

Note: Statistical Inference

Oct 2024

Lecturer:

Typed by: Zhuohua Shen

Contents

1	Preliminary	5
2	Statistical inference fundamentals	7
2.1	Statistical Models	7
2.2	Principles of Data Reduction	8
2.2.1	Sufficiency Principle	8
2.2.2	Likelihood principle	8
3	Multivariate Inference Fundamentals	9
3.1	Random vectors and distributions	9
3.1.1	Multivariate normal distribution	10
3.1.2	The noncentral χ^2 and F distribution	11
3.2	Asymptotic properties	11
3.2.1	Asymptotic distributions of sample means and covariance matrices	11
4	Bayesian Inference	13
5	Structural Equation Model (SEM)	15
5.1	SEM models	15
5.2	Bayesian methods for estimating SEM	17
5.2.1	Bayesian estimation using MCMC	17

Chapter 1

Preliminary

Chapter 2

Statistical inference fundamentals

References: most of the contents are from the undergraduate course STA3020 (by Prof. Jianfeng Mao in 2022-2023 T1, and Prof. Jiasheng Shi in 2023-2024 T2) and postgraduate course STAT5010 (by Kin Wai Keith Chan in 2024-2025 T1), with main textbook Casella and Berger [1]

2.1 Statistical Models

See Chapter 3 of [1]. Suppose $X_i \sim_{\text{iid}} \mathbb{P}_*$, where \mathbb{P}_* refers to the unknown **data generating process** (DGPg), we find $\hat{\mathbb{P}} \approx \mathbb{P}_*$. A **statistical model** is a set of distributions $\mathcal{F} = \{\mathbb{P}_\theta : \theta \in \Theta\}$, where Θ is the **parameter space**. A **parametric model** is the model with $\dim(\Theta) < \infty$, while a **nonparametric model** satisfies $\dim(\Theta) = \infty$.

Definition 2.1.1 (Exponential family). A k -dimensional **exponential family** (EF) $\mathcal{F} = \{f_\theta : \theta \in \Theta\}$ is a model consisting of pdfs of the form

$$f_\theta(x) = c(\theta)h(x) \exp \left\{ \sum_{j=1}^k \eta_j(\theta)T_j(x) \right\} \quad (2.1)$$

where $c(\theta), h(x) \geq 0$, $\Theta = \{\theta : c(\theta) \geq 0, \eta_j(\theta) \text{ being well defined for } 1 \leq j \leq k\}$. Let $\eta_j = \eta_j(\theta)$, the **canonical form** is

$$f_\eta(x) = b(\eta)h(x) \exp \left\{ \sum_{j=1}^k \eta_j T_j(x) \right\}, \quad (2.2)$$

- k -dim **natural exponential family** (NEF): $\mathcal{F}' = \{f_\eta : \eta \in \Xi\}$;
- **natural parameter** $\eta = (\eta_1, \dots, \eta_k)^\top$;
- **natural parameter space**: $\Xi = \{\eta \in \mathbb{R}^k : 0 < b(\eta) < \infty\}$;
- the NEF \mathcal{F}' is of **full rank** if Ξ contains an open set in \mathbb{R}^k ;
- the EF is a **curved exponential family** if $p = \dim(\Theta) < k$.

Properties of EF:

- Let $X \sim f_\eta$, where $\eta \in \Xi$ such that (i) f_η is of the form (2.2) with $B(\eta) = -\log b(\eta)$, and (ii) Ξ contains an open set in \mathbb{R}^k . Then, for $j, j' = 1, \dots, k$, $\mathbb{E}\{T_j(X)\} = \partial B(\eta)/\partial \eta_j$ and $\text{Cov}\{T_j(X), T_{j'}(X)\} = \partial^2 B(\eta)/(\partial \eta_j \partial \eta_{j'})$.
- **Stein's identity**:

Definition 2.1.2 (Location-scale family). Let f be a density.

- A **location-scale family** is given by $\mathcal{F} = \{f_{\mu, \sigma} : \mu \in \mathbb{R}, \sigma \in \mathbb{R}^{++}\}$, where $f_{\mu, \sigma}(x) = f((x - \mu)/\sigma)/\sigma$.
- **location parameter**: μ ; **scale parameter**: σ ; **standard density**: f ;
- A **location family** is $\mathcal{F} = \{f_{\mu, 1} : \mu \in \mathbb{R}\}$.
- A **scale family** is $\mathcal{F} = \{f_{0, \sigma} : \sigma \in \mathbb{R}^{++}\}$

Representation: $X = \mu + \sigma Z$, $Z \sim f_{0,1}(\cdot)$.

- See some examples in Example 3.9, Keith's note 3, and Table 1 in Shi's note L1.
- Transform between location parameter and scale parameter by taking log.

Definition 2.1.3 (Identifiable family). If $\forall \theta_1, \theta_2 \in \Theta$ that

$$\theta_1 \neq \theta_2 \Rightarrow f_{\theta_1}(\cdot) \neq f_{\theta_2}(\cdot),$$

then \mathcal{F} is said to be an **identifiable family**, or equivalently $\theta \in \Theta$ is **identifiable**.

- $p < k$, curved (must).
- $p = k$, of full rank.
- $p > k$, non-identifiable.

2.2 Principles of Data Reduction

Statistics: $T = T(X_{1:n})$, a function of $X_{1:n}$ and free of any unknown parameter.

2.2.1 Sufficiency Principle

Sufficiency principle: If $T = T(X_{1:n})$ is a “sufficient statistics” for θ , then any inference on θ will depend on $X_{1:n}$ only through T .

Definition 2.2.1 (Sufficient, minimal sufficient, ancillary, and complete statistics). Suppose $X_{1:n} \sim \text{iid} \mathbb{P}_\theta$, where $\theta \in \Theta$. Let $T = T(X_{1:n})$ be a statistic. Then T is **sufficient** (SS) for θ

\Leftrightarrow (def) $[X_{1:n} \mid T = t]$ is free of θ for each t .

\Leftrightarrow (technical lemma) $T(x_{1:n}) = T(x'_{1:n})$ implies that $f_\theta(x_{1:n})/f_\theta(x'_{1:n})$ is free of θ .

\Leftrightarrow (Neyman-Fisher factorization theorem) $\forall \theta \in \Theta, x_{1:n} \in \mathcal{X}^n, f_\theta(x_{1:n}) = A(t, \theta)B(x_{1:n})$.

\Leftrightarrow Define $\Lambda(\theta', \theta'' \mid x_{1:n}) := f_{\theta'}(x_{1:n})/f_{\theta''}(x_{1:n})$. $\forall \theta', \theta'' \in \Theta, \exists$ function $C_{\theta', \theta''}$ such that $\Lambda(\theta', \theta'' \mid x_{1:n}) = C_{\theta', \theta''}(t)$, for all $x_{1:n} \in \mathcal{X}^n$ where $t = T(x_{1:n})$.

T is **minimal sufficient** (MSS) for θ

\Leftrightarrow (def) (1) T is a SS for θ ; (2) $T = g(S)$ for any other SS S .

\Leftrightarrow (1) T is a SS for θ ; (2) $S(x_{1:n}) = S(x'_{1:n})$ implies $T(x_{1:n}) = T(x'_{1:n})$ for any SS S .

\Leftrightarrow (Lehmann-Scheffé theorem) $\forall x_{1:n}, x'_{1:n} \in \mathcal{X}^n, f_\theta(x_{1:n})/f_\theta(x'_{1:n})$ is free of $\theta \Leftrightarrow T(x_{1:n}) = T(x'_{1:n})$.

$A = A(X_{1:n})$ is **ancillary** (ANS) if the distribution of A does not depend on θ .

T is **complete** (CS) if $\forall \theta \in \Theta, \mathbb{E}_\theta g(T) = 0$ implies $\forall \theta \in \Theta, \mathbb{P}_\theta\{g(T) = 0\} = 1$.

