

เอกสารและผลงานวิจัยที่เกี่ยวข้อง

เนื้อหาในบทนี้จะแบ่งออกเป็น 5 ส่วนคือ

1. ความเป็นมาและวิธีการตรวจสอบความลำเอียงของข้อสอบ
2. วิธีตรวจสอบความลำเอียงของข้อสอบตามแนวคิดทฤษฎีการตอบข้อสอบ
3. วิธีตรวจสอบความลำเอียงของข้อสอบตามวิธีของแมนเทลและแฮนส์เซล
4. วิธีตรวจสอบความลำเอียงของข้อสอบตามการทดสอบ SIB (SIBTEST)
5. งานวิจัยที่เกี่ยวข้องกับการตรวจสอบความลำเอียงของข้อสอบ

1. ความเป็นมาและวิธีการตรวจสอบความลำเอียงของข้อสอบ

การศึกษาความลำเอียงของข้อสอบหรือแบบสอบนั้น เริ่มมาจากการสังเกตปรากฏการณ์ของผลการสอบเพื่อการคัดเลือก (Selection bias) ซึ่งพบว่าไม่เป็นไปตามสัดส่วน ระดับสติปัญญาหรือโครงสร้างของประชากร ไม่ว่าจะเป็นการคัดเลือกเพื่อการศึกษาต่อเพื่อการบรรจุเข้าทำงาน เพื่อเลื่อนตำแหน่งหรืออื่น ๆ ส่งผลให้เกิดการศึกษาความลำเอียงของแบบสอบ และข้อสอบรายข้อในแบบสอบที่ใช้ (test item bias) ในระยะต่อ ๆ มา

อย่างไรก็ตามแม้แนวคิดในการตรวจสอบความยุติธรรมของการใช้ข้อสอบ เพื่อการคัดเลือกบุคคลจะมีมาตั้งแต่ปี 1951 ดังได้กล่าวมาแล้ว แต่เหตุการณ์ที่ทำให้มีการตื่นตัวเป็นอันมากในเรื่องของความลำเอียงของข้อสอบ ได้แก่ กรณีที่เรียกว่า Defunis and Odegaard ซึ่งเกิดขึ้นในปี 1971 (Breland และ Ironson, 1976: 89) เป็นกรณีที่ Marco Defunis และคณะ ซึ่งถูกปฏิเสธจากโรงเรียนกฎหมายของมหาวิทยาลัยแห่งวอชิงตันได้ฟ้องร้องว่า เขาได้คะแนนการสอบสูงกว่าผู้ได้รับการคัดเลือกบางคนที่มหาวิทยาลัยรับเข้าศึกษา และยื่นฎีกาฟ้องร้อง Charles Odegaard อธิการบดีมหาวิทยาลัยดังกล่าวและคณะเป็นจำเลย เพื่อให้พิจารณาทบทวนการคัดเลือกนักศึกษาใหม่และจากนั้นเป็นต้นมา การพิจารณาตรวจสอบความลำเอียงของข้อสอบระหว่างผู้สอบกลุ่มย่อยเป็นสิ่งที่มีการปฏิบัติกันจนปัจจุบันเป็นเสมือนส่วนหนึ่งของขั้นตอนการ

พัฒนาแบบทดสอบโดยทั่ว ๆ ไป โดยเฉพาะอย่างยิ่งสำหรับการพัฒนาแบบทดสอบมาตรฐาน ตัวอย่างเช่น ศูนย์บริการการทดสอบทางการศึกษา (Educational Testing Service) ของประเทศสหรัฐอเมริกา ซึ่งให้บริการด้านการทดสอบด้วยแบบสอบมาตรฐานทั่วประเทศ สหรัฐอเมริกา ได้มีการตรวจสอบความลำเอียงทั้งก่อนและหลังจากการนำแบบทดสอบไปทดลอง ใช้การตรวจสอบความลำเอียงของข้อสอบก่อนนำไปทดลองใช้จะใช้การตรวจสอบโดยผู้เชี่ยวชาญ ซึ่งประกอบด้วยคณะบุคคลหลายฝ่ายตรวจสอบความไว (sensitivity) ของแบบสอบ โดยพิจารณาถึงรูปแบบของข้อสอบ เนื้อหา ภาพประกอบคำที่ใช้และอื่น ๆ เพื่อไม่ให้มีความลำเอียง เป็นประโยชน์หรือเสียประโยชน์สำหรับผู้สอบกลุ่มย่อยใด ๆ เช่น ในกลุ่มผู้สอบที่ต่างเพศ ต่างเชื้อชาติ และวัฒนธรรม เป็นต้น

เท่าที่ผ่านมาจนถึงปัจจุบัน มีผู้เสนอวิธีการตรวจสอบความลำเอียงของข้อสอบโดยพิจารณาจากค่าสถิติอยู่หลายวิธีดังนี้

1) การพิจารณาความแตกต่างกันในค่าเฉลี่ยระหว่างกลุ่ม กรณีเช่นนี้อาจไม่ได้แสดงถึงความลำเอียงของข้อสอบหรือแบบสอบ เพราะกลุ่มย่อยทั้งสองกลุ่มอาจมีความแตกต่างกันอยู่แล้วในสิ่งที่วัด

2) การวิเคราะห์หัตถ์ประกอบ เป็นการดูความสัมพันธ์ระหว่างคะแนนย่อย ๆ ในแบบสอบสำหรับกลุ่มย่อย ๆ โดยมากใช้เทคนิควิเคราะห์แบบ maximum likelihood ปัญหาอยู่ที่การเปรียบเทียบโครงสร้างระหว่างกลุ่มว่าจะใช้เกณฑ์ใดดีที่สุดรวมทั้งความซับซ้อนของวิธีการและการวิเคราะห์

3) การเปรียบเทียบค่าความยาก เป็นการพิจารณาลำดับที่สัมพันธ์ของค่าความยากที่ได้จากการวิเคราะห์แยกกลุ่ม เป็นการพิจารณาปฏิสัมพันธ์ระหว่างข้อกระทงและกลุ่มผู้สอบ ข้อกระทงใดที่มีลักษณะของความสัมพันธ์สำหรับกลุ่มบางกลุ่มแตกต่างไปจากข้ออื่น ๆ ในแบบสอบเดียวกันจะเป็นข้อสอบที่ถือว่ามีความลำเอียง วิธีการทำได้ทั้งการวิเคราะห์ความแปรปรวนและการใช้แผนภูมิเส้น

ผู้ที่พัฒนาวิธีการนี้ได้สมบูรณ์ที่สุด คือ Angoff ซึ่งแปลงค่า p ให้เป็น ค่า z มาตรฐานเพื่อจัดความสัมพันธ์เชิงเส้นโค้งระหว่างค่า P ของข้อสอบ 2 ชุด แล้วแปลงค่า z เป็นเคลต้า (Δ) มีค่าเฉลี่ย = 13 และค่าส่วนเบี่ยงเบนมาตรฐาน = 4 เพื่อกำจัดค่าที่เป็นลบ แล้วพิจารณาค่าเคลต้าของข้อสอบข้อเดียวกันจากกลุ่มที่ต่างกันโดยใช้แนวเส้นทะแยงที่เป็นเส้นแกนหลัก ข้อสอบจะลำเอียงหรือไม่พิจารณาได้จากระยะทางจากจุดค่าเคลต้าที่ลากไปตั้งฉากกับเส้นแกนหลัก

4) การใช้ทฤษฎีการตอบข้อสอบ เป็นการเปรียบเทียบความน่าจะเป็นในการตอบข้อสอบถูกของกลุ่มผู้สอบข้อสอบ โดยพิจารณาถึงค่าพารามิเตอร์ของข้อสอบ และผู้สอบตามแบบจำลองที่เลือกใช้ซึ่งมีทั้งชนิด 1 2 และ 3 พารามิเตอร์

5) การวิเคราะห์ด้วยค่าโคสแควร์ ซึ่งมีผู้เสนอหลายคน เช่น วิธีของ Scheuneman (1979) เป็นวิธีการที่แปลงค่าความสามารถของสมาชิกในกลุ่มเป็นช่วง ๆ อาจเป็น 3 หรือ 5 ช่วง แล้วใส่ข้อมูลในตาราง 2 ทาง (จำนวนกลุ่ม X การตอบถูก - ผิด) จำนวนตารางเท่ากับจำนวนช่วงคะแนนจากนั้นคำนวณค่าโคสแควร์แต่ละตาราง และเมื่อรวมค่าทุกตาราง จะได้ค่าโคสแควร์ที่มีขึ้นแห่งความเป็นอิสระเท่ากับจำนวนกลุ่ม การจัดกลุ่มตามระดับคะแนนเช่นนี้ ถ้ากลุ่มผู้สอบ 2 กลุ่ม มีคะแนนเฉลี่ยต่างกันแล้วจะเกิดผลที่ดูเป็นความลำเอียงได้

6) ที่แตกแยกออกจาก Loglinear ได้แก่วิธีของ Mantel-Haenszel ที่ศึกษาตัวแปรประเภทแบ่งสองในกลุ่มที่แตกต่างกันด้วยตัวแปรบล็อก และนำมาใช้เป็นตัวบ่งชี้ความลำเอียงของข้อสอบ

7) Linn และ Harnisch, (1981) ได้เสนอวิธีวิเคราะห์ในลักษณะทฤษฎีการตอบข้อสอบแบบ 3 พารามิเตอร์ โดยรวมกลุ่มย่อยเข้าด้วยกัน แล้วแบ่งความสามารถเป็นช่วงดูความแตกต่างระหว่างความน่าจะเป็นของการตอบถูกที่คาดหวัง ซึ่งประมาณได้ด้วยทฤษฎีการตอบข้อสอบกับความน่าจะเป็นของการตอบถูกจากข้อมูลเชิงประจักษ์ พิจารณาเป็นช่วงความสามารถแล้วนำมารวมกันเป็นดัชนีรวมของความลำเอียงของแต่ละข้อ วิธีนี้เรียกว่าวิธีทฤษฎีการตอบข้อสอบแบบเทียบ และใช้ได้ในกรณีที่มีกลุ่มตัวอย่างน้อยทำให้ใช้ทฤษฎีการตอบข้อสอบแบบเต็มรูปไม่ได้

8) วิธีล่าสุดที่เสนอโดย Shealy and Stout (1989) ได้แก่ การทดสอบ SIB (หรือ SIBTEST) เป็นวิธีที่ดัดแปลงจากวิธีทฤษฎีการตอบข้อสอบ สามารถใช้ได้ทั้งเพื่อการตรวจสอบความลำเอียงของข้อสอบรายข้อการตรวจสอบความลำเอียงของกลุ่มข้อสอบหรือแบบสอบทั้งฉบับ

การวิจัยครั้งนี้วิธีการตรวจสอบความลำเอียงของข้อสอบที่ผู้วิจัยนำมาใช้ในการศึกษาได้แก่การใช้ทฤษฎีการตอบข้อสอบ วิธีของ Mantel และ Haenszel และวิธีการทดสอบ SIB ซึ่งจะได้กล่าวอย่างละเอียดต่อไป

2. วิธีตรวจสอบความลำเอียงของข้อสอบตามแนวคิดทฤษฎีการตอบข้อสอบ

2.1 ทฤษฎีการตอบข้อสอบ (Item Response Theory)

ทฤษฎีการตอบข้อสอบเป็นทฤษฎีที่มีกำเนิดมาจากการที่นักวัดผลและนักการศึกษาได้พยายามคิดค้นหาวิธีการแก้ไขข้อบกพร่องของทฤษฎีการทดสอบแบบดั้งเดิม (classical test theory) โดยเริ่มที่งานของ Binet และ Simon ในปี 1916 ที่ได้เริ่มทำการกราฟของความสัมพันธ์ระหว่างผลการสอบและตัวแปรอิสระอื่น ๆ และใช้กราฟนั้นในการพัฒนาแบบสอบ และในปี 1936 Richardson ได้แสดงความสัมพันธ์ระหว่างพารามิเตอร์ตามโมเดลการตอบข้อสอบและพารามิเตอร์ ตามทฤษฎีการวัดดั้งเดิม ซึ่งเป็นวิธีเริ่มแรกในการประมาณค่าพารามิเตอร์ตามทฤษฎีการตอบข้อสอบจนกระทั่งในปี 1943-4 Lawley ได้คิดวิธีการใหม่ในการประมาณค่าพารามิเตอร์ และในปี 1952 Lord ได้เสนอโมเดลปกติสะสม normal ogive model ชนิด 2 พารามิเตอร์ การประมาณค่าพารามิเตอร์ และการนำโมเดลไปใช้ในทางปฏิบัติ ต่อมา Birnbaum (1957-8) ได้เสนอให้ใช้ logistic model แทนและพัฒนาพื้นฐานด้านสถิติสำหรับโมเดลที่เสนอใหม่

ปี 1960 Rasch ได้พัฒนาโมเดลการตอบข้อสอบขึ้น 3 โมเดล ซึ่งได้มีผู้นำไปใช้ในการศึกษาและพัฒนาต่อในภายหลัง เช่น Wright ซึ่งนอกจากจะนำไปใช้ในการศึกษาวิจัยแล้วยังได้ร่วมกับ Panchapakesan (1969) ในการเสนอวิธีการประมาณค่าพารามิเตอร์ และโปรแกรมคอมพิวเตอร์ BICAL เพื่อที่ใช้ในการคำนวณตามโมเดลของ Rasch ด้วย

และหลังจากนั้น มากก็ได้มีนักการศึกษาหลายท่านได้ช่วยกันพัฒนาโปรแกรมคอมพิวเตอร์ สำหรับการวิเคราะห์ตามทฤษฎีขึ้นนี้อีกมากมาย เช่น Lord (1982) ได้พัฒนาโปรแกรมการประมาณค่าพารามิเตอร์ ชื่อ LOGIST และ Mislevy และ Bock (1984) พัฒนาโปรแกรม BILOG เป็นต้น

โมเดลตามทฤษฎีการตอบข้อสอบแสดงถึง ความสัมพันธ์ระหว่างผลการตอบข้อสอบกับลักษณะแฝงหรือความสามารถที่ไม่สามารถสังเกตได้โดยตรง และอธิบายได้ในรูปของฟังก์ชันทางคณิตศาสตร์ โดยมีข้อสรุปว่า (1) ผลการตอบข้อสอบของผู้สอบสามารถทำนาย หรืออธิบายได้ด้วยองค์ประกอบที่เรียกว่า ลักษณะแฝง หรือความสามารถที่มีอยู่ในตัวบุคคล และ (2) ความสัมพันธ์ระหว่างผลการตอบข้อสอบ และความสามารถที่มีอิทธิพลต่อการตอบข้อสอบสามารถอธิบายได้ด้วยฟังก์ชันที่มีลักษณะเป็นโค้งสะสมเพิ่มขึ้นทางเดียว (monotonically increasing function) ซึ่งเรียกว่า ฟังก์ชันลักษณะข้อสอบอันเป็นฟังก์ชันที่แสดงว่า ผู้สอบที่มีคะแนนในลักษณะแฝงที่วัดสูงกว่า จะมีความน่าจะเป็นที่คาดหวังในการตอบข้อสอบได้ถูกต้องสูงกว่าผู้ที่ได้คะแนนการวัดลักษณะแฝงนั้นต่ำกว่า

2.2 ข้อตกลงเบื้องต้นในการใช้ทฤษฎีการตอบข้อสอบ

การวิเคราะห์ข้อสอบโดยใช้ทฤษฎีการตอบข้อสอบ มีข้อตกลงเบื้องต้นเกี่ยวกับแบบสอบดังนี้

2.2.1 ความเป็นเอกมิติ (Unidimensionality) หมายความว่า การตอบข้อสอบทุกข้อได้ถูกต้อง ต้องใช้ความสามารถหรือลักษณะแฝงภายในเพียงลักษณะเดียว หรือข้อสอบทุกข้อในแบบสอบวัดความสามารถแฝงลักษณะเดียว หรือกลุ่มลักษณะแฝงกลุ่มเดียว

2.2.2 Monotonicity

ผู้มีความสามารถในลักษณะแฝงที่วัดสูงกว่า จะมีความน่าจะเป็นในการตอบข้อสอบได้ถูกต้องมากกว่า

2.2.3 Probabilistic models

หมายความว่า เป็นแบบจำลองที่ฟังก์ชันของลักษณะข้อสอบแสดงถึงความสัมพันธ์ระหว่างตัวแปรความสามารถที่สังเกตไม่ได้และตัวแปรการตอบข้อสอบที่สังเกตได้ (ผลการตอบข้อสอบ 0 หรือ 1)

2.2.4 Local independence

การตอบข้อสอบแต่ละข้อเป็นไปอย่างอิสระต่อกัน การตอบข้อสอบข้อหนึ่งไม่มีอิทธิพลต่อการตอบข้อสอบอีกข้อหนึ่ง นั่นคือ

1) การตอบข้อสอบทุกข้อของผู้สอบคนหนึ่ง หรือการตอบข้อสอบข้อหนึ่ง ๆ ของผู้สอบทุกคนมีความเป็นอิสระต่อกัน

2) ในแง่ของสถิติ หมายความว่า ถ้าให้ θ เป็นความสามารถที่มีอิทธิพลต่อการตอบข้อสอบของผู้สอบ U_i เป็นผลการตอบข้อสอบข้อ i ของผู้สอบที่สุ่มเลือกได้ (i มีค่า = $1, 2, \dots, n$) และ $P(U_i | \theta)$ เป็นความน่าจะเป็นในการตอบข้อสอบของผู้สอบที่มีความสามารถ θ $P(U_i = 1 | \theta)$ คือ ความน่าจะเป็นของการตอบข้อสอบถูก และ $P(U_i = 0 | \theta)$ เป็นความน่าจะเป็นของการตอบข้อสอบผิด คุณสมบัติของ local independence แสดงได้ด้วย

$$\begin{aligned} \text{Prob}(U_1, U_2, \dots, U_n | \theta) &= P(U_1 | \theta) P(U_2 | \theta) \dots P(U_n | \theta) \\ &= \prod_{i=1}^n P(U_i | \theta) \end{aligned}$$

ซึ่งให้ความหมายว่า สำหรับผู้สอบใดๆ (หรือผู้สอบที่มีความสามารถตามที่กำหนด) ความน่าจะเป็นของการตอบแบบสอบถามจะเท่ากับผลคูณของความน่าจะเป็นในการตอบข้อสอบแต่ละข้อ เช่น ถ้ารูปแบบผลการตอบข้อสอบ 4 ข้อ เป็น (0, 1, 1, 1) หรือ $U_1 = 0$ $U_2 = 1$ $U_3 = 1$ และ $U_4 = 1$ ตามข้อตกลงเกี่ยวกับ local independence จะได้ว่า

$$\begin{aligned}
 P(U_1 = 0, U_2 = 1, U_3 = 1, U_4 = 1 | \theta) \\
 &= P(U_1 = 0 | \theta) P(U_2 = 1 | \theta) P(U_3 = 1 | \theta) P(U_4 = 1 | \theta) \\
 &= P_1 P_2 P_3 P_4
 \end{aligned}$$

โดยที่ $P_1 = P(U_1 = 1 | \theta)$ และ $Q_1 = 1 - P_1$

3) หลังจากควบคุมค่า θ แล้ว จะให้ความเป็นอิสระอย่างมีเงื่อนไขว่า

$$P_{i,j} | \theta = 0$$

นั่นคือ ผลการตอบข้อสอบ 2 ข้อใด ๆ ของผู้สอบ ไม่มีความสัมพันธ์กัน หรือ Stout (1987) ใช้คำว่า essential independence คือ มีค่าใกล้เคียงศูนย์

Lord (1980: 19) กล่าวว่า "Local independence เป็นคุณสมบัติที่เกิดขึ้นโดยอัตโนมัติจากคุณสมบัติของความเป็นเอกมิติ ไม่ใช่เป็นข้อตกลงเบื้องต้นที่เพิ่มขึ้นมา" ในทางปฏิบัติส่วนมากก่อนจะวิเคราะห์ข้อมูลการทดสอบด้วยโมเดลตามทฤษฎีการตอบข้อสอบ จึงมีการตรวจสอบความเป็นเอกมิติของแบบสอบถามก่อน

2.3 โมเดลที่เป็นที่นิยมในทฤษฎีการตอบข้อสอบ

ลักษณะทั่วไปของทฤษฎีการตอบข้อสอบ คือ โมเดลการตอบข้อสอบซึ่งในทางคณิตศาสตร์แสดงถึงความน่าจะเป็นในการตอบข้อสอบถูกในรูปแบบฟังก์ชันของความสามารถของผู้สอบ การตอบข้อสอบแบบสองประเภท ข้อ j แสดงได้ในรูปของคะแนน

$$X_{ji} = \begin{cases} 1 & \text{ถ้าตอบถูก} \\ 0 & \text{ถ้าตอบผิด} \end{cases}$$

ความสามารถของผู้สอบใช้แทนด้วย θ และความน่าจะเป็นของการตอบข้อสอบข้อ j ถูก คือ

$$P(X_j = 1 | \theta) = P_j(\theta)$$

ดังนั้น ความน่าจะเป็นในการตอบข้อสอบผิด จึงเป็น

$$P(X_j = 0 | \theta) = 1 - P_j(\theta)$$

โมเดลทางทฤษฎีการตอบข้อสอบที่เป็นที่นิยมใช้ในปัจจุบันมีความแตกต่างที่สำคัญ คือ จำนวนพารามิเตอร์ที่ใช้ในการบรรยายถึงข้อสอบ ซึ่งมีอยู่ 3 โมเดล คือ โมเดลชนิด 1, 2 และ 3 พารามิเตอร์ และเป็นโมเดลที่ใช้กับข้อสอบชนิดมีการให้คะแนน 2 ค่า (dichotomous) (Hambleton และคณะ 1991: 12) ดังนี้

1. โมเดลโลจิสติกชนิด 1 พารามิเตอร์ (Rasch, 1960)

โมเดลนี้เป็นโมเดลที่นิยมใช้กันกว้างขวาง ได้ังลักษณะข้อสอบสำหรับโมเดลนี้ แสดงได้ดังสมการ

$$P_i(\theta) = \frac{e^{(\theta - b_i)}}{1 + e^{(\theta - b_i)}} \quad i = 1, 2, \dots, n$$

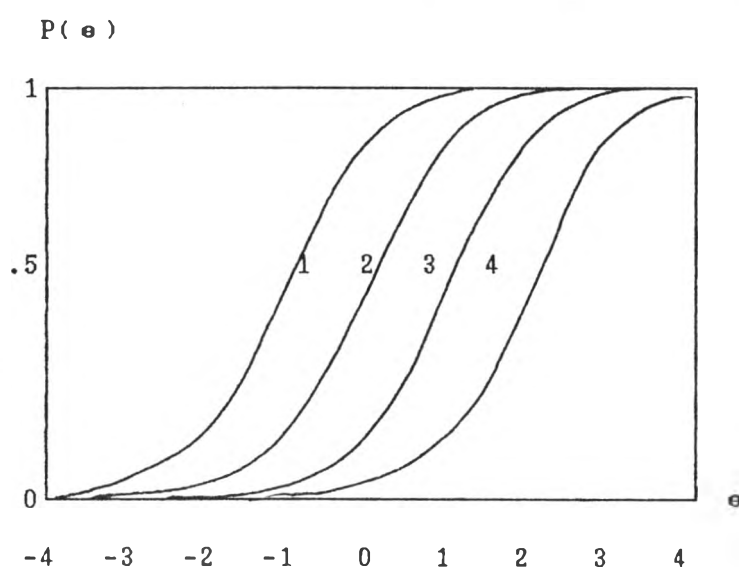
โดยที่

$P_i(\theta)$	คือ ความน่าจะเป็นของผู้สอบที่สัมพันธ์มาที่มีความสามารถ θ จะตอบข้อสอบข้อ i ได้ถูกต้อง
b_i	คือ ค่าพารามิเตอร์ความยากของข้อสอบ
n	คือ จำนวนข้อสอบในแบบสอบ
e	คือ ฐาน (base) ของลอการิทึมธรรมชาติมีค่า 2.718

ค่า b_i คือ ค่าบ่งชี้ความสามารถที่มีความน่าจะเป็นในการตอบข้อสอบถูกเป็น 0.5 ถ้า b_i มีค่าสูง แสดงว่า ข้อสอบนั้นยาก และผู้มีโอกาสจะตอบข้อสอบได้ถูกต้อง 50% ต้องมีความสามารถสูงขึ้น

ถ้าแปลงค่า θ ให้มีค่าเฉลี่ยเป็น 0 และค่าความเบี่ยงเบนมาตรฐานเป็น 1 แล้ว b จะมีค่าประมาณ -2.0 ถึง $+2.0$ (Hambleton และคณะ, 1991: 13) โมเดลนี้ถือว่า ความยากของข้อกระทงเป็นลักษณะข้อสอบลักษณะเดียวที่มีอิทธิพลต่อการตอบข้อสอบ

ภาพที่ 2 เป็นภาพโค้งลักษณะข้อสอบตามโมเดล 1 พารามิเตอร์ 4 ข้อ ซึ่งมีค่า b เป็น $-1, 0, 1$ และ 2 ตามลำดับ สังเกตได้ว่า โค้งตัดแกน $P(\theta)$ ที่ 0 และทุกโค้งมีค่าอำนาจจำแนกเท่ากัน โมเดล 1 พารามิเตอร์นี้ในทางคณิตศาสตร์เทียบได้กับโมเดลของ Rasch



ภาพที่ 2: โค้งลักษณะข้อสอบตามโมเดล 1 พารามิเตอร์

2. โมเดลโลจิสติก ชนิด 2 พารามิเตอร์

Lord เป็นคนแรกที่พัฒนาโมเดลการตอบข้อสอบ ชนิด 2 พารามิเตอร์ ขึ้นในปี 1952 โดยใช้การกระจายแบบโค้งปกติสะสม และ Birnbaum เป็นผู้เสนอให้ใช้ฟังก์ชันโลจิสติกซึ่งมีความสะดวกและง่ายในการคำนวณมากกว่าแทนฟังก์ชันโค้งปกติสะสม

โค้งลักษณะข้อสอบสำหรับโมเดลโลจิสติกแบบ 2 พารามิเตอร์ แสดงได้ในรูปสมการดังนี้

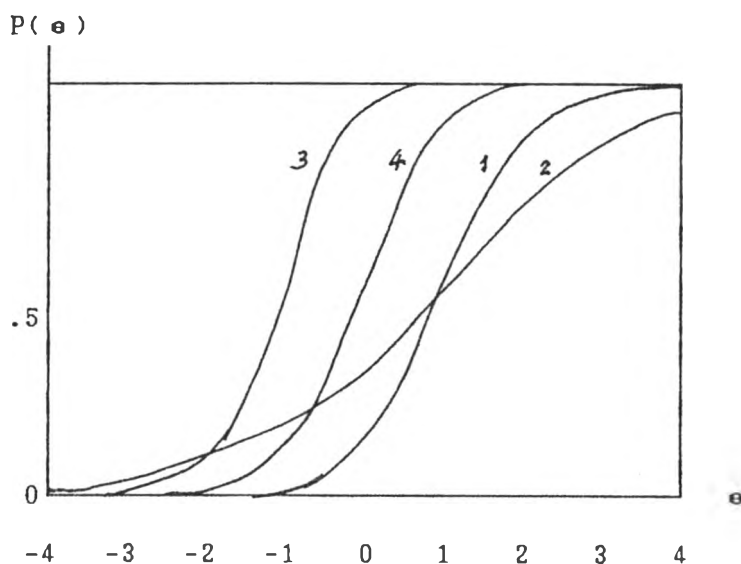
$$P_i(\theta) = \frac{Da_i(\theta - b_i)}{1 + e^{Da_i(\theta - b_i)}} \quad i = 1, 2, \dots, n$$

