

บทที่ 2

เอกสารและงานวิจัยที่เกี่ยวข้อง

ในบทที่ 2 นี้ แบ่งเป็นสองส่วนคือ ส่วนที่หนึ่งกล่าวถึงความรู้พื้นฐานที่เกี่ยวข้องกับข้อมูลสัญญาณได้แก่ ขนาดตัวอย่าง รูปแบบของข้อมูลสัญญาณที่ถูกจำลองขึ้นในการทดลอง เทคนิคที่ใช้ในการประมาณค่าข้อมูลสัญญาณ การแทนค่าข้อมูลสัญญาณด้วยค่าประมาณจากการคำนวณ การประมาณค่าพารามิเตอร์ การจำลองสถานการณ์ ส่วนที่สองเป็นการกล่าวถึงงานวิจัยที่เกี่ยวข้องต่างๆ ที่ได้พัฒนามาเพื่อแก้ปัญหาในเรื่องของการประมาณค่าข้อมูลสัญญาณ

2.1 ความรู้พื้นฐาน

2.1.1. ประชากร (Population)

ประชากร หมายถึง กลุ่มของหน่วยตัวอย่างที่ผู้วิจัยสนใจศึกษา ประชากรอาจเป็นบุคคล องค์กร สัตว์ สิ่งของที่นำมาเป็นหน่วยศึกษาหรือปัญหาในการวิจัย ถ้าประชากรนั้นประกอบด้วยหน่วยตัวอย่างเป็นจำนวนหน่วยที่นับได้แน่นอน หรือสามารถนับจำนวนที่แน่นอนได้ เรียกว่าประชากรจำกัด (finite population) แต่ถ้าประชากรประกอบด้วยหน่วยตัวอย่างจำนวนหน่วยที่นับไม่ได้แน่นอนหรือไม่สามารถนับจำนวนที่แน่นอนได้เรียกว่าประชากรอนันต์ (infinite population) สำหรับงานวิจัยนี้ศึกษาเฉพาะในส่วนของประชากรจำกัดเท่านั้น

2.1.2 ตัวอย่าง (Sample)

ตัวอย่าง หมายถึง ส่วนหนึ่งของประชากรที่ผู้วิจัยสนใจศึกษา ตัวอย่างที่ดีหมายถึง ตัวอย่างที่มีลักษณะต่างๆ ที่สำคัญครบถ้วนเหมือนกับประชากร

ขนาดตัวอย่าง (Sample size)

ขนาดตัวอย่าง คือ จำนวนหน่วยตัวอย่างในการวิจัย เป็นสิ่งสำคัญที่ผู้วิจัยต้องกำหนดให้เหมาะสมและมีความเป็นตัวแทนที่ดีของประชากรที่ทำการศึกษาเพื่อจะช่วยให้ผลการวิจัยมีความน่าเชื่อถือ

การกำหนดขนาดตัวอย่าง

การกำหนดขนาดตัวอย่างควรมีขนาดเท่าใดนั้นผู้วิจัยควรคำนึงถึงสิ่งต่อไป หลายอย่างมาประกอบกัน (Librero, 1985 อ้างถึงใน ธีรรุษิ เอกะกุล, 2543) ดังนี้

(1) ค่าใช้จ่าย เวลา แรงงานและเครื่องมือที่ใช้ในการเก็บรวบรวมข้อมูลตัวอย่างนั้นว่า มีเพียงพอหรือไม่ และคุ้มค่าเพียงใด

(2) ขนาดประชากร ถ้าประชากรมีขนาดใหญ่ มีความจำเป็นต้องเลือกตัวอย่าง ถ้าประชากรมีขนาดเล็กและสามารถที่จะศึกษาได้ควรจะศึกษาจากประชากรทั้งหมด

(3) ความเหมือนกัน ถ้าประชากรมีความเหมือนกันมากความแตกต่างของหน่วยตัวอย่างมีน้อยนั่นคือความแปรปรวนของตัวอย่างมีค่าน้อยใช้ตัวอย่างขนาดเล็กได้ ถ้าประชากรมีลักษณะไม่เหมือนกันความแตกต่างของหน่วยตัวอย่างมีค่ามากทำให้ความแปรปรวนในกลุ่มมากจำเป็นต้องใช้ตัวอย่างขนาดใหญ่เพื่อให้ครอบคลุมคุณลักษณะต่าง ๆ ของประชากร

(4) ความแม่นยำ ถ้าต้องการความแม่นยำในเรื่องที่จะศึกษาค้นคว้า ต้องใช้ตัวอย่างขนาดใหญ่ คือ ยิ่งขนาดตัวอย่างใหญ่มากเท่าไหร่ผลการศึกษายิ่งมีความแม่นยามากขึ้นเท่านั้น

(5) ความคลาดเคลื่อนจากการซักตัวอย่าง ความคลาดเคลื่อนที่ยอมให้เกิดขึ้นได้จากการซักตัวอย่างโดยทั่วไปแล้วมักยอมให้เกิดความคลาดเคลื่อนได้ 1% หรือ 5% (สัดส่วน 0.01 หรือ 0.05) และยังขึ้นอยู่กับความสำคัญของเรื่องที่ต้องการศึกษาด้วยถ้าปัญหามีความสำคัญมากก็ควรให้เกิดความคลาดเคลื่อนน้อยที่สุด เช่น 1% แต่ถ้ามีความสำคัญน้อยก็อาจยอมให้เกิดความคลาดเคลื่อนได้บ้าง เช่น 5% เป็นต้น

(6) ความเชื่อมั่น ผู้วิจัยต้องกำหนดความเชื่อมั่นว่าตัวอย่างที่สุ่มมานั้นมีโอกาสได้ค่าอ้างอิงไม่แตกต่างจากค่าที่แท้จริงของประชากรประมาณเท่าไร เช่น ถ้ากำหนดระดับเชื่อมั่น 95% หมายถึง ค่าอ้างอิงมีโอกาสถูกต้อง 95% มีโอกาสผิดพลาดจากค่าที่แท้จริง 5% นั่นคือ ค่าที่ได้จากกลุ่มตัวอย่าง 95 กลุ่ม จาก 100 กลุ่มที่สุ่มมาจากประชากรเดียวกันจะไม่แตกต่างจากค่าที่แท้จริงของประชากรซึ่งระดับความเชื่อมั่นอาจจะเพิ่มขึ้นเป็น 99% หรือลดลงเหลือ 90%

2.1.3 การประมาณค่าที่เก็บไม่ได้ (Imputation)

การประมาณค่าที่เก็บไม่ได้ คือการประมาณค่าข้อมูลสูญหาย และใช้ค่าที่ประมาณนั้นมาวิเคราะห์เชิงสถิติ Laaksonen (2000) ได้จำแนกวิธีการในการประมาณค่าเป็น 2 กลุ่ม คือ

model-donor imputation คือ การประมาณค่าที่ได้มาจากการตัวแบบ เช่น ตัวประมาณค่าเมื่อประมาณค่าข้อมูลสูญหายด้วยวิธีค่าเฉลี่ย ตัวประมาณข้อมูลสูญหายด้วยค่าทดแทนและตัวประมาณข้อมูลสูญหายด้วยค่าอัตราส่วน

real-donor imputation คือ การประมาณค่าที่ได้มาจากการเช็คของค่าสังเกต เช่น ตัวประมาณข้อมูลสูญหายด้วยค่าโคลเดค ตัวประมาณข้อมูลสูญหายด้วยค่าอห์ดเดค และตัวประมาณข้อมูลสูญหายด้วยค่าไกส์สูด