Properties

- (Transformation) If $T = r(T')$, then (i) T is SS $\Rightarrow T'$ is SS; (ii) T' is CS $\Rightarrow T$ is CS; (iii) r is one-to-one, then if one is SS/MSS/CS, then the another is.
- (**Basu's Lemma**) $X_i \sim \text{iid} \mathbb{P}_\theta$, A is ANS and T is CSS, then $A \perp\!\!\!\perp T$.
- (**Bahadur's theorem**) $X_i \sim \text{iid} \mathbb{P}_\theta$, if an MSS exists, then any CSS is also an MSS.
 - Then if a CSS exists, then any MSS is also a CSS $\Rightarrow \text{CSS} = \text{MSS}$.
 - **All or nothing:** start with MSS T , check whether T is CS. (i) Yes, it is both CSS and MSS, then the set of $\text{MSS} = \text{CSS}$; (ii) No, there is no CSS at all.
- (Exp-family) If $X_i \sim \text{iid} f_\eta$ in (2.2), then $T = (\sum_{i=1}^n T_1(X_i), \dots, \sum_{i=1}^n T_k(X_i))$ is a SS, called **natural sufficient statistic**. If Ξ contains an open set in \mathbb{R}^k (i.e., \mathcal{F}' is of full rank), then T is MSS and CSS.

Proof techniques

- Prove T is not sufficient for θ : show if $\exists x_{1:n}, x'_{1:n} \in \mathcal{X}^n$ and $\theta', \theta'' \in \Theta$, such that $T(x_{1:n}) = T(x'_{1:n})$ and $\Lambda(\theta', \theta'' \mid x_{1:n}) \neq \Lambda(\theta', \theta'' \mid x'_{1:n})$.
- Prove A is an ANS: consider location-scale representation.
- Prove T is a CS: use definition or take $d\mathbb{E}_\theta g(T)/d\theta = 0$.
- Disprove T is CS:
 - Construct an ANS $S(T)$ based on T , then $\mathbb{E}S(T)$ is free of θ , then $g(T) = S(T) - \mathbb{E}S(T)$ is free of θ but $g(T) \neq 0$ w.p.1.
 - (Cancel the 1st moment) Find two unbiased estimators for θ as a function of T . E.g., $X_1, X_2 \sim \text{iid} N(\theta, \theta^2)$, $T = (X_1, X_2)$, $g(T) = X_1 - X_2 \sim N(0, 2\theta^2)$.

Remark 2.2.2. • ANS A is useless on its own, but useful together with other information.

- $\mathbb{P}(A(\mathbf{X}) \mid \theta)$ is free of θ , but for non-SS T , $\mathbb{P}(A(\mathbf{X}) \mid T(\mathbf{X}))$ is not necessarily free of θ .

2.2.2 Likelihood principle

Chapter 3

Multivariate Inference Fundamentals

Reference:

- Robb J. Muirhead - Aspects of multivariate statistical theory [5].
- CUHK STAT4002 - Applied Multivariate Analysis (2023 Spring), by Zhixiang Lin.
- Peng DING - Linear Model and Extensions.
- Ronald Christensen - Plane answers to complex questions: the theory of linear models [2].

3.1 Random vectors and distributions

Definition 3.1.1. Let $\mathbf{x} = (x_1, \dots, x_p)^\top \in \mathbb{R}^p$ be a random vector,

- **Mean** $\mathbb{E}\mathbf{x} = \boldsymbol{\mu} = (\mathbb{E}x_1, \dots, \mathbb{E}x_p)^\top = (\mu_j)$.
 - **Covariance matrix** $\text{Var}(\mathbf{x}) = \text{Cov}(\mathbf{x}) = \Sigma = \mathbb{E}[(\mathbf{x} - \mathbb{E}\mathbf{x})(\mathbf{x} - \mathbb{E}\mathbf{x})^\top] = \mathbb{E}\mathbf{x}\mathbf{x}^\top - \mathbb{E}\mathbf{x}\mathbb{E}\mathbf{x}^\top = (\sigma_{ij})$, $\Sigma \succeq \mathbf{0}$.
 - **Correlation matrix** $R = D^{-1/2}\Sigma D^{-1/2}$, where $D = \text{diag}(\sigma_{11}, \dots, \sigma_{pp})$. We have $R_{ij} = \rho_{ij} = \sigma_{ij}/(\sqrt{\sigma_{ii}}\sqrt{\sigma_{jj}})$.
 - If $\mathbf{y} \in \mathbb{R}^q$ random vector, then $\text{Cov}(\mathbf{x}, \mathbf{y}) = \mathbb{E}[(\mathbf{x} - \mathbb{E}\mathbf{x})(\mathbf{y} - \mathbb{E}\mathbf{y})^\top] = \mathbb{E}\mathbf{x}\mathbf{y}^\top - \mathbb{E}\mathbf{x}\mathbb{E}\mathbf{y}^\top \in \mathbb{R}^{p \times q}$.
- If $\mathbf{Z} = (z_{ij}) \in \mathbb{R}^{p \times q}$ is a random matrix,
- $\mathbb{E}\mathbf{Z} = (\mathbb{E}z_{ij})$.

Proposition 3.1.2. Let $\mathbf{x} \in \mathbb{R}^p$ be a random vector, $\mathbf{a}, \mathbf{b} \in \mathbb{R}^p$ be vectors, $A \in \mathbb{R}^{r_1 \times p}$, $B \in \mathbb{R}^{r_2 \times p}$ be matrices,

- $\mathbb{E}\mathbf{a}^\top \mathbf{x} = \mathbf{a}^\top \mathbb{E}\mathbf{x}$, $\text{Var}(\mathbf{a}^\top \mathbf{x}) = \mathbf{a}^\top \Sigma \mathbf{a}$, and $\text{Cov}(\mathbf{a}^\top \mathbf{x}, \mathbf{b}^\top \mathbf{x}) = \mathbf{a}^\top \Sigma \mathbf{b}$.
- $\mathbb{E}A\mathbf{x} = A\mathbb{E}\mathbf{x}$, $\text{Var}(A\mathbf{x}) = A\Sigma A^\top$, and $\text{Cov}(A\mathbf{x}, B\mathbf{x}) = A\Sigma B^\top$.
- If $\mathbf{y} = A\mathbf{x} + \mathbf{b}$, where $A \in \mathbb{R}^{q \times p}$, $\mathbf{b} \in \mathbb{R}^q$, then $\boldsymbol{\mu}_{\mathbf{y}} = A\boldsymbol{\mu}_{\mathbf{x}} + \mathbf{b}$ and $\Sigma_{\mathbf{y}} = A\Sigma A^\top$.
- $\mathbb{E}(\mathbf{x}^\top A\mathbf{x}) = \text{tr}(A\Sigma) + \boldsymbol{\mu}^\top A\boldsymbol{\mu}$.

Let $\mathbf{Z} \in \mathbb{R}^{p \times q}$ be a random matrix, $B \in \mathbb{R}^{m \times p}$, $C \in \mathbb{R}^{q \times n}$, and $D \in \mathbb{R}^{m \times n}$ constants, then

- $\mathbb{E}(B\mathbf{Z}C + D) = B\mathbb{E}(\mathbf{Z})C + D$.

- The $\Sigma \in \mathbb{R}^{p \times p}$ is a covariance matrix (i.e., $\Sigma = \text{Cov}(\mathbf{x})$ for some random vector $\mathbf{x} \in \mathbb{R}^p$) iff $\Sigma \succeq \mathbf{0}$.
– (\Leftarrow): suppose $r(\Sigma) = r \leq p$, write full rank decomposition $\Sigma = CC^\top$, $C \in \mathbb{R}^{p \times r}$. Let $\mathbf{y} \sim [\mathbf{0}_r, I_r]$, then $\text{Cov}(C\mathbf{y}) = \Sigma$.
- If Σ is not PD, then $\exists \mathbf{a} \neq \mathbf{0}_p$ s.t. $\text{Var}(\mathbf{a}^\top \mathbf{x}) = 0$ so w.p.1., $\mathbf{a}^\top \mathbf{x} = k$, i.e., \mathbf{x} lies in a hyperplane.

Theorem 3.1.3. If $\mathbf{x} \in \mathbb{R}^p$ random, then its distribution is uniquely determined by the distributions of $\mathbf{a}^\top \mathbf{x}$, $\forall \mathbf{a} \in \mathbb{R}^p$.

