Transductive Optimization of Top k Precision

Li-Ping Liu Thomas G. Dietterich

EECS, Oregon State University Corvallis, OR 97330, USA

{liuli@eecs.oregonstate.edu, tgd@oregonstate.edu }

Nan Li Zhi-Hua Zhou

Department of Computer Science & Technology, Nanjing University
Nanjing 210023, China
{lin, zhouzh}@lamda.nju.edu.cn

Abstract

Consider a binary classification problem in which the learner is given a labeled training set, an unlabeled test set, and is restricted to choosing exactly k test points to output as positive predictions. Problems of this kind—transductive precision@k—arise in information retrieval, digital advertising, and reserve design for endangered species. Previous methods separate the training of the model from its use in scoring the test points. This paper introduces a new approach, Transductive Top K (TTK), that seeks to minimize the hinge loss over all training instances under the constraint that exactly k test instances are predicted as positive. The paper presents two optimization methods for this challenging problem. Experiments and analysis confirm the importance of incorporating the knowledge of k into the learning process. Experimental evaluations of the TTK approach show that the performance of TTK matches or exceeds existing state-of-the-art methods on 7 UCI datasets and 3 reserve design problem instances.

1 Introduction

The standard approach to this problem is to first train a classifier on the training data and then threshold the predicted test set scores to obtain the k highest-scoring test instances. Any binary classification algorithm that outputs continuous scores (e.g., an SVM) can be employed in this two-step process. Better results can often be obtained by bipartite ranking algorithms [1–6], which seek to minimize a ranking loss (including ranking losses that put more weight on highly-ranked instances). In addition to precision@k, evaluation measures such as discounted cumulative gain (DCG), normalized discounted cumulative gain (NDCG), and average precision (AP) are often employed to train and evaluate these models.

Recent work focuses even more tightly on the top-ranked instances. The Accuracy At The Top (AATP) algorithm [7] seeks to optimize the ranking quality for a specified top quantile of the training data. Maximizing accuracy on the top quantile is intractable, so AATP optimizes a relaxation of the original objective.

Unlike the ranking problems discussed so far, our problem is transductive, because we have the unlabeled test examples available. There is a substantial body of research on transductive classification [8–10]. The goal of transductive classification is to develop a classifier that will perform well on the entire test set. Most transductive classification algorithms are inspired by either the large margin principle or the clustering principle. The large margin principle asserts that the decision boundary should correctly classify the training examples and pass through a low-density, sparse region of the test data. The clustering principle assumes that the classes form clusters, and that those clusters can be more accurately identified by including test points in the training process.

In the transductive precision@k problem, the goal is to select k positive test points—the classifier or ranker is merely an intermediate step. Vapnik's principle [11] dictates that we should not solve a more difficult problem on the way to solving the problem of interest. Hence, in this paper we jointly train the model and determine the threshold to obtain exactly k predicted positive test instances. Our hypothesis is that the algorithm can take advantage of knowing the value of k to choose a decision boundary that optimizes the precision of the k predicted points. Note that this goal means that neither the large margin principle nor the clustering principle apply directly. A decision boundary that predicts exactly k positives and has high precision on the training data is likely to pass through the region (or cluster) of positive instances rather than through a low-density region. Hence, the intuitions that underlie transductive classification do not provide guidance for solving the transductive precision@k problem.

The paper proceeds as follows. We call the constraint that the model must predict exactly k test instances as positives the k-constraint. We start by identifying a deterministic relation between the precision @k measure and the accuracy of any classifier that satisfies the k-constraint. This suggests that the learning objective should maximize classifier accuracy subject to the k-constraint. We adopt the space of linear decision boundaries and introduce an algorithm we call Transductive optimization of Top k precision (TTK). In the TTK optimization problem, the objective is to minimize the hinge loss on the training set subject to the k-constraint. This optimization problem is very challenging. It can be formulated a Mixed Integer Programming (MIP) problem. For small problems, the global optimum can be found by the branch-and-bound algorithm. To solve larger problems, we design a *feasible direction* method, which we find experimentally to converge very rapidly. An experiment comparing the feasible direction method to the exact MIP solution shows that solutions found by the feasible direction method are nearly optimal. We also present a theoretical analysis of the transductive precision @k problem which shows that one should train different scoring functions for different values of k.

