

Table 2: Information of benchmark datasets.

Dataset	Datasize	d	Source
housing	404	13	UCI Repository
airfoil	1202	5	UCI Repository
concrete	824	8	UCI Repository
powerplant	7654	4	UCI Repository
mpg	313	7	UCI Repository
redwine	1279	11	UCI Repository
whitewine	3918	11	UCI Repository
abalone	3341	10	UCI Repository
diabetes	353	10	[10]

A Experiments details

In this appendix, we explain the detailed setting of experiments. First, we describe the procedure of hyper-parameter tuning during the experiments. Then, we provide detailed information on benchmark datasets.

A.1 Procedure of hyper-parameter tuning

To construct risk \hat{R}_{RA} , we need to tune λ, w_1, w_2 , which is done by minimizing empirically approximated $\text{Err}(w_1, w_2)$ defined in (8). Let \bar{y}, \underline{y} be the 0.99-quantile and 0.01-quantile of F_Y , respectively. Note that we can calculate these quantities since we have access to f_Y . Then, we define $\{y^{(i)}\}_{i=1}^{n_{\text{split}}+1}$ as $y_i = \underline{y} + (i-1)/n_{\text{split}}(\bar{y} - \underline{y})$, by which $\text{Err}(w_1, w_2)$ is approximated as

$$\text{Err}(w_1, w_2) \simeq \sum_{i=1}^{n_{\text{split}}+1} f_Y(y_i) |y_i - w_1 F_Y(y_i) - w_2 (1 - F_Y(y_i))|.$$

We employ w_1, w_2 that minimize the empirical approximation above with $n_{\text{split}} = 1000$ and fix λ to be $(w_1 + w_2)/2$ in all cases.

For the TT method, we employ hypothesis space $\mathcal{H}' = \{h(\mathbf{x}) = F_Y^{-1}(\sigma(\boldsymbol{\theta}^\top \mathbf{x})) \mid \boldsymbol{\theta} \in \mathbb{R}^d\}$, which is slightly different from hypothesis space of liner functions $\mathcal{H} = \{h(\mathbf{x}) = \boldsymbol{\theta}^\top \mathbf{x} \mid \boldsymbol{\theta} \in \mathbb{R}^d\}$, where σ is logistic function $\sigma(x) = 1/(1 + \exp(-x))$. This simplifies the loss function and reduces the computational time. We fix $\lambda = 1/2$ for this risk, which yields the loss

$$\begin{aligned} \hat{R}_{\text{TT}}(h) = & \mathfrak{C} - \frac{1}{n_{\text{U}}} \sum_{\mathbf{x}_i \in \mathcal{D}_{\text{U}}} \left(\frac{1}{2} - \sigma(\boldsymbol{\theta}^\top \mathbf{x}_i) \right) \phi'(\sigma(h(\mathbf{x}_i))) + \phi(\sigma(\boldsymbol{\theta}^\top \mathbf{x}_i)) \\ & - \frac{1}{n_{\text{R}}} \sum_{(\mathbf{x}_i^+, \mathbf{x}_i^-) \in \mathcal{D}_{\text{R}}} \frac{1}{4} \phi'(\sigma(\boldsymbol{\theta}^\top \mathbf{x}_i^+)) - \frac{1}{4} \phi'(\sigma(\boldsymbol{\theta}^\top \mathbf{x}_i^-)). \end{aligned}$$

We minimize this loss with respect to $\boldsymbol{\theta}$.

A.2 Benchmark dataset details

We use eight benchmark datasets from UCI repository [8] and one (diabetes) from Efron et al. [10]. The details of datasets can be found in Table 2. As preprocessing, we excluded all instances that contain missing values, and we encoded a categorical feature in abalone as a one-hot vector.

