# Supplemental file for "Efficient learning of nonparametric directed acyclic graph with statistical guarantee"

Yibo Deng[†], Xin He[†], and Shaogao Lv[‡*]

[†] School of Statistics and Management
Shanghai University of Finance and Economics
[†] School of Statistics and Mathematics
Nanjing Audit University

## Technical proofs

**Proof of Theorem 1.** For any $t = 0, ..., T - 1$, given $\mathcal{S}_t$, it is always true that $E \operatorname{Var}(x_j | \mathbf{x}_{\mathcal{S}_t}) = E \operatorname{Var}(x_j | \mathbf{x}_{\mathrm{pa}_j}) = \sigma_j^2$ for any $j \in \mathcal{A}_t$, due to the fact that $\mathrm{pa}_j \subset \mathcal{S}_t$ if $j \in \mathcal{A}_t$. Moreover, for any $j \in \mathcal{V} \backslash \{\mathcal{S}_t \cup \mathcal{A}_t\}$, by total variance, we have

$$E\big[\operatorname{Var}(x_j | \mathbf{x}_{\mathcal{S}_t})\big] = E\big[E[\operatorname{Var}(x_j | \mathbf{x}_{\mathrm{pa}_j}) | \mathbf{x}_{\mathcal{S}_t}]\big] + E\big[\operatorname{Var}\big(E[x_j | \mathbf{x}_{\mathrm{pa}_j}] | \mathbf{x}_{\mathcal{S}_t}\big)\big]$$
$$= \sigma_j^2 + E\big[\operatorname{Var}\big(E[x_j | \mathbf{x}_{\mathrm{pa}_j}] | \mathbf{x}_{\mathcal{S}_t}\big)\big].$$

This completes the first part of Theorem 1. Additionally, by Assumption 1 in the main text, for any $j, j' \in \mathcal{A}_t$, we have

$$E\big[\operatorname{Var}(x_j | \mathbf{x}_{\mathcal{S}_t})\big] = E\big[\operatorname{Var}(x_{j'} | \mathbf{x}_{\mathcal{S}_t})\big] := \sigma_{t,\min}^2,$$

and for any $k \in \mathcal{V} \backslash \{\mathcal{S}_t \cup \mathcal{A}_t\}$, we have

$$E\left[\operatorname{Var}(x_k|\mathbf{x}_{\mathcal{S}_t})\right] = \sigma_k^2 + E\left[\operatorname{Var}\left(E[x_k|\mathbf{x}_{\mathrm{pa}_k}]|\mathbf{x}_{\mathcal{S}_t}\right)\right] > \sigma_{t,\min}^2 + M_{\max}. \tag{1}$$

Clearly, all the nodes in $\mathcal{A}_t$ can be exactly identified by evaluating the expected conditional variance. This completes the proof. ∎

**Proof of Theorem 3.** Note that the sample variance estimator

$$\widehat{\operatorname{Var}}(x_k) = \frac{1}{n-1} \sum_{i=1}^{n} \left( x_{ik} - \frac{1}{n} \sum_{j=1}^{n} x_{jk} \right)^2 = \frac{1}{\binom{n}{2}} \sum_{i<j} \frac{1}{2}(x_{ik} - x_{jk})^2$$

is a U-statistics with kernel $\frac{1}{2}(x_{ik} - x_{jk})^2$. By the definition of $C_{\mathcal{X}}$ that denotes the diameter of the support $\mathcal{X}$, then we have $\frac{1}{2}(x_{ik} - x_{jk})^2 \leq \frac{1}{2}C_{\mathcal{X}}^2$. Then, by McDiarmid's inequality, for any $\zeta > 0$ and $k \in \mathcal{V}$, there holds

$$P\left(\left|\widehat{\operatorname{Var}}(x_k) - \operatorname{Var}(x_k)\right| > \zeta\right) \leq 2 \exp\left(-\frac{n\zeta^2}{2C_{\mathcal{X}}^4}\right). \tag{2}$$

Moreover, we define the following event

$$\mathcal{E}_0 = \left\{ \max_{k \in \mathcal{V}} \left|\widehat{\operatorname{Var}}(x_k) - \operatorname{Var}(x_k)\right| \leq \frac{M_{\max}}{4} \right\},$$

and use the notation $\mathcal{E}_0^c$ to denote its complementary. By (2), we have

$$P(\mathcal{E}_0^c) \leq 2p \exp\left(-\frac{nM_{\max}^2}{32C_{\mathcal{X}}^4}\right). \tag{3}$$

Note that

$$P(\mathcal{A}_0 \neq \widehat{\mathcal{A}}_0) \leq P(\mathcal{A}_0 \neq \widehat{\mathcal{A}}_0, \ \mathcal{E}_0) + P(\mathcal{E}_0^c)$$

$$\leq P\Big(\exists\, k \in \mathcal{A}_0 \text{ such that} \big|\widehat{\mathrm{Var}}(x_k) - \widehat{\sigma}_{\min}^{(0)}\big| \geq \epsilon_0, \ \mathcal{E}_0\Big)$$

$$+ P\Big(\exists\, k \in \mathcal{V}\backslash\{\mathcal{A}_0\} \text{ such that} \big|\widehat{\mathrm{Var}}(x_k) - \widehat{\sigma}_{\min}^{(0)}\big| < \epsilon_0, \ \mathcal{E}_0\Big) + P(\mathcal{E}_0^c)$$

$$= P_1 + P_2 + P(\mathcal{E}_0^c), \tag{4}$$

where $\widehat{\sigma}_{\min}^{(0)} = \min_{j \in \mathcal{V}} \widehat{\mathrm{Var}}(x_j)$. For ease notation, we denote $k_0 = \operatorname{argmin}_{k \in \mathcal{V}} \widehat{\mathrm{Var}}(x_k)$, and it always holds true that $k_0 \in \mathcal{A}_0$. If not, suppose that $k_0 \in \mathcal{V}\backslash\{\mathcal{A}_0\}$ and for any $j \in \mathcal{A}_0$, under the event $\mathcal{E}_0$ and by Theorem 1 in the main text, we have

$$\widehat{\mathrm{Var}}(x_{k_0}) > \mathrm{Var}(x_{k_0}) - \frac{M_{\max}}{2} > \mathrm{Var}(x_j) + \frac{M_{\max}}{2} > \widehat{\mathrm{Var}}(x_j),$$

which contradicts the definition that $k_0 = \operatorname{argmin}_{k \in \mathcal{V}} \widehat{\mathrm{Var}}(x_k)$.

To bound $P_1$, we notice that under the event $\mathcal{E}_0$, for any $j \in \mathcal{A}_0$, there holds

$$\big|\widehat{\mathrm{Var}}(x_j) - \widehat{\mathrm{Var}}(x_{k_0})\big| = \big|\widehat{\mathrm{Var}}(x_j) - \mathrm{Var}(x_j) + \mathrm{Var}(x_j) - \mathrm{Var}(x_{k_0}) + \mathrm{Var}(x_{k_0}) - \widehat{\mathrm{Var}}(x_{k_0})\big|$$

$$\leq \big|\widehat{\mathrm{Var}}(x_j) - \mathrm{Var}(x_j)\big| + \big|\mathrm{Var}(x_j) - \mathrm{Var}(x_{k_0})\big| + \big|\mathrm{Var}(x_{k_0}) - \widehat{\mathrm{Var}}(x_{k_0})\big|$$

$$\leq \frac{M_{\max}}{4} + 0 + \frac{M_{\max}}{4} = \frac{M_{\max}}{2},$$

where the last inequity follows from Assumption 1 in the main text and the definition of $\mathcal{E}_0$. Thus, by taking $\epsilon_0 = \frac{M_{\max}}{2}$, we have $P_1 = 0$.

