

A Supplementary material

A.1 Proof of Lemma 1:

We first show that the feasible set of (3) is contained in the feasible set of (10). We do this by using the fact that a convex set with a smooth boundary is contained in the halfspace defined by the tangent hyperplane at any point of the boundary of the convex set. Consider a point $(\tilde{w}_\ell, \tilde{x}_\ell)$ on the boundary of the convex set defined by the constraints in (3) and observe that

$$\left\{ (w_\ell, x_\ell) \in \mathbb{R}^2 \mid \begin{matrix} s_\ell w_\ell x_\ell \geq |y_\ell| \\ \text{sign}(w_\ell) = t_\ell \end{matrix} \right\} \subseteq \left\{ (w_\ell, x_\ell) \in \mathbb{R}^2 \mid \begin{pmatrix} s_\ell \tilde{x}_\ell \\ s_\ell \tilde{w}_\ell \end{pmatrix} \cdot \begin{pmatrix} w_\ell - \tilde{w}_\ell \\ x_\ell - \tilde{x}_\ell \end{pmatrix} \geq 0 \right\}. \quad (15)$$

Plugging in $w_\ell = \mathbf{b}_\ell^\top \mathbf{h}$ and $x_\ell = \mathbf{c}_\ell^\top \mathbf{m}$, we have that any feasible (\mathbf{h}, \mathbf{m}) satisfies

$$s_\ell \mathbf{c}_\ell^\top \tilde{\mathbf{m}} \mathbf{b}_\ell^\top \mathbf{h} + s_\ell \mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell^\top \mathbf{m} \geq 2|y_\ell|, \quad \ell = 1, \dots, L,$$

which implies $s_\ell (\mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell^\top \tilde{\mathbf{m}} + \mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell^\top \mathbf{m}) \geq 2|y_\ell|$ for all ℓ . So, the feasible set of (10) contains the feasible set of (3). Lastly, note that among all points $(\mathbf{h}, \mathbf{m}) \in (\tilde{\mathbf{h}}, \tilde{\mathbf{m}}) \oplus S$, only $(\tilde{\mathbf{h}}, \tilde{\mathbf{m}})$ is feasible in (3). So, if $(\tilde{\mathbf{h}}, \tilde{\mathbf{m}})$ solves (10) then $(\tilde{\mathbf{h}}, \tilde{\mathbf{m}})$ solves (3). \square

A.2 Proof of Lemma 2:

Define a one-sided loss function:

$$\mathcal{L}(\mathbf{h}, \mathbf{m}) := \frac{1}{L} \sum_{\ell=1}^L \left[2|y_\ell| - s_\ell \mathbf{c}_\ell^\top \tilde{\mathbf{m}} \mathbf{b}_\ell^\top \mathbf{h} - s_\ell \mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell^\top \mathbf{m} \right]_+,$$

where $(\cdot)_+$ denotes the positive side. The LP in (10) can now be equivalently expressed as

$$(\hat{\mathbf{h}}, \hat{\mathbf{m}}) := \underset{(\mathbf{h}, \mathbf{m}) \in \mathbb{R}^{K+N}}{\text{axrgmin}} \quad \|\mathbf{h}\|_1 + \|\mathbf{m}\|_1 \quad \text{subject to} \quad \mathcal{L}(\mathbf{h}, \mathbf{m}) \leq 0. \quad (16)$$

We want to show that there is no feasible descent direction $(\delta \mathbf{h}, \delta \mathbf{m}) \in \mathcal{D}$ around the true solution $(\tilde{\mathbf{h}}, \tilde{\mathbf{m}})$. Since $(\delta \mathbf{h}, \delta \mathbf{m})$ is a feasible perturbation from the proposed optimal $(\tilde{\mathbf{h}}, \tilde{\mathbf{m}})$, we have from (16)

$$\mathcal{L}(\tilde{\mathbf{h}} + \delta \mathbf{h}, \tilde{\mathbf{m}} + \delta \mathbf{m}) \leq 0. \quad (17)$$

We begin by expanding the loss function $\mathcal{L}(\tilde{\mathbf{h}} + \delta \mathbf{h}, \tilde{\mathbf{m}} + \delta \mathbf{m})$ below

$$\begin{aligned} \mathcal{L}(\tilde{\mathbf{h}} + \delta \mathbf{h}, \tilde{\mathbf{m}} + \delta \mathbf{m}) &= \frac{1}{L} \sum_{\ell=1}^L \left[s_\ell (2y_\ell - \mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell^\top (\tilde{\mathbf{m}} + \delta \mathbf{m}) - \mathbf{c}_\ell^\top \tilde{\mathbf{m}} \mathbf{b}_\ell^\top (\tilde{\mathbf{h}} + \delta \mathbf{h})) \right]_+ \\ &\geq \frac{1}{L} \sum_{\ell=1}^L \left[-s_\ell \mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell^\top \delta \mathbf{m} - s_\ell \mathbf{c}_\ell^\top \tilde{\mathbf{m}} \mathbf{b}_\ell^\top \delta \mathbf{h} \right]_+. \end{aligned} \quad (18)$$

