

RankRC: Large-scale Nonlinear Rare Class Ranking

Aditya Tayal^{a,1,*}, Thomas F. Coleman^{b,1,2}, Yuying Li^{a,1}

^a*Cheriton School of Computer Science, University of Waterloo, Waterloo, ON, Canada N2L 3G1*

^b*Combinatorics and Optimization, University of Waterloo, Waterloo, ON, Canada N2L 3G1*

Abstract

Rare class problems are common in real-world applications across a wide range of domains. Standard classification algorithms are known to perform poorly in these cases, since they focus on overall classification accuracy. In addition, we have seen an explosion of data in recent years, resulting in many large scale rare class problems. In this paper, we consider nonlinear kernel based classification methods expressed as a regularized loss minimization problem. We address challenges associated with both rare class problems and large scale learning, by 1) optimizing area under curve of the receiver of operator characteristic in the training process, instead of classification accuracy and 2) using a rare class kernel representation to achieve an efficient time and space algorithm. We call our algorithm RankRC. We provide heuristic and theoretical justification for the rare class representation, and experimentally illustrate the effectiveness of RankRC in both test performance and computational complexity on several datasets.

Keywords: rare class, large scale, nonlinear kernel, receiver operator characteristic, ranking svm

1. Introduction

In many classification problems samples from one class are extremely rare (the minority class), while the number of samples belonging to the other class are plenty (the majority class). This situation is known as the rare class problem. It is also referred to as an unbalanced or skewed class distribution problem. Rare class problems naturally arise in several application domains, for example, fraud detection, customer churn, intrusion detection, fault detection, credit default, insurance risk and medical diagnosis.

Standard classification methods perform poorly when dealing with unbalanced data, e.g. support vector machines (SVM) [1, 2, 3], decision trees [4, 5, 1, 6], neural networks [1], Bayesian networks [7], and nearest neighbor methods [4, 8]. Most classification algorithms are driven by accuracy (i.e. minimizing error). Since minority examples constitute a small proportion of the data,

*Corresponding author

Email addresses: amtayal@uwaterloo.ca (Aditya Tayal), tfcoleman@uwaterloo.ca (Thomas F. Coleman), yuying@uwaterloo.ca (Yuying Li)

¹All three authors acknowledge funding from the National Sciences and Engineering Research Council of Canada

²This author acknowledges funding from the Ophelia Lazaridis University Research Chair. The views expressed herein are solely from the authors.

29 they have little impact on accuracy or total error. Thus majority examples overshadow the minor-
30 ity class, resulting in models which are heavily biased in recognizing the majority class. Also,
31 errors from different classes are assumed to have the same costs, which is usually not true. In most
32 problems, incorrect classification of the rare class is more expensive, for instance, diagnosing a
33 malignant tumor as benign has more severe consequences than the contrary case.

34 Solutions to the class imbalance problem have been proposed at both the data and algorithm
35 level. At the data level, various resampling techniques are used to balance class distribution,
36 including random under-sampling of majority class instances [9], over-sampling minority class
37 instances with new synthetic data generation [10], and focused resampling, in which samples
38 are chosen based on additional criteria [8]. Although sampling approaches have been showed
39 to achieve success in some applications, they are known to have drawbacks, for instance under-
40 sampling can eliminate useful information, while over-sampling can result in overfitting. At the
41 algorithm level, solutions are proposed by adjusting the algorithm itself. This usually involves ad-
42 justing the costs of the classes to counter the class imbalance (cost-sensitive learning) or adjusting
43 the decision threshold. However, true error costs are often unknown and using an inaccurate cost
44 model can lead to additional bias.

45 In this paper we focus on nonlinear kernel based classification methods expressed as a regu-
46 larized loss minimization problem. In recent years, we have seen an explosion of data, resulting
47 in many large scale rare class problems. For example, detecting unauthorized use of a credit card
48 from millions of transactions. Processing large datasets can be prohibitive for many nonlinear ker-
49 nel algorithms, which scale quadratically to cubically in the number of examples and may require
50 quadratic space as well.

51 To address the challenges associated with rare class problems and large scale learning we
52 propose the following:

- 53 1. Instead of maximizing accuracy (minimizing error), we optimize area under curve (AUC)
54 of the receiver operator characteristic. The AUC overcomes inadequacies of accuracy for
55 unbalanced problems and provides a skew independent measure. It is often used as the eval-
56 uation metric for unbalanced problems and therefore it is appropriate to directly optimize
57 it in the training process. This results in a regularized biclass ranking problem, which is a
58 special case of RankSVM with two ordinal levels [11].
- 59 2. To solve a kernel RankSVM problem in the dual, as originally proposed in [11], requires
60 $O(m^6)$ time and $O(m^4)$ space, where m is the number of data samples. Recently, Chapelle and
61 Keerthi [12] proposed a primal approach to solve RankSVM, which results in $O(m^3)$ time
62 and $O(m^2)$ space, for nonlinear kernels. We propose a modification to kernel RankSVM, that
63 takes specific advantage of the unbalanced nature of the problem, to achieve $O(mm_+)$ time
64 and $O(mm_+)$ space, where m_+ is the number of rare class examples. The idea is to restrict
65 the solution to a linear combination of rare class kernel functions. We call it RankRC, since
66 it enforces a rare class representation. We present heuristic and theoretical justification
67 for this choice. Specifically, we show RankRC is optimal with respect to RankSVM for
68 skewed data when a subset of kernel functions are used. We also draw connections to the
69 Nyström approximation method. Several of our results are general and can be applied to
70 other regularized loss minimization problems.

71 The rest of the paper is organized as follows. Sections 2 and 3 review the AUC measure and
 72 RankSVM. Section 4 develops the RankRC problem. Section 5 outlines the optimization method
 73 used to solve RankRC. Section 6 empirically compares RankRC with other kernel methods on
 74 several datasets. Finally, Section 7 concludes with summary remarks and potential extensions.

75 2. ROC Curve

76 Evaluation metrics play an important role in learning algorithms. They provide ways to
 77 assess performance as well as guide modeling. For classification problems, error rate is the
 78 most commonly used metric. For simplicity, we will consider the two-class case. Let $\mathcal{D} =$
 79 $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$ be a set of m training examples, where $\mathbf{x}_i \in \mathcal{X} \subseteq \mathbb{R}^d$, $y_i \in \{+1, -1\}$.
 80 Denote $f(\mathbf{x})$ as the inductive hypothesis obtained by training on example set \mathcal{D} . Then error rate is
 81 defined as,

$$ErrorRate = \frac{1}{m} \sum_{i=1}^m \mathbb{I}[f(\mathbf{x}_i) \neq y_i], \quad (1)$$

82 where $\mathbb{I}[p]$ denotes the indicator function and is equal to 1 if p is true, 0 if p is false. However, for
 83 highly unbalanced datasets, error rate is not appropriate since it can be biased toward the majority
 84 class [13, 14, 15, 16]. In this paper, we follow convention and set the minority class as positive and
 85 the majority class as negative. Consider a dataset that has 1 percent positive cases and 99 percent
 86 negative ones. A naive solution which assigns every example to be positive will obtain only 1
 87 percent error rate. Indeed, classifiers that always predict the majority class can obtain lower error
 88 rates than those that predict both classes equally well. But clearly these are not useful hypotheses.

89 Classification performance can be represented by a confusion matrix as in Table 1, with m_+
 90 denoting the number of majority examples and m_- the number of minority ones. The proportion
 91 of the two rows reflects class distribution and any performance measure that uses values from both
 92 rows will be sensitive to class skew.

		Predicted		Total
		$f(\mathbf{x}) = +1$	$f(\mathbf{x}) = -1$	
Actual	$y = +1$	True Positives (TP)	False Positives (FP)	m_+
	$y = -1$	False Negatives (FN)	True Negatives (TN)	m_-

Table 1: Add caption

93 The Receiver Operating Characteristic (ROC) can be used to obtain a skew independent mea-
 94 sure [13, 17, 18]. Most classifiers intrinsically output a numerical score and a predicted label is
 95 obtained by thresholding the score. For example, a threshold of zero leads to taking the sign of the
 96 numerical output as the label. Each threshold value generates a confusion matrix with different
 97 quantities of false positives and negatives. The ROC graph is obtained by plotting the true posi-
 98 tive rate (number of true positives divided by m_+) against the false positive rate (number of false

99 positives divided by m_-) as the threshold level is varied (see Figure 1). It depicts the trade-off
 100 between benefits (true positive) and costs (false positives) for different choices of the threshold.
 101 Thus it does not depend on a priori knowledge of the costs associated with misclassification. A
 102 ROC curve that dominates another provides a better solution at any cost point.

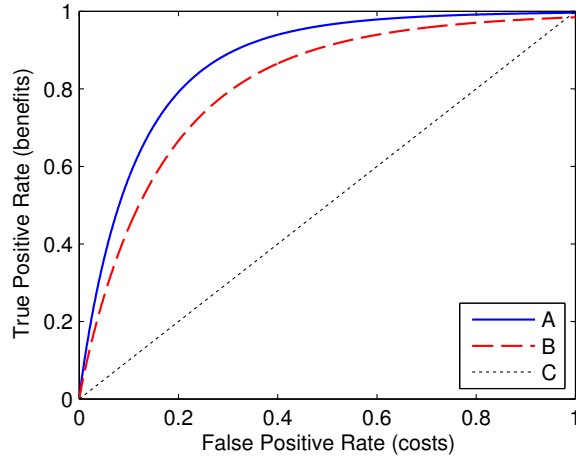


Figure 1: Example ROC curves. Curve A dominates B and curve B dominates C. Curve C has an AUC of 0.5 and indicates a model with no discriminative value.

103 To facilitate comparison, it is convenient to characterize ROC curves using a single measure.
 104 The area under a ROC curve (AUC) can be used for this purpose. It is the average performance of
 105 the model across all threshold levels and corresponds to the Wilcoxon rank statistic [19]. The AUC
 106 can be obtained by forming the ROC curve and using the trapezoid rule to compute area. Also,
 107 given the intrinsic output of a hypothesis, $f(\mathbf{x})$, we can directly compute the AUC by counting
 108 pairwise correct rankings [20]:

$$AUC = \frac{1}{m_+m_-} \sum_{\{i:y_i=+1\}} \sum_{\{j:y_j=-1\}} \mathbb{I}(f(\mathbf{x}_i) > f(\mathbf{x}_j)) . \quad (2)$$

109 Incorporating the AUC in the modeling process leads to a biclass ranking problem, as discussed
 110 in the following section.

