

Fisher Efficient Inference of Intractable Models

Supplementary

A Examples

A.1 Examples of Stein Features $T_{\theta}f(x)$

Example 3. Let $p = \mathcal{N}(0, 1)$, $T_{\mathcal{N}(0,1)}1 = 0$, then $T_{\mathcal{N}(0,1)}x = -x$ and $T_{\mathcal{N}(0,1)}x^2/2 = -x^2 + 1$.

As we see, Stein features with respect to $\mathcal{N}(0, 1)$ using monomials of x are same-order polynomial terms of x which have been widely used as function basis in various function fitting applications.

A.2 Assumption 2 Examples

Example 4. When $f(x) = 0$, by the definition of Stein feature at Section 3.1, $T_{\theta}f(x) \equiv 0$. Our density ratio model does not have any discriminative power and become a constant function 1. We can see $H_{\delta,\delta} = 0$, $H_{\delta,\theta} = 0$ regardless what δ and θ are chosen. Thus, Assumption 2 is not satisfied here. See (14) and (15) in Section B.2 in Appendix for the exact formulas of $H_{\delta,\delta}$ and $H_{\delta,\theta}$.

Example 5. When $f(x) := x$ and $p(x; \theta) := \mathcal{N}(\theta, 1)$, our density ratio model becomes a linear discriminative function (See Example 3). From (14) and (15) we can see, when $\theta = \theta^*$ and $\delta = 0$, $H_{\delta,\delta}^* = -\frac{1}{n_q} \sum_{i=1}^{n_q} (x_q^{(i)} - \theta^*)^2$ which is essentially the negative sample variance and $H_{\delta,\theta}^* = \frac{1}{n_q} \sum_{i=1}^{n_q} \nabla_{\theta}(x_q^{(i)} - \theta) = -1$. Given n_q is sufficiently large, Λ_{\min} and Λ'_{\min} is reasonably small and Λ_{\max} is reasonably large, Assumption 2 should hold at the optimal point $(\theta^*, 0)$ with high probability. We omit the analysis when δ and θ are slightly deviated from their optimal values due to the page limit. Nonetheless, it can be analysed with some extra regularity conditions.

A.3 Example of Asymptotic Efficient Choice of $f(x)$

Example 6. Consider the univariate Gaussian distribution $p(x; \theta) = \exp\{\theta_1 x + \theta_2 x^2\} / Z(\theta)$ for $x \in \mathbb{R}$, $\theta = (\theta_1, \theta_2)$, where $\theta_1 \in \mathbb{R}$, $\theta_2 < 0$, and $Z(\theta)$ is the normalization constant. The score function is $s_1(x; \theta) = x - \frac{1}{Z(\theta)} \partial_{\theta_1} Z(\theta)$, $s_2(x; \theta) = x^2 - \frac{1}{Z(\theta)} \partial_{\theta_2} Z(\theta)$. Let us consider the Stein feature vector for $\mathbf{f}(x) = (x, x^2/2)^\top$, $T_{\theta}\mathbf{f}(x) = (\theta_1 + 2\theta_2 x, 1 + \theta_1 x + 2\theta_2 x^2)^\top$. We know that $\theta_1 Z(\theta) + 2\theta_2 \partial_{\theta_1} Z(\theta) = 0$ and $Z(\theta) + \theta_1 \partial_{\theta_1} Z(\theta) + 2\theta_2 \partial_{\theta_2} Z(\theta) = 0$ (see [11] for details). Thus, $\begin{pmatrix} T_{\theta}f_1(x) \\ T_{\theta}f_2(x) \end{pmatrix} = \begin{pmatrix} 2\theta_2 & 0 \\ \theta_1 & 2\theta_2 \end{pmatrix} \begin{pmatrix} s_1 \\ s_2 \end{pmatrix}$. The coefficient matrix is invertible as long as $\theta_2 \neq 0$. Hence, the DLE with the above \mathbf{f} achieves the asymptotic efficiency bound.

B Proofs

For simplicity, we write all $\sum_{i=1}^{n_q} g(\mathbf{x}_q^{(i)})$ as $\sum_{i=1}^{n_q} g(\mathbf{x}^{(i)})$ from now on as samples always come from dataset X_q . See Table 1 for all defined notations.

B.1 Proof of Lemma 1

Proof. Our proof below is similar to the proof of Lemma 4 in [13]. It can be seen that

$$\begin{aligned} \mathbb{E}_{p_{\theta}}[T_{\theta}f_i(\mathbf{x})] &= \int p(\mathbf{x}; \theta) [\langle \nabla_{\mathbf{x}} \log p(\mathbf{x}; \theta), \nabla_{\mathbf{x}} f_i(\mathbf{x}) \rangle + \text{trace}(\nabla_{\mathbf{x}}^2 f_i(\mathbf{x}))] d\mathbf{x} \\ &= \int \langle \nabla_{\mathbf{x}} p(\mathbf{x}; \theta), \nabla_{\mathbf{x}} f_i(\mathbf{x}) \rangle + p(\mathbf{x}; \theta) \cdot \text{trace}(\nabla_{\mathbf{x}}^2 f_i(\mathbf{x})) d\mathbf{x}. \end{aligned}$$

Table 1: Notations of Symbols

Symbol	Definition
$\ell(\delta, \theta)$	$\frac{1}{n_q} \sum_{i=1}^{n_q} \log r_{\theta}(\mathbf{x}_q^{(i)}; \delta)$, log likelihood ratio
$\nabla \ell(\delta_0, \theta_0)$	$\nabla_{(\delta, \theta)} \ell(\delta_0, \theta_0) _{\delta=\delta_0, \theta=\theta_0}$
$\nabla_{\delta} \ell(\delta_0, \theta_0)$	$\nabla_{\delta} \ell(\delta_0, \theta_0)$
\mathbf{H}	$\nabla_{(\delta, \theta)}^2 \ell(\delta, \theta)$, Hessian of likelihood
$\mathbf{H}_{\delta, \theta}$	$\nabla_{\delta} \nabla_{\theta} \ell(\delta, \theta)$, submatrix of Hessian.
$\text{Ball}(R, \mathbf{x}_0)$	ℓ_2 ball with radius R centered at \mathbf{x}_0
$\ A\ $	ℓ_2 norm of a vector A or the spectral norm of a matrix A
$\mathbf{s}(\mathbf{x}; \theta) \in \mathbb{R}^{\dim(\theta)}$	$\nabla_{\theta} \log p(\mathbf{x}, \theta)$, Score function of p_{θ}
\mathbf{s}	$\mathbf{s}(\mathbf{x}; \theta^*)$

