *liver*, *sonar*, *spambase*, *spect*, *spectf* and *tic-tac-toe*. We treat each class of a data set as positive class in turn. Thus, we actually have 30 tasks³. Since class imbalance may affect cost-sensitive methods [10], a balanced subset is used for each task to focus on the difficulty caused by imprecise costs. 25 cost intervals are used, which are shown in Table 1. 30 times stratified hold-out tests are carried out, with 2/3 data as training set and 1/3 data as test set. All methods use RBF kernels. Parameters are selected for every cost interval by 5-fold cross validation on training data from $\lambda \in \{0.01, 0.1, 1, 10, 100\}$, and kernel parameter σ is selected as $\{1/10, 1/2, 1, 2, 10\}$ times the mean squared distance of training set. There is no guidelines on how to select parameters for standard cost-sensitive methods to handle cost intervals. To test CISVM rigorously, we use the mean risk to select parameters for *all methods*. We use standard SVM-style implementation by invoking MOSEK⁴.

Since the true cost is unknown, we need to evaluate how good a classifier performs in different situations. Here, quadri-section points of an interval are considered in turn as the true cost: $C_{min} + (C_{max} - C_{min}) \times \{0, 0.25, 0.5, 0.75, 1\}$. We call them test points P1-P5.

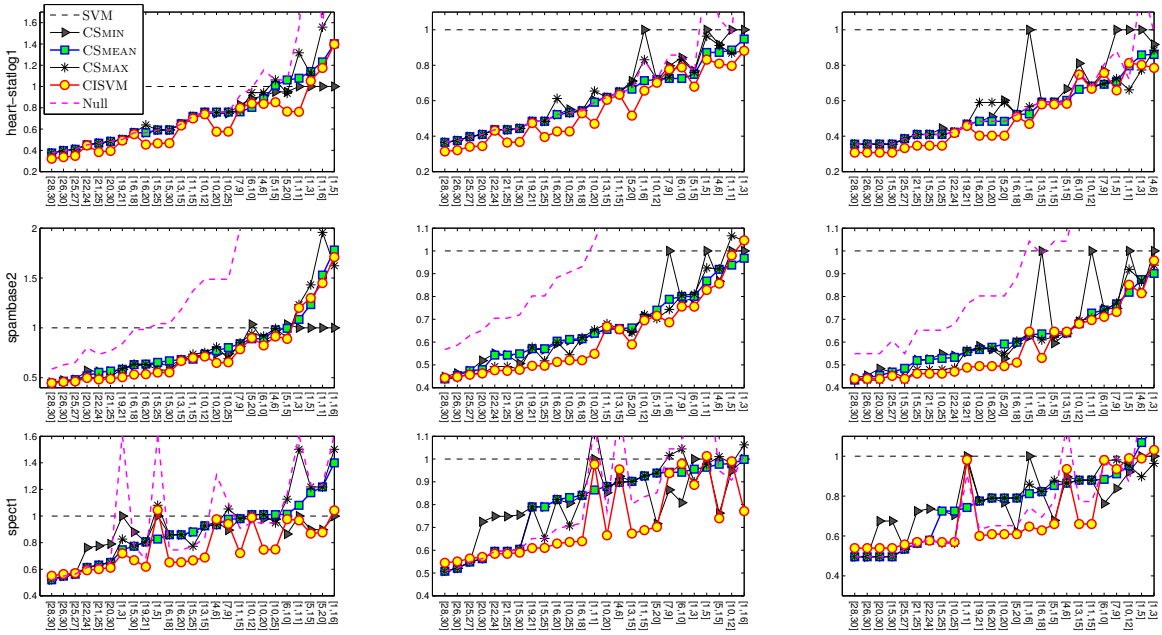
6.2 Results

Due to page limit, we can only show part of the results⁵. Figure 4 shows the performance of each method on two medical data (*heart-statlog1*, *spect1*) and a spam detection data

³Data sets marked by ‘1’ use the first class as negative class and by ‘2’ use the second class as negative class.