111 3. RankSVM

112 The modeling process can usually be expressed as an optimization problem involving a loss
 113 function and a penalty on complexity (e.g. regularization term). For most classification problems,
 114 since the performance measure is error rate, it is natural to consider minimizing the empirical
 115 error rate (1) as the loss function. In practice, $\mathbb{I}[\cdot]$ is often replaced with a convex approximation
 116 such as the hinge loss, logistic loss or exponential loss [21]. Specifically, using the hinge loss,
 117 $\ell_h(z) = \max(0, 1-z)$, with ℓ_2 -regularization leads to the well known support vector machine (SVM)
 118 formulation [22, 23],

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{m} \sum_{i=1}^m \ell_h(y_i \mathbf{w}^T \mathbf{x}_i) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2, \quad (3)$$

119 where $\lambda \in \mathbb{R}_+$ is a parameter that controls complexity and the hypothesis, $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$, is assumed
 120 linear in the input space \mathcal{X} . Since SVMs try to minimize error rate, they can lead to ineffective
 121 class boundaries when dealing with highly skewed datasets, with resulting solutions biased toward
 122 the majority concept [3]. The literature contains several approaches to remedy this problem. Most
 123 prevalent are sampling methods and cost-sensitive learning. However, these approaches explicitly
 124 or implicitly fix the relative costs of misclassification. When the true costs are unknown, this can
 125 lead to suboptimal solutions.

126 Instead of minimizing error rate, we consider optimizing AUC as a natural way to deal with
 127 imbalance. Indeed, if we measure performance using AUC, it is preferable to optimize this quan-
 128 tity directly during the training process. In the AUC formula given in (2), we replace $\mathbb{I}[\cdot]$ with the
 129 hinge loss to obtain a convex ranking loss function. Thus we solve the following regularized loss
 130 minimization problem:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{m_+ m_-} \sum_{\{i: y_i = +1\}} \sum_{\{j: y_j = -1\}} \ell_h(\mathbf{w}^T \mathbf{x}_i - \mathbf{w}^T \mathbf{x}_j) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2. \quad (4)$$

131 Problem (4) is a special case of RankSVM proposed by Herbrich et al. [11] with two ordinal
 132 levels. Like SVM, RankSVM leads to a dual problem which can be expressed in terms of dot-
 133 products between input vectors. This allows us to obtain a non-linear function through the kernel
 134 trick [22], which consists of using a kernel function, $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, that corresponds to a feature
 135 map, $\phi: \mathcal{X} \rightarrow \mathcal{F} \subseteq \mathbb{R}^{d'}$, such that $\forall \mathbf{u}, \mathbf{v} \in \mathcal{X}$, $k(\mathbf{u}, \mathbf{v}) = \phi(\mathbf{u})^T \phi(\mathbf{v})$. Here, k directly computes the
 136 inner product of two vectors in a potentially high-dimensional feature space \mathcal{F} , without the need
 137 to explicitly form the mapping. Consequently, we can replace all occurrences of the dot-product
 138 with k in the dual and work implicitly in space \mathcal{F} .

139 However, since there is a Lagrange multiplier for each constraint associated with the hinge loss,
 140 the dual formulation leads to a problem in $m_+ m_- = O(m^2)$ variables. Assuming the optimization
 141 procedure has cubic complexity in the number of variables, the complexity of the dual method
 142 becomes $O(m^6)$, which is unreasonable for even medium sized datasets.

143 As noted by Chapelle [24], Chapelle and Keerthi [12], we can also solve the primal problem
 144 in the implicit feature space due to the Representer Theorem [25, 26]. This theorem states that the
 145 solution of any regularized loss minimization problem in \mathcal{F} can be expressed as a linear combina-
 146 tion of kernel functions evaluated at the training samples, $k(\mathbf{x}_i, \cdot)$, $i = 1, \dots, m$. Thus, the solution of
 147 (4) in \mathcal{F} can be written as:

$$\mathbf{w} = \sum_{i=1}^m \beta_i k(\mathbf{x}_i, \cdot), \text{ and } f(\mathbf{x}) = \sum_{i=1}^m \beta_i k(\mathbf{x}_i, \mathbf{x}). \quad (5)$$

148 Substituting (5) in (4) we can express the primal problem in terms of β :

$$\min_{\beta \in \mathbb{R}^m} \frac{1}{m_+ m_-} \sum_{\{i:y_i=+1\}} \sum_{\{j:y_j=-1\}} \ell_h \left(\sum_{r=1}^m \beta_r k(\mathbf{x}_r, \mathbf{x}_i) - \sum_{r=1}^m \beta_r k(\mathbf{x}_r, \mathbf{x}_j) \right) + \frac{\lambda}{2} \sum_{i,j=1}^m \beta_i \beta_j k(\mathbf{x}_i, \mathbf{x}_j),$$

149 or more simply,

$$\min_{\beta \in \mathbb{R}^m} \frac{1}{m_+ m_-} \sum_{\{i:y_i=+1\}} \sum_{\{j:y_j=-1\}} \ell_h \left(\beta^T K_i - \beta^T K_j \right) + \frac{\lambda}{2} \beta^T K \beta, \quad (6)$$

150 where $K \in \mathbb{R}^{m \times m}$ is the kernel matrix, $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$, and K_i denotes the i th row of K . To be able
 151 to solve (6) using unconstrained optimization methods such as gradient descent, we require the
 objective to be differentiable. We replace the hinge loss, ℓ_h , with an ϵ -smoothed differentiable

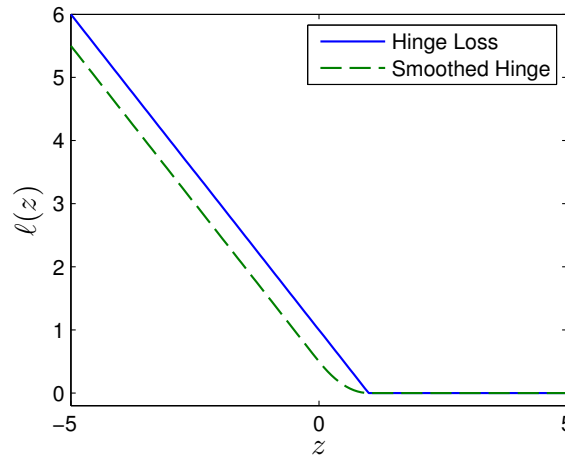


Figure 2: The smoothed hinge is a differentiable approximation of the hinge loss. Here the smoothed hinge is shown with $\epsilon = 0.5$.

152

153 approximation, ℓ_ϵ , defined as,

$$\ell_\epsilon(z) = \begin{cases} (1-\epsilon) - z & \text{if } z < 1 - 2\epsilon \\ \frac{1}{4\epsilon}(1-z)^2 & \text{if } 1 - 2\epsilon \leq z < 1 \\ 0 & \text{if } z \geq 1, \end{cases}$$

154 which transitions from linear cost to zero cost using a quadratic segment (see Figure 2) and pro-
 155 vides similar benefits as the hinge loss. Now we can solve (6) using standard unconstrained opti-
 156 mization techniques. Since there are m variables, Newton's method would for example take $O(m^3)$
 157 operations to converge.

158 RankSVM is popular in the information retrieval community, where linear models are the norm
 159 [e.g. see 27]. For a linear model, with d -dimension input vectors, the complexity of RankSVM
 160 can be reduced to $O(md + m \log m)$ [12]. However, many rare class problems require a nonlinear
 161 function to achieve optimal results. But solving a nonlinear RankSVM in $O(m^3)$ time may not be
 162 practical for mid- to large-sized datasets. Moreover, the method requires $O(m^2)$ space to store the
 163 kernel matrix. We believe this complexity is, in part, the reason why nonlinear RankSVMs are not
 164 commonly used to solve rare class problems.

165 In the next section we propose a modification to nonlinear RankSVMs that takes specific ad-
 166 vantage of the unbalanced nature of the problem to achieve $O(mm_+)$ time and $O(mm_+)$ space, while
 167 not sacrificing performance.

168 4. RankRC: Ranking with Rare Class Representation

169 For highly unbalanced datasets, to make SVM computational feasible for large scale problems,
 170 we propose a rare class (RC) based method. Specifically we propose a RC based method, which
 171 restricts the solution to the form

$$f(\mathbf{x}) = \sum_{\{i: y_i = +1\}} \beta_i k(\mathbf{x}_i, \mathbf{x}), \quad (7)$$

172 so it consists only of kernel function realizations of the minority class.

173 Next we present motivate form (7) by assuming specific properties of the class conditional
 174 distributions and kernel function. Zhu et al. [28] make use of similar assumptions, however, in
 175 their method they attempt to directly estimate the likelihood ratio. In contrast, we use a regularized
 176 loss minimization approach.

177 Recall that the optimal ranking function for a classification problem is the posterior probability,
 178 $P(y = 1|\mathbf{x})$, since it minimizes the Bayes risk for arbitrary costs. From Bayes' Theorem, we have

$$P(y = 1|\mathbf{x}) = \frac{P(y = 1)P(\mathbf{x}|y = 1)}{P(y = 1)P(\mathbf{x}|y = 1) + P(y = -1)P(\mathbf{x}|y = -1)}. \quad (8)$$

179 Any monotonic transformation of (8) also yields equivalent ranking capability. Dividing the nu-
 180 merator and denominator of (8) by $P(y = -1)P(\mathbf{x}|y = -1)$, we note that $P(y = 1|\mathbf{x})$ is a monotonic
 181 transformation of the likelihood ratio, denoted as

$$f(\mathbf{x}) = \frac{P(\mathbf{x}|y = 1)}{P(\mathbf{x}|y = -1)}, \quad (9)$$

182 which is the ranking function we wish to obtain. Now, if we assume that the conditional density,
 183 $P(\mathbf{x}|y = 1)$, is a mixture of m_+ identical spherical normals centered at the rare class examples, we
 184 can write

