

## บทที่ 2

### เอกสารและงานวิจัยที่เกี่ยวข้อง

โดยทั่วไปแล้วเวลาเก็บข้อมูลเพื่อนำมาวิเคราะห์ ข้อมูลที่ได้มามักมีตัวแปรหลายประเภท ดังนั้นในการจัดอันดับโดยใช้การวิเคราะห์ปัจจัย ได้มีการปรับให้สามารถวิเคราะห์ข้อมูลแบบผสมที่มีทั้งตัวแปรแบบต่อเนื่องและตัวแปรเชิงอันดับ โดยทำการวิเคราะห์ปัจจัยผ่านตัวแปรแฝง  $x^*$  แต่ด้วยความซับซ้อนของพารามิเตอร์ทั้งหมดที่ต้องประมาณค่า รวมทั้งตัวแปรแฝงที่ต้องใช้ในการวิเคราะห์ จึงใช้การสุ่มตัวอย่างด้วยวิธีกิบส์และการสุ่มตัวอย่างวิธีฮิตแอนดรันที่อยู่ในกลุ่มวิธีลูกโซ่มาร์คอฟมอนติคาร์โล (Markov chain Monte Carlo) หรือ MCMC ในการจำลองข้อมูล ซึ่งลูกโซ่มาร์คอฟจะเข้าสู่การแจกแจงคงตัว (stationary distribution) ภายใต้อะนุกรมบางประการ แล้วใช้วิธีค่าเฉลี่ยกลุ่ม (Batch mean method) สำหรับการประมาณค่าพารามิเตอร์ และค่าคลาดเคลื่อนมาตรฐาน (SE) ของค่าประมาณ เพื่อนำค่าคลาดเคลื่อนมาตรฐานจากวิธี MCMC ทั้ง 2 มาเปรียบเทียบประสิทธิภาพในการสุ่มเข้าของตัวประมาณ

#### 1. ตัวแบบการวิเคราะห์ปัจจัยสำหรับข้อมูลแบบผสมที่มีทั้งตัวแปรแบบต่อเนื่องและตัวแปรเชิงอันดับ

ให้  $X$  เป็นเมทริกซ์ของค่าสังเกตขนาด  $N \times J$  ซึ่งแต่ละตัวแปรสามารถเป็นได้ทั้งตัวแปรแบบต่อเนื่องและตัวแปรเชิงอันดับที่มี  $C_j > 1$  กลุ่ม โดยให้  $i = 1, 2, \dots, N$  แทนดัชนีของหน่วยที่ให้ค่าสังเกต  $j = 1, 2, \dots, J$  แทนดัชนีของตัวแปร

สมมติให้ ค่าของแต่ละสมาชิกใน  $X$  ได้มาจาก  $X^*$  ซึ่งเป็นเมทริกซ์ของตัวแปรแฝง (Latent variable) ขนาด  $N \times J$  และจุดตัด  $\gamma$

$$x_{ij} = \begin{cases} x_{ij}^* & \text{ถ้า ตัวแปรที่ } j \text{ เป็นตัวแปรแบบต่อเนื่อง} \\ c & \text{ถ้า } x_{ij}^* \in (\gamma_{j(c-1)}, \gamma_{jc}] \text{ และตัวแปรที่ } j \text{ เป็นตัวแปรเชิงอันดับ} \end{cases}$$

โดยที่  $c$  มีค่าเป็นไปได้อย่างตั้งแต่  $1, 2, \dots, C_j$  และให้  $\gamma_{j0} = -\infty, \gamma_{j1} = 0$  และ  $\gamma_{jC_j} = \infty$

