



วรรณคดีที่เกี่ยวข้อง

การเทียบมาตรา (Equating) เป็นคำศัพท์เฉพาะศาสตร์ทางจิตมิติ คำว่าที่เกี่ยวข้องกับการทดสอบและวัดผลทางจิตวิทยา และการศึกษาเบื้องต้นมักไม่ได้อ้างถึง จึงเป็นเรื่องที่บุคคลทั่วไปไม่เข้าใจกัน ส่วนนักวัดผลเองโดยมากไม่ได้ให้ความสนใจเท่าที่ควร (Marco 1981) แต่อย่างไรก็ตามในช่วงเวลาประมาณ ๑๕ ปีมานี้ ได้มีนักวัดผลให้ความสำคัญกับการเทียบมาตรามากขึ้น ความสำคัญ มีการศึกษาวิจัยถึงเทคนิคการเทียบมาตรารูปแบบต่าง ๆ ในสภาพความเป็นจริงหลายลักษณะ นอกจากนี้ ฮอลแลนด์ และ โรบิน (Holland and Rubin 1982) ได้เป็นบรรณาธิการจัดทำคำรา Test Equating ขึ้นเผยแพร่ โดยมุ่งให้เกิดประโยชน์และสร้างเสริมเทคนิควิธีการใหม่ ๆ ให้กว้างขวางและลึกซึ้งยิ่งขึ้น การเสนอวรรณคดีที่เกี่ยวข้องในบทนี้ ได้เสนอรายละเอียด ๔ หัวข้อ ดังนี้ คือ (๑) พัฒนาการทางแนวคิดและรูปแบบการเทียบมาตรา (๒) การออกแบบรวบรวมข้อมูล (๓) วิธีเทียบมาตรารูปแบบอีควิวเปอร์เซนไทล์ (๔) วิธีเทียบมาตรารูปแบบเชิงเส้นตรง (๕) วิธีเทียบมาตรารูปแบบอิงทฤษฎีการตอบข้อสอบ (๖) ความยาวแบบสอบรวม (๗) ความคลาดเคลื่อนของการเทียบมาตรา (๘) การประเมินความเที่ยงพอของการเทียบมาตรา และ (๙) งานวิจัยที่เกี่ยวข้องกับการเปรียบเทียบรูปแบบการเทียบมาตรา

พัฒนาการทางแนวคิดและรูปแบบการเทียบมาตรา

๑. แนวคิดเชิงทฤษฎี

กัลลิกเสน (Gulliksen 1950: 298-304) ให้ความหมายของการเทียบมาตราว่า เป็นวิธีการหาคะแนนที่ได้จากแบบสอบสองชุดของวิชาเดียวกัน ให้เป็นคะแนนสมมูลที่เปรียบเทียบกันได้โดยตรง วิธีการเทียบมาตราที่ใช่กับการสอบทั้งสองชุดของแบบสอบกับกลุ่มคนเพียงกลุ่มเดียว เป็นวิธีง่าย ๆ โดยแปลงคะแนนแต่ละชุดให้เป็นคะแนนมาตรฐาน แล้วนำคะแนนแปลงมาเทียบกันโดยตรง แต่ก่อนอื่นต้องทำการตรวจสอบให้แน่ใจเสียก่อนว่า คะแนนทั้งสองชุดนั้นมีสัมประสิทธิ์ความ

โค้ง และความเข้มของการแจกแจงคล้ายคลึงกัน หรือถ้าคะแนนสองชุดนั้นโค้งทำให้เป็นมาตรฐานใน  
เทอมของเกณฑ์ใดเกณฑ์หนึ่ง ก็ต้องตรวจสอบดูว่า แบบสอบทั้งสองชุดมีค่าสัมประสิทธิ์สหสัมพันธ์กับ  
เกณฑ์เหล่านั้นโดยประมาณ

ฟลานานแกน (Flanagan 1951: 747-748) ได้กล่าวถึงการเทียบมาตรฐานว่า  
เป็นวิธีการที่ทำคะแนนจากแบบสอบต่างชุดให้มีคุณลักษณะที่เปรียบเทียบกันได้ คำว่า "ความสามารถ  
ในการเปรียบเทียบกันได้" มีความหมายเฉพาะที่ว่า เมื่อกำหนดประชากรให้ ถ้าการแจกแจงของ  
คะแนนจริงจากแบบสอบทั้งสองจากกลุ่มตัวอย่างที่เลือกมาขนาดใหญ่ใด ๆ มีลักษณะเหมือนกันแล้ว  
คะแนนดิบของแบบสอบทั้งสองชุดจึงจะสามารถเปรียบเทียบกันได้ หรือถ้าความเที่ยงของแบบสอบสอง  
ชุดนั้นเท่ากันในประชากรนั้นแล้ว ก็สามารถเปรียบเทียบการแจกแจงของค่าที่ได้เช่นกัน จากความ  
หมายดังกล่าวเป็นนิยามเชิงทฤษฎี ในทางปฏิบัติได้นิยามไว้ว่า ในประชากรที่กำหนด ถ้าคะแนนจาก  
แบบสอบสองชุดมีค่าเฉลี่ยเท่ากัน หรือเกือบเท่ากันในทุก ๆ กลุ่มตัวอย่างขนาดใหญ่ใด ๆ การเปรียบเทียบ  
เทียบจะทำได้ อย่างไรก็ตามเมื่อวิเคราะห์ถึงความแท้จริงที่จะให้เกิดลักษณะตามต้นนิยามเป็นเรื่องยาก  
มาก ปัจจัยสำคัญที่องค์การหนึ่ง คือ ความเป็นคู่ขนานของแบบสอบ ฟลานานแกนได้แนะนำว่า ควรเริ่ม  
ตั้งแต่การสร้างแบบสอบให้มีคุณสมบัติของความเป็นคู่ขนานกัน แล้วเลือกวิธีการเทียบมาตรฐานที่เหมาะสม  
ซึ่งเสนอไว้ ๔ วิธี คือ (๑) เทียบค่าเฉลี่ย โดยคำนวณค่าเฉลี่ยของการแจกแจงคะแนนทั้งสองชุด  
ถ้าความแตกต่างของค่าเฉลี่ยอยู่ภายในขอบเขตความแปรผันเชิงสุ่มแล้ว ถือว่าคะแนนสองชุดนั้นเปรียบเทียบ  
เทียบกันได้ แต่ถ้าความแตกต่างมีนัยสำคัญ ให้ใช้วิธีบวกเข้าหรือลบออกเท่าจำนวนที่แตกต่างจาก  
คะแนนชุดหนึ่ง เพื่อให้เกิดคะแนนสมมูลกับอีกชุดหนึ่ง (๒) ใช้เทคนิคของสมการถดถอย ใช้ค่าประมาณ  
ที่ดีที่สุด ("best estimate") ของคะแนนจากแบบสอบชุดหนึ่งซึ่งรู้ค่าของอีกชุดหนึ่ง (๓) ใช้  
คะแนนมาตรฐาน เป็นวิธีปรับคะแนนอย่างคงที่ตลอดการแจกแจง และ (๔) ใช้วิธีคิวเปอร์เซนไทล์  
เป็นวิธีหาค่าคะแนนโดยการเปรียบเทียบความสัมพันธ์ของการแจกแจงคะแนน ซึ่งเป็นวิธีที่เหมาะสมที่สุด

จากมโนทัศน์ดังกล่าว มีผลทำให้การเทียบมาตรฐานในสมัยทศวรรษ ๑๙๕๐ มุ่งเน้น  
เทคนิคการสร้างแบบสอบชุดต่าง ๆ ให้มีคุณสมบัติของความเป็นคู่ขนานมากที่สุด เมื่อแบบสอบเป็น  
คู่ขนานแล้ว คะแนนที่ได้จากทั้งสองชุด ย่อมเป็นคะแนนสมมูลกัน (equivalent scores) หรือ  
ถ้าหากมีความแตกต่าง เกิดขึ้น เช่น ค่าเฉลี่ยหรือส่วนเบี่ยงเบนมาตรฐานต่างกันภายในขอบเขตที่ยอมรับ

จะใช้การปรับค่าความแตกต่างตามที่กลัดกลืนได้แนะนำดังกล่าวมาแล้ว ก็สามารถทำให้คะแนนทั้งสองเป็นคะแนนสมมูลได้ เทคนิคการสร้างแบบสอบคู่ขนานที่กลัดกลืนแนะนำมีดังนี้ (Gulliksen 1950: 207-210)

- (๑) วิเคราะห์ความยากและค่าอำนาจจำแนกรายข้อ แล้วทำจุดลงบนกระดาษกราฟ
- (๒) พิจารณาข้อสอบที่เกาะกลุ่มกันอยู่ว่า แต่ละกลุ่มวัดความสามารถอะไรบ้าง
- (๓) แยกข้อสอบที่เกาะกลุ่ม และวัดสิ่งเดียวกันให้อยู่คนละชุดโดยสุ่ม

วิธีนี้ช่วยให้ผู้พัฒนาแบบสอบสามารถควบคุมความหลากหลายของคำถามจากแบบสอบชุดหนึ่งไปสู่อีกชุดหนึ่งได้ แต่วิธีนี้ ฟลานาแกน (Flanagan 1951: 749-50) พบว่า มีจุดอ่อนอยู่ ๕ ประการ คือ

- (๑) ความยากที่วิเคราะห์จากฉบับเริ่มแรกที่มีจำนวนข้อสอบมาก ๆ นั้นมีค่าแตกต่างจากความยากของข้อสอบข้อเดียวกัน เมื่อกระจายเข้าไปอยู่ในฉบับสมมูล ทั้งนี้เนื่องจากตำแหน่งของการจัดเรียงข้อสอบ และบริบทของแบบสอบมีผลกระทบบ้าง

- (๒) โดยปกติการแบ่งข้อสอบจากฉบับเริ่มแรกไปยังฉบับสมมูล โดยไม่ทำให้สูญเสียลักษณะการสุ่มเนื้อหาในแต่ละชุดนั้นเป็นเรื่องที่ทำได้

- (๓) การปรับปรุงแบบสอบชุดหลัง ๆ ทำได้ยาก เพราะจะต้องคงสภาพความยากคุณสมบัติทั่วไป และจำนวนข้อในแต่ละตอนให้เหมือนฉบับเดิม

- (๔) การดำเนินการสอบที่จัดขึ้นต่างเวลากัน มีความหมายถึงการเปลี่ยนแปลงสถานการณ์ของการทดสอบ ซึ่งมีผลกระทบบ่อการแปรเปลี่ยนของการแจกแจงคะแนนดิบ จะนำมาเปรียบเทียบโดยตรงไม่ได้

- (๕) ความสัมพันธ์ของข้อสอบข้อเดียวกันในแบบสอบค่างชุด ย่อมแตกต่างกัน ซึ่งทำให้ค่าความคงที่ภายในเปลี่ยนแปลงไป ผลที่ตามมาจากการทดสอบจะได้การแจกแจงที่ไม่เหมือนกันด้วย

ดังนั้น การเทียบมาตราของแบบสอบด้วยการจัดกระทำในขั้นตอนการสร้างแบบสอบคู่ขนาน จึงมีข้อจำกัดมาก นอกจากนี้การทดสอบในสภาพการณ์จริงบางกรณี โดยเฉพาะที่เป็นโปรแกรมการสอบระดับชาติ จำเป็นต้องยึดนโยบาย และข้อกำหนดของระเบียบเป็นสำคัญในการปฏิบัติ เทคนิคการเทียบมาตราจึงหันมาเน้นกระบวนการปรับโดยอาศัยวิธีการทางสถิติ ด้วยเหตุนี้การศึกษาค้นคว้า



หาวิธีการเทียบมาตรา จึงมุ่งหาวิธีการทางสถิติที่เหมาะสมเพื่อใช้กับแบบสอบที่มีคุณสมบัติน้อยกว่าที่กำหนดค่านิยามของแบบสอบคู่ขนาน กล่าวโดยสรุปได้ว่า การเทียบมาตราเป็นกระบวนการที่รวมเอาเทคนิคการสร้างแบบสอบอย่างพิถีพิถัน ตลอดจนวิธีการใช้สถิติที่เหมาะสมในการปรับคะแนนที่ได้จากต่างชุดของแบบสอบ เพื่อชดเชยความแตกต่างของความยากที่ต่างกัน (Holland and Rubin 1982: 2) ทั้งนี้ การศึกษาถึงเทคนิคการเทียบมาตราจะเป็นวิธีการทางสถิติมาก โดยถือว่าแบบสอบสองชุดที่กองการเทียบนั้นได้รับการพัฒนาถึงระดับที่วัดความสามารถในสิ่งเดียวกัน และอยู่ภายใต้ข้อจำกัดของสภาพความเป็นจริงที่หลีกเลี่ยงไม่ได้

## ๒. รูปแบบการเทียบมาตรา (Equating models)

### ๒.๑ รูปแบบการเทียบมาตราตามประเพณีนิยม (Traditional models)

แองกอฟ (Angoff 1984: 85-93) แบ่งรูปแบบการเทียบมาตราออกเป็น ๒ รูปแบบ คือ Equating model และ Calibration model ซึ่งมีความแตกต่างในเชิงแนวคิด ดังนี้

Equating model เป็นรูปแบบการแปลงคะแนนที่ยึดนิยามของคะแนนสมมูลเป็นเป้าหมาย (equivalent scores) นิยามของคะแนนสมมูลกล่าวไว้ว่า คะแนนสองจำนวน จำนวนหนึ่งมาจากแบบสอบชุด X ส่วนอีกจำนวนหนึ่งมาจากแบบสอบชุด Y โดยที่ทั้ง X และ Y วัดทั้งชั้นเชิงจิตวิทยาเดียวกันด้วยแบบสอบที่มีความเที่ยงระดับเดียวกัน คะแนนสองจำนวนนั้นนับว่าเป็นคะแนนสมมูล ถ้าค่าก็สมนัยกับค่าแห่งเปอร์เซ็นต์ไคร์ของประชากรที่กำหนดให้ที่ค่าแห่งเดียวกันตามนิยามนี้ ถ้าแบบสอบ ๒ ชุด มีความยากต่างกัน ซึ่งหมายถึง การแจกแจงของคะแนนดิบของแบบสอบสองชุดนั้นย่อมต่างกันด้วยนั้น คะแนนสมมูลที่เกิดจากการใช้วิธีการเทียบมาตรานี้ย่อมเกิดจากการยึดหรือรวมมาตราของคะแนนชุดหนึ่ง เพื่อให้การแจกแจงคะแนนเหมือนกับแจกแจงของอีกชุดหนึ่งนั่นเอง ทั้งนี้ การเทียบมาตราด้วยรูปแบบ Equating นี้มีผลทำให้ผู้สอบแต่ละคนได้คะแนนแปลงเหมือนกันไม่ว่าผู้สอบจะรับการทดสอบด้วยแบบสอบชุดใด ในกรณีที่มีแบบสอบ ๒ ชุด มีลักษณะสำคัญ ๆ ใกล้เคียงกันมากตามทฤษฎีซึ่งเป็นเหตุผลที่จะกำหนดข้อสมมุติฐานที่ว่า การแจกแจงคะแนนดิบของแบบสอบสองชุดมีรูปร่างเหมือนกัน และการแปลงคะแนนจากชุดหนึ่งไปยังอีกชุดหนึ่งสามารถทำได้

โดยเปลี่ยนจุดเริ่มต้นของมาตราและหน่วยการวัดเท่านั้น หรือกล่าวอีกนัยหนึ่งว่า ทำการปรับโมเมนต์ ที่หนึ่งและสองนั่นเอง วิธีการเช่นนี้ คือ การแปลงเชิงเส้นตรง จึงอาจกล่าวสรุปได้ว่า การเทียบ มาตราด้วยการแปลงคะแนนเชิงเส้นตรงให้ค่าใกล้เคียงอย่างมากกับวิธีการอิกวิเปอร์เซนไคล์ เมื่อรูปร่างการแจกแจงคะแนนดิบของแบบสอบทั้งสองคล้ายคลึงกัน ทั้งนี้ ถ้าผู้วิเคราะห์สามารถ กำหนดข้อบกพร่องของความคล้ายคลึงของการแจกแจงได้ ก็สามารถใช่วิธีการเชิงเส้นตรง ซึ่งให้ข้อมูลที่ใจ กว่าอิกวิเปอร์เซนไคล์ตรงที่สามารถหลีกเลี่ยงความลำเอียงที่อาจเกิดขึ้นในขณะทำการปรับเส้นโค้ง แต่อย่างไรก็ตามการแจกแจงของคะแนนจากแบบสอบสองชุด โดยธรรมชาติมีลักษณะการแจกแจง แตกต่างกัน การเทียบมาตราด้วยวิธีอิกวิเปอร์เซนไคล์จึงดู เหมือนให้ความถูกต้องและเหมาะสมกว่า (Angoff 1984: 87)

