

Supplement to “Domain Adaptation with Conditional Transferable Components”

This supplementary material provides the proofs and some details which are omitted in the submitted paper. The equation numbers in this material are consistent with those in the paper.

S1. Proof of Theorem 1

Proof. Combine (3) and (4), we have

$$\sum_{c=1}^C p^{\mathcal{T}}(Y = v_c) p^{\mathcal{T}}(X^{ci}|Y = v_c) = \sum_{c=1}^C p^{\text{new}}(Y = v_c) p^{\mathcal{S}}(X^{ci}|Y = v_c). \quad (15)$$

If the transformation W is non-trivial, there do not exist non-zero $\gamma_1, \dots, \gamma_C$ and ν_1, \dots, ν_C such that $\sum_{c=1}^C \gamma_c p^{\mathcal{T}}(X^{ci}|Y = v_c) = 0$ and $\sum_{c=1}^C \nu_c p^{\mathcal{S}}(X^{ci}|Y = v_c) = 0$. Therefore, we can transform (15) to

$$\sum_{c=1}^C P^{\mathcal{T}}(Y = v_c) P^{\mathcal{T}}(X^{ci}|Y = v_c) - P^{\text{new}}(Y = v_c) p^{\mathcal{S}}(X^{ci}|Y = v_c) = 0. \quad (16)$$

According to \mathbf{A}^{CIC} in Theorem 1, we have $\forall c$,

$$P^{\mathcal{T}}(Y = v_c) P^{\mathcal{T}}(X^{ci}|Y = v_c) - P^{\text{new}}(Y = v_c) p^{\mathcal{S}}(X^{ci}|Y = v_c) = 0. \quad (17)$$

Taking the integral of (17) leads to $P^{\text{new}}(Y = v_c) = P^{\mathcal{T}}(Y = v_c)$, which further implies that $p^{\mathcal{S}}(X^{ci}|Y = v_c) = P^{\mathcal{T}}(X^{ci}|Y = v_c)$. \square

S2. Proof of Lemma 1

Proof.

$$\begin{aligned} \epsilon_{\mathcal{T}}(h) &= \epsilon_{\mathcal{T}}(h) + \epsilon_{\text{new}}(h) - \epsilon_{\text{new}}(h) \\ &\leq \epsilon_{\text{new}}(h) + |\epsilon_{\mathcal{T}}(h) - \epsilon_{\text{new}}(h)| \\ &\leq \epsilon_{\text{new}}(h) + \int |P^{\text{new}}(X^{ci}, Y) - P^{\mathcal{T}}(X^{ci}, Y)| |L(Y, h(X^{ci}))| dX^{ci} dY \\ &\leq \epsilon_{\text{new}}(h) + d_1(p^{\text{new}}(X^{ci}, Y), p^{\mathcal{T}}(X^{ci}, Y)). \end{aligned} \quad (18)$$

\square

S3. Proof of Theorem 2

Proof. In the binary classification problem, because $Y \in \{0, 1\}$ is a discrete variable, we use the Kronecker delta kernel for Y . Then (13) becomes

$$\begin{aligned} &d_k(p^{\text{new}}(X^{ci}, Y), p^{\mathcal{T}}(X^{ci}, Y)) \\ &= \sum_{c=0}^1 \left\| P^{\text{new}}(Y = c) \mu_{p^{\mathcal{S}}(X^{ci}|Y=c)}[\psi(X^{ci})] - P^{\mathcal{T}}(Y = c) \mu_{p^{\mathcal{T}}(X^{ci}|Y=c)}[\psi(X^{ci})] \right\|^2 \\ &= \|\Delta_1\|^2 + \|\Delta_0\|^2 \\ &= \|\Delta_1 + \Delta_0\|^2 - 2\Delta_1^{\top} \Delta_0 \end{aligned}$$

$$\begin{aligned}
 &= \left\| \sum_{c=0}^1 P^{\text{new}}(Y=c) \mu_{p^S}(X^{ci}|Y=c) [\psi(X^{ci})] - \sum_{c=0}^1 P^{\mathcal{T}}(Y=c) \mu_{p^{\mathcal{T}}}(X^{ci}|Y=c) [\psi(X^{ci})] \right\|^2 - 2\Delta_1^\top \Delta_0 \\
 &= \left\| \mu_{p^{\text{new}}}(X^{ci}) [\psi(X^{ci})] - \mu_{p^{\mathcal{T}}}(X^{ci}) [\psi(X^{ci})] \right\|^2 - 2\Delta_1^\top \Delta_0 \\
 &= J^{ci} - 2\Delta_1^\top \Delta_0.
 \end{aligned} \tag{19}$$

Clearly, when $0 < \theta \leq \pi/2$, we have $\Delta_1^\top \Delta_0 \geq 0$. Therefore,

$$d_k(p^{\text{new}}(X^{ci}, Y), p^{\mathcal{T}}(X^{ci}, Y)) \leq J^{ci}. \tag{20}$$