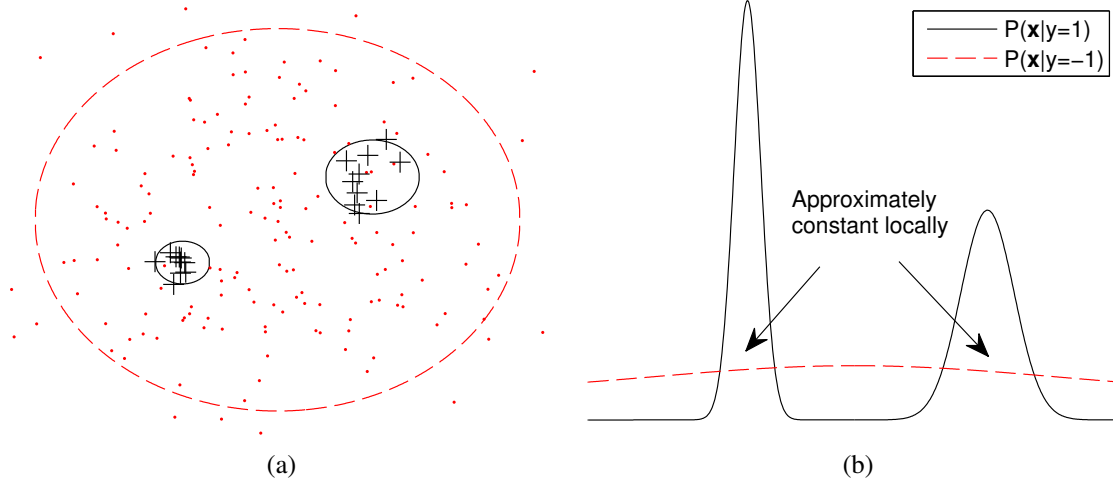


Figure 3: (a) An example of a rare class dataset. Red ‘.’s indicate negative (majority) examples and black ‘+’s indicate positive (minority) examples. (b) The class conditional distributions showing that $P(\mathbf{x}|y = -1)$ is relatively constant in local neighborhood of positive examples.

$$P(\mathbf{x}|y = 1) = \sum_{\{i:y_i=+1\}} a_i \exp \left\{ -\frac{\|\mathbf{x}_i - \mathbf{x}\|^2}{\sigma^2} \right\},$$

185 for some constants a_i . This mixture encompasses a large range of possible distributions from the
 186 m_+ rare examples provided. If we also assume that k denotes a Gaussian kernel function with
 187 width σ , then we have

$$P(\mathbf{x}|y = 1) = \sum_{\{i:y_i=+1\}} a_i k(\mathbf{x}_i, \mathbf{x}).$$

188 Observe that in rare class problems most examples are from the majority class ($y = -1$) and only
 189 a small number are from the rare class ($y = +1$). Therefore it is reasonable to assume $P(\mathbf{x}|y = -1)$
 190 is locally constant in a neighborhood around the minority class examples, see Figure 3 for an
 191 illustration. Let $P(\mathbf{x}|y = -1) \approx c_i$ for each minority example i in the neighbourhood of \mathbf{x}_i .³ Then
 192 (9) can be written as,

$$f(\mathbf{x}) \approx \sum_{\{i:y_i=+1\}} \frac{a_i k(\mathbf{x}_i, \mathbf{x})}{c_i},$$

³We do not make this more precise since we are mainly interested in motivating an approximate form.

193 which corresponds to the rare class representation (7) we have chosen. If the assumptions made
 194 are relaxed, we may still expect the rare class representation to perform reasonably well.

195 For any general regularized loss minimization problem with any loss function $L : \mathbb{R}^m \rightarrow \mathbb{R}$, we
 196 can consider a corresponding rare class regularization problem. Assume that a penalty parameter
 197 $\lambda \in \mathbb{R}_+$ is given, a regularized loss minimization problem can be described as

$$\min_{\mathbf{w} \in \mathbb{R}^{d'}} L(f(\mathbf{x}_1), \dots, f(\mathbf{x}_m)) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2, \quad (10)$$

198 where $f : \mathcal{X} \rightarrow \mathbb{R}$ is a linear hypothesis in feature space, $f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$. Here $L(\cdot)$ is any loss
 199 function including both standard SVM and ranking SVM functions, since SVM-Rank is equivalent
 200 to a 1-class SVM on an enlarged dataset with the set of points $\mathcal{P} = \{\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j) : y_i > y_j, i, j =$
 201 $1, \dots, m\}$. From the Representer Theorem, a solution vector $\mathbf{w} \in \mathcal{S} = \text{span}\{\phi(\mathbf{x}_i) : i = 1, \dots, m\}$ can
 202 be expressed in terms of all the given training points in the feature space.

203 Using the restricted hypothesis (7), we consider the following constrained regularized ranking
 204 problem,

$$\begin{aligned} & \min_{\beta \in \mathbb{R}^{\mathcal{R}}} L(f(\mathbf{x}_1), \dots, f(\mathbf{x}_m)) + \frac{\lambda}{2} \beta^T K_{\mathcal{R}\mathcal{R}} \beta, \\ & \text{subject to } f(\mathbf{x}) = \sum_{i \in \mathcal{R}} \beta_i k(\mathbf{x}_i, \mathbf{x}) \end{aligned} \quad (11)$$

205 where $\mathcal{R} \subseteq \{1, \dots, m\}$. The proposed ranking with a rare class representation, subsequently referred
 206 to as RankRC, is a special case of (11) with $\mathcal{R} = \{i : y_i = 1\}$.

207 In order to see potential advantages of RankRC, we compare the full data set regularized SVM
 208 ranking problem

$$\begin{aligned} & \min_{\beta \in \mathbb{R}^m} \frac{1}{m_+ m_-} \sum_{\{i: y_i=+1\}} \sum_{\{j: y_j=-1\}} \ell_h(f(\mathbf{x}_i) - f(\mathbf{x}_j)) + \frac{\lambda}{2} \beta^T K \beta \\ & \text{subject to } f(\mathbf{x}) = \sum_{i=1}^m \beta_i k(\mathbf{x}_i, \mathbf{x}) \end{aligned} \quad (12)$$

209 with the subset data regularization problem (11).

210 Applying Theorem 2 in [?], we can establish the following theoretical result.

211 **Theorem 1.** *Let $f^*(\mathbf{x})$ be the optimal hypothesis of the full data set SVM-Rank problem (12) under*
 212 *the feature mapping $\phi : \mathcal{X} \rightarrow \mathcal{F}$ and $\tilde{f}^*(\mathbf{x})$ be the optimal hypothesis for the subset data set SVM-*
 213 *Rank problem (11) where $\mathcal{R} \subseteq \{1, \dots, m\}$. Assume there exists $\kappa > 0$ such that $k(\mathbf{x}, \mathbf{x}) \leq \kappa$, where*
 214 *$k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is the kernel map associated with ϕ . Then the following inequality holds for all*
 215 *$\mathbf{x} \in \mathcal{X}$:*

$$|f^*(\mathbf{x}) - \bar{f}^*(\mathbf{x})| \leq \frac{2\kappa}{\lambda} \left(\sum_{\{i:y_i=+1\}} \frac{\mathbb{I}[i \notin \mathcal{R}]}{m_+} + \sum_{\{j:y_j=-1\}} \frac{\mathbb{I}[j \notin \mathcal{R}]}{m_-} \right)^{\frac{1}{2}}, \quad (13)$$

216 where $\mathbb{I}[p]$ denotes the indicator function and is equal to 1 if p is true, 0 if p is false.

217 We note, the difference in hypothesis decreases for larger regularization according to $O(\frac{1}{\lambda})$
 218 and as we include more kernel function realizations in our representation. However, for the rank-
 219 ing loss, the bound decreases asymmetrically depending on whether we include a point from the
 220 positive or negative class. In particular, if the dataset is unbalanced with $m_- \gg m_+$, then $\frac{1}{m_+} \gg \frac{1}{m_-}$,
 221 and the reduction obtained from including a positive class basis is much greater than including
 222 one from the negative class. Hence, for a fixed number of kernel function realizations, the bound
 223 is minimized by first including bases corresponding to the positive or rare class.

224 5. Computational Complexity Comparison between RankRC over SVM-Rank

225 Using a smooth loss function in (??), we have the following RankRC problem in m_+ variables,

$$\min_{\beta \in \mathbb{R}^{m_+}} \frac{1}{m_+ m_-} \sum_{\{i:y_i=+1\}} \sum_{\{j:y_j=-1\}} \ell_\epsilon \left(\beta^T K_{i+} - \beta^T K_{j+} \right) + \frac{\lambda}{2} \beta^T K_{++} \beta. \quad (14)$$

226 Here, K_{i+} denotes i th row of K with column entries corresponding to only the positive class, and
 227 $K_{++} \in \mathbb{R}^{m_+ \times m_+}$ is the square submatrix of K corresponding to the positive class entries. We also
 228 replace ℓ_h with the smooth approximation ℓ_ϵ . To solve (14) we can use several approaches, which
 229 are discussed below.

230 5.1. Linearization

231 Since K_{++} is a positive semi-definite matrix, it has an eigen-decomposition which can be ex-
 232 pressed in the form, $K_{++} = U \Lambda U^T$, with U being an orthonormal matrix (i.e. $U^T U = I$) and Λ a
 233 diagonal matrix containing non-negative eigenvalues of K_{++} . Let $\mathbf{w} = \Lambda^{\frac{1}{2}} U^T \beta$, then

$$\beta = U \Lambda^{\dagger \frac{1}{2}} \mathbf{w}, \quad (15)$$

234 where Λ^\dagger denotes the pseudoinverse of Λ . We can substitute (15) in (14) to obtain the following
 235 problem,

$$\min_{\mathbf{w} \in \mathbb{R}^{m_+}} \frac{1}{m_+ m_-} \sum_{\{i:y_i=+1\}} \sum_{\{j:y_j=-1\}} \ell_\epsilon \left(\mathbf{w}^T \Lambda^{\dagger \frac{1}{2}} U^T K_{i+} - \mathbf{w}^T \Lambda^{\dagger \frac{1}{2}} U^T K_{j+} \right) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2, \quad (16)$$

236 which is a problem in linear space. That is, Problem (16) is equivalent to Problem (4) with data
 237 points given by $\mathbf{x}_i = \Lambda^{\dagger \frac{1}{2}} U^T K_{i+} \in \mathbb{R}^{m_+}$, $i = 1, \dots, m$. Therefore we can use the algorithm described
 238 in Chapelle and Keerthi [12] to solve the linear ranking problem in $O(mm_+ + m \log m) = O(mm_+)$
 239 time. Including the cost of factoring K_{++} , the total time is $O(mm_+ + m_+^3)$. Once we solve for optimal
 240 \mathbf{w} we can use (15) to obtain β for subsequent testing purposes. Also, since we only need entries
 241 $\{K_{ij} : i = 1, \dots, m, y_j = 1\}$, the method only requires $O(mm_+)$ space.

