

## A Proofs

### A.1 Proof of Theorem 1

*Proof of Theorem 1.* The proof just uses assumptions on the transformation function and stability of the training algorithm.

$$\begin{aligned} & \left| \hat{f}^{ta}(x) - f^{ta}(x) \right|^2 \\ &= \left| G\left(\hat{f}^{so}(x), \hat{w}_G(x)\right) - G\left(f^{so}(x), w_G(x)\right) \right|^2 \end{aligned} \quad (2)$$

$$\leq L^2 \left| \hat{f}^{so}(x) - f^{so}(x) \right|^2 + L^2 \left| \hat{w}_G(x) - w_G(x) \right|^2 \quad (3)$$

$$\leq L^2 \left| \hat{f}^{so}(x) - f^{so}(x) \right|^2 + 2L^2 \left| \hat{w}_G(x) - \tilde{w}_G(x) \right|^2 + 2L^2 \left| \tilde{w}_G(x) - w_G(x) \right|^2 \quad (4)$$

$$\leq L^2 \left| \hat{f}^{so}(x) - f^{so}(x) \right|^2 + 2L^2 \left( \sum_{i=1}^{n_{ta}} c_i(X_i^{ta}) \left| W_i - \tilde{W}_i \right| \right)^2 + 2L^2 \left| \tilde{w}_G(x) - w_G(x) \right|^2 \quad (5)$$

where (2) is by the requirement of  $G$ , (3) is by the Lipschitz condition of  $G$ , (4) is because  $(a-b)^2 \leq 2(a-c)^2 + 2(c-b)^2$  and (5) is by our stability assumption of  $\mathcal{A}_{w_G}$ . Now, we are left bounding

$$\begin{aligned} & \left( \sum_{i=1}^{n_{ta}} c_i \left| W_i - \tilde{W}_i \right| \right)^2. \text{ Notice that by the assumption of } H_G, \\ & \left| W_i - \tilde{W}_i \right| = \left| H_G\left(\hat{f}^{so}(X_i^{ta}), (Y_i^{ta})\right) - H_G\left(f^{so}(X_i^{ta}), Y_i^{ta}\right) \right| \leq L \left| \hat{f}^{so}(X_i^{ta}) - f^{so}(X_i^{ta}) \right| \end{aligned} \quad (6)$$

Plugging (6) into (5), we obtain our desired result.  $\square$

### A.2 Proof of Theorem 2

For simplicity, let  $K_h(\cdot) = K(\cdot/h)$  and define the expected regression estimate  $\tilde{f} = \sum_{i=1}^n w_i f(X_i)$ . To prove Theorem 2, we first give some standard supporting lemmas for kernel smoothing.

**Lemma 1 (Lemma 1 of [Kpotufe and Garg, 2013])** *Under the same assumptions in Theorem 2, for all  $x$  with  $\|x\|_2 \leq \Delta_X$ , if  $f$  is  $(\lambda, \alpha)$  Hölder, then, for any  $h > 0$ , we have  $|\tilde{f}(x) - f(x)|^2 \leq \lambda^2 h^{2\alpha}$ .*

**Lemma 2 (Corollary of Lemma 3 and Lemma 7 of [Kpotufe and Garg, 2013])** *Under the same assumptions in Theorem 2, let  $0 < \delta < 1/6$ , for all  $x : \|x\|_2 \leq \Delta_X$  and  $h > 0$ , with probability at least  $1 - \delta$ , we have*

$$|\hat{f}(x) - \tilde{f}(x)|^2 = O\left(\frac{\log(1/\delta)}{nh^d}\right).$$

*Proof of Theorem 2.* we prove Theorem 2 by bounding each corresponding term in Theorem 1. First, by Lemma 1 and Lemma 2, we have for all  $x$ , with probability at least  $1 - \delta$

$$\left| \hat{f}^{so}(x) - f^{so}(x) \right|^2 = O\left(h_{so}^{2\alpha_{so}} + \frac{\log(1/\delta)}{n_{so}h_{so}^d}\right).$$

Specifically, for  $X_1^{ta}, \dots, X_{n_{ta}}^{ta}$ , we have

$$\max_{i=1, \dots, n_{ta}} \left| \hat{f}^{so}(X_i^{ta}) - f^{so}(X_i^{ta}) \right|^2 = O\left(h_{so}^{2\alpha_{so}} + \frac{\log(1/\delta)}{n_{so}h_{so}^d}\right). \quad (7)$$

