

Multiclass Multiple Kernel Learning

Alexander Zien and Cheng Soon Ong

June. 12th, 2007

Existing works focus only on binary classification problem. This work propose an intuitive formulation of the multi-class SVM.

mSVM

$$\text{function : } f_{\mathbf{w},b} = \langle \mathbf{w}, \phi(\mathbf{x}, y) \rangle + b_y$$

$$\text{prediction : } x = \arg \max_{y \in \mathcal{Y}} f_{\mathbf{w},b}(\mathbf{x}, y)$$

$$\text{objective : } \min_{\mathbf{w},b} \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^n \max_{u \neq y_i} \{ \ell(f_{\mathbf{w},b}(\mathbf{x}_i, y_i) - f_{\mathbf{w},b}(\mathbf{x}_i, u)) \}$$

$$\text{hinge loss : } \ell(t) := C \max(0, 1 - t)$$

Multiple Kernel Learning (MKL) Primal

Extend typical multiclass SVM to p feature maps $\phi_k(\mathbf{x}_i, y_i)$:

$$f_{w,b,\beta}(\mathbf{x}, y) = \sum_{k=1}^p \beta_k \langle w_k, \phi_k(x, y) \rangle + b_y$$

$$\underbrace{\sum_{k=1}^p \beta_k = 1, \beta_k \geq 0}$$

feature mapping (kernel) weight

Intuitive Formulation

$$\min_{\beta, w, b, \xi} \quad \frac{1}{2} \sum_{k=1}^p \beta_k \|w_k\|^2 + \sum_{i=1}^n \xi_i$$

$$\text{s.t.} \quad \forall i : \xi_i = \ell(f_{w,b,\beta}(x_i, y_i) - f_{w,b,\beta}(x_i, u))$$

Interpretable and intuitive, but in general not being convex :(

Let $\mathbf{v}_k := \beta_k \mathbf{w}_k$, then

MKL Primal

$$\begin{aligned} \min_{\beta, \mathbf{v}, b, \xi} \quad & \frac{1}{2} \sum_{k=1}^p \frac{\|\mathbf{v}_k\|^2}{\beta_k} + \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & \forall i : \xi_i = \ell(\langle \mathbf{v}, \psi_{iu} \rangle + b_{y_i} - b_u) \\ & \psi_{kiu} = \phi_k(x_i, y_i) - \phi_k(x_i, u) \\ & \psi_{iu} = (\psi_{kiu})_{k=1, \dots, p} \\ & \text{(combine into just one vector)} \end{aligned}$$

This is Convex!

Interpretable and intuitive, but in general not being convex :(

Let $\mathbf{v}_k := \beta_k \mathbf{w}_k$, then

MKL Primal

$$\begin{aligned} \min_{\beta, \mathbf{v}, b, \xi} \quad & \frac{1}{2} \sum_{k=1}^p \frac{\|\mathbf{v}_k\|^2}{\beta_k} + \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & \forall i : \xi_i = \ell(\langle \mathbf{v}, \psi_{iu} \rangle + b_{y_i} - b_u) \\ & \psi_{kiu} = \phi_k(x_i, y_i) - \phi_k(x_i, u) \\ & \psi_{iu} = (\psi_{kiu})_{k=1, \dots, p} \\ & \text{(combine into just one vector)} \end{aligned}$$

This is Convex!

Interpretable and intuitive, but in general not being convex :(

Let $\mathbf{v}_k := \beta_k \mathbf{w}_k$, then

MKL Primal

$$\begin{aligned} \min_{\beta, \mathbf{v}, b, \xi} \quad & \frac{1}{2} \sum_{k=1}^p \frac{\|\mathbf{v}_k\|^2}{\beta_k} + \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & \forall i : \xi_i = \ell(\langle \mathbf{v}, \psi_{iu} \rangle + b_{y_i} - b_u) \\ & \psi_{kiu} = \phi_k(x_i, y_i) - \phi_k(x_i, u) \\ & \psi_{iu} = (\psi_{kiu})_{k=1, \dots, p} \\ & \text{(combine into just one vector)} \end{aligned}$$

This is Convex!