รูปแบบของการรวมตัวระหว่างตัวแปรสังเกตได้  $X$  ถูกจำลองโดยตัวแบบการวิเคราะห์ปัจจัยผ่านตัวแปรแฝง  $X^*$

$$x_i^* = \Lambda \phi_i + \varepsilon_i, \quad i = 1, \dots, N$$

โดยที่  $\phi_i \stackrel{iid}{\sim} N(0, I)$  ,  $i = 1, \dots, N$

และ  $\varepsilon_i \stackrel{iid}{\sim} N(0, \Psi)$  ,  $i = 1, \dots, N$

ซึ่งสามารถเขียนให้อยู่ในรูปของ  $N$  สมการได้ดังนี้

$$X_{N \times J}^* = \Phi_{N \times 2} \Lambda'_{2 \times J} + E_{N \times J}$$

เมื่อ  $X^*$  เป็นเมทริกซ์ของตัวแปรแฝงขนาด  $N \times J$   $\Lambda$  เป็นเมทริกซ์ขนาด  $J \times 2$  ของสัมประสิทธิ์ของปัจจัย เมื่อกำหนดให้มี 1 ปัจจัย,  $\Phi$  เป็นเมทริกซ์ขนาด  $N \times 2$  โดยสมาชิกในหลักที่ 1 เป็น 1 ทั้งหมดและสมาชิกในหลักที่ 2 เป็นคะแนนของปัจจัย และ  $\Psi$  เป็นเมทริกซ์ทแยงมุมที่มีขนาด  $J \times J$  โดยให้  $\Psi_{jj}$  เท่ากับ 1 เมื่อตัวแปรที่  $j$  เป็นตัวแปรอันดับ และ  $\Psi_{jj} \stackrel{iid}{\sim} IG(a_0/2, b_0/2)$  เมื่อตัวแปรที่  $j$  เป็นตัวแปรแบบต่อเนื่อง

ในการวิเคราะห์ปัจจัยเชิงเบสส์สำหรับการแจกแจงก่อน (Prior distribution) ของพารามิเตอร์เกี่ยวข้องเป็นดังต่อไปนี้

$\Lambda$  จะต้องมีสมาชิกอย่างน้อยหนึ่งตัวในหลักที่ 2 ที่ถูกกำหนดให้มีเครื่องหมายเป็นบวกหรือลบเท่านั้น เพื่อกำจัดปัญหาการไม่แปรเปลี่ยนภายใต้การหมุน (Rotational invariance) โดยสมมติให้สมาชิกของ  $\Lambda$  ที่ถูกกำหนดเครื่องหมาย มีการแจกแจงปกติแบบตัดปลายที่ตัดความหนาแน่นบริเวณที่ต่ำกว่าหรือเหนือกว่าศูนย์ออก ขึ้นอยู่กับเครื่องหมายที่กำหนด และสมมติให้สมาชิกของ  $\Lambda$  ทุกตัว ทั้งที่ถูกกำหนดเครื่องหมายและไม่ได้ถูกกำหนดเครื่องหมาย เป็นอิสระซึ่งกันและกัน มีการแจกแจงแบบปกติ มีค่าเฉลี่ยเท่ากับ  $L_{oj}$  และค่าความแม่นยำ (ซึ่งเป็นส่วนกลับของความแปรปรวน) เท่ากับ  $L_{oj}$  สมมติให้  $\gamma$  มีการแจกแจงแบบสม่ำเสมอไม่ตรงแบบ (Improper Uniform) และ  $\Lambda, \phi, \Psi$  และ  $\gamma$  เป็นอิสระซึ่งกันและกัน

ในการจัดอันดับโดยใช้การวิเคราะห์ปัจจัยเชิงเบสส์ เราจะใช้ค่าคาดหวังภายใต้การแจกแจงความน่าจะเป็นภายหลังของคะแนนของปัจจัย ( $\phi_j$ ) มาใช้ในการจัดอันดับ ซึ่งการประมาณค่าคาดหวังดังกล่าวสามารถทำได้จากการจำลองโดยใช้เทคนิคลูกโซ่มาร์คอฟ (MCMC) โดยในงานวิจัยนี้จะใช้การสุ่มตัวอย่างแบบกิบส์ (Gibbs Sampler) และการสุ่มตัวอย่างแบบฮิตแอนด์รัน (Hit-and-run Sampler) แล้วนำผลจากการจำลองมาเฉลี่ยเพื่อมาประมาณค่าคาดหวัง แต่เนื่องจากการประมาณค่าอย่างมีประสิทธิภาพจึงใช้วิธีค่าเฉลี่ยกลุ่ม (batch mean method)