Next, according to Assumption 1 and 2,  $H_G$  is bounded and unbiased and  $w_G$  is bounded, we can view  $\left\{ \left( X_i^{ta}, \tilde{W}_i \right) \right\}_{i=1}^{n_{ta}}$  a training set for function  $w_G$  that  $\tilde{W}_i = w_G(X_i^{ta}) + \epsilon_{w_G}$  where  $\mathbf{E}[\epsilon_{w_G}] = 0$  and  $|\epsilon_{w_G}| \leq 2B$ . Based on this observation, using Lemma 1 and Lemma 2 again, for all  $x : \|x\|_2 \leq \Delta_X$ , we have with probability at least  $1 - \delta$

$$\left| \tilde{w}_G(x) - w_G(x) \right|^2 = O\left(h_{w_G}^{2\alpha_{w_G}} + \frac{\log(1/\delta)}{n_{ta}h_{w_G}^d}\right).$$

Now we are left bounding  $\left\| \mathcal{A}_{w_G}(\mathcal{T}) - \mathcal{A}_{w_G}(\tilde{\mathcal{T}}) \right\|_{\infty}$ . Notice that for  $\mathcal{T}, \tilde{\mathcal{T}}$  in Theorem 1, and for all  $x : \|x\|_2 \leq \Delta_X$ :

$$\begin{aligned} \left| \mathcal{A}_{w_G}(\mathcal{T})(x) - \mathcal{A}_{w_G}(\tilde{\mathcal{T}})(x) \right| &= \frac{\sum_{i=1}^{n_{ta}} K_h(\|x - X_i^{ta}\|_2) (W_i - \tilde{W}_i)}{\sum_{i=1}^{n_{ta}} K_h(\|x - X_i^{ta}\|_2)} \\ &= \frac{\sum_{i=1}^{n_{ta}} K_h(\|x - X_i^{ta}\|_2) |W_i - \tilde{W}_i|}{\sum_{i=1}^{n_{ta}} K_h(\|x - X_i^{ta}\|_2)} \\ &\triangleq \sum_{i=1}^{n_{ta}} c_i |W_i - \tilde{W}_i| \end{aligned}$$

for  $c_i = \frac{K_h(\|x - X_i^{ta}\|_2)}{\sum_{i=1}^{n_{ta}} K_h(\|x - X_i^{ta}\|_2)}$ . Now according to Theorem 1, we only need to bound  $\left( \sum_{i=1}^{n_{ta}} c_i \left| \hat{f}^{so}(X_i^{ta}) - f^{so}(X_i^{ta}) \right| \right)^2$ . With probability at least  $1 - \delta$ , we have:

$$\left( \sum_{i=1}^{n_{ta}} c_i \left| \hat{f}^{so}(X_i^{ta}) - f^{so}(X_i^{ta}) \right| \right)^2 \leq \left( \sum_{i=1}^{n_{ta}} c_i \right)^2 \left( \max_{i=1, \dots, n_{ta}} \left| \hat{f}^{so}(X_i^{ta}) - f^{so}(X_i^{ta}) \right|^2 \right) \quad (8)$$

$$= \max_{i=1, \dots, n_{ta}} \left| \hat{f}^{so}(X_i^{ta}) - f^{so}(X_i^{ta}) \right|^2 \quad (9)$$

$$= O \left( h_{so}^{2\alpha_{so}} + \frac{\log(n_{ta}/\delta)}{n_{so} h_{so}^d} \right), \quad (10)$$

where (8) is because maximum is bigger than other terms, (9) is because  $\sum_{i=1}^{n_{ta}} c_i = 1$  by definition, and (10) is by (7). Putting these all together, using Theorem 1 and choosing the bandwidth according to Theorem 2, we can show for all  $x : \|x\|_2 \leq \Delta_X$

$$\left| f^{ta}(x) - \hat{f}^{ta}(x) \right|^2 = O \left( n_{so}^{\frac{-2\alpha_{so}}{2\alpha_{so}+d}} + n_{ta}^{\frac{-2\alpha_{w_G}}{2\alpha_{w_G}+d}} \right) \log \left( \frac{1}{\delta} \right).$$

Now integrate with respect to  $P_{X^{ta}}$  we obtain our desired result.  $\square$

