

A. Proof of theorem 1

Proof. Let \tilde{H} be the family of hypotheses mapping $\mathcal{X} \times \mathcal{Y}$ to \mathbb{R} defined by $\tilde{H} = \{z = (x, y) \mapsto \rho_h(x, y) : h \in H\}$. Consider the family of functions $\tilde{\mathcal{H}} = \{\Phi_\rho \circ r : r \in \tilde{H}\}$ derived from \tilde{H} , where Φ_ρ is the ρ -margin function loss defined by $\Phi_\rho(x) = 1_{x \leq 0} + \max(0, 1 - x/\rho) 1_{x > 0}$. By the Rademacher complexity bound for functions taking values in $[0, 1]$ (see (Koltchinskii & Panchenko, 2002; Bartlett & Mendelson, 2002)), we can write that with probability at least $1 - \delta$, for all $h \in H$,

$$\mathbb{E} [\Phi_\rho(\rho_h(x, y))] \leq \hat{R}_\rho(h) + 2\hat{\mathfrak{R}}_S(\Phi_\rho \circ \tilde{H}) + 3\sqrt{\frac{\log \frac{2}{\delta}}{2m}}.$$

Since $1_{u \leq 0} \leq \Phi_\rho(u)$ for all $u \in \mathbb{R}$, the generalization error $R(h)$ is a lower bound on the left-hand side, $R(h) = \mathbb{E}[1_{y[h(x') - h(x)] \leq 0}] \leq \mathbb{E}[\Phi_\rho(\rho_h(x, y))]$. Furthermore, since Φ_ρ is $1/\rho$ -Lipschitz, by Talagrand's contraction lemma (Ledoux & Talagrand, 1991), we have $\hat{\mathfrak{R}}_S(\Phi_\rho \circ \tilde{H}) \leq \frac{1}{\rho} \hat{\mathfrak{R}}_S(\tilde{H})$.

For any fixed $y \in \mathcal{Y}$ and any $i \in [1, m]$, define ϵ_i as $2(1_{y=y_i}) - 1$. Since $\epsilon_i \in \{-1, +1\}$, the random variables σ_i and $\sigma_i \epsilon_i$ follow the same distribution. Using this fact and the sub-additivity of sup, $\hat{\mathfrak{R}}_S(\tilde{H})$ can be upper bounded as follows:

$$\begin{aligned} \hat{\mathfrak{R}}_S(\tilde{H}) &= \frac{1}{m} \mathbb{E} \left[\sup_{h \in H} \sum_{i=1}^m \sigma_i \rho_h(x_i, y_i) \right] \\ &= \frac{1}{m} \mathbb{E} \left[\sup_{h \in H} \sum_{i=1}^m \sum_{y \in \mathcal{Y}} \sigma_i \rho_h(x_i, y) 1_{y=y_i} \right] \\ &\leq \frac{1}{m} \sum_{y \in \mathcal{Y}} \mathbb{E} \left[\sup_{h \in H} \sum_{i=1}^m \sigma_i \rho_h(x_i, y) 1_{y=y_i} \right] \\ &= \frac{1}{m} \sum_{y \in \mathcal{Y}} \mathbb{E} \left[\sup_{h \in H} \sum_{i=1}^m \sigma_i \rho_h(x_i, y) \left(\frac{2(1_{y=y_i}) - 1}{2} + \frac{1}{2} \right) \right] \\ &\leq \frac{1}{2m} \sum_{y \in \mathcal{Y}} \mathbb{E} \left[\sup_{h \in H} \sum_{i=1}^m \sigma_i \epsilon_i \rho_h(x_i, y) \right] \\ &\quad + \frac{1}{2m} \sum_{y \in \mathcal{Y}} \mathbb{E} \left[\sup_{h \in H} \sum_{i=1}^m \sigma_i \rho_h(x_i, y) \right] \\ &= \frac{1}{m} \sum_{y \in \mathcal{Y}} \mathbb{E} \left[\sup_{h \in H} \sum_{i=1}^m \sigma_i \rho_h(x_i, y) \right]. \end{aligned}$$

Let $H_{\mathcal{X}}^{(c-1)} = \{\max\{h_1, \dots, h_l\} : h_i \in H_{\mathcal{X}}, i \in [1, c-1]\}$. It is known that the empirical Rademacher of a function class defined as that of the maxima of several hypotheses is upper bounded by the sum of the empirical Rademacher complexities of the sets to which each of these hypotheses belong to (Ledoux & Talagrand, 1991), thus $\hat{\mathfrak{R}}_S(H_{\mathcal{X}}^{(c-1)}) \leq (c-1)\hat{\mathfrak{R}}_S(H_{\mathcal{X}})$.

