# Transfer Learning Through Conditional Quantile Matching

Yikun Zhang [1]   Steven Wilkins-Reeves [2]   Wesley Lee [2]   Aude Hofleitner [2]

## Abstract

We introduce a transfer learning framework for regression that leverages heterogeneous source domains to improve predictive performance in a data-scarce target domain. Our approach learns a conditional generative model separately for each source domain and calibrates the generated responses to the target domain via conditional quantile matching. This distributional alignment step corrects general discrepancies between source and target domains without imposing restrictive assumptions such as covariate or label shift. The resulting framework provides a principled and flexible approach to high-quality data augmentation for downstream learning tasks in the target domain. From a theoretical perspective, we show that an empirical risk minimizer (ERM) trained on the augmented dataset achieves a tighter excess risk bound than the target-only ERM under mild conditions. In particular, we establish new convergence rates for the quantile matching estimator that governs the transfer bias-variance tradeoff. From a practical perspective, extensive simulations and real data applications demonstrate that the proposed method consistently improves prediction accuracy over target-only learning and competing transfer learning methods.

## 1. Introduction

The rapid growth in the volume and heterogeneity of data has created unprecedented opportunities for machine learning methods to improve predictive performances in domains where labeled data are scarce. A prominent paradigm that exploits such opportunities is transfer learning, which aims to improve performances in a target domain by leveraging information from one or more related but heterogeneous source domains (Torrey & Shavlik, 2010;

*Preprint. February 3, 2026.*

Weiss et al., 2016). Commonly, we observe data from $K$ source domains, $\left\{\left(X_i^{(k)}, Y_i^{(k)}\right)\right\}_{i=1}^{n_k} \sim P^{(k)}$ for $k = 1, ..., K$, together with limited data from a target domain $\left\{\left(X_i^{(0)}, Y_i^{(0)}\right)\right\}_{i=1}^{n_0} \sim P^{(0)}$, where $P^{(j)}, j = 0, 1, ..., K$ are probability distributions supported on the product space $\mathcal{X} \times \mathcal{Y}$. The goal is to use the source data to enhance learning in the target domain, especially when $n_0$ is small.

Despite its promise, the success of transfer learning critically depends on how well the source and target distributions align. Distributional discrepancies can lead to negative transfer, in which incorporating source data degrades performances in the target domain (Wang et al., 2019; Zhang et al., 2022). To mitigate this risk, much of the existing literature relies on structural assumptions that constrain the relationship between source and target distributions. Two widely studied assumptions are:

(i) *Covariate Shift* (Shimodaira, 2000): $P^{(k)}(y|x) = P^{(0)}(y|x)$ but $P^{(k)}(x) \neq P^{(0)}(x)$;

(ii) *Label Shift* (Saerens et al., 2002; Nguyen et al., 2016): $P^{(k)}(x|y) = P^{(0)}(x|y)$ but $P^{(k)}(y) \neq P^{(0)}(y)$.

A number of recent extensions have sought to relax these assumptions by introducing invariant or transformed features (Pan et al., 2010; Gong et al., 2016), latent variables (Tsai et al., 2024), or localized shift models (Wilkins-Reeves et al., 2024). Nevertheless, empirical evidence suggests that both covariate and label shift assumptions are often violated in practice (Zhang et al., 2015; Schrouff et al., 2022), or may hold only for a subset of source domains. Moreover, existing methods that attempt to identify transferable sources typically require additional tuning parameters (Bai et al., 2024) or sample splitting of the already limited target data (Tian & Feng, 2023; Wang et al., 2023), which can further limit their effectiveness.

In this paper, we propose a novel transfer learning framework that addresses general distributional shifts between source and target domains from a fundamentally different perspective. Rather than imposing explicit assumptions on the relationship between $P^{(k)}, k = 1, ..., K$ and $P^{(0)}$, we learn conditional generative models from heterogeneous source domains and then calibrate the generated responses to the target domain via quantile matching. This approach enables high-quality synthetic data augmentation for the target
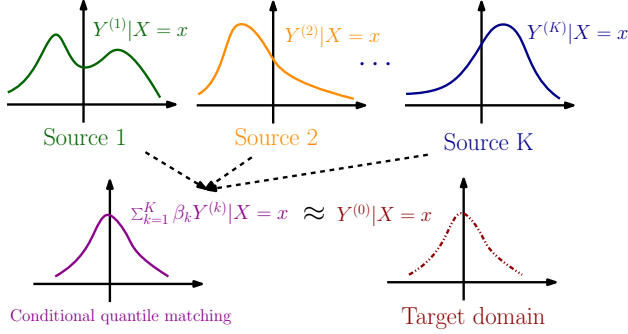
*Figure 1.* Illustration of the conditional quantile matching process.

domain even in challenging settings where both the marginal covariate distributions and the conditional response distributions differ across domains:

$$P^{(i)}(x) \neq P^{(j)}(x) \quad \text{and} \quad P^{(i)}(y|x) \neq P^{(j)}(y|x) \quad (1)$$

for $i \neq j$ and $i, j = 0, ..., K$. Figure 1 illustrates the key idea: synthetic responses generated from multiple source domains are linearly combined and calibrated so that their conditional distribution aligns with that of the target domain. By matching conditional quantiles rather than moments or likelihoods, the proposed framework aligns entire response distributions while remaining computationally tractable. Moreover, the linear structure of the matching step implicitly regularizes the contribution of each source domain, automatically attenuating sources whose conditional responses are poorly aligned with the target.

### 1.1. Contributions and Outline of the Paper

The contributions of this paper are summarized as follows.

1. **Methodology:** We introduce a general transfer learning framework for regression based on conditional quantile matching (TLCQM) in Section 2. The proposed method provides a principled approach to high-quality data augmentation and can be used as a preprocessing step for downstream machine learning models.

2. **Theory:** We establish excess risk bounds on empirical risk minimizers (ERMs) trained on the augmented data in Section 3. Our analysis reveals how transfer bias, generative model error, and quantile matching error jointly govern performance gains. As a key technical contribution, we derive, for the first time, the rate of convergence of the quantile matching estimator introduced in Sgouropoulos et al. (2015).

3. **Experiments:** Through comprehensive experiments on both simulated and real datasets, we demonstrate that the proposed framework consistently improves prediction accuracy over target-only learning and other competing transfer learning methods in Section 4. The

code for our experiments is available at https://github.com/facebookresearch/TLCQM.

### 1.2. Other Related Works

**Transfer learning:** Within the broad literature on transfer learning (Pan & Yang, 2009), our work falls under supervised or transductive transfer learning for regression, also known as multi-domain adaptation (Mansour et al., 2008). We therefore do not review related but distinct lines of work on transfer learning for classification (Cai & Wei, 2021) or unsupervised domain adaptation (Ganin et al., 2016; Long et al., 2015; Chen et al., 2021). Recent theoretical advances in supervised transfer learning for regression have focused on high-dimensional (generalized linear) models (Li et al., 2022; Tian & Feng, 2023; Zhou et al., 2024; 2025), kernel ridge regression (Wang et al., 2023; Wang, 2023; Lin & Reimherr, 2024; 2025), and more general nonparametric settings (Cai & Pu, 2024). The generalized target shift scenario (1) has also been studied through different angles, including location-scale conditional shift models (Zhang et al., 2013), conditional embedding operators from the maximum mean discrepancy in deep learning (Liu et al., 2021), outcome coarsening combined with representation learning (Wu et al., 2023), pairwise similarity-preserving feature extraction to correct distributional discrepancies in $X|Y$ (Taghiyarrenani et al., 2023), and semiparametric regression models with representation learning (He et al., 2024). To the best of our knowledge, no existing work addresses transfer learning under general target shift (1) using quantile matching.

**Quantile matching:** Quantile matching has appeared in various context of the literature, including two-sample testing (Kosorok, 1999), parameter estimation (Dominicy & Veredas, 2013), model diagnostics via quantile-quantile plot (Stuart, 2010), and Bayesian method (Nirwan & Bertschinger, 2020). More broadly, quantile matching can be interpreted as a special case of solving an optimal transport problem between univariate marginal distributions (Mallows, 1972; Villani, 2008).

Finally, we emphasize that our work is distinct from transfer learning for quantile regression (Zhang & Zhu, 2022; Qiao et al., 2024; Bai et al., 2024; Jin et al., 2024), which aims to improve estimation and inference for conditional quantiles in the target domain by leveraging source data. In contrast, our framework focuses on data augmentation through distributional alignment and can be used as a general preprocessing step to enhance a wide range of downstream learning tasks, including but not limited to quantile regression.

### 1.3. Notations

Throughout this paper, we consider a continuous response variable with $\mathcal{Y} \subset \mathbb{R}$ and impose no restrictions on the covariate space $\mathcal{X}$, though our proposed framework can be

readily extended to settings involving multivariate and/or discrete response variables; see Section 5 for details. We adopt standard asymptotic notations: for deterministic sequences $h_n$ and $g_n$ with $g_n > 0$, we write $h_n = O(g_n)$ if $\frac{|h_n|}{g_n} \le C$ for some absolute constant $C > 0$ and all sufficiently large $n$, and $h_n = o(g_n)$ if $\frac{|h_n|}{g_n} \to 0$ as $n \to \infty$. For a random sequence $X_n$ and a deterministic sequence $h_n$, $X_n = o_P(h_n)$ means that $\frac{X_n}{h_n}$ converges to 0 in probability, while $X_n = O_P(h_n)$ indicates that $\frac{X_n}{h_n}$ is bounded in probability as $n \to \infty$. We further write $a_n \lesssim b_n$ (or equivalently $b_n \gtrsim a_n$) if $a_n \le C b_n$ for some constant $C > 0$ and all sufficiently large $n$. When both $a_n \lesssim b_n$ and $a_n \gtrsim b_n$ hold, we denote the asymptotic equivalence by $a_n \asymp b_n$.

## 2. Proposed Framework

In this section, we introduce a general framework for high-quality data augmentation or synthetic data generation in transfer learning for the target domain. A high-level overview of the proposed procedure is summarized in Algorithm 1. The framework consists of learning conditional generative models from source domains, generating synthetic responses at target covariates, and calibrating these responses to the target domain via quantile matching. Among these steps, Steps 1 and 3 are central to the methodology and are detailed in Section 2.1 and Section 2.2, respectively.

### 2.1. Learning Conditional Generative Models via Engression

Depending on data availability and problem structure, a variety of conditional generative models can be employed in Step 1 of Algorithm 1. Here, we consider a neural network-based distributional regression method known as "engression" (Shen & Meinshausen, 2025) as a default choice for its simplicity, flexibility, and empirical effectiveness. Specifically for each source domain $k$, we model the conditional distribution $P^{(k)}(y|x)$ through a measurable function $g^{(k)} : \mathcal{X} \times \mathcal{E} \to \mathcal{Y}$, which maps covariates $x$ and pre-specified noise vectors $\eta \sim P_\eta$ (e.g., Gaussian or uniform) to responses. The induced pushforward measure satisfies $g^{(k)}(x, \cdot)_\# P_\eta = P^{(k)}(\cdot|x)$. Under the engression framework, the population-level estimator is obtained by solving

$$g^{(k)} \in \arg\min_{g \in \mathcal{F}} \mathbb{E}_{(X,Y,\eta) \sim P^{(k)} \times P_\eta} \Big[ |Y - g(X, \eta)| - \frac{1}{2} |g(X, \eta) - g(X, \eta')| \Big]$$

over a certain function class $\mathcal{F}$. In practice, we draw $m$ independent and identically distributed (i.i.d.) samples $\eta_{ij}, j = 1, ..., m$ for each observation $\left(X_i^{(k)}, Y_i^{(k)}\right)$ and

---

**Algorithm 1** Data Augmentation with Conditional Quantile Matching (TLCQM)

**Input:** Target data $\mathcal{D}_T = \left\{ \left( X_i^{(0)}, Y_i^{(0)} \right) \right\}_{i=1}^{n_0}$ and source data $\mathcal{D}_S^{(k)} = \left\{ \left( X_i^{(k)}, Y_i^{(k)} \right) \right\}_{i=1}^{n_k}, k = 1, ..., K$.

**Step 1:** For each source domain $k$, estimate a conditional generative model $\widehat{P}^{(k)}(y|x)$ using the source data $\mathcal{D}_S^{(k)}$.

**Step 2:** For each target covariate vector $X_i^{(0)}$, generate $M$ synthetic responses independently from each learned generative model $\widehat{P}^{(k)}(y|x), k = 1, ..., K$ to obtain $\left\{ \left( X_i^{(0)}, \widehat{\boldsymbol{Y}}_{ij} \right) : i = 1, ..., n_0, j = 1, ..., M \right\}$, where

$$\widehat{\boldsymbol{Y}}_{ij} = \left( \widehat{Y}_{ij}^{(1)}, ..., \widehat{Y}_{ij}^{(K)} \right) \text{ with } \widehat{Y}_{ij}^{(k)} \sim \widehat{P}^{(k)}(y | X_i^{(0)}).$$

**Step 3:** Compute the quantile matching estimator $\widehat{\boldsymbol{\beta}}$ by solving (3).

**Step 4:** Augment the target domain with the synthetic data

$$\mathcal{D}_A = \bigcup_{k=1}^{K} \left\{ \left( X_i^{(k)}, \widehat{\boldsymbol{\beta}}^T \widehat{\boldsymbol{V}}_i^{(k)} \right) \right\}_{i=1}^{n_k},$$

where $\widehat{\boldsymbol{V}}_i^{(k)} = \left( 1, \widehat{Y}_i^{(1,k)}, ..., \widehat{Y}_i^{(K,k)} \right)^T \in \mathbb{R}^{K+1}$ and $\widehat{Y}_i^{(j,k)}$ denotes any predicted value for $Y^{(j)}$ from the $j$-th source generative model evaluated at the covariate $X_i^{(k)}$.

**Step 5 (Optional):** Estimate the density ratios $w_k(x) = \frac{dP^{(0)}(x)}{dP^{(k)}(x)}$ for $k = 1, ..., K$ to correct for covariate shifts.

**Output:** The final augmented data $\mathcal{D}_F = \mathcal{D}_T \cup \mathcal{D}_A$.

---

solve the empirical optimization problem

$$\widehat{g}^{(k)} \in \arg\min_{g \in \mathcal{F}} \frac{1}{n_k} \sum_{i=1}^{n_k} \Bigg[ \frac{1}{m} \sum_{j=1}^{m} |Y_i^{(k)} - g(X_i^{(k)}, \eta_{ij})|$$
$$- \frac{1}{2m(m-1)} \sum_{j=1}^{m} \sum_{j'=1}^{m} |g(X_i^{(k)}, \eta_{ij}) - g(X_i^{(k)}, \eta_{i,j'})| \Bigg].$$

Following Step 2 of Algorithm 1, synthetic responses are then generated at target covariates via $\widehat{Y}_{ij}^{(k)} = \widehat{g}^{(k)}(X_i^{(0)}, \eta_{ij}^{(k)})$, where $\left\{ \eta_{ij}^{(k)} \right\}_{j=1}^{m}$ are i.i.d. noise vectors from the pre-specified distribution and $m > 0$ is chosen sufficiently large (e.g. $m \ge 2000$).

### 2.2. Conditional Quantile Matching

After generating synthetic responses from the learned source-domain generative models $\widehat{P}^{(k)}(y|x), k = 1, ..., K$, the next step is to calibrate these responses to the target domain using the observed data $\mathcal{D}_T = \left\{ \left( X_i^{(0)}, Y_i^{(0)} \right) \right\}_{i=1}^{n_0}$. To this end, we adopt the quantile matching method

(Sgouropoulos et al., 2015) for two major reasons. First, it is computationally more efficient than general optimal transport-based distributional alignment methods (Chernozhukov et al., 2017). Second, the linear structure in (2) implicitly regularizes the contribution of each source domain, enabling automatic identification of transferable sources. At the population level, the quantile matching method seeks a coefficient vector $\boldsymbol{\beta}_* \in \mathbb{R}^{K+1}$ satisfying

$$\boldsymbol{\beta}_* \in \arg\min_{\boldsymbol{\beta} \in \mathbb{R}^{K+1}} S(\boldsymbol{\beta}), \tag{2}$$

where $S(\boldsymbol{\beta}) = \int_0^1 \left[ Q_{Y^{(0)}}(\alpha) - Q_{\boldsymbol{\beta}^T \boldsymbol{V}}(\alpha) \right]^2 d\alpha$ is a squared Mallows' metric (Mallows, 1972) and $Q_\xi(\alpha)$ denotes the $\alpha$-quantile of a random variable $\xi$. In Step 3 of Algorithm 1, we estimate $\boldsymbol{\beta}_*$ by solving

$$\widehat{\boldsymbol{\beta}} \in \arg\min_{\boldsymbol{\beta} \in \mathbb{R}^{K+1}} \sum_{i=1}^{n_0} \sum_{j=1}^{M} \left[ Y_{(i)}^{(0)} - \left( \boldsymbol{\beta}^T \widehat{\boldsymbol{V}} \right)_{(ij)} \right]^2, \tag{3}$$

where $\widehat{\boldsymbol{V}}_{ij} = \left(1, \widehat{\boldsymbol{Y}}_{ij}\right)^T = \left(1, \widehat{Y}_{ij}^{(1)}, ..., \widehat{Y}_{ij}^{(K)}\right)^T \in \mathbb{R}^{K+1}$. Here, $Y_{(1)}^{(0)} \leq \cdots \leq Y_{(n_0)}^{(0)}$ are the order statistics of $Y_1^{(0)}, ..., Y_{n_0}^{(0)}$ and $\left( \boldsymbol{\beta}^T \widehat{\boldsymbol{V}} \right)_{(1)} \leq \cdots \leq \left( \boldsymbol{\beta}^T \widehat{\boldsymbol{V}} \right)_{(n_0 M)}$ are the order statistics of $\boldsymbol{\beta}^T \widehat{\boldsymbol{V}}_{11}, ..., \boldsymbol{\beta}^T \widehat{\boldsymbol{V}}_{n_0 M}$. The number of generated synthetic responses $M$ should be chosen sufficiently large (*e.g.*, $M = 3000$) to ensure convergence of $\widehat{\boldsymbol{\beta}}$ to the global optimum. Although (3) admits no closed-form solution, it can be efficiently solved via an iterative algorithm with established algorithmic convergence guarantees in Sgouropoulos et al. (2015).

