

วรรณคดีที่เกี่ยวข้อง

ในการนำเสนอวรรณคดีที่เกี่ยวข้องกับการวิจัยนี้ ผู้วิจัยได้นำเสนอโดยจำแนก
ออกเป็นตอน ๆ ดังนี้

1. ความหมายของการเทียบมาตรา
2. รูปแบบของการเทียบมาตรา
3. รูปแบบของการเก็บรวบรวมข้อมูล
4. การเทียบค่าความสามารถของผู้สอบ
5. ความคลาดเคลื่อนมาตรฐานของการเทียบมาตรา
6. ความเพียงพอของการเทียบมาตรา
7. งานวิจัยที่เกี่ยวข้องกับการเทียบมาตรา

ความหมายของการเทียบมาตรา

การเทียบมาตรา (Equating) เป็นการศึกษาเชิงประจักษ์ โดยใช้วิธีการทาง
การวัดผลที่มีผู้คิดค้นเพื่อทำให้คะแนนจากแบบสอบต่างฉบับเทียบกันได้ มีผู้ให้นิยามของการ
เทียบมาตราไว้หลายท่าน ดังนี้

กัลลิกเซน (Gulliksen, 1950) กล่าวถึงการเทียบมาตราว่า เป็นวิธีการทำ
คะแนนที่ได้จากแบบสอบสองฉบับที่วัดวิชาเดียวกัน ให้เป็นคะแนนสมมูล (Equivalent
Score) ที่เทียบกันได้โดยตรง โดยเสนอวิธีการให้ผู้สอบกลุ่มเดียวทำแบบสอบสองชุด และ
ใช้วิธีง่าย ๆ คือ แปลงคะแนนแต่ละชุดให้เป็นคะแนนมาตรฐาน แล้วนำคะแนนที่แปลงแล้ว

มาเทียบกันโดยตรง แต่ก่อนอื่นให้ตรวจสอบความเป็นคู่ขนานของแบบสลับโดยใช้สถิติของวิลส์ (Wilks) หรือพิจารณาจากค่าสหสัมพันธ์ระหว่างแบบสลับแต่ละชุดกับเกณฑ์ ถ้ามีค่าเท่ากับ 1 โดยประมาณ ($R_{\text{sc}} = R_{\text{sc}}$) ก็จะสามารถบอกความสามารถของความสัมพันธ์ในการเทียบคะแนนแบบสลับทั้งสองชุด

ฟลานานแกน (Flanagan, 1951:747-748) ให้ความหมายของการเทียบมาตรา ว่า เป็นวิธีการทำคะแนนจากแบบสลับต่างฉบับกัน ให้สามารถนำมาเปรียบเทียบกันได้ คำว่า "ความสามารถในการเปรียบเทียบกันได้" มีความหมายเฉพาะที่ว่า เมื่อกำหนดประชากรให้ ถ้าการแจกแจงของคะแนนจริงจากแบบสลับทั้งสองชุด ซึ่งสลับกับกลุ่มตัวอย่างที่เลือกมาขนาด ใหญ่ใด ๆ มีลักษณะเหมือนกันแล้ว คะแนนดิบจากแบบสลับทั้งสองชุดจึงจะสามารถเปรียบเทียบกันได้ หรือถ้าความเชื่อมั่น (Reliability) ของแบบสลับทั้งสองชุดเท่ากันในประชากรนั้นแล้ว ก็สามารถเปรียบเทียบการแจกแจงของค่าที่ได้เช่นกัน จากความหมายดังกล่าว เป็นนิยามเชิงทฤษฎี ในทางปฏิบัติได้นิยามไว้ว่า ในประชากรที่กำหนด ถ้าคะแนนจากแบบสลับสองชุดมีค่าเฉลี่ยเท่ากัน หรือเกือบเท่ากันในทุก ๆ กลุ่มตัวอย่างขนาดใหญ่ใด ๆ แล้ว การเปรียบเทียบจะทำได้ อย่างไรก็ตาม จะทำให้เกิดลักษณะตามที่นิยามเป็นเรื่องที่ยากมาก สิ่งที่สำคัญคือแบบสลับต้องเป็นคู่ขนานกัน ฟลานานแกนจึงแนะนำว่า ควรเริ่มต้นตั้งแต่การสร้างแบบสลับให้มีคุณสมบัติให้เป็นคู่ขนานกัน แล้วเลือกวิธีเทียบคะแนนที่เหมาะสม ซึ่งเสนอไว้ 4 วิธี คือ

1. ให้ค่าเฉลี่ย โดยคำนวณค่าเฉลี่ยของการแจกแจงคะแนนทั้งสองชุด ถ้าความแตกต่างของค่าเฉลี่ยอยู่ภายในขอบเขตของการแปรผันเชิงสุ่มแล้ว ถือว่าคะแนนทั้งสองชุดนั้นเปรียบเทียบกันได้ แต่ถ้าความแตกต่างมีนัยสำคัญ ให้ใช้วิธีบวกเข้าหรือลบออกเท่าจำนวนที่แตกต่างจากคะแนนชุดที่หนึ่ง เพื่อให้เกิดคะแนนสมมูลกับอีกชุดหนึ่ง
2. ใช้เทคนิคของสมการถดถอย โดยหาค่าประมาณที่ดีที่สุดของคะแนนจากแบบสลับชุดที่หนึ่งซึ่งรู้ค่าของอีกชุดหนึ่ง
3. ให้คะแนนมาตรฐาน ซึ่งเป็นวิธีปรับคะแนนอย่างคงที่ตลอดการแจกแจง
4. ใช้วิธีอิลลิปส์เซนไคล์ ซึ่งเป็นวิธีหาค่าคะแนนโดยการเปรียบเทียบตามสัดส่วนของการแจกแจงคะแนน

มาร์โค (Marco, 1981) ได้อธิบายความหมายของการเทียบมาตรฐานว่า เป็นกระบวนการแปลงคะแนน ที่ได้จากแบบสอบฟอร์มหนึ่ง ให้มีค่าตัวเลขสมมูล (Equivalent) กับตัวเลขที่ได้จากแบบสอบอีกฟอร์มหนึ่ง

แองกอฟ (Angoff, 1971) ให้ความหมายของการเทียบมาตรฐานว่า เป็นการแปลงระบบของหน่วยการวัดของแบบสอบฉบับหนึ่ง ไปสู่ระบบหน่วยการวัดของแบบสอบอีกฉบับหนึ่ง คะแนนที่ผ่านการแปลงแล้วจะให้ความหมายของการเทียบกันโดยตรง แองกอฟได้เสนอวิธีการในการเทียบมาตรฐานไว้ 2 รูปแบบ คือ การเทียบมาตรฐานแบบอควิวเอร์เรนซ์ และการเทียบมาตรฐานเชิงเส้นตรง

ลอร์ด (Lord, 1980) ให้ความหมายของการเทียบมาตรฐานว่าเป็นการแปลงคะแนนจากแบบสอบต่างฟอร์ม ให้มีความหมายให้สับเปลี่ยนกันได้ และเพื่อความเสมอภาคของผู้รับการสอบ

สงบ ลิกหณะ (2522) กล่าวว่า คะแนนจากแบบสอบสองฉบับ วัดสิ่งเดียวกันแต่ไม่จำเป็นต้องเป็นแบบสอบคู่ขนาน จะถือว่าเทียบกันได้ถ้าคะแนนจากแบบสอบทั้งสองมาจากคะแนนจริง (True Score) หรือความสามารถแท้ (True Ability) ที่เท่ากัน

ชูศักดิ์ ทิมภลิต (2527:2) ได้สรุปความหมายของการเทียบมาตรฐานว่า การเทียบมาตรฐานเป็นกิจกรรมที่เกี่ยวโยงกับกิจกรรมสองประการ คือ

1. กระบวนการที่ทำให้แบบสอบ 2 ฉบับใด ๆ มีความทัดเทียมกันหรือเท่ากันในเชิงโครงสร้าง
2. การใช้วิธีการทางสถิติ เพื่อปรับคะแนนที่ได้จากแบบสอบแต่ละฉบับให้อยู่ในมาตรฐานเดียวกัน และเทียบกันได้

จากความหมายข้างต้นพอสรุปได้ว่า การเทียบมาตรฐานเป็น วิธีการทางการวัดผล เพื่อปรับคะแนนที่ได้จากแบบสอบหนึ่งไปยังอีกแบบสอบหนึ่ง ให้สามารถเปรียบเทียบกันได้ โดยถือเอาความสามารถในการตอบข้อสอบของผู้สอบเป็นเกณฑ์

รูปแบบการเทียบมาตรฐาน

รูปแบบการเทียบมาตรฐานมีผู้เสนอไว้หลายรูปแบบดังนี้

1. การเทียบมาตรฐานรูปแบบอิกวิเปอร์เซนไทล์ (Equipercentile Equating) รูปแบบอิกวิเปอร์เซนไทล์เริ่มจากการแจกแจงของคะแนนจากแบบสอบ 2 ฉบับ มีลักษณะคล้ายกัน หรือแตกต่างกันบ้างก็เพียงเล็กน้อย การเทียบมาตรฐาน ณ ตำแหน่งเดียวกันของคะแนน 2 ชุดนั้น ผลของการเทียบมาตรฐานแสดงด้วยกราฟ โดยปกติการเทียบมาตรฐานรูปแบบนี้ให้ภาพการแปลงคะแนนที่สะท้อนถึงระดับความยากง่ายของแบบสอบ 2 ฉบับ ถ้าแบบสอบมีความยากใกล้เคียงกัน เส้นกราฟจะมีลักษณะใกล้เคียงเส้นตรง แต่ถ้ามีความยากแตกต่างกันเส้นกราฟจะเป็นเส้นโค้ง คะแนนสมมูลที่เกิดขึ้นจะถูกยึดหยุ่นตามคะแนนเดิม เพื่อให้คงรักษาให้เหมือนชุดก่อนตามต้องการ (Angoff, 1984)

บรูมและฮอลแลนด์ (Braun and Holland, 1982:19-22 อ้างถึงใน กาวีณี ศรีสุขวัฒนาพันธ์) ได้แสดงถึงความสัมพันธ์ของแบบสอบร่วมในการเทียบมาตรฐานรูปแบบอิกวิเปอร์เซนไทล์ ไว้ดังนี้ ให้กลุ่มตัวอย่าง 2 กลุ่ม คือ P และ Q ที่สุ่มมาต่างอิสระจากประชากรทั้งหมด (T) เขียนเป็นสมการได้ว่า

$$T = fP + (1-f)Q$$

โดย P และ Q ขึ้นกับน้ำหนักของ f และ $(1-f)$ ตามลำดับ การคำนวณ การแจกแจงของตัวอย่างรวม ทำได้โดยการคำนวณการแจกแจงตามเงื่อนไขของกลุ่ม ตัวอย่างกลุ่มที่ 1 (P) และกลุ่มที่ 2 (Q) ก่อน แล้วหาค่าเฉลี่ยที่ถ่วงน้ำหนักด้วย f และ $(1-f)$ ซึ่งได้จากสมการ คือ

$$f = N_1 / (N_1 + N_2)$$

เมื่อ N_1 และ N_2 คือ ขนาดของกลุ่มตัวอย่างที่ 1 และที่ 2 ตามลำดับ ค่าของ f คือค่าเฉลี่ยของกลุ่มตัวอย่างที่ 1 และกลุ่มตัวอย่างที่ 2

ในการเทียบมาตรฐานด้วยรูปแบบอิกวิเปออร์เซนไคล์ มีข้อตกลงเพื่อหาค่าประมาณ การแจกแจงของแบบสโบลบับที่ 1 ($F_T(X)$) และการแจกแจงของแบบสโบลบับที่ 2 ($G_T(X)$) ในประชากรทั้งหมด (T) ซึ่ง $F_T(X)$ และ $G_T(X)$ หาได้ในเทอมของปริมาณ จากข้อมูลที่รวบรวมได้ ฟังก์ชันของการเทียบคะแนนเป็นดังนี้

$$e_x(Y) = F_T^{-1}(G_T(Y))$$

วิธีการเทียบมาตรฐานแบบอิกวิเปออร์เซนไคล์นี้ มีข้อจำกัดที่ต้องคำนึงถึงอยู่หลายประการ คือ

- 1) วิธีนี้ใช้การเขียนเส้นกราฟเพื่อหาคะแนนสมมูลซึ่งจะต้องปรับหรือเกลาเส้นโค้งให้เรียบ ดังนั้นประสบการณ์และการตัดสินใจเกลาเส้นให้เรียบด้วยข้อมูลเหล่านั้น อาจนำไปสู่ความคลาดเคลื่อนอย่างมาก การเทียบมาตรฐานด้วยวิธีนี้จึงเป็นการเทียบไปยังจำนวนเต็มที่อยู่ใกล้ที่สุด (Angoff, 1971)
- 2) รูปแบบอิกวิเปออร์เซนไคล์ มีความไวต่อความแปรปรวนเชิงสุ่มมาก โดยเฉพาะกลุ่มตัวอย่างที่มีขนาดเล็ก (Angoff, 1971; Potthoff, 1982)
- 3) ในกรณีที่มีแบบสองสองชุดมีความแตกต่างกันมาก ผลของการเทียบมาตรฐานจะขาดความคงที่ (Potthoff, 1982)

4) การเทียบคะแนนสมมูลของแบบสอบสองชุด ทำได้เฉพาะในช่วงพิสัยของคะแนนที่มีความถี่ของคะแนนสิ่งเกิดเพียงพอ ส่วนช่วงที่มีความถี่ของคะแนนน้อยจะมีความคลาดเคลื่อนสูงมาก

2. การเทียบมาตรฐานรูปแบบเชิงเส้นตรง (Linear Equating)