242 5.2. Unconstrained Optimization

243 We can also directly solve (14) using standard unconstrained optimization methods. Gradient
 244 only methods, such as steepest descent and nonlinear conjugate gradient do not require estima-
 245 tion of the Hessian. Although this makes each iteration much cheaper, convergence can be slow,
 246 especially near the solution. In contrast Hessian based algorithms, such as Newton's method can
 247 obtain quadratic convergence near the solution, but each iteration can be expensive. In Newton's
 248 method, the p th iterate is updated according to

$$\beta^{(p+1)} = \beta^{(p)} + \mathbf{s},$$

249 where the step, \mathbf{s} , is obtained by minimizing the quadratic Taylor approximation around the current
 250 iterate $\beta^{(p)}$:

$$\min_{\mathbf{s}} \quad \mathbf{s}^T \mathbf{g}^{(p)} + \frac{1}{2} \mathbf{s}^T H^{(p)} \mathbf{s}, \quad (17)$$

251 where $H^{(p)}$ and $\mathbf{g}^{(p)}$ are the Hessian and gradient of the objective at $\beta^{(p)}$, respectively. Problem
 252 (17) has a closed form solution given by

$$\mathbf{s} = - \left(H^{(p)} \right)^{-1} \mathbf{g}^{(p)}.$$

253 Since $H^{(p)}$ is a $m_+ \times m_+$ matrix, this involves $O(m_+^3)$ cost in each iteration. To avoid this, we can
 254 use the truncated Newton method in which $H^{(p)} \mathbf{s} = -\mathbf{g}^{(p)}$ is solved using linear conjugate gradient.
 255 Here, the Hessian is not computed explicitly and the method iteratively approximates the solution
 256 using Hessian-vector products. Since each iteration in the linear conjugate gradient algorithm
 257 leads to a descent direction, we can terminate early while still improving convergence.

258 A drawback of (truncated) Newton's method is that it can be sensitive to the initial point. If the
 259 initial point is not chosen close enough to the solution, the method can be slow to converge, or fail
 260 altogether. Therefore we consider a subspace-trust-region method, which combines the benefit of a
 261 truncated Newton step with steepest descent. In our tests, we found that the subspace-trust-region
 262 method converges with significantly fewer iterations than Newton's method.

263 The idea behind the trust-region method is to solve (17) while constraining the step, \mathbf{s} , to a
 264 neighborhood around the current iterate, in which the approximation is trusted:

$$\begin{aligned} \min_{\mathbf{s}} \quad & \frac{1}{2} \mathbf{s}^T H^{(p)} \mathbf{s} + \mathbf{s}^T \mathbf{g}^{(p)} \\ \text{s.t.} \quad & \|\mathbf{s}\|_2 \leq \Delta^{(p)}. \end{aligned} \quad (18)$$

265 The trust region radius, $\Delta^{(p)}$, is adjusted at each iterate according to standard rules, for example
 266 it is decreased if the solution obtained is worse than the current iterate. Problem (18) can be
 267 solved accurately [e.g see 29], however, the solution uses the full eigen-decomposition of $H^{(p)}$.
 268 To avoid this computation, in the subspace-trust-region method, Problem (18) is restricted to a
 269 two-dimensional subspace spanned by the gradient, $\mathbf{g}^{(p)}$, and an approximate Newton direction,
 270 \mathbf{s}_2 , which can be obtained by solving $H^{(p)}\mathbf{s}_2 = -\mathbf{g}^{(p)}$ using linear conjugate gradient [30]. The idea
 271 behind this choice is to ensure global convergence (via steepest descent direction) and achieve
 272 fast local convergence (via the Newton step). Once the subspace has been computed, solving (18)
 273 costs $O(1)$ time, since in the subspace the problem is only two-dimensional. The implementation
 274 we use is provided in Matlab's optimization toolbox, `fminunc/fmincon`.

275 5.2.1. Computing Gradient and Hessian-Vector Product

We describe how we can compute the gradient and Hessian-vector product of Problem (14) ef-
 276 ficiently. Let $K_{\bullet+} = [K_{ij}]_{i=1,\dots,m,y_j=-1} \in \mathbb{R}^{m \times m_+}$ denote the rectangular submatrix of K with columns
 indexed by the positive class. Consider the expanded matrix

$$A = [K_{i+} - K_{j+}]_{i:y_i=1,j:y_j=-1} \in \mathbb{R}^{m+m_- \times m_+},$$

consisting of the differences of rows in $K_{\bullet+}$ corresponding to all pairwise preferences. In our
 computation we do not explicitly form matrix A , rather we note that A can be expressed as a sparse
 matrix product:

$$A = DK_{\bullet+},$$

276 where $D \in \mathbb{R}^{m+m_- \times m}$ is a sparse matrix that encodes a pairwise preference. That is, if $y_i > y_j$, then
 277 there exists a row r in P such that $D_{ri} = 1, D_{rj} = -1$ and the rest of the row is zero. Let A_r denote
 278 the r th row of A . Then the ranking loss expression in (14) can be written as,

$$\begin{aligned} & \sum_{\{i:y_i=+1\}} \sum_{\{j:y_j=-1\}} \ell_\epsilon \left(\beta^T K_{i+} - \beta^T K_{j+} \right) \\ &= \sum_{r=1}^{m+m_-} \ell_\epsilon \left(\beta^T A_r \right) \\ &= \sum_{r=1}^{m+m_-} \mathbb{I}[r \in \mathcal{L}] \left(1 - \epsilon - \beta^T A_r \right) + \sum_{r=1}^{m+m_-} \mathbb{I}[r \in \mathcal{Q}] \frac{1}{4\epsilon} \left(1 - \beta^T A_r \right)^2, \end{aligned} \quad (19)$$

279 where $\mathcal{L} = \{r : \beta^T A_r < 1 - 2\epsilon\}$ is the set of pairwise differences which are in the linear portion of
 280 ℓ_ϵ , and $\mathcal{Q} = \{r : 1 - 2\epsilon \leq \beta^T A_r < 1\}$ is the set which fall in the quadratic part. Denote $\mathbf{e} \in \mathbb{R}^{m+m_-}$
 281 as a vector of ones. Define $\mathbf{e}^{\mathcal{L}} \in \mathbb{R}^{m+m_-}$ as a binary vector where $\mathbf{e}_r^{\mathcal{L}} = 1$ if $r \in \mathcal{L}$ and $\mathbf{e}_r^{\mathcal{L}} = 0$ if $r \notin \mathcal{L}$.
 282 Also define $I^{\mathcal{Q}} \in \mathbb{R}^{m+m_- \times m+m_-}$ as a diagonal matrix, where $I_{rr}^{\mathcal{Q}} = 1$, if $r \in \mathcal{Q}$, and $I_{rr}^{\mathcal{Q}} = 0$, if $r \notin \mathcal{Q}$.
 283 Then (19) is equivalent to

$$\begin{aligned}
& \left(\mathbf{e}^{\mathcal{L}} \right)^T \left((1-\epsilon)\mathbf{e} - A\boldsymbol{\beta} \right) + \frac{1}{4\epsilon} \left(\mathbf{e} - A\boldsymbol{\beta} \right)^T I^{\mathcal{Q}} \left(\mathbf{e} - A\boldsymbol{\beta} \right) \\
& = \left(\mathbf{e}^{\mathcal{L}} \right)^T \left((1-\epsilon)\mathbf{e} - PK_{\bullet+}\boldsymbol{\beta} \right) + \frac{1}{4\epsilon} \left(\mathbf{e} - PK_{\bullet+}\boldsymbol{\beta} \right)^T I^{\mathcal{Q}} \left(\mathbf{e} - PK_{\bullet+}\boldsymbol{\beta} \right) .
\end{aligned}$$

284 Therefore the objective function in (14) can be expressed as

$$F(\boldsymbol{\beta}) \triangleq \frac{1}{m_+m_-} \left[\left(\mathbf{e}^{\mathcal{L}} \right)^T \left((1-\epsilon)\mathbf{e} - PK_{\bullet+}\boldsymbol{\beta} \right) + \frac{1}{4\epsilon} \left(\mathbf{e} - PK_{\bullet+}\boldsymbol{\beta} \right)^T I^{\mathcal{Q}} \left(\mathbf{e} - PK_{\bullet+}\boldsymbol{\beta} \right) \right] + \frac{\lambda}{2} \boldsymbol{\beta}^T K_{++}\boldsymbol{\beta} . \quad (20)$$

285 We obtain the gradient by taking the derivative of (20) with respect to $\boldsymbol{\beta}$:

$$\begin{aligned}
\mathbf{g} & \triangleq \frac{\partial F}{\partial \boldsymbol{\beta}} = \frac{1}{m_+m_-} \left[- \left(\mathbf{e}^{\mathcal{L}} \right)^T PK_{\bullet+} + \frac{1}{2\epsilon} PK_{\bullet+} I^{\mathcal{Q}} (PK_{\bullet+}\boldsymbol{\beta} - \mathbf{e}) \right] + \lambda K_{++}\boldsymbol{\beta} \\
& = \frac{1}{m_+m_-} \left[- \left(\left(\mathbf{e}^{\mathcal{L}} \right)^T P \right) K_{\bullet+} + \frac{1}{2\epsilon} \left(P \left(K_{\bullet+} \left(I^{\mathcal{Q}} P \right) \left(K_{\bullet+}\boldsymbol{\beta} \right) \right) - P \left(K_{\bullet+} \left(I^{\mathcal{Q}} \mathbf{e} \right) \right) \right) \right] + \lambda K_{++}\boldsymbol{\beta} .
\end{aligned} \quad (21)$$

286 In the last expression we have used brackets to emphasize the order of operations that leads to an
287 efficient implementation and avoids computing $A = PK_{\bullet+}$. It can be verified that the time required
288 is $O(mm_+)$.

289 We obtain the Hessian by taking the derivative of (21) with respect to $\boldsymbol{\beta}$:

$$H \triangleq \frac{\partial^2 F}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} = \frac{1}{2\epsilon m_+m_-} \left(PK_{\bullet+} I^{\mathcal{Q}} PK_{\bullet+} \right) + \lambda K_{++} .$$

290 Note the Hessian requires computing A . However, for the linear conjugate gradient method we
291 only require computing $H\mathbf{s}$ for some vector \mathbf{s} . In this case, we can avoid computing A by using the
292 following order of operations:

$$H\mathbf{s} = \frac{1}{2\epsilon m_+m_-} \left(P \left(K_{\bullet+} \left(I^{\mathcal{Q}} P \right) \left(K_{\bullet+}\mathbf{s} \right) \right) \right) + \lambda K_{++}\mathbf{s} .$$

293 The time required to compute $H\mathbf{s}$ is also $O(mm_+)$.

294 In the subspace-trust-region method we use a maximum of 25 conjugate gradient iterations.
295 We found the solution usually converges in a constant number of trust region iterations. Since
296 each iteration requires $O(mm_+)$ time, the total time required by the algorithm is $O(mm_+)$. Total
297 space is also $O(mm_+)$.