สิ่งที่เพิ่มเข้ามาในสมการคือ D ซึ่งได้แก่ องค์ประกอบของการปรับมาตรการที่ทำให้ฟังก์ชัน
โลจิสติกและฟังก์ชันปกติสะสมมีค่าสัมบูรณ์ของความแตกต่างน้อยกว่า 0.01 สำหรับทุกค่าของ e

สิ่งที่เพิ่มเข้ามาอีกตัว คือ ค่า a_i ซึ่งเป็นค่าพารามิเตอร์อำนาจจำแนก เป็นค่าความ
ชันของโค้งลักษณะข้อสอบที่จุดเปลี่ยนโค้งหรือที่จุด b_i บนมาตรฐานความสามารถ

ค่า a_i มีค่าอยู่ระหว่าง $-\infty$ ถึง $+\infty$ ในทางทฤษฎี แต่ในทางปฏิบัติจะมีค่าอยู่
ระหว่าง -2 ถึง 2 ข้อสอบที่มีค่า a สูง จะเป็นข้อสอบที่มีโค้งการตอบข้อสอบที่ชันมากและสามารถ
จำแนกคนได้ดีกว่าข้อสอบที่มีค่า a ต่ำกว่า

ภาพที่ 3 เป็นภาพโค้งลักษณะข้อสอบโมเดล 2 พารามิเตอร์ ซึ่งข้อ 1 มีค่า $b_1=1.0$
 $a_1=1.0$ ข้อ 2 มีค่า $b_2=1.0$, $a_2=0.5$ ข้อ 3 มีค่า $b_3=-1.0$, $a_3=1.5$ และข้อ 4 มีค่า
 $b_4=0.0$, $a_4=1.2$



ภาพที่ 3: โค้งลักษณะข้อสอบตามโมเดล 2 พารามิเตอร์

(Hambleton, Swaminathan and Rogers, 1991: 16)

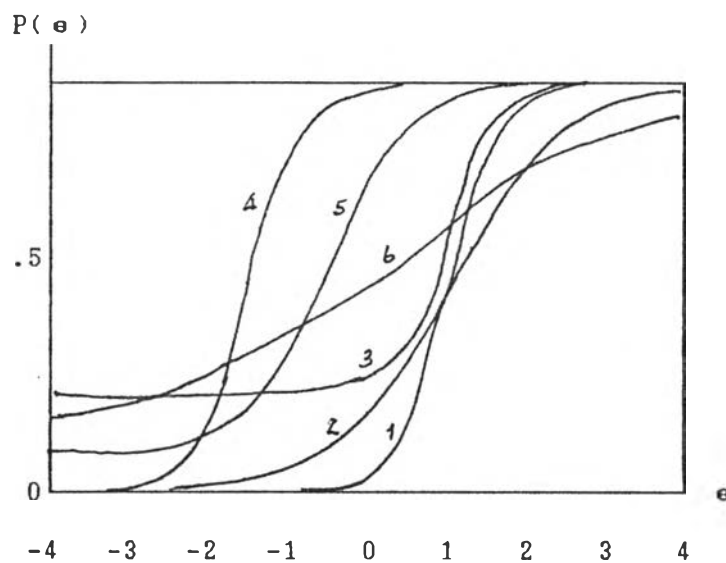


3. โมเดลโลจิสติกแบบ 3 พารามิเตอร์
สมการของฟังก์ชันโลจิสติกแบบ 3 พารามิเตอร์ คือ

$$P_i(\theta) = \frac{c_i + (1 - c_i) e^{Da_i(\theta - b_i)}}{1 + e^{Da_i(\theta - b_i)}} \quad i = 1, 2, \dots, n$$

เมื่อ c_i คือ ความน่าจะเป็นในการตอบข้อสอบเลือกตอบได้ถูกต้องอันเนื่องมาจากการเดา ถ้าตัวเลือกในแบบสอบเรียงไว้แบบสุ่ม และผู้สอบทุกคนเดาตอบ ค่าของ c_i จะเท่ากับ $1/A$ เมื่อ A เป็นจำนวนตัวเลือกในแต่ละข้อ ถ้าผู้สอบแต่ละคนเลือกตอบ โดยมีการกำจัดตัวเลือก 1 ตัว หรือมากกว่าออกก่อน ค่า c_i จะมีค่ามากกว่า $1/A$

ภาพที่ 4 เป็นภาพโค้งลักษณะข้อสอบชนิด 3 พารามิเตอร์ จะเห็นว่า ข้อสอบข้อ 1-3 ยากกว่า ข้อสอบข้อที่ 4-6 และข้อสอบข้อที่ 1, 3 และ 4 มีค่าอำนาจจำแนกสูงกว่าข้อที่ 2, 3 และ 6 และข้อสอบข้อที่ 3, 5 และ 6 แสดงถึงอิทธิพลของค่าการเดา (c_i) ต่อโค้งลักษณะข้อสอบ



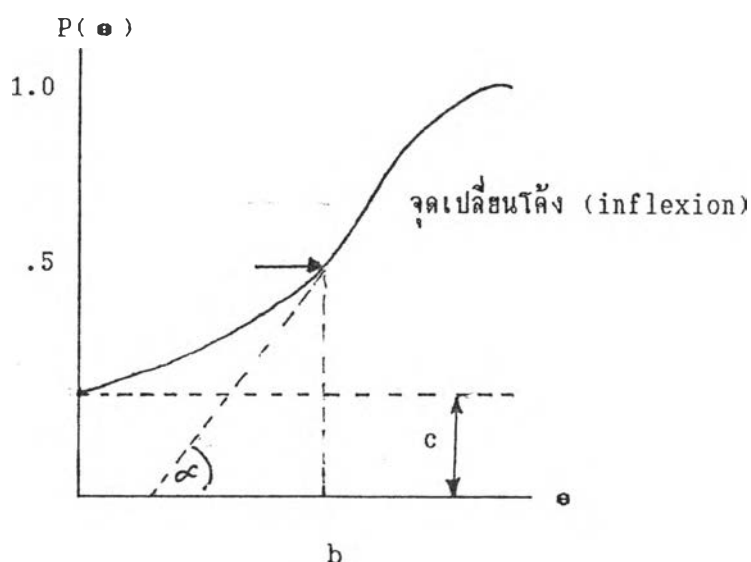
ภาพที่ 4: โค้งลักษณะข้อสอบตามโมเดล 3 พารามิเตอร์

(Hambleton, Swaminathan and Rogers, 1991: 18)

ภาพที่ 5 เป็นภาพแสดงความหมายของค่าพารามิเตอร์ข้อสอบ (Lord, 1980: 13-14) ค่าพารามิเตอร์ c คือ ความน่าจะเป็นที่ผู้สอบที่ไม่มีความสามารถในสิ่งที่วัดโดยสิ้นเชิง ($\theta = -\infty$) จะตอบข้อสอบได้ถูกต้อง เรียกว่า เป็นค่าพารามิเตอร์การเดา หรือระดับคะแนนการตอบถูกโดยบังเอิญ (guessing parameter or pseudo-chance score level) กรณีที่ข้อสอบใดไม่สามารถตอบได้ถูกต้องโดยการเดาแล้ว ค่า $c = 0$

พารามิเตอร์ b เป็นพารามิเตอร์แสดงตำแหน่ง (location parameter) เพราะเป็นตัวแสดงตำแหน่งของโค้งบนมาตรฐานความสามารถเรียกว่า ค่าความยากของข้อสอบ ยิ่งข้อสอบยากมากโค้งจะยิ่งชิดมาทางขวามากจุดเปลี่ยนโค้งของโค้งโลจิสติกอยู่ที่จุดที่ค่า $\theta = b$ กรณีไม่มีการเดา b คือ ระดับความสามารถที่มีความน่าจะเป็นในการตอบข้อสอบถูก = .5 แต่ ถ้ามีการเดา b จะเป็นค่า θ ที่มีความน่าจะเป็นในการตอบถูกอยู่ตรงกึ่งกลางระหว่างค่า c และ 1.00 ($= (1-c)/2$)

ค่าพารามิเตอร์ a เป็นสัดส่วนต่อความชันของโค้งที่จุดเปลี่ยนโค้ง ซึ่งมีค่าเท่ากับ $.425(1-c)$ ดังนั้น a จึงเป็นอำนาจจำแนกของข้อสอบแสดงระดับการตอบข้อสอบที่แปรเปลี่ยนไปภายในระดับความสามารถต่าง ๆ



ภาพที่ 5: ความหมายของค่าพารามิเตอร์ข้อสอบ (Lord, 1980: 14)

2.4 การประมาณค่าความสามารถและพารามิเตอร์ข้อสอบ

การประมาณค่าความสามารถมีวิธีการประมาณที่ใช้กันอยู่หลายวิธี เช่น (Hambleton, Swaminathan และ Rogers, 1991: 46)

1) Joint maximum likelihood procedure พัฒนาขึ้นโดย Lord ในปี 1974 และ 1980 สามารถใช้ประมาณค่าพารามิเตอร์ข้อสอบและความสามารถพร้อม ๆ กัน และสามารถใช้ได้กับทั้งโมเดล 1 2 และ 3 ค่าพารามิเตอร์

2) Marginal maximum likelihood procedure โดย Bock และ Aitkin ปี 1981 มีทั้งชนิดโมเดล 1 2 และ 3 พารามิเตอร์ ประมาณค่าพารามิเตอร์ของข้อสอบ แล้วจึงประมาณค่าพารามิเตอร์ของความสามารถ

3) Conditional maximum likelihood procedure โดย Andersen ปี 1972-3 และ Rasch ปี 1960 ใช้ได้เฉพาะโมเดล 1 พารามิเตอร์ เป็นฟังก์ชันของความเป็นไปได้ โดยมีเงื่อนไขที่คะแนนที่ได้จากการทำแบบสอบถูก

4) Joint and marginal Bayesian estimation procedures โดยมิสเลวี Mislevy ปี 1986 Swaminathan และ Gifford ปี 1982, 1985 และ 1986 ใช้ได้ทั้งกับโมเดล 1 2 และ 3 พารามิเตอร์

5) Heuristic estimation procedure โดย Urry ปี 1974 และ 1978 ใช้ได้ทั้งกับโมเดล 1 2 และ 3 พารามิเตอร์

6) Method based on nonlinear factor analysis procedure โดย McDonald ปี 1967 และ 1989 ใช้ได้กับโมเดล 2 และ 3 พารามิเตอร์ ที่กำหนดค่า c ให้เป็นค่าคงที่

ในที่นี้ผู้วิจัยขอแนะนำเสนอวิธีการประมาณค่าความสามารถและค่าพารามิเตอร์ของข้อสอบ ด้วยวิธีที่นิยมใช้ 3 วิธี ได้แก่ วิธี Joint maximum likelihood procedure ซึ่ง Wingersky, Barton และ Lord (1982) ได้ใช้เป็นฐานในการพัฒนาโปรแกรม LOGIST สำหรับ mainframe computer ซึ่งขณะนี้ยังนิยมใช้กันอยู่และในปัจจุบันศูนย์บริการทางการศึกษา (ETS-Educational Testing Service) กำลังพัฒนารูปแบบใหม่ เพื่อให้สามารถใช้กับคอมพิวเตอร์ส่วนบุคคล (PC-Computer) วิธีที่ 2 ที่จะได้แนะนำเสนอคือ วิธี Joint and marginal Bayesian estimate procedure ซึ่งเป็นส่วนหนึ่งที่ Mislevy และ Bock ได้นำไปพัฒนาโปรแกรม BILOG ในปี 1984 โดยมีทั้งโปรแกรมที่ใช้กับคอมพิวเตอร์ชนิด mainframe และ PC และวิธีสุดท้ายคือ วิธี Marginal Maximum Likelihood (MML) ซึ่งเป็นวิธีที่นำไปใช้ในโปรแกรม BILOG เช่นกัน

2.4.1 วิธีการประมาณด้วย Maximum Likelihood

สมมติให้ $(U_1, U_2, U_3, \dots, U_n)$ เป็นรูปแบบการตอบข้อสอบ n ข้อของผู้สอบ โดยที่ U_j มีค่าเป็น 1 ถ้าตอบข้อสอบถูก และ 0 ถ้าตอบข้อสอบผิด

ภายใต้ข้อตกลงในเรื่องความเป็นอิสระในการตอบข้อสอบแต่ละข้อ (Local independence) จะได้ว่า

$$\begin{aligned} & P(U_1, U_2, \dots, U_j, \dots, U_n \mid \theta) \\ &= P(U_1 \mid \theta) P(U_2 \mid \theta) \dots P(U_j \mid \theta) \dots P(U_n \mid \theta) \end{aligned}$$

หรือเขียนใหม่ได้ในรูป

$$P(U_1, U_2, \dots, U_n \mid \theta) = \prod_{j=1}^n P(U_j \mid \theta)$$

เนื่องจาก U_j มีค่าเป็น 1 หรือ 0 จึงเขียนได้ในรูปฟังก์ชันความเป็นไปได้ ดังนี้

$$P(U_1, U_2, \dots, U_n \mid \theta) = \prod_{j=1}^n P(U_j \mid \theta)^{U_j} [1 - P(U_j \mid \theta)]^{(1-U_j)}$$

$$\begin{aligned} & \text{หรือ} \\ &= \prod_{j=1}^n P_j^{U_j} Q_j^{(1-U_j)} \end{aligned}$$

$$\text{โดยที่ } P_j = P(U_j \mid \theta) \text{ และ } Q_j = 1 - P(U_j \mid \theta)$$

และเมื่อแสดงในรูปฟังก์ชันความน่าจะเป็น (likelihood function) แล้วจะได้ว่า

$$L(u_1, u_2, \dots, u_n \mid \theta) = \prod_{j=1}^n P_j^{u_j} Q_j^{(1-u_j)}$$

และเนื่องจาก P_j และ Q_j เป็นฟังก์ชันของ θ และพารามิเตอร์ของข้อสอบ ฟังก์ชันข้างบน จึงเป็นฟังก์ชันของพารามิเตอร์เหล่านี้ด้วย และจากคุณสมบัติของลอการิทึม ที่ว่า

$$\ln xy = \ln x + \ln y \quad \text{และ}$$

$$\ln x^a = a \ln x$$

ฟังก์ชันข้างบนจึงแสดงได้ในรูปของ log-likelihood ดังนี้

$$\ln L(U | \theta) = \sum_{j=1}^n [u_j \ln P_j + (1-u_j) \ln (1-P_j)]$$

โดยที่ u คือ เวกเตอร์ของการตอบข้อสอบ

ฟังก์ชันความเป็นไปได้ข้างบนมีค่าสูงสุดได้ เมื่อสามารถประมาณค่า θ ที่ใช้แทนค่าในสมการข้างล่าง ซึ่งเรียกว่า สมการ likelihood ได้

$$\frac{d}{d\theta} \ln L(U | \theta) = 0$$

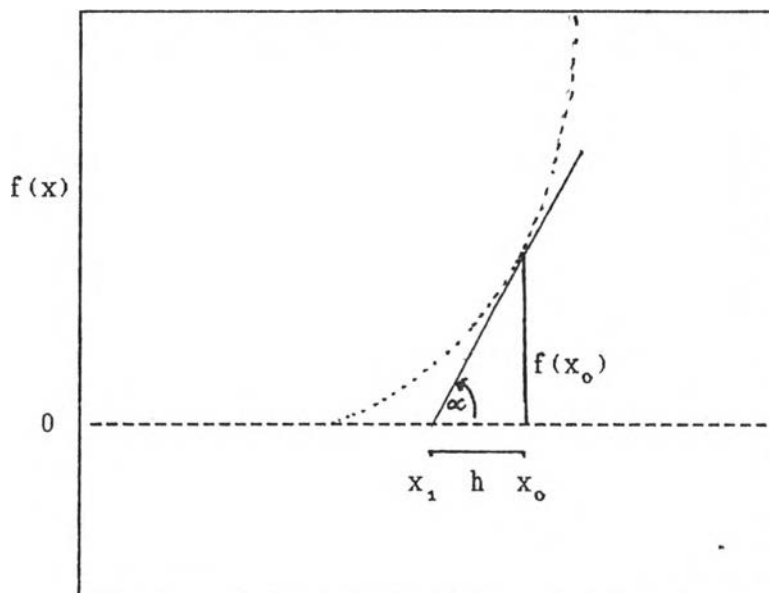
การแก้สมการ likelihood ต้องอาศัยวิธีการของ Newton-Raphson ซึ่งสมมุติว่า ถ้าสมการที่ต้องการแก้คือ

$$f(x) = 0$$

และค่าที่ประมาณได้จากสมการคือ x_0 ในขณะที่ค่าที่ถูกต้องมากกว่าคือ x_1 สามารถแสดงได้ในรูป

$$x_1 = x_0 - h$$

และจากภาพที่ 6
$$h = \frac{f(x_0)}{\tan\alpha}$$



ภาพที่ 6: ภาพแสดงวิธีการของ Newton - Raphson (Hambleton และ Swaminathan, 1985: 80)

และเนื่องจาก $\tan\alpha$ เป็นความชันของ $f(x)$ ที่ x_0 ดังนั้น $\tan\alpha = f'(x_0)$ เมื่อ $f'(x_0)$ เป็น derivative ของฟังก์ชันที่ประเมินได้ที่ x_0 ดังนั้น

$$h = \frac{f(x_0)}{f'(x_0)} \quad \text{และ}$$

$$x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}$$

และเมื่อได้ค่า x_1 แล้ว การประมาณค่า x อื่น ๆ ที่มีความถูกต้องแม่นยำกว่า x_1 ก็ดำเนินการไปตามที่ได้ประมาณค่า x_1 การประมาณค่า x จะมีการทำซ้ำจนกระทั่งความแตกต่างระหว่าง

ค่า x_m และค่า $x_{m+1} (= X_{m+1} - X_m)$ ต่ำกว่าค่าที่กำหนดไว้และค่า X_m ถือเป็นค่าที่ได้จากสมการ likelihood

จากหลักดังกล่าวของ Newton-Raphson

$$\text{เมื่อ } f(x) = \frac{d}{d\theta} \ln L(U|\theta)$$

$$\text{ดังนั้น } f'(x) = \frac{d^2}{d\theta^2} \ln L(U|\theta)$$

และถ้า θ_m เป็นการประมาณค่ารอบที่ m ของค่า θ แล้ว การประมาณค่า θ ในรอบที่ $m+1$ คือ

$$\theta_{m+1} = \theta_m - \left| \frac{d}{d\theta} \ln L(U|\theta) \right|_m / \left| \frac{d^2}{d\theta^2} \ln L(U|\theta) \right|_m$$

และการประมาณซ้ำ จะกระทำต่อไปจนกว่าจะถึงเกณฑ์ที่กำหนด ดังการประมาณค่า X ที่กล่าวข้างบน

2.4.2 วิธี Joint and Marginal Bayesian estimation

วิธี Bayesian เป็นวิธีที่ได้ค่าประมาณของความสามารถที่ถูกต้อง ถ้าทราบข้อมูลเกี่ยวกับความสามารถของผู้สอบมาบ้างแล้ว หรือมีการระบุรูปแบบการกระจายของค่าความสามารถผู้สอบได้ เช่น ระบุว่าผู้สอบเป็นตัวอย่างที่สุ่มมาจากประชากรที่มีการกระจายเป็นปกติ และมีการกระจายของความสามารถเป็นปกติ

$$\theta \sim N(\mu, \sigma^2)$$

ค่า μ และ σ^2 ต้องมีการระบุค่าว่าเท่าไร เช่น Owen (1975 อ้างใน Hambleton และ Swaminathan, 1985: 92) กำหนดให้ $\mu = 0$ และ $\sigma^2 = 1$ ในการศึกษาแบบสอบตัดแปลง (adaptive testing) การระบุการกระจายของความสามารถเป็นการกระจายปกตินับว่าเป็น

วิธีที่ค่อนข้างสะดวก แต่ Birnbaum (1969, cited in Hambleton และ Swaminathan หน้าเดียวกัน) ระบุการกระจายค่าความสามารถเดิมของผู้สอบในรูปของฟังก์ชัน $f(\theta)$ โดยที่

$$f(\theta) = \exp(\theta) / [1 + \exp(\theta)]^2$$

ในบางกรณี การระบุการกระจายของความสามารถเดิม โดยพิจารณาจากข้อมูลเชิงประจักษ์อาจเหมาะสมกว่า เช่น อาจพิจารณาจากการกระจายของคะแนนดิบ หรือคะแนนดิบแปลงตามทฤษฎีของเบย์ส ซึ่งได้กำหนดความสัมพันธ์ของความน่าจะเป็นเชิงเงื่อนไข และความน่าจะเป็นอิสระ (conditional and marginal probabilities) ในรูปของสมการดังนี้

$$P(B|A) = P(A|B) P(B)/P(A) \quad (1)$$

ในบริบทของการประมาณค่าความสามารถ A คือ θ_u และ B คือ U ซึ่งได้แก่ชุดการตอบแบบสอบถาม n ข้อ ดังนั้น สมการ (1) สามารถเขียนใหม่เป็น

$$P(\theta_u|U) = P(U|\theta_u) P(\theta_u)/P(U) \quad (2)$$

และเนื่องจาก θ_u เป็นตัวแปรต่อเนื่อง สมการ (2) จึงเป็นฟังก์ชันความหนาแน่น $P(\theta_u)$ ซึ่งในที่นี้คือ การกระจายเดิมของความสามารถซึ่งในสมการอาจสับสนกับ function การตอบข้อสอบ ดังนั้น ฟังก์ชันความหนาแน่นใน (2) จึงระบุได้ในรูปของ $f(U)$ นั่นคือ

$$f(\theta_u|U) = f(U|\theta_u) f(\theta_u)/f(U) \quad (3)$$

และสำหรับการตอบข้อสอบชุดหนึ่ง ๆ $f(U)$ คือ ค่าคงที่ $f(\theta_u|U)$ เป็นความหนาแน่นภายหลังของ θ_u และ $f(\theta_u)$ เป็นความสามารถเดิม $f(U|\theta_u)$ สามารถแสดงได้ในรูปของฟังก์ชัน likelihood ของค่าที่สังเกตได้ ดังนั้น สมการ (3) จึงเขียนใหม่ได้เป็น

$$f(\theta_u|U) \propto L(U|\theta_u) f(\theta_u) \quad (4)$$

เมื่อ $L(U|\theta_u)$ คือ ฟังก์ชัน likelihood ซึ่งสามารถแปลความหมายได้เป็น ฟังก์ชันภายหลัง \propto ความเป็นไปได้ \times ฟังก์ชันเดิม

$$\text{posterior} \propto \text{likelihood} \times \text{prior} \quad (5)$$

เมื่อจำนวนผู้สอบมี N คน ความหนาแน่นทั้งก่อนและหลังจะเป็นความหนาแน่นร่วมของ $\theta_1, \theta_2, \dots, \theta_n$ ดังนั้น สมการที่ (5) จะขยายออกเป็น

$$f(\theta_1, \theta_2, \dots, \theta_n \mid U_1, \dots, U_n) \propto L(U_1, U_2, \dots, U_n \mid \theta_1, \theta_2, \dots, \theta_n) f(\theta_1, \theta_2, \dots, \theta_n) \quad (6)$$

ฟังก์ชัน likelihood จะได้แก่

$$L(U \mid \theta) = L(U_1, U_2, \dots, U_n \mid \theta_1, \dots, \theta_n) = \prod_{a=1}^N L(U_a \mid \theta_a)$$

$$= \prod_{a=1}^N \prod_{i=1}^n L(U_{ia} \mid \theta) = \prod_{a=1}^N \prod_{i=1}^n P_{ia}^{U_{ia}} (1 - P_{ia})^{1 - U_{ia}} \quad (7)$$

เมื่อได้ค่า $f(\theta_a)$ ออกมาเป็นค่าคงที่แล้วสมการ (4) จะลดลงเป็น

$$f(\theta_a \mid U) \propto L(U \mid \theta_a)$$

ซึ่งได้แก่ ฟังก์ชัน likelihood ในกรณีนี้การประมาณค่าของเบย์ส์จะได้ค่าเท่ากับค่าประมาณค่าด้วย maximum likelihood ความวิธของเบย์ส์แล้ว อาจเริ่มโดยการระบุว่า การกระจายเดิมของ θ เป็นการกระจายปกติมีค่าเฉลี่ยเป็น 0 และความแปรปรวนเป็น 1

$$\theta_a \sim N(0, 1) \quad \text{หรือ}$$

$$f(\theta_a) \propto \exp(-1/2 \theta_a^2) \quad (8)$$

ถ้าเรากำหนดให้ความสามารถในการแจกแจงอย่างเป็นอิสระ ซึ่งเป็นข้อตกลงที่มีเหตุผลเป็นไปได้ แล้วการแจกแจงภายหลังจะเป็นดังนี้

$$f(\theta_1, \theta_2, \dots, \theta_n | U) \propto L(U | \theta_1, \dots, \theta_n) f(\theta_1, \dots, \theta_n) \quad (9)$$

และเนื่องจาก

$$f(\theta_1, \dots, \theta_n) = f(\theta_1) f(\theta_2), \dots, f(\theta_n)$$

$$\begin{aligned} & N \\ & \propto \prod_{a=1}^N \exp((-1/2) \theta_a) \\ & a=1 \end{aligned}$$

$$= \exp((-1/2) \theta_a) \quad (10)$$

ดังนั้น การกระจายภายหลังจึงเป็น

$$f(\theta_1, \theta_2, \dots, \theta_n | U) \propto L(U | \theta_1, \theta_2, \dots, \theta_n) [\exp(-1/2 \theta_a)] \quad (11)$$

จากนั้นประมาณค่าสูงสุดของค่าพารามิเตอร์ที่สอดคล้องกับค่าสูงสุดของฟังก์ชันความหนาแน่นภายหลัง โดยการใช้ค่าล็อกการิทึมเข้าช่วย

$$\ln f(\theta | U) = \text{ค่าคงที่} + \ln L(U | \theta) - \sum_{a=1}^N \theta_a \quad (12)$$

ผลของการแก้สมการ

$$\frac{\partial}{\partial \theta_a} \ln f(\theta | U) = 0 \quad a=1, \dots, n \quad (13)$$



ได้แก่ตัวประมาณค่าฐานนิยมของ e_1, \dots, e_n ของเบย์ส์ สมการ (13) สามารถเขียนได้ในรูป

$$\sum_{i=1}^n k_i (U_{i,n} - P_{i,n}) - e_n = 0 \quad (14)$$

$$\text{เมื่อ } k_i = D_{n,i} (P_{i,n} - c_i) / P_{i,n} (1 - c_i) \quad (15)$$

สำหรับโมเดล 3 พารามิเตอร์

สมการที่ 14 เป็นตัวอธิบายถึงความแตกต่างประการนระหว่างวิธีประมาณค่า MLE และ Bayesian และสามารถประมาณค่า e_n สำหรับผู้ที่ทำข้อสอบผิดหมดหรือถูกหมด ได้จากสมการ

$$\sum_{i=1}^n k_i P_{i,n} = \sum_{i=1}^n k_i U_{i,n} - e_n \quad (16)$$

$$\sum_{i=1}^n k_i P_{i,n} = -e_n \quad (17)$$

และเนื่องจากค่าเฉลี่ยของ $e = 0$ e_n ที่สอดคล้องกับคะแนนรวม = 0 จึงมีค่าเป็นลบ เมื่อพิจารณาสมการ (17) และการพิจารณาค่า e สำหรับผู้ได้คะแนนเต็มก็มีวิธีการเช่นเดียวกัน สมการฐานนิยมซึ่งได้แก่ สมการที่ (14) แก้ได้โดยใช้วิธีการของ Newton-Raphson เช่นเดียวกับที่ใช้ใน MLE

