

Huberized Multiclass Support Vector Machine for Microarray Classification

LI Jun-Tao¹ JIA Ying-Min¹

Abstract This paper proposes a new multiclass support vector machine (SVM) for simultaneous gene selection and microarray classification. Combining the huberized hinge loss function and the elastic net penalty, the proposed SVM can perform automatic gene selection and encourages a grouping effect. The coefficient paths of the proposed SVM are shown to be piecewise linear with respect to the single regularization parameter, based on which the solution path algorithm is developed with low computational complexity. Experiments performed on the leukemia data set are provided to verify the obtained results.

Key words Gene selection, grouping effect, microarray classification, solution path, support vector machine (SVM)

DOI 10.3724/SP.J.1004/2010.00399

The DNA microarray technology^[1-2] is a powerful tool for biological and medical research. To classify microarray gene expression data, support vector machine (SVM) has attracted considerable attention^[3-12]. It is well known that SVM and its various extensions have been successfully applied to two-class classification of microarray data^[3-6]. In recent years, a few attempts have been made to generalize the SVM to multiclass problems^[7-12]. In particular, multiclass support vector machine (MSVM)^[8] that considers all of the classes at once has been developed. However, this method selects genes by using the marginal criterion^[1], which tends to yield redundant genes. To perform joint classification and gene selection, L_1 -norm MSVM^[9] and sup-norm MSVM^[10] have also been developed. Unfortunately, they cannot reveal the mutual information among genes.

The group lasso^[13-14] has been developed to select relevant variables in groups for two-class problems. However, rare results on extending this method to multiclass problems have been reported since it is difficult to construct gene clusters in advance. When used as two-class classifiers, the elastic net penalized methods^[4-5, 15] appear to perform well on microarray data and encourage a grouping effect. The key feature of these methods is the use of the elastic net penalty which not only retains the benefits of the L_1 norm penalty but also tends to generate similar coefficients for highly correlated variables. Taking into account the advantages of the elastic net penalty, this paper is devoted to extending it to multiclass problem and further developing a new MSVM. This is still a challenging problem due to the following two facts:

1) Gene selection of multiclass problems becomes more complex than the binary case. This is because that the MSVM requires to estimate multiple discriminating functions, among which each function has its own subset of important predictors. Moreover, gene selection is a necessary demand for microarray classification, especially grouped gene selection. Hence, the grouping effect for multiclass gene selection should be encouraged.

2) The normal solving algorithms for MSVM, such as quadratic programming^[7-8] and linear programming^[9-10], fail to work since multiple penalties are required simulta-

neously. Furthermore, the piecewise linear regularization path algorithms^[16-17] used for parameter selection are not suitable due to multiple classes and multiple tuning parameters. Hence, new solving algorithm and parameter tuning method should be developed.

This paper attempts to deal with the aforementioned difficulties by developing new statistical learning tool. To this end, we first propose the huberized multiclass support vector machine (HMSVM). Second, the grouping effect of HMSVM described by using 2-norm is presented. Third, we prove that the coefficients of the HMSVM are piecewise linear with respect to the single regularization parameter and give their concrete form. Next, an efficient regularized solution algorithm is developed to compute the optimal coefficients. Finally, we apply our method to leukemia classification and achieve promising results.

1 Problem statement

Assume that the n training pairs $\{(\mathbf{x}_i, y_i), i = 1, \dots, n\}$ are independently and identically distributed according to an unknown probability distribution $P(\mathbf{x}, y)$. For microarray gene expression data with multiple cancer types, \mathbf{x}_i represents the i -th sample and y_i represents the tumor type, which can be coded as $\{1, \dots, K\}$. Let $\mathbf{f} = (f_1, \dots, f_K)$ denote the decision function vector, where $k = 1, \dots, K$. A popular multiclass classifier^[10] can be defined as

$$\phi(\mathbf{x}) = \arg \max_{k=1, \dots, K} f_k(\mathbf{x}) \quad (1)$$

The K -class classification problem is to learn the decision function vector and hence accurately predict the cancer type of a new sample.

There are some popular machine learning methods for the K -class classification problem, e.g., MSVM^[8], 1-norm MSVM^[9], and sup-norm MSVM^[10]. However, all these methods cannot reveal the mutual information among genes. This paper aims to deal with this problem by developing a new MSVM. In the following, we present our notation. Similar to sup-norm MSVM^[10], we use the linear decision functions $f_k(\mathbf{x}) = b_k + \mathbf{w}_k^T \mathbf{x}$, $k = 1, 2, \dots, K$ to build the classifier, and let $\mathbf{b} = (b_1, \dots, b_K)^T$ denote the bias vector, w denote the coefficient matrix, $X = (\mathbf{x}_1; \mathbf{x}_2; \dots; \mathbf{x}_n) = (\mathbf{x}_{(1)}, \mathbf{x}_{(2)}, \dots, \mathbf{x}_{(p)})$ denote the model matrix, where $\mathbf{x}_{(j)} = (x_{1j}, \dots, x_{nj})^T$, $j = 1, \dots, p$. We use $\mathbf{w}_k = (w_{k1}, \dots, w_{kp})^T$ and $\mathbf{w}_{(j)} = (w_{1j}, \dots, w_{Kj})^T$ to represent the k -th row vector and the j -th column vector of w , respectively.

Manuscript received January 6, 2009; accepted March 18, 2009
Supported by National Basic Research Program of China (973 Program) (2005CB321902), National Natural Science Foundation of China (90916024, 60727002, 60774003), the Ph.D. Programs Foundation of Ministry of Education of China (20030006003), and the Commission on Science, Technology, and Industry for National Defense (A2120061303)

1. The Seventh Research Division, Beihang University, Beijing 100191, P. R. China

2 Main results

2.1 Huberized multiclass support vector machine

The elastic net penalty has been successfully applied to binary classification problem^[4-5,15]. Taking into account the advantages of the elastic net penalty, we extend it to K -class classification problem, i.e.,

$$J(\mathbf{b}, \mathbf{w}) = \lambda_2 \sum_{k=1}^K \sum_{j=1}^p w_{kj}^2 + \lambda_1 \sum_{k=1}^K \sum_{j=1}^p |w_{kj}| \quad (2)$$

Similar to binary hybrid huberized SVM^[4], we substitute

$$L_{ki} = \begin{cases} 0, & \text{if } b_k + \mathbf{w}_k^T \mathbf{x}_i < -1 \\ 1 + b_k + \mathbf{w}_k^T \mathbf{x}_i - \frac{\delta}{2}, & \text{if } b_k + \mathbf{w}_k^T \mathbf{x}_i \geq -1 + \delta \\ \frac{(1 + b_k + \mathbf{w}_k^T \mathbf{x}_i)^2}{2\delta}, & \text{otherwise} \end{cases}$$

for hinge loss function $[b_k + \mathbf{w}_k^T \mathbf{x}_i + 1]_+$, where L_{ki} is the abbreviation of $L_{ki}(b_k, \mathbf{w}_k, \mathbf{x}_i)$. Combining the elastic net penalty (2) with the huberized hinge loss function L_{ki} , we propose the following HMSVM

$$\arg \min_{\mathbf{b}, \mathbf{w}} \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K a_{ik} L_{ki} + \lambda_2 \sum_{k=1}^K \sum_{j=1}^p w_{kj}^2 + \lambda_1 \sum_{k=1}^K \sum_{j=1}^p |w_{kj}| \quad (3)$$

subjected to the sum-to-zero constraint^[4]

$$\mathbf{1}^T \mathbf{b} = 0, \quad \mathbf{1}^T \mathbf{w}_{(j)} = 0 \quad (4)$$

where λ_1, λ_2 are regularization parameters, $a_{ik} = I(y_i \neq k)$, and $I(\cdot)$ represents the indicator function. Substituting (4) into (3), we have the following unconstrained convex optimization problem:

$$\arg \min_{\mathbf{b}_-, \mathbf{w}_-} \bar{L}(\lambda_1, \lambda_2, \mathbf{b}_-, \mathbf{w}_-) \quad (5)$$