Now, rewriting $\rho_h(x_i, y)$ explicitly, using again the sub-additivity of sup and observing that $-\sigma_i$ and σ_i are distributed in the same way leads to

$$\begin{aligned} \hat{\mathfrak{R}}_S(\tilde{H}) &\leq \frac{1}{m} \sum_{y \in \mathcal{Y}} \mathbb{E} \left[\sup_{h \in H} \sum_{i=1}^m \sigma_i (h(x_i, y) - \max_{y' \neq y} h(x_i, y')) \right] \\ &\leq \sum_{y \in \mathcal{Y}} \left[\frac{1}{m} \mathbb{E} \left[\sup_{h \in H} \sum_{i=1}^m \sigma_i h(x_i, y) \right] \right. \\ &\quad \left. + \frac{1}{m} \mathbb{E} \left[\sup_{h \in H} \sum_{i=1}^m -\sigma_i \max_{y' \neq y} h(x_i, y') \right] \right] \\ &= \sum_{y \in \mathcal{Y}} \left[\frac{1}{m} \mathbb{E} \left[\sup_{h \in H} \sum_{i=1}^m \sigma_i h(x_i, y) \right] \right. \\ &\quad \left. + \frac{1}{m} \mathbb{E} \left[\sup_{h \in H} \sum_{i=1}^m \sigma_i \max_{y' \neq y} h(x_i, y') \right] \right] \\ &\leq \sum_{y \in \mathcal{Y}} \left[\frac{1}{m} \mathbb{E} \left[\sup_{h \in H_{\mathcal{X}}} \sum_{i=1}^m \sigma_i h(x_i) \right] \right. \\ &\quad \left. + \frac{1}{m} \mathbb{E} \left[\sup_{h \in H_{\mathcal{X}}^{(c-1)}} \sum_{i=1}^m \sigma_i h(x_i) \right] \right] \\ &\leq c \left[\frac{c}{m} \mathbb{E} \left[\sup_{h \in H_{\mathcal{X}}} \sum_{i=1}^m \sigma_i h(x_i) \right] \right] = c^2 \hat{\mathfrak{R}}_S(H_{\mathcal{X}}). \end{aligned}$$

This concludes the proof. \square

B. Proof of lemma 1

Proof. For any $h \in \mathbb{H}_K$ and $x \in \mathcal{X}$, by the reproducing property, we have $h(x) = \langle h, K(x, \cdot) \rangle$. Let $\mathbb{H}_S = \text{span}(\{K(x_i, \cdot) : i \in [1, m]\})$, then, for $i \in [1, m]$, $h(x_i) = \langle h', K(x, \cdot) \rangle$, where h' is the orthogonal projection of h over \mathbb{H}_S . Thus, there exists $\alpha = (\alpha_1, \dots, \alpha_m)^\top \in \mathbb{R}^m$ such that $h' = \sum_{i=1}^m \alpha_i K(x_i, \cdot)$. If $\|h\|_{\mathbb{H}} \leq \Lambda$, then $\alpha^\top \mathbf{K} \alpha = \|h'\|_K^2 \leq \|h\|_K^2 \leq \Lambda^2$ where \mathbf{K} is the kernel matrix of K for the sample S . Conversely, any $\sum_{i=1}^m \alpha_i K(x_i, \cdot)$ with $\alpha^\top \mathbf{K} \alpha \leq \Lambda^2$ is the projection of some $h \in \mathbb{H}_K$ with $\|h\|_K^2 \leq \Lambda^2$.

Thus, for any $y \in \mathcal{Y}$, there exists $\alpha^y = (\alpha_1^y, \dots, \alpha_m^y)^\top \in \mathbb{R}^m$ such that for any $i \in [1, m]$, $h_y(x_i) = \sum_{j=1}^m \alpha_j^y K_\mu(x_i, x_j)$, and $\alpha^{y^\top} \mathbf{K}_\mu \alpha^y \leq \Lambda$ where \mathbf{K}_μ is the kernel matrix associated to K_μ for the sample (x_1, \dots, x_m) . In view of that, we can write

$$\begin{aligned} \hat{\mathfrak{R}}_S(H_{\mathcal{X}}^1) &= \frac{1}{m} \mathbb{E} \left[\sup_{h_y \in H_{\mathcal{X}}^1} \sum_{i=1}^m \sigma_i h_y(x_i) \right] \\ &= \frac{1}{m} \mathbb{E} \left[\sup_{\mu \in \mathcal{M}_1, y \in \mathcal{Y}, \alpha^y} \sigma \mathbf{K}_\mu \alpha^y \right]. \end{aligned}$$