When $\pi/2 < \theta \leq \pi$, we express J^{ci} as

$$\begin{aligned}
 J^{ci} &= \left\| \Delta_1 + \Delta_0 \right\|^2 \\
 &= \left\| \Delta_1 \right\|^2 + \left\| \Delta_0 \right\|^2 + 2 \left\| \Delta_1 \right\| \left\| \Delta_0 \right\| \cos \theta \\
 &= (\left\| \Delta_1 \right\| + \left\| \Delta_0 \right\| \cos \theta)^2 + \left\| \Delta_0 \right\|^2 \sin^2 \theta
 \end{aligned} \tag{21}$$

$$= (\left\| \Delta_0 \right\| + \left\| \Delta_1 \right\| \cos \theta)^2 + \left\| \Delta_1 \right\|^2 \sin^2 \theta. \tag{22}$$

According to (21) and (22), we have $\left\| \Delta_0 \right\|^2 \sin^2 \theta \leq J^{ci}$ and $\left\| \Delta_1 \right\|^2 \sin^2 \theta \leq J^{ci}$. Thus

$$d_k(p^{\text{new}}(X^{ci}, Y), p^{\mathcal{T}}(X^{ci}, Y)) = \left\| \Delta_1 \right\|^2 + \left\| \Delta_0 \right\|^2 \leq 2 \frac{J^{ci}}{\sin^2 \theta}. \tag{23}$$

Combining (20) and (23), we can obtain the results in Theorem 2. \square

S4. Proof of Theorem 3

Proof. We have

$$\begin{aligned}
 \hat{J}^{ci}(\boldsymbol{\beta}, W) &= \left\| \frac{1}{n^S} \psi(W^\top \mathbf{x}^S) \boldsymbol{\beta} - \frac{1}{n^{\mathcal{T}}} \psi(W^\top \mathbf{x}^{\mathcal{T}}) \mathbf{1} \right\|^2 \\
 &= \left\| \frac{1}{n^S} \psi(W^\top \mathbf{x}^S) R^{dis} \boldsymbol{\alpha} - \frac{1}{n^{\mathcal{T}}} \psi(W^\top \mathbf{x}^{\mathcal{T}}) \mathbf{1} \right\|^2 \\
 &= \left\| \left[\frac{1}{n_1} \sum_{i=1}^{n_1} \psi(W_n^\top x_{1i}^S), \dots, \frac{1}{n_C} \sum_{i=1}^{n_C} \psi(W_n^\top x_{Ci}^S) \right] \boldsymbol{\alpha} - \frac{1}{n^{\mathcal{T}}} \psi(W^\top \mathbf{x}^{\mathcal{T}}) \mathbf{1} \right\|^2 \\
 &= \hat{J}^{ci}(\boldsymbol{\alpha}, W),
 \end{aligned} \tag{24}$$

where $x_{ci}^S, c \in \{1, \dots, C\}$ denotes the i -th observation of the c -th class in the source domain.

Define $\Delta = \{\boldsymbol{\alpha} | \boldsymbol{\alpha} \geq 0, \sum_{c=1}^C \boldsymbol{\alpha}_c = 1\}$. We have

$$\begin{aligned}
 &J^{ci}(\boldsymbol{\alpha}_n, W_n) - J^{ci}(\boldsymbol{\alpha}^*, W_n) \\
 &= J^{ci}(\boldsymbol{\alpha}_n, W_n) - \hat{J}^{ci}(\boldsymbol{\alpha}_n, W_n) + \hat{J}^{ci}(\boldsymbol{\alpha}_n, W_n) - \hat{J}^{ci}(\boldsymbol{\alpha}^*, W_n) + \hat{J}^{ci}(\boldsymbol{\alpha}^*, W_n) - J^{ci}(\boldsymbol{\alpha}^*, W_n) \\
 &\quad \text{Since } \boldsymbol{\alpha}_n \text{ is the empirical minimizer and thus } \hat{J}^{ci}(\boldsymbol{\alpha}_n, W_n) \leq \hat{J}^{ci}(\boldsymbol{\alpha}^*, W_n) \\
 &\leq J^{ci}(\boldsymbol{\alpha}_n, W_n) - \hat{J}^{ci}(\boldsymbol{\alpha}_n, W_n) + \hat{J}^{ci}(\boldsymbol{\alpha}^*, W_n) - J^{ci}(\boldsymbol{\alpha}^*, W_n) \\
 &\leq 2 \sup_{\boldsymbol{\alpha} \in \Delta} |J^{ci}(\boldsymbol{\alpha}, W_n) - \hat{J}^{ci}(\boldsymbol{\alpha}, W_n)|.
 \end{aligned} \tag{25}$$