In the experiment section, we first present a small synthetic dataset to show how the TTK algorithm improves the SVM decision boundary. Then, we compare the TTK algorithm with five other algorithms on seven UCI datasets and three reserve design datasets. The results show that the TTK algorithm matches or exceeds the performance of these state-of-the-art algorithms on almost all of these datasets.

2 The TTK model

Let the distribution of the data be \mathcal{D} with support in $\mathcal{X} \times \mathcal{Y}$. In this work, we assume $\mathcal{X} = \mathcal{R}^d$ and only consider the binary classification problem with $\mathcal{Y} = \{-1,1\}$. By sampling from \mathcal{D} independently, a training set $(\mathbf{x},\mathbf{y}) = (x_i,y_i)_{i=1}^n$ and a test set $(\hat{\mathbf{x}},\hat{\mathbf{y}}) = (\hat{x}_j,\hat{y}_j)_{j=1}^m$ are obtained, but the labeling $\hat{\mathbf{y}}$ of the test set is unknown. The problem is to train a classifier and maximize the precision at k on the test set. The hypothesis space is $\mathcal{H} \subset \mathcal{Y}^{\mathcal{X}}$ (functions mapping from \mathcal{X} to \mathcal{Y}). The hypothesis $h \in \mathcal{H}$ is evaluated by the measure precision@k.

When we seek the best classifier from \mathcal{H} for selecting k instances from the test set $\hat{\mathbf{x}}$, we only consider classifiers satisfying the k-constraint, that is, these classifiers must be in the hypothesis space $\mathcal{H}_k(\hat{\mathbf{x}}) = \{h \in \mathcal{H} | \sum_{j=1}^m \mathcal{I}[h(\hat{x}_j) = 1] = k\}$, where $\mathcal{I}[\cdot]$ is 1 if its argument is true and 0 otherwise. All classifiers not predicting k positives on the test set are excluded from \mathcal{H}_k . Note

that with any two-step method with ranking and thresholding, the final classifier is in the hypothesis space $\mathcal{H}_k(\hat{\mathbf{x}})$.

To maximize the precision of $h \in \mathcal{H}_k(\hat{\mathbf{x}})$ on the test set, we essentially need to maximize the classification accuracy of h. This can be seen by the following relation. Let m_- be the number of negative test instances, and let $m_{\rm tp}$, $m_{\rm fp}$ and $m_{\rm tn}$ denote the number of true positives, false positives, and true negatives (respectively) on the test set as determined by h. Then the precision@k of h can be expressed as

$$\rho(h) = \frac{1}{k} m_{\rm tp} = \frac{1}{k} (m_{\rm tn} + k - m_{-}) = \frac{1}{2k} (m_{\rm tp} + m_{\rm tn} + k - m_{-}). \tag{1}$$

Since the number of negative test instances m_{-} is unknown but fixed, there is a deterministic relationship between the accuracy $(m_{\rm tp}+m_{\rm tn})/m$ and the precision@k on the test set. Hence, increasing classification accuracy directly increases the precision. This motivates us to maximize the accuracy of the classifier on the test set while respecting the k-constraint.

In this section, we develop a learning algorithm for linear classifiers and thus $\mathcal{H} = \{h : \mathcal{X} \mapsto \mathcal{Y}, h(x; w, b) = \text{sign}(w^{\top}x + b)\}$. Our learning objective is to minimize the (regularized) hinge loss on the training set, which is a convex upper bound of the zero-one loss. Together with the k-constraint, the optimization problem is

$$\min_{w,b} \frac{1}{2} ||w||_{2}^{2} + C \sum_{i=1}^{n} [1 - y_{i} (w^{\top} x_{i} + b)]_{+},$$

$$s.t. \qquad \sum_{j=1}^{m} \mathcal{I}[w^{\top} \hat{x}_{j} + b > 0] = k ,$$
(2)

where $[\cdot]_+ = \max(\cdot, 0)$ calculates the hinge loss on each instance. Due to the piece-wise constant function in the constraint, the problem is very hard to solve.

