

# MULTICLASS RIDGE-ADJUSTED SLACK VARIABLE OPTIMIZATION USING SELECTED BASIS FOR FAST CLASSIFICATION

Yinan Yu<sup>a,c</sup>, Konstantinos I. Diamantaras<sup>b</sup>, Tomas McKelvey<sup>a</sup>, S.Y. Kung<sup>c</sup>

Chalmers University of Technology <sup>a</sup>  
Gothenburg, Sweden  
{yinan,tomas.mckelvey}@chalmers.se

TEI of Thessaloniki <sup>b</sup>  
Thessaloniki, Greece  
kdiamant@it.teithe.gr

Princeton University <sup>c</sup>  
Princeton, USA  
{yinany, kung}@princeton.edu

## ABSTRACT

Kernel techniques for classification is especially challenging in terms of computation and memory requirement when data fall into more than two categories. In this paper, we extend a binary classification technique called Ridge-adjusted Slack Variable Optimization (RiSVO) to its multiclass counterpart where the label information encoding scheme allows the computational complexity to remain the same to the binary case. The main features of this technique are summarized as follows: (1) Only a subset of data are pre-selected to construct the basis for kernel computation; (2) Simultaneous active training set selection for all classes helps reduce complexity meanwhile improving robustness; (3) With the proposed active set selection criteria, inclusion property is verified empirically. Inclusion property means that once a pattern is excluded, it will no longer return to the active training set and therefore can be permanently removed from the training procedure. This property greatly reduce the complexity. The proposed techniques are evaluated on standard multiclass datasets MNIST, USPS, pendigits and letter which could be easily compared with existing results.

**Index Terms**— RiSVO, kernel, multiclass classification, large scale data, RKHS basis construction

## 1. INTRODUCTION

Multi-class classification problems are of growing importance in a vast range of domains, such as handwriting identification, medical applications, posture recognition, etc, where data fall into  $C$  categories for  $C > 2$ . In the literature, there are mainly two ways of solving such problems:

- (a) Dividing the task into several binary problems using binary classifiers such as SVM [17]. Typical approaches for training and combining the binary classifiers are the One-versus-All (OvA) [12] and All-versus-All (AvA) [16] techniques.
- (b) Direct extensions of binary classifiers using a global objective function. Existing techniques can be found in the literature, such as extensions of SVM [19], Neural Networks [4], Decision Trees [5], K-Nearest Neighbors [3], Naive Bayes [11], etc.

Techniques in category (a) enjoy the benefits of being intuitive and straightforward. However, more training time is

required since more than one binary classifiers are needed, i.e.  $C$  for OvA and  $\frac{C(C-1)}{2}$  for AvA. Nevertheless, classifiers of AvA are in general of much smaller scales than OvA and possibly result in less computations.

On the other hand, (b) has the following features:

- Compared to OvA, intrinsic structure within each class is preserved and results in more balanced classifiers.
- A global optimization problem is formulated and training is carried out simultaneously for all the classes. However, the potential drawback is that this typically results in a huge number of parameters and constraints.
- Outputs for all classes usually have a fair comparison and no normalization is needed.

Due to these properties, we are interested in finding a simultaneous technique without major increase in the complexity compared to its binary version. To this end, we propose an iterative technique extending [20] to handle multiclass case. The reduction of computational complexity is summarized in Table 1.

At each iteration  $k$ , given training size  $N$ , for  $m \ll N_{S,k} \ll N_{S,k-1} \leq N$  and  $N_{S,0} = N$ ,

| Procedure                          | Complexity reduction   |
|------------------------------------|--|
| Active set select.<br>(Sec. 2.2.1) | Kernel eval: $\mathcal{O}(N_{S,k-1}^2) \rightarrow \mathcal{O}(N_{S,k}^2)$   |
| Inclusion prop.<br>(Sec. 2.2.2)    | Training: $\mathcal{O}(N_{S,k-1}^3) \rightarrow \mathcal{O}(N_{S,k}^3)$  |
| RKHS basis<br>select. (Sec. 3)     | Kernel eval: $\mathcal{O}(N_{S,k}^2) \rightarrow \mathcal{O}(N_{S,k}m)$<br>Training: $\mathcal{O}(N_{S,k}^3) \rightarrow \mathcal{O}(m^3)$ |

The paper is organized as follows. First of all, the multi-classification formulation is presented in Section 2. In Section 3, we present a technique to construct the basis matrix for kernel computations by using a subset of the training data. Numerical results are shown in Section 4.