Dual formulation for general loss function

$$\begin{aligned} \min_{\eta, \tilde{\alpha}, \gamma} \quad & \gamma + \sum_i \sum_{u \neq y_i} \eta_{iu} \ell^* \left(-\frac{\tilde{\alpha}_{iu}}{\eta_{iu}} \right) \\ \text{s.t.} \quad & \forall k : \quad \gamma \geq \frac{1}{2} \sum_{i,j} \sum_{u \neq y_i} \sum_{v \neq y_j} \tilde{\alpha}_{iu} \tilde{\alpha}_{jv} \langle \psi_{kiu}, \psi_{k jv} \rangle \\ & \forall i : \quad \forall u \neq y_i : 0 \leq \eta_{iu}, \\ & \forall i : \quad 1 = \sum_{u \neq y_i} \eta_{iu} \\ & \forall v : \quad 0 = \sum_i (1 - \delta_{y_i, v}) \tilde{\alpha}_{iv} - \sum_i \delta_{y_i, v} \sum_{u \neq y_i} \tilde{\alpha}_{iu} \end{aligned}$$

Here ℓ^* is the conjugate function of the loss function ℓ .

$$t_{iu} = \sum_k \langle \mathbf{v}_k, \psi_{kiu} \rangle + b_{y_i} - b_u \quad \text{score difference}$$

$$\xi_i \geq \ell(t_{iu}) \quad \text{Error upper bound}$$

So the Lagrangian is

$$L = \frac{1}{2} \sum_k \frac{1}{\beta_k} \|\mathbf{v}_k\|^2 + \sum_i \xi_i + \gamma \left(\sum_k \beta_k - 1 \right)$$

$$- \sum_k \epsilon_k \beta_k + \sum_i \sum_{u \neq y_i} \eta_{iu} (\ell(t_{iu}) - \xi_i)$$

$$+ \sum_i \sum_{u \neq y_i} \tilde{\alpha}_{iu} \left(t_{iu} - \sum_k \langle \mathbf{v}_k, \psi_{kiu} \rangle - b_{y_i} + b_u \right)$$

$$\frac{\partial L}{\partial \mathbf{v}_k} = \frac{1}{\beta_k} \mathbf{v}_k - \sum_i \sum_{u \neq y_i} \tilde{\alpha}_{iu} \psi_{kiu} = 0$$

$$\Rightarrow \mathbf{w}_k = \frac{1}{\beta_k} \mathbf{v}_k = \sum_i \sum_{u \neq y_i} \tilde{\alpha}_{iu} \psi_{kiu}$$

max min
 $\eta, \tilde{\alpha}, \gamma$ \mathbf{t}

$$\sum_i \sum_{u \neq y_i} (\eta_{iu} \ell(t_i u) + \tilde{\alpha}_{iu} t_{iu}) - \gamma$$

s.t. $\forall k : \gamma \geq \frac{1}{2} \|\mathbf{w}_k\|^2$ (Obtained from $\frac{\partial L}{\partial \beta_k} = 0$)

$\forall i : 1 = \sum_{u \neq y_i} \eta_{iu}$ (Obtained from $\frac{\partial L}{\partial \xi_i} = 0$)

$\forall v : 0 = \sum_i (1 - \delta_{y_i, v}) \tilde{\alpha}_{iv} - \sum_i \delta_{y_i, v} \sum_{u \neq y_i} \tilde{\alpha}_{iu}$ ($\frac{\partial L}{\partial b_v} = 0$)

$$\begin{aligned}
& \max_{\eta, \tilde{\alpha}, \gamma} \min_{\mathbf{t}} \sum_i \sum_{u \neq y_i} (\eta_{iu} \ell(t_i u) + \tilde{\alpha}_{iu} t_{iu}) - \gamma \\
&= \max_{\eta, \tilde{\alpha}, \gamma} \min_{\mathbf{t}} \sum_i \sum_{u \neq y_i} \eta_{iu} \left(\ell(t_i u) + \frac{\tilde{\alpha}_{iu} t_{iu}}{\eta_{iu}} \right) - \gamma \\
&= \max_{\eta, \tilde{\alpha}, \gamma} \min_{\mathbf{t}} \sum_i \sum_{u \neq y_i} -\eta_{iu} \left(-\ell(t_i u) - \frac{\tilde{\alpha}_{iu} t_{iu}}{\eta_{iu}} \right) - \gamma \\
&= \max_{\eta, \tilde{\alpha}, \gamma} \sum_i \sum_{u \neq y_i} -\eta_{iu} \ell^* \left(-\frac{\tilde{\alpha}_{iu}}{\eta_{iu}} \right) - \gamma \\
&\iff \min_{\eta, \tilde{\alpha}, \gamma} \gamma + \sum_i \sum_{u \neq y_i} \eta_{iu} \ell^* \left(-\frac{\tilde{\alpha}_{iu}}{\eta_{iu}} \right)
\end{aligned}$$