Calibration model เป็นรูปแบบการเทียบมาตราที่คงการให้ผลของ คะแนนแปลงสะท้อนลักษณะความยากหรือง่ายกว่าของแบบสอบชุดใหม่ ถ้าพิจารณาย้อนกลับไปที่รูปแบบ Equating วิธีการของอิกวิเปอร์เซนไคล์จะกำหนดให้คะแนนชุดใหม่อยู่ในกรอบของชุดเก่า คะแนน สูงสุดของชุดใหม่ซึ่งวัดระดับความสามารถสูงกว่าชุดเก่าจะถูกกำหนดให้เสมอกัน สภาพเช่นนี้สาร- สนเทศจากการแปลงคะแนนจะถูกจำกัด แองกอฟ (Angoff 1984: 89) ได้เปรียบเทียบการ เทียบมาตราตามรูปแบบ Calibration ว่า เหมือนหนึ่งเทอร์โมมิเตอร์สองชนิดที่สร้างขึ้นเพื่อ วัดอุณหภูมิที่มีขีดจำกัดต่างกัน เช่น วัดอุณหภูมิระหว่าง ๔๐ องศาฟาเรนไฮต์ ถึง ๑๐๐ องศาฟา- เรนไฮต์ กับอีกอันหนึ่งวัดอุณหภูมิระหว่าง ๔๔ องศาฟาเรนไฮต์ กับ ๑๐๔ องศาฟาเรนไฮต์ ทั้งสองอันนี้ต้องการเทียบมาตราเพื่ออ่านอุณหภูมิในหน่วยเดียวกัน รูปแบบการ Calibration จะช่วยให้เทอร์โมมิเตอร์ทั้งสองอันให้ค่าจากการวัดอุณหภูมิที่สะท้อนให้เห็นความแตกต่างอุณหภูมิที่วัดได้ อย่างถูกต้อง ครอนแบค (Cronbach 1970: 111) ได้อธิบายวิธีการ Calibration ว่า เหมือนกับวิธีการทำเครื่องหมายหน่วยวัดบนเครื่องวัดความดันอากาศที่ใช้สารอนีรอยให้มีความ ถูกต้องเช่นเดียวกับเครื่องวัดที่ใช้ปรอท วิธีที่ใช้เทียบมาตราในรูปแบบ Calibration ที่ เหมาะสม คือ การเทียบเชิงเส้นตรง

แองกอฟ (๑๙๘๔: ๕๓) ได้สรุปไว้ว่า การปรับมาตราที่ได้คำนึงถึงความแตกต่างของความยาก และพิสัยความสามารถที่วัดมีอยู่ ๒ รูปแบบ คือ Calibration เป็นการปรับเชิงเส้นตรง หรือใช้คะแนนมาตรฐาน (Z score) ซึ่งเป็นวิธีที่คงการให้คะแนนแปลงสะท้อนให้เห็นความสามารถในการจำแนกค่านที่แบบสอบต้องการ อีกรูปแบบหนึ่ง คือ Equating เป็นการปรับที่ไม่เป็นเส้นตรง (curvilinear) หรืออิกวิเปอร์เซนไคด์ ซึ่งบางครั้งวิธีเชิงเส้นตรงใช้ประมาณค่าได้ก็ และโดยความจริงที่ว่าวิธีเชิงเส้นตรงสมมูลกับอิกวิเปอร์เซนไคด์ เมื่อรูปร่างการแจกแจงของคะแนนดิบเหมือนกัน ซึ่งหมายความว่า นอกจากสองโมเมนต์แรกแล้ว โมเมนต์มาตรฐานที่เหลือของการแจกแจงคะแนนดิบของแบบสอบสองชุดของกลุ่มผู้สอบที่กำหนดให้เหมือนกันหมด ความแตกต่างของการเทียบมาตราสองรูปแบบนี้ จะเห็นได้ชัดเมื่อใช้เทียบคะแนนจากแบบสอบที่มีความยากไม่เท่ากัน ๒ ชุด ไปสู่มาตราใหม่ที่เป็นคนละมาตรากับมาตราคะแนนดิบเดิม จะเห็นความแตกต่างที่สำคัญชัดเจน วิธีการเชิงเส้นตรงจะแสดงลักษณะให้ปรากฏในคะแนนที่มาจากวิธีการแปลง และถ้าแบบสอบมีความยากต่างกัน อาจเกิดเหตุการณ์ที่ว่า คะแนนของชุดหนึ่งอาจสูงเกินกว่าที่คนสอบอีกชุดหนึ่งจะทำได้ถึง ส่วนในกรณีของอิกวิเปอร์เซนไคด์จะปรับความยากเหล่านี้ เพื่อให้คะแนนของแต่ละคนมีคงเดิมโดยไม่สนใจว่า ผู้สอบนั้นจะได้รับการทดสอบด้วยแบบสอบชุดใด

จากการเสนอรูปแบบการเทียบมาตราที่กล่าวมา พบว่า ต่างมีเป้าหมายร่วมกัน คือ ต้องการแปลงหน่วยในระบบมาตราของคะแนนชุดหนึ่ง ไปสู่หน่วยของระบบมาตราคะแนนอีกชุดหนึ่ง โดยให้ตัวเลขภายหลังกระบวนการแปลงสามารถนำมาเปรียบเทียบในเชิงปริมาณได้ สำหรับบทความ และผลงานวิจัยที่ปรากฏในเวลาต่อมา ไม่พบว่ามียุคเน้นความแตกต่างของรูปแบบทั้งสองอีกเลย และการเทียบมาตราที่ปรากฏในระยะต่อมา นี้ ได้ใช้คำว่า Equating เป็นคำทั่วไป ส่วนผลการเทียบจะสะท้อนให้เห็นความยากที่ค้างกันหรือไม่นั้นได้เป็นส่วนหนึ่งที่นำมาพิจารณาถึงสภาพของแบบสอบ และสถานการณ์ของการทดสอบในการวิเคราะห์ เพื่อสร้างคะแนนสมมูลภายใต้การออกแบบในการทดลอง ทั้งนี้ การเทียบมาตราตามประเพณีนิยมมาแต่เดิมอาจสรุปได้ว่ามีสองรูปแบบ คือ รูปแบบอิกวิเปอร์เซนไคด์ และรูปแบบเชิงเส้นตรง (Kolen and Whitney 1982)



## ๒.๒ รูปแบบอิงทฤษฎีการตอบข้อสอบ (Item response theory: IRT models)

เป็นรูปแบบการเทียบมาตราที่พัฒนาต่อมาจากทฤษฎีการตอบที่วัดด้วยการวัดความสามารถจริงที่แฝงอยู่ในแต่ละบุคคล (Latent trait theory) ทฤษฎีว่าด้วยการตอบข้อสอบ เป็นทฤษฎีที่จำลองคำตอบ (response) ของผู้สอบจากการตอบคำถามข้อสอบรายข้อในแบบทดสอบ ทฤษฎีนี้มีประโยชน์ต่อการออกแบบของแบบสอบ ต่อการบรรยายและประเมินข้อสอบ และแบบสอบ ต่อการให้คะแนนคำตอบของผู้สอบรายบุคคลได้อย่างเหมาะสม ต่อการทำนายคะแนนของผู้สอบรายคนและกลุ่มผู้สอบ และต่อการจัดกระทำและแปลความหมายของคะแนนสอบ ซึ่งรวมถึงการเทียบมาตราของแบบสอบต่างชุดด้วย (Lord 1982 (b): 141) ทฤษฎีนี้ได้พัฒนาขึ้นและนำสู่การประยุกต์กับการวัดผลมาไม่น้อยกว่า ๓๐ ปี ในทศวรรษที่ผ่านมา มีงานวิจัยมุ่งเน้นการประยุกต์ทฤษฎีนี้ อย่างจริงจัง จนอาจถือว่าเป็น "ยุคของการปฏิวัติของสาขาการวัดผลทางการศึกษา" (Marco 1977) รูปแบบการเทียบมาตราที่พัฒนาขึ้นภายใต้ทฤษฎีนี้ที่รู้จักกันดี ได้แก่ รูปแบบหนึ่งพารามิเตอร์ (The one parameter model หรือ Rasch model) และรูปแบบโลจิสสามพารามิเตอร์ (The three-parameter logistic model) ซึ่งลอร์ดเป็นผู้นำในการพัฒนาทฤษฎีนี้ และนำไปสู่การประยุกต์กับปัญหาการทดสอบมากมาย รวมถึงการเทียบมาตราคะแนนมาเป็นเวลาเกือบ ๒๐ ปีแล้ว (Marco 1977) สำหรับงานวิจัยนี้ได้เลือกรูปแบบโลจิสสามพารามิเตอร์เพียงอย่างเดียว การเสนอวรรณคดีที่เกี่ยวข้อง จึงจำกัดเฉพาะเรื่องราวของโลจิสสามพารามิเตอร์เท่านั้น

นियามการเทียบมาตรารูปแบบอิงทฤษฎีการตอบข้อสอบ (เรียกย่อ ๆ ว่า IRT) เป็นนियามที่มุ่งเน้นความสำคัญของความเสมอภาค หรือความเป็นธรรมที่บุคคลจะได้รับจากการทดสอบอย่างเท่าเทียมกัน ถึงแม้การสอบนั้นจะใช้แบบสอบต่างชุดกัน (Lord 1980: 195) ความจริงมีโน้ตค้นของความเสมอภาคนี้ แองกอฟ (Angoff 1984: 85) ได้ให้ไว้ในทำนองเดียวกันว่า ผู้ที่รับการทดสอบต้องไม่ได้รับประโยชน์ หรือสูญเสียประโยชน์เป็นพิเศษอันเนื่องมาจากการได้รับการทดสอบด้วยแบบสอบชุดใดชุดหนึ่งที่มีความยากกว่า หรือง่ายกว่ากัน นಿಯามเชิงทฤษฎีในทำนองนี้ เทคนิคการเทียบมาตราที่แล้วมามีข้อจำกัดในทางปฏิบัติอยู่มาก ลอร์ดจึงให้นิยามในเชิงปฏิบัติว่า ความเสมอภาค หมายถึง การที่ผู้สอบทุกคนที่ทุกระดับความสามารถ  $\theta$  มีการแจกแจงความถี่ตาม

เงื่อนไขความสามารถของคะแนน  $x$  ( $f_{x|\theta}$ ) เหมือนกับการแจกแจงความถี่ตามเงื่อนไขความสามารถของคะแนนแปลง  $x(y)$  ซึ่งเขียนเป็นสัญลักษณ์ได้ว่า

$$f_{x|\theta} \equiv f_{x(y)|\theta}$$

เมื่อ  $x(y)$  เป็นฟังก์ชันหนึ่งต่อหนึ่งของ  $y$  และกำหนดให้ความแปรปรวนของการแจกแจงตามเงื่อนไขเท่ากันด้วย (Lord 1980: 195-6) เมื่อวิเคราะห์จากนิยาม พบว่า การเทียบมาตรฐานดังกล่าวจะทำได้เมื่อข้อมูลคะแนนมีความสมบูรณ์ในการระบุความสามารถที่แท้จริงของบุคคลได้ ทั้งนี้เมื่อมีการพัฒนาทฤษฎีการวัดความสามารถแฝง และจำเพาะลงมาถึงการทดสอบข้อสอบ (Item response theory) ซึ่งสามารถหาค่าประมาณความสามารถจริงของบุคคลได้จากการทดสอบข้อสอบ จึงเป็นขั้นการพัฒนาทฤษฎีของการเทียบมาตรฐานอีกระดับหนึ่ง การเทียบมาตรฐานในระดับนี้ เรียกว่า การเทียบมาตรฐานคะแนนจริง (True-score equating) โดยทฤษฎีการเทียบคะแนนจริงมีข้อกำหนด ๓ ประการ (Lord 1980: 199) คือ

(๑) ความเสมอภาค (equity) หมายความว่า ถ้าหากพิจารณาที่ระดับความสามารถ ( $\theta$ ) ใด ๆ การแจกแจงความถี่อย่างมีเงื่อนไขของคะแนนที่แปลงด้วยเทคนิคการเทียบมาตรฐาน IRT จะต้องเหมือนกับการแจกแจงความถี่อย่างมีเงื่อนไขของคะแนนจากแบบสอบที่ทำการเทียบ เขียนในรูปของสมการได้ดังนี้

$$f_{x|\theta} \equiv f_{x(y)|\theta}$$

เมื่อ  $x(y)$  เป็นฟังก์ชันหนึ่งต่อหนึ่งของ  $y$

(๒) ความไม่แปรเปลี่ยนเมื่อเปลี่ยนกลุ่ม (Invariance across groups) คะแนนแปลง  $x(y)$  จะเหมือนกันโดยไม่ขึ้นกับค่าแปรค่าอื่น ๆ ของประชากรที่นำมาสร้างสมการสำหรับการเทียบมาตรฐาน

(๓) ความสมมาตร (Symmetry) หมายถึง การเทียบมาตรฐานคะแนนจะต้องเหมือนกัน ไม่ว่าการเทียบนั้นจะเทียบจากแบบสอบ  $x$  ไปหาแบบสอบ  $y$  หรือแบบสอบ  $y$  ไปหาแบบสอบ  $x$



ข้อกำหนดที่สำคัญนี้ โดยทั่วไปจะหาไม่พบในการแปลงคะแนนด้วยคะแนนสอบ  
ที่มีความคลาดเคลื่อนรวมอยู่ (fallible scores) นอกจากจะต้องทำการเทียบมาตรฐาน  
ด้วยคะแนนจริง (true scores) เท่านั้น

### ๓. สมมติฐานการเทียบมาตรฐาน (Marco 1981)

- นิยามที่ ๑. คะแนนจากแบบสอบ X และแบบสอบ Y นำมาเทียบมาตรฐานกันได้  
ถ้า  $M_{y'} = M_x$  และ  
 $SD_{y'} = SD_x$  สำหรับประชากร P  
เมื่อ M และ SD คือ ค่าเฉลี่ยและส่วนเบี่ยงเบนมาตรฐาน  
 $y'$  คือ คะแนนแปลงที่ได้จากฟังก์ชันการเทียบมาตรฐาน  
 $e_x(Y)$

ในกรณีเช่นนี้ รูปแบบการเทียบมาตรฐานเชิงเส้นตรงสามารถแปลงคะแนน Y ไปสู่  
มาตรฐานของ X อย่างเพียงพอ ด้วยสมการ

$$y' = (SD_x/SD_y) y + M_x - (SD_x/SD_y) M_y$$

โดยทั่วไป  $M_y$  ไม่เท่ากับ  $M_x$  และ  $SD_y$  ไม่เท่ากับ  $SD_x$   
เพราะว่า มีความหมายแตกต่างกันในความยากของคำถามในแบบสอบต่างชุดกัน การแปลงเชิงเส้น  
ตรงจะปรับ  $M_y$  ให้เท่ากับ  $M_x$  และ  $SD_y$  เท่ากับ  $SD_x$  และเป็นการอ้างถึงการจัดการให้  
คะแนนเฉลี่ยและส่วนเบี่ยงเบนมาตรฐานให้เท่ากัน นิยามนี้ไม่มีข้อกำหนดว่า แต่ละรายบุคคลจะต้อง  
มี y เท่ากับ x เพียงแต่กำหนดให้คะแนนเฉลี่ยและส่วนเบี่ยงเบนมาตรฐานจะต้องเท่ากันในประชากร  
ที่ศึกษา และโดยนิยามนี้ไม่มีข้อกำหนดว่า แบบสอบ X และ Y จะต้องวัดในสิ่งเดียวกัน เช่น ความรู้  
หรือความสามารถ หรือทักษะต่าง ๆ

นิยามที่ ๒    คะแนนจากแบบสอบ  $X$  และแบบสอบ  $Y$  นำมาเทียบมาตรากันได้ ถ้า  $y$  และ  $x$  มีตำแหน่งเปอร์เซนต์เท่ากันในการแจกแจงของคะแนน  $X$  และ  $Y$  ในประชากร  $P$

นิยามนี้ตรงกับนิยามรูปแบบอิกวิเปอร์เซนต์ จะเน้นใช้วิธีคำนวณตำแหน่งเปอร์เซนต์ของคะแนนการแจกแจงใน  $X$  และ  $Y$  แล้วหาค่าตำแหน่งคะแนนที่สมมูลกัน ในทางปฏิบัติการแจกแจงของ  $X$  และ  $Y$  จะถือเสมือนหนึ่งเป็นการแจกแจงต่อเนื่องมากกว่าเป็นคะแนนชวาคอน และตำแหน่งเปอร์เซนต์ใช้การคำนวณด้วย linear interpolation การแจกแจงน่าจะปรับเสียก่อนการคำนวณ นิยามนี้ไม่ได้มีข้อกำหนดว่า  $X$  และ  $Y$  ต้องวัดในสิ่งเดียวกัน

นิยามที่ ๓    คะแนนจากแบบสอบ  $X$  และแบบสอบ  $Y$  จะนำมาเทียบมาตรากันได้ เมื่อ  $M'_y = M'_x$  สำหรับประชากรกลุ่มย่อยที่มีความสามารถ ความรู้ และทักษะที่ทำการวัดอยู่ในระดับเดียวกัน นั่นคือ

$$\frac{M'_y|a}{a} = \frac{M'_x|a}{a}$$

$a$  คือ ระดับความสามารถที่กล่าวถึง

ฟังก์ชันการเทียบมาตราที่เป็นไปตามนิยาม ต้องเป็นฟังก์ชันที่ไม่ใช่เส้นตรง (curvilinear) การใช้รูปแบบเชิงเส้นตรงอาจให้การประมาณที่เป็นประโยชน์ได้ ส่วนรูปแบบอิกวิเปอร์เซนต์อาจให้ผลลัพธ์ที่คงเส้นคงวาคำนวณนิยาม หรืออาจไม่เป็นเช่นนั้นก็ได้ แต่ที่ชัดเจน คือ ต้องใช้วิธีการที่สลับซับซ้อนยิ่งขึ้น จึงเพียงพอที่จะให้เกิดการแปลงที่ให้คะแนนสมมูลได้ เช่น วิธีที่อาศัยรูปแบบอิงทฤษฎีการตอบข้อสอบ (Lord 1980, chapter 13) ซึ่งจะให้ค่าประมาณของระดับความสามารถที่อยู่ภายในกลุ่ม นิยามที่ ๓ นี้ ได้กำหนดให้  $X$  และ  $Y$  วัดในสิ่งเดียวกัน วิธีการทางสถิติที่จะให้ผลลัพธ์ที่น่าพอใจตามนิยามก็ต่อเมื่อได้เป็นไปตามข้อตกลงเบื้องต้นของนิยาม ไม่มีข้อกำหนดใด ๆ ที่เกี่ยวกับความคมชัดของแบบสอบ ฉะนั้นจึงสามารถใช้เทียบมาตรากับแบบสอบที่มีความเที่ยงต่างกันได้