Let us relax the equality constraint to an inequality constraint. The optimization problem becomes

$$\min_{w,b} \frac{1}{2} \|w\|_{2}^{2} + C \sum_{i=1}^{n} [1 - y_{i} (w^{\top} x_{i} + b)]_{+},$$

$$s.t. \qquad \sum_{j=1}^{m} \mathcal{I}[w^{\top} \hat{x}_{j} + b > 0] \leq k.$$
(3)

This relaxation generally does not change the solution to the optimization problem. If we neglect the constraint, then the solution that minimizes the objective will be an SVM. In our applications, there are typically significantly more than k positive test points, so the SVM will usually predict more than k positives. In that case, the inequality constraint will be active, and the relaxed optimization problem will give the same solution as the original problem 1 .

Even with the relaxed constraint, the problem is still hard, because the feasible region is non-convex. We first express the problem as a Mixed Integer Program (MIP). Let G be a large constant and η be a binary vector of length m. Then we can write the optimization problem as

The equivalence of (3) and (4) is easy to show, and we omit the proof here. For this MIP, a globally optimal solution can be found for small problem instances via the branch-and-bound method with an off-the-shelf package. We used Gurobi[12].

¹In the extreme case that many data points are nearly identical, the original problem may not have a solution, and the relaxed constraint is satisfied by the "less than" relation. This is unlikely to arise in practice.

For large problem instances, finding the global optimum of the MIP is impractical. We propose to employ a *feasible direction* algorithm [13], which is an iterative algorithm designed for constrained optimization. In each iteration, it first finds a descending direction in the feasible direction cone and then calculates a step size to make a descending step that leads to a improved feasible solution. The feasible direction algorithm fits this problem well. Because the constraint is a polyhedral cone, a step in any direction within the feasible direction cone will generate a feasible solution provided that the step length is sufficiently small. Since our objective is convex but the constraint is highly non-convex, we want to avoid making descending steps along the constraint boundary in order to take bigger steps and avoid local minima.

In each iteration, we first need to find a descending direction. The subgradient $(\nabla w, \nabla b)$ of the objective with respect to (w, b) is calculated as follows. Let $\xi_i = 1 - y_i \ (w^\top x_i + b)$ be the hinge loss on instance i. Then

$$\nabla w = w - C \sum_{i:\xi_i > 0} y_i x_i , \quad \nabla b = -C \sum_{i:\xi_i > 0} y_i.$$
 (5)

We need to project the negative subgradient $(-\nabla w, -\nabla b)$ to a feasible direction to get a feasible descending direction. Let L, E, and R be the sets of test instances predicted to be positive, predicted to be exactly on the decision boundary, and predicted to be negative:

$$L = \{j: w^{\top} \hat{x}_i + b > 0\}, E = \{j: w^{\top} \hat{x}_i + b = 0\}, R = \{j: w^{\top} \hat{x}_i + b < 0\}.$$

With the k-constraint, the feasible direction cone can be written as

$$\mathcal{F} = \left\{ (d_w, d_b) : \sum_{j \in E} \mathcal{I}[\hat{x}_j^\top d_w + d_b > 0] + |L| \le k \right\}.$$
 (6)

We do not project the negative gradient onto \mathcal{F} , both because this is computationally difficult and because it often gives a direction along the boundary of \mathcal{F} . Instead, we find a descending direction by projecting the negative gradient into the null space of a set $B\subseteq E$ of test instances. We first sort the instances in E in descending order according to the value of $-\hat{x}_j^\top \nabla w - \nabla b$. Let $j': 1 \leq j' \leq |E|$ reindex the instances in this order. To construct the set B, we start with $B=\emptyset$ and the initial direction $(d_w,d_b)=-(\nabla w,\nabla b)$. The starting index is $j_0=1$ if |L|=k, and $j_0=2$ if |L|< k. Then with index j' starting from j_0 and increasing, we consecutively put instance j' into B and project (d_w,d_b) into the null space of $\{(\hat{x}_{j^\circ},1):j^\circ\in B\}$. We stop at $j'=j_1$ when all the remaining instances in E have negative inner product with (d_w,d_b) . The final projected direction is denoted by (d_w^*,d_b^*) . The direction (d_w^*,d_b^*) has non-positive inner product with all instances with indices from $j'=j_0$ to j'=|E|, so these instances will not move into the set L when (w,b) moves in that direction. Only when |L|< k, is the first instance allowed to move from E to L. It is easy to check that the final projected direction (d_w^*,d_b^*) is in the feasible cone $\mathcal F$ and that it is a descending direction. This subgradient projection algorithm is summarized in Algorithm 1.