วิธีนี้เป็นวิธีที่ง่ายในขั้นตอนของการแปลงคะแนน โดยมีข้อตกลงว่า แบบสอบสองชุดที่คู่ขนานกันจะมีการแจกแจงคะแนนดิบเป็นอย่างไรก็เหมือนกัน ยกเว้นสำหรับความแตกต่างในค่าเฉลี่ยและความแปรปรวน นั่นคือ นอกจากความแตกต่างในสองโวมเมนต์แรกแล้วโวมเมนต์มาตรฐานของการแจกแจงคะแนนดิบของแบบสอบสองชุดของกลุ่มผู้สอบที่กำหนดให้ จะต้องเหมือนกันหมด (Angoff, 1971)

จากจุดล้นของการเทียบมาตรฐานโคมพิลล์ควิเปอร์เซ็นต์คือ การประมาณค่า $F(x)$ และ $G(y)$ และนั่นโวมเมนต์ที่ง่ายและลาติชการประมาณค่าพารามิเตอร์ที่เป็นไปได้ในแง่การปฏิบัติ คือ โวมเมนต์เชิงเส้นตรง (Linear Model)

$$e_x(y) = ay + b \quad (1)$$

จากโวมเมนต์ของฟังก์ชันการเทียบเชิงเส้นนี้ มีพารามิเตอร์ต้องประมาณค่า 2 ตัว คือ ค่าความชัน (a) และ ค่าคงที่ (b)

ในกรณีของสมการที่ (1) ซึ่งตรงกับฟังก์ชันการเทียบโคมพิลล์ควิเปอร์เซ็นต์ ถ้า

$$F^{-1}(G(y)) = ay + b \quad (2)$$

หรือเช่นเดียวกันถ้า

$$G(y) = F(ay + b) \quad (3)$$

สมการที่ (3) ตรงกับข้อตกลงที่ให้ F และ G อยู่ในรูปของ

$$F(X) = H((x - \mu_x) / \sigma_x) \quad (4)$$

และ

$$G(Y) = H((Y - \mu_y) / \sigma_y) \quad (5)$$

สำหรับฟังก์ชันการกระจายบางอย่าง H ในสมการ (4) และ (5) $\mu_x, \mu_y, \sigma_x, \sigma_y$ คือค่าเฉลี่ยและความแปรปรวนของแบบสอบฟอร์ม X และฟอร์ม Y ถ้าประยุกต์สมการ (4) และ (5) กับสมการ (2) เราจะได้รูปฟอร์มของ a และ b ในสมการ (2) คือ

$$a = \sigma_x / \sigma_y \quad (6)$$

$$b = \mu_x - (\sigma_x / \sigma_y) \mu_y \quad (7)$$

นำไปใช้กับสมการ (1) เราจะได้ฟังก์ชันการเทียบเชิงเส้นตรง คือ

$$e(y) = \mu_x + (\sigma_x / \sigma_y) (y - \mu_y)$$

3. การเทียบมาตรฐานแบบอิงทฤษฎีตอบสนองข้อสอบ

การเทียบมาตรฐานแบบอิงทฤษฎีตอบสนองข้อสอบ ลอร์ด (Lord, 1980) ได้กล่าวว่า แบบสอบสองฉบับใด ๆ ที่เทียบคะแนนกันต้องเป็นแบบสอบที่มีมิติเดียวกัน และเทียบคะแนนจากความสามรถ (e) ที่เท่ากัน ซึ่งจะมีข้อกำหนดที่สำคัญ 3 ประการ ดังนี้

1) ความเสมอภาค (Equity) คือ ถ้าพิจารณาที่ระดับความสามารถ (๑) ใด ๆ การแจกแจงความถี่อย่างมีเงื่อนไขของคะแนนแปลง $x(y)$ หรือ y^* (คะแนนจากแบบสอบ X ที่แปลงมาอยู่สเกลเดียวกับแบบสอบ Y) ที่ ๑ ที่กำหนด ให้อัตลักษณ์การเทียบคะแนน ต้องเหมือนกันกับการแจกแจงความถี่อย่างมีเงื่อนไขของคะแนนจากแบบสอบที่ต้องการเทียบ (x)

2) ความไม่ผันแปรตามกลุ่ม (Invariance Across Groups) คือ คะแนนแปลง $x(y)$ จะเหมือนกันโดยไม่ขึ้นกับตัวแปรอื่น ๆ ของประชากรที่นำมาสร้างสมการเทียบมาตรฐาน

3) ความสมมาตร (Symmetry) คือ คะแนนเทียบมาตรฐานจะต้องเหมือนกันไม่ว่าการเทียบนั้นจะเทียบจากแบบสอบชุดที่ 1 ไปชุดที่ 2 หรือจากแบบสอบชุดที่ 2 ไปชุดที่ 1

การเทียบมาตรฐานโดยวิธีที่ผู้ตอบสนองข้อสอบ มีการเทียบมาตรฐานโดยใช้คะแนนจริงและการเทียบมาตรฐานโดยใช้คะแนนสังเกต ทั้งสองวิธีนี้มีทั้งข้อได้เปรียบเสียเปรียบ กล่าวคือ การเทียบด้วยคะแนนจริงไม่สามารถอธิบายคะแนนที่อยู่ต่ำกว่าระดับการเดาได้ โดยจะให้ความหมายของคะแนนสมมูลเฉพาะคะแนนที่อยู่เหนือค่าเฉลี่ยของการเดา ทั้งถึงแม้จะเป็นการเทียบโดยวิธีใช้คะแนนจริง แต่ยังเป็นค่าที่ประมาณได้จากสูตรในการคำนวณ ดังนั้นจึงยังคงมีความคลาดเคลื่อนอยู่ ส่วนการเทียบโดยวิธีใช้คะแนนสังเกตนั้นเป็นการเทียบคะแนนโดยประมาณ ซึ่งอธิบายคะแนนสมมูลจาก X และ Y ได้ครอบคลุมพิสัยของคะแนนที่สังเกตได้ ลอร์ดได้กล่าวว่า ทั้งสองวิธีนี้มีความสอดคล้องกันมาก แต่การสรุปผลในการอ้างอิงต้องทำอย่างรอบคอบ (Lord, 1980) อย่างไรก็ตามการเทียบมาตรฐานโดยใช้คะแนนสังเกตนี้ ก็นำไปใช้ในงานวิจัยน้อยมาก อาจเป็นเพราะว่า วิธีนี้มีความซับซ้อน และการลงทุนแพงกว่าวิธีการเทียบมาตรฐานโดยใช้คะแนนจริง (Lord and Wingersky, 1984)

ลอร์ด (Lord, 1980) ได้แบ่งรูปแบบการเทียบมาตรฐานโดยอิงทฤษฎีตอบสนองข้อสอบออกเป็นหลายวิธีดังนี้

1. วิธีการเทียบมาตรฐานโดยใช้คะแนนจริง (True-Score Equating)

พิจารณาจากความสัมพันธ์เชิงคณิตศาสตร์ ระหว่างความสามารถและจำนวนข้อที่ทำได้ ซึ่งจะได้คะแนนจริงจากแบบสอบ 2 ฉบับ ดังนี้

$$\xi = \xi(\theta) = \sum_{i=1}^m P_i(\theta_x)$$

$$\eta = \eta(\theta) = \sum_{j=1}^n P_j(\theta_y)$$

- เมื่อ
- ξ แทน คะแนนจริงจากแบบสอบฉบับ X
 - η แทน คะแนนจริงจากแบบสอบฉบับ Y
 - m แทน จำนวนข้อของแบบสอบฉบับ X
 - n แทน จำนวนข้อของแบบสอบฉบับ Y
 - θ_x แทน ค่าความสามารถจากแบบสอบฉบับ X
 - θ_y แทน ค่าความสามารถจากแบบสอบฉบับ Y

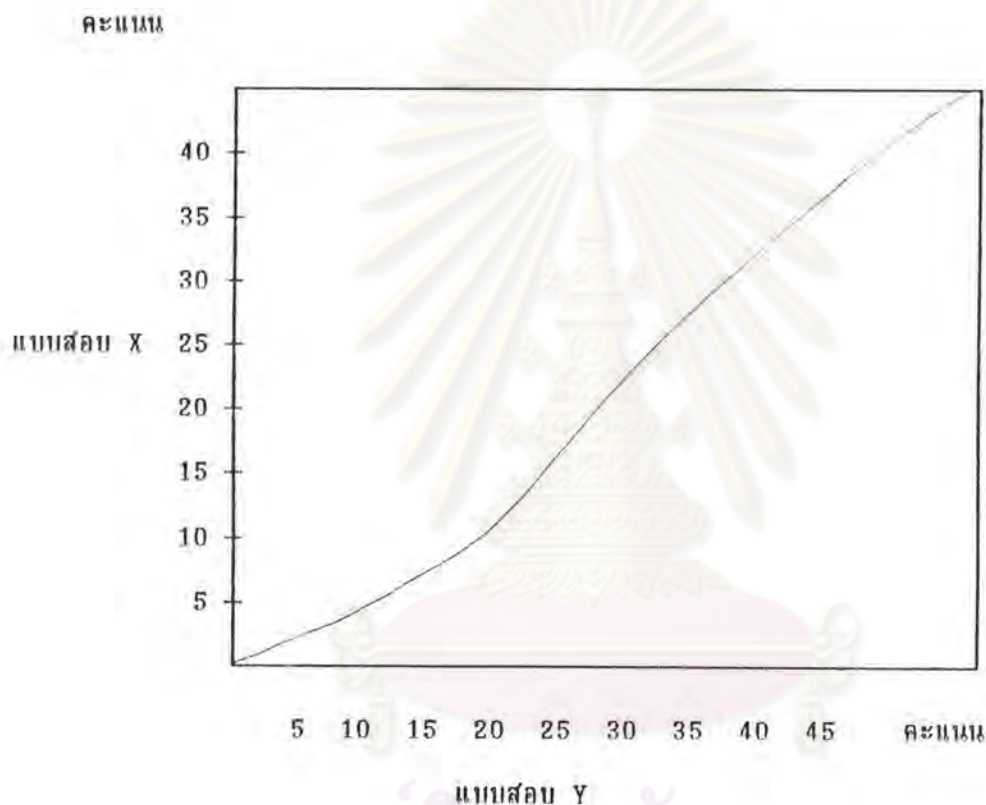
ในทางปฏิบัติ เราจะหาค่า $P_i(\theta)$ และ $P_j(\theta)$ ได้จากการประมาณค่าพารามิเตอร์ของข้อกระทงโดยเลือกที่จะใช้ชนิดหนึ่งพารามิเตอร์, สองพารามิเตอร์ หรือ สามพารามิเตอร์ จากนั้น นำคะแนนจริงของผลการสอบจากแบบสอบฉบับ X และแบบสอบ Y มาหาความสัมพันธ์กัน โดยหาค่าความสามารถ (θ) ที่ระดับเดียวกัน

2. วิธีการเทียบมาตรฐานโดยใช้คะแนนจริงด้วยแบบสอบร่วม (True-score

Equating with an Anchor Test) การให้แบบสอบร่วมทำได้ 2 กรณี คือ ใช้แบบสอบร่วมรวมเข้าเป็นชุดเดียวกับแบบสอบที่ต้องการเทียบ เรียกว่า แบบสอบร่วมภายใน (Internal Anchor Test) ส่วนในกรณีที่จัดแยกเป็นชุดแบบสอบต่างหาก จากแบบสอบที่ต้องการเทียบ เรียกว่า แบบสอบร่วมภายนอก (External Anchor Test) จุดมุ่งหมาย

ของการใช้แบบสอบร่วม คือ เพื่อใช้แบบสอบร่วมเป็นตัวเชื่อมข้อมูลเข้าด้วยกัน ทำให้ค่าพารามิเตอร์ที่วิเคราะห์ออกมาอยู่บนมาตราเดียวกัน ซึ่งถ้าขาดแบบสอบร่วมแล้ว จะไม่สามารถเทียบคะแนนกันได้นอกเสียจากว่า กลุ่มตัวอย่างที่สอบแบบสอบ X และ Y จะมีการกระจายความสามารถเหมือนกัน

ในทางปฏิบัติจะหาค่าความสามารถ (ξ) ของกลุ่มตัวอย่างที่ทำแบบสอบร่วม ไปเทียบมาตรากับแบบสอบฉบับ X และแบบสอบฉบับ Y ดังภาพที่ 1



ภาพที่ 1 การเทียบระดับคะแนนโดยใช้คะแนนจริงจากแบบสอบร่วม

3. วิธีการเทียบมาตราโดยใช้คะแนนดิบด้วยแบบสอบร่วม (Raw-Score Equating with an Anchor) ปัญหาของการเทียบมาตราโดยใช้คะแนนจริง คือ ไม่สามารถทราบคะแนนจริงของแต่ละคนได้ นอกจากใช้วิธีประมาณจากผลการสอบโดยสมการ

$$\xi = \sum_{i=1}^n P_i(\xi)$$

เมื่อ $P_i(\theta)$ เป็นความน่าจะเป็นที่ทำข้อสอบข้อที่ i ถูกที่ระดับความสามารถ θ และ n เป็นจำนวนข้อสอบทั้งหมด ซึ่งค่าที่ได้เป็นค่าประมาณเท่านั้น ยังไม่มีคุณสมบัติเป็นคะแนนจริง ดังนั้นการเทียบมาตรฐานโดยใช้คะแนนจริงจึงคล้ายกับการเทียบมาตรฐานโดยใช้คะแนนดิบชนิดใหม่ นั่นเอง สำหรับการเทียบมาตรฐานโดยใช้คะแนนดิบในข้อนี้อาศัยข้อมูลจากการตอบแบบสอบถาม (Anchor Test)

