

---

# Appendix for "Bootstrap Model Aggregation for Distributed Statistical Learning"

---

**Jun Han**

Department of Computer Science  
Dartmouth College  
jun.han.gr@dartmouth.edu

**Qiang Liu**

Department of Computer Science  
Dartmouth College  
qiang.liu@dartmouth.edu

## 1 Appendix A

We study the asymptotic property of the KL-naive estimator  $\hat{\theta}_{\text{KL}}$ , and prove Theorem 2.

### 1.1 Notations and Assumptions

To simplify the notations for the proofs in the following, we define the following notations.

$$s(\mathbf{x}; \theta) = \log p(\mathbf{x} | \theta); \quad \dot{s}(\mathbf{x}; \theta) = \frac{\partial \log p(\mathbf{x} | \theta)}{\partial \theta}; \quad \ddot{s}(\mathbf{x}; \theta) = \frac{\partial^2 \log p(\mathbf{x} | \theta)}{\partial \theta^2}; \quad (1)$$

$$I(\theta) = \mathbb{E}(\ddot{s}(\mathbf{x}, \theta)); \quad I(\hat{\theta}_k, \theta_{\text{KL}}^*) = \mathbb{E}(\ddot{s}(\mathbf{x}, \theta_{\text{KL}}^*) | \hat{\theta}_k).$$

We start with investigating the theoretical property of  $\hat{\theta}_{\text{KL}}$ .

**Lemma 1.** *Based on Assumption 1, as  $n \rightarrow \infty$ , we have  $\mathbb{E}(\hat{\theta}_{\text{KL}} - \theta_{\text{KL}}^*) = o((dn)^{-1})$ . Further, in terms of estimating the true parameter, we have*

$$\mathbb{E}\|\hat{\theta}_{\text{KL}} - \theta^*\|^2 = O(N^{-1} + (dn)^{-1}). \quad (2)$$

**Proof:** Based on Equation (3) and (4), we know

$$\sum_{k=1}^d \frac{1}{n} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\theta}_{\text{KL}}) - \sum_{k=1}^d \int p(\mathbf{x} | \hat{\theta}_k) \dot{s}(\mathbf{x}; \theta_{\text{KL}}^*) d\mathbf{x} = 0. \quad (3)$$

By the law of large numbers, we can rewrite Equation (3) as

$$\sum_{k=1}^d \int p(\mathbf{x} | \hat{\theta}_k) \dot{s}(\mathbf{x}; \hat{\theta}_{\text{KL}}) d\mathbf{x} - \sum_{k=1}^d \int p(\mathbf{x} | \hat{\theta}_k) \dot{s}(\mathbf{x}; \theta_{\text{KL}}^*) d\mathbf{x} = o_p\left(\frac{1}{n}\right). \quad (4)$$

We also observe that  $\dot{s}(\mathbf{x}; \hat{\theta}_{\text{KL}}) - \dot{s}(\mathbf{x}; \theta_{\text{KL}}^*) = \left[ \int_0^1 \ddot{s}(\mathbf{x}; \theta_{\text{KL}}^* + t(\hat{\theta}_{\text{KL}} - \theta_{\text{KL}}^*)) dt \right] (\theta_{\text{KL}}^* - \hat{\theta}_{\text{KL}})$ . Therefore, Equation (4) can be written as

$$\left[ \sum_{k=1}^d \int p(\mathbf{x} | \hat{\theta}_k) \int_0^1 \ddot{s}(\mathbf{x}; \theta_{\text{KL}}^* + t(\hat{\theta}_{\text{KL}} - \theta_{\text{KL}}^*)) dt d\mathbf{x} \right] (\theta_{\text{KL}}^* - \hat{\theta}_{\text{KL}}) = o_p\left(\frac{1}{n}\right). \quad (5)$$

Under our Assumption 1, the Fish Information matrix  $I(\theta)$  is positive definite in a neighborhood of  $\theta^*$ , then we can find constant  $C_1, C_2$  such that  $C_1 \leq \left\| \int p(\mathbf{x} | \hat{\theta}_k) \int_0^1 \ddot{s}(\mathbf{x}; \theta_{\text{KL}}^* + t(\hat{\theta}_{\text{KL}} - \theta_{\text{KL}}^*)) dt d\mathbf{x} \right\| \leq C_2$ . Therefore, we can get  $\mathbb{E}(\hat{\theta}_{\text{KL}} - \theta_{\text{KL}}^*) = o((dn)^{-1})$ .  $\square$

The following theorem provides the MSE between  $\hat{\theta}_{\text{KL}}$  and  $\theta_{\text{KL}}^*$  and that between  $\hat{\theta}_{\text{KL}}$  and  $\theta^*$ .

