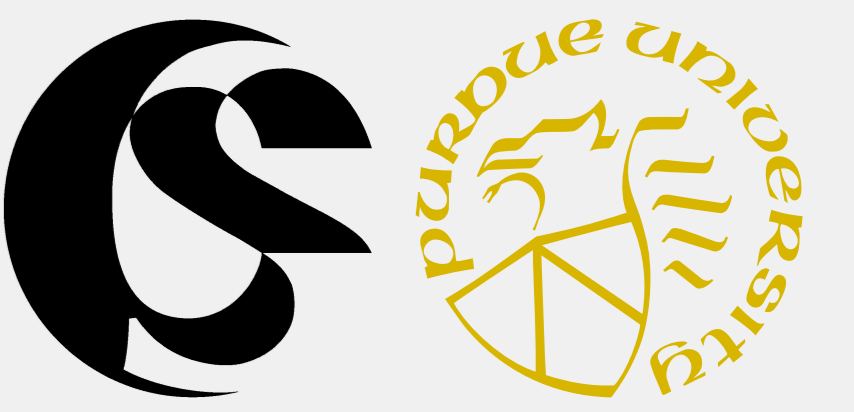


Optimality Implies Kernel Sum Classifiers are Statistically Efficient

Raphael Arkady Meyer Jean Honorio

Computer Science Department, Purdue University



Introduction and Optimization

Background

Statistical Learning Theory + Optimization

- Generalization proofs typically state that all feasible estimators generalize well
- This includes low-accuracy estimators we do not care about
- Proofs often make stringent assumptions on the data distribution
- We combine Optimization and Statistical Learning Theory to prove that optimal estimators generalize well**
- We justify common assumptions made in the Multiple Kernel Learning literature

Multiple Kernel Learning

- Given m kernels k_1, \dots, k_m and dataset $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$
- An estimator picks $\theta_1, \dots, \theta_m$ and α
- Define combined kernel $k_\Sigma(\cdot, \cdot) = \sum_{t=1}^m \theta_t k_t(\cdot, \cdot)$
- Predict with $y(\mathbf{x} | \tilde{\mathbf{K}}_\Sigma, \alpha) = \sum_{i=1}^n \alpha_i k_\Sigma(\mathbf{x}, \mathbf{x}_i)$

Our Approach

- Binary Classification: $y_i \in \{-1, +1\}$
- α is optimal in a Support Vector Machine
- Control generalization error of k_Σ with the error of k_1, \dots, k_m

Optimization-Based Results

Lemma of One Kernel

Let α be the dual-optimal vector for labeled kernel matrix $\tilde{\mathbf{K}}$. Then, by combining the *Stationarity*, *Complementary Slackness*, and *Dual Feasibility* KKT conditions, we find that

$$\|\alpha\|_1 = \alpha^\top \tilde{\mathbf{K}} \alpha$$

Theorem of Two Kernels: Adding Kernels Reduces Complexity

Let α_1 and α_2 be the dual-optimal vectors for labeled kernel matrices $\tilde{\mathbf{K}}_1$ and $\tilde{\mathbf{K}}_2$. Let α_{1+2} be the dual-optimal vector for labeled kernel matrix $\tilde{\mathbf{K}}_{1+2} := \tilde{\mathbf{K}}_1 + \tilde{\mathbf{K}}_2$. Then, following from the *prior lemma*, the *optimality* of α_{1+2} , and some algebra, we have

$$\alpha_{1+2}^\top \tilde{\mathbf{K}}_{1+2} \alpha_{1+2} \leq \frac{1}{3} (\alpha_1^\top \tilde{\mathbf{K}}_1 \alpha_1 + \alpha_2^\top \tilde{\mathbf{K}}_2 \alpha_2)$$

