
Supplemental Information: Graphical Models for Non-Negative Data Using Generalized Score Matching

Shiqing Yu

University of Washington

Mathias Drton

University of Washington

Ali Shojaie

University of Washington

A.2 SCORE MATCHING

The following lemma is used in the proof of Theorem 2.

Lemma A.1. *Assuming that f and g are differentiable a.e., then for all $j = 1, \dots, m$,*

$$\lim_{a \nearrow +\infty, b \searrow 0^+} f(\mathbf{x}_{-j}; a)g(\mathbf{x}_{-j}; a) - f(\mathbf{x}_{-j}; b)g(\mathbf{x}_{-j}; b) = \int_0^\infty f(\mathbf{x}) \frac{\partial g(\mathbf{x})}{\partial x_j} dx_j + \int_0^\infty g(\mathbf{x}) \frac{\partial f(\mathbf{x})}{\partial x_j} dx_j,$$

where $(\mathbf{x}_{-j}; a)$ is the vector obtained by replacing the j -th component of \mathbf{x} by a .

Proof. This is just an analog of Lemma 4 from Hyvärinen (2005) proved by integrating the partial derivatives. \square

Proof of Theorem 2. Recall the following assumptions given in Section 2.3.

- (A1) $p_0(\mathbf{x})h_j(x_j)\partial_j \log p(\mathbf{x}) \rightarrow 0$ as $x_j \nearrow +\infty$ and as $x_j \searrow 0^+$, $\forall \mathbf{x}_{-j} \in \mathbb{R}_+^{m-1}$, $\forall p \in \mathcal{P}_+$,
- (A2) $\mathbb{E}_{p_0} \|\nabla \log p(\mathbf{X}) \circ \mathbf{h}^{1/2}(\mathbf{X})\|_2^2 < +\infty$, $\mathbb{E}_{p_0} \|(\nabla \log p(\mathbf{X}) \circ \mathbf{h}(\mathbf{X}))'\|_1 < +\infty$, $\forall p \in \mathcal{P}_+$,

where

$$\partial_j \log p(\mathbf{x}) \equiv \left. \frac{\partial \log p(\mathbf{y})}{\partial y_j} \right|_{\mathbf{y}=\mathbf{x}}.$$

Without explicitly writing the domains \mathbb{R}_+ or \mathbb{R}_+^m in all integrals, by (4) we have

$$\begin{aligned} J_{\mathbf{h}}(p) &= \frac{1}{2} \int p_0(\mathbf{x}) \left[\|\nabla \log p(\mathbf{x}) \circ \mathbf{h}^{1/2}(\mathbf{x})\|_2^2 - 2(\nabla \log p(\mathbf{x}) \circ \mathbf{h}^{1/2}(\mathbf{x}))^\top (\nabla \log p_0(\mathbf{x}) \circ \mathbf{h}^{1/2}(\mathbf{x})) \right. \\ &\quad \left. + \|\nabla \log p_0(\mathbf{x}) \circ \mathbf{h}^{1/2}(\mathbf{x})\|_2^2 \right] d\mathbf{x} \\ &= \underbrace{\frac{1}{2} \int p_0(\mathbf{x}) \sum_{j=1}^m h_j(x_j) \left(\frac{\partial \log p(\mathbf{x})}{\partial x_j} \right)^2 d\mathbf{x}}_{\equiv A} - \underbrace{\int p_0(\mathbf{x}) \sum_{j=1}^m h_j(x_j) \frac{\partial \log p(\mathbf{x})}{\partial x_j} \frac{\partial \log p_0(\mathbf{x})}{\partial x_j} d\mathbf{x}}_{\equiv B} \\ &\quad + \underbrace{\frac{1}{2} \int p_0(\mathbf{x}) \sum_{j=1}^m h_j(x_j) \left(\frac{\partial \log p_0(\mathbf{x})}{\partial x_j} \right)^2 d\mathbf{x}}_{\equiv C}, \end{aligned}$$

where A will simply appear in the final display as is, C is a constant as it only involves the true pdf p_0 , and we wish to simplify B by integration by parts. We can split the integral into these three parts since A and C are assumed finite in the first part of (A2), and the integrand in B is integrable since $|2ab| \leq a^2 + b^2$. Thus, by linearity and Fubini's theorem, we can write

$$\begin{aligned}
B &= - \sum_{j=1}^m \int p_0(\mathbf{x}) h_j(x_j) \frac{\partial \log p(\mathbf{x})}{\partial x_j} \frac{\partial \log p_0(\mathbf{x})}{\partial x_j} d\mathbf{x} \\
&= - \sum_{j=1}^m \int \left[\int p_0(\mathbf{x}) h_j(x_j) \frac{\partial \log p(\mathbf{x})}{\partial x_j} \frac{\partial \log p_0(\mathbf{x})}{\partial x_j} dx_j \right] d\mathbf{x}_{-j}.
\end{aligned}$$

By the fact that $\frac{\partial \log p_0(\mathbf{x})}{\partial x_j} = \frac{1}{p_0(\mathbf{x})} \frac{\partial p_0(\mathbf{x})}{\partial x_j}$, this can be simplified to

$$B = - \sum_{j=1}^m \int \left[\int \frac{\partial p_0(\mathbf{x})}{\partial x_j} h_j(x_j) \frac{\partial \log p(\mathbf{x})}{\partial x_j} dx_j \right] d\mathbf{x}_{-j}.$$

Then by Lemma A.1 and assumption (A1),

$$\begin{aligned}
B &= - \sum_{j=1}^m \int \left[\lim_{a \nearrow \infty, b \searrow 0^+} [p_0(\mathbf{x}_{-j}; a) h_j(a) \partial_j \log p(\mathbf{x}_{-j}, a) - p_0(\mathbf{x}_{-j}; b) h_j(b) \partial_j \log p(\mathbf{x}_{-j}, b)] \right. \\
&\quad \left. - \int p_0(\mathbf{x}) \frac{\partial (h_j(x_j) \partial_j \log p(\mathbf{x}))}{\partial x_j} dx_j \right] d\mathbf{x}_{-j} \\
&= \sum_{j=1}^m \int \left[\int p_0(\mathbf{x}) \frac{\partial (h_j(x_j) \partial_j \log p(\mathbf{x}))}{\partial x_j} dx_j \right] d\mathbf{x}_{-j}.
\end{aligned}$$

Justified by the second half of (A2), by Fubini-Tonelli and linearity again

$$\begin{aligned}
B &= \sum_{j=1}^m \int p_0(\mathbf{x}) \frac{\partial (h_j(x_j) \partial_j \log p(\mathbf{x}))}{\partial x_j} d\mathbf{x}, \\
&= \sum_{j=1}^m \int h'_j(x_j) \frac{\partial \log p(\mathbf{x})}{\partial x_j} p_0(\mathbf{x}) d\mathbf{x} + \sum_{j=1}^m \int h_j(x_j) \frac{\partial^2 \log p(\mathbf{x})}{\partial x_j^2} p_0(\mathbf{x}) d\mathbf{x}.
\end{aligned}$$