# 3. Theoretical Analysis

This section provides a theoretical comparison of prediction risks between ERMs trained solely on the target domain and those trained on the augmented data $\mathcal{D}_F$ produced by our TLCQM framework (Algorithm 1). Notably, our TLCQM framework is agnostic to the specific learning or optimization procedure used to fit the prediction function and applies to a broad class of machine learning methods.

Let $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$ be a loss function. The population risk of a prediction function $f \in \mathcal{F}$ on the target domain is $R(f) := \mathbb{E}_{P^{(0)}} \left[ \ell \left( Y^{(0)}, f(X^{(0)}) \right) \right] = \mathbb{E} \left[ \ell \left( Y^{(0)}, f(X^{(0)}) \right) \right]$, and we define $f^{(0)} = \arg\min_{f \in \mathcal{F}} R(f)$ as the population risk minimizer over the function class $\mathcal{F}$. Then, the excess risk of any candidate function $f : \mathcal{X} \to \mathcal{Y}$ is defined as $R(f) - R(f^{(0)}) = \mathbb{E} \left[ \ell \left( Y^{(0)}, f(X^{(0)}) \right) \right] - \mathbb{E} \left[ \ell \left( Y^{(0)}, f^{(0)}(X^{(0)}) \right) \right]$.

We analyze empirical risk minimization (ERM) procedures corresponding to different training strategies:

- **(Target-only ERM)** We define $\widehat{f}^{(0)} : \mathcal{X} \to \mathcal{Y}$ as:

$$\widehat{f}^{(0)} = \arg\min_{f \in \mathcal{F}} \frac{1}{n_0} \sum_{i=1}^{n_0} \ell \left( Y_i^{(0)}, f(X_i^{(0)}) \right). \tag{4}$$

- **(TLCQM ERM)** Using the augmented data $\mathcal{D}_F$ from Algorithm 1, the transfer learning prediction function is defined as:

$$\widehat{f}^{(0,tl)} = \arg\min_{f \in \mathcal{F}} \frac{1}{N} \left\{ \sum_{i=1}^{n_0} \ell \left( Y^{(0)}, f(X_i^{(0)}) \right) \right.$$
$$\left. + \sum_{k=1}^{K} \sum_{i=1}^{n_k} \widehat{w}_k(X_i^{(k)}) \cdot \ell \left( \widehat{\boldsymbol{\beta}}^T \widehat{\boldsymbol{V}}_i^{(k)}, f(X_i^{(k)}) \right) \right\}, \tag{5}$$

where $N = \sum_{k=0}^{K} n_k$ is the combined sample size of both source and target domains and $\widehat{w}_k(x)$ is the estimated density ratio $w_k(x) = \frac{dP^{(0)}(x)}{dP^{(k)}(x)}$ for $k = 1, ..., K$.

The density ratio correction in (5) is not necessary when the function class $\mathcal{F}$ for the prediction task is correctly specified (Sugiyama et al., 2007). In practice, however, incorporating importance weights improves robustness to model misspecification and often stabilizes numerical optimization.

## 3.1. Assumptions

**Assumption 3.1** (Basic regularity conditions).

(a) There exists a constant $B_\ell > 0$ such that $|\ell(y, y')| \leq B_\ell$ and $|\ell(y, y') - \ell(y, y'')| \leq B_\ell |y' - y''|$.

(b) The target distribution $P^{(0)}$ is absolutely continuous with respect to each $P^{(k)}$ and $w_k(x) = \frac{dP^{(0)}(x)}{dP^{(k)}(x)} \leq B_w$ for all $k = 1, ..., K$ and a constant $B_w > 0$.

**Assumption 3.2** (Boundedness conditions for excess risk).

(a) The conditional distribution $P^{(k)}(Y^{(k)}|X = x)$ is modeled through $Y^{(k)} = g^{(k)}(x, \eta)$ for $k = 0, 1, ..., K$, where $\eta$ is an independent random vector with a pre-specified distribution. Furthermore, $\sup_{(x,\eta)} \left| g^{(k)}(x, \eta) \right| \leq B_g$ and $\sup_{(x,\eta)} \left| \widehat{g}^{(k)}(x, \eta) \right| \leq B_g$ for a constant $B_g > 0$ and $k = 1, ..., K$.

(b) Let $\mathcal{B} = \arg\min_{\boldsymbol{\beta} \in \mathbb{R}^{K+1}} S(\boldsymbol{\beta})$ as in (2), where $\boldsymbol{V} = \left(1, Y^{(1,0)}, ..., Y^{(K,0)}\right)^T \in \mathbb{R}^{K+1}$ with $Y^{(k,0)} = g^{(k)}(X^{(0)}, \eta^{(k)})$ for $k = 1, ..., K$. For any $\boldsymbol{\beta}_* \in \mathcal{B}$, $\|\boldsymbol{\beta}_*\|_2 \leq B_\beta$ for a constant $B_\beta > 0$.

Assumption 3.1(a) is a standard boundedness and Lipschitz continuity condition on the loss function, which holds, for

example, when $\mathcal{Y}$ is compact or when the loss function $\ell$ is truncated at a certain value $B_\ell > 0$. Assumption 3.1(b) ensures that the density ratios (or importance weights) between the target and source domains are well-controlled and can be consistently estimated, thereby avoiding instability due to extreme covariate shifts (Gretton et al., 2009; Reddi et al., 2015). Assumption 3.2(a) imposes mild boundedness conditions on the conditional generative models and their estimators, which resemble the compactness assumption on $\mathcal{Y}$. Finally, Assumption 3.2(b) restricts the solution set of (2) to be bounded, which can be enforced by adding a regularization term to the objective in (2).

## 3.2. Target-Only Excess Risk

Let $\widehat{\mathrm{Rad}}_n(\mathcal{F}) = \frac{1}{n}\mathbb{E}_\sigma\left[\sup_{f\in\mathcal{F}}\left|\sum_{i=1}^n \sigma_i \cdot f(X_i)\right|\right]$ be the empirical Rademacher complexity of the function class $\mathcal{F}$ with respect to the sample $\{X_i\}_{i=1}^n$, where $\sigma_i$'s are i.i.d. Rademacher random variable, *i.e.*, $\mathbb{P}(\sigma_i = 1) = \mathbb{P}(\sigma_i = -1) = \frac{1}{2}$. Denote $\mathrm{Rad}_n(\mathcal{F}) = \mathbb{E}\left[\widehat{\mathrm{Rad}}_n(\mathcal{F})\right]$.

**Proposition 3.3.** *Under Assumption 3.1(a), the target-only prediction function $\widehat{f}^{(0)}$ in (4) has its excess risk satisfying that, with probability at least $1 - \delta$,*

$$R(\widehat{f}^{(0)}) - R(f^{(0)}) \lesssim \mathrm{Rad}_{n_0}(\mathcal{F}) + \sqrt{\frac{\log(1/\delta)}{n_0}}.$$

Proposition 3.3 is a classical result in statistical learning theory; see, *e.g.*, Theorem 26.5 in Shalev-Shwartz & Ben-David (2014). The upper bound remains valid if $\mathrm{Rad}_{n_0}(\mathcal{F})$ is replaced by its empirical counterpart $\widehat{\mathrm{Rad}}_{n_0}(\mathcal{F})$. In many common settings, $\mathrm{Rad}_{n_0}(\mathcal{F})$ scales as $O\left(\frac{1}{\sqrt{n_0}}\right)$ up to some $\log n_0$ factor, with constants depending critically on the complexity of the functional class $\mathcal{F}$. The representative examples are summarized in Table 3.

## 3.3. Excess Risk Under Our TLCQM framework

**Theorem 3.4.** *Under Assumptions 3.1 and 3.2, the transfer learning prediction function $\widehat{f}^{(0,tl)}$ in (5) has its excess risk satisfying that, with probability at least $1 - \delta$,*

$$R(\widehat{f}^{(0,tl)}) - R(f^{(0)})$$
$$\lesssim \underbrace{\mathrm{Rad}_N(\mathcal{F}) + \sqrt{\frac{K\log(1/\delta)}{N}}}_{\textit{Standard generalization error}} + \underbrace{\frac{1}{N}\sum_{k=1}^K \|\widehat{w}_k - w_k\|_1}_{\textit{Importance weight error}}$$
$$+ \underbrace{\inf_{\boldsymbol{\beta}_*\in\mathcal{B}}\left\|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_*\right\|_1}_{\textit{Quantile matching error}} + \underbrace{\sum_{k=1}^K\left\|\widehat{g}^{(k)} - g^{(k)}\right\|_\infty}_{\textit{Distributional learning error}}$$

$$+ \underbrace{\inf_{\boldsymbol{\beta}_*\in\mathcal{B}}\sqrt{\int_0^1\left[Q_{Y^{(0)}}(\alpha) - Q_{\boldsymbol{\beta}_*^T\boldsymbol{V}}(\alpha)\right]^2 d\alpha}}_{\textit{Transfer bias}}$$

when $N = \sum_{k=0}^K n_k \gg n_0$, where $\|\widehat{w}_k - w_k\|_1 = \sum_{i=1}^{n_k}\left|\widehat{w}_k(X_i^{(k)}) - w_k(X_i^{(k)})\right|$, $\mathcal{B}$ denotes the solution set of (2), and $\left\|\widehat{g}^{(k)} - g^{(k)}\right\|_\infty = \sup_{x\in\mathcal{X}}\mathbb{E}_\eta\left|\widehat{g}^{(k)}(x,\eta) - g^{(k)}(x,\eta)\right|$ for $k = 1, ..., K$.

The proof of Theorem 3.4 is in Section B, where the transfer-learning excess risk is decomposed into five error terms.

- **Standard generalization error:** This term mirrors the classical ERM bound but scales with the augmented sample size $N$ rather than $n_0$. It is strictly smaller than the target-only generalization error whenever $N \gg n_0$.

- **Importance weight error:** This term captures the error induced by estimating density ratios under covariate shift. When the penalized risk minimization (Nguyen et al., 2007) is applied, the rate is $O_P\left(\frac{1}{N}\sum_{k=1}^K\left(n_k^{-\frac{1}{2+\tau}} + n_k^{-\frac{1}{4}}\right)\right)$ up to some $\log n_k$ factors (see Lemma 8 in Reddi et al. 2015). When the kernel mean matching (Gretton et al., 2009) is leveraged, the rate becomes $O_P\left(\frac{1}{N}\sum_{k=1}^K n_k^{-\frac{1}{2}}\right)$ up to some $\log n_k$ factors. In both cases, this term vanishes faster than the target-only excess risk under mild growth of $n_k$ for $k = 1, ..., K$.

- **Quantile matching error:** We derive its rate of convergence in Theorem 3.6. The rate consists of two components. The first component is a stochastic variation term, which scales as the product of "transfer bias" and $O_P\left(n_0^{-\frac{1}{4}}\right)$, together with an additional parametric rate $O_P\left(n_0^{-\frac{1}{2}}\right)$. The second component is a bias term that arises from errors in learning the source conditional distributions and is thus tied to the "distributional learning error" discussed below. Consequently, the quantile matching error decays faster than the target-only excess risk, provided that both the "transfer bias" and the "distributional learning error" vanish at the rate $o_P\left(n_0^{-\frac{1}{2}}\right)$. As we argue below, these conditions are attainable under mild regularity assumptions.

- **Distributional learning error:** This term depends solely on the (expected) $L_\infty$ prediction errors in the source domains. Under standard smoothness assumptions, the optimal nonparametric error rate is $O\left(\sum_{k=1}^K n_k^{-\frac{s}{d+2s}}\right)$, where $s$ is the level of Hölder smoothness of $g^{(k)}$ for $k = 1, ..., K$ (Stone, 1980; 1982). This rate is attainable by most of parametric or nonparametric methods, including sieve meth-

ods (Chen & Christensen, 2013) and neural networks (Imaizumi, 2023). Within the engression framework, when $g^{(k)}$ is quadratic in $x$ and $\eta$, the rate of convergence for $\left\|\widehat{g}^{(k)} - g^{(k)}\right\|_\infty$ is at most $O_P\left(n_k^{-\frac{1}{3}}\right)$ up to a $\log n_k$ factor for $k = 1, ..., K$.

- **Transfer bias:** This term captures the intrinsic approximation error of the population-level quantile matching procedure (2) relative to the target response distribution. Unlike the preceding terms, it reflects a fundamental limitation of transfer learning rather than an estimation error. As shown in (6) of Section B, this bias admits an equivalent representation as $\mathbb{E}_{P^{(k)} \times P_\eta} \left| g^{(0)}(X, \eta) - \beta_0 - \sum_{k=1}^{K} \beta_k \cdot g^{(k)}(X, \eta) \right|$ across the source domains. Importantly, the "transfer bias" vanishes when $P^{(0)}$ lies in the convex hull of $P^{(k)}, k = 1, ..., K$ and there are no covariate shifts, a setting commonly assumed in multi-source transfer learning (see Section 3 in Turrisi et al. 2022). In such cases, the proposed framework achieves asymptotically unbiased transfer learning.

Finally, we emphasize that when the optimal quantile matching estimators are sparse, all those $K$-dependent terms in Theorem 3.4 scale with the number of active source domains rather than the total number of source domains.

### 3.4. Quantile Matching Error

**Assumption 3.5** (Convergence conditions for quantile matching). Let $\mathcal{B} = \arg\min_{\beta \in \mathbb{R}^{K+1}} S(\beta)$ as in (2) and $R_\beta > 0$ be a constant.

(a) The density functions $p_{Y^{(k)}}(y) := p^{(k)}(y), k = 0, 1, ..., K$ and $p_{\beta^T V}(u)$ exist and differentiable for all $\beta$ with $\|\beta - \beta_*\|_2 \le R_\beta$ and $\beta_* \in \mathcal{B}$, where $V = \left(1, Y^{(1,0)}, ..., Y^{(K,0)}\right)$ with $Y^{(k)} \sim P^{(k)}(y|X^{(0)})$.

(b) For any $\beta$ with $\|\beta - \beta_*\|_2 \le R_\beta$ and $\beta_* \in \mathcal{B}$, it holds that $\sup_{\alpha \in [0,1]} \left| p'_{\beta^T V}(Q_{\beta^T V}(\alpha)) \right| < \infty$ and $\inf_{\alpha \in [0,1]} \left| p_{\beta^T V}(Q_{\beta^T V}(\alpha)) \right| > 0$. Furthermore, $\sup_{\alpha \in [0,1]} \left| p'_{Y^{(0)}}(Q_{Y^{(0)}}(\alpha)) \right| < \infty$ and $\sup_{\alpha \in [0,1]} \left| p_{Y^{(0)}}(Q_{Y^{(0)}}(\alpha)) \right| > 0$.

(c) For any $\beta_* \in \mathcal{B}$, there exist a constant $\lambda_{\min} > 0$ such that the smallest eigenvalue of $\nabla^2_\beta S(\beta) \in \mathbb{R}^{(K+1) \times (K+1)}$ is uniformly larger than $\lambda_{\min} > 0$ for all $\beta$ with $\|\beta - \beta_*\|_2 \le R_\beta$.

Assumption 3.5(a,b) is known as the Kiefer condition, which is commonly assumed for the uniform Bahadur-Kiefer bounds for empirical quantile processes (Kiefer, 1970; Kulik, 2007); see also Lemma C.1 in Section C. Assumption 3.5(c) is a standard local convexity condition on the

objective function $S(\beta)$ around its local minima. As we derive in Lemma C.2 that

$$\nabla^2_\beta S(\beta) = 2 \int_0^1 \left\{ \left[ \nabla_\beta Q_{\beta^T V}(\alpha) \right] \left[ \nabla_\beta Q_{\beta^T V}(\alpha) \right]^T \right. \\ \left. - \left[ Q_{Y^{(0)}}(\alpha) - Q_{\beta^T V}(\alpha) \right] \nabla^2_\beta Q_{\beta^T V}(\alpha) \right\}^2 d\alpha,$$

Assumption 3.5(c) is naturally satisfied when the distributions of $Y^{(0)}$ and $\beta_*^T V$ at the optimal matching point $\beta_*$ are close, i.e., $Q_{Y^{(0)}}(\alpha) - Q_{\beta_*^T V}(\alpha) \approx 0$ for all $\alpha \in (0, 1)$. This observation again highlights the fundamental role of the "transfer bias": small population-level mismatch not only improves approximation accuracy but also ensures favorable curvature of the quantile matching objective. This connection mirrors the role of the transfer bias identified in Theorem 3.4 and is central to the theoretical benefits of the proposed TLCQM framework.