Let us rewrite $\mathbb{E}_{p_{\theta}}[T_{\theta} f_i(\mathbf{x})]$ as nested integrals over each component of \mathbf{x} :

$$\begin{aligned} & \mathbb{E}_{p_{\theta}}[T_{\theta} f_i(\mathbf{x})] \\ &= \sum_{j=1}^d \int_{\mathbf{x}_{\setminus j}} \int_{x_j} \partial_{x_j} f_i(\mathbf{x}) \cdot \partial_{x_j} p(\mathbf{x}; \theta) + p(\mathbf{x}; \theta) \cdot \partial_{x_j}^2 f_i(\mathbf{x}) dx_j d\mathbf{x}_{\setminus j}, \end{aligned} \quad (11)$$

$$\begin{aligned} &= \sum_{j=1}^d \int_{\mathbf{x}_{\setminus j}} \underbrace{[p(\mathbf{x}; \theta) \partial_{x_j} f_i(\mathbf{x})]_{x_j \rightarrow -\infty}^{x_j \rightarrow +\infty}}_{0, \text{by assumption}} d\mathbf{x}_{\setminus j} - \int_{\mathbf{x}_{\setminus j}} \int_{x_j} p(\mathbf{x}; \theta) [\partial_{x_j}^2 f_i(\mathbf{x}) - \partial_{x_j}^2 f_i(\mathbf{x})] dx_j d\mathbf{x}_{\setminus j} \\ &= 0. \end{aligned} \quad (12)$$

$$= 0. \quad (13)$$

where $\mathbf{x}_{\setminus j}$ contains all the components in \mathbf{x} except the j -th component. The equality from (11) to (12) is due to one dimensional integration by parts formula. The first term in (12) is zero as the product of $p(\mathbf{x})$ and $\partial_{x_j} f_i(\mathbf{x})$ is assumed to be zero when x_j takes the limit to $+/ - \infty$. Our assumption holds for all i, j , so we can assert $\forall_i \mathbb{E}_{p_{\theta}}[T_{\theta} f_i(\mathbf{x})] = 0$ and $\mathbb{E}_{p_{\theta}}[T_{\theta} \mathbf{f}(\mathbf{x})] = \mathbf{0}$ by its construction. \square

B.2 Derivations of $\nabla_{\delta}^2 \ell(\delta, \theta)$ and $\nabla_{\delta, \theta} \ell(\delta, \theta)$ with $f(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}$

$$\nabla_{\delta}^2 \ell(\delta, \theta) = -\frac{1}{n_q} \sum_{i=1}^{n_q} \frac{[T_{\theta} f(\mathbf{x}^{(i)})]^2}{r_{\theta}^2(\mathbf{x}^{(i)}; \delta)} + 0, \quad (14)$$

$$\nabla_{\delta, \theta} \ell(\delta, \theta) = -\frac{1}{n_q} \sum_{i=1}^{n_q} \frac{T_{\theta} f(\mathbf{x}^{(i)})}{r_{\theta}^2(\mathbf{x}^{(i)}; \delta)} \nabla_{\theta} r_{\theta}(\mathbf{x}^{(i)}; \delta) + \frac{1}{n_q} \sum_{i=1}^{n_q} \frac{1}{r_{\theta}(\mathbf{x}^{(i)}; \delta)} \nabla_{\theta} T_{\theta} f(\mathbf{x}^{(i)}). \quad (15)$$

B.3 Proof of Proposition 2

Proof. First, the definition of Δ_{n_q} gives the boundedness of our ratio, i.e., $\frac{1}{C_{\text{ratio}}} \leq r_{\theta}(\mathbf{x}; \delta) \leq C_{\text{ratio}}, \forall \mathbf{x} \in X_q, \forall \theta \in \Theta$.

Second, $-\mathbf{H}_{\delta, \delta} = \frac{1}{n_q} \sum_{i=1}^{n_q} \frac{1}{r_{\theta}^2(\mathbf{x}^{(i)}; \delta)} \cdot T_{\theta} \psi^{(i)} T_{\theta} \psi^{(i)\top}$, where $T_{\theta} \psi^{(i)}$ is an abbreviation of $T_{\theta} \psi(\mathbf{x}^{(i)})$.

It is a sum over ratio weighted positive semi-definite matrices so we can lower bound its minimum eigenvalue using the lower bound of the ratio:

$$\lambda_{\min}(-\mathbf{H}_{\delta, \delta}) \geq \frac{1}{C_{\text{ratio}}^2} \lambda_{\min} \left(\frac{1}{n_q} \sum_{i=1}^{n_q} T_{\theta} \psi^{(i)} T_{\theta} \psi^{(i)\top} \right) > \frac{\Lambda''_{\min}}{C_{\text{ratio}}^2} > 0, \text{ with high prob.,}$$

due to our assumption. Similarly, we can also upper-bound its maximum eigenvalue

$$\lambda_{\max}(-\mathbf{H}_{\delta, \delta}) \leq C_{\text{ratio}}^2 \lambda_{\max} \left(\frac{1}{n_q} \sum_{i=1}^{n_q} T_{\theta} \psi^{(i)} T_{\theta} \psi^{(i)\top} \right) \leq C_{\text{ratio}}^2 \Lambda''_{\max}, \text{ with high prob.,}$$

Third, $-\mathbf{H}_{\theta,\theta} = \frac{1}{n_q} \sum_{i=1}^{n_q} \frac{1}{r_{\theta}^2(\mathbf{x}^{(i)}; \delta)} \mathbf{J}_{\mathbf{x}} \psi(\mathbf{x}^{(i)}) \mathbf{J}_{\mathbf{x}} \psi(\mathbf{x}^{(i)})^{\top} \delta \delta^{\top} \mathbf{J}_{\mathbf{x}} \psi(\mathbf{x}^{(i)})^{\top} \mathbf{J}_{\mathbf{x}} \psi(\mathbf{x}^{(i)})$. We can see

$$\|\mathbf{H}_{\theta,\theta}\| \leq \frac{C_{\text{ratio}}^2 \cdot \|\delta\|^2}{n_q} \sum_{i=1}^{n_q} \|\mathbf{J}_{\mathbf{x}} \psi(\mathbf{x}^{(i)})\|^4 \leq C_{\text{ratio}}^2 C_2 \cdot \|\delta\|^2 \leq \frac{C_{\text{ratio}}^2 C_2 T}{\sigma(n_q)^2}.$$