Before upper bounding the above defect on the right hand side, we enable some properties of the kernel. Assume that there exists a ψ_{\max} such that for any $x \in \mathcal{X}$, it holds that $-\psi_{\max} \leq \psi(x) \leq \psi_{\max}$ and that $\|\psi_{\max}\|_2 \leq \wedge_2$. Since $\boldsymbol{\alpha} \geq 0$ and $\|\boldsymbol{\alpha}\|_1 = 1$, for any \mathbf{x}^S , it also holds that $\left[\frac{1}{n_1} \sum_{i=1}^{n_1} \psi(W_n^\top x_{1i}^S), \dots, \frac{1}{n_C} \sum_{i=1}^{n_C} \psi(W_n^\top x_{Ci}^S) \right] \boldsymbol{\alpha} \leq \psi_{\max}$.

Now, we have the following Lipschitz property of J^{ci} :

$$|J^{ci}(\boldsymbol{\alpha}, W_n) - \hat{J}^{ci}(\boldsymbol{\alpha}, W_n)|$$

$$\begin{aligned}
 &\leq \left| \max_{\boldsymbol{\alpha}, \mathbf{x}^S} \left[\frac{1}{n_1} \sum_{i=1}^{n_1} \psi(W_n^\top x_{1i}^S), \dots, \frac{1}{n_C} \sum_{i=1}^{n_C} \psi(W_n^\top x_{Ci}^S) \right] \boldsymbol{\alpha} + \max_{\mathbf{x}^S} \frac{1}{n^T} \psi(W_n^\top \mathbf{x}^T) \mathbf{1} \right| \mathbb{E}_{(Y, X) \sim p^S} [\beta(Y) \psi(W_n^\top X)] \\
 &\quad - \mathbb{E}_{X \sim p^T} [\psi(W_n^\top X)] - \left[\frac{1}{n_1} \sum_{i=1}^{n_1} \psi(W_n^\top x_{1i}^S), \dots, \frac{1}{n_C} \sum_{i=1}^{n_C} \psi(W_n^\top x_{Ci}^S) \right] \boldsymbol{\alpha} + \frac{1}{n^T} \psi(W_n^\top \mathbf{x}^T) \mathbf{1} \\
 &\leq 2 |\psi_{\max}|^T \mathbb{E}_{(Y, X) \sim p^S} [\beta(Y) \psi(W_n^\top X)] \\
 &\quad - \mathbb{E}_{X \sim p^T} [\psi(W_n^\top X)] - \left[\frac{1}{n_1} \sum_{i=1}^{n_1} \psi(W_n^\top x_{1i}^S), \dots, \frac{1}{n_C} \sum_{i=1}^{n_C} \psi(W_n^\top x_{Ci}^S) \right] \boldsymbol{\alpha} + \frac{1}{n^T} \psi(W_n^\top \mathbf{x}^T) \mathbf{1}. \tag{26}
 \end{aligned}$$

Then, combining (25) and (26), we have

$$\begin{aligned}
 &J^{ci}(\boldsymbol{\alpha}_n, W_n) - J^{ci}(\boldsymbol{\alpha}^*, W_n) \\
 &\leq 2 \sup_{\boldsymbol{\alpha} \in \Delta} |J^{ci}(\boldsymbol{\alpha}, W_n) - \hat{J}^{ci}(\boldsymbol{\alpha}, W_n)| \\
 &\leq 4 \sup_{\boldsymbol{\alpha} \in \Delta} |\psi_{\max}|^T \mathbb{E}_{(Y, X) \sim p^S} [\beta(Y) \psi(W_n^\top X)] \\
 &\quad - \mathbb{E}_{X \sim p^T} [\psi(W_n^\top X)] - \left[\frac{1}{n_1} \sum_{i=1}^{n_1} \psi(W_n^\top x_{1i}^S), \dots, \frac{1}{n_C} \sum_{i=1}^{n_C} \psi(W_n^\top x_{Ci}^S) \right] \boldsymbol{\alpha} + \frac{1}{n^T} \psi(W_n^\top \mathbf{x}^T) \mathbf{1}. \tag{27}
 \end{aligned}$$

Now, we are going to upper bound the defect:

$$\begin{aligned}
 f^\psi(X, \mathbf{x}^S, \mathbf{x}^T) &\triangleq \mathbb{E}_{(Y, X) \sim p^S} [\beta(Y) \psi(W_n^\top X)] \\
 &\quad - \mathbb{E}_{X \sim p^T} [\psi(W_n^\top X)] - \left[\frac{1}{n_1} \sum_{i=1}^{n_1} \psi(W_n^\top x_{1i}^S), \dots, \frac{1}{n_C} \sum_{i=1}^{n_C} \psi(W_n^\top x_{Ci}^S) \right] \boldsymbol{\alpha} + \frac{1}{n^T} \psi(W_n^\top \mathbf{x}^T) \mathbf{1}. \tag{28}
 \end{aligned}$$

We employ the McDiarmid's inequality to upper bound the ℓ_2 -norm of the defect.