**Theorem 3.6.** *Under Assumptions 3.2 and 3.5, it holds that*

$$\inf_{\beta_* \in \mathcal{B}} \left\| \widehat{\beta} - \beta_* \right\|_1 = O\left( \sqrt{K \sum_{k=1}^{K} \left\| \widehat{g}^{(k)} - g^{(k)} \right\|_\infty} \right) + O_P\left( \sqrt{\frac{K \log \log n_0}{n_0}} \right. $$
$$\left. + \sqrt{K} \left( \frac{\log \log n_0}{n_0} \inf_{\beta_* \in \mathcal{B}} \int_0^1 \left[ Q_{Y^{(0)}}(\alpha) - Q_{\beta_*^T V}(\alpha) \right]^2 d\alpha \right)^{\frac{1}{4}} \right)$$

*up to some Monte Carlo approximation errors $O\left( \frac{1}{\sqrt{M}} \right)$.*
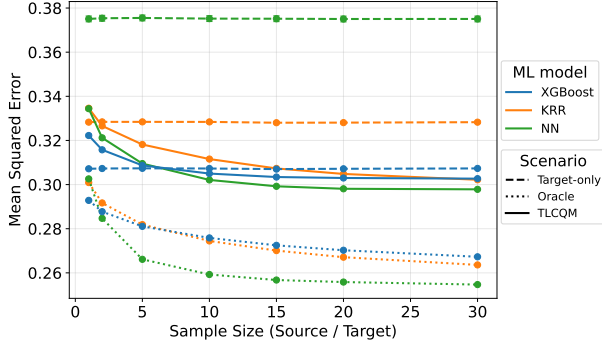
The proof of Theorem 3.6 is in Section C. As discussed after Theorem 3.4, the explicit dependence of the convergence rate on the "transfer bias" $\inf_{\beta_* \in \mathcal{B}} S(\beta_*)$ is essential for the improved performance of our proposed TLCQM framework.

## 4. Numerical Experiments

In this section, we empirically evaluate the performances of several machine learning models under our proposed TLCQM framework in Algorithm 1 and compare them with target-only machine learning models. These experiments are designed to substantiate our theory in Section 3. Furthermore, we benchmark our TLCQM method against several competing transfer learning approaches for regression under distributional shift, including (i) transfer learning with kernel ridge regression and automatic source selection (TKRR; Wang et al. 2023), (ii) deep transfer learning for conditional shift in regression (CDAR; Liu et al. 2021), and (iii) multi-domain adaptation for regression under conditional shift (DARC; Taghiyarrenani et al. 2023). Implementation details for all methods are provided in Section A, and the evaluations are conducted via simulation studies, experiments on a public dataset, and a real-world machine learning application at a technology company.

### 4.1. Simulation Studies

We consider a simulation setting in which both covariate shift and concept shift are present between the source and

*(a)* MSE as a function of the source-to-target sample size ratio for different machine learning models under three data scenarios.

*(b)* MSE as a function of the target-domain sample size for different transfer learning methods.

*Figure 2.* Prediction performances of different machine learning models applied to various data scenarios and other competing transfer learning methods on simulated data.

target domains. For the two candidate source domains, $Y^{(1)} = \sin\left(3\theta^T X^{(1)}\right) + 1 + \epsilon$ and $Y^{(2)} = \cos\left(3\theta^T X^{(2)}\right) + 1 + \epsilon$, where $\theta = \left(1, \frac{1}{2}, ..., \frac{1}{6}\right)^T \in \mathbb{R}^6$ and $X^{(1)}, X^{(2)} \sim \mathcal{N}(\mathbf{1}_6, \boldsymbol{I}_6), \epsilon \sim \mathcal{N}(0, 0.25)$ with $\mathbf{1}_6 = (1, ..., 1)^T \in \mathbb{R}^6$ and $\boldsymbol{I}_6 \in \mathbb{R}^{6\times 6}$ being the identity matrix. For the target domain, $Y^{(0)} = \frac{1}{3}\sin\left(3\theta^T X^{(0)}\right) - 3 + \epsilon$ with $X^{(0)} \sim \mathcal{N}(\mathbf{0}_6, 0.25 \cdot \boldsymbol{I}_6)$ and $\epsilon \sim \mathcal{N}(0, 0.25)$.

We vary the target-domain sample size $n_0$ from 50 to 150 and adjust the ratio $\frac{n_k}{n_0}$ between source and target sample sizes from 1 to 30. To empirically consolidate the improvement of data augmentation by our Algorithm 1, we focus on three widely used machine learning models: XGBoost (Chen & Guestrin, 2016), kernel ridge regression (KRR), and neural network (NN), and apply them to three data scenarios: target-only training, oracle training, and augmented training with our proposed TLCQM framework. Here, the oracle data correspond to samples drawn directly from the target-domain distribution, with a total sample size equal to the combined sizes of the source and target datasets.

Figure 2(a) reports the average mean square errors (MSEs) across 1000 Monte Carlo repetitions for each value of $n_0$. Among all the machine learning methods, the TLCQm-augmented data consistently improves predictive performances relative to the target-only training, with particularly pronounced gains when the sample size ratio between the source and target domains is large (*i.e.*, rich source data and scarce target data). These improvements are most evident for NN and KRR, which are known to benefit more strongly from increased sample sizes than tree-based methods such as XGBoost. We further compare our TLCQM framework with competing transfer learning methods that directly leverage both source and target data. As shown in Figure 2(b), all machine learning models trained on TLCQM-augmented data outperform TKRR, CDAR, and DARC across the considered target sample sizes.

## 4.2. Experiments on Real Public Data

To showcase the effectiveness of our proposed TLCQM framework on real-world data, we next evaluate it on a public dataset `Apartment` from the UCI machine learning repository (Kelly et al., 2019). The response variable $Y$ is defined as the logarithm of the apartment rental price. The covariate vector $X$ comprises 10 features, including the numbers of bathrooms and bedrooms, the apartment size in square feet, and 7 binary indicators capturing the presence of photos, parking, storage, a gym, a pool, and whether cats and dogs are permitted. The source domains consist of observations from three randomly selected states: `IL` ($n_1 = 1036$), `OH` ($n_2 = 4905$), and `WA` ($n_3 = 2612$). The target domain is the state `FL`, which contains 5775 observations in total. To construct the training data for the target domain, we randomly subsample $n_0 \in \{100, 200, 300, 500\}$ observations from `FL` and reserve the remaining observations as the test set.



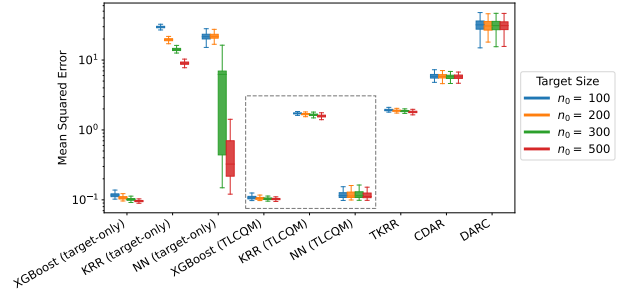*Figure 3.* Performances of different machine learning models applied to target-only and TLCQM scenarios as well as other competing transfer learning methods on the "Apartment" data.

The results, summarized in Figure 3 and Table 1 based on 500 Monte Carlo replications, show that machine learning models trained under our proposed TLCQM framework consistently have smaller MSEs than their target-only counter-

| $n_0$ | 100 | 200 | 300 | 500 |
|---|---|---|---|---|
| XGBoost | 7.16% | 1.59% | -2.28% | -6.80% |
| KRR | 94.20% | 91.41% | 88.43% | 82.53% |
| NN | 99.43% | 99.39% | 96.90% | 87.98% |

*Table 1.* Percentage reduction in MSE achieved by machine learning models trained on TLCQM-augmented data relative to their target-only counterparts.

parts as well as other competing transfer learning methods. The performance gains are particularly pronounced for NN, exceeding those observed for KRR and XGBoost. This pattern aligns with existing findings in the literature that NN typically requires substantially larger sample sizes to surpass tree-based methods (such as XGBoost) on tabular data (Shwartz-Ziv & Armon, 2022).

### 4.3. Case Study: Monthly Active User Prediction

We apply our proposed TLCQM framework (Algorithm 1) to a real-world task involving the prediction of monthly active users (MAUs) for peer apps within a technology company. The dataset consists of country-level MAU measurements collected over a single month, covering two device platforms and five distinct apps. To alleviate the heavy-tailed distribution, we apply a logarithmic transformation $u \mapsto \log(1 + u)$ to all MAU measurements in this case study. Each app is treated as a separate data domain, with approximately 200 observations corresponding to different countries or platforms. The feature vector $X$ includes MAU statistics from other apps as well as a set of country-specific covariates, while the response $Y \in [0, \infty)$ represents the (log-transformed) country-level MAU for each app.

Our goal is to predict the country-level MAU for an app developed by a peer company, for which the third-party MAU estimates are available for a subset of countries. This problem clearly falls into the scenario (1), as both the co-variate distribution and the conditional country-level MAU distribution vary across apps. To emulate this business scenario while enabling systematic evaluation, we adopt a hold-one-out validation, where 4 out of the 5 apps are randomly selected as source domains with complete MAU information, and the remaining app is treated as the target domain, for which MAU statistics are observed only for the aforementioned subset of countries. The proposed TLCQM framework is then employed to augment or impute the missing MAU values for the remaining countries in the target app. Furthermore, to ensure the nonnegativity of the generated MAU values, we impose nonnegative constraints on the coefficients of the quantile matching estimator (3) in our proposed framework, which can be interpreted as an implicit form of regularization tailored to this data application.

The empirical results are summarized in Table 2, which reports the average MSEs for each machine learning model trained on data augmented by our proposed TLCQM framework. For comparison, we also include results from machine learning models trained on the target-only data as well as other competing transfer learning methods that leverage both source and target data, including TKRR, CDAR, and DARC that we mentioned earlier. Since our proposed Algorithm 1 involves generative model learning and synthetic data generation, their reported MSEs exhibit a small degree of Monte Carlo variability as in Table 2. Nevertheless, across all models, the TLCQM-based approaches consistently outperforms the target-only machine learning models and other competing transfer learning methods, even after accounting for this additional randomness. These results demonstrate the robustness and practical relevance of the proposed framework in a business-critical application.

## 5. Discussion

In summary, this work presents a novel transfer learning framework for generating high-quality synthetic regression data that mimics the target-domain distribution. By learning conditional generative models from heterogeneous source-domain data and calibrating the generated samples via conditional quantile matching, the proposed approach yields substantial improvements in downstream prediction performance, supported by both theoretical guarantees and empirical evidence.

Although our analysis focuses on regression with univariate continuous responses, the proposed framework naturally extends to more general settings. For multivariate continuous responses, the quantile matching estimator in (3) can be generalized via a weighted sum of coordinate-wise Mallows' distances. For binary or categorical responses in classification tasks, the objective in (3) reduces to a generalized Brier score. More broadly, the distributional calibration step in our proposed framework could be implemented using alternative metrics or operators, such as isotonic regression (Barlow, 1972), the sliced Wasserstein metric (Bonneel et al., 2015), the Sinkhorn distance (Cuturi, 2013), or maximum mean discrepancy (Cui et al., 2020). A systematic investigation of these extensions is left for future work.

## Acknowledgements

## Impact Statement

This paper presents work whose goal is to advance the field of machine learning. There are many potential societal

| | Target-Only | | | TLCQM | | | TKRR | CDAR | DARC |
|---|---|---|---|---|---|---|---|---|---|
| | XGBoost | NN | KRR | XGBoost | NN | KRR | | | |
| Platform I | 4.596 | 2.910 | 33.281 | **1.545 (0.198)** | **1.689 (0.146)** | **1.556 (0.207)** | 2.204 | 3.183 | 36.380 |
| Platform II | 2.117 | 2.144 | 27.577 | **1.364 (0.034)** | **1.841 (0.015)** | **1.377 (0.034)** | 2.348 | 1.867 | 10.884 |

*Table 2.* (Average) MSEs of different machine learning models trained on target-only and TLCQM-augmented data as well as other competing transfer learning methods for the MAU prediction task across two device platforms. Standard errors of the MSEs based on 30 Monte Carlo experiments are reported in parenthesis, where those values smaller than 0.001 are omitted for brevity.

consequences of our work, none which we feel must be specifically highlighted here.

# References

Awasthi, P., Frank, N., and Mohri, M. On the rademacher complexity of linear hypothesis sets. *arXiv preprint arXiv:2007.11045*, 2020.

Bai, R., Zhang, Y., Yang, H., and Zhu, Z. Transfer learning for high-dimensional quantile regression with distribution shift. *arXiv preprint arXiv:2411.19933*, 2024.

Barlow, R. *Statistical Inference Under Order Restrictions: The Theory and Application of Isotonic Regression.* J. Wiley, 1972.

Bartlett, P. L. and Mendelson, S. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of machine learning research*, 3(Nov):463–482, 2002.

Bartlett, P. L., Foster, D. J., and Telgarsky, M. J. Spectrally-normalized margin bounds for neural networks. *Advances in neural information processing systems*, 30, 2017.

Bartlett, P. L., Harvey, N., Liaw, C., and Mehrabian, A. Nearly-tight vc-dimension and pseudodimension bounds for piecewise linear neural networks. *Journal of Machine Learning Research*, 20(63):1–17, 2019.

Bonneel, N., Rabin, J., Peyré, G., and Pfister, H. Sliced and radon wasserstein barycenters of measures. *Journal of Mathematical Imaging and Vision*, 51(1):22–45, 2015.

Cai, T. T. and Pu, H. Transfer learning for nonparametric regression: Non-asymptotic minimax analysis and adaptive procedure. *arXiv preprint arXiv:2401.12272*, 2024.

Cai, T. T. and Wei, H. Transfer learning for nonparametric classification. *The Annals of Statistics*, 49(1):100–128, 2021.

Chen, T. and Guestrin, C. Xgboost: A scalable tree boosting system. In Krishnapuram, B., Shah, M., Smola, A. J., Aggarwal, C. C., Shen, D., and Rastogi, R. (eds.), *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 13-17, 2016*, pp. 785–794. ACM, 2016.

Chen, X. and Christensen, T. Optimal uniform convergence rates for sieve nonparametric instrumental variables regression. *arXiv preprint arXiv:1311.0412*, 2013.

Chen, X., Wang, S., Wang, J., and Long, M. Representation subspace distance for domain adaptation regression. In Meila, M. and Zhang, T. (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 1749–1759. PMLR, Jul 2021.

Chernozhukov, V., Galichon, A., Hallin, M., and Henry, M. Monge–Kantorovich depth, quantiles, ranks and signs. *The Annals of Statistics*, 45(1):223 – 256, 2017.

Cui, P., Hu, W., and Zhu, J. Calibrated reliable regression using maximum mean discrepancy. *Advances in Neural Information Processing Systems*, 33:17164–17175, 2020.

Cuturi, M. Sinkhorn distances: Lightspeed computation of optimal transport. *Advances in neural information processing systems*, 26, 2013.

Dominicy, Y. and Veredas, D. The method of simulated quantiles. *Journal of Econometrics*, 172(2):235–247, 2013.

Ganin, Y., Ustinova, E., Ajakan, H., Germain, P., Larochelle, H., Laviolette, F., March, M., and Lempitsky, V. Domain-adversarial training of neural networks. *Journal of Machine Learning Research*, 17(59):1–35, 2016.

Gong, M., Zhang, K., Liu, T., Tao, D., Glymour, C., and Schölkopf, B. Domain adaptation with conditional transferable components. In *International Conference on Machine Learning*, pp. 2839–2848. PMLR, 2016.

Gretton, A., Smola, A., Huang, J., Schmittfull, M., Borgwardt, K., Schölkopf, B., et al. Covariate shift by kernel mean matching. *Dataset shift in machine learning*, 3(4): 5, 2009.

Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B., and Smola, A. A kernel two-sample test. *The Journal of Machine Learning Research*, 13(1):723–773, 2012.

He, B., Liu, H., Zhang, X., and Huang, J. Representation transfer learning for semiparametric regression. *arXiv preprint arXiv:2406.13197*, 2024.

Imaizumi, M. Sup-norm convergence of deep neural network estimator for nonparametric regression by adversarial training. *arXiv preprint arXiv:2307.04042*, 2023.

Jin, J., Yan, J., Aseltine, R. H., and Chen, K. Transfer learning with large-scale quantile regression. *Technometrics*, 66(3):381–393, 2024.

Kakade, S. M., Sridharan, K., and Tewari, A. On the complexity of linear prediction: Risk bounds, margin bounds, and regularization. *Advances in neural information processing systems*, 21, 2008.

Kelly, M., Longjohn, R., and Nottingham, K. Apartment for rent classified. UCI Machine Learning Repository, 2019. DOI: https://doi.org/10.24432/C5X623.