298 Finally, we note that we can slightly improve the time required to compute the gradient and
299 Hessian-vector product by first sorting the values of $K_{\bullet+}\boldsymbol{\beta}$ or $K_{\bullet+}\mathbf{s}$. Though this does not improve
300 the *big-O* efficiency, it does reduce the constant factor. We refer the interested reader to [12] for
301 details on a method which can be adapted for the nonlinear RankRC objective (14).

302 6. Experiments

303 In this section we empirically compare RankRC to other methods on several unbalanced prob-
 304 lems. The following methods are compared:

- 305 1. KNN: k-Nearest-Neighbors algorithm. The posterior probability is used as the ranking func-
 306 tion:

$$P(y|\mathbf{x}) = \frac{1}{k} \sum_{i \in \mathcal{K}} \mathbb{I}[y_i = 1],$$

307 where \mathcal{K} is the set of k nearest neighbors in the training dataset.

- 308 2. SVM: This is the standard nonlinear SVM [23], in which the primal problem,

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{m} \sum_{i=1}^m \max(0, 1 - y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b)) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2,$$

309 is solved (in the dual) to obtain the decision function, $f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b = \sum_{i=1}^m \beta_i k(\mathbf{x}_i, \mathbf{x}) + b$,
 310 with $k(\mathbf{x}_i, \mathbf{x}) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x})$.

- 311 3. SVM-w: Weighted SVM [23, 31] in which

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{m} \sum_{i=1}^m \omega_i \max(0, 1 - y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b)) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2,$$

312 is solved, with different weights associated with each class:

$$\omega_i = \begin{cases} \frac{m}{2m_+} & \text{if } y_i = +1 \\ \frac{m}{2m_-} & \text{if } y_i = -1. \end{cases}$$

313 The idea is to penalize misclassification error of minority examples more heavily in order to
 314 reduce the bias towards majority examples.

- 315 4. SVM-RUS: Randomly Under Sample the majority class examples ($y = -1$) to match the
 316 number of minority examples [9]. The resulting dataset, with $2m_+$ points, is used to train a
 317 standard SVM.

- 318 5. SVM-SMT: Uses a Synthetic Minority Oversampling TEchnique (SMOTE) [10] in which
 319 the rare class is over-sampled by creating new synthetic rare class samples according to
 320 each rare class sample and its k nearest neighbors. Each new sample is generated in the
 321 direction of some or all of the nearest neighbors. We oversample to match the number of
 322 majority examples. The resulting dataset, with $2m_-$ points, is used to train a standard SVM.

- 323 6. RANK-SVM: Nonlinear RankSVM problem (6).

- 324 7. RANK-RND: We solve the RankSVM problem constrained to m_+ randomly selected set of
 325 basis function, i.e. Problem (??).
 326 8. RANK-RC: We solve RankSVM constrained to the rare class representation, i.e. Problem
 327 (14).

328 We use LIBSVM [32] to solve the SVM problems (2-5). LIBSVM is a popular and efficient
 329 implementation of the sequential minimal optimization algorithm [33]. We set cache size to 10GB
 330 to minimize cache misses; termination criteria and shrinking heuristics are used in their default
 331 settings. The ranking methods (6-8) are solved using the subspace-trust-region method as outlined
 332 in Section 5. Termination tolerance is set at $1e-6$. For ranking methods, the memory available to
 333 store the kernel matrix is limited to 10GB. Experiments are performed on a Xeon E5620@2.4Ghz
 334 running Linux.

335 All datasets are standardized to zero mean and unit variance before training. Since our fo-
 336 cus is on nonlinear kernels, for all SVM and ranking methods (2-8), we use a Gaussian kernel,
 337 $k(\mathbf{u}, \mathbf{v}) = \exp(-\|\mathbf{u} - \mathbf{v}\|_2^2 / \sigma^2)$ with $\sigma^2 = \frac{1}{m^2} \sum_{i,j=1}^m \|\mathbf{x}_i - \mathbf{x}_j\|_2^2$. The penalty parameter λ is determined
 338 by cross-validation over values $\log_2 \lambda = [-20, -18, \dots, 8, 10]$. For KNN we cross-validate over
 339 $k = [1, 2, \dots, \lceil \min(100, \sqrt{m}) \rceil]$.

340 6.1. Simulated Data

341 We simulate an unbalanced dataset in the following manner. Rare class instances are sampled
 342 from six multivariate normal distributions with equal probability. Their centers, $\mu_i, i = 1, \dots, 6$, are
 343 randomly chosen within a unit cube. The majority class is sampled from $\binom{6}{2} = 15$ multivariate
 344 normal distributions with equal probability. Their centers are chosen along lines connecting all
 345 combinations of two rare class centers, i.e. $t\mu_i + (1-t)\mu_j, i > j$. The parameter $t \in [0, 1]$ is used
 346 to roughly control the degree of class overlap. All covariances are chosen to be spherical, $\sigma^2 I$.
 347 Finally, the imbalance ratio, $\rho = \frac{m_+}{m_-}$, is used to determine the number of samples drawn from
 348 each of the class conditional distributions. An example configuration in 2-dimensional space is
 349 shown in Figure 4. The resulting dataset contains multiple rare-class subconcepts that vary in
 350 discriminative structure.

351 For our experiment we generate data in 5-dimensional space with $\sigma = 0.5$. We set $t = 0.9$,
 352 0.75 , and 0.6 to produce datasets with high, medium and low overlap, respectively. The imbalance
 353 ratio, ρ , is varied from 10% to 40% in 10% increments for each t value. Thus we have a total of
 354 12 datasets. For each dataset we generate 1000 training points, 1000 validation points and 10000
 355 testing points. Results are averaged over 10 trials.

356 Table 2 shows test AUC results using different methods. KNN does not perform as well as SVM
 357 and ranking methods. In general, ranking methods perform better than SVM based methods when
 358 there is greater overlap and higher imbalance (lower ρ). RANK-RND under performs in medium
 359 and low overlap datasets. In comparison, RANK-RC yields statistically similar performance as
 360 RANK-SVM across all datasets. Overall, both RANK-RC and RANK-SVM provide the best models.

361 Figure 5 compares the empirical ranking loss function,

$$\hat{R}_h = \frac{1}{m_+ m_-} \sum_{\{i: y_i = +1\}} \sum_{\{j: y_j = -1\}} \ell_h(f(\mathbf{x}_i) - f(\mathbf{x}_j)) ,$$

Overlap	ρ	KNN	Classification Loss				Ranking Loss				True Bayes
			SVM	SVM-W	SVM-RUS	SVM-SMT	RANK-SVM	RANK-RND	RANK-RC		
High	10%	56.3±0.4	57.3±0.3	59.8±0.3	59.1±0.4	58.9±0.4	61.5±0.3	61.4±0.3	61.4±0.2	66.8	
	20%	55.5±0.4	59.2±0.2	61.2±0.1	61.0±0.4	60.9±0.4	61.6±0.3	61.3±0.4	62.3±0.3		
	30%	57.3±1.0	61.1±0.2	62.2±0.4	62.2±0.1	61.8±0.3	62.3±0.4	62.2±0.3	62.6±0.4		
	40%	59.0±0.7	62.8±0.2	63.2±0.3	63.3±0.2	62.8±0.2	62.7±0.3	62.5±0.2	62.7±0.3		
Medium	10%	54.9±0.6	56.8±0.5	59.5±0.5	58.8±0.3	58.6±0.9	61.4±0.5	60.1±0.5	61.4±0.4	69.5	
	20%	54.5±0.3	60.1±0.4	62.1±0.2	62.1±0.3	62.0±0.5	62.8±0.4	61.5±0.4	63.4±0.2		
	30%	57.6±0.8	63.4±0.1	64.0±0.3	63.3±0.1	63.4±0.4	64.1±0.4	62.4±0.5	64.6±0.5		
	40%	58.0±0.7	65.3±0.2	65.5±0.1	65.4±0.1	65.3±0.2	64.6±0.3	63.0±0.2	64.9±0.3		
Low	10%	55.9±1.0	61.1±0.6	64.0±0.4	63.5±0.2	63.0±0.8	65.3±0.5	62.7±0.6	65.5±0.4	74.4	
	20%	57.2±0.4	64.2±0.3	66.6±0.2	66.2±0.2	66.4±0.2	67.1±0.3	63.5±0.4	67.3±0.2		
	30%	59.9±0.9	69.1±0.4	69.4±0.2	68.9±0.1	69.2±0.2	69.2±0.4	65.8±0.4	69.8±0.4		
	40%	61.6±1.5	71.1±0.1	70.7±0.3	70.5±0.1	71.1±0.1	69.8±0.3	65.8±0.3	69.9±0.3		

Table 2: Comparison of test AUC results for simulated datasets with high ($t = 0.9$), medium ($t = 0.75$) and low ($t = 0.6$) overlap, each with $\rho = 10\%$, 20% , 30% and 40% minority samples. Mean AUC score with standard error over 10 trials are shown. Bolded scores indicate the result is statistically not different than the best performing model using a two-tailed t-test with 99% confidence.

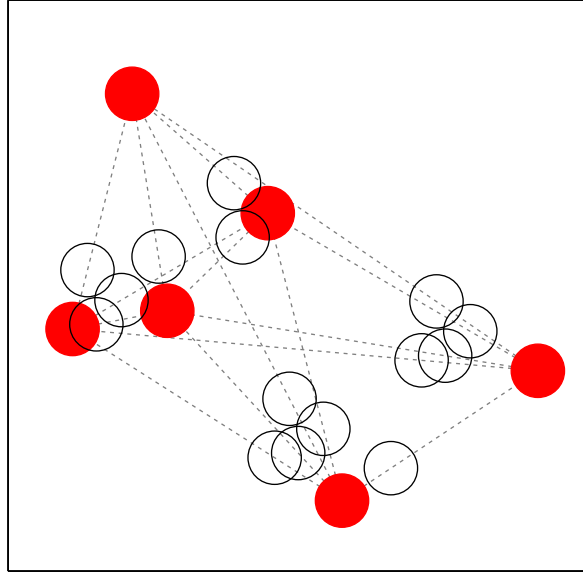


Figure 4: Example configuration of simulated dataset in 2-dimensions with $\sigma = 0.1$ and $t = 0.75$. The red, filled in circles show the locations of the six normal components for the rare class distribution. The black, empty circles show the location of the 15 normal components for the majority class distribution, whose centers lie along the dotted lines connecting all two rare class normal components.