### A.3 Proof of Theorem 3

The proof strategy is similar to that of Theorem 2. Using Theorem 1 we have

$$\begin{aligned} \mathbf{E} \left[ \left| \hat{f}^{ta}(X) - f^{ta}(X) \right|^2 \right] &= O \left( \mathbf{E} \left[ \left| \hat{f}^{so}(X) - f^{so}(X) \right|^2 + |\tilde{w}_G(X) - w_G(X)|^2 + \right. \right. \\ &\quad \left. \left. \left( \sum_{i=1}^{n_{ta}} c_i(X_i^{ta}) \left| \hat{f}^{so}(X_i^{ta}) - f^{so}(X_i^{ta}) \right| \right)^2 \right] \right). \end{aligned}$$

where the expectation is taken over  $P_{x^{ta}}$  and  $\mathcal{T}^{ta}$ . Now we bound three terms on the right hand side separately. By Corollary 3 of Steinwart et al. [2009], we have with probability at least  $1 - \delta$

$$\mathbf{E} \left[ \left| \hat{f}^{so}(X) - f^{so}(X) \right|^2 \right] = O \left( \lambda_{so}^{\beta_{so}} + \frac{\log(1/\delta)}{\lambda_{so}^p n_{so}} \right), \quad (11)$$

where expectation is taken over  $P_x^{ta}$ . Taking union bound over  $X_1^{ta}, \dots, X_{n_{ta}}^{ta}$ , we have

$$\max_{i=1, \dots, n_{ta}} \mathbf{E} \left[ \left| \hat{f}^{so}(X_i^{ta}) - f^{so}(X_i^{ta}) \right|^2 \right] = O \left( \lambda_{so}^{\beta_{so}} + \frac{\log(n_{ta}/\delta)}{\lambda_{so}^p n_{so}} \right). \quad (12)$$

where the expectation is taken over  $\mathcal{T}^{ta}$ . Next, using the exactly same argument as in the Theorem 2, we can view  $\left\{ (X_i^{ta}, \tilde{W}_i) \right\}_{i=1}^{n_{ta}}$  a training set for function  $w_G$  that  $\tilde{W}_i = w_G(X_i^{ta}) + \epsilon_{w_G}$  as

$\widetilde{W}_i = w_G(X_i^{ta}) + \epsilon_{w_G}$  where  $\mathbf{E}[\epsilon_{w_G}] = 0$  and  $|\epsilon_{w_G}| \leq 2B$ . Thus applying Corollary 3 of Steinwart et al. [2009] again, we have with probability at least  $1 - \delta$

$$\mathbf{E} \left[ |\widetilde{w}_G(X) - w_G(X)|^2 \right] = O \left( \lambda_{w_G}^{\beta_{w_G}} + \frac{\log(1/\delta)}{\lambda_{w_G}^p n_{ta}} \right).$$

where expectation is taken over  $P_{x^{ta}}$ . Now we analyze the stability of KRR. We use  $\Phi(x)$  to denotes the feature map corresponding with the given kernel  $K$  so  $K(x, y) = \Phi(x)^\top \Phi(y)$ . Also for simplicity, we denote

$$\Phi_{ta} = (\Phi(x_1^{ta}) \mid \cdots \mid \Phi(x_{n_{ta}}^{ta}))$$

the feature matrix of target domain data. With these notations, we can write

$$\begin{aligned} & \left| \mathcal{A}_{w_G}(\mathcal{T}^{w_G})(x) - \mathcal{A}_{w_G}(\widetilde{\mathcal{T}}^{w_G})(x) \right| \\ &= \left| \begin{pmatrix} W_1 - \widetilde{W}_1 \\ \cdots \\ W_{n_{ta}} - \widetilde{W}_{n_{ta}} \end{pmatrix}^\top (\Phi_{ta}^\top \Phi + n_{ta} \lambda_{w_G} \mathbf{I})^{-1} \Phi_{ta}^\top \Phi(x) \right| \\ &= \left| \left( \Phi_{ta} \begin{pmatrix} W_1 - \widetilde{W}_1 \\ \cdots \\ W_{n_{ta}} - \widetilde{W}_{n_{ta}} \end{pmatrix} \right)^\top (\Phi_{ta} \Phi_{ta}^\top + n_{ta} \lambda_{w_G} \mathbf{I})^{-1} \Phi(x) \right| \\ &\leq \left\| \begin{pmatrix} k^{1/2} |W_1 - \widetilde{W}_1| \\ \cdots \\ k^{1/2} |W_{n_{ta}} - \widetilde{W}_{n_{ta}}| \end{pmatrix} \right\|_2 \left\| (\Phi_{ta} \Phi_{ta}^\top + n_{ta} \lambda_{w_G} \mathbf{I})^{-1} \right\|_{op} k^{1/2} \\ &\leq \sum_{i=1}^{n_{ta}} \frac{k}{n_{ta} \lambda_{w_G}} |W_i - \widetilde{W}_i| \\ &\triangleq \sum_{i=1}^{n_{ta}} c_i |W_i - \widetilde{W}_i|. \end{aligned}$$