The proof uses the fact that a distribution in \mathbb{R}^p is uniquely determined by its ch.f., see Theorem 1.2.2. [5].

Definition 3.1.4. Dataset contains p variables and n observations are represented by $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top$, where the i th row $\mathbf{x}_i^\top = (x_{i1}, \dots, x_{ip})$ is the i th observation vector, $i = 1, \dots, n$.

- (**Sample mean vector**) $\bar{\mathbf{x}} = n^{-1} \sum_{i=1}^n \mathbf{x}_i = (\bar{x}_1, \dots, \bar{x}_p)^\top$, where $\bar{x}_j = n^{-1} \sum_{i=1}^n x_{ij}$.
- (**Sum of squares and cross product (SSCP) matrix**) $\mathbf{A} = \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top$.
- (**Sample covariance matrix**) $\mathbf{S} = (n-1)^{-1} \mathbf{A}$.
- (**Sample correlation matrix**) $\mathbf{R} = D^{-1/2} \mathbf{S} D^{-1/2}$, where $D^{-1/2} = \text{diag}(1/\sqrt{s_{11}}, \dots, 1/\sqrt{s_{pp}})$.

- $\bar{\mathbf{x}} = n^{-1} \mathbf{X}^\top \mathbf{1}_n$, and

$$\begin{aligned} \mathbf{A} &= \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top - n(\bar{\mathbf{x}} - \boldsymbol{\mu})(\bar{\mathbf{x}} - \boldsymbol{\mu})^\top \\ &= (\mathbf{X} - \mathbf{1}_n \bar{\mathbf{x}}^\top)^\top (\mathbf{X} - \mathbf{1}_n \bar{\mathbf{x}}^\top) \succeq \mathbf{0}. \end{aligned}$$

3.1.1 Multivariate normal distribution

Definition 3.1.5 (Original definition of multivariate normal). The random vector $\mathbf{x} \in \mathbb{R}^p$ is said to have an p -variate normal distribution ($\mathbf{x} \sim N_p$) if $\forall \mathbf{a} \in \mathbb{R}^p$, the distribution of $\mathbf{a}^\top \mathbf{x}$ is univariate normal.

Theorem 3.1.6 (Fundamental properties). Let $\mathbf{x} \sim N_p$, we have

1. Both $\boldsymbol{\mu} = \mathbb{E}\mathbf{x}$ and $\Sigma = \text{Cov}(\mathbf{x})$ exist and the distribution of \mathbf{x} is determined by $\boldsymbol{\mu}$ and Σ . Write $\mathbf{x} \sim N_p(\boldsymbol{\mu}, \Sigma)$.
2. (**Representation**) Let $\Sigma \succeq \mathbf{0}_{p \times p}$, $r(\Sigma) = r \leq p$, and $u_{1:r} \sim \text{iid} N(0, 1)$, i.e., $\mathbf{u} \sim N_r(\mathbf{0}_r, I_r)$, then if C is the full rank decomposition of Σ and $\boldsymbol{\mu} \in \mathbb{R}^p$, then $\mathbf{x} = C\mathbf{u} + \boldsymbol{\mu} \sim N_p(\boldsymbol{\mu}, \Sigma)$.
 - Let $\Sigma = HDH^\top$ be the spectral decomposition, then $\mathbf{x} = HD^{1/2}\mathbf{z} + \boldsymbol{\mu}$, where $\mathbf{z} \sim N_p(\mathbf{0}_p, I_p)$.
3. If $\mathbf{x} \sim N_p(\boldsymbol{\mu}, \Sigma)$, then its **ch.f.** $\phi_{\mathbf{x}}(\mathbf{t}) = \exp(i\boldsymbol{\mu}^\top \mathbf{t} - \mathbf{t}^\top \Sigma \mathbf{t}/2)$.
4. (**Density**) $\mathbf{x} \sim N_p(\boldsymbol{\mu}, \Sigma)$ with $\Sigma \succ \mathbf{0}$, then \mathbf{x} has pdf

$$f(\mathbf{x}) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}. \quad (3.1)$$

Note that we guarantee the existence of $N_p(\boldsymbol{\mu}, \Sigma)$ by means of the representation in point 2. By its density, we have MVN kernel: If

$$f(\mathbf{x}) \propto \exp \left\{ -\frac{1}{2} (\mathbf{x}^\top A \mathbf{x} - 2\mathbf{x}^\top B) \right\} = \exp \left\{ -\frac{1}{2} (\mathbf{x} - A^{-1}B)^\top A (\mathbf{x} - A^{-1}B) - B^\top A^{-1}B \right\},$$

then $\mathbf{x} \sim N_p(A^{-1}B, A^{-1})$.

Theorem 3.1.7 (Properties of multivariate normal). If $\mathbf{x} \sim N_p(\boldsymbol{\mu}, \Sigma)$, then we have

1. (**Linearity**) Let $B \in \mathbb{R}^{q \times p}$, $\mathbf{b} \in \mathbb{R}^q$ nonrandom, and $B\Sigma B^\top \succ \mathbf{0}$, then $B\mathbf{x} + \mathbf{b} \sim N_q(B\boldsymbol{\mu} + \mathbf{b}, B\Sigma B^\top)$.
2. (**Linear combinations**) If $\mathbf{x}_k \sim N_p(\boldsymbol{\mu}_k, \Sigma_k) \perp$ for $k = 1, \dots, N$, then for any fixed constants $\alpha_1, \dots, \alpha_N$, $\sum_{k=1}^N \alpha_k \mathbf{x}_k \sim N_p(\sum_{k=1}^N \alpha_k \boldsymbol{\mu}_k, \sum_{k=1}^N \alpha_k^2 \Sigma_k)$.
 - The sample mean $\bar{\mathbf{x}} \sim N_p(\boldsymbol{\mu}, \Sigma/N)$.
3. (**Subset**) The marginal distribution of any subset of $k (< p)$ components of \mathbf{x} is k -variate normal.
4. (**Marginal distribution**) Partition

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}, \quad \boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}, \quad \mathbf{x}_1 \in \mathbb{R}^q, \mathbf{x}_2 \in \mathbb{R}^{p-q}, \Sigma_{12} \in \mathbb{R}^{q \times (p-q)}.$$

Then $\mathbf{x}_1 \sim N_q(\boldsymbol{\mu}_1, \Sigma_{11})$, $\mathbf{x}_1 \perp \mathbf{x}_2$ iff $\Sigma_{12} = \mathbf{0}$.

5. (**Conditional distribution**) Let Σ_{22}^- be a generalized inverse of Σ_{22} (i.e., $\Sigma_{22}\Sigma_{22}^-\Sigma_{22} = \Sigma_{22}$), then
 - (a) $\mathbf{x}_1 - \Sigma_{12}\Sigma_{22}^-\mathbf{x}_2 \sim N_q(\boldsymbol{\mu}_1 - \Sigma_{12}\Sigma_{22}^-\boldsymbol{\mu}_2, \Sigma_{11} - \Sigma_{12}\Sigma_{22}^-\Sigma_{21})$, and $\perp \mathbf{x}_2$.
 - (b) $[\mathbf{x}_1 | \mathbf{x}_2] \sim N_q(\boldsymbol{\mu}_1 + \Sigma_{12}\Sigma_{22}^-(\mathbf{x}_2 - \boldsymbol{\mu}_2), \Sigma_{11} - \Sigma_{12}\Sigma_{22}^-\Sigma_{21})$.
6. (**Cramér**) If $p \times 1$ random vectors $\mathbf{x} \perp \mathbf{y}$ and $\mathbf{x} + \mathbf{y} \sim N_p$, then both $\mathbf{x}, \mathbf{y} \sim N_p$.
7. (**MLE**) of $(\boldsymbol{\mu}, \Sigma)$ is $(\bar{\mathbf{x}}, A/n)$.
8. (**Inverse of Σ and conditional independence**) Denote $\Sigma^{-1} = (\nu^{jk})_{1 \leq j, k \leq p}$. Then $\forall j \neq k$, $\nu^{jk} = 0 \Leftrightarrow x_j \perp x_k | \mathbf{x} \setminus \{x_j, x_k\}$.