Next, we turn to bound $P_2$. Note that for any $k \in \mathcal{V}\backslash\{\mathcal{A}_0\}$, by Theorem 1 in the main text, there holds

$$|\mathrm{Var}(x_k) - \mathrm{Var}(x_{k_0})| \geq M_{\max},$$

and triangle inequality yields that

$$\left|\operatorname{Var}(x_k) - \operatorname{Var}(x_{k_0})\right| \le \left|\operatorname{Var}(x_k) - \widehat{\operatorname{Var}}(x_k)\right| + \left|\widehat{\operatorname{Var}}(x_k) - \widehat{\operatorname{Var}}(x_{k_0})\right| + \left|\widehat{\operatorname{Var}}(x_{k_0}) - \operatorname{Var}(x_{k_0})\right|.$$

Then, under the event $\mathcal{E}_0$, we have

$$\left|\widehat{\operatorname{Var}}(x_k) - \widehat{\operatorname{Var}}(x_{k_0})\right| \ge M_{\max} - \left|\operatorname{Var}(x_k) - \widehat{\operatorname{Var}}(x_k)\right| - \left|\widehat{\operatorname{Var}}(x_{k_0}) - \operatorname{Var}(x_{k_0})\right|$$

$$\ge M_{\max} - \frac{M_{\max}}{4} - \frac{M_{\max}}{4} = \frac{M_{\max}}{2}.$$

Thus, by taking $\epsilon_0 = \frac{M_{\max}}{2}$, there holds $P_2 = 0$.

Clearly, we have

$$P(\mathcal{A}_0 \ne \widehat{\mathcal{A}}_0) \le P(\mathcal{E}_0^c) \le 2p \exp\left(-\frac{nM_{\max}^2}{32C_{\mathcal{X}}^4}\right), \tag{5}$$

by taking $\epsilon_0 = \frac{M_{\max}}{2}$. This completes the proof. ∎

**Lemma S1.** *Suppose that Assumptions 1– 3 in the main text are satisfied. Given the events $\{\widehat{\mathcal{A}}_0 = \mathcal{A}_0, ..., \widehat{\mathcal{A}}_{t-1} = \mathcal{A}_{t-1}\}$ and $\mathcal{J}$ and assume $n\lambda \to \infty$. Then, with probability at least $1 - \delta_n$, for any $j \in \mathcal{V}\backslash\{\mathcal{S}_t\}$, there holds*

$$\|\widehat{f}_j - f_{j,\mathcal{S}_t}^*\|_\infty \le \frac{\kappa_1 C_{j0}}{\lambda\sqrt{n}} \log \frac{2}{\delta_n} + \kappa_1 \lambda^{r-1/2} \|L_{K,t}^{-r} f_{j,\mathcal{S}_t}^*\|_2,$$

*where $C_{j0} = 2\kappa_1 \max\left\{C_{\mathcal{X}} + 2\kappa_1 R, \sqrt{2(2\kappa_1^2 R^2 + \sigma_j^2)}\right\}$.*

**Proof of Lemma S1.** To begin with, we define the sampling operator $S_{\mathbf{x}_{\mathcal{S}_t}} : \mathcal{H}_K \to \mathcal{R}^n$ associated with some copies of $\mathbf{x}_{\mathcal{S}_t} \in \mathcal{X}_t$ as

$$S_{\mathbf{x}_{\mathcal{S}_t}}(f) = \left(f(\mathbf{x}_{1\mathcal{S}_t}), ..., f(\mathbf{x}_{n\mathcal{S}_t})\right)^T,$$

and the adjoint of the sample operator as $S_{\mathbf{x}_{\mathcal{S}_t}}^T : \mathcal{R}^n \to \mathcal{H}_K$ as

$$S_{\mathbf{x}_{\mathcal{S}_t}}^T \mathbf{c} = \sum_{i=1}^{n} c_i K_{\mathbf{x}_{i\mathcal{S}_t}},$$

where $\mathbf{c} = (c_1, ..., c_n)^T \in \mathcal{R}^n$. Note that given the events $\{\widehat{\mathcal{A}}_0 = \mathcal{A}_0, ..., \widehat{\mathcal{A}}_{t-1} = \mathcal{A}_{t-1}\}$ and $\mathcal{J}$, we have $\mathcal{S}_t = \widehat{\mathcal{S}}_t$, and $\|\widehat{f}_j\|_K \leq R$, for any $j \in \mathcal{V}\backslash\{\widehat{\mathcal{S}}_t\}$. Clearly, the solution of (3.2) in Section 3.1 of the main text can be written as

$$\widehat{f}_j = \left(\frac{1}{n}S_{\mathbf{x}_{\mathcal{S}_t}}^T S_{\mathbf{x}_{\mathcal{S}_t}} + \lambda \mathbf{I}_n\right)^{-1} \frac{1}{n} S_{\mathbf{x}_{\mathcal{S}_t}}^T \mathbf{x}_j,$$

where $\mathbf{x}_j = (x_{1j}, ..., x_{nj})^T$ and $\mathbf{I}_n \in \mathcal{R}^{n\times n}$ denotes the identity matrix.

Moreover, we define an immediate function $f_{\lambda,j}$ as

$$f_{\lambda,j} = \underset{f_j \in \mathcal{H}_K}{\operatorname{argmin}} E\big[x_j - f_j(\mathbf{x}_{\mathcal{S}_t})\big]^2 + \lambda\|f_j\|_K^2. \tag{6}$$

Note that solving (6) equals solving the following problem that

$$f_{\lambda,j} = \underset{f_j \in \mathcal{H}_K}{\operatorname{argmin}} \|f_{j,\mathcal{S}_t}^* - f_j\|_{\mathcal{L}^2(\mathcal{X}_{\mathcal{S}_t},\rho_{\mathbf{x}_{\mathcal{S}_t}})}^2 + \lambda\|f_j\|_K^2, \tag{7}$$

by the fact that each node $x_j$ is centered with mean zero and $Ef(\mathbf{x}) = 0$ for all $f \in \mathcal{H}_K$. Thus, the solution of (6) can be derived as

$$f_{\lambda,j} = \left(L_{K,t} + \lambda I\right)^{-1} L_{K,t} f_{j,\mathcal{S}_t}^*,$$

where the integral operator $L_{K,t}$ is defined in Section 4 of the main text.

Simple algebra yields that

$$
\widehat{f}_j - f_{\lambda,j} = \left(\frac{1}{n}S_{\mathbf{x}_{\mathcal{S}_t}}^T S_{\mathbf{x}_{\mathcal{S}_t}} + \lambda \mathbf{I}_n\right)^{-1}\left(\frac{1}{n}S_{\mathbf{x}_{\mathcal{S}_t}}^T \mathbf{x}_j - \frac{1}{n}S_{\mathbf{x}_{\mathcal{S}_t}}^T S_{\mathbf{x}_{\mathcal{S}_t}} f_{\lambda,j} - \lambda f_{\lambda,j}\right)
$$

$$
= \left(\frac{1}{n}S_{\mathbf{x}_{\mathcal{S}_t}}^T S_{\mathbf{x}_{\mathcal{S}_t}} + \lambda \mathbf{I}_n\right)^{-1}\left(\frac{1}{n}\sum_{i=1}^n \left(x_{ij} - f_{\lambda,j}(\mathbf{x}_{i\mathcal{S}_t})\right)K_{\mathbf{x}_{i\mathcal{S}_t}} - L_{K,t}(f_{j,\mathcal{S}_t}^* - f_{\lambda,j})\right).
$$

Thus, we have

$$
\left\|\widehat{f}_j - f_{\lambda,j}\right\|_K \leq \frac{1}{\lambda}\left\|\frac{1}{n}\sum_{i=1}^n \left(x_{ij} - f_{\lambda,j}(\mathbf{x}_{i\mathcal{S}_t})\right)K_{\mathbf{x}_{i\mathcal{S}_t}} - L_{K,t}(f_{j,\mathcal{S}_t}^* - f_{\lambda,j})\right\|_K.
$$