The second equality we used the identity that  $(\Phi^\top \Phi + \lambda \mathbf{I})^{-1} \Phi^\top = \Phi^\top (\Phi \Phi^\top + \lambda \mathbf{I})^{-1}$  for any  $\Phi$  and  $\lambda$ . The first inequality we used sub-multiplicity of operator norm and the assumption  $\|\Phi(x)\|_{\mathcal{H}} \leq k^{1/2}$ . The second inequality we used the fact the lower bound of least eigenvalue of  $(\Phi_{ta} \Phi_{ta}^\top + n_{ta} \lambda_{w_G} \mathbf{I})$  is  $n_{ta} \lambda_{w_G}$ . Therefore, applying Cauchy-Schwartz inequality and using the bound in (12), we have with probability at least  $1 - \delta$ ,

$$\begin{aligned} \mathbf{E} \left[ \left( \sum_{i=1}^{n_{ta}} c_i |\hat{f}^{so}(X_i^{ta}) - f^{so}(X_i^{ta})| \right)^2 \right] &\leq \left( \sum_{i=1}^{n_{ta}} c_i^2 \right) \cdot \left( \sum_{i=1}^{n_{ta}} |\hat{f}^{so}(X_i^{ta}) - f^{so}(X_i^{ta})|^2 \right) \\ &= \sum_{i=1}^{n_{ta}} \frac{k^2}{n_{ta}^2 \lambda_{w_G}^2} \cdot \mathbf{E} \left[ \sum_{i=1}^{n_{ta}} |\hat{f}^{so}(X_i^{ta}) - f^{so}(X_i^{ta})|^2 \right] \\ &\leq \frac{k^2}{\lambda_{w_G}^2} \cdot \max_{i=1, \dots, n_{ta}} \mathbf{E} \left[ |\hat{f}^{so}(X_i^{ta}) - f^{so}(X_i^{ta})|^2 \right] \\ &= O \left( \frac{k^2}{\lambda_{w_G}^2} \left( \lambda_{so}^{\beta_{so}} + \frac{\log(n_{ta}/\delta)}{\lambda_{so}^p n_{so}} \right) \right). \end{aligned}$$

Now putting these all together and choosing  $\lambda_{so}$  and  $\lambda_{w_G}$  according to Theorem 3, we obtain the desired result.  $\square$

#### A.4 Proof of Theorem 4

We first prove a general theorem for cross-validation. This is a standard result for cross-validation and we include the proof for completeness.

**Theorem 5** Let  $\Theta$  be the set of all hypotheses and  $\hat{\theta} = \operatorname{argmin}_{\theta \in \Theta} \sum_{i=1}^{n_{val}} \left( \hat{f}_{\theta}^{ta}(X_i^{val}) - Y_i^{val} \right)^2$  the estimator that minimizes error on the cross-validation set. Then with probability at least  $1 - \delta$ :

$$\mathbf{E} \left[ R \left( \hat{f}_{\hat{\theta}}^{ta} \right) \right] - R \left( f^{ta} \right) = O \left( \mathbf{E} \left[ R \left( \hat{f}_{\theta^*}^{ta} \right) \right] - R \left( f^{ta} \right) + \frac{\log \frac{|\Theta|}{\delta}}{n_{val}} \right),$$

where  $\theta^* = \operatorname{argmin}_{\theta \in \Theta} R \left( \hat{f}_{\theta} \right)$  and the expectation is taken over  $\mathcal{T}^{so}$  and  $\mathcal{T}^{ta}$ .

To prove of Theorem 5, we use the following type of Bernstein's inequality [Craig, 1933]:

**Lemma 3** Let  $X_1, \dots, X_n$  be random variables and suppose that for  $k \geq 3$ :

$$\mathbf{E} [|X_i - \mathbf{E}[X_i]|^k] \leq \frac{\mathbf{Var}[X_i]}{2} k! r^{k-2},$$

for some  $r > 0$ . Then with probability  $> 1 - \delta$ :

$$\frac{1}{n} \sum_{i=1}^n (X_i - \mathbf{E}[X_i]) \leq \frac{\log(1/\delta)}{nt} + \frac{t \mathbf{Var}[X_i]}{2(1-c)},$$

for  $0 \leq tr \leq c < 1$ .