Theorem 4 (McDiarmid's inequality). *Let $X = (X_1, \dots, X_n)$ be an independent and identically distributed sample and X^i a new sample with the i -th example in X being replaced by an independent example $X_i^!$. If there exists $c_1, \dots, c_n > 0$ such that $f : \mathcal{X}^n \rightarrow \mathbb{R}$ satisfies the following conditions:*

$$|f(X) - f(X^i)| \leq c_i, \forall i \in \{1, \dots, n\}. \tag{29}$$

Then for any $X \in \mathcal{X}^n$ and $\epsilon > 0$, the following inequalities hold:

$$\Pr\{|Ef(X) - f(X)| \geq \epsilon\} \leq 2 \exp\left(\frac{-2\epsilon^2}{\sum_{i=1}^n c_i^2}\right). \tag{30}$$

We now check that $f(X, \mathbf{x}^S, \mathbf{x}^T) = \|f^\psi(X, \mathbf{x}^S, \mathbf{x}^T)\|^2$ satisfies the bounded difference property. Let \mathbf{x}_{ci}^S denote the i -th observation belonging to the c -th class. We have

$$\begin{aligned}
 &|f(X, \mathbf{x}_i^S, \mathbf{x}^T) - f(X, \mathbf{x}^S, \mathbf{x}^T)| \\
 &= |(f^\psi(X, \mathbf{x}_i^S, \mathbf{x}^T) + f^\psi(X, \mathbf{x}^S, \mathbf{x}^T))^T (f^\psi(X, \mathbf{x}_i^S, \mathbf{x}^T) - f^\psi(X, \mathbf{x}^S, \mathbf{x}^T))| \\
 &\leq 4 |\psi_{\max}|^T |f^\psi(X, \mathbf{x}_i^S, \mathbf{x}^T) - f^\psi(X, \mathbf{x}^S, \mathbf{x}^T)| \\
 &= 4 |\psi_{\max}|^T \left| \frac{\boldsymbol{\alpha}_c}{n_c} (\psi(W_n^\top \mathbf{x}_{ci}^S) - \psi(W_n^\top \mathbf{x}'_{ci}^S)) \right| \\
 &\leq \frac{8\boldsymbol{\alpha}_c}{n_c} |\psi_{\max}|^T |\psi_{\max}| \leq \frac{8\Lambda_2^2 \boldsymbol{\alpha}_c}{n_c}. \tag{31}
 \end{aligned}$$

Similarly, it holds that

$$|f(X, \mathbf{x}^S, \mathbf{x}_i^T) - f(X, \mathbf{x}^S, \mathbf{x}^T)| \leq \frac{8\Lambda_2^2}{n^S}. \tag{32}$$

Employing McDiarmid's inequality, we have that

$$\Pr\{|f(X, \mathbf{x}^S, \mathbf{x}^T) - E_{\mathbf{x}^S, \mathbf{x}^T} f(X, \mathbf{x}^S, \mathbf{x}^T)| \geq \epsilon\} \leq 2 \exp\left(\frac{-2\epsilon^2}{64 \wedge_2^4 (\sum_{c=1}^C \frac{\alpha_c^2}{n_c} + \frac{1}{n^T})}\right). \quad (33)$$

Combining (27) and (33), we have that for any $\delta > 0$, with probability at least $1 - \delta$,

$$\begin{aligned} & J^{ci}(\boldsymbol{\alpha}_n, W_n) - J^{ci}(\boldsymbol{\alpha}^*, W_n) \\ & \leq 2 \sup_{\boldsymbol{\alpha} \in \Delta} |J^{ci}(\boldsymbol{\alpha}, W_n) - \hat{J}^{ci}(\boldsymbol{\alpha}, W_n)| \\ & \leq 4 \sup_{\boldsymbol{\alpha} \in \Delta} |\psi_{\max}|^\Gamma |f^\psi(X, \mathbf{x}^S, \mathbf{x}^T)| \\ & \quad \text{Using Cauchy-Schwarz inequality} \\ & \leq 4 \sup_{\boldsymbol{\alpha} \in \Delta} \|\psi_{\max}\| \|f^\psi(X, \mathbf{x}^S, \mathbf{x}^T)\| \\ & \leq 4 \wedge_2 \left(E_{\mathbf{x}^S, \mathbf{x}^T} \sup_{\boldsymbol{\alpha} \in \Delta} f(X, \mathbf{x}^S, \mathbf{x}^T) + \wedge_2^2 \sqrt{32 \log \frac{2}{\delta} \left(\sum_{c=1}^C \frac{\alpha_c^2}{n_c} + \frac{1}{n^T} \right)} \right)^{\frac{1}{2}} \\ & \leq 4 \wedge_2 \left(E_{\mathbf{x}^S, \mathbf{x}^T} \sup_{\boldsymbol{\alpha} \in \Delta} f(X, \mathbf{x}^S, \mathbf{x}^T) + 32 \wedge_2^2 \sqrt{\frac{1}{2} \log \frac{2}{\delta} \left(\max_{c \in \{1, \dots, C\}} \frac{1}{n_c} + \frac{1}{n^T} \right)} \right)^{\frac{1}{2}}. \end{aligned} \quad (34)$$