นิยามที่ ๔ คะแนนจากแบบสอบ  $X$  และแบบสอบ  $Y$  นำมาเทียบมารากันได้ เมื่อ  $M_{y|a} = M_x$  ในทุก ๆ กลุ่มย่อยที่มีความสามารถระดับเดียวกันในประชากร  $P$  นั่นคือ

$$M_{y|a} = M_x|a$$

และคะแนนความคลาดเคลื่อนมาตรฐานของการวัด  $y'$  และ  $x$  เท่ากันในประชากร  $P$

นิยามนี้เพิ่มข้อกำหนดจากนิยามที่ ๓ อีกหนึ่งข้อ คือ แบบสอบ  $X$  และ  $Y$  ต้องมีความแม่นยำในการวัดเท่ากัน ซึ่งหมายถึงมีค่าความเที่ยงเท่ากัน ในกรณีนี้แบบสอบ  $X$  และ  $Y$  ควรประกอบด้วยข้อสอบจำนวนที่เท่ากัน

นิยามที่ ๕ คะแนนจากแบบสอบ  $X$  และแบบสอบ  $Y$  นำมาเทียบมารากันได้ เมื่อ  $M_{y|a} = M_x$  และ  $SD_{y|a} = SD_x$  สำหรับทุก ๆ กลุ่มย่อยที่มีความสามารถเดียวกันในประชากร  $P$  นั่นคือ

$$M_{y|a} = M_x|a \quad \text{และ}$$

$$SD_{y|a} = SD_x|a \quad \text{สำหรับทุก ๆ } a$$

นิยามนี้ยากแก่การจัดการให้มีข้อกำหนดที่ระบุไว้ เพราะไม่ใช่เฉพาะแต่ความคลาดเคลื่อนมาตรฐานโดยส่วนรวมจะต้องเท่ากัน แต่หมายถึงความคลาดเคลื่อนมาตรฐานของทุก ๆ กลุ่มย่อยต้องเท่ากันหมดด้วย โดยทางปฏิบัตินิยามนี้ไม่มีทางทำได้เลย นอกจากต้องอาศัยการดำเนินการสอบด้วยคอมพิวเตอร์ และแต่ละคนจะต้องได้รับจำนวนคำถามที่มากพอที่จะให้การประมาณที่แม่นยำเพื่อให้เป็นไปตามนิยาม

นิยามที่ ๖ คะแนนจากแบบสอบ  $X$  และแบบสอบ  $Y$  จะเทียบมารากันได้ ถ้าการแจกแจงความถี่ของ  $y'$  และ  $x$  เหมือนกันในแต่ละกลุ่มตัวอย่างย่อยที่มีความสามารถอยู่ในระดับเดียวกัน





นิยามที่ ๒ นี้ มีข้อกำหนดมากขึ้นกว่านิยามที่ ๕ ในเรื่องของรูปร่างของการแจกแจงของคะแนน  $y$  อย่างมีเงื่อนไขที่ระดับความสามารถใด ๆ ทั้งนี้ยอมหมายถึงว่า  $\mu_{y|a}$  และ  $SD_{y|a}$  จะต้องเข้าคู่กับรูปร่างของการแจกแจงอย่างมีเงื่อนไขของ  $x$  ตลอด (Lord 1980) โคนสกร์ได้เห็นว่า นิยามนี้เป็นไปได้ก็ต่อเมื่อกำหนดในแบบสอบ  $y$  มีความเท่าเทียมหรือสมบูรณ์ในเชิงหน้าที่กับคำถามในแบบสอบ  $x$  เท่ากัน หรือว่าแบบสอบทั้งสองฉบับสามารถให้ผลการวัดเป็นคะแนนสมบูรณ์ (perfectly reliable scores) ในกรณีเช่นนี้  $y = x$  นิยามนี้ดูเหมือนจะไม่มีประโยชน์ แต่เป็นนิยามในเชิงทฤษฎีที่ต้องการให้ความหมายของความไม่แตกต่างกันของการสอบ  $x$  หรือ  $y$  ของบุคคลใด ๆ อย่างแท้จริง

การออกแบบการรวบรวมข้อมูลคะแนน

การเทียบมาตราเป็นวิธีการเชิงประจักษ์ เพื่อกำหนดการแปลงคะแนนที่ได้จากแบบสอบชุดหนึ่ง ไปสู่แบบสอบอีกชุดหนึ่งในความหมายที่เท่าเทียมกัน ลักษณะของกระบวนการเชิงประจักษ์นี้ จึงเกี่ยวข้องกับการออกแบบเพื่อรวบรวมข้อมูล และจัดกระทำในทางสถิติเพื่อระบุการแปลงคะแนนในที่สุด (Marco 1981) แองกอฟ (Angoff 1984: 94-123) ได้บรรยายทั้งสองประเด็นไว้อย่างครอบคลุมในแบบการเทียบมาตรา ๒ แบบ (designs) มาร์โคได้นำมาบรรยายซ้ำ แต่ได้แยกประเด็นของการออกแบบเพื่อรวบรวมข้อมูลกับการจัดกระทำทางสถิติออกจากกัน ทั้งนี้เพราะวิธีการทางสถิติบางวิธีใช้ได้กับแบบแผนการรวบรวมข้อมูลมากกว่าหนึ่งแบบ ซึ่งสรุปสาระสำคัญดังนี้ คือ

แบบแผนที่ ๑ กลุ่มเดียว ทำแบบสอบทั้งสองชุด คือ ชุดใหม่และชุดเก่า การทดสอบอาจจัดทำในวันเดียว หรือ ๒ วันติดต่อกัน แต่หลักการที่สำคัญ คือ ช่วงระยะห่างของการทำแบบสอบ ๒ ชุด ต้องเป็นเวลาสั้น ๆ เพื่อมิให้เกิดประสบการณ์แทรกซ้อนที่มีผลกระทบต่อคะแนน ลำดับการทดสอบชุดใดก่อนไม่ใช่ปัญหา เพราะถือว่าองค์ประกอบต่าง ๆ อันได้แก่ การเรียนรู้ การฝึกฝน ความล้ามีผลต่อคะแนนน้อยมาก

แบบแผนที่ ๒ กลุ่มสองกลุ่ม แต่ละกลุ่มได้รับการทดสอบทั้งสองชุดในลักษณะการจับคู่กับการสอบก่อนหลังสลับให้เกิดความสมดุลกัน แบบแผนที่ ๒ ก็ดัดแปลงมาจากแบบแผนที่ ๑ เป็นแบบแผนที่ที่จะวัดประสิทธิภาพขององค์ประกอบการเรียนรู้ การฝึกฝน และความกล้า เป็นการวัดความกล้าเอียงที่จะเกิดขึ้นกับผลการสอบชุดใดชุดหนึ่ง

แบบแผนที่ ๓ กลุ่มสองกลุ่ม ให้แต่ละกลุ่มทำการสอบเพียงชุดเดียว สิ่งสำคัญของแบบแผนที่ ๓ คือ ต้องมีกลุ่ม ๒ กลุ่ม ที่มีความคล้ายคลึงกันทั้งทางด้านความรู้ ความสามารถ หรือทักษะที่ต้องการวัด เพื่อให้คะแนนสอบที่ได้มาไม่ใช่เป็นผลของความแตกต่างของความสามารถของกลุ่ม แบบแผนที่นี้ได้เคยนำไปใช้กับการทดสอบ the Graduate Management Admission Test, the Law School Admission Test และ the GRE Aptitude Test

แบบแผนที่ ๔ กลุ่มสองกลุ่ม แต่ละกลุ่มทำแบบสอบเพียงชุดเดียว และทำแบบสอบร่วมเหมือนกันอีกส่วนหนึ่ง แบบแผนนี้มีข้อได้เปรียบกว่าแบบแผนที่ ๓ ในประเด็นที่ว่า ได้ข้อสังเกตจากแบบสอบร่วมจากกลุ่มทั้งสอง แบบสอบร่วมอาจเป็นแบบแยกจากแบบสอบทั้งสองชุด (external) หรือเป็นส่วนที่ผนวกเข้าในแบบสอบทั้งสองชุด (internal) ก็ได้ ส่วนที่เป็นแบบสอบร่วมควรประกอบด้วยคำถามที่คล้ายคลึงกับคำถามในแบบสอบที่ต้องการเทียบมาตรฐานทั้งสองชุด ประโยชน์จากแบบสอบร่วมที่มีต่อการเทียบมาตรฐานของแบบสอบสองชุด จะนิยมน้อยขึ้นอยู่กับความสัมพันธ์ของแบบสอบร่วมกับแบบสอบสองชุดนั้น เหตุผลที่แบบแผนที่ ๔ มีข้อได้เปรียบ เพราะคะแนนชุดเก่าและใหม่สามารถรับการปรับเพื่อสะท้อนให้เห็นความแตกต่างที่ปรากฏอยู่ในกลุ่มทั้งสองตามผลของแบบสอบร่วมได้

แบบแผนที่ ๕ กลุ่มไม่ได้จากการสุ่มสองกลุ่ม แต่ละกลุ่มทำแบบสอบเพียงชุดเดียว และทำแบบสอบร่วมเหมือนกันอีกส่วนหนึ่ง เช่นเดียวกับแบบแผนที่ ๔ กลุ่มผู้สอบที่ไม่ได้มาจากการสุ่มมักจะเกิดขึ้นกับการสอบความซื่อสัตย์ ปัจจัยสำคัญของแบบแผนนี้อยู่ที่แบบสอบร่วม ซึ่งต้องให้มีความคล้ายคลึงมากที่สุดกับแบบสอบสองชุดที่ต้องการทำการเทียบมาตรฐาน โดยปกตินี้แล้วไม่มีวิธีการทางสถิติใดที่จะทำการปรับคะแนนในกรณีกลุ่มที่ไม่ได้มาจากการสุ่ม แต่จะทำได้ก็ต่อเมื่อมีแบบสอบร่วมที่มีความเป็นคู่ขนานกับแบบสอบชุดใหม่และเก่า หน้าที่ของแบบสอบร่วม คือ กำหนดจุดอ้างอิงให้เกิดขึ้นกับทั้งสองชุด

ในทำนองก็มีการกำหนดจุดเคี้ยว หรือจุดเยือกแข็งของน้ำ เพื่อใช้กำหนดความราบเรียบเทอร์โมมิเตอร์ ฟาร์เรนไฮต์ให้เท่ากับของสเกลเทอร์โมมิเตอร์แบบเซนติเกรดได้ แบบแผนที่ ๕ นี้ได้นำไปใช้กับการเทียบมาตราในแบบสอบ SAT มาแล้ว

ประเด็นของวิธีการทางสถิติที่ใช้ในการเทียบมาตรา โดยหลักทั่วไป แบบแผนที่สามารถได้ข้อมูลคะแนนสอบทั้งสองชุดจากผู้สอบคนเดียวกันแล้ว การหาคะแนนสมมูลทำได้ด้วยรูปแบบของการเทียบที่ตำแหน่งเปอร์เซ็นต์ และรูปแบบการเทียบเชิงเส้นตรงแบบง่าย คือ เทียบด้วยคะแนนมาตรฐาน เช่น แบบแผนที่ ๑, ๒ และ ๓ ส่วนในกรณีที่ผู้สอบคนเดียวกันทำเพียงชุดเดียว จะมีปัญหาในการประมาณค่าคะแนนสอบของคนทั้งหมดในการทำแบบสอบทั้งสองชุด ปัญหานี้เมื่อใช้แบบแผนที่ ๔ และ ๕ จะช่วยขจัดได้ วิธีทางสถิติที่ใช้ คือ การวิเคราะห์การถดถอย (regression analysis) และการประมาณด้วยความเป็นไปได้สูงสุด (Maximum likelihood) หลังจากได้ค่าประมาณผลสอบของคนทั้งหมดในแบบสอบทั้งสองชุดแล้ว ก็ทำการหาคะแนนสมมูลได้ (Marco 1981)

จากแบบแผนการเก็บรวบรวมข้อมูล และการใช้สถิติในการเทียบมาตรานี้ พบว่า แบบแผนที่สำคัญที่ก่อให้เกิดการประยุคต์ใช้ได้อย่างมาก คือ แบบแผนที่มีการใช้แบบสอบร่วม ซึ่งยังคงใช้อยู่ในโครงการทดสอบของ SAT และสอดคล้องกับเงื่อนไขในสภาพแวดล้อมปัจจุบันด้วย

### วิธีเทียบมาตรารูปแบบอีควิเปอร์เซ็นต์

#### ๑. วิธีการทั่วไป

มโนทัศน์ของการเทียบมาตรารูปแบบอีควิเปอร์เซ็นต์ เริ่มจากการแจกแจงของคะแนนจากแบบสอบชุด X และชุด Y ที่มีลักษณะคล้ายกันหรือถ้ามีความแตกต่างเกิดขึ้นบ้าง ก็มีเพียงเล็กน้อย เช่น ความถี่ที่ค่าเฉลี่ยและส่วนเบี่ยงเบนมาตรฐาน การเทียบหาคะแนนสมมูลทำได้โดยใช้คะแนน ณ ตำแหน่งเปอร์เซ็นต์เดียวกันของคะแนน ๒ ชุดนั้น ผลการเทียบมาตราแสดงด้วยกราฟขั้นตอนการแปลงคะแนนมีดังนี้ คือ เลือกกลุ่มผู้สอบที่มีความสามารถกระจาย มีทั้งเก่ง ปานกลาง และอ่อน แล้วแบ่งเป็นกลุ่มย่อยโดยสุ่ม ๒ กลุ่ม ให้กลุ่มหนึ่งทำแบบสอบ X และอีกกลุ่มทำแบบสอบ Y



ผลการสอบนำมาทำการแจกแจงคะแนน  $X$  และ  $Y$  คำนวณหาจุดกลางเปอร์เซนต์ของแต่ละการแจกแจง อ่านและทำความเข้าใจสำหรับค่า  $X$  และ  $Y$  จากการแจกแจงคะแนนที่สมนัยกับบนกระดานกราฟ ทำจุดบนกระดานประมาณ ๓๐ จุด แล้วลากเส้นเชื่อมจุดจะเกิดเป็นเส้นกราฟ ทำการปรับเส้นให้เรียบ เส้นกราฟนี้จะใช้อ่านค่า  $x$  ที่สมนัยกับ  $y$  หรืออ่านค่า  $y$  ที่สมนัยกับ  $x$  ก็ได้ จากนั้นทำการร่างสำเร็จเพื่ออ่านค่าคะแนนแปลง โดยปกติการเทียบมาครารูปแบบอิกิวเปอร์เซนต์ให้ภาพของการแปลงที่สะท้อนถึงระดับความยากของแบบสอบสองชุด ถ้าแบบสอบสองชุดมีความยากใกล้เคียงกัน เส้นกราฟจะมีลักษณะใกล้เคียงเส้นตรง แต่ถ้าแบบสอบมีความยากต่างกัน เส้นกราฟจะเป็นเส้นโค้ง (curvilinear) คะแนนสมมูลที่เกิดขึ้นจะถูกยึดหรือหย่อนคะแนนเทียบ เพื่อให้คงรักษามาตราให้เหมือนชุดก่อนตามต้องการ (Angoff 1984: 97-101)

๒. การเทียบมาตราโดยใช้แบบสอบรวม (CIE) (CIE)

บรุน และ ฮอลแลนด์ (Braun and Holland 1982: 19-22) ได้แสดงความสัมพันธ์ของแบบสอบรวมในกระบวนการเทียบมาครารูปแบบอิกิวเปอร์เซนต์ไว้ดังนี้ ให้  $T$  เป็นประชากรสังเคราะห์ ได้มาจากประชากร  $P$  และประชากร  $Q$  เขียนเป็นสมการแทนได้ว่า

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

โดยที่  $P$  และ  $Q$  รวมกันด้วยสัดส่วนน้ำหนัก  $f$  และ  $(1-f)$  ตามลำดับ การคำนวณการแจกแจงของ  $T$  ทำได้โดยคำนวณการแจกแจงความเงื่อนไข (condition distributions) ของ  $P$  และ  $Q$  ก่อนแล้วหาค่าเฉลี่ยที่ถ่วงน้ำหนักด้วย  $f$  และ  $(1-f)$  สมการปกติ คือ

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

เมื่อ  $N_1$  และ  $N_2$  เป็นขนาดกลุ่มตัวอย่างของ  $P$  และ  $Q$  ตามลำดับ ค่าของ  $f$  นี้ เป็นค่าเฉลี่ยของ  $P$  และ  $Q$

ในรูปแบบการเขียนมาตราคิวเปอร์เซนไค์ ใช้ชอคคลง (assumption) น้อยที่สุด วัตถุประสงค์ คือ หาค่าประมาณ  $F_T(x)$  และ  $G_T(y)$  ซึ่งเป็นการแจกแจงของแบบ สอบ  $x$  และ  $y$  ในประชากรสังเคราะห์  $T$  นั้นเอง ข้อมูลปริมาณที่คองนำมาวิเคราะห์ให้นิยามไว้ ดังนี้