where $\mathbf{b}_-, \mathbf{w}_-$ denote the vector and matrix formed by the first $K-1$ rows of \mathbf{b} and \mathbf{w} , and \bar{L} is defined as

$$\begin{aligned} \bar{L} = & \frac{1}{n} \left(\sum_{i=1}^n a_{iK} L_{Ki} \left(-\sum_{k'=1}^{K-1} b_{k'}, -\sum_{k'=1}^{K-1} \mathbf{w}_{k'}^T \mathbf{x}_i \right) + \right. \\ & \left. \sum_{i=1}^n \sum_{k=1}^{K-1} a_{ik} L_{ki} \right) + \lambda_2 \left(\sum_{k=1}^{K-1} \sum_{j=1}^p w_{kj}^2 + \sum_{j=1}^p \left(\sum_{k=1}^{K-1} w_{kj} \right)^2 \right) + \\ & \lambda_1 \left(\sum_{k=1}^{K-1} \sum_{j=1}^p |w_{kj}| + \sum_{j=1}^p \left| \sum_{k=1}^{K-1} w_{kj} \right| \right) \end{aligned} \quad (6)$$

Define each region of the k -th classifier as $\mathcal{L}_k = \{i : b_k + \mathbf{w}_k^T \mathbf{x}_i < -1\}$, $\mathcal{E}_k = \{i : -1 \leq b_k + \mathbf{w}_k^T \mathbf{x}_i < -1 + \delta\}$, $\mathcal{R}_k = \{i : b_k + \mathbf{w}_k^T \mathbf{x}_i \geq -1 + \delta\}$ for $k = 1, \dots, K$, define the indices for non-zero w_{kj} as $\mathcal{A}_k = \{j : w_{kj} \neq 0, j = 1, 2, \dots, p\}$ for $k = 1, \dots, K-1$, and define the indices for non-zero $\sum_{k'=1}^{K-1} w_{k'j}$ as $\bar{\mathcal{A}} = \{j : \sum_{k'=1}^{K-1} w_{k'j} \neq 0, j = 1, 2, \dots, p\}$. If we continuously decrease λ_1 or λ_2 or both of them, some of the sets of $\mathcal{L}_k, \mathcal{E}_k, \mathcal{R}_k, \mathcal{A}_k \cap \bar{\mathcal{A}}$ will change. We call this an event, and four types of events may occur: 1) A training point reaches the boundary between \mathcal{L}_k and \mathcal{E}_k ; 2) A training point reaches the boundary between \mathcal{R}_k and \mathcal{E}_k ; 3) An index j leaves $\mathcal{A}_k \cap \bar{\mathcal{A}}$; 4) An index j joins $\mathcal{A}_k \cap \bar{\mathcal{A}}$.

2.2 The grouping effect

Since the correlations between genes sharing the same biological ‘‘pathway’’ can be high, the correlated and relevant genes should be selected or removed together. From the statistical point of view, this can be described as a grouping effect. In the following, we prove that the HMSVM can encourage a grouping effect described by using 2-norm.

Theorem 1. Let $\hat{\mathbf{b}}$ and $\hat{\mathbf{w}}$ denote the optimal solution of the huberized multiclass support machine. If $\mathbf{x}_{(m)}$ and $\mathbf{x}_{(l)}$ are normalized, then

$$\|\hat{\mathbf{w}}_{(m)} - \hat{\mathbf{w}}_{(l)}\|_2 \leq \frac{\sqrt{K}}{\sqrt{n}\lambda_2} \sqrt{2(1-\rho)} \quad (7)$$

where $\rho = \mathbf{x}_{(m)}^T \mathbf{x}_{(l)} = \sum_{i=1}^n x_{im} x_{il}$.

Proof. From the sum-to-zero constraint, it is easy to know that $\hat{\mathbf{b}}_-$ and $\hat{\mathbf{w}}_-$ formed by the first $K-1$ rows of $\hat{\mathbf{b}}$ and $\hat{\mathbf{w}}$ will be the optimal solution of unconstrained convex optimization problem (5). For $k = 1, 2, \dots, K-1, j' = 1, 2, \dots, p$ and any given $1 \leq m, l \leq p$, we construct the following vector and matrix

$$\mathbf{b}_-^* = \hat{\mathbf{b}}_-$$

$$w_{kj'}^* = \begin{cases} \frac{1}{2}(\hat{w}_{km} + \hat{w}_{kl}), & \text{if } j' = m \text{ or } j' = l \\ \hat{w}_{kj'}, & \text{otherwise} \end{cases}$$

By the definition of $\hat{\mathbf{b}}_-, \hat{\mathbf{w}}_-, \mathbf{b}_-^*$, and w_-^* , we have

$$0 \leq \bar{L}(\lambda_1, \lambda_2, \mathbf{b}_-^*, w_-^*) - \bar{L}(\lambda_1, \lambda_2, \hat{\mathbf{b}}_-, \hat{\mathbf{w}}_-) \quad (8)$$

Note that $|L_{ki}(b_k^*, \mathbf{w}_k^*, \mathbf{x}_i) - L_{ki}(\hat{b}_k, \hat{\mathbf{w}}_k, \mathbf{x}_i)| \leq |(\mathbf{w}_k^* - \hat{\mathbf{w}}_k)^T \mathbf{x}_i|$, $|\sum_{k=1}^{K-1} (\hat{w}_{km} - \hat{w}_{kl})(x_{im} - x_{il})| = |(\hat{w}_{Km} - \hat{w}_{Kl})(x_{im} - x_{il})|$. Hence,

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^{K-1} a_{ik} \left[L_{ki}(b_k^*, \mathbf{w}_k^*, \mathbf{x}_i) - L_{ki}(\hat{b}_k, \hat{\mathbf{w}}_k, \mathbf{x}_i) \right] + \\ & \frac{1}{n} \sum_{i=1}^n a_{iK} \left[L_{Ki} \left(-\sum_{k'=1}^{K-1} b_{k'}^*, -\sum_{k'=1}^{K-1} \mathbf{w}_{k'}^T \mathbf{x}_i \right) - \right. \\ & \left. L_{Ki} \left(-\sum_{k'=1}^{K-1} \hat{b}_{k'}, -\sum_{k'=1}^{K-1} \hat{\mathbf{w}}_{k'}^T \mathbf{x}_i \right) \right] \leq \\ & \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^{K-1} \left| (\mathbf{w}_k^* - \hat{\mathbf{w}}_k)^T \mathbf{x}_i \right| + \frac{1}{n} \sum_{i=1}^n \left| \sum_{k=1}^{K-1} (\mathbf{w}_k^* - \hat{\mathbf{w}}_k)^T \mathbf{x}_i \right| = \\ & \frac{1}{2n} \sum_{i=1}^n \sum_{k=1}^{K-1} |(\hat{w}_{km} - \hat{w}_{kl})(x_{im} - x_{il})| + \\ & \frac{1}{2n} \sum_{i=1}^n \left| \sum_{k=1}^{K-1} (\hat{w}_{km} - \hat{w}_{kl})(x_{im} - x_{il}) \right| = \\ & \frac{1}{2n} \sum_{i=1}^n |x_{im} - x_{il}| \cdot \sum_{k=1}^K |\hat{w}_{km} - \hat{w}_{kl}| = \\ & \frac{1}{2n} \|\mathbf{x}_{(m)} - \mathbf{x}_{(l)}\|_1 \cdot \|\hat{\mathbf{w}}_{(m)} - \hat{\mathbf{w}}_{(l)}\|_1 \end{aligned} \quad (9)$$