Fourth, using the fact that $-\mathbf{H}_{\delta,\delta}$ is a positive definite matrix, which we have just proved, we can see

$$\begin{aligned} \lambda_{\min} \left\{ -\mathbf{H}_{\theta,\delta} \mathbf{H}_{\delta,\delta}^{-1} \mathbf{H}_{\delta,\theta} \right\} &= \lambda_{\min} (-\mathbf{H}_{\delta,\delta}^{-1} \mathbf{H}_{\delta,\theta} \mathbf{H}_{\theta,\delta}) \\ &\geq \lambda_{\min} (-\mathbf{H}_{\delta,\delta}^{-1}) \lambda_{\min} (\mathbf{H}_{\delta,\theta} \mathbf{H}_{\theta,\delta}) \\ &= \frac{\lambda_{\min} (\mathbf{H}_{\delta,\theta} \mathbf{H}_{\theta,\delta})}{\lambda_{\max} (-\mathbf{H}_{\delta,\delta})} \geq \frac{\lambda_{\min} (\mathbf{H}_{\delta,\theta} \mathbf{H}_{\theta,\delta})}{C_{\text{ratio}}^2 \Lambda''_{\max}}, \end{aligned}$$

where 2nd line is due to Theorem 7, [21]. So we only need to find a lower bound for $\lambda_{\min} (\mathbf{H}_{\delta,\theta} \mathbf{H}_{\theta,\delta})$. We can write $\mathbf{H}_{\theta,\delta}$ as

$$\mathbf{H}_{\theta,\delta} = \frac{1}{n_q} \sum_{i=1}^{n_q} \frac{1}{r_{\theta}(\mathbf{x}^{(i)}; \delta)} \mathbf{J}_{\mathbf{x}} \psi(\mathbf{x}^{(i)}) \mathbf{J}_{\mathbf{x}} \psi(\mathbf{x}^{(i)})^{\top} \quad (16)$$

$\underbrace{\hspace{10em}}_{\mathbf{A}}$

$$- \frac{1}{n_q} \sum_{i=1}^{n_q} \frac{1}{r_{\theta}^2(\mathbf{x}^{(i)}; \delta)} \mathbf{J}_{\mathbf{x}} \psi(\mathbf{x}^{(i)}) \mathbf{J}_{\mathbf{x}} \psi(\mathbf{x}^{(i)})^{\top} \delta T_{\theta} \psi(\mathbf{x}^{(i)})^{\top} \quad (17)$$

$\underbrace{\hspace{10em}}_{\mathbf{B}}$

Therefore $\mathbf{H}_{\delta,\theta} \mathbf{H}_{\theta,\delta}$ can be written as

$$\mathbf{A} \mathbf{A}^{\top} - \mathbf{A} \mathbf{B}^{\top} - \mathbf{B} \mathbf{A}^{\top} + \mathbf{B} \mathbf{B}^{\top}.$$

Since we are analyzing the minimum eigenvalue, we can safely ignore the last term $\mathbf{B} \mathbf{B}^{\top}$ as it is positive semi-definite. This gives the following inequality:

$$\begin{aligned} \lambda_{\min} \left\{ \mathbf{A} \mathbf{A}^{\top} - \mathbf{A} \mathbf{B}^{\top} - \mathbf{B} \mathbf{A}^{\top} \right\} &\geq \lambda_{\min} \left\{ \mathbf{A} \mathbf{A}^{\top} \right\} + \lambda_{\min} \left\{ -\mathbf{A} \mathbf{B}^{\top} - \mathbf{B} \mathbf{A}^{\top} \right\} \\ &\geq \lambda_{\min} \left\{ \mathbf{A} \mathbf{A}^{\top} \right\} - \|\mathbf{A} \mathbf{B}^{\top} + \mathbf{B} \mathbf{A}^{\top}\| \end{aligned}$$

As \mathbf{A} is a sum of ratio weighted positive semi-definite matrices, we can use the same trick in the second step to lower bound its eigenvalue using the lower bound of the density ratio, eventually, using our assumption on $\lambda_{\min} \left\{ \frac{1}{n_q} \sum_{i=1}^{n_q} \mathbf{J}^{(i)} \mathbf{J}^{(i)\top} \right\} \geq C_3$, we can get,

$$\lambda_{\min}(\mathbf{A}) \geq \frac{C_3}{C_{\text{ratio}}}, \lambda_{\min}(\mathbf{A} \mathbf{A}^{\top}) \geq \lambda_{\min}(\mathbf{A}) \cdot \lambda_{\min}(\mathbf{A}) \geq \frac{C_3^2}{C_{\text{ratio}}^2}.$$

Now we analyze $\|\mathbf{A} \mathbf{B}^{\top} + \mathbf{B} \mathbf{A}^{\top}\|$ which is further upperbounded by $2\|\mathbf{A}\|\|\mathbf{B}\|$.

Similarly to how $\lambda_{\min}(\mathbf{A})$ is bounded, we can upper-bound $\|\mathbf{A}\|$ using the upperbound of the ratio: $\|\mathbf{A}\| \leq C_{\text{ratio}} C_4$. Let us write $\mathbf{B} = \frac{1}{n_q} \sum_{i=1}^{n_q} \frac{1}{r_i^2} \mathbf{J}^{(i)} \mathbf{J}^{(i)\top} \delta T \psi^{(i)\top}$ where $\mathbf{J}^{(i)}$ and r_i are abbreviations of $\mathbf{J}_{\mathbf{x}} \psi(\mathbf{x}^{(i)})$ and $r_{\theta}(\mathbf{x}^{(i)}; \delta)$. It can be seen that

$$\begin{aligned} \|\mathbf{B}\| &\leq \frac{1}{n_q} \sum_{i=1}^{n_q} \left\| \frac{1}{r_i^2} \mathbf{J}^{(i)} \mathbf{J}^{(i)\top} \right\| \cdot \|\delta T \psi^{(i)}\| \leq \frac{1}{n_q} \sum_{i=1}^{n_q} \left\| \frac{1}{r_i^2} \mathbf{J}^{(i)} \mathbf{J}^{(i)\top} \right\| \cdot \|\delta\| \cdot \|T \psi^{(i)}\|, \\ &\leq C_{\text{ratio}}^2 \cdot C_5 T / \sigma(n_q). \end{aligned}$$

Now we can bound

$$\begin{aligned} \lambda_{\min} \left\{ \mathbf{A} \mathbf{A}^{\top} - \mathbf{A} \mathbf{B}^{\top} - \mathbf{B} \mathbf{A}^{\top} + \mathbf{B} \mathbf{B}^{\top} \right\} &\geq \lambda_{\min} \left\{ \mathbf{A} \mathbf{A}^{\top} \right\} - 2\|\mathbf{A}\|\|\mathbf{B}\| \\ &\geq \frac{C_3^2}{C_{\text{ratio}}^2} - C_{\text{ratio}}^3 C_4 \cdot C_5 T / \sigma(n_q) \end{aligned}$$

There exists a constant $N > 0$, such that when $n_q > N$,

$$\begin{aligned}\lambda_{\min} \left\{ -\mathbf{H}_{\theta,\delta} \mathbf{H}_{\delta,\delta}^{-1} \mathbf{H}_{\delta,\theta} \right\} &\geq \frac{\lambda_{\min}(\mathbf{H}_{\delta,\theta} \mathbf{H}_{\theta,\delta})}{C_{\text{ratio}}^2 \Lambda''_{\max}} \geq \frac{C_3^2}{C_{\text{ratio}}^4 \Lambda'_{\max}} - \frac{C_{\text{ratio}} C_4 \cdot C_5 T}{\sigma(n_q) \Lambda''_{\max}} \\ &\geq \frac{C_{\text{ratio}}^2 C_2 T}{\sigma(n_q)^2} \geq \|\mathbf{H}_{\theta,\theta}\|.\end{aligned}$$