ประมาณค่าคลาดเคลื่อนมาตรฐาน (standard error) ของตัวประมาณ และนำค่าดังกล่าวมาใช้ในการเปรียบเทียบประสิทธิภาพของการสุ่มตัวอย่างแบบกิบส์ กับการสุ่มตัวอย่างแบบฮิตแอนดรีน

## 2. ลูกโซ่มาร์คอฟ (Markov Chain)

พิจารณากระบวนการเฟ้นสุ่มแบบไม่ต่อเนื่อง (discrete-time stochastic process)  $\{X_n; n = 0, 1, 2, \dots\}$  กระบวนการนี้เป็นลูกโซ่มาร์คอฟแบบไม่ต่อเนื่อง (discrete-time Markov chain) ถ้า  $m, n \geq 0$  และสถานะ  $i, j \in S, x_u$  และ  $0 \leq u < m$  เป็นจำนวนเต็มที่ไม่เป็นลบ (nonnegative integers)

$$P(X_{m+n} = j | X_m = i, X_u = x_u, 0 \leq u < m) = P(X_{m+n} = j | X_m = i)$$

ลูกโซ่มาร์คอฟแบบไม่ต่อเนื่อง หรือลูกโซ่มาร์คอฟ เป็นกระบวนการเฟ้นสุ่มที่มีคุณสมบัติมาร์คอฟ (Markov property) คือ ความน่าจะเป็นแบบมีเงื่อนไขของอนาคต  $X_{m+n}$  เมื่อกำหนดปัจจุบัน  $X_m$  และอดีต  $X_u$  โดยที่  $0 \leq u < m$  จะขึ้นอยู่กับปัจจุบันเท่านั้น และเป็นอิสระกับอดีต  $X_u$  เรียก  $S$  ว่า สเปซของสถานะ (state space)

### 2.1 ทฤษฎีลูกโซ่มาร์คอฟสำหรับสเปซทั่วไป

สำหรับลูกโซ่มาร์คอฟบนสเปซ  $(S, \mathfrak{S})$  และ  $\{X_n; n = 0, 1, 2, \dots\}$  เป็นลูกโซ่มาร์คอฟ ถ้า

$$\begin{aligned} P(X_{n+1} \in A | X_n, X_{n-1}, \dots, X_0) &= P(X_{n+1} \in A | X_n) \\ &= P(X_n, A) \quad \text{สำหรับ } A \in \mathfrak{S} \end{aligned}$$

ซึ่ง  $P(X_n, A)$  มีคุณสมบัติดังนี้

1. สำหรับแต่ละ  $x \in \mathfrak{S}$  ฟังก์ชัน  $P(x, \cdot)$  เมเชอร์ความน่าจะเป็น (measure probability) บน  $(S, \mathfrak{S})$
2. สำหรับแต่ละ  $A \in \mathfrak{S}$  ฟังก์ชัน  $P(\cdot, A)$  สามารถวัดได้ (measurable)

นิยาม ทฤษฎีบทจำกัด

$$\mu P(A) = \int P(x, A) \mu(dx)$$

เมื่อกำหนดความน่าจะเป็นในการเปลี่ยนสถานะขั้นที่ 1 คือ

$$P^1(x, A) = P(x, A)$$

และความน่าจะเป็นในการเปลี่ยนสถานะขั้นที่  $n$  คือ

$$P^n(x, A) = \int_S P^{n-1}(y, A) P(x, dy)$$

ดังนั้น

$$P(X_n \in A) = \mu P^n(A)$$

เมเชอร์ความน่าจะเป็น  $\pi$  บน  $(S, \mathfrak{S})$  เป็นการแจกแจงคงตัว (stationary distribution) ของ  $P$  ถ้า

$$\pi P = P$$

## 2.2 ทฤษฎีบทการลู่เข้าของลูกโซ่มาร์คอฟในกรณี total variation (Markov Chain Convergence in Total Variation)

ให้  $P$  เป็น  $\varphi$  ที่ irreducible สำหรับบางค่า  $\sigma$  - เมเชอร์ที่หาค่าได้  $\varphi$  บน  $(S, \mathfrak{S})$  ถ้า  $P$  เป็น aperiodic แล้ว

$$\lim_{n \rightarrow \infty} P(X_n \in A | X_0 = x) = \pi(A) \text{ สำหรับทุก } A \in \mathfrak{S}$$