*Proof of Theorem 5:* For a given  $\theta \in \Theta$ , we obtain a corresponding estimated regression function  $\hat{f}_{\theta}$ .

Define  $U_i^{\theta} \triangleq - \left( Y_i^{val} - \hat{f}_{\theta}^{ta}(X_i^{val}) \right)^2 + (Y_i - f^{ta}(X_i^{val}))^2$ . Compute the expectation:

$$\begin{aligned} \mathbf{E} [U_i^{\theta}] &= - \mathbf{E} \left[ -2Y_i^{val} \hat{f}_{\theta}^{ta}(X_i^{val}) + \hat{f}_{\theta}^{ta}(X_i^{val})^2 + 2Y_i^{val} f^{ta}(X_i^{val}) - f^{ta}(X_i^{val})^2 \right] \\ &= - \mathbf{E} \left[ \left( \hat{f}_{\theta}^{ta}(X_i^{val}) - f^{ta}(X_i^{val}) \right)^2 \right] \\ &= R(f^{ta}) - R(\hat{f}_{\theta}^{ta}). \end{aligned}$$

Also, by definition, it is easy to see

$$\frac{1}{n_{val}} \sum_{i=1}^{n_{val}} U_i^{\theta} = \hat{R}(f) - \hat{R}(\hat{f}_{\theta}^{ta}).$$

In order to apply Bernstein's inequality, we must first bound the variance of  $U_i^{\theta}$ :

$$\begin{aligned} \mathbf{var} [U_i^{\theta}] &\leq \mathbf{E} [(U_i^{\theta})^2] \\ &= \mathbf{E} \left[ \left( - \left( Y_i^{val} - \hat{f}_{\theta}^{ta}(X_i^{val}) \right)^2 + \left( Y_i^{val} - f^{ta}(X_i^{val}) \right)^2 \right)^2 \right] \\ &= \mathbf{E} \left[ \left( f^{ta}(X_i) - \hat{f}_{\theta}^{ta} \right)^4 + 4\epsilon_i \left( f^{ta}(X_i^{val}) - \hat{f}_{\theta}^{ta}(X_i^{val}) \right)^3 + 4\epsilon_i^2 \left( f^{ta}(X_i^{val}) - \hat{f}_{\theta}^{ta}(X_i^{val}) \right)^2 \right] \\ &\leq -4\Delta_Y^2 \mathbf{E} [U_i] \end{aligned}$$

where in the last inequality we used the domain of  $Y$  is bounded. Since  $U_i$  is a sum of bounded random variables, the moment condition is satisfied with  $r = 4\Delta_Y^2$ . Now apply Craig-Bernstein inequality to  $U_i^{\theta}$ s, with probability at least  $1 - \delta$ :

$$\frac{1}{n_{val}} \sum_{i=1}^{n_{val}} U_i^{\theta} - \mathbf{E} [U_i^{\theta}] \leq \frac{\log(1/\delta)}{n_{val}t} + \frac{-2t\Delta_Y^2 \mathbf{E} [U_i^{\theta}]}{1-c}.$$

We need to ensure that  $c < 1$ . To do this, let  $c = tr = 4t\Delta_Y^2$  and let  $t < \frac{1}{6\Delta_Y^2}$ , then it is easy to see that  $c < 1$ . For simplicity, define  $a = \frac{2t\Delta_Y^2}{1-c} < 1$ . Now grouping terms we get:

$$\begin{aligned} (1-a)(-\mathbf{E}[U_i^\theta]) + \frac{1}{n_{val}} \sum_{i=1}^{n_{val}} U_i^\theta &\leq \frac{\log(1/\delta)}{n_{val}t} \\ (1-a)(R(\hat{f}_\theta^{ta}) - R(f)) - (\hat{R}(\hat{f}_\theta^{ta}) - \hat{R}(f)) &\leq \frac{\log(1/\delta)}{n_{val}t} \\ R(\hat{f}_\theta^{ta}) - R(f^{ta}) &\leq \frac{1}{1-a} \left( \hat{R}(\hat{f}_\theta) - \hat{R}(f^{ta}) + \frac{\log(1/\delta)}{n_{val}t} \right). \end{aligned}$$