B Estimating density function and cumulative distribution function

In this section, we discuss the case where the true probability density function f_Y is not given. In such a case, we need a slight modification of proposed methods since we have to estimate f_Y from the set of target values $\mathcal{D}_Y = \{y_i\}_{i=1}^{n_Y}$, where n_Y is the size of \mathcal{D}_Y . We first introduce a modification of the RA method and derive an estimation error bound for it. Then, we discuss the same for the TT method as well.

B.1 Modification of the risk approximation method

Although \hat{R}_{RA} does not depend on f_Y or F_Y , we need the information of P_Y when tuning weights w_1, w_2 , which is done by the minimization of Err defined in (8). Since, Err can not be directly calculated without f_Y and F_Y , we propose another quantity $\widehat{\text{Err}}$ below, which substitute expectation over P_Y and CDF function F_Y to empirical mean and the empirical CDF.

$$\widehat{\text{Err}}(w_1, w_2) = \frac{1}{n_Y} \sum_{i=1}^{n_Y} |y_i - w_1 \hat{F}_Y(y_i) - w_2 (1 - \hat{F}_Y(y_i))|,$$

where \hat{F}_Y is the empirical CDF defined as

$$\hat{F}_Y(y) = \frac{1}{n_Y} \sum_{i=1}^{n_Y} \mathbb{1}[y_i \leq y].$$

Note that $\widehat{\text{Err}}$ can be minimized given \mathcal{D}_Y . To show the validity of the method, we establish an estimation error bound involving $\widehat{\text{Err}}$ as follows.

Theorem 5. *Let \mathcal{Y} be bounded in $\mathcal{Y} \subseteq [-L, L]$. Then, for all $w_1, w_2 \in [-L, L]$, we have*

$$|\text{Err}(w_1, w_2) - \widehat{\text{Err}}(w_1, w_2)| \leq O\left(\sqrt{\frac{\log \delta}{n_Y}}\right)$$

with probability $1 - 2\delta$.

Proof. Since the weights are bounded, from Mohri et al. [22, Thm. 10.3], we have

$$\text{Err}(w_1, w_2) \leq \frac{1}{n_Y} \sum_{i=1}^{n_Y} |y_i - w_1 F_Y(y_i) - w_2 (1 - F_Y(y_i))| + O\left(\sqrt{\frac{\log 1/\delta}{m}}\right),$$

with probability $1 - \delta$. Furthermore, from Dvoretzky-Kiefer-Wolfowitz inequality [20], we have

$$\|F_Y(y) - \hat{F}_Y(y)\|_\infty \leq \sqrt{\frac{\log(2/\delta)}{2n_Y}} \quad (11)$$

with probability $1 - \delta$, which yields

$$\frac{1}{n_Y} \sum_{i=1}^{n_Y} |y_i - w_1 F_Y(y_i) - w_2 (1 - F_Y(y_i))| \leq \widehat{\text{Err}} + O\left(\sqrt{\frac{\log 1/\delta}{m}}\right).$$

Therefore, from the union bound, we have

$$|\text{Err}(w_1, w_2) - \widehat{\text{Err}}(w_1, w_2)| \leq O\left(\sqrt{\frac{\log \delta}{n_Y}}\right)$$

with probability $1 - 2\delta$. □

From Theorems 2 and 5, we have

$$R(\hat{h}_{\text{RA}}) \leq R(h^*) + O\left(\sqrt{\frac{\log 1/\delta}{n_U}}\right) + O\left(\sqrt{\frac{\log 1/\delta}{n_R}}\right) + O\left(\sqrt{\frac{\log 1/\delta}{n_Y}}\right) + M\widehat{\text{Err}}(w_1, w_2),$$

with probability $1 - 5\delta$ under the conditions given in these theorems.