2.1.4 ข้อมูลสูญหาย (Missing data)

ข้อมูลสูญหาย คือ ค่าสังเกตที่ต้องการทราบค่าแต่ไม่สามารถทราบค่าได้หรือไม่สามารถเก็บรวบรวมได้ ข้อมูลสูญหายไม่เกิดขึ้นหากวิธีการที่ใช้ในการรวบรวมข้อมูลหรือในการวัดค่ามีประสิทธิภาพดีหรือมีความเหมาะสม ข้อมูลสูญหายส่งผลต่อการวิจัยทั้งในส่วนของการวิเคราะห์และการสรุปผลต่อความ โดยระดับความรุนแรงของผลกระทบนี้ขึ้นอยู่กับองค์ประกอบหลายส่วน แต่ที่สำคัญคือขนาดของข้อมูลสูญหาย ประเภทของข้อมูลสูญหายที่เกิดขึ้น และวิธีการจัดการกับข้อมูลสูญหาย (ปีบัตรณ์ ประสิทธิ์วัฒนารี และสุคนธ์ ประสิทธิ์วัฒนารี, 2549)

2.1.4.1 ขนาดของข้อมูลสูญหาย (Size of missing data)

กรณีที่ไม่มีค่าตอบหรือเกิดการปล่อยว่างไว้สำหรับบางรายการ นักวิจัยจำเป็นต้องตัดสินใจดำเนินการระหว่างการพิจารณาไม่เลือกใช้ข้อมูลในรายที่เกิดปัญหานี้หรือควรแทนค่าสูญหาย (imputed value) ปัจจุบันยังไม่มีเกณฑ์ที่ย่างง่ายสำหรับใช้ในการตัดสินใจ เลือกใช้วิธีการใดวิธีการหนึ่งแต่มีแนวทางหนึ่งที่ถูกนำมาใช้กันอย่างแพร่หลายคือ การพิจารณาจากขนาดของข้อมูลสูญหาย เมื่อจำนวนของหน่วยตัวอย่างที่มีค่าของข้อมูลสูญหายมีจำนวนน้อย (อาจใช้เกณฑ์ $< 5\%$ จากหน่วยตัวอย่างที่มีขนาดใหญ่) อาจเลือกใช้วิธีการย่างง่ายคือใช้วิธีการไม่เลือกข้อมูลที่มีปัญหามาใช้ในการวิเคราะห์ ดังนั้นในการศึกษาจึงควรพยายามลดขนาดของข้อมูลสูญหายให้เหลือจำนวนน้อยที่สุดหรือไม่เหลือเลย การหาทางป้องกันไม่ให้เกิดข้อมูลสูญหาย แนวทางที่จะลดขนาดของข้อมูลสูญหายจำเป็นต้องพิจารณาถึงสาเหตุต่าง ๆ ที่เป็นไปได้ในการเกิดข้อมูลสูญหาย

2.1.4.2 ประเภทของข้อมูลสูญหาย (Missing type of data)

ประสิทธิภาพของขั้นตอนวิธีหรือเทคนิคที่นำมาใช้ในการประมาณค่าข้อมูลสูญหาย นั้นจะดีหรือไม่ส่วนหนึ่งขึ้นอยู่กับประเภทการสูญหายของข้อมูล ถ้าสามารถทราบสาเหตุที่ทำให้เกิดข้อมูลสูญหายอาจจะสามารถเติมเต็มหรือเดาข้อมูลส่วนนั้นได้ไม่ยากนัก ซึ่งในการทำงานจริงไม่อาจจะทราบได้ว่าข้อมูลสูญหายนั้นมีสาเหตุมาจากการและสูญหายในลักษณะใด ดังนั้นในการทดลองต่าง ๆ จึงใช้วิธีจำลองรูปแบบของข้อมูลสูญหายในหลายลักษณะโดยจำแนกออกได้เป็น 3 ประเภท ดังต่อไปนี้

(1) การสูญหายแบบสุ่มอย่างสมบูรณ์ (Missing completely at random หรือ MCAR)

เกิดขึ้นเมื่อความน่าจะเป็นของการสูญหายของข้อมูลที่ตำแหน่งใด ๆ นั้น ไม่มีความสัมพันธ์กับค่าของข้อมูลตัวอื่น ๆ ไม่ว่าจะเป็นข้อมูลที่ทราบค่าหรือข้อมูลที่เกิดการสูญหายด้วยกันก็ตาม หมายถึงความน่าจะเป็นที่จะเกิดการสูญหายของค่าข้อมูลในทุก ๆ ตำแหน่งมีค่าเท่ากัน ในการสุ่มรูปแบบนี้พบว่าวิธีการที่ใช้ในการประมาณค่าข้อมูลนั้นสามารถนำมาประยุกต์ใช้ได้ในทุก ๆ วิธีโดยที่จะไม่ทำให้ข้อมูลเกิดการเบี่ยงเบน

สำหรับสาเหตุที่ทำให้ข้อมูลเกิดการสูญหายแบบสุ่มอย่างสมบูรณ์ อาจเกิดขึ้น เนื่องจากเครื่องมือเสีย อุปกรณ์เกิดข้อบกพร่อง สภาพอากาศแลวร้าย กลุ่มเป้าหมายที่ศึกษาล้มป่วย หรือการนำเข้าข้อมูลไม่ถูกต้อง

สำหรับข้อมูลสูญหายประเภทนี้ จัดเป็นข้อมูลสูญหายที่ก่อให้เกิดปัญหาน้อยที่สุด เพราะว่าข้อมูลสูญหายไม่มีความเกี่ยวข้องต่อผลลัพธ์ของข้อมูล เพราะฉะนั้นจึงสามารถเลือกวิเคราะห์ ข้อมูลในส่วนที่สมบูรณ์ได้

(2) การสูญหายแบบสุ่ม (Missing at random หรือ MAR) การสูญหายของค่าข้อมูล ในกรณีนี้ขึ้นอยู่กับค่าของข้อมูลตัวอื่น ๆ ที่ทราบค่าข้อมูลด้วย แต่ไม่ขึ้นอยู่กับค่าข้อมูลที่สูญหาย ของตัวเองยกตัวอย่าง เช่น หากพบว่าเฉพาะกลุ่มผู้ได้รับการศึกษาน้อยที่ไม่ให้ความร่วมมือในการตอบคำถามเกี่ยวกับทัศนคติในการแพทย์สเปดิค ในลักษณะนี้สามารถกล่าวได้ว่าข้อมูลทัศนคติ ในการแพทย์สเปดิค มีค่าการสูญหายแบบสุ่มทั้งนี้เนื่องจากเป็นค่าสูญหายที่เกิดขึ้นเฉพาะบางส่วน ของตัวแปรระดับการศึกษา