$$\begin{aligned}
& \sum_{k=1}^{K-1} \sum_{j=1}^p (|w_{kj}^*| - |\hat{w}_{kj}|) + \sum_{j=1}^p \left(\left| \sum_{k'=1}^{K-1} w_{k'j}^* \right| - \left| \sum_{k'=1}^{K-1} \hat{w}_{k'j} \right| \right) = \\
& \sum_{k=1}^{K-1} \left(2 \left| \frac{\hat{w}_{k'm} + \hat{w}_{k'l}}{2} \right| - |\hat{w}_{km}| - |\hat{w}_{kl}| \right) + \\
& 2 \left[\sum_{k'=1}^{K-1} \left| \frac{\hat{w}_{k'm} + \hat{w}_{k'l}}{2} \right| - \left| \sum_{k'=1}^{K-1} \hat{w}_{k'm} \right| - \left| \sum_{k'=1}^{K-1} \hat{w}_{k'l} \right| \right] \leq 0 \\
& \sum_{k=1}^{K-1} \sum_{j=1}^p (w_{kj}^{*2} - \hat{w}_{kj}^2) + \sum_{j=1}^p \left[\left(\sum_{k'=1}^{K-1} w_{k'j}^* \right)^2 - \left(\sum_{k'=1}^{K-1} \hat{w}_{k'j} \right)^2 \right] = \\
& - \frac{1}{2} \sum_{k=1}^{K-1} (\hat{w}_{km} - \hat{w}_{kl})^2 - \frac{1}{2} \left(\sum_{k=1}^{K-1} (\hat{w}_{km} - \hat{w}_{kl}) \right)^2 = \\
& - \frac{1}{2} \|\hat{\mathbf{w}}_{(m)} - \hat{\mathbf{w}}_{(l)}\|_2^2
\end{aligned} \tag{10}$$

From (8) ~ (11), we have

$$\begin{aligned}
0 & \leq \frac{1}{n\lambda_2} \|\mathbf{x}_{(m)} - \mathbf{x}_{(l)}\|_1 \cdot \|\hat{\mathbf{w}}_{(m)} - \hat{\mathbf{w}}_{(l)}\|_1 - \\
& \|\hat{\mathbf{w}}_{(m)} - \hat{\mathbf{w}}_{(l)}\|_2^2 \leq \frac{\sqrt{K}}{n\lambda_2} \|\mathbf{x}_{(m)} - \mathbf{x}_{(l)}\|_1 \times \\
& \|\hat{\mathbf{w}}_{(m)} - \hat{\mathbf{w}}_{(l)}\|_2 - \|\hat{\mathbf{w}}_{(m)} - \hat{\mathbf{w}}_{(l)}\|_2^2
\end{aligned} \tag{12}$$

which is equivalent to

$$\|\hat{\mathbf{w}}_{(m)} - \hat{\mathbf{w}}_{(l)}\|_2 \leq \frac{\sqrt{K}}{n\lambda_2} \sum_{i=1}^n |x_{im} - x_{il}| \tag{13}$$

If $\mathbf{x}_{(m)}$ and $\mathbf{x}_{(l)}$ are normalized, then we have

$$\begin{aligned}
\frac{\sqrt{K}}{n\lambda_2} \|\mathbf{x}_{(m)}^\top - \mathbf{x}_{(l)}^\top\|_1 & \leq \frac{\sqrt{K}}{n\lambda_2} \cdot \sqrt{n} \|\mathbf{x}_{(m)} - \mathbf{x}_{(l)}\|_2 = \\
\frac{\sqrt{K}}{\sqrt{n}\lambda_2} \sqrt{2 - 2\mathbf{x}_{(m)}^\top \mathbf{x}_{(l)}} & = \frac{\sqrt{K}}{\sqrt{n}\lambda_2} \sqrt{2(1 - \rho)}
\end{aligned} \tag{14}$$

Substituting (14) into (13) yields (7). \square

For the highly correlated predictors $\mathbf{x}_{(m)}$ and $\mathbf{x}_{(l)}$ ($\rho = 1$), HMSVM tends to assign the same coefficient vectors to them. This means that the highly correlated genes tend to be selected or removed together. According to the terminology of statistical learning^[15], we claim that HMSVM exhibits a grouping effect. For $K = 2$, it can be easily obtained that

$$|\hat{w}_{(1m)} - \hat{w}_{(1l)}| \leq \frac{1}{\sqrt{n}\lambda_2} \sqrt{2(1 - \rho)}$$

This coincides with the known results of binary SVMs^[4-5].

2.3 The piecewise linear regularization solution

Similar to the hybrid huberized SVM^[4], we continuously decrease λ_1 and use the superscript l to index the sets immediately after the l -th event has occurred. Let b_k^l , w_{kj}^l , λ_1^l , λ_2^l be the values of these parameters at the point of entry, m_k be the cardinality of set \mathcal{A}_k^l . For $\lambda_1^l \geq \lambda_1 > \lambda_1^{l+1}$, let

$$\lambda_2^l = \beta - \frac{\beta - \alpha}{\ln(e + \lambda_1^l)} \tag{15}$$

where $0 < \alpha < \beta$ are given constants.

Let $\mathbf{1}_{K-1}$ denote the $(K-1)$ -th-order matrix whose elements are all ones, and I_p denote the p -th-order identity

matrix. Let $\tilde{\mathbf{b}}_- = (\hat{b}_1 - \hat{b}_1^l, \hat{b}_2 - \hat{b}_2^l, \dots, \hat{b}_{K-1} - \hat{b}_{K-1}^l)^\top$, $\tilde{\mathbf{w}}_k = (\hat{w}_{kj_1} - \hat{w}_{kj_1}^l, \hat{w}_{kj_2} - \hat{w}_{kj_2}^l, \dots, \hat{w}_{kj_{m_k}} - \hat{w}_{kj_{m_k}}^l)^\top$, $\tilde{\mathbf{w}}_- = (\tilde{\mathbf{w}}_1^\top, \tilde{\mathbf{w}}_2^\top, \dots, \tilde{\mathbf{w}}_{K-1}^\top)^\top$, $A_{11} = \sum_{i \in \mathcal{E}_K^l} a_{iK} \mathbf{1}_{K-1} + \text{diag}\{\sum_{i \in \mathcal{E}_1^l} a_{i1}, \dots, \sum_{i \in \mathcal{E}_{K-1}^l} a_{iK-1}\}$, $A_{22_1} = (\tilde{A}_{kk'})_{(m_1+m_2+\dots+m_{K-1}) \times (m_1+m_2+\dots+m_{K-1})} + \text{diag}\{\tilde{A}_{11}, \tilde{A}_{22}, \dots, \tilde{A}_{K-1K-1}\}$,

$$\begin{aligned}
\bar{A}_{22_2} & = \begin{bmatrix} 2I_p & I_p & \cdots & I_p \\ I_p & 2I_p & \cdots & I_p \\ \vdots & \vdots & \ddots & \vdots \\ I_p & I_p & \cdots & 2I_p \end{bmatrix} \\
A_{12} & = \begin{bmatrix} \bar{X}_1^\top + \bar{X}_1^\top & \bar{X}_2^\top & \cdots & \bar{X}_{K-1}^\top \\ \bar{X}_1^\top & \bar{X}_2^\top + \bar{X}_2^\top & \cdots & \bar{X}_{K-1}^\top \\ \vdots & \vdots & \ddots & \vdots \\ \bar{X}_1^\top & \bar{X}_2^\top & \cdots & \bar{X}_{K-1}^\top + \bar{X}_{K-1}^\top \end{bmatrix}
\end{aligned}$$

where $\bar{X}_k^\top = (\sum_{i \in \mathcal{E}_k^l} a_{ik} x_{ij_{k1}}, \dots, \sum_{i \in \mathcal{E}_k^l} a_{ik} x_{ij_{km_k}})$, $\bar{X}_k^\top = (\sum_{i \in \mathcal{E}_K^l} a_{iK} x_{ij_{k1}}, \dots, \sum_{i \in \mathcal{E}_K^l} a_{iK} x_{ij_{km_k}})$, $\tilde{\mathbf{X}}^\top(j, k') = (\sum_{i \in \mathcal{E}_K^l} a_{iK} x_{ij} x_{ij_{k'1}}, \dots, \sum_{i \in \mathcal{E}_K^l} a_{iK} x_{ij} x_{ij_{k'm_k}})$,