362 obtained by the ranking methods on four of the training and testing sets as λ is varied. We observe
 363 that the difference between RANK-SVM and the restricted basis models (RANK-RND and RANK-
 364 RC) decreases as λ is increased. Since restricting basis functions also limits the complexity of
 365 the model, the test loss of RANK-RND and RANK-RC is lower than that of RANK-SVM for small
 366 values of λ . However, RANK-RND is unable to achieve the optimal test loss levels at moderate
 367 values of λ (more noticeably in Figures 5c and 5d). In contrast, RANK-RC does not forfeit any test
 368 performance compared to RANK-SVM, while providing additional robustness as λ is reduced.

369 6.2. Real Datasets

370 In this section we compare methods on several unbalanced real datasets obtained from various
 371 sources. Table 3 lists the datasets along with their characteristics. For each dataset, three-quarters
 372 of the observations are used for training and the remaining one-quarter for out-of-sample testing.
 373 Results are averaged over 20 stratified random splits of the data. The model parameter (λ or k) is
 374 tuned by running 10-fold cross-validation on the training set for each split.

375 Table 4 shows the mean test AUC score with standard error for each model. Overall, RANK-
 376 SVM and RANK-RC yield the best performing models with statistically similar results. RANK-RND,
 377 on the other hand, under performs on some datasets, indicating that a random basis set is not as
 378 effective as the rare class basis on unbalanced problems. SVM based methods generally do not
 379 perform as well as ranking methods, except when there appears to be more discriminative patterns
 380 in the data.

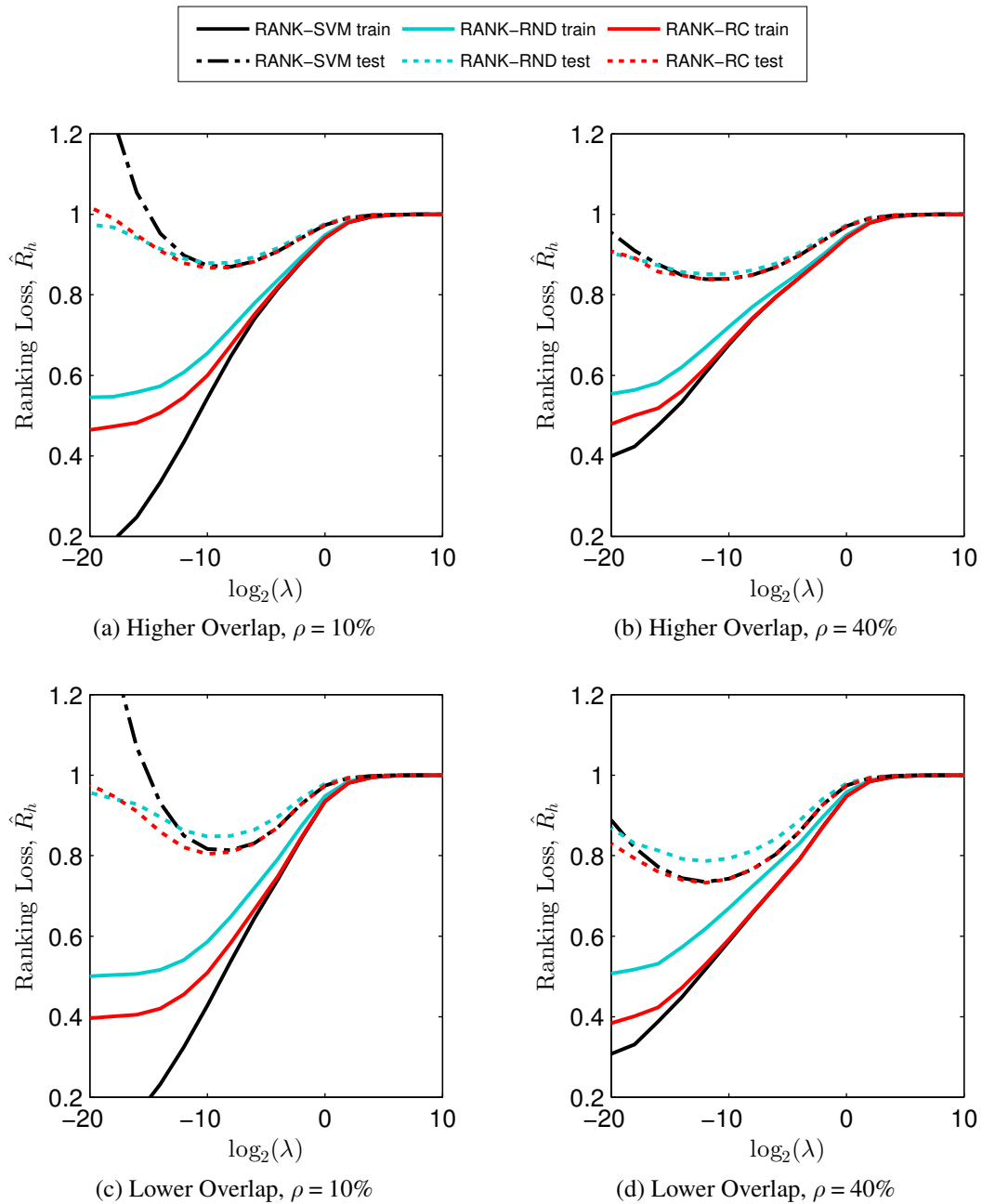


Figure 5: Comparison of empirical train and test ranking loss obtained by the ranking methods on four of the simulated datasets as λ is varied.

381 Table 5 compares the number of support vectors used by the SVM and ranking models. RANK-
382 SVM uses more support vectors than SVM based models. It can use even more support vectors than
383 SVM-SMT, which is trained on an enlarged dataset almost twice the size. This suggests that training
384 RANK-SVM using a working-set type algorithm, which only tracks active support vectors (e.g. as
385 proposed in [24] for standard SVMs), would still run costly in time and space. In comparison,
386 RANK-RND and RANK-RC use significantly fewer support vectors. Moreover, with RANK-RC, test
387 performance is also not compromised.

388 6.3. Intrusion Detection

389 In this section we use the KDD Cup 1999 dataset [36] to evaluate a large-scale unbalanced
390 problem. The objective is to detect network intrusion by distinguishing between legitimate (nor-
391 mal) and illegitimate (attack) connections to a computer network. The dataset is a collection of
392 simulated raw TCP dump data over a period of nine weeks on a local area network. The first seven
393 weeks of data is used for training and the last two for test, providing a total of 4 898 431 training
394 observations and 311 029 test cases. We processed the data to remove duplicate entries (as done in
395 [37]) resulting in 1 074 975 training observations and 77 286 test cases. Each observation contains
396 41 features, three of which are categorical and the rest numerical. The three categorical features
397 are protocol (3 categories), service (70 categories) and flag (11 categories). We represent proto-
398 col using three binary features, where each feature is an indicator for one of the three categories.
399 Service and flag categories are replaced by the frequency in the training sample (i.e. probability)
400 corresponding to the event of observing an attack given the category is present. Thus we obtain a
401 total of 43 features. Finally, as done for all datasets, we standardize each feature to zero mean and
402 unit variance.

403 The attack types are grouped in four categories, DOS (Denial of Service), Probing (Surveil-
404 lance, e.g. port scanning), U2R (user to root), R2L (remote to local). Table 6 shows the distribution
405 of attack types in the training and test sets. Together, the U2R and R2L attacks constitute 4.0%
406 of the test dataset, which is a substantial increase compared to the training set, but still a small
407 fraction. Poor results have been reported in literature for identifying the U2R and R2L attacks
408 [38]. In this experiment, we focus on identifying these attack types by forming a binary classifi-
409 cation problem with the positive class representing either a U2R or R2L attack, and the negative
410 class representing all other connection types. Thus the final training set is highly skewed with only
411 0.098% positive instances.

412 We train using 5%, 10%, 25%, 50%, and 75% of the training data. The remaining training
413 data is used for validation. We are unable to train RANK-SVM, even with just 5% of the data
414 (53 749 samples), since the kernel matrix is too large to store in memory (>10GB). Clearly, this
415 is an example where a large-scale solution is necessary to solve the ranking problem. We do not
416 train SVM-SMT due to the large number of samples as well. We are able to train SVM-W using up
417 to 50% of the data. With more samples SVM-W does not converge, likely due to the large number
418 of support vectors which do not fit in the cache.

419 Figure 6a shows test AUC results obtained by different methods as training data is increased.
420 We observe that SVM and SVM-RUS perform poorly. RANK-RC, RANK-RND and SVM-W produce
421 better results, with RANK-RC performing the best. Figures 6b and 6c compare training time and
422 number of support vectors, respectively, as training data is increased. SVM and SVM-RUS train in

Name	Source	Subject	Features		Samples		
			Original	Derived (d)	m	m_+	ρ
Abalone19	UCI	Life	1N,7C	10	4177	32	0.8%
Mammograph	[34]	Life	6C	6	11183	260	2.3%
Ozone	UCI	Environment	72C	72	2536	73	2.9%
YeastME2	UCI	Life	8C	8	1484	51	3.4%
Wine4	UCI	Chemistry	11C	11	4898	183	3.7%
Oil	[35]	Environment	49C	49	937	41	4.4%
SolarM0	UCI	Nature	10N	32	1389	68	4.9%
Coil	KDD	Business	85C	85	9822	586	6.0%
Thyroid	UCI	Life	21N,7C	52	3772	231	6.1%
Libras	UCI	Physics	90C	90	360	24	6.7%
Scene	LibSVM	Nature	294C	294	2407	177	7.4%
YeastML8	LibSVM	Life	103C	103	2417	178	7.4%
Crime	UCI	Economics	122C	100	1994	150	7.5%
Vowel0	Keel	Computer	10C	10	989	90	9.1%
Euthyroid	UCI	Life	18N,7C	42	3163	293	9.3%
Abalone7	UCI	Life	1N,7C	10	4177	391	9.4%
Satellite	UCI	Nature	36C	36	6435	626	9.7%
Page0	Keel	Computer	10C	10	5472	559	10.2%
Ecoli	UCI	Life	7C	7	336	35	10.4%
Contra2	Keel	Life	9C	9	1473	333	22.6%

Table 3: List of datasets and their characteristics that we use to evaluate methods. Under original features, 'N' is used to denote number of nominal features, 'C', is used to denote number of continuous features. We derive d features by converting nominal features to an indicator representation and use continuous features as is. Under samples, m is the total number of observations, m_+ is the number of rare class observations, and $\rho = \frac{m_+}{m}$ is the percentage of rare class examples.