Kiefer, J. Deviations between the sample quantile process and the sample df. *Nonparametric techniques in statistical inference*, pp. 299–319, 1970.

Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

Koltchinskii, V. *Oracle inequalities in empirical risk minimization and sparse recovery problems: Ecole D'Eté de Probabilités de Saint-Flour XXXVIII-2008*, volume 2033. Springer, 2011.

Kosorok, M. R. Two-sample quantile tests under general conditions. *Biometrika*, 86(4):909–921, 1999.

Kulik, R. Bahadur–kiefer theory for sample quantiles of weakly dependent linear processes. *Bernoulli*, 13(4):1071 – 1090, 2007.

Leboeuf, J.-S., LeBlanc, F., and Marchand, M. Decision trees as partitioning machines to characterize their generalization properties. *Advances in neural information processing systems*, 33:18135–18145, 2020.

Li, S., Cai, T. T., and Li, H. Transfer learning for high-dimensional linear regression: Prediction, estimation and minimax optimality. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 84(1):149–173, 2022.

Lin, H. and Reimherr, M. Smoothness adaptive hypothesis transfer learning. In Salakhutdinov, R., Kolter, Z., Heller, K., Weller, A., Oliver, N., Scarlett, J., and Berkenkamp, F. (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 30286–30316. PMLR, 21–27 Jul 2024.

Lin, H. and Reimherr, M. Model-robust and adaptive-optimal transfer learning for tackling concept shifts in nonparametric regression. *arXiv preprint arXiv:2501.10870*, 2025.

Liu, X., Li, Y., Meng, Q., and Chen, G. Deep transfer learning for conditional shift in regression. *Knowledge-Based Systems*, 227:107216, 2021.

Long, M., Cao, Y., Wang, J., and Jordan, M. Learning transferable features with deep adaptation networks. In *International Conference on Machine Learning*, pp. 97–105. PMLR, 2015.

Mallows, C. L. A note on asymptotic joint normality. *The Annals of Mathematical Statistics*, 43(2):508–515, 1972.

Mansour, Y., Mohri, M., and Rostamizadeh, A. Domain adaptation with multiple sources. *Advances in neural information processing systems*, 21, 2008.

Massart, P. The tight constant in the dvoretzky-kiefer-wolfowitz inequality. *The Annals of Probability*, 18(3): 1269–1283, 1990.

Massart, P. Some applications of concentration inequalities to statistics. *Annales de la Faculté des sciences de Toulouse: Mathématiques*, 9(2):245–303, 2000.

Mohri, M., Rostamizadeh, A., and Talwalkar, A. *Foundations of machine learning*. MIT press, second edition, 2018.

Nguyen, T. D., Christoffel, M., and Sugiyama, M. Continuous target shift adaptation in supervised learning. In *Asian Conference on Machine Learning*, pp. 285–300. PMLR, 2016.

Nguyen, X., Wainwright, M. J., and Jordan, M. Estimating divergence functionals and the likelihood ratio by penalized convex risk minimization. In Platt, J., Koller, D., Singer, Y., and Roweis, S. (eds.), *Advances in Neural Information Processing Systems*, volume 20, 2007.

Nirwan, R.-S. and Bertschinger, N. Bayesian quantile matching estimation. *arXiv preprint arXiv:2008.06423*, 2020.

Pan, S. J. and Yang, Q. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10): 1345–1359, 2009.

Pan, S. J., Tsang, I. W., Kwok, J. T., and Yang, Q. Domain adaptation via transfer component analysis. *IEEE transactions on neural networks*, 22(2):199–210, 2010.

Qiao, S., He, Y., and Zhou, W. Transfer learning for high-dimensional quantile regression with statistical guarantee. *Transactions on Machine Learning Research*, 2024.

Reddi, S., Poczos, B., and Smola, A. Doubly robust covariate shift correction. In *Proceedings of the AAAI conference on artificial intelligence*, volume 29, 2015.

Saerens, M., Latinne, P., and Decaestecker, C. Adjusting the outputs of a classifier to new a priori probabilities: a simple procedure. *Neural Computation*, 14(1):21–41, 2002.

Schrouff, J., Harris, N., Koyejo, S., Alabdulmohsin, I. M., Schnider, E., Opsahl-Ong, K., Brown, A., Roy, S., Mincu, D., Chen, C., et al. Diagnosing failures of fairness transfer across distribution shift in real-world medical settings. *Advances in Neural Information Processing Systems*, 35: 19304–19318, 2022.

Sgouropoulos, N., Yao, Q., and Yastremiz, C. Matching a distribution by matching quantiles estimation. *Journal of the American Statistical Association*, 110(510):742–759, 2015.

Shalev-Shwartz, S. and Ben-David, S. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.

Shen, X. and Meinshausen, N. Engression: extrapolation through the lens of distributional regression. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 87(3):653–677, 2025.

Shimodaira, H. Improving predictive inference under covariate shift by weighting the log-likelihood function. *Journal of Statistical Planning and Inference*, 90(2):227–244, 2000.

Shwartz-Ziv, R. and Armon, A. Tabular data: Deep learning is not all you need. *Information Fusion*, 81:84–90, 2022.

Stone, C. J. Optimal rates of convergence for nonparametric estimators. *The Annals of Statistics*, pp. 1348–1360, 1980.

Stone, C. J. Optimal global rates of convergence for nonparametric regression. *The Annals of Statistics*, pp. 1040–1053, 1982.

Stuart, E. A. Matching methods for causal inference: A review and a look forward. *Statistical Science*, 25(1):1, 2010.

Sugiyama, M., Krauledat, M., and Müller, K.-R. Covariate shift adaptation by importance weighted cross validation. *Journal of Machine Learning Research*, 8(5), 2007.

Taghiyarrenani, Z., Nowaczyk, S., Pashami, S., and Bouguelia, M.-R. Multi-domain adaptation for regression under conditional distribution shift. *Expert systems with applications*, 224:119907, 2023.

Tian, Y. and Feng, Y. Transfer learning under high-dimensional generalized linear models. *Journal of the American Statistical Association*, 118(544):2684–2697, 2023.

Torrey, L. and Shavlik, J. Transfer learning. In *Handbook of research on machine learning applications and trends: algorithms, methods, and techniques*, pp. 242–264. IGI Global Scientific Publishing, 2010.

Truong, L. V. On rademacher complexity-based generalization bounds for deep learning. *arXiv preprint arXiv:2208.04284*, 2022.

Tsai, K., Pfohl, S. R., Salaudeen, O., Chiou, N., Kusner, M., D'Amour, A., Koyejo, S., and Gretton, A. Proxy methods for domain adaptation. In *International Conference on Artificial Intelligence and Statistics*, pp. 3961–3969. PMLR, 2024.

Turrisi, R., Flamary, R., Rakotomamonjy, A., and Pontil, M. Multi-source domain adaptation via weighted joint distributions optimal transport. In *Uncertainty in artificial intelligence*, pp. 1970–1980. PMLR, 2022.

Villani, C. *Optimal transport: old and new*, volume 338. Springer, 2008.

Wang, C., Wang, C., He, X., and Feng, X. Transfer learning for kernel-based regression. *arXiv preprint arXiv:2310.13966*, 2023.

Wang, K. Pseudo-labeling for kernel ridge regression under covariate shift. *arXiv preprint arXiv:2302.10160*, 2023.

Wang, Z., Dai, Z., Póczos, B., and Carbonell, J. Characterizing and avoiding negative transfer. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11293–11302, 2019.

Weiss, K., Khoshgoftaar, T. M., and Wang, D. A survey of transfer learning. *Journal of Big data*, 3(1):9, 2016.

Wilkins-Reeves, S., Chen, X., Ma, Q., Agarwal, C., and Hofleitner, A. Multiply robust estimation for local distribution shifts with multiple domains. In *International Conference on Machine Learning*, pp. 52972–52993. PMLR, 2024.

Wu, Y., Parmigiani, G., and Ren, B. Multi-source domain adaptation for regression. *arXiv preprint arXiv:2312.05460*, 2023.

Yin, D., Kannan, R., and Bartlett, P. Rademacher complexity for adversarially robust generalization. In *International Conference on Machine Learning*, pp. 7085–7094. PMLR, 2019.

Zhang, K., Schölkopf, B., Muandet, K., and Wang, Z. Domain adaptation under target and conditional shift. In *International conference on machine learning*, pp. 819–827. Pmlr, 2013.

Zhang, K., Gong, M., and Schölkopf, B. Multi-source domain adaptation: A causal view. *Proceedings of the AAAI Conference on Artificial Intelligence*, 29(1), 2015.

Zhang, W., Deng, L., Zhang, L., and Wu, D. A survey on negative transfer. *IEEE/CAA Journal of Automatica Sinica*, 10(2):305–329, 2022.

Zhang, Y. and Zhu, Z. Transfer learning for high-dimensional quantile regression via convolution smoothing. *arXiv preprint arXiv:2212.00428*, 2022.

Zhou, D., Li, M., Cai, T., and Liu, M. Model-assisted and knowledge-guided transfer regression for the underrepresented population. *arXiv preprint arXiv:2410.06484*, 2024.

Zhou, D., Liu, M., Li, M., and Cai, T. Doubly robust augmented model accuracy transfer inference with high dimensional features. *Journal of the American Statistical Association*, 120(549):524–534, 2025.

The first row in Table 3 is also known as the finite class lemma or Massart's lemma (Massart, 2000). By Sauer's Lemma, if the VC dimension $\text{VC}(\mathcal{F})$ of a function class $\mathcal{F}$ is less than $n$, then $\text{Rad}_n(\mathcal{F}) = O\left(\sqrt{\frac{\text{VC}(\mathcal{F})\log n}{n}}\right)$.

| Function Class | Upper Bound | References |
|:---:|:---:|:---:|
| Finite hypothesis class $\mathcal{F}$ with $\sup_{x\in\mathcal{X}}|f(x)| \leq B_f$ for any $f \in \mathcal{F}$. | $O\left(B_f\sqrt{\frac{\log|\mathcal{F}|}{n}}\right)$ | Theorem 3.7 in Mohri et al. (2018); see also Massart (2000). |
| Linear predictor class $\mathcal{F} = \{x \mapsto \Upsilon^T x : \|\Upsilon\|_2 \leq B_\gamma, \|x\|_2 < B_x\}$. | $O\left(\frac{B_\gamma B_x}{\sqrt{n}}\right)$ | Theorem 3 in Kakade et al. (2008); see also Theorem 1 in Awasthi et al. (2020). |
| Sparse linear predictor $\mathcal{F} = \{x \mapsto \Upsilon^T x : x \in \mathbb{R}^d, \|\Upsilon\|_0 = s, \|x\|_2 < B_x\}$. | $O\left(B_x\sqrt{\frac{s\log d}{n}}\right)$ | Section 3.1 in Kakade et al. (2008); see also Lemma 7.1 in Koltchinskii (2011). |
| Reproducing kernel Hilbert space $\mathcal{F} = \{f : \mathcal{X} \to \mathcal{H} : \|f\|_{\mathcal{H}} \leq B_f\}$ | $O\left(\frac{B_f}{\sqrt{n}}\right)$ | Lemma 22 in Bartlett & Mendelson (2002). |
| Binary decision tree class with $L$ nodes and $d$ real-valued features | $O\left(\sqrt{\frac{L\log(Ld)\log n}{n}}\right)$ | Corollary 10 in Leboeuf et al. (2020). |
| Neural network with the ReLU activation function, $W$ weights, and $L$ layers. | $O\left(\sqrt{\frac{WL\log W\log n}{n}}\right)$ | Theorem 6 in Bartlett et al. (2019); see also Theorem 5 in Yin et al. (2019), Theorem 14 in Truong (2022), and Theorem 3.3 in Bartlett et al. (2017). |

*Table 3.* Upper bounds on the Rademacher complexity $\text{Rad}_n(\mathcal{F})$ for common function classes

## A. Additional Implementation Details

In this section, we document the practical considerations involved in implementing our proposed TLCQM framework and other experimental setups.

### A.1. Setups for Engression and Conditional Quantile Matching

For the engression method described in Section 2.1, we employ a neural network model with two hidden layers, each consisting of 100 neurons. The latent noise dimension is set to 5, the learning rate to 0.001, and the model is trained for 1000 epochs unless otherwise specified.

For the conditional quantile matching estimator introduced in Section 2.2, we set the number of generative samples for each source domain to $M = 3000$. The iterative algorithm for solving (3) is initialized using the standard least-square estimator

$$\widetilde{\boldsymbol{\beta}} \in \arg\min_{\boldsymbol{\beta}\in\mathbb{R}^{K+1}} \sum_{i=1}^{n_0}\sum_{j=1}^{M}\left[Y_i^{(0)} - \boldsymbol{\beta}^T\widehat{\boldsymbol{V}}_{ij}\right]^2$$

with the same input data. This initialization provides a reasonable starting point for numerical stability and accelerates convergence.

### A.2. Method for Density Ratio Estimation

In Step 5 of Algorithm 1, we estimate the density ratio $w_k(x) = \frac{dP^{(0)}(x)}{dP^{(k)}(x)}$ for $k = 1, ..., K$ using an off-the-shelf approach called kernel mean matching (Gretton et al., 2009; 2012). Let $n_0$ and $n_k, k = 1, ..., K$ denote the sample sizes of the target and source domains, respectively. For each source domain $k$, given a positive definite kernel function $G : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$, we define a matrix $G^{(k)} \in \mathbb{R}^{n_k \times n_k}$ with entries $G_{ij}^{(k)} := G(X_i^{(k)}, X_j^{(k)})$ and a vector $\kappa^{(k)} \in \mathbb{R}^{n_k}$ with entries $\kappa_i^{(k)} = \frac{n_k}{n_0}\sum_{j=1}^{n_0} G(X_i^{(k)}, X_j^{(0)})$. Then, the kernel mean matching estimates the density ratios $\zeta^{(k)} = $

$\left(w_k(X_1^{(k)}), ..., w_k(X_{n_k}^{(k)})\right)^T \in \mathbb{R}^{n_k}$ evaluated on the sample covariates in the $k$-th source domain by solving the following quadratic optimization problem as:

$$\min_{\zeta^{(k)} \in \mathbb{R}^{n_k}} \frac{1}{n_k^2} \zeta^{(k)T} G^{(k)} \zeta^{(k)} - \frac{2}{n_k} \zeta^{(k)T} \kappa^{(k)} \quad \text{subject to } \zeta_i^{(k)} \in [0, B_\zeta] \text{ and } \left| \sum_{i=1}^{n_k} \zeta_i^{(k)} - n_k \right| \le n_k \cdot \xi,$$

where $\xi$ is often chosen to be $O\left(\frac{B_\zeta}{\sqrt{n_k}}\right)$.

### A.3. Hyperparameter Tuning for Standard Machine Learning Models and Comparative Methods

When we apply the standard machine learning methods to datasets under various scenarios, their hyperparameters are tuned via 5-fold cross-validations across a combination of the following candidate choices:

- **XGBoost:** `learning_rate`: $[0.001, 0.01, 0.1]$, `n_estimators`: $[10, 50, 100]$, `max_depth`: $[3, 5]$, `sub_sample`: $[0.8, 1.0]$, `colsample_bytree`: $[0.8, 1.0]$.

- **Kernel Ridge Regression (KRR):** `alpha` (penalty parameter): $\left[3^{-2}, 3^{-1}, ..., 3^6\right] \times \frac{0.1}{n}$, where $n$ is the sample size of the entire training set.

- **Neural network:** `hidden_layer_sizes`: $[(10,), (50,), (100,)]$ and `alpha` (learning rate): $[0.0001, 0.001, 0.01]$ with the activation function as "ReLU", the optimizer as Adam (Kingma & Ba, 2014), and the number of epochs as 1000.

For the transfer learning with kernel ridge regression (TKRR; Wang et al. 2023), we also apply 5-fold cross-validations for selecting the optimal penalty parameter `alpha` within the candidate set $\left[3^{-3}, 3^{-1}, ..., 3^7\right] \times \frac{0.1}{n}$, where $n$ is the sample size of the entire training set. For deep transfer learning for conditional shift in regression (CDAR; Liu et al. 2021), we follow the same setup in their paper, except that we replace the convolution architecture with a 4-layer fully-connected neural network with $d \times 64 \times 16 \times 8$ hidden neurons. Finally, for multi-domain adaptation for regression under conditional shift (DARC; Taghiyarrenani et al. 2023), we set the embedding or latent feature dimension across domains as 8, whose feature extraction model is a 3-layer neural network with $(d + 8) \times 100 \times 100$ hidden neurons.