Finally we analyze $\|\mathbf{H}_{\theta,\delta} \mathbf{H}_{\delta,\delta}^{-1}\|$. $\|\mathbf{H}_{\theta,\delta} \mathbf{H}_{\delta,\delta}^{-1}\| \leq \|\mathbf{H}_{\theta,\delta}\| \cdot \|\mathbf{H}_{\delta,\delta}^{-1}\|$. As $-\mathbf{H}_{\delta,\delta}$ is positive definite, the operator norm of its inverse is the inverse of its minimum eigenvalue, which is upperbounded by $C_{\text{ratio}}^2 / \Lambda''_{\min}$. On the other hand, we can rewrite (16) as $\mathbf{H}_{\theta,\delta} = \frac{1}{n_q} \sum_{i=1}^{n_q} \frac{1}{r_i} \cdot \mathbf{J}^{(i)} \mathbf{J}^{(i)\top} \cdot \left(\text{Iden} - \frac{1}{r_i} \cdot \delta T \psi^{(i)\top} \right)$, so

$$\begin{aligned}\|\mathbf{H}_{\theta,\delta}\| &\leq \frac{1}{n_q} \sum_{i=1}^{n_q} \frac{1}{r_i} \cdot \left\| \mathbf{J}^{(i)} \mathbf{J}^{(i)\top} \cdot \left(\text{Iden} - \frac{1}{r_i} \cdot \delta T \psi^{(i)\top} \right) \right\| \\ &\leq \frac{C_{\text{ratio}}}{n_q} \sum_{i=1}^{n_q} \underbrace{\left\| \mathbf{J}^{(i)} \mathbf{J}^{(i)\top} \right\| \cdot \left\| \text{Iden} - \frac{1}{r_i} \cdot \delta T \psi^{(i)\top} \right\|}_C\end{aligned}$$

From calculation, we know $\|C\| \leq 1 + |(r_i - 1)/r_i| \leq 2 + C_{\text{ratio}}$. Therefore $\|\mathbf{H}_{\theta,\delta}\| \leq \frac{C_{\text{ratio}}^2 + 2C_{\text{ratio}}}{n_q} \sum_{i=1}^{n_q} \|\mathbf{J}^{(i)} \mathbf{J}^{(i)\top}\| \leq (C_{\text{ratio}}^2 + 2C_{\text{ratio}}) C_4$. Therefore $\|\mathbf{H}_{\theta,\delta} \mathbf{H}_{\delta,\delta}^{-1}\|$ is upperbounded by $(C_{\text{ratio}}^4 + 2C_{\text{ratio}}^3) C_4 / \Lambda''_{\min}$.

Refer to [21, 6] for inequalities of eigenvalue of matrix summation and product. \square

B.4 Proof of Theorem 1

Proof. We denote Hessian \mathbf{H} as a block matrix:

$$\mathbf{H} = \nabla^2 \ell(\delta, \theta) = \begin{pmatrix} \mathbf{H}_{11} & \mathbf{H}_{12} \\ \mathbf{H}_{21} & \mathbf{H}_{22} \end{pmatrix} = \begin{pmatrix} \nabla_{\delta}^2 \ell(\delta, \theta) & \nabla_{\delta} \nabla_{\theta} \ell(\delta, \theta) \\ \nabla_{\theta} \nabla_{\delta} \ell(\delta, \theta) & \nabla_{\theta}^2 \ell(\delta, \theta) \end{pmatrix},$$

then Assumption 2 states that for every $\delta \in \Delta_{n_q}$ and $\theta \in \Theta$, $\lambda(\mathbf{H}_{21} \mathbf{H}_{11}^{-1} \mathbf{H}_{12})$ is lower bounded by $2 \|\mathbf{H}_{22}\|$ and $\|\mathbf{H}_{21} \mathbf{H}_{11}^{-1}\|$ is upper bounded.

We can write the optimality condition of (7) and expand them at $(\delta^* \equiv \mathbf{0}, \theta^*)$:

$$\nabla_{\delta} \ell(\hat{\delta}, \hat{\theta}) = \mathbf{0} = \nabla_{\delta} \ell(\delta^*, \theta^*) + \bar{\mathbf{H}}_{11}(\hat{\delta} - \delta^*) + \bar{\mathbf{H}}_{12}(\hat{\theta} - \theta^*) \quad (18)$$

$$\nabla_{\theta} \ell(\hat{\delta}, \hat{\theta}) = \mathbf{0} = \nabla_{\theta} \ell(\delta^*, \theta^*) + \bar{\mathbf{H}}_{21}(\hat{\delta} - \delta^*) + \bar{\mathbf{H}}_{22}(\hat{\theta} - \theta^*), \quad (19)$$

where $\bar{\mathbf{H}}$ is the Hessian evaluated at a $(\bar{\delta}, \bar{\theta})$ which is in between $(\hat{\delta}, \hat{\theta})$ and (δ^*, θ^*) in an *element-wise fashion*. This expansion is basically one-dimensional mean-value theorem applied on *each individual dimension* of $\nabla_{\delta} \ell(\hat{\delta}, \hat{\theta})$ and $\nabla_{\theta} \ell(\hat{\delta}, \hat{\theta})$.

Given (18) and (19) we can solve equations for $\hat{\delta} - \delta^*$ and $\hat{\theta} - \theta^*$.

From (18) we can get

$$\hat{\delta} - \delta^* = \bar{\mathbf{H}}_{11}^{-1} \left[-\nabla_{\delta} \ell(\delta^*, \theta^*) - \bar{\mathbf{H}}_{12}(\hat{\theta} - \theta^*) \right]. \quad (20)$$

Substituting (20) into (19) we get

$$\mathbf{0} = \nabla_{\theta} \ell(\delta^*, \theta^*) - \bar{\mathbf{H}}_{21} \bar{\mathbf{H}}_{11}^{-1} \nabla_{\delta} \ell(\delta^*, \theta^*) + \left[-\bar{\mathbf{H}}_{21} \bar{\mathbf{H}}_{11}^{-1} \bar{\mathbf{H}}_{12} + \bar{\mathbf{H}}_{22} \right] (\hat{\theta} - \theta^*).$$

Rearranging terms, we get

$$\hat{\theta} - \theta^* = \left[\bar{\mathbf{H}}_{21} \bar{\mathbf{H}}_{11}^{-1} \bar{\mathbf{H}}_{12} - \bar{\mathbf{H}}_{22} \right]^{-1} \left(\nabla_{\theta} \ell(\delta^*, \theta^*) - \bar{\mathbf{H}}_{21} \bar{\mathbf{H}}_{11}^{-1} \nabla_{\delta} \ell(\delta^*, \theta^*) \right) \quad (21)$$

$$= \left[-\bar{\mathbf{H}}_{21} \bar{\mathbf{H}}_{11}^{-1} \bar{\mathbf{H}}_{12} + \bar{\mathbf{H}}_{22} \right]^{-1} \bar{\mathbf{H}}_{21} \bar{\mathbf{H}}_{11}^{-1} \nabla_{\delta} \ell(\delta^*, \theta^*). \quad (22)$$

The last line uses the fact that $\nabla_{\theta} \ell(\delta^*, \theta^*) \equiv \mathbf{0}$.

