

Details on Optimizing P1

1 Optimization of P1

For notation convenience, we switch the symbols v_{m+1} and w_{m+1} in the original **P1**. Then **P1** in the submitted paper becomes

$$\begin{aligned}
 \mathbf{P1} \quad & \min_{\mathbf{w}, \mathbf{v}, \boldsymbol{\xi}, \mathbf{d}} \frac{1}{2} \sum_{l=1}^{m+1} d_l \|v_l\|^2 + \frac{1}{2} \sum_{l=1}^m C_l d_l \|w_l\|^2 + \frac{1}{2} \sum_{l=1}^m \sum_{i=1}^n d_l (\xi_i^l)^2 + C \sum_{i=1}^n \xi_i^{m+1} \\
 \text{s.t.} \quad & y_i^l (w_l^T \phi_l(x_i^l)) \geq 1 - \xi_i^l, \quad 1 \leq l \leq m, \quad 1 \leq i \leq n; \\
 & \mathbf{d} \geq \mathbf{0}, \quad \|\mathbf{d}\|_p^p \leq 1, \quad \boldsymbol{\xi} \geq \mathbf{0}, \\
 & y_i^{m+1} \left(\sum_{l=1}^{m+1} d_l (w_l + v_l)^T \phi_l(x_i^l) \right) \geq 1 - \xi_i^{m+1}, \quad \text{when } 1 \leq i \leq n
 \end{aligned} \tag{1}$$

where $\mathbf{0}$ is a column vector whose elements are all 0's. Inequalities between two vectors are taken element-wise, and $C_l, 1 \leq l \leq m$ and C are positive user defined parameters.

To convert **P1** into a convex optimization problem, we can simply replace w_l , v_l ($1 \leq l \leq m+1$), and ξ_i^l ($1 \leq i \leq n$ and $1 \leq l \leq m$) with \hat{w}_l/d_l , \hat{v}_l/d_l , and $\hat{\xi}_i^l/d_l$, respectively. If $d_l = 0$, we define $a/d_l = \infty$ when $a \neq 0$, and $a/d_l = 0$ when $a = 0$. Omitting the *hat* notation for \hat{w} , \hat{v} , and $\hat{\xi}$ for simplicity, **P1** becomes

$$\mathbf{P2} \quad \min_{\mathbf{w}, \mathbf{v}, \boldsymbol{\xi}, \mathbf{d}} \frac{1}{2} \sum_{l=1}^{m+1} \frac{\|v_l\|^2}{d_l} + \frac{1}{2} \sum_{l=1}^m \frac{C_l \|w_l\|^2}{d_l} + \frac{1}{2} \sum_{l=1}^m \sum_{i=1}^n \frac{(\xi_i^l)^2}{d_l} + C \sum_{i=1}^n (\xi_i^{m+1}) \tag{2}$$

$$\text{s.t.} \quad y_i^l (w_l^T \phi_l(x_i^l)) \geq d_l - \xi_i^l, \quad 1 \leq l \leq m; \quad 1 \leq i \leq n \tag{3}$$

$$y_i^{m+1} \left(\sum_{l=1}^{m+1} (w_l + v_l) \phi_l(x_i^l) \right) \geq 1 - \xi_i^{m+1}, \quad 1 \leq i \leq n \tag{4}$$

$$\mathbf{d} \geq \mathbf{0}, \quad \|\mathbf{d}\|_p^p \leq 1, \quad \boldsymbol{\xi} \geq \mathbf{0} \tag{5}$$

Now it is clear from **P2** that the quadratic error weighted by ' \mathbf{d} ' (see Eq. 1.d of the submitted paper) make it easy to formulate a convex problem. Without the ' \mathbf{d} ' in Eq. (1.d) of the submitted paper, the error terms for auxiliary tasks in **P2** become $(\xi_i^l)^2/d_l^2$ which is not convex. It is also important to note that analytically eliminating ' \mathbf{d} ' in **P2** is not so simple as that in MKL (nor its p-norm variant considered by Micchelli and Pontil (2007)) with the technique proposed by Rakotomamonjy et al. (2008), and Micchelli and Pontil (2007), because ' \mathbf{d} ' presents in both the objective function and auxiliary task constraints of **P2**.