เมื่อ  $\pi$  เป็นการแจกแจงคงตัว (stationary distribution) ของ  $P$  ดังนั้น

$$\pi(\cdot) = \pi P(\cdot) = \int_S P(x, \cdot) \pi(dx)$$

### 1. The First Lyapunov Condition (FLC)

สำหรับเซตไม่ว่าง  $B \subset S$  สเกลาร์ที่เป็นบวก (positive scalars)  $a < 1$   $b$  และ  $\delta$  เป็นจำนวนเต็ม  $m \geq 1$  การแจกแจงความน่าจะเป็น  $\varphi$  บน  $S$  และฟังก์ชัน  $V: S \rightarrow [1, \infty]$  จะได้ว่า

- $P(X_m \in \cdot | X_0 \in z) \geq \delta \varphi(\cdot)$  สำหรับทุก  $z \in B$
- $E(V(X_1) | X_0 = z) \leq aV(z) + bI(z \in B)$  สำหรับทุก  $z \in S$

หมายความว่า ถ้าลูกโซ่มาร์คอฟเป็นไปตามเงื่อนไข FLC และมีคุณสมบัติ aperiodic แล้วลูกโซ่มาร์คอฟนี้เป็น ergodic อย่างสม่ำเสมอ (Uniformly ergodic)

## 2. Markov Chain Strong Law of Large Numbers (MCSSLN)

ให้  $\{X_n; n = 0, 1, 2, \dots\}$  เป็นลูกโซ่มาร์คอฟที่ ergodic อย่างสม่ำเสมอ (uniformly ergodic Markov chain) บนสเปซ  $S$  ที่มี  $\pi$  เป็นการแจกแจงคงตัว (stationary distribution) และฟังก์ชัน  $h: S \rightarrow \mathfrak{R}$  ถ้า  $\pi|h| = E_\pi[h(X_0)] = \int_S h(x)\pi(dx) < \infty$  แล้ว

$$\frac{1}{n} \sum_{i=0}^{n-1} h(X_i) \rightarrow \pi h = \int_S h(x)\pi(dx) \text{ เมื่อ } n \rightarrow \infty$$

## 3. Markov Chain Central Limit Theorem (MCCLT)

ให้  $\{X_n; n = 0, 1, 2, \dots\}$  เป็นลูกโซ่มาร์คอฟที่ ergodic อย่างสม่ำเสมอ (uniformly ergodic Markov chain) บนสเปซ  $S$  และฟังก์ชัน  $h: S \rightarrow \mathfrak{R}$  ด้วยค่า  $h^2(z) \leq cV(z)$  สำหรับบางค่า  $c > 0$  และทุกค่า  $z$  จะได้ว่า

$$\sqrt{n} \left( \frac{1}{n} \sum_{i=0}^{n-1} h(X_i) - \pi h \right) \Rightarrow \sigma N(0, 1)$$

เมื่อ  $\pi$  เป็นการแจกแจงความน่าจะเป็นคงตัว (stationary probability distribution) ของ  $X$

$$\text{และ } \sigma^2 = \text{Var}_\pi[h(X_0)] + 2 \sum_{k=1}^{\infty} \text{cov}_\pi[h(X_0), h(X_k)]$$

## 4. Reversibility

ถ้าลูกโซ่มาร์คอฟใดมีคุณสมบัติ reversibility แล้ว จะสามารถออกแบบเมทริกซ์ของความน่าจะเป็นในการเปลี่ยนสถานะ (transition probability matrix) ได้ง่ายขึ้น

ทฤษฎีบท ถ้าความน่าจะเป็นในการเปลี่ยนสถานะ  $P$  มีคุณสมบัติ reversibility ตามความน่าจะเป็น  $\pi$  แล้ว  $\pi$  จะเป็นการแจกแจงคงตัวของ  $P$

$$\begin{aligned} \text{ให้ } A \subset S \text{ แล้ว } \quad \pi P(A) &= \int_S P(x, A)\pi(dx) \\ &= \int_A P(x, S)\pi(dx) \\ &= \int_A 1\pi(dx) \\ &= \pi(A) \end{aligned}$$