B.2 Modification on the target transformation method

Let \tilde{R}_{TT} be the risk which substitute F_Y in R_{TT} to empirical CDF, defined as

$$\begin{aligned}\tilde{R}_{\text{TT}}(h; \lambda) = & \mathfrak{C} - \frac{1}{n_U} \sum_{\mathbf{x}_i \in \mathcal{D}_U} \left((\lambda - \hat{F}_Y(h(\mathbf{x}_i))) \phi'(\hat{F}_Y(h(\mathbf{x}_i))) + \phi(\hat{F}_Y(h(\mathbf{x}_i))) \right) \\ & - \frac{1}{n_R} \sum_{(\mathbf{x}_i^+, \mathbf{x}_i^-) \in \mathcal{D}_R} \left(\frac{1-\lambda}{2} \phi'(\hat{F}_Y(h(\mathbf{x}_i^+))) - \frac{\lambda}{2} \phi'(\hat{F}_Y(h(\mathbf{x}_i^-))) \right).\end{aligned}$$

Using (11), we have

$$|\hat{R}_{\text{TT}}(h) - \tilde{R}_{\text{TT}}(h)| \leq O \left(\sqrt{\frac{\log 1/\delta}{n_Y}} \right)$$

for all $h \in \mathcal{H}$ with probability $1 - \delta$. Let \tilde{h}_{TT} be the minimizer of \tilde{R}_{TT} in hypothesis space \mathcal{H} . Then, under the condition given in Theorem 4, we have

$$R_{\text{TT}}(\tilde{h}_{\text{TT}}) \leq R_{\text{TT}}(h_{\text{TT}}) + O \left(\sqrt{\frac{\log 1/\delta}{n_Y}} \right) + O \left(\sqrt{\frac{\log 1/\delta}{n_R}} \right) + O \left(\sqrt{\frac{\log 1/\delta}{n_U}} \right),$$

with probability $1 - 4\delta$, therefore we have

$$R(\tilde{h}_{\text{TT}}) \leq R(h^*) + 2 \left(\frac{P}{p} \sigma \right)^2 + O \left(\sqrt{\frac{\log 1/\delta}{n_Y}} \right) + O \left(\sqrt{\frac{\log 1/\delta}{n_R}} \right) + O \left(\sqrt{\frac{\log 1/\delta}{n_U}} \right)$$

with probability $1 - 4\delta$, which can be shown by the slight modification of the proof of Theorem 4.

C Proofs

C.1 Proof of Lemma 1

Lemma 1 can be proved as follows.

Proof of Lemma 1. Let $f_{\mathbf{X}^+}$ be the probability density function (PDF) of $P_{\mathbf{X}^+}$. From the definition of \mathbf{X}^+ , we have

$$\begin{aligned}f_{\mathbf{X}^+}(\mathbf{x}) &= \frac{1}{Z} \iiint f_{\mathbf{X},Y}(\mathbf{x}, y) f_{\mathbf{X},Y}(\mathbf{x}', y') \mathbb{1}[y > y'] dy dy' d\mathbf{x}' \\ &= \frac{1}{Z} \int f_{\mathbf{X},Y}(\mathbf{x}, y) \left[\int f_Y(y') \mathbb{1}[y > y'] dy' \right] dy \\ &= \frac{1}{Z} \int f_{\mathbf{X},Y}(\mathbf{x}, y) F_Y(y) dy,\end{aligned}$$

where Z is the normalizing constant and $f_{\mathbf{X},Y}(y)$ is the PDF of $P_{\mathbf{X},Y}$. Now, Z is calculated as

$$\begin{aligned}Z &= \iint f_{\mathbf{X},Y}(\mathbf{x}, y) F_Y(y) dy d\mathbf{x} \\ &= \int f_Y(y) F_Y(y) dy \\ &= \frac{1}{2},\end{aligned}$$

where the last equality holds from the integration by parts. Therefore, we have

$$\begin{aligned}\mathbb{E}_{\mathbf{X}^+} [\phi'(\mathbf{X}^+)] &= \int f_{\mathbf{X}^+}(\mathbf{x}) \phi'(\mathbf{x}) d\mathbf{x} \\ &= \int 2 \left\{ \int f_{\mathbf{X},Y}(\mathbf{x}, y) F_Y(y) dy \right\} \phi'(\mathbf{x}) d\mathbf{x} \\ &= \mathbb{E}_{\mathbf{X},Y} [F_Y(Y) \phi'(\mathbf{x})].\end{aligned}$$

The expectation over $P_{\mathbf{X}^-}$ can be derived in the same way. \square

C.2 Proof of Theorem 2

Here, we show the proof of Theorem 2. First, we show the gap between R and R_{RA} can be bounded as follows.