## 2. MULTICLASS RISVO FORMULATION

### 2.1. Formulation

Given training data  $\mathbf{x}_i$ , a nonlinear mapping  $\varphi : \mathbb{R}^d \rightarrow \mathbb{R}^p$  ( $d \ll p$ ) maps training data  $\mathbf{x}_i$  to a high dimensional feature space. Kernel tricks [15] are usually used when computations

in such high dimensional space are prohibitive. That is, instead of explicitly constructing  $\varphi(\mathbf{x}_i)$ , we replace all inner products by  $k(\mathbf{x}_i, \mathbf{x}_i) = \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_i)$ , where  $k(\cdot, \cdot)$  is a pre-defined kernel function. Without ambiguity, we refer  $\varphi_i$  as the training data for convenience. Given training data  $\varphi_i$  and corresponding class label  $t_i \in \{1, \dots, C\}$ , from the representer theorem, the class prediction  $\hat{y}$  of an unseen data vector  $\varphi$  is defined as:

$$\hat{y} = \arg \max_c \left\{ \sum_{\varphi_i \in G} A(i, c) \varphi_i^T \varphi + b(c) \right\}, \quad (1)$$

where  $A(i, c)$  is the  $(i, c)^{th}$  entry of matrix  $\mathbf{A}$  and  $b(c)$  is the  $c^{th}$  element of the bias vector  $\mathbf{b}$ . Matrix  $\begin{bmatrix} \mathbf{A} \\ \mathbf{b}^T \end{bmatrix}^*$  is referred to as the ‘‘solution matrix’’ learned from the training process. Set  $G$  is a subset of the whole training data to construct the basis for kernel computations. It is defined in Sec. 3.

Similar to RiSVO, the multiclass version is an iterative method works as follows. Let matrix  $\Phi_M = [\varphi_{m_1}, \dots, \varphi_{m_M}]$ ,  $\forall \varphi_{m_i} \in M \subseteq D$ , where  $D$  is the whole training set. At each iteration  $k \geq 1$ , the optimization problem is cast as

$$\begin{bmatrix} \mathbf{A}^k \\ \mathbf{b}^{kT} \end{bmatrix}^* = \arg \min_{\mathbf{A}, \mathbf{b}} \frac{1}{2} \left( E(\boldsymbol{\xi}^T \boldsymbol{\xi} | \varphi_i \in S^k) + \rho \|\Phi_G \mathbf{A}\|^2 \right), \quad (2)$$

where  $S^k \subseteq S^{k-1}$  is the active set defined in Section 2.2.1 with  $S^0 = D$ ,  $\rho$  the ridge parameter [9]  $E(\cdot)$  the expectation operator and  $\|\cdot\|$  denotes the Frobenius norm. Vector  $\boldsymbol{\xi}$  is called the slack vector. It is closely related to the coding scheme of the label information. Note that although  $\boldsymbol{\xi}$  is  $k$  dependent, to simplify the notation, we only specify the step number explicitly for  $S^k$ , and everything else follows.

Now let us define  $\boldsymbol{\xi}$ . For each given training pattern  $\varphi_i \in S^k$ , let  $\mathbf{T}_i$  be a diagonal matrix and the entry  $\mathbf{T}_i(j, j) = \begin{cases} +1 & \text{if } t_i = j \\ -1 & \text{otherwise} \end{cases}$  and

$$\boldsymbol{\lambda}_i^k = \mathbf{T}_i (\mathbf{A}^{(k-1)T} \Phi_G^T \varphi_i + \mathbf{b}^{(k-1)}) \quad (3)$$

The multiclass slack vector  $\boldsymbol{\xi}_i$  is thus defined as:

$$\boldsymbol{\xi}_i = \mathbf{e} - \boldsymbol{\lambda}_i^k \quad (4)$$

where  $\mathbf{e}$  is a all-one vector.

Let us define the kernel matrix  $\mathbf{K}_{AB} = \Phi_A^T \Phi_B$  and  $\mathbf{K}_A = \Phi_A^T \Phi_A$ ,  $\mathbf{E}$  an all-one matrix, and  $N_M$  the cardinality of set  $M$ . Without loss of generality, we drop the iteration number  $k$ , the solution of (2) can be found by taking the derivative to zero, which results in

$$\begin{bmatrix} \mathbf{A}^k \\ \mathbf{b}^{kT} \end{bmatrix}^* = \begin{cases} [\mathbf{K}_S + \rho \mathbf{I} & , & \mathbf{e}]^+ \mathbf{\Pi}, & \text{for } G = S \\ \Omega^+ (\mathbf{K}_{GS} + \mathbf{E}) \mathbf{\Pi}, & \text{otherwise} \end{cases} \quad (5)$$

where

$$\Omega = [(\mathbf{K}_{GS} + \mathbf{E}) \mathbf{K}_{SG} + \rho \mathbf{K}_G, \quad \mathbf{K}_{GSe} + N_S \mathbf{e}] \quad (6)$$

and  $\mathbf{\Pi} = [\mathbf{T}_i \mathbf{e}]_{\varphi_i \in S}$ .