(3) การสูญหายแบบไม่สุ่ม (Not missing at random หรือ NMAR) การสูญหายของ ค่าข้อมูลในกรณีนี้เกิดขึ้นเมื่อความน่าจะเป็นของแ雷ทีนีการสูญหายที่ดำเนินการ ไม่ขึ้นอยู่กับค่า ของข้อมูลตำแหน่งนั้น ๆ หรือในบางกรณีค่าของข้อมูลสูญหายอาจไม่ขึ้นอยู่กับค่าตัวแปรใด ๆ ใน ฐานข้อมูลเลยแต่ขึ้นอยู่กับตัวแปรอื่นที่ไม่ได้ถูกเก็บรวบรวมไว้ในการศึกษารังนั้น เช่น ค่าน้ำหนัก ตัวที่ลดลงขึ้นอยู่กับน้ำหนักตัวตอนเริ่มต้นแต่เนื่องจากตัวแปรน้ำหนักตัวตอนเริ่มต้นไม่ได้ถูก เก็บรวบรวมไว้ในฐานข้อมูล ดังนั้นค่าสูญหายของน้ำหนักตัวที่ลดลงจึงขึ้นอยู่กับตัวแปรภายนอก ฐานข้อมูล ลักษณะข้อมูลสูญหายประเภทนี้จัดเป็นข้อมูลสูญหายที่สามารถส่งผลกระทบอย่าง รุนแรงในการวิเคราะห์ข้อมูล

ในทางปฏิบัติลักษณะของข้อมูลสูญหายประเภทการสูญหายแบบสุ่มอย่าง สมบูรณ์มักไม่ค่อยพบบ่อยนักส่วนใหญ่เป็นข้อมูลสูญหายประเภทการสูญหายแบบสุ่ม ดังนั้น วิธีการเชิงสถิติต่าง ๆ ที่พัฒนาขึ้นมาเพื่อแก้ปัญหาข้อมูลสูญหายมักคำนึงถึงการภายใต้ข้อสมมุติของ การสูญหายแบบสุ่ม

2.1.4.3 วิธีแก้ไขข้อมูลสูญหาย (Methods for treating missing data) จากการศึกษา พบว่าวิธีแก้ไขข้อมูลสูญหายมีหลายวิธี การเลือกใช้วิธีการใดขึ้นอยู่กับลักษณะของข้อมูลสูญหายที่ เกิดขึ้นหากเลือกวิธีการที่ไม่เหมาะสมสมนาใช้อาจทำให้ผลลัพธ์ผิดพลาดได้ วิธีการจัดการกับข้อมูล สูญหายที่มีการนำมาใช้มีดังนี้

(1) การตัดกลุ่มข้อมูลที่สูญหายทิ้งไป (Ignoring and discarding data) เป็น วิธีการที่แบ่งเป็น 2 กลุ่มหลัก ๆ คือ การตัดแ雷ทีสูญหายทิ้งไปทั้งหมด (Complete case analysis) ซึ่ง

นิยมใช้ในงานด้านสถิติ อีกวิธี คือ การตัดเฉพาะเดาที่มีการสูญหายของค่าข้อมูลทิ้ง ซึ่งก่อนจะตัด ตำแหน่งใด ๆ ออกไป ต้องวิเคราะห์ความสัมพันธ์ระหว่างข้อมูลก่อนเสมอ(Afifi & Clark,1996 ถ้าลงใน ลงลักษณ์ วิรัชชัย, 2542)

(2) การตัดข้อมูลแบบลิสท์ไวส์ (Listwise data deletion) เป็นวิธีการจัดการ ข้อมูลสูญหายที่ง่ายมากนั่นคือ ไม่สนใจข้อมูลสูญหายที่เกิดขึ้น โดยวิเคราะห์ข้อมูลจากข้อมูลเฉพาะ ส่วนที่สมบูรณ์ แนวทางนี้มีความเหมาะสมในกรณีที่ข้อมูลสูญหายมีจำนวนน้อยมากหรือผลจากการ วิเคราะห์ข้อมูลนี้ความชัดเจนมากซึ่งวิธีการนี้ก菽กกำหนดให้ใช้เป็นตัวเลือกอัตโนมัติ (by default) สำหรับจัดการกับข้อมูลสูญหายในโปรแกรมคอมพิวเตอร์ทางสถิติทั่วๆ ไปหากไม่เจาะจงเลือกใช้ วิธีการอื่นใดในการจัดการกับข้อมูลสูญหาย

ข้อเสียของการตัดข้อมูลแบบลิสท์ไวส์

(2.1) ไม่เหมาะสมที่จะนำไปใช้กับข้อมูลจริงเนื่องจากเดาของข้อมูลที่ไม่สมบูรณ์ ถูกตัดทิ้งไปทำให้มีจำนวนของข้อมูลที่นำไปวิเคราะห์น้อยลงจนบางครั้งไม่เพียงพอต่อการ นำไปวิเคราะห์

(2.2) ในกรณีที่ต้องการคำนวนหาค่าสถิตินางค่าของข้อมูลนั้น ข้อมูลต้องมีการ สูญหายแบบสุ่มอย่างสมบูรณ์ถ้าข้อมูลมีการสูญหายแบบสุ่มพบว่าค่ากลางที่คำนวนมาได้อาจเกิด การเบี่ยงเบนจากความเป็นจริงได้

(3) การตัดข้อมูลแบบแพร์ไวส์ (Pairwise data deletion) เป็นวิธีการจัดการกับ ข้อมูลสูญหายสำหรับในกรณีที่วิเคราะห์ความสัมพันธ์ระหว่างตัวแปรคู่ โดยวิเคราะห์ข้อมูลเฉพาะ ข้อมูลที่สมบูรณ์ของทั้งสองตัวแปร วิธีการนี้ขอดีในส่วนของการใช้งานข้อมูลได้เต็มที่และมี ประสิทธิภาพ แต่มีข้อด้อยคือต้องซับซ้อนกว่าวิธีการตัดข้อมูลแบบลิสท์ไวส์เล็กน้อยและเสียเวลามากกว่า จึงได้รับความนิยมน้อยกว่า

(4) วิธีการประมาณข้อมูลสูญหายด้วยค่าเฉลี่ยเป็นวิธีการแทนค่าข้อมูล สูญหายด้วยค่าเฉลี่ยของข้อมูลที่ทราบค่าในแต่ละกลุ่มย่อยของตัวแปรอื่น ซึ่งเป็นวิธีที่พัฒนามากจาก การแทนค่าข้อมูลสูญหายด้วยค่าเฉลี่ยของข้อมูลที่ทราบค่าทั้งนี้เนื่องจากข้อมูลนุ่ติว่าค่าของข้อมูล สูญหายควรขึ้นอยู่กับลักษณะของหน่วยตัวอย่าง โดยลักษณะของหน่วยตัวอย่างที่ใกล้เคียงกันควรมี ค่าข้อมูลที่คล้ายคลึงกันวิธีนี้เสนอโดย Wilks ในปี ค.ศ. 1932 โดยกำหนดให้

$$y_{\cdot i} = \begin{cases} y_i & \text{if } i \in R \\ \bar{y}_M & \text{if } i \in R^c \end{cases} \quad (2.1)$$

ตัวประมาณค่าเฉลี่ยประชากรที่ใช้ข้อมูลที่ประมาณค่าข้อมูลสูญหายด้วยสมการ (2.1)

อยู่ในรูป

$$\bar{y}_s = \frac{1}{n} \sum_{i \in s} y_i = \bar{y}_{MI} \quad (2.2)$$

เมื่อ $\bar{y}_{MI} = \frac{1}{r} \sum_{i \in R} y_i$ โดยที่

y_i คือ ค่าของตัวแปรที่สนใจ y ของหน่วยสังเกตที่ i ซึ่งมีทั้งหมด r หน่วย
 y_r คือ ค่าของตัวแปรที่สนใจ y ของข้อมูลที่สมบูรณ์ r หน่วย และค่าของตัวแปรที่สนใจ y ของ
 หน่วยสังเกตที่ i หลังทำการประมาณค่าข้อมูลสูญหาย $k - r$ หน่วย