การดำเนินการเทียบมาตรฐานวิธีนี้ จะเริ่มด้วยการประมาณลักษณะการแจกแจงของความสามารถของกลุ่มผู้สอบรวม $r(\theta)$ ซึ่งหมายถึง ผู้สอบทั้งหมดที่ทำแบบสอบถามการแจกแจง (θ) ในกลุ่มเป็นการประมาณการกระจายของ $r(\theta)$

การประมาณการแจกแจงของคะแนนดิบ X สำหรับกลุ่ม $\phi_x(x)$ ได้จากสมการ

$$\hat{\phi}_x(x) = \frac{1}{N} \sum \phi_x(x/\hat{\theta}_n)$$

เมื่อ $a = 1, 2, 3, \dots, N$ คือ ผู้สอบแต่ละระดับ และ $\phi_x(x/\hat{\theta}_n)$ เป็นการแจกแจงของคะแนนดิบ x สำหรับผู้ที่มีความสามารถ θ

ถ้า $r(\theta)$ มีลักษณะต่อเนื่อง (Continuous) สามารถประมาณด้วยสมการ

$$\hat{\phi}_x(x) = \int_{-\alpha}^{\alpha} \phi_x(x/\theta) r(\theta) d\theta$$

เมื่อ $\phi_x(x/\hat{\theta}_n)$ ได้มาจากการประมาณค่าพารามิเตอร์โดยวิธีทฤษฎีคลอสนองข้อสอบ (Item Parameter) ของข้อสอบในแบบสอบฉบับ X และสมการต่าง ๆ สามารถนำมาประยุกต์ใช้กับแบบสอบฉบับ Y ได้เช่นกัน

เนื่องจากแบบสอบฉบับ X และฉบับ Y เป็นอิสระจากกัน เมื่อความสามารถ (θ) คงที่ การเทียบมาตรฐานต้องคำนวณการกระจายร่วม (Joint Distribution) จากสมการ

$$\hat{\phi}(x, y) = \frac{1}{N} \sum \phi_x(x/\theta_n) \phi_y(y/\theta_n)$$

$$\hat{\phi}(x, y) = \int_{-\alpha}^{\alpha} \phi_x(x/\theta) \phi_y(y/\theta) Y d\theta \quad (8)$$

จากสมการข้างต้นจะเห็นบทบาทของแบบสอร่วมได้ชัดเจน ทำให้สามารถประมาณการแจกแจงร่วม (Joint Distribution) ของ x และ y ถึงแม้จะไม่มีผู้ใดทำแบบสอทั้งสองฉบับ

สมการ (8) เป็นการแจกแจงร่วมกันของตัวแปรสามตัว คือ θ , x และ y เนื่องจาก θ เป็นตัวระบุคะแนนจริง ξ และ η การแจกแจงนี้ถือว่าเป็นการแจกแจงร่วมกันของตัวแปร 4 ตัว คือ ξ , η , x และ y สมการนี้ไม่มีทางเลือกประมาณได้ นอกจากระบุความสัมพันธ์เชิงลิวลิปเปอร์เซนส์ไคส์ ระหว่าง x และ y จากสองสมการสุดท้าย

อย่างไรก็ตาม รูปแบบการเทียบมาตรฐานด้วยทฤษฎีตอบสนองข้อสอบยังแบ่งออกเป็น 3 แบบ ซึ่งแตกต่างกันตามจำนวนพารามิเตอร์ที่ใช้ในแต่ละรูปแบบ ดังนี้ (Hambleton and Swanminathan, 1985)

1. แบบที่ใช้พารามิเตอร์ 3 ตัว (Three-Parameters Logistic Model)
ฟังก์ชันที่ให้กับเส้นลักษณะของข้อสอบ เป็นดังนี้

$$P_i(\theta) = c_i + (1-c_i) \frac{e^{Da_i(\theta-b_i)}}{1 + e^{Da_i(\theta-b_i)}} ; i = 1, 2, 3, \dots, n$$

เมื่อ $P_i(\theta)$ คือ ความน่าจะเป็นในการตอบข้อสอบข้อที่ i ถูกเมื่อผู้สอบมีความสามารถ θ

a_i คือ ค่าอำนาจจำแนกของข้อสอบข้อที่ i

b_i คือ ค่าความยากของข้อสอบข้อที่ i

c_i คือ ค่าการเดาของข้อสอบข้อที่ i

D คือ scaling factor มีค่า 1.702

e คือ ค่าคงที่มีค่าประมาณ 2.7182818...

2. แบบที่ใช้พารามิเตอร์ 2 ตัว (Two-Parameters Logistic Model)

ฟังก์ชันที่ใช้เป็นดังนี้

$$P_i(\theta) = \frac{e^{Da_i(\theta - b_i)}}{1 + e^{Da_i(\theta - b_i)}} \quad ; i=1, 2, 3, \dots, n$$



3. แบบที่ใช้พารามิเตอร์ 1 ตัว (One-Parameter Logistic Model)

ซึ่งราสช์ (Rasch) ได้พัฒนาขึ้นมาใช้ โดยรูปแบบนี้เป็นกรณีพิเศษของแบบ 2 พารามิเตอร์ ที่เบิร์นบอม (Birnbaum) ได้พัฒนาขึ้น ข้อตกลงเบื้องต้นที่เพิ่มขึ้น คือ ค่าอำนาจจำแนกของข้อสอบ (a) เป็นค่าคงที่สำหรับข้อสอบแต่ละข้อ และค่าความยากของข้อสอบแปรเปลี่ยนไป ฟังก์ชันที่ใช้ เป็นดังนี้

$$P_i(\theta) = \frac{e^{D\bar{a}(\theta - b_i)}}{1 + e^{D\bar{a}(\theta - b_i)}} \quad ; i=1, 2, 3, \dots, n$$

รูปแบบในการเก็บรวบรวมข้อมูล

ในการเทียบคะแนนแบบสอบต่างชุดซึ่งสอบกับผู้สอบต่างกลุ่มนั้น การประมาณค่าพารามิเตอร์แยกกันระหว่างแบบสอบแต่ละชุดจะไม่สามารถอยู่ในสเกลเดียวกัน ดังนั้นจึงจำเป็นต้องออกแบบในการเก็บรวบรวมข้อมูลภายใต้รูปแบบที่สามารถจะจัดกระทำทางสถิติ เพื่อหาค่าคงที่สำหรับการแปลงค่าพารามิเตอร์ที่ประมาณค่าได้จากแบบสอบต่างชุดให้อยู่ในสเกลร่วมกัน และสามารถเทียบหาคะแนนสมมูลระหว่างแบบสอบได้ สำหรับรูปแบบในการเก็บรวบรวมข้อมูลเพื่อเทียบคะแนนนั้นพอแบ่งได้เป็น 3 รูปแบบดังนี้ (Lord, 1982 : Hambleton, 1985)

1. รูปแบบที่ใช้กลุ่มผู้สอบกลุ่มเดียวหรือใช้ผู้สอบร่วม โดยกลุ่มเดียวกันทำแบบสอบทั้งสองชุด การประมาณค่าความยาก จึงไม่เกี่ยวข้องกับความสามารถ

ของกลุ่มผู้สอบแต่ลำดับการสอบชุดใดก่อนหลังอาจทำให้เกิดการเวียนรู้ การฝึกฝน และความเหนื่อยล้า ซึ่งจะมีผลกระทบต่อคะแนนผลการสอบได้ ดังนั้นเพื่อจัดความล่าเอียงเหล่านี้ จึงคิดเปลี่ยนรูปแบบนี้ โดยสุ่มผู้สอบเป็นสองกลุ่ม แต่ละกลุ่มทำแบบสอบทั้งสองชุดในลักษณะของการจัดลำดับก่อนหลังต่างกัน โดยวิธีสุ่ม

2. รูปแบบที่ใช้กลุ่มที่เทียบเท่ากัน รูปแบบนี้กลุ่มผู้สอบที่เทียบเท่ากับสองกลุ่ม ซึ่งอาจได้มาโดยการสุ่ม หรือเป็นกลุ่มที่ได้จากการจับคู่ (Matching) ทางด้านความรู้และความสามารถ ทั้งนี้เพื่อให้คะแนนผลการสอบจากแบบสอบแต่ละชุดที่นำมาเทียบคะแนนนั้น ไม่มีผลของความแตกต่างในความสามารถระหว่างกลุ่มผู้สอบ อย่างไรก็ตาม กลุ่มผู้สอบทั้งสองอาจมีความแตกต่างในการแจกแจงความสามารถ ถึงแม้จะมีเพียงเล็กน้อยก็อาจทำให้เกิดความล่าเอียงในการเทียบคะแนนได้

3. รูปแบบที่ใช้แบบสอบร่วม รูปแบบนี้กลุ่มผู้สอบสองกลุ่มที่แตกต่างกัน (อาจเป็นกลุ่มที่ได้มาจากการสุ่มหรือไม่ได้จากการสุ่ม) แต่ละกลุ่มทำแบบสอบเพียง 1 ชุด ซึ่งแต่ละชุดจะประกอบด้วยข้อสอบร่วมจำนวนหนึ่ง โดยอาจสอบรวมหรือสอบแยกต่างหากจากแบบสอบทั้งสองชุดก็ได้ แบบสอบร่วมจะต้องคู่ขนาน หรือเป็นตัวแทนแบบสอบทั้งสองชุดนั้นให้มากที่สุด แบบสอบร่วมจะกำหนดสเกลร่วมของคะแนนจากแบบสอบแต่ละชุดที่สอบโดยผู้สอบทั้งสองกลุ่ม เพื่อนำไปปรับค่าความยากของข้อสอบ หรือค่าความสามารถของผู้สอบจากแบบสอบต่างชุดให้อยู่ในสเกลร่วมกัน

การเทียบค่าประมาณความสามารถของผู้สอบ

การเทียบค่าประมาณความสามารถ จะแตกต่างกันไปตามรูปแบบในการเก็บรวบรวมข้อมูลซึ่งมีรายละเอียดในการเทียบ ดังนี้

1. รูปแบบที่ผู้สอบกลุ่มเดียว เมื่อแบบสอบชุด X และชุด Y สอบกับผู้สอบกลุ่มเดียวกัน มีวิธีการที่เหมาะสม 2 วิธี คือ

วิธีที่ 1 รวมข้อมูลทั้งสองชุดและประมาณค่าพารามิเตอร์ของข้อสอบและค่าความสามารถ ซึ่งจะได้ค่าพารามิเตอร์อยู่ในสเกลของแบบสอบชุด X และแบบสอบชุด Y

วิธีที่ 2 แยกข้อมูลแต่ละชุด (ถ้าจำเป็นต้องประมาณค่าพารามิเตอร์แยกกัน สำหรับแบบสลับแต่ละชุด) แล้วกำหนดค่า θ ให้เหมือนกันในการประมาณค่าพารามิเตอร์แต่ละชุด ด้วยวิธีนี้ค่าประมาณความสามารถจะอยู่บนสเกลร่วมกัน

2. รูปแบบที่ใช้กลุ่มเทียบเท่ากัน กรณีนี้ไม่มีข้อสอบร่วมหรือผู้สอบร่วม จึงจำเป็นต้องประมาณค่าพารามิเตอร์แยกกันสำหรับแบบสลับแต่ละชุด โดยมีลำดับขั้นดังนี้

2.1 กลุ่มผู้สอบที่ได้จากการสุ่ม (จำนวนผู้สอบเท่ากัน) โดยแต่ละกลุ่มทำแบบสลับเพียง 1 ชุด คือ ชุด X หรือชุด Y

2.2 ประมาณค่าพารามิเตอร์แยกกันสำหรับกลุ่มผู้สอบแต่ละกลุ่ม โดยกำหนดค่า θ ให้เหมือนกันในการประมาณค่าทั้งสองกลุ่ม

2.3 เรียงลำดับค่าประมาณ θ และจับคู่ระหว่างค่าประมาณที่ต่ำสุดของ θ_x กับ θ_y ที่ต่ำสุดที่ละคู่เรื่อยไป

2.4 พล็อตกราฟระหว่างค่า θ_x กับค่า θ_y ซึ่งถ้าเป็นไปตามข้อตกลงของทฤษฎีตอบสนองข้อสอบแล้ว ผลที่ได้จะเป็นเส้นตรง อย่างไรก็ตามในทางปฏิบัติให้ค่าประมาณพารามิเตอร์ ดังนั้นผลที่ได้อาจไม่เป็นเส้นตรง โดยเฉพาะส่วนเวลาของค่า θ_x และ θ_y จะมีความคลาดเคลื่อนในการประมาณค่ามากกว่าส่วนอื่น

3. รูปแบบที่ใช้แบบสอบร่วม ในรูปแบบนี้ผู้สอบ N_x ทำแบบสลับชุด X ซึ่งมี n_x ข้อ และแบบสอบร่วม n_y ข้อ เช่นเดียวกับผู้สอบ N_y ทำแบบสลับชุด Y ซึ่งมี n_y ข้อ และข้อสอบร่วมอีก n_x ข้อ มีวิธีการที่เหมาะสม 2 วิธี สำหรับการเทียบมาตรฐาน คือ

วิธีที่ 1 ประมาณค่าแยกกันสำหรับแบบสลับแต่ละชุด คือ ชุด X มีผู้สอบ N_x คน และมีข้อสอบ $(n_x + n_y)$ ข้อ ส่วนชุด Y มีผู้สอบ N_y คน และมีข้อสอบ $(n_y + n_x)$ ข้อ กำหนดค่า θ ให้เหมือนกัน ในการประมาณค่าทั้งสองชุด แล้วหาเส้นความสัมพันธ์ระหว่างแบบสลับโดยผ่านข้อสอบร่วม n_x ข้อ ที่สอบกับกลุ่มผู้สอบทั้งสองกลุ่ม $(N_x + N_y)$ คน ตามวิธีที่อธิบายในข้อ 1 ดังนั้น ความสามารถของผู้สอบจะเทียบกันได้โดยให้เส้นความสัมพันธ์

