

บทที่ 2

เอกสารและงานวิจัยที่เกี่ยวข้อง

ในบทนี้จะกล่าวถึงเอกสารและงานวิจัยที่เกี่ยวข้องต่าง ๆ ที่เป็นความรู้พื้นฐาน โดยผู้วิจัยแบ่งออกเป็น 5 ตอนดังนี้

1. วิธีการจัดการข้อมูลสูญหาย
2. การสุ่มตัวอย่าง
3. การจำลองสถานการณ์ด้วยวิธีการมอนติ คาร์โล
4. ความคลาดเคลื่อนประเภทที่ 1 ประเภทที่ 2 และอำนาจการทดสอบ
5. เกณฑ์การเปรียบเทียบค่าประมาณพารามิเตอร์ระหว่างวิธีการจัดการข้อมูลสูญหายแบบต่าง ๆ
6. งานวิจัยที่เกี่ยวข้อง
7. กรอบแนวคิดในการวิจัย

วิธีการจัดการข้อมูลสูญหาย

ข้อมูลสูญหายเป็นปัญหาที่พบได้โดยทั่วไปในงานวิจัยเชิงสำรวจเพราะว่าจะต้องศึกษากับคนและตัวแปรจำนวนมาก ข้อมูลสูญหายมี 2 ลักษณะ คือ กลุ่มตัวอย่างไม่ตอบแบบสอบถาม เช่น ส่งแบบสอบถามไปแล้วไม่ส่งกลับมา ผู้สัมภาษณ์ไม่พบคนใดคนหนึ่งที่บ้าน และกลุ่มตัวอย่างไม่ตอบข้อความบางข้อ การกลับไปวัดซ้ำคุณลักษณะที่สนใจทำให้สูญเสียทั้งเงินและเวลา จึงทำให้นักสถิติได้คิดวิธีที่จะแทนค่าข้อมูลสูญหายเหล่านั้น

การแก้ปัญหาข้อมูลสูญหายวิธีที่ง่ายที่สุดก็คือการตัดข้อมูลออกไปจากการวิเคราะห์แต่วิธีนี้เป็นวิธีที่ไม่เหมาะสมเพราะกลุ่มตัวอย่างจะหายไปเป็นจำนวนมาก ซึ่งในปัจจุบันมีเทคนิคทางสถิติมากมายที่ใช้ในการจัดการข้อมูลสูญหาย ผู้วิจัยสามารถนำไปประยุกต์ใช้ในการวิจัยเพื่อทำให้การสรุปผลงานวิจัยนั้นได้ถูกต้องยิ่งขึ้น

ปัญหาของข้อมูลสูญหายมี 3 ประเภท คือ

1. โอมิสชั่น (Omissions) จะเกิดขึ้นเมื่อผู้ตอบคำถามตอบไม่ครบในการสำรวจข้อมูล

แต่ครั้ง การสูญหายประเภทนี้จะเกิดขึ้นน้อย ปัญหาที่พบบางงานวิจัยก็คือ ผู้ตอบจะตอบคำถามไม่ครบในการสำรวจข้อมูลอีกครั้งหนึ่ง เหตุผลหลักของการสูญหายจากการสำรวจข้อมูลใหม่ก็คือ การอ่านหนังสือซ้ำ ดังนั้นถ้ารวมทักชะการอ่านไว้ในโมเดลของข้อมูลสูญหายผลกระทบของการสูญหายประเภทนี้ก็ไม่สำคัญ

2. แอททริชัน (Attrition) ผู้ตอบจะไม่ตอบข้อความทั้งหมดในการวิจัยที่มีการวัดซ้ำหลาย ๆ ครั้ง จะเกิดการสูญหายของผู้ตอบในการวัดครั้งแรก แต่จะตอบคำถามเมื่อมีการวัดซ้ำหรือบางคนจะไม่ตอบคำถามไม่ว่าจะวัดกี่ครั้ง ลักษณะของการสูญหายที่บางครั้งตอบบางครั้งไม่ตอบสามารถใช้วิธีการทางสถิติ คือ การถดถอยมาช่วยในการทำนายข้อมูลสูญหายได้ ปัญหาใหญ่ก็คือ ผู้ตอบสูญหายไป แต่ก็มีวิธีการทางสถิติที่สามารถจัดการกับข้อมูลประเภทนี้ได้

3. แพลนมิสซิงเนส (Planned missingness) การสูญหายของข้อมูลประเภทนี้อยู่ภายใต้การควบคุมของนักวิจัยลักษณะของการสูญหายแบบนี้จะเป็น MCAR (Missing completely at random) วิธีการแบบนี้จะทำให้อำนาจการทดสอบน้อย และเป็นอันตรายต่อการสรุปทางสถิติ

ผลที่ตามมาเมื่อมีข้อมูลสูญหาย

นักวิจัยบางคนที่เก็บข้อมูลไม่ครบ ไม่สมบูรณ์ อาจจะไม่กังวลใจมากนักเมื่อเกิดข้อมูลสูญหาย เพราะเมื่อเตรียมข้อมูลเรียบร้อยแล้วเครื่องคอมพิวเตอร์ก็สามารถคำนวณได้ ถึงแม้ว่าจะมีข้อมูลสูญหาย หรือ ตัดข้อมูลที่ไม่สมบูรณ์ออกไป หรือ เก็บเพิ่มเติม หรืออาจจะใช้วิธีเก็บข้อมูลมากกว่าขนาดตัวอย่างที่ต้องการถ้าพบว่าหน่วยข้อมูลใดมีข้อมูลที่ไม่สมบูรณ์ก็ตัดทิ้งไป (สุรินทร์า วีรกุลเทวีณ, 2544. หน้า 17)

ข้อมูลสูญหายทำให้เกิดปัญหาในทางสถิติ 2 ประการ คือ

1. อำนาจในการทดสอบสมมติฐานทางสถิติ
2. ความถูกต้องของค่าพารามิเตอร์ที่ต้องการประมาณค่า

ยกตัวอย่าง เช่น การทำวิจัยเรื่องหนึ่งมีการเก็บข้อมูลจากตัวแปร 10 ตัวแปร เก็บข้อมูลจากกลุ่มตัวอย่าง 100 คน ถ้าเก็บข้อมูลครบจะมีข้อมูลทั้งหมด 1000 ค่า แต่ถ้ามีข้อมูลสูญหายจากผู้ตอบ 25 คน เพียงตัวแปรเดียว ผู้วิจัยตัดข้อมูลออกจะเหลือข้อมูลเพียง 75% จะเห็นว่าขนาดกลุ่มตัวอย่างลดลงซึ่งจะส่งผลต่ออำนาจการทดสอบ ดังนั้นถ้าผู้วิจัยใช้วิธีการตัดข้อมูลออกไปก็จะทำให้เกิดปัญหาในการสรุปผลที่ได้จากการวิเคราะห์ข้อมูลชุดนั้น

ผลที่ตามมาอีกประการหนึ่ง คือ การวิจัยโดยทั่วไปผู้วิจัยต้องการศึกษาจากกลุ่มตัวอย่างแล้วนำค่าสถิติไปอ้างอิงประชากร หรือสรุปไปยังค่าของประชากร เช่น การสอบถามเกี่ยวกับรายได้ ผู้ที่มีรายได้สูงจะไม่ค่อยให้ข้อมูลเกี่ยวกับรายได้ซึ่งถ้าเก็บข้อมูลแล้วนำมาหาค่าเฉลี่ย ข้อมูลจริง ๆ ที่มีค่าสูงจะสูญหายไปทำให้ค่าเฉลี่ยที่คำนวณจากกลุ่มตัวอย่างต่ำกว่าเป็นจริง ทำให้การประมาณค่าเฉลี่ยรายได้ของประชากรเกิดอคติทางลบ (negative bias)

สาเหตุของข้อมูลสูญหาย

สาเหตุของข้อมูลสูญหาย (missing data mechanisms) มี 3 ประเภท คือ

1. ข้อมูลสูญหายแบบสุ่มสมบูรณ์ (missing completely at random : MCAR) เช่น ถ้ากลุ่มตัวอย่างได้รับการคัดเลือกขึ้นมาแบบสุ่มเพื่อให้ได้รับการวัดอย่างใดอย่างหนึ่งจะมีกลุ่มตัวอย่างบางคนที่ไม่สามารถวัดได้อย่างสมบูรณ์ ข้อมูลที่สูญหายไปจะเป็นแบบ MCAR ถึงแม้ว่าการสูญหายจะเกิดขึ้นจากตัวแปรบางตัวแต่ไม่สัมพันธ์กับตัวแปรที่มีข้อมูลสูญหาย ข้อมูลสูญหายยังคงมีลักษณะเป็น MCAR ข้อดีของการสูญหายแบบ MCAR ก็คือ สาเหตุของการสูญหายไม่ต้องนำไปเป็นส่วนหนึ่งของการวิเคราะห์เพื่อควบคุมอคติของการสูญหาย มีวิธีการแบบเก่าในการจัดการข้อมูลสูญหาย เช่น การตัดข้อมูลออกแบบลิสต์ไวส์ ให้ผลการวิเคราะห์ที่ไม่มีอคติ เมื่อข้อมูลสูญหายเป็นแบบ MCAR แต่ก็ยังเป็นวิธีการที่ไม่ดีเพราะอำนาจการทดสอบมีค่าต่ำ
2. ข้อมูลสูญหายแบบมีระบบ (systematic missing data) การสูญหายของข้อมูลเกิดขึ้นจากตัวแปรอื่น ๆ ที่มีความสัมพันธ์กับตัวแปรที่มีข้อมูลสูญหาย สาเหตุของการสูญหายจะไม่เป็นแบบ MCAR ถ้าสาเหตุของการสูญหายวัดได้แล้วนำมาวิเคราะห์จะเรียกว่าเป็น accessible missing data mechanisms ดังนั้นอคติทั้งหมดที่เกี่ยวข้องกับข้อมูลสูญหายจะถูกปรับ ลิทเทิลและรูบิน (Little & Rubin, 1987) เรียกสถานการณ์นี้ว่า ignorable (missing at random)
3. มีสถานการณ์อื่น ๆ ที่เป็นสาเหตุของการสูญหาย เช่น สาเหตุของการสูญหายไม่สามารถวัดได้ และสาเหตุของการสูญหายสัมพันธ์กับตัวแปรที่มีข้อมูลสูญหาย จะเรียกว่าเป็น inaccessible missing data mechanisms (nonignorable mechanisms) สถานการณ์นี้เกิดขึ้นเมื่อค่าของตัวแปรสูญหายเป็นสาเหตุให้เกิดการสูญหาย เช่น คนที่ดื่มสุรามาก ๆ จะหลีกเลี่ยงการตรวจแอลกอฮอล์มากกว่าคนที่ดื่มน้อย หรือเด็กที่ดื้อรั้นจะต่อต้านการตรวจปัสสาวะทั้ง ๆ ที่เขาอาจจะไข้ยาหรือไม่ไข้ยาก็ได้ การสูญหายของข้อมูลก็จะเกิดขึ้นและเป็นสาเหตุที่เข้าถึงไม่ได้ inaccessible mechanism

ลิตเทิลและรูบิน (Little & Rubin, 1987. pp. 6-7) ได้แบ่งวิธีการจัดการข้อมูลสูญหายไว้ 4 วิธี คือ

1. วิธีที่ใช้ข้อมูลสมบูรณ์ (Procedures based on completely recorded units) วิธีนี้เป็นวิธีที่ง่ายที่สุดโดยการตัดข้อมูลที่สูญหายออกไปแล้ววิเคราะห์ข้อมูลที่สมบูรณ์เท่านั้น เช่น การตัดข้อมูลแบบลิสต์ไวส์ (listwise deletion) การตัดข้อมูลแบบแพร์ไวส์ (pairwise deletion) เป็นวิธีที่มีอยู่ในโปรแกรมคอมพิวเตอร์ เช่น โปรแกรม SPSS
2. วิธีการแทนค่า (Imputation based procedures) วิธีนี้ใช้การแทนข้อมูลสูญหายด้วยค่าที่ได้จากวิธีการต่างๆ เช่น การแทนค่าแบบฮอตเดค (hot deck imputation) การแทนค่าโดยใช้ค่าเฉลี่ย (mean imputation) และการแทนค่าโดยวิธีการถดถอย (regression imputation)
3. วิธีการถ่วงน้ำหนัก (Weighting procedures) เป็นวิธีที่ใช้การถ่วงน้ำหนักนำไปปรับข้อมูลที่สูญหาย การถ่วงน้ำหนักมีความสัมพันธ์กับการแทนค่าด้วยค่าเฉลี่ย ถ้าน้ำหนักที่กำหนดไว้เป็นค่าคงที่ในกลุ่มตัวอย่างย่อยแล้วทั้งการแทนค่าด้วยค่าเฉลี่ยของกลุ่มย่อย และการให้น้ำหนักหน่วยที่ตอบเป็นสัดส่วนของการตอบในแต่ละกลุ่มย่อย จะทำให้การแทนค่าข้อมูลสูญหายเหมือนกัน
4. วิธีการที่ได้จากการนิยามโมเดล (model-based procedures) วิธีนี้ได้จากการนิยามโมเดลของข้อมูลสูญหายบางส่วนและใช้หลักการการอ้างถึงเกี่ยวกับไลค์ลิฮูด (likelihood) ในโมเดลด้วยวิธีการประมาณค่าพารามิเตอร์ที่เรียกว่าแมกซิมัมไลค์ลิฮูด (maximum likelihood) โดยใช้วิธีการทำซ้ำ เช่น วิธีอีเอ็ม (EM : expectation maximization) เอฟไอเอ็มแอล (FIML : full information maximization likelihood estimation) และวิธีเอ็มไอ (MI : multiple imputation)

วิธีที่ค่อนข้างจะนิยมใช้ในการจัดการข้อมูลสูญหาย มีอยู่ 7 วิธี คือ

1. การตัดข้อมูลสูญหายออกแบบลิสต์ไวส์ (Listwise deletion) วิธีนี้จะตัดข้อมูลที่มีข้อมูลสูญหายออกไปนำข้อมูลที่สมบูรณ์เท่านั้นไปวิเคราะห์ ข้อดีของวิธีนี้คือเป็นวิธีที่ง่ายในการนำไปใช้แต่ผลลัพธ์ที่ได้จากการวิเคราะห์มีอคติ การจัดการข้อมูลสูญหายแบบนี้ทำให้เกิดปัญหาดังต่อไปนี้

- 1.1 ลดอำนาจการทดสอบทางสถิติ เพราะว่ามีขนาดของกลุ่มตัวอย่างน้อยลง

- 1.2 ไม่สามารถทราบสาเหตุของการสูญหายได้ ถึงแม้ว่าจะให้ค่าประมาณที่ไม่มีอคติเมื่อสาเหตุของการสูญหายเป็น missing completely at random วิธีนี้ดูเหมือนว่าเป็นวิธีที่ดีแต่ก็ไม่แนะนำให้ใช้

2. การตัดข้อมูลสูญหายแบบแพร์ไวด์ (pairwise deletion) วิธีนี้จะตัดหน่วยวิเคราะห์ที่มีข้อมูลสูญหายออกไปเมื่อนำข้อมูลของตัวแปรนั้นมาวิเคราะห์ ผลที่ได้ในการวิเคราะห์ความสัมพันธ์จะได้มาจากกลุ่มตัวอย่างที่แตกต่างกัน วิธีนี้ให้การประมาณค่าความสัมพันธ์ที่สุดเพราะใช้ข้อมูลทั้งหมดที่เก็บรวบรวมมา ถ้าข้อมูลมีการสูญหายแบบสุ่ม (random missing data) ปัญหาที่สำคัญคือ เมตริกซ์สหสัมพันธ์จะไม่เป็นบวกแน่นอน (non-positive definite) ซึ่งจะพบได้ในการวิเคราะห์โดยใช้โปรแกรมลิสเรล (LISREL) ค่าความสัมพันธ์ที่ได้จากการตัดข้อมูลสูญหายแบบแพร์ไวด์จะมีค่ามากกว่าหรือน้อยกว่าค่าความสัมพันธ์ที่ได้จากกลุ่มตัวอย่างทั้งหมดและมีโปรแกรมที่ใช้ในการวิเคราะห์ข้อมูลด้วยวิธีการตัดข้อมูลออกด้วยวิธีแพร์ไวด์ วิธีนี้จะทำให้สูญเสียอำนาจการทดสอบน้อยกว่าวิธีลิสท์ไวด์

3. การแทนค่าข้อมูลสูญหายด้วยค่าเฉลี่ย (mean substitution) เป็นวิธีที่ใช้กันมากในการจัดการข้อมูลสูญหาย วิธีนี้จะแทนข้อมูลสูญหายด้วยค่าเฉลี่ยจากตัวแปรที่มีข้อมูลสมบูรณ์ เช่น ถ้าค่าเฉลี่ยรายได้มีค่าเท่ากับ 52,140 แล้วทุก ๆ คน ที่ไม่ได้รายงานค่าเฉลี่ยจะถูกกำหนดให้มียาได้ เป็น 52,140 กรณีที่ข้อมูลมีการแจกแจงเป็นโค้งปกติใช้ค่าเฉลี่ยจะเหมาะสมแต่ถ้าข้อมูลมีลักษณะการแจกแจงแบบเบ้ควรแทนค่าข้อมูลสูญหายด้วยมัธยฐานจะดีกว่า

ถ้าข้อมูลมีการแจกแจงแบบปกติและเป็นข้อมูลสูญหายแบบสุ่มแล้วการแทนค่าแบบนี้จะไม่มีอคติ อย่างไรก็ตามทุกคนที่ไม่ตอบจะถูกกำหนดให้มีคะแนนเท่ากับค่าเฉลี่ย ซึ่งจริง ๆ แล้วน่าจะมีข้อมูลที่แตกต่างกัน ดังนั้นความแปรปรวนและความแปรปรวนร่วมระหว่างตัวแปรที่มีข้อมูลสูญหายจะมีค่าลดลง ซึ่งทำให้การประมาณค่า R^2 และ β น้อยลง นักวิจัยบางคนแนะนำว่าไม่ควรใช้วิธีนี้มีวิธีการแทนค่าที่มีประโยชน์อย่างชัดเจนมากกว่า เช่น การแทนข้อมูลสูญหายด้วยค่าเฉลี่ยของกลุ่มย่อย (group mean substitution)

4. การแทนค่าข้อมูลสูญหายด้วยค่าเฉลี่ยของกลุ่มย่อย (group mean substitution) วิธีนี้พัฒนามาจากวิธีการแทนค่าข้อมูลสูญหายด้วยค่าเฉลี่ย โดยการแทนค่าด้วยค่าเฉลี่ยจากกลุ่มตัวอย่างที่มีลักษณะคล้ายกับหน่วยข้อมูลที่มีข้อมูลสูญหาย เช่น ถ้ารายได้ของพ่อแม่มีข้อมูลสูญหายจะไม่แทนค่าข้อมูลสูญหายด้วยค่าเฉลี่ยของข้อมูลทั้งหมด แต่จะแทนด้วยค่าเฉลี่ยของกลุ่มย่อยที่มีลักษณะอื่น ๆ เหมือนกันกับตัวแปรที่มีข้อมูลสูญหาย ความแปรปรวนของตัวแปรที่มีข้อมูลสูญหายจะลดน้อยลงแต่ดีกว่าการแทนค่าข้อมูลสูญหายด้วยค่าเฉลี่ยเพราะข้อมูลสูญหายจะถูกแทนค่าด้วยค่าที่ไม่เหมือนกัน แม้กระนั้นก็ตามวิธีนี้ทำให้ความแปรปรวนและความแปรปรวนร่วมลดลงอย่างมีอคติ

วิธีนี้มีข้อดีในการแทนค่าข้อมูลสูญหายได้ถูกต้องมากกว่าแทนค่าด้วยค่าเฉลี่ยทั้งหมด เพราะค่าที่ได้จะขึ้นอยู่กับลักษณะของกลุ่มย่อยที่มีตัวแปรอื่น ๆ เหมือนกันกับหน่วยข้อมูลที่มีข้อมูลสูญหาย และวิธีนี้ยังรักษาความแปรปรวนของตัวแปรที่มีข้อมูลสูญหาย แต่ถ้าลักษณะของตัวแปรแตกต่างจากหน่วยข้อมูลที่มีข้อมูลสูญหายก็จะไม่แทนค่าข้อมูลสูญหาย

5. การแทนค่าข้อมูลสูญหายโดยใช้วิธีการถดถอย (regression imputation) วิธีนี้ใช้การวิเคราะห์การถดถอยเพื่อสร้างสมการทำนายข้อมูลสูญหายจากข้อมูลที่สมบูรณ์ทั้งหมด สุนันทา วีรกุลเทวัญ (2544. หน้า 20) กล่าวว่า ถึงแม้ว่าวิธีการแทนค่าด้วยวิธีการถดถอยจะให้ความแปรปรวนของข้อมูลมีค่ามากกว่าความแปรปรวนที่ได้จากวิธีการแทนค่าข้อมูลสูญหายด้วยค่าเฉลี่ย แต่ก็มีข้อเสียคือจะทำให้ค่าที่เกี่ยวข้องกับความสัมพันธ์ (association) ระหว่างตัวแปรบิดเบือน (distort) ไป เพราะสมการประมาณค่าข้อมูลสูญหายสร้างมาจากความสัมพันธ์ระหว่างตัวแปร

6. การแทนค่าข้อมูลสูญหายด้วยวิธีอีเอ็ม (EM algorithm) เป็นวิธีการแทนค่าข้อมูลสูญหายโดยใช้วิธีการทำซ้ำ ตัวอย่างวิธีการแทนค่าข้อมูลสูญหายด้วยวิธีการอีเอ็ม

ตาราง 1 ข้อมูลการตอบแบบสอบถาม

คนที่	V1	V2	V3	V4
1.	5	4	3	2
2.	-1	3	2	1
3.	2	-1	4	5
4.	-1	2	-1	3
5.	2	2	-1	-1
6.	5	4	3	2
7.	3	2	1	1
8.	3	2	5	-1

ขั้นตอนของวิธีการอีเอ็มมีดังนี้

ขั้นตอนที่ 1 คอลัมน์แรกตัวแปร V1 ข้อมูลสูญหายคือ -1 ของคนที่ 2 สามารถแทนค่าข้อมูลสูญหายนี้โดยใช้ข้อมูลของคนี่ 2 ที่มีอยู่ คือ V2 V3 และ V4 การประมาณค่าทำ

ได้หลายวิธี ถ้าเราสมมุติว่าข้อมูลมีการแจกแจงแบบปกติและมีความสัมพันธ์เชิงเส้นตรง (normality and linearity) ตัวทำนายที่ดีที่สุดคือ การถดถอยเชิงเส้นตรง ซึ่งมีโปรแกรมสำเร็จรูปที่เป็นมาตรฐานใช้ในการประมาณค่าได้จากสมการ

$$V1 = B_0 + B_1 V_2 + B_2 V_3 + B_3 V_4$$

คนที่ 4 มีข้อมูลสูญหายในตัวแปรที่ 1 สามารถประมาณค่าข้อมูลสูญหายได้จากสมการ

$$V1 = B_0 + B_1 V_2 + B_2 V_4$$

กรณีของคนที่ 4 จะไม่ใช้ตัวแปร V_3 เพราะว่าคนที่ 4 ไม่ได้ตอบข้อที่ 3 ดังนั้นเราสามารถทำแบบนี้ไปเรื่อย ๆ ในแต่ละคอลัมน์จนกระทั่งแทนค่าข้อมูลสูญหายได้ทั้งหมด