⁴The code is available at http://lamda.nju.edu.cn/code_CISVM.ashx

⁵More results will be reported in a longer version.



Test Point P1 (the true cost is C_{min}) Test Point P3 (the true cost is C_{μ}) Test Point P5 (the true cost is C_{max})

Figure 4: Some Results of Learning with Cost Intervals. X-axis shows cost intervals. Y-axis is the ratio of the loss of each method against that of SVM. The lower the value, the better the performance. A trivial classifier ‘Null’ which always predict positive is also shown.

Table 2: Win/tie/loss counts of method-in-row vs. method-in-column after t -tests at 95% significance level. Bold/Italic highlights significant better/worse of sign-tests results at 95% significance level.

	SVM	CSMIN	CSMEAN	CSMAX
P1 CSMIN	457/271/22			
P1 CSMEAN	484/162/104	<i>83/487/180</i>		
P1 CSMAX	488/160/102	<i>107/461/182</i>	<i>92/522/136</i>	
P1 CISVM	521/96/133	<i>257/262/231</i>	316/248/186	303/259/188
P2 CSMIN	464/272/14			
P2 CSMEAN	533/175/42	<i>121/514/115</i>		
P2 CSMAX	525/186/39	<i>139/494/117</i>	<i>82/542/126</i>	
P2 CISVM	578/97/75	304/263/183	296/269/185	282/282/186
P3 CSMIN	466/273/11			
P3 CSMEAN	550/180/20	<i>143/509/98</i>		
P3 CSMAX	554/173/23	<i>159/495/96</i>	<i>77/560/113</i>	
P3 CISVM	593/93/64	304/271/175	276/292/182	260/309/181
P4 CSMIN	472/270/8			
P4 CSMEAN	554/180/16	<i>150/503/97</i>		
P4 CSMAX	561/176/13	<i>168/493/89</i>	<i>89/565/96</i>	
P4 CISVM	599/91/60	299/275/176	256/296/198	<i>232/327/191</i>
P5 CSMIN	475/268/7			
P5 CSMEAN	556/178/16	<i>150/503/97</i>		
P5 CSMAX	566/170/14	<i>173/488/89</i>	<i>100/556/94</i>	
P5 CISVM	606/85/59	299/270/181	252/294/204	<i>226/325/199</i>
All CSMIN	2334/1354/62			
All CSMEAN	2677/875/198	<i>647/2516/587</i>		
All CSMAX	2694/865/191	<i>746/2431/573</i>	<i>440/2745/565</i>	
All CISVM	2897/462/391	1463/1341/946	1396/1399/955	1303/1502/945

(*spambase2*) at P1, P3 and P5. Table 2 summaries the results of paired t -tests and sign-tests for each test point separately. Table 3 shows the loss of each method on the first set of data sets for some cost intervals, where we count the number of significantly best performance (SBP) [22] achieved by each method. SBPs are the best performances and there is no significant difference among them according to t -tests results. The SBP counts are further summarized in Table 4 for each cost interval and each test point.

- *At P1.* The true cost is C_{min} . CSMIN is significantly better than CSMEAN and CSMAX. It is not a surprise, since CSMIN learns with the true cost. CSMEAN and CSMAX use C_p larger than the true cost, which obviously causes the risk being overestimated. While CISVM is comparable to

CSMIN, and beats CSMEAN and CSMAX significantly. Note that, P1 is the most far point from CISVM’s cost parameter ($C_p = C_{max}$), but in general it does not prevent CISVM from effectively reducing risk. Some work indicated the best C_p value is usually unequal to the true cost [5, 13]. This might be the reason why CSMIN is not always the best at P1.

- *At P2.* None of the methods learns with the true cost. CSMIN is comparable to CSMEAN and CSMAX. CSMEAN performs better than CSMAX. While CISVM beats them all significantly.

- *At P3.* The true cost is C_{μ} , and CSMEAN learns with the true cost. It is not surprising that CSMEAN is significantly better than CSMIN and CSMAX. The latter two methods underestimate and overestimates the risk, respectively. While CSMEAN is significantly worse than CISVM.