Thus,

$$\begin{aligned}
J_{\mathbf{h}}(p) &= B + A + C \\
&= \int_{\mathbb{R}_+^m} p_0(\mathbf{x}) \sum_{j=1}^m \left[h'_j(x_j) \frac{\partial \log p(\mathbf{x})}{\partial x_j} + h_j(x_j) \frac{\partial^2 \log p(\mathbf{x})}{\partial x_j^2} + \frac{1}{2} h_j(x_j) \left(\frac{\partial \log p(\mathbf{x})}{\partial x_j} \right)^2 \right] d\mathbf{x} + C,
\end{aligned}$$

where C is a constant that does not depend on p . □

Proof of Theorem 3. By definition $J_{\mathbf{h}}(p_{\boldsymbol{\theta}}) \geq 0$ and $J_{\mathbf{h}}(p_{\boldsymbol{\theta}_0}) = 0$, so $\boldsymbol{\theta}_0$ minimizes $J_{\mathbf{h}}(p_{\boldsymbol{\theta}})$. Conversely, suppose $J_{\mathbf{h}}(p_{\boldsymbol{\theta}}) = 0$ for some $\boldsymbol{\theta}_1 \in \Theta$. By assumption $p_0(\mathbf{x}) > 0$ almost surely (hereafter a.s.) and $h_j^{1/2}(\mathbf{x}) > 0$ a.s. for all $j = 1, \dots, m$. Therefore, we must have $\nabla \log p_{\boldsymbol{\theta}_1}(\mathbf{x}) = \nabla \log p_0(\mathbf{x})$ a.s., or equivalently, $p_{\boldsymbol{\theta}_1}(\mathbf{x}) = \text{const} \times p_0(\mathbf{x})$ for all almost every $\mathbf{x} \in \mathbb{R}_+^m$. Since $p_{\boldsymbol{\theta}_1}$ and p_0 are both continuous probability density functions, we necessarily have $p_{\boldsymbol{\theta}_1}(\mathbf{x}) = p_0(\mathbf{x})$ for all $\mathbf{x} \in \mathbb{R}_+^m$, which implies $\boldsymbol{\theta}_1 = \boldsymbol{\theta}_0$ by the identifiability assumption. The last claim follows by the law of large numbers, and is an analog of Corollary 3 in Hyvärinen (2005). □

A.3 EXPONENTIAL FAMILIES

Consider the case where $\{p_{\boldsymbol{\theta}} : \boldsymbol{\theta} \in \Theta \subseteq \mathbb{R}^r\}$ contains exponential families with densities

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\theta}^\top \mathbf{t}(\mathbf{x}) - \psi(\boldsymbol{\theta}) + b(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}_+^m.$$

Then the empirical generalized \mathbf{h} -score matching loss becomes

$$\hat{J}_{\mathbf{h}}(p_{\boldsymbol{\theta}}) = \frac{1}{2} \boldsymbol{\theta}^\top \Gamma(\mathbf{x}) \boldsymbol{\theta} - \mathbf{g}(\mathbf{x})^\top \boldsymbol{\theta} + \text{const},$$

where

$$\mathbf{\Gamma}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m h_j(X_j^{(i)}) \mathbf{t}'_j(\mathbf{X}^{(i)}) \mathbf{t}'_j(\mathbf{X}^{(i)})^\top \in \mathbb{R}^{r \times r} \quad \text{and} \quad (\text{A.1})$$

$$\mathbf{g}(\mathbf{x}) = -\frac{1}{n} \sum_{i=1}^n \left[h_j(X_j^{(i)}) b'_j(\mathbf{X}^{(i)}) \mathbf{t}'_j(\mathbf{X}^{(i)}) + h_j(X_j^{(i)}) \mathbf{t}''_j(\mathbf{X}^{(i)}) + h'_j(X_j^{(i)}) \mathbf{t}'_j(\mathbf{X}^{(i)}) \right] \in \mathbb{R}^r. \quad (\text{A.2})$$

Proof of (6). For exponential families, under the assumptions the empirical loss $\hat{J}_{\mathbf{h}}(p_{\boldsymbol{\theta}})$ becomes (up to an additive constant)

$$\begin{aligned} & \hat{J}_{\mathbf{h}}(p_{\boldsymbol{\theta}}) \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m \left[h'_j(X_j^{(i)}) \frac{\partial \log p_{\boldsymbol{\theta}}(\mathbf{X}^{(i)})}{\partial X_j^{(i)}} + h_j(X_j^{(i)}) \frac{\partial^2 \log p_{\boldsymbol{\theta}}(\mathbf{X}^{(i)})}{\partial (X_j^{(i)})^2} + \frac{1}{2} h_j(X_j^{(i)}) \left(\frac{\partial \log p_{\boldsymbol{\theta}}(\mathbf{X}^{(i)})}{\partial X_j^{(i)}} \right)^2 \right] \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m \left[h'_j(X_j^{(i)}) (\boldsymbol{\theta}^\top \mathbf{t}'_j(\mathbf{X}^{(i)}) + b'_j(\mathbf{X}^{(i)})) + h_j(X_j^{(i)}) (\boldsymbol{\theta}^\top \mathbf{t}''_j(\mathbf{X}^{(i)}) + b''_j(\mathbf{X}^{(i)})) \right. \\ & \quad \left. + \frac{1}{2} h_j(X_j^{(i)}) (\boldsymbol{\theta}^\top \mathbf{t}'_j(\mathbf{X}^{(i)}) + b'_j(\mathbf{X}^{(i)}))^2 \right] \\ &= \frac{1}{n} \left\{ \frac{1}{2} \boldsymbol{\theta}^\top \left[\sum_{i=1}^n \sum_{j=1}^m h_j(X_j^{(i)}) \mathbf{t}'_j(\mathbf{X}^{(i)}) \mathbf{t}'_j(\mathbf{X}^{(i)})^\top \right] \boldsymbol{\theta} \right. \\ & \quad \left. + \left[\sum_{i=1}^n h_j(X_j^{(i)}) b'_j(\mathbf{X}^{(i)}) \mathbf{t}'_j(\mathbf{X}^{(i)}) + h_j(X_j^{(i)}) \mathbf{t}''_j(\mathbf{X}^{(i)}) + h'_j(X_j^{(i)}) \mathbf{t}'_j(\mathbf{X}^{(i)}) \right]^\top \boldsymbol{\theta} \right\} + \text{const}, \end{aligned}$$

which is quadratic in $\boldsymbol{\theta}$. Let

$$\mathbf{\Gamma}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m h_j(X_j^{(i)}) \mathbf{t}'_j(\mathbf{X}^{(i)}) \mathbf{t}'_j(\mathbf{X}^{(i)})^\top, \quad (\text{A.3})$$

$$\mathbf{g}(\mathbf{x}) = -\frac{1}{n} \sum_{i=1}^n \left[h_j(X_j^{(i)}) b'_j(\mathbf{X}^{(i)}) \mathbf{t}'_j(\mathbf{X}^{(i)}) + h_j(X_j^{(i)}) \mathbf{t}''_j(\mathbf{X}^{(i)}) + h'_j(X_j^{(i)}) \mathbf{t}'_j(\mathbf{X}^{(i)}) \right]. \quad (\text{A.4})$$

Then we can write $\hat{J}_{\mathbf{h}}(p_{\boldsymbol{\theta}}) = \frac{1}{2} \boldsymbol{\theta}^\top \mathbf{\Gamma}(\mathbf{x}) \boldsymbol{\theta} - \mathbf{g}(\mathbf{x})^\top \boldsymbol{\theta} + \text{const}$. \square