วิธีนี้มีความคลาดเคลื่อนแบบสุ่มรวมอยู่ในค่าที่ถูกทำนาย ถ้าสมการการถดถอยให้ผลการทำนายที่ดีก็จะมี ความคลาดเคลื่อนน้อย แต่ถ้าสมการถดถอยให้ผลการทำนายที่ไม่ดี จะมีความคลาดเคลื่อนมาก ขนาดของความคลาดเคลื่อนจะเกี่ยวข้องกับสมการทำนายว่าจะทำนายได้ดีหรือไม่ วิธีที่ง่ายที่จะพิจารณาความคลาดเคลื่อนก็คือพิจารณาจากค่าที่เหลือ (residual) ของตัวแปร V_1 เมื่อ V_1 ไม่มีค่าสูญหาย แล้วเลือกความคลาดเคลื่อนขึ้นมา 1 ค่าอย่างสุ่มเพื่อนำไปรวมหรือลบออกจากค่าที่ถูกแทน

จากขั้นตอนที่กล่าวมาจะได้ชุดข้อมูล (data matrix) ที่ไม่มีข้อมูลสูญหาย แล้วจึงคำนวณความแปรปรวนร่วมจากชุดข้อมูลดังกล่าวได้

ขั้นตอนที่ 2 พิจารณาคนที่ 4 แทนที่จะใช้ตัวแปร V_2 และ V_4 เพื่อทำนาย V_1 จะแทนค่าตัวแปร V_3 ก่อน แล้วจึงสร้างสมการเพื่อทำนาย V_1 ดังนั้นในขั้นตอนที่ 2 จะสร้างสมการพยากรณ์ใหม่โดยใช้ข้อมูลที่สมบูรณ์จากขั้นตอนที่ 1 การประมาณค่าข้อมูลสูญหายในขั้นตอนที่ 2 จะประมาณค่าได้ดีกว่าเพราะว่าใช้ข้อมูลมากขึ้นหลังจากขั้นตอนที่ 2 ก็จะมีข้อมูลใหม่ และเมตริกซ์ความแปรปรวนร่วม สามารถทดสอบได้ว่าขั้นตอนที่ 1 และขั้นตอนที่ 2 เพิ่มขึ้นอย่างมีนัยสำคัญหรือไม่ โดยการเปรียบเทียบความแปรปรวนร่วมในขั้นตอนที่ 1 และขั้นตอนที่ 2 ถ้าแตกต่างกันอย่างมีนัยสำคัญขั้นตอนที่ 2 ก็จะนำไปใช้ประโยชน์ได้

ขั้นตอนที่ 3 สร้างสมการพยากรณ์ใหม่ก็จะได้เมตริกซ์ความแปรปรวนร่วมใหม่ถ้าเมตริกซ์ในขั้นตอนที่ 3 ยังแตกต่างกันอย่างมีนัยสำคัญจากขั้นตอนที่ 2 ก็จะทำขั้นตอนที่ 4 วิธีการประมาณค่าจะจบเมื่อเมตริกซ์ความแปรปรวนร่วมทั้งสองครั้งไม่แตกต่างกัน และไม่สามารถดำเนินการต่อไปได้อีก ดังนั้นก็จะได้ข้อมูลสูญหายโดยวิธีการทำซ้ำ

วิธีการประมาณค่าด้วยวิธีอีเอ็ม มีประเด็นที่นักสถิติพิจารณาก็คือความคลาดเคลื่อนมาตรฐานไม่ตรง (invalid) คือ ไม่ได้คำนวณอย่างตรงไปตรงมาในการวิเคราะห์ขั้นสุดท้าย ดังนั้นจะมีวิธีการอื่น ๆ ที่มีความตรงในการประมาณค่าความคลาดเคลื่อนมาตรฐาน วิธีการประมาณค่าข้อมูลสูญหายวิธีนี้เป็นวิธีที่นักสถิติแนะนำในการแก้ปัญหาข้อมูลสูญหายเพราะใช้วิธีการถดถอยซึ่งน่าจะแทนค่าข้อมูลสูญหายได้ถูกต้องมากยิ่งขึ้น

7. การแทนค่าข้อมูลสูญหายด้วยวิธีเอ็มไอ (multiple imputation) วิธีนี้พัฒนามาจากวิธีอีเอ็ม (EM algorithm) แต่การทำซ้ำจะทำให้เกิดกลุ่มที่แตกต่างของค่าที่ถูกแทน ข้อมูลที่แตกต่างกันนี้ใช้ในการประมาณค่าความคลาดเคลื่อนในโมเดล วิธีนี้จะมีประโยชน์เพราะว่าการประมาณค่าความคลาดเคลื่อนมาตรฐานพิสูจน์ได้โดยวิธีการทางคณิตศาสตร์แต่ขั้นตอนการทำงานต้องใช้เวลามาก

การจัดการข้อมูลสูญหายด้วยวิธีเอ็มไอแตกต่างจากวิธีอีเอ็ม 2 ประการคือ

1. การบันทึกความแปรปรวนที่สูญหายไป ความแปรปรวนที่ได้จากการประมาณค่าด้วยวิธีการถดถอยจะสูญหายไปเพราะว่าค่าที่ถูกแทนแต่ละตัวได้จากการประมาณค่าโดยปราศจากความคลาดเคลื่อน ถึงแม้ว่าไม่มีข้อมูลสูญหายนักวิจัยก็ทราบว่าการประมาณค่าที่ได้จากวิธีการถดถอยจะประมาณค่าได้สูงกว่าหรือต่ำกว่าค่าจริง ความแตกต่างระหว่างค่าจริงกับค่าที่ได้จากการทำนายเป็นความคลาดเคลื่อน และมีเหตุผลที่จะสมมุติว่า การแจกแจงความคลาดเคลื่อนของข้อมูล ที่ไม่สูญหายสามารถอธิบายการแจกแจงความคลาดเคลื่อนของข้อมูลสูญหายได้ ดังนั้นวิธีที่จะเก็บความแปรปรวนของค่าที่ถูกแทนคือสุ่มหน่วยข้อมูลขึ้นมา 1 กรณี จากการแจกแจงความคลาดเคลื่อนของข้อมูลที่มีค่าสมบูรณ์และรวมไปกับค่าที่ถูกแทนแต่ละตัว

2. เนื่องจากส่วนของความแปรปรวนอื่น ๆ ในการแทนค่าด้วยวิธีการถดถอยจะสูญหายไป เพราะค่าที่ถูกแทนจากการประมาณค่าเพียงครั้งเดียวของเมตริกซ์ความแปรปรวนร่วมซึ่งประมาณค่าตัวของมันเองด้วยความคลาดเคลื่อน ดังนั้น การแทนค่าข้อมูลสูญหายด้วยวิธีเอ็มไอจะสร้างเมตริกซ์ความแปรปรวนร่วมโดยใช้วิธีบูทสทราฟ (Bootstrap) และประมาณค่าเมตริกซ์ความแปรปรวนร่วมโดยใช้วิธีการอีเอ็ม

จุดแข็งและจุดอ่อนของวิธีการจัดการข้อมูลสูญหาย

วิธีการจัดการข้อมูลสูญหายที่กล่าวมามีจุดแข็งและจุดอ่อนดังต่อไปนี้

ตาราง 2 จุดอ่อนจุดแข็งของวิธีการจัดการข้อมูลสูญหาย

วิธีการ	จุดแข็ง	จุดอ่อน
การตัดข้อมูลออกแบบลิสต์ไวส์	<ul style="list-style-type: none"> - ทำให้ได้เมตริกซ์ความสัมพันธ์ที่เป็นบวกแน่นอน - การวิเคราะห์ใช้กลุ่มตัวอย่างกลุ่มเดียวกัน - ใช้ง่ายและมีอยู่ในโปรแกรมสำเร็จรูปทั่วไป 	<ul style="list-style-type: none"> - มีอำนาจการทดสอบน้อย - ลดความแปรปรวนของข้อมูล
การตัดข้อมูลออกแบบแพร์ไวส์	<ul style="list-style-type: none"> - ใช้ข้อมูลทั้งหมดในการวิเคราะห์ - มีอำนาจการทดสอบมากกว่าการตัดข้อมูลแบบลิสต์ไวส์ - ใช้ง่ายและมีอยู่ในโปรแกรมโปรแกรมสำเร็จรูปทั่วไป 	<ul style="list-style-type: none"> - ทำให้เกิดเมตริกซ์ความสัมพันธ์ที่ไม่เป็นบวกแน่นอน - ค่าความสัมพันธ์ได้มาจากกลุ่มตัวอย่างคนละกลุ่ม - ให้การประมาณค่าดีกว่าวิธีลิสต์ไวส์ถ้าข้อมูลสูญหายแบบสุ่ม
การแทนค่าข้อมูลสูญหายด้วยค่าเฉลี่ย	<ul style="list-style-type: none"> - ง่ายในการคำนวณ 	<ul style="list-style-type: none"> - ลดความแปรปรวนของข้อมูลมีผลต่อความแปรปรวนร่วม และความสัมพันธ์

ตาราง 2 (ต่อ)

วิธีการ	จุดแข็ง	จุดอ่อน
การแทนค่าข้อมูลสูญหายด้วยค่าเฉลี่ยของกลุ่มย่อย	<ul style="list-style-type: none"> - ง่ายในการคำนวณและมีเหตุผลในการแทนค่า - รักษาความแปรปรวนของข้อมูลมีคติน้อยกว่าการแทนค่าข้อมูลสูญหายด้วยค่าเฉลี่ย 	<ul style="list-style-type: none"> - ลดความแปรปรวนของข้อมูล
การแทนค่าข้อมูลสูญหายโดยใช้การวิเคราะห์การถดถอย	<ul style="list-style-type: none"> - ทำนายค่าสูญหายได้ดี 	<ul style="list-style-type: none"> - ให้ผลการทำนายที่มากเกินไป - ทำให้เกิดปัญหา multicollinearity - ลดความแปรปรวนของข้อมูล
การแทนค่าข้อมูลสูญหายด้วยวิธีอีเอ็ม	<ul style="list-style-type: none"> - มีอำนาจการทดสอบสูงและการแก้ปัญหาที่ดีเพราะใช้วิธีการทำซ้ำ - ใช้ข้อมูลทั้งหมดวิเคราะห์ 	<ul style="list-style-type: none"> - ใช้โปรแกรมที่เฉพาะ
การแทนค่าข้อมูลสูญหายด้วยวิธีเอ็มไอ	<ul style="list-style-type: none"> - มีอำนาจการทดสอบสูงและการแก้ปัญหาที่ดีเพราะใช้วิธีการทำซ้ำ - ประมาณค่าความคลาดเคลื่อนมาตรฐานได้ดี 	<ul style="list-style-type: none"> - ใช้โปรแกรมที่เฉพาะ - การวิเคราะห์ใช้เวลานาน

การสุ่มตัวอย่าง

การสำรวจเพื่อหาข้อมูลทางสถิติระดับใหญ่ซึ่งเรียกว่า (Large Scale Statistical) เช่น การสำรวจแรงงานของประชาชน การสำรวจเพื่อหาข้อมูลทางเศรษฐกิจและสังคม การสำรวจผล

สัมฤทธิ์ทางการเรียนของนักเรียน ฯลฯ แบ่งระเบียบวิธีการสำรวจ (Survey methodology) ออกเป็น 2 วิธีใหญ่ ๆ (นิยม ปุราคำ, 2517) คือ

1. การสำรวจด้วยวิธีการแจงนับอย่างครบถ้วน หรือการสำมะโน (Complete enumeration) ตัวอย่างการสำรวจลักษณะเช่นนี้ เช่น การเก็บรวบรวมเกี่ยวกับจำนวนและคุณลักษณะของประชากรและเคหสถาน ซึ่งเรียกกันว่า สำมะโนประชากรและเคหะ ซึ่งรัฐบาลจัดทำโดยจัดส่งพนักงานแจงนับออกไปยังครัวเรือนทุกครัวเรือน เพื่อสอบถามข้อมูลต่าง ๆ เกี่ยวกับประชากรทุกคนและบ้านที่กลงในแบบสำมะโนที่กำหนดขึ้น การเก็บโดยวิธีนี้ต้องลงทุนคือใช้กำลังคนและงบประมาณมาก ซึ่งอาจทำให้เกิดปัญหาในการหาปัจจัยเหล่านี้ให้เพียงพอแก่ความต้องการของการสำรวจนั้น ถึงแม้ว่าจะหางบประมาณมาได้แต่ก็มีปัญหาเรื่องเวลาเพราะการสำรวจแบบนี้ต้องใช้เวลาอันยาวนานข้อมูลที่ได้จึงไม่ทันสมัย ไม่ทันต่อความต้องการของผู้ที่ต้องการข้อมูล ข้อเสียของการสำรวจสถิติโดยการสำมะโนมีข้อเสียอยู่หลายอย่าง คือ ใช้กำลังคนมาก ใช้งบประมาณมาก ใช้เวลามาก ยากในการควบคุมการทำงานของผู้ที่เก็บรวบรวมข้อมูลเป็นผลให้ข้อมูลที่ได้อาจผิดพลาดได้มาก

2. การสำรวจโดยการแจงนับเพียงบางส่วน หรือ การสุ่มสำรวจ (Partial enumeration method or Sampling survey) วิธีนี้ถูกพัฒนาขึ้นมาเพื่อแก้ปัญหาและข้อเสียของวิธีการสำรวจแบบแรกที่ได้กล่าวมาแล้ว โดยอาศัยทฤษฎีทางคณิตศาสตร์สถิติ (Mathematical statistics) เกี่ยวกับความน่าจะเป็น (Probability) การแจกแจงตัวแปรสุ่ม (Distribution of random variables) ประกอบกับการวางแผนการสำรวจ (Survey design) และการประมาณผล (Estimation) วิธีนี้จะเลือกตัวแทนของประชากรที่เราต้องการศึกษาขึ้นมาจำนวนหนึ่งแล้วใช้ข้อมูลที่ได้จากตัวอย่าง ประมาณตัวเลขที่ต้องการโดยอาศัยหลักการประมาณผลทางสถิติที่เหมาะสม ทำให้ได้ค่าประมาณของตัวเลขที่ใกล้เคียงกับค่าที่ได้จากประชากรทั้งหมด ข้อดีของวิธีนี้คือใช้กำลังคนน้อย ใช้งบประมาณน้อย ใช้เวลาในงานสนามและการประมวลผลน้อย สามารถควบคุมคุณภาพของข้อมูลได้ ทำให้ศึกษาและวิเคราะห์ปัญหาได้ลึกและสมบูรณ์กว่าการศึกษาข้อมูลจำนวนมาก

การสำรวจข้อมูลจากกลุ่มตัวอย่างเป็นเครื่องมือที่สำคัญที่สุดในการหาข้อมูลสถิติต่าง ๆ ของรัฐบาลและเอกชน เช่น การสำรวจข้อมูลทางด้านการเกษตร อุตสาหกรรม สาธารณสุข การคมนาคม การศึกษา และข้อมูลทางเศรษฐกิจและสังคมอื่น ๆ รวมทั้งการหยั่งเสียงประชามติ การวิจัยตลาด มีข้อควรระวังในการศึกษาข้อมูลจากกลุ่มตัวอย่างดังนี้

2.1 ถ้าการสุ่มตัวอย่างนั้นไม่ได้รับการออกแบบวิธีการสุ่มตัวอย่าง และนำไปปฏิบัติอย่างถูกต้องแล้วผลที่ได้จะมีความคลาดเคลื่อน และนำไปสู่การสรุปที่ผิดพลาดได้

2.2 การสุ่มตัวอย่างต้องอาศัยเทคนิคที่เฉพาะของแต่ละวิธี และมีข้อจำกัดในแต่ละวิธีจะต้องคำนึงถึงความเหมาะสมกับเนื้อเรื่องที่ต้องทำการศึกษาในด้านด้วย

2.3 การเก็บรวบรวมข้อมูลจากกลุ่มตัวอย่างจะต้องเตรียมรายละเอียดเกี่ยวกับประชากร เช่น บัญชีรายชื่อ การจำแนกแยกประเภทประชากร หน่วยการสุ่ม วิธีการสุ่ม เป็นต้น

2.4 ถ้าใช้ขนาดตัวอย่างไม่มากพอจะมีปัญหาในแง่ของการสรุปอ้างอิงไปยังประชากรตลอดจนไม่สามารถประมาณค่าตัวเลขในระดับย่อย ๆ ได้

มโนทัศน์พื้นฐานที่สำคัญ

มโนทัศน์ที่สำคัญที่เกี่ยวกับการสุ่มตัวอย่างมีดังนี้ (ศิริชัย กาญจนวาสี, 2533)

1. ประชากร (Population) หมายถึง กลุ่มของสิ่งต่าง ๆ ทั้งหมดที่สนใจศึกษาซึ่งอาจจะเป็นคน สิ่งของหรือเหตุการณ์ต่าง ๆ ค่าที่คำนวณได้จากข้อมูลทุกหน่วยของประชากร เรียกว่าค่าพารามิเตอร์ (Parameter) จึงเป็นค่าที่แท้จริงและถือว่าเป็นค่าที่ถูกต้อง เช่น ค่าเฉลี่ยเลขคณิต (μ) ส่วนเบี่ยงเบนมาตรฐาน (σ) สัดส่วน (π) สหสัมพันธ์ (ρ) เป็นต้น

2. กลุ่มตัวอย่าง (Sample) หมายถึง ส่วนหนึ่งของประชากรที่ถูกเลือกหรือสุ่มเพื่อใช้ศึกษาแทนประชากร ค่าที่คำนวณได้จากข้อมูลกลุ่มตัวอย่างเรียกว่า ค่าสถิติ (Statistics) จึงเป็นค่าที่เป็นจริงเฉพาะกลุ่มตัวอย่างนั้น โดยทั่วไปเป็นค่าที่มีความคลาดเคลื่อนจากค่าพารามิเตอร์

3. การสุ่มตัวอย่าง (Sampling Techniques) หมายถึง การกระทำเพื่อให้ได้มาซึ่งกลุ่มตัวอย่างเพื่อใช้ศึกษาแทนประชากร การออกแบบการสุ่มตัวอย่างที่ดีประกอบด้วยวิธีการสุ่มตัวอย่างที่ดี คือไม่ลำเอียง ปราศจากอคติ และขนาดของกลุ่มตัวอย่างที่เหมาะสม (ขนาดกลุ่มตัวอย่างที่ถูกต้องตามหลักสถิติ คือ มีความเป็นตัวแทนของประชากร ตลอดจนเป็นขนาดกลุ่มตัวอย่างที่พอเหมาะตามหลักปฏิบัติ คือ มีความเป็นไปได้ในการติดตามเก็บรวบรวมข้อมูล)

4. การสรุปอ้างอิงทางสถิติ (Statistical Inference) หมายถึง การวิเคราะห์ทางสถิติเพื่อใช้ค่าสถิติที่คำนวณได้จากกลุ่มตัวอย่างไปทำการประมาณค่า (Estimation) หรือทดสอบสมมติฐาน (Hypothesis testing) เกี่ยวกับค่าพารามิเตอร์ของประชากร

ขั้นตอนการสุ่มตัวอย่าง

ในทางปฏิบัติแล้วผู้วิจัยอาจวางแผนในการที่จะเลือกสุ่มกลุ่มตัวอย่างที่ดีได้ตามลำดับขั้น

ดังนี้



1. พิจารณาความมุ่งหมายของการวิจัย หรือการเลือกกลุ่มตัวอย่าง ผู้วิจัยจะต้อง**สำนักหอสมุด** วิเคราะห์จุดมุ่งหมายของการวิจัยอย่างละเอียดและให้ชัดเจนในจุดมุ่งหมาย หรือปัญหาของการวิจัยเพื่อเป็นการหาลักษณะของกลุ่มตัวอย่างหรือประชากรที่ถูกต้องที่สุดว่าจะศึกษากับใคร มีคุณสมบัติเช่นไร

2. ให้คำจำกัดความของประชากรที่จะใช้เพื่อสนองจุดมุ่งหมายของการวิจัยนั้น ผู้วิจัยจะต้องให้ความหมายของประชากรได้ว่าหมายถึงใคร ครอบคลุมกว้างขวางหรือมีขอบเขตแคไหน มีคุณลักษณะที่จำกัดอย่างไร เช่น อาจหมายถึงนักเรียนที่เรียนชั้นมัธยมศึกษาปีที่ 3 ทั่วประเทศ หรือเฉพาะนักเรียนที่เรียนในโรงเรียนสาธิตของมหาวิทยาลัยหรือวิทยาลัยครูเท่านั้น เป็นต้น

3. การกำหนดหน่วยตัวอย่าง (Sampling unit) ก่อนที่จะเลือกกลุ่มตัวอย่างผู้วิจัยจะต้องกำหนดหน่วยตัวอย่างก่อน โดยพิจารณาจากคุณลักษณะของประชากรที่เราจะศึกษาส่วนใหญ่แล้วหน่วยของตัวอย่งนี้มักจะหมายถึง หน่วยหรือส่วนที่ย่อยที่สุดของประชากรในการเลือกแต่ละครั้ง ซึ่งการเลือกหน่วยที่ย่อยที่สุดนี้ บางครั้งก็ทำได้ง่ายบางครั้งก็ทำได้ยาก ทั้งนี้ขึ้นอยู่กับหน่วยของตัวอย่งนั้นว่าชัดเจนและสะดวกในการเลือกแคไหน ฉะนั้น ก่อนเลือกกลุ่มตัวอย่างผู้วิจัยจะต้องกำหนดหน่วยของตัวอย่างให้ได้เสียก่อน

4. กำหนดขอบข่ายของประชากร การกำหนดขอบข่ายของประชากร ก็คือ การรวบรวมรายชื่อหรือทำบัญชีหน่วยของตัวอย่าง (Sampling unit) หรือเป็นการหาขนาดของประชากรนั่นเอง

5. ประมาณขนาดของตัวอย่าง (Estimated sample size) ดังที่กล่าวมาแล้วว่าการเลือกกลุ่มตัวอย่างนั้นต้องคำนึงถึงว่ากลุ่มตัวอย่างที่ได้นั้นจะต้องเป็นตัวแทนที่ดีของประชากร ฉะนั้น การเลือกกลุ่มตัวอย่างมาศึกษานั้นจะต้องให้มีขนาดหรือจำนวน ที่เป็นสัดส่วนพอเหมาะ กับขนาดของประชากร

6. เลือกวิธีการสุ่มตัวอย่างที่เหมาะสม (Sampling techniques) การเลือกกลุ่มตัวอย่างที่ดีและเหมาะสมนั้นขึ้นอยู่กับผู้วิจัยมีความรู้เกี่ยวกับประชากรมากน้อยเพียงใด ซึ่งถ้าเราทราบเรื่องของประชากรที่เราจะวิจัยมาก โอกาสที่เราจะเลือกกลุ่มตัวอย่างได้ตรงกับวัตถุประสงค์ของการวิจัยก็มากขึ้น