Let $\psi_t(s) := (s)_+ - (s - t)_+$. Using the fact that $\psi_t(s) \leq (s)_+$, and that for every $\alpha, t \geq 0$, and $s \in \mathbb{R}$, $\psi_{\alpha t}(s) = t\psi_\alpha(\frac{s}{t})$, we have

$$\begin{aligned} \frac{1}{L} \sum_{\ell=1}^L \left[-s_\ell \mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell^\top \delta \mathbf{m} - s_\ell \mathbf{c}_\ell^\top \tilde{\mathbf{m}} \mathbf{b}_\ell^\top \delta \mathbf{h} \right]_+ &\geq \frac{1}{L} \sum_{\ell=1}^L \psi_{\tau \|(\delta \mathbf{h}, \delta \mathbf{m})\|_2} \left(-s_\ell \mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell^\top \delta \mathbf{m} - s_\ell \mathbf{c}_\ell^\top \tilde{\mathbf{m}} \mathbf{b}_\ell^\top \delta \mathbf{h} \right) \\ &= \|(\delta \mathbf{h}, \delta \mathbf{m})\|_2 \cdot \frac{1}{L} \sum_{\ell=1}^L \psi_\tau \left(-s_\ell \left\langle (\mathbf{c}_\ell^\top \tilde{\mathbf{m}} \mathbf{b}_\ell, \mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell), \frac{(\delta \mathbf{h}, \delta \mathbf{m})}{\|(\delta \mathbf{h}, \delta \mathbf{m})\|_2} \right\rangle \right) \\ &= \|(\delta \mathbf{h}, \delta \mathbf{m})\|_2 \left[\frac{1}{L} \sum_{\ell=1}^L \mathbb{E} \psi_\tau \left(-s_\ell \left\langle (\mathbf{c}_\ell^\top \tilde{\mathbf{m}} \mathbf{b}_\ell, \mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell), \frac{(\delta \mathbf{h}, \delta \mathbf{m})}{\|(\delta \mathbf{h}, \delta \mathbf{m})\|_2} \right\rangle \right) - \right. \\ &\quad \left. \frac{1}{L} \sum_{\ell=1}^L \left(\mathbb{E} \psi_\tau \left(-s_\ell \left\langle (\mathbf{c}_\ell^\top \tilde{\mathbf{m}} \mathbf{b}_\ell, \mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell), \frac{(\delta \mathbf{h}, \delta \mathbf{m})}{\|(\delta \mathbf{h}, \delta \mathbf{m})\|_2} \right\rangle \right) - \psi_\tau \left(-s_\ell \left\langle (\mathbf{c}_\ell^\top \tilde{\mathbf{m}} \mathbf{b}_\ell, \mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell), \frac{(\delta \mathbf{h}, \delta \mathbf{m})}{\|(\delta \mathbf{h}, \delta \mathbf{m})\|_2} \right\rangle \right) \right) \right]. \end{aligned} \quad (19)$$

The proof mainly relies on lower bounding the right hand side above uniformly over all $(\delta \mathbf{h}, \delta \mathbf{m}) \in \mathcal{D}$. To this end, define a centered random process $\mathcal{R}(\mathbf{B}, \mathbf{C})$ as follows

$$\mathcal{R}(\mathbf{B}, \mathbf{C}) := \sup_{(\delta \mathbf{h}, \delta \mathbf{m}) \in \mathcal{D}} \frac{1}{L} \sum_{\ell=1}^L \left[\mathbb{E} \psi_{\tau} \left(-s_{\ell} \left\langle (\mathbf{c}_{\ell}^{\top} \tilde{\mathbf{m}} \mathbf{b}_{\ell}, \mathbf{b}_{\ell}^{\top} \tilde{\mathbf{h}} \mathbf{c}_{\ell}), \frac{(\delta \mathbf{h}, \delta \mathbf{m})}{\|(\delta \mathbf{h}, \delta \mathbf{m})\|_2} \right\rangle \right) - \psi_{\tau} \left(-s_{\ell} \left\langle (\mathbf{c}_{\ell}^{\top} \tilde{\mathbf{m}} \mathbf{b}_{\ell}, \mathbf{b}_{\ell}^{\top} \tilde{\mathbf{h}} \mathbf{c}_{\ell}), \frac{(\delta \mathbf{h}, \delta \mathbf{m})}{\|(\delta \mathbf{h}, \delta \mathbf{m})\|_2} \right\rangle \right) \right],$$

and an application of bounded difference inequality McDiarmid [1989] yields that $\mathcal{R}(\mathbf{B}, \mathbf{C}) \leq \mathbb{E} \mathcal{R}(\mathbf{B}, \mathbf{C}) + t\tau/\sqrt{L}$ with probability at least $1 - e^{-2Lt^2}$. It remains to evaluate $\mathbb{E} \mathcal{R}(\mathbf{B}, \mathbf{C})$, which after using a simple symmetrization inequality van der Vaart and Wellner [1997] yields