Weyl's inequality states:

$$\lambda_{\min}(A + B) \geq \lambda_{\min}(A) + \lambda_{\min}(B).$$

As $\bar{\delta} \in \Delta_{n_q}$ and $\bar{\theta} \in \Theta$, \bar{H} is regulated by Assumption 2. Since

$$\lambda_{\min}(-\bar{H}_{21} \bar{H}_{11}^{-1} \bar{H}_{12}) \geq \Lambda_{\min}$$

and

$$\lambda_{\min}(\bar{H}_{22}) \geq -\|\bar{H}_{22}\| \geq -\frac{\Lambda_{\min}}{2}$$

which are assumed by Assumption 2, we have

$$\lambda_{\min}(-\bar{H}_{21} \bar{H}_{11}^{-1} \bar{H}_{12} + \bar{H}_{22}) \geq \Lambda_{\min}/2 > 0.$$

Denote $-\bar{H}_{21} \bar{H}_{11}^{-1} \bar{H}_{12} + \bar{H}_{22}$ as \bar{H}/\bar{H}_{22} (it is actually the Schur Complement of \bar{H}). Using Holder's inequality, we get

$$\begin{aligned} \|\hat{\theta} - \theta^*\| &\leq \left\| [\bar{H}/\bar{H}_{22}]^{-1} \right\| \left\| \bar{H}_{21} \bar{H}_{11}^{-1} \right\| \|\nabla_{\delta} \ell(\delta^*, \theta^*)\| \\ &\leq \frac{\left\| \bar{H}_{21} \bar{H}_{11}^{-1} \right\|}{\lambda_{\min}[\bar{H}/\bar{H}_{22}]} \cdot \|\nabla_{\delta} \ell(\delta^*, \theta^*)\| \leq \frac{2\Lambda_{\max}}{\Lambda_{\min}} \cdot \|\nabla_{\delta} \ell(\delta^*, \theta^*)\|. \end{aligned} \quad (23)$$

Further, we have $\mathbb{E}_q[T_{\theta^*} \mathbf{f}(x)] = \mathbb{E}_{p_{\theta^*}}[T_{\theta^*} \mathbf{f}(x)] = \mathbf{0}$. The first equality is due to Assumption 1 and the second equality is given by Stein identity.

Therefore, $\nabla_{\delta} \ell(\delta^*, \theta^*) = \frac{1}{n_q} \sum_{i=1}^{n_q} T_{\theta^*} \mathbf{f}(x^{(i)}) - \mathbf{0} = \frac{1}{n_q} \sum_{i=1}^{n_q} T_{\theta^*} \mathbf{f}(x^{(i)}) - \mathbb{E}_q[T_{\theta^*} \mathbf{f}(x)]$, which converges to 0 in ℓ_2 norm in probability due to Assumption 3. This gives the convergence in probability of $\|\hat{\theta} - \theta^*\|$. Finite sample convergence rate can be given if the convergence rate of $\|\nabla_{\delta} \ell(\delta^*, \theta^*)\|$ is known.

Now we show the consistency of $\hat{\delta}$. From (20) we can see that

$$\hat{\delta} - \delta^* = -\bar{H}_{11}^{-1} \nabla_{\delta} \ell(\delta^*, \theta^*) - \bar{H}_{11}^{-1} \bar{H}_{12} (\hat{\theta} - \theta^*),$$

and due to Holder's inequality, we get

$$\begin{aligned} \|\hat{\delta} - \delta^*\| &= \left\| -\bar{H}_{11}^{-1} \right\| \|\nabla_{\delta} \ell(\delta^*, \theta^*)\| + \left\| \bar{H}_{11}^{-1} \bar{H}_{12} \right\| \|\hat{\theta} - \theta^*\| \\ &\leq \frac{1}{\Lambda'_{\min}} \|\nabla_{\delta} \ell(\delta^*, \theta^*)\| + \Lambda_{\max} \|\hat{\theta} - \theta^*\|. \end{aligned} \quad (24)$$

Combine (24) with (23) we get

$$\|\hat{\delta} - \delta^*\| \leq \frac{2\Lambda_{\max}^2 \Lambda'_{\min} + \Lambda_{\min}}{\Lambda_{\min} \Lambda'_{\min}} \cdot \|\nabla_{\delta} \ell(\delta^*, \theta^*)\|$$

Again, due to Assumption 3, $\|\nabla_{\delta} \ell(\delta^*, \theta^*)\| \xrightarrow{\mathbb{P}} \mathbf{0}$. This completes the proof. \square

B.5 Proof of Theorem 2

Proof. Due to Assumption 4, it can be seen that $\bar{H} \xrightarrow{\mathbb{P}} \mathbb{E}_q[\bar{H}]$. Moreover, as $\bar{\theta} \xrightarrow{\mathbb{P}} \theta^*$ and $\bar{\delta} \xrightarrow{\mathbb{P}} \mathbf{0}$ (proved in Theorem 1), we can see $\mathbb{E}_q[\bar{H}] \xrightarrow{\mathbb{P}} \mathbb{E}_q[H^*]$ due to continuous mapping. Thus $\bar{H} = \mathbb{E}_q[H^*] + o_p(1)$. From now on, for simplicity, let us denote $-\mathbb{E}_q[H^*]$ as \mathbf{I}^2 .

² \mathbf{I} for ‘‘information matrix’’. Do not confuse with the identify matrix which is denoted as \mathbf{I}_{den} in this paper

We again write the optimality condition of (7) and apply asymptotic expansion at $(\delta^* \equiv \mathbf{0}, \theta^*)$:

$$\nabla_{\delta} \ell(\hat{\delta}, \hat{\theta}) = \mathbf{0} = \nabla_{\delta} \ell(\delta^*, \theta^*) + (-\mathbf{I}_{11} + o_p(1))(\hat{\delta} - \delta^*) + (-\mathbf{I}_{12} + o_p(1))(\hat{\theta} - \theta^*) \quad (25)$$

$$\nabla_{\theta} \ell(\hat{\delta}, \hat{\theta}) = \mathbf{0} = \nabla_{\theta} \ell(\delta^*, \theta^*) + (-\mathbf{I}_{21} + o_p(1))(\hat{\delta} - \delta^*) + (-\mathbf{I}_{22} + o_p(1))(\hat{\theta} - \theta^*). \quad (26)$$

Note we have replaced all \bar{H} with $-\mathbf{I} + o_p(1)$, and $o_p(1)$ will be ignored in future algebraic calculations.