For notation simplicity, we denote $\xi_i = \left(x_{ij} - f_{\lambda,j}(\mathbf{x}_{i\mathcal{S}_t})\right)K_{\mathbf{x}_{i\mathcal{S}_t}}$, which satisfies

$$
E[\xi_i] = L_{K,t}(f_{j,\mathcal{S}_t}^* - f_{\lambda,j}), \ \ \|\xi_i\|_K \leq \frac{\kappa_1}{2}\left(C_{\mathcal{X}} + 2\|f_{\lambda,j}\|_\infty\right)
$$

$$
\text{and } E[\|\xi_i\|_K^2] \leq \kappa_1^2 \int \left(x_j - f_{\lambda,j}(\mathbf{x}_{\mathcal{S}_t})\right)^2 d\rho_{\mathbf{x}_{\mathcal{S}_t\cup\{j\}}},
$$

and then by Lemma 2 of Smale and Zhou (2007), for any $\delta_n \in (0,1)$, with probability at least $1 - \delta_n$, we have

$$
\|\widehat{f}_j - f_{\lambda,j}\|_K \leq \frac{\kappa_1\left(C_{\mathcal{X}} + 2\|f_{\lambda,j}\|_\infty\right)\log(2/\delta_n)}{\lambda n} + \frac{\kappa_1}{\lambda}\sqrt{\frac{2\int \left(x_j - f_{\lambda,j}(\mathbf{x}_{\mathcal{S}_t})\right)^2 d\rho_{\mathbf{x}_{\mathcal{S}_t\cup\{j\}}}\log(2/\delta_n)}{n}}.
$$

It is clear that by pluging $f_j = 0$ into (7), we have $\|f_{\lambda,j}\|_K \leq \frac{\kappa_1\|f_{j,\mathcal{S}_t}^*\|_K}{\lambda^{1/2}}$, and thus we have $\|f_{\lambda,j}\|_\infty \leq \kappa_1\|f_{\lambda,j}\|_K < \frac{\kappa_1^2\|f_{j,\mathcal{S}_t}^*\|_K}{\lambda^{1/2}}$. To bound $\int \left(x_j - f_{\lambda,j}(\mathbf{x}_{\mathcal{S}_t})\right)^2 d\rho_{\mathbf{x}_{\mathcal{S}_t\cup\{j\}}}$, simple calculation yields that for any $f \in \mathcal{H}_K$,

$$
\int \left(x_j - f(\mathbf{x}_{\mathcal{S}_t})\right)^2 d\rho_{\mathbf{x}_{\mathcal{S}_t\cup\{j\}}} - \int \left(x_j - f_{j,\mathcal{S}_t}^*(\mathbf{x}_{\mathcal{S}_t})\right)^2 d\rho_{\mathbf{x}_{\mathcal{S}_t\cup\{j\}}} = \|f - f_{j,\mathcal{S}_t}^*\|_{\mathcal{L}^2(\mathcal{X}_{\mathcal{S}_t}, \rho_{\mathbf{x}_{\mathcal{S}_t}})}^2. \tag{8}
$$

By taking $f = 0$, there holds

$$\int \left(x_j - f_{j,\mathcal{S}_t}^*(\mathbf{x}_{\mathcal{S}_t})\right)^2 d\rho_{\mathbf{x}_{\mathcal{S}_t \cup \{j\}}} + \|0 - f_{j,\mathcal{S}_t}^*\|_{\mathcal{L}^2(\mathcal{X}_{\mathcal{S}_t}, \rho_{\mathbf{x}_{\mathcal{S}_t}})}^2 = E[f_j^*(\mathbf{x}_{\mathrm{pa}_j}) + n_j]^2 \leq \kappa_1^2 \|f_{j,\mathcal{S}_t}^*\|_K^2 + \sigma_j^2,$$

where the last equality follows from the generating scheme of Model 1 in the main text and the last inequality follows from Assumption 3 in the main text. Moreover, we notice that from (7) and by the definition of $f_{\lambda,j}$, there holds

$$\|f_{j,\mathcal{S}_t}^* - f_{\lambda,j}\|_{\mathcal{L}^2(\mathcal{X}_{\mathcal{S}_t}, \rho_{\mathbf{x}_{\mathcal{S}_t}})}^2 + \lambda \|f_{\lambda,j}\|_K^2 \leq \|f_{j,\mathcal{S}_t}^* - 0\|_{\mathcal{L}^2(\mathcal{X}_{\mathcal{S}_t}, \rho_{\mathbf{x}_{\mathcal{S}_t}})}^2 + \lambda \|0\|_K^2 \leq \kappa_1^2 \|f_{j,\mathcal{S}_t}^*\|_K^2,$$

and then, by plugging $f = f_{\lambda,j}$ into (8), we have

$$\int \left(x_j - f_{\lambda,j}(\mathbf{x}_{\mathcal{S}_t})\right)^2 d\rho_{\mathbf{x}_{\mathcal{S}_t \cup \{j\}}} = \int \left(x_j - f_{j,\mathcal{S}_t}^*(\mathbf{x}_{\mathcal{S}_t})\right)^2 d\rho_{\mathbf{x}_{\mathcal{S}_t \cup \{j\}}} + \|f_{\lambda,j} - f_{j,\mathcal{S}_t}^*\|_2^2 \leq 2\kappa_1^2 \|f_{j,\mathcal{S}_t}^*\|_K^2 + \sigma_j^2,$$

Therefore, with probability $1 - \delta_n$, we have

$$\|\widehat{f}_j - f_{\lambda,j}\|_K \leq \frac{\kappa_1(C_\mathcal{X} + 2\kappa_1^2\|f_{j,\mathcal{S}_t}^*\|_K/\lambda^{1/2})\log(2/\delta_n)}{\lambda n} + \frac{\kappa_1}{\lambda}\sqrt{\frac{2(2\kappa_1^2\|f_{j,\mathcal{S}_t}^*\|_K^2 + \sigma_j^2)\log(2/\delta_n)}{n}}$$
$$\leq \frac{C_{j0}\log(2/\delta_n)}{\lambda\sqrt{n}}, \tag{9}$$

where $C_{j0} = 2\kappa_1 \max\left\{C_\mathcal{X} + 2\kappa_1\|f_{j,\mathcal{S}_t}^*\|_K, \sqrt{2(2\kappa_1^2\|f_{j,\mathcal{S}_t}^*\|_K^2 + \sigma_j^2)}\right\}$ and the last inequality follows from the fact that $n\lambda \to \infty$.

Then, we turn to bound $\|f_{\lambda,j} - f_{j,\mathcal{S}_t}^*\|_K$ following similar treatments as in Smale and Zhou (2005). Specifically, for the integral operator $L_{K,t}$ defined in Section 4 of the main text with normalized eigenpairs $\{(\mu_k, \phi_k)\}_{k=1}^\infty$, we have

$$L_{K,t}^{1/2}\phi_i = \sum_{j\geq 1} \mu_j^{1/2}\langle\phi_i, \phi_j\rangle_2 \phi_j = \mu_i^{1/2}\phi_i \in \mathcal{H}_K,$$

and

$$\|\mu_i^{1/2}\phi_i\|_K = \Big(\sum_{j\geq 1} \frac{\langle \mu_i^{1/2}\phi_i, \phi_j\rangle_2^2}{\mu_j}\Big)^{1/2} = \langle\phi_i,\phi_i\rangle_2 = 1.$$

Thus by Assumption 2 of the main text, there exists some function $h_{j,t} = \sum_{i\geq 1}\langle h_{j,t}, \phi_i\rangle_2\phi_i \in \mathcal{L}^2(\mathcal{X}_{\mathcal{S}_t}, \rho_{\mathbf{x}_{\mathcal{S}_t}})$ such that $f_{j,\mathcal{S}_t}^* = L_{K,t}^r h_{j,t} = \sum_{i\geq 1}\mu_i^r\langle h_{j,t}, \phi_i\rangle_2\phi_i \in \mathcal{H}_K$.