## B. Proof of Theorem 3.4

*Theorem* 3.4. Under Assumptions 3.1 and 3.2, the transfer learning prediction function $\widehat{f}^{(0,tl)}$ in (5) has its excess risk satisfying that, with probability at least $1 - \delta$,

$$R(\widehat{f}^{(0,tl)}) - R(f^{(0)}) \lesssim \underbrace{\text{Rad}_N(\mathcal{F}) + \sqrt{\frac{K \log(1/\delta)}{N}}}_{\text{Standard generalization error}} + \underbrace{\frac{1}{N} \sum_{k=1}^{K} \|\widehat{w}_k - w_k\|_1}_{\text{Importance weight error}} + \underbrace{\inf_{\boldsymbol{\beta}_* \in \mathcal{B}} \left\|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_*\right\|_1}_{\text{Quantile matching error}}$$

$$+ \underbrace{\sum_{k=1}^{K} \left\|\widehat{g}^{(k)} - g^{(k)}\right\|_\infty}_{\text{Distributional learning error}} + \underbrace{\inf_{\boldsymbol{\beta}_* \in \mathcal{B}} \sqrt{\int_0^1 \left[Q_{Y^{(0)}}(\alpha) - Q_{\boldsymbol{\beta}_*^T \boldsymbol{V}}(\alpha)\right]^2 d\alpha}}_{\text{Transfer bias}}$$

when $N = \sum_{k=0}^{K} n_k \gg n_0$, where $\|\widehat{w}_k - w_k\|_1 = \sum_{i=1}^{n_k} \left|\widehat{w}_k(X_i^{(k)}) - w_k(X_i^{(k)})\right|$, $\mathcal{B}$ denotes the solution set of (2), and $\left\|\widehat{g}^{(k)} - g^{(k)}\right\|_\infty = \sup_{x \in \mathcal{X}} \mathbb{E}_\eta \left|\widehat{g}^{(k)}(x, \eta) - g^{(k)}(x, \eta)\right|$ for $k = 1, ..., K$.

*Proof of Theorem 3.4.* Denote the objective function in (5) by

$$\widehat{R}_N(f; \widehat{\boldsymbol{w}}, \widehat{\boldsymbol{\beta}}, \widehat{\boldsymbol{g}}) := \frac{1}{N} \left\{ \sum_{i=1}^{n_0} \ell\left(Y^{(0)}, f(X_i^{(0)})\right) + \sum_{k=1}^{K} \sum_{i=1}^{n_k} \widehat{w}_k(X_i^{(k)}) \cdot \ell\left(\widehat{\boldsymbol{\beta}}^T \widehat{\boldsymbol{V}}_i^{(k)}, f(X_i^{(k)})\right) \right\}.$$

Given $\widehat{f}^{(0,tl)}$ in (5), we can decompose an upper bound of its excess risk into three terms as follows:

$$R(\widehat{f}^{(0,tl)}) - R(f^{(0)})$$

$$\leq R(\widehat{f}^{(0,tl)}) - \widehat{R}_N(\widehat{f}^{(0,tl)}; \widehat{w}, \widehat{\beta}, \widehat{g}) + \widehat{R}_N(f^{(0)}; \widehat{w}, \widehat{\beta}, \widehat{g}) - R(f^{(0)})$$

$$= R(\widehat{f}^{(0,tl)}) - \widehat{R}_N(\widehat{f}^{(0,tl)}; w, \widehat{\beta}, \widehat{g}) + \widehat{R}_N(f^{(0)}; w, \widehat{\beta}, \widehat{g}) - R(f^{(0)})$$

$$\quad + \widehat{R}_N(\widehat{f}^{(0,tl)}; w, \widehat{\beta}, \widehat{g}) - \widehat{R}_N(\widehat{f}^{(0,tl)}; \widehat{w}, \widehat{\beta}, \widehat{g}) + \widehat{R}_N(f^{(0)}; \widehat{w}, \widehat{\beta}, \widehat{g}) - \widehat{R}_N(f^{(0)}; w, \widehat{\beta}, \widehat{g})$$

$$= \mathbb{E}\left[\ell\left(Y^{(0)}, \widehat{f}^{(0,tl)}(X^{(0)})\right)\right] - \frac{1}{N}\left\{\sum_{i=1}^{n_0} \ell\left(Y_i^{(0)}, \widehat{f}^{(0,tl)}(X_i^{(0)})\right) + \sum_{k=1}^{K}\sum_{i=1}^{n_k} w_k(X_i^{(k)}) \cdot \ell\left(Y_{ki}^{(0)}, \widehat{f}^{(0,tl)}(X_i^{(k)})\right)\right\}$$

$$+ \frac{1}{N}\left\{\sum_{i=1}^{n_0} \ell\left(Y_i^{(0)}, f^{(0)}(X_i^{(0)})\right) + \sum_{k=1}^{K}\sum_{i=1}^{n_k} w_k(X_i^{(k)}) \cdot \ell\left(Y_{ki}^{(0)}, f^{(0)}(X_i^{(k)})\right)\right\} - \mathbb{E}\left[\ell\left(Y^{(0)}, f^{(0)}(X^{(0)})\right)\right]$$

$$+ \widehat{R}_N(\widehat{f}^{(0,tl)}; w, \widehat{\beta}, \widehat{g}) - \widehat{R}_N(\widehat{f}^{(0,tl)}; \widehat{w}, \widehat{\beta}, \widehat{g}) + \widehat{R}_N(f^{(0)}; \widehat{w}, \widehat{\beta}, \widehat{g}) - \widehat{R}_N(f^{(0)}; w, \widehat{\beta}, \widehat{g})$$

$$+ \frac{1}{N}\sum_{k=1}^{K}\sum_{i=1}^{n_k} w_k(X_i^{(k)})\left[\ell\left(Y_{ki}^{(0)}, \widehat{f}^{(0,tl)}(X_i^{(k)})\right) - \ell\left(\widehat{\beta}^T\widehat{V}_i^{(k)}, \widehat{f}^{(0,tl)}(X_i^{(k)})\right)\right]$$

$$+ \frac{1}{N}\sum_{k=1}^{K}\sum_{i=1}^{n_k} w_k(X_i^{(k)})\left[\ell\left(\widehat{\beta}^T\widehat{V}_i^{(k)}, f^{(0)}(X_i^{(k)})\right) - \ell\left(Y_{ki}^{(0)}, f^{(0)}(X_i^{(k)})\right)\right]$$

$$\leq \underbrace{2\sup_{f\in\mathcal{F}}\left|\widehat{R}_N(f) - R(f)\right|}_{\textbf{Term I}} + \underbrace{\widehat{R}_N(\widehat{f}^{(0,tl)}; w, \widehat{\beta}, \widehat{g}) - \widehat{R}_N(\widehat{f}^{(0,tl)}; \widehat{w}, \widehat{\beta}, \widehat{g}) + \widehat{R}_N(f^{(0)}; \widehat{w}, \widehat{\beta}, \widehat{g}) - \widehat{R}_N(f^{(0)}; w, \widehat{\beta}, \widehat{g})}_{\textbf{Term II}}$$

$$+ \underbrace{\frac{1}{N}\sum_{k=1}^{K}\sum_{i=1}^{n_k} w_k(X_i^{(k)})\left[\ell\left(Y_{ki}^{(0)}, \widehat{f}^{(0,tl)}(X_i^{(k)})\right) - \ell\left(\beta_*^T V_i^{(k)}, \widehat{f}^{(0,tl)}(X_i^{(k)})\right)\right]}_{\textbf{Term III}}$$

$$+ \underbrace{\frac{1}{N}\sum_{k=1}^{K}\sum_{i=1}^{n_k} w_k(X_i^{(k)})\left[\ell\left(\beta_*^T V_i^{(k)}, f^{(0)}(X_i^{(k)})\right) - \ell\left(Y_{ki}^{(0)}, f^{(0)}(X_i^{(k)})\right)\right]}_{\textbf{Term III}}$$

$$+ \underbrace{\frac{1}{N}\sum_{k=1}^{K}\sum_{i=1}^{n_k} w_k(X_i^{(k)})\left[\ell\left(\beta_*^T V_i^{(k)}, \widehat{f}^{(0,tl)}(X_i^{(k)})\right) - \ell\left(\widehat{\beta}^T\widehat{V}_i^{(k)}, \widehat{f}^{(0,tl)}(X_i^{(k)})\right)\right]}_{\textbf{Term IV}}$$

$$+ \underbrace{\frac{1}{N}\sum_{k=1}^{K}\sum_{i=1}^{n_k} w_k(X_i^{(k)})\left[\ell\left(\widehat{\beta}^T\widehat{V}_i^{(k)}, f^{(0)}(X_i^{(k)})\right) - \ell\left(\beta_*^T V_i^{(k)}, f^{(0)}(X_i^{(k)})\right)\right]}_{\textbf{Term IV}},$$

where $Y_{ki}^{(0)}, i = 1, ..., n_k$ are random samples from the conditional distribution $P^{(0)}(Y|X = X_i^{(k)}) = g^{(0)}(X_i^{(k)}, \eta_i^{(k)})$ for some independent noise vector $\eta_i^{(k)}$ and $V_i^{(k)} = \left(1, Y_i^{(1,k)}, ..., Y_i^{(K,k)}\right)^T \in \mathbb{R}^{K+1}$ with $Y_i^{(j,k)} = g^{(j)}(X_i^{(k)}, \eta_{ij}^{(k)})$ and some independent noise vector $\eta_{ij}^{(k)}$ for $j = 1, ..., K$. Additionally, $\beta_*$ is the projection of $\widehat{\beta}$ to the solution set $\mathcal{B} = \arg\min_{\beta\in\mathbb{R}^{K+1}} \int_0^1 \left[Q_{Y^{(0)}}(\alpha) - Q_{\beta^T V}(\alpha)\right]^2 d\alpha$ in (2).

• **Term I:** Similar to Proposition 3.3, we know that

$$2\sup_{f\in\mathcal{F}}\left|\widehat{R}_N(f) - R(f)\right| \leq 4B_\ell \cdot \text{Rad}_N(\mathcal{F}) + 2B_\ell\sqrt{\frac{2\log(2/\delta)}{N}}$$

with probability at least $1 - \delta$.

- **Term II:** We compute that

$$\textbf{Term II} = \frac{1}{N} \sum_{k=1}^{K} \sum_{i=1}^{n_k} \left[ \widehat{w}_k(X_i^{(k)}) - w_k(X_i^{(k)}) \right] \ell \left( \widehat{\boldsymbol{\beta}}^T \widehat{\boldsymbol{V}}_i^{(k)}, \widehat{f}^{(0,tl)}(X_i^{(k)}) \right)$$

$$+ \frac{1}{N} \sum_{k=1}^{K} \sum_{i=1}^{n_k} \left[ \widehat{w}_k(X_i^{(k)}) - w_k(X_i^{(k)}) \right] \ell \left( \widehat{\boldsymbol{\beta}}^T \widehat{\boldsymbol{V}}_i^{(k)}, f^{(0)}(X_i^{(k)}) \right)$$

$$\le \frac{2B_\ell}{N} \sum_{k=1}^{K} \| \widehat{w}_k - w_k \|_1 \,,$$

where the last inequality follows from the boundedness of the loss function under Assumption 3.1(a) and $\| \widehat{w}_k - w_k \|_1 = \sum_{i=1}^{n_k} \left| \widehat{w}_k(X_i^{(k)}) - w_k(X_i^{(k)}) \right|$ for $k = 1, ..., K$.

- **Term III:** Define $\widehat{D}_{N-n_0}(f) := \frac{1}{N-n_0} \sum_{k=1}^{K} \sum_{i=1}^{n_k} w_k(X_i^{(k)}) \left[ \ell \left( Y_{ki}^{(0)}, f(X_i^{(k)}) \right) - \ell \left( \boldsymbol{\beta}_*^T \boldsymbol{V}_i^{(k)}, f(X_i^{(k)}) \right) \right]$ and

$$D(f) := \mathbb{E} \left[ \widehat{D}_{N-n_0}(f) \right] = \mathbb{E}_{X \sim P^{(0)}} \left[ \ell \left( Y^{(0)}, f(X) \right) - \ell \left( \boldsymbol{\beta}_*^T \boldsymbol{V}^{(0)}, f(X) \right) \right],$$

where $\boldsymbol{V}^{(0)} = \left( 1, Y^{(1,0)}, ..., Y^{(K,0)} \right)^T \in \mathbb{R}^{K+1}$ with $Y^{(j,0)} = g^{(j)}(X^{(0)}, \eta_j^{(0)})$. Then, we can upper bound **Term III** as:

$$\textbf{Term III} = \frac{N - n_0}{N} \left[ \widehat{D}_{N-n_0}(\widehat{f}^{(0,tl)}) - \widehat{D}_{N-n_0}(f^{(0)}) \right]$$

$$\le \frac{2(N - n_0)}{N} \sup_{f \in \mathcal{F}} \left| \widehat{D}_{N-n_0}(f) - D(f) \right|$$

$$+ \frac{1}{N} \sum_{k=1}^{K} n_k \cdot \mathbb{E}_{X \sim P^{(0)}} \left[ \ell \left( Y^{(0)}, \widehat{f}^{(0,tl)}(X) \right) - \ell \left( \boldsymbol{\beta}_*^T \boldsymbol{V}^{(0)}, \widehat{f}^{(0,tl)}(X) \right) \right]$$

$$+ \frac{1}{N} \sum_{k=1}^{K} n_k \cdot \mathbb{E}_{X \sim P^{(0)}} \left[ \ell \left( Y^{(0)}, f^{(0)}(X) \right) - \ell \left( \boldsymbol{\beta}_*^T \boldsymbol{V}^{(0)}, f^{(0)}(X) \right) \right]$$

$$\overset{(i)}{\le} \frac{4B_\ell(N - n_0)}{N} \cdot \text{Rad}_{N-n_0}(\mathcal{F}) + 2B_\ell \cdot \frac{\sqrt{2(N - n_0) \log(2/\delta)}}{N}$$

$$+ \frac{2}{N} \sum_{k=1}^{K} n_k B_\ell \cdot \mathbb{E}_{X \sim P^{(0)}} \left[ \inf_\gamma \int_{\mathcal{Y} \times \mathcal{Y}} |y_p - y_q| \, d\gamma(y_p, y_q) \right]$$

$$= \frac{4B_\ell(N - n_0)}{N} \cdot \text{Rad}_{N-n_0}(\mathcal{F}) + 2B_\ell \cdot \frac{\sqrt{2(N - n_0) \log(2/\delta)}}{N}$$

$$+ \frac{2(N - n_0)B_\ell}{N} \cdot \mathbb{E}_{X^{(0)} \sim P^{(0)}} \left[ W_1 \left( P^{(0)}(Y^{(0)}|X^{(0)}), P(\boldsymbol{\beta}_*^T \boldsymbol{V}^{(0)}|X^{(0)}) \right) \right]$$

$$\overset{(ii)}{\le} \frac{4B_\ell(N - n_0)}{N} \cdot \text{Rad}_{N-n_0}(\mathcal{F}) + 2B_\ell \sqrt{\frac{2 \log(2/\delta)}{N}}$$

$$+ \frac{4(N - n_0)B_\ell B_g}{N} \sqrt{\int_0^1 \left[ Q_{Y^{(0)}}(\alpha) - Q_{\boldsymbol{\beta}_*^T \boldsymbol{V}}(\alpha) \right]^2 d\alpha}$$

with probability at least $1 - \delta$, where (i) follows from Proposition 3.3 and the Kantorovich-Rubinstein duality with $\gamma$ being a coupling between the conditional distributions $P^{(0)}(y|X)$ and $P(\boldsymbol{\beta}_*^T \boldsymbol{V}^{(0)}|X)$, as well as (ii) leverages the boundedness of $Y^{(k)}, k = 0, 1, ..., K$ in Assumption 3.2(a) with Cauchy-Schwarz inequality. Here, $W_1 \left( P(Y|X), P(Z|X) \right)$ is the Wasserstein-1 distance between the conditional distributions $P(Y|X), P(Z|X)$, and the Wasserstein-2 distance in (ii) above satisfies

$$\mathbb{E}_{X^{(0)} \sim P^{(0)}} \left[ W_2 \left( P^{(0)}(Y^{(0)}|X^{(0)}), P(\boldsymbol{\beta}_*^T \boldsymbol{V}^{(0)}|X^{(0)}) \right) \right] = \sqrt{\int_0^1 \left[ Q_{Y^{(0)}}(\alpha) - Q_{\boldsymbol{\beta}_*^T \boldsymbol{V}}(\alpha) \right]^2 d\alpha}.$$

As a result, we derive that

$$\textbf{Term III} \leq \frac{4B_\ell(N-n_0)}{N} \cdot \mathrm{Rad}_{N-n_0}(\mathcal{F}) + 2B_\ell\sqrt{\frac{2\log(2/\delta)}{N}}$$
$$+ \frac{4(N-n_0)B_\ell B_g}{N}\sqrt{\int_0^1 \left[Q_{Y^{(0)}}(\alpha) - Q_{\boldsymbol{\beta}^T\boldsymbol{V}}(\alpha)\right]^2 d\alpha}$$

with probability at least $1 - \delta$.