**Theorem 2.** Based on Assumption 1, as  $n \rightarrow \infty$ ,  $\mathbb{E}\|\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*\|^2 = O(\frac{1}{nd})$ . Further, in terms of estimating the true parameter, we have

$$\mathbb{E}\|\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}^*\|^2 = O(N^{-1} + (dn)^{-1}). \quad (6)$$

**Proof:** According to the Equation (4),

$$\hat{\boldsymbol{\theta}}_{\text{KL}} = \arg \max_{\boldsymbol{\theta} \in \Theta} \sum_{k=1}^d \frac{1}{n} \sum_{j=1}^n s(\tilde{\mathbf{x}}_j^k; \boldsymbol{\theta}). \quad (7)$$

Then the first order derivative of Equation (7) with respect to  $\boldsymbol{\theta}$  at  $\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}_{\text{KL}}$  is zero,

$$\sum_{k=1}^d \frac{1}{n} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_{\text{KL}}) = 0. \quad (8)$$

By Taylor expansion of Equation (8), we get

$$\sum_{k=1}^d \frac{1}{n} \sum_{j=1}^n (\dot{s}(\tilde{\mathbf{x}}_j^k; \boldsymbol{\theta}_{\text{KL}}^*) + \ddot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_{\text{KL}})(\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*)) + o_p(\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*) = 0.$$

By the law of large numbers,  $\frac{1}{n} \sum_{j=1}^n \ddot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_{\text{KL}}) = I(\hat{\boldsymbol{\theta}}_k, \boldsymbol{\theta}_{\text{KL}}^*) + o_p(\frac{1}{n})$ . Under our Assumption 1,  $I(\boldsymbol{\theta})$  is positive definite in a neighborhood of  $\boldsymbol{\theta}^*$ . Since  $\hat{\boldsymbol{\theta}}_k$  are in the neighborhood of  $\boldsymbol{\theta}^*$ ,  $I(\hat{\boldsymbol{\theta}}_k, \boldsymbol{\theta}_{\text{KL}}^*)$  is positive definite, for  $k = 1 \in [d]$ . Then we have

$$\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^* = \left( \sum_{k=1}^d I(\hat{\boldsymbol{\theta}}_k, \boldsymbol{\theta}_{\text{KL}}^*) \right)^{-1} \sum_{k=1}^d \frac{1}{n} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^k; \boldsymbol{\theta}_{\text{KL}}^*) + o_p(\frac{1}{n}) = 0. \quad (9)$$

By the central limit theorem,  $\frac{1}{\sqrt{n}} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^k; \boldsymbol{\theta}_{\text{KL}}^*)$  converges to a normal distribution. By some simple calculation, we have

$$\text{Cov}(\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*, \hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*) = \frac{1}{n} \left( \sum_{k=1}^d I(\hat{\boldsymbol{\theta}}_k, \boldsymbol{\theta}_{\text{KL}}^*) \right)^{-1} \sum_{k=1}^d \text{Var}(\dot{s}(\mathbf{x}; \boldsymbol{\theta}_{\text{KL}}^*) \mid \hat{\boldsymbol{\theta}}_k) \left( \sum_{k=1}^d I(\hat{\boldsymbol{\theta}}_k, \boldsymbol{\theta}_{\text{KL}}^*) \right)^{-1}. \quad (10)$$

According to our Assumption 1, we already know  $I(\hat{\boldsymbol{\theta}}_k, \boldsymbol{\theta}_{\text{KL}}^*)$  is positive definite,  $C_1 \leq \|I(\hat{\boldsymbol{\theta}}_k, \boldsymbol{\theta}_{\text{KL}}^*)\| \leq C_2$ . We have  $(\sum_{k=1}^d I(\hat{\boldsymbol{\theta}}_k, \boldsymbol{\theta}_{\text{KL}}^*))^{-1} = O(\frac{1}{d})$  and  $\sum_{k=1}^d \text{Var}(\dot{s}(\mathbf{x}; \boldsymbol{\theta}_{\text{KL}}^*) \mid \hat{\boldsymbol{\theta}}_k) = O(d)$ . Therefore,  $\mathbb{E}\|\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*\|^2 = \text{trace}(\text{Cov}(\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*, \hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*)) = O(\frac{1}{nd})$ . Because the MSE between the exact KL estimator  $\boldsymbol{\theta}_{\text{KL}}^*$  and the true parameter  $\boldsymbol{\theta}^*$  is  $O(N^{-1})$  as shown in Liu and Ihler (2014), the MSE between  $\hat{\boldsymbol{\theta}}_{\text{KL}}$  and the true parameter  $\boldsymbol{\theta}^*$  is

$$\mathbb{E}\|\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}^*\|^2 \approx \mathbb{E}\|\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*\|^2 + \mathbb{E}\|\boldsymbol{\theta}_{\text{KL}}^* - \boldsymbol{\theta}^*\|^2 = O(N^{-1} + (dn)^{-1}).$$

We complete the proof of this theorem.  $\square$

## 2 Appendix B

In this section, we analyze the MSE of our proposed estimator  $\hat{\boldsymbol{\theta}}_{\text{KL}-C}$  and prove Theorem 3.