In this design, we have the following considerations. When |L| < k, the instance that has the largest inner product with the negative subgradient is not used to constrain the projected direction, since this can avoid a very hard constraint on the projection. We allow at most one instance to move from E to L to reduce the chance that (w,b) hits the boundary. In the projecting iterations, instances with large inner products are selected first to reduce the number of projections.

Once a descending direction is chosen, we perform a line search to determine the step size. We first find the maximum step size α that guarantees the feasibility of the descending step. That is, no points in R will cross the decision boundary and enter S with the step length α .

$$\alpha = \min_{j \in R : \hat{x}_{j}^{\top} d_{w}^{\star} + d_{b} > 0} \frac{-(\hat{x}_{j}^{\top} w + b)}{\hat{x}_{j}^{\top} d_{w}^{\star} + d_{b}}.$$
 (7)

Then we do a line search in $[0, 0.5\alpha]$ to find the best step length α^* . Note that the objective function is a convex piece-wise quadratic function, so we only need to check these elbow points plus a

²By "inner product" between a direction (d_w, d_b) and an instance x, we mean $x^{\top} d_w + d_b$.

Algorithm 1 Find a descending feasible direction

```
Input: subgradient (\nabla w, \nabla b), instance set \{\hat{x}_j : j \in E\}, size |L|, k

Output: descending feasible direction (d_w^\star, d_b^\star)

Sort instances in E in descending order according to -\hat{x}_j^\top \nabla w - \nabla b

Initialize (d_w, d_b) = -(\nabla w, \nabla b)

Initialize B = \emptyset

j_0 = \min(k - |L|, 1) + 1

for j' = j_0 to |E| do

if \exists j'' : j' \leq j'' \leq |E|, \hat{x}_{j''}^\top d_w + d_b > 0 then

B = B \cup \{j'\}

project (d_w, d_b) into the null space of \{(\hat{x}_{j^\circ}, 1) : j^\circ \in B\}

else

break

end for

d_w^\star = d_w, d_b^\star = d_b
```

minimum between two elbow points to find the best step length. We omit the details. The shrinkage 0.5 of α reduces the chance of (w, b) hitting the boundary.

We initialize w by training a standard linear SVM (although any linear model can be used) and then initialize b to satisfy the k-positive constraint. This gives us a pair (w,b) that is a feasible solution to (3). Then (w,b) is updated in each iteration according to $(w,b) := (w,b) + \alpha^{\star}(d_w^{\star},d_b^{\star})$ until convergence.

We set the maximum number of iterations, T, to 500; the algorithm typically requires only 200-300 iterations to converge. In each iteration, the two most expensive calculations are computing the subgradient and projecting the negative subgradient. The first calculation requires O(nd) operations, and the second one takes $O(d^3)$ operations, since there are usually no more than (d+1) instances in the set E. The overall running time is the time of training an initial model plus $O(T(nd+d^3))$.

Though motivated differently, the AATP algorithm solves a similar optimization problem. The AATP objective is equivalent to ours—but applied to the training set. Their constraint is that the top q quantile of training instances must receive positive scores and all others, negative scores. The AATP authors assume that the decision boundary must go though a training instance, so their relaxation of the optimization problem is constrained to require one instance to be on the decision boundary, classify at least quantile q training examples as positive, and minimize the training set hinge loss. This tends to find solutions that classify more than quantile q of the instances as positive. We will see below that when we compute the exact optimum using an MIP solver, that solution often has many instances on the decision boundary. This demonstrates that the AATP relaxation is quite loose.