Proof of Theorem 4. Recall that $\hat{J}_{\mathbf{h}}(p_{\boldsymbol{\theta}}) = \frac{1}{2} \boldsymbol{\theta}^\top \mathbf{\Gamma} \boldsymbol{\theta} - \mathbf{g}^\top \boldsymbol{\theta} + \text{const}$. The minimizer of $\hat{J}_{\mathbf{h}}(p_{\boldsymbol{\theta}})$ is thus available in the unique closed form $\hat{\boldsymbol{\theta}} \equiv \mathbf{\Gamma}(\mathbf{x})^{-1} \mathbf{g}(\mathbf{x})$ as long as $\mathbf{\Gamma}$ is invertible (C1). Since $\mathbf{\Gamma}$ and \mathbf{g} are sample averages, by Khinchin's weak law of large numbers we have $\mathbf{\Gamma} \xrightarrow{p} \mathbb{E}_{p_0} \mathbf{\Gamma} \equiv \mathbf{\Gamma}_0$ and $\mathbf{g} \xrightarrow{p} \mathbb{E}_{p_0} \mathbf{g} \equiv \mathbf{g}_0$, where existence of $\mathbf{\Gamma}_0$ and \mathbf{g}_0 is assumed in (C2). Since $J_{\mathbf{h}}(p_{\boldsymbol{\theta}}) = \mathbb{E}[\hat{J}_{\mathbf{h}}(p_{\boldsymbol{\theta}})] = \mathbb{E}[\frac{1}{2} \boldsymbol{\theta}^\top \mathbf{\Gamma}(\mathbf{x}) \boldsymbol{\theta} - \mathbf{g}(\mathbf{x})^\top \boldsymbol{\theta}] = \frac{1}{2} \boldsymbol{\theta}^\top \mathbf{\Gamma}_0 \boldsymbol{\theta} - \mathbf{g}_0^\top \boldsymbol{\theta}$ and we know $\boldsymbol{\theta}_0$ minimizes $J_{\mathbf{h}}(p_{\boldsymbol{\theta}})$ by definition, by first-order condition we must have $\mathbf{\Gamma}_0 \boldsymbol{\theta}_0 = \mathbf{g}_0$. Then by Lindeberg-Lévy central limit theorem (recall that $\mathbf{g}(\mathbf{x})$ and $\mathbf{\Gamma}(\mathbf{x})$ are sample averages)

$$\sqrt{n}(\mathbf{g}(\mathbf{x}) - \mathbf{\Gamma}(\mathbf{x}) \boldsymbol{\theta}_0) \rightarrow_d \mathcal{N}_m(\mathbf{0}, \boldsymbol{\Sigma}_0),$$

where $\boldsymbol{\Sigma}_0 \equiv \mathbb{E}_{p_0}[(\mathbf{\Gamma}(\mathbf{x}) \boldsymbol{\theta}_0 - \mathbf{g}(\mathbf{x}))(\mathbf{\Gamma}(\mathbf{x}) \boldsymbol{\theta}_0 - \mathbf{g}(\mathbf{x}))^\top]$, as long as $\boldsymbol{\Sigma}_0$ exists (C2). Then by Slutsky's theorem,

$$\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \equiv \sqrt{n}(\mathbf{\Gamma}(\mathbf{x})^{-1}(\mathbf{g}(\mathbf{x}) - \mathbf{\Gamma}(\mathbf{x}) \boldsymbol{\theta}_0)) \rightarrow_d \mathcal{N}_r(\mathbf{0}, \mathbf{\Gamma}_0^{-1} \boldsymbol{\Sigma}_0 \mathbf{\Gamma}_0^{-1}),$$

as long as $\mathbf{\Gamma}_0$ is invertible (C2).

For the second half of the theorem, (C2) $\mathbb{E}_{p_0} \mathbf{\Gamma}(\mathbf{x}) < \infty$ and $\mathbb{E}_{p_0} \mathbf{g}(\mathbf{x}) < \infty$ implies $\mathbb{E}_{p_0} |\mathbf{\Gamma}(\mathbf{x})| < \infty$ and $\mathbb{E}_{p_0} |\mathbf{g}(\mathbf{x})| < \infty$, so by strong law of large numbers (and a union bound on at most k^2 null sets)

$$\mathbf{\Gamma}(\mathbf{x}) \rightarrow_{\text{a.s.}} \mathbf{\Gamma}_0, \quad \mathbf{g}(\mathbf{x}) \rightarrow_{\text{a.s.}} \mathbf{g}_0.$$

Then outside a null set,

$$\hat{\boldsymbol{\theta}} \equiv \boldsymbol{\Gamma}(\mathbf{x})^{-1} \mathbf{g}(\mathbf{x}) \rightarrow_{\text{a.s.}} \boldsymbol{\Gamma}_0^{-1} \mathbf{g}_0 = \boldsymbol{\theta}_0.$$

□

Proof for Example 5. For the family of univariate truncated Gaussian distributions with unknown mean parameter μ and known variance parameter σ^2 , we have

$$p_\theta(x) \propto \exp(\theta t(x) + b(x)), \quad \theta \equiv \frac{\mu}{\sigma^2}, \quad t(x) \equiv x, \quad b(x) = -\frac{x^2}{2\sigma^2}.$$

We choose to estimate $\theta \equiv \mu/\sigma^2$. Then by (A.1) and (A.2),

$$\begin{aligned} \hat{\mu}_h &= \sigma^2 \hat{\theta} \equiv \sigma^2 \boldsymbol{\Gamma}(\mathbf{x})^{-1} \mathbf{g}(\mathbf{x}) \\ &= -\sigma^2 \left[\sum_{i=1}^n h(X_i) t'(X_i)^2 \right]^{-1} \left[\sum_{i=1}^n h(X_i) b'(X_i) t'(X_i) + h(X_i) t''(X_i) + h'(X_i) t'(X_i) \right] \\ &= -\sigma^2 \left[\sum_{i=1}^n h(X_i) \right]^{-1} \left[\sum_{i=1}^n -h(X_i) \frac{X_i}{\sigma^2} + h'(X_i) \right]. \end{aligned}$$

By Theorem 4,

$$\sqrt{n}(\hat{\mu}_h - \mu_0) \rightarrow_d \mathcal{N} \left(0, \frac{\sigma^4 \mathbb{E}_0 \left[h(X) \frac{\mu_0 - X}{\sigma^2} + h'(X) \right]^2}{\mathbb{E}_0^2[h(X)]} \right) \sim \mathcal{N} \left(0, \frac{\mathbb{E}_0 \left[h(X)(\mu_0 - X) + \sigma^2 h'(X) \right]^2}{\mathbb{E}_0^2[h(X)]} \right).$$

By integration by parts, (suppressing the dependence of p_{μ_0} on μ_0)

$$\begin{aligned} &\mathbb{E}_0[h(X)h'(X)(X - \mu_0)] \\ &= \int_0^\infty h'(x)h(x)(x - \mu_0)p(x) dx = \int_0^\infty h(x)(x - \mu_0)p(x) dh(x) \\ &= h^2(x)(x - \mu_0)p(x) \Big|_0^\infty - \int h(x) dh(x)(x - \mu_0)p(x) \\ &= - \int h^2(x)p(x) dx - \int h(x)h'(x)(x - \mu_0)p(x) dx + \int h^2(x) \frac{(x - \mu_0)^2}{\sigma^2} p(x) dx, \end{aligned}$$

where the last step follows from the assumptions $\lim_{x \searrow 0^+} h(x) = 0$ and $\lim_{x \nearrow +\infty} h^2(x)(x - \mu_0)p_{\mu_0}(x) = 0$. So

$$\mathbb{E}_0[h(X)h'(X)(X - \mu_0)] = \frac{\mathbb{E}[h^2(X)((X - \mu_0)^2/\sigma^2 - 1)]}{2}. \quad (\text{A.5})$$

The asymptotic variance thus becomes

$$\begin{aligned} &\frac{\mathbb{E}_0 \left[h(X)(\mu_0 - X) + \sigma^2 h'(X) \right]^2}{\mathbb{E}_0^2[h(X)]} \\ &= \frac{\mathbb{E}_0 \left[h^2(X)(X - \mu_0)^2 - 2\sigma^2 h^2(X) \left((X - \mu_0)^2/\sigma^2 - 1 \right) / 2 + \sigma^4 h'^2(X) \right]}{\mathbb{E}_0^2[h(X)]} \\ &= \frac{\mathbb{E}_0[\sigma^2 h^2(X) + \sigma^4 h'^2(X)]}{\mathbb{E}_0^2[h(X)]}. \end{aligned}$$

We note that the Cramér-Rao lower bound is $\frac{\sigma^4}{\text{var}(X - \mu_0)}$, which follows from taking the second derivative of $\log p_{\mu_0}$ with respect to μ_0 . □

