

# Efficient Multiple Kernel Learning (2)

Lei Tang

Arizona State University

May. 29th, 2007

For simplicity, we assume all Kernel matrix have the same trace (without scaling problem). Given  $G_1, G_2, \dots, G_p$ , the weight for each kernel is  $\beta_k$ , and  $\sum \beta_k = 1$ .

## Goal

$$\begin{aligned} \min_{\beta} \max_{\alpha} \quad & \alpha^T e - \frac{1}{2} \alpha^T \text{diag}(y) \left( \sum_{k=1}^p \beta_k G_k \right) \text{diag}(y) \alpha \\ \text{s.t.} \quad & \alpha^T y = 0 \\ & C \geq \alpha \geq 0 \\ & \beta^T e = 1, \beta \geq 0 \end{aligned}$$

$$\begin{aligned}
& \min_{\beta} \max_{\alpha} \alpha^T e - \frac{1}{2} \alpha^T \text{diag}(y) \left( \sum_{k=1}^p \beta_k G_k \right) \text{diag}(y) \alpha \\
= & \max_{\alpha} \min_{\beta} \alpha^T e - \frac{1}{2} \alpha^T \text{diag}(y) \left( \sum_{k=1}^p \beta_k G_k \right) \text{diag}(y) \alpha \\
= & \max_{\alpha} \alpha^T e - \max_{\sum_k \beta_k = 1} \frac{1}{2} \sum_{k=1}^p \beta_k \left( \alpha^T \text{diag}(y) G_k \text{diag}(y) \alpha \right) \\
= & \max_{\alpha} \alpha^T e - \max_i \left( \frac{1}{2} \alpha^T \text{diag}(y) G_k \text{diag}(y) \alpha \right)
\end{aligned}$$

So the problem can be reformulated as

$$\begin{aligned}
 \max_{\alpha} \quad & \alpha^T e - \frac{1}{2}t \\
 \text{s.t.} \quad & t \geq \alpha^T y = 0 \\
 & C \geq \alpha \geq 0 \\
 & t \geq \frac{1}{2} \alpha^T \text{diag}(y) G_k \text{diag}(y) \alpha \quad \forall k
 \end{aligned}$$

This is a QCQP problem, which can be solved by general optimization package. But it does not scale up!!

Let

$$S_k(\alpha) = \alpha^T e - \frac{1}{2} \alpha^T \text{diag}(y) G_k \text{diag}(y) \alpha$$

As  $\sum \beta_k = 1$ , the objective becomes

$$\min_{\beta} \max_{\alpha} \sum_{k=1}^p \beta_k S_k(\alpha)$$

So the goal is equivalent to

$$\begin{aligned} \min_{\beta} \max_{\alpha} \quad & \sum_{k=1}^p \beta_k S_k(\alpha) \\ \text{s.t.} \quad & \alpha^T y = 0 \\ & C \geq \alpha \geq 0 \\ & \beta^T e = 1, \beta \geq 0 \end{aligned}$$

$$\begin{aligned}
\min_{\beta} \max_{\alpha} \quad & \sum_{k=1}^p \beta_k S_k(\alpha) \\
\text{s.t.} \quad & \alpha^T y = 0 \\
& C \geq \alpha \geq 0 \\
& \beta^T e = 1, \beta \geq 0
\end{aligned}$$

Assume  $\alpha^*$  is the optimal, let  $\theta^* := \sum_{k=1}^p \beta_k S_k(\alpha)$  would be the maximal,

$$\sum_{k=1}^p \beta_k S_k(\alpha) \leq \theta^*$$

$$\begin{aligned} \min \quad & \theta \\ \text{s.t.} \quad & \beta \geq 0, \sum_k \beta_k = 1 \\ & \sum_{k=1}^p \beta_k S_k(\alpha) \leq \theta \\ & \text{for all } \alpha \in \mathcal{R}^N, 0 \leq \alpha \leq C, \alpha^T y = 0 \end{aligned}$$

Note that  $\theta$  and  $\beta$  are linearly constrained with infinitely many constraints for all possible  $\alpha$ .

$$\begin{aligned} \max \quad & \theta \\ \text{s.t.} \quad & \beta \geq 0, \sum_k \beta_k = 1, \sum_{k=1}^p \beta_k S_k(\alpha) \geq \theta \\ & \text{for all } \alpha \in \mathcal{C} \end{aligned}$$

## Column Generation

1. *Restricted master problem*: Compute optimal  $(\beta, \theta)$  for a restricted subsets of constraints (**Typical Linear Programming problem**)
2. Add the constraints that maximize the constraint violation for the given intermediate solution  $(\beta, \theta)$ . (**Exactly an SVM dual formulation with a combined kernel**)

Idea is exactly the same as *cutting plane*



$$\begin{aligned} \max \quad & \theta \\ \text{s.t.} \quad & \beta \geq 0, \sum_k \beta_k = 1, \sum_{k=1}^p \beta_k S_k(\alpha) \geq \theta \\ & \text{for all } \alpha \in \mathcal{C} \end{aligned}$$