วิธีที่ 2 ใช้โปรแกรม LOGIST ประมาณค่าพารามิเตอร์ของข้อสอบและค่าความสามารถของผู้สอบพร้อมกันทั้งสองชุด ดังนี้

1) จัดข้อมูลรวบรวบว่าผู้สอบ $N_x + N_y$ คน ทำแบบสอบ $(n_x + n_y + n_z)$ ข้อ

2) จัดข้อสอบ n_z ข้อ ซึ่งผู้สอบ N_x คน ไม่ได้ทำแบบสอบให้เป็นข้อสอบที่ผู้สอบทำไม่ถึง (not reached) ส่วนข้อสอบ n_x ข้อ ซึ่งผู้สอบ N_y คน ไม่ได้ทำข้อสอบก็จะจัด ให้เป็นข้อสอบที่ผู้สอบทำไม่ถึงเช่นกัน

3) ประมาณค่าพารามิเตอร์ความสามารถ สำหรับผู้สอบ $(N_x + N_y)$ คน และพารามิเตอร์ของข้อสอบ สำหรับข้อสอบ $(n_x + n_y + n_z)$ ข้อ พร้อมกันครั้งเดียวด้วยวิธีนี้ ค่าประมาณพารามิเตอร์ที่ได้จากแบบสอบต่างชุดจะอยู่ในสเกลเดียวกันทั้งหมด และค่าประมาณความสามารถระหว่างผู้สอบแบบสอบต่างชุดสามารถนำมาเปรียบเทียบได้

ความคลาดเคลื่อนมาตรฐานของการเทียบมาตรฐาน

การประเมินความคลาดเคลื่อนของการเทียบมาตรฐานที่สมบูรณ์ คือ การวิเคราะห์แหล่งความลำเอียง (bias) และความแปรผัน (variability) แล้วนำเสนอให้เห็นเป็นภาพเดียวกัน เพื่อช่วยในการตัดสินใจว่า ในการนำการเทียบมาตรฐานด้วยรูปแบบนั้น ๆ ไปใช้ มีความเหมาะสมหรือไม่อย่างไร (Braun and Hooland, 1982)

1. ความลำเอียงของการเทียบมาตรฐาน (Bias error)

ความลำเอียงเป็นองค์ประกอบหนึ่งของความคลาดเคลื่อน ซึ่งจะตั้งพิจารณากันอย่างรอบคอบ เมื่อต้องการนำผลของการเทียบมาตรฐานไปใช้ บรุนและฮอลแลนด์ (1982) ได้เสนอทวิเคราะห์ไว้ว่า ความลำเอียงทางสถิติ มีความหมาย 2 ประการ คือ

1) ความแตกต่างของค่าเฉลี่ยของค่าตัวประมาณตลอดจนการทำารสุ่มซ้ำ ๆ จากประชากรเดียวกัน กับค่าของประชากรที่ถูกประมาณ

2) ความแตกต่างของฟังก์ชันของการเทียบมาตรฐานโดยประมาณ (Estimated equating function) และค่าที่แท้จริงของฟังก์ชันการเทียบมาตรฐาน (equating function)

แหล่งของความลำเอียงของการเทียบมาตรฐานที่สำคัญมี 2 แหล่ง คือ

2.1 ความผันแปรของประชากร (Population variability) เนื่องจากการเทียบมาตรฐานไม่ได้จัดกระทำกับประชากรกลุ่มเดียวตลอด ฟังก์ชันของการเทียบมาตรฐานที่สร้างขึ้นกับเพื่อใช้กับประชากรหนึ่ง อาจนำไปใช้กับประชากรหนึ่งไม่ได้ เช่น ถ้า $X_p^*(Y)$ เทียบมาตรฐานจาก Y ไปยังมาตรฐาน X ในประชากร P และ $Y_q^*(Z)$ เทียบมาตรฐานจาก Z ไปยังมาตรฐาน Y ในประชากร Q แล้ว จะสรุปว่า ฟังก์ชันเชิงซ้อน $X_p^*(Y_q^*(Z))$ เทียบมาตรฐานจาก Z ไปยัง X ในประชากร P หรือ Q ย่อมไม่ได้ ในแบบแผนการเทียบมาตรฐานที่ใช้แบบสอบรวมได้แก้ไขปัญหานี้ โดยการสร้างประชากรสังเคราะห์ (Synthesis population) ซึ่งเป็นผลรวมของประชากร P และ Q ตามสัดส่วน แต่ยังไม่รับประกันผลว่า จะไม่มีความลำเอียงในทุก ๆ ประชากรที่แปรเปลี่ยนไป

2.2 ความคลาดเคลื่อนของรูปแบบ (Model error) ความคลาดเคลื่อนนี้ เกิดขึ้นจากการจัดการกระทำข้อสมมุติฐานที่ว่าตัวรูปการแจกแจงให้ต่างกันอย่างผิดพลาด คั้งแบ่งเป็น 2 พวก คือ

- 1) ความคลาดเคลื่อนของรูปแบบที่ทดสอบได้ (Testable model error) เช่น ความเป็นเส้นตรงของฟังก์ชันการถดถอย สามารถทดสอบได้จากข้อมูลจำนวนมาก ซึ่งจะเห็นผลได้ว่า มีความเพี้ยนพลหรือไม่
- 2) ความคลาดเคลื่อนรูปแบบที่ทดสอบไม่ได้ (Nontestable model errors) เป็นกรณีที่ไม่สามารถหาข้อมูลมาทดสอบเพื่อพิสูจน์ว่า สมมุติฐานที่ว่าไว้เพี้ยนพลหรือไม่

2. ความแปรผันเชิงสุ่ม และความคลาดเคลื่อนมาตรฐานของการเทียบมาตรฐาน

วิธีการเทียบมาตรฐานทุกวิธี ไม่ว่าจะอยู่ในรูปแบบการเทียบมาตรฐานใด เมื่อกลุ่มตัวอย่างผู้สอบเป็นกลุ่มที่สุ่มมาจากประชากรเดี่ยว หรือหลายประชากร ย่อมมีความผันแปรเชิงสุ่มเกิดขึ้น จึงนิยามให้เทคนิคการประมาณค่าความคลาดเคลื่อนเชิงสุ่ม (Estimate sampling error) กับวิธีการเทียบมาตรฐานต่าง ๆ ความคลาดเคลื่อนเชิงสุ่มมีสมมุติฐานเบื้องต้นว่า กลุ่มตัวอย่างมาจากการสุ่ม และใช้ความคลาดเคลื่อนมาตรฐานของการเทียบมาตรฐาน (Standard error of equating : SEE) เป็นการวัดความแปรผันประเภทนี้

ความคลาดเคลื่อนมาตรฐานโดยอิงรูปแบบทฤษฎีตอบสนองข้อสอบ ค่าประมาณคะแนนจำนวนข้อที่ตอบถูก (Number-right score) ξ ของผู้สอบที่ทำแบบสอบ X และ η ของผู้ที่ทำแบบสอบ Y มีค่าเท่ากับฟังก์ชันคุณลักษณะที่ประเมินที่ระดับความสามารถ (θ) เดียวกัน ให้ความหมายของการเท่ากันในเชิงคะแนนสมมูล ในทางปฏิบัติจึงนำความสัมพันธ์เชิงฟังก์ชัน ξ และ η มาใช้เทียบคะแนน X และ Y จากคะแนนที่สอบได้จริงของคนทั้งสองกลุ่ม การประมาณเพื่อให้เทียบ η ไปกับ ξ ตามฟังก์ชัน จะต้องใช้ค่าประมาณของประชากรข้อสอบ ค่าประมาณเหล่านี้ คือ ที่มาของความคลาดเคลื่อนเชิงสุ่ม (Sampling error) ในการเทียบมาตรา (Lord, 1981)

ความเพียงพอของการเทียบมาตรา

วิธีการเทียบมาตราแต่ละวิธี ประกอบด้วยตัวรูปแบบของการเทียบมาตรา (Model) ซึ่งมีข้อตกลงที่ว่าด้วยสมมติฐานเบื้องต้น (Assumptions) ของแต่ละรูปแบบ และประกอบด้วยวิธีการออกแบบ (Designs) เพื่อจัดเก็บข้อมูลให้เป็นไปตามข้อตกลงต่าง ๆ ถ้าหากทุกอย่างเป็นไปตามเงื่อนไข เหล่านี้แล้ว ผลของการเทียบมาตราจะมีความถูกต้อง (Accurate) และความแม่นยำ (Precise) ตามทฤษฎี แต่ความแท้จริงของการทดสอบมักไม่ได้เป็นไปตามอุดมการณ์ เพราะมีหลายสิ่งหลายอย่างอยู่นอกเหนือการควบคุม เช่น โปรแกรมการทดสอบระดับชาติ ซึ่งมีทั้งนโยบาย และกฎเกณฑ์เป็นตัวกำหนด จึงไม่สามารถควบคุมการทดสอบให้เป็นไปตามข้อตกลงเชิงทฤษฎี ตัวแบบสอบเองก็จำเป็นต้องมีการเปลี่ยนแปลงไป ดังนั้นข้อมูลที่น่ามาใช้ศึกษา จึงไม่สามารถระบุว่าเป็นตัวอย่างประชากรใดอย่างชัดเจน และโดยความเป็นจริงเป็นการจัดกระทำกับประชากรมากกว่า (Population quantities) ไม่ใช่ค่าประมาณ (Sample estimates) (Braun and Holland, 1982:10)

ด้วยสภาพความเป็นจริงดังกล่าว ทำให้การเทียบมาตราที่จัดกระทำอยู่นั้นอยู่ในภาวะที่มีเงื่อนไขน้อยกว่าความพอดีตามข้อตกลงในแต่ละรูปแบบ ดังนั้น การเทียบมาตราที่ได้พัฒนาขึ้น จึงจำเป็นต้องมีการตรวจสอบความเพียงพอของรูปแบบ (The adequacy of equating models) (Petersen, Marco and Stewart, 1982:71) วิธีการประเมินความเพียงพอมีผู้เสนอแนวคิดและวิธีการปฏิบัติไว้ดังนี้

1. ดัชนีตรวจสอบความเพียงพลของเจเกอร์ (Jaeger)

เจเกอร์ (Jaeger, 1981:26) ได้เสนอดัชนี 5 ตัว เพื่อตรวจสอบความเพียงพลของการใช้รูปแบบเชิงเส้นตรงว่า เทคนิควิธีที่ได้นำมาใช้ในการเทียบมาตรฐานเพียงพลกับการปรับความแตกต่างระหว่างการแจกแจงของคะแนนจากแบบสอบหรือไม่ หรือจำเป็นต้องมีรูปแบบอื่นที่มีความสลับซับซ้อนมากขึ้น ดัชนี 5 ตัว มีดังนี้ คือ

1.1 ดัชนีความคล้ายคลึงของการแจกแจงคะแนนสะสมของแบบสอบเก่าและชุดใหม่ โดยพิจารณาปรับความแตกต่างระหว่างค่าเฉลี่ยและส่วนเบี่ยงเบนมาตรฐาน การทดสอบความเหมือนของการแจกแจงใช้ทดสอบด้วย The Kolmogorov-smirnov two-sample test. (Smirnov, 1948 cited by Jaeger, 1981:27)

1.2 รูปแบบของการแจกแจงคะแนนดิบกับคะแนนแปลง (Shape of the raw score transformation) ดัชนีตัวนี้บอกความเพียงพลของการเทียบมาตรฐานด้วยรูปแบบเชิงเส้น ซึ่งสามารถอธิบายความแตกต่างในการแจกแจงคะแนนดิบทั้งสองชุดในการเทียบมาตรฐานด้วยรูปแบบเชิงเส้น ถ้าสามารถอธิบายความแตกต่างในการแจกแจงคะแนนดิบทั้งสองชุดอย่างเพียงพล ก็สามารถยอมรับได้ว่า การแปลงคะแนนดิบจากแบบสอบชุดใหม่ไปยังแบบสอบชุดเก่าเป็นเส้นตรงอย่างแน่นอน

1.3 ความคงเส้นคงวาของผลลัพธ์ของการเทียบมาตรฐานตามรูปแบบเชิงเส้นกับการเทียบที่ตำแหน่งเปอร์เซ็นต์ไทล์ (Consistency of linear and equipercentile equating results) การวิเคราะห์นี้ อาศัยข้อตกลงที่เป็นสมมติฐานเบื้องต้นที่ว่า ถ้ารูปแบบเชิงเส้นตรงมีความเพียงพลแล้ว ฟังก์ชันของรูปแบบการเทียบมาตรฐานที่ตำแหน่งเปอร์เซ็นต์ไทล์จะแปรผันไปโดยสุ่มรอบ ๆ ฟังก์ชันของรูปแบบเชิงเส้นตรงที่สมนัยกัน

1.4 ความคล้ายคลึงของการแจกแจงความยากของข้อสอบ (Similarity of item difficulty distributions) โดยอาศัยหลักการที่ว่า การที่ใช้รูปแบบเชิงเส้นตรงมีความเพียงพลอย่างแท้จริงกับแบบสอบที่มีคุณสมบัติเป็นคู่ขนานกัน ถ้ามีความเบี่ยงเบนจากความเป็คู่ขนานมากเท่าใด แสดงว่า ต้องการรูปแบบการเทียบมาตรฐานที่ซับซ้อนขึ้น เพราะการแจกแจงของแบบสอบที่ไม่ใช่คู่ขนานจะมีความแตกต่างเกิดขึ้นในระดับโมเมนต์ที่สูงขึ้น