3 Analysis

Before presenting experiments, we first argue that different values of k require us, in general, to train different models. We work with the population distribution \mathcal{D} instead of with samples, and we assume linear models. Suppose the distributions of positive instances and negative instances have probability measures μ_+ and μ_- defined on \mathcal{R}^d . The total distribution is a mixture of the two distributions, and it has measure $\mu = \lambda \mu_+ + (1 - \lambda)\mu_-$ with $\lambda \in (0,1)$. The classifier (w,b) defines a positive region $R_{w,b} = \{x \in \mathcal{R}^d, w^\top x + b > 0\}$. Assume $\mu_+(R_{w,b})$ and $\mu_-(R_{w,b})$ are both differentiable with respect to (w,b). If we consider classifier that classify fraction q of the instances as positive, then $\mu(R_{w,b}) = q$. The precision of the classifier will be $\lambda \mu_+(R_{w,b}) / q$. The optimal classifier is therefore

$$(w^*, b^*) = \underset{(w,b)}{\arg \max} \lambda \mu_+(R_{w,b})$$

 $s.t. \quad \lambda \mu_+(R_{w,b}) + (1 - \lambda)\mu_-(R_{w,b}) = q.$ (8)

If we change q, we might hope that we do not need to modify w^* but instead can just change b^* . However, this is unlikely to work.

Theorem 1 If (w^*, b_1) and (w^*, b_2) are two optimal solutions for (8) with two different quantile values q_1 and q_2 , then

$$\exists s_1, t_1, s_2, t_2, \in \mathbf{R}, \qquad s_1 \frac{\partial \mu_+(R_{w^*, b_1})}{\partial (w^*, b_1)} = t_1 \frac{\partial \mu_-(R_{w^*, b_1})}{\partial (w^*, b_1)},$$

$$\partial \mu_+(R_{w^*, b_1}) \qquad \partial \mu_-(R_{w^*, b_1})$$

$$(9)$$

$$s_2 \frac{\partial \mu_+(R_{w^*,b_2})}{\partial (w^*,b_2)} = t_2 \frac{\partial \mu_-(R_{w^*,b_2})}{\partial (w^*,b_2)}.$$
 (10)

The proof follows directly from the KKT conditions. Note that (9) and (10) are two vector equations. When b_1 is changed into be b_2 , the vectors of partial derivatives, $\partial \mu_+(R_{w^*,b_1})/\partial (w^*,b_1)$ and $\partial \mu_-(R_{w^*,b_1})/\partial (w^*,b_1)$ must change their directions in the same way to maintain optimality. This will only be possible for very special choices of μ_+ and μ_- . This suggests that (w^*,b^*) should be optimized jointly to achieve each target quantile value q.

4 Experimental Tests

4.1 An illustrative synthetic dataset

We begin with a simple synthetic example to provide some intuition for how the TTK algorithm improves the SVM decision boundary, see Figure 1. The dataset consists of 40 training and 40 test instances. The training and testing sets each contain 22 positive and 18 negative instances. Our goal is to select k=4 positive test instances. The bold line is the decision boundary of the SVM. It is an optimal linear classifier both for overall accuracy and for precision@k for k=24. However, when we threshold the SVM score to select 4 test instances, this translates the decision boundary to the dashed line, which gives very poor precision of 0.5. This dashed line is the starting point of the TTK algorithm. After making feasible direction descent steps, TTK finds the solution shown by the dot-dash-dot line. The k test instances selected by this boundary are all positive. Notice that if k=24, then the SVM decision boundary gives the optimal solution. This provides additional intuition for why the TTK algorithm should be rerun whenever we change the desired value of k.