Lemma 2. *For all $h \in \mathcal{H}$, such that $|\phi'(h(\mathbf{x}))| \leq M$ for all $\mathbf{x} \in \mathcal{X}$, we have*

$$|R(h) - R_{\text{RA}}(h; \lambda; w_1, w_2)| \leq M \text{Err}(w_1, w_2)$$

for all $\lambda \in \mathbb{R}$.

Proof. From Lemma 1 and the fact $\mathbb{E}_{\mathbf{X}} [\phi'(\mathbf{X})] = \frac{1}{2} \mathbb{E}_{\mathbf{X}^+} [\phi'(\mathbf{X}^+)] + \frac{1}{2} \mathbb{E}_{\mathbf{X}^-} [\phi'(\mathbf{X}^-)]$, we have

$$\begin{aligned} & |R(h) - R_{\text{RA}}(h; \lambda, w_1, w_2)| \\ &= |\mathbb{E}_{\mathbf{X}, Y} [Y \phi'(h(\mathbf{X}))] - w_1 \mathbb{E}_{\mathbf{X}^+} [\phi'(h(\mathbf{X}^+))] - w_2 \mathbb{E}_{\mathbf{X}^-} [\phi'(h(\mathbf{X}^-))]| \\ &= \left| \int f_{\mathbf{X}, Y}(\mathbf{x}, y) \phi'(h(\mathbf{x})) \{y - 2w_1 F_Y(y) - 2w_2(1 - F_Y(y))\} dy d\mathbf{x} \right| \\ &\leq \int f_{\mathbf{X}, Y}(\mathbf{x}, y) |\phi'(h(\mathbf{x}))| |y - 2w_1 F_Y(y) - 2w_2(1 - F_Y(y))| dy d\mathbf{x} \\ &\leq M \int f_Y(y) |y - 2w_1 F_Y(y) - 2w_2(1 - F_Y(y))| dy \\ &\leq M \text{Err}(w_1, w_2). \end{aligned}$$

□

Now, Theorem 2 can be derived as follows.

Proof of Theorem 2. Let \tilde{d}, \tilde{d}' be the pseudo-dimensions defined as

$$\begin{aligned} \tilde{d} &= \text{Pdim}(\{\mathbf{x} \rightarrow \phi'(h(\mathbf{x})) \mid h \in \mathcal{H}\}), \\ \tilde{d}' &= \text{Pdim}(\{\mathbf{x} \rightarrow h(\mathbf{x})\phi'(h(\mathbf{x})) - \phi(h(\mathbf{x})) \mid h \in \mathcal{H}\}), \end{aligned}$$

where $\text{Pdim}(\mathcal{F})$ denotes the pseudo-dimension of the functional space \mathcal{F} . From the assumptions in Theorem 2, using the discussion in Mohri et al. [22, Theorem 10.6], each of following bound holds with probability $1 - \delta$ for all $h \in \mathcal{H}$.