1.5 ความคล้ายคลึงของค่าอำนาจจำแนกของข้อสอบ (Similarity of item discrimination distributions) มีเหตุผลท่านเองเดียวกับดัชนีตัวที่ 4

2. ดัชนีความแตกต่าง (Discrepancy indices)

ปีเตอร์สัน มาร์โค และ สตีเวอร์ท (Petersen, marco and Stewart, 1982) ได้เสนอวิธีประเมินความเพียงพอของรูปแบบการเทียบมาตรฐานโดยพิจารณาความแตกต่างระหว่างคะแนนแปลง (An estimated criterion score : t') ซึ่งเป็นผลจากการเทียบมาตรฐาน กับคะแนนเกณฑ์ (Criterion score : t) ที่สัมพันธ์กัน ถ้าความแตกต่างมีค่าน้อย มีความหมายว่า ความคลาดเคลื่อนที่เกิดขึ้นจากการใช้รูปแบบการเทียบมาตรฐานนั้นน้อยด้วย รูปแบบการเทียบมาตรฐานดังกล่าวย่อมมีความเหมาะสมกับสถานการณ์ที่ใช้

ดัชนีความแตกต่างที่นำมาเปรียบเทียบระหว่างรูปแบบและสถานการณ์ที่ต่าง ๆ กัน คือ กำลังสองของค่าเฉลี่ยของความแตกต่างที่ถ่วงน้ำหนักด้วยความแปรปรวนของคะแนนเกณฑ์ (The weighed mean-square difference) เป็นค่ามาตรฐาน ค่าที่คำนวณออกมานี้ เรียกว่า ความแปรปรวนรวม (The total error) ซึ่งมีสูตรดังนี้

$$\text{total error} = \sum f_j d_j^2 / nS_c^2 \dots\dots\dots (9)$$

$$\text{เมื่อ } d_j = (t - t')$$

$$n = \text{จำนวนคะแนนที่ใช้}$$

$$S_c^2 = \text{ความแปรปรวนของคะแนน } t$$

3. ดัชนีเปรียบเทียบเปอร์เซ็นต์ (The percentile comparison index) เป็นมาตรการวัดความไม่สอดคล้องระหว่างการแจกแจงของคะแนนในแบบสอบชุด X กับแบบสอบชุด Y ที่ได้แปลงไปอยู่ในมาตราคะแนนของ X แล้ว ตามวิธีที่ได้พัฒนาขึ้น ดัชนีเปรียบเทียบเปอร์เซ็นต์ คือ ค่าเฉลี่ยกำลังสองของความแตกต่าง (The mean-squared difference) ที่ได้จากการแจกแจงของคะแนนต่าง ๆ ของเกณฑ์ X กับคะแนนแปลง X^* ที่แปลงมาจาก Y ด้วยวิธีเทียบมาตราที่ระบุไว้ ณ ที่ตำแหน่งเปอร์เซ็นต์เดียวกัน ดัชนีได้เสนอโดย โคลเลน (Kolen, 1982) ได้แนะนำให้ใช้ข้อมูลจากกลุ่มตัวอย่างสอบทานผล ซึ่งได้สุ่มมาจากประชากรเดียวกันกับกลุ่มตัวอย่างที่ได้พัฒนาตารางคะแนนแปลง สูตรการคำนวณมีดังนี้ คือ

$$C = \frac{\sum (X_i - X_i^*)^2}{nk} \dots\dots\dots (10)$$

- เมื่อ n คือ จำนวนของคะแนนดิบของกลุ่มสอบทานผล
 k คือ จำนวนข้อสอบในแบบสอบรวมที่ใช้

ค่า C ที่ได้ ถ้ามีค่าน้อยจะให้ความหมายว่า รูปแบบการเทียบมาตราที่นำมาสร้างตารางคะแนนแปลงนี้มีความเหมาะสมและเพียงพอที่จะให้ผลการแปลงคะแนนอย่างคงเส้นคงวา

วิธีการประเมินความเพียงพอที่ได้กล่าวมานี้ อาจจำแนกเป็นสองพวก พวกแรกเป็นการประเมินก่อนดำเนินการ เช่น ดัชนีความคล้ำยคลึงของการแจกแจงคะแนนสะสมของแบบสอบสองชุด ดัชนีความคล้ำยคลึงของการแจกแจงของค่าความยากของข้อสอบ เป็นต้น พวกหลังเป็นการประเมินผลของการเทียบมาตรา ซึ่งอาศัยคะแนนเกณฑ์ที่เลือกสรรแล้วเป็นหลักในการเทียบหาความแตกต่าง สำหรับดัชนีที่เสนอโดย ปีเตอร์สันและคณะนั้น (1982) คะแนนเกณฑ์ที่ใช้ คือ ผลการแปลงคะแนนด้วยรูปแบบอิงทฤษฎีตอบสนองข้อสอบ

ค่าดัชนีที่คำนวณมีลักษณะเป็นหน่วยมาตรฐานแล้ว สามารถนำค่าเหล่านี้ที่ได้จากการใช้รูปแบบที่ต่างกัน ตลอดจนสถานการณ์ที่ได้ข้อมูลที่แตกต่างกันมาเปรียบเทียบกันโดยตรงได้ ส่วนดัชนีที่เสนอโดย โคลเลน (Kolen, 1982) ได้ใช้ข้อมูลจากผู้สอบเองเป็นเกณฑ์ในการหาความแตกต่าง ข้อมูลเหล่านี้ได้จากการออกแบบด้วยการใช้กลุ่มตัวอย่างสอบทานผล ซึ่งผู้สอบในกลุ่มตัวอย่างนี้ได้รับการทดสอบด้วยแบบสอบทั้งสองชุด ดังนั้น การใช้คะแนนของตนเองเป็นเกณฑ์จึงมีความเป็นอิสระ ไม่ขึ้นกับกระบวนการแปลงคะแนนอื่น ๆ เช่นวิธีที่เสนอโดยปีเตอร์สันและคณะ ด้วยเหตุผลดังกล่าว ผู้วิจัยจึงได้เลือกวิธีการประเมินความเพียงพอจากการวิเคราะห์ผลในกลุ่มตัวอย่างสอบทานผล แต่การหาค่าดัชนีนี้ได้ดัดแปลงจากสูตรของ โคลเลน ไปให้ตามแนวความคิดของปีเตอร์สันและคณะ คือ ใช้ค่าความแปรปรวนเป็นตัวถ่วงน้ำหนักเพื่อให้ค่าที่ได้มีหน่วยเป็นมาตรฐาน

สูตรที่ดัดแปลงและนำมาใช้ในการศึกษาคั้งนี้ คือ

$$C = \frac{\sum (X - X_1^m)^2}{nS_x^2} \dots\dots\dots (11)$$

- เมื่อ X_1 คือ คะแนนเกณฑ์หรือคะแนนจากการสอบชุด X ของคนที่ 1
 X_1^m คือ คะแนนที่ได้เทียบด้วยตารางคะแนนที่สัมพันธ์กันของคนี่ 1
 n คือ จำนวนคนในกลุ่มตัวอย่างสอบทานผลที่นำมาวิเคราะห์
 S_x^2 คือ ค่าความแปรปรวนของคะแนน X

ศูนย์วิทยทรัพยากร
 จุฬาลงกรณ์มหาวิทยาลัย

งานวิจัยที่เกี่ยวข้องกับการเทียบมาตรฐาน

งานวิจัยที่เกี่ยวข้องกับการเทียบมาตรฐานโดยอิงทฤษฎีตอบสนองข้อสอบในต่างประเทศ สแกจส์และลิทซ์ซิล (Skaggs Gary and Lissitz Robertes .1986:498-501) ได้สรุปไว้ดังนี้

ตารางที่ 1 การเทียบมาตรฐานตามแนวคิดโดยอิงทฤษฎีตอบสนองข้อสอบ

ปีที่ศึกษา	ผู้ศึกษา	แบบสอบ	กลุ่ม ตัวอย่าง	วิธีการ	ตัวแปรที่เทียบมาตรฐาน	วิธีประเมิน
1968	Wright	LSAT	นักเรียน กฎหมาย (980)	ราส์ซ	แบบสอบฉบับง่ายกับฉบับ ยาก และความสามารถ ของกลุ่มต่ำกับกลุ่มสูง	- เปรียบเทียบความ แตกต่างมาตรฐาน ของแบบสอบ - ลักษณะโค้งความ สัมพันธ์ b
1968	Ander- son, Kearney และ Everett	แบบสอบ เช่าวี- ปัญญา	ทหาร (610 875	ราส์ซ	ความยากของแบบสอบ จากกลุ่มตัวอย่าง 2 กลุ่ม	ความสัมพันธ์ b

ตารางที่ 1 (ต่อ) การเทียบมาตรฐานตามแนวตั้งโดยอิงทฤษฎีตอบสนองข้อสอบ

ปีการศึกษา	ผู้ศึกษา	แบบสอบ	กลุ่มตัวอย่าง	วิธีการ	ตัวแปรที่เทียบมาตรฐาน	วิธีประเมิน
1975	Tinsley และ Dawis	คำอุปมาอุปมัย	นักเรียนมัธยมศึกษา, ลูกจ้างชั้นดี	ราล์ฟ	ความยากของข้อสอบจากกลุ่มตัวอย่างที่ต่างกัน	ความสัมพันธ์ b
1978	Slinde และ Linn	แบบสอบผลสัมฤทธิ์ทางคณิตศาสตร์	นักศึกษาปีที่ 1 (390-810)	ราล์ฟ	เทียบมาตรฐานในแนวตั้งจากกลุ่มตัวอย่างก่อนที่มีความสามารถต่างกัน	ความแตกต่างมาตรฐานของความสามารถ
1979	Slinde และ Linn	ความเข้าใจทางการอ่านและคำศัพท์	นักเรียนระดับ 5 (510-570)	ราล์ฟ	เทียบมาตรฐานในแนวตั้งจากกลุ่มตัวอย่างก่อนที่มีความสามารถต่างกัน	ความแตกต่างมาตรฐานของความสามารถ
1979	Marco, Peterson และ Stewart	SAT-ภาษา	นักเรียนมัธยม (ส่วนมาก) (1580)	ราล์ฟ, 3 พารา มิเตอร์, อีควิเปอร์ เซนไต์ล် และเท็ง-เส้นตรง	เทียบมาตรฐานในแนวตั้ง, กลุ่มตัวอย่างที่สัมพันธ์ คล้ายกัน, แบบสอบร่วมภาษาออกกับภาษาใน, ความยากของแบบสอบร่วม, IRT กับอีควิเปอร์ เซนไต์ล်, เปรียบเทียบรูปแบบการเทียบมาตรฐาน	ค่าความคลาดเคลื่อนกำลังสองและดัชนีความลำเอียง

ตารางที่ 1 (ต่อ) การเทียบมาตรฐานตามแนวตั้งโดยอิงทฤษฎีจุดประสงค์ของข้อสอบ

ปีการศึกษา	ผู้ศึกษา	แบบสอบ	กลุ่ม ตัวอย่าง	วิธีการ	ตัวแปรที่เทียบมาตรฐาน	วิธีประเมิน
1980	Lord และ Hoover	ITBS, ความเข้าใจ- จิตทางคณิต ศาสตร์	นักเรียน ระดับ 6 - 8 (245-300)	ราส์ช	เทียบมาตรฐานในแนวตั้ง จากตัวอย่างที่มีระดับ ความสามารถต่างกัน	เปรียบเทียบ จากกราฟ
1981	Guskey	ITBS: ความ เข้าใจใน การอ่าน	นักเรียน ที่มีความ สามารถ สูงใน ระดับ 6-8	ราส์ช	เทียบมาตรฐานในแนวตั้ง: เปรียบเทียบสเกลความ สามารถโดยราส์ชกับ คะแนนเกรดที่เทียบเท่า กัน	เปรียบเทียบ จากกราฟของ คะแนนเดิมกับ คะแนนแปลง
1981	Divgi	MAT: การอ่าน	นักเรียน ระดับ 6-8 (5500)	ราส์ช	เทียบมาตรฐานในแนวตั้ง ของแบบสอบก่อน	ความแตกต่าง มาตรฐาน, ความล้มเหลว จากฟังก์ชันของ คะแนนรวม

ตารางที่ 1 (ต่อ) การเทียบมาตรฐานตามแนวคิดอิงทฤษฎีผู้ตอบสนองข้อสอบ

ปีที่ศึกษา	ผู้ศึกษา	แบบสอบ	กลุ่มตัวอย่าง	วิธีการ	ตัวแปรที่เทียบมาตรฐาน	วิธีประเมิน
1981	Kolen	ITED: คำศัพท์เกี่ยวกับความยาวของสระ	นักเรียนระดับ 9-12 (1580-1925)	รอสส์, 2, 3 พารามิเตอร์, อีควิเปอร์เซ็นต์, เส้นตรง	เทียบมาตรฐานในแนวนอนและแนวดิ่ง, เปรียบเทียบโมเดลการเทียบมาตรฐาน, คำศัพท์กับความยาวของสระ	Cross-Validation, ค่าเฉลี่ยกำลังสองของความคลาดเคลื่อน
1981	Peter-son, Cook, Stocking	SAT: ภาษาและคณิตศาสตร์	นักเรียนมัธยมศึกษา (2670)	3 พารามิเตอร์, เส้นตรง, อีควิเปอร์เซ็นต์	เปรียบเทียบรูปแบบการเทียบมาตรฐาน, เปรียบเทียบเมตริกซ์การแปลงของ 3 พารามิเตอร์, ทดสอบภาษากับคณิต	ความคลาดเคลื่อนกำลังสองเฉลี่ยและดัชนีความล่าช้าเชิง, แบบสอบเทียบมาตรฐาน
1981	Cook, Dubar, Eignor	SAT/PSAT ภาษาและคณิตศาสตร์	นักเรียนมัธยมศึกษา (2000)	3 พารามิเตอร์, เส้นตรง, อีควิเปอร์เซ็นต์	เปรียบเทียบรูปแบบการเทียบมาตรฐาน, ทดสอบภาษาและคณิตศาสตร์	ความคลาดเคลื่อนกำลังสองเฉลี่ย-การเทียบมาตรฐานโดยให้ 3 พารามิเตอร์เป็นเกณฑ์