**Theorem 3.** Under Assumptions 1, we have

$$\text{as } n \rightarrow \infty, \quad n\mathbb{E}\|\hat{\boldsymbol{\theta}}_{\text{KL}-C} - \boldsymbol{\theta}_{\text{KL}}^*\|^2 < n\mathbb{E}\|\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*\|^2.$$

Since  $\tilde{\boldsymbol{\theta}}_k$  is the MLE of data  $\{\tilde{\mathbf{x}}_j^k\}_{j=1}^n$ , then we have

$$(\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k) = -I(\hat{\boldsymbol{\theta}}_k)^{-1} \frac{1}{n} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k) + o_p(\frac{1}{n}). \quad (11)$$

Then  $\mathbb{E}(\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k) = o(\frac{1}{n})$ . According to Theorem (2), when  $\mathfrak{B}_k$  is a constant matrix, for  $k \in [d]$ ,

$$\mathbb{E}(\hat{\boldsymbol{\theta}}_{\text{KL}-C} - \boldsymbol{\theta}_{\text{KL}}^*) = \mathbb{E}(\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*) + \sum_{k=1}^d \mathfrak{B}_k \mathbb{E}(\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k) = o(\frac{1}{n}).$$

Notice that  $\frac{1}{n} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^r; \hat{\boldsymbol{\theta}}_r)$  and  $\frac{1}{n} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^t; \hat{\boldsymbol{\theta}}_t)$  are independent when  $r \neq t$ . According to Equation (9), we know  $\sum_{k=1}^d \frac{1}{n} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^k; \boldsymbol{\theta}_{\text{KL}}^*)$  and  $\frac{1}{n} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k)$  are correlated to each other for  $k \in [d]$ ,

$$\begin{aligned} \text{Cov}((\hat{\boldsymbol{\theta}}_{\text{KL}-C} - \boldsymbol{\theta}_{\text{KL}}^*), (\hat{\boldsymbol{\theta}}_{\text{KL}-C} - \boldsymbol{\theta}_{\text{KL}}^*)) &= \text{Cov}(\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*, \hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*) \\ &+ 2 \sum_{k=1}^d \mathfrak{B}_k \text{Cov}(\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*, \tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k)^T + \sum_{k=1}^d \mathfrak{B}_k \text{Cov}((\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k), (\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k)) \mathfrak{B}_k^T. \end{aligned}$$

When  $\mathbf{B}_k = -(\text{Cov}(\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k))^{-1} \text{Cov}(\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*, \tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k)$ , we have

$$\begin{aligned} \text{Cov}(\hat{\boldsymbol{\theta}}_{\text{KL}-C} - \boldsymbol{\theta}_{\text{KL}}^*, \hat{\boldsymbol{\theta}}_{\text{KL}-C} - \boldsymbol{\theta}_{\text{KL}}^*) &= \text{Cov}(\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*, \hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*) - \\ &\sum_{k=1}^d \text{Cov}(\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k, \tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k)^{-1} \text{Cov}(\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*, \tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k) \text{Cov}(\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*, \tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k)^T. \end{aligned} \quad (12)$$

We know  $\mathbb{E}\|\hat{\boldsymbol{\theta}}_{\text{KL}-C} - \boldsymbol{\theta}_{\text{KL}}^*\|^2 = \text{trace}(\text{Cov}(\hat{\boldsymbol{\theta}}_{\text{KL}-C} - \boldsymbol{\theta}_{\text{KL}}^*, \hat{\boldsymbol{\theta}}_{\text{KL}-C} - \boldsymbol{\theta}_{\text{KL}}^*))$ ,  $\mathbb{E}\|\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*\|^2 = \text{trace}(\text{Cov}(\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*, \hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*))$ . The second term of Equation (12) is a positive definite matrix, therefore we have  $n\mathbb{E}\|\hat{\boldsymbol{\theta}}_{\text{KL}-C} - \boldsymbol{\theta}_{\text{KL}}^*\|^2 < n\mathbb{E}\|\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*\|^2$  as  $n \rightarrow \infty$ . We complete the proof of this theorem.  $\square$

**Theorem 4.** Under Assumption 1, when  $N > n \times d$ , we have  $E\|\hat{\boldsymbol{\theta}}_{\text{KL}-C} - \boldsymbol{\theta}_{\text{KL}}^*\|^2 = O(\frac{1}{dn^2})$  as  $n \rightarrow \infty$ . Further, in terms of estimating the true parameter, we have

$$\mathbb{E}\|\hat{\boldsymbol{\theta}}_{\text{KL}-C} - \boldsymbol{\theta}^*\|^2 = O(N^{-1} + (dn^2)^{-1}).$$