การได้ตัวอย่างเพื่อให้เป็นตัวแทนที่ดีของประชากร ก็จำเป็นที่จะต้องคำนึงถึงวิธีการสุ่มตัวอย่างและขนาดของกลุ่มตัวอย่างที่เหมาะสม สำหรับวิธีการสุ่มตัวอย่างสามารถแบ่งได้เป็น 2 ประเภท (สมชัย วงษ์นายะ, 2533. หน้า 12-21) คือ

1. การสุ่มแบบไม่ใช้ความน่าจะเป็น (Nonprobability Sampling)
2. การสุ่มแบบใช้ความน่าจะเป็น (Probability Sampling)

การสุ่มแบบไม่ใช้ความน่าจะเป็น (Nonprobability Sampling)

เป็นการเลือกกลุ่มตัวอย่างที่ไม่ทราบว่าสมาชิกแต่ละหน่วยในประชากรมีโอกาสได้รับเลือกเท่าไร สามารถแบ่งออกได้เป็น 4 ประเภทย่อย ๆ คือ

1. วิธีการสุ่มแบบบังเอิญ (Accidental Sampling) เป็นวิธีการเลือกกลุ่มตัวอย่างโดยเลือกเก็บข้อมูลเฉพาะหน่วยของประชากรที่พบ
2. วิธีการสุ่มตามวัตถุประสงค์ (Purposive Sampling) เป็นวิธีการเลือกกลุ่มตัวอย่างที่ผู้วิจัยตัดสินใจเลือกหน่วยต่าง ๆ ของประชากรตามที่ตนเองเห็นสมควร โดยเลือกให้สอดคล้องกับวัตถุประสงค์ของการวิจัย
3. วิธีการสุ่มตามโควตา (Quota Sampling) เป็นวิธีการสุ่มกลุ่มตัวอย่างที่ผู้วิจัยกำหนดลักษณะและจำนวนหน่วยตัวอย่างตามที่ตนเองเห็นสมควร เช่น กำหนดว่าเป็นนักเรียนชายและนักเรียนหญิงกลุ่มละเท่าไร หลังจากนั้นก็ใช้การสุ่มแบบบังเอิญจนได้จำนวนครบตามต้องการ
4. วิธีการสุ่มแบบเชือกก่อนหิมะ (Snowball Sampling) หรือเรียกว่าวิธีการสุ่มแบบลูกโซ่ (Chain Sampling) เป็นวิธีการสุ่มกลุ่มตัวอย่างที่ผู้วิจัยเลือกกลุ่มตัวอย่างที่มีลักษณะที่ต้องการมาจำนวนหนึ่งที่มีลักษณะตรงตามวัตถุประสงค์ของการวิจัย หลังจากนั้นก็ให้ผู้ถูกเก็บรวบรวมข้อมูลเสนอรายชื่อคนอื่น ๆ ที่มีลักษณะดังกล่าวเป็นขั้นที่ 2 ต่อไปก็ทำเช่นเดิมอีกเป็นหลาย ๆ ขั้นต่อเนื่องกันไปเป็นลูกโซ่ จนกว่าจะได้กลุ่มตัวอย่างครบตามที่ต้องการ

การสุ่มแบบใช้ความน่าจะเป็น (Probability Sampling)

เป็นวิธีการเลือกกลุ่มตัวอย่างโดยทราบว่าสมาชิกแต่ละหน่วยในประชากรมีโอกาสได้รับเลือกเท่าไร ในกรณีนี้ที่แต่ละหน่วยในประชากรมีโอกาสได้รับเลือกเป็นกลุ่มตัวอย่างโดยเท่าเทียมกันวิธีการเลือกแบบนี้จะเรียกว่าวิธีการสุ่มตัวอย่าง ซึ่งเป็นวิธีที่ยอมรับของนักวิจัย เพราะสามารถควบคุมอคติของผู้วิจัยที่จะเลือกกลุ่มตัวอย่างตามใจชอบเพื่อให้ผลการวิจัยเป็นไปตามความคาดหวังหรือตามสมมุติฐานของการวิจัยที่กำหนดไว้ นอกจากนี้วิธีการสุ่มนี้ก็ดำเนินการตามหลักของทฤษฎีความน่าจะเป็นสามารถใช้สถิติอนุมานได้

วิธีการสุ่มแบบใช้ความน่าจะเป็น สามารถแบ่งออกได้เป็น 5 วิธี คือ

1. วิธีการสุ่มอย่างง่าย (Simple Random Sampling)
2. วิธีการสุ่มแบบมีระบบ (Systematic Sampling)
3. วิธีการสุ่มแบบแบ่งชั้น (Stratified Random Sampling)
4. วิธีการสุ่มตามกลุ่ม (Cluster Sampling)
5. วิธีการสุ่มแบบหลายขั้นตอน (Multi-Stage Sampling)

สำหรับรายละเอียดของวิธีการสุ่มแบบใช้ความน่าจะเป็นสรุปได้ดังตารางที่ 3



ตาราง 3 รายละเอียดของวิธีการสุ่มแบบใช้ความน่าจะเป็น

วิธีการสุ่ม	ข้อดี/ข้อเสีย
<p>1. วิธีการสุ่มอย่างง่าย</p> <p>เป็นวิธีการสุ่มตัวอย่างที่สมาชิกแต่ละหน่วยในประชากรมีโอกาสได้รับเลือกเท่า ๆ กัน สำหรับวิธีการสุ่มตัวอย่างเพื่อให้ได้กลุ่มตัวอย่างตามการสุ่มอย่างง่ายสามารถทำได้หลายวิธี ได้แก่</p> <p>1.1 ใช้วิธีจับฉลาก</p> <p>1.2 ใช้ตารางเลขสุ่ม</p> <p>1.3 ใช้คอมพิวเตอร์</p>	<p>ข้อดี</p> <ol style="list-style-type: none"> 1. วางแผนการสุ่มได้ง่ายและสะดวกในการสุ่ม 2. เหมาะกับประชากรที่มีลักษณะที่คล้ายคลึงกัน (Homogeneity) 3. การคำนวณหาขนาดของกลุ่มตัวอย่างตลอดจนค่าสถิติต่าง ๆ ใช้สูตรที่ไม่ยุ่งยาก <p>ข้อเสีย</p> <ol style="list-style-type: none"> 1. ทำให้เสียค่าใช้จ่ายในการวิจัยสูง เช่น การเดินทางไปเก็บข้อมูลจากกลุ่มตัวอย่างที่อยู่กระจัดกระจาย 2. เสียเวลาและไม่สะดวกในการเก็บรวบรวมข้อมูล 3. ในกรณีที่สมาชิกในประชากรมีลักษณะเป็นวิวิธพันธ์ (Heterogeneity) อาจทำให้ผลการวิจัยคลาดเคลื่อนได้มา
<p>2. วิธีการสุ่มแบบมีระบบ</p> <p>เป็นวิธีการสุ่มตัวอย่างจากส่วนต่าง ๆ ของประชากร ซึ่งเรียกว่าช่วงของการสุ่ม ส่วนละ 1 หน่วย สำหรับวิธีการสุ่มแบบมีระบบสามารถแบ่งออกได้เป็น 4 ประเภท คือ (สมชัย วงษ์นายะ, 2533. หน้า 15 อ้างอิงจาก โนรี ใจใส, 2531. หน้า 84-95)</p>	<p>ข้อดี</p> <ol style="list-style-type: none"> 1. สะดวกและง่ายในการสุ่มตัวอย่าง 2. ถ้าหน่วยต่าง ๆ ในประชากรเรียงกันอย่างสุ่ม ทำให้การสุ่มตัวอย่างวิธีนี้ให้ผลใกล้เคียงกับวิธีการสุ่มอย่างง่าย

ตาราง 3 (ต่อ)

วิธีการสุ่ม	ข้อดี/ข้อเสีย
<p>2.1 การสุ่มตัวอย่างแบบมีระบบเชิงเส้น (Linear Systematic Sampling)</p> <p>จะใช้ในกรณีที่ช่วงของการสุ่ม (k) ซึ่งได้มาจากการนำจำนวนกลุ่มตัวอย่าง (n) ไปหารจำนวนประชากร (N) มีค่าเป็นเลขจำนวนเต็ม แต่ถ้าช่วงของการสุ่มไม่เป็นเลขจำนวนเต็มจะใช้วิธีการสุ่มแบบมีระบบประเภทต่อไป</p>	<p>ข้อเสีย</p> <p>1. ในกรณีที่สมาชิกในประชากรถูกจัดเรียงไว้ด้วยระบบอย่างใดอย่างหนึ่งบางครั้งเมื่อใช้การสุ่มแบบมีระบบจะทำให้ได้กลุ่มตัวอย่างที่มีลักษณะคล้าย ๆ กัน ไม่เป็นตัวแทนที่ดีของประชากร</p>
<p>2.2 การสุ่มตัวอย่างแบบมีระบบชนิดวงกลม (Circular Systematic Sampling)</p> <p>เป็นวิธีการสุ่มตัวอย่างที่แก้ปัญหาของการสุ่มแบบมีระบบในกรณีที่ $N \neq nk$ โดยมีวิธีการสุ่มตัวอย่างคือ สุ่มหน่วยแรกระหว่าง 1 ถึง N สมมติให้ได้ r จะได้หน่วยตัวอย่างทั้งหมดคือ $(r+jk)$ ถ้า $r+jk < N$ เมื่อ $j=0,1,2,\dots,(n-1)$ $(r+jk-N)$ ถ้า $r+jk > N$ เมื่อ $j=0,1,2,\dots,(n-1)$</p>	
<p>2.3 การสุ่มตัวอย่างแบบมีระบบชนิดใช้ช่วงการสุ่มที่เป็นเศษส่วน</p> <p>เป็นวิธีการสุ่มที่แก้ปัญหาของการสุ่มแบบมีระบบในกรณีที่ $N \neq nk$ โดยมีวิธีการสุ่มตัวอย่างคือ สุ่มหน่วยแรกระหว่าง 1 ถึง k สมมติได้ r จะได้หน่วยตัวอย่างทั้ง n หน่วย โดยพิจารณาจาก $r+jk$ เมื่อ $j=0,1,2,\dots,(n-1)$ ซึ่งมีหลักว่า จะได้หน่วยตัวอย่างที่ i เมื่อ $(i-1) < r \leq i$</p>	

ตาราง 3 (ต่อ)

วิธีการสุ่ม	ข้อดี/ข้อเสีย
<p>2.4 การสุ่มตัวอย่างแบบมีระบบชนิดสุ่มทั้งหน้าและหลัง</p> <p>เป็นวิธีการสุ่มที่แก้ปัญหของวิธีการสุ่มแบบมีระบบในกรณีที่มี $N \neq nk$ โดยมีวิธีการสุ่มคือสุ่มหน่วยเริ่มต้นระหว่าง 1 ถึง N สมมุติได้ r หน่วยตัวอย่าง n หน่วย จะประกอบด้วย $r, r \pm k, r \pm 2k, r \pm 3k, \dots, r \pm (n-1)k \leq N$</p>	
<p>3. วิธีการสุ่มแบบระดับชั้น</p> <p>เป็นวิธีการสุ่มตัวอย่างที่มีการแบ่งประชากรที่จะทำการศึกษาออกเป็นกลุ่ม ๆ ซึ่งเรียกว่าระดับชั้นหรือชั้นภูมิ (Stratum) ตามตัวแปรที่ใช้ในการวิจัยโดยจำแนกให้หน่วยตัวอย่างภายในกลุ่มมีลักษณะคล้ายคลึงกันมากที่สุด (Homogeneous) และต่างกลุ่มกันมีลักษณะที่ต่างกัน (Heterogeneous) แล้วสุ่มตัวอย่างจากแต่ละระดับชั้น ในการสร้างระดับชั้นนั้น จำเป็นที่จะต้องใช้ตัวแปรแบ่งระดับชั้นที่เหมาะสม สำหรับตัวแปรที่ใช้กำหนดระดับชั้นที่ดีที่สุดก็คือตัวแปรที่เรากำลังต้องการศึกษานั้นเอง ในทางปฏิบัติหลักในการหาตัวแปรกำหนดระดับชั้น ก็คือหาตัวแปรที่มีความสัมพันธ์อย่างสูงกับตัวแปรที่สนใจ (สมชัย วงษ์นายะ, 2533. หน้า 17 อ้างอิงจาก สุชาติดา กิระนันท์, 2525)</p>	<p>ข้อดี</p> <ol style="list-style-type: none"> 1. ทำให้สามารถเสนอผลการวิจัยได้ตามจุดมุ่งหมายของการวิจัยที่ต้องการเสนอผลตามตัวแปรอิสระต่าง ๆ ที่ใช้แบ่งระดับชั้น เช่น ตามอาชีพ วุฒิกการศึกษา เป็นต้น 2. ในกรณีที่ใช้ตัวแปรแบ่งระดับชั้นที่มีความสัมพันธ์สูงกับตัวแปรที่ศึกษาทำให้วิธีการสุ่มตัวอย่างตามระดับชั้นได้กลุ่มตัวอย่างที่เป็นตัวแทนที่ดีของประชากรและทำให้ได้ค่าสถิติต่าง ๆ ที่มีความแม่นยำมากขึ้น <p>ข้อเสีย</p> <ol style="list-style-type: none"> 1. ทำให้มีงานเพิ่มขึ้น ทั้งในขั้นการวางแผนการสุ่มตัวอย่าง การเก็บรวบรวมข้อมูล และการวิเคราะห์ข้อมูล 2. ถ้าใช้จำนวนชั้นภูมิมากเกินไป อาจทำให้ไม่สามารถจำกัดขนาดของกลุ่มตัวอย่างให้เป็นไปตามที่ต้องการได้

ตาราง 3 (ต่อ)

วิธีการสุ่ม	ข้อดี/ข้อเสีย
	<p>3. ต้องใช้สูตรสำหรับการคำนวณหาขนาดของกลุ่มตัวอย่างและสูตรการคำนวณค่าสถิติต่าง ๆ ที่มีความยุ่งยาก</p>
<p>4. วิธีการสุ่มตามกลุ่ม</p> <p>วิธีการสุ่มตัวอย่างวิธีนี้จะแบ่งประชากรออกเป็นกลุ่ม ๆ (Clusters) โดยให้ประชากรในกลุ่มย่อยแต่ละกลุ่มมีลักษณะคล้ายคลึงกัน (Homogeneous) แต่สมาชิกภายในกลุ่มย่อยมีลักษณะแตกต่างกัน (Heterogeneous) แล้วใช้วิธีการสุ่มอย่างง่าย สุ่มมาเพียง 1 กลุ่ม หรือหลายกลุ่มให้ได้ขนาดของกลุ่มตัวอย่างตามที่กำหนดไว้</p>	<p>ข้อดี</p> <ol style="list-style-type: none"> 1. เป็นการประหยัดเวลาและค่าใช้จ่ายในการเก็บรวบรวมข้อมูล 2. ไม่จำเป็นต้องมีกรอบตัวอย่างที่สมบูรณ์ <p>ข้อเสีย</p> <ol style="list-style-type: none"> 1. มีความยุ่งยากในการแบ่งกลุ่มให้แต่ละกลุ่มมีความคล้ายคลึงกันและสมาชิกภายในกลุ่มมีความแตกต่างกันในกรณีที่กลุ่มไม่ได้แบ่งอยู่ตามธรรมชาติ 2. มีความยุ่งยากในการคำนวณหาขนาดของกลุ่มตัวอย่างและการคำนวณค่าสถิติต่าง ๆ
<p>5. การสุ่มแบบหลายขั้นตอน</p> <p>เป็นการสุ่มตัวอย่างที่พัฒนามาจากวิธีการสุ่มตามกลุ่ม จะมีการสุ่มตั้งแต่ 2 ขั้นขึ้นไป จะมีการสุ่มกลุ่มหรือระดับชั้นใหญ่ขึ้นมาก่อนแล้วจึงสุ่มกลุ่มหรือระดับชั้นย่อยลงไปจนถึงหน่วยที่ต้องการศึกษา</p>	<p>ข้อดี</p> <ol style="list-style-type: none"> 1. เป็นการประหยัดเวลาและค่าใช้จ่ายในการเก็บรวบรวมข้อมูล 2. ใช้ได้กับประชากรที่มีจำนวนมากจนไม่สามารถหาบัญชีรายชื่อได้ 3. ไม่จำเป็นต้องมีกรอบตัวอย่างที่สมบูรณ์ของทุก ๆ ชั้น <p>ข้อเสีย</p> <ol style="list-style-type: none"> 1. มีความยุ่งยากในการคำนวณหาขนาดของกลุ่มตัวอย่างและการคำนวณหาค่าสถิติต่าง ๆ เนื่องจากมีการสุ่มหลายครั้ง

ตาราง 3 (ต่อ)

วิธีการสุ่ม	ข้อดี/ข้อเสีย
	<p>ข้อเสีย</p> <p>2. ต้องใช้ขนาดของกลุ่มตัวอย่างมากกว่าวิธีการสุ่มอื่น ๆ จึงจะได้ค่าประมาณใกล้เคียงกัน</p>

การวางแผนการสุ่มตัวอย่างประชากรแต่ละวิธีมีข้อดีและข้อเสียอยู่ในตัวของมันเอง ซึ่งจะต้องพิจารณาให้รอบคอบและเลือกวิธีที่ดีที่สุด เหมาะสมที่สุด เพื่อให้ได้ตัวแทนของประชากรที่สอดคล้องเหมาะสมกับสภาพที่ทำการศึกษ (ดวงใจ ปวีณอภิชาติ, 2535. หน้า 22)

การจำลองสถานการณ์ด้วยวิธีการมอนติ คาร์โล

ความเป็นมาของวิธีมอนติ คาร์โล

เทคนิควิธีมอนติ คาร์โล เป็นสาขาหนึ่งของคณิตศาสตร์ใช้คอมพิวเตอร์ช่วยในการจำลองสถานการณ์ (Simulation) โดยอาศัยตัวเลขสุ่ม (Random Number) มาสร้างตัวแปรให้เหมือนกับสถานการณ์จริงและมีการทดลองซ้ำหลาย ๆ ครั้ง เพื่อให้ได้ค่าที่แน่นอนที่จะใช้เป็นข้อสรุปหรืออธิบายปรากฏการณ์ต่าง ๆ ในสถานการณ์จริง (ต่าย เที่ยงฉี, 2534. หน้า 62-68) หรือช่วยหาคำตอบในเรื่องราวต่าง ๆ ที่ยังไม่แน่ใจในผลที่จะเกิดขึ้น

เทคนิควิธีมอนติ คาร์โล ได้มีการใช้มานานแล้วแต่ในสมัยก่อน ๆ ไม่ได้เรียกว่า มอนติ คาร์โล โดยนำมาใช้พัฒนาทฤษฎีความน่าจะเป็น (Probability Theory) ในราวปี ค.ศ.1753 จอร์ส หลุยส์ เลคเลอร์ และบุฟฟอง (Georges Louis Leclere and Comte de Buffon) ทำการทดลองหาค่า (π) โดยการโยนเข็มที่มีความยาว k หน่วย อย่างสุ่มลงมาบนพื้นราบที่มีเส้นขนานอยู่ โดยให้ระยะห่างระหว่างเส้นขนานแต่ละเส้นห่างกัน d หน่วย และกำหนดให้ $d > k$ จะมีความน่าจะเป็นที่เข็มจะตัดกับเส้นขนาน $P = 2k/\pi d$ ซึ่งถ้าความน่าจะเป็น (P) เป็นค่าสุ่ม ก็จะหาค่า π ได้ ในราวปี ค.ศ. 1908 กอสเซท (W.S.Gosset) ได้ศึกษาการแจกแจงความถี่ของค่าสูงของนักโทษอาชญากรรมจำนวน 3,000 คน โดยเทียบกับการแจกแจงความถี่ของกลุ่มตัวอย่างที่สุ่มมาครั้งละ 4 คน จำนวน 750 กลุ่มตัวอย่าง ผลการศึกษาพบว่า การแจกแจงความถี่ทั้งสองมี

ลักษณะเหมือนกัน กอสมิท ได้ตั้งชื่อการแจกแจงความถี่ที่ค้นพบนี้ว่า การแจกแจงค่าที่ (t-distribution) ซึ่งถือได้ว่าเป็นจุดเริ่มต้นของเทคนิควิธีมอนติ คาร์โล (Monte Carlo Method) เทคนิคมอนติ คาร์โล ได้รับการพัฒนาอย่างจริงจังในปี ค.ศ.1944 ช่วงสงครามโลกครั้งที่ 2 อูลาม และอน นิวแมน (Ulam and Von Neumann) เป็นผู้ตั้งชื่อ มอนติ คาร์โล เป็นรหัสรับของงานที่ทำใน ลอส อลามอส (Los Alamos) ได้นำเทคนิคนี้มาหาผลของการแพร่อย่างสุ่มของนิวตรอน (Neutron diffusion) ในวัสดุเชื้อเพลิงซึ่งเป็นการทดลองทางคณิตศาสตร์เพื่อหาผลของค่าตอบก่อนที่จะทำการทดลองจริง ซึ่งทำให้ไม่เกิดอันตรายและช่วยประหยัดค่าใช้จ่ายก่อนการทดลองจริง หลังจากนั้นเทคนิคมอนติ คาร์โล ได้มีการนำมาใช้อย่างกว้างขวางทั้งทางด้านฟิสิกส์ คณิตศาสตร์ สถิติ และการวิจัย นับได้ว่าเทคนิควิธีมอนติ คาร์โล มีประโยชน์อย่างมากในการขยายความรู้เชิงทฤษฎี เช่น การนำมาศึกษาค่าความคลาดเคลื่อนทางสถิติ เปรียบเทียบประสิทธิภาพการทดสอบชนิดต่าง ๆ เป็นต้น

ขั้นตอนของระเบียบวิธีมอนติ คาร์โล

หลักการที่สำคัญของวิธีการมอนติ คาร์โล ก็คือ การนำเอาตัวเลขสุ่ม (Random Number) มาประยุกต์แก้ปัญหาต่าง ๆ มีขั้นตอนที่สำคัญดังนี้

1. สร้างตัวเลขสุ่ม (Generate random number) ระยะเวลาทำได้โดยอาศัยเครื่องมือทางกายภาพ เช่น ล้อรูเล็ต ลูกเต๋า ไฟ กระจายเขียนเบอร์ เป็นต้น แต่ได้ตัวเลขสุ่มไม่มาก ต่อมาจึงหันมาใช้เครื่องมืออิเล็กทรอนิกส์ เช่น เครื่องสร้างตัวเลขสุ่มที่สร้างขึ้นโดยบริษัท แรนต์ (Rand) โดยสร้างตัวเลขสุ่มจากเครื่องกำเนิดพัลส์ (Pulse) ซึ่งสามารถสร้างตัวเลขสุ่มได้เป็นล้านตัว