ตารางที่ 1 (ต่อ) การเทียบมาตรฐานตามแนวตั้งโดยอิงทฤษฎีตอบสนองข้อสอบ

ปีที่ศึกษา	ผู้ศึกษา	แบบสอบ	กลุ่ม ตัวอย่าง	วิธีการ	ตัวแปรที่เทียบมาตรฐาน	วิธีประเมิน
1981	Kolen, Whitney	GED Test	ผู้ใหญ่ (200)	3 พารา- มิเตอร์, ราล์ฟ, เส้นตรง, อิกวิเปอร์ เซนไคล์	เปรียบเทียบรูปแบบการ เทียบมาตรฐานด้วยการ สลับฉบับกัน	ความคลาดเคลื่อน กำลังสองเฉลี่ย Cross-Vali- dation
1981	Cowell	TOELT- SLEP	นักเรียน มัธยม- ศึกษา ($n_1=290-320$ $n_2=2070-3170$)	3 พารา- มิเตอร์, ราล์ฟ, เส้นตรง	กลุ่มตัวอย่างใหญ่และ เล็ก, เปรียบเทียบรูปแบบการ เทียบมาตรฐาน ด้วยการสลับฉบับกัน	สถิติโดยสรุป ต่าง ๆ-วิธีการ เปรียบเทียบ รายคู่
1981	Patience	TTED-วิธี เขียน	นักเรียน ระดับ 9 -10 (1000)	ราล์ฟ, 2, 3 พารา- มิเตอร์, อิกวิเปอร์ เซนไคล์	เปรียบเทียบการเทียบ มาตรฐานแนวตั้ง	ความสัมพันธ์ ระหว่างมาตร คะแนนของผู้- สอบด้วยจุดเริ่ม ต้น, เกณฑ์บนพื้น ฐานวิธีเทียบ มาตรฐานแบบอิกวิ เปอร์เซนไคล์ กับคะแนนรวม

ตารางที่ 1 (ต่อ) การเทียบมาตรฐานตามแนวคิดโดยอิงทฤษฎีตอบสนองข้อสอบ

ปีที่ศึกษา	ผู้ศึกษา	แบบสอบ	กลุ่ม ตัวอย่าง	วิธีการ	ตัวแปรที่เทียบมาตรฐาน	วิธีประเมิน
1982	Modu	CB-Adv. ฟิลิกส์กาส- ภาพ	นักเรียน ม.ต้น	3 พารา- มิเตอร์, เส้นตรง อิลควิเปอร์ เช่นไคล์	เทียบมาตรฐานของฟอร์มที่ เลือกภายใต้เงื่อนไข Multidimensiona- lity	เปรียบเทียบ กราฟของฟังก์ชัน การเทียบมาตรฐาน
1983	Holmes, Doody, Bogan	CTBS-ศัพท์ ความเข้าใจ ในการ อ่าน	นักเรียน ระดับ 4 5, 8, 9	3 พารา- มิเตอร์	วิธีการ 3 วิธีของการ สร้างเมตริกซ์ธรรมชาติ	ค่าความคลาด- เคลื่อนกำลัง สองประสิทธิ์ไป สู่กลุ่มตามผล
1983	Doody, Bogan, Yen	การแปล	(2000)	3 พารา- มิเตอร์	-เทียบมาตรฐานในแนวคิด -เมตริกซ์ข้อมูลมิติเดียว และหลายมิติ	เปรียบเทียบ โดยประมาณ ความสามารถ จากแบบสอบข้อสอบ
1983	Loyd	ITBS ความคิด รวบยอด ทางคณิต- ศาสตร์	นักเรียน ระดับ 3-8 (250)	ราล์ฟ, 3 พารา- มิเตอร์	เปรียบเทียบรูปแบบ การเทียบมาตรฐาน, แบบ สอบร่วมภายในและภายใน นอก, การเทียบมาตรฐาน ในแนวคิด	เทียบโดยใช้ กราฟ

ตารางที่ 1 (ต่อ) การเทียบมาตรฐานแนวตั้งโดยอิงทฤษฎีตอบสนองข้อสอบ

ปีที่ศึกษา	ผู้ศึกษา	แบบสอบ	กลุ่มตัวอย่าง	วิธีการ	ตัวแปรที่เทียบมาตรฐาน	วิธีประเมิน
1984	Cook, Eignor, Taft	ชีววิทยา	นักเรียนมัธยมศึกษา (2400-3900)	3 พารามิเตอร์, เชิงเส้น-ตรง, อิกวิเปอร์เซนไคล์	IRT และแบบดั้งเดิม เทียบบนฐานของกลุ่มตัวอย่างที่ต่างกัน, เปรียบเทียบโมเดล การเทียบมาตรฐาน	ฟังก์ชันของการเปรียบเทียบ การเทียบมาตรฐาน
1985	Harris, Kolen	ACT- การใช้คณิตศาสตร์	นักเรียนมัธยมศึกษา (3870-3970)	3 พารามิเตอร์, อิกวิเปอร์เซนไคล์, เชิงเส้น-ตรง	เปรียบเทียบโมเดลการเทียบมาตรฐาน, กลุ่มตัวอย่างที่มีความสามารถสูง-ต่ำ	เปรียบเทียบจากกลุ่มที่แตกต่าง - ความคลาดเคลื่อน กำลังสองและดัชนีความล่าช้า

สไลด์และลินด์ (Slinde and Linn, 1977-1979) ได้ทำการตรวจสอบปัญหาของการเทียบมาตรฐานแนวตั้ง (Vertical Equating) ของแบบสอบสองชุดที่สร้างขึ้นเพื่อใช้กับประชากรที่มีความสามารถต่างกันด้วยวิธีอ้อม จากการศึกษาได้ให้ข้อเสนอแนะว่า การใช้รูปแบบการเทียบมาตรฐาน 3 รูปแบบ คือ เชิงเส้นตรง อิกวิเปอร์เซนไคล์ และอิงทฤษฎีตอบสนองข้อสอบชนิดหนึ่งพารามิเตอร์ รูปแบบโลจิสต์อาจมีข้อจำกัดในกระบวนการเทียบมาตรฐานแนวตั้ง โดยเฉพาะอย่างยิ่งเมื่อทำการเทียบมาตรฐานโดยใช้แบบสอบสองชุดที่มีความแตกต่างกันมาก และกลุ่มตัวอย่าง

สองกลุ่มที่มีระดับความสามารถต่างกันมาก จากการศึกษานี้ได้ให้ข้อสังเกตว่า ถ้าใช้รูปแบบอิงทฤษฎีตอบสนองข้อสอบชนิดสามพารามิเตอร์ของโลจิสต์ อาจให้ผลการเทียบมาตรฐานที่ต่ำกว่าในสถานการณ์เช่นนั้น

มาร์โค ปีเตอร์สัน และสตีเวอร์ท (Marco, Petersen and Stewart, 1979) ได้ประเมินและเปรียบเทียบวิธีการเทียบมาตรฐานเชิงเส้นตรง วิธีการเทียบมาตรฐานโดยวิธีเปอร์เซ็นต์ และวิธีการใช้โด่งลักษณะข้อสอบโดยวิธีแบบสอบ SAT ผลปรากฏว่าแบบสอบที่มีความยากแตกต่างกัน วิธีการเทียบมาตรฐานโดยวิธีโด่งลักษณะข้อสอบให้ผลดีที่สุด และวิธีเทียบมาตรฐานเชิงเส้นตรงให้ผลน้อยที่สุด และเมื่อปีเตอร์สันและคณะทำซ้ำอีก ในปี ค.ศ. 1983 โดยวิเคราะห์ซ้ำกับปี ค.ศ. 1979 ผลคงปรากฏเหมือนเดิม แต่เมื่อโคเลน (Kolen, 1981) ได้ศึกษารูปแบบการเทียบมาตรฐานระหว่างรูปแบบดั้งเดิม 2 วิธี คือ รูปแบบเชิงเส้นตรงและรูปแบบอควิเปอร์เซ็นต์ กับรูปแบบการเทียบมาตรฐานโดยวิธีโด่งลักษณะข้อสอบที่ศึกษาทั้ง 1, 2, 3 พารามิเตอร์ แต่ละวิธีเทียบด้วยค่าประมาณจากค่าจริง (Estimated true score equating) และเทียบด้วยการประมาณจากค่าสังเกต (Estimated observed score equating) ให้ข้อมูลจากโครงการ IOWA Test of Education Development (ITED) ค.ศ. 1978 ศึกษา 2 รูปแบบ คือ เทียบจากฟอร์มที่มีความยากเท่าเทียมกัน และเทียบจากฟอร์มที่มีความยากแตกต่างกัน กลุ่มตัวอย่างเป็นนักเรียนระดับที่ 1 คือ ระดับเกรด 9 และ 10 ระดับ 2 คือ ระดับเกรด 11 และ 12 จำนวนทั้งหมด 10,728 คน จากโรงเรียนในรัฐไอโอวา 34 แห่ง การศึกษานี้ใช้เกณฑ์ Cross-Validation criterion ผลการวิจัยพบว่า การเทียบมาตรฐานทั้งหมดแตกต่างกันอย่างมีนัยสำคัญทางสถิติที่ระดับ .01 และ .05 การเทียบมาตรฐานโดยวิธีอควิเปอร์เซ็นต์เหมาะสมที่สุดกับแบบสอบที่มีความยากแตกต่างกัน ส่วนการเทียบมาตรฐานโดยวิธีโด่งลักษณะข้อสอบพบว่า กรณีใช้พารามิเตอร์เดียว ทั้งเป็นวิธีที่ไม่ถูกต้องสำหรับแบบสอบที่ต่างกัน ซึ่งเนื่องมาจากเกิดการเดาเกิดขึ้นมาก ส่วนกรณีใช้สามพารามิเตอร์ มีปัญหาที่เกิดขึ้นอย่างหลีกเลี่ยงไม่ได้ คือ การประมาณค่าไม่ถูกต้องของ Lower asymptote parameter

คูก ดันบาร์ และไอเนอร์ (Cook, Dunbar and Eignor, 1981) ได้ศึกษาผลการเทียบมาตรฐานแบบใช้ทฤษฎีตอบสนองข้อสอบ เปรียบเทียบกับการเทียบมาตรฐานแบบดั้งเดิม 2 วิธี คือ การเทียบมาตรฐานเชิงเส้นตรง และการเทียบมาตรฐานแบบอควิเปอร์เซน-ไสต์ การเทียบมาตรฐานโดยใช้ทฤษฎีตอบสนองข้อสอบให้โมเดลสามพารามิเตอร์ ซึ่งประมาณค่าพารามิเตอร์โดยให้โปรแกรม LOGIST ด้วยวิธีความเป็นไปได้สูงสุด (Maximum likelihood procedure) สำหรับการประมาณค่าที่อยู่ต่ำกว่าระดับการเดา ใช้วิธีการเชิงเส้นตรง (Linear interpolation) ในการเทียบมาตรฐานศึกษาทั้งสองกรณี คือ มีแบบสอบรวม และไม่มีแบบสอบรวม โดยให้ประชากรที่ไม่ได้สุ่มมาจากประชากรเดียวกัน ทั้งนี้เพื่อศึกษาความเป็นไปได้ (Feasibility) ของวิธีการที่ใช้ทฤษฎีตอบสนองข้อสอบ ในการนำไปแก้ปัญหาทางปฏิบัติ การเปรียบเทียบพิจารณา 2 ประการ คือ ความสอดคล้องสัมพันธ์ (Relative agreement) ระหว่างวิธีที่ใช้ทฤษฎีตอบสนองข้อสอบกับวิธีดั้งเดิมแต่ละวิธี โดยพิจารณาจากการพล็อตกราฟ ลักษณะประการหนึ่งคือ ค่าความแตกต่างสำหรับการแจกแจงคะแนนรวม และแต่ละส่วนของการแจกแจง (สูงกว่า 20 เปอร์เซ็นต์, ตรงกลาง 60 เปอร์เซ็นต์ และต่ำกว่า 20 เปอร์เซ็นต์ของการแจกแจง) ซึ่งคำนวณจากค่าที่ถ่วงน้ำหนักของค่าเฉลี่ยกำลังสองของความแตกต่างระหว่างค่าประมาณจากการเทียบมาตรฐานแต่ละวิธีที่ศึกษากับคะแนนเกณฑ์ เกณฑ์ที่ใช้ศึกษาในที่นี้ คือ คะแนนแปลงด้วยการเทียบมาตรฐานที่ใช้ทฤษฎีตอบสนองข้อสอบชนิดสามพารามิเตอร์ ซึ่งงานวิจัยที่มีมาก่อน (Kolen, 1981 ; Slind and Linn, 1977) แนะนำว่ามีความเหมาะสมในการเทียบมาตรฐานในกรณีที่มีแบบสอบมีความแตกต่างกัน สอดคล้องกับกลุ่มตัวอย่างที่มีความสามารถต่างกัน