Take union bound over  $\Theta$ , and consider  $\hat{f}_\theta$ :

$$R(\hat{f}_\theta^{ta}) - R(f^{ta}) \leq \frac{1}{1-a} \left( \hat{R}(\hat{f}_\theta^{ta}) - \hat{R}(f^{ta}) + \frac{\log(|\Theta|/\delta)}{n_{val}t} \right).$$

Now, recall that  $\hat{f}_\theta^{ta}$  is the minimizer for  $\hat{R}$  among all estimators induced by  $\Theta$ , we have

$$R(\hat{f}_\theta^{ta}) - R(f^{ta}) \leq \frac{1}{1-a} \left( \hat{R}(\hat{f}_{\theta^*}^{ta}) - \hat{R}(f^{ta}) + \frac{\log(|\Theta|/\delta)}{n_{val}t} \right).$$

Now taking expectation over  $\mathcal{T}^{val}$  then over  $\mathcal{T}^{so}$  and  $\mathcal{T}^{ta}$  we obtain the desired result.  $\square$

Now we are ready to prove Theorem 4. Since  $\bar{\mathcal{G}}$  is an  $\epsilon$ -cover of  $\mathcal{G}$ , there exists  $G' \in \bar{\mathcal{G}}$  such that  $\|G' - G^*\|_\infty \leq \epsilon$ . For any  $x$ ,

$$\begin{aligned} &\left| f^{ta}(x) - \hat{f}_{G'}^{ta}(x) \right| \\ &= \left| G^*(f^{so}(x), w_{G^*}(x)) - G'(\hat{f}^{so}(x), \hat{w}_{G'}(x)) \right| \\ &\leq \left| G^*(f^{so}(x), w_{G^*}(x)) - G^*(\hat{f}^{so}(x), \hat{w}_{G^*}(x)) \right| + \left| G^*(\hat{f}^{so}(x), \hat{w}_{G^*}(x)) - G'(\hat{f}^{so}(x), \hat{w}_{G^*}(x)) \right| \\ &\quad + \left| G'(\hat{f}^{so}(x), \hat{w}_{G^*}(x)) - G'(\hat{f}^{so}(x), \hat{w}_{G'}(x)) \right| \end{aligned} \quad (13)$$

where  $\hat{w}_{G^*} = \mathcal{A}_{w_G}(\{X_i^{ta}, W_i^*\})$  and  $W_i^* = H_{G'}(\hat{f}^{so}(X_i^*), Y_i^*) + w_{G^*}(X_i^{ta}) - w_{G'}(X_i^{ta})$ , i.e. an un-biased estimated of  $w_{G^*}(X_i^*)$ . We can bound three terms in (13) separately. The first term is just the difference between estimator based on  $G^*$  and the true  $f^{ta}$ , so after taking expectation it becomes the excess risk of  $\hat{f}_{G^*}^{ta}$ . By our construction of  $\epsilon$ -cover of  $\mathcal{G}$ , the second term is smaller than  $\epsilon$ . For the third term, notice that by Lipschitz assumption on  $G$ s and our assumptions on  $G$ s in  $\mathcal{G}$  in the theorem 4, we have:

$$\begin{aligned} &\left| G'(\hat{f}^{so}(x), \hat{w}_{G^*}(x)) - G'(\hat{f}^{so}(x), \hat{w}_{G'}(x)) \right| \\ &\leq L(|\hat{w}_{G^*}(x) - \hat{w}_{G'}(x)|) \\ &\leq L^2 \sum_{i=1}^{n_{ta}} c_i \|G^* - G'\|_\infty \\ &= O\left(\sum_{i=1}^{n_{ta}} c_i \epsilon\right). \end{aligned}$$

Now we have shown  $R(\hat{f}_{G'}^{ta}) - R(f^{ta}) = O\left(R(\hat{f}_{G^*}^{ta}) - R(f^{ta})\right)$ . Let  $\bar{G}_* = \operatorname{argmin}_{G \in \bar{\mathcal{G}}} R(\hat{f}_G)$ , the best transformation function in  $\bar{\mathcal{G}}$ . By the optimality of  $\bar{G}_*$ , we have  $R(\hat{f}_{\bar{G}_*}^{ta}) - R(f^{ta}) = O\left(R(\hat{f}_{G^*}^{ta}) - R(f^{ta})\right)$ . Applying Theorem 5 with our assumptions on  $\epsilon$  and  $n_{val}$  we know  $R(\hat{f}_{\bar{G}_*}^{ta}) - R(f^{ta}) = O\left(R(\hat{f}_{\bar{G}_*}^{ta}) - R(f^{ta})\right)$ . Combing these facts we have  $R(\hat{f}_{\bar{G}_*}^{ta}) - R(f^{ta}) = O\left(R(\hat{f}_{G^*}^{ta}) - R(f^{ta})\right)$ .