For point 3, each component of a random vector is (marginally) normal does not imply that the vector has a multivariate normal distribution. Counterexample: let $U_1, U_2, U_3 \sim \text{iid} N(0, 1)$, $Z \perp U_{1:3}$. Define

$$X_1 = \frac{U_1 + ZU_3}{\sqrt{1+Z^2}}, \quad X_2 = \frac{U_2 + ZU_3}{\sqrt{1+Z^2}}.$$

Then $[X_1 | Z] \sim N(0, 1)$, free of Z , so $X_1 \sim N(0, 1)$, and $X_2 \sim N(0, 1)$. But (X_1, X_2) not normal. The converse is true if the components of \mathbf{x} are all independent and normal, or if \mathbf{x} consists of independent subvectors, each of which is normally distributed.

For the proof of point 5, we use the lemma: if $\Sigma \succeq \mathbf{0}$, then $\ker(\Sigma_{22}) \subset \ker(\Sigma_{12})$, and $\text{range}(\Sigma_{21}) \subset \text{range}(\Sigma_{22})$. So $\exists B \in \mathbb{R}^{q \times (p-q)}$ satisfying $\Sigma_{12} = B\Sigma_{22}$.

Quadratic form of MVN

Proposition 3.1.8 (Variance). If $\mathbf{x} \sim N(\boldsymbol{\mu}, \Sigma)$, then

$$\begin{aligned} \text{Var}(\mathbf{x}^\top A \mathbf{x}) &= 2\text{tr}(A\Sigma A\Sigma) + 4\boldsymbol{\mu}^\top A\Sigma A\boldsymbol{\mu}, \\ \text{Cov}(\mathbf{x}^\top A_1 \mathbf{x}, \mathbf{x}^\top A_2 \mathbf{x}) &= 2\text{tr}(A_1\Sigma A_2\Sigma) + 4\boldsymbol{\mu}^\top A_1\Sigma A_2\boldsymbol{\mu}. \end{aligned}$$

See Problem B 2.14 in Peng DING.

Theorem 3.1.9 (Independence of quadratic form). Assume $\mathbf{x} \sim N_p(\boldsymbol{\mu}, \Sigma)$.

1. For two symmetric matrix $A, B \in \mathbb{R}^{p \times p}$, $\mathbf{x}^\top A \mathbf{x} \perp \mathbf{x}^\top B \mathbf{x}$ iff

$$\Sigma A \Sigma B \Sigma = \mathbf{0}, \quad \Sigma A \Sigma B \boldsymbol{\mu} = \Sigma B \Sigma A \boldsymbol{\mu} = \mathbf{0}, \quad \boldsymbol{\mu}^\top A \Sigma B \boldsymbol{\mu} = 0.$$

If $\Sigma \succ \mathbf{0}$, then iff $A \Sigma B = \mathbf{0}$.

2. If $\Sigma \succ \mathbf{0}$, $A \in \mathbb{R}^{p \times p}$ symmetric, and $B \in \mathbb{R}^{r \times p}$, then $\mathbf{x}^\top A \mathbf{x} \perp B \mathbf{x}$ iff $B \Sigma A = \mathbf{0}$.

See Theorem B.11 in Peng DING, and Problem 1.22–1.23 in [5].

Theorem 3.1.10. If $\mathbf{x}, \mathbf{x}_{1:N} \sim_{\text{iid}} N_p(\boldsymbol{\mu}, \Sigma)$, where Σ is nonsingular, then

- $(\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \sim \chi_p^2$,
- $\mathbf{x}^\top \Sigma^{-1} \mathbf{x} \sim \chi_p^2(\boldsymbol{\mu}^\top \Sigma^{-1} \boldsymbol{\mu})$,
- partition

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}, \quad \boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}, \quad \mathbf{x}_1, \boldsymbol{\mu}_1 \in \mathbb{R}^k, \quad \Sigma_{11} \in \mathbb{R}^{k \times k}, \quad \text{then}$$

$$Q = (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) - (\mathbf{x}_1 - \boldsymbol{\mu}_1)^\top \Sigma_{11}^{-1} (\mathbf{x}_1 - \boldsymbol{\mu}_1) \sim \chi_{p-k}^2.$$

- $N(\bar{\mathbf{x}}_N - \boldsymbol{\mu})^\top \Sigma^{-1} (\bar{\mathbf{x}}_N - \boldsymbol{\mu}) \sim \chi_p^2$,
- the Mahalanobis distance $d_i^2 = (\mathbf{x}_i - \bar{\mathbf{x}}_N)^\top \mathbf{S}^{-1} (\mathbf{x}_i - \bar{\mathbf{x}}_N) \xrightarrow{d} \chi_p^2$.

If $r(\Sigma) = k \leq p$, then

1.

$$(\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^- (\mathbf{x} - \boldsymbol{\mu}) \sim \chi_k^2,$$

$$(\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^+ (\mathbf{x} - \boldsymbol{\mu}) \sim \chi_k^2.$$

2. If Σ is idempotent with $r(\Sigma) = k$, then $(\mathbf{x} - \boldsymbol{\mu})^\top (\mathbf{x} - \boldsymbol{\mu}) \sim \chi_k^2$.

- If $\boldsymbol{\mu} \in \text{Col}(\Sigma)$, then $\mathbf{x}^\top \mathbf{x} \sim \chi_k^2(\boldsymbol{\mu}^\top \boldsymbol{\mu})$.

Theorem 3.1.11. If $\mathbf{x} \sim N_p(\boldsymbol{\mu}, \Sigma)$,

1. if Σ is nonsingular, $B \in \mathbb{R}^{p \times p}$ is symmetric, then $\mathbf{x}^\top B \mathbf{x} \sim \chi_k^2(\delta)$ iff $B \Sigma$ is idempotent (equiv., $B \Sigma B = B$), in which case $k = r(B)$ and $\delta = \boldsymbol{\mu}^\top B \boldsymbol{\mu}$;
2. for $A \in \mathbb{R}^{p \times p}$, $\mathbf{x}^\top A \mathbf{x} \sim \chi_{r(A \Sigma)}^2(\boldsymbol{\mu}^\top A \boldsymbol{\mu})$ if

$$(1) \Sigma A \Sigma A \Sigma = \Sigma A \Sigma, \quad (2) \boldsymbol{\mu}^\top A \Sigma A \boldsymbol{\mu} = \boldsymbol{\mu}^\top A \boldsymbol{\mu}, \quad (3) \Sigma A \Sigma A \boldsymbol{\mu} = \Sigma A \boldsymbol{\mu}.$$

3.1.2 The noncentral χ^2 and F distribution

3.2 Asymptotic properties

3.2.1 Asymptotic distributions of sample means and covariance matrices

Refer to section 1.2.2, [5].

Theorem 3.2.1 (CLT for sample means). Let $\mathbf{x}_{1:n} \sim_{\text{iid}} [\boldsymbol{\mu}, \Sigma]$, then

$$\sqrt{n}(\bar{\mathbf{x}}_n - \boldsymbol{\mu}) = \frac{1}{\sqrt{n}} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu}) \xrightarrow{d} N_p(\mathbf{0}_p, \Sigma).$$

Theorem 3.2.2 (CLT for sample covariance matrices). Let $\mathbf{x}_{1:n} \sim_{\text{iid}} [\boldsymbol{\mu}, \Sigma]$ with finite fourth moments, SSCP matrix $\mathbf{A} = \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top$, and $\mathbf{S} = (n-1)^{-1} \mathbf{A}$. Let $V = \text{Cov}[\text{vec}((\mathbf{x}_1 - \boldsymbol{\mu})(\mathbf{x}_1 - \boldsymbol{\mu})^\top)]$, then

$$\frac{1}{\sqrt{n}}(\text{vec}(\mathbf{A}) - n \cdot \text{vec}(\Sigma)) \xrightarrow{d} N_{p^2}(\mathbf{0}, V),$$

$$\sqrt{n-1}(\text{vec}(\mathbf{S}) - \text{vec}(\Sigma)) \xrightarrow{d} N_{p^2}(\mathbf{0}, V).$$

Note that $V \in \mathbb{R}^{p^2 \times p^2}$ is singular as the LHS vectors above have repeated elements.