$$\begin{aligned} & \left| \mathbb{E}_{\mathbf{X}^+} [\phi'(h(\mathbf{X}^+))] - \frac{1}{n_{\text{R}}} \sum_{\mathbf{x}_i^+ \in \mathcal{D}_{\text{R}}^+} \phi'(h(\mathbf{x}_i^+)) \right| \leq M \sqrt{\frac{2\tilde{d} \log \frac{en_{\text{R}}}{d}}{n_{\text{R}}}} + M \sqrt{\frac{\log \frac{1}{\delta}}{2n_{\text{R}}}}, \\ & \left| \mathbb{E}_{\mathbf{X}^-} [\phi'(h(\mathbf{X}^-))] - \frac{1}{n_{\text{R}}} \sum_{\mathbf{x}_i^- \in \mathcal{D}_{\text{R}}^-} \phi'(h(\mathbf{x}_i^-)) \right| \leq M \sqrt{\frac{2\tilde{d} \log \frac{en_{\text{R}}}{d}}{n_{\text{R}}}} + M \sqrt{\frac{\log \frac{1}{\delta}}{2n_{\text{R}}}}, \\ & \left| \mathbb{E}_{\mathbf{X}} [g(\mathbf{X})] - \frac{1}{n_{\text{U}}} \sum_{\mathbf{x}_i \in \mathcal{D}_{\text{U}}^+} g(\mathbf{x}_i) \right| \leq m \sqrt{\frac{2\tilde{d}' \log \frac{en_{\text{U}}}{d'}}{n_{\text{U}}}} + m \sqrt{\frac{\log \frac{1}{\delta}}{2n_{\text{U}}}}, \end{aligned}$$

where $g(\mathbf{x}) = h(\mathbf{x})\phi'(h(\mathbf{x})) + \phi(h(\mathbf{x}))$. From the uniform bound, we have

$$\begin{aligned} & |R_{\text{RA}}(h; w_1, w_2) - \hat{R}_{\text{RA}}(h; \lambda, w_1, w_2)| \\ &\leq \left(\left| w_1 - \frac{\lambda}{2} \right| + \left| w_2 - \frac{\lambda}{2} \right| \right) \left(M \sqrt{\frac{2\tilde{d} \log \frac{en_{\text{R}}}{d}}{n_{\text{R}}}} + M \sqrt{\frac{\log \frac{1}{\delta}}{2n_{\text{R}}}} \right) \\ &\quad + (m + \lambda M) \left(\sqrt{\frac{2\tilde{d}' \log \frac{en_{\text{U}}}{d'}}{n_{\text{U}}}} + \sqrt{\frac{\log \frac{1}{\delta}}{2n_{\text{U}}}} \right) \end{aligned}$$

with probability $1 - 3\delta$ for all $h \in \mathcal{H}$. Hence, with probability $1 - 3\delta$, we have

$$\begin{aligned}
& R(\hat{h}_{\text{RA}}) - R(h^*) \\
& \leq R_{\text{RA}}(\hat{h}_{\text{RA}}; \lambda, w_1, w_2) - R_{\text{RA}}(h^*; \lambda, w_1, w_2) + |R(h^*) - R_{\text{RA}}(h^*; \lambda, w_1, w_2)| \\
& \quad + |R(\hat{h}_{\text{RA}}) - R_{\text{RA}}(\hat{h}_{\text{RA}}; \lambda, w_1, w_2)| \\
& \leq (R_{\text{RA}}(\hat{h}_{\text{RA}}; \lambda, w_1, w_2) - \hat{R}_{\text{RA}}(h^*; \lambda, w_1, w_2)) \\
& \quad - (R_{\text{RA}}(h^*; \lambda, w_1, w_2) - \hat{R}_{\text{RA}}(h^*; \lambda, w_1, w_2)) + 2M\text{Err}(w_1, w_2) \\
& \leq (R_{\text{RA}}(\hat{h}_{\text{RA}}; \lambda, w_1, w_2) - \hat{R}_{\text{RA}}(\hat{h}_{\text{RA}}; \lambda, w_1, w_2)) \\
& \quad - (R_{\text{RA}}(h^*; \lambda, w_1, w_2) - \hat{R}_{\text{RA}}(h^*; \lambda, w_1, w_2)) + 2M\text{Err}(w_1, w_2) \\
& \leq O\left(\sqrt{\frac{\log 1/\delta}{n_{\text{U}}}}\right) + O\left(\sqrt{\frac{\log 1/\delta}{n_{\text{R}}}}\right) + 2M\text{Err}(w_1, w_2),
\end{aligned}$$

where the second inequality holds from the fact $\hat{R}_{\text{RA}}(\hat{h}_{\text{RA}}; \lambda, w_1, w_2) \leq \hat{R}_{\text{RA}}(\hat{h}^*; \lambda, w_1, w_2)$ and Lemma 2. \square