3. วิธี Marginal Maximum Likelihood (MML)

การประมาณค่าพารามิเตอร์ด้วย MML สามารถใช้ได้กับแบบสอบทั้งฮาวและสั้น และเริ่มจากรูปแบบการตอบข้อสอบของผู้สอบ

$$X = (x_1, x_2, \dots, x_n)$$

และสำหรับผู้ที่มีความสามารถ e ค่าความเป็นไปได้แสดงได้ดังนี้

$$P(X|\theta) = \prod_{j=1}^n [P_j(\theta)]^{x_j} [1-P_j(\theta)]^{1-x_j}$$

ค่าความน่าจะเป็นที่ได้เป็นความน่าจะเป็นที่มีเงื่อนไขว่าทราบค่าความสามารถ = θ ซึ่งต่างจากความน่าจะเป็นของรูปแบบ x จากผู้ไม่ทราบค่า θ ซึ่งสุ่มมาจากประชากรที่มีการกระจายของค่า θ อยู่ในรูปฟังก์ชันความหนาแน่นต่อเนื่อง (continuous density $g(\theta)$) อยู่ในรูปของความน่าจะเป็นที่ไม่มีเงื่อนไข และแสดงได้ในรูปของการอินทิเกรต

$$P(x) = \int_{-\infty}^{+\infty} P(x|\theta) g(\theta) d\theta$$

ค่านี้เรียกว่าค่าความน่าจะเป็นโดยปราศจากเงื่อนไขของ X และเนื่องจากค่า θ ถูก integrate ออกไป ค่านี้จึงเป็นฟังก์ชันของพารามิเตอร์ข้อสอบเท่านั้น และสามารถหาค่าได้จากสูตร gaussian quadrature ดังนี้

$$\bar{P}(x) = \sum_{k=1}^q p(x|x_k) A(X_k)$$

เมื่อ X คือ quadrature point

$A(x_k)$ คือ น้ำหนักทางบวกซึ่งสอดคล้องกับฟังก์ชัน $g(x)$

ซึ่งหาค่าได้จากตารางใน Stroud และ

Sechrest, 1966 (อ้างถึงใน Mislevy และ

Bock, 1990 : 1- 7)

และในวิธีของ MML ค่าพารามิเตอร์ของข้อสอบจะพิจารณาใช้ค่าที่ทำให้สมการข้างล่างมีค่าสูงสุด

$$\log L_M = \sum_{l=1}^s r_l \log \bar{P}(X_l)$$

เมื่อ $\bar{P}(X_l)$ ประมาณได้ดังสูตร gaussian quadrature r_l คือ ความถี่ของรูปแบบการตอบข้อสอบรูปแบบ X_l จากจำนวนผู้เข้าสอบทั้งหมด (N) และ S เป็นจำนวนรูปแบบที่เป็นไปได้ทั้งหมดจากข้อมูลการตอบข้อสอบ

กรณี 3 พารามิเตอร์ ค่าสูงสุดของสมการข้างบนสำหรับข้อสอบข้อ j หาได้จากสมการ likelihood

$$\sum_{k=1}^q \frac{\bar{r}_{jk} - \bar{N}_k P_j(X_k)}{P_j(X_k) [1 - P_j(X_k)]} \cdot \frac{\partial P_j(X_k)}{\partial \begin{bmatrix} c_j \\ a_j \\ b_j \end{bmatrix}} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\text{เมื่อ } \bar{r}_{jk} = \sum_{l=1}^s r_{lj} P(X_l | X_k) A(X_k) / \bar{P}(X_k)$$

$$\bar{N}_k = \sum_{l=1}^s r_{lj} P(X_l | X_k) A(X_k) / \bar{P}(X_k)$$

เป็นค่าคาดหวังภายหลังของจำนวนข้อสอบที่ถูกและจำนวนข้อสอบที่ทำที่จุด X_k (x_{1j} คือ คะแนน 0, 1 สำหรับข้อสอบ j ในรูปแบบการตอบ l)

จากนั้นใช้ขั้นตอนการคำนวณ EM และวิธีของ Newton-Gauss (= Fisher-scoring) แก้สมการ

2.5 การตรวจสอบความลำเอียงของข้อสอบตามทฤษฎีการตอบข้อสอบ

นิยามของความลำเอียงของข้อสอบในแง่ของทฤษฎีการตอบข้อสอบก็คือ ข้อสอบจะลำเอียงถ้าฟังก์ชันการตอบข้อสอบข้อนั้น ๆ ของประชากรย่อย 2 กลุ่ม ไม่เป็นเส้นเดียวกัน

จากนิยามดังกล่าวการตรวจสอบความลำเอียงของข้อสอบตามแนวความคิดของทฤษฎีการตอบข้อสอบสามารถทำได้โดยเปรียบเทียบฟังก์ชันการตอบข้อสอบของผู้สอบกลุ่มย่อยที่ต้องการศึกษา และการเปรียบเทียบอาจกระทำได้นี้ (Hulin Drasgow และ Parsons, 1983)

1) ทฤษฎีการตอบข้อสอบโมเดล 1 พารามิเตอร์ ซึ่งกำหนดให้ค่าการเดา = 0 และค่าอำนาจจำแนกของทุกข้อ = 1.00 มีค่าความยาก (b) แปรผันไปตามกลุ่ม เริ่มด้วยการทดสอบความเหมาะสมของแบบจำลองกับข้อมูลทั้ง 2 กลุ่ม เป็นรายข้อด้วย H_1 ซึ่งมีการกระจายแบบ X^2 มี $df = N-1$ ตัดข้อที่ไม่เหมาะสมออก แล้ววิเคราะห์ความล่าเอียงรายข้อด้วย Z_1

$$H_1 = \sum_{j=1}^N \frac{[U_{i,j} - \hat{P}_1(\theta_j)]^2}{[\hat{P}_1(\theta_j) \hat{Q}_1(\theta_j)]}$$

$U_{i,j}$ = คะแนนจากการตอบข้อสอบ แบบ 0-1 ของคนที่ j ในข้อที่ i

$$P_1(\theta) = \frac{1}{1 + \exp[-D(\theta - b_1)]}$$

$$Z_1 = \frac{\hat{b}_{1A} - \hat{b}_{1B}}{SE^2 b_{1A} + SE^2 b_{1B}}$$

$$SE^2_{b_1} = \sqrt{I_{a1} / (I_{a1} I_{b1} - I_{ab1}^2)}$$

ค่า Z_1 สูงแสดงว่า ข้อสอบนั้นยากสำหรับกลุ่มหนึ่งมากกว่าอีกกลุ่มหนึ่ง

หรืออาจทดสอบด้วยค่า t ซึ่งคำนวณได้จากสูตร

$$t = \frac{\hat{b}_{1A} - \hat{b}_{1B}}{SE^2_{b_{1A}} + SE^2_{b_{1B}}}$$

ในการนี้ Draba (Cited in Berk, 1982 : 138) เสนอให้ใช้ค่า $t > 2.4$ เป็นเกณฑ์การตัดสินความล่าช้าของข้อสอบ

2) ทฤษฎีการตอบข้อสอบโมเดล 2 พารามิเตอร์ กำหนดค่าการเดา = 0 ส่วนค่าความยากและค่าอำนาจจำแนกผันแปรไปตามกลุ่มย่อย วิธีการเริ่มด้วยการประมาณค่าพารามิเตอร์แยกตามกลุ่มย่อย A และ B ด้วยวิธี maximum likelihood แล้วเปรียบเทียบโค้งการตอบข้อสอบ ด้วยค่า F

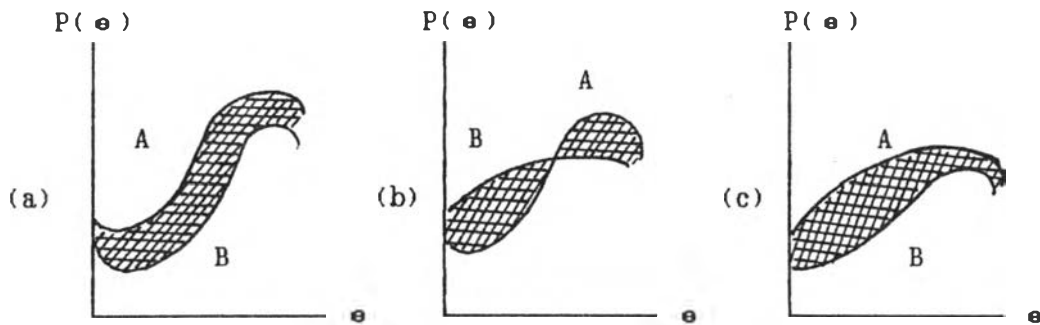
$$F = \frac{SSE(\text{pooled}) - [SSE(A) - SSE(B)]}{SSE(A) + SSE(B)} \cdot \frac{[J_A - J_B - 4]}{2}$$

ค่าขึ้นแห่งความเป็นอิสระ คือ 2, $(J_A + J_B - 4)$ เมื่อ J เป็นจำนวนช่วงความสามารถ การทดสอบค่า F จะบอกว่าเส้นถดถอยรวมสามารถอธิบายความแปรปรวนของสัดส่วนการตอบของเส้นถดถอยจากแต่ละกลุ่มได้เท่ากันหรือไม่ ถ้าเท่ากัน (ไม่มีนัยสำคัญ) ก็ไม่จำเป็นต้องใช้เส้นถดถอยแยกกันเพราะความสัมพันธ์ระหว่างความสามารถ และการตอบข้อสอบของทั้งสองกลุ่มเป็นเช่นเดียวกัน และข้อสอบทำหน้าที่ในการวัดความสามารถของคน 2 กลุ่ม นี้ได้เหมือนกัน

3) ทฤษฎีการตอบข้อสอบโมเดล 3 พารามิเตอร์ เป็นแบบที่ปล่อยให้ค่าการเดาค่าความยาก และค่าอำนาจจำแนกของข้อสอบแปรเปลี่ยนไปตามกลุ่ม

วิธีนี้เริ่มจากการประมาณค่าพารามิเตอร์รวมของทุกคนทุกกลุ่ม แล้วแปลงเป็นค่ามาตรฐานบนค่าประมาณความยาก จากนั้นแยกประมาณค่าความยากและอำนาจจำแนกตามกลุ่มย่อย โดยใช้ค่าการเดาที่ประมาณได้จากการรวมทุกกลุ่ม การทำค่าความยากให้เป็นมาตรฐานต้องตัดข้อที่มีค่าความยากสูงมากหรือต่ำมากออก

การวิเคราะห์ความล่าช้าของข้อสอบตามโมเดล 3 พารามิเตอร์นั้น Linn, Levine, Hastings และ Wardrop, (1981 cited in Hulin, Drasgows และ Parsons, 1983: 176) ได้พัฒนาการทดสอบความล่าช้าจากความแตกต่างระหว่างโค้งลักษณะข้อสอบโดยตรง โดยพิจารณาจากพื้นที่ระหว่างโค้ง 2 โค้งเป็นตัวแสดงความล่าช้า โดยมีลักษณะของพื้นที่อยู่ 3 ลักษณะ ดังภาพที่ 7



ภาพที่ 7: โด็งการตอบข้อสอบที่มีความลำเอียง 3 ลักษณะสำหรับกลุ่ม A และ B

รูป (a) แสดงถึงความลำเอียงที่คงที่ (consistent bias) ข้อสอบมีความลำเอียงต่อกลุ่ม B อย่างสม่ำเสมอตลอดความต่อเนื่องของความสามารถเป็นความลำเอียงอย่างสม่ำเสมอ (uniform bias)

รูป (b) แสดงว่า ข้อสอบมีความยากคล้ายคลึงกันสำหรับทั้ง 2 กลุ่ม แต่มีอำนาจจำแนกสำหรับกลุ่ม B น้อยกว่ากลุ่ม A ทำให้มีความลำเอียงต่อกลุ่ม B ในพวกที่มีความสามารถสูงและลำเอียงต่อกลุ่ม A ที่มีความสามารถต่ำ

รูป (c) มีจุดตัดของโด็งและแกนความเป็นไปได้ในการตอบถูกไม่เท่ากัน สำหรับ 2 กลุ่ม แต่ความยากและอำนาจจำแนกใกล้เคียงกัน มีความลำเอียงต่อกลุ่ม B ในผู้ที่มีความสามารถต่ำ แต่จะลดลงในผู้ที่มีความสามารถสูงขึ้น ลักษณะความลำเอียงเช่นนี้อาจเป็นเพราะสมาชิกในกลุ่ม A ทำข้อสอบถูกเพราะการเดาก็เป็นได้

Shepard และคณะ (1985 : 81) กล่าวว่า เมื่อพบว่าโด็งลักษณะข้อสอบของกลุ่มอ้างอิงและกลุ่มเป้าหมายต่างกันมากพอที่จะเป็นหลักฐานว่า ข้อสอบที่ศึกษาอยู่ไม่ได้มีลักษณะแฝงเดียวกันสำหรับกลุ่มทั้งสองแล้ว จะใช้ดัชนีชี้ถึงระดับของความลำเอียง ดังนี้

1. พื้นที่ ชนิดไม่มีเครื่องหมาย (unsigned Area) เป็นค่าสัมบูรณ์ของพื้นที่ระหว่างโด็งทั้งสอง
2. ผลบวกกำลังสองแบบที่ 1 (SOS1) เป็นค่าผลบวกของความแตกต่างกำลังสองทุกค่าของความสามารถที่เป็นไปได้
3. ผลบวกกำลังสองแบบที่ 2 (SOS2) คล้าย SOS1 แต่ให้น้ำหนักความแตกต่างกำลังสองด้วยส่วนกลับของความคลาดเคลื่อนของความแปรปรวน ในการประมาณค่าความเป็นไปได้ ดังนั้นถ้ามีการประมาณโด็งอย่างไม่ดีในบริเวณใดส่วนหนึ่งแล้ว จะทำให้ลดประสิทธิผลต่อความแตกต่างระหว่างกลุ่มสอง

4. ไคสแควร์ ($IRT-\chi^2$) เป็นการทดสอบนัยสำคัญโดยการเปรียบเทียบความแตกต่างของค่าความชาก และอำนาจจำแนกระหว่างกลุ่มหรือม ๆ กัน โดย Lord (1980)

5. พื้นที่ชนิดมีเครื่องหมาย เหมือนพื้นที่ชนิดไม่มีเครื่องหมาย แต่มีการให้เครื่องหมายไว้เพื่อแสดงให้เห็นว่า กลุ่มใดเสียประโยชน์ และกลุ่มใดได้ประโยชน์ ถ้าต้องการตอบข้อสอบตัดกัน แสดงว่ามีความลำเอียงต่างกัน สำหรับผู้มีความสามารถต่างกัน และค่าสัมบูรณ์ของพื้นที่ชนิดมีเครื่องหมาย จะน้อยกว่าชนิดไม่มีเครื่องหมาย

6. ผลบวกกำลังสอง แบบที่ 3 (SOS3) เป็นผลบวกกำลังสองของความแตกต่างชนิดมีเครื่องหมายของ SOS1 โดยการใช้เครื่องหมายของความแตกต่างแต่ละตัวที่ยังไม่ยกกำลังสอง

7. ผลบวกกำลังสอง แบบที่ 4 (SOS4) เป็นดัชนีคู่ขนานกับ SOS2 แต่เป็นผลบวกที่มีการให้น้ำหนัก และให้เครื่องหมาย

การวัดพื้นที่ของความแตกต่างระหว่างโค้งลักษณะข้อสอบของกลุ่มตัวอย่าง เพื่อพิจารณาถึงขนาดของความลำเอียงของข้อสอบนั้นมีอยู่ 2 ลักษณะ คือ การวัดช่วงเปิด และช่วงปิด (open and closed intervals) กรณีของการวัดช่วงปิดนั้น จะกำหนดขอบเขตของการวัดให้อยู่ในช่วงความสามารถที่กำหนด ในขณะที่การวัดช่วงเปิดนั้นวัดภายในช่วงความสามารถทั้งหมด เพื่อให้ได้พื้นที่ที่แน่นอน (exact area) ระหว่างโค้ง 2 โค้ง

ภายใต้การวัดแบบช่วงปิดนั้น สามารถวัดพื้นที่ได้ 2 ลักษณะ คือ พื้นที่ชนิดมีเครื่องหมายและไม่มีเครื่องหมาย (signed and unsigned areas) โดยมีสูตรการคำนวณทั่ว ๆ ไป ดังนี้

จากฟังก์ชันการตอบข้อสอบ แบบ 3 พารามิเตอร์

$$P(\theta) = c + (1-c) P^*(\theta)$$

$$\text{โดยที่ } P^*(\theta) = \{1 + \exp[Da(\theta - b)]\}^{-1}$$

a, b และ c เป็นพารามิเตอร์ที่แสดงลักษณะของข้อสอบ และ D เป็นตัวคงที่ปกติกำหนดให้ = 1.7

พื้นที่ระหว่างความสามารถที่กำหนดช่วงไว้ 2 จุด บนมาตรฐานความสามารถได้จาก

$$S(\theta_1, \theta_2) = \int_{\theta_1}^{\theta_2} P(\theta) d\theta = c(\theta_2 - \theta_1) + (1-c)(Da)^{-1} \ln \frac{\{1 - \exp[Da(\theta_2 - b)]\}}{\{1 + \exp[Da(\theta_1 - b)]\}}$$

ในการศึกษาความลำเอียงของข้อสอบนั้น จะมีพารามิเตอร์ของ 2 กลุ่ม คือ กลุ่มอ้างอิง ได้แก่ a_R , b_R , และ c_R และกลุ่มเปรียบเทียบ ได้แก่ a_F , b_F , และ c_F พื้นที่ชนิดมีเครื่องหมาย (Closed-interval Signed Area-CSA) ได้จาก

$$CSA = \int_{e_1}^{e_2} [P_R(e) - P_F(e)] de = S_R(e_1, e_2) - S_F(e_1, e_2)$$

และพื้นที่ชนิดไม่มีเครื่องหมาย (Closed-interval Unsigned Area-CUA) หรือ Φ ได้จาก

$$CUA = \int_{e_1}^{e_2} |P_R(e) - P_F(e)| de$$

Laksana และ Coffman (1980: 11) เสนอให้ใช้ค่า $CUA > .20$ เป็นเกณฑ์ของความลำเอียงของข้อสอบ

การวิเคราะห์ความลำเอียงของข้อสอบด้วยทฤษฎีการตอบข้อสอบนั้น Subkoviak, Mack, Ironson และ Craig (1984) พบว่าจากการเปรียบเทียบกับวิธีอื่น ๆ ได้แก่ ค่าความยากแปลง และโคสแควร์แบบ 5 ช่วง พบว่าการวิเคราะห์ด้วยทฤษฎีการตอบข้อสอบแบบ 3 พารามิเตอร์ ให้ความถูกต้องมากที่สุด Shepard Camilli และ Williams (1985) ค้นพบสอดคล้องกันในเรื่องที่ว่าวิธีที่ดีที่สุดกว่าวิธี χ^2 และค่าความยากแปลงแต่ 3 คนหลังนี้ ใช้ทฤษฎีการตอบข้อสอบเทียมนำไปประยุกต์ใช้ Baghi และ Ferrara (1990) พบว่าในขนาดตัวอย่าง 750 คน สามารถใช้ MH แทน IRT ชนิด 3 พารามิเตอร์ ได้ Kim และ Kohen (1991) พบว่าการเปรียบเทียบความแตกต่างของพื้นที่ใต้โค้งด้วยดัชนีไม่มีเครื่องหมายดีกว่าแบบมีเครื่องหมายรวม ทั้งมีความคลาดเคลื่อนต่ำกว่าแบบมีเครื่องหมาย และระบุว่าการลำเอียงโดยที่ไม่ลำเอียงจริง น้อยกว่าวิธีอื่น ๆ คิดเป็นร้อยละ 5 ชื่อ Lautenschlager และ Park (1988) พบว่าวิธีของ Linn (1980) ดีกว่าวิธีของ Warm (1978) Rudner Getson และ Knight 1980) พบว่าวิธีการตรวจสอบความลำเอียงด้วยทฤษฎีการตอบข้อสอบโมเดล 3 พารามิเตอร์ ไม่มีผลกระทบจากความยาวของข้อสอบ

การใช้โปรแกรมการถดถอยแบบโลจิสต์ โดย Smaminathan และ Rogers (1990) จากข้อมูลจำลอง พบว่า สามารถตรวจค้นความลำเอียงของข้อสอบได้ร้อยละ 50-75 แล้วแต่ความยาวของข้อสอบ และขนาดของกลุ่มตัวอย่าง ในกรณีที่เป็นการลำเอียงอย่างไม่

สม่ำเสมอ แต่สามารถตรวจค้นความล่าเอียงอย่างสม่ำเสมอได้ใกล้เคียงกับวิธี MH แต่ค่าใช้จ่ายมากกว่าถึง 3-4 เท่า Rogers (1989) ได้ใช้การทดลองแบบโลจิสต์เพื่อตรวจค้นความล่าเอียงของข้อสอบแบบไม่สม่ำเสมอ พบว่า ใช้ได้ดีกับกลุ่มตัวอย่างขนาด 250 คนขึ้นไป และจะตรวจค้นความล่าเอียงอย่างสม่ำเสมอได้เท่ากับ MH แต่ตรวจค้นได้ดีกว่าในกรณีของความล่าเอียงแบบไม่สม่ำเสมอ

สำหรับการวิจัยครั้งนี้ ผู้วิจัยได้เลือกศึกษา วิธี IRT-2 พารามิเตอร์ โดยเลือกพัฒนาเกณฑ์ของ พื้นที่ความแตกต่างแบบเปิด (open interval) เพื่อให้ได้พื้นที่แน่นอน (exact area) โดยใช้ทั้งพื้นที่ชนิดที่มีเครื่องหมาย (SA) และพื้นที่ชนิดไม่มีเครื่องหมาย (UA)

3. วิธีตรวจสอบความล่าเอียงของข้อสอบด้วยวิธีของแมนเทลและแฮนส์เชล

เทคนิคแมนเทล-แฮนส์เชล (MH) เป็นเทคนิคที่ Mantel และ Haenszel ได้เสนอขึ้นใช้ตั้งแต่ปี 1959 แต่ Holland และ Thayer เพิ่งนำเสนอมาใช้เพื่อการตรวจสอบความล่าเอียงของข้อสอบในปี 1986 และในปีนั้นก็ได้มีการตรวจสอบความล่าเอียงของข้อสอบด้วยเทคนิค MH โดย Dorans Hambleton Rogers และ Arrasmith Holland และ Thayer, Wright Raju Bode และ Larsen (1989 : 2)

เทคนิค MH เป็นเทคนิคที่มีความคล้ายคลึงกับวิธีไคสแควร์ที่เสนอโดย Scheuneman (1979) Marascuilo และ Slaughter (1981) และ Mellenberg (1982) โดยเป็นการเปรียบเทียบผลการตอบข้อสอบของกลุ่มผู้สอบ 2 กลุ่ม หรือที่ฮอลแลนด์เรียกว่ากลุ่มอ้างอิงและกลุ่มเปรียบเทียบ (reference and focal group) โดยที่กลุ่มแรกจะใช้อ้างอิงถึงกลุ่มที่คาดว่าจะได้ประโยชน์จากข้อสอบและกลุ่มหลังเป็นกลุ่มเสียประโยชน์จากข้อสอบในกรณีที่ข้อสอบมีความล่าเอียง โดยมีการตรวจสอบทุก ๆ ระดับคะแนนรวมจากการสอบ ข้อสอบใดที่ผู้สอบทั้ง 2 กลุ่ม ทำได้เท่า ๆ กัน จะเป็นข้อสอบที่ถือได้ว่าไม่มีความล่าเอียงต่อกลุ่มใดกลุ่มหนึ่ง วิธีการในการตรวจสอบมีดังนี้

ถ้าให้ $N_{r,j}$ และ $N_{f,j}$ แทนจำนวนผู้สอบในกลุ่มอ้างอิงและกลุ่มเปรียบเทียบที่มีคะแนนรวมในช่วงคะแนน j และให้ N_j ซึ่ง $= N_{r,j} + N_{f,j}$ เป็นจำนวนผู้สอบรวมทั้ง 2 กลุ่ม ที่ได้คะแนนรวมจากการสอบอยู่ในช่วงคะแนน j แล้วนำมาเขียนในรูปตาราง 2 ทาง แสดงผลการตอบถูก (1) และผิด (0) ดังตารางที่ 1

ตารางที่ 1 ข้อมูลผลการสอบของกลุ่มตัวอย่างกลุ่ม R และ F ที่มีคะแนนรวมอยู่ในช่วงคะแนน j

		คะแนนที่ได้จากข้อสอบที่ต้องการตรวจสอบความล่าเอียง		
		1	0	รวม
กลุ่มผู้สอบ	R	A_{j}	B_{j}	$N_{R,j}$
	F	$C_{F,j}$	$D_{F,j}$	$N_{F,j}$
	รวม	$N_{1,j}$	$N_{0,j}$	N_{j}

เมื่อ A_{j} , B_{j} , C_{j} และ D_{j} เป็นความถี่ที่สังเกตได้ของการตอบถูก (1) และผิด (0) ในช่วงคะแนน j และแสดงในรูปสัดส่วนได้ ดังตารางที่ 2

ตารางที่ 2 สัดส่วนการตอบข้อสอบของกลุ่มประชากร 2 กลุ่ม

กลุ่ม	ผลการตอบ	1	0	รวม
	R	$P_{R,j}$	$Q_{R,j}$	1
F	$P_{F,j}$	$Q_{F,j}$	1	

โดยที่ $P_{R,j}$ คือ สัดส่วนของกลุ่มอ้างอิงที่อยู่ในช่วงความสามารถ j ที่ตอบข้อสอบถูก

$P_{F,j}$ คือ สัดส่วนของกลุ่มสนใจที่อยู่ในช่วงความสามารถ j ที่ตอบข้อสอบถูก

$Q_{R,j}$ คือ $1 - P_{R,j}$

$Q_{F,j}$ คือ $1 - P_{F,j}$



จากตารางที่ 1 จะได้ว่า ในช่วงระดับคะแนน j A_{j1} และ B_{j1} คือ จำนวนผู้สอบกลุ่ม
อ้างอิงที่ตอบข้อสอบข้อที่ถูกต้อง และไม่ถูกต้องตามลำดับ ในขณะที่ C_{j1} และ D_{j1} คือ จำนวนผู้สอบ
กลุ่มเปรียบเทียบที่ตอบข้อสอบข้อที่ถูกต้องและไม่ถูกต้องตามลำดับเช่นกัน

ตามเทคนิควิธี MH ที่เสนอโดย Holland และ Thayer นั้น ข้อมูลการตอบข้อสอบ
แต่ละข้อจะถูกนำมาแสดงในรูปตาราง 2 ทางนี้ เป็นจำนวนตารางเท่ากับคะแนนรวมที่มีผู้สอบทำ
ได้บวกด้วย 1 จากนั้นจะคำนวณค่าความเป็นไปได้ในรูปของค่าสัดส่วนการตอบข้อสอบถูก-ผิด
ระหว่างกลุ่มโดยกำหนดให้เป็นค่า α ดังนี้

$$\alpha_{MH} = \frac{\sum A_{j1} D_{j1} / N_{j1}}{\sum B_{j1} C_{j1} / N_{j1}}$$

ค่า α_{MH} จะมีค่าระหว่าง 0 และ 1 โดยที่ $\alpha_{MH} = 1$ จะแสดงถึงสมมุติฐานศูนย์
ที่ว่าไม่มีความแตกต่างหรือไม่มีความลำเอียง ซึ่งแสดงว่ากลุ่มตัวอย่างทั้ง 2 กลุ่ม ทำข้อสอบได้
ผลเช่นเดียวกันในแต่ละระดับคะแนนรวม ค่า $A_{j1} D_{j1} / N_{j1}$ และ $B_{j1} C_{j1} / N_{j1}$ มีค่าเท่ากันทั้ง