ข้อมูลที่ใช้ในการศึกษานี้ เป็นผลการสอบสองครั้งที่ดำเนินการแล้ว โดย The Collage Board Admissions Testing Program ซึ่งใช้แบบสอบ 2 ชุด ที่มีความแตกต่างในความยาวและความยาก แต่ละชุดมีข้อสอบ 2 ตอน คือ ข้อสอบภาษาและคณิตศาสตร์ ให้กับกลุ่มตัวอย่างที่ไม่ได้เท่าเทียมกันโดยการสุ่ม

ผลการศึกษาพบว่า วิธีการเทียบมาตรฐานแบบดั้งเดิม กับวิธีเทียบมาตรฐานโดยใช้ทฤษฎีตอบสนองข้อสอบมีความสอดคล้องกันมาก เนื่องจากโดยความจริงแล้ว การแจกแจงคะแนนดิบจากแบบสอบทั้งสองชุดมีรูปร่างคล้ายคลึงกันมาก จึงทำให้รูปแบบเชิงเส้นตรง

และเชิงเส้นโค้งหรือลิวเปอร์ เช่น ไซล์มีผลใกล้เคียงกัน โดยการเทียบเชิงเส้นโค้งมีผลใกล้เคียงกับการเทียบแบบใช้ทฤษฎีตอบสนองข้อสอบมากกว่าการเทียบเชิงเส้นตรง ล่าง-ไรก็ตาม ความแตกต่างที่เกิดขึ้นที่ส่วนปลายของการเทียบแบบดั้งเดิมทุกวิธีที่ต่างไปจากการเทียบที่ใช้ทฤษฎีตอบสนองข้อสอบ เป็นผลจากข้อมูลส่วนนั้นน้อยและหาชานมาก ในทางทฤษฎีจึงแนะนำให้ใช้การเทียบที่ใช้ทฤษฎีตอบสนองข้อสอบเพราะไม่มีผลเนื่องจากขาดแคลนข้อมูลในส่วนปลายของการแจกแจง ล่างไรก็ตาม การเทียบที่ใช้ทฤษฎีตอบสนองข้อสอบก็ไม่สามารถเตรียมการเทียบคะแนนในส่วนที่อยู่ปลายด้านล่างที่ต่ำกว่าระดับการคาดได้ ซึ่งจำเป็นต้องใช้วิธีอื่นประมาณค่าเพิ่มเติม นอกจากนี้ผลการวิจัยพบว่า การเทียบมาตราโดยใช้ทฤษฎีตอบสนองข้อสอบโดยตรงเมื่อไม่มีแบบสลับร่วม และเมื่อกลุ่มตัวอย่างไม่ได้เป็นกลุ่มที่เท่าเทียมกันโดยการใช้วิธีนั้น ให้ผลแตกต่างจากการเทียบแบบมีแบบสลับร่วมเพียงเล็กน้อย

เมื่อโคเลนและวิทนี (Kolen and Whitney, 1982) ได้เปรียบเทียบความถูกต้องของการเทียบมาตราโดยวิธีแบบสลับร่วม 4 วิธี คือ วิธีลิวเปอร์เช่นไซล์ วิธีเชิงเส้นตรง พารามิเตอร์ตัวเดียว และพารามิเตอร์สามตัว ใต้แบบสอบ General Education Development (GED) ซึ่งเป็นแบบสอบผลสัมฤทธิ์ เพื่อใช้ในการตัดสินให้ประกาศนียบัตรที่ต้องการเทียบความวุฒาระดับเตรียมอุดมศึกษาทั่วประเทศ แบบสอบมี 12 ชุด โดยวิธีฟอร์มที่ 12 เป็นแบบสลับร่วม ที่เหลืออีก 11 ฟอร์ม เป็นแบบสอบที่ใช้เทียบมาตรา ผู้เข้าสอบแต่ละคนต้องทำแบบสอบ 2 ฟอร์ม ในจำนวน 11 ฟอร์ม ส่วนฟอร์มที่ 12 ต้องทำทุกคน กลุ่มตัวอย่างมาจากการสอบปี 1980 จำนวนมากกว่า 800,000 คน ทำการสุ่มแบบแบ่งชั้นหลายชั้นจากประเภทของโรงเรียน เขตทางภูมิศาสตร์ สถานภาพทางสังคม สดักทำขั้วได้นักเรียนที่เป็นตัวอย่างโรงเรียนละ 22 คน จำนวนคนในแต่ละชุดของแบบสอบประมาณ 200 คน ประมาณค่าพารามิเตอร์โดยวิธีโปรแกรม LOGIST แล้วสร้างตารางสมมูลระหว่างแบบสอบอีก 11 ชุด การเทียบผลการวิเคราะห์โดยวิธี Cross-Validation ผลการศึกษาพบว่า วิธีการเทียบมาตราโดยวิธีลิวเปอร์เช่นไซล์และวิธีโค้งลักษณะข้อสอบได้ผลไม่เป็นที่ยอมรับ ขณะเดียวกันวิธีการเทียบมาตราโดยวิธีพารามิเตอร์ตัวเดียวให้ผลเพียงพอเทียบเท่ากับรูปแบบสามพารามิเตอร์ ในการประเมินผลการศึกษา

ครั้งนี้มีการเปรียบเทียบอื่น ๆ ซึ่งสรุปได้ว่า ความถูกต้องของการเทียบมาตรฐานขึ้นอยู่กับองค์ประกอบหลายอย่าง เช่น ลักษณะของแบบสอบ รูปแบบของการเทียบมาตรฐาน ขนาดของกลุ่มตัวอย่าง เป็นต้น

ฮัทเทน (Hutten, 1982) ได้ศึกษาความเหมาะสมของข้อมูลจริงกับรูปแบบทฤษฎีตอบสนองข้อสอบ 2 รูปแบบ คือ รูปแบบของราส์ช กับรูปแบบ 3 พารามิเตอร์ ในการประมาณค่าพารามิเตอร์ ความสามารถ และความยากของข้อสอบ สำหรับกลุ่มตัวอย่างขนาดเล็ก (250 คน) และแบบสอบสั้น (20 ข้อ) ผลการศึกษาพบว่า รูปแบบของราส์ช และรูปแบบ 3 พารามิเตอร์ มีความเหมาะสมกับข้อมูลประมาณร้อยละ 80 ทั้งสองรูปแบบเมื่อเปรียบเทียบคะแนนที่คำนวณกับคะแนนที่วัดได้ โดยใช้สถิติ Kolmogorov-Smirnov แล้ว รูปแบบของราส์ชเหมาะสมกับข้อมูลทั้งหมดดีกว่ารูปแบบ 3 พารามิเตอร์ ผู้วิจัยสนับสนุนให้ใช้รูปแบบ 3 พารามิเตอร์ กับกลุ่มตัวอย่างที่มีขนาดตั้งแต่ 1,000 คนขึ้นไป เพื่อจะได้ความแม่นยำในการประมาณค่าพารามิเตอร์ต่าง ๆ ผลการวิจัยนี้สอดคล้องกับที่ดักลาสส์ (Douglass, 1981) ได้ศึกษาเปรียบเทียบการใช้รูปแบบทฤษฎีตอบสนองข้อสอบกับผลการสอบในระดับห้องเรียน โดยใช้รูปแบบทฤษฎีตอบสนองข้อสอบประมาณค่าพารามิเตอร์ของข้อสอบ 100 ข้อ ที่เป็นข้อสอบปลายปีของ 4 ปีการศึกษา จากกลุ่มตัวอย่างระดับมหาวิทยาลัยจำนวน 594-1,082 คน จัดข้อสอบแบ่งเป็น 4 ชุด จากนั้นทำการสุ่มข้อสอบชุดละ 43-53 ข้อ สุ่มผู้สอบจำนวน 200, 500 และ 800 คน มาวิเคราะห์ ปรากฏว่าการวิเคราะห์ข้อสอบด้วยรูปแบบของราส์ชมีความคงที่ของค่าพารามิเตอร์ดีกว่ารูปแบบ 3 พารามิเตอร์ และเมื่อมีการเทียบคะแนนของแบบสอบข้ามกลุ่มตัวอย่าง รูปแบบของราส์ชก็มีความคงที่ดีกว่าด้วย

ลอร์ด และวิงเกอร์สกี (Lord and Wingersky, 1984) ได้เปรียบเทียบการเทียบมาตรฐานระหว่างรูปแบบทฤษฎีตอบสนองข้อสอบที่เทียบโดยใช้คะแนนจริง และรูปแบบทฤษฎีตอบสนองข้อสอบที่เทียบคะแนนสังเกตโดยใช้วิธีอควิเพลร์ เชนไคล์ จากแบบสอบ SAT ด้านภาษา จำนวน 6 ชุด เทียบเข้าสู่อันในลักษณะไอ้ โดทท์แบบสอบร่วม ทั้งนี้แบบสอบชุดแรกและชุดสุดท้ายเป็นชุดเดียวกัน การออกแบบรวบรวมข้อมูลดำเนินการดังภาพ

ภาพที่ 2 การออกแบบรวบรวมข้อมูลของลอร์ดและวังเกอร์สกี

V4	fe
	fe X2

X2	fm
	fm Y3

Y3	fw
	fw B3

B3	fk
	fk Y2

Y2	fu
	fu Z5

Z5	et.
	et. V4

จากภาพ แบบสอบ V4 เทียบไปสู่แบบสอบ X2 โดสที่ใช้แบบสอบร่วม fe ส่วนแบบสอบ X2 เทียบไปสู่แบบสอบ Y3 โดสที่ใช้แบบสอบร่วม fm ซึ่งจะได้ผลการเทียบจาก V4 ไป Y3 กระทำการเทียบคะแนนในลักษณะเดียวกันไปยังแบบสอบชุดอื่น สุดท้ายจะได้ผลการเทียบ V4 กับ V4 ซึ่งเป็นการเทียบแบบสอบกับตัวเอง และใช้เป็นเกณฑ์ในการเปรียบเทียบความคงเส้นคงวาของวิธีการเทียบมาตรฐานทั้ง 2 วิธี ทั้งนี้เกณฑ์ที่แตกต่างกันจากที่โคเลน (Kolen, 1981) ได้เคยใช้ศึกษาความคงเส้นคงวาของวิธีเทียบมาตรฐานต่างๆ โคลนใช้เกณฑ์จากกลุ่มทานผลซึ่งเป็นกลุ่มอิสระอีกกลุ่มหนึ่ง โดสผู้วิจัยกล่าวไว้ว่า เกณฑ์ที่โคเลนใช้ ไม่ใช่เกณฑ์ที่มีคุณสมบัติสำหรับการเลือกวิธีการเทียบมาตรฐานที่ดีที่สุด เพราะวิธีเทียบมาตรฐานที่ไม่ถูกต้อง อาจให้ผลการเทียบคงเส้นคงวามากกว่าวิธีการที่ถูกต้อง เกณฑ์การเทียบแบบสอบกับตัวเองที่ใช้ในการศึกษานี้ อาจไม่เที่ยงตรงในเชิงทฤษฎี และไม่น่าพอใจในการปฏิบัติสำหรับการเทียบแบบสอบสองฉบับที่แตกต่างกัน แต่เป็นเกณฑ์ที่มีประโยชน์ต่อการเปลี่ยนวิธีการเทียบมาตรฐานใหม่ หากพบว่าคะแนนที่แปลงแล้วในขั้นสุดท้ายไม่เป็นผลอย่างเดียวกันกับคะแนนเดิมของตนเอง

ผลการวิจัยยังไม่มีการตัดสินใจเชิงทฤษฎีที่ชัดเจนในการเลือกวิธีการเทียบมาตรฐานแบบทฤษฎีตอบสนองข้อสอบที่ใช้คะแนนจริง กับแบบทฤษฎีตอบสนองข้อสอบที่ใช้คะแนนสังเกต โดสอิวเปอร์เซ็นต์ เนื่องจากทั้งสองวิธีนี้มีความสอดคล้องกัน

โดรานส์ และคิงส์ตัน (Dorrans and Kingston, 1985) ได้ศึกษาผลการผ่านเป็นเกี่ยวกับความมีมิติเดียว ในการประมวลค่าพารามิเตอร์ของข้อสอบ และพารามิเตอร์ความสามารถ โดสใช้การเทียบมาตรฐานแบบสามพารามิเตอร์กับแบบสอบวัดความถนัดชุด GRE ด้านภาษา จำนวน 4 ชุด แต่ละชุดวัด 2 ด้าน คือ ความเข้าใจในการอ่าน และหลักภาษา การออกแบบการเทียบมาตรฐานมี 2 แบบ คือ เทียบคะแนนรวมทั้งชุดซึ่งมีความเป็นวิวิธพันธ์ (Heterogeneous) และเทียบคะแนนรวมแต่ละด้านซึ่งมีความเป็นเอกพันธ์ (Homogeneous) โดยการออกแบบการรวมรวมข้อมูล เมื่อใช้กลุ่มที่เท่าเทียมกัน และเมื่อใช้แบบสอบร่วม จำนวนกลุ่มตัวอย่างมีตั้งแต่ 2,579 ถึง 4,351 คน

ผลการวิจัยพบว่า การเทียบมาตรฐานแบบทฤษฎีตอบสนองข้อสอบ เมื่อผ่านข้อตกลงเกี่ยวกับแบบสลับต้องมิตีเดีวั้นนั้น อาจมีผลกระทบ (impact) ต่อการเทียบมาตรฐาน แต่ไม่อาจยืนยันได้ เนื่องจากมีความคล้ายคลึงกัน ระหว่างผลการเทียบเมื่อก่อนข้อสอบเป็นเอกพจน์และวิวิธพันธ์ ดังนั้นผู้วิจัยจึงแนะนำว่า การเทียบมาตรฐานโดยวิธีทฤษฎีตอบสนองข้อสอบมีความแกร่งเพื่อส่งผลต่อการผ่านการมีมิตีเดีว