Therefore, we have

$$f_{\lambda,j} - f_{j,\mathcal{S}_t}^* = (L_{K,t} + \lambda I)^{-1}\big(-\lambda f_{j,\mathcal{S}_t}^*\big) = -\sum_{i\geq 1}\frac{\lambda}{\lambda + \mu_i}\mu_i^r\langle h_{j,t}, \phi_i\rangle_2\phi_i,$$

and then, the $K$-norm of $f_{\lambda,j} - f_{j,\mathcal{S}_t}^*$ can be bounded as

$$
\begin{aligned}
\big\|f_{\lambda,j} - f_{j,\mathcal{S}_t}^*\big\|_K^2 &= \sum_{i\geq 1}\Big(\frac{\lambda}{\lambda + \mu_i}\mu_i^{r-1/2}\langle h_{j,t}, \phi_i\rangle_2\Big)^2 \\
&= \lambda^{2r-1}\sum_{i\geq 1}\Big(\frac{\lambda}{\lambda + \mu_i}\Big)^{3-2r}\Big(\frac{\mu_i}{\lambda + \mu_i}\Big)^{2r-1}\langle h_{j,t}, \phi_i\rangle_2^2 \\
&\leq \lambda^{2r-1}\sum_{i\geq 1}\langle h_{j,t}, \phi_i\rangle_2^2 = \lambda^{2r-1}\|h_{j,t}\|_2^2 = \lambda^{2r-1}\|L_{K,t}^{-r}f_{j,\mathcal{S}_t}^*\|_2^2. \qquad (10)
\end{aligned}
$$

Combining (9) and (10), under the events $\{\widehat{\mathcal{A}}_0 = \mathcal{A}_0, ..., \widehat{\mathcal{A}}_{t-1} = \mathcal{A}_{t-1}\}$ and $\mathcal{J}$, with probability at least $1 - \delta_n$, we have

$$
\begin{aligned}
\|\widehat{f}_j - f_{j,\mathcal{S}_t}^*\|_K &\leq \|\widehat{f}_j - f_{\lambda,j}\|_K + \|f_{\lambda,j} - f_{j,\mathcal{S}_t}^*\|_K \\
&\leq \frac{C_{j0}}{\lambda\sqrt{n}}\log\frac{2}{\delta_n} + \lambda^{r-1/2}\|L_{K,t}^{-r}f_{j,\mathcal{S}_t}^*\|_2.
\end{aligned}
$$

Moreover, we notice that $\|\widehat{f}_j - f_{j,\mathcal{S}_t}^*\|_\infty \leq \kappa_1\|\widehat{f}_j - f_{j,\mathcal{S}_t}^*\|_K$ by the reproducing property and the requirement that $\|f_{j,\mathcal{S}_t}^*\|_K^2 \leq R/2$ in Section 4 of the main text. This completes the proof. ∎

**Proof of Theorem 4.** Given the event $\{\widehat{\mathcal{A}}_0 = \mathcal{A}_0, ..., \widehat{\mathcal{A}}_{t-1} = \mathcal{A}_{t-1}\}$, we have $\mathcal{S}_t = \widehat{\mathcal{S}}_t$. Then, for

any $j \in \mathcal{V} \backslash \{\mathcal{S}_t\}$, there holds

$$
\left| \widehat{E\mathrm{Var}}(x_j|\mathbf{x}_{\mathcal{S}_t}) - E\,\mathrm{Var}(x_j|\mathbf{x}_{\mathcal{S}_t}) \right|
$$

$$
= \left| \frac{1}{n}\sum_{i=1}^{n}(x_{ij})^2 - \frac{1}{n}\sum_{i=1}^{n}\left(\widehat{f}_j(\mathbf{x}_{i\mathcal{S}_t})\right)^2 - E[x_j^2] + E\left[E[x_j|\mathbf{x}_{\mathcal{S}_t}]^2\right] \right|
$$

$$
\leq \left| E[x_j^2] - \frac{1}{n}\sum_{i=1}^{n}(x_{ij})^2 \right| + \left| E\left[E[x_j|\mathbf{x}_{\mathcal{S}_t}]^2\right] - \frac{1}{n}\sum_{i=1}^{n}\left(\widehat{f}_j(\mathbf{x}_{i\mathcal{S}_t})\right)^2 \right|. \tag{11}
$$

To bound the first term of (11), we notice that each $x_j$ is required to be centered with mean zero in Section 2 of the main text, which implies that zero belong to the support $\mathcal{X}$, and then $x_{ij}^2$ are bounded by $\frac{C_{\mathcal{X}}^2}{4}$ from the definition of $C_{\mathcal{X}}$, which denotes the diameter of the support $\mathcal{X}$. Then by the Hoeffding's inequality, for any $\frac{\varsigma}{2} > 0$, there holds

$$
P\left( \left| E[x_j^2] - \frac{1}{n}\sum_{i=1}^{n}(x_{ij})^2 \right| > \frac{\varsigma}{2} \right) \leq 2\exp\left( -\frac{8n\varsigma^2}{C_{\mathcal{X}}^4} \right). \tag{12}
$$

Next, the second term of (11) can be decomposed as

$$
\left| E\left[E[x_j|\mathbf{x}_{\mathcal{S}_t}]^2\right] - \frac{1}{n}\sum_{i=1}^{n}\left(\widehat{f}_j(\mathbf{x}_{i\mathcal{S}_t})\right)^2 \right|
$$

$$
\leq \left| E\left[E[x_j|\mathbf{x}_{\mathcal{S}_t}]^2 - \widehat{f}_j^2(\mathbf{x}_{\mathcal{S}_t})\right] \right| + \left| E\left[\widehat{f}_j^2(\mathbf{x}_{\mathcal{S}_t})\right] - \frac{1}{n}\sum_{i=1}^{n}\left(\widehat{f}_j(\mathbf{x}_{i\mathcal{S}_t})\right)^2 \right|
$$

$$
= \Delta_1 + \Delta_2,
$$

and thus it suffices to bound $\Delta_1$ and $\Delta_2$ sequentially under the events $\{\widehat{\mathcal{A}}_0 = \mathcal{A}_0, ..., \widehat{\mathcal{A}}_{t-1} = \mathcal{A}_{t-1}\}$ and $\mathcal{J}$.

To bound $\Delta_1$, we notice that

$$
\begin{aligned}
\Delta_1 &= \left| E\big[f^*_{j,\mathcal{S}_t}(\mathbf{x}_{\mathcal{S}_t})\big(f^*_{j,\mathcal{S}_t}(\mathbf{x}_{\mathcal{S}_t}) - \widehat{f}_j(\mathbf{x}_{\mathcal{S}_t})\big)\big] + E\big[\widehat{f}_j(\mathbf{x}_{\mathcal{S}_t})\big(f^*_{j,\mathcal{S}_t}(\mathbf{x}_{\mathcal{S}_t}) - \widehat{f}_j(\mathbf{x}_{\mathcal{S}_t})\big)\big] \right| \\
&\leq \|f^*_{j,\mathcal{S}_t} - \widehat{f}_j\|_\infty \left| \int |f^*_{j,\mathcal{S}_t}(\mathbf{x}_{\mathcal{S}_t})| d\rho_{\mathbf{x}_{\mathcal{S}_t}} + \int |\widehat{f}_j(\mathbf{x}_{\mathcal{S}_t})| d\rho_{\mathbf{x}_{\mathcal{S}_t}} \right| \\
&\leq 2\kappa_1 \max\{\|\widehat{f}_j\|_K, \|f^*_{j,\mathcal{S}_t}\|_K\} \|f^*_{j,\mathcal{S}_t} - \widehat{f}_j\|_\infty \leq 2\kappa_1 R \|f^*_{j,\mathcal{S}_t} - \widehat{f}_j\|_\infty,
\end{aligned}
$$

where the last inequality follows from the reproducing property of RKHS, the requirement that $\|f^*_{j,\mathcal{S}_t}\|_K \leq R/2$ in Section 4 of the main text and and under the event $\mathcal{J}$ that $\|\widehat{f}_j\|_K \leq R$. Then, by Lemma S1, with probability at least $1 - \delta_n/2$, we have

$$
\Delta_1 \leq 2\kappa_1^2 R \Big( \frac{C_{j0}}{\lambda\sqrt{n}} \log \frac{4}{\delta_n} + \lambda^{r-1/2} \|L_{K,t}^{-r} f^*_{j,\mathcal{S}_t}\|_2 \Big). \tag{13}
$$