$$\tilde{A}_{kk} = \begin{bmatrix} \sum_{i \in \mathcal{E}_k^l} a_{ik} x_{ij_{k1}} x_{ij_{k1}} & \cdots & \sum_{i \in \mathcal{E}_k^l} a_{ik} x_{ij_{k1}} x_{ij_{km_k}} \\ \sum_{i \in \mathcal{E}_k^l} a_{ik} x_{ij_{k2}} x_{ij_{k1}} & \cdots & \sum_{i \in \mathcal{E}_k^l} a_{ik} x_{ij_{k2}} x_{ij_{km_k}} \\ \vdots & \ddots & \vdots \\ \sum_{i \in \mathcal{E}_k^l} a_{ik} x_{ij_{km_k}} x_{ij_{k1}} & \cdots & \sum_{i \in \mathcal{E}_k^l} a_{ik} x_{ij_{km_k}} x_{ij_{km_k}} \end{bmatrix}$$

and $m_k \times m_k'$ submatrix of \tilde{A}

$$\tilde{A}_{kk'} = \begin{bmatrix} \tilde{\mathbf{X}}^\top(j = j_{k1}, k' = k') \\ \tilde{\mathbf{X}}^\top(j = j_{k2}, k' = k') \\ \vdots \\ \tilde{\mathbf{X}}^\top(j = j_{km_k}, k' = k') \end{bmatrix}$$

for $k = 1, \dots, K-1$, $k' = 1, \dots, K-1$. Let $A_{22} = \frac{1}{n\delta} A_{22_1} + 2\lambda_2 A_{22_2}$, where A_{22_2} is constructed by orderly selecting the j_{k1} -th, j_{k2} -th, \dots , j_{km_k} -th rows and columns of the k -th p rows and columns in A_{22_2} .

Theorem 2. If the regularization parameters satisfy (15), then the coefficient paths of HMSVM are piecewise linear with respect to the regularization parameter λ_1 , i.e.,

$$\begin{cases} \hat{b}_k = \hat{b}_k^l + (\lambda_1 - \lambda_1^l) \bar{a}_{k0} \\ \hat{w}_{kj} = \hat{w}_{kj}^l + (\lambda_1 - \lambda_1^l) \bar{a}_{kj}, \text{ for } j \in \mathcal{A}_k^l \cap \bar{\mathcal{A}}^l \end{cases} \tag{16}$$

hold for $\lambda_1^l \geq \lambda_1 > \lambda_1^{l+1}$, $k = 1, \dots, K-1$, and $j = 1, \dots, p$, where \bar{a}_{k0} and \bar{a}_{1j} are respectively the k -th and $(K+j)$ -th elements of vector $A_l^{-1} \mathbf{1}^\alpha$, \bar{a}_{kj} is the $(K-1+m_1+\dots+m_{k-1}+k+j)$ -th elements of vector $A_l^{-1} \mathbf{1}^\alpha$ for $k \geq 2$, and A_l and $\mathbf{1}^\alpha$ are defined as

$$\begin{aligned}
A_l & = \begin{bmatrix} A_{11} & A_{12} \\ A_{12}^\top & A_{22} \end{bmatrix} \\
\mathbf{1}^\alpha & = \begin{bmatrix} 0_{K-1} \\ \text{sgn}((A_{22_2} - I) \tilde{\mathbf{w}}_-) + \text{sgn}(\hat{\mathbf{w}}_-^l) \end{bmatrix}
\end{aligned}$$

Proof. Note that (5) is an unconstrained convex optimization problem. Hence, for $b_k \neq 0, w_{kj} \neq 0, \sum_{k=1}^{K-1} b_k \neq 0$, and $\sum_{k=1}^{K-1} w_{kj} \neq 0$, we have

$$\begin{cases} \left. \frac{\partial \bar{L}(\lambda_1, \lambda_2, \mathbf{b}_-, w_-)}{\partial b_k} \right|_{\mathbf{b}_- = \hat{\mathbf{b}}_-, w_- = \hat{w}_-} = 0 \\ \left. \frac{\partial \bar{L}(\lambda_1, \lambda_2, \mathbf{b}_-, w_-)}{\partial w_{kj}} \right|_{\mathbf{b}_- = \hat{\mathbf{b}}_-, w_- = \hat{w}_-} = 0 \end{cases} \quad (17)$$

where $k = 1, \dots, K-1, j = 1, \dots, p$. Note that sets $\mathcal{E}_k^l, \mathcal{R}_k^l, \mathcal{L}_k^l, \mathcal{A}_k^l, \bar{\mathcal{A}}$ and the regularization parameter λ_2 will not change for $\lambda_1^l \geq \lambda_1 > \lambda_1^{l+1}$. Hence, we have

$$\begin{aligned} & \sum_{i \in \mathcal{E}_k^l} a_{ik} \frac{1 + \hat{b}_k + \hat{\mathbf{w}}_k^T \mathbf{x}_i}{n\delta} + \\ & \sum_{i \in \mathcal{E}_K^l} \frac{a_{iK}}{n\delta} \left(-1 + \sum_{k'=1}^{K-1} (\hat{b}_{k'} + \hat{\mathbf{w}}_{k'}^T \mathbf{x}_i) \right) + \\ & \frac{1}{n} \left(\sum_{i \in \mathcal{R}_k^l} a_{ik} - \sum_{i \in \mathcal{R}_K^l} a_{iK} \right) = 0 \end{aligned} \quad (18)$$

$$\begin{aligned} & \sum_{i \in \mathcal{E}_k^l} a_{ik} x_{ij} \frac{1 + \hat{b}_k + \hat{\mathbf{w}}_k^T \mathbf{x}_i}{n\delta} + 2\lambda_2^l \left(\hat{w}_{kj} + \sum_{k'=1}^{K-1} \hat{w}_{k'j} \right) + \\ & \sum_{i \in \mathcal{E}_K^l} \frac{a_{iK} x_{ij}}{n\delta} \left(-1 + \sum_{k'=1}^{K-1} (\hat{b}_{k'} + \hat{\mathbf{w}}_{k'}^T \mathbf{x}_i) \right) + \\ & \frac{1}{n} \left(\sum_{i \in \mathcal{R}_k^l} a_{ik} x_{ij} - \sum_{i \in \mathcal{R}_K^l} a_{iK} x_{ij} \right) + \\ & \lambda_1 \left(\text{sgn}(\hat{w}_{kj}) + \text{sgn} \left(\sum_{k'=1}^{K-1} \hat{w}_{k'j} \right) \right) = 0 \end{aligned} \quad (19)$$

where $k = 1, \dots, K-1, j \in \mathcal{A}_k^l \cap \bar{\mathcal{A}}^l$. Analogously, for $\lambda_1 = \lambda_1^l$, we also have

$$\begin{aligned} & \sum_{i \in \mathcal{E}_k^l} a_{ik} \frac{1 + \hat{b}_k^l + \hat{\mathbf{w}}_k^{lT} \mathbf{x}_i}{n\delta} + \frac{1}{n} \left(\sum_{i \in \mathcal{R}_k^l} a_{ik} - \sum_{i \in \mathcal{R}_K^l} a_{iK} \right) + \\ & \sum_{i \in \mathcal{E}_K^l} \frac{a_{iK}}{n\delta} \left(-1 + \sum_{k'=1}^{K-1} (\hat{b}_{k'}^l + \hat{\mathbf{w}}_{k'}^{lT} \mathbf{x}_i) \right) = 0 \\ & \sum_{i \in \mathcal{E}_k^l} a_{ik} x_{ij} \frac{1 + \hat{b}_k^l + \hat{\mathbf{w}}_k^{lT} \mathbf{x}_i}{n\delta} + 2\lambda_2^l \left(\hat{w}_{kj}^l + \sum_{k'=1}^{K-1} \hat{w}_{k'j}^l \right) + \\ & \sum_{i \in \mathcal{E}_K^l} \frac{a_{iK} x_{ij}}{n\delta} \left(-1 + \sum_{k'=1}^{K-1} (\hat{b}_{k'}^l + \hat{\mathbf{w}}_{k'}^{lT} \mathbf{x}_i) \right) + \\ & \frac{1}{n} \left(\sum_{i \in \mathcal{R}_k^l} a_{ik} x_{ij} - \sum_{i \in \mathcal{R}_K^l} a_{iK} x_{ij} \right) + \\ & \lambda_1^l \left(\text{sgn}(\hat{w}_{kj}^l) + \text{sgn} \left(\sum_{k'=1}^{K-1} \hat{w}_{k'j}^l \right) \right) = 0 \end{aligned} \quad (21)$$