4.2 Effectiveness of Optimization

One way to compare different algorithms is to see how well they optimize the training and test surrogate loss functions. We trained a standard SVM, AATP, TTK (MIP) and TTK (feasible direction) on three UCI^3 data sets: diabetes, ionosphere and sonar. For the SVM and AATP methods, we fit them to the training data and then obtain a top-k prediction by adjusting the intercept term k. We set k to select 5% of the test instances. Table 1 reports the regularized hinge loss on the training set and the hinge loss on the test set. The hyper-parameter k is set to 1 for all methods. The results show that TTK with either solver obtains much lower losses than the competing methods. From the difference between the third (MIP) and the fourth (feasible direction) columns, we can also see that the feasible direction method finds near-optimal solutions.

To understand and compare the behavior of AATP and TTK, we performed a non-transductive experiment (by making the training and test sets identical). We measured the number of training instances that fall on the decision boundary and the fraction of training instances classified as positive (see Table 2). The optimal solution given by the MIP solver always puts multiple instances on the decision boundary, whereas the AATP method always puts a single instance on the boundary. The MIP always exactly achieves the desired k, whereas AATP always classifies many more than k instances as positive. This shows that the AATP assumption that the decision boundary should pass through exactly one training instance is wrong.

4.3 Precision evaluation on real-world datasets

In this subsection, we evaluate our TTK method on ten datasets. Seven datasets {diabetes, iono-sphere, sonar, spambase, splice} from UCI repository and {german-numer, svmguide3} from the

³https://archive.ics.uci.edu/ml/datasets.html

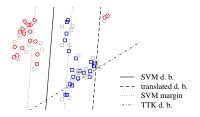


Figure 1: TTK improves the SVM decision boundary ("d. b."). Square/circle: positive/negative instance, colored/gray: training/testing instance. k = 4.

Table 1: Training and test loss attained by different methods

dataset	SVM	AATP	TTK_{MIP}	TTK_{FD}	
diabetes					
train obj.	311 ± 25	265 ± 23	224 ± 7	226 ± 7	
test loss	323 ± 20	273 ± 24	235 ± 6	235 ± 5	
ionosphere					
train obj.	325 ± 46	474 ± 88	127 ± 4	136 ± 4	
test loss	338 ± 44	488 ± 85	146 ± 5	150 ± 7	
sonar					
train obj.	167 ± 52	166 ± 41	20 ± 8	30 ± 10	
test loss	216 ± 22	213 ± 30	103 ± 19	105 ± 24	

Table 2: AATP and TTK solution statistics: Number of instances on the decision boundary ("# at d.b.") and fraction of instances predicted as positive ("fraction +")

dataset, dimension, ratio of positives	AATP		TTK_{MIP}	
dataset, difficultion, ratio of positives	# at d.b.	fraction +	# at d.b.	fraction +
diabetes, $d = 8, n_{+}/n = 0.35$	1	0.12	5	0.05
ionosphere, $d = 33, n_{+}/n = 0.64$	1	0.53	21	0.05
sonar, $d = 60, n_+/n = 0.47$	1	0.46	40	0.05

LIBSVM web site⁴, are widely studied binary classification datasets. The other three datasets, NY16, NY18 and NY88, are three species distribution datasets extracted from a large eBird dataset [14]; each of them has 634 instances and 38 features. The eBird dataset contains a large number of checklists of bird counts reported from birders around the world. Each checklist is associated with the latitude and longitude of the observation and a set of 38 features describing the habitat. We chose a subset of the data consisting of checklists of three species from New York state in June of 2012. To correct for spatial sampling bias, we formed spatial cells by imposing a grid over New York and combining all checklists reported within each grid cell. This gives 634 cells (instances). Each instance is labeled with whether a species was present or absent in the corresponding cell.

We compare the TTK algorithm with 5 other algorithms. The SVM algorithm [15] is the baseline. The Transductive SVM [8] compared here (denoted by TSVM) uses the UniverSVM [16] implementation, which optimizes its objective with the convex-concave procedure. SVMperf [17] can optimize multiple ranking measures and it is parameterized here to optimize precision@k. Two algorithms, Accuracy At The Top (AATP) [7] and TopPush [6], are specially designed for top precision optimization. The proposed TTK objective is solved using both the MIP solver and the feasible direction method (denoted TTK_{MIP} and TTK_{FD}, respectively). Each algorithm is run 10 times on 10 random splits of each dataset. Each of these algorithms requires setting the regularization parameter C. This was done by performing five 2-fold internal cross-validation runs within each training set and selecting the value of C from the set $\{0.01, 0.1, 1, 10, 100\}$ that maximized precision on the top 5% of the (cross-validation) test points. With the chosen value of C, the algorithm was then run on the full training set (and unlabeled test set) and the precision on the top 5% was measured. The achieved precision values were then averaged across the 10 independent runs.