C.3 Proof of Theorem 1

Theorem 1 can be shown as follows.

Proof of Theorem 1. The variance of \hat{R}_{RA} denoted as $\text{Var} \left[\hat{R}_{\text{RA}}(h; \lambda, w_1, w_2) \right]$ can be expressed as

$$\text{Var} \left[\hat{R}_{\text{RA}}(h; \lambda, w_1, w_2) \right] = \left(w_1 - \frac{\lambda}{2} \right)^2 \frac{\sigma_+^2}{n_{\text{R}}} + \left(w_2 - \frac{\lambda}{2} \right)^2 \frac{\sigma_-^2}{n_{\text{R}}}$$

when $n_{\text{U}} \rightarrow \infty$. By solving the above quadratic optimization problem, we have

$$\arg \min_{\lambda} \text{Var} \left[\hat{R}_{\text{RA}}(h; \lambda, w_1, w_2) \right] = \frac{2(w_1\sigma_+^2 + w_2\sigma_-^2)}{\sigma_+^2 + \sigma_-^2}.$$

\square

C.4 Proof of Theorem 3

We can construct a simple example satisfies the conditions in Theorem 3 as follows.

Proof. Let $f_{\mathbf{X},Y}, \tilde{f}_{\mathbf{X},Y}$ be the PDF of $P_{\mathbf{X},Y}, \tilde{P}_{\mathbf{X},Y}$, respectively. If we consider $\mathcal{X} = [-1, 1]$ and $\mathcal{Y} = [0, 4]$ and these PDF to be

$$\begin{aligned}
f_{\mathbf{X},Y}(x, y) &= \begin{cases} \frac{1}{6} & (y \in [0, 2] \cup [3, 4]), \\ 0 & (\text{otherwise}), \end{cases} \\
\tilde{f}_{\mathbf{X},Y}(x, y) &= \begin{cases} \frac{1}{8} & (x \in [-1, 0), y \in [0, 1]), \\ \frac{1}{4} & (x \in [-1, 0), y \in [1, 2]), \\ \frac{1}{8} & (x \in [-1, 0), y \in [3, 4]), \\ \frac{5}{24} & (x \in [0, 1], y \in [0, 1]), \\ \frac{1}{12} & (x \in [0, 1], y \in [1, 2]), \\ \frac{5}{24} & (x \in [0, 1], y \in [3, 4]), \\ 0 & (\text{otherwise}). \end{cases}
\end{aligned}$$

Then, by the simple calculation, we can see that they have the same PDF $f_{\mathbf{X}}(x), f_Y(y), f_{\mathbf{X}^+, \mathbf{X}^-}(x^+, x^-)$, each represents the PDF of $P_{\mathbf{X}}, P_Y, P_{\mathbf{X}^+, \mathbf{X}^-}$, respectively, which are

$$\begin{aligned}
f_{\mathbf{X}}(x) &= 0.5, \\
f_Y(y) &= \begin{cases} \frac{1}{3} & (y \in [0, 2] \cup [3, 4]), \\ 0 & (\text{otherwise}), \end{cases} \\
f_{\mathbf{X}^+, \mathbf{X}^-}(x^+, x^-) &= 0.25.
\end{aligned}$$