Therefore, the computational complexity scales as  $\mathcal{O}(m^3)$ , where  $m$  is the size of the basis  $\Phi_G$  in the Reproducing Kernel Hilbert Space (RKHS) (cf. Sec. 3).

## 2.2. Active set selection and inclusion property

Two key ideas of multiclass RiSVO algorithm are **active set selection** and the **inclusion property**, where the former improves the classification performance and the latter reduces computational complexity. Both properties are verified empirically in Sec. 4.

### 2.2.1. Rules for Active Set Selection

Starting from  $S^0 = D$ , at each iteration  $k \geq 1$ , we select the active set  $S^k \subset S^{k-1}$  and the iterations are terminated when  $S^k = S^{k-1}$ . The underlying idea is to improve robustness by discarding (i) **the well classified** patterns and (ii) the ‘‘**outliers**’’, i.e. the selected patterns need to fall between two marginal hyperplanes. This idea is used by LASVM [2] in binary classification, whereas we extend it to a multiclass version, where the inclusion property is preserved. More precisely, for all  $\varphi_i \in S^{k-1}$  and  $j \in \{1, \dots, C\}$ , the selection criteria at step  $k$  are:

- $C_1^k = \{\varphi_i : \gamma_l < \boldsymbol{\lambda}_i^k(t_i) < \gamma_u\}$
- $C_2^k = \{\varphi_i : \zeta_l < E(\boldsymbol{\lambda}_i^k(j) | j \neq t_i) < \zeta_u\}$

where  $\boldsymbol{\lambda}_i^k$  is defined in Equation (3) and  $\boldsymbol{\lambda}_i^k(l)$  denotes the  $l^{th}$  entry. Parameters  $\gamma_l, \gamma_u, \zeta_l$  and  $\zeta_u$  are user defined marginal boundaries, where  $\zeta_l > \gamma_l, \zeta_u < \gamma_u$ .

Set  $C_1^k$  means that for a training pattern  $\varphi_i$  from class  $c$  to be selected, it has to fall between two hyperplanes  $\boldsymbol{\lambda}_i^k(t_i) = \gamma_l$  and  $\boldsymbol{\lambda}_i^k(t_i) = \gamma_u$  for dimension  $t_i$ . Similarly,  $C_2^k$  indicates that the averaged position of  $\varphi_i$  on the other  $C - 1$  dimensions should fall between  $E(\boldsymbol{\lambda}_i^k(j) | j \neq t_i) = \zeta_l$  and  $E(\boldsymbol{\lambda}_i^k(j) | j \neq t_i) = \zeta_u$ . Since  $t_i$  is the true label of data  $\varphi_i$ , the condition  $\zeta_l > \gamma_l, \zeta_u < \gamma_u$  indicates that we care more about the correct dimension  $t_i$ .

The active set at iteration  $k$  for multiclass RiSVO is hence defined as

$$S^k = C_1^k \cup C_2^k. \quad (7)$$

### 2.2.2. Inclusion Property

To guarantee and speed up convergence, we always search among  $\varphi_i \in S^{k-1}$  to identify  $S^k$ . This requires the validity of the Inclusion Property. That is, even we search among all training data  $\varphi_i \in D$  at every iteration  $k$ , the following property holds.

$$S^K \subseteq S^{K-1} \subsetneq \dots \subsetneq S^1 \subsetneq S^0 \quad (8)$$

where maximum iteration  $K$  is defined by user or the minimum integer when  $S^K = S^{K-1}$ .

Equation (8) is called the *Inclusion Property*. **Inclusion property means that at each iteration  $k$ , the active set  $S^k$  can be identified by searching within the previous active set  $S^{k-1}$ .** It means that once a pattern is removed from the active set, it will no longer satisfy Eq. (7) at future iterations. The sufficient condition for this property to hold will be presented in the full paper due to the limited scope. In this paper, it is only verified empirically.

A summary of the implementation for Multiclass RiSVO can be found in *Algorithm Multiclass RiSVO* and the code will be available on the authors' website shortly.