From Equation (4), we know

$$\sum_{k=1}^d \frac{1}{n} \sum_{j=1}^n \frac{\partial \log p(\tilde{\mathbf{x}}_j^k | \hat{\boldsymbol{\theta}}_{\text{KL}})}{\partial \boldsymbol{\theta}} = 0. \quad (13)$$

By Taylor expansion, Equation (13) can be rewritten as

$$\sum_{k=1}^d [\frac{1}{n} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k) + \ddot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k)(\hat{\boldsymbol{\theta}}_{\text{KL}} - \hat{\boldsymbol{\theta}}_k)] + O_p(\|\hat{\boldsymbol{\theta}}_{\text{KL}} - \hat{\boldsymbol{\theta}}_k\|^2) = 0. \quad (14)$$

$\|\hat{\boldsymbol{\theta}}_{\text{KL}} - \hat{\boldsymbol{\theta}}_k\|^2 \leq \|\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*\|^2 + \|\boldsymbol{\theta}_{\text{KL}}^* - \hat{\boldsymbol{\theta}}_k\|^2$ . As we know from Liu and Ihler (2014), we have

$$\|\boldsymbol{\theta}_{\text{KL}}^* - \hat{\boldsymbol{\theta}}_k\|^2 \leq \|\boldsymbol{\theta}_{\text{KL}}^* - \boldsymbol{\theta}^*\|^2 + \|\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}_k\|^2 = O_p(\frac{d}{N}), \quad (15)$$

When  $N > n \times d$ , we have  $\|\hat{\boldsymbol{\theta}}_{\text{KL}} - \hat{\boldsymbol{\theta}}_k\|^2 = O_p(\frac{1}{nd})$ . And it is also easy to derive

$$\hat{\boldsymbol{\theta}}_{\text{KL}} - \hat{\boldsymbol{\theta}}_k = \hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^* + \boldsymbol{\theta}_{\text{KL}}^* - \boldsymbol{\theta}^* + \boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}_k = o_p(\frac{1}{N}) + o_p(\frac{1}{N}) + o_p(\frac{d}{N}) = o_p(\frac{1}{nd} + \frac{d}{N}), \quad (16)$$

where  $\boldsymbol{\theta}_{\text{KL}}^* - \boldsymbol{\theta}^* = o_p(\frac{1}{N})$  has been proved in Liu and Ihler's paper(2014). According to the law of large numbers,  $\frac{1}{n} \sum_{j=1}^n \ddot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k) = I(\hat{\boldsymbol{\theta}}_k) + o_p(\frac{1}{n})$ , then we have

$$(\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*) = -(\sum_{k=1}^d I(\hat{\boldsymbol{\theta}}_k))^{-1} \sum_{k=1}^d \frac{1}{n} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k) + O_p(\frac{1}{nd}). \quad (17)$$

Notie that  $\frac{1}{n} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^r; \hat{\boldsymbol{\theta}}_r)$  and  $\frac{1}{n} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^t; \hat{\boldsymbol{\theta}}_t)$  are independent when  $r \neq t$ . Therefore from (11) and (17), the covariance matrix of  $n(\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*)$  and  $n(\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k)$  is

$$\text{Cov}(n(\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*), n(\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k)) = n\left(\sum_{k=1}^d I(\hat{\boldsymbol{\theta}}_k)\right)^{-1} + \left(\sum_{k=1}^d I(\hat{\boldsymbol{\theta}}_k)\right)^{-1} O(1),$$

for  $k \in [d]$ . According to Assumption 1, we know  $\sum_{k=1}^d I(\hat{\boldsymbol{\theta}}_k) = O(d)$ . Then we will have

$$\text{Cov}(n(\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*), n(\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k)) = n\left(\sum_{k=1}^d I(\hat{\boldsymbol{\theta}}_k)\right)^{-1} + O\left(\frac{1}{d}\right), \text{ for } k \in [d]. \quad (18)$$

According to Theorem 2 and Equation (10), by the law of large numbers, it is easy to derive

$$\text{Cov}(n(\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*), n(\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*)) = n\left(\sum_{k=1}^d I(\hat{\boldsymbol{\theta}}_k)\right)^{-1} + o(1).$$

$$\begin{aligned} \text{Cov}(n(\hat{\boldsymbol{\theta}}_{\text{KL}-C} - \boldsymbol{\theta}_{\text{KL}}^*), n(\hat{\boldsymbol{\theta}}_{\text{KL}-C} - \boldsymbol{\theta}_{\text{KL}}^*)) &= \text{Cov}(n(\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*), n(\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*)) \\ &+ 2 \sum_{k=1}^d \mathfrak{B}_k \text{Cov}(n(\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*), n(\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k))^\top + \sum_{k=1}^d \mathfrak{B}_k \text{Cov}(n(\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k), n(\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k)) \mathfrak{B}_k^\top, \end{aligned} \quad (19)$$

where  $\mathfrak{B}_k$  is defined in (8),

$$\mathfrak{B}_k = -\left(\sum_{k=1}^d I(\hat{\boldsymbol{\theta}}_k)\right)^{-1} I(\hat{\boldsymbol{\theta}}_k), \quad k \in [d].$$