However, the conditional expectation $\mathbb{E}_{Y|\mathbf{X}=x}[Y]$ defined on $P_{\mathbf{X},Y}$ is

$$\mathbb{E}_{Y|\mathbf{X}=x}[Y] = \frac{11}{6},$$

while the conditional expectation $\tilde{\mathbb{E}}_{Y|\mathbf{X}=x}[Y]$ defined on $\tilde{P}_{\mathbf{X},Y}$ is

$$\tilde{\mathbb{E}}_{Y|\mathbf{X}=x}[Y] = \begin{cases} \frac{7}{4} & (x \in [-1, 0)), \\ \frac{23}{12} & (x \in [0, 1]). \end{cases}$$

□

C.5 Proof of Theorem 4

The Theorem 4 can be shown as follows.

Proof of Theorem 4. We first show that under the conditions, we have

$$\|h_{\text{true}}(\mathbf{x}) - h_{\text{TT}}(\mathbf{x})\|_{\infty} \leq \frac{\sigma P}{p}.$$

Since $(F_Y(y))' = f_Y(y) \leq P$ and $(F_Y^{-1}(q))' = 1/f_Y(F^{-1}(q)) \leq 1/p$ for any $y \in \mathcal{Y}$ and $q \in [0, 1]$, $F_Y(y), F_Y^{-1}(q)$ are $P, 1/p$ -Lipschitz continuous, respectively. Therefore, we have

$$\begin{aligned} h_{\text{TT}}(\mathbf{x}) &= F_Y^{-1}(\mathbb{E}_{Y|\mathbf{X}=\mathbf{x}}[F_Y(Y)]) \\ &= F_Y^{-1}(\mathbb{E}_{\epsilon}[F_Y(h_{\text{true}}(\mathbf{x}) + \epsilon)]) \\ &\leq F_Y^{-1}(F_Y(h_{\text{true}}(\mathbf{x}) + \sigma P)) \\ &\leq h_{\text{true}}(\mathbf{x}) + \frac{\sigma P}{p} \end{aligned}$$

for all $\mathbf{x} \in \mathcal{X}$. With the same discussion, we have $|h_{\text{TT}}(\mathbf{x}) - h_{\text{true}}(\mathbf{x})| \leq \frac{\sigma P}{p}$. Therefore, we have

$$\|h_{\text{true}}(\mathbf{x}) - h_{\text{TT}}(\mathbf{x})\|_{\infty} \leq \frac{\sigma P}{p}.$$

Now, if $\phi(x) = x^2$, which means $R(h) = \mathbb{E}_{\mathbf{X},Y}[(h(\mathbf{X}) - Y)^2]$, we have

$$\begin{aligned} R(\hat{h}_{\text{TT}}) &= \mathbb{E}_{\mathbf{X},Y}[(\hat{h}_{\text{TT}}(\mathbf{x}) - Y)^2] \\ &= \mathbb{E}_{\mathbf{X},Y}[(\hat{h}_{\text{TT}}(\mathbf{x}) - \hat{h}_{\text{true}}(\mathbf{x}) + \hat{h}_{\text{true}}(\mathbf{x}) - Y)^2] \\ &= \mathbb{E}_{\mathbf{X}}[(\hat{h}_{\text{TT}}(\mathbf{x}) - \hat{h}_{\text{true}}(\mathbf{x}))^2] + \mathbb{E}_{\mathbf{X},Y}[(\hat{h}_{\text{true}}(\mathbf{x}) - Y)^2] \\ &\quad + 2\mathbb{E}_{\mathbf{X},Y}[(\hat{h}_{\text{TT}}(\mathbf{x}) - \hat{h}_{\text{true}}(\mathbf{x}))(\hat{h}_{\text{true}}(\mathbf{x}) - Y)] \\ &= \mathbb{E}_{\mathbf{X}}[(\hat{h}_{\text{TT}}(\mathbf{X}) - \hat{h}_{\text{true}}(\mathbf{X}))^2] + \mathbb{E}_{\mathbf{X},Y}[(h_{\text{true}}(\mathbf{X}) - Y)^2] \\ &\leq R(h_{\text{true}}) + 2\mathbb{E}_{\mathbf{X}}[(\hat{h}_{\text{TT}}(\mathbf{x}) - h_{\text{TT}}(\mathbf{x}))^2] + 2\mathbb{E}_{\mathbf{X}}[(h_{\text{true}}(\mathbf{x}) - h_{\text{TT}}(\mathbf{x}))^2]. \end{aligned}$$