## B Regression Calibration for Measurement Error Problem

Given,  $f^{so}$ , in this section we provide a standard technique to obtain an unbiased estimate of  $w_G(X_i^{ta})$ s. Since we assume

$$Y^{ta} = f^{ta}(X^{ta}) + \epsilon^{ta},$$

the measurement error model corresponds to *classical error model* in Carroll et al. [2006]. Regression calibration is a widely used and reasonably well investigated method for measurement error problem. The algorithm is as follows (we have adapted the general algorithm to our HTL problem):

- Compute an estimate of  $f^{ta}(X_i^{ta})$ :  $\tilde{f}^{ta}(X_i^{ta})$ . Note that directly using  $Y_i^{ta}$  is one of the option for  $\tilde{f}^{ta}(X_i^{ta})$ .
- Compute  $G_{f^{so}(X_i^{ta})}^{-1}(\tilde{f}^{ta}(X_i^{ta}))$ .
- Calibrate our previous computed value by applying some function  $F$ :

$$\tilde{W}_i = F\left(G_{f^{so}(X_i^{ta})}^{-1}(\tilde{f}^{ta}(X_i^{ta}))\right)$$

where  $F$  depends on  $G$  and the specific distribution on noise.

Now we consider the loglinear mean model as a concrete example. Suppose

$$G(f^{so}(x), w_G(x)) = \beta f^{so}(x) \log(w_G(x))$$

where  $\beta$  is some constant. Further, we assume  $\epsilon^{ta} \sim \mathcal{N}(0, \sigma^2)$  Now we apply the regression calibration algorithm.

- First we choose  $Y_i^{ta}$  as our estimate for  $\tilde{f}^{ta}(X_i^{ta})$ .
- Second, by our choice of  $G$ :

$$G_{f^{so}(X_i^{ta})}^{-1}(Y_i^{ta}) = \exp\left(\frac{Y_i^{ta}}{\beta f^{so}(X_i^{ta})}\right)$$

- Last, for our choice of  $G$  and assumption of  $\epsilon^{ta}$ , the corresponding  $F$  and final estimate of  $w_G(X_i^{ta})$  is

$$\begin{aligned} \tilde{W}_i &= F\left(G_{f^{so}(X_i^{ta})}^{-1}(\tilde{f}^{ta}(X_i^{ta}))\right) \\ &= \exp\left(\log\left(G_{f^{so}(X_i^{ta})}^{-1}(\tilde{f}^{ta}(X_i^{ta}))\right) + \sigma^2 (f^{so}(X_i^{ta}))^2\right) \\ &= \exp\left(\frac{Y_i^{ta}}{\beta f^{so}(X_i^{ta})} + \sigma^2 (f^{so}(X_i^{ta}))^2\right). \end{aligned}$$

The estimator for  $w_G(X_i^{ta})$  depends on some distribution specific parameters which may be unknown, like  $\sigma^2$  in the previous example. In such cases, we may replace these parameters by our estimates. For example, in the previous Gaussian noise case, suppose for each  $X_i^{ta}$ , we have multiple observations  $\{Y_{ij}\}_{j=1}^{n_i}$ . Then we can estimate  $\sigma^2$  by

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^{n_{ta}} \sum_{j=1}^{n_i} (Y_{ji}^{ta} - \bar{Y}_i^{ta})^2}{\sum_{i=1}^{n_{ta}} (n_i - 1)}$$

where  $\bar{Y}_i^{ta} = \frac{\sum_{j=1}^{n_i} Y_{ij}}{n_i}$ .

Here we only provide one method for measurement error problem. There are other techniques such simulation extrapolation and likelihood method which may be also applicable in many situations. The choice of method depends on specific transformation  $G$  and assumptions on the distribution of the noise. Again, interested readers are referred to Carroll et al. [2006] for details.