2 กลุ่มทำให้ $\alpha_{MH} = 1$

อย่างไรก็ตามในทางปฏิบัติแล้ว แม้ว่าจะมีการควบคุมให้ผู้สอบ 2 กลุ่ม มีความสามารถ
ในสิ่งที่ต้องการวัดเท่ากัน แต่เป็นไปได้ที่ผลการตอบข้อสอบของผู้สอบทั้ง 2 กลุ่ม จะเป็น
เช่นเดียวกันดังกล่าว Mantel และ Haenszel (1959) จึงเสนอค่าสถิติไคสแควร์เพื่อทดสอบค่า
ที่ได้ว่าจะมีความแตกต่างจาก 1.00 อย่างมีนัยสำคัญหรือไม่ที่ระดับชั้นแห่งความเป็นอิสระ = 1
ดังนี้

$$MH-CHISQ = \frac{\sum A_{j1} - E(A_{j1}) - 0.5)^2}{\sum \text{Var}(A_{j1})}$$

โดยที่ค่าผลรวมเป็นค่าที่ได้จากการรวมทุกชั้นของคะแนนรวม และค่า $E(A_{j1})$ มีค่า
ดังนี้

$$E(A_{j1}) = (N_{r1}) (N_{1j}) / N_{11}$$

$$\text{Var}(A_{.j}) = \frac{N_{R.j}N_{E.j}N_{I.j}N_{O.j}}{N_{.j}(N_{.j}-1)}$$

และสมมติฐานศูนย์แสดงว่าข้อสอบไม่ลำเอียง ได้แก่

$$H_0 : P_{R.j} = P_{F.j} \quad \text{สำหรับทุกชั้นคะแนน } j$$

สมมติฐานศูนย์นี้เป็นสมมติฐานของความเป็นอิสระอย่างมีเงื่อนไขของสมาชิกกลุ่ม และคะแนนที่ได้จากการตอบข้อสอบที่ต้องการตรวจสอบความลำเอียง และภายใต้สมมติฐานศูนย์จะได้ค่าคาดหวังของการตอบข้อสอบ ดังนี้

$$E(A_{.j}) = n_{R.j} m_{1.j} / T_{.j}$$

$$E(B_{.j}) = n_{R.j} m_{0.j} / T_{.j}$$

$$E(C_{.j}) = n_{F.j} m_{1.j} / T_{.j}$$

$$E(D_{.j}) = n_{F.j} m_{0.j} / T_{.j}$$

สมมติฐานอื่นของแมนเทลและแฮนส์เชล ได้แก่

$$H_1 : \frac{p_{R.j}}{q_{R.j}} = \alpha \frac{p_{F.j}}{q_{F.j}} \quad j = 1, \dots, K \quad \text{เมื่อ } \alpha \text{ มีค่า } \neq 1$$

เมื่อ α มีค่า = 1 ซึ่งสอดคล้องกับสมมติฐานศูนย์ จะได้ว่า

$$H_0 : \frac{p_{R.j}}{q_{R.j}} = \frac{p_{F.j}}{q_{F.j}}$$

และค่าประมาณของ $\alpha_{MH} = \frac{p_{R.j}}{q_{R.j}} \cdot \frac{p_{F.j}}{q_{F.j}} = \frac{p_{R.j}q_{F.j}}{p_{F.j}q_{R.j}}$ สำหรับทุก $j = 1, \dots, k$

นอกจากนั้น Holland และ Thayer ได้เสนอแนะให้แปลงค่า α_{MH} ให้เป็นเคลต้า (Δ) เพื่อให้เป็นค่าที่มีค่าเฉลี่ยเป็น 0 ดังนี้

$$\Delta_{MH} = \frac{-4}{1.7} \ln(\alpha_{MH}) = -2.35 \ln(\alpha_{MH})$$

ค่า Δ_{MH} มีค่าระหว่าง -2.6 ถึง 2.6 และใช้เป็นดัชนีแสดงความล่าเอียง ข้อสอบที่ไม่ล่าเอียง จะมีค่า $\alpha_{MH} = 1$ หรือ $\Delta_{MH} = 0$

เท่าที่ผ่านมา มีผู้นำวิธี MH มาศึกษา เช่น Clauser และคณะ (1991) ใช้ MH วิเคราะห์ความล่าเอียงอย่างสม่ำเสมอ พบว่า MH จะใช้ได้ดีในกรณีที่ข้อสอบมีค่าอำนาจจำแนกสูง แต่จะไม่สามารถวิเคราะห์ข้อสอบที่มีค่าความยากสูงมากได้ ซึ่งสอดคล้องกับข้อค้นพบของ Mazor และคณะ (1991) Sudweeks และ Tolman (1990) พบว่า ข้อสอบที่ล่าเอียงมักเป็นข้อสอบยาก Baghi และ Ferrara (1990) พบว่า เมื่อใช้ขนาดตัวอย่าง 750 คนขึ้นไป MH ใช้แทน IRT-3 ได้ Swaminathan และ Rogers (1990) ซึ่งศึกษาจากข้อมูลจำลองพบว่า MH วิเคราะห์ได้ดีกว่าการถดถอยแบบโลจิสต์เล็กน้อย โดยตรวจค้นได้ถูกต้อง ร้อยละ 75 ในกรณีใช้กลุ่มตัวอย่าง 250 คน และตรวจค้นได้ ร้อยละ 100 กรณีกลุ่มตัวอย่าง 500 คน กรณีที่ล่าเอียงอย่างสม่ำเสมอและกรณีไม่สม่ำเสมอที่ติดกันที่ปลายข้างใดข้างหนึ่ง แต่ MH มีค่าใช้จ่ายน้อยกว่าโลจิสต์ประมาณ 3-4 เท่า Hambleton และคณะ (1986) พบว่า MH ให้ค่าใกล้เคียงกับ IRT ทั้งที่ใช้ค่าความแตกต่างของค่าเฉลี่ยกำลังสองและการตรวจสอบความแตกต่างของพื้นที่รวมได้โด่ง แต่ MH มีค่าใช้จ่ายต่ำกว่าและใช้เวลาน้อยกว่า Lincre (1986) พบว่า MH ใช้ได้ดีในสถานการณ์จำลองทุกสถานการณ์พอ ๆ กับโปรแกรม PROX ของ Rasch แต่ควรใช้เกณฑ์ในการคำนวณความล่าเอียงและความคลาดเคลื่อนหลาย ๆ เกณฑ์ และใช้ตัวประมาณค่าความคลาดเคลื่อนมาตรฐานมากกว่า 1 ตัว Perlman และคณะ (1988) พบว่า MH มีปัญหาด้านความเที่ยง เมื่อจำนวนกลุ่มตัวอย่างน้อยกว่า 660 คน Thissen Steinberg และ Wainer (1988) พบว่า MH ให้ผลการวิเคราะห์คล้าย IRT-LR และอาจใช้วิเคราะห์ความล่าเอียงก่อนใช้ IRT-LR

ในการศึกษาครั้งนี้ ผู้วิจัยได้นำดัชนี α_{MH} มาใช้เพื่อพัฒนาเกณฑ์ตัดสินใจข้อสอบล่าเอียง

4. วิธีตรวจสอบความลำเอียงของข้อสอบตามการทดสอบ SIB (SIBTEST)

การทดสอบ SIB เป็นวิธีการที่มีแนวคิดบนการตรวจสอบความลำเอียงของแบบสอบชนิดพหุมิติ (multidimensional) โดยมีพื้นฐานบนทฤษฎีการตอบข้อสอบแบบพหุมิติ แต่การทดสอบ SIB แตกต่างจาก IRT ตรงที่เป็นการใช้การทดสอบค่าสถิติไร้พารามิเตอร์ (non-parametric) และมีการคำนวณที่ง่ายไม่ซับซ้อน ใช้ได้สำหรับการตรวจสอบความลำเอียงอย่างสม่ำเสมอและมีทิศทางเดียว (unidirectional)

การทดสอบ SIB เป็นวิธีการตรวจสอบความลำเอียงของข้อสอบที่เสนอโดย Shealy และ Stout (Shealy, 1989; Shealy และ Stout, 1991 cited in Shealy และ Stout, 1992: 2) และรูปแบบอยู่ที่การวิเคราะห์ข้อมูลที่มีข้อตกลงว่ามีมิติการวัด 2 มิติ มิติหนึ่งเป็นความสามารถหรือลักษณะแฝงเป้าหมายที่ต้องการวัด และอีกมิติเป็นลักษณะแฝงแทรกซ้อนที่ไม่ต้องการวัด

Shealy และ Stout (1992) เริ่มจากกลุ่มผู้สอบ 2 กลุ่ม คือ กลุ่มอ้างอิงและกลุ่มเปรียบเทียบ การตอบข้อสอบของผู้สอบที่สุ่มมาจากกลุ่มผู้สอบคือ U โดยที่ $U = (U_1, \dots, U_n)$ เมื่อ U_i มีค่าเป็น 0 ถ้าตอบข้อสอบข้อ i ผิดและเป็น 1 เมื่อตอบถูก

สำหรับโมเดล IRT แล้ว โดยทั่ว ๆ ไป U จะเกิดได้จาก 2 องค์ประกอบ คือ (1) พารามิเตอร์ความสามารถของผู้สอบจำนวน d มิติ และ (2) ฟังก์ชันการตอบข้อสอบแต่ละข้อ ซึ่งเป็นตัวกำหนดความน่าจะเป็นของการตอบข้อสอบได้ถูก

ในที่นี้ d คือ 2 เพราะเราจะพิจารณาความสามารถแทรกซ้อนเพิ่มเข้ามาจากความสามารถเป้าหมายที่ต้องการวัดเวกเตอร์ความสามารถในที่นี้คือ (θ, η) เมื่อ θ คือ ความสามารถเป้าหมายและ η คือ ความสามารถสอดแทรก ส่วนฟังก์ชันการตอบข้อสอบสำหรับข้อ i คือ $P_i(\theta, \eta)$ โดยที่ข้อสอบทุกข้อจะประกอบด้วย θ ในขณะที่บางข้อจะมี η อยู่ด้วย สำหรับข้อที่ประกอบด้วย θ อย่างเดียว ฟังก์ชันการตอบข้อสอบ คือ $P_i(\theta)$ ค่า $P_i(\theta, \eta)$ จะสูงขึ้นใน θ และใน η เมื่อข้อสอบข้อ i ต้องการความสามารถทั้ง 2 อย่าง และจะสูงขึ้นใน θ ถ้าต้องการเฉพาะความสามารถ θ

นิยามของความลำเอียงของข้อสอบในการทดสอบ SIB คือ "ความลำเอียงของข้อสอบจะเกิดขึ้นต่อกลุ่มเปรียบเทียบที่ระดับความสามารถ θ ถ้าฟังก์ชันการตอบข้อสอบสำหรับความสามารถเป้าหมาย (target ability) ของกลุ่มอ้างอิงมีค่ามากกว่าฟังก์ชันการตอบข้อสอบของกลุ่มเปรียบเทียบ" (Stout และ Shealy, 1993) ซึ่งแสดงได้ในรูปของ

$$T_R(\theta) > T_F(\theta)$$

$$\begin{aligned} \text{เมื่อ } T_R(\theta) &= E_R[P(\theta, \eta) \mid \theta] \\ &= \int p(\theta, \eta) f_R(\eta \mid \theta) d\eta \\ T_F(\theta) &= E_F[P(\theta, \eta) \mid \theta] \\ &= T_R(\theta) \\ \theta &= (\theta, \eta) \end{aligned}$$

โดยที่ R คือ กลุ่มอ้างอิง F คือ กลุ่มเปรียบเทียบ

θ = ความสามารถเป้าหมาย

η = ความสามารถส่อคนแทรกซึ่งอาจเป็นพหุมิติก็ได้

$T_R(\theta)$ = ฟังก์ชันการตอบข้อสอบของกลุ่มอ้างอิง

$T_F(\theta)$ = ฟังก์ชันการตอบข้อสอบของกลุ่มเปรียบเทียบ

ความลำเอียงของข้อสอบจะเกิดขึ้นเมื่อ

$$\eta_F \mid \theta \neq \eta_R \mid \theta$$

นิยามความลำเอียงในแง่ของแบบสอบทั้งฉบับ เมื่อให้ U เป็นการตอบข้อสอบแบบสุ่มของผู้สอบคนหนึ่ง $h(U)$ คือ คะแนนจากแบบสอบ โดยที่ $h(U) = \sum U_i$ แล้วจะได้ว่าความลำเอียงของข้อสอบจะเกิดต่อกลุ่มเปรียบเทียบในระดับความสามารถ θ ถ้า

$$\begin{aligned} T_R(\theta) &> T_F(\theta) \\ \text{เมื่อ } T_R(\theta) &\stackrel{\text{def}}{=} E_R[h(U) \mid \theta] \\ T_F(\theta) &\stackrel{\text{def}}{=} E_F[h(U) \mid \theta] \end{aligned}$$

โดยที่ปริมาณของความลำเอียงที่ระดับ θ คือ $\beta(\theta) \stackrel{\text{def}}{=} T_R(\theta) - T_F(\theta)$ และดัชนีของความลำเอียง คือ

$\beta_u = \int B(\theta) f_F(\theta) d\theta$ ซึ่งเป็นความลำเอียงเฉลี่ยต่อผู้สอบแต่ละคนในกลุ่ม
ที่สนใจศึกษา

สมมติฐานเพื่อการทดสอบความลำเอียงของข้อสอบหรือแบบสอบของ SIBTEST

$$H_0: \beta = 0 \quad H_1: |\beta| > 0$$

การทดสอบทำโดยแบ่งแบบสอบออกเป็นแบบสอบย่อย 2 ฉบับ คือ V เป็นแบบสอบ
ที่มีความตรง และ S แบบสอบที่ต้องการศึกษา และคะแนนของแต่ละแบบสอบย่อยคือ

$$Y = \sum_{S} U_i$$

$$X = \sum_{V} U_i$$

สำหรับการตรวจสอบความลำเอียงของข้อสอบรายข้อ S คือ ข้อสอบแต่ละข้อที่ตรวจสอบ
และ V คือ ข้อสอบที่เหลือในแต่ละรอบของการตรวจสอบ

\bar{Y}_{rk} คือ คะแนนเฉลี่ยที่ได้จากแบบสอบย่อยที่ศึกษาสำหรับผู้เข้าสอบทุกคนที่มีคะแนนรวม
ในแบบสอบย่อยที่ถือว่าเป็นแบบสอบที่มีความตรง ($X = k$)

และความลำเอียงวัดได้จาก $(\bar{Y}_{rk} - \bar{Y}_{Fk}) S$ ทั้งนี้ เพราะ $\bar{Y}_{rk} - \bar{Y}_{Fk}$
เป็นความแตกต่างของการตอบของกลุ่มย่อยที่มีความสามารถระดับเดียวกัน ($= k$) บน ความ
สามารถเป้าหมาย และกรณีข้อสอบไม่มีความลำเอียงค่า $\bar{Y}_{rk} - \bar{Y}_{Fk} = 0$

สำหรับค่าสถิติที่ประมาณค่าความลำเอียง β ได้จาก

$$\beta = \sum_{k=0}^n P_k (\bar{Y}_{rk} - \bar{Y}_{Fk})$$

เมื่อ $P_k =$ สัดส่วนของผู้สอบกลุ่มที่สนใจศึกษาที่ตอบแบบสอบย่อยที่ถือว่าไม่ลำเอียงได้ถูก

k ข้อ

สถิติทดสอบของ SIBTEST คือ

$$T = \frac{\hat{\beta}}{\sigma(\hat{\beta})} \quad \text{ซึ่ง} \quad T \sim N(0,1) \quad \text{เมื่อ} \quad \beta = 0$$

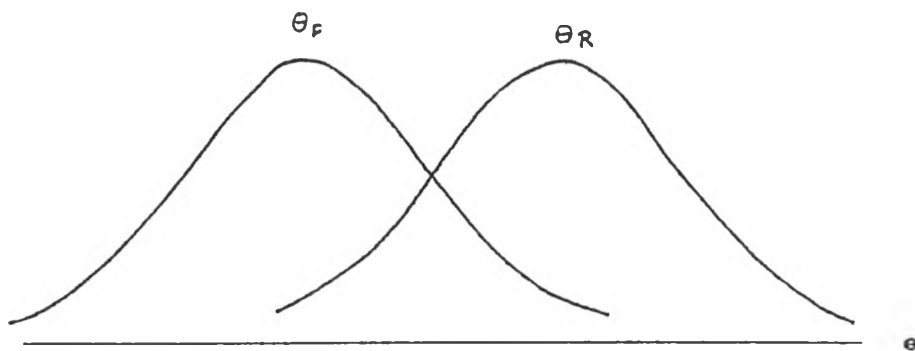
โดยที่

$$\sigma(\hat{\beta}) = \left[\sum_{k=0}^n P_k \left[(1/J_{Rk}) \sigma^2(y|k,R) + (1/J_{Fk}) \sigma^2(y|k,F) \right] \right]^{1/2}$$

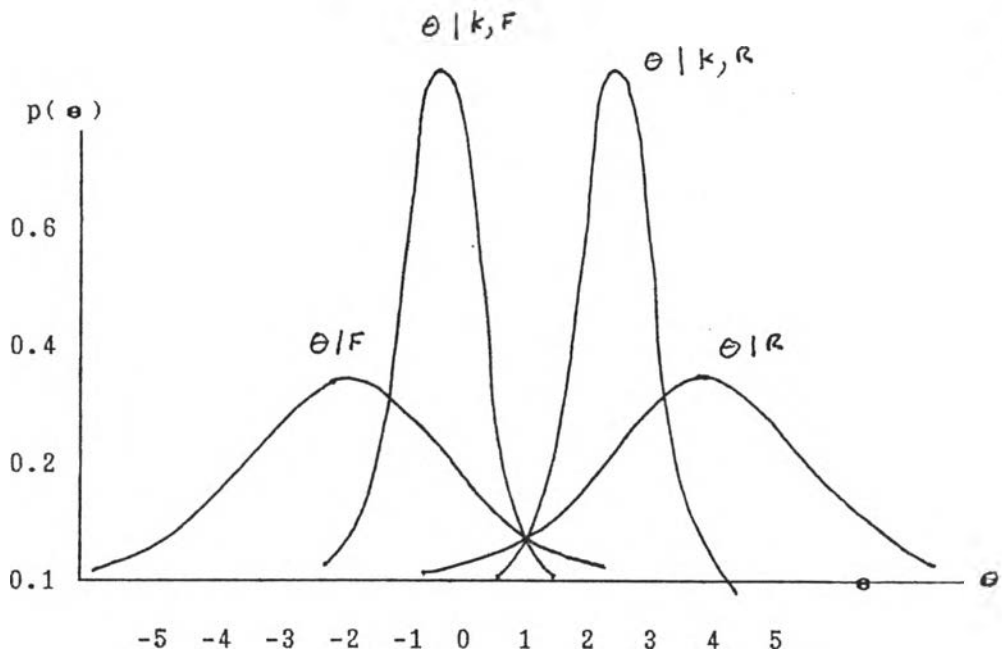
เมื่อ $\sigma^2(y|k,g)$ เป็นความแปรปรวนของคะแนนของกลุ่มอ้างอิงและกลุ่มเปรียบเทียบ ในช่วงคะแนน k

อย่างไรก็ตามเมื่อมีความแตกต่างระหว่างความสามารถเป้าหมาย ค่าสถิติ T จะสูงกว่าปกติ ทำให้เกิดการระบุว่าข้อสอบลำเอียงทั้ง ๆ ที่ไม่ลำเอียง

จากภาพข้างล่าง ซึ่งแสดงความแตกต่างระหว่างกลุ่มประชากรย่อยในความสามารถเป้าหมาย (θ)



เมื่อแบ่งระดับความสามารถตามแบบสอบย่อยที่ถือว่าไม่ลำเอียงจะปรากฏปัญหาการมีค่าสถิติสูงกว่าปกติ ดังภาพ



นั่นคือ ที่ระดับความสามารถ k เราคาดได้ว่า แม้ข้อสอบไม่ลำเอียง ผู้สอบในกลุ่มสนใจศึกษา จะมีความสามารถต่ำกว่าผู้สอบในกลุ่มอ้างอิง นั่นคือ

$$\bar{Y}_{Rk} - \bar{Y}_{Fk} > 0 \quad \text{ทุก ๆ ค่าของ } k$$

ดังนั้น จึงใช้ค่าแก้ด้วยการใช้ค่าถดถอยตามทฤษฎีการตอบข้อสอบชนิด 3 พารามิเตอร์ เพื่อกำจัดค่าเพื่อของค่าสถิติ โดยการแปลงค่า

$$\bar{Y}_{Rk} - \bar{Y}_{Fk} \quad \text{ให้เป็น} \quad \bar{Y}_{Rk}^* - \bar{Y}_{Fk}^*$$

โดยที่

$$Y_{gk} = \bar{Y}_{gk} + \hat{M}_{gk} [\hat{V}(k) - \hat{V}_g(k)]; \quad g \text{ คือ } R \text{ หรือ } F$$

$$\bar{Y}_{g, (k+1)} - \bar{Y}_{g, (k-1)}$$

$$\text{และ } \hat{M}_{gk} = \frac{\bar{Y}_{g, (k+1)} - \bar{Y}_{g, (k-1)}}{V_g(k+1) - V_g(k-1)}$$

$$V(k) = (1/2)[V_R(k) + V_F(k)] = (1/2)(T_V(\theta_{Rk}) + T_V(\theta_{Fk}))$$

$$V(k) = T_V(\theta_k)$$

การทดสอบนัยสำคัญทางสถิติ เมื่อ $H_0: \beta = 0$ ได้แก่

$$Z_{\beta} = \frac{\beta}{SE(\beta)}$$

และจะปฏิเสธสมมติฐานศูนย์ที่ว่า $\beta = 0$

เมื่อ $\beta > Z_{\alpha}$ โดยที่ค่า Z_{α} กำหนดไว้จาก $PN(0,1) > Z_{\alpha} = \alpha$

การศึกษาวิจัยครั้งนี้ ผู้วิจัยได้เลือกพัฒนาเกณฑ์ การตัดสินใจข้อสอบล่าเอียงของดัชนี B_{S1B} จากข้อมูลเชิงประจักษ์ เพื่อนำไปใช้แทนการทดสอบนัยสำคัญ ของค่า Z ที่ได้จากโปรแกรมการวิเคราะห์

5. งานวิจัยที่เกี่ยวข้องกับการตรวจสอบความล่าเอียงของข้อสอบ

งานวิจัยเกี่ยวกับการตรวจค้นความล่าเอียงของข้อสอบเท่าที่มีผู้ทำมานั้น แบ่งเป็น งานวิจัยภายในประเทศไทยและงานวิจัยต่างประเทศ ดังนี้

5.1 งานวิจัยในประเทศ

ชี้ชัย เผ่าพงษ์ (2526) ได้ศึกษาความล่าเอียงระหว่างผู้สอบเพศชายและหญิง ของแบบทดสอบมาตรฐานวัดความถนัดทางการเรียนด้านคณิตศาสตร์ และภาษาในระดับมัธยมศึกษาตอนต้น จำนวน 2 ฉบับ ด้วยวิธีวิเคราะห์ความแตกต่างของโค้งลักษณะข้อสอบ จากทฤษฎีการตอบข้อสอบ 3 พารามิเตอร์ กลุ่มตัวอย่างที่ใช้เป็นนักเรียนชั้นมัธยมศึกษาปีที่ 3 ในปีการศึกษา 2524 จากทุกภาคภูมิศาสตร์ของประเทศไทย ได้แก่ ภาคเหนือ ภาคกลาง ภาคใต้ ภาคตะวันออก และภาคตะวันออกเฉียงเหนือ

ผลการวิจัยพบว่า แบบทดสอบวัดความถนัดทางการเรียนด้านคณิตศาสตร์ มีข้อ กระทั่งล่าเอียงต่อกลุ่มนักเรียนชายและหญิงในระดับปานกลางขึ้นไป จำนวน 5 ข้อ เนื้อหาของ ข้อกระทั่งเป็นเรื่องเกี่ยวกับร้อยละ การหาปริมาตรและการหาความยาวเส้นรอบรูปสามเหลี่ยม อย่างละ 1 ข้อ และเป็นเนื้อหาเกี่ยวกับโจทย์ปัญหาทางคณิตศาสตร์ 2 ข้อ ข้อสอบที่พบว่า ล่าเอียงทั้ง 5 ข้อ ล่าเอียงต่อนักเรียนชาย 4 ข้อ อีก 1 ข้อ ล่าเอียงต่อนักเรียนหญิงในช่วง ความสามารถแรก และล่าเอียงต่อนักเรียนชายในช่วงความสามารถต่อมา

สำหรับการวิเคราะห์ความลำเอียงในแบบทดสอบวัดความถนัดทางภาษาเกี่ยวกับการอ่านเข้าใจ มีข้อสอบที่พบว่าการอ่านเข้าใจในระดับปานกลางขึ้นไป 9 ข้อ เป็นข้อที่มีเนื้อหาวัดความเข้าใจเกี่ยวกับการอ่านคำประพันธ์ และบทร้อยกรองอย่างละ 1 ข้อ อีก 7 ข้อ เป็นเนื้อหาเกี่ยวกับความเข้าใจในการอ่าน ในจำนวนข้อสอบทั้ง 9 ข้อนี้ เป็นข้อสอบที่มีความลำเอียงต่อกลุ่มนักเรียนชาย 1 ข้อ ต่อนักเรียนหญิง 6 ข้อ และอีก 2 ข้อ มีความลำเอียงต่อนักเรียนชายในระดับความสามารถแรก ๆ และลำเอียงต่อนักเรียนหญิงในระดับความสามารถที่สูงขึ้น

ทัศนีย์ นีรมนตรี (2530) ได้ศึกษาถึงความลำเอียงของแบบสอบวิชาคณิตศาสตร์ของโครงการตรวจสอบคุณภาพการศึกษานักเรียนมัธยมศึกษาปีที่ 6 ปีการศึกษา 2526 ด้วยวิธีวิเคราะห์ 3 วิธีคือ วิธีกำหนดจุดค่าเฉลี่ย ($D + .75 Z_u$) (2) วิธีการตอบข้อสอบแบบ 3 พารามิเตอร์ ($\phi > .40$) และ (3) วิธีการทดสอบความแตกต่างระหว่างกลุ่มด้วยสถิติไคสแควร์ในโมเดลล็อกลิเนียร์ 2 โมเดล คือ โมเดลที่ไม่มีพารามิเตอร์ผลร่วมระหว่างระดับคะแนนกับกลุ่ม และโมเดลที่ไม่มีพารามิเตอร์ของผลหลักที่เกิดจากกลุ่ม (เปรียบเทียบค่า χ^2 ที่คำนวณได้กับค่า χ^2 ในตาราง) โดยเปรียบเทียบจำนวนข้อกระทงที่มีความลำเอียงระหว่างกลุ่มนักเรียนในกรุงเทพมหานครกับกลุ่มนักเรียนในภาคภูมิศาสตร์ 5 ภาค ได้แก่ ภาคกลาง ภาคเหนือ ภาคใต้ ภาคตะวันออก และภาคตะวันออกเฉียงเหนือ