### 2.3 การวิเคราะห์สถานะเสถียรภาพของการจำลอง (Steady State Analysis of a Simulation)

ในการจำลองกระบวนการเพิ่มสุ่มที่เป็นลูกโซ่มาร์คอฟ  $\{X_n; n = 0, 1, 2, \dots\}$  เมื่อ  $X_n$  เป็นเวกเตอร์ของตัวแปรสุ่ม ซึ่ง  $X_n \in S$  และ  $S$  เป็นสเปซของสถานะ (state space) โดยเงื่อนไขทั่วไป การแจกแจงของ  $X_n$  จะเข้าสู่การแจกแจงคงตัว  $\pi$  ซึ่งในการวิเคราะห์เชิงเบย์ก็คือการแจกแจงภายหลัง และเป็นอิสระกับจุดเริ่มต้น  $x_0$  ดังนี้

$$p(X_n \in A | X_0 = x_0) \rightarrow \pi(A)$$

จาก MCSLLN จะได้ว่า ค่าเฉลี่ยในระยะยาว (long run average) จะสามารถประมาณค่า  $\pi h$  ดังนี้

$$\frac{1}{N} \sum_{i=0}^{N-1} h(X_i) \rightarrow \pi h \text{ สำหรับ } N \text{ ที่มีค่ามาก}$$

จาก MCCLT จะได้ว่า การแจกแจงของค่าเฉลี่ยในระยะยาว (distribution of long run average) จะประมาณ (approximate) ด้วยการแจกแจงแบบปกติ (normal distribution) ดังนี้

$$\frac{1}{N} \sum_{i=0}^{N-1} h(X_i) \sim N(\pi h, \sigma^2) \text{ สำหรับ } N \text{ ที่มีค่ามาก}$$

$$\text{เมื่อ } \sigma^2 = \text{Var}_\pi [h(X_0)] + 2 \sum_{k=1}^{\infty} \text{cov}_\pi [h(X_0), h(X_k)]$$

### 3. เทคนิคลูกโซ่มาร์คอฟมอนติคาร์โล (Markov Chain Monte Carlo)

ในทฤษฎีบทของเบย์สกรณีที่มีการแจกแจงภายหลัง (Posterior distribution) ไม่ได้อยู่ในรูปแบบสามัญนั้น การคำนวณค่าคาดหวังของพารามิเตอร์จากฟังก์ชันความหนาแน่นน่าจะเป็นภายหลังเป็นปัญหาสำคัญของการสร้างแบบจำลองเบย์ เทคนิคลูกโซ่มาร์คอฟมอนติคาร์โล (Markov Chain Monte Carlo) หรือ MCMC จึงได้รับการเสนอเพื่อแก้ปัญหาดังกล่าว โดย MCMC เป็นเทคนิคที่ใช้ในการสุ่มตัวอย่างของตัวแปรสุ่มจากฟังก์ชันความหนาแน่นน่าจะเป็นภายหลัง

MCMC เป็นการสร้างลำดับของเวกเตอร์สุ่ม  $\{\bar{\theta}_0, \bar{\theta}_1, \dots, \bar{\theta}_n\}$  โดยที่การสุ่ม  $\bar{\theta}_{j+1}$  จากการแจกแจง  $p(\bar{\theta}_{j+1} | \bar{\theta}_j)$  เมื่อ  $j \geq 0$  ซึ่งหมายถึง  $\bar{\theta}_{j+1}$  ถูกสุ่มโดยขึ้นอยู่กับ  $\bar{\theta}_j$  เพียงอย่างเดียวไม่ได้ขึ้นอยู่กับ  $\{\bar{\theta}_0, \bar{\theta}_1, \dots, \bar{\theta}_{j-1}\}$  นั่นคือ  $p(\bar{\theta}_{j+1} | \bar{\theta}_j)$  เป็น Transition Kernel ซึ่งเป็นความ

น่าจะเป็นของการเปลี่ยนสถานะจากสถานะที่  $j$  ไปยังสถานะที่  $j+1$  โดยเงื่อนไขทั่วไป เมื่อ  $j$  มีค่ามากขึ้น การแจกแจงของ  $\bar{\theta}_j$  จะเข้าสู่การแจกแจงคงตัว  $\pi$  (stationary distribution) ในกรณีการวิเคราะห์เชิงเบย์  $\pi$  ก็คือการแจกแจงภายหลัง