- *At P4.* None of the methods learns with the true cost. CSMIN has the most inappropriate cost parameter among all methods, and the results show it is significantly worse than the others. CSMEAN and CSMAX are comparable. While CISVM is significantly better than CSMIN and CSMEAN, and comparable to or slightly better than CSMEAN.

- *At P5.* The true cost is C_{max} . Both CSMAX and CISVM learn with the true cost. Again, CSMIN is the worst one, due to its under-bias. CSMAX does not achieve to perform better than CSMEAN, but CISVM does. CISVM is comparable to or slightly better than CSMAX.

- Over all test points, CISVM beats all other methods significantly. The SBPs of CISVM are over 20% more than others (Table 4). CSMAX is better than CSMIN yet worse than CSMEAN. The reason might be that, though CSMEAN could not achieve the minimal worst case risk, it is more robust than CSMAX. The results of CSMIN show that min-

Table 5: Grouped Results of t -tests: CISVM vs. Others. Results in bold is significant in sign-tests.

		<i>Width</i> = 3	<i>Width</i> = 5	<i>Width</i> = 10	<i>Width</i> = 15
size		450/600/450	300/300/300	300/300/150	150/300/150
mean	CSMIN	33/43/23	44/32/25	45/32/23	61/15/23
	CSMEAN	38/38/24	34/38/28	36/35/29	48/25/27
[1, 10]	CSMAX	40/37/23	40/34/26	41/32/27	45/25/31
mean	CSMIN	39/34/27	40/31/29	41/32/26	36/32/31
	CSMEAN	40/35/25	44/29/27	44/32/23	48/26/26
(10, 20]	CSMAX	40/36/24	40/35/25	36/43/21	38/39/23
mean	CSMIN	30/48/22	35/42/23	45/33/23	39/33/27
	CSMEAN	25/49/26	28/48/23	25/52/23	37/40/23
(20, 30]	CSMAX	22/52/26	23/52/26	19/51/30	21/49/30

imizing the best case risk performs well only when the true cost is very close to the minimal cost.

To study how CISVM’s performance will change with different cost intervals, we split 25 cost intervals into 12 groups according to different width and means and show the t -tests results for each group in Table 5. Since there are different numbers of intervals contained in each group, we normalize the win/tie/lose counts for the ease of comparison, with total counts shown in the row *size*. It shows that, when the mean cost is larger than 20, the performance of CISVM decreases, yet it is still comparable to or better than other methods. While when the mean cost is smaller than 20, the performance of CISVM is very similar for intervals with different width and means.

From the experiments we conclude that: (1) The success of CISVM implies that learning with cost intervals should be solved specially, and standard cost-sensitive methods could not handle it well. (2) CISVM is generally the best method over the whole interval. Though we consider only 5 test points, it is very likely that CISVM is the best over the whole interval since it is better than or comparable to the methods learning with the true costs. Even at P1, where CISVM makes the most wrong guess of the true risk, CISVM still beats other methods sometimes. (3) Generally, CISVM is better than CSSVM with cost parameter being any value in a cost interval.

6.3 The Effect of Using Cost Distributions

To see whether more cost information will help improve performance, we consider two typical types of distributions, uniform and normal distributions, in our experiments: $U(1, 3)$, $U(1, 6)$, $U(1, 11)$, $U(5, 7)$, $U(5, 10)$, $U(5, 15)$, $N(5, 0.5)$, $N(5, 1)$, $N(5, 2)$, $N(10, 1)$, $N(10, 2)$, $N(10, 4)$, $N(20, 1)$, $N(20, 2)$, $N(20, 4)$. There are 15 cost distributions altogether. For CODIS, we use example-dependent cost-sensitive SVM [3] as learning algorithm A and set $T = 9$. We compare CODIS with CSSVM taking $E[c]$ as the true cost on the first 15 data sets used above. The settings for the 30 trails experiment and parameter selection are as same as that in the previous section.