Alternatively, we can also bound **Term III** as:

$$\textbf{Term III} = \frac{1}{N}\sum_{k=1}^K\sum_{i=1}^{n_k} w_k(X_i^{(k)})\left[\ell\left(Y_{ki}^{(0)}, \widehat{f}^{(0,tl)}(X_i^{(k)})\right) - \ell\left(\boldsymbol{\beta}_*^T\boldsymbol{V}_i^{(k)}, \widehat{f}^{(0,tl)}(X_i^{(k)})\right)\right]$$
$$+ \frac{1}{N}\sum_{k=1}^K\sum_{i=1}^{n_k} w_k(X_i^{(k)})\left[\ell\left(\boldsymbol{\beta}_*^T\boldsymbol{V}_i^{(k)}, f^{(0)}(X_i^{(k)})\right) - \ell\left(Y_{ki}^{(0)}, f^{(0)}(X_i^{(k)})\right)\right]$$
$$\leq \frac{2B_\ell B_w}{N}\sum_{k=1}^K\sum_{i=1}^{n_k}\left|\boldsymbol{\beta}_*^T\boldsymbol{V}_i^{(k)} - Y_{ki}^{(0)}\right|$$
$$= \frac{2B_\ell B_w}{N}\sum_{k=1}^K\sum_{i=1}^{n_k}\left[\left|\boldsymbol{\beta}_*^T\boldsymbol{V}_i^{(k)} - Y_{ki}^{(0)}\right| - \mathbb{E}\left|\boldsymbol{\beta}_*^T\boldsymbol{V}_i^{(k)} - Y_{ki}^{(0)}\right|\right] + \frac{2B_\ell B_w}{N}\sum_{k=1}^K\sum_{i=1}^{n_k}\mathbb{E}\left|\boldsymbol{\beta}_*^T\boldsymbol{V}_i^{(k)} - Y_{ki}^{(0)}\right|$$
$$\overset{(iii)}{\leq} \frac{4B_\ell B_w B_\beta B_g(N-n_0)(K+1)}{N}\sqrt{\frac{\log(2/\delta)}{N-n_0}} + \frac{2B_\ell B_w}{N}\sum_{k=1}^K n_k \cdot \inf_{\boldsymbol{\beta}\in\mathcal{B}}\mathbb{E}\left|\boldsymbol{\beta}^T\boldsymbol{V}^{(k)} - Y^{(0,k)}\right|$$
$$\leq 4B_\ell B_w B_\beta B_g(K+1)\sqrt{\frac{\log(2/\delta)}{N}}$$
$$+ \frac{2B_\ell B_w}{N}\sum_{k=1}^K n_k \cdot \inf_{\boldsymbol{\beta}\in\mathcal{B}}\mathbb{E}_{(X,\eta)\sim P^{(k)}\times P_\eta}\left|g^{(0)}(X,\eta) - \beta_0 - \sum_{k=1}^K \beta_k \cdot g^{(k)}(X,\eta)\right|$$

with probability at least $1 - \delta$, where the inequality (iii) follows from Hoeffding's inequality with

$$\boldsymbol{\beta}_*^T\boldsymbol{V}_i^{(k)} = \beta_{*0} + \sum_{j=1}^K \beta_{*j}\cdot g^{(j)}(X_i^{(k)}, \eta_i^{(k)}) \leq B_\beta + KB_\beta B_g \leq (K+1)B_\beta B_g$$

when $B_g > 1$. In this case, the "transfer bias" term is expressed as

$$\sum_{k=1}^K n_k \cdot \inf_{\boldsymbol{\beta}\in\mathcal{B}}\mathbb{E}_{(X,\eta)\sim P^{(k)}\times P_\eta}\left|g^{(0)}(X,\eta) - \beta_0 - \sum_{k=1}^K \beta_k \cdot g^{(k)}(X,\eta)\right|, \tag{6}$$

whose magnitude relies on how well the target distribution $g^{(0)}(X,\eta)$ or $P^{(0)}(X,Y)$ can be approximated by the convex hull of $g^{(k)}(X,\eta), k = 1,...,K$ or $P^{(k)}(X,Y), k = 1,...,K$ with intercept in source domains under the covariate distribution in each source domain.

- **Term IV:** By Assumption 3.1, we know that

$$\textbf{Term IV} \leq \frac{2B_\ell B_w}{N}\sum_{k=1}^K\sum_{i=1}^{n_k}\left|\widehat{\boldsymbol{\beta}}^T\widehat{\boldsymbol{V}}_i^{(k)} - \boldsymbol{\beta}_*^T\boldsymbol{V}_i^{(k)}\right|$$
$$\leq \underbrace{\frac{2B_\ell B_w}{N}\sum_{k=1}^K\sum_{i=1}^{n_k}\left|\left(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_*\right)^T\boldsymbol{V}_i^{(k)}\right|}_{\textbf{Term IVa}} + \underbrace{\frac{2B_\ell B_w}{N}\sum_{k=1}^K\sum_{i=1}^{n_k}\left|\boldsymbol{\beta}_*^T\left(\widehat{\boldsymbol{V}}_i^{(k)} - \boldsymbol{V}_i^{(k)}\right)\right|}_{\textbf{Term IVb}}$$

$$+ \underbrace{\frac{2B_\ell B_w}{N} \sum_{k=1}^{K} \sum_{i=1}^{n_k} \left| \left( \widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_* \right)^T \left( \widehat{\boldsymbol{V}}_i^{(k)} - \boldsymbol{V}_i^{(k)} \right) \right|}_{\textbf{Term IVc}},$$

where we recall that $\widehat{\boldsymbol{V}}_i^{(k)} = \left( 1, \widehat{Y}_i^{(1,k)}, ..., \widehat{Y}_i^{(K,k)} \right)^T \in \mathbb{R}^{K+1}$ with $\widehat{Y}_i^{(j,k)} = \frac{1}{M} \sum_{m=1}^{M} \widehat{g}^{(j)}(X_i^{(k)}, \eta_{im}^{(k)})$ and $\boldsymbol{V}_i^{(k)} = \left( 1, Y_i^{(1,k)}, ..., Y_i^{(K,k)} \right)^T \in \mathbb{R}^{K+1}$ with $Y_i^{(j,k)} = g^{(j)}(X_i^{(k)}, \eta_{ij}^{(k)})$ for $j = 1, ..., K$.

Note that **Term IVc** will be dominated by the maximum of **Term IVa** and **Term IVb** when they are small, so we will focus on the first two terms.

For **Term IVa**, we have that

$$\frac{2B_\ell B_w}{N} \sum_{k=1}^{K} \sum_{i=1}^{n_k} \left| \left( \widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_* \right)^T \boldsymbol{V}_i^{(k)} \right| \leq \frac{2B_\ell B_w B_g (N - n_0)}{N} \left\| \widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_* \right\|_1$$

under Assumption 3.2(a).

For **Term IVb**, we also have that

$$\frac{2B_\ell B_w}{N} \sum_{k=1}^{K} \sum_{i=1}^{n_k} \left| \boldsymbol{\beta}_*^T \left( \widehat{\boldsymbol{V}}_i^{(k)} - \boldsymbol{V}_i^{(k)} \right) \right|$$

$$\leq \frac{2B_\ell B_w}{N} \sum_{k=1}^{K} \sum_{i=1}^{n_k} \left[ \left| \boldsymbol{\beta}_*^T \left( \widehat{\boldsymbol{V}}_i^{(k)} - \boldsymbol{V}_i^{(k)} \right) \right| - \mathbb{E}_\eta \left| \boldsymbol{\beta}_*^T \left( \widehat{\boldsymbol{V}}_i^{(k)} - \boldsymbol{V}_i^{(k)} \right) \right| + \mathbb{E}_\eta \left| \boldsymbol{\beta}_*^T \left( \widehat{\boldsymbol{V}}_i^{(k)} - \boldsymbol{V}_i^{(k)} \right) \right| \right]$$

$$\overset{(iv)}{\leq} \frac{2B_\ell B_w (N - n_0)}{N} \sqrt{\frac{2K B_g B_\beta \log(2/\delta)}{N - n_0}} + \frac{2B_\ell B_w B_\beta (N - n_0)}{N} \sum_{k=1}^{K} \left\| \widehat{g}^{(k)} - g^{(k)} \right\|_\infty$$

$$\leq 2B_\ell B_w \sqrt{\frac{2K B_g B_\beta \log(2/\delta)}{N}} + \frac{2B_\ell B_w B_\beta (N - n_0)}{N} \sum_{k=1}^{K} \left\| \widehat{g}^{(k)} - g^{(k)} \right\|_\infty,$$

with probability at least $1 - \delta$, where (iv) follows from the Hoeffding's inequality under Assumption 3.2. Here, $\left\| \widehat{g}^{(k)} - g^{(k)} \right\|_\infty = \sup_{x \in \mathcal{X}} \mathbb{E}_\eta \left| \widehat{g}^{(k)}(x, \eta) - g^{(k)}(x, \eta) \right|$ is the expected $L_\infty$ prediction error over the noise vector. As a result, we derive that

$$\textbf{Term IV} \leq \frac{2B_\ell B_w B_g (N - n_0)}{N} \left\| \widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_* \right\|_1 + 2B_\ell B_w \sqrt{\frac{2K B_g B_\beta \log(2/\delta)}{N}}$$

$$+ \frac{2B_\ell B_w B_\beta (N - n_0)}{N} \sum_{k=1}^{K} \left\| \widehat{g}^{(k)} - g^{(k)} \right\|_\infty$$

with probability at least $1 - \delta$.

Combining all the results in **Terms I, II, III, IV**, we conclude that with probability at least $1 - \delta$,

$$R(\widehat{f}^{(0,tl)}) - R(f^{(0)})$$

$$\leq \text{Rad}_N(\mathcal{F}) + \sqrt{\frac{\log(1/\delta)}{N}} + \frac{1}{N} \sum_{k=1}^{K} \| \widehat{w}_k - w_k \|_1 + \frac{4B_\ell(N - n_0)}{N} \cdot \text{Rad}_{N-n_0}(\mathcal{F}) + 2B_\ell \sqrt{\frac{2 \log(2/\delta)}{N}}$$

$$+ \frac{4(N - n_0) B_\ell B_g}{N} \sqrt{\int_0^1 \left[ Q_{Y^{(0)}}(\alpha) - Q_{\boldsymbol{\beta}_*^T \boldsymbol{V}}(\alpha) \right]^2 d\alpha} + \frac{2B_\ell B_w B_g (N - n_0)}{N} \left\| \widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_* \right\|_1$$

$$+ 2B_\ell B_w \sqrt{\frac{2K B_g B_\beta \log(2/\delta)}{N}} + \frac{2B_\ell B_w B_\beta (N - n_0)}{N} \sum_{k=1}^{K} \left\| \widehat{g}^{(k)} - g^{(k)} \right\|_\infty$$

$$\overset{(v)}{\lesssim} \mathrm{Rad}_N(\mathcal{F}) + \sqrt{\frac{K \log(1/\delta)}{N}} + \frac{1}{N} \sum_{k=1}^{K} \|\widehat{w}_k - w_k\|_1 + \sqrt{\int_0^1 \left[Q_{Y^{(0)}}(\alpha) - Q_{\boldsymbol{\beta}_*^T \boldsymbol{V}}(\alpha)\right]^2 d\alpha}$$

$$+ \inf_{\boldsymbol{\beta}_* \in \mathcal{B}} \left\|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_*\right\|_1 + \sum_{k=1}^{K} \left\|\widehat{g}^{(k)} - g^{(k)}\right\|_\infty,$$

where (v) follows from the condition that $N = \sum_{k=0}^{K} n_k \gg n_0$. $\qquad\qquad\square$

## C. Proof of Theorem 3.6

Before proving Theorem 3.6, we first introduce several notations and supporting lemmas. For any random variable $Z$, we denote the quantile function corresponding to its empirical distribution of $\{Z_1, ..., Z_n\}$ by $\alpha \mapsto Q_{n,Z}(\alpha)$, *i.e.*,

$$Q_{n,Z}(\alpha) = \inf \{z \in \mathbb{R} : F_{n,Z}(z) \geq \alpha\}$$

with $\alpha \in (0,1)$ and $F_{n,Z}(z) = \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}(Z_i \leq z)$, where $\mathbb{1}(A)$ is an indicator function of the event $A$. Then, ignoring the Monte Carlo approximation in (3), we denote

$$\widehat{S}_{n_0}(\boldsymbol{\beta}) = \frac{1}{n_0} \sum_{i=1}^{n_0} \left[Y_{(i)}^{(0)} - \left(\boldsymbol{\beta}^T \widehat{\boldsymbol{V}}\right)_{(i)}\right]^2 = \frac{1}{n_0} \sum_{j=1}^{n_0} \left[Q_{n_0,Y^{(0)}}\left(\frac{j}{n_0}\right) - Q_{n_0,\boldsymbol{\beta}^T \widehat{\boldsymbol{V}}}\left(\frac{j}{n_0}\right)\right]^2,$$

$$S_{n_0}(\boldsymbol{\beta}) = \frac{1}{n_0} \sum_{i=1}^{n_0} \left[Y_{(i)}^{(0)} - (\boldsymbol{\beta}^T \boldsymbol{V})_{(i)}\right]^2 = \frac{1}{n_0} \sum_{j=1}^{n_0} \left[Q_{n_0,Y^{(0)}}\left(\frac{j}{n_0}\right) - Q_{n_0,\boldsymbol{\beta}^T \boldsymbol{V}}\left(\frac{j}{n_0}\right)\right]^2. \tag{7}$$

**Lemma C.1.** *Under Assumption 3.5, it holds that*

$$\sup_{\alpha \in [0,1]} \left|\sqrt{n} \cdot p_{\boldsymbol{\beta}^T \boldsymbol{V}}(Q_{\boldsymbol{\beta}^T \boldsymbol{V}}(\alpha)) \left[Q_{n,\boldsymbol{\beta}^T \boldsymbol{V}}(\alpha) - Q_{\boldsymbol{\beta}^T \boldsymbol{V}}(\alpha)\right] + \sqrt{n} \left[F_{n,\boldsymbol{\beta}^T \boldsymbol{V}}(Q_{\boldsymbol{\beta}^T \boldsymbol{V}}(\alpha)) - \alpha\right]\right| = O_P\left(\frac{(\log n)^{\frac{1}{2}} (\log \log n)^{\frac{1}{4}}}{n^{\frac{1}{4}}}\right),$$

$$\sup_{\alpha \in [0,1]} \left|\sqrt{n} \cdot p_{Y^{(0)}}(Q_{Y^{(0)}}(\alpha)) \left[Q_{n,Y^{(0)}}(\alpha) - Q_{Y^{(0)}}(\alpha)\right] + \sqrt{n} \left[F_{n,Y^{(0)}}(Q_{Y^{(0)}}(\alpha)) - \alpha\right]\right| = O_P\left(\frac{(\log n)^{\frac{1}{2}} (\log \log n)^{\frac{1}{4}}}{n^{\frac{1}{4}}}\right).$$

This is a classical result established by Kiefer (1970) for i.i.d. samples and Kulik (2007) for some weakly dependent quantile processes, and we thus refer readers to these two classical papers for its proof.

**Lemma C.2.** *Under Assumption 3.5(a,b), it holds that*

$$\nabla_{\boldsymbol{\beta}} Q_{\boldsymbol{\beta}^T \boldsymbol{V}}(\alpha) = \mathbb{E}\left[\boldsymbol{V} | \boldsymbol{\beta}^T \boldsymbol{V} = Q_{\boldsymbol{\beta}^T \boldsymbol{V}}(\alpha)\right],$$

$$\nabla_{\boldsymbol{\beta}}^2 Q_{\boldsymbol{\beta}^T \boldsymbol{V}}(\alpha) = \mathrm{Cov}\left[\widehat{\boldsymbol{V}} \widehat{\boldsymbol{V}}^T | \boldsymbol{\beta}^T \widehat{\boldsymbol{V}} = Q_{\boldsymbol{\beta}^T \boldsymbol{V}}(\alpha)\right] - \frac{2 p'_{\boldsymbol{\beta}^T \boldsymbol{V}}(Q_{\boldsymbol{\beta}^T \boldsymbol{V}}(\alpha))}{p_{\boldsymbol{\beta}^T \boldsymbol{V}}(Q_{\boldsymbol{\beta}^T \boldsymbol{V}}(\alpha))} \cdot \mathbb{E}\left[\boldsymbol{V} | \boldsymbol{\beta}^T \boldsymbol{V} = Q_{\boldsymbol{\beta}^T \boldsymbol{V}}(\alpha)\right] \mathbb{E}\left[\boldsymbol{V} | \boldsymbol{\beta}^T \boldsymbol{V} = Q_{\boldsymbol{\beta}^T \boldsymbol{V}}(\alpha)\right]^T,$$

$$\nabla_{\boldsymbol{\beta}}^2 S(\boldsymbol{\beta}) = 2 \int_0^1 \left\{\left[\nabla_{\boldsymbol{\beta}} Q_{\boldsymbol{\beta}^T \boldsymbol{V}}(\alpha)\right] \left[\nabla_{\boldsymbol{\beta}} Q_{\boldsymbol{\beta}^T \boldsymbol{V}}(\alpha)\right]^T - \left[Q_{Y^{(0)}}(\alpha) - Q_{\boldsymbol{\beta}^T \boldsymbol{V}}(\alpha)\right] \nabla_{\boldsymbol{\beta}}^2 Q_{\boldsymbol{\beta}^T \boldsymbol{V}}(\alpha)\right\}^2 d\alpha,$$

*where $S(\boldsymbol{\beta}) = \int_0^1 \left[Q_{Y^{(0)}}(\alpha) - Q_{\boldsymbol{\beta}^T \boldsymbol{V}}(\alpha)\right]^2 d\alpha$ as in (2) and $\boldsymbol{V} \in \mathbb{R}^{K+1}$ can be replaced by any random vector $\boldsymbol{Z} \in \mathbb{R}^{K+1}$ with coordinatewise compact supports and differentiable density function $u \mapsto p_{\boldsymbol{\beta}^T \boldsymbol{Z}}(u)$.*

*Proof of Lemma C.2.* First, we compute the gradient of $G(\boldsymbol{\beta}) := \mathbb{E}\left[\mathbb{1}(\boldsymbol{\beta}^T \boldsymbol{V} \leq z)\right] = F_{\boldsymbol{\beta}^T \boldsymbol{V}}(z)$. For any direction $\boldsymbol{h} \in \mathbb{R}^{K+1}$, we know that

$$\lim_{t \to 0} \frac{G(\boldsymbol{\beta} + t\boldsymbol{h}) - G(\boldsymbol{\beta})}{t} = \lim_{t \to 0} \mathbb{E}\left[\frac{\mathbb{1}\left(z < \boldsymbol{\beta}^T \boldsymbol{V} \leq z - t\boldsymbol{h}^T \boldsymbol{V}\right)}{t}\right]$$

$$\overset{(i)}{=} -\mathbb{E}\left[p_{\boldsymbol{\beta}^T \boldsymbol{V}}(z) \cdot \boldsymbol{h}^T \boldsymbol{V}\right]$$

19

$$= -p_{\boldsymbol{\beta}^T\boldsymbol{V}}(z) \cdot \mathbb{E}\left[\boldsymbol{h}^T\boldsymbol{V}\big|\boldsymbol{\beta}^T\boldsymbol{V} = z\right],$$

where (i) follows from the Lebesgue differentiation theorem under Assumption 3.5(a,b) as $\lim_{t\to 0}\frac{\mathbb{1}\left(z<\boldsymbol{\beta}^T\boldsymbol{V}\leq z-t\boldsymbol{h}^T\boldsymbol{V}\right)}{t} = -p_{\boldsymbol{\beta}^T\boldsymbol{V}}(z) \cdot \boldsymbol{h}^T\boldsymbol{V}$. Hence, $\nabla_{\boldsymbol{\beta}}G(\boldsymbol{\beta}) = -p_{\boldsymbol{\beta}^T\boldsymbol{V}}(z) \cdot \mathbb{E}\left[\boldsymbol{V}\big|\boldsymbol{\beta}^T\boldsymbol{V} = z\right]$ given the definition of directional derivatives.