$$\mathbb{E} \mathcal{R}(\mathbf{B}, \mathbf{C}) \leq 2 \mathbb{E} \sup_{(\delta \mathbf{h}, \delta \mathbf{m}) \in \mathcal{D} \cap \mathcal{B}} \frac{1}{L} \sum_{\ell=1}^L \varepsilon_{\ell} \psi_{\tau} \left(-s_{\ell} \left\langle (\mathbf{c}_{\ell}^{\top} \tilde{\mathbf{m}} \mathbf{b}_{\ell}, \mathbf{b}_{\ell}^{\top} \tilde{\mathbf{h}} \mathbf{c}_{\ell}), \frac{(\delta \mathbf{h}, \delta \mathbf{m})}{\|(\delta \mathbf{h}, \delta \mathbf{m})\|_2} \right\rangle \right), \quad (20)$$

where $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_L$ are independent Rademacher random variables. Using the fact that $\psi_t(s)$ is a contraction: $|\psi_t(\alpha_1) - \psi_t(\alpha_2)| \leq |\alpha_1 - \alpha_2|$ for all $\alpha_1, \alpha_2 \in \mathbb{R}$, we have from the Rademacher contraction inequality Ledoux and Talagrand [2013] that

$$\begin{aligned} \mathbb{E} \sup_{(\delta \mathbf{h}, \delta \mathbf{m}) \in \mathcal{D}} \frac{1}{L} \sum_{\ell=1}^L \varepsilon_{\ell} \psi_{\tau} \left(-s_{\ell} \left\langle (\mathbf{c}_{\ell}^{\top} \tilde{\mathbf{m}} \mathbf{b}_{\ell}, \mathbf{b}_{\ell}^{\top} \tilde{\mathbf{h}} \mathbf{c}_{\ell}), \frac{(\delta \mathbf{h}, \delta \mathbf{m})}{\|(\delta \mathbf{h}, \delta \mathbf{m})\|_2} \right\rangle \right) \\ \leq \mathbb{E} \sup_{(\delta \mathbf{h}, \delta \mathbf{m}) \in \mathcal{D}} \frac{1}{L} \sum_{\ell=1}^L -\varepsilon_{\ell} s_{\ell} \left\langle (\mathbf{c}_{\ell}^{\top} \tilde{\mathbf{m}} \mathbf{b}_{\ell}, \mathbf{b}_{\ell}^{\top} \tilde{\mathbf{h}} \mathbf{c}_{\ell}), \frac{(\delta \mathbf{h}, \delta \mathbf{m})}{\|(\delta \mathbf{h}, \delta \mathbf{m})\|_2} \right\rangle \\ = \mathbb{E} \sup_{(\delta \mathbf{h}, \delta \mathbf{m}) \in \mathcal{D}} \frac{1}{L} \sum_{\ell=1}^L \varepsilon_{\ell} \left\langle (\mathbf{c}_{\ell}^{\top} \tilde{\mathbf{m}} \mathbf{b}_{\ell}, \mathbf{b}_{\ell}^{\top} \tilde{\mathbf{h}} \mathbf{c}_{\ell}), \frac{(\delta \mathbf{h}, \delta \mathbf{m})}{\|(\delta \mathbf{h}, \delta \mathbf{m})\|_2} \right\rangle, \end{aligned} \quad (21)$$

where the last equality is the result of the fact that multiplying Rademacher random variables with signs does not change the distribution. In addition, using the facts that $t\mathbf{1}(s \geq t) \leq \psi_t(s)$, and that random vectors $\{(\mathbf{c}_{\ell}^{\top} \tilde{\mathbf{m}} \mathbf{b}_{\ell}, \mathbf{b}_{\ell}^{\top} \tilde{\mathbf{h}} \mathbf{c}_{\ell})\}_{\ell=1}^L$ are identically distributed and the distribution is symmetric, it follows

$$\begin{aligned} \tau \mathbb{P} \left(-s_{\ell} \left\langle (\mathbf{c}_{\ell}^{\top} \tilde{\mathbf{m}} \mathbf{b}_{\ell}, \mathbf{b}_{\ell}^{\top} \tilde{\mathbf{h}} \mathbf{c}_{\ell}), \frac{(\delta \mathbf{h}, \delta \mathbf{m})}{\|(\delta \mathbf{h}, \delta \mathbf{m})\|_2} \right\rangle \geq \tau \right) &= \tau \mathbb{P} \left(\left\langle (\mathbf{c}_{\ell}^{\top} \tilde{\mathbf{m}} \mathbf{b}_{\ell}, \mathbf{b}_{\ell}^{\top} \tilde{\mathbf{h}} \mathbf{c}_{\ell}), \frac{(\delta \mathbf{h}, \delta \mathbf{m})}{\|(\delta \mathbf{h}, \delta \mathbf{m})\|_2} \right\rangle \geq \tau \right) \\ &= \tau \mathbb{E} \left[\mathbf{1} \left(\left\langle (\mathbf{c}_{\ell}^{\top} \tilde{\mathbf{m}} \mathbf{b}_{\ell}, \mathbf{b}_{\ell}^{\top} \tilde{\mathbf{h}} \mathbf{c}_{\ell}), \frac{(\delta \mathbf{h}, \delta \mathbf{m})}{\|(\delta \mathbf{h}, \delta \mathbf{m})\|_2} \right\rangle \geq \tau \right) \right] \leq \mathbb{E} \psi_{\tau} \left(\left\langle (\mathbf{c}_{\ell}^{\top} \tilde{\mathbf{m}} \mathbf{b}_{\ell}, \mathbf{b}_{\ell}^{\top} \tilde{\mathbf{h}} \mathbf{c}_{\ell}), \frac{(\delta \mathbf{h}, \delta \mathbf{m})}{\|(\delta \mathbf{h}, \delta \mathbf{m})\|_2} \right\rangle \right). \end{aligned} \quad (22)$$