Subtracting (20) from (18) gives

$$\begin{aligned} & \sum_{i \in \mathcal{E}_k^l} a_{ik} (\hat{b}_k - \hat{b}_k^l + (\hat{\mathbf{w}}_k - \hat{\mathbf{w}}_k^l)^T \mathbf{x}_i) + \\ & \sum_{i \in \mathcal{E}_K^l} a_{iK} \sum_{k'=1}^{K-1} [\hat{b}_{k'} - \hat{b}_{k'}^l + (\hat{\mathbf{w}}_{k'} - \hat{\mathbf{w}}_{k'}^l)^T \mathbf{x}_i] = 0 \end{aligned} \quad (22)$$

Note that $\text{sgn}(\hat{w}_{kj}) = \text{sgn}(\hat{w}_{kj}^l)$ and $\text{sgn}(\sum_{k'=1}^{K-1} \hat{w}_{k'j}) = \text{sgn}(\sum_{k'=1}^{K-1} \hat{w}_{k'j}^l)$. Subtracting (21) from (19) gives

$$\begin{aligned} & \sum_{i \in \mathcal{E}_k^l} a_{ik} x_{ij} \frac{\hat{b}_k - \hat{b}_k^l + (\hat{\mathbf{w}}_k - \hat{\mathbf{w}}_k^l)^T \mathbf{x}_i}{n\delta} + \\ & \sum_{i \in \mathcal{E}_K^l} a_{iK} x_{ij} \sum_{k'=1}^{K-1} \frac{\hat{b}_{k'} - \hat{b}_{k'}^l + (\hat{\mathbf{w}}_{k'} - \hat{\mathbf{w}}_{k'}^l)^T \mathbf{x}_i}{n\delta} + \\ & 2\lambda_2^l (\hat{w}_{kj} - \hat{w}_{kj}^l + \sum_{k'=1}^{K-1} (\hat{w}_{k'j} - \hat{w}_{k'j}^l)) + \\ & (\lambda_1 - \lambda_1^l) \left(\text{sgn}(\hat{w}_{kj}^l) + \text{sgn} \left(\sum_{k'=1}^{K-1} \hat{w}_{k'j}^l \right) \right) = 0 \end{aligned} \quad (23)$$

Note that

$$\begin{aligned} & \sum_{i \in \mathcal{E}_k^l} a_{ik} (\hat{b}_k - \hat{b}_k^l) + \sum_{i \in \mathcal{E}_K^l} a_{iK} \sum_{k'=1}^{K-1} (\hat{b}_{k'} - \hat{b}_{k'}^l) = \\ & \sum_{i \in \mathcal{E}_k^l} a_{ik} (\hat{b}_k - \hat{b}_k^l) + \sum_{k'=1}^{K-1} \sum_{i \in \mathcal{E}_K^l} a_{iK} (\hat{b}_{k'} - \hat{b}_{k'}^l) \\ & \sum_{i \in \mathcal{E}_k^l} a_{ik} (\hat{\mathbf{w}}_k - \hat{\mathbf{w}}_k^l)^T \mathbf{x}_i + \sum_{i \in \mathcal{E}_K^l} a_{iK} \sum_{k'=1}^{K-1} (\hat{\mathbf{w}}_{k'} - \hat{\mathbf{w}}_{k'}^l)^T \mathbf{x}_i = \\ & \sum_{k'=1}^{K-1} \sum_{j' \in \mathcal{A}_{k'}^l} \sum_{i \in \mathcal{E}_K^l} a_{iK} x_{ij'} (\hat{w}_{k'j'} - \hat{w}_{k'j'}^l) + \\ & \sum_{j \in \mathcal{A}_k^l} \sum_{i \in \mathcal{E}_k^l} a_{ik} x_{ij} (\hat{w}_{kj} - \hat{w}_{kj}^l) \end{aligned}$$

Hence, system of linear equations in (22) can be rewritten as

$$A_{11} \tilde{\mathbf{b}}_- + A_{12} \tilde{\mathbf{w}}_- = 0 \quad (24)$$

Note also that

$$\begin{aligned} & \sum_{i \in \mathcal{E}_k^l} a_{ik} x_{ij} (\hat{b}_k - \hat{b}_k^l) + \sum_{i \in \mathcal{E}_K^l} a_{iK} x_{ij} \sum_{k'=1}^{K-1} (\hat{b}_{k'} - \hat{b}_{k'}^l) = \\ & \sum_{i \in \mathcal{E}_k^l} a_{ik} x_{ij} (\hat{b}_k - \hat{b}_k^l) + \sum_{k'=1}^{K-1} \sum_{i \in \mathcal{E}_K^l} a_{iK} x_{ij} (\hat{b}_{k'} - \hat{b}_{k'}^l) \\ & \sum_{i \in \mathcal{E}_k^l} a_{ik} x_{ij} (\hat{\mathbf{w}}_k - \hat{\mathbf{w}}_k^l)^T \mathbf{x}_i + \sum_{i \in \mathcal{E}_K^l} a_{iK} x_{ij} \sum_{k'=1}^{K-1} (\hat{\mathbf{w}}_{k'} - \hat{\mathbf{w}}_{k'}^l)^T \mathbf{x}_i = \\ & \sum_{j'' \in \mathcal{A}_k^l} \sum_{i \in \mathcal{E}_k^l} a_{ik} x_{ij} x_{ij''} (\hat{w}_{kj''} - \hat{w}_{kj''}^l) + \\ & \sum_{k'=1}^{K-1} \sum_{j' \in \mathcal{A}_{k'}^l} \sum_{i \in \mathcal{E}_K^l} a_{iK} x_{ij} x_{ij'} (\hat{w}_{k'j'} - \hat{w}_{k'j'}^l) \end{aligned}$$

$$A_{22_2} \tilde{\mathbf{w}}_- = \begin{bmatrix} \hat{w}_{1j_{11}} - \hat{w}_{1j_{11}}^l + \sum_{k'=1}^{K-1} (\hat{w}_{k'j_{11}} - \hat{w}_{k'j_{11}}^l) \\ \vdots \\ \hat{w}_{2j_{21}} - \hat{w}_{2j_{21}}^l + \sum_{k'=1}^{K-1} (\hat{w}_{k'j_{21}} - \hat{w}_{k'j_{21}}^l) \\ \vdots \\ A_{22_\alpha} \end{bmatrix}$$

where $A_{22_\alpha} = \sum_{k'=1}^{K-1} (\hat{w}_{k'j_{K-1m_{K-1}}} - \hat{w}_{k'j_{K-1m_{K-1}}}^l) + \hat{w}_{K-1j_{K-1m_{K-1}}} - \hat{w}_{K-1j_{K-1m_{K-1}}}^l$. Hence, system of linear equations in (23) can be rewritten as

$$\frac{1}{n\delta} A_{12}^T \tilde{\mathbf{b}}_- + 2\lambda_2 A_{22_2} \tilde{\mathbf{w}}_- = (\lambda_1 - \lambda_1^l) (\text{sgn}((A_{22_2} - I) \hat{\mathbf{w}}_-^l) + \text{sgn}(\hat{\mathbf{w}}_-^l)) \quad (25)$$