We now get an asymptotic version of (22):

$$\begin{aligned} \sqrt{n_q} (\hat{\theta} - \theta^*) &\rightsquigarrow -(\mathbf{I}_{21} \mathbf{I}_{11}^{-1} \mathbf{I}_{12} - \mathbf{I}_{22})^{-1} \mathbf{I}_{21} \mathbf{I}_{11}^{-1} \nabla_{\delta} \ell(\delta^*, \theta^*) \cdot \sqrt{n_q} \\ &= -(\mathbf{I}_{21} \mathbf{I}_{11}^{-1} \mathbf{I}_{12})^{-1} \mathbf{I}_{21} \mathbf{I}_{11}^{-1} \nabla_{\delta} \ell(\delta^*, \theta^*) \cdot \sqrt{n_q} \end{aligned}$$

The last equality is due to $\mathbf{I}_{22} \equiv \mathbf{0}$.

Noticing that $\mathbf{I}_{11}^{-1} \nabla_{\delta} \ell(\delta^*, \theta^*) \cdot \sqrt{n_q}$ is a sum of independent random variables with zero mean and covariance $-\mathbf{I}_{11}^{-1}$. Applying CLT on $\mathbf{I}_{11}^{-1} \nabla_{\delta} \ell(\delta^*, \theta^*) \cdot \sqrt{n_q}$ yields

$$\mathbf{I}_{11}^{-1} \nabla_{\delta} \ell(\delta^*, \theta^*) \rightsquigarrow \mathcal{N}(\mathbf{0}, -\mathbf{I}_{11}^{-1}),$$

thus

$$\sqrt{n_q} (\hat{\theta} - \theta^*) \rightsquigarrow \mathcal{N}[\mathbf{0}, (-\mathbf{I}_{21} \mathbf{I}_{11}^{-1} \mathbf{I}_{12})^{-1}].$$

□

B.6 Proof of Lemma 3

Proof. Let us shorten the Stein feature vector $T_{\theta} f(x)$ as $t(x; \theta) \in \mathbb{R}^b$ and t as $t(x; \theta^*)$. We start by computing each factors in the variance. Since $r_{\theta}(x; \delta^*) = 1$ holds for all θ , we have $\nabla_{\theta} r_{\theta}(x; \delta^*) = \mathbf{0}$. Then, we have

$$\begin{aligned} -\mathbb{E}_q[\mathbf{H}_{\delta, \delta}^*] &= -\mathbb{E}_q[\nabla_{\delta}^2 \log r_{\theta^*}(x; \delta^*)] \\ &= \mathbb{E}_q\left[\frac{1}{r(x; \delta^*, \theta^*)^2} t t^{\top}\right] = \mathbb{E}_q[t t^{\top}] \in \mathbb{R}^{b \times b}, \\ \mathbb{E}_q[\mathbf{H}_{\theta, \delta}^*] &= \mathbb{E}_q\left[\frac{1}{r} \nabla_{\theta} t(x; \theta^*)^{\top} - \frac{1}{r^2} \nabla_{\theta} r_{\theta^*}(x; \delta^*) t(x; \theta^*)^{\top}\right] \\ &= \mathbb{E}_q[\nabla_{\theta} t(x; \theta^*)^{\top}] \in \mathbb{R}^{\dim(\theta) \times b}. \end{aligned}$$

Since the equality $\mathbb{E}_{p_{\theta}}[t(x; \theta)] = \mathbf{0}$ holds for all θ , we have $\nabla_{\theta} \mathbb{E}_{p_{\theta}}[t(x; \theta)] = \mathbf{0}$. Exchangeability of the integration and the derivative yields

$$\nabla_{\theta} \mathbb{E}_{p_{\theta}}[t(x; \theta)] = \mathbb{E}_{p_{\theta}}[s(x; \theta) t(x; \theta)^{\top}] + \mathbb{E}_{p_{\theta}}[\nabla_{\theta} t(x; \theta)^{\top}] = \mathbf{0}.$$

As a result, we obtain

$$\mathbb{E}_q[\mathbf{H}_{\theta, \delta}^*] = -\mathbb{E}_q[s t^{\top}].$$

□

B.7 Proof of Theorem 5

Proof. Use Taylor series to expand $\mathbb{E}_q[\ell(\hat{\delta}, \hat{\theta})]$ on (θ^*, δ^*) , we get

$$\begin{aligned} \mathbb{E}_q[\ell(\hat{\delta}, \hat{\theta})] &= \mathbb{E}_q[\ell(\delta^*, \theta^*)] + \nabla_{\delta} \mathbb{E}_q[\ell(\delta^*, \theta^*)]^{\top} [\hat{\delta} - \delta^*] + \nabla_{\theta} \mathbb{E}_q[\ell(\delta^*, \theta^*)]^{\top} [\hat{\theta} - \theta^*] \\ &\quad + \frac{1}{2} [\hat{\eta} - \eta^*]^{\top} \nabla_{\eta}^2 \mathbb{E}_q[\ell(\bar{\eta})] [\hat{\eta} - \eta^*] \\ &= 0 + 0 + 0 + \frac{1}{2} [\hat{\eta} - \eta^*]^{\top} \nabla_{\eta}^2 \mathbb{E}_q[\ell(\bar{\eta})] [\hat{\eta} - \eta^*] \end{aligned} \quad (27)$$

where we denote $\boldsymbol{\eta} := \begin{bmatrix} \boldsymbol{\delta} \\ \boldsymbol{\theta} \end{bmatrix}$ for short and $\bar{\boldsymbol{\eta}}$ is defined in between $\hat{\boldsymbol{\eta}}$ and $\boldsymbol{\eta}^*$ in an element-wise fashion. The second equality is due to $\boldsymbol{\delta}^* = \mathbf{0}$ and $\mathbb{E}_q [\nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*)] = \mathbf{0}$, which is given by Stein identity. Similarly we can expand

$$\ell(\hat{\boldsymbol{\delta}}, \hat{\boldsymbol{\theta}}) = \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*)^\top [\hat{\boldsymbol{\delta}} - \boldsymbol{\delta}^*] + \frac{1}{2} [\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}^*]^\top \nabla_{\boldsymbol{\eta}}^2 \ell(\bar{\boldsymbol{\eta}}) [\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}^*], \quad (28)$$

where $\bar{\boldsymbol{\eta}}$ is similarly defined as $\bar{\boldsymbol{\eta}}$. It can be seen that $\nabla_{\boldsymbol{\eta}}^2 \ell(\bar{\boldsymbol{\eta}}) \xrightarrow{\mathbb{P}} -\mathbf{I}$ and $\nabla_{\boldsymbol{\eta}}^2 \mathbb{E}_q [\ell \bar{\boldsymbol{\eta}}] \xrightarrow{\mathbb{P}} -\mathbf{I}$ due to Assumption 4 and our consistency results. Taking the difference between (27) and (28) after multiplying n_q yields

$$n_q \mathbb{E}_q [\ell(\hat{\boldsymbol{\delta}}, \hat{\boldsymbol{\theta}})] - n_q \ell(\hat{\boldsymbol{\delta}}, \hat{\boldsymbol{\theta}}) = -n_q \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*)^\top [\hat{\boldsymbol{\delta}} - \boldsymbol{\delta}^*] + o_p(1).$$