Plugging (22), and (21) in (19), we have

$$\begin{aligned} \frac{1}{L} \sum_{\ell=1}^L \left[-s_{\ell} \left\langle (\mathbf{c}_{\ell}^{\top} \tilde{\mathbf{m}} \mathbf{b}_{\ell}, \mathbf{b}_{\ell}^{\top} \tilde{\mathbf{h}} \mathbf{c}_{\ell}), \frac{(\delta \mathbf{h}, \delta \mathbf{m})}{\|(\delta \mathbf{h}, \delta \mathbf{m})\|_2} \right\rangle \right]_+ &\geq \\ \tau \|(\delta \mathbf{h}, \delta \mathbf{m})\|_2 \mathbb{P} \left(\left\langle (\mathbf{c}_{\ell}^{\top} \tilde{\mathbf{m}} \mathbf{b}_{\ell}, \mathbf{b}_{\ell}^{\top} \tilde{\mathbf{h}} \mathbf{c}_{\ell}), \frac{(\delta \mathbf{h}, \delta \mathbf{m})}{\|(\delta \mathbf{h}, \delta \mathbf{m})\|_2} \right\rangle \geq \tau \right) & \\ - \|(\delta \mathbf{h}, \delta \mathbf{m})\|_2 \left(2 \mathbb{E} \sup_{(\delta \mathbf{h}, \delta \mathbf{m}) \in \mathcal{D}} \frac{1}{L} \sum_{\ell=1}^L \varepsilon_{\ell} \left\langle (\mathbf{c}_{\ell}^{\top} \tilde{\mathbf{m}} \mathbf{b}_{\ell}, \mathbf{b}_{\ell}^{\top} \tilde{\mathbf{h}} \mathbf{c}_{\ell}), \frac{(\delta \mathbf{h}, \delta \mathbf{m})}{\|(\delta \mathbf{h}, \delta \mathbf{m})\|_2} \right\rangle + \frac{t\tau}{\sqrt{L}} \right) & \end{aligned}$$

Combining this with (17) and (18), we obtain the final result

$$\begin{aligned} \|(\delta \mathbf{h}, \delta \mathbf{m})\|_2 \left[\tau \mathbb{P} \left(\left\langle (\mathbf{c}_{\ell}^{\top} \tilde{\mathbf{m}} \mathbf{b}_{\ell}, \mathbf{b}_{\ell}^{\top} \tilde{\mathbf{h}} \mathbf{c}_{\ell}), \frac{(\delta \mathbf{h}, \delta \mathbf{m})}{\|(\delta \mathbf{h}, \delta \mathbf{m})\|_2} \right\rangle \geq \tau \right) \right. \\ \left. - \left(2 \mathbb{E} \sup_{(\delta \mathbf{h}, \delta \mathbf{m}) \in \mathcal{D}} \frac{1}{L} \sum_{\ell=1}^L \varepsilon_{\ell} \left\langle (\mathbf{c}_{\ell}^{\top} \tilde{\mathbf{m}} \mathbf{b}_{\ell}, \mathbf{b}_{\ell}^{\top} \tilde{\mathbf{h}} \mathbf{c}_{\ell}), \frac{(\delta \mathbf{h}, \delta \mathbf{m})}{\|(\delta \mathbf{h}, \delta \mathbf{m})\|_2} \right\rangle + \frac{t\tau}{\sqrt{L}} \right) \right] \leq 0. \end{aligned}$$

Using the definitions in (13), and (14), we can write

$$\|(\delta \mathbf{h}, \delta \mathbf{m})\|_2 \left(\tau \mathbf{p}_{\tau}(\mathcal{D}) - \frac{(2\mathfrak{C}(\mathcal{D}) + t\tau)}{\sqrt{L}} \right) \leq 0.$$

It is clear that choosing $L \geq \left(\frac{2\mathfrak{C}(\mathcal{D}) + t\tau}{\tau \mathfrak{p}_\tau(\mathcal{D})} \right)^2$ implies

$$\|(\delta \mathbf{h}, \delta \mathbf{m})\|_2 \leq 0,$$

which directly means that $(\delta \mathbf{h}, \delta \mathbf{m}) = (0, 0)$. Recall that $\mathcal{S} \subset \mathcal{N}$, and $\mathcal{D} \perp \mathcal{N}$, where \mathcal{S} is defined in (11), this implies that the minimizer $(\tilde{\mathbf{h}}, \tilde{\mathbf{m}})$ of the LP (10) resides in the set $(\tilde{\mathbf{h}}, \tilde{\mathbf{m}}) \oplus \mathcal{S}$. This completes the proof of Lemma 2.