ผลการศึกษาพบว่าเมื่อแยกข้อสอบออกตามค่าความยากที่วิเคราะห์ด้วยวิธีดั้งเดิมพบว่า ข้อที่ลำเอียงมีจำนวน 43 ข้อ และไม่ลำเอียง 17 ข้อ เหมือนกันทุกคู่ทุกภาคและเมื่อเปรียบเทียบจำนวนข้อที่ลำเอียงระหว่างกลุ่มนักเรียนในกรุงเทพมหานคร และกลุ่มนักเรียนภาคอื่น ๆ พบว่าในการใช้ IRT-3 พารามิเตอร์ พบข้อกระทงที่มีความลำเอียงจำนวนมากที่สุด มีข้อกระทงที่ลำเอียงซ้ำกันระหว่างกลุ่มนักเรียนในกรุงเทพมหานครกับทุก ๆ ภาค แต่จำนวนไม่เท่ากันในแต่ละวิธี วิธีที่ 1 และ 3 มีข้อลำเอียงซ้ำกันมากที่สุดระหว่างกรุงเทพมหานครและภาคตะวันออกเฉียงเหนือ ส่วนวิธีที่ 2 มีจำนวนข้อที่ลำเอียงซ้ำกันมากที่สุดระหว่างกรุงเทพมหานครและภาคตะวันออกเฉียงเหนือ

การเปรียบเทียบผลการวิเคราะห์ทั้ง 3 วิธี กับความลำเอียงที่เกิดขึ้นในภาคเดียวกันพบข้อกระทงที่ลำเอียงซ้ำกัน ข้อกระทงที่ลำเอียงส่วนใหญ่เป็นข้อที่ง่ายสำหรับนักเรียนในกรุงเทพมหานครมากกว่ากลุ่มนักเรียนในภาคอื่น ๆ สำหรับวิธีวิเคราะห์ที่ 1 แต่เป็นข้อที่มีความลำเอียงในเกณฑ์ต่ำสำหรับวิธีที่ 2 และเป็นข้อที่ลำเอียงอย่างสม่ำเสมอในการวิเคราะห์ด้วยวิธีที่ 3

สรศักดิ์ อมรัตน์ศักดิ์ (2531) ได้ศึกษาถึงความลำเอียงของข้อสอบ โดยพิจารณาถึงสัมประสิทธิ์สหสัมพันธ์ระหว่าง วิธีวิเคราะห์ความลำเอียงของข้อสอบ 4 วิธี คือ วิธีวิเคราะห์ความแปรปรวน วิธีแปลงค่าความยากของข้อสอบ ($D + 3 Sd$) วิธีโค้งลักษณะข้อสอบ



โมเดล 1 พารามิเตอร์ (ทดสอบความแตกต่างระหว่างค่าความยากด้วย t เกณฑ์ที่แสดงว่าข้อสอบ
ลำเอียง คือ $t \geq 2.4$) และ 3 พารามิเตอร์ (ϕ, γ, δ) และเปรียบเทียบความแตกต่าง
ของผลการคัดเลือกก่อนและหลังการศึกษา ความลำเอียงของข้อสอบในด้านจำนวนผู้ได้รับการ
คัดเลือกสัดส่วนชาย : หญิงที่ได้รับการคัดเลือกและความเที่ยงของแบบสอบ ผลการวิจัยสรุปได้ว่า

1. วิธีโด่งลักษณะข้อสอบโมเดล 3 พารามิเตอร์ ค้นพบข้อสอบลำเอียงได้มาก
ที่สุด รองลงมา ได้แก่ วิธีวิเคราะห์ความแปรปรวน วิธีที่ตรวจค้นข้อสอบที่ลำเอียงได้น้อยที่สุด
ได้แก่ วิธีแปลงค่าความยากของข้อสอบ

2. วิธีทั้ง 4 วิธี มีความสัมพันธ์กันทางบวกอย่างมีนัยสำคัญที่ระดับ .001 โดยมี
ค่าระหว่าง .7535-.9921

3. การใช้คะแนนดิบและคะแนนรวมแบบอื่น ๆ อีก 5 วิธี มีจำนวนผู้ได้รับการ
คัดเลือกแตกต่างกันประมาณร้อยละ 4 ถึง 24 ส่วนการใช้คะแนนมาตรฐานที่ปกปิดรวมกับคะแนน
แปลงแบบอื่น ๆ อีก 4 วิธี มีจำนวนผู้ได้รับการคัดเลือกแตกต่างกันร้อยละ 4-23

4. เมื่อตัดข้อสอบที่มีความลำเอียงออก พบว่าสัดส่วนหญิงและชายที่ได้รับการ
คัดเลือกมีความใกล้เคียงกัน และค่าความเที่ยงของแบบสอบลดลงเล็กน้อย

สุวัฒน์ สุกมลสันต์ (2534) ได้วิเคราะห์ความลำเอียงของข้อสอบภาษาอังกฤษ
เพื่อคัดเลือกเข้าศึกษาในมหาวิทยาลัย โดยใช้ข้อมูลการตอบแบบทดสอบภาษาอังกฤษเข้า
มหาวิทยาลัยชุด กข และ ชุด กขค ปี 2531-2533 ซึ่งมีข้อสอบชุดละ 100 ข้อ ตัวแปรที่ศึกษา
ความลำเอียงคือ เพศ และภาคภูมิศาสตร์ของผู้สอบ ซึ่งแยกออกเป็น 5 ภาคตามภูมิภาคนาของ
ผู้สอบ ได้แก่ ภาคกลาง ภาคตะวันออก ภาคตะวันออกเฉียงเหนือ ภาคเหนือ และภาคใต้
เทคนิควิธีที่นำมาวิเคราะห์ความลำเอียงมี 3 วิธี ได้แก่ วิธีกำหนดจุดค่าเดลต้า (Delta-Plot
Method) วิธีไคสแควร์ชนิดที่แบ่งความสามารถของผู้สอบเป็น 3 ระดับ ได้แก่ กลุ่มความสามารถ
ระดับต่ำ (ผู้ได้คะแนนรวม 0-40 คะแนน) กลุ่มความสามารถระดับปานกลาง (ผู้ได้คะแนนรวม
41-70 คะแนน) และระดับสูง (ผู้ได้คะแนนรวม 71-100 คะแนน) วิธีที่ 3 ที่นำมาใช้วิเคราะห์
ความลำเอียงของข้อสอบได้แก่ วิธีการวัดพื้นที่ความแตกต่างระหว่างโด่งลักษณะข้อสอบที่วิเคราะห์
ตามทฤษฎีการตอบข้อสอบแบบ 3 พารามิเตอร์

เกณฑ์การตัดสินความลำเอียงของข้อสอบสำหรับวิธีต่าง ๆ ที่ สุวัฒน์ สุกมลสันต์
กำหนดขึ้นคือ

1. วิธีกำหนดจุดค่าเดลต้า ข้อสอบลำเอียงคือ ข้อที่มีระยะตั้งฉาก จากค่า
อันดับเดลต้าไปยังเส้นแกนหลัก (d) มากกว่า $= + 1.964 Sd$.

2. วิธีไคสแควร์ ข้อสอบลำเอียง คือ ข้อสอบที่ผู้สอบต่างกลุ่มที่อยู่ในระดับ
คะแนนเดียวกัน มีสัดส่วนการตอบถูกหรือผิดแตกต่างกันอย่างมีนัยสำคัญที่ระดับ $\alpha = 0.017$

3. วิธีการวัดพื้นที่ความแตกต่างระหว่างโด่งลักษณะของข้อสอบ ข้อสอบล่าเอียง คือ ข้อที่มีพื้นที่ความแตกต่างมากกว่า 0.40 โดยแบ่งเป็น 3 ระดับ

3.1 พื้นที่ > 0.70 แสดงว่าล่าเอียงมาก

3.2 พื้นที่ระหว่าง 0.40-0.70 แสดงว่าล่าเอียงปานกลาง

3.3 พื้นที่ที่มีค่ามากกว่า 0.00 แต่น้อยกว่า 0.40 แสดงว่าล่าเอียงน้อย

ผลการวิเคราะห์พบว่าแบบทดสอบภาษาอังกฤษ ฉบับ กข และ กขค ปี 2531-2533 มีความล่าเอียงต่อเพศ 7-28 ข้อ และ 4-41 ข้อ โดยมีแนวโน้มที่ล่าเอียงต่อเพศชายมากกว่าเพศหญิง และล่าเอียงต่อภาคภูมิศาสตร์ 6-45 ข้อ และ 5-43 ข้อ โดยมีค่าความล่าเอียงต่อผู้สอบจากภาคอื่นมากกว่าภาคกลางประมาณ 2-3 เท่า

การวิเคราะห์เปรียบเทียบผลของการวิเคราะห์โดยใช้วิธีวิเคราะห์ต่างกัน พบว่า มีจำนวนข้อสอบที่ระบุว่าล่าเอียงต่างกันอย่างมีนัยสำคัญ แต่ละวิธีให้ผลที่มีความสัมพันธ์กัน อย่างไม่มีนัยสำคัญ ผลการวิเคราะห์ด้วยวิธีทฤษฎีการตอบข้อสอบ เมื่อใช้เกณฑ์พื้นที่มากกว่า 0.40 แสดงความล่าเอียงของข้อสอบ พบว่ามีจำนวนข้อสอบที่ล่าเอียงมากที่สุด รองลงมา ได้แก่ วิธีไคสแควร์ และวิธีกำหนดจุดค่าเคลต้า ซึ่งวิธีสุดท้ายนี้พบข้อสอบล่าเอียงน้อยกว่า 2 วิธีแรกประมาณ 3-4 เท่า

จะเห็นว่า การศึกษาถึงความล่าเอียงของข้อสอบในประเทศไทยมีการศึกษา อยู่ 2 ลักษณะ คือ การศึกษาความล่าเอียงระหว่างผู้สอบเพศชายและหญิง และระหว่างผู้สอบที่มีภูมิลำเนาอยู่ต่างภาคภูมิศาสตร์ วิธีที่นำมาใช้ ได้แก่ วิธีทฤษฎีการตอบข้อสอบโมเดล 1 และ 3 พารามิเตอร์ วิธีกำหนดจุดค่าเคลต้า วิธีไคสแควร์ วิธีวิเคราะห์ความแปรปรวน และวิธีแปลงค่าความชากของข้อสอบ ผลการศึกษาพบว่า ในขณะที่ ชัยชัย เผ่าพงษ์ พบว่าข้อสอบวัดความถนัดด้านคณิตศาสตร์ที่ล่าเอียงส่วนมากล่าเอียงต่อเพศชาย และข้อสอบวัดความถนัดด้านภาษาที่ล่าเอียงส่วนมากล่าเอียงต่อเพศหญิงนั้น สุวัฒน์ สุกมลสันต์ พบว่าข้อสอบล่าเอียงด้านภาษามีแนวโน้มล่าเอียงต่อเพศชายมากกว่าหญิง

ในด้านความล่าเอียงระหว่างภาคพบว่ามีความสอดคล้องกันระหว่างการศึกษาของทัศนีย์ พิรมนตรี ที่พบว่าข้อสอบล่าเอียงมักเป็นข้อสอบที่ง่ายสำหรับผู้สอบในกรุงเทพมหานคร ในขณะที่ สุวัฒน์ สุกมลสันต์ พบว่าข้อสอบที่ล่าเอียงมีความล่าเอียงต่อผู้สอบภาคอื่นมากกว่าภาคกลาง โดยเฉลี่ยแล้วล่าเอียงต่อผู้สอบภาคตะวันออกเฉียงเหนือ ภาคเหนือ และภาคใต้เป็นจำนวนใกล้เคียงกันและมากกว่าที่ล่าเอียงต่อผู้สอบภาคกลางโดยเฉลี่ยแล้วล่าเอียงต่อผู้สอบภาคตะวันออกเฉียงเหนือ ภาคเหนือและภาคใต้เป็นจำนวนใกล้เคียงกันและมากกว่าที่ล่าเอียงต่อผู้สอบภาคกลาง

ข้อค้นพบอีกประการคือ นอกจากแต่ละวิธีที่ใช้จะให้ผลการค้นพบจำนวนข้อสอบล่าเอียง ไม่เท่ากันและผลการศึกษาของสัทธน์ สุกมลสันต์ พบว่า มีความสัมพันธ์กันระหว่างการค้นพบแต่ละวิธีอย่างไม่มีนัยสำคัญ นอกจากนั้นเกณฑ์ตัดสินความล่าเอียง ด้วยดัชนี ϕ หรือ (SA) ยังมีความแตกต่างกันไปด้วย

กล่าวได้ว่า การศึกษาเรื่องความล่าเอียงของข้อสอบในประเทศไทยยังมีน้อยมาก วิธีการ MH และ SIBTEST ยังไม่มีการนำมาศึกษาวิจัย และยังไม่พบว่ามีการศึกษาเปรียบเทียบด้านความยาวของข้อสอบ และการพัฒนาเกณฑ์การตัดสินความล่าเอียง นอกจากนั้น ข้อค้นพบความล่าเอียงของข้อสอบคณิตศาสตร์และภาษาอังกฤษระหว่างผู้สอบต่างประเทศยังไม่สอดคล้องกัน ผู้วิจัยจึงหวังว่าการศึกษาของผู้วิจัยในครั้งนี้จะช่วยให้ข้อมูลเพิ่มเติมมากขึ้นในส่วนที่กล่าวมา และจะช่วยให้ได้มีการศึกษาวิจัยในเรื่องนี้มากขึ้นในอนาคต

5.2 งานวิจัยต่างประเทศ

งานวิจัยเกี่ยวกับความล่าเอียงของข้อสอบในต่างประเทศนั้น ผู้วิจัยจะนำเสนอเป็น 2 ประเด็น คือ (1) การตรวจค้นข้อสอบล่าเอียงของข้อมูลจำลอง และ (2) การตรวจค้นข้อสอบล่าเอียงจากข้อมูลจริง

5.2.1 การตรวจค้นข้อสอบล่าเอียงจากข้อมูลจำลอง

การตรวจค้นข้อสอบล่าเอียงจากข้อมูลจำลองส่วนมากเป็นวิธีการศึกษาเปรียบเทียบเทคนิคการตรวจสอบต่าง ๆ ที่ผู้เสนอขึ้นมาใช้ดังนี้

Rudner Getson และ Knight (1980) ได้เปรียบเทียบเทคนิคการตรวจค้นความล่าเอียงของข้อสอบจากข้อมูล 3 วิธี ได้แก่ เทคนิคความยากแปลงชนิดแกนหลัก (Transformed Item Difficulties Major Axis = TID - MA) เทคนิคความยากแปลงเส้นแกน 45 องศา (Transformed Item Difficulties-45 Line = TID-45) เทคนิคทฤษฎีการตอบข้อสอบ 3 พารามิเตอร์ (ICC-3) และ 1 พารามิเตอร์ ชนิดค่าสถิติความเหมาะสม (ICC-1F) ชนิดความแตกต่างในความง่ายของข้อสอบ (ICC-1E) เทคนิคไคสแควร์แบบ 5 ช่วง (CHI-5) และเทคนิคไคสแควร์ชนิดหลายช่วง (CHI-N) ผลจากการวิเคราะห์ข้อมูลพบว่า เมื่อใช้ความยาวของข้อสอบหลายขนาดไม่พบว่ามีเทคนิคใดที่มีผลจากความยาวของข้อสอบเป็นพิเศษ เมื่อพิจารณาจากค่าสหสัมพันธ์ของผลการตรวจค้นกับข้อมูลที่สร้างขึ้น โดยพิจารณาจากปริมาณของความล่าเอียงของข้อสอบ พบว่า วิธีที่ถูกต้องมากที่สุด ได้แก่ ICC-3 CHI-5 TID-45 ($r = .80 .73$ และ $.68$) และ ที่ถูกต้องน้อยที่สุด ได้แก่ ICC-1F ($r = .55$) และทุกวิธียกเว้น ICC-1F สามารถตรวจค้นความล่าเอียงได้ในกรณีที่มีความยากต่างกันเท่านั้น มีเพียง ICC-1F ที่มีความไวต่อความล่าเอียงในอำนาจจำแนกของข้อสอบด้วย

เทคนิค CHI-5 มีประสิทธิภาพเท่า ๆ กับ ICC-3 และพบว่าถ้าเพิ่มจำนวนช่วงความสามารถมากขึ้น จะลดความถูกต้องในการตรวจค้นลง ซึ่งอาจเป็นเพราะช่วงความสามารถในแต่ละข้อนั้นต่างกันไปได้

สรุปว่า ถ้าไม่คำนึงถึงความยาวของแบบสอบ และขนาดหรือธรรมชาติของความล่าเอียงแล้ว ICC-3 และ ไคสแควร์แบบ 5 ช่วง เป็นวิธีที่ควรใช้วิธีค่าความยากแปลงก็อาจนำมาใช้แทนได้ ถ้าจะตรวจค้นในด้านความต่างของค่าความยากของข้อสอบ

Van Der Flier Mellenbergh Ader และ Win (1984) ได้เปรียบเทียบวิธีการตรวจค้นความล่าเอียงของข้อสอบจากข้อมูลจำลองด้วยวิธีที่มีเงื่อนไข วิธีโลจิกแบบตัดแปลงจากวิธีของ Lord ที่เสนอการใช้ทฤษฎีการตอบข้อสอบ 3 พารามิเตอร์ แบบเป็นขั้นตอนโดย

ขั้นที่ 1 ประมาณค่าพารามิเตอร์ แล้วคำนวณค่าสถิติที่แสดงความล่าเอียงของข้อสอบทุกข้อ

ขั้นที่ 2 ตัดข้อสอบที่ล่าเอียงออก แล้วประมาณค่าพารามิเตอร์ใหม่ จากนั้นจึงนำค่าพารามิเตอร์ที่ได้คำนวณค่าสถิติความล่าเอียงของข้อสอบทุกข้อ (รวมทั้งข้อสอบที่ล่าเอียงด้วย) ข้อที่มีค่าไคสแควร์สูงสุดในรอบนั้นและมีนัยสำคัญ ถือว่าเป็นข้อสอบที่ล่าเอียง ข้ออื่น ๆ นอกนั้น ถือว่าเป็นข้อที่ไม่ล่าเอียงแม้ว่า ค่าไคสแควร์ จะมีนัยสำคัญ

ขั้นที่ 3 ทำขั้นที่ 1 และ 2 อีกครั้ง การทำซ้ำจะจบลงถ้าทำซ้ำเท่ากับจำนวนรอบที่กำหนดไว้ หรือค่าไคสแควร์สูงสุดของรอบนั้นไม่มีนัยสำคัญ

ในแต่ละรอบของการคำนวณ ค่าตัวเลขค่าไคสแควร์จากข้อสอบทุกข้อรวมทั้งข้อที่ตรวจค้นว่ามีความล่าเอียงในรอบก่อน ๆ

การคำนวณแต่ละรอบ จำนวนข้อสอบที่ขจัดออกไปจะมากกว่าจำนวนข้อสอบในรอบที่แล้ว 1 ข้อ

ข้อสอบที่ตัดออกจากการคำนวณคะแนนรวมในการคำนวณรอบที่แล้ว ไม่จำเป็นต้องถูกตัดออกไปในรอบต่อไป ถ้าค่าไคสแควร์ ของค่านั้นยอมรับได้ในรอบต่อมา ด้วยวิธีนี้คะแนนรวมจะเป็นอิสระจากข้อสอบล่าเอียง และข้อสอบทุกข้อจะได้รับการตรวจสอบโดยมีระดับความสามารถที่ได้จากแบบสอบที่ไม่ล่าเอียง

ข้อดีของการทำซ้ำคือจะแน่ใจว่าได้ข้อสอบที่ล่าเอียงจริง ๆ จากการตรวจสอบเพราะจะระบุเฉพาะข้อที่มีค่าไคสแควร์ที่สูงสุด และมีนัยสำคัญเท่านั้น

Lautenschlager และ Park (1988) ได้ใช้ข้อมูลจำลองศึกษาเปรียบเทียบวิธีการเชื่อมโยงค่าพารามิเตอร์ (parameter linking method) ซึ่งมี การแปลงค่าพารามิเตอร์ที่ประมาณได้จากทฤษฎีการตอบข้อสอบ ให้เป็นค่ามาตรฐานเพื่อเปรียบเทียบกันได้

ข้ามกลุ่ม โดยมีวิธีที่ Linn และคณะ (1980 ; cited in Lautenschlager และ Park, 1988) และวิธีของ Warm (1978) cited in Lautenschlager และ Park 1988)

วิธีของ Warm ง่ายที่สุดจากการใช้ความคิดรวบยอดที่ว่า ความยากของข้อสอบของทุกกลุ่มจะเหมือนกันได้ด้วยการแปลงค่าเชิงเส้นตรง โดยได้แปลงค่าพารามิเตอร์สัมพันธ์โดยใช้สูตรดังนี้

$$b_1 = \frac{SD_{b_1}}{SD_{b_2}} b_2 + (M_{b_1} - \frac{SD_{b_1}}{SD_{b_2}} M_{b_2})$$

$$a_1 = a_2 \frac{SD_{b_2}}{SD_{b_1}}$$

b_1 คือ ค่าพารามิเตอร์ b บนมาตราของกลุ่ม 1

b_2 คือ ค่าพารามิเตอร์ b บนมาตราของกลุ่ม 2

a_1 คือ ค่าพารามิเตอร์ a บนมาตราของกลุ่ม 1

a_2 คือ ค่าพารามิเตอร์ a บนมาตราของกลุ่ม 2

M_{b_1}, SD_{b_1} คือ ค่าเฉลี่ยและค่าความเบี่ยงเบนมาตรฐานของค่า b บนมาตราของกลุ่ม 1

M_{b_2}, SD_{b_2} คือ ค่าเฉลี่ยและค่าความเบี่ยงเบนมาตรฐานของค่า b บนมาตราของกลุ่ม 2

วิธีการของ Linn นั้น จะประมาณความถูกต้องของ b_1 ด้วยค่าเฉลี่ยและค่าความเบี่ยงเบนมาตรฐานที่ได้ให้น้ำหนักไว้ น้ำหนักของแต่ละข้อจะได้แก่ส่วนกลับของความแปรปรวนในการสุ่มของ b_1 ที่ประมาณค่าได้สูงสุด ดังนั้นข้อที่มีค่าความแปรปรวนในการสุ่มของ b_1 มากที่สุดจะมีน้ำหนักน้อยที่สุดในการพิจารณาค่าคงที่ของการเชื่อมโยงค่าพารามิเตอร์

จากการวิเคราะห์ข้อมูลจำลอง ด้วยการทำทดสอบค่า χ^2 ของ Lord พบว่าภายใต้สถานการณ์ค่าเอียงแบบทิศทางเดียว (แต่ละข้อค่าเอียงเข้าข้างกลุ่มใดกลุ่มหนึ่งเพียงกลุ่มเดียว) นั้น ยากที่จะตรวจค้นจำนวนข้อสอบที่ถูกระบุตรงข้ามกับความเป็นจริง (ข้อสอบค่าเอียงแต่ตรวจพบว่าเป็นข้อสอบไม่ค่าเอียงหรือข้อสอบไม่ค่าเอียงแต่ตรวจพบว่าเป็นข้อที่ค่าเอียง)

ทำให้ยากที่จะเลือกว่าควรใช้วิธีใด แต่ในการศึกษาครั้งนี้ พอลจะพูดได้ว่าวิธีของ Linn และคณะ ดูจะดีกว่า และยังพบว่าข้อสอบที่มีระดับความลำเอียงน้อย มักจะถูกระบุจากการตรวจค้นว่าเป็น ข้อสอบไม่ลำเอียง

Lautenschlager และ Park กล่าวว่า การตรวจสอบข้อสอบลำเอียง ในอนาคต น่าจะปรับปรุงโดยใช้ทฤษฎีการตอบข้อสอบแบบพหุมิติ (multidimensional IRT) แทนเมื่อแบบสอบขาดคุณสมบัติด้านเอกมิติ

Linacre (1988) ได้ตรวจสอบความถูกต้องของวิธี MH และวิธี PROX ของ Rasch ในการตรวจสอบความลำเอียงของข้อสอบที่ได้จากการจำลองข้อมูลขึ้นโดยใช้ตัวประมาณค่าความคลาดเคลื่อนมาตรฐานหลาย ๆ วิธี สำหรับวิธี MH ผลการวิเคราะห์ พบว่า ทั้งสองวิธีใช้ได้ดีในสถานการณ์จำลองทุกสถานการณ์ ส่วนตัวประมาณค่าความคลาดเคลื่อนมาตรฐานของ MH มีคุณสมบัติแตกต่างกันไปและควรมีการแปลผลอย่างระมัดระวัง เมื่อใช้วิธี MH ผู้วิเคราะห์ ควรจะคำนวณขนาดของความลำเอียงของข้อสอบ และความแปรปรวนของความคลาดเคลื่อนด้วยการใช้เกณฑ์หลาย ๆ เกณฑ์ และใช้ตัวประมาณค่าความคลาดเคลื่อนมาตรฐานมากกว่า 1 ตัว สำหรับวิธี PROX จะให้การตรวจสอบการประมาณที่ได้จากวิธี MH

Thissen, Steinberg และ Wainer (1988) ได้ศึกษาความแตกต่างของโครงสร้างการตอบข้อสอบระหว่างกลุ่ม ด้วยทฤษฎีการตอบข้อสอบ เพื่อศึกษาถึงสภาพการณ์ในการทดสอบสมมุติฐาน ความเท่ากันของค่าพารามิเตอร์ด้วยค่าไคสแควร์แบบอัตราส่วน likelihood (IRT-LR) โดยใช้ข้อมูลจำลอง และได้ทำการเปรียบเทียบการทดสอบสมมุติฐานความเท่ากันของค่าพารามิเตอร์ ด้วย IRT-LR และวิธีของ MH พบว่า IRT-LR สามารถตรวจค้นความลำเอียงได้ในเมื่อค่าความแตกต่างของความยากระหว่างกลุ่มมีค่า 0.3 และความแตกต่างในด้านการเดามีค่า 0.1 จำนวนตัวอย่างในกลุ่มน้อย (=500คน) โดยจำนวนข้อสอบร่วมจะน้อย หรือมากก็ไม่มีผลกระทบ ค่าความแตกต่างที่น้อยกว่านี้ไม่สามารถใช้ IRT-LR ตรวจค้นได้ว่าลำเอียงหรือไม่ แต่อาจทำได้ถ้าจำนวนตัวอย่างมากกว่านี้

การระบุความลำเอียงของข้อสอบ ใช้วิธีเปรียบเทียบความแตกต่างระหว่างกลุ่มในข้อสอบที่ต้องการตรวจค้นความลำเอียงกับข้อสอบร่วม ซึ่งมีข้อตกลงว่าไม่มีความลำเอียงอย่างไรก็ตามควรได้ศึกษาเพื่อหาวิธีเลือกข้อสอบร่วม และการทดสอบความลำเอียงของแบบสอบทั้งฉบับ

สำหรับผลการตรวจค้นความลำเอียงของข้อสอบด้วยวิธี IRT-LR และ MH ให้ผลคล้ายคลึงกัน ข้อเสนอนี้คือ MH อาจเป็นประโยชน์ในการประเมินแบบสอบยาว ๆ และอาจใช้เป็นเครื่องมือกลั่นกรองก่อนจะใช้ IRT-LR ต่อไป

Oshima (1989) ได้ศึกษาผลของข้อมูลหลายมิติที่มีต่อการตรวจค้นความล่าช้า โดยจำลองข้อมูล 42 ชุดจากข้อสอบ 40 ข้อ มีโครงสร้างความเป็นพหุมิติ 2 ชนิด คือ ชนิดที่มีลักษณะแฝงที่เด่นอยู่ลักษณะเดียว และที่มีลักษณะแฝงเด่น 2 ลักษณะ เกณฑ์ที่ใช้ในการตัดสินได้จากการคำนวณค่าดัชนี SA UA SSOS (ผลบวกกำลังสองที่คิดเครื่องหมาย) USOS (ผลบวกกำลังสองที่ไม่คิดเครื่องหมาย) เกณฑ์การตัดสินความล่าช้าคือ ค่าที่มากกว่า 2 ส่วนเบี่ยงเบนมาตรฐานของ SA UA SSOS และ USOS ที่คำนวณได้จากกลุ่มอ้างอิง 2 กลุ่ม ผลจากการศึกษารังนี้พบความสัมพันธ์กับผลการศึกษาของ Oshima (1989) ที่พบว่าค่าที่ได้จากวิธีนี้ไม่คงที่ ขึ้นอยู่กับกลุ่มผู้สอบ ค่าที่นำมาใช้ในครั้งนี้นี้จึงใช้ค่าเฉลี่ยจากการทำซ้ำ 5 ครั้ง