Chapter 4

Bayesian Inference

Chapter 5

Structural Equation Model (SEM)

Reference:

- CUHK STAT5020 - Topics in multivariate analysis (2025 Spring), by Xin Yuan SONG.
- Sik-Yum Lee and Xin-Yuan Song - Basic and advanced Bayesian structural equation modeling: With applications in the medical and behavioral sciences [4].

5.1 SEM models

Goal: to examine the relationships among the variables of interest.

Approach: group observed variables to form latent variables. Advantages:

- Reduce the number of variables compared to direct regression.
- As highly correlated observed variables are grouped into latent variables, the problem induced by multicollinearity is alleviated.
- It gives better assessments on the interrelationships of latent constructs.

Definition 5.1.1 (Linear SEMs). Assume that the observed variables $\mathbf{y}_p \sim \text{iid} N_p$ with mean $\boldsymbol{\mu}_p$. Let $\boldsymbol{\omega}_q$ be latent variables. $\boldsymbol{\omega}_q = (\boldsymbol{\eta}_{q_1}^\top, \boldsymbol{\xi}_{q_2}^\top)^\top$, where $\boldsymbol{\eta}_{q_1}$ is the key outcome latent variables, and $\boldsymbol{\xi}_{q_2}$ is the explanatory latent variables. Define

$$\begin{aligned} \text{(Measurement equation)} \quad & \mathbf{y}_p = \boldsymbol{\mu}_p + \boldsymbol{\Lambda}_{p \times q} \boldsymbol{\omega}_q + \boldsymbol{\epsilon}_p, \\ \text{(Structural equation)} \quad & \boldsymbol{\eta}_{q_1} = \boldsymbol{\Gamma}_{q_1 \times q_2} \boldsymbol{\xi}_{q_2} + \boldsymbol{\delta}_{q_1}, \end{aligned} \quad (5.1)$$

where $\boldsymbol{\Lambda}$ is the unknown factor loading matrix, $\boldsymbol{\Gamma}$ is the unknown matrix of regression coefficients, and $\boldsymbol{\epsilon}$ and $\boldsymbol{\delta}$ are measurement (residual) errors. An extension:

$$\boldsymbol{\eta}_{q_1} = \boldsymbol{\Pi}_{q_1 \times q_1} \boldsymbol{\eta}_{q_1} + \boldsymbol{\Gamma}_{q_1 \times q_2} \boldsymbol{\xi}_{q_2} + \boldsymbol{\delta}_{q_1},$$

where $\boldsymbol{\Pi}$ is a matrix of unknown coefficients such that $I_{q_1} - \boldsymbol{\Pi}$ is nonsingular and the diagonal elements of $\boldsymbol{\Pi}$ are zero.

Assumption 5.1.2 (Standard linear SEMs). For $i = 1, \dots, n$,

- (A1) $\boldsymbol{\epsilon}_i \sim \text{iid} N[\mathbf{0}_p, \boldsymbol{\Psi}_\epsilon]$, where $\boldsymbol{\Psi}_\epsilon \in \mathbb{R}^{p \times p}$ is diagonal.
- (A2) $\boldsymbol{\xi}_i \sim \text{iid} N[\mathbf{0}_{q_2}, \boldsymbol{\Phi}]$, where $\boldsymbol{\Phi}$ is a general.
- (A3) $\boldsymbol{\delta}_i \sim \text{iid} N[\mathbf{0}_{q_1}, \boldsymbol{\Psi}_\delta]$, where $\boldsymbol{\Psi}_\delta$ is diagonal.
- (A4) $\boldsymbol{\delta}_i \perp \boldsymbol{\xi}_i$, and $\boldsymbol{\epsilon}_i \perp \boldsymbol{\omega}_i, \boldsymbol{\delta}_i$.

These assumptions imply that

$$\boldsymbol{\eta}_i \sim \text{iid} N_{q_1}(\mathbf{0}_{q_1}, \boldsymbol{\Gamma} \boldsymbol{\Phi} \boldsymbol{\Gamma}^\top + \boldsymbol{\Psi}_\delta), \quad \boldsymbol{\omega}_i \sim \text{iid} N_q \left(\mathbf{0}_q, \boldsymbol{\Sigma}_\omega = \begin{bmatrix} \boldsymbol{\Gamma} \boldsymbol{\Phi} \boldsymbol{\Gamma}^\top + \boldsymbol{\Psi}_\delta & \boldsymbol{\Gamma} \boldsymbol{\Phi} \\ \boldsymbol{\Phi} \boldsymbol{\Gamma}^\top & \boldsymbol{\Phi} \end{bmatrix} \right), \quad \mathbf{y}_i \sim \text{iid} N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}(\boldsymbol{\theta}) = \boldsymbol{\Lambda} \boldsymbol{\Sigma}_\omega \boldsymbol{\Lambda}^\top + \boldsymbol{\Psi}_\epsilon).$$

There is an identifiability issue relevant to all SEMs: The measurement equation is identified if $\forall \boldsymbol{\theta}_1, \boldsymbol{\theta}_2$, $\text{MeaEq}(\boldsymbol{\theta}_1) = \text{MeaEq}(\boldsymbol{\theta}_2)$ implies $\boldsymbol{\theta}_1 = \boldsymbol{\theta}_2$. The structural equation is identified if $\forall \boldsymbol{\theta}_1, \boldsymbol{\theta}_2$, $\text{StEq}(\boldsymbol{\theta}_1) = \text{StEq}(\boldsymbol{\theta}_2)$ implies $\boldsymbol{\theta}_1 = \boldsymbol{\theta}_2$. The SEM is identified if both of its MeaEq and StEq are identified.

- A simple and common method is using a $\boldsymbol{\Lambda}$ with the non-overlapping structure, e.g.,

$$\boldsymbol{\Lambda}^\top = \begin{bmatrix} 1 & \lambda_{21} & \lambda_{31} & \lambda_{41} & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & \lambda_{62} & \lambda_{72} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & \lambda_{93} & \lambda_{10,3} \end{bmatrix}$$

where 1's are fixed to introduce a scale to latent variables.

Example 5.1.3. Study the kidney disease of type 2 diabetic patients. We observe: plasma creatine (PCr), urinary albumin creatinine ratio (ACR), systolic blood pressure (SBP), diastolic blood pressure (DBP), body mass index (BMI), waist hip ratio (WHR), glycated hemoglobin (HbA1c), fasting plasma glucose (FPG). Group

- {PCr, ACR}: ‘kidney disease (KD)’
- {SBP, DBP}: ‘blood pressure (BP)’
- {BMI, WHR}: ‘obesity (OB)’
- {HbA1c, FPG}: ‘glycemic control (GC)’

$\mathbf{y} = (\text{PCr}, \text{ACR}, \text{SBP}, \text{DBP}, \text{BMI}, \text{WHR})^\top$, $\boldsymbol{\omega} = (\text{KD}, \text{BP}, \text{OB})^\top$, $\boldsymbol{\eta} = \text{KD}$, $\boldsymbol{\xi} = (\text{BP}, \text{OB})^\top$, $p = 6$, $q = 3$, $q_1 = 1$, $q_2 = 2$. Then the measurement equation is

$$\begin{bmatrix} \text{PCr} \\ \text{ACR} \\ \text{SBP} \\ \text{DBP} \\ \text{BMI} \\ \text{WHR} \end{bmatrix} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \\ \mu_4 \\ \mu_5 \\ \mu_6 \end{bmatrix} + \begin{bmatrix} \lambda_{11} & 0 & 0 \\ \lambda_{21} & 0 & 0 \\ 0 & \lambda_{32} & 0 \\ 0 & \lambda_{42} & 0 \\ 0 & 0 & \lambda_{53} \\ 0 & 0 & \lambda_{63} \end{bmatrix} \begin{bmatrix} \text{KD} \\ \text{BP} \\ \text{OB} \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \epsilon_4 \\ \epsilon_5 \\ \epsilon_6 \end{bmatrix}.$$

We know that KD is only linked with PCr and ACR, ... The structural equation can be defined as:

$$\text{KD} = \gamma_1 \text{BP} + \gamma_2 \text{OB} + \delta.$$

Suppose we wish to study the effects of $\boldsymbol{\xi}$ on $\boldsymbol{\eta} = (\text{KD}, \eta_A)^\top$, $q_1 = 2$,

$$\begin{pmatrix} \text{KD} \\ \eta_A \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ \pi & 0 \end{pmatrix} \begin{pmatrix} \text{KD} \\ \eta_A \end{pmatrix} + \begin{pmatrix} \gamma_1 & \gamma_2 \\ \gamma_3 & \gamma_4 \end{pmatrix} \begin{pmatrix} \text{BP} \\ \text{OB} \end{pmatrix} + \begin{pmatrix} \delta \\ \delta_A \end{pmatrix}.$$

Definition 5.1.4 (More SEM models).