$F_P(x v)$	เป็นฟังก์ชันการแจกแจงความถี่ของ $x$ เมื่อกำหนด $v$ ; ( $v$ in $P$ )
$F_Q(x v)$	" " " $x$ " " $v$ ; ( $v$ in $Q$ )
$G_P(y v)$	" " " $y$ " " $v$ ; ( $v$ in $P$ )
$G_Q(y v)$	" " " $y$ " " $v$ ; ( $v$ in $Q$ )
$K_P(v)$	เป็นฟังก์ชันการแจกแจงของ $v$ ในประชากรย่อย $P$
$K_Q(v)$	" " " $v$ " " " $Q$
$K_T(v)$	" " " $v$ ในประชากรสังเคราะห์ $T$

ดังนั้น จากนิยาม  $K_T(v) = fK_P(v) + (1-f)K_Q(v)$  แต่  $F_Q(x|v)$  และ  $G_P(y|v)$  ไม่สามารถประมาณได้จากข้อมูลจริงเลย จึงจำเป็นต้องระบุเป็นชอคคลงที่ไม่มี การทดสอบ ดังนี้

$$F_T(x) = \int F_P(x|v) dK_P(v) f + \int F_Q(x|v) dK_Q(v) (1-f)$$

และ

$$G_T(y) = \int G_P(y|v) dK_P(v) f + \int G_Q(y|v) dK_Q(v) (1-f)$$

การอินทิเกรตของสมการข้างต้นี้ ความจริงแล้วเป็นยอรวม เพราะทุก ๆ การ แจกแจงเป็นจำนวนเลขไม่คองเนื่อง แต่การเขียนอย่างนี้รักุมมากกว่า และเพื่อให้เกิดความง่ายขึ้น จึงกำหนดชอคคลงว่า การแจกแจงของ  $x$  เมื่อกำหนด  $v = v$  เหมือนกันในประชากรย่อย  $P$  และ  $Q$  และการแจกแจงของ  $y$  เมื่อกำหนด  $v = v$  เหมือนกันทั้งใน  $P$  แล้ว  $Q$  ฉะนั้นรูปสมการ ที่ง่ายขึ้น คือ

$$\begin{aligned} \text{ถ้า} \quad F_P(x|\nu) &= F_Q(x|\nu) \\ \text{และถ้า} \quad G_P(y|\nu) &= G_Q(y|\nu) \quad \text{แล้ว} \\ F_T(x) &= \int F_P(x|\nu) dK_T(\nu) \\ \text{และ} \\ G_T(y) &= \int G_Q(y|\nu) dK_T(\nu) \end{aligned}$$

ซึ่งหมายความว่า  $F_T(x)$  และ  $G_T(y)$  หาได้ในเทอมของปริมาณโดยประมาณจากข้อมูลที่รวบรวมได้ พังชั้นเทียบมาคร่าจะเป็นดังนี้

$$\hat{G}_X(Y) = \hat{F}_T^{-1}(\hat{G}_T(Y))$$

### ๓. ประเด็นอภิปรายในการใช้รูปแบบอิกวิเปอร์เซนไทล์

การใช้รูปแบบอิกวิเปอร์เซนไทล์ มีประเด็นที่เป็นข้ออภิปรายเสมอ คือ เทคนิคของการเกลลาหรือปรับเส้นให้เรียบ แองกอฟ (Angoff 1984: 11-12) ได้ให้รายละเอียดการเกลลาเส้นด้วยมือ โดยให้เขียนจุดสร้างกราฟลงบนกระดาษกราฟชนิดที่เป็นความน่าจะเป็นปกติ (normal probability paper) ซึ่งจะทำให้การแจกแจงที่เป็นโค้งปกติมีภาพเป็นเส้นตรง และง่ายต่อการเกลลาด้วยมือ อุปกรณ์ที่ช่วยในการลากเส้น ได้แก่ กระจุกงู หลักการเกลลาโดยทั่วไป คือ ให้ลากเส้นผ่านจุดต่าง ๆ ที่ไม่อยู่ในลู่เดียวกันในลักษณะให้เฉลี่ยความไกลเคียงกันทั้งสองข้าง ความชำนาญของผู้ปฏิบัติจะสามารถมองเห็นความผิดปกติของเส้นได้ นอกจากการเกลลาด้วยมือแล้ว มีเทคนิคการเกลลาอีกกลุ่มหนึ่ง คือ เทคนิควิธีเชิงวิเคราะห์ซึ่งจะทำการวิเคราะห์ก่อนเขียนกราฟ วิธีเชิงวิเคราะห์เป็นวิธีเกลลาเส้นอย่างเป็นปรนัย ไม่เกิดความลำเอียงทั้ง เช่นวิธีไข่มื้อ วิธีเชิงวิเคราะห์มีหลายวิธี เช่นวิธีของ เคอร์ตัน และ คูทท์ ใช้วิธีการทำให้จุดต่าง ๆ ที่อยู่กันเป็นชุดเฉลี่ยให้ใหม่แนวใหม่เป็นพาราโบลิกหรือคิวบิก (Cureton and Tukey 1951 cited by Angoff 1984: 12) ถ้าการแจกแจงเชิงทฤษฎีกับการแจกแจงเชิงประจักษ์สอดคล้องกัน แองกอฟได้แนะนำให้ใช้วิธีการ (the negative hypergeometric (Angoff 1984: 13)) นอกจากนี้มีวิธีการปรับเกลลาเทียบคะแนนเชิงเส้นตรง (smoothing linearly interpolated equated score) ซึ่งลินเช และ ปรีชาค





(Lindsay and Prichard 1971 cited by Kolen 1984) เสนอไว้ แนววิธีนี้มีข้อจำกัดสองประการ คือ ผลการเทียบให้คะแนนแปลงที่รากคุณสมบัติของความสมมาตร และไม่สามารถระบุความคลาดเคลื่อนมาตรฐานของคะแนนเทียบเมื่อเปลี่ยนจุด โคลเลน (Kolen 1984) จึงได้ทำการศึกษากันไรจุดอ่อนโดยใช้เทคนิค Cubic splines ซึ่งเป็นวิธีทางคณิตศาสตร์ ตามปกติ Cubic splines เป็นเครื่องมือที่ใช้ปรับเส้นโค้งธรรมชาติที่ใช้ในการเขียนแบบ ปัจจุบันได้พัฒนาไปใช้กับสถานการณ์ต่าง ๆ ที่ต้องการทำเส้นโค้งให้ติดกับข้อมูล การประเมินประสิทธิภาพของการทดลองใช้วิธีนี้ตัดสินโดยเปรียบเทียบผลลัพธ์จากการปรับอีควิเปอร์เซนไทล์กับผลจากการเทียบมาตราด้วยวิธีอื่น ๆ จากการวิเคราะห์กลุ่มสอบทานผล ข้อมูลได้จากกลุ่มตัวอย่างสุ่มสมมูลประมาณ ๓๐๐๐ คนต่อการทำแบบสอบ ๑ ชุด จากแบบสอบทั้งหมด ๔ ชุด ผลการศึกษา พบว่า การเทียบมาตราโดยทำการปรับเส้นด้วยวิธีเชิงวิเคราะห์ตามเทคนิคของ cubic splines เป็นวิธีที่มีความเพียงพอสูงสุดในการเทียบมาตราของแบบสอบ AAP โดยเฉพาะอย่างยิ่งกับชุดแบบสอบที่มีความคล้ายคลึงน้อยที่สุด

ลอร์ด (Lord 1982a) ยอมรับว่า การปรับหรือเกลารเส้นโค้งของการเทียบมาตราในรูปแบบอีควิเปอร์เซนไทล์ช่วยลดความคลาดเคลื่อนเชิงสุ่มได้ส่วนหนึ่ง แต่มีโอกาสทำให้เกิดความลำเอียงได้เช่นกัน และมีความเห็นว่า การระบุความคลาดเคลื่อนของคะแนนแปลงที่ปรับแล้วเป็นงานหนักเกินความจำเป็น

รูปแบบอีควิเปอร์เซนไทล์ มีข้อจำกัดที่ต้องคำนึงอยู่หลายประการ คือ

(๑) ในการเขียนเส้นกราฟที่แสดงคะแนนสมมูลที่ยังไม่ได้ปรับหรือเกลารให้เปรียบเทียบ จะมีความคลาดเคลื่อนมาก แต่เทคนิคการปรับเส้นเท่าที่ทำได้มานั้นยังไม่มียุติวิธีใดที่ให้หลักประกันผลได้อย่างน่าเชื่อถือ ลักษณะเช่นนี้ยังก่อให้เกิดความลำเอียงที่ยังคงหาทางปรับปรุงอีก (Angoff 1984: 97; Potthoff 1982: 209)

(๒) รูปแบบอีควิเปอร์เซนไทล์ มีความไวต่อความแปรปรวนเชิงสุ่มมาก โดยเฉพาะขนาดตัวอย่างที่มีขนาดเล็ก (Angoff 1984: 97; Potthoff 1982: 210)

(๓) ในกรณีที่แบบสอบสองชุดมีความเที่ยงต่างกันมาก ผลของการเทียบมาตราจากความคงที่ (Potthoff 1982: 210)

(๔) การเทียบคะแนนของผู้สอบของแบบสอบสองชุดทำได้เฉพาะในช่วงพิสัยของคะแนนที่มีความถี่ของคะแนนสังเกตเท่านั้น ส่วนที่อยู่นอกพิสัยดังกล่าวจะมีความคลาดเคลื่อนของการเทียบมากกว่าสูงมาก (Angoff 1984: 97)

### วิธีเทียบมาตรารูปแบบเชิงเส้นตรง

รูปแบบการเทียบมาตราเชิงเส้นตรง เป็นรูปแบบที่ได้รับความนิยมมากมาเป็นเวลานาน ได้ใช้ในโครงการทดสอบระดับชาติหลายโครงการ เช่น การทดสอบความถนัดเชิงวิชาการ (SAT) (Donlon and Angoff 1971: 32) เริ่มด้วยการเทียบมาตราคะแนนชุดต่าง ๆ ที่ใช้เวลาคือ ๆ มากกับแบบสอบฉบับที่แบบที่ใช้เมื่อมิถุนายน ๑๙๖๐ ตัวอย่างการเทียบมาตราที่จัดกระทำเมื่อปี ค.ศ. ๑๙๖๐ ได้ดำเนินการดังนี้ คือ ออกแบบเพื่อการเทียบมาตราด้วยการใช้ข้อสอบจำนวนหนึ่งที่เป็นของชุดที่ใช้เมื่อ เมษายน ๑๙๖๐ เข้าไปรวมอยู่ในชุดใหม่ที่ใช้เมื่อมิถุนายน ๑๙๖๐ คะแนนเฉลี่ยและคะแนนเบี่ยงเบนมาตรฐานของกลุ่มข้อสอบจำนวนนี้ที่รวมกัน (the carried over items หรือ common items) ได้ปรากฏว่า แยกค่างกันในกลุ่มผู้สอบในเดือน เมษายน และมิถุนายน ความแตกต่างของค่าสถิติทั้งสองนี้ได้นำไปใช้ปรับค่าสถิติอื่น ๆ ของกลุ่มผู้สอบเหล่านั้นที่รับการทดสอบจากแบบสอบต่างชุดกัน หลังจากการปรับได้ทำให้เกิดคะแนนเฉลี่ยและส่วนเบี่ยงเบนมาตรฐานของคะแนนดิบที่ปรับแล้วของชุดใหม่ที่เท่ากับคะแนนเฉลี่ยและส่วนเบี่ยงเบนมาตรฐานของคะแนนมาตราของชุดเก่า ด้วยวิธีการเช่นนี้ทำให้เกิดการแปลงคะแนนเชิงเส้นตรง (a linear conversion) จากคะแนนดิบไปสู่คะแนนมาตราสำหรับแบบสอบชุดใหม่ได้ ในการสอบเมื่อเดือน เมษายน ค.ศ. ๑๙๖๐ ได้ดำเนินการเทียบด้วยวิธีการเดียวกัน

#### ๑. วิธีเทียบโดยใช้แบบสอบรวม (CIE)

ผู้ที่คิดการเทียบมาตราคะแนนโดยใช้แบบสอบรวมในรูปแบบเชิงเส้นตรง คือ เลดยาร์ด ทักเกอร์ (Ledyard Tucker) และได้ตั้งชื่อวิธีการเทียบมาตรานี้ว่า Tucker equating เพื่อเป็นเกียรติแก่ผู้ริเริ่ม (Angoff 1961 cited by Donlon and Angoff 1971: 38; Braun and Holland 1982: 23) วิธีนี้เริ่มต้นด้วยการตั้งข้อตกลงเบื้องต้นว่า การเทียบมาตราคะแนนทั้งสองขั้นเป็นเส้นตรงดังนี้

$$e_X(Y) = \mu_X + \frac{\sigma_X}{\sigma_Y} (Y - \mu_Y) \text{ -----(๑)}$$

เมื่อ  $\mu_X$   $\sigma_X$   $\mu_Y$  และ  $\sigma_Y$  เป็นค่าเฉลี่ยและส่วนเบี่ยงเบนมาตรฐานของประชากร T ที่คนทำแบบสอบถาม X และ Y พิจารณาเฉลี่ยของประชากรในแบบสอบถาม X คือ

$$\mu_X = \int E_P(X|\nu) dK_P(\nu) f + \int E_Q(X|\nu) dK_Q(\nu) (1-f) \text{ -----(๒)}$$

และมีความแปรปรวนของ X คือ

$$\begin{aligned} \sigma_X^2 = & \int \text{Var}_P(X|\nu) dK_P(\nu) f + \int \text{Var}_Q(X|\nu) dK_Q(\nu) (1-f) \\ & + \int E_P^2(X|\nu) dK_P(\nu) f + \int E_Q^2(X|\nu) dK_Q(\nu) (1-f) - \mu_X^2 \end{aligned} \text{ -----(๓)}$$

สำหรับทั้งค่าเฉลี่ยและความแปรปรวนของประชากรในแบบสอบถาม Y เป็นในทำนองเดียวกัน ซึ่งสัญลักษณ์ต่าง ๆ มีความหมายแทนสิ่งต่อไปนี้ คือ

P คือ ประชากร P ที่ทำแบบสอบถาม X และแบบสอบถาม Y

Q คือ ประชากร Q ที่ทำแบบสอบถาม Y และแบบสอบถาม X

T คือ ประชากรสังเคราะห์มาจากการรวม P และ Q

f เป็นสัดส่วนน้ำหนักของประชากร P ใน T

1-f เป็นสัดส่วนน้ำหนักของประชากร Q ใน T

E และ Var คือ ค่าคาดหวังและความแปรปรวนของประชากร

$K_P(\nu)$  และ  $K_Q(\nu)$  เป็นฟังก์ชันการแจกแจงของแบบสอบถาม  $\nu$  ใน P และ Q ตามลำดับ

เทอมต่าง ๆ ในสมการ (๒) และ (๓) ที่สามารถหาค่าประมาณจากข้อมูล คือ

$$E_P(X|\mathbf{v}) = \int x dF_P(x|\mathbf{v})$$

$$\text{Var}_P(X|\mathbf{v}) = \int [x - E_P(X|\mathbf{v})]^2 dF_P(x|\mathbf{v})$$

และเป็นทำนองเดียวกันในเทอม  $E_Q(X|\mathbf{v})$  และ  $\text{Var}_Q(X|\mathbf{v})$

จากสมการข้างต้นจะพบว่า ค่า  $\mu_X$  และ  $\sigma_X^2$  หาได้ก็ต่อเมื่อได้กำหนดเงื่อนไขเป็นข้อตกลงเบื้องต้น ๒ ประการก่อน คือ

(๑) ค่าประมาณเฉลี่ยและค่าความแปรปรวนอย่างมีเงื่อนไขของ  $X$  เมื่อกำหนด  $\mathbf{v}$  ในประชากรย่อย  $P$  เหมือนกับในประชากร  $Q$  และในกรณี  $Y$  เมื่อกำหนด  $\mathbf{v}$  ในประชากรย่อย  $Q$  ย่อมเหมือนกับในประชากร  $P$  เขียนในรูปสูตรดังนี้

$$\text{ถ้า } E_P(X|\mathbf{v}) = E_Q(X|\mathbf{v}) \quad \text{และ}$$

$$\text{ถ้า } \text{Var}_P(X|\mathbf{v}) = \text{Var}_Q(X|\mathbf{v}) \quad \text{แล้ว}$$

$$\mu_X = \int E_P(X|\mathbf{v}) dK_T(\mathbf{v}) \quad \text{และ} \quad \text{-----} \quad (๔)$$

$$\sigma_X^2 = \int \text{Var}_P(X|\mathbf{v}) dK_T(\mathbf{v}) + \int E_P^2(X|\mathbf{v}) dK_T(\mathbf{v}) - \mu_X^2 \quad \text{---} \quad (๕)$$

สำหรับค่า  $\mu_Y$  และ  $\sigma_Y^2$  หาได้ในทำนองเดียวกัน (๔) และ (๕)