According to Equation (11), we know  $\text{Cov}(n(\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k), n(\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k)) = n(I(\hat{\boldsymbol{\theta}}_k))^{-1} + o(1)$ . By some simple calculation, we know that  $n^2 \text{Cov}(\hat{\boldsymbol{\theta}}_{\text{KL}-C} - \boldsymbol{\theta}_{\text{KL}}^*, \hat{\boldsymbol{\theta}}_{\text{KL}-C} - \boldsymbol{\theta}_{\text{KL}}^*) = O\left(\frac{1}{d}\right)$ . Therefore, under the Assumption 1, when  $N > n \times d$ , we get the following result,

$$\mathbb{E}\|\hat{\boldsymbol{\theta}}_{\text{KL}-C} - \boldsymbol{\theta}_{\text{KL}}^*\|^2 = \text{trace}(\text{Cov}(\hat{\boldsymbol{\theta}}_{\text{KL}-C} - \boldsymbol{\theta}_{\text{KL}}^*, \hat{\boldsymbol{\theta}}_{\text{KL}-C} - \boldsymbol{\theta}_{\text{KL}}^*)) = O\left(\frac{1}{dn^2}\right).$$

We know  $\mathbb{E}\|\boldsymbol{\theta}_{\text{KL}}^* - \boldsymbol{\theta}^*\|^2 = O(N^{-1})$  from Liu and Ihler (2014). Then we have

$$\mathbb{E}\|\hat{\boldsymbol{\theta}}_{\text{KL}-C} - \boldsymbol{\theta}^*\|^2 \approx \mathbb{E}\|\hat{\boldsymbol{\theta}}_{\text{KL}-C} - \boldsymbol{\theta}_{\text{KL}}^*\|^2 + \mathbb{E}\|\boldsymbol{\theta}_{\text{KL}}^* - \boldsymbol{\theta}^*\|^2 = O(N^{-1} + (dn^2)^{-1}).$$

The proof of this theorem is complete.  $\square$

### 3 Appendix C

In this section, we analyze the asymptotic property of  $\hat{\boldsymbol{\theta}}_{\text{KL}-W}$  and prove Theorem 5. We show the MSE between  $\hat{\boldsymbol{\theta}}_{\text{KL}-W}$  and  $\boldsymbol{\theta}_{\text{KL}}^*$  is much smaller than the MSE between the KL-naive estimator  $\hat{\boldsymbol{\theta}}_{\text{KL}}$  and  $\boldsymbol{\theta}_{\text{KL}}^*$ .

**Lemma 5.** *Under Assumption 1, as  $n \rightarrow \infty$ ,  $\tilde{\eta}(\boldsymbol{\theta})$  is a more accurate estimator of  $\eta(\boldsymbol{\theta})$  than  $\hat{\eta}(\boldsymbol{\theta})$ , i.e.,*

$$n\text{Var}(\tilde{\eta}(\boldsymbol{\theta})) \leq n\text{Var}(\hat{\eta}(\boldsymbol{\theta})), \quad \text{for any } \boldsymbol{\theta} \in \Theta. \quad (20)$$

By Taylor expansion,

$$\frac{p(\mathbf{x}|\hat{\boldsymbol{\theta}}_k)}{p(\mathbf{x}|\tilde{\boldsymbol{\theta}}_k)} = 1 + (\log p(\mathbf{x}|\hat{\boldsymbol{\theta}}_k) - \log p(\mathbf{x}|\tilde{\boldsymbol{\theta}}_k)) + O_p(\|\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k\|^2), \quad (21)$$

we will have

$$\tilde{\eta}(\boldsymbol{\theta}) = \sum_{k=1}^d \left[ \frac{1}{n} \sum_{j=1}^n (1 + (s(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k) - s(\tilde{\mathbf{x}}_j^k; \tilde{\boldsymbol{\theta}}_k))) s(\tilde{\mathbf{x}}_j^k; \boldsymbol{\theta}) + O_p(\|\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k\|^2) \right],$$