If A_l has full rank, then (16) can be obtained by solving the combined systems of linear equations (24) and (25). \square

It should be noted that if A_l does not have full rank, then the solution paths are not unique, and more care has to be taken^[16]. According to the sum-to-zero constraint, it can be easily obtained that $\hat{b}_K = -\sum_{k=1}^{K-1} \hat{b}_k - (\lambda_1 - \lambda_1^l) \sum_{k=1}^{K-1} \bar{a}_{k0}$ and $\hat{w}_{Kj} = -\sum_{k=1}^{K-1} \hat{w}_{kj} - (\lambda_1 - \lambda_1^l) \sum_{k=1}^{K-1} \bar{a}_{kj}$. Substituting (16) into $f_k(\mathbf{x}) = b_k + \mathbf{w}_k^T \mathbf{x}$ gives $f_k(\mathbf{x}_i) = f_k^l(\mathbf{x}_i) + (\lambda_1 - \lambda_1^l) (\bar{a}_{k0} + \sum_{j \in \mathcal{A}_k^l \cap \bar{\mathcal{A}}^l} x_{ij} \bar{a}_{kj})$ and $f_K(\mathbf{x}_i) = -\sum_{k=1}^{K-1} f_k^l(\mathbf{x}_i) - (\lambda_1 - \lambda_1^l) \sum_{k=1}^{K-1} (\bar{a}_{k0} + \sum_{j \in \mathcal{A}_k^l \cap \bar{\mathcal{A}}^l} x_{ij} \bar{a}_{kj})$, where $k = 1, \dots, K-1$.

2.4 Algorithm

Similar to the binary solution path algorithms^[4-6, 16], our algorithm starts from $\lambda_1 \rightarrow \infty$ and we let $\hat{w}^0 = 0_{K \times p}$, $\hat{\mathbf{b}}_-^0 = \arg \min_{\mathbf{b}_-} \frac{1}{n} (\sum_{i=1}^n \sum_{k=1}^{K-1} a_{ik} L_k(b_k, 0, 0) + \sum_{i=1}^n a_{iK} L_K(-\sum_{k=1}^{K-1} b_{k'}, 0, 0))$, $\hat{b}_K^0 = -\sum_{k=1}^{K-1} \hat{b}_k^0$. According to the obtained \hat{b}_k^0 , we determine the sets \mathcal{E}_k^0 , \mathcal{L}_k^0 , and \mathcal{R}_k^0 for $k = 1, \dots, K$. Furthermore, we let $\lambda_1^0 = \max_{j=1}^p \max_{k=1}^{K-1} (|\sum_{i \in \mathcal{E}_k^0} \frac{a_{ik} x_{ij}}{n\delta} (\hat{b}_k^0 + \sum_{k'=1}^{K-1} \hat{b}_{k'}^0)|)$, $\mathcal{A}_k^0 = \arg \max_{j \in \{1, \dots, p\}} (|\sum_{i \in \mathcal{E}_k^0} \frac{a_{ik} x_{ij}}{n\delta} (\hat{b}_k^0 + \sum_{k'=1}^{K-1} \hat{b}_{k'}^0)|)$, $\lambda_2^0 = \beta - (\beta - \alpha) / (\ln(e + \lambda_1^0))$, and determine $\bar{\mathcal{A}}^0$ according to \mathcal{A}_k^0 . This completes the initialization.

After the l -th event has occurred, the most important problem of our algorithm is to determine the step size for the event which will occur first. Note that the first event will occur when $f_k(\mathbf{x}_i)$ reaches -1 . Hence, the step size for the first event can be determined by

$$d_1 = \min \{d_{11}, d_{12}\} \quad (26)$$

where $d_{11} = \min_{k=1, \dots, k-1} \min_{i \in \mathcal{E}_k^l \cup \mathcal{L}_k^l} |(-1 - f_k^l(\mathbf{x}_i)) / (\bar{a}_{k0} + \sum_{j \in \mathcal{A}_k^l \cap \bar{\mathcal{A}}^l} x_{ij} \bar{a}_{kj})|$, $d_{12} = \min_{i \in \mathcal{E}_K^l \cup \mathcal{L}_K^l} |(-1 + \sum_{k=1}^{K-1} f_k^l(\mathbf{x}_i)) / (\sum_{k=1}^{K-1} (\bar{a}_{k0} + \sum_{j \in \mathcal{A}_k^l \cap \bar{\mathcal{A}}^l} x_{ij} \bar{a}_{kj}))|$. Let $d_{21} = \min_{k=1, \dots, k-1} \min_{i \in \mathcal{E}_k^l \cup \mathcal{R}_k^l} |(-1 + \delta - f_k^l(\mathbf{x}_i)) / (\bar{a}_{k0} + \sum_{j \in \mathcal{A}_k^l \cap \bar{\mathcal{A}}^l} x_{ij} \bar{a}_{kj})|$, $d_{22} = \min_{i \in \mathcal{E}_K^l \cup \mathcal{R}_K^l} |(-1 + \delta + \sum_{k=1}^{K-1} f_k^l(\mathbf{x}_i)) / (\sum_{k=1}^{K-1} (\bar{a}_{k0} + \sum_{j \in \mathcal{A}_k^l \cap \bar{\mathcal{A}}^l} x_{ij} \bar{a}_{kj}))|$. Analogously, the step size for the second event can be determined by

$$d_2 = \min \{d_{21}, d_{22}\} \quad (27)$$

Note that the third event will occur when a non-zero parameter \hat{w}_{kj} reduces to zero or the sum of the first $K-1$

parameters in $\hat{\mathbf{w}}_{(j)}$ becomes zero. Hence, the step size for the third event can be determined by

$$d_3 = \min_{j \in \mathcal{A}_k^l \cap \bar{\mathcal{A}}^l} \min_{k \in \{1, \dots, K\}} \left| \frac{\hat{w}_{kj}^l}{\bar{a}_{kj}} \right| \quad (28)$$

where $\hat{w}_{Kj}^l = \sum_{k'=1}^{K-1} \hat{w}_{k'j}^l$ and $\bar{a}_{Kj} = \sum_{k'=1}^{K-1} \bar{a}_{k'j}$.

For $k = 1, \dots, K-1$, $j = 1, \dots, p$, we define $C_{kj} = \sum_{i \in \mathcal{E}_k^l} a_{iK} x_{ij} (-1 + \sum_{k'=1}^{K-1} (\hat{b}_{k'} + \hat{\mathbf{w}}_k^T \mathbf{x}_i)) / (n\delta) + \sum_{i \in \mathcal{E}_k^l} a_{iK} x_{ij} (1 + \hat{b}_k + \hat{\mathbf{w}}_k^T \mathbf{x}_i) / (n\delta) + (\sum_{i \in \mathcal{R}_k^l} a_{iK} x_{ij} - \sum_{i \in \mathcal{R}_k^l} a_{iK} x_{ij}) / n + 2\lambda_2^l (\hat{w}_{kj} + \sum_{k'=1}^{K-1} \hat{w}_{k'j})$. Note that $|C_{kj}| \leq 2\lambda_1$, for $j \notin \mathcal{A}_k^l \cap \bar{\mathcal{A}}^l$; $C_{kj} = 0$ for $j \in \mathcal{A}_k^l \cap \bar{\mathcal{A}}^l$ and $\text{sgn}(\hat{w}_{kj}) = -\text{sgn}(\sum_{k'=1}^{K-1} \hat{w}_{k'j})$; $C_{kj} = -2\lambda_1 \text{sgn}(\hat{w}_{kj})$, for $j \in \mathcal{A}_k^l \cap \bar{\mathcal{A}}^l$ and $\text{sgn}(\hat{w}_{kj}) = \text{sgn}(\sum_{k'=1}^{K-1} \hat{w}_{k'j})$. Hence, after $|C_{kj}|$ ($j \notin \mathcal{A}_k^l \cap \bar{\mathcal{A}}^l$) meets the decreasing $2\lambda_1$ or reaches zero, the fourth event will occur if we keep moving λ_1 in the same direction. For $j \notin \mathcal{A}_k^l \cap \bar{\mathcal{A}}^l$, there are three cases: 1) $j \notin \mathcal{A}_k^l$ and $j \notin \bar{\mathcal{A}}^l$; 2) $j \in \mathcal{A}_k^l$ and $j \notin \bar{\mathcal{A}}^l$; 3) $j \notin \mathcal{A}_k^l$ and $j \in \bar{\mathcal{A}}^l$. Note that $\hat{w}_{kj} = 0$ and $\sum_{k'=1}^{K-1} \hat{w}_{k'j} = 0$ for $j \notin \mathcal{A}_k^l$ and $j \notin \bar{\mathcal{A}}^l$. Hence, for case 1), we have $C_{kj} = C_{kj}^l + C_{\mathcal{E}kj} (\lambda_1 - \lambda_1^l)$, where $C_{\mathcal{E}kj} = \sum_{i \in \mathcal{E}_k^l} \frac{a_{iK} x_{ij}}{n\delta} (\bar{a}_{k0} + \sum_{j'' \in \mathcal{A}_k^l} x_{ij''} \bar{a}_{kj''}) + \sum_{i \in \mathcal{E}_K^l} \frac{a_{iK} x_{ij}}{n\delta} \sum_{k'=1}^{K-1} (\bar{a}_{k'0} + \sum_{j' \in \mathcal{A}_{k'}^l} x_{ij'} \bar{a}_{k'j'})$. The step size d_{411} , which makes $|C_{kj}|$ ($j \notin \mathcal{A}_k^l \cap \bar{\mathcal{A}}^l$) reach zero, can be calculated by