(๒) เป็นข้อตกลงเบื้องต้นตามรูปแบบ (model assumption) เพราะเป็นการยอมรับการทำ  $E_P(X|\mathbf{v}) = E_Q(X|\mathbf{v})$   $\text{Var}_P(X|\mathbf{v})$  และ  $\text{Var}_Q(X|\mathbf{v})$  ให้มีรูปที่ง่ายขึ้น ข้อตกลงเบื้องต้นเหล่านี้เป็นแบบฉบับที่กำหนดในการวิเคราะห์การถดถอยว่า ค่าคาดหมายอย่างมีเงื่อนไขเป็นเส้นตรง และความแปรปรวนอย่างมีเงื่อนไขเป็นค่าคงที่ ดังเช่น



$$E_P(X|v) = (a_{X|V.P})v + b_{X|V.P}$$

$$E_Q(Y|v) = (a_{Y|V.Q})v + b_{Y|V.Q}$$

$$\text{Var}_P(X|v) = \sigma_{X|V.P}^2$$

$$\text{Var}_Q(Y|v) = \sigma_{Y|V.P}^2$$

ผลจากการกำหนดข้อตกลงเบื้องต้น ๒ ประการที่กล่าวมานี้ สรุปเป็นหลักของการเทียบมาตรฐานวิธีการของทักเกอร์ (Tucker equating) ได้ดังนี้

$$\text{ถ้า } E_P(X|v) = E_Q(X|v) = (a_{X|V.P})v + b_{X|V.P}$$

$$\text{และ } \text{Var}_P(X|v) = \text{Var}_Q(X|v) = \sigma_{X|V.P}^2 \quad \text{แล้วจะได้}$$

$$\mu_X = (a_{X|V.P}) \mu_{V.T} + b_{X|V.P} \quad \text{----- (๖)}$$

และ

$$\sigma_X^2 = \sigma_{X|V.P}^2 + (a_{X|V.P})^2 \sigma_{V.T}^2 \quad \text{----- (๗)}$$

ในทำนองเดียวกัน สามารถหาค่า  $\mu_Y$  และ  $\sigma_Y^2$  จาก (๖) และ (๗)

ดังนี้

$$\mu_Y = (a_{Y|V.Q}) \mu_{V.T} + b_{Y|V.Q} \quad \text{----- (๘)}$$

และ

$$\sigma_Y^2 = \sigma_{Y|V.Q}^2 + (a_{Y|V.Q})^2 \sigma_{V.T}^2 \quad \text{----- (๙)}$$

ในการคำนวณค่าต่าง ๆ ตามข้อตกลงเบื้องต้นนี้ ให้ใช้ค่าจากกลุ่มตัวอย่าง เพื่อหาค่าประมาณกำลังสองน้อยที่สุด (least squares estimates) ของ  $a_{X|V.P}$

$b_{X|V.P}$   $\sigma_{X|V.P}^2$   $a_{Y|V.Q}$   $b_{Y|V.Q}$  และ  $\sigma_{Y|V.P}^2$  ข้อมูลทั้งหมดเมื่อนำมารวมเข้าด้วยกันจะประมาณค่า  $\mu_{V.T}$  และ  $\sigma_{V.T}^2$  จากนั้นนำค่าประมาณเหล่านี้แทนในสมการ (๖)

จะได้ฟังก์ชันของการเทียบมาตรฐานเชิงเส้นตรงของ Y ไปสู่ X (Braun and Holland

1982: 23-24)

การหาค่าประมาณค่าเฉลี่ยและความแปรปรวนของประชากรของคะแนนจาก  
 ชุค X และชุก Y ด้วยวิธีการของทักเกอร์ของประมาณการคอมของกลุ่มตัวอย่าง Q ในการทำแบบสอบ  
 X และพิจารณาผู้สอบ ๒ กลุ่ม ที่มีความแตกต่างกันตามธรรมชาติ ต่อมาลอว์ (Lord 1955)  
 ได้คิดวิธีการหาค่าประมาณค่าเฉลี่ยและความแปรปรวนของประชากร โดยพิจารณาจากผู้สอบเมื่อได้  
 รับการสุ่มมาจากประชากรเดียวกัน จะประมาณค่าการคอมหรือคะแนนของทั้งประชากรในแบบสอบ  
 ทั้งสองชุก แทนที่จะประมาณเพียงกลุ่ม Q ในการคอมชุก X เท่านั้น วิธีการของลอว์ที่นำมาใช้หา  
 ค่าประมาณ คือ วิธีการความเป็นไปได้สูงสุด (the maximum likelihood method)  
 ได้สมการที่ใช้หาค่าประมาณดังนี้

$$\hat{\mu}_X = M_{X\alpha} + b_{XV\alpha} (\hat{\mu}_V - M_{V\alpha}) \quad \text{--- (๑๐)}$$

$$\hat{\mu}_Y = M_{Y\beta} + b_{YV\beta} (\hat{\mu}_V - M_{V\beta}) \quad \text{--- (๑๑)}$$

$$\hat{\sigma}_X^2 = S_{X\alpha}^2 + b_{XV\alpha}^2 (\hat{\sigma}_V^2 - S_{V\alpha}^2) \quad \text{--- (๑๒)}$$

$$\hat{\sigma}_Y^2 = S_{Y\beta}^2 + b_{YV\beta}^2 (\hat{\sigma}_V^2 - S_{V\beta}^2) \quad \text{--- (๑๓)}$$

เมื่อ

$$\hat{\mu}_V = M_{Vt}$$

$$\hat{\sigma}_V^2 = S_{Vt}^2$$

และ

$$t = \alpha + \beta$$

ค่าประมาณเหล่านี้นำไปแทนในสมการที่ (๑) จะเป็นสมการเทียบมาตรา

เชิงเส้นตรง

ในการศึกษารังนี้ ผู้วิจัยได้ใช้วิธีการสุ่มกลุ่ม ๒ กลุ่ม ให้ทำแบบสอบถาม x หรือ y กลุ่มละแบบสอบถามเดียว และทั้ง ๒ กลุ่ม ได้ทำแบบสอบถามร่วม ๆ ด้วย การเลือกวิธีการทางสถิติมาใช้ หากค่าประมาณประชากรต่าง ๆ จึงได้พิจารณาเลือกใช้วิธีการความเป็นไปได้สูงสุด (the maximum likelihood) ตามที่ลอว์คได้เสนอไว้ (Angoff 1984: 105; Lord 1955: 198)

## ๒. ประเด็นอภิปรายในการใช้รูปแบบเชิงเส้นตรง

โดยทั่วไปการเทียบมาตราเชิงเส้นตรง เป็นรูปแบบที่ให้ทั้งความสะดวก และความ เป็นปรนัย แต่ก็ยังมีข้อจำกัดในการใช้อยู่มาก โดยเฉพาะแบบสอบถามที่มีความยากต่างกัน พอททอฟ (Potthoff) ได้แนะนำทางเลือกอื่นเมื่อพบว่า รูปแบบเชิงเส้นตรงยังไม่เพียงพอ ดังนี้ (Potthoff 1982: 209-210)

(๑) ปรับปรุงวิธีการให้คะแนน เพื่อให้คะแนนไม่เกาะกลุ่มมากเกินไป ซึ่งเป็นสาเหตุหนึ่งที่ทำให้การเทียบมาตราเชิงเส้นตรงไม่ได้ผลน่าพอใจ

(๒) ควรใช้วิธีการตรวจสอบที่เป็นปรนัย เพื่อช่วยในการตัดสินใจว่า ควรใช้รูปแบบเชิงเส้นตรงหรือไม่ ในกรณีที่วิธีการเทียบที่มีข้อบกพร่องเบื้องต้นของความเที่ยงเท่ากัน อาจทดสอบ ความแตกต่างของการแจกแจงของคะแนนภายหลังการเทียบดูว่า ยังคงเหมือนกันหรือไม่ ถ้าพบว่า ไม่สามารถปฏิเสธสมมุติฐานศูนย์ ก็แสดงว่ารูปแบบเชิงเส้นตรงใช้ได้แล้ว ไม่จำเป็นต้องหาวิธีการอื่น ๆ ที่ไม่ใช่เชิงเส้นตรง

(๓) ถ้าพบว่ารูปแบบเชิงเส้นตรงไม่เหมาะในการเทียบมาตราแล้ว ให้พิจารณา รูปแบบอื่น แต่หากพบว่ารูปแบบอื่นก็ยังให้ผลไม่ดีไปกว่าเชิงเส้นตรง อาจตัดสินใจใช้เชิงเส้นตรง แต่ควรมีรายงานประกอบเพื่อช่วยให้ผู้นำผลของคะแนนแปลงไปใช้ได้ เห็นข้อจำกัดที่ปรากฏอยู่ สิ่งที่ควรมีในรายงาน คือ

ก. ระบุลักษณะเฉพาะของแบบสอบถาม คะแนนแปลง และให้ข้อสังเกตเพื่อเตือนถึงผลกระทบของคะแนนต่อการเทียบในช่วงพิสัยต่าง ๆ

ข. น่าจะเสนอพิสัยคะแนนแปลงในช่วงความเชื่อมั่นที่เหมาะสม โดยเฉพาะ ในกรณีที่แบบสอบถามสองชุดมีความเที่ยงไม่เท่ากัน ช่วงความเชื่อมั่นย่อมต่างกัน

(๔) ในบางกรณีพบว่า รูปแบบเชิงเส้นตรงให้ประสิทธิภาพดีมาก ยกเว้นคะแนนที่กระจายอยู่ที่ปลายสุดทั้งสองข้าง อาจดำเนินการเทียบมาตรฐานเสริมในส่วนนี้ ปัญหาที่เกิดขึ้นจากแบบสอบ ๒ ชุด อาจมีความยาก (หรือง่าย) ไม่เพียงพอที่จะจำแนกได้ก็ในส่วนปลายที่สูงสุดหรือต่ำสุด การแก้ไข คือ นำแบบสอบ ๒ ชุดนี้ไปทำการสอบกับคนที่มีระดับความสามารถในระดับความยากนั้น ๆ เพิ่มเติม

วิธีการเทียบมาตรฐานด้วยรูปแบบอิงทฤษฎีการตอบสนอง

ลอร์ด (Lord 1977) ได้แนะนำการพิจารณาเลือกวิธีเทียบมาตรฐานคะแนนของแบบสอบสองชุดไว้ว่า ถ้าไม่รู้แน่ว่าแบบสอบ ๒ ฉบับ ที่กำลังจะทำการเทียบมาตรฐาน มีความเป็นคู่ขนานกันหรือไม่ วิธีการที่ปลอดภัยที่สุด คือ ให้จัดกระทำกับแบบสอบนั้นในฐานะไม่เป็นคู่ขนานกัน และใช้วิธีการประเภทที่ไม่ใช่โมเดลเส้นตรง การออกแบบเพื่อศึกษาการเทียบมาตรฐานโดยอิงทฤษฎีการตอบสนอง สามารถนำมาใช้กับคน ๒ กลุ่ม ที่ต่างกัน และใช้กับแบบสอบ ๒ ชุด ที่ไม่เป็นคู่ขนานกัน ซึ่งวิธีการเทียบมาตรฐานแบบดั้งเดิมไม่สามารถทำได้เหมาะสมนัก ในกรณีเช่นนี้การเทียบมาตรฐานแบบสอบรวมเข้ามาช่วยพิจารณา และปรับความแตกต่างของความสามารถ ๒ กลุ่ม แต่ละกลุ่มของผู้สอบจะสขเพียงฉบับเดียวกับแบบสอบรวม ตัวอย่างการเทียบมาตรฐานด้วยแบบแบบนี้ เช่น มาร์โก (Marco 1977 cited by Lord 1977) ได้วางแผนศึกษาการเทียบมาตรฐานสอบวัดผลสัมฤทธิ์วิชาแคลคูลัส ฉบับเพื่อการจัดชั้นเรียนในระดับวิทยาลัย ซึ่งมี ๔๕ ข้อ (College Board Advanced placement program: AP) และฉบับโปรแกรมการสอบระดับวิทยาลัย (The College Level Examination Program: CLEP) ซึ่งง่ายกว่าฉบับแรกเล็กน้อย มี ๕๐ ข้อ และมีแบบสอบรวม (Anchor test) ๑๗ ข้อ แทรกอยู่ในแต่ละฉบับ กลุ่มตัวอย่างผู้สอบชุด AP จำนวน ๑๒๐๐ คน เมื่อ พฤษภาคม ๑๙๗๒ ส่วนกลุ่มตัวอย่างผู้สอบชุด CLEP จำนวน ๑๒๐๘ คน เมื่อ พฤษภาคม ๑๙๗๔ กระบวนการเทียบมาตรฐานเริ่มด้วยการใช้โปรแกรมคอมพิวเตอร์ประมาณค่าประชากรข้อสอบ และความสามารถจากการทำงานของเครื่อง ๑ ครั้ง ได้ค่าประมาณพร้อม ๆ กันทั้งหมด ทั้งที่เป็นค่าประมาณความสามารถ มี ๒๒๒๒ ตัว (ทั้ง ๒ กลุ่มรวมกัน และหักออก ๒๔๓ คน ที่สอบไม่สมบูรณ์) และค่าประมาณข้อสอบ ๗๕ ข้อ (45+50-17) เป็นค่า a b และ c สำหรับทุกข้อ รายละเอียดการใช่แบบสอบรวมในรูปแบบนี้ จะแยกกล่าวเป็น ๒ ส่วน



คือ การเทียบคะแนนจริง และคะแนนสังเกต

๑. การเทียบมาตราคะแนนจริง

พิจารณาจากความสัมพันธ์เชิงคณิตศาสตร์ระหว่างความสามารถ และจำนวนคะแนน  
นับความข้อที่ทำได้ คะแนนจริง  $\xi$  จากแบบสอบ  $X$  และ  $\eta$  จากแบบสอบ  $Y$  เป็นดังต่อไปนี้

$$\xi = \xi(\theta) = \sum_{i=1}^{n_X} P_i(\theta) \quad \text{----- (๑๔)}$$

$$\eta = \eta(\theta) = \sum_{j=1}^{n_Y} P_j(\theta) \quad \text{----- (๑๕)}$$

เมื่อ  $n_X$  และ  $n_Y$  เป็นจำนวนข้อในแบบสอบชุด  $X$  และชุด  $Y$  ตามลำดับ สมการทั้งสองเป็น  
สมการคณิตศาสตร์ ดังนั้น ถ้าแทนค่าความสามารถ  $\theta$  ลงในสมการทั้งสองก็จะได้คะแนนจริงที่สมมูล  
กันมาตราเดียวกัน หรืออีกวิธีหนึ่งทำสมการทั้งสองให้เท่ากัน โดยกำจัดค่า  $\theta$  ให้หายไป ก็จะ  
สามารถหาคะแนนสมมูลของ  $X$  และ  $Y$  ได้ การหาค่าที่เทียบมาตราเดียวกันนี้ ใช้คำนวณจากค่า  
 $P_i(\theta)$  และ  $P_j(\theta)$  ซึ่งใช้ค่าประมาณแทน คือ  $\hat{P}_i(\theta)$  และ  $\hat{P}_j(\theta)$  ที่ได้จากการแทน  
ค่า  $\hat{a}$ ,  $\hat{b}$  และ  $\hat{c}$  (Lord 1977)

การใช้โปรแกรม LOGIST จะวิเคราะห์ข้อมูลการตอบของทุกคน และข้อสอบทุกข้อ  
พร้อม ๆ กัน โดยถือว่าข้อสอบชุดหนึ่งที่ยุ่สอบไม่ได้ทำนั้นเป็นส่วนที่ยุ่สอบทำไม่ถึง (not reached)  
วิธีนี้ทำให้การหาค่าประมาณประชากรปรากฏออกมามาตราเดียวกันหมด นับว่าเป็นวิธีที่มีประสิทธิภาพ  
โดยมีสมมุติฐานเบื้องต้นว่า ความสามารถ  $\theta$  ใด ๆ ของผู้สอบไม่มีความแตกต่างกันจากแบบสอบ  
หนึ่งไปยังอีกแบบสอบหนึ่ง (Lord 1980: 20)

การใช้แบบสอบร่วมทำได้ ๒ กรณี คือ ให้แบบสอบร่วมผนวกเข้าเป็นชุดเดียวกับ  
แบบสอบที่ต้องการเทียบ ซึ่งเรียกว่า แบบสอบร่วมภายใน (internal anchor test) เช่น  
ตัวอย่างแบบสอบร่วม ๑๓ ข้อ ในวิชาแคลคูลัสที่กล่าวถึงข้างต้น ส่วนในกรณีที่จัดแยกเป็นชุดแบบสอบร่วม

ต่างหากจากแบบสอบที่ทำการเทียบ เรียกว่า แบบสอบร่วมภายนอก (external anchor test) จุดมุ่งหมายของแบบสอบร่วม คือ เชื่อมข้อมูลค่าคอมเข้าด้วยกัน เพื่อให้ค่าพารามิเตอร์ที่วิเคราะห์ออกมาอยู่บนมาตราเดียวกัน ถ้าหากข้อสอบร่วมจำนวนนี้แล้ว จะไม่สามารถเทียบกันได้ นอกเสียจากว่ากลุ่มผู้สอบ X และ Y จะมีการแจกแจงของความสามารถเหมือนกัน (Lord 1980: 202)

๒. การเทียบมาตราคะแนนสังเกต (Observed score equating)

จากการเทียบมาตราด้วยคะแนนจริงคงที่กล่าวมาแล้ว ทำให้เห็นปัญหาในทางปฏิบัติ คือ การประมาณคะแนนของผู้สอบจากค่าคอม  $\hat{\xi} = \sum_1 \hat{p}_1(\hat{\theta})$  แท้จริง คือ คะแนนที่มีความคลาดเคลื่อนอีกชุดหนึ่งนั่นเอง ไม่มีคุณสมบัติของคะแนนจริง ซึ่งลอร์ด (Lord 1980 : 196-198, 202) ได้พิสูจน์ให้เห็นแล้วว่า การเทียบมาตราจะทำให้ได้ข้ออย่างเข้มงวดไม่ได้ นอกเสียจากว่าแบบสอบ ๒ ฉบับนี้คงที่คู่ขนานกันจริง ลอร์ดได้เสนอการเทียบด้วยคะแนนสังเกต ซึ่งเป็นวิธีการเทียบมาตราโดยประมาณทำได้ในทางปฏิบัติที่สถานการณ์การสอบไม่สามารถจัดให้ผู้สอบทำแบบสอบได้ทั้งสองชุด กระบวนการเริ่มด้วยการประมาณการแจกแจงความสามารถของคนในกลุ่มผู้สอบ  $\hat{r}(\theta)$  (หมายถึง a, หรือ b หรือ c) โดยใช้การแจกแจงจริงของ  $\hat{\theta}_a$  ในกลุ่มเป็นตัวอย่างประมาณ  $r(\theta)$  การแจกแจงของคะแนนดิบหาได้จากสมการต่อไปนี้

$$\hat{\phi}_x(x) = \frac{1}{N} \sum_{a=1}^N \phi_x(x | \hat{\theta}_a) \quad \text{เมื่อ } a = 1, 2, \dots, N$$