- **SEMs with fixed covariates**

$$\begin{aligned} \mathbf{y}_p &= \mathbf{A}_{p \times r_1} \mathbf{c}_{r_1} + \mathbf{\Lambda}_{p \times q} \boldsymbol{\omega}_q + \boldsymbol{\epsilon}_p, \\ \boldsymbol{\eta}_{q_1} &= \mathbf{B}_{q_1 \times r_2} \mathbf{d}_{r_2} + \mathbf{\Pi}_{q_1 \times q_1} \boldsymbol{\eta}_{q_1} + \mathbf{\Gamma}_{q_1 \times q_2} \boldsymbol{\xi}_{q_2} + \boldsymbol{\delta}_{q_1}, \end{aligned} \quad (5.2)$$

where \mathbf{A}, \mathbf{B} are a matrix of unknown coefficients, \mathbf{c}, \mathbf{d} are a vector of fixed covariates (known).

- **Nonlinear SEMs**

$$\begin{aligned} \mathbf{y}_p &= \boldsymbol{\mu}_p + \mathbf{\Lambda}_{p \times q} \boldsymbol{\omega}_q + \boldsymbol{\epsilon}_p, \\ \boldsymbol{\eta}_{q_1} &= \mathbf{\Pi}_{q_1 \times q_1} \boldsymbol{\eta}_{q_1} + \mathbf{\Gamma}_{q_1 \times t} \mathbf{F}_t(\boldsymbol{\xi}_{q_2}) + \boldsymbol{\delta}_{q_1}, \end{aligned} \quad (5.3)$$

where $\mathbf{F}(\boldsymbol{\xi}) = (f_1(\boldsymbol{\xi}), \dots, f_t(\boldsymbol{\xi}))^\top$ with nonzero, known, and linearly independent differentiable functions f_1, \dots, f_t , $t \geq q_2$.

- **Nonlinear SEMs with fixed covariates**

$$\begin{aligned} \mathbf{y}_p &= \mathbf{A}_{p \times r_1} \mathbf{c}_{r_1} + \mathbf{\Lambda}_{p \times q} \boldsymbol{\omega}_q + \boldsymbol{\epsilon}_p, \\ \boldsymbol{\eta}_{q_1} &= \mathbf{B}_{q_1 \times r_2} \mathbf{d}_{r_2} + \mathbf{\Pi}_{q_1 \times q_1} \boldsymbol{\eta}_{q_1} + \mathbf{\Gamma}_{q_1 \times t} \mathbf{F}_t(\boldsymbol{\xi}_{q_2}) + \boldsymbol{\delta}_{q_1}. \end{aligned} \quad (5.4)$$

A simple extension of the StEq is

$$\boldsymbol{\eta}_{q_1} = \mathbf{\Pi}_{q_1 \times q_1} \boldsymbol{\eta}_{q_1} + \mathbf{\Lambda}_{\omega} \mathbf{G}_t(\mathbf{d}, \boldsymbol{\xi}) + \boldsymbol{\delta}_{q_1}$$

where $\mathbf{G}(\mathbf{d}, \boldsymbol{\xi}) = (g_1(\mathbf{d}, \boldsymbol{\xi}), \dots, g_t(\mathbf{d}, \boldsymbol{\xi}))^\top$ is a vector-valued function with nonzero, known, and linearly independent differentiable functions.

Let $\boldsymbol{\Lambda}_k^\top$ be the k th row of $\boldsymbol{\Lambda}$, and $\boldsymbol{\Lambda}_k^\top = (\boldsymbol{\Lambda}_{k\boldsymbol{\eta}}^\top, \boldsymbol{\Lambda}_{k\boldsymbol{\xi}}^\top)$ be a partition correspondings to the partition of $\boldsymbol{\omega} = (\boldsymbol{\eta}^\top, \boldsymbol{\xi}^\top)^\top$. For model (5.3),

$$\begin{aligned} \mathbb{E}(\boldsymbol{\xi}) &= \mathbf{0}_{q_2}, \quad \mathbb{E}(\boldsymbol{\eta}) = [(\mathbf{I}_{q_1} - \mathbf{\Pi})^{-1} \mathbf{\Gamma}] \mathbb{E}(\mathbf{F}(\boldsymbol{\xi})), \\ \mathbb{E}(y_k) &= \mu_k + \boldsymbol{\Lambda}_{k\boldsymbol{\eta}}^\top [(\mathbf{I}_{q_1} - \mathbf{\Pi})^{-1} \mathbf{\Gamma}] \mathbb{E}(\mathbf{F}(\boldsymbol{\xi})). \end{aligned}$$

For model (5.4), let \mathbf{A}_k^\top be the k th row of \mathbf{A} ,

$$\mathbb{E}(y_k) = \mathbf{A}_k^\top \mathbf{c} + \boldsymbol{\Lambda}_{k\boldsymbol{\eta}}^\top \mathbb{E}(\boldsymbol{\eta}) = \mathbf{A}_k^\top \mathbf{c} + \boldsymbol{\Lambda}_{k\boldsymbol{\eta}}^\top [(\mathbf{I}_{q_1} - \mathbf{\Pi})^{-1} \mathbf{\Lambda}_{\omega}] \mathbb{E}(\mathbf{G}(\mathbf{d}, \boldsymbol{\xi})).$$

In developing the Bayesian methods for analyzing SEMs, we usually assign fixed known values to the hyperparameters in the conjugate prior distributions. Consider the modified model (5.4):

$$\begin{aligned} \mathbf{y}_i &= \boldsymbol{\mu} + \boldsymbol{\Lambda} \boldsymbol{\omega}_i + \boldsymbol{\epsilon}_i, \\ \boldsymbol{\eta}_i &= \mathbf{B} \mathbf{d}_i + \boldsymbol{\Pi} \boldsymbol{\eta}_i + \boldsymbol{\Gamma} \mathbf{F}(\boldsymbol{\xi}_i) + \boldsymbol{\delta}_i = \boldsymbol{\Lambda}_\omega \mathbf{G}(\boldsymbol{\omega}_i) + \boldsymbol{\delta}_i, \end{aligned}$$

where $\boldsymbol{\Lambda}_\omega = (\mathbf{B}, \boldsymbol{\Pi}, \boldsymbol{\Gamma}) \in \mathbb{R}^{q_1 \times (r_2 + q_1 + t)}$, and $\mathbf{G}(\boldsymbol{\omega}_i) = (\mathbf{d}_i^\top, \boldsymbol{\eta}_i^\top, \mathbf{F}(\boldsymbol{\xi}_i)^\top)^\top \in \mathbb{R}^{r_2 + q_1 + t}$. Assumption 5.1.2 is satisfied. $\boldsymbol{\xi}_i \sim \text{iid} \mathcal{N}[\mathbf{0}_{q_2}, \boldsymbol{\Phi}]$, $\boldsymbol{\epsilon}_i \sim \text{iid} \mathcal{N}[\mathbf{0}_p, \boldsymbol{\Psi}_\epsilon]$, $\boldsymbol{\Psi}_\epsilon = \text{diag}(\psi_{\epsilon k})$, and $\boldsymbol{\delta}_i \sim \text{iid} \mathcal{N}[\mathbf{0}_{q_1}, \boldsymbol{\Psi}_\delta]$, $\boldsymbol{\Psi}_\delta = \text{diag}(\psi_{\delta k})$.