Substitute $(\hat{\boldsymbol{\delta}} - \boldsymbol{\delta}^*)$ with (20) we get

$$n_q \mathbb{E}_q [\ell(\hat{\boldsymbol{\delta}}, \hat{\boldsymbol{\theta}})] - n_q \ell(\hat{\boldsymbol{\delta}}, \hat{\boldsymbol{\theta}}) = n_q \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*)^\top [\bar{\mathbf{H}}_{11}^{-1} \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*) + \bar{\mathbf{H}}_{11}^{-1} \bar{\mathbf{H}}_{12} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)] + o_p(1).$$

Substitute $(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)$ using (22), we get

$$\begin{aligned} n_q \mathbb{E}_q [\ell(\hat{\boldsymbol{\delta}}, \hat{\boldsymbol{\theta}})] - n_q \ell(\hat{\boldsymbol{\delta}}, \hat{\boldsymbol{\theta}}) &= n_q \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*)^\top \bar{\mathbf{H}}_{11}^{-1} \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*) \\ &\quad - n_q \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*)^\top \bar{\mathbf{H}}_{11}^{-1} \bar{\mathbf{H}}_{12} [\bar{\mathbf{H}} / \bar{\mathbf{H}}_{22}]^{-1} \bar{\mathbf{H}}_{21} \bar{\mathbf{H}}_{11}^{-1} \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*) + o_p(1) \end{aligned} \quad (29)$$

Replacing submatrices of $\bar{\mathbf{H}}_{a,b}$ using submatrices of $-\mathbf{I}_{a,b}$ in (29) and using the fact that $\mathbf{I}_{22} \equiv \mathbf{0}$ (due to $\boldsymbol{\delta}^* = \mathbf{0}$),

$$\begin{aligned} n_q \mathbb{E}_q [\ell(\hat{\boldsymbol{\delta}}, \hat{\boldsymbol{\theta}})] - n_q \ell(\hat{\boldsymbol{\delta}}, \hat{\boldsymbol{\theta}}) &= -\sqrt{n_q} \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*)^\top \mathbf{I}_{11}^{-1} \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*) \sqrt{n_q} \\ &\quad + \sqrt{n_q} \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*)^\top \mathbf{I}_{11}^{-1} \mathbf{I}_{12} [\mathbf{I}_{21} \mathbf{I}_{11}^{-1} \mathbf{I}_{12}]^{-1} \mathbf{I}_{21} \mathbf{I}_{11}^{-1} \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*) \sqrt{n_q} + o_p(1) \end{aligned} \quad (30)$$

Taking the expectation,

$$\begin{aligned} n_q \mathbb{E} \left\{ \mathbb{E}_q [\ell(\hat{\boldsymbol{\delta}}, \hat{\boldsymbol{\theta}})] - \ell(\hat{\boldsymbol{\delta}}, \hat{\boldsymbol{\theta}}) \right\} &= -\text{trace}(\mathbf{I}_{11} \mathbf{I}_{11}^{-1}) + \text{trace}(\mathbf{I}_{11}^{-1} \mathbf{I}_{12} [\mathbf{I}_{21} \mathbf{I}_{11}^{-1} \mathbf{I}_{12}]^{-1} \mathbf{I}_{21}) + o_p(1) \\ &= -\text{rank}(\mathbf{I}_{11}) + \text{rank}(\mathbf{I}_{21} \mathbf{I}_{11}^{-1} \mathbf{I}_{12}) + o_p(1). \end{aligned}$$

In the case when $\mathbf{I}_{11} \in \mathbb{R}^{b \times b}$, $\mathbf{I}_{12} \in \mathbb{R}^{b \times \dim(\boldsymbol{\theta})}$ are full-rank and $\dim(\boldsymbol{\theta}) \leq b$, $\text{rank}(\mathbf{I}_{11}) = b$ and $\text{rank}(\mathbf{I}_{21} \mathbf{I}_{11}^{-1} \mathbf{I}_{12}) = \dim(\boldsymbol{\theta})$, which completes the proof. \square

B.8 The Asymptotic Distribution of $2n_q \ell(\hat{\boldsymbol{\delta}}, \hat{\boldsymbol{\theta}})$

We show $2n_q \ell(\hat{\boldsymbol{\delta}}, \hat{\boldsymbol{\theta}})$ follows a χ^2 distribution based on previously assumed assumptions.

Theorem 6. Suppose Assumption 1, 2, 3 and 4 holds, $\mathbb{E}_q [\mathbf{H}_{\boldsymbol{\delta}, \boldsymbol{\delta}}^*]$ is invertible and $\mathbb{E}_q [\mathbf{H}_{\boldsymbol{\delta}, \boldsymbol{\theta}}^*]$ are full-rank and $\dim(\boldsymbol{\theta}) < b$, then $2n_q \ell(\hat{\boldsymbol{\delta}}, \hat{\boldsymbol{\theta}}) \rightsquigarrow \chi^2(b - \dim(\boldsymbol{\theta}))$.

Proof. First we expand $2n_q \ell(\hat{\boldsymbol{\delta}}, \hat{\boldsymbol{\theta}})$ using mean value theorem:

$$2n_q \ell(\hat{\boldsymbol{\delta}}, \hat{\boldsymbol{\theta}}) = 2n_q \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*)^\top d\boldsymbol{\delta} + n_q d\boldsymbol{\delta} \bar{\mathbf{H}}_{11} d\boldsymbol{\delta} + n_q d\boldsymbol{\delta} \bar{\mathbf{H}}_{12} d\boldsymbol{\theta} + n_q d\boldsymbol{\theta} \bar{\mathbf{H}}_{21} d\boldsymbol{\delta} + n_q d\boldsymbol{\theta} \bar{\mathbf{H}}_{22} d\boldsymbol{\theta} \quad (31)$$

where $d\mathbf{t}$ is short for $\hat{\mathbf{t}} - \mathbf{t}^*$. Note $\ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*) = 0$. Now we analyze each term.