การสร้างหรือเลือกใช้ตัวเลขสุ่มดังกล่าวกับเครื่องคอมพิวเตอร์ยังมีปัญหา 2 ประการคือ เป็นการยากที่จะทำให้คอมพิวเตอร์สามารถเรียกใช้ได้เมื่อมีความต้องการ และยากที่จะทำให้เครื่องมือดังกล่าวสร้างตัวเลขสุ่มชุดเดิม เมื่อต้องการใช้เปรียบเทียบวิธีการต่าง ๆ ภายใต้เงื่อนไขของระบบเลขสุ่มชุดเดียวกัน หรือถ้าจะเก็บเลขสุ่มเหล่านี้ไว้ในหน่วยความจำหรือจานแม่เหล็กก็จะทำให้สูญเสียหน่วยความจำหรือเสียเวลาในการค้นหา ฉะนั้นการสร้างตัวเลขสุ่มในคอมพิวเตอร์จึงนิยมสร้างตัวเลขสุ่มเทียม (Pseudo random number) โดยอาศัยสูตรทางคณิตศาสตร์

ในปัจจุบันมีโปรแกรมสำหรับสร้างตัวเลขสุ่มในเครื่องคอมพิวเตอร์ เช่น ในภาษาเบสิก (BASIC) มีคำสั่งเรียกใช้ตัวเลขสุ่มได้คือ RANDOMIZED และ RND ในภาษาฟอร์แทน

(FORTRAN) ก็มีคำสั่งเรียกตัวเลขสุ่มได้ คือ RANDUM ส่วนภาษาฟอกโปร (FOXPRO) มีคำสั่งเรียกตัวเลขสุ่มคือ RAND()

คุณสมบัติของตัวเลขสุ่มที่ดี

1. ตัวเลขสุ่มที่ได้ต้องมีลักษณะการกระจายความน่าจะเป็นแบบสม่ำเสมอ (Uniform Distribution)

2. ตัวเลขสุ่มที่ได้ต้องเป็นอิสระต่อกัน

3. อนุกรมของตัวเลขสุ่มที่ได้ต้องสามารถสร้างซ้ำเดิมได้

4. อนุกรมตัวเลขสุ่มที่ได้ต้องไม่ซ้ำเดิมในช่วงที่ต้องการใช้ตัวเลขสุ่ม

5. ต้องใช้เวลาน้อยในการสร้างตัวเลขสุ่ม

6. ต้องใช้หน่วยความจำในคอมพิวเตอร์น้อย

2. นำตัวเลขสุ่มมาประยุกต์ใช้กับปัญหาต่าง ๆ เป็นการนำตัวเลขสุ่มไปสร้างตัวแปรตามลักษณะการแจกแจงของปัญหาที่จะศึกษาเพื่อเป็นข้อมูลของปัญหานั้น เช่น สร้างตัวเลขสุ่มแล้วนำตัวเลขสุ่มไปสร้างเป็นคะแนนการสอบของผู้เรียน แต่บางครั้งตัวแปรของปัญหาที่จะศึกษาไม่ได้สร้างจากตัวเลขสุ่มโดยตรง แต่ใช้ตัวเลขสุ่มเป็นพื้นฐานก็ได้

3. ทำการทดลองซ้ำหลาย ๆ ครั้ง การศึกษาด้วยวิธีมอนติ คาร์โล ต้องมีการทดลองซ้ำหลาย ๆ ครั้งเพื่อลดความคลาดเคลื่อนของคำตอบที่จะได้ และสามารถสรุปเป็นความน่าจะเป็นของการเกิดเหตุการณ์ในปัญหานั้น ๆ

จุดเด่นของการใช้เทคนิคมอนติ คาร์โล

เทคนิคมอนติ คาร์โล ใช้ตัวเลขสุ่มเป็นพื้นฐานในการสร้างตัวแปรของปัญหาโดยอาศัยทฤษฎี สูตร หรือกฎเกณฑ์ต่าง ๆ ที่มีอยู่ และมีการทดลองซ้ำหลาย ๆ ครั้ง เพื่อลดความคลาดเคลื่อนต่าง ๆ ซึ่งนับว่ามีประโยชน์ที่สำคัญดังนี้

1. สามารถควบคุมตัวแปรแทรกซ้อนและสามารถสังเกตได้อย่างสมบูรณ์ และทำการทดลองซ้ำภายใต้สภาพแวดล้อม (Context) เดิมหลาย ๆ ครั้งได้ ส่วนในการทดลองจริงนั้นทำไม่ได้เพราะไม่สามารถรักษาสภาพแวดล้อมให้เหมือนเดิมทุกอย่างได้เมื่อเวลาเปลี่ยนไป

2. ถ้ามีสูตรหรือกฎเกณฑ์ต่าง ๆ ที่ถูกต้องรองรับในการสร้างตัวแปรของปัญหาในการทดลองแล้วจะให้ผลที่ถูกต้องแม่นยำกว่าเมื่อใช้ทดลองในสถานการณ์จริง เพราะสามารถลดตัวแปรแทรกซ้อนเชิงจิตวิทยาได้

3. สิ้นเปลืองเวลา แรงงาน และค่าใช้จ่ายน้อยกว่าเมื่อเทียบกับการทดลองในสถานการณ์จริง

ความคลาดเคลื่อนประเภทที่ 1 ความคลาดเคลื่อนประเภทที่ 2 และอำนาจการทดสอบ

การทดสอบสมมติฐานทางสถิติเป็นการศึกษาข้อมูลจากกลุ่มตัวอย่างเพื่ออ้างอิงหรือสรุปไปยังประชากร เทคนิคที่เป็นมาตรฐานในการทดสอบสมมติฐาน ก็คือ การเลือกระดับนัยสำคัญ (α) เพื่อกำหนดค่าวิกฤต (critical values) สุ่มตัวอย่างขึ้นมาศึกษา คำนวณค่าสถิติ และยอมรับหรือปฏิเสธสมมติฐาน การตัดสินใจในการทดสอบสมมติฐานจะถูกหรือผิดก็ได้ การตัดสินใจจะถูกต้องในการยอมรับสมมติฐานสุญถ้าสมมติฐานสุญนั้นเป็นความจริง แต่ถ้าปฏิเสธสมมติฐานสุญที่เป็นจริง เป็นการตัดสินใจที่ไม่ถูกต้องจะเรียกว่า ความคลาดเคลื่อนประเภทที่ 1 การตัดสินใจจะถูกต้องอีกเช่นเดียวกันในการปฏิเสธสมมติฐานสุญที่ผิดและตัดสินใจไม่ถูกต้องในการยอมรับสมมติฐานสุญที่ผิด ซึ่งจะเกิดความคลาดเคลื่อนประเภทที่ 2 (Type II error) สรุปได้ตามภาพที่ 1

		สมมติฐานสุญ	
		ถูกต้อง	ไม่ถูกต้อง
การตัดสินใจ	ยอมรับ	ตัดสินใจถูก ($1-\alpha$)	ตัดสินใจผิด ความคลาดเคลื่อน ประเภทที่ 2, β
	ปฏิเสธ	ความคลาดเคลื่อน ประเภทที่ 1, α	ตัดสินใจถูก ($1-\beta$)

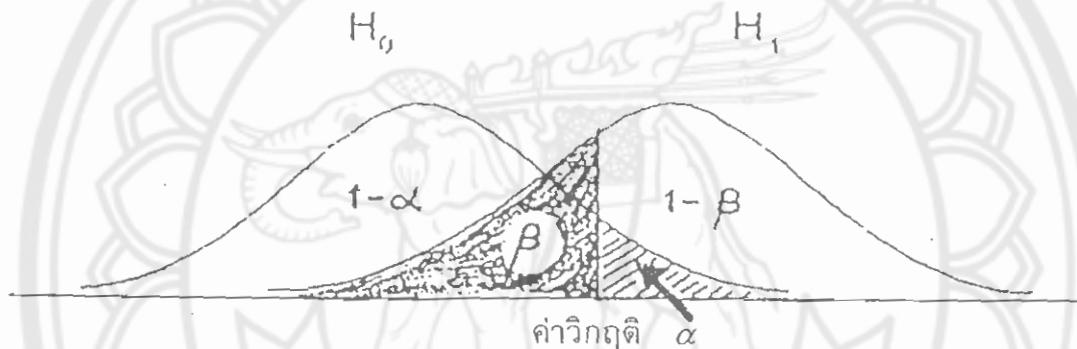
ภาพ 1 รูปแบบของความคลาดเคลื่อนประเภทที่ 1 และประเภทที่ 2

การทดสอบสมมติฐานสุญ (H_0) ในงานวิจัยเชิงปริมาณเพื่อปฏิเสธสมมติฐานสุญและยอมรับสมมติฐานอื่น (H_1) ซึ่งการตัดสินใจที่น่าเชื่อถือและถูกต้องมากที่สุดคือ การตัดสินใจปฏิเสธสมมติฐานสุญและสมมติฐานสุญไม่ถูกต้อง จากภาพที่ 1 ก็คือค่าของ $1-\beta$ เรียก $1-\beta$

นี้ว่าเป็นอำนาจของการทดสอบ (power of test) ดังนั้นอำนาจของการทดสอบ ก็คือ โอกาสหรือความน่าจะเป็นที่จะปฏิเสธสมมติฐานศูนย์เมื่อสมมติฐานศูนย์ที่ตั้งไว้ไม่ถูกต้อง หรืออาจกล่าวได้อีกนัยหนึ่งว่า อำนาจของการทดสอบ คือ โอกาสหรือความน่าจะเป็นที่จะปฏิเสธสมมติฐานศูนย์ เมื่อสมมติฐานอื่นถูกต้อง หรือยอมรับสมมติฐานอื่นเมื่อสมมติฐานอื่นที่ตั้งไว้ถูกต้อง (สำราญ กำจัดภัย, 2542. หน้า 70-74)

ความสัมพันธ์ระหว่าง α กับอำนาจของการทดสอบ

ความสัมพันธ์ระหว่าง α กับอำนาจของการทดสอบแสดงให้เห็นตามภาพที่ 2

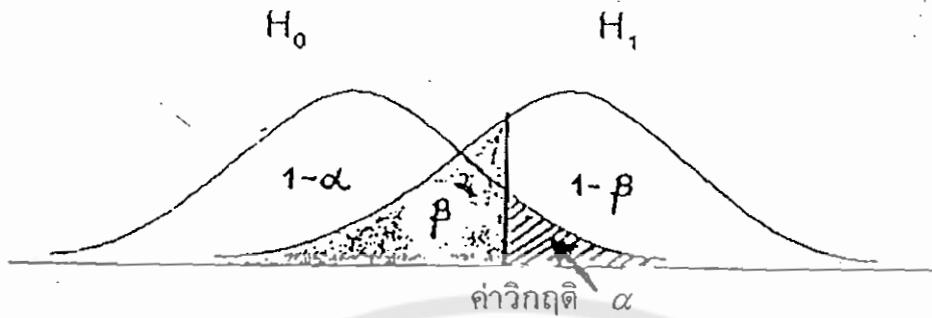


ภาพ 2 แสดงความสัมพันธ์ระหว่าง α กับอำนาจของการทดสอบ

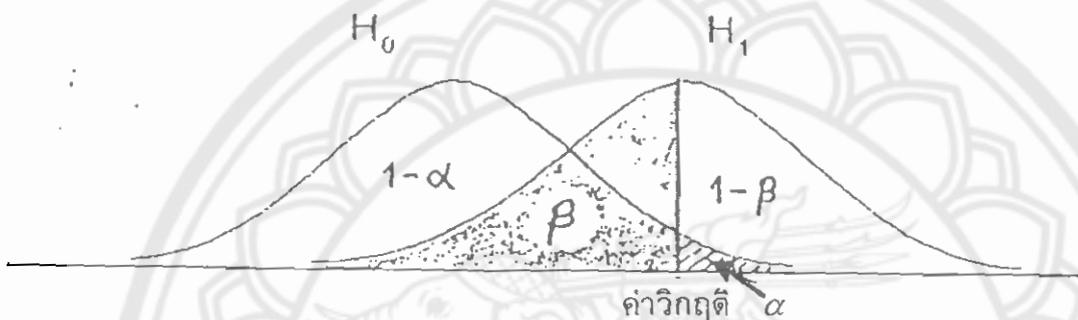
จากภาพ α คือ พื้นที่ใต้โค้งแทนความน่าจะเป็นของการตัดสินใจผิดพลาดชนิด

Type I error β คือ พื้นที่ใต้โค้งแทนความน่าจะเป็นของการตัดสินใจผิดพลาดชนิด Type II error และ $1 - \beta$ คือ พื้นที่ใต้โค้งแทนความน่าจะเป็นในการปฏิเสธสมมติฐานศูนย์เมื่อสมมติฐานศูนย์นั้นไม่ถูกต้อง เรียกว่า อำนาจของการทดสอบ

ถ้าลด α ให้มีพื้นที่น้อยลงจากภาพที่ 3 (ก) เป็น 3 (ข) ทำให้ค่าวิกฤตเลื่อนไปทางขวา (ทางบวก) ส่งผลให้พื้นที่ของ β เพิ่มขึ้น และในขณะเดียวกันพื้นที่ของ $1 - \beta$ จะลดน้อยลง ดังนั้นสรุปได้ว่า ถ้าลดโอกาสเสี่ยงต่อความผิดพลาดชนิด α จะส่งผลให้ความผิดพลาดชนิด β เพิ่มขึ้น และในขณะเดียวกันอำนาจของการทดสอบจะน้อยลง ซึ่งแสดงให้เห็นได้ดังภาพ



ภาพ 3 (ก) แสดงความสัมพันธ์ระหว่าง α กับอำนาจของการทดสอบ

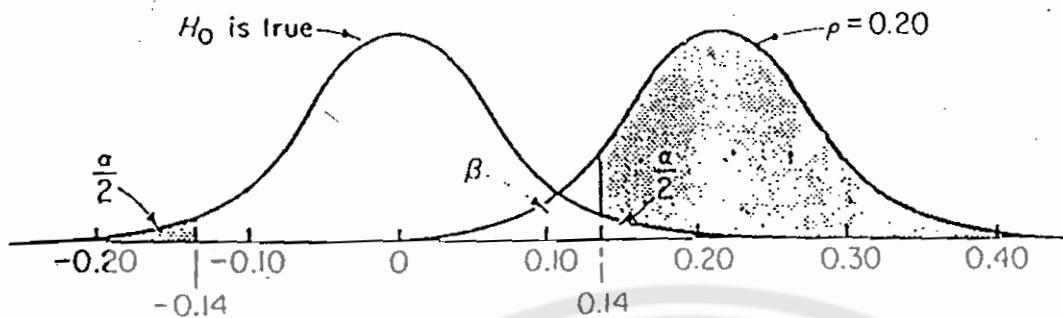


ภาพ 3 (ข) แสดงความสัมพันธ์ระหว่าง α กับอำนาจของการทดสอบ

ความสัมพันธ์ระหว่าง β กับอำนาจของการทดสอบ

ดังได้กล่าวแล้วว่าในการทดสอบสมมติฐาน สมมติฐานใดสมมติฐานหนึ่งจะต้องเป็นจริง นั่นก็คือ หลังจากเก็บรวบรวมข้อมูลจากกลุ่มตัวอย่างแล้วการตัดสินใจอย่างใดอย่างหนึ่งต้องเกิดขึ้น ยอมรับสมมติฐานศูนย์ ปฏิเสธสมมติฐานอื่น หรือยอมรับสมมติฐานอื่น ปฏิเสธสมมติฐานศูนย์ การตัดสินใจที่ไม่ถูกต้อง คือ การยอมรับสมมติฐานศูนย์เมื่อสมมติฐานศูนย์นั้นไม่เป็นความจริง เป็นความคลาดเคลื่อนประเภทที่ 2 ให้สัญลักษณ์ β เช่น ผู้วิจัยต้องการทดสอบสมมติฐาน $H_0: \rho = 0$ สุ่มตัวอย่างขึ้นมาศึกษาจำนวน 200 คน ต้องการปฏิเสธ $H_0: \rho = 0$ เมื่อสมมติฐานศูนย์นี้เกิดขึ้นจริงเพียง 5 ใน 100 ครั้ง จึงกำหนดระดับนัยสำคัญที่ต้องการทดสอบเท่ากับ .05 ($\alpha = .05$) ผู้วิจัยกำหนดค่าวิกฤตซึ่งจะนำไปสู่การปฏิเสธสมมติฐานศูนย์ ดังภาพที่ 4 ค่าของ r มากกว่าหรือ เท่ากับ .14 หรือน้อยกว่าหรือเท่ากับ -.14 จะปฏิเสธสมมติฐานศูนย์ แสดงว่า สมมติฐานศูนย์ไม่ถูกต้อง ซึ่งเราทราบว่ถ้าสมมติฐานศูนย์เป็นจริง $\rho = 0$ จะมีโอกาสปฏิเสธ 5 ใน 100 ครั้ง ($\alpha = .05$)

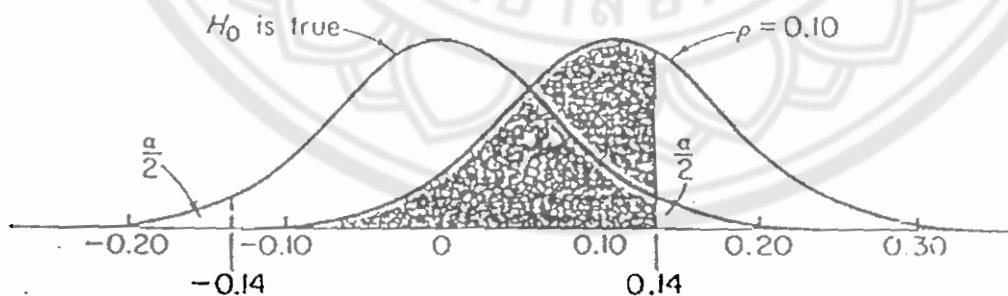
อย่างไรก็ตามถ้าในความเป็นจริงแล้ว $\rho = .20$ กรณีนี้ $H_0: \rho = 0$ ควรจะถูกปฏิเสธ และสรุปได้ว่าความสัมพันธ์แตกต่างจากศูนย์ ความน่าจะเป็นในการปฏิเสธสมมติฐานศูนย์เมื่อไม่เป็นความจริง ก็คือ อำนาจของการทดสอบดังส่วนที่แรเงาในภาพที่ 4



ภาพ 4 แสดงอำนาจการทดสอบ $H_0: \rho = 0$ และ $H_0: \rho \neq 0$ ($\rho = .20$) จำนวนกลุ่มตัวอย่าง 200 และ $\alpha = .05$

พื้นที่ที่สูงกว่าค่าวิกฤติ คือ ค่าสหสัมพันธ์ทุก ๆ ค่าจาก .14 ไปถึง 1.00 ดังนั้น อำนาจของการทดสอบสมมติฐานเพื่อปฏิเสธ H_0 เมื่อ $\rho = .20$ ก็คือพื้นที่เหนือกว่า .14 ภายใต้โค้งด้านขวามือ ซึ่งมีพื้นที่ประมาณ 82% ส่วนค่า β คือ ความน่าจะเป็นในการยอมรับสมมติฐานสูญที่ไม่ถูกต้องมีค่าเท่ากับ .18 สรุปได้ว่า ในการทดสอบสมมติฐานครั้งนี้ถ้า $\rho = .20$ แล้วการทดสอบนี้ จะมีความคลาดเคลื่อนประเภทที่ 2 เกิดขึ้น โดยมีอำนาจของการทดสอบเท่ากับ .82

ถ้าความสัมพันธ์เปลี่ยนจาก .20 เป็น .10 อำนาจการทดสอบจะเป็นเท่าไร เมื่อค่าระดับนัยสำคัญ และจำนวนกลุ่มตัวอย่างเหมือนเดิม ($\alpha = .05, n = 200$) ค่าวิกฤติของการทดสอบยังเหมือนเดิม คือ -1.00 ถึง $-.14$ และ $.14$ ถึง 1.00 การแจกแจงการสุ่มของค่า r กรณีกลุ่มตัวอย่างเท่ากับ 200 ยังไม่เปลี่ยนแปลงจากที่กล่าวมา ดังภาพที่ 5



ภาพ 5 แสดงอำนาจการทดสอบ $H_0: \rho = 0$ และ $H_0: \rho \neq 0$ ($\rho = .10$) จำนวนกลุ่มตัวอย่าง 200 และ $\alpha = .05$

จากการวัดพื้นที่ได้โค้งด้านขวามือเหนือค่าวิกฤต .14 แสดงให้เห็นว่าอำนาจการทดสอบ $H_0: \rho = 0$ เมื่อจริง ๆ แล้ว $\rho = .10$ มีค่าเท่ากับ .29 ($\beta = .71$) แสดงให้เห็นว่าถ้าค่าความสัมพันธ์น้อยลงอำนาจของการทดสอบจะลดลง

การทดสอบสมมติฐาน $H_0: \rho = 0$ และ $H_1: \rho \neq 0$ ด้วย $\alpha = .10$ และ $n=200$ อำนาจการทดสอบจะมากกว่าหรือน้อยกว่ากรณีการทดสอบที่ $\alpha = .05$ เมื่อในความเป็นจริงแล้ว $\rho = .20$ จากการวัดพื้นที่ได้โค้งปกติด้านขวาของค่าวิกฤตตั้งแต่ค่า $t = .12$ (เมื่อระดับนัยสำคัญเปลี่ยนไปค่าวิกฤตจะเปลี่ยนไปด้วย) แสดงให้เห็นว่าอำนาจของการทดสอบ $H_0: \rho = 0$ เมื่อจริง ๆ แล้ว $\rho = .20$ ที่ $\alpha = .10$ มีค่าเท่ากับ .87 ($\beta = .13$) แสดงให้เห็นว่าถ้าระดับนัยสำคัญสูงขึ้นค่าอำนาจของการทดสอบก็จะสูงขึ้นด้วย

อำนาจของการทดสอบจะเพิ่มขึ้นเมื่อค่าความสัมพันธ์แตกต่างจากศูนย์มากขึ้น ค่าความสัมพันธ์นั้นนักวิจัยไม่สามารถควบคุมได้ แต่อย่างไรก็ตามนักวิจัยสามารถควบคุมขนาดของกลุ่มตัวอย่างและระดับนัยสำคัญได้ กล่าวได้ว่า เมื่อ $\rho \neq 0$ อำนาจของการทดสอบจะเพิ่มขึ้นเมื่อขนาดของกลุ่มตัวอย่างและค่านัยสำคัญเพิ่มมากขึ้น

องค์ประกอบที่มีผลกระทบต่ออำนาจของการทดสอบ

บุญเรียง ขจรศิลป์ (2533, หน้า 81-82) ได้กล่าวว่า องค์ประกอบที่มีผลกระทบต่ออำนาจของการทดสอบมีดังนี้