\bar{y}_{MI} คือ ค่าเฉลี่ยของตัวแปรที่สนใจ y จากข้อมูลที่สมบูรณ์ทั้งหมด r หน่วย
 เป็นตัวประมาณค่าที่ไม่เออนเอียง (Unbiased) และมีค่าความแปรปรวน

$$V(\bar{y}_{MI}) = \left(\frac{1}{r} - \frac{1}{N} \right) S_y^2 \quad (2.3)$$

(5) วิธีการประมาณข้อมูลสูญหายด้วยค่าทดแทน เป็นวิธีการประมาณค่าข้อมูล
 สูญหายโดยการสร้างสมการทดแทนระหว่างตัวแปรที่สนใจ y กับตัวแปรช่วง x โดยใช้ข้อมูลที่สมบูรณ์
 ทั้งหมด r หน่วยแล้วประมาณค่าข้อมูลสูญหายด้วยสมการทดแทน สุนันทา วีรกุลเทวัญ (2544) กล่าว
 ว่า ถึงแม้ว่าวิธีการ แทนค่าด้วยวิธีการประมาณข้อมูลสูญหายด้วยค่าทดแทน ทำให้ความแปรปรวน
 ของข้อมูลมีค่าน้อยกว่าความแปรปรวนที่ได้จากการแทนค่าข้อมูลสูญหายด้วยค่าเฉลี่ยแต่มีข้อเสีย
 คือทำให้ค่าที่เกี่ยวข้องกับความสัมพันธ์ (Association) ระหว่างตัวแปรบิดเบี้ยน (Distort) เพราะ
 สมการทดแทนที่ใช้ประมาณค่าข้อมูลสูญหายสร้างมาจากความสัมพันธ์ระหว่างตัวแปร

(6) วิธีการประมาณข้อมูลสูญหายด้วยค่าใกล้สุด เป็นวิธีการแทนข้อมูลสูญหายโดย
 ใช้ข้อมูลจากการเลือกหน่วยตัวอย่างที่มีลักษณะคล้ายคลึงกันมากที่สุดกับหน่วยตัวอย่างที่เกิดค่า
 สูญหาย จากนั้นแทนค่าข้อมูลสูญหายด้วยค่าของหน่วยตัวอย่างที่คล้ายคลึงกัน

(7) วิธีการประมาณข้อมูลสูญหายด้วยค่าอothเดค เป็นวิธีการแทนข้อมูลสูญหาย
 โดยใช้ข้อมูลจากผู้ให้ข้อมูล (Donor) ในงานวิจัยหรืองานสำรวจเรื่องเดียวกันที่เพิ่งสำรวจเสร็จสิ้น
 ใหม่ ๆ (เรียกว่า hot) ไปเป็นข้อมูลทดแทนให้แก่รายการที่สูญหาย (Roth & Switzer, 1995)

(8) วิธีการประมาณข้อมูลสูญหายด้วยค่าโคลเดค เป็นวิธีการแทนข้อมูลสูญหาย
 โดยใช้ข้อมูลจากผู้ให้ข้อมูล ในงานวิจัยหรืองานสำรวจจากแฟ้มอื่นแต่เป็นงานสำรวจเรื่องเดียวกัน
 ในการสำรวจครั้งก่อน (เรียกว่า cold) ไปเป็นข้อมูลทดแทนให้แก่รายการที่สูญหาย

(9) วิธีการประมาณข้อมูลสัญญาด้วยค่าคอมพ�ไนซ์ เป็นวิธีการแทนข้อมูลสัญญาด้วยค่าประมาณที่นำเสนอโดย Singh and Horn (2000) โดยกำหนดให้

$$y_{r,i} = \begin{cases} \alpha \frac{n}{r} y_i + (1-\alpha)\hat{b}x_i, & \text{if } i \in R \\ (1-\alpha)\hat{b}x_i, & \text{if } i \in R^c \end{cases} \quad (2.4)$$

เมื่อ $\hat{b} = \frac{\sum_{i=1}^r y_i}{\sum_{i=1}^r x_i}$ และ α คือ ค่าที่เหมาะสมที่สูงเลือก

ตัวประมาณค่าเฉลี่ยประชากรที่ใช้ข้อมูลที่ปรับมาณค่าข้อมูลสัญญาด้วยสมการ (2.4) อุปในรูป

$$\bar{y}_{COMP} = \alpha \bar{y}_r + (1-\alpha) \bar{y}_r \frac{\bar{x}_n}{\bar{x}_r} \quad (2.5)$$

เมื่อ

\bar{x}_r คือ ค่าเฉลี่ยของตัวแปรช่วง x ที่มีค่าข้อมูลของตัวแปรที่สนใจ y สมบูรณ์ทั้งหมด r หน่วย

\bar{x}_n คือ ค่าเฉลี่ยของตัวแปรช่วง x จากข้อมูลที่สมบูรณ์ทั้งหมด n หน่วย

\bar{y}_r คือ ค่าเฉลี่ยของตัวแปรที่สนใจ y จากข้อมูลที่สมบูรณ์ทั้งหมด r หน่วย
เป็นตัวประมาณค่าที่เออนเอียง โดยมีความเออนเอียง (Biasedness) ดังนี้

$$b(\bar{y}_{COMP}) = (1-\alpha) \left(\frac{1}{r} - \frac{1}{n} \right) \bar{Y} (C_x^2 - \rho_{yx} C_y C_x) \quad (2.6)$$

และมีค่าคลาดเคลื่อนกำลังสองเฉลี่ย

$$MSE(\bar{y}_{COMP}) = \left(\frac{1}{r} - \frac{1}{N} \right) \bar{Y}^2 C_r^2 + \left(\frac{1}{r} - \frac{1}{n} \right) \bar{Y}^2 \left[(1-\alpha)^2 C_x^2 - 2(1-\alpha)\rho_{yx} C_y C_x \right] \quad (2.7)$$

และมีค่าคลาดเคลื่อนกำลังสองเฉลี่ยที่เหมาะสม

$$MSE(\bar{y}_{COMP})_{opt} = M(\bar{y}_{RAT}) - \left(\frac{1}{r} - \frac{1}{n} \right) \left(1 - \rho_{yx} \frac{C_y}{C_x} \right)^2 \bar{Y}^2 C_x^2 \quad (2.8)$$

2.1.5 การประมาณค่าพารามิเตอร์

การประมาณค่าพารามิเตอร์ มี 2 ชนิด คือ

การประมาณค่าแบบค่าเดียว (Point estimation) คือ การประมาณค่าพารามิเตอร์ของประชากรโดยใช้ค่าสถิติเพียงค่าเดียว คำนวณได้จากตัวอย่าง โดยทั่วไป การประมาณค่าพารามิเตอร์ต่าง ๆ จะทำได้โดยการเลือกตัวอย่างมาจำนวน n หน่วยเก็บรวบรวมข้อมูลจากหน่วยตัวอย่างที่เลือกมาได้จากนั้นคำนวณหาค่าของตัวสถิติแล้วนำไปใช้สำหรับการประมาณค่าของประชากร