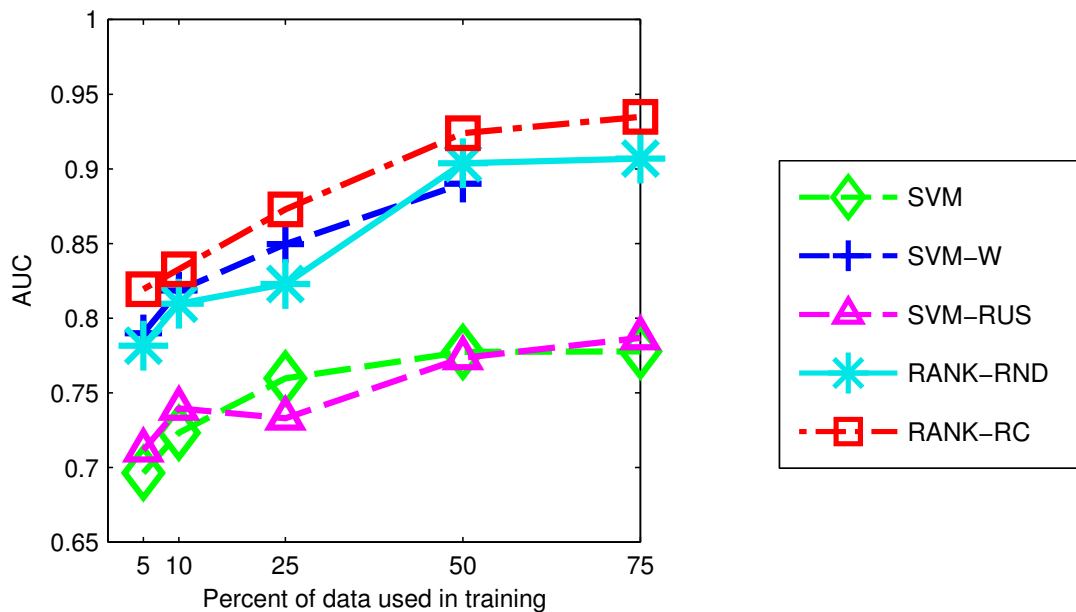
Dataset	KNN	Classification Loss						Ranking Loss			
		SVM	SVM-W	SVM-RUS	SVM-SMT	RANK-SVM	RANK-RND	RANK-RC			
Abalone19	55.7±2.2	54.9±3.1	64.3±1.3	74.1±1.5	67.4±1.1	81.0±1.2	79.1±1.1	81.4±1.1			
Mammograph	80.7±0.4	88.4±0.4	90.1±0.7	92.8±0.4	91.3±0.4	93.7±0.4	93.9±0.4	94.4±0.3			
Ozone	66.4±2.2	85.0±1.1	85.5±0.7	86.4±0.9	85.9±0.8	89.4±0.9	88.7±0.8	90.1±0.9			
YeastME2	69.3±1.7	81.8±0.5	85.5±0.7	87.5±0.7	86.8±1.1	90.8±0.8	89.0±0.9	89.4±1.1			
Wine4	61.9±0.5	74.9±0.7	71.6±0.9	79.1±0.8	78.9±0.8	83.5±0.6	79.6±0.6	82.7±0.7			
Oil	71.6±1.7	91.1±0.9	88.0±1.4	90.6±0.8	90.6±1.0	92.5±0.9	89.2±1.0	91.7±0.8			
SolarM0	58.9±1.6	55.4±0.7	63.1±1.3	71.5±0.8	73.1±0.4	78.5±0.5	77.2±0.8	77.5±0.8			
Coil	53.9±0.3	59.2±0.8	62.9±0.4	68.8±0.4	67.5±0.5	70.0±0.4	69.8±0.4	72.3±0.2			
Thyroid	73.9±0.6	94.8±0.4	93.4±0.5	94.8±0.3	94.4±0.4	95.7±0.4	91.3±0.5	95.7±0.3			
Libras	87.4±2.1	96.8±0.9	96.7±0.9	96.4±0.8	96.8±0.9	97.6±0.8	95.4±1.0	94.8±1.1			
Scene	59.3±0.8	67.3±0.8	75.4±0.9	74.8±0.6	74.0±0.9	77.1±0.7	76.4±0.8	77.5±0.6			
YeastML8	54.5±0.7	57.1±0.8	59.6±0.6	57.9±0.5	59.7±0.4	61.5±0.5	60.2±0.7	62.0±0.5			
Crime	71.8±1.5	87.6±0.7	87.3±0.6	90.1±0.3	90.8±0.3	92.3±0.3	91.2±0.3	91.6±0.3			
Vowel0	100.0±0.0	100.0±0.0	100.0±0.0	99.8±0.0	100.0±0.0	100.0±0.0	98.4±0.1	100.0±0.0			
Euthyroid	75.8±0.8	95.0±0.4	95.0±0.4	94.6±0.4	95.0±0.4	95.2±0.4	90.7±0.4	94.1±0.4			
Abalone7	78.2±2.1	56.1±3.2	76.3±0.5	77.4±0.3	74.4±0.2	87.0±0.3	86.5±0.3	87.1±0.3			
Satellite	83.8±0.3	94.8±0.1	94.6±0.1	94.3±0.1	95.1±0.1	95.3±0.1	94.3±0.1	95.1±0.1			
Page0	90.4±0.4	98.4±0.1	98.1±0.1	98.1±0.1	98.2±0.1	98.6±0.1	95.6±0.1	98.4±0.1			
Ecoli	75.6±2.0	94.6±0.7	93.7±0.7	94.1±0.6	93.2±0.6	94.1±0.6	93.4±0.9	94.5±0.7			
Contra2	60.5±0.9	66.9±0.8	70.2±0.5	70.5±0.6	70.6±0.8	73.2±0.5	72.6±0.5	73.4±0.4			

21

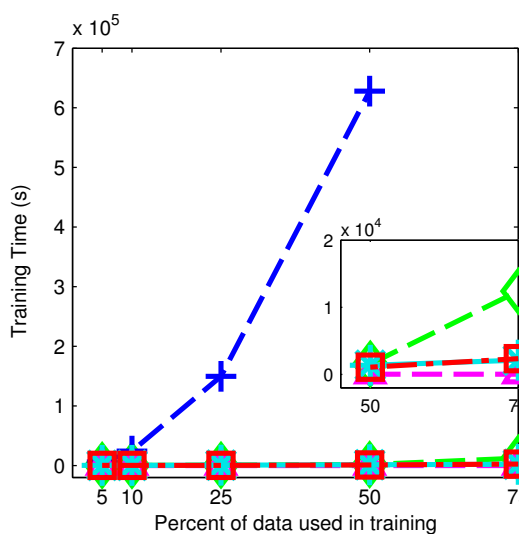
Table 4: Comparison of test AUC results for real datasets (listed in Table 3). Mean AUC score with standard error over 20 trials are shown. Each trial uses one-quarter data for out-of-sample testing. Bolded scores indicate the result is statistically not different than the best performing model using a two-tailed t-test with 99% confidence.

Dataset	Classification Loss					Ranking Loss				
	SVM	SVM-W	SVM-RUS	SVM-SMT	RANK-SVM	RANK-RND	RANK-SVM	RANK-RND	RANK-RC	
Abalone19	117	1555	32	2644	2979	24	2979	24	24	
Mammograph	306	2152	119	2987	7252	195	7252	195	193	
Ozone	206	746	61	1165	995	55	995	55	55	
YeastME2	105	429	41	594	1066	38	1066	38	38	
Wine4	458	2109	179	3166	3550	136	3550	136	136	
Oil	84	311	32	340	488	31	488	31	31	
SolarM0	183	845	90	1114	1042	51	1042	51	51	
Coil	1754	5560	704	8178	7284	435	7284	435	435	
Thyroid	332	723	137	1020	2739	168	2739	168	172	
Libras	35	93	22	124	197	18	197	18	18	
Scene	603	1171	213	1566	1748	132	1748	132	133	
YeastML8	833	1669	257	1562	1804	133	1804	133	133	
Crime	308	631	115	867	1326	112	1326	112	112	
Vowel0	35	37	25	45	730	68	730	68	67	
Euthyroid	389	673	177	1002	2303	216	2303	216	219	
Abalone7	713	1391	274	2076	3079	291	3079	291	292	
Satellite	773	1158	301	1526	4734	466	4734	466	469	
Page0	322	570	145	924	4012	415	4012	415	416	
Ecoli	47	68	18	115	248	26	248	26	26	
Contra2	560	843	396	912	1096	249	1096	249	249	

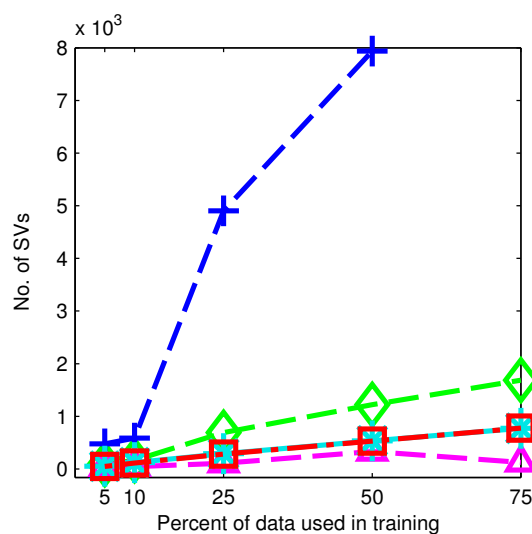
Table 5: Average number of support vectors used by the SVM and ranking models over 20 trials. For ranking models, support vectors are counted as the number of non-zero coefficients associated with kernel functions.



(a)



(b)



(c)

Figure 6: Comparison of (a) test AUC score, (b) training time in seconds, and (c) number of support vectors, for the intrusion detection problem as percent of data used for training is increased from 5% to 75%. In our experiment setup, we were unable to train RANK-SVM due to the large size of the dataset. Also, for more than 50% of data, SVM-W did not converge after more than 72 hours of training.

423 reasonable time, though they do not produce good models. On the other hand, SVM-W quickly
424 becomes very expensive. RANK-RC and RANK-RND scale well, while able to produce effective
425 models. RANK-RC and RANK-RND also use significantly fewer support vectors than SVM-W.