Now, by the Cauchy-Schwarz inequality, the supremum $\sup_{\alpha^y} \sigma^\top \mathbf{K}_\mu \alpha^y$ is reached for $\mathbf{K}_\mu^{-1/2} \alpha^y$ collinear

with $\mathbf{K}_\mu^{1/2}\boldsymbol{\sigma}$, which gives $\sup_{\boldsymbol{\alpha}^y \in \mathcal{A}} \boldsymbol{\sigma}^\top \mathbf{K}_\mu \boldsymbol{\alpha}^y = \Lambda \sqrt{\boldsymbol{\sigma}^\top \mathbf{K}_\mu \boldsymbol{\sigma}}$. Thus,

$$\begin{aligned} \widehat{\mathfrak{R}}_S(H_{\mathcal{X}}^1) &= \frac{\Lambda}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{\boldsymbol{\mu} \in \mathcal{M}_1} \sqrt{\boldsymbol{\sigma}^\top \mathbf{K}_\mu \boldsymbol{\sigma}} \right] \\ &= \frac{\Lambda}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{\boldsymbol{\mu} \in \mathcal{M}_1} \sqrt{\boldsymbol{\mu} \cdot \mathbf{u}_{\boldsymbol{\sigma}}} \right], \end{aligned}$$

which concludes the proof. \square

C. Proof of lemma 3

Proof. It first helps to rewrite the optimization in the following equivalent form:

$$\min_{\lambda \geq 0, t} t \quad \text{subject to: } \forall k \in [1, p], t \geq \mathbf{u}_{\boldsymbol{\sigma}, k} + \lambda(\tilde{\gamma}_k - \gamma_0).$$

The Lagrangian L associated to this problem can be defined for any $t \in \mathbb{R}$, $\lambda \geq 0$, and $\boldsymbol{\beta} \in \mathbb{R}^p$, $\boldsymbol{\beta} \geq 0$, by

$$L(t, \lambda, \boldsymbol{\beta}) = t + \sum_{k=1}^p \beta_k (\mathbf{u}_{\boldsymbol{\sigma}, k} + \lambda(\tilde{\gamma}_k - \gamma_0) - t).$$

By the KKT conditions, the following holds:

$$\frac{\partial L}{\partial t} = 1 - \sum_{k=1}^p \beta_k^* = 0 \iff \sum_{k=1}^p \beta_k^* = 1, \quad (10)$$

$$\frac{\partial L}{\partial \lambda} = \sum_{k=1}^p \beta_k^* (\tilde{\gamma}_k - \gamma_0) = 0, \quad (11)$$

$$\forall k : \beta_k^* > 0, \mathbf{u}_{\boldsymbol{\sigma}, k} + \lambda^* (\tilde{\gamma}_k - \gamma_0) = t^*. \quad (12)$$

Here λ^* , t^* , and $\boldsymbol{\beta}^*$ denote the optimal values for the dual problem $\max_{\boldsymbol{\beta} \geq 0} \min_{\lambda, t} L(t, \lambda, \boldsymbol{\beta})$. Note that Slater's condition holds for the convex primal problem, which implies that strong duality holds and that λ^* and t^* are also optimal solutions to the primal problem. The conditions (10) and (11) follow from the fact that $\nabla L = \mathbf{0}$ at the optimum and (12) follows from the complementary slackness condition that guarantees for all k , $\beta_k^* (\mathbf{u}_{\boldsymbol{\sigma}, k} + \lambda^* (\tilde{\gamma}_k - \gamma_0) - t^*) = 0$. Note that the condition in (10) as well as the constraint $\beta_k \geq 0$, which is imposed on dual variables that correspond to inequality constraints, implies $\boldsymbol{\beta}^* \in \Delta_1$.

We first consider the subset of solutions to the dual optimization that are equal to $T_1 = \max_{k: \tilde{\gamma}_k \geq \gamma_0} \mathbf{u}_{\boldsymbol{\sigma}, k}$. Note that any feasible point in this simpler optimization is also a feasible point in the original problem. Now, consider the dual optimization, which can be simplified by removing t when we impose $\boldsymbol{\beta} \in \Delta_1$: $\max_{\boldsymbol{\beta} \in \Delta_1} \min_{\lambda} \sum_{k=1}^p \beta_k (\mathbf{u}_{\boldsymbol{\sigma}, k} + \lambda(\tilde{\gamma}_k - \gamma_0))$. Whenever it is the case that $\lambda^* = 0$ or that $\tilde{\gamma}_k = \gamma_0$ for all k where $\beta_k > 0$, we can further simplify

the objective to $\max_{\boldsymbol{\beta} \in \Delta_1} \sum_{k=1}^p \beta_k \mathbf{u}_{\boldsymbol{\sigma}, k} = \max_k \mathbf{u}_{\boldsymbol{\sigma}, k}$. Note this implies that $\lambda^* = 0$ iff $\tilde{\gamma}_{k_{\max}} \geq \gamma_0$ (where $k_{\max} = \operatorname{argmax}_k \mathbf{u}_{\boldsymbol{\sigma}, k}$), since all constraints of the original problem in lemma 2 must be satisfied at the optimum. Thus, T_1 is found as the solution to the dual problem whenever $\lambda^* = 0$ or $\tilde{\gamma}_k = \gamma_0$ for all k where $\beta_k > 0$.