Since  $s(\mathbf{x}; \hat{\boldsymbol{\theta}}_k) - s(\mathbf{x}; \tilde{\boldsymbol{\theta}}_k) = \dot{s}(\mathbf{x}; \hat{\boldsymbol{\theta}}_k)(\hat{\boldsymbol{\theta}}_k - \tilde{\boldsymbol{\theta}}_k)$ , according to equation (11), we have

$$\tilde{\eta}(\boldsymbol{\theta}) = \hat{\eta}(\boldsymbol{\theta}) - \sum_{k=1}^d \frac{1}{n} \sum_{j=1}^n s(\tilde{\mathbf{x}}_j^k; \boldsymbol{\theta}) \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k)(\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k) + O_p(\|\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k\|^2),$$

Then according to equation (11), we have

$$\hat{\eta}(\boldsymbol{\theta}) = \tilde{\eta}(\boldsymbol{\theta}) - \sum_{k=1}^d \mathbb{E}(s(\tilde{\mathbf{x}}_j^k; \boldsymbol{\theta}) \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k) \mid \hat{\boldsymbol{\theta}}_k) I(\hat{\boldsymbol{\theta}}_k)^{-1} \frac{1}{n} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k) + O_p\left(\frac{d}{n}\right),$$

Denote  $\hat{\xi}(\boldsymbol{\theta}) = -\sum_{k=1}^d \mathbb{E}(s(\tilde{\mathbf{x}}_j^k; \boldsymbol{\theta}) \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k) \mid \hat{\boldsymbol{\theta}}_k) I(\hat{\boldsymbol{\theta}}_k)^{-1} \frac{1}{n} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k)$ . According to Henmi et al. (2007),  $\hat{\xi}(\boldsymbol{\theta})$  is the orthogonal projection of  $\hat{\eta}(\boldsymbol{\theta})$  onto the linear space spanned by the score vector component for each  $\hat{\boldsymbol{\theta}}_k$ , where  $k \in [d]$ . Then we will have  $\text{Var}(\hat{\eta}(\boldsymbol{\theta})) = \text{Var}(\tilde{\eta}(\boldsymbol{\theta})) + \text{Var}(\hat{\xi}(\boldsymbol{\theta}))$ . Therefore,  $n\text{Var}(\tilde{\eta}(\boldsymbol{\theta})) \leq n\text{Var}(\hat{\eta}(\boldsymbol{\theta}))$ .

**Theorem 6.** Under the Assumption 1, for any  $\{\hat{\boldsymbol{\theta}}_k\}$ , we have that

$$\text{as } n \rightarrow \infty, \quad n\mathbb{E}\|\hat{\boldsymbol{\theta}}_{\text{KL-W}} - \boldsymbol{\theta}_{\text{KL}}^*\|^2 \leq n\mathbb{E}\|\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*\|^2.$$

**Proof:** From Equation (10), we know

$$\sum_{k=1}^d \frac{1}{n} \sum_{j=1}^n \frac{p(\tilde{\mathbf{x}}_j^k \mid \hat{\boldsymbol{\theta}}_k)}{p(\tilde{\mathbf{x}}_j^k \mid \tilde{\boldsymbol{\theta}}_k)} \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_{\text{KL-W}}) = 0.$$

Since  $\frac{p(\mathbf{x} \mid \hat{\boldsymbol{\theta}}_k)}{p(\mathbf{x} \mid \tilde{\boldsymbol{\theta}}_k)} = \exp\{\log p(\mathbf{x} \mid \hat{\boldsymbol{\theta}}_k) - \log p(\mathbf{x} \mid \tilde{\boldsymbol{\theta}}_k)\} = 1 + (\log p(\mathbf{x} \mid \hat{\boldsymbol{\theta}}_k) - \log p(\mathbf{x} \mid \tilde{\boldsymbol{\theta}}_k)) + O_p(\|\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k\|^2)$ , we have

$$\sum_{k=1}^d \frac{1}{n} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_{\text{KL-W}}) - \sum_{k=1}^d \left[ \frac{1}{n} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_{\text{KL-W}}) \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k)^T (\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k) + O_p(\|\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k\|^2) \right] = 0. \quad (22)$$

From the asymptotic property of MLE, we know  $\mathbb{E}\|\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k\|^2 = \frac{1}{n} \text{trace}(I(\hat{\boldsymbol{\theta}}_k))$ . Therefore, we know  $\|\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k\|^2 = O_p(\frac{1}{n})$  and  $\sum_{k=1}^d \|\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k\|^2 = O_p(\frac{d}{n})$ .