Second, we compute the Hessian of $G(\boldsymbol{\beta}) := \mathbb{E}\left[\mathbb{1}(\boldsymbol{\beta}^T\boldsymbol{V} \leq z)\right] = F_{\boldsymbol{\beta}^T\boldsymbol{V}}(z)$. For any two directions $\boldsymbol{h}_1, \boldsymbol{h}_2 \in \mathbb{R}^{K+1}$, we know that

$$
\begin{aligned}
D^2G(\boldsymbol{\beta})[\boldsymbol{h}_1, \boldsymbol{h}_2] &= -\lim_{t\to 0}\frac{DG(\boldsymbol{\beta}+t\boldsymbol{h}_2)[\boldsymbol{h}_1] - DG(\boldsymbol{\beta})[\boldsymbol{h}_1]}{t} \\
&= -\lim_{t\to 0}\frac{p_{\boldsymbol{\beta}+t\boldsymbol{h_2^T}\boldsymbol{V}}(z) \cdot \mathbb{E}\left[\boldsymbol{h}_1^T\boldsymbol{V}\big|(\boldsymbol{\beta}+t\boldsymbol{h}_2)^T\boldsymbol{V} = z\right] - p_{\boldsymbol{\beta}^T\boldsymbol{V}}(z) \cdot \mathbb{E}\left[\boldsymbol{h}_1^T\boldsymbol{V}\big|\boldsymbol{\beta}^T\boldsymbol{V} = z\right]}{t} \\
&\overset{(ii)}{=} p'_{\boldsymbol{\beta}^T\boldsymbol{V}}(z) \cdot \mathbb{E}\left[\boldsymbol{h}_1^T\boldsymbol{V}|\boldsymbol{\beta}^T\boldsymbol{V} = z\right]\mathbb{E}\left[\boldsymbol{h}_2^T\boldsymbol{V}|\boldsymbol{\beta}^T\boldsymbol{V} = z\right] - p_{\boldsymbol{\beta}^T\boldsymbol{V}}(z) \cdot \mathrm{Cov}\left[(\boldsymbol{h}_1^T\boldsymbol{V})(\boldsymbol{h}_2^T\boldsymbol{V})|\boldsymbol{\beta}^T\boldsymbol{V} = z\right],
\end{aligned}
$$

where (ii) follows from the fact that $\frac{d}{dt}\mathbb{E}\left[\boldsymbol{h}_1^T\boldsymbol{V}\big|(\boldsymbol{\beta}+t\boldsymbol{h}_2)^T\boldsymbol{V} = z\right]\Big|_{t=0} = \mathrm{Cov}\left(\boldsymbol{h}_1^T\boldsymbol{V}, \boldsymbol{h}_2^T\boldsymbol{V}|\boldsymbol{\beta}^T\boldsymbol{V} = z\right)$. Hence, $\nabla_{\boldsymbol{\beta}}^2 G(\boldsymbol{\beta}) = p'_{\boldsymbol{\beta}^T\boldsymbol{V}}(z) \cdot \mathbb{E}\left[\boldsymbol{V}|\boldsymbol{\beta}^T\boldsymbol{V} = z\right]\mathbb{E}\left[\boldsymbol{V}|\boldsymbol{\beta}^T\boldsymbol{V} = z\right]^T - p_{\boldsymbol{\beta}^T\boldsymbol{V}}(z) \cdot \mathbb{E}\left[\boldsymbol{V}\boldsymbol{V}^T|\boldsymbol{\beta}^T\boldsymbol{V} = z\right]$ by the definition of second-order directional derivatives.

Now, under Assumption 3.5(a), we know that $F_{\boldsymbol{\beta}^T\boldsymbol{V}}(Q_{\boldsymbol{\beta}^T\boldsymbol{V}}(\alpha)) = \alpha$. Taking the gradient $\nabla_{\boldsymbol{\beta}}$ on both sides of the equality gives us that

$$p_{\boldsymbol{\beta}^T\boldsymbol{V}}(Q_{\boldsymbol{\beta}^T\boldsymbol{V}}(\alpha)) \cdot \nabla_{\boldsymbol{\beta}}Q_{\boldsymbol{\beta}^T\boldsymbol{V}}(\alpha) + \nabla_{\boldsymbol{\beta}}F_{\boldsymbol{\beta}^T\boldsymbol{V}}(u)\big|_{u=Q_{\boldsymbol{\beta}^T\boldsymbol{V}}(\alpha)} = 0. \tag{8}$$

Therefore, we conclude that $\nabla_{\boldsymbol{\beta}}Q_{\boldsymbol{\beta}^T\boldsymbol{V}}(\alpha) = \mathbb{E}\left[\boldsymbol{V}\big|\boldsymbol{\beta}^T\boldsymbol{V} = Q_{\boldsymbol{\beta}^T\boldsymbol{V}}(\alpha)\right]$. Furthermore, taking an extra gradient $\nabla_{\boldsymbol{\beta}}$ on both sides of (8) yields that

$$p'_{\boldsymbol{\beta}^T\boldsymbol{V}}(Q_{\boldsymbol{\beta}^T\boldsymbol{V}}(\alpha))\left[\nabla_{\boldsymbol{\beta}}Q_{\boldsymbol{\beta}^T\boldsymbol{V}}(\alpha)\right]\left[\nabla_{\boldsymbol{\beta}}Q_{\boldsymbol{\beta}^T\boldsymbol{V}}(\alpha)\right]^T + p_{\boldsymbol{\beta}^T\boldsymbol{V}}(Q_{\boldsymbol{\beta}^T\boldsymbol{V}}(\alpha)) \cdot \nabla_{\boldsymbol{\beta}}^2 Q_{\boldsymbol{\beta}^T\boldsymbol{V}}(\alpha) + \nabla_{\boldsymbol{\beta}}^2 F_{\boldsymbol{\beta}^T\boldsymbol{V}}(u)\big|_{u=Q_{\boldsymbol{\beta}^T\boldsymbol{V}}(\alpha)} = 0.$$

Therefore, we conclude that

$$\nabla_{\boldsymbol{\beta}}^2 Q_{\boldsymbol{\beta}^T\boldsymbol{V}}(\alpha) = \mathrm{Cov}\left[\boldsymbol{V}\boldsymbol{V}^T|\boldsymbol{\beta}^T\boldsymbol{V} = Q_{\boldsymbol{\beta}^T\boldsymbol{V}}(\alpha)\right] - \frac{2p'_{\boldsymbol{\beta}^T\boldsymbol{V}}(Q_{\boldsymbol{\beta}^T\boldsymbol{V}}(\alpha))}{p_{\boldsymbol{\beta}^T\boldsymbol{V}}(Q_{\boldsymbol{\beta}^T\boldsymbol{V}}(\alpha))} \cdot \mathbb{E}\left[\boldsymbol{V}|\boldsymbol{\beta}^T\boldsymbol{V} = Q_{\boldsymbol{\beta}^T\boldsymbol{V}}(\alpha)\right]\mathbb{E}\left[\boldsymbol{V}|\boldsymbol{\beta}^T\boldsymbol{V} = Q_{\boldsymbol{\beta}^T\boldsymbol{V}}(\alpha)\right]^T.$$

Finally, by direct calculations, we have that

$$\nabla_{\boldsymbol{\beta}}^2 S(\boldsymbol{\beta}) = 2\int_0^1 \left\{\left[\nabla_{\boldsymbol{\beta}}Q_{\boldsymbol{\beta}^T\boldsymbol{V}}(\alpha)\right]\left[\nabla_{\boldsymbol{\beta}}Q_{\boldsymbol{\beta}^T\boldsymbol{V}}(\alpha)\right]^T - \left[Q_{Y^{(0)}}(\alpha) - Q_{\boldsymbol{\beta}^T\boldsymbol{V}}(\alpha)\right]\nabla_{\boldsymbol{\beta}}^2 Q_{\boldsymbol{\beta}^T\boldsymbol{V}}(\alpha)\right\}^2 d\alpha.$$

The proof is thus completed. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

**Lemma C.3.** *Let $\{a_i\}_{i=1}^n$ and $\{b_i\}_{i=1}^n$ be any two sequences of real numbers. Then,*

$$\sum_{i=1}^n\left|a_{(i)} - b_{(i)}\right| \leq \sum_{i=1}^n|a_i - b_i|, \quad \sum_{i=1}^n\left|a_{(i)} - b_{(i)}\right|^2 \leq \sum_{i=1}^n(a_i - b_i)^2$$

*with $\{a_{(i)}\}_{i=1}^n$ and $\{b_{(i)}\}_{i=1}^n$ being the order statistics of $\{a_i\}_{i=1}^n$ and $\{b_i\}_{i=1}^n$, i.e., $a_{(1)} \leq \cdots \leq a_{(n)}$ and $b_{(1)} \leq \cdots \leq b_{(n)}$.*

*Proof of Lemma C.3.* The second inequality was proved in Lemma 1 of (Sgouropoulos et al., 2015), so we only prove the first inequality $\sum_{i=1}^n\left|a_{(i)} - b_{(i)}\right| \leq \sum_{i=1}^n|a_i - b_i|$ by induction.

When $n = 2$, it suffices to show that $|a_{(1)} - b_{(1)}| + |a_{(2)} - b_{(2)}| \leq |a_{(1)} - b_{(2)}| + |a_{(2)} - b_{(1)}|$.

● **Case I:** The ranges of $a_1, a_2$ and $b_1, b_2$ have no overlap. Without loss of generality, we assume that $b_1 \leq b_2 \leq a_1 \leq a_2$. Then,

$$|a_{(1)} - b_{(1)}| + |a_{(2)} - b_{(2)}| = a_1 - b_2 + b_2 - b_1 + a_2 - a_1 + a_1 - b_2$$

$$= b_2 - b_1 + a_2 - a_1 + 2(a_1 - b_2)$$
$$= |a_{(1)} - b_{(2)}| + |a_{(2)} - b_{(1)}|.$$

• **Case II:** The ranges of $a_1, a_2$ and $b_1, b_2$ have an overlap. Without loss of generality, it could be either $b_1 \leq a_1 \leq a_2 \leq b_2$ or $b_1 \leq a_1 \leq b_2 \leq a_2$. If $b_1 \leq a_1 \leq a_2 \leq b_2$, then

$$
\begin{aligned}
|a_{(1)} - b_{(1)}| + |a_{(2)} - b_{(2)}| &= a_1 - b_1 + b_2 - a_2 \\
&\leq a_1 - b_1 + b_2 - a_2 + 2(a_2 - a_1) \\
&= b_2 - a_1 + a_2 - b_1 \\
&= |a_{(1)} - b_{(2)}| + |a_{(2)} - b_{(1)}|.
\end{aligned}
$$

If $b_1 \leq a_1 \leq b_2 \leq a_2$, then

$$
\begin{aligned}
|a_{(1)} - b_{(1)}| + |a_{(2)} - b_{(2)}| &= a_1 - b_1 + a_2 - b_2 \\
&\leq a_1 - b_1 + b_2 - a_2 + 2(b_2 - a_1) \\
&= b_2 - a_1 + a_2 - b_1 \\
&= |a_{(1)} - b_{(2)}| + |a_{(2)} - b_{(1)}|.
\end{aligned}
$$

Hence, the inequality holds when $n = 2$.

Now, support that $\sum_{i=1}^{k} |a_{(i)} - b_{(i)}| \leq \sum_{i=1}^{k} |a_i - b_i|$, and we need to prove this inequality when $n = k + 1$. Without loss of generality, we assume that $a_{k+1} = a_{(1)}$ and $b_j = b_{(1)}$. If $j = k + 1$, then the inequality naturally holds for $n = k + 1$. When $j \neq k + 1$, it follows the proof above for the case when $n = 2$, because

$$|a_{(1)} - b_{(1)}| + |a_j - b_{k+1}| = |a_{k+1} - b_j| + |a_j - b_{k+1}| \leq |a_j - b_j| + |a_{k+1} - b_{k+1}|.$$

As a result,

$$
\begin{aligned}
\sum_{i=1}^{k+1} |a_i - b_i| &\geq |a_{(1)} - b_{(1)}| + |a_j - b_{k+1}| + \sum_{1 \leq i \leq k, i \neq j} |a_i - b_i| \\
&\geq |a_{(1)} - b_{(1)}| + \sum_{i=2}^{k+1} |a_{(i)} - b_{(i)}|,
\end{aligned}
$$

where the last inequality follows from the induction hypothesis for $n = k$. The proof is thus completed. $\qquad\square$

**Lemma C.4.** *Under Assumption 3.2(a) and $\left\| \widehat{g}^{(k)} - g^{(k)} \right\|_{\infty} = \sup_{x \in \mathcal{X}} \mathbb{E}_{\eta} \left| \widehat{g}^{(k)}(x, \eta) - g^{(k)}(x, \eta) \right| \leq 1$, it holds that*

$$\sup_{\boldsymbol{\beta}: \|\boldsymbol{\beta}\|_2 \leq C_\beta} \left| \widehat{S}_{n_0}(\boldsymbol{\beta}) - S_{n_0}(\boldsymbol{\beta}) \right| \leq 8 B_g C_\beta^2 \sqrt{K+1} \sum_{k=1}^{K} \left\| \widehat{g}^{(k)} - g^{(k)} \right\|_{\infty}$$

*for any finite constant $C_\beta > 0$, where $\widehat{S}_{n_0}(\boldsymbol{\beta})$ and $S_{n_0}(\boldsymbol{\beta})$ are defined in (7).*

*Proof of Lemma C.4.* For any fixed $\boldsymbol{\beta} \in \mathbb{R}^{K+1}$, we know that

$$
\begin{aligned}
\left| \widehat{S}_{n_0}(\boldsymbol{\beta}) - S_{n_0}(\boldsymbol{\beta}) \right| &= \left| \frac{1}{n_0} \sum_{i=1}^{n_0} \left[ Y_{(i)}^{(0)} - \left( \boldsymbol{\beta}^T \widehat{\boldsymbol{V}} \right)_{(i)} \right]^2 - \frac{1}{n_0} \sum_{i=1}^{n_0} \left[ Y_{(i)}^{(0)} - \left( \boldsymbol{\beta}^T \boldsymbol{V} \right)_{(i)} \right]^2 \right| \\
&\leq \frac{1}{n_0} \sum_{i=1}^{n_0} \left[ \left( \boldsymbol{\beta}^T \boldsymbol{V} \right)_{(i)} - \left( \boldsymbol{\beta}^T \widehat{\boldsymbol{V}} \right)_{(i)} \right]^2 + \frac{2}{n_0} \sum_{i=1}^{n_0} \left| Y_{(i)}^{(0)} - \left( \boldsymbol{\beta}^T \boldsymbol{V} \right)_{(i)} \right| \left| \left( \boldsymbol{\beta}^T \boldsymbol{V} \right)_{(i)} - \left( \boldsymbol{\beta}^T \widehat{\boldsymbol{V}} \right)_{(i)} \right| \\
&\overset{(i)}{\leq} \frac{1}{n_0} \sum_{i=1}^{n_0} \left[ \boldsymbol{\beta}^T \boldsymbol{V}_i - \boldsymbol{\beta}^T \widehat{\boldsymbol{V}}_i \right]^2 + \frac{2 B_g (1 + C_\beta \sqrt{K+1})}{n_0} \sum_{i=1}^{n_0} \left| \boldsymbol{\beta}^T \boldsymbol{V}_i - \boldsymbol{\beta}^T \widehat{\boldsymbol{V}}_i \right|
\end{aligned}
$$