Proof for Example 6. For the family of univariate truncated Gaussian distributions with known mean parameter μ and unknown variance parameter $\sigma^2 > 0$, we have

$$p_\theta(x) \propto \exp(\theta t(x) + b(x)), \quad \theta \equiv \frac{1}{\sigma^2}, \quad t(x) \equiv -(x - \mu)^2/2, \quad b(x) = 0.$$

We estimate $\theta \equiv 1/\sigma^2$. By (A.1) and (A.2),

$$\begin{aligned} \hat{\theta} &\equiv \Gamma(\mathbf{x})^{-1} g(\mathbf{x}) \\ &= - \left[\sum_{i=1}^n h(X_i) t'(X_i)^2 \right]^{-1} \left[\sum_{i=1}^n h(X_i) b'(X_i) t'(X_i) + h(X_i) t''(X_i) + h'(X_i) t'(X_i) \right] \\ &= \left[\sum_{i=1}^n h(X_i) (X_i - \mu)^2 \right]^{-1} \left[\sum_{i=1}^n h(X_i) + h'(X_i) (X_i - \mu) \right]. \end{aligned}$$

By Theorem 4, $\sqrt{n}(\hat{\theta} - \theta) \rightarrow_d \mathcal{N}(0, \varsigma^2)$, where

$$\begin{aligned} \varsigma^2 &\equiv \frac{\mathbb{E}_0 [h(X)((X - \mu)^2/\sigma_0^2 - 1) - h'(X)(X - \mu)]^2}{\mathbb{E}_0^2[h(X)(X - \mu)^2]} \\ &= \frac{1}{\mathbb{E}_0^2[h(X)(X - \mu)^2]} \left(\mathbb{E}_0[h^2(X)(X - \mu)^4/\sigma_0^4 - 2h^2(X)(X - \mu)^2/\sigma_0^2 + h^2(X) + h'^2(X)(X - \mu)^2 \right. \\ &\quad \left. - 2h(X)h'(X)(X - \mu)^3/\sigma_0^2 + 2h(X)h'(X)(X - \mu) \right). \end{aligned}$$

By integration by parts, (suppressing the dependence of $p_{\sigma_0^2}$ on σ_0^2)

$$\begin{aligned} &\mathbb{E}_0[h(X)h'(X)(X - \mu)^3] \\ &= \int_0^\infty h'(x)h(x)(x - \mu)^3 p(x) dx = \int_0^\infty h(x)(x - \mu)^3 p(x) dh(x) \\ &= h^2(x)(x - \mu)^3 p(x) \Big|_0^\infty - \int h(x) dh(x)(x - \mu)^3 p(x) \\ &= - \int h(x)h'(x)(x - \mu)^3 p(x) dx - 3 \int h^2(x)(x - \mu)^2 p(x) dx + \int h^2(x) \frac{(x - \mu)^4}{\sigma_0^2} p(x) dx, \end{aligned}$$

where the last step follows from the assumptions $\lim_{x \searrow 0^+} h(x) = 0$ and $\lim_{x \nearrow +\infty} h^2(x)(x - \mu)^3 p_{\sigma_0^2}(x) = 0$. Combining this with (A.5) we get

$$\sqrt{n}(\hat{\theta} - \theta) \rightarrow_d \mathcal{N}(0, \varsigma^2) \sim \mathcal{N}\left(0, \frac{2\mathbb{E}_0[h^2(X)(X - \mu)^2/\sigma_0^2] + \mathbb{E}_0[h'^2(X - \mu)^2]}{\mathbb{E}_0^2[h(X)(X - \mu)^2]}\right),$$

and so by the delta method, for $\hat{\sigma}_k^2 \equiv \hat{\theta}^{-1}$,

$$\sqrt{n}(\hat{\sigma}_h^2 - \sigma_0^2) \rightarrow_d \mathcal{N}\left(0, \frac{2\sigma_0^6 \mathbb{E}_0[h^2(X)(X - \mu)^2] + \sigma_0^8 \mathbb{E}_0[h'^2(X - \mu)^2]}{\mathbb{E}_0^2[h(X)(X - \mu)^2]}\right).$$

We note that the Cramér-Rao lower bound is $\frac{4\sigma_0^8}{\text{var}(X - \mu)^2}$, which follows from taking the second derivative of $\log p_{\sigma_0^2}$ with respect to σ_0^2 . \square

A.4 REGULARIZED GENERALIZED SCORE MATCHING

We first verify assumptions (A1)–(A2) in the case of truncated Gaussian distributions.

Lemma A.2 (Assumptions for truncated Gaussian). *Consider the non-centered truncated Gaussian distribution with density*

$$\log p_0(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_0)^\top \mathbf{K}_0(\mathbf{x} - \boldsymbol{\mu}_0) + \text{const}$$

with unknown positive definite inverse covariance parameter \mathbf{K}_0 and unknown mean parameter $\boldsymbol{\mu}_0$. Then assuming $0 \leq h_j \leq M_j$, $\lim_{x_j \searrow 0^+} h_j(x_j) = 0$ and $|h'_j| \leq M'_j$, assumptions (A1)–(A2) for score matching are satisfied for any proposed parameters $\mathbf{K} \succ \mathbf{0}$ and $\boldsymbol{\mu}$. Taking $\boldsymbol{\mu} \equiv \boldsymbol{\mu}_0 \equiv \mathbf{0}$ the assumptions also hold in the centered setting. Choosing $m = 1$ gives the univariate case.

Proof of Lemma A.2. Consider $p \sim \text{TN}(\boldsymbol{\mu}, \mathbf{K})$, with \mathbf{k}_j the j -th column of \mathbf{K} . Let $M \equiv \max_j M_j$ and $M' \equiv \max_j M'_j$.

(A1) For any fixed $\mathbf{x}_{-j} \in \mathbb{R}_+^{m-1}$ and any $p \in \mathcal{P}_+$ with parameters \mathbf{K} and $\boldsymbol{\mu}$,

$$\begin{aligned} \lim_{x_j \nearrow \infty} h_j(x_j) p_0(\mathbf{x}) \partial_j \log p(\mathbf{x}) &\propto \lim_{x_j \nearrow \infty} h_j(x_j) \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_0)^\top \mathbf{K}_0(\mathbf{x} - \boldsymbol{\mu}_0)\right) \mathbf{k}_j^\top (\mathbf{x} - \boldsymbol{\mu}) \\ &= \lim_{x_j \nearrow \infty} h_j(x_j) \exp\left(C_1 + C_2 x_j - \frac{1}{2} \kappa_{0,jj} x_j^2\right) (C_3 + C_4 x_j) \end{aligned}$$

for some constants C_1, C_2, C_3 , and C_4 depending on \mathbf{x}_{-j} , \mathbf{K}_0 , \mathbf{K} , $\boldsymbol{\mu}_0$ and $\boldsymbol{\mu}$. Since $\kappa_{0,jj} > 0$ and we assumed h_j to be bounded, the limit equals to 0 for all j and \mathbf{x}_{-j} .

Similarly,

$$\begin{aligned} \lim_{x_j \searrow 0^+} h_j(x_j) p_0(\mathbf{x}) \partial_j \log p(\mathbf{x}) &\propto \lim_{x_j \searrow 0^+} h_j(x_j) \exp\left(C_1 + C_2 x_j - \frac{1}{2} \kappa_{0,jj} x_j^2\right) (C_3 + C_4 x_j) \\ &= \exp(C_1) C_3 \lim_{x_j \searrow 0^+} h_j(x_j) = 0 \end{aligned}$$

if and only if we assume $\lim_{x_j \searrow 0^+} h_j(x_j) = 0$.