Table 3 shows the performance of the algorithms. For datasets with more than 1000 instances, the AATP and $TTK_{\rm MIP}$ algorithms do not finish within a practical amount of time, so results are not reported for these algorithms on those datasets. This is indicated in the table by "NA". The results for each pair of algorithms are compared by a paired-differences t-test at the p < 0.05 significance level. If one algorithm is not significantly worse than any of the other algorithms, then it is regarded as one the best and its performance is shown in bold face. Wins, ties and losses of of $TTK_{\rm MIP}$ and $TTK_{\rm FD}$ with respect to all other algorithms are reported in the last two rows of Table 3.

On each of the six small datasets, the performance of $TTK_{\rm MIP}$ matches or exceeds that of the other algorithms. The $TTK_{\rm FD}$ method does almost as well—it is among the best algorithms on 8 of the 10 datasets. It loses once to SVMperf (on symguide3) and once to AATP (on ionosphere). None of the

⁴http://www.csie.ntu.edu.tw/ cjlin/libsvmtools/datasets/binary.html

Table 3: Mean Precision (\pm 1 standard deviation	n) of classifiers when %5 of testing instances are
predicted as positives.	

dataset	SVM	TSVM	SVMperf	TopPush	AATP	TTK_{MIP}	TTK_{FD}
diabetes	.86±.08	.86±.09	$.69 \pm .20$.80±.10	.68±.28	.85±.10	.86±.08
ionosphere	.76±.13	$.80 \pm .17$	$.82 \pm .22$	$.71 \pm .16$	$1.00 \pm .00$.97±.05	$.84 {\pm} .15$
sonar	.96±.08	$.98 \pm .06$	$.85 \pm .16$	$.88 \pm .13$	$.90 \pm .11$.96±.08	$1.00\pm.00$
german-numer	.70±.08	$.72 \pm .08$	$.56 \pm .17$	$.63 \pm .12$	NA.	NA.	$.71 \pm .06$
splice	$1.00 \pm .00$	$1.00\pm.00$	$1.00\pm.01$	$1.00\pm.00$	NA.	NA.	$\boldsymbol{1.00 {\pm}.00}$
spambase	.97±.02	$.97 \pm .02$	$.98 \pm .01$	$.96 \pm .02$	NA.	NA.	$.98 \pm .01$
svmguide3	.86±.07	$.85 \pm .07$	$.91 \pm .04$	$.83 \pm .07$	NA.	NA.	$.87 \pm .06$
NY16	.64±.08	$.64 \pm .09$	$.65\pm.12$	$.62 \pm .10$	$.62 \pm .08$.68±.07	$.70\pm.09$
NY18	$.44 {\pm} .11$	$.45 \pm .10$	$.36 \pm .07$	$.43 {\pm} .13$	$.46 \pm .12$.46±.08	$.47 \pm .12$
NY88	.40±.08	$.33 \pm .12$	$.37 \pm .15$	$.34 \pm .08$	$.31 \pm .09$.40±.09	$.42 \pm .07$
TTK _{MIP} w/t/l	1/5/0	2/4/0	2/4/0	1/5/0	2/4/0		
$TTK_{\rm FD}$ w/t/l	3/7/0	4/6/0	3/6/1	7/3/0	4/1/1		

other methods performs as well. By comparing TTK_{FD} with SVM, we see that the performance is improved on almost all datasets, so the TTK_{FD} method can be viewed as a safe treatment of the SVM solution. As expected, the transductive SVM does not gain much advantage from the availability of the testing instances, because it seeks to optimize accuracy rather than precision@k. The TopPush algorithm is good at optimizing the precision of the very top instance. But when more positive instances are needed, the TopPush algorithm does not perform as well as TTK.