Similar to the derivation of equation (9), according to equation (11), we have the following equation,

$$\begin{aligned} \hat{\boldsymbol{\theta}}_{\text{KL-W}} - \boldsymbol{\theta}_{\text{KL}}^* &= \left( \sum_{k=1}^d I(\hat{\boldsymbol{\theta}}_k, \boldsymbol{\theta}_{\text{KL}}^*) \right)^{-1} \sum_{k=1}^d \frac{1}{n} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^k; \boldsymbol{\theta}_{\text{KL}}^*) - \\ &\quad \left( \sum_{k=1}^d I(\hat{\boldsymbol{\theta}}_k, \boldsymbol{\theta}_{\text{KL}}^*) \right)^{-1} \sum_{k=1}^d \mathbb{E}(\dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_{\text{KL-W}})^T \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k) \mid \hat{\boldsymbol{\theta}}_k) \frac{1}{n} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k) = O_p\left(\frac{d}{n}\right). \end{aligned}$$

Then we have,

$$\begin{aligned} \hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^* &= \hat{\boldsymbol{\theta}}_{\text{KL-W}} - \boldsymbol{\theta}_{\text{KL}}^* \\ &\quad - \left( \sum_{k=1}^d I(\hat{\boldsymbol{\theta}}_k, \boldsymbol{\theta}_{\text{KL}}^*) \right)^{-1} \sum_{k=1}^d \mathbb{E}(\dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_{\text{KL-W}})^T \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k) \mid \hat{\boldsymbol{\theta}}_k) \frac{1}{n} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k) = O_p\left(\frac{d}{n}\right). \end{aligned}$$

According to Henmi et al.(2007), we know the second term of above equation is the orthogonal projection of  $(\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*)$  onto the linear space spanned by the score component for each  $\hat{\boldsymbol{\theta}}_k$ , for  $k \in [d]$ . Then

$$n\mathbb{E}\|\hat{\boldsymbol{\theta}}_{\text{KL-W}} - \boldsymbol{\theta}_{\text{KL}}^*\|^2 \leq n\mathbb{E}\|\hat{\boldsymbol{\theta}}_{\text{KL}} - \boldsymbol{\theta}_{\text{KL}}^*\|^2.$$

We complete the proof of this theorem.  $\square$

**Theorem 7.** Under the Assumptions 1, when  $N > n \times d$ ,  $\mathbb{E}\|\hat{\boldsymbol{\theta}}_{\text{KL-W}} - \boldsymbol{\theta}_{\text{KL}}^*\|^2 = O(\frac{1}{dn^2})$ . Further, its MSE for estimating the true parameter  $\boldsymbol{\theta}^*$  is

$$\mathbb{E}\|\hat{\boldsymbol{\theta}}_{\text{KL-W}} - \boldsymbol{\theta}^*\|^2 = O(N^{-1} + (dn^2)^{-1}).$$

According to Equation (22),

$$\sum_{k=1}^d \frac{1}{n} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_{\text{KL-W}}) - \sum_{k=1}^d \frac{1}{n} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_{\text{KL-W}}) \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k)^T (\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k) = O_p\left(\frac{d}{n}\right).$$

Approximating the first term of the above equation by Taylor expansion, we will get

$$\begin{aligned} \sum_{k=1}^d \frac{1}{n} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_{\text{KL-W}}) &= \sum_{k=1}^d \left[ \frac{1}{n} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k) \right. \\ &\quad \left. + \sum_{k=1}^d \frac{1}{n} \sum_{j=1}^n \ddot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k) (\hat{\boldsymbol{\theta}}_{\text{KL-W}} - \hat{\boldsymbol{\theta}}_k) + O_p(\|\hat{\boldsymbol{\theta}}_{\text{KL-W}} - \hat{\boldsymbol{\theta}}_k\|^2) \right]. \end{aligned} \quad (23)$$

Since  $\|\hat{\boldsymbol{\theta}}_{\text{KL-W}} - \hat{\boldsymbol{\theta}}_k\|^2 \leq \|\hat{\boldsymbol{\theta}}_{\text{KL-W}} - \boldsymbol{\theta}_{\text{KL}}^*\|^2 + \|\boldsymbol{\theta}_{\text{KL}}^* - \hat{\boldsymbol{\theta}}_k\|^2$ , according to equation (15), then  $\|\hat{\boldsymbol{\theta}}_{\text{KL-W}} - \hat{\boldsymbol{\theta}}_k\|^2 = O_p(\|\hat{\boldsymbol{\theta}}_{\text{KL-W}} - \boldsymbol{\theta}_{\text{KL}}^*\|^2 + \frac{d}{N})$ . We can easily derive  $\dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_{\text{KL-W}}) = \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k) + O_p(\hat{\boldsymbol{\theta}}_{\text{KL-W}} - \hat{\boldsymbol{\theta}}_k)$  for  $k \in [d]$ . When  $N > n \times d$ , we will have

$$\begin{aligned} \sum_{k=1}^d \frac{1}{n} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k) &+ \sum_{k=1}^d \frac{1}{n} \sum_{j=1}^n \ddot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k) (\hat{\boldsymbol{\theta}}_{\text{KL-W}} - \hat{\boldsymbol{\theta}}_k) \\ &- \sum_k \frac{1}{n} \sum_{j=1}^n \dot{s}(\mathbf{x}_j^k; \hat{\boldsymbol{\theta}}_k) \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k)^T (\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k) + O_p(\|\hat{\boldsymbol{\theta}}_{\text{KL-W}} - \boldsymbol{\theta}_{\text{KL}}^*\|^2) = O\left(\frac{d}{n}\right). \end{aligned} \quad (24)$$