(A2) For any $p \in \mathcal{P}_+$ with parameters \mathbf{K} and $\boldsymbol{\mu}$,

$$\begin{aligned} \mathbb{E}_{p_0} \|\nabla \log p(\mathbf{X}) \circ \mathbf{h}^{1/2}(\mathbf{X})\|_2^2 &\leq M \mathbb{E}_{p_0} \|\nabla \log p(\mathbf{X})\|_2^2 = M \text{tr}(\mathbb{E}_{p_0} [(\mathbf{K}(\mathbf{X} - \boldsymbol{\mu}))(\mathbf{K}(\mathbf{X} - \boldsymbol{\mu}))^\top]) \\ &= M \text{tr}(\mathbf{K} \mathbb{E}_{p_0} [(\mathbf{X} - \boldsymbol{\mu}_0 + (\boldsymbol{\mu}_0 - \boldsymbol{\mu}))(\mathbf{X} - \boldsymbol{\mu}_0 + (\boldsymbol{\mu}_0 - \boldsymbol{\mu}))^\top] \mathbf{K}^\top) \\ &= M \text{tr}(\mathbf{K}(\mathbf{K}_0^{-1} + (\boldsymbol{\mu}_0 - \boldsymbol{\mu})(\boldsymbol{\mu}_0 - \boldsymbol{\mu})^\top) \mathbf{K}) < +\infty \end{aligned}$$

since M , \mathbf{K} , \mathbf{K}_0 , $\boldsymbol{\mu}$, $\boldsymbol{\mu}_0$ are all finite constants. We also have

$$\begin{aligned} \mathbb{E}_{p_0} \|(\nabla \log p(\mathbf{X}) \circ \mathbf{h}(\mathbf{X}))'\|_1 &= \sum_{i=1}^m \mathbb{E}_{p_0} |h'_i(X_i) \partial_i \log p(\mathbf{X}) + h_i(X_i) \partial_i^2 \log p(\mathbf{X})| \\ &\leq \sum_{i=1}^m \mathbb{E}_{p_0} |h'_i(X_i) \partial_i \log p(\mathbf{X})| + \mathbb{E}_{p_0} |h_i(X_i) \partial_i^2 \log p(\mathbf{X})| \\ &\leq \sum_{i=1}^m M' \mathbb{E}_{p_0} |\mathbf{k}_i^\top (\mathbf{X} - \boldsymbol{\mu})| + M \kappa_{ij} \\ &\leq \sum_{i=1}^m M' |\mathbf{k}_i|^\top \mathbb{E}_{p_0} \mathbf{X} + M' |\mathbf{k}_i^\top \boldsymbol{\mu}| + M \text{tr}(\mathbf{K}) < +\infty. \end{aligned}$$

Hence, (A1) and (A2) are both satisfied. \square

Our analysis of the regularized generalized \mathbf{h} -score matching estimator follows the proof for the following theorem from Lin et al. (2016), restated below. In our definition and implementation we choose to optimize over all symmetric matrices, but we adopt the following theorem in whose proof the symmetry condition is not explicitly imposed, in order to decouple the columns of \mathbf{K} and to highlight the scaling.

Theorem A.3 (Analog of Theorem 1 from Lin et al. (2016)). Recall that $S_0 \equiv S(\mathbf{K}_0) \equiv \{(i, j) : \kappa_{0,ij} \neq 0\}$. Suppose $\mathbf{\Gamma}_{0,S_0,S_0}$ is invertible and satisfies the irrepresentability condition (10) with incoherence parameter $\alpha \in (0, 1]$. Assume that

$$\|\mathbf{\Gamma}(\mathbf{x}) - \mathbf{\Gamma}_0\|_\infty < \epsilon_1, \quad \|\mathbf{g}(\mathbf{x}) - \mathbf{g}_0\|_\infty < \epsilon_2, \quad (\text{A.6})$$

with $d_{\mathbf{K}_0}\epsilon_1 \leq \alpha/(6c_{\mathbf{\Gamma}_0})$. If

$$\lambda > \frac{3(2-\alpha)}{\alpha} \max\{c_{\mathbf{K}_0}\epsilon_1, \epsilon_2\},$$

then the following statements hold:

- (a) The regularized generalized \mathbf{h} -score matching estimator $\hat{\mathbf{K}}$ in (9) is unique, with support $\hat{S} \equiv S(\hat{\mathbf{K}}) \subseteq S_0$, and satisfies

$$\|\hat{\mathbf{K}} - \mathbf{K}_0\|_\infty \leq \frac{c_{\mathbf{\Gamma}_0}}{2-\alpha} \lambda.$$

- (b) If

$$\min_{1 \leq j < k \leq m} |\mathbf{K}_{0,jk}| > \frac{c_{\mathbf{\Gamma}_0}}{2-\alpha} \lambda,$$

then $\hat{S} = S_0$ and $\text{sign}(\hat{\mathbf{K}}_{jk}) = \text{sign}(\mathbf{K}_{0,jk})$ for all $(j, k) \in S_0$.

This is a deterministic result, and the improvement of our generalized estimator over the one in Lin et al. (2016) is in its asymptotic guarantees, as in Theorem 10. We present a corollary to this theorem, as seen in the second and third inequalities in Theorem 10 (a).

Corollary A.1. Suppose the same assumptions under Theorem A.3 hold. Then $\hat{\mathbf{K}}$ satisfies

$$\begin{aligned} \|\hat{\mathbf{K}} - \mathbf{K}_0\|_F &\leq \frac{c_{\mathbf{\Gamma}_0}}{2-\alpha} \lambda \sqrt{|S_0|} \leq \frac{c_{\mathbf{\Gamma}_0}}{2-\alpha} \lambda \sqrt{d_{\mathbf{K}_0} m}, \\ \|\hat{\mathbf{K}} - \mathbf{K}_0\|_2 &\leq \frac{c_{\mathbf{\Gamma}_0}}{2-\alpha} \lambda \min(\sqrt{|S_0|}, d_{\mathbf{K}_0}). \end{aligned}$$

Proof of Corollary A.1. By Theorem A.3, under assumptions in that theorem, the support of $\hat{\mathbf{K}}$ is a subset of the true support of \mathbf{K}_0 , and $\|\hat{\mathbf{K}} - \mathbf{K}_0\|_\infty \leq \frac{c_{\mathbf{\Gamma}_0}}{2-\alpha} \lambda$. Since \mathbf{K}_0 has $|S_0|$ nonzero entries,

$$\|\hat{\mathbf{K}} - \mathbf{K}_0\|_F = \left[\sum_{\mathbf{K}_{0,jk} \neq 0} (\hat{\mathbf{K}}_{jk} - \mathbf{K}_{0,jk})^2 \right]^{1/2} \leq \sqrt{|S_0|} \|\hat{\mathbf{K}} - \mathbf{K}_0\|_\infty \leq \frac{c_{\mathbf{\Gamma}_0}}{2-\alpha} \lambda \sqrt{|S_0|}.$$

Similarly, by the definition of matrix ℓ_∞ - ℓ_∞ norm,

$$\|\hat{\mathbf{K}} - \mathbf{K}_0\|_2 \leq \|\hat{\mathbf{K}} - \mathbf{K}_0\|_\infty = \max_{j=1, \dots, m} \sum_{k=1}^m |\hat{\mathbf{K}}_{jk} - \mathbf{K}_{0,jk}| \leq \frac{c_{\mathbf{\Gamma}_0}}{2-\alpha} \lambda d_{\mathbf{K}_0}.$$

The result follows by also noting that $\|\hat{\mathbf{K}} - \mathbf{K}_0\|_2 \leq \|\hat{\mathbf{K}} - \mathbf{K}_0\|_F$. \square