อย่างไรก็ตาม เมื่อจำนวนข้อสอบที่มีลักษณะพหุมิติ (ข้อสอบที่วัดลักษณะแฝงเด่นมากกว่า 1 ลักษณะ) ในแบบสอบมีจำนวนไม่เท่ากัน โดยเมื่อมีจำนวนข้อสอบที่มีลักษณะพหุมิติอยู่ในแบบสอบร้อยละ 5 10 และ 15 ข้อ ดัชนี SA ที่ใช้เป็นเกณฑ์จะมีค่าเท่ากับ +.173 +.205 +.201 ดัชนี UA เท่ากับ .262 .277 .293 ดัชนี SSOS เท่ากับ +1.13 +1.23 +1.55 และดัชนี USOS เท่ากับ 1.56 1.60 และ 2.03 ตามลำดับ

ผลการศึกษาพบว่า เมื่อผู้สอบ 2 กลุ่ม มีค่าเฉลี่ยของความสามารถแทรกซ้อน (Nuisance ability - σ_2) เท่ากัน พบว่าการระบุว่าข้อสอบล่าช้าทั้งหมดที่จำลองขึ้นให้ไม่ล่าช้า ทั้งนี้พบทั้งในข้อสอบลักษณะเอกมิติ และพหุมิติ และเมื่อค่าเฉลี่ยของ σ_2 ต่างกัน จะสามารถระบุข้อสอบที่ล่าช้าได้อย่างถูกต้องในอัตราสูง

การค้นพบข้อสอบล่าช้าร้อยละ 100 สำหรับทุกดัชนี เกิดขึ้นเมื่อค่าเฉลี่ยของ σ_2 ของกลุ่มผู้สอบทั้ง 2 กลุ่มมีค่าเท่ากัน และมีจำนวนข้อสอบพหุมิติ อยู่ร้อยละ 5 ข้อ (มีความเป็นเอกมิติมากกว่า) อัตราการค้นพบสูงเพราะค่าเฉลี่ยของดัชนีมีค่าสูง ทั้งสำหรับข้อสอบพหุมิติและข้อสอบที่จำลองให้ล่าช้า อย่างไรก็ตามเมื่อจำนวนข้อสอบล่าช้า ในแบบสอบเพิ่มขึ้น สัดส่วนของข้อสอบล่าช้าที่ค้นพบลดลง

โดยสรุปแล้วดัชนีจาก IRT สามารถตรวจสอบข้อสอบพหุมิติที่สร้างให้ล่าช้า และที่สร้างให้ไม่ล่าช้าได้อย่างถูกต้อง แต่อำนาจในการตรวจสอบลดลง ถ้าจำนวนข้อสอบล่าช้า ในแบบสอบมีมากขึ้น

Oshima เสนอว่า เนื่องจากพบว่าเมื่อใส่ข้อสอบพหุมิติเข้าไปในแบบสอบร้อยละ 20 ทำให้การทดสอบความเป็นเอกมิติถูกปฏิเสธมากกว่าครึ่งหนึ่งของการทดสอบ ดังนั้นจึงควรทดสอบ ความเป็นเอกมิติของแบบสอบทุกครั้ง เมื่อจะวิเคราะห์ด้วย IRT แบบสอบที่ไม่เป็นเอกมิติ ไม่ควรใช้ IRT แบบเอกมิติ แต่อาจพิจารณาใช้ IRT แบบพหุมิติแทน

Swaminathan และ Rogers (1990) ได้ตรวจค้นความล่าเอียงของข้อสอบด้วยวิธีการถดถอยแบบโลจิสต์ และวิธี MH โดยมีขนาดของตัวอย่าง 2 ขนาด ได้แก่ 250 และ 500 คน และความยาวของแบบสอบ 3 ขนาด ได้แก่ 40 60 และ 80 ข้อ จำลองข้อมูลการตอบข้อสอบด้วยโปรแกรม DATAGEN แบบ 3 พารามิเตอร์ ข้อมูลการจำลองความล่าเอียงใช้ 2 พารามิเตอร์โดยให้อำนาจจำแนกของ 2 กลุ่มเท่ากัน ส่วนค่าความชากแปรผันไปเพื่อให้เกิดระดับความล่าเอียงต่าง ๆ ตามต้องการ แต่ในลักษณะการล่าเอียงแบบไม่สม่ำเสมอ (Nonuniform) นั้นให้ค่าความชากของ 2 กลุ่มเท่ากัน แต่ค่าอำนาจจำแนกแปรไปผลการวิเคราะห์พบว่าการตรวจค้นความล่าเอียงของข้อสอบ 2 วิธี โดยใช้การทดสอบ χ^2 มีชั้นความเป็นอิสระ (df) เท่ากับ 2 สำหรับวิธีทดสอบโลจิสต์และการทดสอบ MH- χ^2 ให้ผลใกล้เคียงกัน โดย MH ดีกว่าเล็กน้อย คือ มีการตรวจค้นได้ถูกต้องร้อยละ 75 กรณีใช้ตัวอย่าง 250 คน และตรวจค้นได้ร้อยละ 100 ในตัวอย่าง 500 คน ในทุกความยาวแบบสอบและในกรณีความล่าเอียงแบบสม่ำเสมอ (Uniform) ในกรณีของการล่าเอียงแบบไม่สม่ำเสมอ นั้น MH ไม่สามารถตรวจค้นได้ ส่วนการถดถอยแบบโลจิสต์ตรวจค้นได้ถูกต้องประมาณร้อยละ 50 ในกรณีตัวอย่างน้อยและข้อสอบสั้น และถูกต้องประมาณร้อยละ 75 ในแบบสอบยาวและตัวอย่างขนาดใหญ่แต่กรณีมีความล่าเอียงแบบไม่สม่ำเสมอ โดยที่โค้งตัดกันที่ปลายสุดของโค้งไม่ว่าจะเป็นด้านความสามารถสูงหรือต่ำก็ตาม MH จะดีกว่า

ในการตรวจค้นผิดคือระบุว่าล่าเอียงทั้ง ๆ ที่ไม่ล่าเอียงนั้น MH ดีกว่า โดยตรวจค้นผิดประมาณร้อยละ 1 ส่วนการถดถอยแบบโลจิสต์ผิดประมาณร้อยละ 1-6

วิธีการถดถอยแบบโลจิสต์แม้จะยืดหยุ่นกว่า และมีความเป็นทั่วไปมากกว่า แต่ค่าใช้จ่ายมากกว่าประมาณ 3-4 เท่า

ประโยชน์ของการถดถอยแบบโลจิสต์ อยู่ตรงที่เป็นวิธีการที่มีพื้นฐานเป็นแบบจำลองที่สามารถเพิ่มตัวแปรความสามารถเข้าไปในแบบจำลองได้ เช่นเดียวกับองค์ประกอบอื่น ๆ ที่มีความเกี่ยวข้องและสำคัญซึ่งทำให้เข้าใจในธรรมชาติของความล่าเอียงได้ดีกว่า

Clauser และคณะ (1991) ได้ศึกษาถึงอิทธิพลต่าง ๆ ที่มีผลต่อวิธีการตรวจค้นความล่าเอียงของข้อสอบด้วย MH ข้อมูลที่ใช้เป็นข้อมูลจำลองขึ้นเพื่อใช้ในการศึกษาการตรวจค้นความล่าเอียงแบบสม่ำเสมอ ข้อมูลที่จำลองขึ้นเป็นแบบสอบจำนวน 75 ข้อ มีความล่าเอียงอยู่ 16 ข้อ มีระดับของความแตกต่างในค่าความชากหลายระดับระหว่างกลุ่มอ้างอิงและกลุ่มเป้าหมายซึ่งมีจำนวนกลุ่มละ 1,500 คน มีการวิเคราะห์ระดับความชากที่ต่างกันหลายระดับกับระดับความชากของแบบสอบ 5 ระดับ อำนาจจำแนก 4 ระดับ และการกระจายของความสามารถ 2 ลักษณะ ผลการวิเคราะห์แสดงอิทธิพลจากความแตกต่างในระดับความชากรวมทั้งค่าอำนาจจำแนกและค่าความชาก ซึ่งทำให้ได้ข้อเสนอนี้ว่าค่าสถิติ Mantel Haenszel

นี้ใช้ได้ สำหรับข้อกระทงที่มีค่าอำนาจจำแนกสูงและมีแนวโน้มจะไม่สามารถตรวจค้นความ
ล่าเอียงในข้อสอบที่มีค่าความยากสูง

Mazor และคณะ (1991) ได้ใช้วิธีการตรวจค้นความล่าเอียงแบบ
MH ศึกษาข้อมูลที่จำลองขึ้นโดยใช้กลุ่มตัวอย่างขนาด 2,000 1,000 500 200 และ 100 คน
เพื่อเปรียบเทียบอัตราการตรวจค้นโดยใช้ข้อสอบที่จำลองขึ้น 5 ชุด ชุดละ 75 ข้อ ตรวจสอบ
ความล่าเอียงด้วยการทดสอบค่า $MH-\chi^2$ $df = 1$ ผลการศึกษาพบว่าเมื่อใช้จำนวนผู้สอบ 2,000
คน วิธีนี้ตรวจค้นความล่าเอียงได้ผิดพลาดร้อยละ 30 และเมื่อใช้จำนวนผู้สอบ 500 คนลงไป
พบว่าตรวจค้นความล่าเอียงผิดพลาดร้อยละ 50 ข้อ กระทงที่ไม่สามารถตรวจค้นความล่าเอียงได้
ได้แก่ข้อกระทงที่ยากมากข้อกระทงที่มีค่าความยากต่างกันเล็กน้อยสำหรับคน 2 กลุ่ม และเป็นข้อ
กระทงที่มีค่าอำนาจจำแนกต่ำ

Park และ Lautenschlager (1990) ได้ตรวจค้นความล่าเอียง
ของข้อสอบโดยใช้ทฤษฎีการตอบข้อสอบที่ได้ดัดแปลงจากวิธีของ Lord (1980) ที่เรียกว่า วิธี
'ปลอดจากความล่าเอียง' (Purification) ซึ่งมีวิธีการดังนี้

1) วิเคราะห์แบบสอบทั้งฉบับ โดยการประมาณค่าพารามิเตอร์รวม
แล้วแปลงค่าความยาก (b) ให้เป็นค่ามาตรฐาน จากนั้นกำหนดค่าการเดา (c_1) ตามที่ได้จาก
การวิเคราะห์แล้วประมาณค่าความยากและอำนาจจำแนก (a_1) ใหม่อีกครั้ง โดยแยกตามกลุ่มที่
ศึกษาแล้วแปลงค่า b_1 เป็นค่ามาตรฐานอีกครั้ง จากนั้นเปรียบเทียบฟังก์ชันการตอบข้อสอบเพื่อ
ความล่าเอียง

2) ประมาณค่าระดับความสามารถโดยรวมทุกกลุ่มเข้าด้วยกันใหม่ใช้
แต่เฉพาะข้อสอบที่ไม่ล่าเอียง จะได้ข้อสอบที่มีความเป็นเอกพันธ์

3) กำหนดค่า e เท่ากับที่ประมาณได้ในข้อ 2 แล้วประมาณค่า a_1
และ b_1 ของข้อสอบทุกข้อรวมข้อที่มีความล่าเอียงไว้ด้วย

4) เปรียบเทียบฟังก์ชันการตอบข้อสอบ หรือค่าพารามิเตอร์อีกครั้ง
ด้วยค่าไคสแควร์

วิธีนี้เรียกว่าปลอดจากความล่าเอียง เพราะพยายามที่จะประมาณค่า
พารามิเตอร์โดยตัดข้อสอบล่าเอียงออกจากแบบสอบก่อน

วิธีที่นำมาใช้ในการศึกษาดังนี้เป็นวิธีที่ปรับปรุงจากวิธีของลอร์ดได้แก่

1. วิธีดัดแปลงจาก Drasgow ที่ได้เสนอไว้ เป็นการตรวจค้นความ
ล่าเอียงโดยการทำซ้ำ ๆ หลายรอบ จนกว่าจะได้ชุดข้อสอบที่ล่าเอียงที่ซ้ำกันในการทำซ้ำ 2 รอบ
สุดท้าย วิธีนี้ใช้ค่าไคสแควร์ต่ำสุดของ Divgi (1985) แทนวิธีของ Stocking และ Lord
(1983) วิธีการนี้เรียกว่า M-DIL (Modified-Drasgow's Iterative Linking
Procedure) มีวิธีดังนี้

- 1) ใช้วิธีการเชื่อมโยงมาตรฐานของ สติคกิ้งและลอร์ด ประมาณค่าพารามิเตอร์บนมาตรฐานที่เท่ากันโดยประมาณ
- 2) คำนวณค่าความล่าเอียง
- 3) เชื่อมโยงค่าพารามิเตอร์ที่ประมาณได้อีกครั้ง โดยใช้แค่ข้อที่ไม่ล่าเอียง
- 4) คำนวณค่าความล่าเอียงของทุกข้ออีกครั้ง
- 5) ทำซ้ำ ๆ จนกว่าจะได้ข้อสอบล่าเอียงชุดเดียวกันใน 2 รอบสุดท้าย

2. วิธีดัดแปลงจากวิธี 'ปลอดล่าเอียง' ของ Lord (1980) ซึ่ง Park (1988) ได้เสนอขึ้นเรียกว่า M-LTP (Modified-Lord Test Purification) ซึ่งมีขั้นตอนดังนี้

- 1) คำนวณค่า θ จากทุกกลุ่มรวมกัน
- 2) จากค่า θ ของแต่ละคนที่คำนวณได้ใน 1) นำมาคำนวณค่าพารามิเตอร์ของข้อสอบ แล้วคำนวณค่าสถิติความล่าเอียง
- 3) คำนวณค่า θ ใหม่อีกครั้ง โดยใช้เฉพาะข้อสอบที่ไม่ล่าเอียง
- 4) คำนวณค่าพารามิเตอร์ของข้อสอบอีกครั้งโดยใช้ค่า θ ที่คำนวณได้จาก 3)
- 5) คำนวณค่าสถิติความล่าเอียง
- 6) ทำซ้ำจนกว่าจะได้ข้อสอบซ้ำกัน ในการทำซ้ำ 2

รอบสุดท้าย

3. ใช้วิธีผสมระหว่าง M-DIL และวิธีปลอดล่าเอียงเรียกว่า ILAP (Iterative Parameter Linking and θ Scale Purification) ซึ่งจะเริ่มด้วย M-DIL จนถึงรอบสุดท้ายแล้วเริ่มประมาณค่า θ ใหม่อีกครั้งแยกตามกลุ่มเฉพาะข้อสอบไม่ล่าเอียง จากนั้น จึงคำนวณค่าพารามิเตอร์ข้อสอบแยกตามกลุ่ม แล้วใช้วิธีการของ Divgi ประมาณค่าคงที่ของการเชื่อมโยง โดยใช้เฉพาะข้อที่ไม่ล่าเอียง แล้วคำนวณค่าสถิติความล่าเอียงใหม่ แล้วตามด้วย M-DIL ต่อ



ข้อมูลศึกษา เป็นข้อมูลจำลองด้วยวิธีการจำลองของ Sympon (1978) ภาษาใต้สภาพการณ์ล่าเอียงชนิดทิศทางเดียว และทิศทางผสม พบว่าการวิเคราะห์ 2 วิธีแรก 1 ให้ผลเท่ากัน แต่ M-DIL ต้องทำซ้ำในจำนวนรอบที่มากกว่า M-LTP และถ้าใช้วิธีสุดท้ายแล้วจะสามารถลดจำนวนข้อสอบที่ล่าเอียง แต่ถูกระบุจากการตรวจค้นว่าไม่ ล่าเอียงลง ได้มากกว่า 2 วิธีแรก

Cohen และ Kim (1993) ศึกษาเปรียบเทียบการใช้ทฤษฎีการตอบข้อสอบ 2 พารามิเตอร์ โดยใช้สถิติทดสอบ χ^2 ของลอร์ดและสถิติทดสอบพื้นที่ความแตกต่างระหว่างโค้งลักษณะข้อสอบชนิดคิดเครื่องหมายและไม่คิดเครื่องหมาย อันได้แก่ สถิติ Z ของราชู โดยใช้ข้อมูลจำลองตามทฤษฎีการตอบข้อสอบ 2 พารามิเตอร์ จำลองการตอบข้อสอบ 2 ขนาด ข้อ 20 และ 60 ข้อ จำนวนผู้สอบ 100 และ 500 คน และมีจำนวนข้อสอบล่าเอียงในแบบทดสอบจำนวนร้อยละ 0 10 และ 20

ผลการวิเคราะห์พบว่า

1. แบบสอบยิ่งยาวการระบุข้อสอบล่าเอียงผิดพลาด (FP) จะมากขึ้น และถ้ามีจำนวนข้อสอบล่าเอียงอยู่ในแบบสอบน้อย จำนวนข้อสอบที่ถูกระบุว่าล่าเอียงผิดพลาดจะยิ่งมากขึ้น

องค์ประกอบที่มีอิทธิพลต่อการระบุข้อสอบล่าเอียงผิดพลาด (FP) อีกประการ คือ ระดับนัยสำคัญ (α) ค่า α ยิ่งมาก FP จะเกิดมาก (ที่ .05 จะเกิด FP มากกว่าที่ .01) ผลดังนี้เกิดในทุกสถานการณ์ที่ศึกษา สำหรับ Z แต่เกิดในบางสถานการณ์สำหรับ χ^2

2. แบบสอบที่ยาวกว่า มีร้อยละของข้อสอบล่าเอียงมากกว่า จะมีข้อสอบที่ถูกระบุว่าไม่ล่าเอียงผิดพลาด (FN) มากกว่า แต่จะเกิดขึ้นกับการทดสอบด้วย Z มากกว่า

กรณีที่จำนวนผู้สอบเท่ากับ 100 คน พบว่าความแตกต่างในการระบุข้อสอบไม่ล่าเอียงผิดพลาด สำหรับ Z และ χ^2 ไม่แตกต่างกันมากเท่ากรณีใช้ตัวอย่าง 500 คนซึ่งมีการระบุไม่ล่าเอียงผิดพลาด สำหรับ χ^2 น้อยกว่า Z โดยเฉพาะอย่างยิ่งในกรณีที่มีข้อสอบล่าเอียงในแบบสอบมากที่สุดคือ ร้อยละ 20

Cohen และ Kim วิจัยพบว่า การที่ค่าสถิติทดสอบ χ^2 ให้ผลการตรวจสอบข้อสอบล่าเอียงดีกว่า Z ของพื้นที่ที่คิดเครื่องหมายอาจเป็นเพราะ Z (ESA) ใช้การเปรียบเทียบเฉพาะค่าความยาก (b) ในขณะที่ χ^2 ใช้ทั้งค่าความยากและอำนาจจำแนก (b และ a) แต่ไม่สามารถอธิบายความแตกต่างได้ในกรณีใช้ Z(H) ซึ่งทดสอบพื้นที่ที่ไม่คิดเครื่องหมายเพราะ Z(H) ก็ใช้การทดสอบค่าพารามิเตอร์ทั้ง 2 ตัวเช่นกัน

5.2.2 การตรวจค้นข้อสอบล่าเอียงจากข้อมูลจริง

การตรวจค้นข้อสอบล่าเอียงจากข้อมูลจริงเท่าที่ผ่านมาใช้ข้อมูลทั้งที่ได้จากแบบสอบมาตรฐานและแบบสอบที่ผู้วิจัยแต่ละคนพัฒนาขึ้นเองดังนี้

Shoener (1984) ได้เปรียบเทียบการตรวจค้นความล่าเอียงของข้อสอบด้วยวิธีการทางสถิติ ได้แก่ IRT 3 พารามิเตอร์ ซึ่งใช้การทดสอบ χ^2 และการตัดสินใจด้วยผู้ตัดสินจำนวน 24 คน โดยใช้แบบสอบอิงเกณฑ์กับนักเรียนชั้นมัธยมศึกษาชาย และหญิงจำนวน 1,064 คน ที่มีภูมิหลังทางวัฒนธรรมต่าง ๆ กัน การตัดสินใจใช้ผู้ตัดสินที่มีความรู้ทางด้านคณิตศาสตร์และหลักสูตรคณิตศาสตร์เป็นอย่างดีและใช้แบบฟอร์มการให้คะแนนเป็นช่วง (rating form) ผู้ตัดสินเป็นผิวดำ 8 คน ผิวดำขาว 8 คน และเชื้อสายสเปน 8 คน เป็นชาย และหญิงเท่า ๆ กัน การวิเคราะห์ได้พิจารณา ถึงความสอดคล้องของการตัดสินใจของกลุ่มผู้ตัดสิน 3 กลุ่ม และความสอดคล้องระหว่างวิธีการทางสถิติและการตัดสินใจของผู้ตัดสิน โดยใช้ค่าสถิติของ Kappa ผลจากการวิเคราะห์พบว่าความสอดคล้องระหว่างการตัดสินใจของผู้เชี่ยวชาญต่างกลุ่มไม่มีนัยสำคัญสำหรับทั้งความล่าเอียงด้านวัฒนธรรมและทางเพศ พบความสอดคล้องกันอย่างมีนัยสำคัญในตัวบ่งชี้บางตัวที่คำนวณขึ้นเพื่อรวมการให้คะแนนของผู้ตัดสิน แต่ไม่พบความสอดคล้องระหว่างวิธีการตรวจค้นความล่าเอียงทั้ง 2 วิธี การตัดข้อกระทงที่พบทางสถิติว่าล่าเอียงออกไปไม่ทำให้ลำดับของคนในกลุ่มย่อยและรวมทั้งหมดเปลี่ยนไป และทั้งแบบสอบฉบับเดิมและแบบสอบที่ตัดข้อกระทงที่ล่าเอียงออกแล้วมีความสัมพันธ์อย่างมีนัยสำคัญกับแบบสอบวิธีคณิตศาสตร์อิงกลุ่มที่เป็นแบบสอบมาตรฐาน

Subkoviak Mack Ironson และ Craig (1984) ได้เปรียบเทียบวิธีการตรวจค้นความล่าเอียงของข้อสอบ 3 วิธีคือ (1) ICC ชนิด 3 พารามิเตอร์ (2) ค่าความยากแปลง (TID) และ (3) ไคสแควร์ซึ่งมี 2 วิธี ได้แก่ วิธีที่เสนอโดย Scheuneman (1979) พิจารณาเฉพาะตัวเลือกถูก (CHIS) และที่เสนอโดย Camilli (1979) พิจารณาทั้งตัวเลือกถูกและผิด (CHIC)

วิธีการแปลงค่าความยาก TID ทำโดยคำนวณค่าความยากของข้อสอบตามกลุ่มผิวดำและผิวขาว จากนั้นแปลงค่าความยาก (p) เป็นค่า z ซึ่งมีค่าเท่ากับเปอร์เซ็นต์ไทล์ที่ $1-p$ ในการกระจายแบบปกติ แล้วแปลงค่า z เป็นค่าเคลต้าโดยที่ $= 4Z + 13$ ค่าระยะทางระหว่างจุดจบของค่า z ทั้งสองกลุ่ม และเส้นแกนหลักเป็นค่าแสดงความล่าเอียงและจะใช้เครื่องหมายบวก ถ้าเป็นการล่าเอียงเข้าข้างผิวดำ และใช้ค่าลบถ้าเข้าข้างผิวขาว

วิธี CHIS จะแบ่งคะแนนรวมเป็น 5 ช่วง คำนวณความถี่คาดหวัง (E) ในการตอบถูกในแต่ละช่วงตามกลุ่ม แล้วเปรียบเทียบความถี่คาดหวัง กับความถี่จากการสังเกต (O) ด้วยค่าสถิติไคสแควร์ รวมทุกช่วงคะแนนค่าที่ได้จะเป็นการวัดความล่าเอียงโดยไม่มี

เครื่องหมาย การวัดโดยคิดเครื่องหมายหรือมีทิศทางทำโดย ให้เครื่องหมายบวกกรณีล่าเอียงเข้าข้างกลุ่มผิวดำ และลบเมื่อเข้าข้างกลุ่มผิวขาว โดยให้เครื่องหมายในแต่ละช่วงคะแนน

วิธี CHIC ก็เช่นเดียวกับ CHIS แต่จะวิเคราะห์ตัวเลือกผิดด้วย เครื่องมือที่ใช้เป็นแบบสอบเลือกค่าเหมือนแบบเลือกตอบ 4 ตัวเลือกจำนวน 50 ข้อ มีอยู่ 10 ข้อ ที่เป็นข้อสอบวัดค่าสแลงสำหรับผิวดำอีก 40 ข้อ เป็นคำศัพท์มาตรฐาน ให้รหัส 10 ข้อแรก เป็น 1 คือข้อที่ล่าเอียงอีก 40 ข้อให้รหัส 0 เป็นข้อไม่ล่าเอียงผลการวิเคราะห์มีดังนี้

ค่าสหสัมพันธ์ระหว่างผลการตรวจค้นของแต่ละวิธีการ กับข้อสอบที่สร้างขึ้นให้รหัสว่าล่าเอียงหรือไม่ล่าเอียง พบว่าค่าสหสัมพันธ์ระหว่าง CHIS และ CHIC มีค่าสูงสุด (ชนิดไม่คิดเครื่องหมาย = .978 ชนิดคิดเครื่องหมาย = .970) รองลงมาคือค่าสหสัมพันธ์ระหว่าง TID และ CHIC ถ้ามีจำนวนตัวอย่างน้อย เช่น 300 คน หรือน้อยกว่า แทนที่จะใช้ IRT ควรเลือกค่า z และทำให้ค่าถูกต้องใกล้เคียงกับ z คือ ไคสแควร์ส่วนค่าเฉลี่ยนั้น ในกรณีสุดโต่งจะมีการชี้ถึงความล่าเอียงตรงข้ามกับความเป็นจริง ดัชนีเคลด้าที่เหลือ (residual delta index) ซึ่งได้ลองนำมาศึกษานั้นพบว่าใช้ได้เกือบเท่า ๆ ไคสแควร์ ซึ่งควรมีการแก้ไขต่อไปเพื่อนำมาใช้แทน ไคสแควร์ได้

Shepard Camilli และ Williams (1985) ได้ศึกษาความตรงของเทคนิคการตรวจค้นความล่าเอียงด้วยวิธีไคสแควร์ ดัชนีเคลด้าของแองกอฟ และทฤษฎีการตอบข้อสอบแบบเติม โดยใช้ทั้งข้อมูลจริงและข้อมูลจำลองพบว่าผลจากการศึกษาด้วยข้อมูล 2 ชนิด มีความสอดคล้องกัน คือ พบว่าทฤษฎีการตอบข้อสอบแบบเติมเป็นวิธีที่ดีที่สุด โดยวิเคราะห์จากทั้งค่าสัมประสิทธิ์สหสัมพันธ์ และร้อยละของความสอดคล้อง ไคสแควร์ให้ความถูกต้องใกล้เคียงกัน แต่เคลด้าใช้ได้ไม่พอเพียง และในกรณีที่ล่าเอียงมาก ๆ ยังระบุได้ตรงข้ามกับความ เป็นจริง เช่นระบุว่าเข้าข้างกลุ่มผิวดำ ทั้ง ๆ ที่จริงเป็นข้อล่าเอียงเข้าข้างผิวขาว แต่อย่างไรก็ตามดัชนีดัดแปลงของ Angoff (Modified Angoff Index) ซึ่งนำค่าความแตกต่างของค่าความยาก มาดกอบบนค่าสหสัมพันธ์แบบพอยท์ไบซีเรียลของข้อกระทงและใช้ค่าที่เหลือเป็นดัชนีให้การตรวจค้นดีเกือบเท่ากับไคสแควร์ แต่ไม่ควรระวังถ้าจะใช้เพื่อการตัดข้อกระทงทั้ง