	$n_{ta} = 10$	$n_{ta} = 20$	$n_{ta} = 40$	$n_{ta} = 80$	$n_{ta} = 160$	$n_{ta} = 320$
Only Target KS	0.005 ± 0.001	0.003 ± 0.001	0.003 ± 0.001	0.003 ± 0.000	0.002 ± 0.000	0.002 ± 0.000
Only Target KRR	<b>0.001 ± 0.001</b>	<b>0.001 ± 0.000</b>	<b>0.000 ± 0.000</b>	<b>0.000 ± 0.000</b>	<b>0.000 ± 0.000</b>	<b>0.000 ± 0.000</b>
Only Source KS	0.031 ± 0.012	0.031 ± 0.012	0.031 ± 0.012	0.031 ± 0.012	0.031 ± 0.012	0.031 ± 0.012
Only Source KRR	0.016 ± 0.013	0.016 ± 0.013	0.016 ± 0.013	0.016 ± 0.013	0.016 ± 0.013	0.016 ± 0.013
Combined KS	0.023 ± 0.017	0.029 ± 0.011	0.017 ± 0.013	0.007 ± 0.007	0.002 ± 0.000	0.002 ± 0.000
Combined KRR	0.006 ± 0.008	0.009 ± 0.010	0.002 ± 0.002	0.001 ± 0.000	0.001 ± 0.000	0.001 ± 0.000
CDM	0.004 ± 0.002	0.007 ± 0.001	0.004 ± 0.002	0.001 ± 0.000	0.001 ± 0.000	0.012 ± 0.002
Offset KS	0.003 ± 0.001	0.002 ± 0.001	0.002 ± 0.000	0.002 ± 0.000	0.002 ± 0.000	0.001 ± 0.000
Offset KRR	0.002 ± 0.001	<b>0.001 ± 0.000</b>	<b>0.000 ± 0.000</b>	<b>0.000 ± 0.000</b>	<b>0.000 ± 0.000</b>	<b>0.000 ± 0.000</b>
Scale KS	0.004 ± 0.002	0.003 ± 0.001	0.002 ± 0.001	0.002 ± 0.000	0.002 ± 0.000	0.002 ± 0.000
Scale KRR	<b>0.001 ± 0.000</b>	<b>0.001 ± 0.000</b>	<b>0.000 ± 0.000</b>	<b>0.000 ± 0.000</b>	<b>0.000 ± 0.000</b>	<b>0.000 ± 0.000</b>

Table 3: 1 standard deviation intervals for the mean squared errors of various algorithms when transferring from kin-8nh to kin-8fm. The values in bold are the best errors for each  $n_{ta}$ .

## C Additional Experimental Results

### C.1 Synthetic data

This section gives details of the synthetic data. For both experiments, we use  $n_{so} = 10000$  samples from the source domain, and  $n_{ta} = 100$  samples from the target domain. We put Gaussian noise on the labels:  $\epsilon^{so} \sim \mathcal{N}(0, 0.01)$ ,  $\epsilon^{ta} \sim \mathcal{N}(0, 0.01)$ ; and we use KS with a gaussian kernel for estimating  $f^{so}$  and  $w_G$ .

Figure 1b shows the offset example in Section 3, where we consider

$$f^{so}(x) = \sqrt{x(1-x)} \sin\left(\frac{2.1\pi}{x+0.05}\right), f^{ta}(x) = f^{so}(x) + x.$$

We used the transformation function  $G(a, b) = a + b$ . The bandwidths of the kernels were chosen by cross validation. For estimating  $f^{so}$ , the chosen bandwidth is  $h_{so} = 10^{-8}$ , and for estimating  $w_G$ , the chosen value is  $h_{w_G} = 10^{-5}$ . Figure 1c shows the scale example in Section 3, where we consider the same source regression function and  $f^{ta}(x) = 5f^{so}(x)$ . We tested the transformation function  $G(a, b) = ab$ . Bandwidth parameters were again chosen by cross validation:  $h_{so} = 10^{-7}$  for estimating  $f^{so}$ , and  $h_{w_G} = 5 \times 10^{-4}$  for estimating  $w_G$ . The plots show that by using our proposed transfer learning framework with an appropriate transformation function, we can estimate the target regression function better, especially in regions where  $f^{ta}$  is not smooth.

### C.2 Transferring from kin-8nh to kin-8fm

Now we briefly discuss the results of the second transfer task with the robotic arm data described in Section 6. The source domain is kin-8nh and the target domain is kin-8fm. The results are shown in Table 3. Here we see the effects of trying to transfer to an “easy” domain. We do not gain any advantage by using the transfer algorithm, except for the smallest value of  $n_{ta}$  - even here the gain is minimal. However, it should be noted that using transfer learning does not negatively affect performance. And we point out that in a dataset where the smoothness conditions are unknown, we would use cross-validation to decide whether or not to use the source data.