To bound $\Delta_2$, we notice that

$$
\begin{aligned}
\Delta_2 &= \left| \int \widehat{f}_j(\mathbf{x}_{\mathcal{S}_t}) \langle \widehat{f}_j, K_{\mathbf{x}_{\mathcal{S}_t}} \rangle_K d\rho_{\mathbf{x}_{\mathcal{S}_t}} - \frac{1}{n} \sum_{i=1}^n \widehat{f}_j(\mathbf{x}_{i\mathcal{S}_t}) \langle \widehat{f}_j, K_{\mathbf{x}_{i\mathcal{S}_t}} \rangle_K \right| \\
&= \left| \Big\langle \widehat{f}_j, \int \widehat{f}_j(\mathbf{x}_{\mathcal{S}_t}) K_{\mathbf{x}_{\mathcal{S}_t}} d\rho_{\mathbf{x}_{\mathcal{S}_t}} \Big\rangle_K - \frac{1}{n} \Big\langle \widehat{f}_j, S^T_{\mathbf{x}_{\mathcal{S}_t}} S_{\mathbf{x}_{\mathcal{S}_t}} \widehat{f}_j \Big\rangle_K \right| \\
&= \left| \Big\langle \widehat{f}_j, \int \widehat{f}_j(\mathbf{x}_{\mathcal{S}_t}) K_{\mathbf{x}_{\mathcal{S}_t}} d\rho_{\mathbf{x}_{\mathcal{S}_t}} - \frac{1}{n} S^T_{\mathbf{x}_{\mathcal{S}_t}} S_{\mathbf{x}_{\mathcal{S}_t}} \widehat{f}_j \Big\rangle_K \right| \\
&= \left| \Big\langle \widehat{f}_j, \big(L_{K,t} - \frac{1}{n} S^T_{\mathbf{x}_{\mathcal{S}_t}} S_{\mathbf{x}_{\mathcal{S}_t}}\big) \widehat{f}_j \Big\rangle_K \right| \leq \|\widehat{f}_j\|_K \Big\| L_{K,t} - \frac{1}{n} S^T_{\mathbf{x}_{\mathcal{S}_t}} S_{\mathbf{x}_{\mathcal{S}_t}} \Big\|_{HS},
\end{aligned}
$$

where $S^T_{\mathbf{x}_{\mathcal{S}_t}}$ and $S_{\mathbf{x}_{\mathcal{S}_t}}$ denote the sampling operators defined in Lemma S1, and $\|\cdot\|_{HS}$ denotes the norm endowed with a Hilbert space $HS(K)$ containing all the Hilbert-Schmidt operators on $\mathcal{H}_K$ and satisfying $\|T\|_K \leq \|T\|_{HS}$ for any $T \in HS(K)$. Note that under the event $\mathcal{J}$, we have $\|\widehat{f}_j\|_K \leq R$. Moreover, by Lemma 18 of Rosasco et al. (2013), with probability at least $1 - \delta_n/2$,

we have

$$\left\| L_{K,t} - \frac{1}{n} S_{\mathbf{x}_{S_t}}^T S_{\mathbf{x}_{S_t}} \right\|_{HS} \le \frac{2\sqrt{2}\kappa_1^2}{\sqrt{n}} \log \frac{4}{\delta_n}.$$

Clearly, with probability at least $1 - \delta_n/2$, we have $\Delta_2 \le \frac{2R\sqrt{2}\kappa_1^2}{\sqrt{n}} \log \frac{4}{\delta_n}$.

Combining the upper bounds of $\Delta_1$ and $\Delta_2$, with probability at least $1 - \delta_n$, there holds

$$\left| E\left[ E[x_j|\mathbf{x}_{S_t}]^2 \right] - \frac{1}{n} \sum_{i=1}^{n} \left( \widehat{f}_j(\mathbf{x}_{iS_t}) \right)^2 \right|$$

$$\le 2\kappa_1^2 R\left( \frac{C_{j0}}{\lambda\sqrt{n}} \log \frac{4}{\delta_n} + \lambda^{r-1/2} \| L_{K,t}^{-r} f_{j,S_t}^* \|_2 \right) + \frac{2R\sqrt{2}\kappa_1^2}{\sqrt{n}} \log \frac{4}{\delta_n}$$

$$\le 2\kappa_1^2 R\left( \frac{C_{j0} + \sqrt{2}}{\lambda\sqrt{n}} \log \frac{4}{\delta_n} + \lambda^{r-1/2} \| L_{K,t}^{-r} f_{j,S_t}^* \|_2 \right).$$

Then, by taking $\lambda = n^{-\frac{1}{2r+1}}$, for any $\delta_n \in (0,1)$, with probability at least $1 - \delta_n$ there holds

$$\left| E\left[ E[x_j|\mathbf{x}_{S_t}]^2 \right] - \frac{1}{n} \sum_{i=1}^{n} \left( \widehat{f}_j(\mathbf{x}_{iS_t}) \right)^2 \right| \le C_{jt} n^{-\frac{2r-1}{2(2r+1)}} \log \frac{4}{\delta_n},$$

where $C_{jt} = 6\kappa_1^2 R \max \left\{ C_{j0}, \sqrt{2}, \| L_{K,t}^{-r} f_{j,S_t}^* \|_2 \right\}$. Correspondingly, for any $\zeta > 0$, we have

$$P\left( \left| E\left[ E[x_j|\mathbf{x}_{S_t}]^2 \right] - \frac{1}{n} \sum_{i=1}^{n} \left( \widehat{f}_j(\mathbf{x}_{iS_t}) \right)^2 \right| > \frac{\zeta}{2} \mid \widehat{\mathcal{A}}_0 = \mathcal{A}_0, ..., \widehat{\mathcal{A}}_{t-1} = \mathcal{A}_{t-1}, \mathcal{J} \right)$$

$$\le 4 \exp\left( -\frac{\zeta}{2C_{jt}} n^{\frac{2r-1}{2(2r+1)}} \right). \tag{14}$$

Combining (12) and (14), for any $\zeta > 0$, there holds

$$P\left( \left| E\mathrm{Var}(x_j|\mathbf{x}_{S_t}) - \widehat{E\mathrm{Var}}(x_j|\mathbf{x}_{S_t}) \right| > \zeta \mid \widehat{\mathcal{A}}_0 = \mathcal{A}_0, ..., \widehat{\mathcal{A}}_{t-1} = \mathcal{A}_{t-1}, \mathcal{J} \right)$$

$$\le 2 \exp\left( -\frac{8n\zeta^2}{C_{\mathcal{X}}^4} \right) + 4 \exp\left( -\frac{\zeta}{2C_{jt}} n^{\frac{2r-1}{2(2r+1)}} \right). \tag{15}$$

This completes the proof of the first part of Theorem 4.