The definition of conjugate function of f is defined as

$$f^*(y) = \sup_{x \in \text{dom}(f)} (y^T x - f(x))$$

$$\begin{aligned}
 \min_{\gamma, \tilde{\alpha}, \eta} \quad & \gamma + \sum_i \sum_{u \neq y_i} \eta_{iu} \ell^* \left(-\frac{\tilde{\alpha}_{iu}}{\eta_{iu}} \right) \\
 \text{s.t.} \quad & \forall k : \quad \gamma \geq \frac{1}{2} \sum_{i,j} \sum_{u \neq y_i} \sum_{v \neq y_j} \tilde{\alpha}_{iu} \tilde{\alpha}_{jv} \langle \psi_{kiu}, \psi_{k jv} \rangle \\
 & \forall i : \quad \forall u \neq y_i : 0 \leq \eta_{iu}, \\
 & \forall i : \quad 1 = \sum_{u \neq y_i} \eta_{iu} \\
 & \forall v : \quad 0 = \sum_i (1 - \delta_{y_i, v}) \tilde{\alpha}_{iv} - \sum_i \delta_{y_i, v} \sum_{u \neq y_i} \tilde{\alpha}_{iu}
 \end{aligned}$$

Dual Formulation for Hinge Loss

When the loss $\ell(t) = C \cdot \max(0, 1 - t)$,

$$\ell^*(\nu) = \begin{cases} \nu & -C \leq \nu \leq 0 \\ \infty & \text{otherwise} \end{cases}$$

So in order to make the dual solvable, we must have

$$\begin{aligned} & -C \leq -\frac{\tilde{\alpha}_{iu}}{\eta_{iu}} \leq 0 \\ \text{As } \forall i: & 1 = \sum_{u \neq y_i} \eta_{iu} \\ \text{and } \forall i: & \forall u \neq y_i : 0 \leq \eta_{iu} \\ \text{We have } & \sum_{u \neq y_i} \tilde{\alpha}_{iu} \leq C \\ & \forall u \neq y_i : 0 \leq \tilde{\alpha}_{iu} \end{aligned} \quad \begin{aligned} & \gamma + \sum_i \sum_{u \neq y_i} \eta_{iu} \ell^* \left(-\frac{\tilde{\alpha}_{iu}}{\eta_{iu}} \right) \\ & = \gamma + \sum_i \sum_{u \neq y_i} \eta_{iu} \left(-\frac{\tilde{\alpha}_{iu}}{\eta_{iu}} \right) \\ & = \gamma - \sum_i \sum_{u \neq y_i} \tilde{\alpha}_{iu} \end{aligned}$$

Dual Formulation for Hinge Loss

When the loss $\ell(t) = C \cdot \max(0, 1 - t)$,

$$\ell^*(\nu) = \begin{cases} \nu & -C \leq \nu \leq 0 \\ \infty & \text{otherwise} \end{cases}$$

So in order to make the dual solvable, we must have

$$-C \leq -\frac{\tilde{\alpha}_{iu}}{\eta_{iu}} \leq 0$$

As $\forall i: \quad 1 = \sum_{u \neq y_i} \eta_{iu}$

and $\forall i: \quad \forall u \neq y_i: 0 \leq \eta_{iu}$

We have $\sum_{u \neq y_i} \tilde{\alpha}_{iu} \leq C$

$$\forall u \neq y_i: 0 \leq \tilde{\alpha}_{iu}$$

$$\begin{aligned} & \gamma + \sum_i \sum_{u \neq y_i} \eta_{iu} \ell^* \left(-\frac{\tilde{\alpha}_{iu}}{\eta_{iu}} \right) \\ &= \gamma + \sum_i \sum_{u \neq y_i} \eta_{iu} \left(-\frac{\tilde{\alpha}_{iu}}{\eta_{iu}} \right) \\ &= \gamma - \sum_i \sum_{u \neq y_i} \tilde{\alpha}_{iu} \end{aligned}$$