ในการสร้าง MCMC นั้นจะมีขั้นตอนวิธี (Algorithm) หลายแบบ เช่น การสุ่มตัวอย่างแบบเมโทรโพลิส-เฮสติงส์ (Metropolis-Hastings Sampling), การสุ่มตัวอย่างแบบกิบส์ (Gibbs Sampling) และ การสุ่มตัวอย่างแบบฮิตแอนดร์น ซึ่งมีรายละเอียดดังต่อไปนี้

### 3.1 การสุ่มตัวอย่างแบบเมโทรโพลิส-เฮสติงส์ (Metropolis-Hastings Sampling)

สมมติว่าต้องการสุ่มพารามิเตอร์  $\bar{\theta}$  จากการแจกแจง  $\pi$  ซึ่งก็คือการแจกแจงภายหลัง (posterior distribution) ของการวิเคราะห์เชิงเบย์ ในการกระบวนการ (Algorithm) ของเมโทรโพลิส-เฮสติงส์มีขั้นตอนดังต่อไปนี้

ขั้นที่ 1 กำหนดพารามิเตอร์เริ่มต้น  $\bar{\theta}_0$  โดยที่  $\pi > 0$

ขั้นที่ 2 สุ่มพารามิเตอร์  $\bar{\theta}^{(can)}$  จาก Proposal Distribution  $q(\bar{\theta}_j | \bar{\theta}_{j+1})$  ซึ่งเป็นการนำจะเป็นในการเปลี่ยนค่าพารามิเตอร์จาก  $\bar{\theta}_j$  ไปเป็น  $\bar{\theta}_{j+1}$

ขั้นที่ 3 คำนวณค่า  $\alpha$  ซึ่งเป็นอัตราส่วนของฟังก์ชันความหนาแน่น  $\pi$  ของพารามิเตอร์ที่สุ่มขึ้นมาใหม่  $\bar{\theta}^{(can)}$  เทียบกับฟังก์ชันความหนาแน่นของพารามิเตอร์  $\bar{\theta}_j$

$$\alpha = \min \left( \frac{\pi(\bar{\theta}^{(can)}) q(\bar{\theta}_j | \bar{\theta}^{(can)})}{\pi(\bar{\theta}_j) q(\bar{\theta}^{(can)} | \bar{\theta}_j)}, 1 \right)$$

ขั้นที่ 4 สุ่ม  $u$  จากการแจกแจงยูนิฟอร์ม (Uniform Distribution)  $U(0,1)$

ยอมรับ  $\bar{\theta}^{(can)}$  ถ้า  $u < \alpha$  และให้  $\bar{\theta}_{j+1} = \bar{\theta}^{(can)}$

หรือ ปฏิเสธ  $\bar{\theta}^{(can)}$  ถ้า  $u \geq \alpha$  และให้  $\bar{\theta}_{j+1} = \bar{\theta}_j$

ขั้นที่ 5 วนซ้ำในขั้นที่ 2 ใหม่

โดยเงื่อนไขทั่วไป การแจกแจงของ  $\bar{\theta}_j$  จะเข้าสู่การแจกแจงคงตัว  $\pi$  (stationary distribution) ซึ่งสามารถนำไปประมาณค่าคาดหวังภายใต้การแจกแจงภายหลังได้

### 3.2 การสุ่มตัวอย่างแบบกิบส์ (Gibbs Sampling)

การสุ่มตัวอย่างแบบกิบส์เป็นกรณีพิเศษของ Single-component Metropolis Hasting Method กล่าวคือจะสุ่ม  $\bar{\theta}^{(can)}$  จาก Full Condition Distribution แทนการสุ่มจาก Proposal Distribution นั่นคือ

$$q(\bar{\theta}_j^{(can)} | \bar{\theta}_j, \bar{\theta}_{-j}) = \pi(\bar{\theta}_j^{(can)} | \bar{\theta}_{-j})$$

$$\text{และ } q(\bar{\theta}_j | \bar{\theta}_j^{(can)}, \bar{\theta}_{-j}) = \pi(\bar{\theta}_j | \bar{\theta}_{-j})$$