สำหรับงานวิจัยที่เกี่ยวกับการเทียบมาตรฐานในประเทศไทย มีผู้ศึกษาเกี่ยวกับความถูกต้องในการเทียบมาตรฐาน โดย ชูชีพ พงศ์สมบูรณ์ (2528) ได้ศึกษาเปรียบเทียบรูปแบบการเทียบมาตรฐาน 3 วิธี คือ การเทียบเชิงเส้นตรง การเทียบแบบอิกวิเปอร์เซ็นต์ไคล์ และการเทียบโดยการใช้โด่งลักษณะข้อสอบ จากรูปแบบที่ใช้แบบทดสอบร่วมกับรูปแบบที่ใช้ผู้สอบร่วม เพื่อต้องการเปรียบเทียบประสิทธิภาพและความคงที่ของการเทียบมาตรฐานทั้ง 3 วิธี ใช้กลุ่มตัวอย่างเป็นนักเรียนชั้นมัธยมศึกษาปีที่ 6 ปีการศึกษา 2526 ทั่วประเทศสุ่มมาจำนวน 3,721 คน ที่เข้าสอบโครงการตรวจคุณภาพการศึกษาของกรมวิศการ กระทรวงศึกษาธิการ ในจำนวน 60 ข้อ แบ่งเป็น 2 ฉบับ ฉบับละ 38 ข้อ มีแบบสอบร่วมชนิดแบบสอบร่วมภาษาใน จำนวน 20 ข้อ ผลการเปรียบเทียบได้ว่าประสิทธิภาพของการเทียบมาตรฐานระหว่างรูปแบบทดสอบร่วมกับรูปแบบที่ใช้ผู้สอบร่วมในแต่ละวิธีมีประสิทธิภาพไม่แตกต่างกัน ความคงที่ของวิธีการเทียบมาตรฐานโดยรูปแบบที่ใช้ผู้สอบร่วมได้ผลว่า วิธีการเทียบโดยใช้อิกวิเปอร์เซ็นต์ไคล์มีความคงที่มากกว่าการเทียบเชิงเส้นตรง แต่มีความคงที่ผล ๆ กับการเทียบโดยการใช้โด่งลักษณะข้อสอบ ความคงที่ของวิธีการเทียบมาตรฐานในรูปแบบที่ใช้แบบสอบร่วมได้ผลเหมือนกับรูปแบบที่ใช้ผู้สอบร่วม

ส่วนภาวิณี ศรีสุวัตตานันท์ (2528) ได้เปรียบเทียบผลการใช้รูปแบบการเทียบมาตรฐาน 3 รูปแบบ คือ รูปแบบอิกวิเปอร์เซ็นต์ไคล์ รูปแบบเชิงเส้นตรง รูปแบบ IRT แบบสามพารามิเตอร์ โดยที่ใช้แบบสอบร่วมภาษาในต่างกัน 3 ขนาด คือ ขนาดร้อยละ 20 (7 ข้อ) ร้อยละ 40 (14 ข้อ) ร้อยละ 60 (21 ข้อ) โดยให้ข้อสอบทั้งหมด 35 ข้อ

การประเมินผลของการเทียบมาตรฐาน 2 ลักษณะ คือ ค่าความคลาดเคลื่อนมาตรฐานของการเทียบมาตรฐาน (SEE) ซึ่งเป็นการวิเคราะห์ในเชิงทฤษฎีของรูปแบบ และดัชนีการเปรียบเทียบความแตกต่าง (Index:C) ที่ได้จากวิเคราะห์กลุ่มตัวอย่างสอบทานผล ซึ่งเป็นการวิเคราะห์ความเพิงพอลของการเทียบมาตรฐานโดยตรง กับผลการเทียบสุดท้าย ส่วนประชากรและกลุ่มตัวอย่างแยกเป็น 2 กรณี คือ 1) กรณีสอบเลือก ประชากรคือ ผลการสอบรายบุคคลผู้มีคุณสมบัติขั้นต่ำตามกำหนดระเบียบการสอบคัดเลือกที่กำหนดไว้ 15,875 คน สุ่มตัวอย่างสำหรับการเทียบมาตรฐาน 2 กลุ่ม กลุ่มละ 1,500 คน สำหรับสอบทานผลกลุ่มละ 1,500 คน 2) กรณีแบบสอบผลสัมฤทธิ์ ประชากรคือ ผลการสอบรายบุคคลในวิชาที่กำหนดในระดับอุดมศึกษาจำนวน 12,383 คน สุ่มตัวอย่างเหมือนกรณีแรก แบบสอบมีกรณีละ 2 ชุด โดยที่ชุดละ 35 ข้อ มีข้อสอบรวม 21 ข้อ, 14 ข้อ, และ 7 ข้อ โดยยึดแบบสอบรวมมีค่าความเชื่อมั่นเท่าเทียมกันในแต่ละคู่ที่ใช้รูปแบบการเทียบมาตรฐาน 3 รูปแบบ สร้างตารางแปลงคะแนนสมมูลจากแบบสอบชุดที่ 1 ไปกับแบบสอบชุดที่ 2 แล้ววิเคราะห์ความคลาดเคลื่อนมาตรฐานของการเทียบมาตรฐาน ค่าดัชนีประสิทธิภาพสัมพันธ์ การวิเคราะห์กลุ่มสอบทานผล โดยคำนวณค่าดัชนีเปรียบเทียบความแตกต่าง (C) ผลการศึกษาพบว่า การเทียบระดับคะแนน 9 รูปแบบ ได้ผลเป็นที่ยอมรับว่า ความยาวของแบบสอบรวมมีผลต่อความแม่นยำ และความเพิงพอลของการเทียบมาตรฐาน แต่ในกรณีแบบสอบคัดเลือกมีผลไม่ชัดเจน ในด้านมิติเกี่ยวกับรูปแบบการเทียบมาตรฐานที่ได้ผลดีตามลำดับ คือ รูปแบบอิกวิเปอร์เทนไคส์ รูปแบบ IRT และเชิงเส้นตรง กรณีการสอบวัดผลสัมฤทธิ์ รูปแบบที่ได้ผลคือ รูปแบบเชิงเส้นตรง รูปแบบอิกวิเปอร์เทนไคส์ และรูปแบบ IRT

ต่อมา เรวดี อินทสระ (2530) ได้เปรียบเทียบความคลาดเคลื่อนของการเทียบมาตรฐานและความเพิงพอลจากการใช้รูปแบบอิงทฤษฎีการตอบสนองข้อสอบ กับรูปแบบการใช้เทคนิคการวิเคราะห์องค์ประกอบ จากกลุ่มตัวอย่างจำนวน 2,823 คน ใช้แบบทดสอบผลสัมฤทธิ์ทางคณิตศาสตร์ จำนวน 150 ข้อ วิเคราะห์โดยใช้โปรแกรม LOGIST 5 นำค่าพารามิเตอร์แยกเป็น 2 ฉบับ ฉบับละ 45 ข้อ ทั้งสองมีข้อสอบรวม 15 ข้อ ไปทดสอบกับกลุ่มตัวอย่าง แล้วนำมาเทียบระดับคะแนนและทานผลการเทียบมาตรฐานได้ผลว่า

เมื่อเทียบมาตรฐานจากฉบับที่ 2 ไปฉบับที่ 1 เปรียบเทียบความสมมูลของคะแนนทั้งสองได้ว่า ช่วงคะแนน 1 - 18 รูปแบบ IRT มีคะแนนสมมูลสูงกว่ารูปแบบเทคนิคการวิเคราะห์องค์ประกอบ ช่วงคะแนน 19 - 23 ทั้งสองรูปแบบมีคะแนนสมมูลใกล้เคียงกัน ช่วงคะแนน 24 - 45 รูปแบบเทคนิควิเคราะห์องค์ประกอบมีคะแนนสมมูลสูงกว่ารูปแบบ IRT เมื่อตรวจสอบผลได้ว่าการทานผลการเทียบระดับคะแนนแบบอิงทฤษฎีตอบสนองข้อสอบ มีดัชนีความเพียงพออยู่ในระดับความพอใจ ส่วนรูปแบบการให้เทคนิคการวิเคราะห์องค์ประกอบ มีสัดส่วนความคลาดเคลื่อนมาตรฐานทั้งสองฉบับสูงกว่ารูปแบบอิงทฤษฎีตอบสนองข้อสอบ

อาจารย์ กาญจนกิจโสภณ (2531) ได้ศึกษาการสร้างแบบทดสอบและตารางเทียบมาตรฐานตามแนวคิด ในวิชาคณิตศาสตร์ เรื่องสมการและลสมการ โดยวิธีที่เทียบคะแนนแบบราส์สำหรับนักเรียนชั้นมัธยมศึกษาตอนต้น กลุ่มตัวอย่างที่ใช้เทียบคะแนนเป็นนักเรียนชั้นมัธยมศึกษาปีที่ 1, 2 และ 3 จำนวนระดับชั้นละ 490, 482 และ 415 คน ตามลำดับ ผู้วิจัยได้นำแบบสอบระดับชั้นละ 90 ข้อ ไปทดสอบกับกลุ่มตัวอย่างอื่น แล้วคัดข้อสอบที่มีคุณภาพ เรียงตามความยากที่ต้องการระดับชั้นละ 40, 30 และ 30 ข้อ โดยแต่ละชั้นจะมีข้อสอบร่วมระหว่างชั้นอยู่ 10 ข้อ หลังจากนั้นนำไปทดสอบกับกลุ่มตัวอย่างที่ใช้เทียบคะแนน และสร้างตารางเทียบคะแนน เพื่อใช้ศึกษาพัฒนาการของนักเรียนกลุ่มเดิมที่ผ่านขึ้นไปเรียนในชั้นที่สูงขึ้นว่า นักเรียนมีความก้าวหน้าเป็นไปตามปกติหรือเบี่ยงเบนไปจากปกติ

จากการศึกษาเอกสารและงานวิจัยที่ผ่านมาพบว่า ไม่มีวิธีการเทียบมาตรฐานวิธีใดเพียงวิธีเดียวที่ให้ผลที่ดีที่สุดในทุกสถานการณ์ ดังนั้นในทางปฏิบัติจึงจำเป็นต้องศึกษาข้อตกลงต่าง ๆ ของวิธีการเทียบมาตรฐานแต่ละวิธี เช่น ลักษณะการแจกแจงของคะแนนที่ได้จากแบบสอบต่างชุดกัน ลักษณะของแบบสอบ ลักษณะของกลุ่มตัวอย่าง เป็นต้น เพื่อเป็นแนวทางในการตัดสินใจเลือกวิธีการในการเทียบมาตรฐานที่เหมาะสม ในแต่ละสถานการณ์ที่ต้องการเทียบ จึงพอสรุปได้ดังนี้

1. เมื่อแบบสอบมีความยากต่างกัน แต่กลุ่มตัวอย่างเป็นกลุ่มสมบูรณ์แล้ว วิธีการของรูปแบบอิคิวเปอร์เซนต์ไคล์ ซึ่งในความสัมพันธ์ของการแปลงที่ไม่เป็นเส้นตรงให้ความเหมาะสมต่อการแปลงได้ แต่ถ้าไม่เป็นกลุ่มสมบูรณ์ บางครั้งอาจเกิดจากการแจกแจงคะแนนที่ต่างมาก จนทำให้เกิดความคลาดเคลื่อนที่ทำให้ผลการเทียบหาความเพิงผลของการสร้างคะแนนสมบูรณ์
2. เมื่อแบบสอบมีความคล้ายคลึงในระดับความยาก และกลุ่มตัวอย่างเป็นกลุ่มที่สมบูรณ์ ควรใช้วิธีการของรูปแบบเชิงเส้นตรง ซึ่งสะดวกและให้ความแม่นยำมากกว่า
3. การเทียบมาตราโดยทฤษฎีตั้งเดิม ในกรณีแบบสอบมีความยากแตกต่างกันและกลุ่มตัวอย่างมีความสามารถต่างกันมากแล้ว การใช้วิธีอิคิวเปอร์เซนต์ไคล์ดีกว่าการเทียบเชิงเส้นตรง
4. เมื่อนำแบบสอบที่มีความยากแตกต่างกันมาเทียบ กลุ่มตัวอย่างไม่สมบูรณ์รูปแบบทฤษฎีตอบสนองข้อสอบสามพารามิเตอร์น่าจะให้ผลเป็นที่น่าพอใจกว่าวิธีอื่น
5. เมื่อกลุ่มตัวอย่างค่อนข้างน้อย และแบบสอบวัดในสิ่งเดียวกันแล้ว วิธีการเชิงเส้นตรง และวิธีใช้ทฤษฎีตอบสนองข้อสอบชนิดหนึ่งพารามิเตอร์ให้ผลเป็นที่น่าพอใจกว่าวิธีอิคิวเปอร์เซนต์ไคล์

จากข้อสรุปข้างต้น จะเห็นว่า ในกรณีที่นำแบบสอบที่มีความยากแตกต่างกันมาสอบกับกลุ่มตัวอย่างที่มีระดับความสามารถแตกต่างกัน แล้วนำผลที่ได้มาเทียบมาตรากัน การเทียบมาตราโดยวิธีรูปแบบอิงทฤษฎีตอบสนองข้อสอบ แบบ 3 พารามิเตอร์ น่าจะให้ผลเป็นที่น่าพอใจกว่าวิธีอื่น ดังนั้นผู้วิจัยจึงทำการศึกษาคุณภาพของการเทียบมาตราโดยวิธีรูปแบบอิงทฤษฎีการตอบสนองข้อสอบ แบบ 3 พารามิเตอร์

จุฬาลงกรณ์มหาวิทยาลัย