Dual Formulation for Hinge Loss

When the loss $\ell(t) = C \cdot \max(0, 1 - t)$,

$$\ell^*(\nu) = \begin{cases} \nu & -C \leq \nu \leq 0 \\ \infty & \text{otherwise} \end{cases}$$

So in order to make the dual solvable, we must have

$$\begin{aligned} & -C \leq -\frac{\tilde{\alpha}_{iu}}{\eta_{iu}} \leq 0 \\ \text{As } \forall i: & \quad 1 = \sum_{u \neq y_i} \eta_{iu} \\ \text{and } \forall i: & \quad \forall u \neq y_i : 0 \leq \eta_{iu} \\ \text{We have } & \quad \sum_{u \neq y_i} \tilde{\alpha}_{iu} \leq C \\ & \quad \forall u \neq y_i : 0 \leq \tilde{\alpha}_{iu} \end{aligned} \quad \begin{aligned} & \gamma + \sum_i \sum_{u \neq y_i} \eta_{iu} \ell^* \left(-\frac{\tilde{\alpha}_{iu}}{\eta_{iu}} \right) \\ & = \gamma + \sum_i \sum_{u \neq y_i} \eta_{iu} \left(-\frac{\tilde{\alpha}_{iu}}{\eta_{iu}} \right) \\ & = \gamma - \sum_i \sum_{u \neq y_i} \tilde{\alpha}_{iu} \end{aligned}$$

Dual for Hinge Loss (Folded formulation)

$$\min_{\tilde{\alpha}, \gamma} \quad \gamma - \sum_i \sum_{u \neq y_i} \tilde{\alpha}_{iu}$$

$$\text{s.t.} \quad \forall k : \gamma \geq \frac{1}{2} \|w_k(\tilde{\alpha})\|^2$$

$$\sum_{u \neq y_i} \tilde{\alpha}_{iu} \leq C$$

$$\forall u \neq y_i : 0 \leq \tilde{\alpha}_{iu}$$

$$\forall v : \quad 0 = \sum_i (1 - \delta_{y_i, v}) \tilde{\alpha}_{iv} - \sum_i \delta_{y_i, v} \sum_{u \neq y_i} \tilde{\alpha}_{iu}$$

$$\alpha_{iu} = \begin{cases} -\tilde{\alpha}_{iu} & \text{if } u \neq y_i \\ \sum_{v \neq y_i} \tilde{\alpha}_{iv} & \text{if } u = y_i \end{cases}$$

Unfolded Formulation

$$\min_{\alpha, \gamma} \quad \gamma - \sum_i \alpha_{i, y_i}$$

$$\text{s.t.} \quad \forall k : \gamma \geq \frac{1}{2} \|w_k(\alpha)\|^2$$

$$\alpha \in \mathcal{S} := \left\{ \alpha \left| \begin{array}{l} \forall i : 0 \leq \alpha_{i, y_i} \leq C \\ \forall i : \forall u \neq y_i : \alpha_{iu} \leq 0 \\ \forall i : \sum_{u \in \mathcal{Y}} \alpha_{iu} = 0 \\ \forall u \in \mathcal{Y} : \sum_i \alpha_{iu} = 0 \end{array} \right. \right\}$$

SILP formulation

$$\max_{\beta, \theta} \quad \theta$$

$$s.t. \quad \forall \alpha \in \mathcal{S} : \theta \leq \frac{1}{2} \sum_k \beta_k \|w_k(\alpha)\|^2 - \sum_i \alpha_{i,y_i}$$

- Both version of formulation are QCQPs.
- When there's only one single kernel, it reduces to the dual of mSVM.
- When there are only 2 classes, the formulation reduces to Gert's formulation of kernel selection
- Unfolded formulation can be easily transformed into semi-infinite linear programming.