11

Next, we define the following event

$$\mathcal{E}_{1t} = \Big\{ \max_{j \in \mathcal{V} \backslash \{\mathcal{S}_t\}} \big| E\mathrm{Var}(x_j|\mathbf{x}_{\mathcal{S}_t}) - \widehat{E\mathrm{Var}}(x_j|\mathbf{x}_{\mathcal{S}_t}) \big| \leq \frac{M_{\max}}{4} \Big\},$$

and use the notation $\mathcal{E}_{1t}^c$ to denote its complementary. Then, by (15), we have

$$P(\mathcal{E}_{1t}^c \mid \widehat{\mathcal{A}}_0 = \mathcal{A}_0, ..., \widehat{\mathcal{A}}_{t-1} = \mathcal{A}_{t-1}, \mathcal{J})$$

$$\leq 2(p - |\mathcal{S}_t|) \exp\Big( - \frac{n M_{\max}^2}{2C_{\mathcal{X}}^4} \Big) + 4(p - |\mathcal{S}_t|) \exp\Big( - \frac{M_{\max} n^{\frac{2r-1}{2(2r+1)}}}{8C_{jt}} \Big). \tag{16}$$

Note that

$$P(\widehat{\mathcal{A}}_t \neq \mathcal{A}_t \mid \widehat{\mathcal{A}}_0 = \mathcal{A}_0, ..., \widehat{\mathcal{A}}_{t-1} = \mathcal{A}_{t-1}, \mathcal{J})$$

$$\leq P(\widehat{\mathcal{A}}_t \neq \mathcal{A}_t, \ \mathcal{E}_{1t} \mid \widehat{\mathcal{A}}_0 = \mathcal{A}_0, ..., \widehat{\mathcal{A}}_{t-1} = \mathcal{A}_{t-1}, \mathcal{J}) + P(\mathcal{E}_{1t}^c \mid \widehat{\mathcal{A}}_0 = \mathcal{A}_0, ..., \widehat{\mathcal{A}}_{t-1} = \mathcal{A}_{t-1}, \mathcal{J})$$

$$\leq P\big( \exists \, j \in \mathcal{A}_t \text{ such that} |\widehat{E\mathrm{Var}}(x_j|\mathbf{x}_{\mathcal{S}_t}) - \widehat{\sigma}_{\min}^{(t)}| \geq \epsilon_t, \ \mathcal{E}_{1t} \mid \widehat{\mathcal{A}}_0 = \mathcal{A}_0, ..., \widehat{\mathcal{A}}_{t-1} = \mathcal{A}_{t-1}, \mathcal{J}\big)$$

$$\quad + P\big( \exists \, j \in \mathcal{V} \backslash \{\mathcal{S}_t \cup \mathcal{A}_t\} \text{ such that } |\widehat{E\mathrm{Var}}(x_j|\mathbf{x}_{\mathcal{S}_t}) - \widehat{\sigma}_{\min}^{(t)}| < \epsilon_t, \ \mathcal{E}_{1t} \mid \widehat{\mathcal{A}}_0 = \mathcal{A}_0, ..., \widehat{\mathcal{A}}_{t-1} = \mathcal{A}_{t-1}, \mathcal{J}\big)$$

$$\quad + P(\mathcal{E}_{1t}^c \mid \widehat{\mathcal{A}}_0 = \mathcal{A}_0, ..., \widehat{\mathcal{A}}_{t-1} = \mathcal{A}_{t-1}, \mathcal{J})$$

$$= P_3 + P_4 + P(\mathcal{E}_{1t}^c \mid \widehat{\mathcal{A}}_0 = \mathcal{A}_0, ..., \widehat{\mathcal{A}}_{t-1} = \mathcal{A}_{t-1}, \mathcal{J}). \tag{17}$$

Note that following the similar treatments as that of $P_1$ and $P_2$ in the proof of Theorem 3 in the main text and by taking $\epsilon_t = \frac{M_{\max}}{2}$, we have $P_3 = 0$ and $P_4 = 0$. Finally, the bound (17) reduces to

$$P(\widehat{\mathcal{A}}_t \neq \mathcal{A}_t \mid \widehat{\mathcal{A}}_0 = \mathcal{A}_0, ..., \widehat{\mathcal{A}}_{t-1} = \mathcal{A}_{t-1}, \mathcal{J})$$

$$\leq P(\mathcal{E}_{1t}^c \mid \widehat{\mathcal{A}}_0 = \mathcal{A}_0, ..., \widehat{\mathcal{A}}_{t-1} = \mathcal{A}_{t-1}, \mathcal{J})$$

$$\leq 2(p - |\mathcal{S}_t|) \exp\Big( - \frac{n M_{\max}^2}{2C_{\mathcal{X}}^4} \Big) + 4(p - |\mathcal{S}_t|) \exp\Big( - \frac{M_{\max} n^{\frac{2r-1}{2(2r+1)}}}{8C_{jt}} \Big).$$

This completes the proof. ∎

**Proof of Lemma 1.** Given the event $\{\widehat{\mathcal{A}}_0 = \mathcal{A}_0, ..., \widehat{\mathcal{A}}_t = \mathcal{A}_t\}$, $t \geq 1$, we have $\mathcal{S}_t = \widehat{\mathcal{S}}_t$. At first, for some positive constant $C_3$, we define the following event

$$\mathcal{E}_{2t} = \left\{ \max_{j \in \mathcal{A}_t, k \in \mathcal{S}_t} \left| \|\widehat{g}_{jk}\|_n^2 - \|g_{jk}^*\|_2^2 \right| \leq C_3 n^{-\frac{2r-1}{2(2r+1)}} \log \left( 4|\mathcal{S}_t| \max\{n, |\mathcal{S}_t|\} \right) \right\},$$

and use the notation $\mathcal{E}_{2t}^c$ to denote its complementary.

We notice that

$$P\left( \{\mathcal{E}_j \neq \widehat{\mathcal{E}}_j : j \in \widehat{\mathcal{A}}_t\} \mid \mathcal{A}_0 = \widehat{\mathcal{A}}_0, ..., \mathcal{A}_t = \widehat{\mathcal{A}}_t, \mathcal{J} \right)$$

$$\leq P\left( \{\mathcal{E}_j \neq \widehat{\mathcal{E}}_j : j \in \widehat{\mathcal{A}}_t\}, \mathcal{E}_{2t} \mid \mathcal{A}_0 = \widehat{\mathcal{A}}_0, ..., \mathcal{A}_t = \widehat{\mathcal{A}}_t, \mathcal{J} \right)$$

$$+ P\left( \mathcal{E}_{2t}^c \mid \mathcal{A}_0 = \widehat{\mathcal{A}}_0, ..., \mathcal{A}_t = \widehat{\mathcal{A}}_t, \mathcal{J} \right). \tag{18}$$

Note that by the definition that $\widehat{\mathcal{E}}_j = \{ k \to j, \|\widehat{g}_{jk}\|_n^2 > v_n^{(t)}, \text{for any } k \in \widehat{\mathcal{S}}_t \}$ and by Assumption 4 of the main text, for the first term of (18), there holds

$$P\left( \{\mathcal{E}_j \neq \widehat{\mathcal{E}}_j : j \in \widehat{\mathcal{A}}_t\}, \mathcal{E}_{2t} \mid \mathcal{A}_0 = \widehat{\mathcal{A}}_0, ..., \mathcal{A}_t = \widehat{\mathcal{A}}_t, \mathcal{J} \right)$$

$$\leq P\left( \exists\, k \in \mathrm{pa}_j \text{ such that } \|\widehat{g}_{jk}\|_n^2 \leq v_n^{(t)}, \mathcal{E}_{2t} \mid \mathcal{A}_0 = \widehat{\mathcal{A}}_0, ..., \mathcal{A}_t = \widehat{\mathcal{A}}_t, \mathcal{J} \right)$$

$$+ P\left( \exists\, k \in \widehat{\mathrm{pa}}_j \text{ such that } \|g_{jk}^*\|_2^2 = 0, \mathcal{E}_{2t} \mid \mathcal{A}_0 = \widehat{\mathcal{A}}_0, ..., \mathcal{A}_t = \widehat{\mathcal{A}}_t, \mathcal{J} \right)$$

$$= P_5 + P_6, \tag{19}$$

where $\widehat{\mathrm{pa}}_j = \{ k : k \to j \in \widehat{\mathcal{E}}_j \}$.