Now we seek an expression T_2 that is equal to the optimum of the dual optimization in the cases not accounted for by T_1 . That is, we consider the case $\lambda^* > 0$ and the existence of at least one k such that $\beta_k^* > 0$ and $\tilde{\gamma}_k \neq \gamma_0$. In order for condition (11) to be satisfied in this case, there must be at least two coordinates β_k^* and $\beta_{k'}^*$ that are non-zero for k and k' that satisfy $\tilde{\gamma}_k < \gamma_0 < \tilde{\gamma}_{k'}$. This is because both a negative coefficient, i.e. $(\tilde{\gamma}_k - \gamma_0)$, as well as a positive coefficient, i.e. $(\tilde{\gamma}_{k'} - \gamma_0)$, must be present in order for 0 to be found as convex combination. Now, fix any two coordinates k and k' that are non-zero in $\boldsymbol{\beta}^*$ and that satisfy $\tilde{\gamma}_k \leq \gamma_0 \leq \tilde{\gamma}_{k'}$ with $\tilde{\gamma}_k \neq \tilde{\gamma}_{k'}$ (there exists *at least* two such coordinates by the argument just discussed). From condition (12) we know that

$$\begin{aligned} \mathbf{u}_{\boldsymbol{\sigma}, k} + \lambda^* (\tilde{\gamma}_k - \gamma_0) &= \mathbf{u}_{\boldsymbol{\sigma}, k'} + \lambda^* (\tilde{\gamma}_{k'} - \gamma_0) \\ \iff \lambda^* &= \frac{\mathbf{u}_{\boldsymbol{\sigma}, k} - \mathbf{u}_{\boldsymbol{\sigma}, k'}}{\tilde{\gamma}_k - \tilde{\gamma}_{k'}} > 0. \end{aligned} \quad (13)$$

Plugging this back into (12) we find an expression for the optimal objective value t^* :

$$t^* = \alpha_{k, k'} \mathbf{u}_{\boldsymbol{\sigma}, k} + (1 - \alpha_{k, k'}) \mathbf{u}_{\boldsymbol{\sigma}, k'},$$

where $0 \leq \alpha_{k, k'} = \frac{\tilde{\gamma}_{k'} - \gamma_0}{\tilde{\gamma}_{k'} - \tilde{\gamma}_k} \leq 1$. However, we still do not know which k and k' are active at the optimum. We do know that for all k , $t^* \geq \mathbf{u}_{\boldsymbol{\sigma}, k} + \lambda^* (\tilde{\gamma}_k - \gamma_0)$ since all constraints must hold at the optimum, which also implies, for all k and k' ,

$$\begin{aligned} t^* &\geq \alpha_{k, k'} (\mathbf{u}_{\boldsymbol{\sigma}, k} + \lambda^* (\tilde{\gamma}_k - \gamma_0)) \\ &\quad + (1 - \alpha_{k, k'}) (\mathbf{u}_{\boldsymbol{\sigma}, k'} + \lambda^* (\tilde{\gamma}_{k'} - \gamma_0)) \\ &= \alpha_{k, k'} \mathbf{u}_{\boldsymbol{\sigma}, k} + (1 - \alpha_{k, k'}) \mathbf{u}_{\boldsymbol{\sigma}, k'} \\ &\quad + \lambda^* \underbrace{(\alpha_{k, k'} (\tilde{\gamma}_k - \gamma_0) + (1 - \alpha_{k, k'}) (\tilde{\gamma}_{k'} - \gamma_0))}_{= 0}. \end{aligned}$$

Thus, we can maximize over all feasible choice of k and k' in order to find the value at which the above inequality is tight:

$$t^* = \max_{(k, k') \in J_p} \alpha_{k, k'} \mathbf{u}_{\boldsymbol{\sigma}, k} + (1 - \alpha_{k, k'}) \mathbf{u}_{\boldsymbol{\sigma}, k'},$$

which gives us the expression T_2 for the optimum in the intersection case. Finally, taking the maximum over T_1 and T_2 completes the lemma. \square