- Prior (conjugate) for $\boldsymbol{\theta}_y = (\boldsymbol{\mu}, \boldsymbol{\Lambda}, \boldsymbol{\Psi}_\epsilon)$: let $\boldsymbol{\Lambda}_k^\top$ be the k th row of $\boldsymbol{\Lambda}$,

$$\begin{aligned} \psi_{\epsilon k} &\sim \text{IG}(\alpha_{0\epsilon k}, \beta_{0\epsilon k}), \quad [\boldsymbol{\Lambda}_k \mid \psi_{\epsilon k}] \sim \mathcal{N}_q(\boldsymbol{\Lambda}_{0k}, \psi_{\epsilon k} \mathbf{H}_{0y k}), \quad k = 1, \dots, p, \\ \boldsymbol{\mu} &\sim \mathcal{N}_p(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0). \end{aligned}$$

- Prior (conjugate) for $\boldsymbol{\theta}_\omega = (\boldsymbol{\Lambda}_\omega, \boldsymbol{\Psi}_\delta, \boldsymbol{\Phi})$: let $\boldsymbol{\Lambda}_{\omega k}^\top$ be the k th row of $\boldsymbol{\Lambda}_\omega$, $\boldsymbol{\Psi}_\delta$,

$$\begin{aligned} \boldsymbol{\Phi} &\sim \text{IW}_{q_2}(\mathbf{R}_0^{-1}, \rho_0), \quad \text{or } \boldsymbol{\Phi}^{-1} \sim \text{W}_{q_2}(\mathbf{R}_0, \rho_0) \\ \psi_{\delta k} &\sim \text{IG}(\alpha_{0\delta k}, \beta_{0\delta k}), \quad [\boldsymbol{\Lambda}_{\omega k} \mid \psi_{\delta k}] \sim \mathcal{N}_{r_2 + q_1 + t}(\boldsymbol{\Lambda}_{0\omega k}, \psi_{\delta k} \mathbf{H}_{0\omega k}), \quad k = 1, \dots, q_1. \end{aligned}$$

- Assume the prior $\boldsymbol{\theta}_y \perp \boldsymbol{\theta}_\omega$.

Hyperparameter selection: If we have good prior information about a parameter – select the prior distribution with a small variance. E.g.,

- if $\boldsymbol{\Lambda}_k \approx \boldsymbol{\Lambda}_{0k}$, then $\mathbf{H}_{0y k} = 0.5 \mathbf{I}_q$. If not, select the prior with a larger variance;
- since $\epsilon_{ik} \sim \mathcal{N}(0, \psi_{\epsilon k})$, if the variation is small, $\psi_{\epsilon k}$ is small, then choose small $\mathbb{E}(\psi_{\epsilon k}) = \beta_{0\epsilon k} / (\alpha_{0\epsilon k} - 1)$ and $\text{Var}(\psi_{\epsilon k}) = \beta_{0\epsilon k}^2 / \{(\alpha_{0\epsilon k} - 1)^2(\alpha_{0\epsilon k} - 2)\}$;
- if $\boldsymbol{\Phi} \approx \boldsymbol{\Phi}_0$, since $\mathbb{E}(\boldsymbol{\Phi}) = \mathbf{R}_0^{-1} / (\rho_0 - q_2 - 1)$, choose $\mathbf{R}_0^{-1} = (\rho_0 - q_2 - 1) \boldsymbol{\Phi}_0$.

If the sample size is large, can use a portion of the data to estimate $\boldsymbol{\Lambda}_{0k}$, $\boldsymbol{\Lambda}_{0\omega k}$ and $\boldsymbol{\Phi}_0$. If the sample size is moderate, can use the same data twice.

Noninformative prior (Jeffrey): If information is not available and the sample size is small,

$$\begin{aligned} \mathbb{P}(\boldsymbol{\Lambda}, \boldsymbol{\Psi}_\epsilon) &\propto \mathbb{P}(\psi_{\epsilon 1}, \dots, \psi_{\epsilon p}) \propto \prod_{k=1}^p \psi_{\epsilon k}^{-1}, \quad \mathbb{P}(\boldsymbol{\Lambda}_\omega, \boldsymbol{\Psi}_\delta) \propto \mathbb{P}(\psi_{\delta 1}, \dots, \psi_{\delta q_1}) \propto \prod_{k=1}^{q_1} \psi_{\delta k}^{-1}, \\ \mathbb{P}(\boldsymbol{\Phi}) &\propto |\boldsymbol{\Phi}|^{-(q_2+1)/2}. \end{aligned}$$

5.2.1 Bayesian estimation using MCMC

Model: Linear SEM with fixed covariates (5.2) without intercept:

$$\begin{aligned} \mathbf{y}_i &= \boldsymbol{\Lambda} \boldsymbol{\omega}_i + \boldsymbol{\epsilon}_i, \\ \boldsymbol{\eta}_i &= \mathbf{B} \mathbf{d}_i + \boldsymbol{\Pi} \boldsymbol{\eta}_i + \boldsymbol{\Gamma} \boldsymbol{\xi}_i + \boldsymbol{\delta}_i = \boldsymbol{\Lambda}_\omega \mathbf{v}_i + \boldsymbol{\delta}_i, \end{aligned}$$

where $\boldsymbol{\Lambda}_\omega = (\mathbf{B}, \boldsymbol{\Pi}, \boldsymbol{\Gamma}) \in \mathbb{R}^{q_1 \times (r_2 + q_1 + q_2)}$, and $\mathbf{v}_i = (\mathbf{d}_i^\top, \boldsymbol{\eta}_i^\top, \boldsymbol{\xi}_i^\top)^\top \in \mathbb{R}^{r_2 + q_1 + q_2}$. That is, assume $\boldsymbol{\mu} = \mathbf{0}_p$ and $\mathbf{F}(\boldsymbol{\xi}_i) = \boldsymbol{\xi}_i$.

Denote data $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_n) = (\mathbf{Y}_1, \dots, \mathbf{Y}_p)^\top \in \mathbb{R}^{p \times n}$, $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_n) = (\mathbf{V}_1, \dots, \mathbf{V}_{r_2 + q_1 + q_2})^\top$, $\boldsymbol{\Xi}_k = (\eta_{1k}, \dots, \eta_{nk})^\top$ for $k = 1, \dots, q_1$, matrix of latent variables $\boldsymbol{\Omega} = (\boldsymbol{\omega}_1, \dots, \boldsymbol{\omega}_n) \in \mathbb{R}^{q \times n}$, $\boldsymbol{\Omega}_1 = (\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_n)$, $\boldsymbol{\Omega}_2 = (\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_n)$, and

$$\boldsymbol{\theta} = (\boldsymbol{\Lambda}, \mathbf{B}, \boldsymbol{\Pi}, \boldsymbol{\Gamma}, \boldsymbol{\Phi}, \boldsymbol{\Psi}_\epsilon, \boldsymbol{\Psi}_\delta) = (\underbrace{\boldsymbol{\Lambda}, \boldsymbol{\Psi}_\epsilon}_{\boldsymbol{\theta}_y}, \underbrace{\mathbf{B}, \boldsymbol{\Pi}, \boldsymbol{\Gamma}, \boldsymbol{\Phi}, \boldsymbol{\Psi}_\delta}_{\boldsymbol{\theta}_\omega}).$$