A.3 Proof of Theorem 2:

In light of Lemma 2, the proof of Theorem 2 comes down to computing the Rademacher complexity $\mathfrak{C}(\mathcal{D})$ defined in (13), and the tail probability estimate $\mathfrak{p}_\tau(\mathcal{D})$ defined in (14) of the set of descent directions \mathcal{D} defined in (12).

Upper Bound on Rademacher Complexity: We will start by evaluating $\mathfrak{C}(\mathcal{D})$

$$\begin{aligned} \mathfrak{C}(\mathcal{D}) &= \mathbb{E} \sup_{(\delta \mathbf{h}, \delta \mathbf{m}) \in \mathcal{D}} \frac{1}{\sqrt{L}} \sum_{\ell=1}^L \varepsilon_\ell \left\langle \mathbf{c}_\ell^\top \tilde{\mathbf{m}} \mathbf{b}_\ell, \mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell \right\rangle, \frac{(\delta \mathbf{h}, \delta \mathbf{m})}{\|(\delta \mathbf{h}, \delta \mathbf{m})\|_2} \Bigg\rangle \\ &\leq \mathbb{E} \left\| \frac{1}{\sqrt{L}} \sum_{\ell=1}^L \varepsilon_\ell \left(\mathbf{c}_\ell^\top \tilde{\mathbf{m}} \mathbf{b}_\ell|_{\Gamma_h}, \mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell|_{\Gamma_m} \right) \right\|_2 \cdot \sup_{(\delta \mathbf{h}, \delta \mathbf{m}) \in \mathcal{D}} \left\| \frac{(\delta \mathbf{h}_{\Gamma_h}, \delta \mathbf{m}_{\Gamma_m})}{\|(\delta \mathbf{h}, \delta \mathbf{m})\|_2} \right\|_2 \\ &\quad + \mathbb{E} \left\| \frac{1}{\sqrt{L}} \sum_{\ell=1}^L \varepsilon_\ell \left(\mathbf{c}_\ell^\top \tilde{\mathbf{m}} \mathbf{b}_\ell|_{\Gamma_h^c}, \mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell|_{\Gamma_m^c} \right) \right\|_\infty \cdot \sup_{(\delta \mathbf{h}, \delta \mathbf{m}) \in \mathcal{D}} \left\| \frac{(\delta \mathbf{h}_{\Gamma_h^c}, \delta \mathbf{m}_{\Gamma_m^c})}{\|(\delta \mathbf{h}, \delta \mathbf{m})\|_2} \right\|_1. \end{aligned} \quad (23)$$

First note that on set \mathcal{D} (12), we have

$$\left\| \frac{(\delta \mathbf{h}_{\Gamma_h^c}, \delta \mathbf{m}_{\Gamma_m^c})}{\|(\delta \mathbf{h}, \delta \mathbf{m})\|_2} \right\|_1 \leq \sqrt{S_1 + S_2} \left\| \frac{(\delta \mathbf{h}_{\Gamma_h}, \delta \mathbf{m}_{\Gamma_m})}{\|(\delta \mathbf{h}, \delta \mathbf{m})\|_2} \right\|_2 \leq \sqrt{S_1 + S_2}.$$

As for the remaining terms, we begin by writing

$$\begin{aligned} \mathbb{E} \left\| \frac{1}{\sqrt{L}} \sum_{\ell=1}^L \varepsilon_\ell \left(\mathbf{c}_\ell^\top \tilde{\mathbf{m}} \mathbf{b}_\ell|_{\Gamma_h}, \mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell|_{\Gamma_m} \right) \right\|_2 &\leq \sqrt{\mathbb{E} \left\| \frac{1}{\sqrt{L}} \sum_{\ell=1}^L \varepsilon_\ell \left(\mathbf{c}_\ell^\top \tilde{\mathbf{m}} \mathbf{b}_\ell|_{\Gamma_h}, \mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell|_{\Gamma_m} \right) \right\|_2^2} \\ &= \sqrt{\frac{1}{L} \sum_{\ell=1}^L \mathbb{E} \left(|\mathbf{c}_\ell^\top \tilde{\mathbf{m}}|^2 \|\mathbf{b}_\ell|_{\Gamma_h}\|_2^2 + |\mathbf{b}_\ell^\top \tilde{\mathbf{h}}|^2 \|\mathbf{c}_\ell|_{\Gamma_m}\|_2^2 \right)} \\ &= \sqrt{\|\tilde{\mathbf{m}}\|_2^2 S_1 + \|\tilde{\mathbf{h}}\|_2^2 S_2}, \end{aligned}$$