D. Theorem 3

Theorem 3. Fix $\rho > 0$ and let $p' = \text{Card}(I_p) \leq p$ and $p'' = \text{Card}(J_p) < p^2$. Then, for any $\delta > 0$, with probability at least $1 - \delta$ over a sample of size m , the following multi-class classification generalization bound holds for all $h \in H^1$:

$$R(h) \leq \widehat{R}_\rho(h) + \frac{2c^2\Lambda}{m\rho} \sqrt{T_{\gamma_0} + m\lambda_{\max} \sqrt{\frac{\log(p' + p'')}{2}}} + 3\sqrt{\frac{\log \frac{2}{\delta}}{2m}}, \quad (14)$$

where $T_{\gamma_0} = \max(\max_{k \in I_p} \text{Tr}[\mathbf{K}_k], \max_{(k,k') \in J_p} \text{Tr}[\mathbf{K}_{k,k'}])$ and $\lambda_{\max} = \max(\max_{k \in I_p} \|\mathbf{K}_k\|_2, \max_{(k,k') \in J_p} \|\mathbf{K}_{k,k'}\|_2)$, with $\mathbf{K}_{k,k'} = \alpha_{k,k'} \mathbf{K}_k + (1 - \alpha_{k,k'}) \mathbf{K}_{k'}$ and $\alpha_{k,k'} = \frac{\tilde{\gamma}_{k'} - \gamma_0}{\tilde{\gamma}_{k'} - \tilde{\gamma}_k}$.

Proof. Let $M_p = \{(k, k') \in [0, p] \times [1, p] : (\tilde{\gamma}_k \leq \gamma_0 \leq \tilde{\gamma}_{k'}) \wedge (\tilde{\gamma}_k \neq \tilde{\gamma}_{k'})\}$ as in the proof of Theorem 2. By lemmas 1-3 and Jensen's inequality, we can write:

$$\begin{aligned} \widehat{\mathfrak{R}}_S(H_\chi^1) &\leq \frac{\Lambda}{m} \mathbb{E} \left[\sqrt{\max_{(k,k') \in M_p} \sigma^\top \mathbf{K}_{k,k'} \sigma} \right] \\ &\leq \frac{\Lambda}{m} \sqrt{\mathbb{E} \left[\max_{(k,k') \in M_p} \sigma^\top \mathbf{K}_{k,k'} \sigma \right]}. \end{aligned}$$

Note that for any k, k' , $\sigma^\top \mathbf{K}_{k,k'} \sigma \leq \|\sigma\|^2 \|\mathbf{K}_{k,k'}\|_2 = m \|\mathbf{K}_{k,k'}\|_2$. Now, for any $t \in \mathbb{R}$, by the convexity of exp and Jensen's inequality, we have

$$\begin{aligned} e^{t \mathbb{E}[\max_{(k,k') \in M_p} \sigma^\top \mathbf{K}_{k,k'} \sigma]} &\leq \mathbb{E}[e^{t \max_{(k,k') \in M_p} \sigma^\top \mathbf{K}_{k,k'} \sigma}] \\ &= \mathbb{E}[\max_{(k,k') \in M_p} e^{t \sigma^\top \mathbf{K}_{k,k'} \sigma}] \\ &\leq \mathbb{E}[\sum_{(k,k') \in M_p} e^{t \sigma^\top \mathbf{K}_{k,k'} \sigma}]. \end{aligned}$$

For any $(k, k') \in M_p$, we have $\mathbb{E}[\sigma^\top \mathbf{K}_{k,k'} \sigma] = \text{Tr}[\mathbf{K}_{k,k'}]$. Thus, by Hoeffding's inequality, the following holds

$$\begin{aligned} \mathbb{E}[e^{t \sigma^\top \mathbf{K}_{k,k'} \sigma}] &= e^{t \text{Tr}[\mathbf{K}_{k,k'}]} \mathbb{E}[e^{t(\sigma^\top \mathbf{K}_{k,k'} \sigma - \text{Tr}[\mathbf{K}_{k,k'}])}] \\ &\leq e^{t \text{Tr}[\mathbf{K}_{k,k'}]} e^{t^2 \lambda_{\max}^2 m^2 / 8}. \end{aligned}$$

Therefore, we can write

$$\begin{aligned} e^{t \mathbb{E}[\max_{(k,k') \in M_p} \sigma^\top \mathbf{K}_{k,k'} \sigma]} &\leq (p' + p'') e^{t \max_{k,k'} \text{Tr}[\mathbf{K}_{k,k'}]} e^{t^2 \lambda_{\max}^2 m^2 / 8}. \end{aligned}$$

Taking the log of both sides and rearranging gives

$$\begin{aligned} &\mathbb{E}[\max_{(k,k') \in M_p} \sigma^\top \mathbf{K}_{k,k'} \sigma] \\ &\leq \max_{(k,k') \in M_p} \text{Tr}[\mathbf{K}_{k,k'}] + \frac{\log(p' + p'')}{t} + t \lambda_{\max}^2 m^2 / 8. \end{aligned}$$

Choosing $t = \sqrt{8(\log(p' + p'')) / (\lambda_{\max}^2 m^2)}$ to minimize the upper bound gives

$$\begin{aligned} &\mathbb{E}[\max_{(k,k') \in M_p} \sigma^\top \mathbf{K}_{k,k'} \sigma] \\ &\leq \max_{(k,k') \in M_p} \text{Tr}[\mathbf{K}_{k,k'}] + m \lambda_{\max} \sqrt{\frac{\log(p' + p'')}{2}}. \end{aligned}$$

