

A. Proof of Theorem 3

Although $f_{\mathcal{P}}^*$ is a feasible solution, it is not a local optimum for $\theta \in [0, 1]$ and $s \leq 0$ because

$$\alpha_i \leq C\theta \quad \text{for } i \in \tilde{\mathcal{I}} \cap \mathcal{O}, \quad (12a)$$

$$\alpha_i \geq C \quad \text{for } i \in \tilde{\mathcal{O}} \cap \mathcal{I}, \quad (12b)$$

violate the KKT conditions (7) for $\tilde{\mathcal{P}}$. These feasibility and sub-optimality indicates that

$$J_{\tilde{\mathcal{P}}}(f_{\tilde{\mathcal{P}}}^*; \theta) < J_{\mathcal{P}}(f_{\mathcal{P}}^*; \theta), \quad (13)$$

we arrive at (9).

Q.E.D.

B. Proof of Theorem 4

Sufficiency: If (10e) is true, i.e., if there are NO instances with $y_i f_{\mathcal{P}}^*(\mathbf{x}_i) = s$, then any convex problems defined by different partitions $\tilde{\mathcal{P}} \neq \mathcal{P}$ do not have feasible solutions in the neighborhood of $f_{\mathcal{P}}^*$. This means that if $f_{\mathcal{P}}^*$ is a conditionally optimal solution, then it is locally optimal. (10a)-(10d) are sufficient for $f_{\mathcal{P}}^*$ to be conditionally optimal for the given partition \mathcal{P} . Thus, (10) is sufficient for $f_{\mathcal{P}}^*$ to be locally optimal.

Necessity: From Theorem 3, if there exists an instance such that $y_i f_{\mathcal{P}}^*(\mathbf{x}_i) = s$, then $f_{\mathcal{P}}^*$ is a feasible but not locally optimal. Then (10e) is necessary for $f_{\mathcal{P}}^*$ to be locally optimal. In addition, (10a)-(10d) are also necessary for local optimality, because of every local optimal solutions are conditionally optimal for the given partition \mathcal{P} . Thus, (10) is necessary for $f_{\mathcal{P}}^*$ to be locally optimal.

Q.E.D.

C. Implementation of D-step

In D-step, we work with the following convex problem

$$f_{\tilde{\mathcal{P}}}^* := \operatorname{argmin}_{f \in \operatorname{pol}(\mathcal{P}; s)} J_{\tilde{\mathcal{P}}}(f; \theta). \quad (14)$$

where, $\tilde{\mathcal{P}}$ is updated from \mathcal{P} as (8).

Let us define a partition $\Pi := \{\mathcal{R}, \mathcal{E}, \mathcal{L}, \tilde{\mathcal{I}}', \tilde{\mathcal{O}}', \hat{\mathcal{O}}''\}$ of \mathbb{N}_n such that

$$i \in \mathcal{R} \Rightarrow y_i f(\mathbf{x}_i) > 1, \quad (15a)$$

$$i \in \mathcal{E} \Rightarrow y_i f(\mathbf{x}_i) = 1, \quad (15b)$$

$$i \in \mathcal{L} \Rightarrow s < y_i f(\mathbf{x}_i) < 1, \quad (15c)$$

$$i \in \tilde{\mathcal{I}}' \Rightarrow y_i f(\mathbf{x}_i) = s \text{ and } i \in \tilde{\mathcal{I}}, \quad (15d)$$

$$i \in \tilde{\mathcal{O}}' \Rightarrow y_i f(\mathbf{x}_i) = s \text{ and } i \in \tilde{\mathcal{O}}, \quad (15e)$$

$$i \in \hat{\mathcal{O}}'' \Rightarrow y_i f(\mathbf{x}_i) < s. \quad (15f)$$

If we write the conditionally optimal solution as

$$f_{\tilde{\mathcal{P}}}^*(x) := \sum_{j \in \mathbb{N}_n} \alpha_j^* y_j K(x, \mathbf{x}_j), \quad (16)$$

$\{\alpha_j^*\}_{j \in \mathbb{N}_n}$ must satisfy the following KKT conditions

$$y_i f_{\tilde{\mathcal{P}}}^*(\mathbf{x}_i) > 1 \Rightarrow \alpha_i^* = 0 \quad (17a)$$

$$y_i f_{\tilde{\mathcal{P}}}^*(\mathbf{x}_i) = 1 \Rightarrow \alpha_i^* \in [0, C], \quad (17b)$$

$$s < y_i f_{\tilde{\mathcal{P}}}^*(\mathbf{x}_i) < 1 \Rightarrow \alpha_i^* = C \quad (17c)$$

$$y_i f_{\tilde{\mathcal{P}}}^*(\mathbf{x}_i) = s, i \in \tilde{\mathcal{I}}' \Rightarrow \alpha_i^* \geq C, \quad (17d)$$

$$y_i f_{\tilde{\mathcal{P}}}^*(\mathbf{x}_i) = s, i \in \tilde{\mathcal{O}}' \Rightarrow \alpha_i^* \leq C\theta, \quad (17e)$$

$$y_i f_{\tilde{\mathcal{P}}}^*(\mathbf{x}_i) < s, i \in \hat{\mathcal{O}}'' \Rightarrow \alpha_i^* = C\theta. \quad (17f)$$