Connection to other Regularizer

$$\begin{aligned} \min_{\beta, \mathbf{w}, \mathbf{b}, \xi, \mathbf{s}} \quad & \frac{1}{2} \sum_{k=1}^p \beta_k \|\mathbf{w}_k\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & \xi_i = \max_{u \neq y_i} s_{iu}, \quad s_{iu} \geq 0 \\ & \sum_{k=1}^p \beta_k \langle \mathbf{w}_k, \Psi_{kiu} \rangle + b_{y_i} - b_u \geq 1 - s_{iu} \end{aligned}$$

Intuitive MKL: Equation (3); Non-Convex

substitute $\mathbf{v}_k := \beta_k \mathbf{w}_k$

$$\begin{aligned} \min_{\beta, \mathbf{v}, \mathbf{b}, \xi, \mathbf{s}} \quad & \frac{1}{2} \sum_{k=1}^p \frac{1}{\beta_k} \|\mathbf{v}_k\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & \xi_i \geq s_{iu}, \quad s_{iu} \geq 0 \\ & \sum_{k=1}^p \langle \mathbf{v}_k, \Psi_{kiu} \rangle + b_{y_i} - b_u \geq 1 - s_{iu} \end{aligned}$$

Equation (4); Convex

$$\begin{aligned} \min_{\beta, \mathbf{w}, \mathbf{b}, \xi, \mathbf{s}} \quad & \frac{1}{2} \left(\sum_{k=1}^p \beta_k \|\mathbf{w}_k\| \right)^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & \xi_i = \max_{u \neq y_i} s_{iu}, \quad s_{iu} \geq 0 \\ & \sum_{k=1}^p \beta_k \langle \mathbf{w}_k, \Psi_{kiu} \rangle + b_{y_i} - b_u \geq 1 - s_{iu} \end{aligned}$$

Generalized from Sonnenburg et. al. [20]; Non-Convex

substitute $\mathbf{v}_k := \beta_k \mathbf{w}_k$

$$\begin{aligned} \min_{\mathbf{v}, \mathbf{b}, \xi, \mathbf{s}} \quad & \frac{1}{2} \left(\sum_{k=1}^p \|\mathbf{v}_k\| \right)^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & \xi_i = \max_{u \neq y_i} s_{iu}, \quad s_{iu} \geq 0 \\ & \sum_{k=1}^p \langle \mathbf{v}_k, \Psi_{kiu} \rangle + b_{y_i} - b_u \geq 1 - s_{iu} \end{aligned}$$

Generalized from Bach et. al. [1]; Convex

Theorem 1

(1)

$$\begin{aligned} \min_{\alpha} \quad & \gamma - \sum_i \alpha_{iy_i} \\ \text{s.t.} \quad & \forall i : 0 \leq \alpha_{iy_i} \leq C, \quad \forall i : \forall u \neq y_i : \alpha_{iu} \leq 0 \\ & \forall i : \sum_{u \in \mathcal{Y}} \alpha_{iu} = 0 \quad \forall u : \sum_i \alpha_{iu} = 0 \\ & \forall k : \gamma \geq \frac{1}{2} \sum_{i, j, u, v} \alpha_{iu} \alpha_{jv} \langle \Phi_k(\mathbf{x}_i, u), \Phi_k(\mathbf{x}_j, v) \rangle \end{aligned}$$

Common Lagrange Dual; Equation (9)

Time Complexity Comparison

	Examples	Classes	Kernels
QP, unfolded	2.5	1.8	–
QCQP, unfolded (9)	3.0	2.0	2.3
SILP, unfolded (11)	2.4	1.7	1.1
SILP, compact	2.6	2.2	1.0

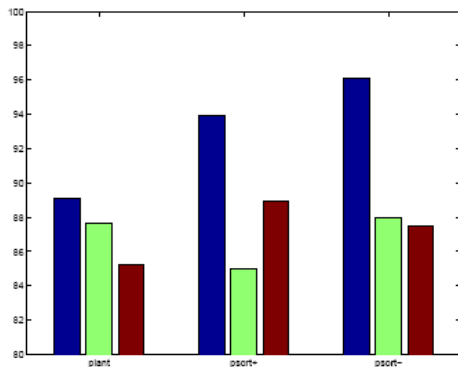


Figure 3: Protein Subcellular Localization results. The bars within each group correspond to different methods: (left, blue) MKL; (center, green) unweighted sum of kernels; (right, red) current state of the art.

- 1 Contribution: Generalize the kernel selection problem to any kinds of convex loss functions and multiple classes.
- 2 If just focus on multi-class SVM, deriving from the dual formulation is much easier.
- 3 Possible interesting direction: kernel selection in y space, kernel selection for structured output.
- 4 Whether or not all the classes should share the same feature space (Mentioned in the paper, but has been done by us).
- 5 Can we remove the weight β in the regularizer? Different formulation and variants might need more understanding.