---

### Algorithm Multiclass RiSVO

---

- Initialization: let  $S^0 = D$ .
  - Define maximum iteration  $K$ .
  - **Repeat:** iteration  $k$ 
    - Update  $\begin{bmatrix} \mathbf{A}^k \\ \mathbf{b}^{kT} \end{bmatrix}^*$  according to Equation (5)
    - Compute  $\xi_i$  and  $\lambda_i^k$  from Equations (4) and (3)
    - Identify  $S^k$  according to Equations (7) and (8).
    - Stop if  $k = K$  or  $S^k = S^{k-1}$
- 

### 3. RKHS BASIS SELECTION FOR LARGE SCALE DATA

In a linear classifier, we typically have the label estimation

$$\hat{y} = \arg \max_c \{ \mathbf{w}_c^T \varphi + b(c) \}, \quad (9)$$

where  $\mathbf{w}_c$  is the  $c^{\text{th}}$  column of the decision matrix  $\mathbf{W}$ . The relation between Equation (9) and (1) can be explained by the representer theorem [15], where the solution vector can be written as a linear combination of the training data in the Reproducing Kernel Hilbert Space (RKHS), i.e.  $\mathbf{W} = \sum_{\varphi_i \in G} \varphi_i \mathbf{A}(i, \cdot)$ . In another word, the set of vectors  $G = \{ \varphi_1, \varphi_2, \dots, \varphi_m \}$  form a *basis*. Ideally speaking, all the training data could be used to construct this basis. However, from Equation (5), we know that the order of the computational complexity depends on the size of the basis  $m$ . Therefore, the computational complexity can be reduced by constructing the basis using a subset of the training data. In the literature, this is often referred to as the kernel approximation problem. The review of such techniques is out of the scope. Related topics can be found in [14, 1, 7, 18, 10, 13]. In the context of finding subspace basis, we develop a selection technique based on subspace projection. By removing data with respect to small (normalized) projection distance, the subspace directions not ‘‘informative’’ enough are therefore rejected. These directions are sometimes considered dominated by noise. To fully explore the subspace concept, the technique is analyzed from three perspectives: (1) numerical, (2) basis reduction and (3) computation. Note that the ‘‘redundant’’ samples are only discarded for the basis construction, but still kept for training due to the rich label information.

Let us fix some notations before we proceed. Some might be different from previous sections due to different purposes. Let the basis matrix  $\Phi_m = [\varphi(\mathbf{x}_1), \varphi(\mathbf{x}_2), \dots, \varphi(\mathbf{x}_m)]$ . Define  $\mathbf{K}_m = \Phi_m^T \Phi_m$  and  $K_m(x_i, x_j) = \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)$  the  $(i, j)^{\text{th}}$  entry of the  $m \times m$  matrix. In the empirical space, when a new pattern  $\varphi(\mathbf{x}_r)$  is included as a part of the basis, the matrix  $\Phi_m$  becomes  $\Phi_{m+1}$  and the corresponding kernel matrix is denoted as  $\mathbf{K}_{m+1}$ . Finally, let vector  $\mathbf{k}_m = [K_m(\mathbf{x}_1, \mathbf{x}_r) \ \dots \ K_m(\mathbf{x}_m, \mathbf{x}_r)]^T$ .

### 3.1. Numerical perspective

First of all, numerical issues are discussed. Let  $\tilde{\Phi}_m = [\varphi(\mathbf{x}_{g_1}), \varphi(\mathbf{x}_{g_2}), \dots, \varphi(\mathbf{x}_{g_m})]$ ,  $g_i \in \mathbb{N}$  be an orthogonal basis in the intrinsic space and  $\angle(\mathbf{A}, \mathbf{B})$  denote the principal angle [8] between matrices  $\mathbf{A}$  and  $\mathbf{B}$ .

**Definition 1** (Highly Redundant Pattern). *Given a new pattern  $\mathbf{x}_r$ , let  $\Phi_{m+1} = [\tilde{\Phi}_m \ \varphi(\mathbf{x}_r)]$ . Pattern  $\mathbf{x}_r$  is called Highly Redundant (w.r.t.  $\tilde{\Phi}_m$ ), if and only if  $|\mathbf{K}_{m+1}| = 0$ , i.e.  $\angle(\Phi_{m+1}, \tilde{\Phi}_m) = 0$ .  $\square$*

The definition tells us that by including a highly redundant pattern, the kernel matrix becomes rank deficient which causes computational issues. Therefore, highly redundant patterns should be discarded during basis construction. However,  $|\mathbf{K}_{m+1}|$  in Definition 1 is generally non-zero. In order to have a meaningful numerical definition, relaxation needs to be introduced, i.e. practically speaking,  $\mathbf{x}_r$  is considered Highly Redundant if  $\mathbf{K}_{m+1}$  is numerically rank deficient.