We can obtain the dual of **P2** following a standard method as that used in the non-sparse MKL (Kloft et al., 2009). However, directly solving the dual problem may create some numerical problems as discussed by Kloft et al. (2009). Hence, we employ the cutting plane algorithm to solve **P2**. The dual objective will still be used for checking the stopping condition, i.e., the relative duality gap. Fixing \mathbf{d} , we can solve the partial Lagrangian w.r.t. \mathbf{w} , \mathbf{v} , and $\boldsymbol{\xi}$. For brevity, we provide the semi-infinite programming formulation directly:

$$\mathbf{P3} \quad \min_{\mathbf{d}, \rho} \quad \rho \quad \text{s.t.} \quad \mathbf{d} \geq \mathbf{0}, \quad \|\mathbf{d}\|_p^p \leq 1,$$

$$\text{and} \quad \rho \geq \left(d_1 \mathbf{e}^T, d_2 \mathbf{e}^T, \dots, d_m \mathbf{e}^T, \mathbf{e}^T \right) \boldsymbol{\alpha} - \frac{1}{2} \left(\sum_{l=1}^{m+1} d_l Q_l \right)$$

$$\text{for all } \boldsymbol{\alpha} \text{ satisfying } \mathbf{0} \leq \boldsymbol{\alpha}, \text{ and } \boldsymbol{\alpha}^{m+1} \leq C$$

where $\boldsymbol{\alpha}$ is the Lagrange multiplier vector such that $\boldsymbol{\alpha} = \left((\boldsymbol{\alpha}^l)_{(l=1)}^{(m+1)} \right)$, and for each l , $\boldsymbol{\alpha}^l = \left((\alpha_i^l)_{(i=1)}^{(n)} \right)$, and $Q_l = \left(\boldsymbol{\alpha}^l \right)^T \left(\frac{1}{C_l} K_l \circ \left(\mathbf{y}^l (\mathbf{y}^l)^T \right) + \mathbf{I} \right) \boldsymbol{\alpha}^l + \bar{C}_l \left(\boldsymbol{\alpha}^{m+1} \right)^T \left(K_l \circ \left(\mathbf{y}^{m+1} (\mathbf{y}^{m+1})^T \right) \right) \boldsymbol{\alpha}^{m+1}$

$+ \frac{2}{\bar{C}_l} (\alpha^l)^T \left(K_l \circ (\mathbf{y}^l (\mathbf{y}^{m+1})^T) \right) \alpha^{m+1}$, where $\bar{C}_l = 1/C_l + 1$, for $1 \leq l \leq m$.

When $l = m + 1$,

$Q_{m+1} = (\alpha^{m+1})^T \left(K_{m+1} \circ (\mathbf{y}^{m+1} (\mathbf{y}^{m+1})^T) \right) \alpha^{m+1}$, where “ \circ ” denotes the element-wise product between two matrices.

P3 can be solved by the cutting plane algorithm with a standard QP (quadratic programming) to find its most violated constraint. For the *restricted master problem* (Sonnenburg et al., 2006) of this algorithm, we use CVX¹ directly, rather than solve an approximate problem (Kloft et al., 2009). Based on Slater’s condition, we can use the relative duality gap as a stopping criterion and we set a threshold of 10^{-2} in our experiment.

2 Derivation for the Dual of P2

The dual of **P2** is used in checking the stopping condition. Suppose the Lagrangian multipliers for (3) and (4) in **P2** are $((\alpha_i^l)_{(l,i)=(1,1)}^{(m+1,n)})$. Multipliers for the three terms in (5) are $((\gamma_l)_{l=1}^{m+1})$, $(1/p)\beta$, and $((\eta_i^l)_{(l,i)=(1,1)}^{(m+1,n)})$ respectively.

Let’s define $D = \{a/b : a > b, a \text{ is a positive even number, and } b \text{ is a positive odd number.}\}$

Proposition 1 *If $p \in D$, the Wolfe dual of P2 is the following:*

$$\mathbf{D1} \quad \max_{\alpha} \quad e^T \alpha^{m+1} - \left(\sum_{l=1}^{m+1} (Q'_l)^{\frac{p}{p-1}} \right)^{\frac{p-1}{p}} \quad (6)$$

$$\text{s.t.} \quad \mathbf{0} \leq \alpha, \text{ and } \alpha^{m+1} \leq C \quad (7)$$

for $1 \leq l \leq m$, we let

$$Q'_l = \max \left(\frac{1}{2} Q_l - e^T \alpha^l, 0 \right) \quad (8)$$

where

$$\begin{aligned} Q_l = & (\alpha^l)^T \left(\frac{1}{C_l} K_l \circ (\mathbf{y}^l (\mathbf{y}^l)^T) + \mathbf{I} \right) \alpha^l \\ & + \bar{C}_l (\alpha^{m+1})^T \left(K_l \circ (\mathbf{y}^{m+1} (\mathbf{y}^{m+1})^T) \right) \alpha^{m+1} \\ & + \frac{2}{C_l} (\alpha^l)^T \left(K_l \circ (\mathbf{y}^l (\mathbf{y}^{m+1})^T) \right) \alpha^{m+1} \end{aligned} \quad (9)$$

where $\bar{C}_l = 1/C_l + 1$ and “ \circ ” represents the element-wise product between two matrices. When $l = m + 1$, $Q'_{m+1} = (1/2)Q_{m+1}$ where

$$Q_{m+1} = (\alpha^{m+1})^T \left(K_{m+1} \circ (\mathbf{y}^{m+1} (\mathbf{y}^{m+1})^T) \right) \alpha^{m+1}$$