การประมาณค่าแบบเป็นช่วง (Interval estimation) คือ การประมาณค่าของพารามิเตอร์ของประชากรออกมานเป็นช่วงซึ่งผลที่ได้มาจากการประมาณแบบค่าเดียว ช่วงของการประมาณค่า พารามิเตอร์ดังกล่าวเรียกว่าช่วงความเชื่อมั่น (confidence interval) สำหรับงานวิจัยนี้ ศึกษาเฉพาะการประมาณค่าแบบค่าเดียวเท่านั้น

2.1.6 การจำลองสถานการณ์ (Simulation) เป็นการรวบรวมวิธีการต่าง ๆ ที่ใช้ในการจำลองสถานการณ์จริงหรือพฤติกรรมของระบบต่าง ๆ มาไว้บนคอมพิวเตอร์ เพื่อศึกษาการให้ผลของกิจกรรมในรูปแบบต่าง ๆ โดยมีการเก็บข้อมูลและวิเคราะห์หารูปแบบที่ถูกต้องจากโปรแกรมคอมพิวเตอร์เพื่อปรับปรุงในอนาคต (Kelton, Sadowski, and Sturrock, 2003)

เนื่องจากการปฏิบัติงานจริง ไม่สามารถทำการทดลองหรือปรับเปลี่ยนกระบวนการทำงานได้จนกว่าจะมองเห็นถึงประโยชน์ที่ได้รับ ดังนั้นการจำลองสถานการณ์จะช่วยให้สามารถวิเคราะห์สภาพที่เป็นอยู่ในปัจจุบันของระบบ ช่วยหาแนวทางหรือทางเลือก (alternative) ที่เหมาะสมก่อนนำไปใช้กับสถานการณ์หรือการปฏิบัติงานจริง ซึ่งจะช่วยให้ลดความเสี่ยงในการเกิดความผิดพลาดหรือความล้มเหลว ได้ นอกจากนี้ยังช่วยให้ประยุกต์ใช้จ่ายและเวลาอีกด้วย (Maria, 1997)

คุณลักษณะสำคัญของแบบจำลองสถานการณ์

แบบจำลองสถานการณ์มีคุณลักษณะที่แตกต่างไปจากแบบจำลองชนิดอื่น ๆ ทั่วไปดังนี้

(1) แบบจำลองสถานการณ์ทุกชนิดต้องมีการตรวจสอบความถูกต้องก่อนเป็นอันดับแรก เพื่อไม่ให้เกิดข้อผิดพลาด มีการตรวจสอบทั้งตรรกะ (Logic) และการคำนวณว่าถูกต้องหรือไม่

(2) มีการตรวจสอบว่าผลที่ได้ต้องอยู่ในขอบเขตของผลลัพธ์ที่คาดคะเนไว้ และ

แบบจำลองนั้นทำงานอย่างถูกต้องโดยสามารถนำผลลัพธ์นั้นมาวิเคราะห์ได้

(3) ลดความเบี่ยงเบน โดยใช้ค่าสัมมิคิยาภันเพื่อลดความแปรผัน และเพิ่มความถูกต้องเมื่อเปรียบเทียบกับองค์ประกอบที่ต่างกัน ได้

(4) มีลักษณะเป็นการเลียนแบบสถานการณ์จริงมากกว่าเป็นการนำเสนอสถานการณ์จริง

(5) มีลักษณะเป็นการคาดการณ์สถานการณ์จริงที่จะเกิดขึ้นภายใต้เงื่อนไขต่าง ๆ กัน

(6) เป็นแบบจำลองที่ใช้กับปัญหาที่มีความซับซ้อนสูง

ประโยชน์ของการจำลองสถานการณ์

การที่ระบบสนับสนุนการตัดสินใจมีการใช้แบบจำลองสถานการณ์เพื่อเลียนแบบสถานการณ์ปัญหาต่าง ๆ ก่อให้เกิดประโยชน์ดังต่อไปนี้

(1) แบบจำลองสถานการณ์เป็นทฤษฎีที่มีการใช้งานเพื่อคาดการณ์เหตุการณ์ในอนาคตอย่างตรงไปตรงมา

(2) แบบจำลองสถานการณ์สามารถทำงานที่มีเวลาเข้าไปเกี่ยวข้องเป็นจำนวนมาก ๆ ได้ตี

(3) แบบจำลองสถานการณ์ค่อนข้างเป็นการอธิบายให้เห็นเป็นรูปร่างมากกว่าการใช้เป็นเครื่องมือธรรมชาติ

(4) ผู้สร้างระบบการตัดสินใจสามารถติดต่อกับผู้ใช้ได้เพื่อรับรู้เรื่องราวเกี่ยวกับปัญหาได้อย่างลึกซึ้ง

(5) สามารถสร้างแบบจำลองสถานการณ์ที่มาจากการบันทึกของผู้ใช้ได้

(6) แบบจำลองถูกสร้างขึ้นเฉพาะเหตุการณ์ที่มาจากการบันทึกของผู้ใช้ได้

(7) แบบจำลองสถานการณ์สามารถจัดการกับปัญหาได้มากหมายหลากหลายชนิด เช่น การจัดการกับคลังสินค้า (Inventory) และการจัดการทรัพยากรบุคคล (Human Resource) อีกทั้งสามารถทำหน้าที่ในเชิงบริหารระดับสูงอีกด้วย เช่น การวางแผนการต่าง ๆ ในระยะยาว ดังนั้นผู้ใช้สามารถนำแบบจำลองสถานการณ์มาใช้ชักการสิ่งต่าง ๆ ได้ทุกช่วงเวลา

(8) สามารถทดลองป้อนตัวแปรที่แตกต่างกันตามแต่ละเหตุการณ์ลงในแบบจำลองเพื่อดูผลลัพธ์ที่เป็นทางเลือกต่าง ๆ จากนั้นจึงเลือกทางเลือกที่ดีที่สุดเพียงอย่างเดียว

(9) โดยทั่วไปแล้วแบบจำลองชนิดนี้มักจะนำมาใช้เพื่อรับรู้ปัญหาของเหตุการณ์จริงที่มีความซับซ้อนกล่าวคือหากเป็นปัญหาง่าย ๆ ก็ไม่จำเป็นต้องใช้แบบจำลองชนิดนี้ ตัวอย่างเช่น อาจใช้แบบจำลองสถานการณ์เพื่อการกระจายความน่าจะเป็นจริง ๆ หากกว่าการนำมาใช้เพื่อการประมาณการธรรมชาติซึ่งการกระจายนั้นเรากระจายโดยนำทฤษฎีมาใช้ด้วย

(10) สามารถใช้แบบจำลองสถานการณ์ เพื่อเป็นเครื่องมือวัดประสิทธิภาพของตัวแปรได้จำนวนมากและยังสามารถสะท้อนกลับมาถึงผู้ตัดสินใจได้โดยตรง

2.2 งานวิจัยที่เกี่ยวข้อง

Buck (1960) (อ้างถึงใน Little & Rubin, 1987) เสนอเงื่อนไขของวิธีการประมาณข้อมูลสูญหายด้วยค่าเฉลี่ย ซึ่งอยู่ในรูปของวิธีการประมาณข้อมูลสูญหายด้วยค่าคงคลาย โดยสร้างสมการ ลดคลายระหว่างตัวแปรที่สนใจ y กับตัวแปรช่วย x โดยใช้ข้อมูลที่สมบูรณ์ทั้งหมด r หน่วยแล้ว ประมาณค่าข้อมูลสูญหายด้วยสมการลดคลาย $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$ นำค่า \hat{y}_i ที่ได้ไปแทนที่ข้อมูลสูญหาย ของตัวแปรที่สนใจ y โดยกำหนดให้

$$y_{i,i} = \begin{cases} y_i, & \text{if } i \in R \\ \hat{\beta}_0 + \hat{\beta}_1 x_i, & \text{if } i \in R^c \end{cases} \quad (2.9)$$