Hambleton และคณะ (1986) ได้เปรียบเทียบวิธีการตรวจค้นความล่าเอียงของข้อสอบ 4 วิธี ได้แก่ วิธีของ Mantel-Haenszel วิธีการลงจุดค่าความยาก วิธีค่าความแตกต่างของค่าเฉลี่ยกำลังสอง (the root meansquared difference) และวิธีของพินท์รวม ซึ่ง 2 วิธีหลังนี้เป็นวิธีการที่ใช้ทฤษฎีการตอบข้อสอบ

แบบสอบที่ใช้ในการศึกษา เป็นแบบสอบความสามารถในการอ่านของ คลีฟแลนด์มีข้อสอบจำนวน 75 ข้อ โดยใช้กลุ่มตัวอย่างชาย 451 คน หญิง 486 คน คะแนนจุดตัดในการแปลความหมายค่าสถิติที่แสดงความล่าเอียงได้จากการจำลองข้อมูลขึ้น ผลปรากฏว่าวิธีทั้ง 4 วิธี ให้ผลการตรวจค้นข้อสอบล่าเอียงใกล้เคียงกัน ปัญหาด้านวิธีวิทยาอยู่ที่ความไม่แจ่มชัดใน

การกำหนดจุดตัด ความคลาดเคลื่อนชนิดที่หนึ่ง และการประมาณค่าพารามิเตอร์ที่ไม่ดี ข้อค้นพบเด่น ๆ อยู่ที่ความสำคัญของการเลือกใช้ช่วงคะแนนความสามารถในการวัดความล่าเอียง และพบว่า วิธี MH เป็นวิธีเลือกใช้แทนวิธีทฤษฎีการตอบข้อสอบได้อย่างประหยัดกว่าทั้งเงินและเวลา

Doolittle และ Cleary (1987) ได้ตรวจค้นความล่าเอียงระหว่างเพศของข้อสอบในแบบวัดผลสัมฤทธิ์ทางคณิตศาสตร์โดยใช้ตัวอย่างจำนวน 8 กลุ่ม ที่ทำแบบสอบ ACT Assessment Mathematic Test รวม 8 ฟอรม แต่ละกลุ่มมีความเท่าเทียมกันและมีจำนวนประมาณ 1,300-1,400 คน ตัวอย่างเป็นหญิงประมาณร้อยละ 55 วิเคราะห์ด้วยดัชนีของ ลินน์ และ ฮาร์นิสซ์ (Linn and Harnisch, 1981) ซึ่งเป็นโมเดลโลจิสติกแบบ 3 พารามิเตอร์ ถ้าดัชนี Z ที่ได้เป็นบวกแสดงว่าง่าย สำหรับกลุ่มเปรียบเทียบถ้าเป็นลบแสดงว่ายาก สำหรับกลุ่มเปรียบเทียบ ผลการวิเคราะห์พบว่าดัชนี Z มีค่าเป็นลบ สำหรับข้อกระทงที่วัดด้านเรขาคณิตและการให้เหตุผลเชิงพีชคณิต และเลขคณิตในทุกฟอรมแสดงว่าง่ายสำหรับชายมากกว่าหญิงที่เหลือน้อยมากเป็นบวก

ผลจากการวิเคราะห์ด้วยการวิเคราะห์ความแปรปรวน พบผลเช่นเดียวกับการคำนวณค่า Z

Perlman และคณะ (1988) ได้ศึกษาความคงที่ของวิธีการประมาณค่าความล่าเอียงของข้อสอบ 4 วิธี ได้แก่ วิธีค่าความยากแปลงวิธีค่าความยากแปลงที่นำมาตัดแปลงโดย Shepard วิธีการวิเคราะห์ค่าที่เหลือนแบบ 1 พารามิเตอร์ของ Rasch และวิธี MH กลุ่มตัวอย่างที่ใช้มี 30 กลุ่ม โดยมีขนาด 600 จนถึง 2,000 คน ซึ่งสุ่มมาจากนักเรียนชั้น 9 จำนวน 54,896 คน ซึ่งมีทั้งผิวขาว ผิวดำ และชนชาติที่พูดภาษาสเปนที่ทำแบบทดสอบทักษะความสามารถขั้นต่ำของรัฐชิคาโก ซึ่งประกอบด้วยข้อสอบเลือกตอบ 46 ข้อ จากธนาคารข้อสอบจำนวน 1,000 ข้อ ผลจากการศึกษาพบว่าความเที่ยงของตัวบ่งชี้ความล่าเอียงจะมีปัญหา ถ้าใช้จำนวนตัวอย่างน้อยกว่า 666 คน ไม่มีวิธีใดที่ให้ผลการตรวจค้นความล่าเอียงได้มากกว่า หรือน้อยกว่าวิธีอื่น ๆ อย่างคงที่ข้ามขนาดตัวอย่าง

Baeza (1989) ได้ศึกษาพฤติกรรมกรรมการตอบข้อสอบของชาวอเมริกันอินเดียนและอเมริกันคอเคเซียน เพื่อพิจารณาระดับของความล่าเอียงทั้งภายในและภายนอกข้อสอบโดยใช้ตัวแปรด้านเศรษฐกิจ สังคม และระดับคะแนนเฉลี่ยในชั้นมัธยมปลายเป็นตัวแปรจับคู่กลุ่มตัวอย่างนอกเหนือไปจากเพศ โรงเรียนที่เรียนอยู่ และปีการศึกษาที่เรียน

วิธีการที่ใช้ในการตรวจค้นความล่าเอียง ได้แก่ วิธี MH และวิธีการจับคู่กลุ่มตัวอย่างของ McNemar ตัวแปรจับคู่สำหรับ MH ได้แก่ คะแนนรวมจากแบบสอบ ส่วนตัวแปรจับคู่สำหรับ McNemar ได้แก่ ตัวแปรด้านเศรษฐกิจสังคมหรือระดับคะแนนเฉลี่ย

ผลการวิเคราะห์พบว่าวิธี MH ซึ่งใช้เกณฑ์ภายในตรวจสอบที่ค่าเฉลี่ยเพียงเล็กน้อยส่วน McNemar ซึ่งใช้เกณฑ์ภายนอกนั้นพบว่าเมื่อใช้ตัวแปรจับคู่เป็นสถานภาพทางเศรษฐกิจสังคม พบว่าข้อสอบมีระดับความลำเอียงต่ำกว่าการจับคู่ด้วยระดับคะแนนเฉลี่ย แสดงว่าระดับคะแนนเฉลี่ยของคน 2 กลุ่มนี้มีระดับของความหมายไม่เท่ากัน

Baghi และ Ferrara (1989) ได้เปรียบเทียบผลการสืบค้นความลำเอียงด้วยวิธีการแตกต่างกัน 3 วิธี ได้แก่ การใช้ทฤษฎีการตอบข้อสอบ วิธีการใช้ค่าเฉลี่ย (ค่าความยากแปลง) และวิธี MH โดยศึกษาในแง่ของผลกระทบจากขนาดของกลุ่มตัวอย่างในผลการตรวจค้นความลำเอียง และระดับของความสัมพันธ์ระหว่างค่าสถิติจากทั้ง 3 วิธีนี้

ข้อมูลที่ใช้ในการวิเคราะห์ เป็นข้อมูลที่สุ่มจากนักเรียนระดับ 9 จำนวน 50,000 คน ที่ตอบแบบสอบถามทักษะความเป็นพลเมืองดี ประกอบด้วยข้อสอบแบบเลือกตอบ 45 ข้อ ประเมินความรู้ และ ทักษะ รวม 3 ด้าน ได้แก่ด้านรัฐบาลประชาธิปไตย ด้านการเมือง และ พฤติกรรมทางการเมือง และ ด้านหลักการ สิทธิ และความรับผิดชอบ จำนวนกลุ่มตัวอย่างที่สุ่มมาใช้มี 4 ขนาด ได้แก่ ขนาด 1,000 คน 750 คน 500 คน และ 200 คน

ผลจากการวิเคราะห์ไม่พบว่า มีข้อสอบข้อใดที่มีความลำเอียงทั้งใน การเปรียบเทียบระหว่างผู้สอบพิวค่าและพิวขาว และผู้สอบหญิง และชาย การแสดงค่าประมาณความยากด้วยแผนภูมิเส้นแสดงความสัมพันธ์ระหว่างค่าความยากต่างกลุ่มอยู่ในรูปเส้นตรงเกือบสมบูรณ์ ความสอดคล้องกันระหว่างวิธีการตรวจค้นความลำเอียงของ Rasch และค่าความยากแปลงมีความสัมพันธ์กันสูงมากทุกขนาดของกลุ่มตัวอย่างจากค่าสัมประสิทธิ์สหสัมพันธ์แบบลำดับ และ ทั้ง 2 วิธี มีความสอดคล้องสูงสุดกับวิธีการตอบข้อสอบ 3 พารามิเตอร์ ในการตรวจค้นความลำเอียงและไม่ลำเอียงของข้อสอบ

Hambleton และ Rogers (1989) ได้ศึกษาเปรียบเทียบวิธีการตรวจค้นข้อสอบลำเอียง 2 วิธี ได้แก่ MH และ IRT โดยใช้ข้อมูลการตอบแบบสอบ New Mexico High School Proficiency Exam (NMHSPE) จำนวน 150 ข้อ ของนักเรียนชาวอเมริกัน พิวขาว 8,000 คน และชาวอเมริกันพื้นเมือง 2,600 คน จากผู้สอบทั้งหมด 23,000 คน

NMHSPE เป็นแบบสอบชนิดเลือกตอบ 4 ตัวเลือก วัดทักษะการดำเนินชีวิต 5 ด้าน คือ (1) ความรู้เกี่ยวกับแหล่งทรัพยากรของชุมชน (2) ความรู้ด้านผู้บริโภค (3) ความรู้เกี่ยวกับรัฐบาลและกฎหมาย (4) ความรู้ด้านสุขภาพทางกายและจิต และ (5) ความรู้ด้านอาชีพ

กลุ่มตัวอย่างที่นำมาศึกษา ซึ่งได้แก่ ชาวอเมริกันผิวขาวและอเมริกันพื้นเมืองนั้นได้สุ่มมาศึกษาอย่างละ 2,000 คน และแบ่งเป็น 2 กลุ่ม ๆ ละ 1,000 คน ตามลำดับคือ การศึกษาความล่าเอียง 2 คู่ เช่นนี้เพื่อพิจารณาถึงความคงที่ในการระบุข้อสอบล่าเอียง

ในการศึกษาความล่าเอียง ผู้วิจัยได้ใช้ข้อสอบเพียง 75 ข้อ โดยตัดข้อสอบที่ง่ายมาก ($p > .70$) และค่าอำนาจจำแนกต่ำ ($r < .10$) วิธี MH วิเคราะห์โดยใช้ 76 ระดับคะแนน ส่วนวิธีทฤษฎีการตอบข้อสอบใช้การวัดพื้นที่ระหว่างโค้งลักษณะข้อสอบ (ICC) พิสัยของความสามารถที่นำมาใช้ในการคำนวณพื้นที่กำหนดค่าใช้ 2 ค่าเบี่ยงเบนมาตรฐานสูงกว่าความสามารถเฉลี่ยของชาวอเมริกันผิวขาว และ 2 ค่าเบี่ยงเบนมาตรฐานต่ำกว่าความสามารถเฉลี่ยของชาวอเมริกันพื้นเมือง ความคงที่ของการตรวจสอบความล่าเอียง ได้แก่ ร้อยละของข้อสอบที่ถูกระบุล่าเอียงหรือไม่ล่าเอียงซ้ำกันในการวิเคราะห์ 2 ครั้ง

เกณฑ์ตัดสินความล่าเอียง สำหรับวิธี IRT - 3 พารามิเตอร์ได้จากการวิเคราะห์ดัชนี SA ระหว่างผู้สอบที่เป็นอเมริกันพื้นเมือง 2 กลุ่ม และใช้ค่าสูงสุดที่ได้เป็นเกณฑ์ซึ่งมีค่า = .468 และด้วยวิธีเดียวกันเกณฑ์ตัดสินความล่าเอียง สำหรับวิธี MH คือ $MH-\chi^2 = 6.64$

ผลการวิเคราะห์พบว่า วิธี MH มีความสอดคล้องของการตรวจสอบร้อยละ 80 ส่วน ICC ร้อยละ 73 และพบว่า ร้อยละ 47 ของข้อสอบที่ถูกระบุล่าเอียงในการวิเคราะห์คู่ที่ 1 เป็นข้อสอบล่าเอียงในการวิเคราะห์คู่ที่ 2 ด้วย สำหรับ MH และร้อยละ 61 สำหรับ ICC ในทางกลับกันพบว่า ร้อยละ 64 ของข้อสอบล่าเอียงในการวิเคราะห์คู่ที่ 2 เป็นข้อสอบล่าเอียงในการวิเคราะห์คู่ที่ 1 สำหรับ MH และร้อยละ 56 สำหรับ ICC ผลดังกล่าวแสดงให้เห็นว่าความคงที่ของการวิเคราะห์ทั้ง MH และ ICC อยู่ในระดับปานกลางซึ่งอาจเป็นเพราะจำนวนตัวอย่างค่อนข้างน้อย (MH มีข้อสอบที่ระบุล่าเอียงในการวิเคราะห์ทั้ง 2 ครั้ง 7 ข้อ ICC 14 ข้อ และพื้นที่แสดงความแตกต่างของ ICC มีค่าระหว่าง -2.7 และ 1.5)

สรุปได้ว่า ทั้ง MH และ ICC มีความไม่เที่ยงในระดับหนึ่งในการใช้ตรวจสอบข้อสอบล่าเอียงโดยที่ MH มีความคงที่ประมาณร้อยละ 80 และ ICC ประมาณ 73

Reynold (1989) ได้ใช้วิธีการ 4 วิธีในการตรวจค้น

ความล่าเอียงของข้อสอบ CTBS/S ฉบับภาษาสเปน โดยใช้ข้อมูลจากนักเรียน 895 และ 911 คน วิธีการที่ใช้ในการตรวจค้นมี 4 วิธี ได้แก่

1. ค่าความยากแปลงของแองกอล์ฟ
2. การทดสอบไคสแควร์
3. เปรียบเทียบความแตกต่างในค่าความยาก และการทดสอบความเหมาะสมด้วยโมเดลของ Rasch
4. ทฤษฎีการตอบข้อสอบ 3 พารามิเตอร์ (พิจารณาจากพื้นที่และความแตกต่างในค่าสถิติที่ใช้ทดสอบความเหมาะสม)

ตรวจสอบความตรงเชิงโครงสร้างในการวัดความล่าเอียงของข้อสอบ โดยการพิจารณาค่าสหสัมพันธ์ระหว่างวิธีต่าง ๆ ร้อยละของความสอดคล้องระหว่างวิธีต่าง ๆ ที่ใช้และการวิเคราะห์ตัวประกอบ

ผลจากการวิเคราะห์พบว่าเคลด้า และไคสแควร์ มีความสัมพันธ์เกือบจะสมบูรณ์วิธีที่ 1 และ 2 กับวิธีที่ 3 และ 4 มีความสัมพันธ์กันน้อย สนับสนุนโดยผลจากการวิเคราะห์ร้อยละของความสอดคล้อง การวิเคราะห์ตัวประกอบแสดงว่า เคลด้า ไคสแควร์ และ IRT แบบ 3 พารามิเตอร์ วัดโครงสร้าง (construct) เดียวกัน คือ ความล่าเอียงด้านภาษา และความสัมพันธ์ระหว่างเคลด้าและไคสแควร์นั้นเห็นได้ชัดแม้แต่ในกรณีที่มีความล่าเอียงน้อย

Baghi และ Ferrara (1990) ได้เปรียบเทียบผลจากการสืบค้นความล่าเอียงของข้อสอบด้วยทฤษฎีการตอบข้อสอบแบบ 3 พารามิเตอร์ (IRT) และสถิติไคสแควร์แบบ MH (MHCS) โดยใช้ข้อมูลจากการบริหารแบบสอบทักษะความเป็นพลเมืองดีของรัฐแมรี่แลนด์ ให้แก่นักเรียนระดับ 9 จำนวน 50,000 คน ในเดือนมกราคม-กุมภาพันธ์ 1988 จำนวนตัวอย่างที่ใช้ในการสืบค้นด้วย MHCS มี 4 ขนาด ได้แก่ 1,000 คน 750 คน 500 คน และ 200 คน ส่วนจำนวนตัวอย่างในการสืบค้นด้วย IRT ใช้กลุ่มละ 1,000 คน ในแต่ละกลุ่มเปรียบเทียบ การวิเคราะห์ประกอบด้วย (1) การพิจารณาความคงที่ของค่าสถิติ MHCS เมื่อใช้กลุ่มตัวอย่างขนาดต่าง ๆ (2) ความคงที่ของดัชนีความล่าเอียง (MH Alpha) ในวิธี MHCS ซ้ำในช่วงคะแนน (3) ความสัมพันธ์ระหว่างดัชนีความล่าเอียงแบบ IRT และค่า MH Alpha และ (4) ความสอดคล้องระหว่างวิธีทั้งหมดในการสืบค้นข้อกระทงที่ล่าเอียง

ผลปรากฏว่าเมื่อใช้วิธี IRT สืบค้นความล่าเอียงของข้อสอบระหว่างกลุ่มหญิง และชายพบข้อสอบล่าเอียง 4 ข้อ และพบข้อสอบล่าเอียง 3 ข้อในการเปรียบเทียบระหว่างผู้สอบพิวชาวและพิวค่า ซึ่งมีความน่าจะเป็นในการตอบข้อสอบได้ถูกต้องไม่เท่ากันอย่างน้อยสำคัญ และในขณะที่ผลจากการหาค่าสหสัมพันธ์ระหว่างดัชนีทั้งสองตัวแสดงว่าไม่สอดคล้องกันในการบ่งชี้ข้อสอบที่ล่าเอียงนั้น ความสอดคล้องในผลรวมแสดงว่าเทคนิค MHCS เป็นวิธีที่ใช้แทน IRT แบบ 3 พารามิเตอร์ได้อย่างเพียงพอถ้าจำนวนตัวอย่างมีขนาดอย่างน้อย 750 คน

Sudweeks และ Tolman (1990) ได้เปรียบเทียบผลจากการใช้วิธี MH และการตัดสินใจด้วยผู้เชี่ยวชาญในการตรวจค้นความล่าช้าของข้อสอบวิชาวิทยาศาสตร์ระหว่างนักเรียนชายและหญิงระดับ 5 จำนวน 926 คน เครื่องมือที่เก็บรวบรวมข้อมูลเป็นแบบสอบอิงเกณฑ์ชนิดเลือกตอบที่สร้างขึ้นตามมาตรฐาน และวัตถุประสงค์ของหลักสูตรวิทยาศาสตร์ร่วมในรัฐยูทาห์ ผลจากการวิเคราะห์พบว่า วิธีการทั้งสองให้ผลต่างกัน และข้อสอบหลายข้อที่ถูกระบุด้วยวิธีการ 2 วิธีนี้ว่า ล่าช้าจะเป็นข้อสอบที่ยาก

Cohen Kim และ Subkoviak (1991) ได้ศึกษาอิทธิพลของการทราบการกระจายในการตรวจค้นการทำหน้าที่เบี่ยงเบนของข้อสอบสำหรับกลุ่มผู้สอบผิดค่า และผู้ชวที่สร้างขึ้นอย่างมีเจตนาให้ล่าช้า โดยใช้โปรแกรม LOGIST ซึ่งใช้วิธี Joint Maximum Likelihood Estimation (JMLE) และโปรแกรม BILOG โดยที่โปรแกรม BILOG ใช้การประมาณค่าพารามิเตอร์ 3 สถานการณ์ คือ

1. MMLE (Marginal Maximum Likelihood Estimation) เป็น BILOG ที่ไม่มีการกำหนดค่าพารามิเตอร์ ค่าอำนาจจำแนก (a) จะอยู่ในช่วง $0 - M_u = 1.13$
 $a = .604$ เป็น BILOG A

2. MMAP (Marginal Maximum a Posteriori Estimation) กำหนดค่า a ในรูป lognormal ค่าเฉลี่ย = 0 (ไม่มี FLOAT option) เป็น BILOG B

3. MMAP กำหนดค่า a ในรูป lognormal เช่นกัน แต่มีข้อตกลงว่าการกระจายของค่า (ค่าแปลงของ a) มีการกระจายเป็นปกติ มีค่าความเบี่ยงเบนมาตรฐาน = 0.5 (มี FLOAT option) เป็น BILOG C

ส่วนการวัดความแตกต่างของพื้นที่ระหว่างโค้งการตอบข้อสอบนั้น ทำ 2 ลักษณะ คือ อย่างคิดเครื่องหมาย (SA) และไม่คิดเครื่องหมาย (UA) การประมาณค่าพารามิเตอร์ใช้แบบ 2 พารามิเตอร์ โดยการกำหนดค่า c ให้เป็นค่าคงที่

พบจากค่าสัมประสิทธิ์สหสัมพันธ์แบบเพียร์สันว่า การวัดพื้นที่ โดยใช้ BILOG B และ BILOG C มีความสัมพันธ์กันสูงมาก ทั้งในดัชนีแบบคิดเครื่องหมาย และไม่มีเครื่องหมาย ส่วนระหว่าง LOGIST และ BILOG ก็สูงเช่นกัน แต่ความคงที่ในรูปแบบของความคล้ายคลึงจะน้อยกว่า ส่วนความคล้ายคลึงกันระหว่าง LOGIST กับ BILOG B และ BILOG C ลดลงกว่าเล็กน้อย เป็นผลที่สอดคล้องกับการศึกษาที่คนอื่น ๆ เคยพบ ในการวัดพื้นที่ของความเบี่ยงเบนของข้อสอบ

BILOG B และ C จะดีกว่าวิธีอื่น ๆ และมีความไวในการตรวจค้นมากเมื่อค่า b อยู่ในช่วง +5 มากกว่าในช่วง +3 ยกเว้น BILOG B ที่เท่า ๆ กันทั้ง 2 ช่วง ส่วน MMLE จะดีกว่า JMLE มาก ในการตรวจค้นความเบี่ยงเบน ในกรณีที่ค่าสัมบูรณ์ของ b อยู่

ในช่วง +3 และ ต่ำกว่าเล็กน้อยในกรณีค่าสัมบูรณ์ของ b น้อยกว่า หรือเท่ากับ 5 ในการวัดแบบมีเครื่องหมาย จำนวนตัวอย่างที่ใช้ก็ไม่มีข้อผิดพลาดนักต่อการตรวจค้นความเบี่ยงเบนของข้อสอบสรุปแล้ว การใช้ BILOG แบบ B และ C ในการตรวจค้นความลำเอียงของข้อสอบที่สร้างขึ้นมาโดยทราบค่า หรือ กำหนดค่าความลำเอียงไว้ก่อนแล้วดีกว่าวิธีอื่น

Kim และ Cohen (1991) ได้ศึกษาความลำเอียงของข้อสอบที่สร้างขึ้นอย่างจงใจให้ลำเอียง โดยการเปรียบเทียบผลการคำนวณพื้นที่ความแตกต่างระหว่างโค้งที่ได้จากทฤษฎีการตอบข้อสอบของกลุ่มตัวอย่างผิวขาว และผิวดำ 2 วิธี ได้แก่ การคำนวณพื้นที่ภายใต้ความสามารถที่มีขอบเขตจำกัด และที่ไม่มีขอบเขตจำกัด

จำนวนกลุ่มตัวอย่างนี้ใช้กลุ่มละ 1,000 คน สุ่มมาจากเดิมที่มีผิวขาว 1,021 คน และผิวดำ 1008 คน ทั้งนี้เพราะขีดจำกัดในการใช้โปรแกรม BILOG ด้วยไมโครคอมพิวเตอร์ เครื่องมือเก็บรวบรวมข้อมูล คือ แบบสอบคำศัพท์แบบเลือกตอบ ชนิด 4 ตัวเลือก จำนวน 50 ข้อ เป็นคำศัพท์ภาษาอังกฤษมาตรฐาน 40 ข้อ อีก 10 ข้อ เป็นคำศัพท์ที่มีความง่ายสำหรับกลุ่มผิวดำ มากกว่ากลุ่มผิวขาว

การวิเคราะห์ความเป็นเอกมิติ (Unidimension) ด้วยการวิเคราะห์ตัวประกอบหลัก (principal components analysis) พบว่า ตัวประกอบแรกที่ได้ โดยไม่มีการหมุนแกน และใช้ค่าสหสัมพันธ์แบบเตตตราคอร์ริคแสดงถึงความเป็นเอกมิติได้พอเพียงในการใช้แบบจำลองทฤษฎีการตอบข้อสอบ สำหรับกลุ่มผิวขาวโดยอธิบายความแปรปรวนได้ร้อยละ 25 ในขณะที่ค่อนข้างน้อยสำหรับผิวดำ คือ อธิบายได้ร้อยละ 16

ใช้ทฤษฎีการตอบข้อสอบแบบ 2 พารามิเตอร์และแบบ 3 พารามิเตอร์ ชนิดบังคับให้ $c_r = c_f = c = .23$

ใช้วิธีการวัดพื้นที่ 4 วิธี จากพารามิเตอร์ข้อสอบที่เทียบมาตรฐานแล้ว คือ ESA ซึ่งได้แก่ พื้นที่แบบปิดคิดเครื่องหมาย (Exact Signed Area) ได้จากการคำนวณพื้นที่ตลอดช่วงความสามารถโดยไม่จำกัดขอบเขต และ EUA คือพื้นที่แบบปิดไม่คิดเครื่องหมาย (Exact Unsigned Area) กรณีที่เป็นพื้นที่ที่คิดเครื่องหมาย กลุ่มอ้างอิง (R) คือ กลุ่มผิวดำดังนั้นค่าที่เป็นบวกของ CSA หรือ ESA จะแสดงว่าข้อสอบข้อนั้นลำเอียงเข้าข้างกลุ่มผิวดำ

ข้อกระทงที่สร้างขึ้นเพื่อให้เกิดความลำเอียงเข้าข้างผิวดำ 10 ข้อ ให้ค่าเป็น 1 อีก 40 ข้อ ที่เหลือเป็น 0 แล้วหาค่าสหสัมพันธ์ แบบ point biserial ระหว่างค่าดังกล่าวกับการวัด 4 วิธี เพื่อดูผลการตรวจค้นแต่ละวิธี และหาค่าสหสัมพันธ์ระหว่างวิธีต่าง ๆ เพื่อดูความคล้ายคลึงกัน