ซึ่งทำให้

$$\begin{aligned} \alpha &= \min \left( \frac{\pi(\bar{\theta}_j^{(can)} | \bar{\theta}_{-j}) q(\bar{\theta}_j | \bar{\theta}_j^{(can)}, \bar{\theta}_{-j})}{\pi(\bar{\theta}_j | \bar{\theta}_{-j}) q(\bar{\theta}_j^{(can)} | \bar{\theta}_j, \bar{\theta}_{-j})}, 1 \right) \\ &= \min \left( \frac{\pi(\bar{\theta}_j^{(can)} | \bar{\theta}_{-j}) p(\bar{\theta}_j | \bar{\theta}_{-j})}{\pi(\bar{\theta}_j | \bar{\theta}_{-j}) p(\bar{\theta}_j^{(can)} | \bar{\theta}_{-j})}, 1 \right) = 1 \end{aligned}$$

การสุ่ม  $\bar{\theta}_j^{(can)}$  จาก Full Conditional Distribution จึงถูกยอมรับในทุกๆ ครั้ง ดังนั้น การสุ่มพารามิเตอร์ในการสุ่มตัวอย่างแบบกิบส์จึงเป็นการสุ่มจาก Full Conditional Distribution เป็นดังนี้

$$\pi(\bar{\theta}_j | \bar{\theta}_{-j}) = \frac{\pi(\bar{\theta})}{\int \pi(\bar{\theta}) d\theta_j} \text{ โดยที่ } \bar{\theta}_{-j} = \{\theta_1, \theta_2, \dots, \theta_{j-1}, \theta_{j+1}, \dots, \theta_n\}$$

โดยเงื่อนไขทั่วไป การแจกแจงของ  $\bar{\theta}_j$  จะลู่เข้าสู่การแจกแจงคงตัว  $\pi$  (stationary distribution) ซึ่งสามารถนำไปประมาณค่าคาดหวังภายใต้การแจกแจงภายหลังได้

### 3.3 การสุ่มตัวอย่างแบบฮิตแอนด์รัน (Hit-and-run sampler)

ให้  $S \subset \mathcal{R}^n$  และให้การแจกแจงเป้าหมาย (Target distribution) มีฟังก์ชันความน่าจะเป็น  $\pi$  บน  $S$

สร้าง  $\bar{\theta}_n$  ด้วยกระบวนการดังนี้

ขั้นที่ 1 กำหนดพารามิเตอร์เริ่มต้น  $\bar{\theta}_0 \in S$

ขั้นที่ 2 ที่  $\bar{\theta}_n = \bar{\theta} \in S$  สุ่มเลือกทิศทาง  $\bar{d}$  บนพื้นผิวทรงกลมรัศมี 1 หน่วย

ลากเส้นตรง  $L$  ผ่าน  $\bar{\theta}$  และตัดกับ  $S$  โดยให้  $L_n = S \cap \{l : l = \bar{\theta} + r\bar{d}, r \in \mathcal{R}\}$



ขั้นที่ 3 สุ่มเลือกจุด  $\bar{\theta}_{n+1}$  บนเส้นตรง  $L_n$  ด้วยการแจกแจงแบบ  $\pi$  โดยมีเงื่อนไขบน  $L_n$

ขั้นที่ 4 วนซ้ำในขั้นที่ 2 ใหม่

โดยเงื่อนไขทั่วไป การแจกแจงของ  $\bar{\theta}_j$  จะเข้าสู่การแจกแจงคงตัว  $\pi$  (stationary distribution) ซึ่งสามารถนำไปประมาณค่าคาดหวังภายใต้การแจกแจงภายหลังได้

จากการสุ่มซ้ำของลูกโซ่มาร์คอฟจึงมีค่าความแปรปรวนร่วม (covariance) เพราะ  $\bar{\theta}_j$  แต่ละตัวไม่เป็นอิสระกัน ดังนั้น งานวิจัยนี้จึงใช้วิธีค่าเฉลี่ยกลุ่ม (Batch Means method) สำหรับแก้ปัญหาค่าความแปรปรวนร่วม