ถ้า  $\hat{r}(\theta)$  เป็นค่าต่อเนื่อง จะประมาณค่า  $\hat{\phi}(x)$  จากสมการ

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

ทั้งนี้  $\phi_x(x | \hat{\theta}_a)$  ใช้ประมาณค่าประชากรข้อสอบในฉบับ x ในทำนองเดียวกัน ประมาณค่าประชากรของแบบสอบ Y และเนื่องจากแบบสอบ X และ Y ต่างมีการแจกแจงเป็นอิสระกัน เมื่อกำหนดให้  $\theta$  คงที่ ก็สามารถประมาณส่วนที่เป็นการกระจายร่วมกันของคะแนน X และ Y สำหรับกลุ่มที่ถูกประมาณโดย

$$\hat{\phi}(x,y) = \frac{1}{N} \sum_{a=1}^N \hat{\phi}_x(x|\hat{\theta}_a) \hat{\phi}_y(y|\hat{\theta}_a)$$

หรือโดย

$$\hat{\phi}(x,y) = \int_{-\infty}^{\infty} \hat{\phi}_x(x|\theta) \hat{\phi}_y(y|\theta) r(\theta) d\theta \text{ ---- (๑๖)}$$

จากที่กล่าวมานี้ จะเห็นพาทของแบบสอบร่วมซิกเจน ทำให้สามารถประมาณการแจกแจงร่วม (joint distribution) ของ X และ Y ถึงแม้จะไม่มีผู้ใดสอบทั้ง ๒ ฉบับ

สมการ ๑๖ เป็นการแจกแจงร่วมกันของตัวแปรสามตัว คือ  $\theta, x$  และ  $y$  เนื่องจาก  $\theta$  เป็นตัวระบุคะแนนจริง  $\xi$  และ  $\eta$  การแจกแจงนี้จึงนับว่าเป็นการแจกแจงร่วมกันของตัวแปร ๔ ตัว คือ  $\xi, \eta, x$  และ  $y$  สมการนี้ไม่มีทางอื่นใดจะประมาณได้ นอกจากจากความสัมพันธ์เชิงอิกวิเปอร์เซนไคล์ระหว่าง  $x$  และ  $y$  จากสมการสุดท้ายข้างต้น ๒ สมการ วิธีการเป็นการเทียบมาตรฐานคะแนนของแบบสอบโดยประมาณ

๓. ประเด็นอภิปรายเกี่ยวกับการเทียบมาตรฐานแบบ IRT

จากทฤษฎีการเทียบคะแนนดิบของแบบสอบสองชุดด้วยวิธีการเทียบคะแนนจริง และวิธีเทียบคะแนนสังเกต หรือการเทียบอิกวิเปอร์เซนไคล์ใน IRT ทั้งที่กล่าวมาแล้วนั้น มีข้อได้เปรียบและเสียเปรียบ วิธีเทียบคะแนนจริง  $\xi$  และ  $\eta$  ให้ความหมายของคะแนนสมมุติเฉพาะคะแนนที่อยู่เหนือค่าเฉลี่ยของการเคา ส่วนวิธีการใช้การประมาณคะแนนสังเกต ก็มีจุดอ่อนที่เป็นเพียงการเทียบมาตรฐานอย่างประมาณเท่านั้น ข้อข้อคำถามว่า วิธีใดจะให้ผลการเทียบดีกว่ากันนั้น ลอร์ดได้กล่าวว่า ทั้งสองวิธีมีความสอดคล้องกันมาก แต่การสรุปเปรียบเทียบอ้างอิงทั่วไปต้องทำอย่างพิถีพิถัน (Lord 1980: 203)

โคเลน (Kolen 1981) ได้ทดลองใช้วิธีอิกวิเปอร์เซนไคล์ IRT พบว่า จากตัวอย่างแบบสอบและสถานการณ์ผู้สอบเฉพาะในการวิจัย วิธีนี้เป็นวิธีที่ดี โดยใช้เกณฑ์ประเมินความคงเส้นคงวาในการวิเคราะห์หักกลุ่มตัวอย่างสอบทานผล

ปี ๑๘๘๔ ลอร์ด และ วินเกอส์กี (Lord and Wingersky 1984)

ได้ศึกษาเปรียบเทียบผลการเทียบมาตรฐานรูปแบบ IRT ด้วยวิธีคะแนนจริงและคะแนนสังเกตพบว่า ทั้งสองวิธีแทบจะไม่ได้มีข้อแตกต่างกันเลย และได้วิจารณ์ไว้ว่า วิธีคะแนนสังเกตมีความซับซ้อนในการคำนวณและลงทุนด้านอื่น ๆ มากกว่าวิธีคะแนนจริง

### ความยาวของแบบสอบรวม

จากการเสนอแบบแผนการจัดการเทียบมาตรฐานที่ใช้แบบสอบรวมอย่างละเอียดของแองกอฟ (Angoff 1984: 106-107) ได้แนะนำให้ใช้แบบสอบรวมที่มีความยาวเพียงพอที่จะก่อให้เกิดข้อสนเทศเพื่อนำไปปรับความแตกต่างระหว่างกลุ่มผู้สอบอย่างมีประสิทธิภาพ จำนวนข้อที่เหมาะสมคือ ต้องไม่น้อยกว่า ๒๐ ข้อ หรือไม่น้อยกว่าร้อยละ ๒๐ ของจำนวนข้อในแบบสอบที่ต้องการเทียบมาตรฐาน แล้วแต่จำนวนไหนจะมากกว่ากัน ขณะเดียวกันไรท์ (Wright 1979: 98) ได้กล่าวว่าควรเป็นข้อสอบที่วัดเรื่องเดียวกันกับแบบสอบที่ต้องการศึกษา และมีเพียง ๑๐ ข้อ ก็เพียงพอแล้ว ที่กล่าวมานี้ยังขาดผลเชิงประจักษ์ในการสนับสนุน ต่อมาปี ค.ศ. ๑๘๘๕ บูดেসคู (Budescu 1985: 13-20) ได้ศึกษาวิเคราะห์เชิงทฤษฎี ส่วนคลื่น และ โคลเลน (Klein and Kolen 1985) ได้ศึกษาเชิงประจักษ์ ซึ่งจะเสนอต่อไป

๑. การวิเคราะห์เชิงทฤษฎีของบูเดสคู (๑๘๘๒, ๑๓ - ๒๐) พบว่า ตัวแปรสำคัญที่ส่งผลถึงประสิทธิภาพของการเทียบมาตรฐาน คือ ค่าสัมประสิทธิ์สหสัมพันธ์ระหว่างแบบสอบรวมกับแบบสอบสองชุดที่ต้องการเทียบมาตรฐาน และในเชิงทฤษฎีวัดนอกรีต (classical test theory) ตัวแปรความยาวของแบบสอบรวมมีความสัมพันธ์กับความเที่ยง อธิบายได้ดังต่อไปนี้

กลุ่ม  $\alpha$  ทำแบบสอบชุด X กลุ่ม B ทำแบบสอบชุด Y และทั้งสองกลุ่มทำแบบสอบรวม V เหมือนกัน ค่าเฉลี่ยกับส่วนเบี่ยงเบนมาตรฐานของแบบสอบทั้งสองชุด และของแบบสอบรวมประมาณได้ความถ่วงกับกลุ่มดังนี้ ให้  $M_{Y\alpha}$   $S_{X\alpha}$   $M_{V\alpha}$   $S_{V\alpha}$  เป็นค่าสถิติที่ได้จากกลุ่ม  $\alpha$  และ  $M_{YB}$   $S_{YB}$   $M_{VB}$   $S_{VB}$  เป็นค่าสถิติที่ได้จากกลุ่ม B ในลักษณะเดียวกัน ค่าเฉลี่ยและส่วนเบี่ยงเบนมาตรฐานของแบบสอบรวม V ของกลุ่มตัวอย่างรวม T ( $\alpha + \beta$ ) ได้โดยตรงจากผลการสอบ ในกรณีที่กลุ่มผู้สอบไม่ได้มาจากการสุ่ม วิธีการประมาณค่าจึงใช้วิธีของทักเกอร์



(Tucker equating method) แต่ถ้าวการจักดำเนินการสอบได้ใช้กลุ่มผู้สอบอย่างสุ่ม การประมาณค่าจะใช้วิธีเป็นไปก็สูงสุด (Maximum likelihood) สูตรที่ใช้มีความแตกต่างกันเล็กน้อย สำคัญของวิธีการใช้แบบสอบรวมอยู่ที่การประมาณค่าเฉลี่ยและความแปรปรวนของชุด X และ Y ในกลุ่มตัวอย่างรวม ซึ่งหาได้ด้วยการใช้ค่าต่าง ๆ ที่รู้แล้ว คือ ค่าสหสัมพันธ์ระหว่างคะแนนของแบบสอบรวม V กับคะแนนจากแบบสอบ X และ Y และค่าความแตกต่างระหว่างคะแนนในแบบสอบรวม V ของกลุ่มตัวอย่างทั้ง ๒ ดังสูตร

$$\hat{M}_{Xt} = M_{X\alpha} + \frac{r_{XV\alpha} S_{X\alpha}}{S_{V\alpha}} (M_{Vt} - M_{V\alpha}) \quad \text{-----} \quad (๑๗)$$

และ 
$$\hat{S}_{Xt}^2 = S_{X\alpha}^2 + \frac{r_{XV\alpha}^2 S_{X\alpha}^2}{S_{V\alpha}^2} (S_{Vt}^2 - S_{V\alpha}^2) \quad \text{-----} \quad (๑๘)$$

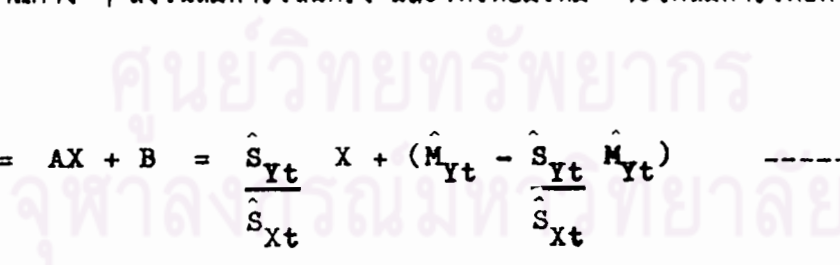
ทำนองเดียวกัน สูตรที่ใช้สำหรับแบบสอบชุด Y มีดังนี้

$$\hat{M}_{Yt} = M_{Y\beta} + \frac{r_{YV\beta} S_{Y\beta}}{S_{V\beta}} (M_{Vt} - M_{V\beta}) \quad \text{-----} \quad (๑๙)$$

$$\hat{S}_{Yt}^2 = S_{Y\beta}^2 + \frac{r_{YV\beta}^2 S_{Y\beta}^2}{S_{V\beta}^2} (S_{Vt}^2 - S_{V\beta}^2) \quad \text{-----} \quad (๒๐)$$

เมื่อแทนค่าประมาณต่าง ๆ ลงในสมการเส้นตรง และจัดเทอมใหม่ จะได้สมการเพื่อการเทียบมาตรฐานดังนี้

$$Y = AX + B = \frac{\hat{S}_{Yt}}{\hat{S}_{Xt}} X + (M_{Yt} - \frac{\hat{S}_{Yt}}{\hat{S}_{Xt}} M_{Xt}) \quad \text{-----} \quad (๒๑)$$



จากสมการข้างต้นเห็นได้ชัดแล้วว่า ถ้าสัมประสิทธิ์สหสัมพันธ์ระหว่างแบบสอบถามกับชุด X หรือ Y เป็นศูนย์ แบบสอบถามที่นำมาใช้ในงานนี้ไม่มีความหมายเลย ครั้นถ้าสัมประสิทธิ์สหสัมพันธ์ดังกล่าวมีค่าเพิ่มสูงขึ้น ๆ จะมีผลต่อการเพิ่มคุณภาพของตัวประมาณประชากรรวมอย่างสูงขึ้นตามกัน (monotonic increasing) การที่จะให้ได้ว่าซึ่งความสัมพันธ์สูง ๆ จะก่อให้เกิดการสร้างและจัดชุดแบบสอบถามอย่างพิถีพิถันที่ให้มีลักษณะคู่ขนานกับแบบสอบถามทั้งสองชุด และที่ก่อให้เกิดการเน้นในงานของบุคค (๑๘๘๕: ๑๕) คือ ความยาวของแบบสอบถาม ซึ่งการวิเคราะห์ต่อไปนี้จะแสดงให้เห็นความสำคัญของความยาวของแบบสอบถาม ทั้งนี้

กลุ่ม  $\alpha$  ทำชุด X และ V ให้ได้คะแนนรวม Z

ดังนั้น Z ของแต่ละคน คือ  $Z = X + V$  ----- (๒๒)

ให้ความเที่ยงของแบบสอบถาม (X + V) เป็น  $r_{ZZ}$  แสดงในเทอมของพารามิเตอร์ของ X และ V ดังนี้ (Horst 1968: 281 cited by Budescu 1985: 16)

$$r_{ZZ} = \frac{(S_x^2 r_{XX} + S_v^2 r_{VV} + 2 r_{XV} S_x S_v)}{(S_x^2 + S_v^2 + 2 r_{XV} S_x S_v)} \text{ ---- (๒๓)}$$

เมื่อ  $r_{XX}$  และ  $r_{VV}$  เป็นค่าประมาณความเที่ยงของแบบสอบถาม X และแบบสอบถาม V

ให้  $p$  ( $0 < p < 1$ ) แทนสัดส่วนของข้อสอบชุด X ที่อยู่ใน Z และ  $q = (1 - p)$  เป็นสัดส่วนของแบบสอบถาม V ในแบบสอบถาม Z ถ้า ๒ แบบสอบ X และ V ต่างเป็นคู่ขนานกัน ความเที่ยงเขียนในรูปของฟังก์ชันของ  $r_{ZZ}$  โดยอาศัยสูตร Spearman - Brown ดังนี้



$$r_{xx} = \frac{p r_{ZZ}}{(1 - q r_{ZZ})} \quad \text{และ}$$

$$r_{vv} = \frac{q r_{ZZ}}{(1 - p r_{ZZ})} \quad \text{-----} \quad (๒๔)$$

ในทำนองเดียวกัน ถ้า  $s_z^2$  เป็นความแปรปรวนของแบบสอบถาม z ความแปรปรวนของส่วนประกอบสามารถเขียนในรูปของฟังก์ชันของความยาวสัมพันธ์ได้ดังนี้

$$s_x^2 = s_z^2 p (1 - q r_{ZZ}) \quad \text{และ}$$

$$s_v^2 = s_z^2 q (1 - p r_{ZZ}) \quad \text{-----} \quad (๒๕)$$

แทนสมการ (๒๔) และ (๒๕) ใน (๒๓) จะพบว่า  $r_{xv}$  เป็นฟังก์ชันของความเที่ยงของแบบสอบถามรวม z และความยาวของส่วนประกอบ x และ v ดังนี้

$$r_{xv} = r_{ZZ} \left[ \frac{p q}{(1 - p r_{ZZ})(1 - q r_{ZZ})} \right]^{1/2} \quad \text{-----} \quad (๒๖)$$

นิพจน์ในสูตร (๒๖) เป็นอิสระจากค่าเฉลี่ยและความแปรปรวนรวมของแบบสอบถามที่เกี่ยวข้องทั้ง ๒ ดังนั้น สามารถเอาไปใช้กับกลุ่มผู้สอบอีกกลุ่ม (B) ได้เหมือนกัน สรุปจากการวิเคราะห์ได้ว่า โขจรธรรมชาติแล้วสหสัมพันธ์ระหว่างแบบสอบถาม x กับแบบสอบถามรวมเป็นฟังก์ชันเพิ่มขึ้นตามความเที่ยงของแบบสอบถามรวม ( $r_{ZZ}$ ) เมื่อความเที่ยงของแบบสอบถามรวม ( $r_{ZZ}$ ) คงที่ ค่าสหสัมพันธ์ระหว่างแบบสอบถาม x กับแบบสอบถามรวม ( $r_{xv}$ ) มีค่าสูงสุดที่ความยาวแบบสอบถามทั้งสองเท่ากัน ( $p = q = .5$ ) ถ้า  $r_{xv}$  มีลักษณะสมมาตรและเพิ่มขึ้นตามผลคูณของ p และ q จากจุดสำคัญนี้ทำให้ได้ค่าอ้างอิง คือ ประสิทธิภาพสูงสุดของการเทียบมาตราเมื่อจัดสัดส่วนของแบบสอบถามที่ต้องการเทียบกับแบบสอบถามให้มีสัดส่วนเท่ากัน ในกรณีเช่นนี้จะได้สหสัมพันธ์ระหว่างแบบสอบถามที่ต้องการเทียบกับแบบสอบถามรวม ( $r_{xv}$ ) มีค่าเท่ากับ  $r_{ZZ} / (2 - r_{ZZ})$  (๑๙๕๕: ๑๖ - ๑๗)