Now we are going to upper bound the term $E_X \sup_{\boldsymbol{\alpha} \in \Delta} f(X, \mathbf{x}^S, \mathbf{x}^T)$. Let

$$g_n(\mathbf{x}^S, \mathbf{x}^T) \triangleq \left[\frac{1}{n_1} \sum_{i=1}^{n_1} \psi(W_n^\top x_{1i}^S), \dots, \frac{1}{n_C} \sum_{i=1}^{n_C} \psi(W_n^\top x_{Ci}^S) \right] \boldsymbol{\alpha} - \frac{1}{n^T} \psi(W_n^\top \mathbf{x}^T) \mathbf{1} \quad (35)$$

and

$$g(X) \triangleq \mathbb{E}_{(Y, X) \sim p^S} [\boldsymbol{\beta}(Y) \psi(W_n^\top X)] - \mathbb{E}_{X \sim p^T} [\psi(W_n^\top X)]. \quad (36)$$

We have that

$$\begin{aligned} & E_{\mathbf{x}^S, \mathbf{x}^T} \sup_{\boldsymbol{\alpha} \in \Delta} \|f^\psi(X, \mathbf{x}^S, \mathbf{x}^T)\|^2 \\ & = E_{\mathbf{x}^S, \mathbf{x}^T} \sup_{\boldsymbol{\alpha} \in \Delta} \|g(X) - g_n(\mathbf{x}^S, \mathbf{x}^T)\|^2 \\ & = E_{\mathbf{x}^S, \mathbf{x}^T} \sup_{\boldsymbol{\alpha} \in \Delta} \|E_{\mathbf{x}'^S, \mathbf{x}'^T} g_n(\mathbf{x}'^S, \mathbf{x}'^T) - g_n(\mathbf{x}^S, \mathbf{x}^T)\|^2 \\ & \leq E_{\mathbf{x}^S, \mathbf{x}^T, \mathbf{x}'^S, \mathbf{x}'^T} \sup_{\boldsymbol{\alpha} \in \Delta} \|g_n(\mathbf{x}'^S, \mathbf{x}'^T) - g_n(\mathbf{x}^S, \mathbf{x}^T)\|^2. \end{aligned} \quad (37)$$

where $\mathbf{x}'^S, \mathbf{x}'^T$ are ghost samples which are i.i.d. with $\mathbf{x}^S, \mathbf{x}^T$, respectively.

Since $\mathbf{x}^j, \mathbf{x}'^j, j = S, T$ are i.i.d. samples, $\sum_{i=1}^{n_c} \psi(W_n^\top \mathbf{x}_{ci}^j) - \psi(W_n^\top \mathbf{x}'_{ci}^j)$ has a symmetric property, which means it has an even density function. Thus, $\sum_{i=1}^{n_c} \psi(W_n^\top \mathbf{x}_{ci}^j) - \psi(W_n^\top \mathbf{x}'_{ci}^j)$ and $\sum_{i=1}^{n_c} \sigma_{ci} (\psi(W_n^\top \mathbf{x}_{ci}^j) - \psi(W_n^\top \mathbf{x}'_{ci}^j))$ has the same distribution, where σ_{ci} are independent variables uniformly distributed from $\{-1, 1\}$. Then, we have

$$E_{\mathbf{x}^S, \mathbf{x}^T, \mathbf{x}'^S, \mathbf{x}'^T} \sup_{\boldsymbol{\alpha} \in \Delta} \|g_n(\mathbf{x}'^S, \mathbf{x}'^T) - g_n(\mathbf{x}^S, \mathbf{x}^T)\|^2 = E_{\mathbf{x}^S, \mathbf{x}^T, \mathbf{x}'^S, \mathbf{x}'^T, \boldsymbol{\sigma}} \sup_{\boldsymbol{\alpha} \in \Delta} \|g_n(\mathbf{x}'^S, \mathbf{x}'^T, \boldsymbol{\sigma}) - g_n(\mathbf{x}^S, \mathbf{x}^T, \boldsymbol{\sigma})\|^2, \quad (38)$$

where

$$g_n(\mathbf{x}^S, \mathbf{x}^T, \boldsymbol{\sigma}) \triangleq \left[\frac{1}{n_1} \sum_{i=1}^{n_1} \sigma_{1i} (\psi(W_n^\top \mathbf{x}_{ci}^S)) \dots \frac{1}{n_C} \sum_{i=1}^{n_C} \sigma_{Ci} (\psi(W_n^\top \mathbf{x}_{ci}^S)) \right] \boldsymbol{\alpha} - \frac{1}{n^T} \sum_{i=1}^{n^T} \sigma_{Ti} \psi(W_n^\top \mathbf{x}_i^T). \quad (39)$$