and the second term in (23) is

$$\begin{aligned} \mathbb{E} \left\| \frac{1}{\sqrt{L}} \sum_{\ell=1}^L \varepsilon_\ell \left(\mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell|_{\Gamma_m^c}, \mathbf{c}_\ell^\top \tilde{\mathbf{m}} \mathbf{b}_\ell|_{\Gamma_h^c} \right) \right\|_\infty &\leq \sqrt{\mathbb{E} \left\| \frac{1}{\sqrt{L}} \sum_{\ell=1}^L \varepsilon_\ell \left(\mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell|_{\Gamma_m^c}, \mathbf{c}_\ell^\top \tilde{\mathbf{m}} \mathbf{b}_\ell|_{\Gamma_h^c} \right) \right\|_\infty^2} \\ &\leq \sqrt{2e \log(K + N) \cdot \frac{1}{L} \sum_{\ell=1}^L \mathbb{E} \max \left\{ |\mathbf{c}_\ell^\top \tilde{\mathbf{m}}|^2 \|\mathbf{b}_\ell|_{\Gamma_h^c}\|_\infty^2, |\mathbf{b}_\ell^\top \tilde{\mathbf{h}}|^2 \|\mathbf{c}_\ell|_{\Gamma_m^c}\|_\infty^2 \right\}} \\ &\leq \sqrt{2e \log(K + N) \mathbb{E} \max \left\{ |\mathbf{b}^\top \tilde{\mathbf{h}}|^2 \|\mathbf{c}|_{\Gamma_m^c}\|_\infty^2, |\mathbf{c}^\top \tilde{\mathbf{m}}|^2 \|\mathbf{b}|_{\Gamma_h^c}\|_\infty^2 \right\}} \\ &\leq C \sqrt{\max \{ \|\tilde{\mathbf{h}}\|_2^2, \|\tilde{\mathbf{m}}\|_2^2 \} \log^2(K + N)}, \end{aligned}$$

where the second inequality by the application of Lemma 5.2.2 in Akritas et al. [2016], and the final equality is due to the fact that $\|\mathbf{c}|_{\Gamma_m^c}\|_\infty^2$, and $\|\mathbf{b}|_{\Gamma_h^c}\|_\infty^2$ are subexponential and using Lemma 3 in van de Geer and Lederer [2013].

Plugging the bounds above back in (23), we obtain the upper bound on the Rademacher complexity given below

$$\mathfrak{C}(\mathcal{D}) \leq C \sqrt{(\|\tilde{\mathbf{m}}\|_2^2 + \|\tilde{\mathbf{h}}\|_2^2)(S_1 + S_2) \log^2(K + N)}. \quad (24)$$

Tail Probability: To apply the result in Lemma 2, we also need to evaluate

$$\mathbf{p}_\tau(\mathcal{D}) = \inf_{(\delta\mathbf{h}, \delta\mathbf{m}) \in \mathcal{D}} \mathbb{P} \left(\left\langle (\mathbf{c}_\ell^\top \tilde{\mathbf{m}} \mathbf{b}_\ell, \mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell), \frac{(\delta\mathbf{h}, \delta\mathbf{m})}{\|(\delta\mathbf{h}, \delta\mathbf{m})\|_2} \right\rangle \geq \tau \right). \quad (25)$$

It suffice to estimate the probability $\mathbb{P}(|\mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell^\top \delta\mathbf{m} + \mathbf{b}_\ell^\top \delta\mathbf{h} \mathbf{c}_\ell^\top \tilde{\mathbf{m}}| \geq \tau)$. Using Paley-Zygmund inequality, we obtain

$$\begin{aligned} \mathbb{P} \left(|\mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell^\top \delta\mathbf{m} + \mathbf{b}_\ell^\top \delta\mathbf{h} \mathbf{c}_\ell^\top \tilde{\mathbf{m}}|^2 \geq \frac{1}{2} \mathbb{E} |\mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell^\top \delta\mathbf{m} + \mathbf{b}_\ell^\top \delta\mathbf{h} \mathbf{c}_\ell^\top \tilde{\mathbf{m}}|^2 \right) \\ \geq \frac{1}{4} \cdot \frac{(\mathbb{E} |\mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell^\top \delta\mathbf{m} + \mathbf{b}_\ell^\top \delta\mathbf{h} \mathbf{c}_\ell^\top \tilde{\mathbf{m}}|^2)^2}{\mathbb{E} |\mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell^\top \delta\mathbf{m} + \mathbf{b}_\ell^\top \delta\mathbf{h} \mathbf{c}_\ell^\top \tilde{\mathbf{m}}|^4}. \end{aligned}$$

By the norm equivalence of Gaussian random variables, we have that $(\mathbb{E} |\mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell^\top \delta\mathbf{m} + \mathbf{b}_\ell^\top \delta\mathbf{h} \mathbf{c}_\ell^\top \tilde{\mathbf{m}}|^4)^{1/4} \leq c(\mathbb{E} |\mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell^\top \delta\mathbf{m} + \mathbf{b}_\ell^\top \delta\mathbf{h} \mathbf{c}_\ell^\top \tilde{\mathbf{m}}|^2)^{1/2}$, this implies that

$$\mathbb{P}(|\mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell^\top \delta\mathbf{m} + \mathbf{b}_\ell^\top \delta\mathbf{h} \mathbf{c}_\ell^\top \tilde{\mathbf{m}}|^2 \geq \frac{1}{2} \mathbb{E} |\mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell^\top \delta\mathbf{m} + \mathbf{b}_\ell^\top \delta\mathbf{h} \mathbf{c}_\ell^\top \tilde{\mathbf{m}}|^2) \geq \frac{1}{4} \cdot \frac{1}{c^4}. \quad (26)$$