Theorem of Many Kernels: Adding Many Kernels Greatly Reduces Complexity

Let $\alpha_1, \dots, \alpha_m$ be the dual-optimal vectors for labeled kernel matrices $\tilde{\mathbf{K}}_1, \dots, \tilde{\mathbf{K}}_m$. Let α_Σ be the dual optimal vector for labeled kernel matrix $\tilde{\mathbf{K}}_\Sigma := \sum_{t=1}^m \tilde{\mathbf{K}}_t$. Then, by repeatedly applying the *prior lemma*, we find

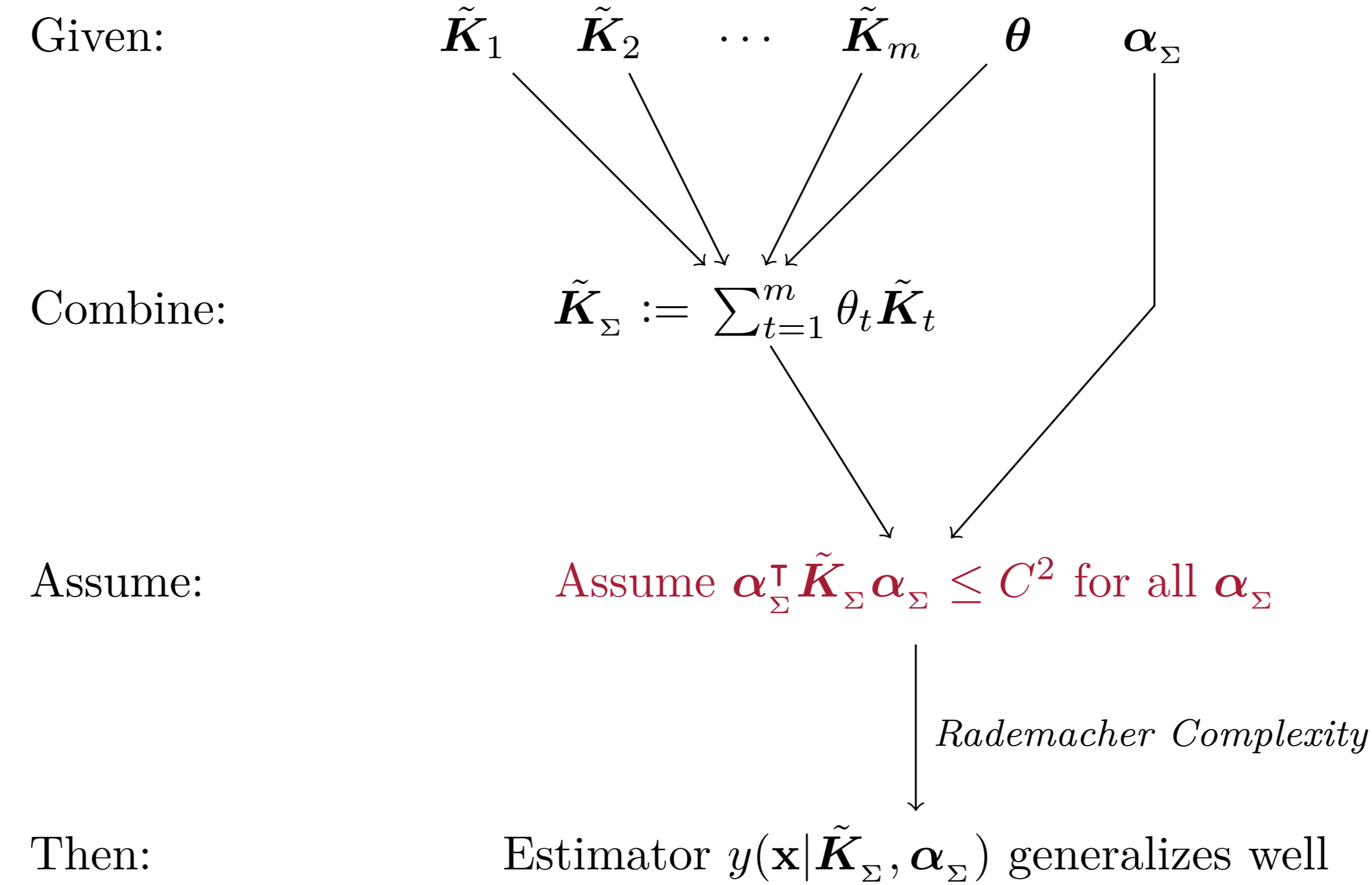
$$\alpha_\Sigma^\top \tilde{\mathbf{K}}_\Sigma \alpha_\Sigma \leq 3m^{-\log_2(3)} \sum_{t=1}^m \alpha_t^\top \tilde{\mathbf{K}}_t \alpha_t$$