According to Talagrand contraction Lemma, we have

$$\begin{aligned}
 & E_{\mathbf{x}^S, \mathbf{x}^T, \mathbf{x}'^S, \mathbf{x}'^T, \boldsymbol{\sigma}} \sup_{\boldsymbol{\alpha} \in \Delta} \left\| g_n(\mathbf{x}'^S, \mathbf{x}'^T, \boldsymbol{\sigma}) - g_n(\mathbf{x}^S, \mathbf{x}^T, \boldsymbol{\sigma}) \right\|^2 \\
 & \leq 2E_{\mathbf{x}^S, \mathbf{x}^T, \mathbf{x}'^S, \mathbf{x}'^T, \boldsymbol{\sigma}} \sup_{\boldsymbol{\alpha} \in \Delta} |\psi_{\max}|^T |g_n(\mathbf{x}'^S, \mathbf{x}'^T, \boldsymbol{\sigma}) - g_n(\mathbf{x}^S, \mathbf{x}^T, \boldsymbol{\sigma})| \\
 & \leq 4E_{\mathbf{x}^S, \mathbf{x}^T, \mathbf{x}'^S, \mathbf{x}'^T, \boldsymbol{\sigma}} \sup_{\boldsymbol{\alpha} \in \Delta} |\psi_{\max}|^T |g_n(\mathbf{x}^S, \mathbf{x}^T, \boldsymbol{\sigma})| \\
 & = 4E_{\mathbf{x}^S, \mathbf{x}^T, \mathbf{x}'^S, \mathbf{x}'^T, \boldsymbol{\sigma}} \sup_{\boldsymbol{\alpha} \in \Delta} |\psi_{\max}|^T \\
 & \left\langle [\boldsymbol{\alpha}^T, -1]^T, \left[\frac{1}{n_1} \sum_{i=1}^{n_c} \sigma_{1i} (\psi(W_n^T \mathbf{x}_{ci}^S)), \dots, \frac{1}{n_C} \sum_{i=1}^{n_c} \sigma_{Ci} (\psi(W_n^T \mathbf{x}_{ci}^S)), \frac{1}{n^T} \sum_{i=1}^{n^T} \sigma_{Ti} \psi(W_n^T \mathbf{x}_i^T) \right]^T \right\rangle. \quad (40)
 \end{aligned}$$

Let

$$\mathbf{v} \triangleq \left[\frac{1}{n_1} \sum_{i=1}^{n_c} \sigma_{1i} (\psi(W_n^T \mathbf{x}_{ci}^S)), \dots, \frac{1}{n_C} \sum_{i=1}^{n_c} \sigma_{Ci} (\psi(W_n^T \mathbf{x}_{ci}^S)), \frac{1}{n^T} \sum_{i=1}^{n^T} \sigma_{Ti} \psi(W_n^T \mathbf{x}_i^T) \right]^T. \quad (41)$$

Since $\|[\boldsymbol{\alpha}^T, -1]^T\|_2 \leq 2$, using Cauchy-Schwarz inequality again, we have

$$\begin{aligned}
 & E_{\mathbf{x}^S, \mathbf{x}^T, \mathbf{x}'^S, \mathbf{x}'^T, \boldsymbol{\sigma}} \sup_{\boldsymbol{\alpha} \in \Delta} \left\| g_n(\mathbf{x}'^S, \mathbf{x}'^T, \boldsymbol{\sigma}) - g_n(\mathbf{x}^S, \mathbf{x}^T, \boldsymbol{\sigma}) \right\|^2 \\
 & \leq 4E_{\mathbf{x}^S, \mathbf{x}^T, \mathbf{x}'^S, \mathbf{x}'^T, \boldsymbol{\sigma}} \sup_{\boldsymbol{\alpha} \in \Delta} |\psi_{\max}|^T \langle [\boldsymbol{\alpha}^T, -1]^T, \mathbf{v} \rangle \\
 & \leq 8E_{\mathbf{x}^S, \mathbf{x}^T, \mathbf{x}'^S, \mathbf{x}'^T, \boldsymbol{\sigma}} |\psi_{\max}|^T \sqrt{\mathbf{v}^T \mathbf{v}} \\
 & \leq 8E_{\mathbf{x}^S, \mathbf{x}^T, \mathbf{x}'^S, \mathbf{x}'^T, \boldsymbol{\sigma}} |\psi_{\max}|^T \sqrt{E_{\boldsymbol{\sigma}} \mathbf{v}^T \mathbf{v}} \\
 & = 8E_{\mathbf{x}^S, \mathbf{x}^T, \mathbf{x}'^S, \mathbf{x}'^T, \boldsymbol{\sigma}} |\psi_{\max}|^T \sqrt{\sum_{c=1}^C \frac{1}{n_c^2} \sum_{i=1}^{n_c} (\psi(W_n^T \mathbf{x}_{ci}^S))^2 + \frac{1}{(n^T)^2} \sum_{i=1}^{n^T} (\psi(W_n^T \mathbf{x}_i^T))^2} \\
 & \leq 8|\psi_{\max}|^T |\psi_{\max}| \sqrt{\sum_{c=1}^C \frac{1}{n_c} + \frac{1}{n^T}} \\
 & \leq 8 \wedge_2^2 \sqrt{\sum_{c=1}^C \frac{1}{n_c} + \frac{1}{n^T}}. \quad (42)
 \end{aligned}$$