Finally, a simple calculation shows that $\mathbb{E} |\mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell^\top \delta\mathbf{m} + \mathbf{b}_\ell^\top \delta\mathbf{h} \mathbf{c}_\ell^\top \tilde{\mathbf{m}}|^2 \geq \min\{\|\tilde{\mathbf{h}}\|_2^2, \|\tilde{\mathbf{m}}\|_2^2\} (\|\delta\mathbf{m}\|_2^2 + \|\delta\mathbf{h}\|_2^2)$.

$$\begin{aligned} \mathbb{E} |\mathbf{b}_\ell^\top \tilde{\mathbf{h}} \mathbf{c}_\ell^\top \delta\mathbf{m} + \mathbf{b}_\ell^\top \delta\mathbf{h} \mathbf{c}_\ell^\top \tilde{\mathbf{m}}|^2 &= \mathbb{E}_b \mathbb{E}_c \tilde{\mathbf{h}}^\top \mathbf{b}_\ell \mathbf{b}_\ell^\top \tilde{\mathbf{h}} \delta\mathbf{m}^\top \mathbf{c}_\ell \mathbf{c}_\ell^\top \delta\mathbf{m} + \delta\mathbf{h}^\top \mathbf{b}_\ell \mathbf{b}_\ell^\top \delta\mathbf{h} \tilde{\mathbf{m}}^\top \mathbf{c}_\ell \mathbf{c}_\ell^\top \tilde{\mathbf{m}} \\ &\quad + 2 \mathbb{E}_b \mathbb{E}_c \delta\mathbf{h}^\top \mathbf{b}_\ell \mathbf{b}_\ell^\top \tilde{\mathbf{h}} \delta\mathbf{m}^\top \mathbf{c}_\ell \mathbf{c}_\ell^\top \tilde{\mathbf{m}} \\ &= \mathbb{E}_b (\|\delta\mathbf{m}\|_2^2 \tilde{\mathbf{h}}^\top \mathbf{b}_\ell \mathbf{b}_\ell^\top \tilde{\mathbf{h}} + \|\tilde{\mathbf{m}}\|_2^2 \delta\mathbf{h}^\top \mathbf{b}_\ell \mathbf{b}_\ell^\top \delta\mathbf{h} + 2 \delta\mathbf{m}^\top \tilde{\mathbf{m}} \delta\mathbf{h}^\top \mathbf{b}_\ell \mathbf{b}_\ell^\top \tilde{\mathbf{h}}) \\ &= \|\delta\mathbf{m}\|_2^2 \|\tilde{\mathbf{h}}\|_2^2 + \|\tilde{\mathbf{m}}\|_2^2 \|\delta\mathbf{h}\|_2^2 + 2 \delta\mathbf{m}^\top \tilde{\mathbf{m}} \delta\mathbf{h}^\top \tilde{\mathbf{h}} \\ &= \|\delta\mathbf{m}\|_2^2 \|\tilde{\mathbf{h}}\|_2^2 + \|\tilde{\mathbf{m}}\|_2^2 \|\delta\mathbf{h}\|_2^2 + 2(\delta\mathbf{h}^\top \tilde{\mathbf{h}})^2 \\ &\geq \|\delta\mathbf{m}\|_2^2 \|\tilde{\mathbf{h}}\|_2^2 + \|\tilde{\mathbf{m}}\|_2^2 \|\delta\mathbf{h}\|_2^2, \\ &\geq \min\{\|\tilde{\mathbf{h}}\|_2^2, \|\tilde{\mathbf{m}}\|_2^2\} (\|\delta\mathbf{m}\|_2^2 + \|\delta\mathbf{h}\|_2^2), \end{aligned}$$

where the last equality follows using the fact $(\delta\mathbf{h}, \delta\mathbf{m}) \in \mathcal{D} \subset \mathcal{N}_\perp$, and hence $\mathcal{D} \perp \mathcal{N}$, which implies that $\delta\mathbf{h}^\top \tilde{\mathbf{h}} = \delta\mathbf{m}^\top \tilde{\mathbf{m}}$. Normalizing by $\|(\delta\mathbf{h}, \delta\mathbf{m})\|_2$, and comparing with (25) directly shows that $\tau^2 = \min\{\|\tilde{\mathbf{h}}\|_2^2, \|\tilde{\mathbf{m}}\|_2^2\}$, and $\mathbf{p}_\tau(\mathcal{D}) \geq \frac{1}{8c^4}$. Plugging these results and the Rademacher complexing bound in (24) in Lemma 2 proves Theorem 2. \square