จากการวิเคราะห์ที่กล่าวมาข้างต้น บูเคส (๑๙๕๕: ๑๗) ให้นำมาสร้างเป็น  
ดัชนีของประสิทธิภาพสัมพัทธ์ของการเทียบมาตรา เรียกว่า ความพร่องสัมพัทธ์ (relative  
deficiency) ตามนิยามของสมการ (๒๗) ดังนี้

$$R.D(q) = 1 - \frac{r_{xv}(q)}{r_{xv}(0.5)} = 1 - \left[ \frac{(2 - r_{zz})^2 pq}{(1 - p r_{zz})(1 - q r_{zz})} \right]^{1/2} \quad \text{-- (๒๗)}$$

มาตรการนี้จะแสดงให้เห็นถึงการเพิ่มประสิทธิภาพในเชิงเปรียบเทียบกับค่าอ้างอิงสูงสุด เมื่อ  
 $p = q = .5$  ดังนั้น กระบวนการเทียบมาตราใด ๆ ที่มีค่าความพร่องสัมพัทธ์ต่ำกว่าหรือเท่ากับ  
๐.๑๐ จัดว่าเป็นการเทียบมาตราที่สามารถให้ผลที่น่าพอใจ (๑๙๕๕: ๑๗)

ในกรณีสูตรสำเร็จของแองกอสที่ใช่ ๒๐% ถ้าจะให้มีประสิทธิภาพสัมพัทธ์ในระดับที่  
น่าพอใจแล้ว จะต้องเป็นแบบสอบที่มีค่าความเที่ยงสูงกว่า ๐.๘๖๖ (๑๙๕๕: ๑๗)

กล่าวโดยสรุป การวิเคราะห์ของบูเคสเป็นการวิเคราะห์เชิงทฤษฎีที่ชี้ให้เห็นความ  
สัมพันธ์ของความยาวแบบสอบร่วมกับประสิทธิภาพของการเทียบมาตรา ในเชิงปฏิบัติไม่สามารถให้  
ข้อกำหนดทั่วไปที่จะดีปฏิบัติกัน ทั้งนี้เพราะการปฏิบัติจะต้องคำนึงถึงตัวแปรแวดล้อมอื่น ๆ อีก เช่น  
เวลา การลงทุน และข้อจำกัดในทางปฏิบัติอื่น ๆ อีก แต่สิ่งที่ควรพิจารณา คือ ให้ประสิทธิภาพของ  
การเทียบมาตรานั้น ๆ อยู่ในระดับของการยอมรับเฉพาะกรณี (๑๙๕๕: ๑๘)

๒. การศึกษาเชิงประจักษ์ของคลิน และ โคลเลน (๑๙๕๕) คลิน และ โคลเลน ได้  
เลือกศึกษากรณีการเทียบมาตราที่กลุ่มผู้สอบไม่ได้มาจากการสุ่ม ซึ่งวิธีการใช้แบบสอบร่วมเป็นสิ่ง  
จำเป็นอย่างยิ่งที่ให้พื้นฐานเพื่อการประมาณความแตกต่างของกลุ่มผู้สอบสองกลุ่ม ประเด็นของแบบ-  
สอบร่วมที่เราทำการศึกษา คือ เรื่องของความยาวของแบบสอบร่วม โดยได้ทำการตรวจสอบผล  
กระทบของความยาวต่อการเทียบมาตรา ดังนี้

ข้อมูลเป็นผลการสอบเพื่อสอบรับประกาศนียบัตรวิชาชีพ แบบสอบเป็นชนิดเลือกตอบ  
มี ๒๕๐ ข้อ ได้ทำการสอบเมื่อปี ค.ศ. ๑๙๕๔ ออกแบบการวิจัยด้วยการแบ่งแบบสอบเป็น ๒ ชุด และ  
แบบสอบร่วม ความยาวต่าง ๆ กัน ๕ ขนาด ในแต่ละขนาดประกอบด้วยข้อสอบที่มีความสมมูลย์ทั้งใน



ก้านเนื้อหา ความยากเฉลี่ย และอำนาจจำแนกเฉลี่ยของแบบสอบทั้งฉบับ แบบสอบรวมที่ใช้จึงมีขนาด ๒๐ ข้อ ๔๐ ข้อ ๖๐ ข้อ ๘๐ ข้อ และ ๑๐๐ ข้อ แต่ละแบบสอบรวมจึงเป็นชั้นเซทของฉบับที่ยาวกว่า และทุกฉบับต่างก็เป็นชั้นเซทของจำนวน ๒๕๐ ข้อ

การจัดเตรียมกลุ่มตัวอย่างผู้สอบ ผู้วิจัยมีจุดประสงค์จะศึกษาผลการเทียบมาตราแบบทักเกอร์ (Tucker equating) ซึ่งเป็นรูปแบบเชิงเส้นตรง (Angoff 1971: 579-583 cited by Klein and Kolen 1985: 4) กับกลุ่มผู้สอบที่มีความสามารถคล้ายคลึงกัน และกับกลุ่มผู้สอบที่มีความสามารถแตกต่างกัน ผู้วิจัยจึงได้แบ่งคนสอบจำนวน ๑,๗๕๕ คน ออกเป็น ๒ กลุ่ม ไม่ซ้ำซ้อนกันตามเกณฑ์ของความสามารถคล้ายคลึงและแตกต่างกันพื้นฐาน (meadin)

ผลการเทียบมาตราที่เกิดกับการใช้แบบสอบรวมขนาดความยาว ๕ ขนาด กับกลุ่มผู้สอบ ๒ ลักษณะ รวมมีผลการเทียบมาตรา ๑๐ ชุด จากนั้นได้คำนวณค่าความคลาดเคลื่อน (root-mean-squared error: RMSE) ของแต่ละกรณีด้วยสูตร

$$\widehat{RMSE} = \left[ \frac{\sum_i N_i (x'_i - x_i)^2}{\sum_i x_i x'_i} \right]^{1/2} \quad \text{--- (๒๔)}$$

- เมื่อ  $x_i$  เป็นคะแนนแถวที่  $i$  ของแบบสอบทั้งหมด
- $N_i$  เป็นจำนวนผู้สอบที่ได้คะแนน  $i$  จากแบบสอบทั้งหมด
- $x'_i$  เป็นคะแนนสมมูลของ  $x_i$  ที่ประมาณด้วยวิธีการของทักเกอร์ในรูปแบบเชิงเส้นตรง

นอกจากนี้ยังได้ประมาณค่าเฉลี่ยของความคลาดเคลื่อนของการเทียบมาตรา ซึ่งเป็นตัวที่ทำให้เกิด  $\widehat{RMSE}$  โดยใช้สมการ

$$\widehat{BIAS} = \bar{x}' - \bar{x} \quad \text{--- (๒๕)}$$

- เมื่อ  $\bar{x}'$  เป็นค่าเฉลี่ยของค่าประมาณของคะแนนดิบสมมูลย์ และ
- $\bar{x}$  เป็นค่าเฉลี่ยของคะแนนดิบ

เปรียบเทียบค่าดัชนีความคลาดเคลื่อน (RMSE) กับส่วนเบี่ยงเบนมาตรฐานของแบบสอบ ซึ่งเท่ากับ ๒๐.๓๒ พบว่า ความคลาดเคลื่อนของการเทียบมาตราสำหรับแบบสอบรวมที่มีขนาด ๒๐ ข้อ และ ๕๐ ข้อ ของกลุ่มความสามารถแตกต่างกันมีค่าสูง

คลินและโคเลนอภิปรายไว้ว่า กรณีแบบสอบที่นำมาศึกษาได้จัดกระทำให้เป็นแบบสอบเทียบมาตรากับแบบสอบรวม แล้วทำการเทียบมาตรากับแบบสอบเดิม เมื่อแบบสอบรวมมีความสมนัยกับแบบสอบทั้งในค่าความครอบคลุมของเนื้อหา ความยาวเฉลี่ย และค่าอำนาจจำแนกเฉลี่ย ผลการวิจัยไม่สามารถสรุปว่า ความยาวของแบบสอบรวมมีความสัมพันธ์อย่างมีระบบกับค่าประมาณความคลาดเคลื่อนมาตรฐานของการเทียบมาตรา ความแม่นยำของการเทียบมาตราที่ใช้แบบสอบรวมขนาด ๒๐ ๕๐ ๖๐ และ ๑๐๐ ในกลุ่มผู้สอบที่คล้ายกัน สำหรับกลุ่มผู้สอบที่มีความสามารถแตกต่างกัน พบว่า ค่าประมาณความคลาดเคลื่อนพุ่งสูงขึ้นอย่างมากเมื่อจำนวนแบบสอบรวมลดลงจาก ๒๐ ข้อ ลงเหลือ ๕๐ ข้อ และจาก ๕๐ ข้อ ลงเหลือ ๒๐ ข้อ

กล่าวโดยสรุป เมื่อแบบสอบที่นำมาเทียบมาตรา มีความคล้ายคลึงกันมาก เช่น ให้กรณีที่ทำการศึกษา คือ เป็นข้อสอบชุดเดียวกัน และกลุ่มผู้สอบมีระดับความสามารถคล้ายคลึงกัน การเพิ่มความยาวให้กับแบบสอบรวมไม่ได้เพิ่มความแม่นยำของการเทียบมาตราอย่างชัดเจน แต่สำหรับกลุ่มผู้สอบที่มีความสามารถแตกต่างกันแล้ว การเพิ่มความยาวของแบบสอบรวมมีผลต่อการเพิ่มความแม่นยำของการเทียบมาตราอย่างชัดเจน จากกรณีการศึกษาครั้งนี้ พบว่า เมื่อแบบสอบรวมมีความยาวน้อยกว่าร้อยละ ๒๕ ของแบบสอบทั้งหมดแล้ว ผลของความคลาดเคลื่อนของการเทียบมาตราจะอยู่ในระดับที่มากจนไม่เป็นที่ยอมรับได้ และที่คงย้ำอีกครั้ง คือ เป็นกรณีที่แบบสอบรวมมีความสมนัยกับแบบสอบทั้งหมดเป็นอย่างมากในค่าความครอบคลุมของเนื้อหา ความยาว และอำนาจจำแนก การศึกษานี้ไม่มีหลักฐานที่จะชี้ว่าความยาวของแบบสอบรวมเป็นตัวแปรสำคัญเพียงตัวแปรเดียว

จากผลการศึกษาของบูเคส (Budescu) คลิน และ โคลน (Klein and Kolen) ที่เกี่ยวกับความยาวของแบบสอบรวมนี้ พบว่า ตัวแปรที่มีความสัมพันธ์และร่วมส่งผลกระทบต่อประสิทธิภาพของการเทียบมาตราที่สำคัญ คือ คุณภาพของแบบสอบรวม และจากผลการศึกษาของ

คลีน และ โคลเนน ยังพบประเด็นของความไม่คงที่ของการเทียบมาตรา เมื่อเปลี่ยนลักษณะของกลุ่มผู้สอบอีก

สำหรับการศึกษาที่กำลังดำเนินการอยู่ในครั้งนี้ ผู้วิจัยได้จัดทำเกี่ยวกับตัวแปรความยาวของแบบสอบรวมใหม่ต่างกัน ๓ ระดับ คือ คิดเป็นอัตราร้อยละเทียบกับแบบสอบที่ทำการเทียบมาตราเป็น ๒๐ ๔๐ และ ๖๐ โดยมีกลุ่มผู้สอบที่เป็นไปคาบสมรรมชาติของการสอบจริง ๒ สถานการณ์ คือ สถานการณ์การสอบแข่งขันเพื่อคัดเลือกเข้าทำงาน และสถานการณ์การสอบเพื่อสัมฤทธิ์ผลรายวิชาในระดับปริญญาตรี โดยผู้วิจัยได้ตั้งสมมุติฐานจากหลักของแองกอฟ ทวีเคราะห์ของบูเทส และผลการวิจัยส่วนหนึ่งของคลีน และ โคลเนนว่า แบบสอบรวมที่ยาวกว่าย่อมได้ผลการเทียบมาตราที่ดีกว่า

### ความคลาดเคลื่อนของการเทียบมาตรา

การประเมินความคลาดเคลื่อนของการเทียบมาตราที่สมบูรณ์ คือ การวิเคราะห์แหล่งความลำเอียง (bias) และความแปรผัน (variability) แล้วนำมาเสนอให้เห็นเป็นภาพเดียวกัน เพื่อช่วยในการตัดสินใจว่า การเทียบมาตราด้วยวิธีนั้น ๆ ควรทำหรือไม่อย่างไร (Braun and Holland 1982: 32) รายละเอียดที่จะเสนอต่อไปนี้ ได้แบ่งเป็น ๒ ตอน คือ เรื่องของความแปรปรวน ซึ่งมีทฤษฎีค่อนข้างหนักแน่นแล้ว และเรื่องของความลำเอียง ซึ่งเป็นบทวิเคราะห์ที่ยังต้องการวิจัยสนับสนุนอีกมาก

#### ๑. ความแปรผันเชิงสุ่ม และความคลาดเคลื่อนมาตรฐานของการเทียบมาตรา

วิธีการเทียบมาตราทุกวิธี ไม่ว่าจะอยู่ในรูปแบบการเทียบมาตราใด เมื่อกลุ่มตัวอย่างผู้สอบเป็นกลุ่มสุ่มจากประชากรเดียว หรือหลายประชากร ย่อมมีความแปรผันเชิงสุ่มเกิดขึ้น จึงนิยมใช้เทคนิคการประมาณค่าความคลาดเคลื่อนเชิงสุ่ม (estimated sampling errors) กับ การเทียบมาตราต่าง ๆ ความคลาดเคลื่อนเชิงสุ่มนี้จึงมีสมมุติฐานเบื้องต้นว่า ตัวอย่างมาจากการสุ่ม และใช้ความคลาดเคลื่อนมาตรฐานของการเทียบมาตรา (Standard error of equating ย่อว่า SEE) เป็นมาตรการวัดความแปรผันประเภทนี้

๑.๑ ความคลาดเคลื่อนมาตรฐานของการเทียบมาตราในรูปแบบเชิงเส้นตรง

ลอร์ด (Lord 1950 cited by Angoff 1984: 96) ได้สร้างสูตรสำหรับคำนวณหาความคลาดเคลื่อนมาตรฐานของรูปแบบเชิงเส้นตรงหลายสูตร ตามแบบแผนการรวบรวมข้อมูล ซึ่งได้มีการนำไปใช้อย่างกว้างขวางมานาน เช่น การเทียบมาตรากับแบบสอบ SAT ตั้งแต่ทศวรรษที่ ๑๙๕๐ (Donlon and Angoff 1971: 33) บรุน และ ฮอลแลนด์ (Braun and Holland 1982: 33-35) ได้นำมาเสนอให้อยู่ในรูปแบบทั่วไป ดังนี้

ให้  $n_X$  คือ จำนวนคนสอบที่ได้จากการสุ่มจากประชากร  $P$  และทำแบบสอบ  $X$

$n_Y$  คือ จำนวนคนสอบที่ได้จากการสุ่มจากประชากร  $P$  เช่นกัน แต่ทำแบบสอบ  $Y$

$n_X$  และ  $n_Y$  ไม่จำเป็นต้องเท่ากัน

$x_i$  คือ คะแนนของคนที่  $i$  ในกลุ่มที่ทำแบบสอบ  $X$

$y_j$  คือ คะแนนของคนที่  $j$  ในกลุ่มที่ทำแบบสอบ  $Y$

$\{X_i\}$  และ  $\{Y_j\}$  ต่างมีการแจกแจงอย่างอิสระทั่วทั้งชั้น  $F(x)$

และ  $G(y)$  ตามลำดับ ตัวอย่าง  $X$  และ  $Y$  เป็นตัวอย่างอิสระจากกัน สัจสุลักษณะโมเมนต์ที่ ๑ และ ๒ ของการแจกแจงใน  $X$  และ  $Y$  มีดังนี้

$$\mu_X = E(X) \qquad \mu_Y = E(Y) \qquad \text{--- (๓๐)}$$

$$\sigma_X^2 = \text{Var}(X) \qquad \sigma_Y^2 = \text{Var}(Y) \qquad \text{--- (๓๑)}$$

ฟังก์ชันเทียบมาตราของ  $Y$  ไปยัง  $X$  ในรูปเชิงเส้นตรง คือ

$$x^*(y) = \mu_X + \frac{\sigma_X}{\sigma_Y} (y - \mu_Y)$$



ซึ่ง  $y$  จะแปรผันไปตลอดคะแนนดิบของ  $Y$  เขียนให้อยู่ในรูปของค่าประมาณประชากร ดังนี้

$$\hat{x}^*(y) = \bar{X} + \frac{s_x}{s_y} (y - \bar{Y})$$

เมื่อ  $\bar{X}$  และ  $s_x^2$  เป็นค่าเฉลี่ยและความแปรปรวนของ  $\{X_i\}$   
ความแปรปรวนแอสิมโทติก (Asymptotic variance) ของ  $\hat{x}^*(y)$  ได้มาจากทฤษฎีบทดังต่อไปนี้

$$\begin{aligned} \text{Var}(\hat{X}^*(y)) &= \sigma_X^2 \left\{ \frac{1}{n_X} + \frac{1}{n_Y} + \left( \frac{\text{Skew}(X)}{n_X} + \frac{\text{Skew}(Y)}{n_Y} \right) Z(y) \right. \\ &+ \left. \left( \frac{2 + \text{kurt}(X)}{4n_X} + \frac{2 + \text{kurt}(Y)}{4n_Y} \right) Z^2(y) \right\} \\ &+ (\text{higher order terms}) \quad \text{----- (๓๒)} \end{aligned}$$

$$\begin{aligned} \text{เมื่อ} \quad \text{Skew}(X) &= \mu_{3X} / \sigma_X^3 \\ \text{kurt}(X) &= (\mu_{4X} / \sigma_X^4) - 3 \\ \mu_{3X} &= E(X - \mu_X)^3 \\ \mu_{4X} &= E(X - \mu_X)^4 \end{aligned}$$