Proof of Theorem 10. By Theorem A.3 it suffices to prove that for any $\tau > 3$, we can bound $\|\mathbf{\Gamma}(\mathbf{x}) - \mathbf{\Gamma}_0\|_\infty$ by some ϵ_1 and $\|\mathbf{g}(\mathbf{x}) - \mathbf{g}_0\|_\infty$ by some ϵ_2 , uniformly with probability $1 - m^{3-\tau}$. Recall from Section 4.2 that the j^{th} block of $\mathbf{\Gamma} \in \mathbb{R}^{m^2 \times m^2}$ has (k, ℓ) -th entry

$$\frac{1}{n} \sum_{i=1}^n X_k^{(i)} X_\ell^{(i)} h_j(X_j^{(i)}),$$

and the entry in $\mathbf{g} \in \mathbb{R}^{m^2}$ (obtained by linearizing a $m \times m$ matrix) corresponding to (j, k) with $j \neq k$, is

$$\frac{1}{n} \sum_{i=1}^n X_k^{(i)} h'_j(X_j^{(i)}),$$

while the entry for (j, j) is

$$\frac{1}{n} \sum_{i=1}^n X_j^{(i)} h'_j(X_j^{(i)}) + \frac{1}{n} \sum_{i=1}^n h_j(X_j^{(i)}).$$

Denote $M \equiv \max_j \sup_{x>0} h_j(x)$ and $M' \equiv \max_j \sup_{x>0} h'_j(x)$, and let $c_{\mathbf{X}} \equiv 2 \max_j (2\sqrt{\Sigma_{jj}} + \sqrt{e}\mathbb{E}_0 X_j)$. Using results for sub-gaussian random variables from Lemma A.6 below and Hoeffding's inequality, we have for any $t_1, t_{2,1}, t_{2,2} > 0$,

$$\begin{aligned} \mathbb{P} \left(\left| \frac{1}{n} \sum_{i=1}^n X_k^{(i)} X_\ell^{(i)} h_j(X_j^{(i)}) - \mathbb{E}_0 X_k X_\ell h_j(X_j) \right| > t_1 \right) &\leq 2 \exp \left(- \min \left(\frac{nt_1^2}{2M^2 c_{\mathbf{X}}^4}, \frac{nt_1}{2M c_{\mathbf{X}}^2} \right) \right), \\ \mathbb{P} \left(\left| \frac{1}{n} \sum_{i=1}^n X_k^{(i)} h'_j(X_j^{(i)}) - \mathbb{E}_0 X_k h'_j(X_j) \right| \geq t_{2,1} \right) &\leq 2 \exp \left(- \frac{nt_{2,1}^2}{2M'^2 c_{\mathbf{X}}^2} \right), \\ \mathbb{P} \left(\left| \frac{1}{n} \sum_{i=1}^n h_j(X_j^{(i)}) - \mathbb{E}_0 h_j(X_j) \right| \geq t_{2,2} \right) &\leq 2 \exp \left(-2nt_{2,2}^2/M^2 \right). \end{aligned}$$

Choosing

$$\begin{aligned} \epsilon_1 &\equiv M c_{\mathbf{X}}^2 \max \left\{ \frac{2(\log m^\tau + \log 6)}{n}, \sqrt{\frac{2(\log m^\tau + \log 6)}{n}} \right\}, \\ \epsilon_{2,1} &\equiv \sqrt{2} M' c_{\mathbf{X}} \sqrt{\frac{\log m^{\tau-1} + \log 6}{n}}, \quad \epsilon_{2,2} \equiv M \sqrt{\frac{\log m^{\tau-2} + \log 6}{2n}}, \end{aligned}$$

and taking union bounds over m^3 , m^2 , and m events, respectively, we have

$$\begin{aligned} \mathbb{P} \left(\sup_{j,k,\ell} \left| \frac{1}{n} \sum_{i=1}^n X_k^{(i)} X_\ell^{(i)} h_j(X_j^{(i)}) - \mathbb{E}_0 X_k X_\ell h_j(X_j) \right| \geq \epsilon_1 \right) &\leq \frac{1}{3m^{\tau-3}}, \\ \mathbb{P} \left(\sup_{j,k} \left| \frac{1}{n} \sum_{i=1}^n X_k^{(i)} h'_j(X_j^{(i)}) - \mathbb{E}_0 X_k h'_j(X_j) \right| \geq \epsilon_{2,1} \right) &\leq \frac{1}{3m^{\tau-3}}, \\ \mathbb{P} \left(\sup_j \left| \frac{1}{n} \sum_{i=1}^n h_j(X_j^{(i)}) - \mathbb{E}_0 h_j(X_j) \right| \geq \epsilon_{2,2} \right) &\leq \frac{1}{3m^{\tau-3}}. \end{aligned}$$

Hence, with probability at least $1 - m^{3-\tau}$, $\|\Gamma(\mathbf{x}) - \Gamma_0\|_\infty < \epsilon_1$ and $\|\mathbf{g}(\mathbf{x}) - \mathbf{g}_0\|_\infty < \epsilon_2 \equiv \epsilon_{2,1} + \epsilon_{2,2}$. Consider any $\tau > 3$, and let

$$\begin{aligned} c_2 &\equiv \frac{6}{\alpha} c_{\Gamma_0}, \quad n \geq \max \{ 2M^2 c_{\mathbf{X}}^4 c_{\mathbf{K}_0}^2 d_{\mathbf{K}_0}^2 (\tau \log m + \log 6), 2M c_{\mathbf{X}}^2 c_2 d_{\mathbf{K}_0} (\tau \log m + \log 6) \}, \\ \lambda &> \frac{3(2-\alpha)}{\alpha} \max \{ c_{\mathbf{K}_0} \epsilon_1, \epsilon_2 \} \\ &\equiv \frac{3(2-\alpha)}{\alpha} \max \left\{ M c_{\mathbf{K}_0} c_{\mathbf{X}}^2 \frac{2(\log m^\tau + \log 6)}{n}, \right. \\ &\quad \left. M c_{\mathbf{K}_0} c_{\mathbf{X}}^2 \sqrt{\frac{2(\log m^\tau + \log 6)}{n}}, \sqrt{2} M' c_{\mathbf{X}} \sqrt{\frac{\log m^{\tau-1} + \log 6}{n}} + M \sqrt{\frac{\log m^{\tau-2} + \log 6}{2n}} \right\}. \end{aligned}$$

Then $d_{\mathbf{K}_0} \epsilon_1 \leq \alpha/(6c_{\Gamma_0})$ and the results follow from Theorem A.3. \square

We now present the definition of sub-Gaussian and sub-exponential norms and variables as well as lemmas required for the proof above.

Definition A.4 (Sub-Gaussian and Sub-Exponential Variables). The *sub-gaussian* ($r = 2$) and *sub-exponential* ($r = 1$) norms of a random variable are defined as

$$\|X\|_{\psi_r} \equiv \sup_{q \geq 1} q^{-1/r} (\mathbb{E}|X|^{rq})^{1/(rq)} \equiv \sup_{q \geq 1} q^{-1/r} \|X\|_{rq}.$$

If $\|X\|_{\psi_2} < \infty$ we say X is *sub-gaussian*; if $\|X\|_{\psi_1} < \infty$ we call X *sub-exponential*. For a *zero-mean* sub-gaussian random variable X also define the *sub-gaussian parameter*

$$\tau(X) = \inf\{\tau \geq 0 : \mathbb{E} \exp(tX) \leq \exp(\tau^2 t^2 / 2), \forall t \in \mathbb{R}\}.$$

Note that the definition of sub-gaussian norm here allows the variable to be non-centered, and is different from the one in Vershynin (2010), which uses $\|X\|_q$ in the definition. Instead, it coincides with θ_2 in Buldygin and Kozachenko (2000). The definition of the sub-gaussian parameter is the same as in Buldygin and Kozachenko (2000), and the definition of the sub-exponential norm is as in Vershynin (2010).