เมื่อ

y_i คือ ค่าของตัวแปรที่สนใจ y ของข้อมูลที่สมบูรณ์ r หน่วย และค่าของตัวแปรที่สนใจ y ของ หน่วยสังเกตุที่ i หลังทำการประมาณค่าข้อมูลสูญหาย $k - r$ หน่วย
 y_i คือ ค่าของตัวแปรที่สนใจ y ของหน่วยสังเกตุที่ i ซึ่งมีทั้งหมด r หน่วย
 x_i คือ ค่าของตัวแปรช่วย x ของหน่วยสังเกตุที่ i โดยที่ค่าของตัวแปรที่สนใจ y เป็นค่าข้อมูลสูญหาย Lee, Rancourt, and Sarndal (1994) ศึกษาวิธีการประมาณข้อมูลสูญหายด้วยค่าอัตราส่วน และวิธีการประมาณข้อมูลสูญหายด้วยค่าไกลสุด ซึ่งวิธีการประมาณข้อมูลสูญหายด้วยค่าอัตราส่วน เป็นวิธีการประมาณค่าพารามิเตอร์เดียวซึ่งศึกษาประชากรจำกัดขนาด N โดยที่ r เป็นจำนวน หน่วยของข้อมูลที่สมบูรณ์ โดยการหักตัวอย่างแบบสุ่มเชิงเดียวชนิดไม่คืนที่ (Simple random sampling without replacement หรือ SRSWOR) จากตัวอย่างที่สุ่มมา n หน่วย จากประชากร $N = 100$ แล้วตัดข้อมูลโดยสุ่ม จากนั้นแทนค่าข้อมูลสูญหายด้วยค่าของตัวประมาณค่าแบบ อัตราส่วนโดยใช้ข้อมูลที่สมบูรณ์ทั้งหมด r หน่วย กำหนดให้

$$y_{i,i} = \begin{cases} y_i, & \text{if } i \in R \\ \hat{b}x_i, & \text{if } i \in R^c \end{cases} \quad \text{เมื่อ} \quad \hat{b} = \frac{\sum_{i=1}^r y_i}{\sum_{i=1}^r x_i} \quad (2.10)$$

ตัวประมาณค่าเฉลี่ยประชากรที่ใช้ข้อมูลที่ประมาณค่าข้อมูลสูญหายด้วยสมการ (2.10)
อยู่ในรูป

$$\bar{y}_{RAT} = \bar{y}_r \frac{\bar{x}_n}{\bar{x}_r} \quad (2.11)$$

เมื่อ

\bar{x}_r คือ ค่าเฉลี่ยของตัวแปรช่วย x ที่มีค่าข้อมูลของตัวแปรที่สนใจ y สมบูรณ์ทั้งหมด r หน่วย
 \bar{x}_n คือ ค่าเฉลี่ยของตัวแปรช่วย x จากข้อมูลที่สมบูรณ์ทั้งหมด n หน่วย

\bar{y}_r คือ ค่าเฉลี่ยของตัวแปรที่สนใจ y จากข้อมูลที่สมบูรณ์ทั้งหมด r หน่วย เป็นตัวประมาณค่าที่เอนเอียง โดยมีความเอนเอียง (biasedness) ดังนี้

$$b(\bar{y}_{RAT}) = \left(\frac{1}{r} - \frac{1}{N} \right) \bar{Y} (C_x^2 - \rho_{yx} C_y C_x) \quad (2.12)$$

และมีค่าคาดคะเนก่อนกำลังสองเฉลี่ย

$$MSE(\bar{y}_{RAT}) = \left(\frac{1}{n} - \frac{1}{N} \right) S_y^2 + \left(\frac{1}{r} - \frac{1}{n} \right) \left\{ S_y^2 + \left(\frac{\bar{Y}}{\bar{X}} \right)^2 S_x^2 - 2 \left(\frac{\bar{Y}}{\bar{X}} \right) S_{xy} \right\} \quad (2.13)$$

ส่วนวิธีการประมาณข้อมูลสูญหายด้วยค่าใกล้สุด เป็นวิธีการประมาณค่าพารามิเตอร์เดียว กำหนดให้

$$y_{i_r} = \begin{cases} y_i & \text{if } i \in R \\ y_{NNI} & \text{if } i \in R^c \end{cases} \quad (2.14)$$

เมื่อ $y_{NNI} = y_l$ ถ้า $\min_{1 \leq l \leq r} |x_i - x_l|$ สำหรับ $l \in r$ และ $l \neq i$

และ y_{NNI} คือ ค่าของตัวแปรที่สนใจ y เป็นค่าข้อมูลสูญหายที่ถูกประมาณค่าด้วยวิธีการประมาณข้อมูลสูญหายด้วยค่าใกล้สุด

นอกจากนี้พิจารณาเปรียบเทียบประสิทธิภาพของการประมาณค่าข้อมูลสูญหายทั้ง 3 วิธี คือวิธีการประมาณข้อมูลสูญหายด้วยค่าอัตราส่วน วิธีการประมาณข้อมูลสูญหายด้วยค่าใกล้สุด และ วิธีการประมาณข้อมูลสูญหายด้วยค่ามัลติเพลท (Multiple imputation) โดยใช้เทคนิค蒙ติคาร์โล ซักตัวอย่างขนาด 30% ของขนาดประชากรและซักตัวอย่างด้วยอัตราค่าข้อมูลสูญหาย 30% ของขนาด ตัวอย่าง ซึ่งค่าของตัวแปรช่วง x มีค่าเฉลี่ยเท่ากับ 48 และค่าความแปรปรวนเท่ากับ 768 และมีการแจกแจงแบบแกมมา (Gamma-distribution) โดยที่ตัวแปรช่วง x มีความสัมพันธ์ในระดับ 0.75 กับตัวแปรที่สนใจ y ที่มีการแจกแจงแบบแกมมา เช่นกันและมีค่าเฉลี่ยเท่ากับ $\mu(x) = a + bx + cx^2$ และค่าความแปรปรวนเท่ากับ $\sigma^2(x) = d^2 x^{2g}$ กำหนดให้ค่า a, b, c, d และ g คือค่าคงตัว สามารถเขียนการแจกแจงแบบแกมมาได้ดังนี้ $\frac{1}{\Gamma(\alpha)\beta^\alpha} x^{\alpha-1} \exp(-x/\beta)$ โดยที่ $x > 0$ และค่า $\mu(x) = a + bx + cx^2 = \alpha\beta$ และค่า $\sigma^2(x) = d^2 x^{2g} = \alpha\beta^2$ โดยใช้เกณฑ์ในการพิจารณาเปรียบเทียบ 3 วิธี คือรากที่สองของค่าคาดคะเนก่อนกำลังสองเฉลี่ย ค่าความเอนเอียงสัมพัทธ์ สัมบูรณ์ (Absolute relative bias) และอัตราความครอบคลุม (Coverage rate) ได้ข้อสรุปว่าวิธีการประมาณข้อมูลสูญหายด้วยค่าอัตราส่วนเป็นวิธีที่เหมาะสมสมเมื่อพิจารณาจากเกณฑ์ของวิธีค่าความเอนเอียงสัมพัทธ์สัมบูรณ์