Furthermore, if we assume that $\alpha_t^\top \tilde{\mathbf{K}}_t \alpha_t \leq B^2$, then

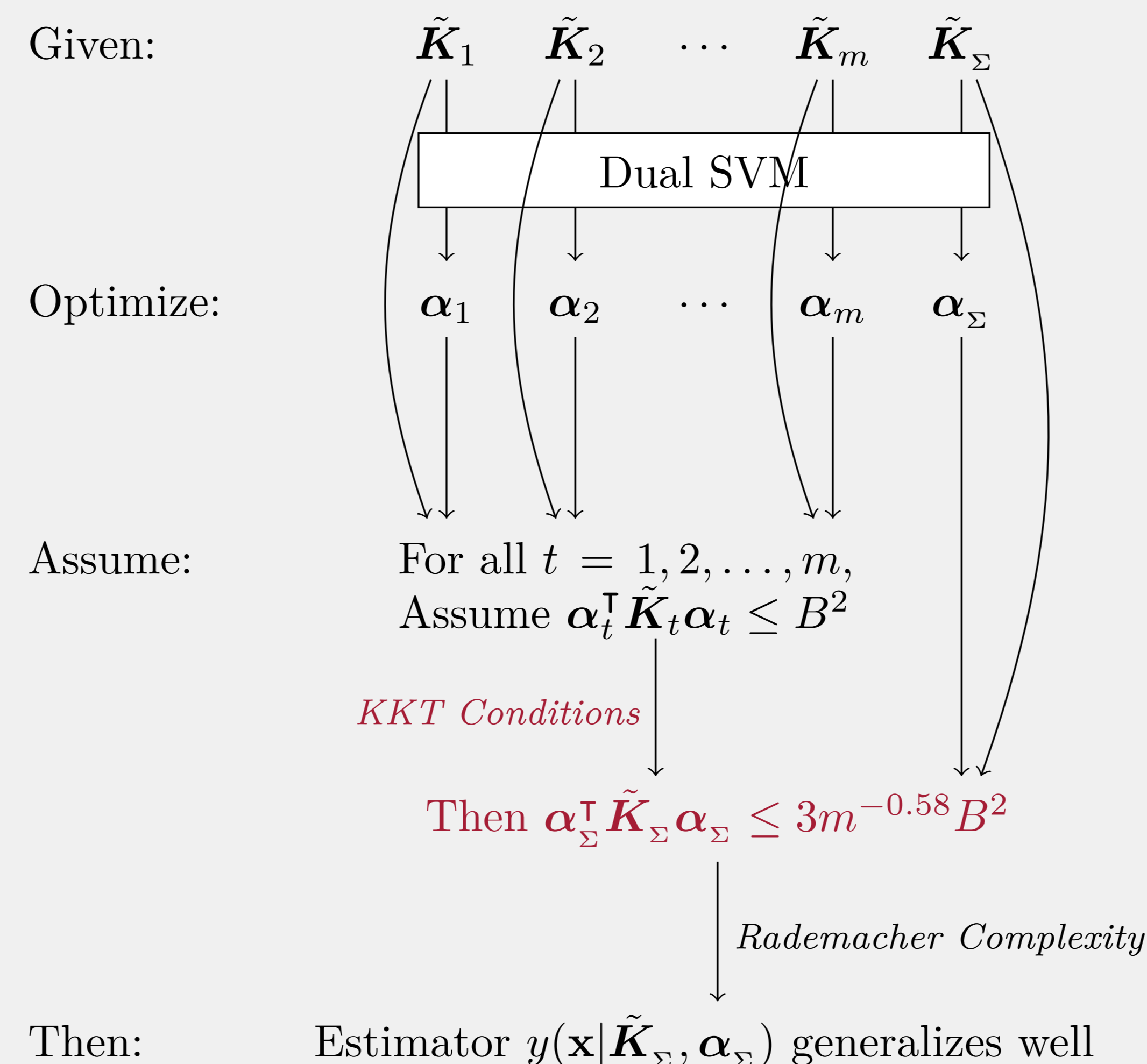
$$\alpha_\Sigma^\top \tilde{\mathbf{K}}_\Sigma \alpha_\Sigma \leq 3m^{-\log_2(3/2)} B^2$$