At the end, combining (34), (37) and (42), with probability at least $1 - \delta$, we have

$$\begin{aligned}
 & J^{ci}(\boldsymbol{\alpha}_n, W_n) - J^{ci}(\boldsymbol{\alpha}^*, W_n) \\
 & \leq 4 \wedge_2 \left(E_{\mathbf{x}^S, \mathbf{x}^T} \sup_{\boldsymbol{\alpha} \in \Delta} f(X, \mathbf{x}^S, \mathbf{x}^T) + 32 \wedge_2 \sqrt{\frac{1}{2} \log \frac{2}{\delta} \left(\max_{c \in \{1, \dots, C\}} \frac{1}{n_c} + \frac{1}{n^T} \right)} \right)^{\frac{1}{2}} \\
 & = 4 \wedge_2 \left(8 \wedge_2^2 \sqrt{\sum_{c=1}^C \frac{1}{n_c} + \frac{1}{n^T}} + 32 \wedge_2^2 \sqrt{\frac{1}{2} \log \frac{2}{\delta} \left(\max_{c \in \{1, \dots, C\}} \frac{1}{n_c} + \frac{1}{n^T} \right)} \right)^{\frac{1}{2}} \\
 & \leq 8 \wedge_2^2 \left(2 \sqrt{\sum_{c=1}^C \frac{1}{n_c} + \frac{1}{n^T}} + 8 \sqrt{\frac{1}{2} \log \frac{2}{\delta} \left(\max_{c \in \{1, \dots, C\}} \frac{1}{n_c} + \frac{1}{n^T} \right)} \right)^{\frac{1}{2}}.
 \end{aligned}$$

The proof ends. \square

S5. Derivatives used in Sec. 2.5

The gradient of \hat{J}^{ct} w.r.t. \tilde{K}^S , $\tilde{K}^{\tau,S}$, and K^τ is

$$\frac{\partial \hat{J}^{ct}}{\partial \tilde{K}^S} = \frac{1}{n^S \sigma^2} \beta \beta^\top, \quad \frac{\partial \hat{J}^{ct}}{\partial \tilde{K}^{\tau,S}} = -\frac{2}{n^S n^\tau} \mathbf{1} \beta^\top, \quad \text{and} \quad \frac{\partial \hat{J}^{ct}}{\partial K^\tau} = \frac{1}{n^{\tau^2}} \mathbf{1} \mathbf{1}^\top.$$

The gradient of $\text{Tr}[\hat{\Sigma}_{YY|X^{ct}}]$ w.r.t. \tilde{K}^S is

$$\frac{\partial \text{Tr}[\hat{\Sigma}_{YY|X^{ct}}]}{\partial \tilde{K}^S} = -\varepsilon (\tilde{K}^S + n^S \varepsilon I)^{-1} L (\tilde{K}^S + n^S \varepsilon I)^{-1}.$$

Using the chain rule, we further have the gradient of \hat{J}^{ct} w.r.t. the entries of W , \mathbf{G} , and \mathbf{H} :

$$\frac{\partial \hat{J}^{ct}}{\partial W_{pq}} = \text{Tr} \left[\left(\frac{\partial \hat{J}^{ct}}{\partial \tilde{K}^S} \right)^\top (\mathbf{D}_{pq}^1) \right] - \text{Tr} \left[\left(\frac{\partial \hat{J}^{ct}}{\partial \tilde{K}^{\tau,S}} \right)^\top (\mathbf{D}_{pq}^2) \right] + \text{Tr} \left[\left(\frac{\partial \hat{J}^{ct}}{\partial K^\tau} \right)^\top (\mathbf{D}_{pq}^3) \right], \quad (43)$$

$$\frac{\partial \hat{J}^{ct}}{\partial \mathbf{G}_{pq}} = \text{Tr} \left[\left(\frac{\partial \hat{J}^{ct}}{\partial \tilde{K}^S} \right)^\top (\mathbf{E}_{pq}^1) \right] - \text{Tr} \left[\left(\frac{\partial \hat{J}^{ct}}{\partial \tilde{K}^{\tau,S}} \right)^\top (\mathbf{E}_{pq}^2) \right], \quad (44)$$