Chaimongkol and Suwattee (2005 c) ศึกษาวิจัยเรื่อง Weighted Nearest Neighbor and Regression Imputation ซึ่งเป็นวิธีการประมาณค่าข้อมูลสูญหายด้วยวิธีการประมาณข้อมูลสูญหายด้วยค่าใกล้สุดและวิธีการประมาณข้อมูลสูญหายด้วยค่าถดถอยมาร่วมกัน เพื่อนำมาใช้ในการจัดการกับข้อมูลสูญหายที่เกิดจากการไม่ตอบเชิงบางคำถatement หรือเชิงบางตัวแปร โดยมีข้อสมมุติว่าความแปรผันของ x และ y มีความสัมพันธ์กันในรูปแบบ $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$ โดยที่ β_0 และ β_1 คือค่าคงตัวและ $\varepsilon_i \sim NID(0, \sigma^2)$ โดยกำหนดให้

$$y_{i \cdot} = \begin{cases} y_i & \text{if } i \in R \\ \hat{y}_i^* = w_i y_i^* + (1 - w_i) \hat{y}_i & \text{if } i \in R^c \end{cases} \quad (2.15)$$

เมื่อ y_i^* คือ การประมาณค่าข้อมูลสูญหายของตัวแปรที่สนใจ y ด้วยวิธีการประมาณข้อมูลสูญหายด้วยค่าใกล้สุด

\hat{y}_i คือ การประมาณค่าข้อมูลสูญหายของตัวแปรที่สนใจ y ด้วยวิธีการประมาณข้อมูลสูญหายด้วยค่าถดถอย

$$\text{และ } w_i = \frac{\sum_{i=1}^r (x_i - \bar{x})^2 + r(x_i - \bar{x})^2}{(r+1)\sum_{i=1}^r (x_i - \bar{x})^2 + r(x_i - \bar{x})^2} \quad \text{เป็นค่าที่ใช้ในการถ่วงน้ำหนัก}$$

ตัวประมาณค่าเฉลี่ยประชากรที่ใช้ข้อมูลที่ประมาณค่าข้อมูลสูญหายด้วยสมการ (2.15)

อยู่ในรูป

$$\bar{y}_{WNR} = \frac{1}{n} \sum_{i=1}^n y_{i \cdot} \quad (2.16)$$

และมีค่าความแปรปรวน

$$V(\bar{y}_{WNR}) = \sigma^2 \left\{ \frac{\sum_{i=1}^r (x_i - \bar{x})^2 + r(x_i - \bar{x})^2}{(r+1)\sum_{i=1}^r (x_i - \bar{x})^2 + r(x_i - \bar{x})^2} \right\} \quad (2.17)$$

โดยที่

$$\sigma^2 = \left\{ \frac{\sum_{i=1}^r (y_i - \hat{y}_i^*)^2}{(r-1)\bar{x} \left\{ 1 - \left(\frac{1}{r} \right) (CV_{x_r})^2 \right\}} \right\} \quad (2.18)$$

$$\text{และ } CV_{x_r} = \frac{S_{x_r}}{\bar{x}_r}$$

ในงานวิจัยนี้พิจารณาเปรียบเทียบประสิทธิภาพของการประมาณค่าข้อมูลสัญญาณทั้ง 3 วิธี คือ วิธีการประมาณข้อมูลสัญญาณด้วยค่าไกล์สุด วิธีการประมาณข้อมูลสัญญาณด้วยค่าถดถอย และวิธีการประมาณข้อมูลสัญญาณด้วยค่าถดถอยของค่าเข้าไกล์สุดแบบถ่วงน้ำหนักโดยใช้เกณฑ์ในการพิจารณาเปรียบเทียบ 4 วิธีคือรากที่สองของค่าคลาดเคลื่อนกำลังสองเฉลี่ย ค่าความแอนเอียงสัมพัทธ์ สัมบูรณ์ ค่าอัตราส่วนความสำเร็จในการประมาณค่า (the ratio for the performances of the imputations) และอัตราความครอบคลุม โดยใช้เทคนิคmondicar โอลจัลลงข้อมูลประชากรจำกัดขนาด $N = 400$ สุ่มตัวอย่างขนาด n มาจำนวน 10%, 20% และ 30% ของขนาดประชากร ด้วยวิธีการซักตัวอย่างแบบง่ายไม่แทนที่ อัตราค่าของข้อมูลสัญญาณที่คาดหวัง $E(m/n)$ ที่ใช้ในการพิจารณา คือ 10%, 20% และ 30% โดยมีค่าสัมประสิทธิ์สหสัมพันธ์ระหว่างตัวแปร x และ y คือ 0.882 ได้ข้อสรุปว่าวิธีการประมาณข้อมูลสัญญาณด้วยค่าถดถอยของค่าเข้าไกล์สุดแบบถ่วงน้ำหนัก เป็นวิธีที่เหมาะสมเมื่อพิจารณาจากเกณฑ์ของวิธีรากที่สองของค่าคลาดเคลื่อนกำลังสองเฉลี่ย

Chaimongkol and Suwattee (2005 a) ศึกษาวิจัยเรื่อง The Effect of Weighted Nearest

Neighbor-Regression Imputation on Simple Linear Regression Analysis เปรียบเทียบผลกระทบของวิธีการประมาณค่าข้อมูลสัญญาณ 4 วิธี คือ วิธีการประมาณข้อมูลสัญญาณด้วยค่าเฉลี่ย วิธีการประมาณข้อมูลสัญญาณด้วยค่าถดถอยโดยสุ่ม วิธีการประมาณข้อมูลสัญญาณด้วยค่าของทศศก และวิธีการประมาณข้อมูลสัญญาณด้วยค่าถดถอยของค่าเข้าไกล์สุดแบบถ่วงน้ำหนัก ในการวิเคราะห์ถดถอยอย่างง่ายโดยใช้เกณฑ์ในการพิจารณาเปรียบเทียบ 4 วิธี คือ ค่าสัมประสิทธิ์สหสัมพันธ์ ค่าสัมประสิทธิ์การตัดสินใจ ค่าเบี่ยงเบนสัมบูรณ์เฉลี่ย (Mean absolute deviation) และช่วงความเชื่อมั่นของค่าเฉลี่ย ซึ่งมีอัตราค่าของข้อมูลสัญญาณที่คาดหวังที่ใช้ในการพิจารณาคือ 10% 15% และ 20% สรุปว่าวิธีการประมาณข้อมูลสัญญาณด้วยค่าถดถอยของค่าเข้าไกล์สุดแบบถ่วงน้ำหนัก เป็นวิธีที่เหมาะสม เมื่อพิจารณาจากเกณฑ์ของค่าสัมประสิทธิ์สหสัมพันธ์ ค่าสัมประสิทธิ์การตัดสินใจดีกว่าวิธีอื่น ๆ ค่าเบี่ยงเบนสัมบูรณ์เฉลี่ยต่ำสุดและมีช่วงความเชื่อมั่นของค่าเฉลี่ยแคบกว่าวิธีอื่น ๆ ภายใต้ตัวแบบจำลองที่ศึกษา