1. โอกาสที่จะเกิดความคลาดเคลื่อนชนิดที่ 2 (β) อำนาจของการทดสอบแปรผกผันกับโอกาสที่จะเกิดความคลาดเคลื่อนชนิดที่ 2 ถ้าโอกาสที่จะเกิดความคลาดเคลื่อนชนิดที่ 2 มากขึ้น อำนาจของการทดสอบก็จะน้อยลง และถ้าโอกาสที่จะเกิดความคลาดเคลื่อนชนิดที่ 2 น้อยลง อำนาจของการทดสอบจะมากขึ้น
2. โอกาสที่จะเกิดความคลาดเคลื่อนชนิดที่ 1 (α) อำนาจของการทดสอบแปรตามกับโอกาสที่จะเกิดความคลาดเคลื่อนชนิดที่ 1 ถ้าโอกาสที่จะเกิดความคลาดเคลื่อนชนิดที่ 1 มากขึ้น อำนาจของการทดสอบก็จะมากขึ้น และถ้าโอกาสที่จะเกิดความคลาดเคลื่อนชนิดที่ 1 น้อยลง อำนาจของการทดสอบจะน้อยลง การทดสอบสมมติฐานแบบหางเดียวจะมีอำนาจของการทดสอบมากกว่าการทดสอบแบบสองหาง
3. ความแปรปรวนของกลุ่มตัวอย่าง อำนาจของการทดสอบจะแปรผกผันกับความแปรปรวนของกลุ่มตัวอย่าง ถ้ากลุ่มตัวอย่างมีความแปรปรวนน้อยอำนาจของการทดสอบจะมาก
4. ขนาดของความแตกต่างระหว่างค่าพารามิเตอร์ของกลุ่มประชากรภายใต้สมมติฐานศูนย์กับสมมติฐานอื่น อำนาจของการทดสอบจะแปรตามขนาดของความแตกต่างระหว่างค่า

พารามิเตอร์ของกลุ่มประชากรภายใต้สมมติฐานสูญญกับสมมติฐานอื่น ถ้ายิ่งแตกต่างกันมาก อำนาจของการทดสอบจะยิ่งมากขึ้น

การกำหนดระดับนัยสำคัญ (α) ในการทำวิจัย

ในการทดสอบสมมติฐานทางสถิติความคิดของนักวิจัยส่วนใหญ่จะมีอยู่ 2 ค่า ค่าแรก กำหนด α (significance level) ไว้ล่วงหน้า ค่าที่สองไม่ได้กำหนด α ไว้ล่วงหน้าแต่จะรายงานค่า p (p-value) หรือค่าสถิติ (statistical value) ที่คำนวณได้แทน (สุชาติดา บวรภักดีวงศ์, 2541. หน้า 16-19)

นักวิจัยค่าที่หนึ่งจะกำหนด α ไว้ก่อน ค่า α ที่นิยมใช้กันมากคือ .05 .01 และ .10 ตามลำดับ แล้วนำค่า p ที่คำนวณได้ไปเปรียบเทียบกับค่า α ที่กำหนด หรือนำค่าสถิติที่คำนวณได้ไปเปรียบเทียบกับค่าวิกฤตที่เปิดจากตารางที่ระดับ α ที่กำหนดไว้ ถ้าค่า p ที่คำนวณได้มีค่าน้อยกว่า α หรือค่าสถิติที่คำนวณได้มีค่าตกอยู่ในบริเวณวิกฤตผลการทดสอบจะปฏิเสธสมมติฐานสูญญที่ตั้งไว้ มิฉะนั้นแล้วจะยอมรับสมมติฐานสูญญ นักวิจัยในกลุ่มนี้จะให้ความสำคัญกับ α ในตำแหน่งที่ไม่ต่อเนื่องคือ .05 .01 หรือ .10 มากเกินไป จึงมีนักสถิติหลายคนตั้งคำถามว่าทำไม α ที่ระดับ .05 จึงมีความสำคัญมากกว่า α ที่ระดับ .06 หรือ .07 มากมายนัก?

นักวิจัยในค่าที่สองให้ความสำคัญกับ α ทุกค่าเท่ากันหมด คือมอง α เป็นค่าต่อเนื่อง นักวิจัยกลุ่มนี้จะรายงานค่า p ที่คำนวณได้เพื่อให้ผู้อ่านตัดสินใจเองว่า ที่ α ระดับใดจะปฏิเสธ H_0 แต่โดยส่วนใหญ่ผู้วิจัยจะสรุปผลการทดสอบว่าจะปฏิเสธ H_0 ในช่วงใด ตัวอย่างเช่น ถ้า p ที่คำนวณได้มีค่าเป็น .006 จะสรุปว่าปฏิเสธ H_0 ที่ $p < .01$ ถ้า p ที่คำนวณได้มีค่าเป็น .015 ก็ จะปฏิเสธ H_0 ที่ $p < .02$ ถ้า p ที่คำนวณได้มีค่าเป็น .045 ก็ จะปฏิเสธ H_0 ที่ $p < .05$ เป็นต้น

แนวทางในการกำหนดขนาดของ α มีดังนี้

1. การนำผลการทดสอบไปใช้หลังการทดลองหรือผลสืบเนื่องในทางปฏิบัติที่จะตามมา เช่น ถ้าผู้วิจัยต้องการทดสอบประสิทธิภาพของยาชนิดหนึ่งว่าสามารถรักษาโรคได้จริงหรือไม่ ดังนั้น

Type I error rate = α = โอกาสที่จะปฏิเสธ H_0 ที่เป็นจริง

Type II error rate = β = โอกาสที่จะยอมรับ H_0 ที่เป็นเท็จ

Power of the test = $1 - \beta$ = โอกาสที่จะปฏิเสธ H_0 ที่เป็นเท็จ

ผู้วิจัยอาจเริ่มต้นด้วยการถามตัวเองว่า ความคลาดเคลื่อนประเภทไหนที่ควรให้ความสำคัญ หรือควรควบคุมมากกว่า α ในที่นี้จะหมายถึง โอกาสที่จะปฏิเสธยาที่รักษาโรคได้ β หมายถึง โอกาสที่จะยอมรับยาที่ไม่สามารถรักษาโรคได้ ถ้าผู้วิจัยเห็นว่างานวิจัยนี้ควรเน้น

β เพราะการยอมรับยาที่ไม่มีประสิทธิภาพในการรักษาโรค น่าจะมีความเสี่ยงมากกว่าการปฏิเสธยาที่มีประสิทธิภาพ ถ้าเป็นเช่นนี้ก็ควรให้ β มีค่าเล็ก (เช่น $\beta = .001$) และ α มีค่าใหญ่ (เช่น $\alpha = .10$) เป็นต้น ดังนั้นนักวิจัยจะต้องตัดสินใจว่าจะควบคุมความคลาดเคลื่อนประเภทใด ถ้าต้องการให้ β เล็ก ก็ต้องกำหนด α ขนาดใหญ่ แต่ถ้าต้องการให้ α เล็ก ก็ต้องกำหนดให้ β ใหญ่

2. วัตถุประสงค์ของการวิจัย ในบางครั้งผู้วิจัยต้องการตรวจสอบว่า ผลงานวิจัยนั้น ยืนยัน (confirm) ทฤษฎี (theory) หรือตัวแบบ (model) ที่มีคนคิดไว้แล้วหรือไม่ การทำงานวิจัยเพื่อสนับสนุนทฤษฎี หรือ model ที่มีอยู่แล้ว อาจใช้ α ขนาดใหญ่ได้ แต่ถ้าต้องการสร้างทฤษฎีใหม่ หรือ model ใหม่ ควรกำหนด α ขนาดเล็ก

3. อำนาจของการทดสอบ (power of the test) อำนาจการทดสอบ คือ โอกาสที่จะปฏิเสธสมมติฐานที่ผิด ที่ผู้วิจัยต้องการให้มีค่าสูง ๆ ในทางทฤษฎีอำนาจการทดสอบจะแปรผันตรงกับขนาดตัวอย่าง ดังนั้น ถ้าขนาดตัวอย่างใหญ่ขึ้น อำนาจการทดสอบจะสูงขึ้นด้วย ในขณะที่เดียวกันจะมีผลทำให้ α สูงขึ้นด้วย อันเป็นผลเนื่องมาจากค่า standard error ซึ่งจะแปรผกผันกับขนาดตัวอย่างอันจะทำให้ α มีค่ามากเมื่อขนาดตัวอย่างใหญ่ ถ้าขนาดตัวอย่างเล็กโอกาสที่จะยอมรับสมมติฐานสูญจะมาก ในทางปฏิบัติ α ขนาดเล็ก .01, .001 ควรใช้กับขนาดตัวอย่างใหญ่ และ α ขนาดใหญ่ .10, .15 ควรใช้กับขนาดตัวอย่างเล็ก

4. ระดับการควบคุมตัวแปรต่าง ๆ ใน design การทำวิจัยที่สามารถควบคุมตัวแปรแทรกซ้อนและตัวแปรที่กำลังศึกษาได้อย่างค่อนข้างเต็มที่ ดังนั้น α ที่มีขนาดใหญ่ก็น่าจะยอมรับได้ สำหรับ design ที่ไม่สามารถควบคุมตัวแปรและปัจจัยได้เต็มที่เหมือนการทดลองในห้องแล็บ (Lab) เพื่อให้ผู้วิจัยมั่นใจในผลการวิจัยควรใช้ α ขนาดเล็ก

5. ความถูกต้องของผลการทดสอบ (robustness of the test) ในการทดสอบทางสถิติเมื่อข้อมูลมีลักษณะบางอย่างไม่เป็นไปตามข้อตกลงเบื้องต้น (assumptions) เช่น การแจกแจงของประชากรไม่เป็นโค้งปกติ หรือ ค่าความแปรปรวนของกลุ่มประชากรไม่เท่ากันทุกกลุ่ม ไม่ได้ทำให้โอกาสที่จะเกิดความคลาดเคลื่อนประเภทที่ 1 (Type I error) หรือความคลาดเคลื่อนประเภทที่ 2 (Type II error) มากขึ้น จะเรียกว่าการทดสอบนั้นมีความคงทนหรือมีความถูกต้อง (robustness of the test) (บุญเรียง ขจรศิลป์, 2525. หน้า 29 ; ไปรมา พจนทิมล, 2526. หน้า 1 อ้างอิงจาก Sawat Pratoomraj, 1970. p. 1) ดังนั้นถ้าต้องการให้ผลการทดสอบมีความถูกต้อง (robust) ควรใช้ α ขนาดเล็ก หมายความว่า ถ้าลักษณะของข้อมูลมีความสอดคล้องกับข้อตกลงเบื้องต้น (assumption) ของตัวสถิติทดสอบ α ขนาดใหญ่ก็น่าจะยอมรับได้

ในทางกลับกัน ถ้าลักษณะของข้อมูลมีความสอดคล้องน้อยหรือไม่สอดคล้องกับ assumption ก็ควรใช้ α ขนาดเล็ก

6. การทดสอบ 1 หาง หรือ 2 หาง (one-tailed or two-tailed test) การทดสอบจะเป็น 1 หาง หรือ 2 หาง ขึ้นอยู่กับคำถามของการวิจัย (research question) ของงานวิจัยนั้น ๆ ถ้าสนใจ ทดสอบในทิศทางน้อยกว่าอย่างเดียว (one-tailed) หรือมากกว่าอย่างเดียว (one-tailed) หรือไม่สนใจทิศทาง (two-tailed) การทดสอบแบบหางเดียวจะมีโอกาสปฏิเสธสมมติฐานศูนย์ได้ง่ายกว่า การทดสอบแบบสองหาง

ขั้นตอนต่าง ๆ ที่ควรพิจารณาก่อนทำวิจัย

1. ผู้วิจัยต้องมีความเข้าใจงานวิจัยนั้น ๆ อย่างแท้จริง เข้าใจถึงปัจจัยสาเหตุ ปัจจัยแทรกซ้อน รวมถึงตัวแปรต่าง ๆ ที่เกี่ยวข้อง
2. ควรมีความเข้าใจในหลักการของตัวสถิติทดสอบได้เป็นอย่างดีโดยเฉพาะข้อตกลงเบื้องต้นของตัวสถิตินั้น ๆ ว่ามีอะไรบ้าง ข้อมูลสอดคล้องหรือไม่สอดคล้องกับข้อตกลงเบื้องต้น และจะทำการทดสอบได้อย่างไร
3. ผู้วิจัยสามารถเสนอแนะได้ว่างานวิจัยนั้นควรใช้ α ที่ระดับใด เพราะเหตุใด
4. การออกแบบการวิจัยควรคำนึงถึงปัจจัยต่าง ๆ ที่มีผลต่ออำนาจของการทดสอบ เช่น ขนาดของตัวอย่าง เทคนิคการสุ่มตัวอย่าง การควบคุมตัวแปรแทรกซ้อน คุณภาพของเครื่องมือที่ใช้ และสถิติทดสอบ เป็นต้น

เกณฑ์การเปรียบเทียบค่าประมาณพารามิเตอร์ระหว่างวิธีการจัดการข้อมูลสูญหายแบบต่าง ๆ

ในการศึกษาวิจัยครั้งนี้ผู้วิจัยต้องการตรวจสอบความแม่นยำ และอำนาจการทดสอบที่ได้จากวิธีการจัดการข้อมูลสูญหายแบบอิมพิวเทชันแบบเอ็ม และแบบลิสต์ไวส์ โดยนำตัวแปรที่เกี่ยวข้องเข้ามาศึกษาด้วยคือ วิธีการสุ่มตัวอย่าง ความสัมพันธ์ระหว่างตัวแปร และจำนวนข้อมูลสูญหาย ความแม่นยำที่ใช้ในการศึกษา คือ ความใกล้เคียงกันระหว่างค่าสถิติกับค่าพารามิเตอร์ ซึ่งก็คือความคลาดเคลื่อน เป็นลักษณะการศึกษาค่าประมาณพารามิเตอร์แบบหนึ่ง ถ้ามีความแม่นยำสูง ค่าสถิติกับค่าพารามิเตอร์ใกล้เคียงกันมากการประมาณค่าพารามิเตอร์ย่อมมีความถูกต้อง แต่ถ้าความแม่นยำต่ำค่าสถิติกับค่าพารามิเตอร์ก็จะแตกต่างกันมาก การประมาณค่าพารามิเตอร์จะมีความผิดพลาด

การศึกษาการประมาณค่าพารามิเตอร์ของประชากรมี 2 แบบ คือ การประมาณค่าแบบจุด (Point estimation) และการประมาณค่าแบบช่วง (Interval Estimation) การประมาณค่าแบบจุดเป็นการประมาณค่าประชากรด้วยค่าเพียงค่าเดียว เช่น ใช้ \bar{X} ประมาณค่า μ , s^2 ประมาณค่า σ^2 หรือ r_{xy} ประมาณค่า ρ เป็นต้น วิธีการประมาณค่าอีกแบบหนึ่ง คือ การประมาณค่าแบบช่วง เป็นการประมาณค่าพารามิเตอร์ของประชากรใดประชากรหนึ่งด้วยช่วงค่าตัวเลขช่วงหนึ่ง โดยมีคุณสมบัติว่า ค่าของประชากรที่แท้จริงจะตกอยู่ในช่วงค่าที่ประมาณนี้ด้วยความเชื่อมั่นระดับหนึ่ง โดยใช้การประมาณค่าแบบจุดและการแจกแจงความน่าจะเป็นของตัวประมาณนั้นในการคำนวณ การประมาณค่าทั้งสองแบบจะเหมาะสมกับกรณีการใช้งานที่ต่างกัน ในกรณีที่ตัวประมาณค่า (Estimator) ซึ่งสามารถนำมาประมาณค่าพารามิเตอร์ได้หลายวิธีได้มีการกำหนด คุณสมบัติของวิธีการประมาณค่าไว้ 4 ประการดังนี้ (ดวงใจ ปวีณอภิชาติ, 2535. หน้า 35-37 อ้างอิงจาก Hay, 1963; Yamane, 1967)

1. ความไม่เอนเอียง (Unbiasness)

ถ้า $\hat{\theta}$ เป็นตัวประมาณค่าที่ไม่เอนเอียงของ θ แล้วจะได้ว่า

$$E(\hat{\theta}) = \theta$$

หมายความว่า ค่าคาดหวัง (Expected value) ของตัวประมาณค่า $\hat{\theta}$ มีค่าเท่ากับค่าของพารามิเตอร์ตัวที่ต้องการประมาณค่า หรือ ค่าเฉลี่ยของค่าสถิติของกลุ่มตัวอย่างที่เป็นไปได้ทั้งหมด (All Possible Sample) มีค่าเท่ากับค่าพารามิเตอร์ที่ต้องการประมาณค่า

ส่วนค่าเอนเอียง (Bias) ของตัวประมาณค่า $\hat{\theta}$ คือ ค่าความแตกต่างระหว่างค่าที่คาดหวังของ $\hat{\theta}$ กับค่าพารามิเตอร์ θ

$$\text{Bias}(\hat{\theta}) = E(\hat{\theta}) - \theta$$

2. ความคงเส้นคงวา

การประมาณค่าพารามิเตอร์ θ ด้วย $\hat{\theta}$ เมื่อขนาดของกลุ่มตัวอย่างใหญ่ขึ้นตัวประมาณค่า $\hat{\theta}$ มีค่าเข้าใกล้ θ มากขึ้นด้วย เขียนเป็นประโยคทางคณิตศาสตร์ได้ดังนี้

$$\lim_{n \rightarrow \infty} P(|\hat{\theta} - \theta| < \varepsilon) = 1$$

เมื่อ \mathcal{E} คือ ค่าความคลาดเคลื่อนที่ยอมให้เกิดขึ้นได้ในการหาค่าพารามิเตอร์ด้วยตัวประมาณค่า $\hat{\theta}$ และ P คือ ความน่าจะเป็นหรือโอกาสที่จะเกิดเหตุการณ์นั้นขึ้น

3. ความมีประสิทธิภาพ (Efficiency)

หมายถึง การพิจารณาถึงความถูกต้องแม่นยำในการประมาณค่าของตัวประมาณค่าหนึ่ง ๆ เกณฑ์ที่ใช้พิจารณาความมีประสิทธิภาพของตัวประมาณค่า คือ ความแปรปรวนที่เปรียบเทียบกับกลุ่มของตัวประมาณค่าที่ไม่ลำเอียงด้วยกัน

โดยทั่วไปจะนิยามประสิทธิภาพของตัวประมาณค่าใด ๆ ว่าเป็นอัตราส่วนระหว่างค่าความแปรปรวนของตัวประมาณค่าที่มีประสิทธิภาพนั้นกับค่าความแปรปรวนของตัวประมาณค่าตัวอื่น ๆ กล่าวคือ ถ้าค่าความแปรปรวนของ $\hat{\theta}_i$ หรือ $\text{Var}(\hat{\theta}_i)$ น้อยกว่าความแปรปรวนของ $\hat{\theta}_j$ หรือ $\text{Var}(\hat{\theta}_j)$ เมื่อทั้ง $\hat{\theta}_i$ และ $\hat{\theta}_j$ เป็นตัวประมาณค่าที่ไม่ลำเอียงแล้วจะได้ว่า

$$E_r = \frac{\text{Var}(\hat{\theta}_i)}{\text{Var}(\hat{\theta}_j)} ; 0 < E_r < 1$$

เมื่อค่า E_r ที่ได้จากการคำนวณน้อยกว่า 1 แสดงว่า $\hat{\theta}_i$ เมื่อเปรียบเทียบกับ $\hat{\theta}_j$ แล้ว $\hat{\theta}_i$ มีประสิทธิภาพมากกว่า (Absolutely efficient) เพราะค่าความแปรปรวนของ $\hat{\theta}_i$ เล็กกว่าความแปรปรวนของ $\hat{\theta}_j$

การพิจารณาคุณภาพของตัวประมาณค่าพารามิเตอร์ $\hat{\theta}$ ลักษณะที่ควรพิจารณาอีกลักษณะหนึ่งคือค่าต่าง ๆ ของ $\hat{\theta}$ มีความแตกต่างจากค่าที่ต้องการประมาณค่าเพียงใด ถ้าค่า $\hat{\theta}$ แตกต่างจาก θ น้อยก็แสดงว่าโอกาสที่ค่าประมาณจะตกอยู่ใกล้ค่าของ θ มาก ตัวประมาณค่าที่มีลักษณะเช่นนี้จะดีกว่าตัวประมาณค่าที่มีค่าเป็นไปได้กระจายออกจากค่าที่ต้องการประมาณมาก เกณฑ์ที่ใช้เปรียบเทียบคือค่าเฉลี่ยความคลาดเคลื่อนกำลังสอง (Mean Square Error) ในทางปฏิบัติเราไม่สามารถที่จะหาตัวประมาณค่าที่มีค่าเฉลี่ยความคลาดเคลื่อนกำลังสองสำหรับทุก ๆ ค่าของ θ ได้ จากการพิจารณาตัวประมาณค่าที่ไม่เอนเอียง พบว่า ค่าเฉลี่ยกำลัง

สองจะเท่ากับความแปรปรวนจึงสามารถหาตัวประมาณค่าที่ดีที่สุดในกลุ่มตัวประมาณค่าที่ไม่เอนเอียงคือตัวประมาณค่าที่มีค่าความแปรปรวนต่ำสุด

4. ความพอเพียง (Sufficiency)

ตัวประมาณค่า $\hat{\theta}$ จะเป็นตัวประมาณค่าที่มีความพอเพียง ถ้าตัวประมาณค่าให้สารสนเทศที่ก่อให้เกิดประโยชน์ได้ทั้งหมดที่ต้องการเกี่ยวกับพารามิเตอร์ที่ต้องการประมาณค่า เช่น \bar{X} เป็นตัวประมาณค่าที่พอเพียงของ μ หมายความว่า ไม่มีตัวประมาณค่าของ μ ตัวอื่น เช่น ฐานนิยม (Mode) มัธยฐาน (Median) ที่จะสามารถให้สารสนเทศเกี่ยวกับ μ เพิ่มขึ้นได้อีกดีกว่า \bar{X}

การศึกษาวិธีการสุ่มตัวอย่างตัวบ่งชี้คุณภาพในการประมาณค่าพารามิเตอร์เมื่อตัวประมาณค่านั้นเป็นค่าที่ไม่ลำเอียง (Unbiasedness) สามารถเลือกใช้เกณฑ์ในการบ่งชี้คุณภาพของวิธีการสุ่มตัวอย่างได้ 2 ประเภท คือ (สมชัย วงษ์นายะ, 2533 อ้างอิงจาก อภิชาติ พงษ์ศรีหิคุลชัย, 2530. หน้า 90-107 ; Jaeger, 1983. pp. 33-35 ; Bickel & Doksum, 1977. pp. 120-130)

1. ความถูกต้อง (Accuracy) เป็นการวัดความใกล้เคียงกันระหว่างค่าประมาณพารามิเตอร์กับค่าพารามิเตอร์ เช่น ใช้ค่าความคลาดเคลื่อนกำลังสองเฉลี่ยของค่าพารามิเตอร์ (Mean Square Error of the Estimate of the Parameter) ค่าความคลาดเคลื่อนสัมบูรณ์เฉลี่ยของค่าประมาณพารามิเตอร์ (Mean Absolute Error of the Estimate of the Parameter) เกณฑ์ความคลาดเคลื่อนกำลังสองเฉลี่ยของค่าประมาณพารามิเตอร์จะคล้ายกับเกณฑ์ค่าเฉลี่ยของกำลังสองของความแตกต่างระหว่างค่าประมาณพารามิเตอร์กับค่าพารามิเตอร์ ส่วนค่าความคลาดเคลื่อนสัมบูรณ์ของค่าประมาณพารามิเตอร์จะคล้ายกับเกณฑ์ส่วนเบี่ยงเบนเฉลี่ยของค่าประมาณพารามิเตอร์