Plugging in this upper bound on the Rademacher complexity of H_χ^1 in the learning guarantee of theorem 1 concludes the proof. \square

E. Proof of lemma 4

Proof. We introduce $M_1(\gamma_0) = \max_{k \in I_p(\gamma_0)} \text{Tr}[\mathbf{K}_k]$ and $M'_1(\gamma_0) = \max_{(k,k') \in J_p(\gamma_0)} \text{Tr}[\mathbf{K}_{k,k'}]$, where we write $I_p(\gamma_0)$ and $J_p(\gamma_0)$ to make the dependency of these sets on γ_0 explicit. With these definitions, we can write $T_{\gamma_0} = \max(M_1(\gamma_0), M'_1(\gamma_0))$.

Notice that if $\text{Tr}[K_k] \leq \text{Tr}[K_{k'}]$, then $\text{Tr}[K_{kk'}] \leq \text{Tr}[K_{k'}]$ since $\text{Tr}[K_{kk'}]$ is a convex combination of $\text{Tr}[K_k]$ and $\text{Tr}[K_{k'}]$. Thus, if (k, k') is the maximizing pair of indices defining $M'_1(\gamma_0)$ and $\text{Tr}[K_k] \leq \text{Tr}[K_{k'}]$, we have $M'_1(\gamma_0) = \text{Tr}[K_{kk'}] \leq \text{Tr}[K_{k'}] \leq \max_{\tilde{\gamma}_{k'} \geq \gamma_0} \text{Tr}[K_{k'}] = M_1(\gamma_0)$. By contraposition, if $M_1(\gamma_0) < M'_1(\gamma_0)$, then $\text{Tr}[K_{k'}] < \text{Tr}[K_k]$. In view of that, we can rewrite $T_{\gamma_0} = \max(M_1(\gamma_0), M_2(\gamma_0))$, where $M_2(\gamma_0) = \max_{(k,k') \in N_p(\gamma_0)} \text{Tr}[\mathbf{K}_{k,k'}]$, where $N_p(\gamma_0) = \{(k, k') \in [1, p]^2 : (\tilde{\gamma}_k \leq \gamma_0 \leq \tilde{\gamma}_{k'}) \wedge (\tilde{\gamma}_k \neq \tilde{\gamma}_{k'}) \wedge (\text{Tr}[K_{k'}] < \text{Tr}[K_k])\}$.

Now, let $\lambda_0 \geq \gamma_0$, we will show that $T_{\lambda_0} \leq T_{\gamma_0}$. First note that if $T_{\lambda_0} = M_1(\lambda_0)$, then, since by definition of M_1 , $M_1(\lambda_0) \leq M_1(\gamma_0) \leq T_{\gamma_0}$, this shows immediately that $T_{\lambda_0} \leq T_{\gamma_0}$.

Otherwise, $T_{\lambda_0} = M_2(\lambda_0)$. Let (l, l') be the maximizing indices in the definition of $M_2(\lambda_0)$. Then, since $\text{Tr}[K_{l'}] - \text{Tr}[K_l] < 0$ holds, we can write:

$$\begin{aligned} T_{\lambda_0} &= \frac{\tilde{\gamma}_{l'} - \lambda_0}{\tilde{\gamma}_{l'} - \tilde{\gamma}_l} \text{Tr}[K_l] + \frac{\lambda_0 - \tilde{\gamma}_l}{\tilde{\gamma}_{l'} - \tilde{\gamma}_l} \text{Tr}[K_{l'}] \\ &= \frac{\tilde{\gamma}_{l'} \text{Tr}[K_l] - \tilde{\gamma}_l \text{Tr}[K_{l'}]}{\tilde{\gamma}_{l'} - \tilde{\gamma}_l} + \frac{\lambda_0 (\text{Tr}[K_{l'}] - \text{Tr}[K_l])}{\tilde{\gamma}_{l'} - \tilde{\gamma}_l} \\ &\leq \frac{\tilde{\gamma}_{l'} \text{Tr}[K_l] - \tilde{\gamma}_l \text{Tr}[K_{l'}]}{\tilde{\gamma}_{l'} - \tilde{\gamma}_l} + \frac{\gamma_0 (\text{Tr}[K_{l'}] - \text{Tr}[K_l])}{\tilde{\gamma}_{l'} - \tilde{\gamma}_l} \\ &= \frac{\tilde{\gamma}_{l'} - \gamma_0}{\tilde{\gamma}_{l'} - \tilde{\gamma}_l} \text{Tr}[K_l] + \frac{\gamma_0 - \tilde{\gamma}_l}{\tilde{\gamma}_{l'} - \tilde{\gamma}_l} \text{Tr}[K_{l'}]. \end{aligned}$$