Singh et al. (2010) ศึกษาวิจัยเรื่อง Estimation of population mean using imputation techniques in sample surveys โดยเสนอตัวประมาณค่าข้อมูลสัญญาณ 2 วิธี วิธีแรกคือการแทนค่าข้อมูลสัญญาณ โดยใช้เฉพาะข้อมูลของตัวแปรช่วย x ของหน่วยตัวอย่างที่ข้อมูลของตัวแปรที่สนใจ y สมบูรณ์ซึ่งมีทั้งหมด r หน่วย ($r < n$) โดยกำหนดให้

$$y_{i,i} = \begin{cases} y_i & \text{if } i \in R \\ \frac{1}{(n-r)} \bar{y}_r \left[\frac{n\{(A+C)\bar{X} + fB\bar{x}_r\}}{\{(A+fB)\bar{X} + C\bar{x}_r\}} - r \right] & \text{if } i \in R^c \end{cases} \quad (2.19)$$

เมื่อ $A = (d-1)(d-2)$, $B = (d-1)(d-4)$, $C = (d-2)(d-3)(d-4)$ และ $f = \frac{n}{N}$

ตัวประมาณค่าเฉลี่ยประชากรที่ใช้ข้อมูลที่ประมาณค่าข้อมูลสูญหายด้วยสมการ (2.19) อยู่ในรูป

$$T_{d1} = \bar{y}_r \left[\frac{(A+C)\bar{X} + fB\bar{x}_r}{(A+fB)\bar{X} + C\bar{x}_r} \right] \quad (2.20)$$

เป็นตัวประมาณค่าที่เอนเอียง โดยมีความเอนเอียง (biasedness) ดังนี้

$$b(T_{d1}) = \left(\frac{1}{r} - \frac{1}{N} \right) \bar{Y} \varphi \left[\rho_{YX} C_Y C_X - \varphi_2 C_X^2 \right] \quad (2.21)$$

เมื่อ $\varphi = \varphi_1 - \varphi_2$, $\varphi_1 = \frac{fB}{A+fB+C}$ และ $\varphi_2 = \frac{C}{A+fB+C}$

โดยที่ C_Y คือ สัมประสิทธิ์การแปรผันของข้อมูลของตัวแปรที่สนใจ Y

C_X คือ สัมประสิทธิ์การแปรผันของข้อมูลของตัวแปรช่วง X

และ ρ_{YX} คือ สัมประสิทธิ์สหสัมพันธ์ระหว่างตัวแปรที่สนใจ Y และตัวแปรช่วง X
และมีค่าคาดคะเนดังนี้

$$MSE(T_{d1}) = \left(\frac{1}{r} - \frac{1}{N} \right) \bar{Y}^2 \left[C_Y^2 + \varphi^2 C_X^2 + 2\varphi \rho_{YX} C_Y C_X \right] \quad (2.22)$$

ซึ่งพบว่า $MSE(T_{d1})$ มีค่าต่ำสุดก็ต่อเมื่อ

$$\varphi = \left[\frac{fB-C}{A+fB+C} \right] = -\rho_{YX} \frac{C_Y}{C_X} \quad (2.23)$$

จากสมการ (2.23) เมื่อทราบค่า ρ_{YX} C_Y และ C_X สามารถหาค่า d ที่เหมาะสมได้ และจะได้ค่า
คาดคะเนดังนี้

$$MSE(T_{d1})_{opt} = \left(\frac{1}{r} - \frac{1}{N} \right) (1 - \rho_{YX}^2) S_Y^2 \quad (2.24)$$

เมื่อ S_Y^2 คือ ความแปรปรวนของข้อมูลของตัวแปรที่สนใจ Y

สำหรับวิธีที่สองที่ Singh et al. เสนอ คือการแทนค่าข้อมูลสูญหายโดยใช้ข้อมูลของ
ตัวแปรช่วง x จากทั้ง n หน่วย โดยกำหนดให้

$$y_{\cdot i} = \begin{cases} y_i & \text{if } i \in R \\ \frac{1}{(n-r)} \bar{y}_r \left[\frac{n \{(A+C)\bar{X} + fB\bar{x}_n\}}{\{(A+fB)\bar{X} + C\bar{x}_n\}} - r \right] & \text{if } i \in R^c \end{cases} \quad (2.25)$$

ตัวประมาณค่าเฉลี่ยประชากรที่ใช้ข้อมูลที่ประมาณค่าข้อมูลสัญญาด้วยสมการ (2.25) อยู่ในรูป

$$T_{d2} = \bar{y}_r \left[\frac{(A+C)\bar{X} + fB\bar{x}_n}{(A+fB)\bar{X} + C\bar{x}_n} \right] \quad (2.26)$$

ความเอนเอียง (biasedness) คือ

$$b(T_{d2}) = \left(\frac{1}{n} - \frac{1}{N} \right) \bar{Y} \varphi \left[\rho_{yx} C_y C_x - \varphi_2 C_x^2 \right] \quad (2.27)$$

มีค่าคลาดเคลื่อนกำลังสองเฉลี่ย

$$MSE(T_{d2}) = \bar{Y}^2 \left[\left(\frac{1}{r} - \frac{1}{N} \right) C_y^2 + \left(\frac{1}{n} - \frac{1}{N} \right) [\varphi^2 C_x^2 + 2\varphi\rho_{yx} C_y C_x] \right] \quad (2.28)$$

และค่าคลาดเคลื่อนกำลังสองเฉลี่ยที่เหมาะสมดังนี้

$$MSE(T_{d2})_{opt} = \left[\left(\frac{1}{r} - \frac{1}{N} \right) - \left(\frac{1}{n} - \frac{1}{N} \right) \rho_{yx}^2 \right] S_y^2 \quad (2.29)$$

พิจารณาเปรียบเทียบประสิทธิภาพของตัวประมาณค่าที่ใช้ประมาณค่า \bar{Y} ได้แก่ ตัวประมาณค่าที่ใช้วิธีการประมาณข้อมูลสัญญาด้วยค่าเฉลี่ย วิธีการประมาณข้อมูลสัญญาด้วยค่าอัตราส่วนและวิธีการประมาณข้อมูลสัญญาด้วยค่าคอมโพรไมซ์ ตัวประมาณค่า T_{d1} และ T_{d2} โดยใช้เกณฑ์ในการพิจารณาเปรียบเทียบ คือ ประสิทธิภาพสัมพห์ (Relative efficiency) ได้ข้อสรุปว่าตัวประมาณค่า T_{d1} เป็นวิธีที่ดีที่สุดสามารถอธิบายได้อย่างมีประสิทธิภาพเมื่อศึกษาด้วยการจำลอง

จากการศึกษาเอกสารและงานวิจัยที่เกี่ยวข้องจะเห็นได้ว่าตัวประมาณค่า T_{d1} มีประสิทธิภาพดีที่สุดเมื่อเปรียบเทียบกับวิธีการประมาณข้อมูลสัญญาด้วยค่าเฉลี่ย วิธีการประมาณข้อมูลสัญญาด้วยค่าอัตราส่วน วิธีการประมาณข้อมูลสัญญาด้วยค่าคอมโพรไมซ์ และตัวประมาณค่า T_{d2} ผู้วิจัยจึงสนใจที่จะพัฒนาตัวประมาณค่า T_{d1} ให้มีประสิทธิภาพสูงยิ่งขึ้น