From the proof in Section B.7 we know

$$\begin{aligned} 2n_q \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*)^\top d\boldsymbol{\delta} &= 2n_q \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*)^\top \mathbf{I}_{11}^{-1} \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*) \\ &\quad - 2n_q \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*)^\top \mathbf{I}_{11}^{-1} \mathbf{I}_{12} [\mathbf{I}_{21} \mathbf{I}_{11}^{-1} \mathbf{I}_{12}]^{-1} \mathbf{I}_{21} \mathbf{I}_{11}^{-1} \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*) + o_p(1). \end{aligned} \quad (32)$$

With the help of (20) and (22) and a few algebra we can see that

$$\begin{aligned} n_q d\boldsymbol{\delta} \mathbf{I}_{11} d\boldsymbol{\delta} &= n_q \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*)^\top \mathbf{I}_{11}^{-1} \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*) \\ &\quad - n_q \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*)^\top \mathbf{I}_{11}^{-1} \mathbf{I}_{12} [\mathbf{I}_{21} \mathbf{I}_{11}^{-1} \mathbf{I}_{12}]^{-1} \mathbf{I}_{21} \mathbf{I}_{11}^{-1} \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*) + o_p(1). \end{aligned} \quad (33)$$

Similar calculations also show $n_q d\boldsymbol{\delta} \bar{\mathbf{H}}_{12} d\boldsymbol{\theta} = n_q d\boldsymbol{\theta} \bar{\mathbf{H}}_{21} d\boldsymbol{\delta} = o_p(1)$ and $n_q d\boldsymbol{\theta} \bar{\mathbf{H}}_{22} d\boldsymbol{\theta} = o_p(1)$. Combine (31), (32) and (33), we can see that

$$\begin{aligned} 2n_q \ell(\hat{\boldsymbol{\delta}}, \hat{\boldsymbol{\theta}}) &= n_q \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*)^\top \mathbf{I}_{11}^{-1} \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*) \\ &\quad - n_q \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*)^\top \mathbf{I}_{11}^{-1} \mathbf{I}_{12} [\mathbf{I}_{21} \mathbf{I}_{11}^{-1} \mathbf{I}_{12}]^{-1} \mathbf{I}_{21} \mathbf{I}_{11}^{-1} \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*) + o_p(1) \\ &= \sqrt{n_q} \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*)^\top \mathbf{I}_{11}^{-1} \left\{ \text{Iden} - \mathbf{I}_{12} [\mathbf{I}_{21} \mathbf{I}_{11}^{-1} \mathbf{I}_{12}]^{-1} \mathbf{I}_{21} \mathbf{I}_{11}^{-1} \right\} \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*) \sqrt{n_q} + o_p(1), \end{aligned}$$

where Iden is identify matrix. Denote $\text{Iden} - \mathbf{I}_{12} [\mathbf{I}_{21} \mathbf{I}_{11}^{-1} \mathbf{I}_{12}]^{-1} \mathbf{I}_{21} \mathbf{I}_{11}^{-1}$ as A . One can verify that $\nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*)^\top \mathbf{I}_{11}^{-1} A$ has covariance $\mathbf{I}_{11}^{-1} A^3$. By checking the eigenvalues of A^4 , it can be seen that $\text{rank}(A) = db - \dim(\boldsymbol{\theta})$ and assuming \mathbf{I}_{11}^{-1} is full rank, $\text{rank}(\mathbf{I}_{11}^{-1} A) = db - \dim(\boldsymbol{\theta})$. Therefore $\sqrt{n_q} \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*)^\top \mathbf{I}_{11}^{-1} A$ is asymptotically a degenerated multivariate normal variable with covariance matrix $\mathbf{I}_{11}^{-1} A$.

We can rewrite $\sqrt{n_q} \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*)^\top \mathbf{I}_{11}^{-1} A \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*) \sqrt{n_q}$ as

$$\sqrt{n_q} \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*)^\top \mathbf{I}_{11}^{-1} A [\mathbf{I}_{11}^{-1} A]^+ \mathbf{I}_{11}^{-1} A \nabla_{\boldsymbol{\delta}} \ell(\boldsymbol{\delta}^*, \boldsymbol{\theta}^*) \sqrt{n_q},$$

where T^+ is the pseudoinverse. This quadratic form has a χ^2 distribution with degree of freedom $\text{rank}(\mathbf{I}_{11}^{-1} A) = db - \dim(\boldsymbol{\theta})$. \square

B.9 Proof of Proposition 3

Proof. We convert the SDRE problem (5) as the following equivalent problem:

$$\max_{\boldsymbol{\delta}, \boldsymbol{\epsilon}} \sum_{i=1}^{n_q} \log \epsilon_i \quad \text{s.t.} \quad \forall i \in \{1 \dots n_q\}, \quad \boldsymbol{\delta}^\top T_{\boldsymbol{\theta}} \mathbf{f}(\mathbf{x}_q^{(i)}) + 1 = \epsilon_i.$$

Let us introduce Lagrangian multipliers $\mu_1 \dots \mu_{n_q}$ over all the constraints. We can write the Lagrangian:

$$\min_{\boldsymbol{\mu}} \max_{\boldsymbol{\delta}, \boldsymbol{\epsilon}} \sum_{i=1}^{n_q} (\log \epsilon_i) - \sum_{i=1}^{n_q} \mu_i \left(\boldsymbol{\delta}^\top T_{\boldsymbol{\theta}} \mathbf{f}(\mathbf{x}_q^{(i)}) + 1 - \epsilon_i \right) \quad (34)$$

Solve the inner max problem with respect to $\boldsymbol{\epsilon}$,

$$\max_{\boldsymbol{\epsilon}} \sum_{i=1}^{n_q} \log \epsilon_i + \mu_i \epsilon_i = \sum_{i=1}^{n_q} [-(\log -\mu_i) - 1], \quad (35)$$

when $\epsilon_i = -\frac{1}{\mu_i}$. This also implies the relationship between the dual parameter μ_i and the primal parameter $\boldsymbol{\delta}$: $r_{\boldsymbol{\theta}}(\mathbf{x}_q^{(i)}; \boldsymbol{\delta}) = \boldsymbol{\delta}^\top T_{p\boldsymbol{\theta}} \mathbf{f}(\mathbf{x}_q^{(i)}) + 1 = \epsilon_i = -\frac{1}{\mu_i}$.

The inner optimization with respect to $\boldsymbol{\delta}$, i.e., $\max_{\boldsymbol{\delta}} - \sum_{i=1}^{n_q} \mu_i \boldsymbol{\delta}^\top T_{\boldsymbol{\theta}} \mathbf{f}(\mathbf{x}_q^{(i)})$ is a linear programming and is only bounded when $\sum_{i=1}^{n_q} \mu_i T_{\boldsymbol{\theta}} \mathbf{f}(\mathbf{x}_q^{(i)}) = \mathbf{0}$ and achieves the optimal value 0.

Substituting the optimal values of these two maximization results into the Lagrangian and adding constraint $\sum_{i=1}^{n_q} \mu_i T_{\boldsymbol{\theta}} \mathbf{f}(\mathbf{x}_q^{(i)}) = \mathbf{0}$ gives the Lagrangian dual (9). Moreover, the primal problem in (5) is concave, we can verify the Slater's condition holds at $\boldsymbol{\delta} = \mathbf{0}, \boldsymbol{\epsilon} = \mathbf{1}$ thus the strong duality holds. \square

³Some calculations show $A^\top \mathbf{I}_{11}^{-1} A = \mathbf{I}_{11}^{-1} A$.

⁴ $\text{eig}(\text{Iden} - T) = 1 - \text{eig}(T)$ and $\text{eig}(ST) = \text{eig}(TS)$.