Proposition 5.2.1. *The above model has the following posterior distributions:*

1. Conditional distribution $\mathbb{P}(\boldsymbol{\Omega} \mid \mathbf{Y}, \boldsymbol{\theta}) = \prod_{i=1}^n \mathbb{P}(\boldsymbol{\omega}_i \mid \mathbf{y}_i, \boldsymbol{\theta}) \propto \prod_{i=1}^n \mathbb{P}(\boldsymbol{\omega}_i \mid \boldsymbol{\theta}) \mathbb{P}(\mathbf{y}_i \mid \boldsymbol{\omega}_i, \boldsymbol{\theta})$, where

$$\begin{aligned} [\boldsymbol{\omega}_i \mid \boldsymbol{\theta}] &\sim \mathcal{N}_q(\boldsymbol{\mu}_{\omega i}, \boldsymbol{\Sigma}_\omega), \quad [\mathbf{y}_i \mid \boldsymbol{\omega}_i, \boldsymbol{\theta}] \sim \mathcal{N}_p(\boldsymbol{\Lambda} \boldsymbol{\omega}_i, \boldsymbol{\Psi}_\epsilon), \\ [\boldsymbol{\omega}_i \mid \mathbf{y}_i, \boldsymbol{\theta}] &\sim \mathcal{N}_q(\boldsymbol{\Sigma}^{*-1}(\boldsymbol{\Sigma}_\omega^{-1} \boldsymbol{\mu}_{\omega i} + \boldsymbol{\Lambda}^\top \boldsymbol{\Psi}_\epsilon^{-1} \mathbf{y}_i), \boldsymbol{\Sigma}^{*-1}) \end{aligned}$$

where

$$\begin{aligned} \boldsymbol{\Pi}_0 &= \mathbf{I}_{q_1} - \boldsymbol{\Pi}, \quad \boldsymbol{\mu}_{\omega i} = \begin{pmatrix} \boldsymbol{\Pi}_0^{-1} \mathbf{B} \mathbf{d}_i \\ \mathbf{0}_{q_2} \end{pmatrix}, \quad \boldsymbol{\Sigma}_\omega = \begin{bmatrix} \boldsymbol{\Pi}_0^{-1}(\boldsymbol{\Gamma} \boldsymbol{\Phi} \boldsymbol{\Gamma}^\top + \boldsymbol{\Psi}_\delta) \boldsymbol{\Pi}_0^{-\top} & \boldsymbol{\Pi}_0^{-1} \boldsymbol{\Gamma} \boldsymbol{\Phi} \\ \boldsymbol{\Phi} \boldsymbol{\Gamma}^\top \boldsymbol{\Pi}_0^{-\top} & \boldsymbol{\Phi} \end{bmatrix}, \\ \boldsymbol{\Sigma}^* &= \boldsymbol{\Sigma}_\omega^{-1} + \boldsymbol{\Lambda}^\top \boldsymbol{\Psi}_\epsilon^{-1} \boldsymbol{\Lambda}. \end{aligned}$$

2. Assume all elements of $\boldsymbol{\Lambda}_k$ and $\boldsymbol{\Lambda}_\omega$ are unknown, conditional distribution $\mathbb{P}(\boldsymbol{\theta} \mid \mathbf{Y}, \boldsymbol{\Omega}) = \mathbb{P}(\boldsymbol{\theta}_y \mid \mathbf{Y}, \boldsymbol{\Omega}) \mathbb{P}(\boldsymbol{\theta}_\omega \mid \boldsymbol{\Omega})$

$\mathbf{Y}, \mathbf{\Omega}$), where we can show $[\mathbf{\Lambda}_k, \psi_{\epsilon k} \mid \mathbf{Y}, \mathbf{\Omega}] \perp\!\!\!\perp$ and $[\mathbf{\Lambda}_{\omega k}, \psi_{\delta k} \mid \mathbf{\Omega}] \perp\!\!\!\perp$, and

$$\begin{aligned}\mathbb{P}(\boldsymbol{\theta}_{\mathbf{y}} \mid \mathbf{Y}, \mathbf{\Omega}) &\propto \prod_{k=1}^p \mathbb{P}(\mathbf{\Lambda}_k, \psi_{\epsilon k} \mid \mathbf{Y}, \mathbf{\Omega}), \\ \mathbb{P}(\boldsymbol{\theta}_{\omega} \mid \mathbf{Y}, \mathbf{\Omega}) &\propto \left[\prod_{k=1}^{q_1} \mathbb{P}(\mathbf{\Lambda}_{\omega k}, \psi_{\delta k} \mid \mathbf{\Omega}) \right] \mathbb{P}(\boldsymbol{\Phi} \mid \mathbf{\Omega}_2)\end{aligned}$$

where

$$\begin{aligned}\mathbb{P}(\mathbf{\Lambda}_k, \psi_{\epsilon k}^{-1} \mid \mathbf{Y}, \mathbf{\Omega}) &\propto N_q(\mathbf{a}_k, \psi_{\epsilon k} \mathbf{A}_k) \cdot \text{Ga}(n/2 + \alpha_0 \epsilon_k, \beta_{\epsilon k}), \\ \mathbb{P}(\mathbf{\Lambda}_{\omega k}, \psi_{\delta k}^{-1} \mid \mathbf{\Omega}) &\propto N_{r_2+q_1+q_2}(\mathbf{a}_{\omega k}, \psi_{\delta k} \mathbf{A}_{\omega k}) \cdot \text{Ga}(n/2 + \alpha_0 \delta_k, \beta_{\delta k}), \\ [\boldsymbol{\Phi} \mid \mathbf{\Omega}_2] &\sim \text{IW}_{q_2}[(\mathbf{\Omega}_2 \mathbf{\Omega}_2^\top + \mathbf{R}_0^{-1}), n + \rho_0].\end{aligned}$$

where

$$\begin{aligned}\mathbf{A}_k &= (\mathbf{H}_{0yk}^{-1} + \mathbf{\Omega} \mathbf{\Omega}^\top)^{-1}, \quad \mathbf{a}_k = \mathbf{A}_k (\mathbf{H}_{0yk}^{-1} \mathbf{\Lambda}_{0k} + \mathbf{\Omega} \mathbf{Y}_k), \\ \beta_{\epsilon k} &= \beta_{0\epsilon k} + \frac{1}{2} (\mathbf{Y}_k^\top \mathbf{Y}_k - \mathbf{a}_k^\top \mathbf{A}_k^{-1} \mathbf{a}_k + \mathbf{\Lambda}_{0k}^\top \mathbf{H}_{0yk}^{-1} \mathbf{\Lambda}_{0k}), \\ \mathbf{A}_{\omega k} &= (\mathbf{H}_{0\omega k}^{-1} + \mathbf{V}_k \mathbf{V}_k^\top)^{-1}, \quad \mathbf{a}_{\omega k} = \mathbf{A}_{\omega k} (\mathbf{H}_{0\omega k}^{-1} \mathbf{\Lambda}_{0\omega k} + \mathbf{V}_k \boldsymbol{\Xi}_k), \\ \beta_{\delta k} &= \beta_{0\delta k} + \frac{1}{2} (\boldsymbol{\Xi}_k^\top \boldsymbol{\Xi}_k - \mathbf{a}_{\omega k}^\top \mathbf{A}_{\omega k}^{-1} \mathbf{a}_{\omega k} + \mathbf{\Lambda}_{0\omega k}^\top \mathbf{H}_{0\omega k}^{-1} \mathbf{\Lambda}_{0\omega k}).\end{aligned}$$

Remark 5.2.2. For general nonlinear SEMs with fixed covariates (5.4), we can define $\mathbf{u} = [\mathbf{c}^\top, \boldsymbol{\omega}^\top]^\top \in \mathbb{R}^{r_1+q}$ and use similar procedure to derive the full conditional distributions. But by the nonlinear structure $\mathbf{G}(\boldsymbol{\omega})$, $\mathbb{P}(\mathbf{\Omega} \mid \mathbf{Y}, \boldsymbol{\theta})$ may not have closed form like normal, while $\mathbb{P}(\boldsymbol{\theta} \mid \mathbf{Y}, \mathbf{\Omega})$ is not affected and keeps normal-Gamma. To handle fixed parameters, see Appendix 3.3 in [4] and Sec 4.3.1 and Appendix 4.3 in [3]. Also see STAT5020 HW1 Q3.

Bibliography

- [1] G. Casella and R. L. Berger. *Statistical inference*, volume 2. Duxbury Pacific Grove, CA, 2002. [2](#), [2.1](#)
- [2] R. Christensen et al. *Plane answers to complex questions*, volume 35. Springer, 2002. [3](#)
- [3] S.-Y. Lee. *Structural equation modeling: A Bayesian approach*. John Wiley & Sons, 2007. [5.2.2](#)
- [4] S.-Y. Lee and X.-Y. Song. *Basic and advanced Bayesian structural equation modeling: With applications in the medical and behavioral sciences*. John Wiley & Sons, 2012. [5](#), [5.2.2](#)
- [5] R. J. Muirhead. *Aspects of multivariate statistical theory*. John Wiley & Sons, 1982. [3](#), [3.1](#), [3.1.1](#), [3.2.1](#)