A.4 Evaluation of the Projection Operator

Given a point $(\mathbf{x}', \mathbf{w}', \boldsymbol{\xi}') \in \mathbb{R}^{3L}$, in this section we focus on deriving a closed-form expression for $\text{proj}_{\mathcal{C}}((\mathbf{x}', \mathbf{w}', \boldsymbol{\xi}'))$, where

$$\mathcal{C} = \{(\mathbf{x}, \mathbf{w}, \boldsymbol{\xi}) \in \mathbb{R}^{3L} \mid s_\ell(\xi_\ell + x_\ell)w_\ell \geq |y_\ell|, \ t_\ell w_\ell \geq 0, \ \ell = 1, \dots, L\}$$

is the convex feasible set of (6). It is straightforward to see that the resulting projection program decouples into L convex programs in \mathbb{R}^3 as

$$\arg \min_{x \in \mathbb{R}, w \in \mathbb{R}, \xi \in \mathbb{R}} \frac{1}{2} \left\| \begin{pmatrix} x \\ w \\ \xi \end{pmatrix} - \begin{pmatrix} x'_\ell \\ w'_\ell \\ \xi'_\ell \end{pmatrix} \right\|_2^2 \quad \text{s.t.} \quad |y_\ell| - s_\ell x w - s_\ell \xi w \leq 0, \quad -t_\ell w \leq 0. \quad (27)$$

Throughout this derivation we assume that $|y_\ell| > 0$ (derivation of the projection for the case y_ℓ is easy) and as a result of which the second constraint $-t_\ell w \leq 0$ is never active (because then $w = 0$ and the first constraint requires that $|y_\ell| \leq 0$). We also consistently use the fact that t_ℓ and s_ℓ are signs and nonzero.

Forming the Lagrangian as

$$\mathcal{L}(x, w, \xi, \mu_1, \mu_2) = \frac{1}{2} \left\| \begin{pmatrix} x \\ w \\ \xi \end{pmatrix} - \begin{pmatrix} x'_\ell \\ w'_\ell \\ \xi'_\ell \end{pmatrix} \right\|_2^2 + \mu_1 (|y_\ell| - s_\ell x w - s_\ell \xi w) - \mu_2 (t_\ell w),$$

along with the primal constraints, the KKT optimality conditions are

$$\frac{\partial \mathcal{L}}{\partial x} = x - x'_\ell - \mu_1 s_\ell w = 0, \quad (28)$$

$$\frac{\partial \mathcal{L}}{\partial w} = w - w'_\ell - \mu_1 s_\ell x - \mu_1 s_\ell \xi - \mu_2 t_\ell = 0, \quad (29)$$

$$\frac{\partial \mathcal{L}}{\partial \xi} = \xi - \xi'_\ell - \mu_1 s_\ell w = 0, \quad (30)$$

$$\mu_1 \geq 0, \quad \mu_1 (|y_\ell| - s_\ell x w - s_\ell \xi w) = 0, \quad (31)$$

$$\mu_2 \geq 0, \quad \mu_2 (t_\ell w) = 0. \quad (32)$$

We now proceed with the possible cases.

Case 1. $\mu_1 = \mu_2 = 0$:

In this case we have $(x, w, \xi) = (x'_\ell, w'_\ell, \xi'_\ell)$ and this result would only be acceptable when $|y_\ell| - s_\ell x'_\ell w'_\ell - s_\ell \xi'_\ell w'_\ell \leq 0$ and $t_\ell w'_\ell \geq 0$.

Case 2. $\mu_1 = 0, t_\ell w = 0$:

In this case the first feasibility constraint of (27) requires that $|y_\ell| \leq 0$, which is not possible when $|y_\ell| > 0$.

Case 3. $|y_\ell| - s_\ell x w - s_\ell \xi w = 0, t_\ell w = 0$:

Similar to the previous case, this cannot happen when $|y_\ell| > 0$.

Case 4. $\mu_2 = 0, |y_\ell| - s_\ell x w - s_\ell \xi w = 0$:

In this case we have

$$|y_\ell| = s_\ell x w + s_\ell \xi w.$$

Now combining this observation with (28) and (30) yields

$$|y_\ell| = s_\ell (x'_\ell + \mu_1 s_\ell w) w + s_\ell (\xi'_\ell + \mu_1 s_\ell w) w, \quad (33)$$

and therefore

$$\mu_1 = \frac{|y_\ell| - s_\ell (x'_\ell + \xi'_\ell) w}{2w^2}. \quad (34)$$

Similarly, (29) yields

$$w = w'_\ell + \mu_1 s_\ell (x'_\ell + \mu_1 s_\ell w) + \mu_1 s_\ell (\xi'_\ell + \mu_1 s_\ell w). \quad (35)$$

Knowing that $w \neq 0$, μ_1 can be eliminated between (33) and (35) to generate the following forth order polynomial equation in terms of w :

$$2w^4 - 2w'_\ell w^3 + s_\ell |y_\ell| (x'_\ell + \xi'_\ell) w - y_\ell^2 = 0.$$

After solving this 4-th order polynomial equation (e.g., the root command in MATLAB) we pick the real root w which obeys

$$t_\ell w \geq 0, \quad |y_\ell| - s_\ell (x'_\ell + \xi'_\ell) w \geq 0. \quad (36)$$

Note that the second inequality in (36) warrants nonnegative values for μ_1 thanks to (34). After picking the right root, we can explicitly obtain μ_1 using (35) and calculate the solutions x and ξ using (28) and (30). Technically, in using the ADMM scheme for each ℓ we solve a forth-order polynomial equation and find the projection.