Lemma A.5 (Properties of Sub-Gaussian and Sub-Exponential Variables). *Then*

- 1) For any X and $r = 1, 2$, $\|X - \mathbb{E}X\|_{\psi_r} \leq 2\|X\|_{\psi_r}$ and $\|X\|_{\psi_r} \leq \|X - \mathbb{E}X\|_{\psi_r} + |\mathbb{E}X|$, as long as the expectation and norms are finite.
- 2) (Buldygin and Kozachenko, 2000) $\tau(X)$ is a norm on the space of all zero-mean sub-gaussian variables; in particular, $\tau(X + Y) \leq \tau(X) + \tau(Y)$ as long as the quantities are defined and finite. If X is zero-mean sub-gaussian, then $\text{var}(X) \leq \tau^2(X)$, $\|X\|_{\psi_2} \leq 2\tau(X)/\sqrt{e}$, $\tau(X) \leq \sqrt{e}\|X\|_{\psi_2}$. If X_1, \dots, X_n are i.i.d. zero-mean sub-gaussian, $\tau\left(\frac{1}{n} \sum_{i=1}^n X_i\right) \leq \frac{1}{\sqrt{n}}\tau(X_i)$.
- 3) If random variables X_1 and X_2 (not necessarily independent) are sub-gaussian with $\|X_1\|_{\psi_2} \leq K_1$ and $\|X_2\|_{\psi_2} \leq K_2$, then $X_1 X_2$ is sub-exponential with $\|X_1 X_2\|_{\psi_1} \leq K_1 K_2$.
- 4) (Buldygin and Kozachenko, 2000) If X is zero-mean sub-gaussian,

$$\mathbb{E}|X|^q \leq 2(q/e)^{q/2} \tau^q(X)$$

for any $q > 0$.

- 5) (Buldygin and Kozachenko, 2000) If X_1, \dots, X_n are independent zero-mean sub-gaussian variables, then for any $\epsilon > 0$,

$$\begin{aligned} \mathbb{P}(|X_1| \geq \epsilon) &\leq 2 \exp\left(-\frac{\epsilon^2}{2\tau^2(X_1)}\right), \\ \mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n X_i\right| > \epsilon\right) &\leq 2 \exp\left(-\frac{n\epsilon^2}{2 \max_i \tau^2(X_i)}\right). \end{aligned}$$

- 6) (Vershynin, 2010) If X_1, \dots, X_n are independent zero-mean sub-exponential random variables with $K \geq \max_i \|X_i\|_{\psi_1}$, then for any $\epsilon > 0$,

$$\begin{aligned} \mathbb{P}(|X_1| \geq \epsilon) &\leq 2 \exp\left(-\min\left(\frac{\epsilon^2}{8e^2 K^2}, \frac{\epsilon}{4eK}\right)\right), \\ \mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n X_i\right| \geq \epsilon\right) &\leq 2 \exp\left(-\min\left(\frac{n\epsilon^2}{8e^2 K^2}, \frac{n\epsilon}{4eK}\right)\right). \end{aligned}$$

Proof. 1) For $r = 1, 2$, by triangle inequality, $\|X - \mathbb{E}X\|_{\psi_r} \leq \|X\|_{\psi_r} + \|\mathbb{E}X\|_{\psi_r} = \|X\|_{\psi_r} + |\mathbb{E}X| \leq \|X\|_{\psi_r} + \mathbb{E}|X| \leq 2\|X\|_{\psi_r}$, where in the last step we used the definition of $\|\cdot\|_{\psi_r}$ with $q = 1$ for $r = 1$ and $\mathbb{E}|X| \leq (\mathbb{E}|X|^2)^{1/2}$ with $q = 2$ for $r = 2$. On the other hand, $\|X\|_{\psi_r} \leq \|X - \mathbb{E}X\|_{\psi_r} + \|\mathbb{E}X\|_{\psi_r} = \|X - \mathbb{E}X\|_{\psi_r} + |\mathbb{E}X|$.

- 2) These follow from Theorems 1.2 and 1.3 and Lemmas 1.2 and 1.7 from Buldygin and Kozachenko (2000), and $\sqrt[4]{3.1}e^{9/16}/\sqrt{2} \approx 1.6467 \leq 1.6487 \approx \sqrt{e}$.
- 3) By Hölder's inequality (or Cauchy-Schwarz),

$$\begin{aligned} \|X_1 X_2\|_{\psi_1} &= \sup_{q \geq 1} q^{-1} (\mathbb{E}|X_1 X_2|^q)^{1/q} = \sup_{q \geq 1} q^{-1} (\mathbb{E}|X_1^q X_2^q|)^{1/q} \\ &\leq \sup_{q \geq 1} q^{-1} \left[(\mathbb{E}|X_1|^{2q})^{1/2} (\mathbb{E}|X_2|^{2q})^{1/2} \right]^{1/q} \\ &\leq \sup_{q \geq 1} \left[q^{-1/2} (\mathbb{E}|X_1|^{2q})^{1/2q} \right] \sup_{q \geq 1} \left[q^{-1/2} (\mathbb{E}|X_2|^{2q})^{1/2q} \right] \\ &= \|X_1\|_{\psi_2} \|X_2\|_{\psi_2} \leq K_1 K_2. \end{aligned}$$

4) This is Lemma 1.4 from Buldygin and Kozachenko (2000).

5) This is Theorem 1.5 from Buldygin and Kozachenko (2000).

6) This follows from Corollary 5.17 from Vershynin (2010). \square

Lemma A.6. *Suppose \mathbf{X} follows a truncated normal distribution on \mathbb{R}_+^m with parameters $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma} = \mathbf{K}^{-1} \succ \mathbf{0}$. Let $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(n)}$ be i.i.d. copies of \mathbf{X} , with j -th component of the i -th copy being $X_j^{(i)}$. Then*

1. For $j = 1, \dots, p$, $\tau(X_j - \mathbb{E}X_j) \leq \sqrt{\Sigma_{jj}}$. That is, the sub-gaussian parameter of any marginal distribution of \mathbf{X} , after centering, is bounded by the square root of its corresponding diagonal entry in the covariance parameter $\boldsymbol{\Sigma}$. Then for any $\epsilon > 0$,

$$\mathbb{P}\left(\left|\frac{1}{n}\sum_{i=1}^n X_j^{(i)} - \mathbb{E}X_j\right| > \epsilon\right) \leq 2 \exp\left(-\frac{n\epsilon^2}{2\Sigma_{jj}}\right).$$

In particular, if h_0 is a function bounded by M_0 , then for any $\epsilon > 0$,

$$\begin{aligned} \mathbb{P}\left(\left|\frac{1}{n}\sum_{i=1}^n X_j^{(i)} h_0(X_k^{(i)}) - \mathbb{E}X_j h_0(X_k)\right| \geq \epsilon\right) &\leq 2 \exp\left(-\frac{n\epsilon^2}{8M_0^2(2\sqrt{\Sigma_{jj}} + \sqrt{\epsilon}\mathbb{E}X_j)^2}\right), \\ \tau\left(\frac{1}{n}\sum_{i=1}^n X_j^{(i)} h_0(X_k^{(i)}) - \mathbb{E}X_j h_0(X_k)\right) &\leq \frac{2M_0}{\sqrt{n}}(2\sqrt{\Sigma_{jj}} + \sqrt{\epsilon}\mathbb{E}X_j), \\ \left\|\frac{1}{n}\sum_{i=1}^n X_j^{(i)} h_0(X_k^{(i)}) - \mathbb{E}X_j h_0(X_k)\right\|_{\psi_2} &\leq \frac{4M_0}{\sqrt{\epsilon n}}(2\sqrt{\Sigma_{jj}} + \sqrt{\epsilon}\mathbb{E}X_j). \end{aligned}$$