การประมาณค่าอัตราความคลาดเคลื่อนของการตรวจค้นความล่าเอียงกระทำโดยการมีข้อตกลงว่า การวัดพื้นที่แต่ละวิธีมีการกระจายเป็นปกติ ค่าเฉลี่ยและค่าความเบี่ยงเบนมาตรฐานของการวัดแต่ละวิธี ประมาณจากข้อสอบ 40 ข้อ กำหนดค่าวิกฤตแบบทางเดียวเป็น 2 ค่า คือ .05 และ .01 ข้อที่มีพื้นที่ใหญ่กว่าค่าวิกฤตจะเป็นข้อสอบที่ล่าเอียงการจักประเภทข้อสอบเช่นนี้ทำให้เกิดความคลาดเคลื่อน 2 ชนิด ได้แก่ ความคลาดเคลื่อนชนิดระบุว่าไม่ล่าเอียง ทั้ง ๆ ที่ล่าเอียง (FN = False Negatives) และระบุว่าล่าเอียง ทั้ง ๆ ที่ไม่ล่าเอียง (FP = False Positives) และเมื่อใช้ $p = .05$ พบว่ามี FP มากกว่าที่ .01

ผลจากการวิเคราะห์ข้อมูลพบว่า เกือบไม่มีความแตกต่างระหว่างการคำนวณพื้นที่แบบช่วงเปิด และแบบช่วงปิดในการตรวจค้นความล่าเอียงของข้อสอบ ส่วนการใช้พื้นที่แบบไม่คิดเครื่องหมายไม่ว่าจะเป็นแบบช่วงเปิดหรือปิดมีความไวต่อการตรวจค้นความล่าเอียงของข้อสอบมากกว่าแบบคิดเครื่องหมาย นอกจากนี้พบว่าการคำนวณพื้นที่ช่วงเปิดนั้นพื้นที่ชนิดมีเครื่องหมาย และไม่คิดเครื่องหมาย มีความสัมพันธ์กันมากกว่าการคำนวณพื้นที่ช่วงปิด

ในด้านของความคลาดเคลื่อนในการตรวจค้นความล่าเอียง พบว่าวิธี 3 พารามิเตอร์ โดยใช้พื้นที่แบบคิดเครื่องหมาย มีอัตราความคลาดเคลื่อนต่ำกว่าวิธีอื่น ๆ ไม่มีความคลาดเคลื่อนเลขในระดับ .01 และ มีความคลาดเคลื่อนในการระบุข้อสอบว่าล่าเอียงโดยไม่ล่าเอียงจริง จำนวน 2 ข้อ ที่ระดับ .05 ส่วนวิธี 3 พารามิเตอร์ แบบไม่คิดเครื่องหมาย มีอัตราความคลาดเคลื่อนสูงกว่าวิธีอื่น ๆ เล็กน้อย

Ryan (1991) ได้ตรวจค้นความล่าเอียงของข้อสอบด้วยวิธีการที่เสนอโดย Mantel และ Haenszel (1959) โดยใช้ตัวอย่างกลุ่มอ้างอิง 500 คน และกลุ่มเปรียบเทียบ 100 คน เพื่อตรวจสอบความคงที่ของค่าประมาณความล่าเอียงด้วย MH ข้ามกลุ่มตัวอย่างและขนาดตัวอย่างต่าง ๆ กัน และเพื่อศึกษาว่า MH มีความแกร่งต่อบริบทของข้อสอบเพียงไร โดยมีข้อสอบรวมจำนวน 40 ข้อ และข้อสอบที่ใช้สลับกันอีกจำนวน 35 ข้อ กลุ่มตัวอย่างเป็นชายผิวขาวและผิวดำจำนวน 5015 คน และ 670 คน ตามลำดับ โดยมีเกณฑ์การแยกประเภทข้อสอบจากผลการวิเคราะห์เป็น 3 ประเภท

ประเภท A ข้อสอบที่กำหนดหน้าที่ได้ตรง มีค่า $\alpha_{MH} < 1.00$ และมีค่าไม่ต่างจาก 0 อย่างมีนัยสำคัญ

ประเภท B ข้อสอบที่พอจะนำมาใช้ได้ แต่ควรเลือกข้อที่มีค่า α_{MH} ค่าคือ $1 = \alpha_{MH} < 1.5$ และไม่ต่างจาก 1 อย่างมีนัยสำคัญ

ประเภท C เลือกมาใช้ถ้าตรงตามวัตถุประสงค์เฉพาะที่ระบุไว้ อย่างมีนัยสำคัญว่าต้องการข้อสอบลักษณะนี้ได้แก่ข้อที่มีค่า $\alpha_{MH} > 1.5$ และมากกว่า 1 อย่างมีนัยสำคัญ



ขั้นตอนในการวิเคราะห์ขั้นแรกมีการวิเคราะห์ 5 แบบ ในการสร้างฐานเรียกว่า ชุด 1 ทำโดยการสุ่มตัวอย่างผิวขาว 670 คน เป็นกลุ่มเป้าหมาย ที่เหลือเป็นกลุ่มอ้างอิงแล้วสุ่มกลุ่มผิวขาวมาอีก 4 กลุ่ม เป็นกลุ่มอ้างอิง ได้แก่ P_1 P_2 P_3 P_4 และวิเคราะห์ข้อสอบรวม 40 ข้อด้วย MH จับคู่การวิเคราะห์เป็น P/P P/P_1 P/P_2 P/P_3 P/P_4 คะแนนที่ใช้แบ่งระดับความสามารถ คือ คะแนนรวมจากข้อสอบ 40 ข้อ

การวิเคราะห์ชุด 2 ทำในลักษณะเดียวกับ ชุด 1 แต่ใช้กลุ่มเป้าหมายเป็นผิวดำวิเคราะห์ด้วย MH 5 คู่ ได้แก่ BW BW_1 BW_2 BW_3 BW_4 เป็นการวิเคราะห์เพื่อดูความคงที่ของค่า α_{MH} ข้ามกลุ่มและขนาดตัวอย่าง

ชุด 3 เป็นการวิเคราะห์ผลของบริบทของข้อสอบ BW_1 BW_2 BW_3 และ BW_4 โดยใช้ข้อสอบ 75 ข้อจากฟอร์ม 1 ฟอร์ม 2 ฟอร์ม 3 ฟอร์ม 4 ใช้คะแนนรวมเป็นเกณฑ์แบ่งช่วงความสามารถ (แต่ละฟอร์มมีข้อสอบรวม 40 ข้อ)

และชุด 4 เป็นการวิเคราะห์เช่นเดียวกับชุด 3 แต่ใช้ข้อสอบ 4 ฟอร์มที่ไม่มีข้อสอบรวม (มีข้อสอบ 35 ข้อ) และใช้คะแนนรวมจาก 35 ข้อ เป็นเกณฑ์แบ่งระดับความสามารถ

ผลการวิเคราะห์พบว่า ความลำเอียงสำหรับกลุ่มผิวขาวด้วยกัน 5 คู่ นั้นมีความสัมพันธ์กันต่ำมาก ไม่มีความแตกต่างอันเนื่องมาจากความลำเอียง แต่มีความแตกต่างจากความคลาดเคลื่อนในการสุ่มตัวอย่างที่มีความซ้ำซ้อนกันมากกว่า จึงอาจเป็นการลำเอียงแบบไม่แท้ก็ได้

เกณฑ์การแบ่งช่วงระดับความสามารถ ซึ่งเป็นบริบทที่ศึกษาไม่มีผลต่อค่า α_{MH} ไม่ว่าจะใช้เกณฑ์คะแนนรวม 35 ข้อหรือ 75 ข้อ ก็ได้ค่าสหสัมพันธ์ของการเปรียบเทียบแต่ละคู่เท่า ๆ กัน

ผลจากการจำแนกประเภทข้อสอบมีเพียง 1 ข้อ ที่มีความสอดคล้องกับทุกการวิเคราะห์ คือ เป็นประเภท B ทุกครั้ง นอกนั้นมีความแปรปรวนไปตามการวิเคราะห์ Harris และ Carlton (1993) ได้ศึกษาลักษณะกว้าง ๆ ของข้อสอบเพื่อชี้ให้เห็นจุดเด่นและจุดด้อย สำหรับผู้สอบหญิงและชายโดยแบ่งลักษณะที่ศึกษาเป็นประเด็นที่วัด รูปแบบและเนื้อหาของข้อสอบ

ข้อสอบที่นำมาศึกษาเป็นข้อมูลการตอบข้อสอบในแบบสอบ SAT-M จำนวน 6 ฟอร์ม ซึ่งเป็นตอนที่ทดสอบเฉพาะคณิตศาสตร์ ประกอบด้วยข้อสอบ 60 ข้อ เป็นการแก้ปัญหาทางคณิตศาสตร์ทั่วไป 40 ข้อ และการเปรียบเทียบปริมาณตัวเลข 20 ข้อ ผู้สอบเป็น

นักเรียนชั้นมัธยมทั้งต้นและปลาย เป็นชาย 181,228 คน หญิง 198,668 คน ทุกคนรายงานด้วยตนเองว่า ภาษาอังกฤษเป็นภาษาที่แต่ละคนใช้ได้ดีที่สุด จำนวนกลุ่มตัวอย่างที่ใช้ใน SAT-M พอร์ม 1-6 เป็นชาย 6,329-74,283 คน หญิง 6,712-83,945 คน

ผู้วิจัย ได้ระบุลักษณะเดิมของข้อสอบโดยใช้อ้อมลูจากการศึกษาในครั้งก่อน ๆ ของผู้อื่น เช่น Bleistein, 1986; Donlon, Ekstrom และ Lockheed, 1979; Doolittle และ Cleary, 1987; McPeck และ Wild, 1987; Schmitt และ Dorans, 1987; Wendler และ Carlton, 1987 เป็นต้น สิ่งที่ระบุได้แก่ ประเด็นที่วัด (นั่นคือ เนื้อหาหลักได้แก่ เลขคณิต พีชคณิต เรขาคณิต หรือทั่ว ๆ ไป) รูปแบบข้อสอบ (มีการใช้ตัวเลือกที่ว่า "ไม่สามารถตัดสินได้หรือไม่") และเนื้อหาของข้อสอบ (มีการอ้างอิงด้านเพศหรือไม่)

ตรวจสอบด้วย วิธี MH และพิจารณาความแตกต่างของค่าเคลต้า ซึ่งใช้แสดงถึงความยากของข้อสอบใน SAT อยู่แล้ว โดยที่ค่าเคลต้า = 1.00 แสดงว่าข้อสอบมีความยาก สำหรับกลุ่มหนึ่งมากกว่าอีกกลุ่มเท่ากับ 1 เคลต้า หรือเท่ากับความแตกต่างร้อยละ 10 โดยประมาณ

ค่าเคลต้าที่เป็นลบแสดงว่ายากสำหรับหญิงมากกว่าชาย ค่าบวกแสดงว่ายากสำหรับชายมากกว่าหญิง

ผลการวิเคราะห์โดยจับคู่ชายและหญิงในด้านความสามารถ ซึ่งได้แก่คะแนนรวมจากแบบสอบแล้ว ได้ค่าเคลต้ามี่ค่า -1.86 ถึง 1.27 ค่าเฉลี่ย -.01 และค่าความเบี่ยงเบนมาตรฐาน = .51 และใช้การวิเคราะห์ความแปรปรวนทดสอบนัยสำคัญของความแตกต่างเฉลี่ยระหว่างข้อที่มีและไม่มีลักษณะดังที่ระบุไว้ พบว่า โดยทั่ว ๆ ไป ถ้าเนื้อหาหลักเป็นเรขาคณิตหรือเลขคณิตแล้วชายจะทำได้ดีกว่าหญิง แต่กรณีเป็นเนื้อหาทั่ว ๆ ไป (เช่น โครงสร้างของระบบจำนวน เซต ฯลฯ) หญิงจะทำได้ดีกว่าชาย ถ้าเป็นเลขคณิตและพีชคณิต หญิงจะทำได้ดีกว่าชาย แต่ถ้าเป็นเลขคณิตและเรขาคณิตแล้ว ชายจะทำได้ดีกว่าหญิง การพบว่าหญิงทำได้ดีกว่าในกรณีของพีชคณิตและชายจะทำได้ดีกว่าในเรขาคณิต (โดยที่หญิงและชายมีความสามารถรวมเท่ากัน) นั้นสอดคล้องกับการค้นพบของ Doolittle และ Cleary (1987)

นอกจากนั้นพบว่า ชายจะทำได้ดีกว่าในข้อสอบที่ต้องการกระบวนการทางสมองในระดับสูงกว่า ในข้อที่เป็นการประยุกต์นำมาใช้ในชีวิตประจำวัน ส่วนหญิงทำได้ดีกว่ากรณีที่มีตัวแปร (เช่น x , a , b) ในปัญหาและในตัวเลือก กรณีเป็นปัญหาตรงไปตรงมาจากหลักสูตรหรือหนังสือเรียนโดยไม่มีการนำมาประยุกต์ และในกรณีมีการอ้างอิงถึงบุคคลโดยไม่มีการแสดงถึงเพศหรือสถานภาพทางเชื้อชาติ

Raju และคณะ (1993) ได้ศึกษาเปรียบเทียบผลการทดสอบค่า Z ของ Raju ทั้งที่ได้จากพื้นที่ชนิดมีและไม่คิดเครื่องหมาย การทดสอบค่า χ^2 ของ Lord (ใช้โมเดล 2 พารามิเตอร์) และการทดสอบ MH- χ^2 โดยใช้เกณฑ์ตัดสิน ดังนี้

$$\chi^2 \text{ ของ Lord} = 13.80$$

Z มีค่าระหว่าง ± 3.30 และ

$$\text{MH-}\chi^2 \text{ มีค่า} = 10.83$$

ข้อมูลที่ใช้ในการตรวจค้นความล่าเอียงของข้อสอบ คือ ผลการตอบแบบสอบ Gater-Mac Ginitie Reading Tests (GMRT) เฉพาะแบบสอบย่อยที่วัดคำศัพท์ 45 ข้อ เป็นแบบเลือกตอบ 5 ตัวเลือก กับนักเรียนระดับ 10 และ 12 ในปี 1987 จำนวน 839 คน ในจำนวนนี้มีนักเรียนผิวดำ 245 คน ผิวดำ 436 เป็นหญิง 440 และชาย 399 คน

ผลจากการศึกษา ในกรณีการตรวจสอบความล่าเอียงระหว่างนักเรียนผิวดำและผิวดำ พบว่า มีข้อสอบล่าเอียง 1 ข้อ สำหรับการทดสอบด้วย χ^2 และ Z ทั้งพื้นที่ชนิดคิดเครื่องหมายและไม่คิดเครื่องหมาย ได้แก่ ข้อ 41 ส่วน MH พบข้อสอบล่าเอียง 2 ข้อ ได้แก่ ข้อ 14 และ 27 ทั้ง 3 ข้อ เป็นข้อที่มีความแตกต่างกันมากสำหรับผู้สอบผิวดำและผิวดำ

ส่วนผลการตรวจสอบความล่าเอียงระหว่างชายและหญิง 3 วิธีแรก พบข้อสอบล่าเอียงซ้ำกัน 4 ข้อ ส่วน MH พบ 5 ข้อ ซ้ำกับ 3 วิธีแรก 4 ข้อ (ร้อยละ 80) ได้แก่ ข้อ 2, 18, 23, 33 และข้อที่ไม่ซ้ำกับ 3 วิธีแรกคือ ข้อ 41 ทุกข้อมีค่าความชุกสูงกว่าข้อที่ไม่ล่าเอียง

การพบว่า MH ค้นพบข้อสอบล่าเอียงได้มากกว่าวิธีที่ใช้ทฤษฎีการตอบข้อสอบเช่นนี้ค้านกับการค้นพบของ สแวมมินาซาน และโรเจอร์ส (1990) ที่พบว่าวิธีที่ใช้ทฤษฎีการตอบข้อสอบค้นพบข้อสอบล่าเอียงได้มากกว่า MH

ข้อควรสังเกต คือ การทดสอบ Z ของราชูอาจให้ผลถูกต้องกว่าการกำหนดเกณฑ์ตัดสินตามอำเภอใจ (arbitrary) สำหรับพื้นที่ความแตกต่างระหว่างโค้งเพราะค่าความคลาดเคลื่อนมาตรฐานในการวัดเป็นตัวหนึ่งเข้ามาเกี่ยวข้องทำให้ข้อสอบข้อที่มีความแตกต่างของพื้นที่น้อยกว่าเป็นข้อสอบล่าเอียง และข้อสอบที่มีค่าความแตกต่างของพื้นที่มากกว่าเป็นข้อสอบไม่ล่าเอียงได้

Roussos และ Stout (1993) ได้ศึกษาผลของการใช้จำนวนตัวอย่างขนาดเล็ก ความแตกต่างกันของค่าพารามิเตอร์ข้อสอบรวมทั้งความแตกต่างในค่าเฉลี่ยของความสามารถของกลุ่มอ้างอิง และกลุ่มเปรียบเทียบเพื่อศึกษาถึงผลที่ต่อความคลาดเคลื่อนชนิดที่ 1 (Type I error) ซึ่งเกิดจากการปฏิเสธสมมติฐานศูนย์ที่เป็นจริงหรือในกรณีนี้ คือ การยอมรับว่าข้อสอบล่าเอียงทั้ง ๆ ที่ข้อสอบไม่ล่าเอียง

การศึกษาที่ 1 ใช้กลุ่มตัวอย่าง และกลุ่มเปรียบเทียบกลุ่มละ 100 200 500 และ 1000 คน เลือกข้อสอบ 1 ข้อ จากแบบสอบ ASVAB เลือกข้อที่มีค่าพารามิเตอร์ของข้อสอบคือ a b และ c ใกล้เคียงกับค่าเฉลี่ยของแบบสอบมากที่สุด คือ 1.32 0.03 และ 0.25 (ค่าเฉลี่ยคือ 1.22 0.09 และ 0.02) และเป็นข้อสอบที่ไม่ลำเอียงความสามารถ (θ) ของกลุ่มตัวอย่างทั้ง 2 กลุ่มถูกจำลองให้มีการกระจายเป็นปกติมีค่าความแปรปรวนเป็น 1.00 ความแตกต่างระหว่างค่าเฉลี่ยของการแจกแจงของ 2 กลุ่ม ถูกจำลองให้มีค่า 0.0 0.5 และ 1.0 แล้วคำนวณค่าความสามารถเฉลี่ยของแต่ละกลุ่มให้จุดกลางของค่าเฉลี่ยอันดับระดับความยากเฉลี่ยของแบบสอบ จำลองการตอบ 400 ครั้ง และวิเคราะห์ 400 ครั้ง ในแต่ละสถานการณ์

ผลการใช้ SIBTEST และ MH ตรวจสอบความลำเอียงของข้อสอบที่เลือกขึ้นมาดังกล่าว พบว่านอกจากกรณีกลุ่มตัวอย่างขนาด 100 คน และค่าความแตกต่างของความสามารถ (d_T) = 0.0 แล้ว อัตราการปฏิเสธสมมติฐานศูนย์ หรืออัตราความคลาดเคลื่อนประเภทที่ 1 ของวิธีทั้งสองแตกต่างกันอย่างไม่มีนัยสำคัญ ทั้งยังมีความแตกต่างจากค่าระดับนัยสำคัญที่กำหนด คือ 0.05 อย่างไม่มีนัยสำคัญเช่นกัน สำหรับความแตกต่างที่พบในจำนวนตัวอย่าง 100 คน และ $d_T = 0.0$ นั้น อัตราความคลาดเคลื่อนประเภทที่ 1 ของทั้งสองวิธีต่ำกว่า 0.05 อย่างไรก็ตาม SIBTEST มีอัตราความคลาดเคลื่อนประเภทที่ 1 สูงกว่า MH เล็กน้อย สำหรับการวิเคราะห์ความลำเอียงโดยใช้กลุ่มตัวอย่างไม่เกิน 1000 คน และมีแนวโน้มว่าอัตราความคลาดเคลื่อนจะสูงขึ้นเมื่อกลุ่มตัวอย่างมากขึ้นทั้ง 2 วิธี

การศึกษาที่ 2 ใช้แบบสอบเดิม แต่เลือกค่า a b และ c สำหรับข้อสอบที่นำมาวิเคราะห์ให้แตกต่างกันโดย

$$a = 0.4 \quad 1.0 \quad \text{และ} \quad 2.5$$

$$b = -1.5 \quad -0.5 \quad 0.0 \quad 0.5 \quad \text{และ} \quad 1.5$$

$$c = 0.20$$

จำลองการกระจายของความสามารถเป้าหมายให้เป็นปกติ d_T มีค่า 0.0 และ 1.0 นอกจากนั้นการจำลองทุกอย่างเป็นเช่นเดียวกับการศึกษาที่ 1 แต่ใช้จำนวนผู้สอบขนาดใหญ่ คือ 500 1,000 และ 3,000 คน และแต่ละสถานการณ์จำลองการตอบ 100 ครั้ง และตรวจสอบโดยใช้ SIBTEST และ MH พบว่าเมื่อค่า a สูงขึ้นและ b ค่า อัตราการปฏิเสธสมมติฐานจะสูงกว่าระดับที่กำหนดไว้คือ 0.05 ทั้ง 2 วิธี และ MH จะมีอัตราการปฏิเสธสูงกว่า SIBTEST และยิ่งจำนวนกลุ่มตัวอย่างใหญ่ขึ้น อัตราการปฏิเสธจะยิ่งสูงขึ้น ในขนาดกลุ่มตัวอย่าง 3,000 คน ที่ค่า a = 2.5 และ b = -1.5 อัตราการปฏิเสธของ MH สูงถึง 1.00 และ SIBTEST = 0.91 ในกรณีที่ $d_T = 1.0$ และ c = 0.20 ในภาพรวมแล้ว SIBTEST มีอัตรา

การปฏิเสธต่ำกว่า MH แต่เมื่อใช้ค่า $d_T = 0.0$ ไม่มีค่า c ผลปรากฏว่า อัตราการปฏิเสธของ ทั้ง 2 วิธีอยู่ในระดับใกล้เคียงกับ 0.05 แม้ว่า SIBTEST จะสูงกว่า MH แต่ไม่มีนัยสำคัญ

สรุปว่าเมื่อใช้จำนวนตัวอย่างน้อย ไม่เกิน 1,000 คน MH และ SIBTEST ให้ผลในด้านความคลาดเคลื่อนประเภทที่ 1 ใกล้เคียงกัน และอยู่ภายในระดับ 0.05 ในขณะที่พบว่า SIBTEST มีอัตราความคลาดเคลื่อนประเภทที่ 1 ต่ำกว่า MH ในกรณีที่จำนวน ตัวอย่างมีขนาดใหญ่ (มากกว่า 1,000 คน) และค่าเฉลี่ยของการแจกแจงของกลุ่มอ้างอิง และกลุ่มเปรียบเทียบมีความแตกต่างกัน รวมทั้งข้อสอบมีความยากค่า และค่าอำนาจจำแนกสูง นอกจากผลทางด้านอัตราการปฏิเสธแล้ว ยังพบว่าตัวประมาณความล่าเอียงคือ β และ Δ ที่มีค่า ต่ำกว่าเกณฑ์หรือมีค่าแสดงความล่าเอียงเพียงเล็กน้อยก็อาจมีนัยสำคัญได้ ($\Delta = 1.5$ และ 1.0 และ $\beta = 0.10$ หรือ 0.05)

การศึกษาเกี่ยวกับการตรวจค้นความล่าเอียงของข้อสอบในต่างประเทศ เท่าที่ผ่านมาจะเห็นว่า มีการใช้ทั้งข้อมูลจริงและข้อมูลจำลอง ข้อมูลจำลองที่นำมาใช้ในการศึกษา เป็นข้อมูลที่จำลองขึ้นตามทฤษฎีการตอบข้อสอบทั้งสิ้น ซึ่งเป็นจุดที่น่าวิจารณ์ว่าเป็นผลให้พบว่า การ ตรวจสอบข้อสอบล่าเอียงด้วย IRT มีความถูกต้องกว่าวิธีอื่น แม้จะยังมีความคลาดเคลื่อนในการ ตรวจสอบทั้งชนิด FN และ FP ก็ตาม

ข้อค้นพบที่สอดคล้องกับอีกประการ คือ ส่วนมากพบความสัมพันธ์กันใน ระดับสูงระหว่างวิธี IRT ชนิด 2 และ 3 พารามิเตอร์ และ MH แม้จะมีการพบข้อขัดแย้งกันอยู่ เช่นกันเมื่อใช้ข้อมูลจริงในด้านจำนวนข้อสอบล่าเอียงที่ตรวจพบ ซึ่ง Swaminathan และ Rogers (1990) พบว่า IRT ตรวจพบข้อสอบล่าเอียงได้มากกว่า MH ในขณะที่ Raju และคณะ (1993) พบว่า MH ตรวจพบได้มากกว่า IRT รวมทั้งการพบความไม่คงที่ของการตรวจพบข้อสอบล่าเอียง เมื่อใช้กลุ่มตัวอย่างต่างกันอันเนื่องมาจากความไวของค่าสถิติทดสอบต่อจำนวนกลุ่มตัวอย่าง และ ปัญหาการใช้เกณฑ์เพื่อใช้ตัดสินว่าข้อสอบล่าเอียง ซึ่งส่วนมากกำหนดเกณฑ์แตกต่างกันตามอำเภอใจ (arbitrary) (Hambleton 1986 Thissen Steinberg และ Wainer 1988 Linacre 1988 Perlman และคณะ 1988 Hambleton และ Rogers 1989 Baghi และ Ferrara 1990 Swaminathan และ Rogers 1990 และ Raju และคณะ 1993

จากข้อค้นพบดังกล่าว ทำให้ผู้วิจัยสนใจในการพัฒนาเกณฑ์ที่จะใช้ตัดสิน ข้อสอบล่าเอียงขึ้นจากข้อมูลเชิงประจักษ์ โดยเลือกวิธีการที่ผู้วิจัยเห็นว่ามีความคิด ในเรื่องของ ความล่าเอียงเหมาะสม เป็นแนวเดียวกันและมีผลการศึกษาที่ใกล้เคียงกัน 3 วิธี คือ วิธี IRT-2 พารามิเตอร์ วิธี MH และวิธี SIBTEST โดยนำดัชนีที่ได้จาก 3 วิธีมาใช้ในการพัฒนาเกณฑ์ 4 ตัว ได้แก่ SA UA α_{MH} และ β_{SIB} ซึ่งจากเอกสารรายงานที่เกี่ยวข้องดังกล่าว สรุปการใช้เกณฑ์ที่ผ่านมาของวิธีตรวจค้นความล่าเอียงดังตารางที่ 3

จากข้อค้นพบดังกล่าว ทำให้ผู้วิจัยสนใจในการพัฒนาเกณฑ์ที่จะใช้ตัดสินข้อสอบล่าเอียงขึ้นจากข้อมูลเชิงประจักษ์ โดยเลือกวิธีการที่ผู้วิจัยเห็นว่ามีความคิด ในเรื่องของความล่าเอียงเหมาะสม เป็นแนวเดียวกันและมีผลการศึกษาที่ใกล้เคียงกัน 3 วิธี คือ วิธี IRT-2 พารามิเตอร์ วิธี MH และวิธี SIBTEST โดยนำดัชนีที่ได้จาก 3 วิธีนี้มาใช้ในการพัฒนาเกณฑ์ 4 ตัว ได้แก่ SA UA α_{MH} และ β_{SIB}

ตารางที่ 3 สรุปเกณฑ์ที่พบว่ามี การนำมาใช้ตัดสินข้อสอบล่าเอียง

วิธี	เกณฑ์การตัดสิน
IRT	<ol style="list-style-type: none"> 1. การทดสอบนัยสำคัญของค่า χ^2 ของ Lord 2. $\chi^2 > 13.80$ UA (หรือ ϕ) $> .20$ UA $> .40$ 3. SA $> .468$ 4. การทดสอบค่า Z (Raju)
MH	<ol style="list-style-type: none"> 1. ทดสอบนัยสำคัญของ MH-χ^2 2. MH-$\chi^2 > 6.64$ 3. $\alpha_{MH} > 1.00$ 4. MH-$\chi^2 > 10.83$
SIB-TEST	ทดสอบนัยสำคัญของค่า Z