2. ความแม่นยำ (Precision) เป็นการวัดในกรณีที่ไม่ทราบค่าพารามิเตอร์ของประชากร เช่น ใช้ค่าความแปรปรวนของการสุ่ม (Sampling Variance) ค่าความคลาดเคลื่อนมาตรฐาน (Standard Error) ค่าความแม่นยำสัมพัทธ์ (Relative Precision) ค่าสัมประสิทธิ์ความแปรปรวน (Coefficient of Variance)

จากการศึกษางานวิจัยเกี่ยวกับวิธีการจัดการข้อมูลสูญหาย พบว่า มีการใช้เกณฑ์ที่แตกต่างกันไปในการตัดสินว่าวิธีการจัดการข้อมูลสูญหายวิธีใดดีกว่ากัน เช่น งานวิจัยของ รุท

(Roth, 1994, pp. 540-541) ใช้เกณฑ์ความถูกต้อง (Accuracy) โดยนิยามจากขนาดของการกระจายจากค่าพารามิเตอร์ เช่น ความถูกต้องอาจนิยามเป็นการกระจายของความแตกต่างของค่าสัมประสิทธิ์สหสัมพันธ์ที่ได้จากกลุ่มตัวอย่างกับค่าสัมประสิทธิ์สหสัมพันธ์ของประชากร ยังเกอร์และฟอร์ไซท์ (Younger & Forsyth, 1998, p. 203) ได้ศึกษาเปรียบเทียบวิธีการแทนค่าข้อมูลสูญหาย 4 วิธี ในการพยากรณ์ที่มีตัวแปร 2 ตัว เกณฑ์ที่ใช้คือ ค่าเฉลี่ยความแตกต่างของค่าสัมประสิทธิ์การถดถอย (Beta weight) ที่ได้จากการแทนค่าด้วยวิธีการจัดการข้อมูลสูญหายกับค่าสัมประสิทธิ์การถดถอยของประชากร คอเมรย์และฮินส์ (Kromrey & Hines, 1994, p. 577) ได้ศึกษาผลของข้อมูลสูญหายแบบสุ่มในการวิเคราะห์การถดถอยที่มีตัวแปรทำนาย 2 ตัวแปร ประสิทธิภาพของการจัดการข้อมูลสูญหายแบบธรรมดา 5 วิธี พิจารณาจากความแตกต่างของค่าความสัมพันธ์พหุคูณ (R^2) และค่าสัมประสิทธิ์การถดถอยมาตรฐานที่ได้จากวิธีการจัดการข้อมูลสูญหายแบบต่าง ๆ กับค่าที่คำนวณจากข้อมูลสมบูรณ์ ในปี ค.ศ. 1999 ราจเมเคอร์ (Raaijmakers, 1999, p. 734) ได้ศึกษาประสิทธิผลของวิธีการจัดการข้อมูลสูญหายที่แตกต่างกันในการสำรวจโดยใช้แบบสอบถามแบบลิเคอร์ท เกณฑ์ที่ใช้ตัดสินคือ ค่าเฉลี่ยและส่วนเบี่ยงเบนมาตรฐานของความแตกต่างของคะแนนที่ได้จากวิธีการจัดการข้อมูลสูญหายและคะแนนจากข้อมูลสมบูรณ์และเปรียบเทียบค่าสถิติ R^2 , B และค่าสถิติที ที่ได้จากวิธีการจัดการข้อมูลสูญหายและคะแนนจากข้อมูลสมบูรณ์

สำหรับในประเทศไทย ชะไมพร ธรรมวัฒน์ไพศาล (2522, หน้า 124-126) ได้ศึกษาวิธีการประมาณค่าข้อมูลสูญหายในการวิเคราะห์การถดถอยวิธีต่าง ๆ 6 วิธี ด้วยกัน เกณฑ์ที่ใช้ในการเปรียบเทียบ คือ ค่าความแปรปรวนร่วมที่สามารถอธิบายได้ หรือ (R^2) เป็นดัชนีในการตัดสินว่าวิธีประมาณค่าใดสามารถประมาณค่าได้ใกล้เคียงกับค่าที่สูญหายมากกว่ากัน พรศิริ หมื่นไชยศรี (2529, หน้า 75-76) ได้ศึกษาวิธีการจัดการข้อมูลสูญหายในลักษณะเดียวกันกับที่ ชะไมพร ธรรมวัฒน์ไพศาล แต่ใช้เกณฑ์ค่าเฉลี่ยความคลาดเคลื่อนกำลังสอง (Mean Square Error) ส่วน ถวัลย์ จันทร์เพ็ง (2531) ได้ศึกษาเปรียบเทียบความแม่นยำของวิธีประมาณค่าข้อมูลสูญหาย 3 วิธี คือ วิธีค่าเฉลี่ยจากกลุ่มตัวอย่าง วิธีใช้สมการถดถอย และวิธีใช้ค่าเฉลี่ยระหว่างค่าเฉลี่ยจากกลุ่มตัวอย่างและสมการถดถอย เกณฑ์ที่ใช้คือ ค่าเฉลี่ยของผลต่างระหว่างค่าของข้อมูลสูญหายกับค่าที่ประมาณได้จากวิธีการประมาณค่า และค่าเฉลี่ยกำลังสองของผลต่างระหว่างค่าของข้อมูลสูญหายกับค่าที่ประมาณได้จากวิธีการประมาณ และในปี พ.ศ. 2538 วารุณี ศรีบำรุงศักดิ์ (2538, หน้า 5) ได้ศึกษาการพยากรณ์ด้วยวิธีการถดถอยเชิงเส้นพหุเมื่อตัวแปรตามมีค่าสูญหาย เกณฑ์ที่ใช้ตัดสินคือ ค่าความคลาดเคลื่อนระหว่างค่าพยากรณ์ของ

ตัวแปรตามกับค่าจริงในรูปของค่ารากที่สองของค่าเฉลี่ยของความคลาดเคลื่อนกำลังสอง (The Square Root of Mean Squares Error : RMSE) วิธีใดให้ค่า RMSE ต่ำกว่าจะเป็นวิธีการประมาณค่าที่ดีกว่า

เกณฑ์ที่ใช้ตัดสินวิธีการจัดการข้อมูลสูญหายวิธีใดดีกว่ากันสรุปได้ดังนี้

1. ขนาดของการกระจายจากค่าพารามิเตอร์
2. ค่าเฉลี่ยและส่วนเบี่ยงเบนมาตรฐานของความแตกต่างของค่าสถิติที่ได้จากการแทนค่าด้วยวิธีการจัดการข้อมูลสูญหายแบบต่าง ๆ กับค่าพารามิเตอร์หรือค่าที่คำนวณจากข้อมูลสมบูรณ์
3. ความแปรปรวนร่วมที่สามารถอธิบายได้ (R^2)
4. ค่าเฉลี่ยความคลาดเคลื่อนกำลังสอง (Mean square error)
5. ค่าเฉลี่ยของผลต่างระหว่างค่าของข้อมูลสูญหายกับค่าที่ประมาณได้จากวิธีประมาณค่าข้อมูลสูญหายและค่าเฉลี่ยกำลังสองของผลต่างระหว่างค่าของข้อมูลสูญหายกับค่าที่ประมาณได้จากวิธีประมาณค่าข้อมูลสูญหาย
6. ค่าความคลาดเคลื่อนระหว่างค่าพยากรณ์ของตัวแปรตามกับค่าจริงในรูปของค่ารากที่สองของค่าเฉลี่ยของความคลาดเคลื่อนกำลังสอง (The Square Root of Mean Squares Error : RMSE)

จะเห็นว่าโดยส่วนใหญ่จะใช้เกณฑ์ค่าเฉลี่ยความคลาดเคลื่อนกำลังสอง (Mean Square Error) การใช้เกณฑ์ค่าเฉลี่ยความคลาดเคลื่อนกำลังสอง หรือค่าเฉลี่ยความแตกต่างกำลังสองดีกว่าค่าเฉลี่ยความคลาดเคลื่อน หรือค่าเฉลี่ยความแตกต่างระหว่างค่าสถิติที่ได้จากการจัดการข้อมูลสูญหายแบบต่าง ๆ กับค่าที่ได้จากข้อมูลสมบูรณ์ เพราะว่า เมื่อนำมาแยกกำลังสองจะไม่มีค่าติดลบซึ่งทำให้การแปลความหมายง่ายขึ้น และถ้าตัวประมาณค่าไม่เอนเอียง เราจะพบว่าค่าเฉลี่ยความคลาดเคลื่อนกำลังสองจะเท่ากับความแปรปรวนของตัวประมาณค่านั้น และใช้ค่าความแปรปรวนเป็นเกณฑ์ในการพิจารณาตัดสินว่าตัวประมาณค่าดีหรือไม่ดีได้ โดยพิจารณาจากตัวประมาณค่าที่มีความแปรปรวนต่ำสุด (ดวงใจ ปวีณอภิชาติ, 2535. หน้า 37) งานวิจัยที่ใช้เกณฑ์ค่าเฉลี่ยความคลาดเคลื่อนกำลังสอง ดังเช่นงานวิจัยของ รุท (Roth, 1994. pp. 540-541) พรศิริ หมื่นไชยศรี (2529, หน้า 75-76) วารุณี ตรีบำรุงศักดิ์ (2538, หน้า 5) ส่วนงานวิจัยของ ยังเกอร์ และฟอร์ไรท์ (1998, p. 203) และถวัลย์ จันทร์เพ็ง (2531) ใช้เกณฑ์ค่าเฉลี่ยความคลาดเคลื่อนกำลังสอง และค่าเฉลี่ยความคลาดเคลื่อน

ในการวิจัยนี้ผู้วิจัยได้นิยาม ความแม่นยำ หมายถึง ความใกล้เคียงกันระหว่างค่าสถิติ และค่าพารามิเตอร์ซึ่งมีลักษณะเป็นความคลาดเคลื่อนเช่นเดียวกับที่กล่าวมา ดังนั้นผู้วิจัยจึงใช้เกณฑ์ในการเสนอผลการเปรียบเทียบว่าวิธีการจัดการข้อมูลสูญหาย วิธีสุ่มตัวอย่าง ความสัมพันธ์ระหว่างตัวแปร และจำนวนข้อมูลสูญหายที่แตกต่างกัน วิธีที่ดีที่สุด คือ ค่าเฉลี่ยกำลังสองของความคลาดเคลื่อนระหว่างค่าสถิติที่ได้จากวิธีการจัดการข้อมูลสูญหาย วิธีสุ่มตัวอย่าง ความสัมพันธ์ระหว่างตัวแปร และจำนวนข้อมูลสูญหายที่แตกต่างกันกับค่าพารามิเตอร์

เหตุผลที่ผู้วิจัยใช้เกณฑ์เดียวเพราะว่าความคลาดเคลื่อนของค่าสถิติกับค่าพารามิเตอร์ เมื่อนำไปยกกำลังสองค่าที่ได้จะมีค่าเป็นบวกทั้งหมดไม่มีค่าติดลบทำให้การแปลความหมายง่ายขึ้น ค่าที่น้อยและเข้าใกล้ศูนย์จะถือว่ามีความแม่นยำสูง และถ้าค่าเฉลี่ยกำลังสองของความคลาดเคลื่อนของวิธีการจัดการข้อมูลสูญหาย 2 วิธีเท่ากัน วิธีที่ให้อำนาจการทดสอบสูงกว่าจะถือว่าเป็นวิธีที่ดีที่สุด ด้วยเหตุผลที่ว่า ข้อมูลสูญหายทำให้เกิดปัญหา 2 ประการ (Roth, 1994, pp. 538-539) คือ การประมาณค่าพารามิเตอร์มีอคติและอำนาจการทดสอบลดลง การตรวจสอบด้วยค่าเฉลี่ยกำลังสองของความคลาดเคลื่อนเป็นการตรวจสอบด้านการประมาณค่าพารามิเตอร์ ดังนั้นเมื่อค่าเฉลี่ยกำลังสองของความคลาดเคลื่อนเท่ากัน จึงใช้อำนาจการทดสอบเป็นเกณฑ์ในการตัดสินอีกเกณฑ์หนึ่ง

งานวิจัยที่เกี่ยวข้อง

งานวิจัยที่ศึกษาเกี่ยวกับการจัดการข้อมูลสูญหายมีดังต่อไปนี้

ชะไมพร ธรรมวัฒน์ไพศาล (2522, หน้า 124-126) ได้ศึกษาเกี่ยวกับวิธีการประมาณค่าข้อมูลสูญหายในการวิเคราะห์การถดถอยวิธีต่าง ๆ 6 วิธีด้วยกัน คือ วิธีกำลังสองน้อยที่สุด วิธีอันดับศูนย์หรือวิธีใช้ค่าเฉลี่ยจากกลุ่มตัวอย่าง วิธีอันดับศูนย์ตัดแปลง วิธีถดถอยอันดับหนึ่งหรือวิธีใช้สมการถดถอย วิธีถดถอยสองชั้น และวิธีผสมหรือวิธีใช้ค่าเฉลี่ยระหว่างค่าเฉลี่ยจากกลุ่มตัวอย่างและสมการถดถอย จากข้อมูลที่เก็บมาจัดกระทำให้สูญหายแบบสุ่ม เกณฑ์ที่ใช้ในการเปรียบเทียบคือ ค่าความแปรปรวนร่วมที่สามารถอธิบายได้หรือ (R^2) ใช้เป็นดัชนีในการตัดสินใจว่าวิธีประมาณค่าใดสามารถประมาณค่าได้ใกล้เคียงกับค่าที่สูญหายมากกว่ากัน ถ้า (R^2) มีค่าเพิ่มมากขึ้น ผลที่ได้จากงานวิจัยพบว่า วิธีที่ให้ค่า (R^2) สูงกว่าวิธีอื่น ๆ มีอยู่ 3 วิธีเรียงตามลำดับจากมากไปน้อยคือ วิธีถดถอยสองชั้น วิธีถดถอยอันดับหนึ่งหรือวิธีใช้สมการถดถอย และวิธีกำลังสองน้อยที่สุด

พรศิริ หมั่นไชยศรี (2529. หน้า 75-76) ได้ศึกษาเปรียบเทียบวิธีการประมาณค่าข้อมูลที่สูญหายในการวิเคราะห์ตัวแปรพหุคูณ 4 วิธี คือ วิธีใช้ค่าเฉลี่ยจากกลุ่มตัวอย่าง วิธีวิเคราะห์การถดถอยพหุคูณเชิงเส้น วิธีวิเคราะห์ความถดถอยพหุคูณเชิงเส้นดัดแปลง และวิธีวิเคราะห์ส่วนประกอบหลัก ใช้เทคนิคมอนติคาร์โลซิมูเลชัน และใช้เกณฑ์ในการเปรียบเทียบคือ ค่าเฉลี่ยความคลาดเคลื่อนกำลังสอง (mean square error) จากกลุ่มตัวอย่าง 30, 50, 70 100 และ 200 จำนวนตัวแปรเท่ากับ 3, 5, 7 และ 10 ตัวแปร ค่าสัมประสิทธิ์สหสัมพันธ์ระหว่างตัวแปรเท่ากับ .10, .20,90 และกำหนดสัดส่วนข้อมูลที่สูญหายของแต่ละตัวแปรมีค่าเท่ากับ 10% จากการศึกษาพบว่าวิธีประมาณค่าข้อมูลที่สูญหายในการวิเคราะห์ตัวแปรพหุคูณทั้ง 4 วิธี ให้ค่าเฉลี่ยความคลาดเคลื่อนกำลังสองไม่แตกต่างกันอย่างมีนัยสำคัญที่ระดับ .05 ไม่ว่าจะเป็สถานการณ์ใดก็ตามที่มีข้อมูลสูญหายเกิดขึ้น

ถวัลย์ จันทร์เพ็ง (2531) ได้ศึกษาเปรียบเทียบความแม่นยำของวิธีประมาณค่าข้อมูลที่สูญหาย 3 วิธี คือ วิธีใช้ค่าเฉลี่ยจากกลุ่มตัวอย่าง วิธีใช้สมการถดถอย และวิธีใช้ค่าเฉลี่ยระหว่างค่าเฉลี่ยจากกลุ่มตัวอย่างและสมการถดถอย จากกลุ่มตัวอย่างขนาด 5, 10 และ 15 กำหนดจำนวนข้อมูลสูญหายครั้งละ 1 ค่า และ 2 ค่า โดยใช้ข้อมูลที่มีลักษณะการแจกแจงแบบปกติและในกรณีที่ใช้สมการพยากรณ์ช่วยในการประมาณค่าได้กำหนดค่าสัมประสิทธิ์สหสัมพันธ์ของตัวแปรเกณฑ์กับตัวพยากรณ์เท่ากับ 0.2, 0.4 และ 0.6 ทำการทดลองด้วยเทคนิคมอนติคาร์โลซิมูเลชัน โดยการจำลองด้วยเครื่องคอมพิวเตอร์ ผลที่ได้พบว่า วิธีที่ใช้ค่าเฉลี่ยจากกลุ่มตัวอย่างประมาณค่าได้ไม่ดีเท่ากับวิธีใช้สมการถดถอย และวิธีใช้ค่าเฉลี่ยระหว่างค่าเฉลี่ยจากกลุ่มตัวอย่างและสมการถดถอย

วารุณี ตรีบำรุงศักดิ์ (2538) ได้ศึกษาการพยากรณ์ด้วยวิธีการถดถอยเชิงเส้นพหุ เมื่อตัวแปรตามมีค่าสูญหาย วิธีที่ใช้ประมาณค่าตัวแปรตามเมื่อมีข้อมูลสูญหาย คือ วิธีค่าเฉลี่ย วิธีสมการถดถอย วิธีอีเอ็ม (EM algorithm) และวิธีของฮันท์ (Hunt's Method) ใช้กลุ่มตัวอย่างขนาด 10, 20, 30, 50, และ 70 ค่าเบี่ยงเบนมาตรฐานของความคลาดเคลื่อน 5, 10, 15, 20 และ 25 สัดส่วนการสูญหายของตัวแปรตาม 10%, 20%, 30%, 40%, 50%, 60% และ 70% ข้อมูลที่ใช้ในการวิจัยได้จากการจำลองด้วยเทคนิคมอนติคาร์โล และหารากที่สองของค่าเฉลี่ยของความคลาดเคลื่อนกำลังสองของค่าพยากรณ์ ผลการศึกษาพบว่า วิธีการของฮันท์เป็นวิธีการที่ดีเมื่อกลุ่มตัวอย่างมีขนาดเล็ก ความคลาดเคลื่อนน้อย และสัดส่วนการสูญหายมาก แต่ถ้าความคลาดเคลื่อนสูง วิธีค่าเฉลี่ยจะเป็นวิธีการที่ดีในทุกสัดส่วนการสูญหายของตัวแปรตาม แต่ถ้การสูญหายมีขนาดใหญ่ขึ้นวิธีสูญหายจะเหมาะสมเกือบทุกกรณี

งานวิจัยเกี่ยวกับข้อมูลสูญหายในประเทศไทยยังมีน้อยโดยส่วนใหญ่จะศึกษากับข้อมูลในสถานการณ์จำลอง ในการวิเคราะห์ตัวแปรพหุคูณ และเป็นงานวิจัยเชิงทดลอง แต่ไม่พบงานวิจัยที่ศึกษากับงานวิจัยเชิงสำรวจที่มีข้อมูลสูญหายค่อนข้างมาก

ชานและตัน (ถวัลย์ จันทน์เพ็ง, 2531. หน้า 22 อ้างอิงจาก Chan & Dunn, 1972. pp. 473-477) ได้ศึกษาเปรียบเทียบวิธีการประมาณค่าข้อมูลสูญหายในการวิเคราะห์จำแนกประเภทจำนวน 5 วิธี โดยใช้เทคนิคมอนติคาร์โลซิมูเลชัน คือ 1. ศึกษาเมื่อไม่มีข้อมูลสูญหายเลย 2. วิธีตัดตัวอย่างที่มีข้อมูลสูญหายออก 3. วิธีใช้ค่าเฉลี่ยจากกลุ่มตัวอย่าง 4. วิธีใช้สมการถดถอย และ 5. วิธีวิเคราะห์ส่วนประกอบหลัก โดยศึกษาข้อมูลที่มีลักษณะการแจกแจงแบบปกติ 2 ตัวแปรและประชากรทั้งสองกลุ่มมีความแปรปรวนเท่ากัน เกณฑ์ในการพิจารณาคือร้อยละของการจำแนกผิด ผลการวิจัยพบว่า ค่า (R^2) ของวิธีใช้สมการถดถอยจะสูงกว่าวิธีใช้ค่าเฉลี่ย วิธีวิเคราะห์ส่วนประกอบหลัก และวิธีตัดตัวอย่างที่มีข้อมูลสูญหายออก และถ้าจำนวนตัวแปรมากขึ้นวิธีใช้สมการถดถอยจะ ให้ค่า (R^2) เพิ่มสูงขึ้น

ฟิงค์ไบเนอร์ (Finkbeiner, 1979. pp. 416-420) ศึกษาเปรียบเทียบวิธีประมาณค่าข้อมูลสูญหายในการวิเคราะห์ตัวแปรพหุคูณ 6 วิธีด้วยกัน คือ 1. วิธี maximum likelihood (ML) 2. วิธีใช้ค่าเฉลี่ยจากกลุ่มตัวอย่าง (MR) 3. วิธีที่ใช้เฉพาะตัวอย่างที่ไม่มีข้อมูลสูญหาย (CD) 4. วิธีใช้สมการถดถอย (REG) 5. วิธีวิเคราะห์องค์ประกอบหลัก (PC) และ 6. วิธีที่ใช้เฉพาะคู่ตัวอย่างสมบูรณ์ (CPO) โดยใช้เทคนิคมอนติคาร์โลซิมูเลชันในการจำลองข้อมูล ขนาดกลุ่มตัวอย่างเท่ากับ 64 จำนวนการสูญหายของข้อมูล 2 รูปแบบ ทดลองซ้ำ 50 ครั้ง เกณฑ์ที่ใช้ในการเปรียบเทียบคือ ค่าเฉลี่ยจากกลุ่มตัวอย่างและการกระจายของค่าพารามิเตอร์ และใช้ค่าเฉลี่ยความคลาดเคลื่อนกำลังสองของวิธีที่มีค่าน้อยที่สุด จากการศึกษาพบว่า ค่าเฉลี่ยจากกลุ่มตัวอย่างและค่าส่วนเบี่ยงเบนมาตรฐานของ ทุกวิธีไม่แตกต่างกัน แต่ค่าเฉลี่ยความคลาดเคลื่อนกำลังสองแตกต่างกันเรียงตามลำดับจากน้อยไปมากคือ ML, MR, CPO และ REG สรุปได้ว่าวิธีที่สามารถประมาณค่าข้อมูลสูญหายได้ดีที่สุดคือวิธี maximum likelihood