For the bound of $P_5$, by taking $v_n^{(t)} = \frac{C_2}{2} n^{-\frac{2r-1}{2(2r+1)}} \left( \log \left( 4|\mathcal{S}_t| \max\{n, |\mathcal{S}_t|\} \right) \right)^\beta$ and Assumption 3 in the main text, we have

$$\left| \|\widehat{g}_{jk}\|_n^2 - \|g_{jk}^*\|_2^2 \right| \geq \|g_{jk}^*\|_2^2 - \|\widehat{g}_{jk}\|_n^2 > 2v_n^{(t)} - v_n^{(t)} = v_n^{(t)},$$

which contradicts with $\mathcal{E}_{2t}$. Precisely, under $\mathcal{E}_{2t}$, we have

$$\max_{j\in\mathcal{A}_t, k\in\mathcal{S}_t}\left|\|\widehat{g}_{jk}\|_n^2 - \|g_{jk}^*\|_2^2\right| \leq C_3 n^{-\frac{2r-1}{2(2r+1)}}\log\left(4|\mathcal{S}_t|\max\{n, |\mathcal{S}_t|\}\right),$$

and then the different rates of convergence lead to the contradiction. Thus, when $n$ is sufficiently large, $P_5 = 0$. To bound $P_6$, it is obvious that $\left|\|\widehat{g}_{jk}\|_n^2 - \|g_{jk}^*\|_2^2\right| > \nu_n^{(t)}$, which contradicts with $\mathcal{E}_{2t}$ again, which yields that $P_6 = 0$.

Now, we turn to bound $P\left(\mathcal{E}_{2t}^c \mid \mathcal{A}_0 = \widehat{\mathcal{A}}_0, ..., \mathcal{A}_t = \widehat{\mathcal{A}}_t, \mathcal{J}\right)$. At first, we define the sample operators for gradients $\widehat{D}_{t,l} : \mathcal{H}_K \to \mathcal{R}^n$ as $(\widehat{D}_{t,l}f)_i = \langle f, \partial_l K_{\mathbf{x}_{i\mathcal{S}_t}}\rangle_K$, their adjoint operators $\widehat{D}_{t,l}^* : \mathcal{R}^n \to \mathcal{H}_K$ as $\widehat{D}_{t,l}^*\mathbf{c} = \frac{1}{n}\sum_{i=1}^n \partial_l K_{\mathbf{x}_{i\mathcal{S}_t}}c_i$, and define the integral operators for gradients $D_{t,l} : \mathcal{H}_K \to \mathcal{L}^2(\mathcal{X}_{\mathcal{S}_t}, \rho_{\mathbf{x}_{\mathcal{S}_t}})$ as $D_{t,l}f = \langle f, \partial_l K_{\mathbf{x}_{\mathcal{S}_t}}\rangle_K$, $D_{t,l}^* : \mathcal{L}^2(\mathcal{X}_{\mathcal{S}_t}, \rho_{\mathbf{x}_{\mathcal{S}_t}}) \to \mathcal{H}_K$ as $D_{t,l}^*f = \int \partial_l K_{\mathbf{x}_{\mathcal{S}_t}}f(\mathbf{x}_{\mathcal{S}_t})d\rho_{\mathbf{x}_{\mathcal{S}_t}}$. Then, we have

$$D_{t,l}^*D_{t,l}f_{j,\mathcal{S}_t}^* = \int \partial_l K_{\mathbf{x}_{\mathcal{S}_t}}g_{jl}^*(\mathbf{x}_{\mathcal{S}_t})d\rho_{\mathbf{x}_{\mathcal{S}_t}} \text{ and } \widehat{D}_{t,l}^*\widehat{D}_{t,l}f_{j,\mathcal{S}_t}^* = \frac{1}{n}\sum_{i=1}^n \partial_l K_{\mathbf{x}_{i\mathcal{S}_t}}g_{jl}^*(\mathbf{x}_{i\mathcal{S}_t}).$$

Note that $D_{t,l}^*D_{t,l}$ and $\widehat{D}_{t,l}^*\widehat{D}_{t,l}$ are the Hilbert-Schmidt operators belonging to a Hilbert space endowed with norm $\|\cdot\|_{HS}$.

Moreover, we notice that for any $j \in \mathcal{A}_t$ and $k \in \mathcal{S}_t$

$$\begin{aligned}
\left|\|\widehat{g}_{jk}\|_n^2 - \|g_{jk}^*\|_2^2\right| &= \left|\frac{1}{n}\sum_{i=1}^n\left(\widehat{g}_{jk}(\mathbf{x}_{i\mathcal{S}_t})\right)^2 - \int\left(g_{jk}^*(\mathbf{x}_{\mathcal{S}_t})\right)^2 d\rho_{\mathbf{x}_{\mathcal{S}_t}}\right| \\
&= \left|\left\langle\widehat{f}_j, \frac{1}{n}\sum_{i=1}^n\widehat{g}_{jk}(\mathbf{x}_{i\mathcal{S}_t})\partial_k K_{\mathbf{x}_{i\mathcal{S}_t}}\right\rangle_K - \left\langle f_{j,\mathcal{S}_t}^*, \int g_{jk}^*(\mathbf{x}_{\mathcal{S}_t})\partial_k K_{\mathbf{x}_{\mathcal{S}_t}}d\rho_{\mathbf{x}_{\mathcal{S}_t}}\right\rangle_K\right| \\
&= \left|\left\langle\widehat{f}_j - f_{j,\mathcal{S}_t}^*, \widehat{D}_{t,k}^*\widehat{D}_{t,k}(\widehat{f}_j - f_{j,\mathcal{S}_t}^*)\right\rangle_K + \langle\widehat{D}_{t,k}^*\widehat{D}_{t,k}f_{j,\mathcal{S}_t}^*, \widehat{f}_j - f_{j,\mathcal{S}_t}^*\rangle_K \right. \\
&\quad \left. + \langle f_{j,\mathcal{S}_t}^*, \widehat{D}_{t,k}^*\widehat{D}_{t,k}(\widehat{f}_j - f_{j,\mathcal{S}_t}^*)\rangle_K + \langle f_{j,\mathcal{S}_t}^*, (\widehat{D}_{t,k}^*\widehat{D}_{t,k} - D_{t,k}^*D_{t,k})f_{j,\mathcal{S}_t}^*\rangle_K\right| \\
&\leq \|\widehat{f}_j - f_{j,\mathcal{S}_t}^*\|_K^2\|\widehat{D}_{t,k}^*\widehat{D}_{t,k}\|_{HS} + 2\|\widehat{f}_j - f_{j,\mathcal{S}_t}^*\|_K\|f_{j,\mathcal{S}_t}^*\|_K\|\widehat{D}_{t,k}^*\widehat{D}_{t,k}\|_{HS} \\
&\quad + \|\widehat{D}_{t,k}^*\widehat{D}_{t,k} - D_{t,k}^*D_{t,k}\|_{HS}\|f_{j,\mathcal{S}_t}^*\|_K^2.
\end{aligned}$$

14

By Assumption 3 in the main text, direct calculation yields that

$$\max_{k \in \mathcal{S}_t} \|\widehat{D}_{t,k}^* \widehat{D}_{t,k}\|_{HS} = \max_{k \in \mathcal{S}_t} \|\partial_k K_{\mathbf{x}_{\mathcal{S}_t}}\|_K^2 \leq \kappa_2^2.$$

Moreover, by Lemma 18 of Rosasco et al. (2013), for any $\delta_n \in (0,1)$, with probability at least $1 - \delta_n/2$, we have

$$\max_{k \in \mathcal{S}_t} \|\widehat{D}_{t,k}^* \widehat{D}_{t,k} - D_{t,k}^* D_{t,k}\|_{HS} \leq \frac{2\sqrt{2}\kappa_2^2}{\sqrt{n}} \log \frac{4|\mathcal{S}_t|}{\delta_n},$$

given $\{\widehat{\mathcal{A}}_0 = \mathcal{A}_0, ..., \widehat{\mathcal{A}}_t = \mathcal{A}_t\}$. Thus, by taking $\delta_n = (\max\{n, |\mathcal{S}_t|\})^{-1}$, with probability at least $1 - \delta_n/2$, there holds

$$\max_{k \in \mathcal{S}_t} \|\widehat{D}_k^* \widehat{D}_k - D_k^* D_k\|_{HS} \leq \frac{2\sqrt{2}\kappa_2^2}{\sqrt{n}} \log \left(4|\mathcal{S}_t|\max\{n, |\mathcal{S}_t|\}\right).$$