#### 4. วิธีค่าเฉลี่ยกลุ่ม (Batch means method)

เป็นวิธีการหนึ่งที่ใช้ประมาณค่าเฉลี่ยประชากร โดยใช้การคำนวณ Monte Carlo standard error (MCSE) สมมติว่าเราสนใจประมาณค่าเฉลี่ย  $E_\pi(g(\bar{\theta}))$  เมื่อ  $\bar{\theta}$  มีการแจกแจง  $\pi$  ซึ่งกระบวนการของค่าเฉลี่ยกลุ่ม (Batch means) มีดังนี้

ขั้นที่ 1 จำลองลูกโซ่มาร์คอฟ  $\{\bar{\theta}_n\}$  จำนวน  $n = ab$  รอบ โดยที่  $a$  และ  $b$  เป็นจำนวนเต็ม  $a$  เป็นจำนวน batch และแต่ละ batch มีขนาด  $b$

ขั้นที่ 2 หาค่าเฉลี่ยของตัวอย่าง ใน batch ที่  $k$  จากสูตร  $\bar{Y}_k = \frac{\sum_{i=(k-1)b+1}^{kb} g(\bar{\theta}_i)}{b}$

ขั้นที่ 3 กำหนดให้  $\hat{\sigma}_g^2 = \frac{b}{a-1} \sum_{k=1}^a (\bar{Y}_k - \bar{g}_n)^2$  โดยที่  $\bar{g}_n = \frac{1}{a} \sum_{k=1}^a \bar{Y}_k$  จะได้ว่า

ค่าประมาณของ Monte Carlo standard error คือ  $\frac{\hat{\sigma}_g}{\sqrt{n}}$

ขั้นที่ 4 คำนวณหาช่วงความเชื่อมั่นของ  $E_\pi(g(\bar{\theta}))$  ได้จากสูตร

$$\bar{g}_n \pm t_{\frac{\alpha}{2}, (a-1)} \frac{\hat{\sigma}_g}{\sqrt{n}}$$

เอกสารและงานวิจัยที่เกี่ยวข้อง



Chen และ Schmeiser (1993) ได้ทำการเปรียบเทียบประสิทธิภาพการสุ่มตัวอย่างแบบฮิตแอนด์รัน และการสุ่มตัวอย่างแบบกิบส์ พบว่าการสุ่มตัวอย่างแบบฮิตแอนด์รันมีประสิทธิภาพมากกว่าสำหรับในกรณีที่เป็นการแจกแจงแบบปกติของ 2 ตัวแปร

Quinn (2004) ได้นำเสนอการใช้ตัวแบบการวิเคราะห์ปัจจัยเชิงเบส (Bayesian factor analysis) ซึ่งได้ปรับให้เหมาะกับการวิเคราะห์ข้อมูลแบบผสมที่มีทั้งตัวแปรแบบต่อเนื่องและตัวแปรเชิงอันดับ และใช้วิธีการคำนวณที่วางอยู่บนพื้นฐานของการสุ่มตัวอย่างแบบกิบส์ และวิธีของ Cowles (1993) ซึ่งวิธีดังกล่าวจัดอยู่ในกลุ่มวิธีลูกโซ่มาร์คอฟ มอนติคาร์โล (Markov chain Monte Carlo) หรือ MCMC ในการประมาณค่าพารามิเตอร์

Lovasz และ Vempala (2006) ได้ศึกษาการสุ่มตัวอย่างแบบฮิตแอนดรันบนบริเวณคอนเวกซ์ พบว่า การสุ่มตัวอย่างแบบฮิตแอนดรันเป็นหนึ่งในวิธีที่เร็วที่สุด ในการสร้างจุดตัวอย่างซึ่งลู่อู่เข้าสู่การแจกแจงแบบสม่ำเสมอบน Log-concave

จึงเป็นที่น่าสนใจว่า หากแทนคำนวณวิธีการสุ่มตัวอย่างแบบกิบส์ด้วยการสุ่มตัวอย่างแบบฮิตแอนดรันในตัวแบบการวิเคราะห์ปัจจัยเชิงเบสที่ Quinn ได้นำเสนอ ผลที่ได้ยังคงสอดคล้องกับงานวิจัยของ Chen และ Schmeiser และงานวิจัยของ Lovasz และ Vempala หรือไม่