### 3.2. Basis reduction perspective

For large data set, the purpose of basis construction is not only to avoid numerical issue, but also to reduce the size of the basis in order to scale down the computations with preserved information. To this end, Empirical Redundancy is introduced.

**Definition 2.** *Given a pattern  $\mathbf{x}_r$ , let  $\Phi_{m+1} = [\tilde{\Phi}_m | \varphi(\mathbf{x}_r)]$ . Pattern  $\mathbf{x}_r$  is called empirically non-redundant (w.r.t.  $\tilde{\Phi}_m$ ), if and only if  $\frac{\pi}{2} \geq \angle(\Phi_{m+1}, \tilde{\Phi}_m) > \frac{\pi}{2} - \theta_\epsilon$  for small  $\theta_\epsilon \geq 0$ .  $\square$*

This leads us to Theorem 1.

**Theorem 1.** *Given  $\tilde{\Phi}_m$  and pattern  $\mathbf{x}_r$ ,  $\mathbf{x}_r$  is empirically non-redundant w.r.t.  $\tilde{\Phi}_m$  iff*

$$\left| \frac{\mathbf{k}_{m+1}^T \tilde{\mathbf{K}}_m^{-1} \mathbf{k}_{m+1}}{K_{m+1}(\mathbf{x}_r, \mathbf{x}_r)} \right| < \epsilon. \quad (10)$$

for some small number  $\epsilon \geq 0$ .  $\square$

*Proof.* We denote  $\tilde{\varphi}(\mathbf{x}_r)$  the projection of  $\varphi(\mathbf{x}_r)$  onto subspace  $\tilde{\Phi}_m$ , and  $(\tilde{\Phi}_m^T \tilde{\Phi}_m)^{-1}$ . Thus the angle between  $\varphi(\mathbf{x}_r)$  and  $\tilde{\varphi}(\mathbf{x}_r)$  is computed as

$$\theta = \arg \cos \frac{\varphi(\mathbf{x}_r)^T \tilde{\varphi}(\mathbf{x}_r)}{\|\varphi(\mathbf{x}_r)\| \|\tilde{\varphi}(\mathbf{x}_r)\|}$$

$$\begin{aligned} \text{where } \varphi(\mathbf{x}_r)^T \tilde{\varphi}(\mathbf{x}_r) &= \mathbf{k}_{m+1}^T \tilde{\mathbf{K}}_m^{-1} \mathbf{k}_{m+1} \\ \|\varphi(\mathbf{x}_r)\| &= \sqrt{K_{m+1}(\mathbf{x}_r, \mathbf{x}_r)} \\ \|\tilde{\varphi}(\mathbf{x}_r)\| &= \sqrt{\mathbf{k}_{m+1}^T \tilde{\mathbf{K}}_m^{-1} \mathbf{k}_{m+1}} \end{aligned}$$

This gives us  $\cos^2 \theta = \left| \frac{\mathbf{k}_{m+1}^T \tilde{\mathbf{K}}_m^{-1} \mathbf{k}_{m+1}}{K_{m+1}(\mathbf{x}_r, \mathbf{x}_r)} \right|$ . Since  $\mathbf{x}_r$  empirically non-redundant, i.e.  $\frac{\pi}{2} \geq \angle(\Phi_{m+1}, \tilde{\Phi}_m) \geq \frac{\pi}{2} - \theta_\epsilon$ , we have  $\cos^2 \theta < \epsilon$  for some small  $\epsilon > 0$ . Thus it follows the conclusion.  $\square$

Due to the lack of space, the proof will be given in the full paper. The theorem tells us that as a data selection rule for basis construction, update  $\tilde{\Phi}_{m+1} = \Phi_{m+1}$  is applicable if Equation (10) holds true; otherwise,  $\tilde{\Phi}_{m+1} = \tilde{\Phi}_m$ .