At the beginning of the D-step, $f_{\tilde{\mathcal{P}}}^*(\mathbf{x}_i)$ violates the KKT conditions by

$$\Delta f_i := y_i \begin{bmatrix} \mathbf{K}_{i, \Delta_{\mathcal{I} \rightarrow \mathcal{O}}} & \mathbf{K}_{i, \Delta_{\mathcal{O} \rightarrow \mathcal{I}}} \end{bmatrix} \begin{bmatrix} \alpha_{\Delta_{\mathcal{I} \rightarrow \mathcal{O}}}^{(\text{bef})} - \mathbf{1}C\theta \\ \alpha_{\Delta_{\mathcal{O} \rightarrow \mathcal{I}}}^{(\text{bef})} - \mathbf{1}C \end{bmatrix}.$$

where $\alpha^{(\text{bef})}$ is the corresponding α at the beginning of the D-step, while $\Delta_{\mathcal{I} \rightarrow \mathcal{O}}$ and $\Delta_{\mathcal{O} \rightarrow \mathcal{I}}$ denote the difference in $\tilde{\mathcal{P}}$ and \mathcal{P} defined as

$$\Delta_{\mathcal{I} \rightarrow \mathcal{O}} := \{i \in \mathcal{I} \mid y_i f_{\mathcal{P}}(\mathbf{x}_i) = s\},$$

$$\Delta_{\mathcal{O} \rightarrow \mathcal{I}} := \{i \in \mathcal{O} \mid y_i f_{\mathcal{P}}(\mathbf{x}_i) = s\}.$$

Then, we consider the following another parametrized problem with a parameter $\mu \in [0, 1]$:

$$f_{\tilde{\mathcal{P}}}(\mathbf{x}_i; \mu) := f_{\tilde{\mathcal{P}}}(\mathbf{x}_i) + \mu \Delta f_i \quad \forall i \in \mathbb{N}_n.$$

In order to always satisfy the KKT conditions for $f_{\tilde{\mathcal{P}}}(\mathbf{x}_i; \mu)$, we solve the following linear system

$$\mathbf{Q}_{\mathcal{A}, \mathcal{A}} \begin{bmatrix} \alpha_{\mathcal{E}} \\ \alpha_{\tilde{\mathcal{I}}'} \\ \alpha_{\tilde{\mathcal{O}}'} \end{bmatrix} = \begin{bmatrix} \mathbf{1} \\ s \\ s \end{bmatrix} - \mathbf{Q}_{\mathcal{A}, \mathcal{L}} \mathbf{1}C - \mathbf{Q}_{\mathcal{A}, \hat{\mathcal{O}}''} \mathbf{1}C\theta \\ - \begin{bmatrix} \mathbf{Q}_{\mathcal{A}, \Delta_{\mathcal{I} \rightarrow \mathcal{O}}} & \mathbf{Q}_{\mathcal{A}, \Delta_{\mathcal{O} \rightarrow \mathcal{I}}} \end{bmatrix} \begin{bmatrix} \alpha_{\Delta_{\mathcal{I} \rightarrow \mathcal{O}}}^{(\text{bef})} - \mathbf{1}C\theta \\ \alpha_{\Delta_{\mathcal{O} \rightarrow \mathcal{I}}}^{(\text{bef})} - \mathbf{1}C \end{bmatrix} \mu,$$

where $\mathcal{A} := \{\mathcal{E}, \tilde{\mathcal{I}}', \tilde{\mathcal{O}}'\}$. This linear system can also be solved by using the piecewise-linear parametric programming while the scalar parameter μ is continuously moved from 1 to 0.

In this parametric problem, we can show that $f_{\tilde{\mathcal{P}}}^*(\mathbf{x}_i; \mu) = f_{\tilde{\mathcal{P}}}^*(\mathbf{x}_i)$ if $\mu = 1$ and $f_{\tilde{\mathcal{P}}}^*(\mathbf{x}_i; \mu) = f_{\mathcal{P}}^*(\mathbf{x}_i)$ if $\mu = 0$ for all $i \in \mathbb{N}_n$.

Since the number of elements in $\Delta_{\mathcal{I} \rightarrow \mathcal{O}}$ and $\Delta_{\mathcal{O} \rightarrow \mathcal{I}}$ are typically small, the D-step can be efficiently implemented by a technique used in the context of incremental learning (Cauwenberghs & Poggio, 2001).