Table 6 shows the results for each cost distribution separately. The results of SVM and CSSVM for normal distributions with the same mean value are all the same since they just use the expected cost. It can be found that CODIS can reduce 60% more risks than CSSVM. Paired t -tests and a following sign-tests both at significance level 95% are also carried out. It shows that CODIS is significantly better than other methods for any cost distribution on any data set in all 225 cases (15×15). Moreover, CODIS is very robust to different cost distributions since its results are very similar. For a specific cost distribution $N(10, 2)$, its 98.8% confidence interval is [5,15] which can be handled by CISVM. The expected cost is 10 which is used as the true cost in CSSVM. In this case, CSSVM is equivalent to CSMEAN. The results

can be found in Table 3 at P3 since the true cost is the mean value. It shows CODIS is much better than CISVM, and CISVM is much better than CSSVM. All the results indicate that when more information can be obtained in addition to cost intervals, a great advantage can be obtained by using the additional information appropriately.

7. CONCLUSION

In many real-world applications, it is generally difficult to have precise costs, while relatively easier to get cost intervals. This paper reports the first study on learning with cost intervals. We propose the CISVM method and experiments show that it is significantly superior to standard cost-sensitive SVMs. We also consider the situation where distributions are known about costs in addition to cost intervals, and propose the CODIS approach which is able to exploit cost distributions to reduce 60% more risks than standard cost-sensitive SVM.

8. ACKNOWLEDGEMENT

This research was supported by the NSFC (60635030), the JiangsuSF (BK2008018), and the 973 Program (2010CB327903).

9. REFERENCES

- [1] A. Asuncion and D. Newman. UCI machine learning repository, 2007.
- [2] A. P. Bradley. The use of the area under the ROC curve in the evaluation of machine learning algorithms. *Patt. Recogn.*, 30(7):1145–1159, 1997.
- [3] U. Brefeld, P. Geibel, and F. Wyszotzki. Support vector machines with example dependent costs. In *Proc. ECML*, pages 23–34, 2003.
- [4] N. V. Chawla, D. A. Cieslak, L. O. Hall, and A. Joshi. Automatically countering imbalance and its empirical relationship to cost. *DMKD*, 17(2):225–252, 2008.
- [5] M. Ciraco, M. Rogalewski, and G. Weiss. Improving classifier utility by altering the misclassification cost ratio. In *Proc. Intl. Workshop on Utility-Based Data Mining*, pages 46–52, 2005.
- [6] P. Domingos. MetaCost: A general method for making classifiers cost-sensitive. In *Proc. SIGKDD*, pages 155–164, 1999.
- [7] C. Drummond and R. C. Holte. Cost curves: An improved method for visualizing classifier performance. *Mach. Learn.*, 65(1):95–130, 2006.
- [8] C. Elkan. The foundations of cost-sensitive learning. In *Proc. IJCAI*, pages 973–978, 2001.
- [9] M. J. Gardner and D. G. Altman. Confidence intervals rather than p values: Estimation rather than hypothesis testing. *Stat. in Medic.*, 292:746–750, 1986.
- [10] X.-Y. Liu and Z.-H. Zhou. The influence of class imbalance on cost-sensitive learning: An empirical study. In *Proc. IEEE ICDM*, pages 970–974, 2006.
- [11] A. C. Lozano and N. Abe. Multi-class cost-sensitive boosting with p-norm loss functions. In *Proc. SIGKDD*, pages 506–514, 2008.
- [12] H. Masnadi-Shirazi and N. Vasconcelos. Asymmetric boosting. In *Proc. ICML*, pages 609–616, 2007.
- [13] V. S. Sheng and C. X. Ling. Thresholding for making classifiers cost-sensitive. In *Proc. AAAI*, 2006.
- [14] S. J. Stolfo, W. Fan, W. Lee, A. Prodromidis, and P. K. Chan. Cost-based modeling for fraud and intrusion detection: Results from the JAM project. *DARPA Information Survivability Conference and Exposition*, 2:1130–1144, 2000.
- [15] K. M. Ting. An instance-weighting method to induce cost-sensitive trees. *IEEE TKDE*, 14(3):659–665, 2002.