$$d_{411} = \left| \frac{C_{kj}^l}{C_{\mathcal{E}kj}} \right| \quad (29)$$

The step size d_{412} , which makes $|C_{kj}|$ ($j \notin \mathcal{A}_k^l \cap \bar{\mathcal{A}}^l$) meet the $2\lambda_1$, can be calculated by

$$|C_{kj}^l - C_{\mathcal{E}kj} d_{412}| = |C_{kj}| = 2\lambda_1 = 2(\lambda_1^l - d_{412}) \quad (30)$$

If $\text{sgn}(C_{kj}^l) = \text{sgn}(C_{\mathcal{E}kj})$ and $d_{412} \leq |C_{kj}^l| / |C_{\mathcal{E}kj}|$, then

$$d_{412} = \frac{2\lambda_1^l - |C_{kj}^l|}{2 - |C_{\mathcal{E}kj}|}$$

If $\text{sgn}(C_{kj}^l) = \text{sgn}(C_{\mathcal{E}kj})$ and $d_{412} > |C_{kj}^l| / |C_{\mathcal{E}kj}|$, then

$$d_{412} = \frac{2\lambda_1^l + |C_{kj}^l|}{2 + |C_{\mathcal{E}kj}|}$$

Otherwise, we have

$$d_{412} = \frac{2\lambda_1^l - |C_{kj}^l|}{2 + |C_{\mathcal{E}kj}|}$$

It should be noted that the inequality $d_{412} \leq |C_{kj}^l| / |C_{\mathcal{E}kj}|$ (or $d_{412} > |C_{kj}^l| / |C_{\mathcal{E}kj}|$) is equivalent to the computable inequality $\lambda_1^l |C_{\mathcal{E}kj}| \leq |C_{kj}^l|$ (or $\lambda_1^l |C_{\mathcal{E}kj}| > |C_{kj}^l|$). Hence, the step size d_{41} , which determines the fourth event for case 1), can be determined by

$$d_{41} = \min_{j \in \mathcal{A}_k^l, j \notin \bar{\mathcal{A}}^l} \min_{k=1, \dots, K-1} \min \{d_{411}, d_{412}\} \quad (31)$$

Analogously, we can get the step sizes d_{42} and d_{43} , which determine the fourth event for case 2) and 3), respectively. It should be noted that the calculation procedure is the same as the case 1) except the different value of $C_{\mathcal{E}kj}$. Note that $\hat{w}_{kj} \neq 0$ and $\sum_{k'=1}^{K-1} \hat{w}_{k'j} = 0$ for

$j \in \mathcal{A}_k^l$ and $j \notin \bar{\mathcal{A}}^l$. Hence, for case 2), we have $C_{\mathcal{E}kj} = 2\lambda_2^l \bar{a}_{kj} + \sum_{i \in \mathcal{E}_k^l} \frac{a_{ik} x_{ij}}{n\delta} (\bar{a}_{k0} + \sum_{j'' \in \mathcal{A}_k^l} x_{ij''} \bar{a}_{kj''}) + \sum_{i \in \mathcal{E}_K^l} \frac{a_{iK} x_{ij}}{n\delta} \sum_{k'=1}^{K-1} (\bar{a}_{k'0} + \sum_{j'' \in \mathcal{A}_{k'}^l} x_{ij''} \bar{a}_{k'j''})$. Note also that $\hat{w}_{kj} = 0$ and $\sum_{k'=1}^{K-1} \hat{w}_{k'j} \neq 0$ for $j \notin \mathcal{A}_k^l$ and $j \in \bar{\mathcal{A}}^l$. Hence, for case 3), we have $C_{\mathcal{E}kj} = 2\lambda_2^l \sum_{k'=1}^{K-1} \bar{a}_{k'j} + \sum_{i \in \mathcal{E}_k^l} \frac{a_{ik} x_{ij}}{n\delta} (\bar{a}_{k0} + \sum_{j'' \in \mathcal{A}_k^l} x_{ij''} \bar{a}_{kj''}) + \sum_{i \in \mathcal{E}_K^l} \frac{a_{iK} x_{ij}}{n\delta} \sum_{k'=1}^{K-1} (\bar{a}_{k'0} + \sum_{j'' \in \mathcal{A}_{k'}^l} x_{ij''} \bar{a}_{k'j''})$. To sum up, the step size for the fourth event can be determined by

$$d_4 = \min\{d_{41}, d_{42}, d_{43}\} \quad (32)$$

The algorithm that computes the whole solution path $\hat{\mathbf{b}}, \hat{w}$ proceeds as follows:

Step 1. Initialization: calculate $\hat{w}^0, \hat{\mathbf{b}}^l, \lambda_1^0, \lambda_2^0, \mathcal{E}_k^0, \mathcal{L}_k^0, \mathcal{R}_k^0, \mathcal{A}_k^0$, and $\bar{\mathcal{A}}^0$, where $k = 1, \dots, K$.

Step 2. Find λ_1^{l+1} and λ_2^{l+1} : let $\lambda_2^l = \beta - (\beta - \alpha)/(\ln(e + \lambda_1^l))$, calculate $\hat{\mathbf{b}}, \hat{w}$ and $f_k(\mathbf{x}_i)$, and determine the step size $d = \min\{d_1, d_2, d_3, d_4\}$.

Step 3. If the generalized correlation reduces to zero or a pre-specified maximum iteration number is reached, then stop the algorithm.

Step 4. Otherwise, let $l = l + 1, \lambda_1^{l+1} = \lambda_1^l - d, \lambda_2^{l+1} = \beta - (\beta - \alpha)/(\ln(e + \lambda_1^{l+1}))$ and update $\hat{\mathbf{b}}, \hat{w}^l, \mathcal{E}_k^l, \mathcal{L}_k^l, \mathcal{R}_k^l, \mathcal{A}_k^l$ and $\bar{\mathcal{A}}^l$. Then goto the Step 2.