¹Available at: <http://stanford.edu/~boyd/cvx>

Proof. Taking derivative w.r.t. to the Lagrangian $\mathcal{L}(\alpha, \gamma, \beta, \eta)$ of **P2** we have

$$\frac{\partial \mathcal{L}}{\partial w_l} = C_l \frac{w_l^T}{d_l} - \sum_{i=1}^n \alpha_i^l y_i^l \phi_l(x_i^l) - \sum_{i=1}^n \alpha_i^{m+1} y_i^{m+1} \phi_l(x_i^l), 1 \leq l \leq m \quad (10)$$

$$\frac{\partial \mathcal{L}}{\partial v_l} = \frac{v_l^T}{d_l} - \sum_{i=1}^n \alpha_i^{m+1} y_i^{m+1} \phi_l(x_i^l), 1 \leq l \leq m+1 \quad (11)$$

$$\frac{\partial \mathcal{L}}{\partial \xi_i^l} = \frac{\xi_i^l}{d_l} - \alpha_i^l - \eta_i^l, 1 \leq l \leq m; 1 \leq i \leq n \quad (12)$$

$$\frac{\partial \mathcal{L}}{\partial \xi_i^{m+1}} = C - \alpha_i^{m+1} - \eta_i^{m+1}, 1 \leq i \leq n \quad (13)$$

$$\frac{\partial \mathcal{L}}{\partial d_l} = -\frac{C_l}{2} \frac{\|w_l\|^2}{d_l^2} - \frac{\|v_l\|^2}{d_l^2} - \frac{1}{2} \sum_{i=1}^n \frac{(\xi_i^l)^2}{d_l^2} + e^T \alpha^l + \beta d_l^{p-1} - \gamma_l, 1 \leq l \leq m \quad (14)$$

$$\frac{\partial \mathcal{L}}{\partial d_{m+1}} = -\frac{1}{2} \frac{\|v_{m+1}\|^2}{d_{m+1}^2} + \beta d_{m+1}^{p-1} - \gamma_{m+1}, \quad (15)$$

Setting all these to zero, and plugging all back into the Lagrangian, we have

$$\mathcal{L}(\alpha, \gamma, \beta, \eta) = e^T \alpha^{m+1} - \frac{1}{p} \beta - \frac{p-1}{p} \beta^{\frac{-1}{p-1}} \left(\sum_{l=1}^{m+1} (G_l)^{\frac{p}{p-1}} \right) \quad (16)$$

where for $1 \leq l \leq m$

$$G_l = Q_l'' + \gamma_l + \eta_l^T \alpha^l + \frac{1}{2} \eta_l^T \eta_l \quad (17)$$

and

$$G_{m+1} = Q_{m+1}'' + \gamma_{m+1} \quad (18)$$

where $Q_l'' = \frac{1}{2} Q_l - e^T \alpha^l$, $1 \leq l \leq m$ and $Q_{m+1}'' = Q'_{m+1}$.

We have $p \in D$ and this ensures the following: when solving for d_l from Eq.(14) after replacing $\|w_l\|^2/d_l^2$, $\|v_l\|^2/d_l^2$ and $(\xi_i^l)^2/d_l^2$ by Eq.(11), (12), and (13), we always obtain a real solution for d_l . Also, since $p \in D$, $(G_l)^{\frac{p}{p-1}}$ is always a non-negative value. Thus, to maximize $\mathcal{L}(\alpha, \gamma, \beta, \eta)$: (a). if $Q_l'' \geq 0$, then γ_l and $\eta_l = \mathbf{0}$; (b). if $Q_l'' \leq 0$, then $\gamma_l = -Q_l''$, and still $\eta_l = \mathbf{0}$. Hence, to maximize Eq.(16), we can eliminate γ and η , and replace all the Q_l'' with Q' . Further, following Kloft et al. (2009), let's take derivative of (16) w.r.t β and ignoring that it is non-negative. Now, at the maximum of the Lagrangian,

$$\beta = \left(\sum_{l=1}^{m+1} (Q_l')^{\frac{p}{p-1}} \right)^{\frac{p-1}{p}} \quad (19)$$

and β here is always non-negative. This concludes the proof. \blacksquare

References

- M. Kloft, U. Brefeld, S. Sonnenburg, P. Laskov, K.R. Müller, , and A. Zien. Efficient and accurate lp-norm multiple kernel learning. In *Advances in Neural Information Processing Systems 22*, pp. 997–1005, 2009.
- C.A. Micchelli, and M. Pontil. Feature space perspectives for learning the kernel. *Machine Learning*, 66(2): 297–319, 2007.
- A. Rakotomamonjy, F. Bach, Y. Grandvalet, and S. Canu. SimpleMKL. *Journal of Machine Learning Research*, 9:2491–2521, 2008.
- S. Sonnenburg, G. Rätsch, C. Schäfer, and B. Schölkopf. Large scale multiple kernel learning. *The Journal of Machine Learning Research*, 7:1565, 2006.