5 Summary

This paper introduced and studied the transductive precision@k problem, which is to train a model on a labeled training set and an unlabeled test set and then select a fixed number k of positive instances from the testing set. Most existing methods first train a scoring function and then adjust a threshold to select the top k test instances. We show that by learning the scoring function and the threshold together, we are able to achieve better results.

We presented the TTK method. The TTK objective is the same as the SVM objective, but TTK imposes the constraint that the learned model must select exactly k positive instances from the testing set. This constraint guarantees that the final classifier is optimized for its target task. The optimization problem is very challenging, since it involves a set selection problem. We introduced two algorithms for solving it. First, we formulated it as a mixed integer program and solved it exactly via the branch-and-bound method. Second, we designed a feasible direction algorithm that is able to scale to larger datasets but that attempts to find a good solution. We compared both TTK algorithms to several state-of-the-art methods on ten datasets. The results indicate that the performance of the TTK methods matches or exceeds all of the other algorithms on most of these datasets.

Our analysis and experimental results show that the TTK objective is a step in the right direction. However, we believe that the performance can be further improved if we can minimize a tighter (possibly non-convex) bound on the zero-one loss. In the future, we will extend the TTK formulation to problems with multiple labels so that it can be applied to reserve design problems involving multiple species.

References

- [1] C. Burges, T. Shaked, E. Renshaw, A. Lazier, M. Deeds, N. Hamilton, and G. Hullender. Learning to rank using gradient descent. In *Proceedings of the 22nd International Conference on Machine Learning*, 2005.
- [2] C. Rudin. The p-norm push: A simple convex ranking algorithm that concentrates at the top of the list. *Journal of Machine Learning Research*, 10, 2009.
- [3] N. Usunier, D. Buffoni, and P. Gallinari. Ranking with ordered weighted pairwise classification. In *Proceedings of the 26th International Conference on Machine Learning*, 2009.
- [4] S. Agarwal. The infinite push: a new support vector ranking algorithm that directly optimizes accuracy at the absolute top of the list. In *Proceedings of the SIAM International Conference on Data Mining*, 2011.
- [5] A. Rakotomamonjy. Sparse support vector infinite push. In *Proceedings of the 29th International Conference on Machine Learning*, 2012.
- [6] N. Li, R. Jin, and Z.-H. Zhou. Top rank optimization in linear time. In Advances in Neural Information Processing Systems 27, 2014.
- [7] S. Boyd, C. Cortes, M. Mohri, and A. Radovanovic. Accuracy at the top. In Advances in Neural Information Processing Systems 25, 2012.
- [8] T. Joachims. Transductive inference for text classification using support vector machines. In *Proceedings* of the 16th International Conference on Machine Learning, 1999.
- [9] V. Sindhwani and S. S. Keerthi. Large scale semi-supervised linear syms. In Proceedings of the 29th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, 2006.
- [10] D. Pechyony. Theory and Practice of Transductive Learning. PhD thesis, Department of Computer Science, Technion-Israel Institute of Technology, 2008.
- [11] V. N. Vapnik. Statistical Learning Theory. Wiley-Interscience, 1998.
- [12] Inc. Gurobi Optimization. Gurobi optimizer reference manual, 2015. URL http://www.gurobi.
- [13] M. S. Bazaraa, H. D Sherali, and C. M. Shetty. *Nonlinear Programming: Theory and Algorithms*. Wiley Publishing, 3rd edition, 2006.
- [14] B. L. Sullivan, C. L. Wood, M. J. Iliff, R. E. Bonney, D. Fink, and S. Kelling. ebird: a citizen-based bird observation network in the biological sciences. *Biological Conservation*, 142, 2009.
- [15] B. Schölkopf and A. J. Smola. Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond. Adaptive computation and machine learning. MIT Press, 2002.
- [16] F. Sinz and M. Roffilli. Universym, 2012. http://mloss.org/software/view/19/.
- [17] T. Joachims. A support vector method for multivariate performance measures. In *Proceedings of the 22nd International Conference on Machine Learning*, 2005.