3 Experiments on leukemia data

The aim of the leukemia benchmark^[2] is to form a decision rule capable of distinguishing between acute myeloid leukemia (AML) and acute lymphoblastic leukemia (ALL). This data set contains 38 samples for training and 34 samples for testing. The samples were assayed using Affymetrix Hgu6800 chips and data on the expression of 7129 genes (Affymetrix probes) are available. The original data set is available at <http://www.broad.mit.edu/cgi-bin/cancer/datasets.cgi>. The identification of the two cancer types was based on their origins, lymphoid (lymph or lymphatic tissue related), and myeloid (bone marrow related), respectively. Similar to MSVM^[8], we modify the leukemia data set as a three-class problem by dividing ALL into B-cell ALL and T-cell ALL. This study examines 72 samples from three types of acute leukemia with 38 samples in B-cell ALL, 9 samples in T-cell ALL, and 25 samples in AML.

In the first experiment, we train and test HMSVM using the original separation for the training and testing data. In the original training data set, there are 19 samples in B-cell ALL, 8 samples in T-cell ALL, and 11 samples in AML. The remaining 34 samples are used to test the prediction accuracy. To make the computation more manageable, we use the pre-processing steps proposed by Dudoit^[1] and select the most significant 3571 genes as the predictors. Let $\beta = 10, \alpha = 0.5, \delta = 0.05$. We compute the entire regularization solution paths according to the algorithm in Subsection 2.4. Fig. 1 shows the curve of correlation between the two regularization parameters. It is shown that λ_2 is the piecewise constant function with respect to λ_1 . Fig. 2 shows the curve of prediction error. It is shown that the prediction error is piecewise linear with respect to λ_1 and the optimal mode is given when the regularization parameter λ_1 is selected in the interval $[2.280905, 2.951294]$. The corresponding prediction error is 0.02941176, i.e., only one sample in the test data set is misclassified.

In the second experiment, we compare HMSVM with

several competitors for multiclass classification, including one-versus-all (OVA) classifier, L_2 -norm MSVM and L_1 -norm MSVM. We test their averaged prediction accuracy using a randomly splitting approach: the original training and testing data are combined together and we randomly split it into 38 and 34 samples for training and testing, respectively. The entire process is repeated 50 times and Table 1 gives the averaged test error and the averaged number of selected genes. In L_2 -norm MSVM, 100 most important genes are selected to build the classifier by adopting the marginal criterion^[1]. Although L_2 -norm MSVM still works well by setting appropriate Gaussian kernel parameter σ and regularization parameter λ_2 , it does not select genes in a satisfactory way. L_1 -norm MSVM seems to overcome this disadvantage and improve the prediction accuracy. However, the number of the selected genes is upper bounded by the sample size and the selected genes vary largely with the randomly splitting training set. The HMSVM gives the best averaged test error and selects the moderate amount of genes. It is likely that the grouping effect in the process of automatic gene selection contributes much to the improved prediction accuracy.

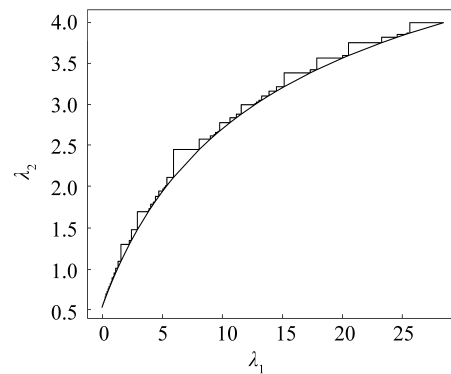


Fig. 1 Correlation between λ_2 and λ_1

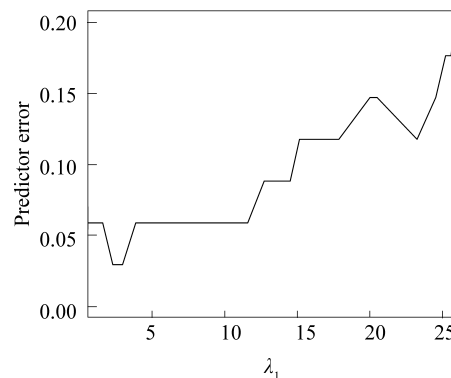


Fig. 2 Curve of the prediction error

4 Conclusions

The huberized multiclass SVM has been proposed in this paper, and the corresponding properties have been studied. It is shown that the HMSVM can encourage a grouping effect for multiclass gene selection. Based on a reasonable correlation of the two regularization parameters, the optimal coefficients are shown to be piecewise linear with re-

spect to the single regularization parameter and an efficient regularized path algorithm is developed. We have applied HMSVM to the leukemia data and achieved promising results.

Table 1 Classification results for the leukemia data

Method	Averaged test error (%)	Averaged number of selected genes
OVA	6.24	26.73
2-norm MSVM	4.20	100
1-norm MSVM	3.76	20.92
HMSVM	2.94	62.37

References

- Dudoit S, Fridlyand J, Speed T P. Comparison of discrimination methods for the classification of tumors using gene expression data. *Journal of the American Statistical Association*, 2002, **97**(457): 77–87
- Golub T R, Slonim D K, Tamayo P, Huard C, Gaasenbeek M, Mesirov J P. Molecular classification of cancer: class discovery and class prediction by gene expression monitoring. *Science*, 1999, **286**(5439): 531–537
- Guyon I, Weston J, Barnhill S, Vapnik V. Gene selection for cancer classification using support vector machines. *Machine Learning*, 2002, **46**(1-3): 389–422
- Wang L, Zhu J, Zou H. Hybrid huberized support vector machines for microarray classification and gene selection. *Bioinformatics*, 2008, **24**(3): 412–419
- Wang L, Zhu J, Zou H. The doubly regularized support vector machine. *Statistica Sinica*, 2006, **16**(2): 589–615
- Zhu J, Rosset S, Hastie T, Tibshirani R. 1-norm support vector machines. *Advances in Neural Information Processing Systems*. New York: MIT Press, 2004. 49–56
- Lee Y, Cui Z. Characterizing the solution path of multicategory support vector machines. *Statistica Sinica*, 2006, **16**(2): 391–409
- Lee Y, Lin Y, Wahba G. Multicategory support vector machines: theory, and application to the classification of microarray data and satellite radiance data. *Journal of the American Statistical Association*, 2004, **99**(465): 67–81
- Wang L F, Shen X T. On l_1 -norm multi-class support vector machines: methodology and theory. *Journal of the American Statistical Association*, 2007, **102**(478): 583–594
- Zhang H H, Liu Y, Wu Y, Zhu J. Variable selection for the multicategory SVM via adaptive sup-norm regularization. *Electronic Journal of Statistics*, 2008, **2**(1): 149–167
- Crammer K, Singer Y. On the algorithmic implementation of multiclass kernel-based vector machines. *Journal of Machine Learning Research*, 2001, **2**(3): 265–292
- Li J, Jia Y, Du J, Yu F. Gene selection of multiple cancer types via huberized multi-class support vector machine. In: Proceedings of the 48th IEEE Conference on Decision and Control and the 28th Chinese Control Conference. Shanghai, China: IEEE, 2009. 1520–1525
- Ma S, Song X, Huang J. Supervised group lasso with applications to microarray data analysis. *BMC Bioinformatics* [Online], available: <http://www.biomedcentral.com/1471-2105/8/60>, March 15, 2009
- Yuan M, Lin Y. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society, Series B*, 2006, **68**(1): 49–67
- Zou H, Hastie T. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society, Series B*, 2005, **67**(2): 301–320
- Hastie T, Rosset S, Tibshirani R, Zhu J. The entire regularization path for the support vector machine. *Journal of Machine Learning Research*, 2004, **5**(12): 1391–1415
- Rosset S, Zhu J. Piecewise linear regularized solution paths. *Annals of Statistics*, 2007, **35**(3): 1012–1030



LI Jun-Tao Ph.D. candidate at the Seven Research Division, Beihang University. His research interest covers intelligent control, statistical learning, data mining, and machine-learning-based bioinformatics. Corresponding author of this paper. E-mail: juntaol@mail@yahoo.com.cn



JIA Ying-Min Professor at Beihang University. His research interest covers multivariable systems, robust control, adaptive control, and intelligent control and their applications in vehicle systems and industrial processes. E-mail: ymjia@buaa.edu.cn