426 7. Conclusion

427 In this paper, we use a ranking loss function to tackle the problem of learning from unbal-
428 anced datasets. Minimizing biclass ranking loss is equivalent to maximizing the AUC measure,
429 which overcomes the inadequacies of accuracy, used by conventional classification algorithms.
430 The resulting regularized loss minimization problem corresponds to a biclass RankSVM problem.
431 We modify RankSVM to take advantage of the rare class situation by restricting the solution to a
432 linear combination of rare class kernel functions (RankRC). This allows us to solve the nonlinear
433 ranking problem in $O(mm_+)$ time and $O(mm_+)$ space, thus enabling us to solve problems which are
434 too large for kernel RankSVM. We provided heuristic and theoretical justification for this choice
435 and experimentally illustrated the effectiveness of RankRC, in both test performance and training
436 time.

437 Below we list a few extensions/variants one may consider using the rare class representation:

- 438 1. Regularization: In problem (14) we can use an ℓ_1 -regularizer, $\|\beta\|_1$, instead of $\beta^T K_{++}\beta$.
439 This would lead to sparser solutions [39] and could be solved using coordinate descent
440 methods [40].
- 441 2. Loss function: We can replace the loss function with other variants of ranking loss. The
442 AUC concentrates uniformly across all threshold levels. We can use weighted AUC [41] or
443 the p-norm push [42] to emphasize specific portions of the AUC curve. Also, we can use
444 list based ranking methods to optimize other criteria such as F_1 -score or Precision/Recall
445 breakeven point [43]. The rare-class representation allows us to learn a nonlinear function
446 for unbalanced datasets with more complex loss functions, in reasonable time and space.
- 447 3. Stochastic Learning: For very large datasets, the $m \times m_+$ kernel submatrix may be too large
448 to fit in memory. In this case, we can store $K_{++} \in \mathbb{R}^{m_+ \times m_+}$ and cycle (randomly) through
449 majority class examples updating the $\beta \in \mathbb{R}^{m_+}$ vector via gradient descent using an adaptive
450 learning rate [44]. Unlike standard stochastic gradient descent, in each iteration we use the
451 full set of minority examples and a single (or small subset) of majority samples to perform
452 the update. This should lead to faster convergence while using only $O(m_+m_+)$ space.

453 In summary, the rare class representation offers significant benefits to learn nonlinear models
454 for large-scale rare class problems.

455 References

- 456 [1] N. Japkowicz, S. Stephen, The class imbalance problem: A systematic study, *Intell. Data Anal.* 6 (2002)
457 429–449.
- 458 [2] B. Raskutti, A. Kowalczyk, Extreme re-balancing for svms: a case study, *SIGKDD Explor. Newsl.* 6 (2004)
459 60–69.
- 460 [3] G. Wu, E. Y. Chang, Class-boundary alignment for imbalanced dataset learning, in: *ICML, 2003*, pp. 49–56.

- 461 [4] G. E. A. P. A. Batista, R. C. Prati, M. C. Monard, A study of the behavior of several methods for balancing
462 machine learning training data, *SIGKDD Explor. Newsl.* 6 (2004) 20–29.
- 463 [5] N. V. Chawla, N. Japkowicz, A. Kotcz, Editorial: special issue on learning from imbalanced data sets, *SIGKDD*
464 *Explor. Newsl.* 6 (2004) 1–6.
- 465 [6] G. M. Weiss, Mining with rarity: a unifying framework, *SIGKDD Explor. Newsl.* 6 (2004) 7–19.
- 466 [7] K. Ezawa, M. Singh, S. W. Norton, Learning goal oriented bayesian networks for telecommunications risk
467 management, in: *ICML*, 1996, pp. 139–147.
- 468 [8] J. Zhang, I. Mani, KNN Approach to Unbalanced Data Distributions: A Case Study Involving Information
469 Extraction, in: *ICML*, 2003.
- 470 [9] M. Kubat, S. Matwin, Addressing the curse of imbalanced training sets: One-sided selection, in: *ICML*, 1997,
471 pp. 179–186.
- 472 [10] N. V. Chawla, K. W. Bowyer, L. O. Hall, W. P. Kegelmeyer, Smote: Synthetic minority over-sampling technique,
473 *Journal of Artificial Intelligence Research* 16 (2002) 321–357.
- 474 [11] R. Herbrich, T. Graepel, K. Obermayer, Large Margin Rank Boundaries for Ordinal Regression, MIT Press,
475 2000.
- 476 [12] O. Chapelle, S. S. Keerthi, Efficient algorithms for ranking with svms, *Inf. Retr.* 13 (2010) 201–215.
- 477 [13] F. Provost, T. Fawcett, R. Kohavi, The case against accuracy estimation for comparing induction algorithms, in:
478 *In Proceedings of the Fifteenth International Conference on Machine Learning*, 1997, pp. 445–453.
- 479 [14] M. A. Maloof, Learning when data sets are imbalanced and when costs are unequal and unknown, in: *ICML*,
480 2003.
- 481 [15] Y. Sun, M. S. Kamel, A. K. C. Wong, Y. Wang, Cost-sensitive boosting for classification of imbalanced data,
482 *Pattern Recogn.* 40 (2007) 3358–3378.
- 483 [16] H. He, E. A. Garcia, Learning from imbalanced data, *IEEE Trans. on Knowl. and Data Eng.* 21 (2009) 1263–
484 1284.
- 485 [17] A. P. Bradley, The use of the area under the roc curve in the evaluation of machine learning algorithms, *Pattern*
486 *Recognition* 30 (1997) 1145–1159.
- 487 [18] C. E. Metz, Basic principles of ROC analysis., *Seminars in nuclear medicine* 8 (1978) 283–298.
- 488 [19] J. A. Hanley, B. J. Mcneil, The meaning and use of the area under a receiver operating characteristic (ROC)
489 curve., *Radiology* 143 (1982) 29–36.
- 490 [20] E. R. DeLong, D. M. DeLong, D. L. Clarke-Pearson, Comparing the Areas under Two or More Correlated
491 Receiver Operating Characteristic Curves: A Nonparametric Approach, *Biometrics* 44 (1988) 837–845.
- 492 [21] P. L. Bartlett, M. I. Jordan, J. D. McAuliffe, Convexity, Classification, and Risk Bounds, *Journal of the American*
493 *Statistical Association* 101 (2006) 138–156.
- 494 [22] B. E. Boser, I. M. Guyon, V. N. Vapnik, A training algorithm for optimal margin classifiers, in: *Proceedings of*
495 *the fifth annual workshop on Computational learning theory, COLT '92, ACM*, 1992, pp. 144–152.
- 496 [23] V. N. Vapnik, *Statistical Learning Theory*, 1st ed., Wiley, 1998.
- 497 [24] O. Chapelle, Training a support vector machine in the primal, *Neural Comput.* 19 (2007) 1155–1178.
- 498 [25] G. Kimeldorf, G. Wahba, A correspondence between Bayesian estimation of stochastic processes and smoothing
499 by splines, *Ann. Math. Statist.* 41 (1970) 495–502.
- 500 [26] B. Schölkopf, R. Herbrich, A. J. Smola, A generalized representer theorem, in: *Proceedings of the 14th Annual*
501 *Conference on Computational Learning Theory and and 5th European Conference on Computational Learning*
502 *Theory*, 2001, pp. 416–426.
- 503 [27] T. Joachims, Optimizing search engines using clickthrough data, in: *KDD*, 2002, pp. 133–142.
- 504 [28] M. Zhu, W. Su, H. A. Chipman, LAGO: A Computationally Efficient Approach for Statistical Detection, *Tech-*
505 *nometrics* 48 (2006).
- 506 [29] S. Boyd, L. Vandenberghe, *Convex Optimization*, Cambridge University Press, 2004.
- 507 [30] M. A. Branch, T. F. Coleman, Y. Li, A subspace, interior, and conjugate gradient method for large-scale bound-
508 constrained minimization problems, *SIAM J. Scientific Computing* 21 (1999) 1–23.
- 509 [31] E. E. Osuna, R. Freund, F. Girosi, *Support Vector Machines: Training and Applications*, Technical Report, MIT,
510 1997.
- 511 [32] C.-C. Chang, C.-J. Lin, LIBSVM: A library for support vector machines, *ACM Trans-*

- 512 actions on Intelligent Systems and Technology 2 (2011) 27:1–27:27. Software available at
513 <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- 514 [33] J. C. Platt, Advances in kernel methods, MIT Press, 1999, pp. 185–208.
- 515 [34] K. Woods, C. Doss, K. Bowyer, J. Solka, C. Preibe, P. Keglmeyer, Comparative evaluation of pattern recognition
516 techniques for detection of microcalcifications in mammography, International Journal of Pattern Recognition
517 and Artificial Intelligence 7 (1993) 1417–1436.
- 518 [35] M. Kubat, R. Holte, S. Matwin, Machine learning for the detection of oil spills in satellite radar images, Machine
519 Learning 30 (1998) 195–215.
- 520 [36] KDD Cup 1999, <http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html>, Accessed:
521 2013-08-31.
- 522 [37] M. Tavallaei, E. Bagheri, W. Lu, A. A. Ghorbani, A detailed analysis of the kdd cup 99 data set, in: Pro-
523 ceedings of the Second IEEE international conference on Computational intelligence for security and defense
524 applications, CISDA'09, 2009, pp. 53–58.
- 525 [38] M. Sabhnani, Application of machine learning algorithms to kdd intrusion detection dataset within misuse
526 detection context, in: ICML, 2003, pp. 209–215.
- 527 [39] R. Tibshirani, Regression shrinkage and selection via the lasso, J. Roy. Statist. Soc. Ser. B 58 (1996) 267–288.
- 528 [40] J. Friedman, T. Hastie, R. Tibshirani, Regularization paths for generalized linear models via coordinate descent,
529 Journal of Statistical Software (2009).
- 530 [41] C. G. Weng, J. Poon, A new evaluation measure for imbalanced datasets, in: Seventh Australasian Data Mining
531 Conference (AusDM 2008), volume 87, 2008, pp. 27–32.
- 532 [42] C. Rudin, The p-norm push: A simple convex ranking algorithm that concentrates at the top of the list, J. Mach.
533 Learn. Res. 10 (2009) 2233–2271.
- 534 [43] T. Joachims, A support vector method for multivariate performance measures, in: ICML, 2005, pp. 377–384.
- 535 [44] L. Bottou, O. Bousquet, The tradeoffs of large scale learning, in: NIPS, 2008, pp. 161–168.

	Training		Test	
Normal	812808	75.6%	47913	62.0%
DOS	247266	23.0%	23568	30.5%
Probing	13850	1.3%	2677	3.5%
U2R	52	0.005%	215	0.278%
R2L	999	0.093%	2913	3.769%
Total	1074975	100%	77286	100%

Table 6: Distribution of connection types in training and test sets for the intrusion detection problem.