$\frac{1}{n} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k) = I(\hat{\boldsymbol{\theta}}_k) + o_p(\frac{1}{n})$  and we also know that  $\frac{1}{n} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k) \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k)^T = I(\hat{\boldsymbol{\theta}}_k) + o_p(1)$ . From (16), we know  $\boldsymbol{\theta}_{\text{KL}}^* - \hat{\boldsymbol{\theta}}_k = o_p(\frac{d}{N}) = o_p(\frac{1}{n})$ . When  $N > n \times d$ , we have

$$\begin{aligned} \sum_{k=1}^d \frac{1}{n} \sum_{j=1}^n \dot{s}(\tilde{\mathbf{x}}_j^k; \hat{\boldsymbol{\theta}}_k) &+ \sum_{k=1}^d I(\hat{\boldsymbol{\theta}}_k) (\hat{\boldsymbol{\theta}}_{\text{KL-W}} - \boldsymbol{\theta}_{\text{KL}}^*) \\ &+ \sum_{k=1}^d \frac{1}{n} I(\hat{\boldsymbol{\theta}}_k) (\tilde{\boldsymbol{\theta}}_k - \hat{\boldsymbol{\theta}}_k) + O_p(\|\hat{\boldsymbol{\theta}}_{\text{KL-W}} - \boldsymbol{\theta}_{\text{KL}}^*\|^2) = O\left(\frac{d}{n}\right). \end{aligned} \quad (25)$$

Based on the Equation (11), the first term and the third term of Equation (25) are cancelled. By some simple calculation, we will get

$$n^2 (\hat{\boldsymbol{\theta}}_{\text{KL-W}} - \boldsymbol{\theta}_{\text{KL}}^*)^T \left( \sum_{k=1}^d I(\hat{\boldsymbol{\theta}}_k) \right) \left( \sum_{k=1}^d I(\hat{\boldsymbol{\theta}}_k) \right) (\hat{\boldsymbol{\theta}}_{\text{KL-W}} - \boldsymbol{\theta}_{\text{KL}}^*) = O_p(d). \quad (26)$$

This indicates,  $\text{Cov}(n(\sum_{k=1}^d I(\hat{\boldsymbol{\theta}}_k))(\hat{\boldsymbol{\theta}}_{\text{KL-W}} - \boldsymbol{\theta}_{\text{KL}}^*), n(\sum_{k=1}^d I(\hat{\boldsymbol{\theta}}_k))(\hat{\boldsymbol{\theta}}_{\text{KL-W}} - \boldsymbol{\theta}_{\text{KL}}^*)) = O(d)$  as  $n \rightarrow \infty$ . We know  $n^2 \mathbb{E} \|\hat{\boldsymbol{\theta}}_{\text{KL-W}} - \boldsymbol{\theta}_{\text{KL}}^*\|^2 = \text{trace}(\text{Cov}(n(\hat{\boldsymbol{\theta}}_{\text{KL-W}} - \boldsymbol{\theta}_{\text{KL}}^*), n(\hat{\boldsymbol{\theta}}_{\text{KL-W}} - \boldsymbol{\theta}_{\text{KL}}^*)))$ . According to Assumption 1,  $I(\hat{\boldsymbol{\theta}}_k)$  is positive definite and then  $\text{trace}(\sum_{k=1}^d I(\hat{\boldsymbol{\theta}}_k)) = O(d)$ . Therefore, we have

$$\mathbb{E} \|\hat{\boldsymbol{\theta}}_{\text{KL-W}} - \boldsymbol{\theta}_{\text{KL}}^*\|^2 = O\left(\frac{d}{d^2 n^2}\right) = O\left(\frac{1}{dn^2}\right).$$

We know  $\mathbb{E} \|\boldsymbol{\theta}_{\text{KL}}^* - \boldsymbol{\theta}^*\|^2 = O(N^{-1})$  from Liu and Ihler (2014). Then we have

$$\mathbb{E} \|\hat{\boldsymbol{\theta}}_{\text{KL-W}} - \boldsymbol{\theta}^*\|^2 \approx \mathbb{E} \|\hat{\boldsymbol{\theta}}_{\text{KL-W}} - \boldsymbol{\theta}_{\text{KL}}^*\|^2 + \mathbb{E} \|\boldsymbol{\theta}_{\text{KL}}^* - \boldsymbol{\theta}^*\|^2 = O(N^{-1} + (dn^2)^{-1}).$$

The proof of this theorem is complete.  $\square$