### 3.3. Computational perspective

However, with the increasing size of the basis, the computation of  $\mathbf{K}_m^{-1}$  in Equation (10) becomes costly. By using *matrix inversion lemma* [6],  $\mathbf{K}_m^{-1}$  is updated in an iterative fashion and there are only vector operations involved in the computations.

$$\mathbf{K}_m^{-1} = \begin{bmatrix} \mathbf{K}_{m-1}^{-1} & 0 \\ 0 & 0 \end{bmatrix} + \frac{1}{\mu} \begin{bmatrix} -\mathbf{K}_{m-1}^{-1} \mathbf{k}_m \\ 1 \end{bmatrix} \begin{bmatrix} -\mathbf{k}_m^T \mathbf{K}_{m-1}^{-1} & 1 \end{bmatrix}$$

where  $\mu = K_m(\mathbf{x}_r, \mathbf{x}_r) - \mathbf{k}_m^T \tilde{\mathbf{K}}_{m-1}^{-1} \mathbf{k}_m$ . The implementation of updating the basis is summarized in *Algorithm Basis Construction (B.C.)*.

---

#### Algorithm Basis Construction (B.C.)

---

- Initialization:  
Select the basis size  $m$  and a small number  $0 < \epsilon \ll 1$ . Randomly select  $m_0$  patterns to initialize the basis matrix  $\Phi_0 = [\varphi(\mathbf{x}_{g_1}), \dots, \varphi(\mathbf{x}_{g_{m_0}})]$ . Compute  $\mathbf{K}_0$ , where  $K_0(i, j) = K(\mathbf{x}_{g_i}, \mathbf{x}_{g_j})$ . Let  $q = 0, r = 1, \mathbf{P}_0 = \mathbf{K}_0^{-1}$ .

- Step 1: For every new pattern  $\mathbf{x}_r$ , compute:

$$\mathbf{p}_r = \mathbf{P}_q \mathbf{k}, \quad \delta_r = \frac{\mathbf{k}^T \mathbf{p}_r}{K(\mathbf{x}_r, \mathbf{x}_r)}$$

where  $\mathbf{k} \in \mathbb{R}^{m_q \times 1}$  and  $\mathbf{k}(i) = K(\mathbf{x}_{g_i}, \mathbf{x}_r)$ .

- Step 2: If  $|\delta_r| < \epsilon$ , select the pattern  $\mathbf{x}_r$  and update:

$$\Phi_{q+1} = [\Phi_q \varphi(\mathbf{x}_{g_r})], \quad q = q + 1, \quad m_q = m_0 + q.$$

$$\mathbf{P}_q = \begin{bmatrix} \mathbf{P}_{q-1} & 0 \\ 0 & 0 \end{bmatrix} + \frac{1}{\mu} \begin{bmatrix} -\mathbf{p}_r \\ 1 \end{bmatrix} \begin{bmatrix} -\mathbf{p}_r^T & 1 \end{bmatrix}$$

$$\text{where } \mu = K(\mathbf{x}_r, \mathbf{x}_r) - \mathbf{k}^T \mathbf{p}_r \quad (11)$$

Otherwise ignore the pattern.

- Step 3: if  $m_q < m$ , set  $r \leftarrow r + 1$  and go to Step 1; otherwise, stop.
- 

## 4. EXPERIMENTAL RESULTS

### 4.1. Multiclass RiSVO

In this section, multiclass RiSVO is evaluated on the data sets MNIST, USPS, UCI letters and UCI pendigits. The parameters are chosen by cross-validation using the training patterns and testing error rates are presented in this section. The results and chosen parameters are summarized in Table 1 which could be compared with other existing techniques [21, 22, 23]. For example, SVM gives 2.03% testing error rate on dataset pendigits and 5.75% on letters, whereas our method achieve 0.6% and 2.56%, respectively. If we compare the

number of support vectors, SVM has 1200 on pendigits and 8400 on letters, where multiclass RiSVO uses 1500 and 4000 basis vectors for representing the solution matrix. On the other hand, multiclass RiSVO does not overperform SVM on dataset MNIST and pendigits. Our guess is that RiSVO is essentially a subspace approach, which is not designed for very sparse images without any preprocessing. However, if we take into consideration the computational complexity, for  $C$  classes, multiclass RiSVO still scales as  $\mathcal{O}(m^3)$ . The convergence of three data sets are shown in Fig. 1. The advantages with respect to the computational efficiency are summarized.

- a. From computational perspective, multiclass RiSVO scales as  $\mathcal{O}(m^3)$ , which remains the same as binary RiSVO, i.e. the order of complexity does not depend on the number of classes.
- b. Alg. *Basis Construction (B.C.)* offers a fast and numerically robust tool of reducing the basis size with preserved information.
- c. By using the proposed active set selection criteria, the inclusion property is empirically verified on all datasets. That is, when a data set is excluded at step  $k_0$ , it does not come back to the active set for  $k > k_0$ .
- d. Parallelizability for big data computations. It fits the framework of Map-Reduce as presented in [20], where the data set is divided into  $M$  chunks and computations are carried out independently. The final result of multiclass RiSVO is obtained by various methods, such as averaging over all  $M$  processors.