Since $\|h_{\text{TT}}(\mathbf{x}) - h_{\text{true}}(\mathbf{x})\|_{\infty} \leq \frac{\sigma P}{p}$, we have

$$\mathbb{E}_{\mathbf{X}}[(h_{\text{true}}(\mathbf{X}) - h_{\text{TT}}(\mathbf{X}))^2] \leq \left(\frac{\sigma P}{p}\right)^2.$$

Furthermore, using the characteristic of expectation, if $\phi(x) = x^2$, which means $R_{\text{TT}}(h) = \mathbb{E}_{\mathbf{X},Y}[(F_Y(h(\mathbf{X})) - F_Y(Y))^2]$, we have

$$\begin{aligned} R_{\text{TT}}(\hat{h}_{\text{TT}}) &= \mathbb{E}_{\mathbf{X},Y}[(F_Y(\hat{h}_{\text{TT}}(\mathbf{X})) - F_Y(Y))^2] \\ &= \mathbb{E}_{\mathbf{X},Y}[(F_Y(\hat{h}_{\text{TT}}(\mathbf{X})) - F_Y(h_{\text{TT}}(\mathbf{X})))^2] + \mathbb{E}_{\mathbf{X},Y}[(F_Y(Y) - F_Y(h_{\text{TT}}(\mathbf{X})))^2] \\ &= \mathbb{E}_{\mathbf{X},Y}[(F_Y(\hat{h}_{\text{TT}}(\mathbf{X})) - F_Y(h_{\text{TT}}(\mathbf{X})))^2] + R_{\text{TT}}(h_{\text{TT}}). \end{aligned}$$

Since $(F_Y(y))' \geq p$, we have

$$\begin{aligned}\mathbb{E}_{\mathbf{X}} \left[(\hat{h}_{\text{TT}}(\mathbf{X}) - h_{\text{TT}}(\mathbf{X}))^2 \right] &\leq \frac{1}{p^2} \mathbb{E}_{\mathbf{X}, Y} \left[(F_Y(\hat{h}_{\text{TT}}(\mathbf{X})) - F_Y(h_{\text{TT}}(\mathbf{X})))^2 \right] \\ &= \frac{1}{p^2} \left(R_{\text{TT}}(\hat{h}_{\text{TT}}) - R_{\text{TT}}(h_{\text{TT}}) \right) \\ &\leq O \left(\sqrt{\frac{\log 1/\delta}{n_{\text{U}}}} \right) + O \left(\sqrt{\frac{\log 1/\delta}{n_{\text{R}}}} \right)\end{aligned}$$

with probability $1 - 3\delta$, where the last inequality holds from the same discussion as in Theorem 2. Note that $|\phi'(F_Y(h(\mathbf{x})))|$, $|F_Y(h(\mathbf{x}))\phi'(F_Y(h(\mathbf{x}))) - \phi(F_Y(h(\mathbf{x})))|$ are bounded since $F_Y(h(\mathbf{x})) \in [0, 1]$ by definition. Combining these inequalities, we can see that

$$R(\hat{h}_{\text{TT}}) \leq R(h_{\text{true}}(\mathbf{x})) + 2 \left(\frac{\sigma P}{p} \right)^2 + O \left(\sqrt{\frac{\log 1/\delta}{n_{\text{U}}}} \right) + O \left(\sqrt{\frac{\log 1/\delta}{n_{\text{R}}}} \right)$$

with probability $1 - 3\delta$. □