2. For $j, k, \ell \in \{1, \dots, p\}$, if h_0 is a function bounded by M_0 , then with $c_{\mathbf{X}} \equiv 2 \max_j(2\sqrt{\Sigma_{jj}} + \sqrt{\epsilon}\mathbb{E}X_j)$,

$$\|X_j X_k h_0(X_\ell) - \mathbb{E}X_j X_k h_0(X_\ell)\|_{\psi_1} \leq \frac{M_0}{2e} c_{\mathbf{X}}^2. \quad (\text{A.7})$$

In particular, for any $\epsilon > 0$,

$$\mathbb{P}\left(\left|\frac{1}{n}\sum_{i=1}^n X_j^{(i)} X_k^{(i)} h_0(X_\ell^{(i)}) - \mathbb{E}X_j X_k h_0(X_\ell)\right| > \epsilon\right) \leq 2 \exp\left(-\min\left(\frac{n\epsilon^2}{2M_0^2 c_{\mathbf{X}}^4}, \frac{n\epsilon}{2M_0 c_{\mathbf{X}}^2}\right)\right).$$

Proof of Lemma A.6. 1. Without loss of generality choose $j = 1$. By the definition of sub-gaussian parameters, we need to show that for all $t \in \mathbb{R}$,

$$\mathbb{E} \exp(tX_1) \leq \exp(t^2 \Sigma_{11}/2 + t\mathbb{E}X_1),$$

which is equivalent to

$$t^2 \Sigma_{11}/2 + t\mathbb{E}X_1 - \log \mathbb{E} \exp(tX_1) \geq 0 \quad \forall t \in \mathbb{R}. \quad (\text{A.8})$$

Since the left-hand side of (A.8) equals 0 at $t = 0$, it suffices to show that its derivative

$$t\Sigma_{11} + \mathbb{E}X_1 - \frac{d \log \mathbb{E} \exp(tX_1)}{dt} = t\Sigma_{11} + \mathbb{E}X_1 - \frac{\frac{d\mathbb{E} \exp(tX_1)}{dt}}{\mathbb{E} \exp(tX_1)} \quad (\text{A.9})$$

is non-negative on $(0, \infty)$ and non-positive on $(-\infty, 0)$. By properties of moment-generating functions, $\frac{d\mathbb{E} \exp(tX_1)}{dt}$ evaluated at $t = 0$ equals $\mathbb{E}X_1$, so (A.9) equals 0 at $t = 0$. It in turn suffices to show the derivative of (A.9), namely

$$\Sigma_{11} - \frac{d^2 \log \mathbb{E} \exp(tX_1)}{dt^2} \quad (\text{A.10})$$

is non-negative in $t \in \mathbb{R}$.

By Tallis (1961), denoting the first column of Σ as Σ_1 , the moment-generating function of the marginal distribution of X_1 is

$$\frac{\int_{\mathbb{R}_+^p - \mu - t\Sigma_1} \exp\left(-\frac{1}{2}\mathbf{x}^\top \Sigma^{-1}\mathbf{x}\right) d\mathbf{x}}{\int_{\mathbb{R}_+^p - \mu} \exp\left(-\frac{1}{2}\mathbf{x}^\top \Sigma^{-1}\mathbf{x}\right) d\mathbf{x}} \exp\left(t\mu_1 + \frac{1}{2}t^2\Sigma_{11}^2\right).$$

(A.10) thus becomes

$$-\frac{d^2}{dt^2} \log \int_{\mathbb{R}_+^p - \mu - t\Sigma_1} \exp\left(-\frac{1}{2}\mathbf{x}^\top \Sigma^{-1}\mathbf{x}\right) d\mathbf{x}.$$

Showing this is non-negative in $t \in \mathbb{R}$ is equivalent to showing that the integral itself is log-concave in t . But

$$\int_{\mathbb{R}_+^p - \mu - t\Sigma_1} \exp\left(-\frac{1}{2}\mathbf{x}^\top \Sigma^{-1}\mathbf{x}\right) d\mathbf{x} = \int_{\mathbb{R}^p} \exp\left(-\frac{1}{2}\mathbf{x}^\top \Sigma^{-1}\mathbf{x}\right) \mathbf{1}_{\mathbb{R}_+^p - \mu}(\mathbf{x} + t\Sigma_1) d\mathbf{x}$$

with $\exp\left(-\frac{1}{2}\mathbf{x}^\top \Sigma^{-1}\mathbf{x}\right)$ log-concave in \mathbf{x} and $\mathbf{1}_{\mathbb{R}_+^p - \mu}(\mathbf{x} + t\Sigma_1)$ log-concave in (\mathbf{x}, t) since $\mathbb{R}_+^p - \mu$ is a convex set (half-space). Here $\mathbf{1}_S(\cdot)$ is the indicator function of a set S . Since log-concavity is closed under multiplication and integration over \mathbb{R}^p , the integral is indeed log-concave, and our proof of the bound on the sub-gaussian parameter of $X_j - \mathbb{E}X_j$ is complete. The tail bound follows from 5) of Lemma A.5.

Now by 1) and 2) of Lemma A.5,

$$\|X_j\|_{\psi_2} \leq 2\sqrt{\Sigma_{jj}/e} + \mathbb{E}X_j.$$

If h_0 is a function bounded by M_0 , then by definition

$$\|X_j h_0(X_k)\|_{\psi_2} \leq M_0 \left(2\sqrt{\Sigma_{jj}/e} + \mathbb{E}X_j\right).$$

By 1) and 2) of Lemma A.5 again,

$$\begin{aligned} \tau(X_j h_0(X_k) - \mathbb{E}X_j h_0(X_k)) &\leq \sqrt{e} \|X_j h_0(X_k) - \mathbb{E}X_j h_0(X_k)\|_{\psi_2} \\ &\leq 2\sqrt{e} \|X_j h_0(X_k)\|_{\psi_2} \\ &\leq 2M_0(2\sqrt{\Sigma_{jj}} + \sqrt{e}\mathbb{E}X_j). \end{aligned}$$

The tail bound thus follows from the first inequality using 5) of Lemma A.5. By 2),

$$\begin{aligned} \tau\left(\frac{1}{n} \sum_{i=1}^n X_j^{(i)} h_0(X_k^{(i)}) - \mathbb{E}X_j h_0(X_k)\right) &\leq \frac{2M_0}{\sqrt{n}} (2\sqrt{\Sigma_{jj}} + \sqrt{e}\mathbb{E}X_j), \\ \left\| \frac{1}{n} \sum_{i=1}^n X_j^{(i)} h_0(X_k^{(i)}) - \mathbb{E}X_j h_0(X_k) \right\|_{\psi_2} &\leq \frac{4M_0}{\sqrt{en}} (2\sqrt{\Sigma_{jj}} + \sqrt{e}\mathbb{E}X_j). \end{aligned}$$

2. By the proof of 1) of this lemma, $\|X_j\|_{\psi_2} \leq 2\sqrt{\Sigma_{jj}/e} + \mathbb{E}X_j$, and by 3) of Lemma A.5,

$$\|X_j X_k\|_{\psi_1} \leq (2\sqrt{\Sigma_{jj}/e} + \mathbb{E}X_j)(2\sqrt{\Sigma_{kk}/e} + \mathbb{E}X_k) \leq \max_j \left(2\sqrt{\Sigma_{jj}/e} + \mathbb{E}X_j\right)^2.$$

Since h_0 is a function bounded by M_0 , by definition

$$\|X_j X_k h_0(X_\ell)\|_{\psi_1} \leq M_0 \max_j \left(2\sqrt{\Sigma_{jj}/e} + \mathbb{E}X_j\right)^2.$$

Then by 1) of Lemma A.5 again,

$$\|X_j X_k h_0(X_\ell) - \mathbb{E}X_j X_k h_0(X_\ell)\|_{\psi_1} \leq 2M_0 \max_j \left(2\sqrt{\Sigma_{jj}/e} + \mathbb{E}X_j\right)^2.$$

The tail bound then follows from 6) of Lemma A.5.

□