When $\|\widehat{f}_j - f_{j,\mathcal{S}_t}^*\|_K$ is sufficiently small and by taking $\lambda = n^{-\frac{1}{2r+1}}$, with probability at least $1 - \delta_n$, we have

$$\max_{j \in \mathcal{A}_t, k \in \mathcal{S}_t} \left| \|\widehat{g}_{jk}\|_n^2 - \|g_{jk}\|_2^2 \right|$$

$$\leq \max\{\kappa_2^2, \kappa_2^2 \|f_{j,\mathcal{S}_t}^*\|_K, \|f_{j,\mathcal{S}_t}^*\|_K^2\} \left( 3 \max_{j \in \mathcal{A}_t, k \in \mathcal{S}_t} \|\widehat{f}_j - f_{j,\mathcal{S}_t}^*\|_K + \max_{k \in \mathcal{S}_t} \|\widehat{D}_k^* \widehat{D}_k - D_k^* D_k\|_{HS} \right)$$

$$\leq \max\{\kappa_2^2, \kappa_2^2 \|f_{j,\mathcal{S}_t}^*\|_K, \|f_{j,\mathcal{S}_t}^*\|_K^2\} \left( 3 \max_{j \in \mathcal{A}_t, k \in \mathcal{S}_t} \{C_{j0} + \|L_{K,t}^{-r} f_{j,\mathcal{S}_t}^*\|_2\} n^{-\frac{2r-1}{2(2r+1)}} \log \frac{4|\mathcal{S}_t|}{\delta_n} + \frac{2\sqrt{2}\kappa_2^2}{\sqrt{n}} \log \frac{4|\mathcal{S}_t|}{\delta_n} \right)$$

$$\leq C_3 n^{-\frac{2r-1}{2(2r+1)}} \log \left(4|\mathcal{S}_t|\max\{n, |\mathcal{S}_t|\}\right),$$

where $C_3 = 3 \max\{\kappa_2^2, \kappa_2^2 \|f_{j,\mathcal{S}_t}^*\|_K, \|f_{j,\mathcal{S}_t}^*\|_K^2\} \max_{j \in \mathcal{A}_t, k \in \mathcal{S}_t} \{3C_{j0}, 3\|L_{K,t}^{-r} f_{j,\mathcal{S}_t}^*\|_2, 2\sqrt{2}\kappa_2^2\}$. Thus,

we have

$$P\big(\mathcal{E}_{2t}^c|\widehat{\mathcal{A}}_0 = \mathcal{A}_0, ..., \widehat{\mathcal{A}}_t = \mathcal{A}_t, \mathcal{J}\big) \leq \frac{1}{\max\{n, |\mathcal{S}_t|\}},$$

Finally, by (18) and (19), we have

$$P\Big(\{\mathcal{E}_j = \widehat{\mathcal{E}}_j : j \in \widehat{\mathcal{A}}_t\} \mid \mathcal{A}_0 = \widehat{\mathcal{A}}_0, ..., \mathcal{A}_t = \widehat{\mathcal{A}}_t, \mathcal{J}\Big) \geq 1 - \frac{1}{\max\{n, |\mathcal{S}_t|\}}.$$

This completes the proof. ∎

**Proof for Theorem 5.** Note that

$$P(\widehat{\mathcal{G}} \neq \mathcal{G}) = P(\widehat{\mathcal{G}} \neq \mathcal{G}, \mathcal{J}) + P(\widehat{\mathcal{G}} \neq \mathcal{G}, \mathcal{J}^c).$$

For $P(\widehat{\mathcal{G}} \neq \mathcal{G}, \mathcal{J}^c)$, we have

$$P(\widehat{\mathcal{G}} \neq \mathcal{G}, \mathcal{J}^c) \leq P(\mathcal{J}^c) \leq Tp \max_{1 \leq t \leq T-1, \, j \in \mathcal{V}\backslash\{\mathcal{S}_t\}} P\big(\|\widehat{f}_j\|_K > R\big). \tag{20}$$

Note that by Theorem 1 and Lemma 3 of Smale and Zhou (2007), for any $t$ and $j \in \mathcal{V}\backslash\{\mathcal{S}_t\}$, we have $P\big(\|\widehat{f}_j\|_K > R\big) \leq \frac{1}{n}$ if the sample size satisfies $n \geq \left(\frac{C_4}{R}\log 2n\right)^{\frac{2(2r+1)}{2r-1}}$ for some positive constant $C_4$ and the $K$-norm of the target function is upper bounded by $R/2$ as assumed in Section 4 of the main text, and thus $P(\widehat{\mathcal{G}} \neq \mathcal{G}, \mathcal{J}^c) \to 0$ as $n \to \infty$.

For $P(\widehat{\mathcal{G}} \neq \mathcal{G}, \mathcal{J})$, we have

$$P(\widehat{\mathcal{G}} \neq \mathcal{G}, \mathcal{J}) \leq P\big( \cup_{t=0}^{T-1} \{\widehat{\mathcal{A}}_t \neq \mathcal{A}_t\} \cup \{\widehat{\mathcal{E}} \neq \mathcal{E}\}, \mathcal{J} \big)$$

$$\leq P(\widehat{\mathcal{A}}_0 \neq \mathcal{A}_0) + \sum_{t=1}^{T-1} P(\widehat{\mathcal{A}}_t \neq \mathcal{A}_t \mid \widehat{\mathcal{A}}_0 = \mathcal{A}_0, ..., \widehat{\mathcal{A}}_{t-1} = \mathcal{A}_{t-1}, \mathcal{J}) +$$

$$P(\widehat{\mathcal{E}} \neq \mathcal{E} \mid \widehat{\mathcal{A}}_0 = \mathcal{A}_0, ..., \widehat{\mathcal{A}}_{T-1} = \mathcal{A}_{T-1}, \mathcal{J})$$

$$\leq P(\widehat{\mathcal{A}}_0 \neq \mathcal{A}_0) + \sum_{t=1}^{T-1} P(\widehat{\mathcal{A}}_t \neq \mathcal{A}_t \mid \widehat{\mathcal{A}}_0 = \mathcal{A}_0, ..., \widehat{\mathcal{A}}_{t-1} = \mathcal{A}_{t-1}, \mathcal{J}) +$$

$$\sum_{t=1}^{T-1} P\Big( \{\mathcal{E}_j = \widehat{\mathcal{E}}_j : j \in \widehat{\mathcal{A}}_t\} \mid \mathcal{A}_0 = \widehat{\mathcal{A}}_0, ..., \mathcal{A}_t = \widehat{\mathcal{A}}_t, \mathcal{J} \Big).$$

Clearly, combining with Theorem 4 and Lemma 1 in the main text, we have $P(\widehat{\mathcal{G}} \neq \mathcal{G}, \mathcal{J}) \to 0$ as $n \to \infty$. This completes the proof. ∎

# References

Rosasco, L., S. Villa, S. Mosci, M. Santoro, and A. Verri (2013). Nonparametric sparsity and regularization. *Journal of Machine Learning Research* *14*, 1665–1714.

Smale, S. and D. Zhou (2005). Shannon sampling ii: Connections to learning theory. *Applied and Computational Harmonic Analysis* *19(3)*, 285–302.

Smale, S. and D. Zhou (2007). Learning theory estimates via integral operators and their approximations. *Constructive Approximation* *26(2)*, 153–172.