$$\frac{\partial \hat{J}^{ct}}{\partial \mathbf{H}_{pq}} = \text{Tr} \left[\left(\frac{\partial \hat{J}^{ct}}{\partial \tilde{K}^S} \right)^\top (\mathbf{F}_{pq}^1) \right] - \text{Tr} \left[\left(\frac{\partial \hat{J}^{ct}}{\partial \tilde{K}^{\tau,S}} \right)^\top (\mathbf{F}_{pq}^2) \right], \quad (45)$$

and the gradient of $\text{Tr}[\hat{\Sigma}_{YY|X^{ct}}]$ w.r.t. the entries of W , \mathbf{G} , and \mathbf{H} :

$$\frac{\partial \text{Tr}[\hat{\Sigma}_{YY|X^{ct}}]}{\partial W_{pq}} = \text{Tr} \left[\left(\frac{\partial \text{Tr}[\hat{\Sigma}_{YY|X^{ct}}]}{\partial \tilde{K}^S} \right)^\top (\mathbf{D}_{pq}^1) \right], \quad (46)$$

$$\frac{\partial \text{Tr}[\hat{\Sigma}_{YY|X^{ct}}]}{\partial \mathbf{G}_{pq}} = \text{Tr} \left[\left(\frac{\partial \text{Tr}[\hat{\Sigma}_{YY|X^{ct}}]}{\partial \tilde{K}^S} \right)^\top (\mathbf{E}_{pq}^1) \right], \quad (47)$$

$$\frac{\partial \text{Tr}[\hat{\Sigma}_{YY|X^{ct}}]}{\partial \mathbf{H}_{pq}} = \text{Tr} \left[\left(\frac{\partial \text{Tr}[\hat{\Sigma}_{YY|X^{ct}}]}{\partial \tilde{K}^S} \right)^\top (\mathbf{F}_{pq}^1) \right], \quad (48)$$

where

$$[\mathbf{D}_{pq}^1]_{ij} = -\frac{\tilde{k}^s(x_i^s, x_j^s)}{\sigma^2} \left[\sum_{k=1}^D w_{kq} (a_{qi} x_{ik}^s - a_{qj} x_{jk}^s) (a_{qi} x_{ip}^s - a_{qj} x_{jp}^s) + (a_{qi} x_{ip}^s - a_{qj} x_{jp}^s) (b_{qi} - b_{qj}) \right],$$

$$[\mathbf{D}_{pq}^2]_{ij} = -\frac{\tilde{k}^{t,s}(x_i^t, x_j^s)}{\sigma^2} \left[\sum_{k=1}^D w_{kq} (x_{ik}^t - a_{qj} x_{jk}^s) (x_{ip}^t - a_{qj} x_{jp}^s) + a_{qj} x_{jp}^s b_{qj} \right],$$

$$[\mathbf{D}_{pq}^3]_{ij} = -\frac{\tilde{k}^t(x_i^t, x_j^t)}{\sigma^2} \left[\sum_{k=1}^D w_{kq} (x_{ik}^t - x_{jk}^t) (x_{ip}^t - x_{jp}^t) \right],$$

$$[\mathbf{E}_{pq}^1]_{ij} = -\frac{\tilde{k}^s(x_i^s, x_j^s)}{\sigma^2} (x_{jq}^{ct} - x_{iq}^{ct}) (x_{jq}^s R_{jp}^{dis} - x_{iq}^s R_{ip}^{dis}),$$

$$[\mathbf{E}_{pq}^2]_{ij} = -\frac{\tilde{k}^{t,s}(x_i^t, x_j^s)}{\sigma^2} x_{jq}^s R_{jp}^{dis} (x_{jq}^{ct} - x_{iq}^t),$$

$$[\mathbf{F}^1_{pq}]_{ij} = -\frac{\tilde{k}^s(x_i^s, x_j^s)}{\sigma^2}(x_{jq}^{ct} - x_{iq}^{ct})(R_{jp}^{dis} - R_{ip}^{dis}),$$

$$[\mathbf{F}^2_{pq}]_{ij} = -\frac{\tilde{k}^{t,s}(x_i^t, x_j^s)}{\sigma^2}R_{jp}^{dis}(x_{jq}^{ct} - x_{iq}^t).$$

The derivative of J^{reg} w.r.t. \mathbf{G} and \mathbf{H} is

$$\frac{\partial J^{reg}}{\partial \mathbf{G}} = \frac{2\lambda_S}{n^S}R^{dis\top}(\mathbf{A}^\top - \mathbf{1}_{n^S \times d}), \text{ and}$$

$$\frac{\partial J^{reg}}{\partial \mathbf{H}} = \frac{2\lambda_L}{n^S}R^{dis\top}\mathbf{B}^\top.$$