**Notes on selecting hyperparameters  $\gamma_l, \gamma_u, \zeta_l, \zeta_u$ :** These parameter are chosen from  $\gamma_l, \gamma_u \in \{\pm 1, \pm 1.1, \pm 1.2\}$  and  $\zeta_l, \zeta_u \in \{\pm 0.5, \pm 0.6, \pm 0.7\}$ . For every pattern  $\varphi$ , by setting  $\gamma_l < \zeta_l$  and  $\gamma_u > \zeta_u$ , we are assigning more weights to the “correct” position  $l_i$ .

| $m$             | 100   | 700  | 1000 | 1300 | 1600 | 1900 |
|-----------------|-------|------|------|------|------|------|
| <i>Random</i>   | 10.12 | 4.20 | 3.57 | 3.38 | 3.0  | 2.99 |
| <i>Alg.B.C.</i> | 9.78  | 3.98 | 3.51 | 3.33 | 2.95 | 2.92 |

**Table 2.** Testing error comparison between random basis selection and Alg. Basis Construction (B.C.).

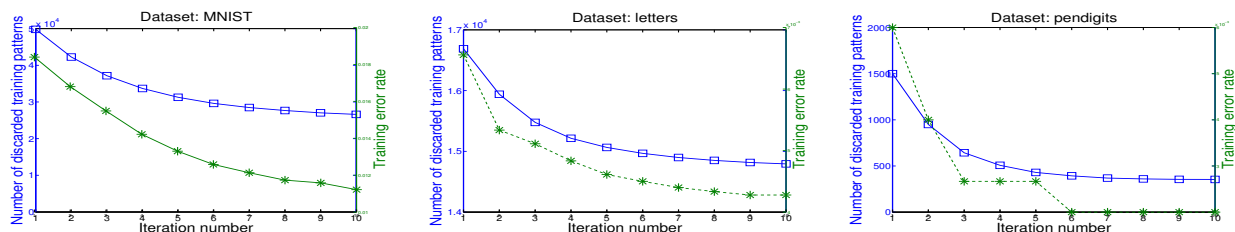
### 4.2. RKHS basis construction

By reducing the size of the basis,  $m$ , the computations of Multiclass RiSVO is reduced, since the complexity mainly depends on  $m$ . However, a random data selection might degrade the performance. A brief comparison between random selection and Alg. Basis Construction (B.C.) in terms of classification error on data set MNIST is shown in Table 2. Note that Alg. Basis Construction (B.C.) scales as  $\mathcal{O}(m^2 N_D)$ .

**Notes on selecting  $\epsilon$ :** according to Section 3.1,  $\epsilon$  **shouldn't be chosen too large, i.e.  $\epsilon \rightarrow 1$** . Otherwise at the next iteration,  $\delta_r \rightarrow 0$  and  $\frac{1}{\mu} \rightarrow +\infty$ . The problem will thus become ill-conditioned. Therefore,  $0 < \epsilon \ll 1$  (c.f. Algorithm Basis Construction (B.C.) and Sec.3.1).

| Data               |    |                    | Parameters        |                 |            |                        |                      |        | $P_{err}$ |
|--------------------|----|--------------------|-------------------|-----------------|------------|------------------------|----------------------|--------|-----------|
| Name               | C  | Training data      | Kernel            | Basis $\Phi(G)$ |            | RiSVO (Eq.2, 7)        |                      |        |           |
|                    |    |                    |                   | $m$             | $\epsilon$ | $(\gamma_l, \gamma_u)$ | $(\zeta_l, \zeta_u)$ | $\rho$ |           |
| MNIST[21]          | 10 | $784 \times 60000$ | poly, $d = 2$     | 3500            | 0.5        | $(-1.1, 1.1)$          | $(-0.5, 0.5)$        | 0.001  | 2.33%     |
| USPS[22]           | 10 | $256 \times 3035$  | poly, $d = 2$     | 3035            | -          | $(-1.1, 1.1)$          | $(-0.5, 0.5)$        | 0.001  | 4.81%     |
| UCI Letters[23]    | 26 | $16 \times 20000$  | rbf, $\sigma = 5$ | 4000            | 0.3        | $(-1.2, 1.2)$          | $(-0.7, 0.7)$        | 0.001  | 2.56%     |
| UCI pendigits [23] | 10 | $16 \times 1500$   | rbf, $\sigma = 1$ | 1500            | -          | $(-1, 1)$              | $(-0.5, 0.5)$        | 0.001  | 0.60%     |

**Table 1.** The data sets are from standard multiclass databases, which can be found in the corresponding references.



**Fig. 1.** Convergence of number of discarded patterns and classification error on the training set for three data bases.