21

$$\leq C_\beta^2 \sum_{k=1}^K \left\|\widehat{g}^{(k)} - g^{(k)}\right\|_\infty^2 + 2B_g C_\beta (1 + C_\beta \sqrt{K+1}) \sum_{k=1}^K \left\|\widehat{g}^{(k)} - g^{(k)}\right\|_\infty$$

$$\leq 4B_g C_\beta (1 + C_\beta \sqrt{K+1}) \sum_{k=1}^K \left\|\widehat{g}^{(k)} - g^{(k)}\right\|_\infty$$

$$\leq 8B_g C_\beta^2 \sqrt{K+1} \sum_{k=1}^K \left\|\widehat{g}^{(k)} - g^{(k)}\right\|_\infty,$$

where (i) applies Lemma C.3. Since the bound is independent of $\boldsymbol{\beta}$, taking the supremum over $\boldsymbol{\beta}$ for all $\|\boldsymbol{\beta}\|_2 \leq C_\beta$ leads to the final result. $\qquad\square$

*Theorem* 3.6. Under Assumptions 3.2 and 3.5, it holds that

$$\inf_{\boldsymbol{\beta}_* \in \mathcal{B}} \left\|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_*\right\|_1 = O_P\left( \sqrt{K} \left( \frac{\log \log n_0}{n_0} \inf_{\boldsymbol{\beta}_* \in \mathcal{B}} \int_0^1 \left[ Q_{Y^{(0)}}(\alpha) - Q_{\boldsymbol{\beta}_*^T \boldsymbol{V}}(\alpha) \right]^2 d\alpha \right)^{\frac{1}{4}} + \sqrt{\frac{K \log \log n_0}{n_0}} \right)$$

$$+ O\left( \sqrt{K \sum_{k=1}^K \left\|\widehat{g}^{(k)} - g^{(k)}\right\|_\infty} \right)$$

up to some Monte Carlo approximation errors $O\left(\frac{1}{\sqrt{M}}\right)$.

*Proof of Theorem 3.6.* By Theorem 2 in Sgouropoulos et al. (2015), we know that $\widehat{\boldsymbol{\beta}}$ will eventually fall into the $R_\beta$ ball around the solution set $\mathcal{B}$ and $S(\widehat{\boldsymbol{\beta}}) \leq S(\boldsymbol{\beta}_*) + r$ for some constant $r > 0$ with probability tending to 1 as $n_0 \to \infty$. Additionally, we ignore the Monte Carlo approximation error $O\left(\frac{1}{\sqrt{M}}\right)$ in (3) for sufficiently large $M > 0$ and focus on the sums in (7).

Under Assumption 3.5(c) and Taylor's expansion, we can always choose $\boldsymbol{\beta}_* \in \mathcal{B}$ such that $S(\widehat{\boldsymbol{\beta}}) - S(\boldsymbol{\beta}_*) \geq \frac{\lambda_{\min}}{2} \left\|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_*\right\|_2^2$ when $n_0$ is sufficiently large. Hence, with probability tending to 1 as $n_0 \to \infty$, we have that

$$\left\|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_*\right\|_1 \leq \sqrt{K+1} \left\|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_*\right\|_2 \leq \sqrt{\frac{2(K+1)}{\lambda_{\min}} \left[ S(\widehat{\boldsymbol{\beta}}) - S(\boldsymbol{\beta}_*) \right]}, \tag{9}$$

and it suffices to prove the rate of convergence for the excess risk $S(\widehat{\boldsymbol{\beta}}) - S(\boldsymbol{\beta}_*)$. To this end, we compute that

$$S(\widehat{\boldsymbol{\beta}}) - S(\boldsymbol{\beta}_*) = S(\widehat{\boldsymbol{\beta}}) - \widehat{S}_{n_0}(\widehat{\boldsymbol{\beta}}) + \underbrace{\widehat{S}_{n_0}(\widehat{\boldsymbol{\beta}}) - \widehat{S}_{n_0}(\boldsymbol{\beta}_*)}_{\leq 0} + \widehat{S}_{n_0}(\boldsymbol{\beta}_*) - S(\boldsymbol{\beta}_*)$$

$$\leq S(\widehat{\boldsymbol{\beta}}) - S_{n_0}(\widehat{\boldsymbol{\beta}}) + S_{n_0}(\widehat{\boldsymbol{\beta}}) - \widehat{S}_{n_0}(\widehat{\boldsymbol{\beta}}) + \widehat{S}_{n_0}(\boldsymbol{\beta}_*) - S_{n_0}(\boldsymbol{\beta}_*) + S_{n_0}(\boldsymbol{\beta}_*) - S(\boldsymbol{\beta}_*)$$

$$\leq 2 \sup_{\boldsymbol{\beta}: S(\boldsymbol{\beta}) \leq S(\boldsymbol{\beta}_*) + r} |S_{n_0}(\boldsymbol{\beta}) - S(\boldsymbol{\beta})| + 2 \sup_{\boldsymbol{\beta}: \|\boldsymbol{\beta} - \boldsymbol{\beta}_*\|_2 \leq R_\beta} \left|\widehat{S}_{n_0}(\boldsymbol{\beta}) - S_{n_0}(\boldsymbol{\beta})\right| \tag{10}$$

$$\overset{(i)}{\leq} 2 \sup_{\boldsymbol{\beta}: S(\boldsymbol{\beta}) \leq S(\boldsymbol{\beta}_*) + r} |S_{n_0}(\boldsymbol{\beta}) - S(\boldsymbol{\beta})| + 8B_g C_\beta^2 \sqrt{K+1} \sum_{k=1}^K \left\|\widehat{g}^{(k)} - g^{(k)}\right\|_\infty,$$

where (i) follows from Lemma C.4. It remains to derive the rate of convergence for $\sup_{\boldsymbol{\beta}: \|\boldsymbol{\beta} - \boldsymbol{\beta}_*\|_2 \leq R_\beta} |S_{n_0}(\boldsymbol{\beta}) - S(\boldsymbol{\beta})|$. For any fixed $\boldsymbol{\beta}$ with $\|\boldsymbol{\beta} - \boldsymbol{\beta}_*\|_2 \leq R_\beta$, we know that

$$|S_{n_0}(\boldsymbol{\beta}) - S(\boldsymbol{\beta})|$$

$$= \left| \frac{1}{n_0} \sum_{j=1}^{n_0} \left[ Q_{n_0, Y^{(0)}}\left(\frac{j}{n_0}\right) - Q_{n_0, \boldsymbol{\beta}^T \boldsymbol{V}}\left(\frac{j}{n_0}\right) \right]^2 - \int_0^1 \left[ Q_{Y^{(0)}}(\alpha) - Q_{\boldsymbol{\beta}^T \boldsymbol{V}}(\alpha) \right]^2 d\alpha \right|$$

$$\leq \left| \frac{1}{n_0} \sum_{j=1}^{n_0} \left[ Q_{n_0, Y^{(0)}} \left( \frac{j}{n_0} \right) - Q_{n_0, \boldsymbol{\beta}^T \boldsymbol{V}} \left( \frac{j}{n_0} \right) \right]^2 - \frac{1}{n_0} \sum_{j=1}^{n_0} \left[ Q_{Y^{(0)}} \left( \frac{j}{n_0} \right) - Q_{\boldsymbol{\beta}^T \boldsymbol{V}} \left( \frac{j}{n_0} \right) \right]^2 \right|$$

$$+ \left| \frac{1}{n_0} \sum_{j=1}^{n_0} \left[ Q_{Y^{(0)}} \left( \frac{j}{n_0} \right) - Q_{\boldsymbol{\beta}^T \boldsymbol{V}} \left( \frac{j}{n_0} \right) \right]^2 - \int_0^1 \left[ Q_{Y^{(0)}}(\alpha) - Q_{\boldsymbol{\beta}^T \boldsymbol{V}}(\alpha) \right]^2 d\alpha \right|$$

$$\leq \frac{2}{n_0} \sum_{j=1}^{n_0} \left[ Q_{n_0, Y^{(0)}} \left( \frac{j}{n_0} \right) - Q_{Y^{(0)}} \left( \frac{j}{n_0} \right) \right]^2 + \frac{2}{n_0} \sum_{j=1}^{n_0} \left[ Q_{n_0, \boldsymbol{\beta}^T \boldsymbol{V}} \left( \frac{j}{n_0} \right) - Q_{\boldsymbol{\beta}^T \boldsymbol{V}} \left( \frac{j}{n_0} \right) \right]^2$$

$$+ \left| \frac{2}{n_0} \sum_{j=1}^{n_0} \left[ Q_{Y^{(0)}} \left( \frac{j}{n_0} \right) - Q_{\boldsymbol{\beta}^T \boldsymbol{V}} \left( \frac{j}{n_0} \right) \right] \left[ Q_{Y^{(0)}} \left( \frac{j}{n_0} \right) - Q_{n, Y^{(0)}} \left( \frac{j}{n_0} \right) + Q_{n, \boldsymbol{\beta}^T \boldsymbol{V}} \left( \frac{j}{n_0} \right) - Q_{\boldsymbol{\beta}^T \boldsymbol{V}} \left( \frac{j}{n_0} \right) \right] \right|$$

$$+ \left| \frac{1}{n_0} \sum_{j=1}^{n_0} \left[ Q_{Y^{(0)}} \left( \frac{j}{n_0} \right) - Q_{\boldsymbol{\beta}^T \boldsymbol{V}} \left( \frac{j}{n_0} \right) \right]^2 - \int_0^1 \left[ Q_{Y^{(0)}}(\alpha) - Q_{\boldsymbol{\beta}^T \boldsymbol{V}}(\alpha) \right]^2 d\alpha \right|$$

$$\stackrel{(ii)}{=} \frac{2}{n_0} \sum_{j=1}^{n_0} \left[ \frac{F_{n, Y^{(0)}}(Q_{Y^{(0)}}(j/n_0)) - j/n_0}{p_{Y^{(0)}}(Q_{Y^{(0)}}(j/n_0))} + \frac{r_{n_0}}{\sqrt{n_0}} \right]^2 + \frac{2}{n_0} \sum_{j=1}^{n_0} \left[ \frac{F_{n, \boldsymbol{\beta}^T \boldsymbol{V}}(Q_{\boldsymbol{\beta}^T \boldsymbol{V}}(j/n_0)) - j/n_0}{p_{\boldsymbol{\beta}^T \boldsymbol{V}}(Q_{\boldsymbol{\beta}^T \boldsymbol{V}}(j/n_0))} + \frac{r_{n_0}}{\sqrt{n_0}} \right]^2$$

$$+ \sqrt{\frac{1}{n_0} \sum_{j=1}^{n_0} \left| Q_{Y^{(0)}} \left( \frac{j}{n_0} \right) - Q_{\boldsymbol{\beta}^T \boldsymbol{V}} \left( \frac{j}{n_0} \right) \right|^2}$$

$$\times \sqrt{\frac{2}{n_0} \sum_{j=1}^{n_0} \left| \frac{F_{n, Y^{(0)}}(Q_{Y^{(0)}}(j/n_0)) - j/n_0}{p_{Y^{(0)}}(Q_{Y^{(0)}}(j/n_0))} + \frac{F_{n, \boldsymbol{\beta}^T \boldsymbol{V}}(Q_{\boldsymbol{\beta}^T \boldsymbol{V}}(j/n_0)) - j/n_0}{p_{\boldsymbol{\beta}^T \boldsymbol{V}}(Q_{\boldsymbol{\beta}^T \boldsymbol{V}}(j/n_0))} + \frac{r_{n_0}}{\sqrt{n_0}} \right|^2} + O\left( \frac{1}{n_0^2} \right)$$

$$\stackrel{(iii)}{=} \frac{4}{n_0} \sum_{i=1}^{n_0} \left[ u_{n_0} + \frac{r_{n_0}}{\sqrt{n_0}} \right]^2 + 2 \left[ u_{n_0} + \frac{r_{n_0}}{\sqrt{n_0}} \right] \sqrt{\int_0^1 \left[ Q_{Y^{(0)}}(\alpha) - Q_{\boldsymbol{\beta}^T \boldsymbol{V}}(\alpha) \right]^2 d\alpha} + O\left( \frac{1}{n_0^2} \right)$$

$$= O_P \left( \sqrt{\frac{S(\boldsymbol{\beta}) \log \log n_0}{n_0}} \right),$$

where (ii) leverages Lemma C.1 and Cauchy-Schwarz inequality with $r_n = O_P \left( \frac{(\log n)^{\frac{1}{2}} (\log \log n)^{\frac{1}{4}}}{n^{\frac{1}{4}}} \right)$ as well as the fact that the Riemann sum approximation has its error bound as

$$\left| \frac{1}{n_0} \sum_{i=1}^{n_0} \left[ Q_{Y^{(0)}} \left( \frac{j}{n_0} \right) - Q_{\boldsymbol{\beta}^T \boldsymbol{V}} \left( \frac{j}{n_0} \right) \right]^2 - \int_0^1 \left[ Q_{Y^{(0)}}(\alpha) - Q_{\boldsymbol{\beta}^T \boldsymbol{V}}(\alpha) \right]^2 d\alpha \right| = O\left( \frac{1}{n_0^2} \right)$$

under Assumption 3.5 and midpoint rule, and (iii) utilizes the Dvoretzky-Kiefer-Wolfowitz inequality (Massart, 1990) with $u_n = O_P \left( \sqrt{\frac{\log \log n_0}{n_0}} \right)$ and

$$\mathbb{P} \left( \sup_{\alpha \in [0,1]} \left| F_{n, Y^{(0)}}(Q_{Y^{(0)}}(\alpha)) - \alpha \right| > C \right) \leq 2e^{-2nC^2} \quad \text{and} \quad \mathbb{P} \left( \sup_{\alpha \in [0,1]} \left| F_{n, \boldsymbol{\beta}^T \boldsymbol{V}}(Q_{\boldsymbol{\beta}^T \boldsymbol{V}}(\alpha)) - \alpha \right| > C \right) \leq 2e^{-2nC^2}.$$

Furthermore, the above display implies that

$$\sup_{\boldsymbol{\beta}: S(\boldsymbol{\beta}) \leq S(\boldsymbol{\beta}_*) + r} |S_{n_0}(\boldsymbol{\beta}) - S(\boldsymbol{\beta})| = O_P \left( \sqrt{\frac{(S(\boldsymbol{\beta}_*) + r) \log \log n_0}{n_0}} \right). \tag{11}$$

Hence, if we assume that $r_n = S(\widehat{\boldsymbol{\beta}}) - S(\boldsymbol{\beta}_*)$, then by (10) and (11), we obtain that

$$r_n = S(\widehat{\boldsymbol{\beta}}) - S(\boldsymbol{\beta}_*) \leq 2 \sup_{\boldsymbol{\beta}: S(\boldsymbol{\beta}) \leq S(\boldsymbol{\beta}_*) + r} |S_{n_0}(\boldsymbol{\beta}) - S(\boldsymbol{\beta})| + 8 B_g C_\beta^2 \sqrt{K+1} \sum_{k=1}^{K} \left\| \widehat{g}^{(k)} - g^{(k)} \right\|_\infty$$

$$\lesssim \sqrt{\frac{(S(\boldsymbol{\beta}_*) + r_n) \log \log n_0}{n_0}} + \sqrt{K} \sum_{k=1}^{K} \left\| \widehat{g}^{(k)} - g^{(k)} \right\|_\infty$$

in probability. Solving for $r_n$ in this (probabilistic) inequality yields that

$$r_n \lesssim \sqrt{\frac{S(\boldsymbol{\beta}_*) \log \log n_0}{n_0}} + \frac{\log \log n_0}{n_0} + \sqrt{K} \sum_{k=1}^{K} \left\| \widehat{g}^{(k)} - g^{(k)} \right\|_\infty.$$

Therefore, we conclude from (9) and the above calculations that

$$\inf_{\boldsymbol{\beta}_* \in \mathcal{B}} \left\| \widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_* \right\|_1 = O_P \left( \sqrt{K} \left( \frac{\log \log n_0}{n_0} \int_0^1 \left[ Q_{Y^{(0)}}(\alpha) - Q_{\boldsymbol{\beta}_*^T \boldsymbol{V}}(\alpha) \right]^2 d\alpha \right)^{\frac{1}{4}} + \sqrt{\frac{K \log \log n_0}{n_0}} \right)$$
$$+ O \left( \sqrt{K \sum_{k=1}^{K} \left\| \widehat{g}^{(k)} - g^{(k)} \right\|_\infty} \right).$$

The result thus follows. $\qquad \square$