## Column Generation

1. *Restricted master problem*: Compute optimal  $(\beta, \theta)$  for a restricted subsets of constraints (**Typical Linear Programming problem**)
2. Add the constraints that maximize the constraint violation for the given intermediate solution  $(\beta, \theta)$ . (**Exactly an SVM dual formulation with a combined kernel**)

Idea is exactly the same as *cutting plane*

## Other tricks to accelerate calculation

- Step 2 is not necessary to be exact if  $\beta$  is still far away from optimal.
- *Chunking*: maintain a smaller number of optimization variables.
- *Warm Start* for step 1 to solve LP
- Specific data structure to maintain the gram matrix (tries), save memory. Kernel caching.
- Parallel processing

Works up to 20 kernels and one million examples

## Goal

$$\begin{aligned} \min_{\beta} \max_{\alpha} \quad & \alpha^T e - \frac{1}{2} \alpha^T \text{diag}(y) \left( \sum_{k=1}^p \beta_k G_k \right) \text{diag}(y) \alpha \\ \text{s.t.} \quad & \alpha^T y = 0 \\ & C \geq \alpha \geq 0 \\ & \beta^T e = 1, \beta \geq 0 \end{aligned}$$

Let

$$J(\beta) = \begin{cases} \max_{\alpha} & \alpha^T e - \frac{1}{2} \alpha^T \text{diag}(y) \left( \sum_{k=1}^p \beta_k G_k \right) \text{diag}(y) \alpha \\ \text{s.t.} & \alpha^T y = 0 \\ & C \geq \alpha \geq 0 \end{cases}$$

The problem can be reformulated as

$$\min_{\beta} J(\beta) \text{ such that } \sum_{k=1}^p \beta_k = 1, \beta \geq 0$$

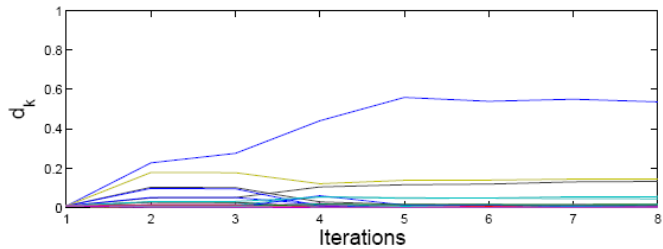
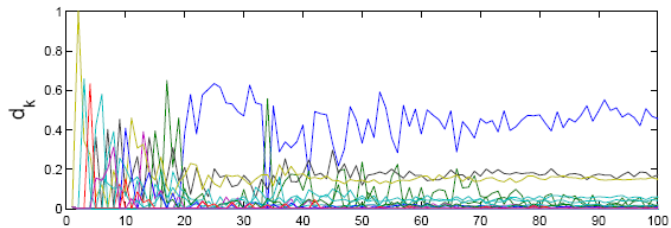
Let  $\alpha^*$  be the optimal value for  $J(\beta)$ , then

$$J(\beta) = \alpha^{*T} e - \frac{1}{2} \alpha^{*T} \text{diag}(y) \left( \sum_{k=1}^p \beta_k G_k \right) \text{diag}(y) \alpha^*$$
$$\frac{\partial J(\beta)}{\partial \beta_k} = -\frac{1}{2} \alpha^{*T} \text{diag}(y) G_k \text{diag}(y) \alpha^*$$

- For fixed  $\beta$ , calculate the optimal  $J(\beta)$ ; (Exactly an SVM problem)
- Calculate a decent direction for  $\beta$  such that the constraint is satisfied.
- Update  $\beta$  by finding the step length using some line search method (e.g. Armijo's rule). This involves multiple evaluations of  $J(\beta)$  ( A single kernel SVM) with some small variants of  $\beta$ . This could be speed up by initializing the SVM with  $\alpha^*$ .

# Cutting Plane v.s. Steepest Decent

- Cutting plane gradually add constraints whereas steepest decent involves fixed number of constraints in each iteration.
- Both involves calculation of SVM in each iteration. Steepest Decent involves a little more calculation while computing the step length.
- Convergence Analysis: Both converge. But Cutting planes method are unstable, especially when the number of lower-bounding affine functions is small.



Credit  $M = 208$

Algorithm	# Kernel	Accuracy	Time (s)
SILP	$10.4 \pm 1.9$	$86.4 \pm 1.4$	$124 \pm 28$
Our	$16.8 \pm 5.4$	$86.3 \pm 1.4$	$32.7 \pm 9.7$

Ionosphere  $M = 442$

Algorithm	# Kernel	Accuracy	Time (s)
SILP	$19.9 \pm 2.1$	$92.2 \pm 1.5$	$152 \pm 39$
Our	$30.5 \pm 7.0$	$92.3 \pm 1.4$	$18.1 \pm 5.8$

Sonar  $M = 793$

Algorithm	# Kernel	Accuracy	Time (s)
SILP	$29.1 \pm 3.5$	$79.2 \pm 4.6$	$383 \pm 113$
Our	$46.0 \pm 7.6$	$78.6 \pm 4.2$	$21.8 \pm 16.9$



- ① Large Scale Multiple Kernel Learning, JMLR, 2006
- ② More Efficiency in Multiple Kernel Learning, ICML, 2007