## 5. CONCLUSION

In this paper, we present a multi-classification technique aiming for high efficiency and robustness. The technique is based on three important features. (1) Active set selection for improving robustness; (2) The empirical validity of Inclusion property for reducing complexity; and (3) Selection  $G \subsetneq D$  to construct the basis for representing the solution matrix. This can be viewed as a preselection of the support vectors from a subspace projection viewpoint. As a future direction, we will provide further empirical and theoretical analysis on multi-class RiSVO.

## 6. REFERENCES

- [1] Bach F. and Jordan M. I., Predictive low-rank decomposition for kernel methods, in *Proceeding of ICML*, 2005.
- [2] Bordes A., Ertekin S., Weston J., and Bottou L., Fast Kernel Classifiers with Online and Active Learning, *Journal of Machine Learning Research*, 6:1579-1619, 2005.
- [3] Bay S. D., *Combining Nearest Neighbor Classifiers Through Multiple Feature Subsets*, Proc. of the 15th Intl Conf on Machine Learning, pp. 37-45, Morgan Kaufmann Publishers Inc., USA, 1998
- [4] Bishop C. M., *Neural Networks for Pattern Recognition* Oxford University Press, 1995.
- [5] Breiman L., Friedman J., Olshen R. A., and Stone C. J., *Classification and Regression Trees*, 1984.
- [6] Diamantaras K. I. and Kung SY, *Principal Component Neural Networks: Theory and Applications*, 1996.
- [7] Drineas P., and Mahoney M. W., On the Nyström method for approximating a Gram matrix for improved kernel-based learning. *JMLR*, 6:2153 – 2175, 2005.
- [8] Golub G.H., Van Loan, C. F., *Matrix Computation*, 3rd edition, The Johns Hopkins University Press, 1996.
- [9] Hoerl A. E. and Kennard R. W. *Ridge regression: Biased estimation for nonorthogonal problems*, *Technometrics*, 42(1):80-86, 1970.
- [10] Rahimi, A. and Recht B., Random features for large-scale kernel machines. In *Proceeding of NIPS*, 2007.
- [11] Rennie J.D.M., *Improving multi-class text classification with naive Bayes*, Tech. Report, MIT, 2001.
- [12] Rifkin R. and Klautau A., *In Defense of One-Vs-All Classification*, *JMLR*, pp. 101-141, 2004.
- [13] B. Savas and I.S. Dhillon, Clustered low rank approximation of graphs in information science applications. In *SDM*, 2011.
- [14] Schölkopf B., Knirsch P, Smola A. and Burges C., *Fast Approximation of Support Vector Kernel Expansions, and an Interpretation of Clustering as Approximation in Feature Spaces*, *Mustererkennung, Informatik aktuell*, pp 125-132, 1998.
- [15] Schölkopf, B. and Smola, A. J., *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond* 1st Edition, The MIT Press, Dec. 2001.
- [16] Hastie T. and Tibshirani R. *Classification by pairwise coupling*, Editors: Jordan M. I., Kearns M. J., and Solla S. A., *Advances in Neural Information Processing Systems*, vol. 10. The MIT Press, 1998
- [17] Vapnik, V., *Statistical Learning Theory*, 1st Edition, Wiley-Interscience, September, 1998.
- [18] Wang H., Hu Z. and Zhao Y., *Kernel Principal Component Analysis for Large Scale Data Set*, *Lecture Notes in Computer Science*, vol. 4113, pp 745-756, 2006.
- [19] Weston J. and Watkins C., *Support vector machines for multiclass pattern recognition*, Proc. of the 7th European Symposium on Artificial Neural Networks, 1999.
- [20] Yu Y., Diamantaras K. I., McKelvey T., and Kung S.Y., *Ridge-adjusted Slack Variable Optimization for Supervised Classification*, *IEEE Intl Workshop on Machine Learning for Signal Processing*, UK, 2013.
- [21] <http://yann.lecun.com/exdb/mnist>
- [22] <http://www.kernel-machines.org/data.html>
- [23] <http://www.ics.uci.edu/mllearn/MLRepository.html>