If $(l, l') \in N_p(\gamma_0)$, the right-hand side is upper bounded by $M_2(\gamma_0) \leq T_{\gamma_0}$. Otherwise, if (l, l') is not in $N_p(\gamma_0)$, since by definition of $N_p(\gamma_0)$, $\tilde{\gamma}_l \geq \lambda_0 \geq \gamma_0$, this can only be because $\gamma_0 \leq \tilde{\gamma}_l$. But in that case both $\tilde{\gamma}_l$ and $\tilde{\gamma}_{l'}$ are greater than or equal to γ_0 . Then, by definition of $M_1(\gamma_0)$, $\text{Tr}[K_{ll'}] \leq \max(\text{Tr}[K_l], \text{Tr}[K_{l'}]) \leq$

$M_1(\gamma_0) \leq T_{\gamma_1}$. Thus, the inequality $T_{\lambda_0} \leq T_{\gamma_0}$ holds in all cases. \square

F. Alternative M³K optimization

Here we present an alternative formulation of the M³K algorithm, which results in a quadratically constrained linear program. Such a problem can be solved with any standard second order cone programming (SOCP) solver. We find this formulation can be faster to solve than the SILP formulation for smaller size problems, especially in the case of fewer classes.

First, consider the dual formulation with the margin constraint appearing instead as an additional penalty term in the objective with regularization parameter Γ . For every choice of γ_0 in the constraint version of the optimization there exists a choice of Γ that results in an equivalent optimization problem.

$$\begin{aligned} \min_{\boldsymbol{\mu}, \boldsymbol{\gamma}} \max_{\boldsymbol{\alpha}} \quad & \sum_{i=1}^m \boldsymbol{\alpha}_i^\top \mathbf{e}_{y_i} - \frac{C}{2} \sum_{i,j=1}^m \boldsymbol{\alpha}_i^\top \boldsymbol{\alpha}_j \sum_{k=1}^p \mu_k K_k(x_i, x_j) \\ & - \Gamma \sum_{i=1}^m \gamma_i \\ \text{s.t. } \forall i, \quad & \boldsymbol{\alpha}_i \leq \mathbf{e}_{y_i}, \quad \boldsymbol{\alpha}_i^\top \mathbf{1} = 0 \\ & \forall i, \forall y \neq y_i, \quad \gamma_i \leq \boldsymbol{\mu}^\top \boldsymbol{\eta}(x_i, y_i, y) \\ & \boldsymbol{\mu} \geq \mathbf{0}, \quad \boldsymbol{\mu}^\top \mathbf{1} = \Lambda. \end{aligned}$$

The objective is linear in $\boldsymbol{\mu}$ and $\boldsymbol{\gamma}$ and concave in $\boldsymbol{\alpha}$ and both $\boldsymbol{\mu}$ and $\boldsymbol{\alpha}$ are drawn from convex compact sets. The $\boldsymbol{\gamma}$ is drawn from a closed convex set, however, it is unbounded from below. One can add a lower bound on $\boldsymbol{\gamma}$ that has no effect on the optimal solution in order to achieve compactness and, thus, the minimax theorem applies and we permute the min and the max and solve for $\boldsymbol{\mu}$ and $\boldsymbol{\gamma}$ first. Focusing on the terms that depend on $\boldsymbol{\mu}$ and $\boldsymbol{\gamma}$, we have:

$$\begin{aligned} \min_{\boldsymbol{\mu}, \boldsymbol{\gamma}} \quad & -\frac{C}{2} \boldsymbol{\mu}^\top \mathbf{u} - \Gamma \sum_{i=1}^m \gamma_i \\ \text{s.t. } \forall i, \forall y \neq y_i, \quad & \gamma_i \leq \boldsymbol{\mu}^\top \boldsymbol{\eta}(x_i, y_i, y) \\ & \boldsymbol{\mu} \geq \mathbf{0}, \quad \boldsymbol{\mu}^\top \mathbf{1} = \Lambda, \end{aligned}$$

where $u_k = \sum_{i,j=1}^m \boldsymbol{\alpha}_i^\top \boldsymbol{\alpha}_j K_k(x_i, x_j)$. Consider the partial Lagrangian for the constraints on $\boldsymbol{\gamma}$, which introduces dual variables $\boldsymbol{\beta} \geq \mathbf{0}$:

$$\begin{aligned} L = -\frac{C}{2} \boldsymbol{\mu}^\top \mathbf{u} - \Gamma \sum_{i=1}^m \gamma_i \\ + \sum_{i=1}^m \sum_{y \neq y_i} \beta_{i,y} (\gamma_i - \boldsymbol{\mu}^\top \boldsymbol{\eta}(x_i, y_i, y)). \end{aligned}$$

At the optimum, it is guaranteed that

$$\forall i, \quad \frac{\partial L}{\partial \gamma_i} = -\Gamma + \sum_{y \neq y_i} \beta_{i,y} = 0 \iff \sum_{y \neq y_i} \beta_{i,y} = \Gamma.$$