คอมเรย์และฮินส์ (Kromrey & Hines, 1991. pp. 13-16) ได้ศึกษาผลของข้อมูลสูญหายแบบสุ่มในการวิเคราะห์การถดถอยที่มีตัวแปรทำนาย 2 ตัว ประสิทธิภาพของการจัดการข้อมูลสูญหายแบบธรรมดา 5 วิธี พิจารณาจากความแตกต่างของค่าความสัมพันธ์พหุคูณ (R^2) และค่าสัมประสิทธิ์การถดถอยมาตรฐาน ขนาดกลุ่มตัวอย่างที่ใช้ในการศึกษา คือ 50, 100 และ 200 สุ่มจากข้อมูลจริง แล้วสร้างข้อมูลสูญหายในแต่ละกลุ่มตัวอย่าง แล้วเปรียบเทียบค่าประมาณพารามิเตอร์ที่ได้จากกลุ่มตัวอย่างที่ไม่มีข้อมูลสูญหาย ผลการศึกษาพบว่า วิธีการแทน

ค่า 3 วิธี คือ การแทนค่าข้อมูลสูญหายด้วยค่าเฉลี่ย การแทนค่าข้อมูลสูญหายโดยใช้วิธีการวิเคราะห์การถดถอยอย่างง่ายและพหุคูณ ประมาณค่าความสัมพันธ์พหุคูณ (R^2) และสัมประสิทธิ์การถดถอยอย่างมีอคติ วิธีการจัดการข้อมูลสูญหายแบบตัดออก คือ ลิสต์ไวส์และแพร์ไวส์ (Listwise and Pairwise deletion) ประมาณค่าพารามิเตอร์ได้ถูกต้องเมื่อข้อมูลสูญหายเท่ากับ 60%

วิททาและไคเซอร์ (Witta & Kaizer, 1991) ได้ศึกษาวิธีการจัดการข้อมูลสูญหายโดยใช้วิธีการจัดการข้อมูลสูญหายเหมือนกันกับที่คอมพิวเตอร์และฮาร์ดแวร์ พบว่า มีความแตกต่างระหว่างค่าเฉลี่ยประชากรและค่าเฉลี่ยที่ได้จากการแทนค่าข้อมูลสูญหายด้วยวิธีค่าเฉลี่ยอย่างมีนัยสำคัญทางสถิติที่ระดับ .01 วิธีการอื่น ๆ ไม่แตกต่างกัน ดังนั้น วิธีการแทนค่าข้อมูลสูญหายด้วยค่าเฉลี่ยไม่เหมาะสมในการจัดการข้อมูลสูญหาย

มาร์แคนโทนีโอ (Marcantonio, 1992. <http://thailis.uni.net.th/dao/detail.nsp>) ได้ศึกษาวิธีการจัดการข้อมูลสูญหายหลายวิธีจากวิธีแบบเก่าจนถึงวิธีการแบบใหม่ เช่น วิธีอีเอ็ม (Em algoritm) โดยพิจารณาประสิทธิภาพของสัมประสิทธิ์การถดถอยภายใต้เงื่อนไขที่แตกต่างกัน คือ จำนวนข้อมูลสูญหาย ขนาดของกลุ่มตัวอย่าง ประเภทของการสูญหาย และการแจกแจงแบบปกติหลายตัวแปร ความสัมพันธ์ระหว่างตัวแปรมีค่าจาก .10 ถึง .55 ตัวแปรตามที่ศึกษา ได้แก่ สัมประสิทธิ์การถดถอยของคะแนนดิบ ความคลาดเคลื่อนกำลังสองเฉลี่ย และความแปรปรวนของสัมประสิทธิ์การถดถอยคะแนนดิบ ใช้ค่าสัมประสิทธิ์จากการตัดข้อมูลสูญหายออกแบบลิสต์ไวส์เป็นตัวเปรียบเทียบ ผลการวิจัยพบว่า ประสิทธิภาพสัมพัทธ์ที่ดีที่สุดขึ้นอยู่กับข้อมูลสูญหายเป็นแบบสุ่ม (missing at random) วิธีอีเอ็มจะดีในทุกสถานการณ์ เมื่อสาเหตุของการสูญหายเป็นแบบสุ่มอย่างสมบูรณ์ (missing completely at random) วิธีการถดถอยจะดี วิธีการตัดข้อมูลสูญหายออกแบบลิสต์ไวส์ไม่มีประสิทธิภาพมากที่สุดเมื่อข้อมูลสูญหายเป็นแบบสุ่ม วิธีการถดถอยจะไม่มีประสิทธิภาพและไม่มีประสิทธิภาพมากขึ้นเมื่อขนาดกลุ่มตัวอย่างมากขึ้น ผลของความเป็นปกติหลายตัวแปรมีค่าน้อย แสดงว่าประสิทธิภาพของสัมประสิทธิ์ต่าง ๆ มีผลน้อยจากความเป็นปกติของข้อมูล

สุนันทา วีรกุลเทวีญ (Viragoontavan, 2000. <http://thailis.uni.net.th/dao/detail.nsp>) ได้ศึกษาเปรียบเทียบประสิทธิภาพสัมพัทธ์ (relative effectiveness) ของวิธีการจัดการข้อมูลสูญหาย 6 วิธี คือ การตัดข้อมูลสูญหายออกแบบลิสต์ไวส์ การแทนค่าด้วยค่าเฉลี่ยของกลุ่ม การแทนค่าด้วยวิธีการถดถอย การแทนค่าด้วยวิธีฮอทเดค (hot-deck imputation) การแทนค่าวิธีเอ็มไอด้วยโปรแกรม SOLAR และการแทนค่าวิธีเอ็มไอด้วยโปรแกรม NORM ข้อมูลได้จาก

สถานการณ์จำลองตามเงื่อนไขต่อไปนี้ ความสัมพันธ์ระหว่างตัวแปร (ต่ำ ปานกลาง สูง) ขนาดกลุ่มตัวอย่าง (100 200 500) และสัดส่วนข้อมูลสูญหาย (.05 .10 .20) จัดกระทำข้อมูล สมบูรณ์ให้เป็นข้อมูลสูญหายแบบสุ่มตามสัดส่วนที่กำหนดแล้วจัดกระทำ ข้อมูลสูญหายด้วยวิธีการทั้ง 6 วิธี นำข้อมูลไปวิเคราะห์เชิงจำแนกแล้วประเมินประสิทธิภาพสัมพัทธ์โดยพิจารณาอัตรา (hit rate) และอำนาจการจำแนก (discriminating power) ของสมการจำแนกสมการแรก ผลการวิจัยพบว่า วิธีการจัดการข้อมูลสูญหายแบบเอ็มไอทั้งสองวิธีมีประสิทธิภาพมากที่สุด การตัดข้อมูลสูญหายออกแบบลิสต์ไวส์มีประสิทธิภาพต่ำสุด วิธีการจัดการข้อมูลสูญหายทั้งหมดเมื่อข้อมูลมีความสัมพันธ์กันน้อยประมาณค่าได้ถูกต้องมากกว่าเมื่อข้อมูลมีความสัมพันธ์กันสูง การประมาณค่าอัตราและอำนาจการจำแนกได้ถูกต้องมากขึ้นเมื่อจำนวนกลุ่มตัวอย่างมากขึ้นและจำนวนข้อมูลสูญหายน้อย

เอ็นเดอร์ (Enders, 2001, pp. 713-740) ได้ศึกษาความสามารถของการประมาณค่าแบบเอฟไอเอ็มแอล (Full information maximum likelihood : FIML) ในการวิเคราะห์การถดถอยพหุคูณเมื่อมีข้อมูลสูญหาย ตัวแปรที่ศึกษามี 4 ตัวแปร คือ วิธีการจัดการข้อมูลสูญหาย จำนวนข้อมูลสูญหาย ขนาดของกลุ่มตัวอย่าง และขนาดของความสัมพันธ์ระหว่างตัวแปร ใช้เทคนิคอนติ คาร์โล ซิมูเลชัน จำลองข้อมูลที่มีรูปแบบของการสูญหายแตกต่างกัน 3 แบบ คือ การสูญหายแบบสุ่มสมบูรณ์ (missing completely at random) การสูญหายแบบสุ่ม (missing at random) และการสูญหายไม่เป็นแบบสุ่ม (nonrandom pattern) ตัวแปรตามที่ศึกษา คือ สัมประสิทธิ์การถดถอย สัมประสิทธิ์การตัดสินใจ (R^2) และประสิทธิภาพ นอกจากนี้ เอ็นเดอร์ ยังได้ศึกษาปฏิสัมพันธ์ของ วิธีการจัดการข้อมูลสูญหาย จำนวนข้อมูลสูญหาย และขนาดของความสัมพันธ์ระหว่างตัวแปร ผลการวิจัยพบว่า การประมาณค่าข้อมูลสูญหายแบบเอฟไอเอ็มแอล ดีกว่าการตัดข้อมูลสูญหายออกแบบลิสต์ไวส์ แพร่ไวส์ การแทนค่าข้อมูลด้วยค่าเฉลี่ยและการประมาณค่าข้อมูลสูญหายแบบเอฟไอเอ็มแอลมีอคติน้อย และมีปฏิสัมพันธ์ระหว่างวิธีการจัดการข้อมูลสูญหาย จำนวนข้อมูลสูญหาย และขนาดของความสัมพันธ์ระหว่างตัวแปร

ซู (Zhou, 2002) ได้ศึกษาเปรียบเทียบวิธีการแทนค่าข้อมูลสูญหาย 4 วิธี ในการประมาณค่าพารามิเตอร์ของแบบวัดแบบลิเคอร์ท โดยมีวัตถุประสงค์ที่ต้องการประเมินประสิทธิภาพของวิธีการแทนค่าข้อมูลสูญหาย 4 วิธี ภายใต้เงื่อนไขของข้อมูลที่หลากหลายและสาเหตุของการสูญหายแบบ MCAR (Missing completely at random) ใช้การวิจัยเชิงทดลองแบบแฟคทอเรียล ($4 \times 3 \times 3$) มีวิธีการจัดการข้อมูลสูญหาย 4 วิธี เปอร์เซ็นต์ของกลุ่มตัวอย่างที่มีข้อมูลสูญหาย 3 ระดับ และเปอร์เซ็นต์ของข้อคำถามที่มีข้อมูลสูญหาย 3 ระดับ เช่นเดียวกัน ข้อมูลที่ใช้ศึกษา

เป็นข้อมูลทุติยภูมิ (secondary data) สุ่มมา 500 คน จากคนทั้งหมด 2,058 คน ผลการวิจัยพบว่า ความแม่นยำที่ได้จากวิธีการแทนค่า 4 วิธี มีความแตกต่างกันอย่างมีนัยสำคัญทางสถิติ โดยไม่พิจารณาเงื่อนไขของข้อมูล วิธีการจัดการข้อมูลสูญหายแบบ IMS (Item mean substitution) เลวที่สุด PMS (Person mean substitution) เลวเป็นลำดับที่สองในการประมาณค่าพารามิเตอร์ ส่วนวิธี EM algorithm และ SRI (Stochastic regression imputation) ให้การประมาณค่าความแม่นยำได้ถูกต้องมากขึ้น วิธีการแทนค่าข้อมูลสูญหายแบบ EM algorithm ดีที่สุดในการประมาณค่าเฉลี่ยของข้อคำถาม (item mean) ความสัมพันธ์ระหว่างข้อคำถาม (correlations between item) และความสัมพันธ์ระหว่างมาตรวัด (correlations between scale) วิธีการแทนค่าข้อมูลสูญหายแบบ SRI ดีที่สุดในการประมาณค่าส่วนเบี่ยงเบนมาตรฐานของข้อคำถาม และความเที่ยง

จากการศึกษางานวิจัยในประเทศไทยและต่างประเทศ พบว่า งานวิจัยเกี่ยวกับข้อมูลสูญหายส่วนใหญ่จะศึกษาโดยใช้เทคนิคมอนติ คาร์โล ซิมูเลชัน จำลองข้อมูล ส่วนมากจะศึกษากับกลุ่มตัวอย่างขนาดเล็ก แต่ไม่พบว่ามีการศึกษาเปรียบเทียบในกรณีที่มีการสุ่มตัวอย่างแตกต่างกัน เพราะการสุ่มตัวอย่างด้วยวิธีการที่แตกต่างกันจะประมาณค่าพารามิเตอร์ได้แตกต่างกัน เมื่อใช้วิธีการจัดการข้อมูลสูญหายแตกต่างกันก็อาจจะทำให้การประมาณค่าพารามิเตอร์แตกต่างกันด้วย นอกจากนั้นยังมีตัวแปรอื่น ๆ เข้ามาเกี่ยวข้องทำให้ผลของวิธีการจัดการข้อมูลสูญหายแตกต่างกัน เช่น จำนวนข้อมูลสูญหาย ความสัมพันธ์ระหว่างตัวแปร ดังนั้นจึงควรที่จะศึกษาปฏิสัมพันธ์ระหว่างวิธีการจัดการ ข้อมูลสูญหายกับตัวแปรอื่น ๆ ด้วย ซึ่งจะทำให้ข้อค้นพบที่ได้สามารถนำไปประยุกต์ใช้ได้อย่างมีประสิทธิภาพ

กรอบแนวคิดที่ใช้ในการวิจัย

จากการศึกษาเอกสารและงานวิจัยที่เกี่ยวข้องกับการจัดการข้อมูลสูญหาย มีวิธีการจัดการข้อมูลสูญหายหลายวิธีทั้งวิธีธรรมดา เช่น การตัดข้อมูลสูญหายออกแบบลิสต์ไวส์ แพร์ไวส์ ที่มีใช้ในโปรแกรมคอมพิวเตอร์ทั่วไป และวิธีที่ต้องใช้แนวคิดทางสถิติขั้นสูงประมาณค่าพารามิเตอร์ด้วยการทำซ้ำ และพบว่าความแตกต่างของวิธีการต่าง ๆ จะน้อยลงด้วยปัจจัยต่อไปนี้ (Raaijmakers, 1999. p. 728)

1. ขนาดของกลุ่มตัวอย่างมากขึ้น
2. จำนวนเปอร์เซ็นต์ของการสูญหายน้อย

3. ตัวแปรสูญหายน้อย
4. การลดลงในระดับความสัมพันธ์ระหว่างตัวแปร

และจากการศึกษาของคอมเรย์ และฮินส์ (Kromrey & Hines, 1994. p. 575) กล่าวว่า มีวิธีการจัดการข้อมูลสูญหายเป็นจำนวนมากการเลือกใช้จึงเกิดความสับสนว่าวิธีใดดีที่สุดกับ ลักษณะของข้อมูล และจำนวนข้อมูลสูญหาย ดังนั้นการศึกษาวิธีการจัดการข้อมูลสูญหายจะต้อง ศึกษา 3 ลักษณะ คือ

1. ศึกษาว่าวิธีการจัดการข้อมูลสูญหายมีผลกระทบต่อการประมาณค่าพารามิเตอร์หรือไม่
2. เมื่อใช้วิธีการสุ่มที่แตกต่างกันจะมีผลกระทบต่อการประมาณค่าพารามิเตอร์หรือไม่
3. ข้อมูลที่นำมาศึกษาควรจะเป็นข้อมูลในสถานการณ์จริง

จะเห็นว่าการศึกษาวิธีการจัดการข้อมูลสูญหายมีตัวแปรอื่นๆ เข้ามาเกี่ยวข้องหลาย ตัวแปร แต่เมื่อพิจารณาจากขั้นตอนการทำวิจัยแล้ว ผู้วิจัยจะต้องสุ่มตัวอย่างประชากรขึ้นมา ศึกษา ข้อมูลที่สุ่มขึ้นมา มีลักษณะความสัมพันธ์แตกต่างกันไป และจำนวนข้อมูลสูญหายก็จะแตกต่างกัน ซึ่งขึ้นอยู่กับข้อคำถามและความตั้งใจของผู้ตอบ ตัวแปรดังกล่าวน่าจะมีผลร่วมกันในการประมาณค่าพารามิเตอร์ ซึ่งสอดคล้องกับงานวิจัยของ เอ็นเดอร์ (Enders, 2001. pp. 713-740) ที่ศึกษาปฏิสัมพันธ์ของวิธีการจัดการข้อมูลสูญหาย จำนวนข้อมูลสูญหาย และขนาดของ ความสัมพันธ์ระหว่างตัวแปร ผลการวิจัยพบว่า มีปฏิสัมพันธ์ระหว่างวิธีการจัดการข้อมูลสูญหาย จำนวนข้อมูลสูญหาย และขนาดของความสัมพันธ์ระหว่างตัวแปร ดังนั้นในการศึกษาวิจัยครั้งนี้ ผู้วิจัยจึงกำหนดตัวแปรอิสระ 4 ตัวแปร ได้แก่

1. วิธีการสุ่มตัวอย่าง
2. ความสัมพันธ์ระหว่างตัวแปร
3. จำนวนข้อมูลสูญหาย
4. วิธีการจัดการข้อมูลสูญหาย

โดยมีรายละเอียดในการกำหนดระดับของแต่ละตัวแปรดังนี้

1. วิธีการสุ่มตัวอย่าง จากการศึกษางานวิจัยที่ศึกษาเปรียบเทียบค่าประมาณ พารามิเตอร์จากแบบแผนการสุ่มตัวอย่างต่างแบบของ นิเวศ คำรัตน์ (2534) โดยศึกษาเปรียบเทียบค่าประมาณมัชฌิมเลขคณิตของอายุประชากรตามคุณสมบัติของตัวประมาณค่าที่ดี 3 ด้าน

โดยใช้การสุ่มตัวอย่างแบบง่าย แบบมีระบบ แบบแบ่งชั้น แบบสองชั้น และแบบสามชั้น พบว่าวิธีการสุ่มตัวอย่างแบบแบ่งชั้นให้ค่าประมาณมัชฌิมเลขคณิตของประชากรมีประสิทธิภาพมากที่สุด สมชัย วงษ์นายะ (2533) ได้ศึกษาเปรียบเทียบค่าประมาณพารามิเตอร์ จากแบบแผนการสุ่มตัวอย่างต่างแบบรวม 7 วิธี ได้แก่ การสุ่มตัวอย่างแบบง่าย แบบมีระบบ แบบแบ่งชั้นที่ใช้ตัวแปรจำแนกชั้นต่างกัน คือ ขนาดโรงเรียน คุณภาพของโรงเรียน อำเภอ และการสุ่มตัวอย่างแบบแบ่งชั้น 2 ระยะ โดยใช้ตัวแปรจำแนกชั้นอำเภอ ขนาดโรงเรียน และตัวแปรจำแนกชั้นอำเภอ คุณภาพของโรงเรียน ผลการศึกษา พบว่า วิธีการสุ่มตัวอย่างแบบมีระบบให้ค่าเฉลี่ยของกำลังสองของความแตกต่างระหว่างค่าประมาณพารามิเตอร์กับค่าพารามิเตอร์ และค่าส่วนเบี่ยงเบนเฉลี่ยของค่าประมาณพารามิเตอร์น้อยที่สุด ดวงใจ ปวีณอภิชาติ (2535) ศึกษาเปรียบเทียบค่าประมาณพารามิเตอร์ของแผนการสุ่มแบบแบ่งชั้น ที่มีตัวแปรจำแนกชั้นภูมิภาค และวิธีการกำหนดขนาดของกลุ่มตัวอย่างย่อยที่แตกต่างกัน พบว่า วิธีการสุ่มที่ใช้ขนาดของโรงเรียนเป็นตัวแปรจำแนกชั้นและกำหนดขนาดกลุ่มตัวอย่างย่อยแบบนี้มีประสิทธิภาพมากที่สุด ในการประมาณค่ามัชฌิมเลขคณิต และประมาณค่าความแปรปรวน และสุกัญญรัตน์ คงงาม (2539) ได้ศึกษาเปรียบเทียบคุณสมบัติของตัวประมาณค่าพารามิเตอร์ที่ได้จากกลุ่มตัวอย่างสุ่มแบบหลายชั้นตอนที่แตกต่างกัน 3 วิธี คือ วิธีสุ่มแบบแบ่งชั้น วิธีสุ่มแบบแบ่งชั้น 2 ระยะ และวิธีสุ่มแบบกลุ่ม ใช้การสุ่มตัวอย่างย่อยที่แตกต่างกัน 2 วิธี คือ วิธีสุ่มแบบง่ายกับแบบมีระบบ ผลการวิจัยพบว่าเมื่อใช้วิธีการสุ่มตัวอย่างในหน่วยใหญ่แบบแบ่งชั้น และแบบแบ่งชั้น 2 ระยะ ใช้วิธีสุ่มตัวอย่างย่อยแบบมีระบบให้ค่าประมาณด้านความไม่เอนเอียง ความคงเส้นคงวา และความมีประสิทธิภาพสูงกว่าค่าประมาณจากวิธีการสุ่มตัวอย่างย่อยแบบง่าย แต่เมื่อสุ่มด้วยวิธีสุ่มแบบกลุ่มพบว่าวิธีการสุ่มตัวอย่างย่อยแบบง่ายจะให้ค่าประมาณที่มีคุณสมบัติในการประมาณค่าพารามิเตอร์สูงกว่าค่าประมาณจากวิธีการสุ่มตัวอย่างย่อยแบบมีระบบ จากการศึกษางานวิจัยเกี่ยวกับการสุ่มตัวอย่าง พบว่า การสุ่มตัวอย่างต่างกันประมาณค่าพารามิเตอร์ได้แตกต่างกัน ดังนั้นในการศึกษาวิจัยครั้งนี้ผู้วิจัยจึงกำหนดวิธีการสุ่มตัวอย่างดังต่อไปนี้