Main Theorem and Context

Template of Prior Works



Our Optimality Assumption



Conclusions

Learning Theory Results

Support Vector Machines Styles

- We consider the standard SVM and a L2-penalized SVM for nonseparable data:

$$\min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|_2^2 + \frac{C}{2} \|\xi\|_2^2 \quad \max_{\alpha, \xi} \|\alpha\|_1 - \frac{1}{2} \alpha^\top \tilde{\mathbf{K}} \alpha - \frac{1}{2} \|\xi\|_2^2$$

s.t. $y_i \mathbf{w}^\top \phi(\mathbf{x}_i) \geq 1 - \xi_i \quad \forall i \in [n]$ s.t. $0 \leq \alpha_i \leq C \xi_i \quad \forall i \in [n]$
 $\xi_i \geq 0 \quad \forall i \in [n]$

(a) Primal SVM Problem

(b) Dual SVM Problem

Figure 1. Primal and Dual SVM Problems. The L2 penalties are in gray.

- We prove statistical efficiency for standard SVM and $C = \frac{1}{2}$ in the L2-SVM

Ways to Combine Kernels Together

- Our core theorem complements existing Rademacher Complexity proofs
- Generalization error is bounded by the Rademacher Complexity $\hat{\mathcal{R}}(\mathcal{F})$:

$$\hat{\mathcal{R}}(\mathcal{F}) := \mathbb{E}_{\sigma \sim \{\pm 1\}^n} \left[\sup_{h \in \mathcal{F}} \left(\frac{1}{n} \sum_{i=1}^n \sigma_i h(\mathbf{x}_i) \right) \right]$$

- Different proofs consider different ways to combine kernels:

- Kernel Sums: If all $\theta_t = 1$, then

$$\hat{\mathcal{R}}(\mathcal{F}) = O\left(\frac{BRm^{0.208}}{\sqrt{n}}\right)$$

- Kernel Subsets: If all $\theta_t \in \{0, 1\}$, then

$$\hat{\mathcal{R}}(\mathcal{F}) = O\left(\frac{BRm^{0.208} \cdot \sqrt{\ln(m)}}{\sqrt{n}}\right)$$

- Convex Combinations*: If we have $\theta_t \in \{0\} \cup \left[\frac{10}{m}, 1\right]$ and $\sum_{t=1}^m \theta_t = 1$, then

$$\hat{\mathcal{R}}(\mathcal{F}) = O\left(\frac{BRm \sqrt{\ln(m)}}{\sqrt{n}}\right)$$

Table of Constants

Variable	Meaning
n	Number of Samples
i, j	Index of a Sample
m	Number of Kernels
t	Index of a Kernel
$\tilde{\mathbf{K}}$	Labeled Kernel Matrix (i.e. $[\tilde{\mathbf{K}}]_{i,j} := y_i y_j k(\mathbf{x}_i, \mathbf{x}_j)$)
α_t	Dual Solution Vector for SVM with $\tilde{\mathbf{K}}_t$
B^2	Upper Bound for all $\alpha_t^\top \tilde{\mathbf{K}}_t \alpha_t$
R^2	Upper Bound for all $k_t(\mathbf{x}_i, \mathbf{x}_i) = \ \phi(\mathbf{x}_i)\ _2^2$