Enforcing this constraint explicitly gives the following simplified objective, where $\boldsymbol{\gamma}$ is no longer appears (since $\Gamma \sum_i \gamma_i - \sum_{i,y \neq y_i} \beta_{i,y} \gamma_i = 0$):

$$\begin{aligned} \min_{\boldsymbol{\mu}} \max_{\boldsymbol{\beta} \geq \mathbf{0}} \quad & -\frac{C}{2} \boldsymbol{\mu}^\top \mathbf{u} - \sum_{i=1}^m \sum_{y \neq y_i} \boldsymbol{\mu}^\top \beta_{i,y} \boldsymbol{\eta}(x_i, y_i, y) \\ \text{s.t. } \forall i, \quad & \sum_{y \neq y_i} \beta_{i,y} = \Gamma \\ & \boldsymbol{\mu} \geq \mathbf{0}, \quad \boldsymbol{\mu}^\top \mathbf{1} = \Lambda. \end{aligned}$$

Note the minimax theorem applies once again (the objective is linear in $\boldsymbol{\mu}$ and $\boldsymbol{\beta}$ which are both drawn from convex compact sets) and we can first solve the minimization over $\boldsymbol{\mu}$ for a fixed $\boldsymbol{\beta}$. Since the objective and the constraints are linear in $\boldsymbol{\mu}$ it is sufficient to consider a solution that places all of the possible mass on a single coordinate of $\boldsymbol{\mu}$. Thus, an equivalent objective is: $\min_{\boldsymbol{\mu}} \boldsymbol{\mu}^\top (-\frac{C}{2} \mathbf{u} - \sum_{i=1}^m \sum_{y \neq y_i} \beta_{i,y} \boldsymbol{\eta}(x_i, y_i, y)) = \min_k \Lambda (-\frac{C}{2} u_k - \sum_{i=1}^m \sum_{y \neq y_i} \beta_{i,y} \eta_k(x_i, y_i, y))$. A final reformulation, via the introduction of the variable t , simplifies the objective further:

$$\begin{aligned} \max_{\boldsymbol{\beta}, t} \quad & \Lambda t \\ \text{s.t. } \forall k, \quad & t \leq -\frac{C}{2} u_k - \sum_{i=1}^m \sum_{y \neq y_i} \beta_{i,y} \eta_k(x_i, y_i, y) \\ & \forall i, \quad \sum_{y \neq y_i} \beta_{i,y} = \Gamma, \quad \boldsymbol{\beta} \geq \mathbf{0}. \end{aligned}$$

Plugging this solution for $\boldsymbol{\mu}$ and $\boldsymbol{\gamma}$ back into the original optimization gives the final overall quadratically constraint linear program:

$$\begin{aligned} \max_{\boldsymbol{\alpha}, \boldsymbol{\beta}, t} \quad & \sum_{i=1}^m \boldsymbol{\alpha}_i^\top \mathbf{e}_{y_i} + \Lambda t \\ \text{s.t. } \forall k, \quad & t \leq -\frac{C}{2} \sum_{i,j=1}^m \boldsymbol{\alpha}_i^\top \boldsymbol{\alpha}_j K_k(x_i, x_j) \\ & - \sum_{i=1}^m \sum_{y \neq y_i} \beta_{i,y} \eta_k(x_i, y_i, y) \\ & \forall i, \quad \sum_{y \neq y_i} \beta_{i,y} = \Gamma, \quad \boldsymbol{\beta} \geq \mathbf{0} \\ & \forall i, \quad \boldsymbol{\alpha}_i \leq \mathbf{e}_{y_i}, \quad \boldsymbol{\alpha}_i^\top \mathbf{1} = 0. \end{aligned}$$

Table 4. Performance of several algorithms on the *caltech101* dataset for varying numbers of training points per class (PPC). The dataset consists of 102 classes and 48 kernels.

PPC	UNIF	BINMKL	OBSC	UFO	M ³ K
5	46.0 ± 0.9	54.0 ± 0.7	52.5 ± 0.6	47.9 ± 0.7	51.4 ± 1.2
10	57.9 ± 0.8	65.9 ± 0.9	65.1 ± 1.0	62.9 ± 0.7	66.0 ± 1.1
15	64.6 ± 1.0	71.8 ± 0.5	71.4 ± 0.6	70.6 ± 0.6	72.6 ± 0.9
20	68.4 ± 0.9	75.4 ± 1.1	75.7 ± 1.2	73.5 ± 1.1	76.0 ± 0.8
25	71.9 ± 1.6	77.1 ± 1.2	78.2 ± 0.8	75.5 ± 1.3	79.3 ± 1.0

G. Caltech 101 performance

Here we present the numerical values that are displayed in figure 2.