1.1 วิธีการสุ่มแบบแบ่งชั้น มีขั้นตอนการสุ่มดังนี้

1.1.1 แบ่งประชากรออกเป็นระดับชั้น 3 ระดับชั้น

1.1.2 กำหนดขนาดกลุ่มตัวอย่างในแต่ละระดับชั้นแบบนี้แมน

1.1.3 สุ่มตัวอย่างจากประชากรที่ได้จำแนกไว้ในข้อที่ 1.1.1 ด้วยขนาดที่กำหนด

ไว้ในแต่ละระดับชั้นในข้อที่ 1.1.2 โดยใช้วิธีการสุ่มแบบเป็นระบบ

1.2 วิธีการสุ่มแบบกลุ่ม มีขั้นตอนการสุ่มดังนี้

1.2.1 แบ่งประชากรออกเป็นกลุ่มทั้งหมด 14 กลุ่ม

1.2.2 สุ่มกลุ่มออกมาครั้งหนึ่งโดยการสุ่มอย่างง่าย

1.2.3 กำหนดขนาดกลุ่มตัวอย่างในแต่ละกลุ่มแบบนีย์แมน

1.2.4 สุ่มตัวอย่างตามกลุ่มที่สุ่มได้ในข้อที่ 1.2.2 ตามขนาดที่กำหนดไว้ในข้อที่

1.2.3 โดยการสุ่มอย่างง่าย

1.3 วิธีการสุ่มแบบหลายขั้นตอน มีขั้นตอนการสุ่มดังนี้

1.3.1 แบ่งประชากรออกเป็นกลุ่มทั้งหมด 14 กลุ่ม

1.3.2 สุ่มกลุ่มออกมาครั้งหนึ่งโดยการสุ่มอย่างง่าย

1.3.3 ในแต่ละกลุ่มแบ่งประชากรออกเป็นระดับชั้น 3 ระดับชั้น

1.3.4 กำหนดขนาดกลุ่มตัวอย่างในแต่ละระดับชั้นแบบนีย์แมน

1.3.5 สุ่มตัวอย่างจากประชากรที่จำแนกไว้ในข้อที่ 1.3.3 ด้วยขนาดที่กำหนดไว้

ในข้อที่ 1.3.4 โดยวิธีการสุ่มแบบเป็นระบบ

2. ความสัมพันธ์ระหว่างตัวแปร ความสัมพันธ์ระหว่างตัวแปรมีบทบาทสำคัญที่จะทำให้วิธีการจัดการข้อมูลสูญหายให้ผลแตกต่างกัน ดังเช่นงานวิจัยของ คอมเรย์และฮินส์ (Kromrey and Hines, 1994. p. 575) ที่ได้ศึกษาเกี่ยวกับวิธีการจัดการข้อมูลสูญหายแล้วสรุปไว้ว่า การเลือกใช้วิธีการจัดการข้อมูลสูญหายจะต้องพิจารณาถึงโครงสร้างของข้อมูลและระดับของข้อมูลสูญหาย ราจเมเคอร์ (Raaijmakers, 1999. p. 728) กล่าวว่า มีวิธีการจัดการข้อมูลสูญหายหลายวิธีและพบว่าความแตกต่างของวิธีการต่าง ๆ จะน้อยลงเมื่อความสัมพันธ์ระหว่างตัวแปรน้อยลง สอดคล้องกับงานวิจัยของ สุนันทา วีระกุลเทวีญ (Viragoontavan, 2000) ที่ได้ศึกษาประสิทธิภาพสัมพัทธ์ (relative effectiveness) ของวิธีการจัดการข้อมูลสูญหาย 6 วิธี คือ การตัดข้อมูลสูญหายออกแบบลิสท์ไวส์ การแทนค่าด้วยค่าเฉลี่ยของกลุ่ม การแทนค่าด้วยวิธีการถดถอย การแทนค่าด้วยวิธีการฮอทเดค (hot-deck imputation) การแทนค่าด้วยวิธีเอ็มไอใช้โปรแกรม SOLAR และการแทนค่าด้วยวิธีเอ็มไอใช้โปรแกรม NORM พบว่า วิธีการจัดการข้อมูลสูญหายทั้งหมดเมื่อข้อมูลมีความสัมพันธ์กันน้อยประมาณค่าได้ถูกต้องมากกว่าเมื่อข้อมูลมีความสัมพันธ์กันมาก ซึ่งข้อค้นพบของนักวิจัยทั้ง 2 ท่านดังกล่าวขัดแย้งกับงานวิจัยของ รุท (Roth, 1994. p. 537-560) ที่พบว่า ถ้าความสัมพันธ์ระหว่างตัวแปรมีค่ามาก วิธีการจัดการข้อมูลสูญหายแบบการถดถอยจะดีกว่าเมื่อข้อมูลมีความสัมพันธ์ระหว่างตัวแปรน้อย แต่ทั้งนี้ก็ยังขึ้นอยู่กับตัวแปรอื่น ๆ ด้วย

ดังนั้นนักวิจัยที่ศึกษาวิธีการจัดการข้อมูลสูญหายจะนำลักษณะความสัมพันธ์ระหว่างตัวแปรเข้ามาศึกษาด้วย (พรศิริ หมั่นไชยศรี, 2529 ; ถวัลย์ จันทร์เพ็ง, 2531 ; Marcantino, 1992 ; Viragoontavan, 2000 ; Enders, 2001) ลักษณะความสัมพันธ์ระหว่างตัวแปรที่นำมาศึกษาในการวิจัยครั้งนี้ผู้วิจัยกำหนดขนาดความสัมพันธ์โดยพิจารณาจากความแปรปรวนร่วม (common variance หรือ shared variance) หรือเรียกว่า สัมประสิทธิ์การทำนาย (coefficient of determination, r^2) เป็นการวัดขนาดของความแปรปรวนของตัวแปรหนึ่งซึ่งอธิบายได้จากความแปรปรวนของอีกตัวแปรหนึ่ง ความสัมพันธ์ที่มีค่า .30 จะมีค่า $r^2 = .09$ แสดงว่ามีความแปรปรวนร่วมกัน (command variance) เท่ากับ 9% หมายความว่า มีความสัมพันธ์กันน้อย ความสัมพันธ์ที่มีค่า .50 จะมีค่า $r^2 = .25$ มีความแปรปรวนร่วมกันเท่ากับ 25% แสดงว่า มีความสัมพันธ์กันระดับปานกลาง และความสัมพันธ์ที่มีค่า .70 จะมีค่า $r^2 = .49$ มีความแปรปรวนร่วมกันเท่ากับ 49% แสดงว่า มีความสัมพันธ์กันมาก (Grimm, 1992. p. 377) และนอกจากจะพิจารณาความสัมพันธ์จากค่าสัมประสิทธิ์การทำนายแล้ว ผู้วิจัยยังพิจารณาขนาดความสัมพันธ์ให้สอดคล้องกับงานวิจัยที่ศึกษาเกี่ยวกับวิธีการจัดการข้อมูลสูญหายที่ผ่านมา ดังเช่นงานวิจัยของ รุท (Roth, 1994. p. 543) ที่พบว่า วิธีการจัดการ ข้อมูลสูญหายแบบการถดถอยจะมีประสิทธิภาพสูงเมื่อความสัมพันธ์ระหว่างตัวแปรมีค่ามาก แต่จะมีประสิทธิภาพต่ำเมื่อความสัมพันธ์ระหว่างตัวแปรมีค่าน้อย ($r=.20$ ถึง $r=.30$ หรือต่ำกว่า) ขนาดของความสัมพันธ์สูงในการวิจัยที่ผ่านมา นั้นกำหนดไว้ถึง .90 แต่ทั้งนี้ในสภาพความเป็นจริง ขนาดของความสัมพันธ์ที่สูงถึง .90 อาจจะเป็นไปไม่ได้ ผู้วิจัยจึงกำหนดขนาดของความสัมพันธ์สูงไว้เท่ากับ .70 ($r^2 = .49$) ซึ่งเป็นระดับต่ำสุดของความสัมพันธ์มากตามตารางสรุปขนาดความสัมพันธ์ของกริมม์ ซึ่งสอดคล้องกับงานวิจัยของ ถวัลย์ จันทร์เพ็ง (2534) ที่กำหนดความสัมพันธ์มากที่สุดไว้เท่ากับ .60 และใกล้เคียงกับงานวิจัยของ มาร์แคนโทนีโอ (Marcantonio, 1992) ที่ได้กำหนดขนาดความสัมพันธ์มากที่สุดไว้เท่ากับ .55

ดังนั้นในการศึกษาวิจัยครั้งนี้ผู้วิจัยจึงกำหนดขนาดความสัมพันธ์ระหว่างตัวแปรเป็น

3 ระดับ คือ

- 2.1 ความสัมพันธ์สูง ($r=.70$)
- 2.2 ความสัมพันธ์ปานกลาง ($r=.50$)
- 2.3 ความสัมพันธ์ต่ำ ($r=.30$)

3. จำนวนข้อมูลสูญหาย วิธีการจัดการข้อมูลสูญหายจะดีหรือไม่ดีขึ้นอยู่กับตัวแปรอีกตัวหนึ่ง คือ จำนวนข้อมูลสูญหาย ถ้ามีข้อมูลสูญหายจำนวนน้อยประมาณ 5% การตัดหน่วยตัวอย่างออกไปดูเหมือนว่าจะมีเหตุผลในการแก้ปัญหาข้อมูลสูญหาย แต่ถ้ามีการสูญหายมากการตัดข้อมูลออกจะไม่มีประสิทธิภาพข้อมูลที่เหลืออยู่จะไม่ใช่ตัวแทนของประชากรซึ่งมีเป้าหมายในการอ้างอิง (Schafer, 1997. p. 1 ; Little & Rubin, 1987. p. 5) แต่ก็ยังมีงานวิจัยของ เรมอนด์และโรเบิร์ต (Roth, 1994. p. 542 citing Raymond & Robert, 1987) ที่กล่าวว่า ถ้าใช้วิธีการถดถอยจัดการข้อมูลสูญหายจะมีความเหมาะสมเมื่อมีข้อมูลสูญหายมากกว่า 20% ส่วนเฟรดและลี (Fred & Lii, 1998. <http://www.findarticles.com>) ได้ศึกษาประเมินความถูกต้องของวิธีการจัดการข้อมูลสูญหายหลายวิธีภายใต้รูปแบบของการสูญหายที่แตกต่างกันโดยกำหนดจำนวนของข้อมูลสูญหายแบบสุ่มมีค่าต่ำสุดเป็น 10% และสูงสุดเท่ากับ 20% นอกจากนี้มีนักวิจัยคนอื่น ๆ ได้กำหนดจำนวนข้อมูลสูญหายแตกต่างออกไป เช่น รุท (Roth, 1994. p. 551) ได้กล่าวว่า การเลือกวิธีการจัดการข้อมูลสูญหายจะมีความสำคัญเมื่อจำนวนข้อมูลสูญหายอยู่ระหว่าง 15-20% และจะมีความสำคัญมากที่สุดเมื่อจำนวนข้อมูลสูญหายเท่ากับ 30-40% ที่ระดับนี้การเลือกใช้วิธีการจัดการข้อมูลสูญหายจะให้ผลลัพธ์ที่แตกต่างกัน จะเห็นว่าจำนวนข้อมูลสูญหายมีค่าอยู่ระหว่าง 5-40%

เมื่อพิจารณาอัตราการตอบแบบสอบถามกลับมาเป็นข้อมูลสมบูรณ์ที่ผู้วิจัยสามารถวิเคราะห์ได้ก็จะอยู่ระหว่าง 80-90% (เซาเวร์ อินโย, 2541. หน้า 194) จึงถือ่าใช้ได้ ถ้าน้อยกว่านี้ก็ไม่มีคามหมายในการสรุปผลการศึกษาคั้งนั้น สมาคมการศึกษาของอเมริกัน (NEA) พบว่าข้อมูลจากแบบสอบถามที่จะถือได้ว่าเป็นตัวแทนของประชากรได้นักวิจัยต้องได้รับแบบสอบถามกลับคืนมามากกว่าร้อยละ 90 และอัตราการตอบกลับของแบบสอบถามต่ำสุดซึ่งทำให้ตัวประมาณค่าไม่มีความลำเอียง คือ อัตราการตอบกลับตั้งแต่ร้อยละ 95 (วีระยุทธ ชาตะกาญจน์, 2538) จากงานวิจัยและเอกสารดังกล่าวข้อมูลสูญหายจะอยู่ระหว่าง 5% ถึง 20% อยู่ในช่วงที่ผู้วิจัยยอมรับและสามารถนำไปวิเคราะห์ต่อไปได้

ในการวิจัยครั้งนี้ใช้เทคนิคการจำลองสถานการณ์ขึ้นมาศึกษา เพื่อให้ได้ข้อค้นพบที่เป็นประโยชน์มากยิ่งขึ้นในเรื่องของจำนวนข้อมูลสูญหาย ผู้วิจัยจึงกำหนดจำนวนข้อมูลสูญหายให้มีจำนวนมากกว่า 20% โดยกำหนดจำนวนข้อมูลสูญหายสูงสุดเท่ากับ 30% ซึ่งจำนวนข้อมูลสูญหายดังกล่าวสอดคล้องกับงานวิจัยของ รุท (Roth) ที่กล่าวว่า การเลือกวิธีการจัดการข้อมูลสูญหายมีความสำคัญมากที่สุดเมื่อข้อมูลสูญหายมีค่าอยู่ระหว่าง 30-40%

ดังนั้นในการศึกษาวิจัยครั้งนี้ผู้วิจัยจึงกำหนดจำนวนข้อมูลสูญหายออกเป็น

4 ระดับ คือ

3.1 5%

3.2 10%

3.3 20%

3.4 30%

4. วิธีการจัดการข้อมูลสูญหาย วิธีการจัดการข้อมูลสูญหายเป็นตัวแปรที่ผู้วิจัยสนใจนำมาศึกษา การศึกษาวิธีการจัดการข้อมูลสูญหายในเมืองไทยยังมีน้อยมาก ในต่างประเทศมีการศึกษาค้นคว้ามากมาย มีทั้งบทความ ตำรา และโปรแกรมที่ช่วยจัดการข้อมูลสูญหาย จากการศึกษาเอกสารและงานวิจัยที่เกี่ยวข้อง พบว่า วิธีการจัดการข้อมูลสูญหายที่ใช้กันมากที่สุด ได้แก่ การตัดข้อมูลสูญหายออกแบบลิสต์ไวส์ การตัดข้อมูลสูญหายออกแบบแพร์ไวส์ การแทนค่าข้อมูลสูญหายด้วยค่าเฉลี่ย การแทนค่าข้อมูลสูญหายโดยวิธีการถดถอย การแทนค่าข้อมูลสูญหายด้วยวิธีอีเอ็ม (EM algorithm) และเมื่อนำวิธีการจัดการข้อมูลสูญหายต่าง ๆ มาเปรียบเทียบกัน พบว่า เทคนิควิธีการที่ค่อนข้างมีเหตุผลในการประมาณค่าพารามิเตอร์ได้ถูกต้องแม่นยำคือการทำซ้ำโดยเฉพาะวิธีอีเอ็ม ดังเช่นงานวิจัยของ ฟิงไบเนอร์ (Finkbeiner, 1979, pp. 416-420) มาร์แคนโทนีโอ (Marcantonio, 1992) และยังมีงานวิจัยที่ใช้การประมาณค่าด้วยการทำซ้ำลักษณะเช่นเดียวกับวิธีอีเอ็ม ได้แก่ งานวิจัยของ สุพันธ์ วิรุณเทวัญ (Viragoontavan, 2000) ใช้วิธีการแทนค่าข้อมูลสูญหายด้วยวิธีเอ็มไอโดยใช้โปรแกรม SOLAR และโปรแกรม NORM พบว่า วิธีการจัดการข้อมูลสูญหายแบบเอ็มไอมีประสิทธิภาพมากที่สุด

วิธีการที่ควรนำมาใช้ในการจัดการข้อมูลสูญหายคือวิธีอีเอ็ม (EM algorithm) เพราะประมาณค่าพารามิเตอร์ได้ถูกต้องแม่นยำและมีโปรแกรมสำเร็จรูปช่วยในการวิเคราะห์ แต่วิธีการดังกล่าวมีจุดบกพร่องดังต่อไปนี้ การแทนค่าข้อมูลสูญหายครั้งแรกในสมการ $\hat{Y} = a + bX$ ใช้ค่าสถิติได้มาจากกลุ่มตัวอย่างที่มีข้อมูลสมบูรณ์ตัดหน่วยตัวอย่างที่มีข้อมูลสูญหายออกไป การประมาณค่าพารามิเตอร์ด้วยค่าเหล่านี้จึงมีอคติ (bias) หมายความว่าค่าที่ได้อาจจะไม่ใช่ค่าที่แท้จริงของประชากร ประกอบกับการประมาณค่าแบบนี้ใช้วิธี Maximun likelihood ซึ่งใช้เทคนิคการทำซ้ำ แนวคิดการวิเคราะห์ใช้สถิติขั้นสูง และไม่สามารถนำไปรวมกับโปรแกรมคอมพิวเตอร์ทั่ว ๆ ไปได้

ดังนั้นผู้วิจัยจึงเสนอวิธีการประมาณค่าข้อมูลสูญหาย ขั้นตอนแรกประมาณค่าข้อมูลสูญหายด้วยวิธีการถดถอยอย่างง่าย โดยการประมาณค่าพารามิเตอร์ด้วยวิธีการทำซ้ำ นำค่าพารามิเตอร์ที่ได้ไปแทนในสมการ $\hat{Y} = a + bX$ แล้วจึงคัดเลือกสมการทำนายที่มีความคลาดเคลื่อนน้อยที่สุด จากการทำซ้ำ 1,000 ครั้ง เพื่อให้สมการทำนายข้อมูลสูญหายได้อย่างถูกต้องแม่นยำ เรียกว่าวิธีนี้ว่า วิธีการจัดการข้อมูลสูญหายแบบอีพีเอสเอสอี

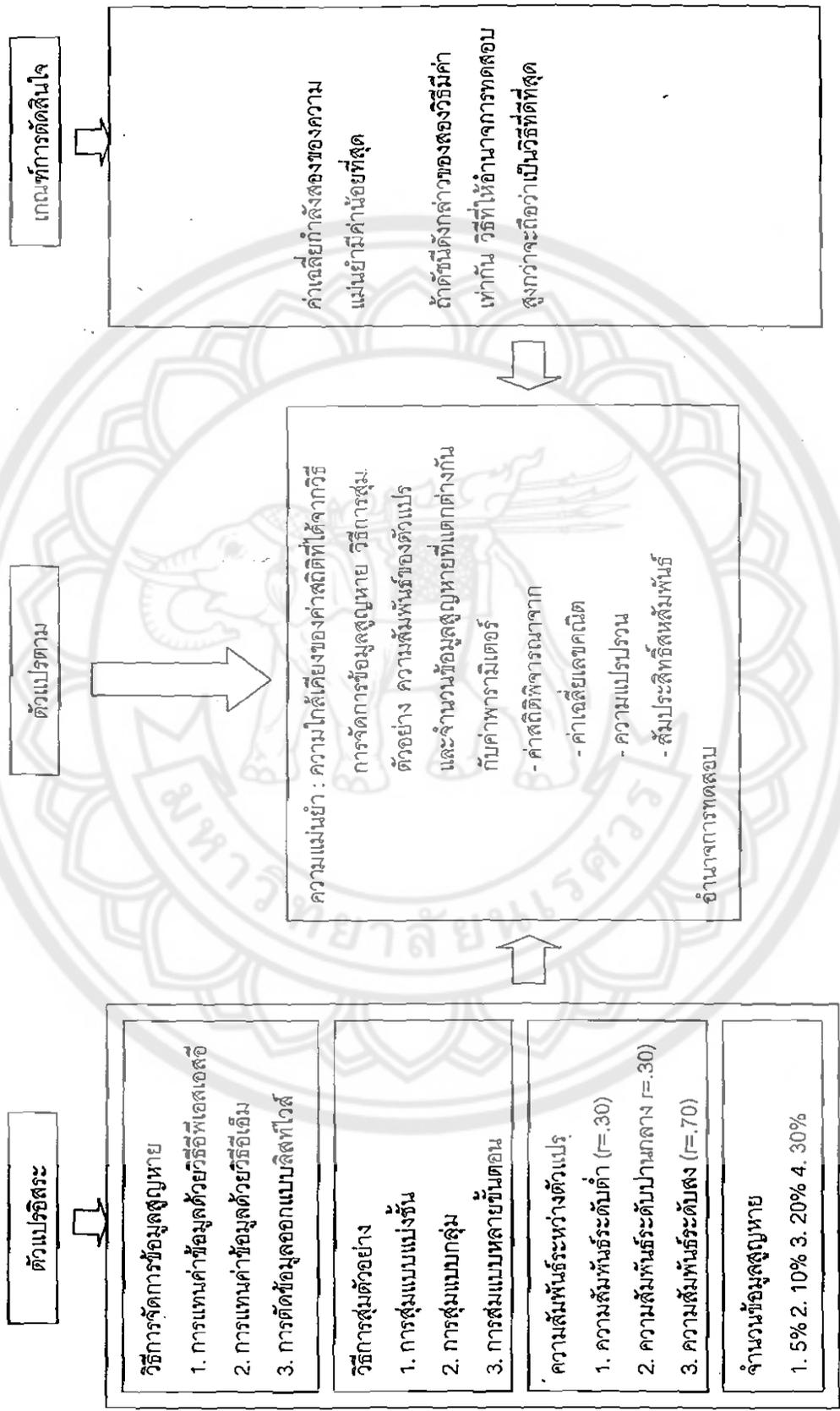
การวิจัยครั้งนี้ผู้วิจัยมีวัตถุประสงค์เพื่อเปรียบเทียบว่าวิธีการจัดการข้อมูลสูญหายแบบใดดีกว่ากันโดยพิจารณาจากความแม่นยำและอำนาจการทดสอบ ความแม่นยำพิจารณาจากความใกล้เคียงกันของค่าสถิติกับค่าพารามิเตอร์โดยมีเกณฑ์ตัดสินความแม่นยำดังนี้ คือ ค่าเฉลี่ยกำลังสองของความคลาดเคลื่อนมีค่าน้อยที่สุด ส่วนอำนาจการทดสอบ หมายถึง ความน่าจะเป็นในการปฏิเสธสมมติฐานสูญเมื่อสมมติฐานนั้นไม่เป็นความจริง ดังนั้นจึงนำวิธีการจัดการข้อมูลสูญหายโดยการตัดข้อมูลสูญหายออกแบบลิสต์ไวส์มาเป็นวิธีพื้นฐานในการเปรียบเทียบ สอดคล้องกับการศึกษาของ มาร์แคนโทนีโอ (Marcantonio, 1992) ที่ใช้วิธีการตัดข้อมูลสูญหายออกแบบลิสต์ไวส์เป็นตัวเปรียบเทียบ เพราะวิธีการตัดข้อมูลสูญหายออกแบบลิสต์ไวส์จะตัดหน่วยตัวอย่างที่มีรายการคำตอบใดคำตอบหนึ่งที่เป็นข้อมูลสูญหายออกไปจากการวิเคราะห์ จึงเป็นวิธีที่ทำแล้วข้อมูลที่เหลืออยู่เป็นข้อมูลที่สมบูรณ์ การนำไปเปรียบเทียบกับวิธีการแทนค่าแบบอื่น ๆ ทำให้เห็นความแตกต่างได้อย่างชัดเจนขึ้นว่าการตัดออกไปหรือการแทนค่าวิธีใดดีกว่ากัน

ในการวิจัยครั้งนี้ผู้วิจัยจึงกำหนดวิธีการจัดการข้อมูลสูญหาย 3 วิธี คือ

- 4.1 การตัดข้อมูลสูญหายออกแบบลิสต์ไวส์
- 4.2 การแทนค่าข้อมูลสูญหายด้วยวิธีอีเอ็ม
- 4.3 การแทนค่าข้อมูลสูญหายด้วยวิธีอีพีเอสเอสอี

รายละเอียดตัวแปรที่ใช้ในการวิจัยสรุปได้ดังแผนภาพต่อไปนี้

จากการศึกษาเอกสารและงานวิจัยที่เกี่ยวข้อง ผู้วิจัยได้กำหนดกรอบแนวคิดที่ใช้ในการวิจัยดังนี้



ภาพ 6 กรอบแนวคิดที่ใช้ในการวิจัย