สำหรับเทอมต่าง ๆ ของ  $Y$  หาได้ในทำนองเดียวกัน และ

$$Z(y) = (y - \mu_Y) / \sigma_Y$$

สำหรับเทอมต่าง ๆ ที่มีกำลังสูงขึ้นเป็นกำลังที่เป็น + ของ  $n_X^{-1}$  และ  $n_Y^{-1}$   
จากทฤษฎีบทนี้ ได้อาศัยบทแทรก ๒ บท ช่วยทำให้สูตรชัดเจนขึ้น

บทแทรกที่ ๑. ถ้า X และ Y มีการแจกแจงปกติแล้ว

$$\text{Skew}(X) = 0 \quad \text{Skew}(Y) = 0 \quad \text{และ}$$

$$\text{kurt}(X) = 0 \quad \text{kurt}(Y) = 0$$

ดังนั้น  $\text{Var}(x_{(y)}^*) = (\sigma_x^2 / n_h)(2 + Z^2(y)) + (\text{higher order terms})$

เมื่อ  $n_h = [\frac{1}{2}(n_X^{-1} + n_Y^{-1})]^{-1}$  คือ ตัวกลางฮาร์โมนิกของ  $n_X$  และ  $n_Y$

บทแทรกที่ ๒. ถ้า X และ Y มีการแจกแจงเป็นปกติ และ

$$n_X = n_Y = N/2 \quad \text{แล้ว}$$

$$\text{Var}(x_{(y)}^*) = \frac{2\sigma_x^2}{N} (2 + Z^2(y)) + (\text{higher order terms}) \quad \text{--- (๓๓)}$$

โดยทั่วไปสภาพข้อมูลมักไม่เป็นไปตามบทแทรกที่ ๑ คือ การแจกแจงของคะแนนย่อมมีลักษณะเบ้ แต่เหมาะกับการจัด  $n_X$  ไม่เท่ากับ  $n_Y$  ความคลาดเคลื่อนมาตรฐานของการเทียบมาตรา คือ ปรากฏที่ ๒ ของ  $\text{Var}(x_{(y)}^*)$  หรือ  $\text{SEE}_{x^*}^2$

สูตรสำหรับหาความคลาดเคลื่อนมาตรฐาน เมื่อข้อมูลแตกต่างจากที่กล่าวมานี้ จะแตกต่างไปตามเงื่อนไข ซึ่งสรุปได้ดังนี้

(๑) เมื่อกลุ่มสุ่มสองกลุ่มทำแบบสอบถามละชุด (เหมือนที่กล่าวมาข้างต้น)

หาค่าประมาณความคลาดเคลื่อนจากรากที่สองของสูตรต่อไปนี้ (Lord 1950 cited by Angoff 1984: 97)

$$\text{SEE}_{x^*}^2 = \frac{2s_x^2}{N_t} (Z_y^2 + 2) \quad \text{--- (๓๔)}$$

เมื่อ  $N_t = n_X + n_Y$

$$Z_y = (X - n_Y) / s_Y$$

(๒) เมื่อกลุ่มผู้สอบสองกลุ่ม แต่ละกลุ่มทำแบบสอบทั้งสองชุด แต่จัดค่าในการทดสอบให้ผู้สอบรับการทดสอบ X หรือ Y ก่อนหลังในลักษณะสลับ สูตรที่ใช้ คือ (Lord 1950 cited by Angoff 1984: 103)

$$SEE_{X^*}^2 = \frac{\sigma_X^2 (1 - r_{XY}) Z_Y^2 (1 + r_{XY}) + 2}{N_t} \quad \text{--- (๓๕)}$$

เมื่อ  $r_{XY}$  คือ สัมประสิทธิ์สหสัมพันธ์ของคะแนนชุด X และ Y พิจารณาจากสูตรนี้จะเห็นว่า ถ้า X และ Y มีความสัมพันธ์สูง SEE จะมีค่าต่ำลง เช่น ถ้า  $r_{XY}$  มีค่า .๘๐ ค่า  $SEE_{X^*}^2$  ที่  $Z_Y = 0$  จะมีค่าเป็น ๐/๑๐ ของ  $SEE_{X^*}^2$  ใน (๑) ซึ่งอธิบายได้ในแง่ของจำนวนผู้สอบว่า กรณีที่เทียบมาคร่าวๆ (๑) จะต้องใช้จำนวนผู้สอบมากกว่า (๒) เป็น ๑๐ จึงจะได้ผลที่มีความคมชัดเท่ากัน

(๓) เมื่อกลุ่มผู้สอบสองกลุ่มทำแบบสอบกลุ่มละชุด (X หรือ Y) และทำแบบสอบรวมเหมือนกัน (Angoff 1984: 104) ค่าประมาณความคลาดเคลื่อนหาได้จากสูตร (Lord 1950 cited by Angoff 1984: 106) ดังนี้

$$SEE_{X^*}^2 = \frac{2\sigma_X^2 (1 - \hat{r}^2)(1 + \hat{r}^2) Z_Y^2 + 2}{N_t} \quad \text{--- (๓๖)}$$

โดยมีข้อตกลงว่า

$$\hat{r} = \frac{b_{xv}\alpha \hat{\sigma}_v}{\hat{\sigma}_X} = \frac{b_{yv}\beta \hat{\sigma}_v}{\hat{\sigma}_Y}$$

เมื่อ  $b_{xv}\alpha$  และ  $b_{yv}\beta$  คือ สัมประสิทธิ์ถดถอย และ  $\hat{\sigma}_v$   $\hat{\sigma}_X$  และ  $\hat{\sigma}_Y$  คือ ค่าประมาณส่วนเบี่ยงเบนมาตรฐานของ V X และ Y ตามลำดับ

พิจารณาจากสูตร ถ้า  $\hat{r} = 0$  ค่าประมาณความแปรปรวนคลาดเคลื่อน จะมีขนาดเท่ากับหาได้จาก (๑) ถ้า  $\hat{r} = .7$  ค่าที่ได้เป็นเพียงครึ่งหนึ่งที่  $Z_Y = 0$

๐.๒ ความคลาดเคลื่อนมาตรฐานของการเทียบมารูปแบบฮิวเปอร์เซนไคล์

การศึกษาความคลาดเคลื่อนมาตรฐานของการเทียบมารูปแบบฮิวเปอร์เซนไคล์ เป็นงานที่ค่อนข้างใหม่ของลอร์ด (Lord 1982) ในการพัฒนาสูตรใ้กำหนดให้กลุ่มตัวอย่าง ๒ กลุ่ม ทำแบบสอบถามอิสระ และสร้างข้อตกลงให้คะแนนสอบเป็นตัวแทนของ การหาความแปรปรวนแอสิมโทติกมีดังนี้

ให้  $F(X)$  และ  $G(Y)$  แทนการแจกแจงความถี่สะสมของคะแนน  $X$  และ  $Y$  ในประชากร  $P$  ให้กลุ่มตัวอย่าง  $N_2$  ทำแบบสอบถาม  $Y$  แล้วหาสัดส่วน  $q$  ที่มีคะแนนต่ำกว่าค่าคงที่  $y_0$  ที่เลือกมาตามความสนใจ ในกลุ่มตัวอย่าง  $N_1$  ทำแบบสอบถาม  $X$  ให้  $N_{1q}$  เป็นสถิติอันดับในกลุ่มตัวอย่างนี้ โดยที่  $x^*$  สมมูลกับ  $y_0$  ดังนั้น แต่ละค่าของ  $y_0$  ที่สนใจ จะพบความแปรปรวนเชิงสัมประสิทธิ์ของ  $x^*$  สำหรับแต่ละค่าของ  $y_0$  ถ้าให้  $q$  คงที่  $x^*$  จะมีการแจกแจงอย่างปกติด้วยค่าเฉลี่ย  $\mu_{x^*|q}$  มีความสัมพันธ์เป็น  $F(\mu_{x^*|q}) = q$  ----- (๓๗)

และความแปรปรวน  $\sigma_{x^*|q}^2 = pq / N_1 [g(\mu_{x^*|q})]^2$  ----- (๓๘)

เมื่อ  $p \equiv 1 - q$

และ  $f(x)$  คือ ความน่าจะเป็นของความหนาแน่นที่  $x$

แต่อย่างไรก็ตาม  $q$  เป็นสัดส่วนตัวอย่าง จึงต้องหา  $\text{Var } x^*$  เมื่อ  $q$  ไม่คงที่จาก

$$\text{Var } x^* = \text{Var}(\mu_{x^*|q}) + E(\sigma_{x^*|q}^2) \text{ ----- (๓๙)}$$

เมื่อ  $\text{Var}$  และ  $E$  แทนความแปรปรวนและค่าคาดหวังของ  $x^*$  สำหรับทุกค่าของ  $q$  ที่  $y_0$  คงที่

$$\frac{d}{dq} \mu_{x^*|q} = \frac{1}{g(\mu_{x^*|q})}$$

เนื่องจาก  $\text{Var } q = PQ / N_2$

เพราะฉะนั้น  $\text{Var}(\mu_{x^*|q}) = PQ / N_2 g_0^2$  ----- (๔๐)





เมื่อ  $q$  เป็นค่าประชากรของ  $q$  ซึ่งนิยามไว้ว่า

$$Q \equiv G(y_0) \quad \text{-----} \quad (๔๑)$$

$$P \equiv 1 - Q \quad \text{และ}$$

$$g_0 \equiv g(x_0)$$

$$G(x_0) \equiv Q \quad \text{-----} \quad (๔๒)$$

สมการที่ (๔๐) ขึ้นอยู่กับ  $n_2$  ไม่ใช่  $n_1$  เพราะว่า ความผันแปรใน  
 เกิดขึ้นจากความผันแปรของ  $q$  ซึ่งเป็นสถิติที่หนึ่งที่มีความแปรปรวน  $PQ / n_2$  เพื่อ  
 ประเมิน  $E \sigma_{x^*|q}^2$  ให้เขียน (๓๔) โดยขยายอนุกรมและไม่คงสนใจเทอมที่มีกำลังสูง ๆ ดังนี้

$$\begin{aligned} E \sigma_{x^*|q}^2 &= E \frac{Pq}{N_1 g_0^2 \left(1 - \frac{g - g_0}{g_0}\right)^2} \\ &= E \frac{Pq}{N_1 g_0^2} \left(1 - 2 \frac{g - g_0}{g_0} + \dots\right) \\ &= E \frac{P - P^2}{N_1 g_0^2} = \frac{1}{N_1 g_0^2} (P - \text{Var } p - P^2) \\ &= \frac{PQ}{N_1 g_0^2} \quad \text{-----} \quad (๔๓) \end{aligned}$$

ศูนย์วิทยพัชกร  
 จุฬาลงกรณ์มหาวิทยาลัย

เมื่อ  $\sigma \equiv \sigma(\mu_{x^*|q})$  แทน (๕๐) และ (๕๓) ลงใน (๕๔) ได้สูตร  
สุดท้าย

$$\text{Var } x^* \quad \equiv \quad \frac{pq}{g_0^2} \left( \frac{1}{N_1} + \frac{1}{N_2} \right) \text{-----}(๕๕)$$

$$\text{SEE}^2 \quad \equiv \quad \frac{pq}{g_0^2} \left( \frac{1}{N_1} + \frac{1}{N_2} \right) \text{-----}(๕๕)$$

๖.๓ ความคลาดเคลื่อนมาตรฐานของทฤษฎีการตอบข้อสอบ (IRT)

ในทฤษฎีการตอบข้อสอบ (IRT) ค่าประมาณคะแนนจำนวนข้อที่ตอบถูก (number-right scores)  $\xi$  ของผู้สอบที่ทำแบบสอบ X และ  $\eta$  ของผู้สอบที่ทำแบบสอบ Y มีค่าเท่ากับฟังก์ชันคุณลักษณะที่ประเมินที่ระดับความสามารถ ( $\theta$ ) เดียวกัน ให้ความหมายของความเท่ากันในเรื่องคะแนนสมมูล ในทางปฏิบัติจึงนำความสัมพันธ์เชิงฟังก์ชันของ  $\xi$  และ  $\eta$  มาใช้เทียบมาตรฐานคะแนน X และ Y จากคะแนนที่สอบได้จริงของคนทั้งสองกลุ่ม การประมาณเพื่อให้เทียบ  $\eta$  ไปยัง  $\xi$  ตามฟังก์ชันจะต้องใช้ค่าประมาณของประชากรข้อสอบ ค่าประมาณเหล่านี้ คือ ที่มาของความคลาดเคลื่อนเชิงสุ่ม (sampling error) ในการเทียบมาตรา (Lord 1981: 2)

ลอร์ด (Lord 1981) ได้สังเคราะห์สูตรสำหรับหาความคลาดเคลื่อนมาตรฐาน (the asymptotic standard error) ของการเทียบมาตราคะแนนจริงในรูปแบบของทฤษฎีการตอบข้อสอบแบบสามพารามิเตอร์ของโลจิส ออกแบบการเทียบมาตราโดยให้กลุ่มผู้สอบ ๒ กลุ่ม ทำการสอบจากแบบสอบคนละฉบับ (X หรือ Y) และทำแบบสอบร่วมชนิดภายนอก (๘) เหมือนกัน การเทียบมาตราคะแนน  $\eta$  ไปยัง  $\xi$  ใช้วิธีการทำให้ค่าความสามารถอยู่คงที่ในระดับเดียวกัน ในรูปแบบดังกล่าวมีงาน ๒ ชั้นคอน คือ ชั้นแรกเทียบ  $\omega$  ไปยัง  $\xi$  ในผู้สอบกลุ่มแรก ชั้นที่สองเทียบ  $\eta$  ไปหา  $\omega$  ในผู้สอบกลุ่มที่สอง การทำงาน ๒ ชั้นคอนนี้ จะทำให้เทียบ  $\eta$  ไปยัง  $\xi$  ได้ การหาค่าประมาณคะแนนจริงจากฟังก์ชันใช้แทนค่าประมาณประชากรข้อสอบที่ได้จากการประมาณด้วยวิธีการความน่าจะเป็นสูงสุด (Maximum likelihood estimates) ซึ่งโดยอาศัย

โปรแกรมสำเร็จรูป LOGIST ช่วยในคำนวณค่าความ สมการคะแนนจริงมีดังนี้

$$\hat{\xi} = \sum_g \hat{P}_{g1}(\theta_1) \text{ ----- (๔๖)}$$

$$\hat{\omega} = \sum_g \hat{P}_{g2}(\theta_2) \text{ ----- (๔๗)}$$

$$\hat{\omega} = \sum_g \hat{P}_{g3}(\theta_3) \text{ ----- (๔๘)}$$

$$\hat{\eta} = \sum_g \hat{P}_{g4}(\theta_4) \text{ ----- (๔๙)}$$

สมการทั้งสี่นี้ แสดงให้เห็นว่า  $\hat{\eta}$  เป็นฟังก์ชันของค่าประมาณประชากรข้อสอบทุกค่ารวมกันตามการกำหนดค่า  $\xi$

จากการหาอนุพันธ์ของข้อสอบในแต่ละพารามิเตอร์ a b c ของข้อสอบแต่ละข้อ เมื่อกำหนดให้  $\xi$  คงที่ ทำให้หาค่าประมาณคะแนนจริงของ  $\eta$  และได้สมการความแปรปรวน  $\eta$  ดังนี้

$$\text{Var } \hat{\eta} = \sum_{p=1}^4 \sum_{g=1}^{n_p} \sum_{r=1}^3 \sum_{s=1}^3 \eta'_{rgp} \eta'_{sgp} \text{COV}(\hat{\xi}_{rgp}, \hat{\xi}_{sgp}) \text{ ----- (๕๐)}$$

ค่าที่ได้เป็นค่าประมาณอันดับค่า สัจคุณลักษณะข้อสอบมีดังนี้

p คือ จำนวนฟังก์ชันของการประมาณคะแนนจริง ซึ่งมี ๔ ฟังก์ชัน

- คือ
- ๑ หมายถึง กลุ่มที่หนึ่งทำแบบสอบ X
  - ๒ หมายถึง กลุ่มที่หนึ่งทำแบบสอบ ω
  - ๓ หมายถึง กลุ่มที่สองทำแบบสอบ ω
  - ๔ หมายถึง กลุ่มที่สองทำแบบสอบ Y

r และ s หมายถึง พารามิเตอร์ข้อสอบ มี ๓ ตัว คือ

- ๑ หมายถึง a      ๒ หมายถึง b      และ ๓ หมายถึง c

n = 1 ถึง n<sub>p</sub> หมายถึง ลำดับข้อสอบในแบบสอบแต่ละฉบับ

$\eta'_{rGP} = \partial \eta / \partial t_{rGP}$  คือ ค่าอนุพันธ์ของคะแนนจริง  $\eta$  ที่เทียบกับ พหุคูณของข้อสอบ

$Cov(\hat{t}_{rGP}, \hat{t}_{sGP})$  หมายถึง ความแปรปรวนร่วมของค่าประมาณ สูงสุดของพหุคูณ เมื่อกำหนดค่า  $\theta$  คงที่

จากการศึกษาความคลาดเคลื่อนมาตรฐานของการเทียบมาตรา IRT ของ Lord (1981) ได้ใช้ตัวอย่างคะแนนจากการสอบจริงในโครงการทดสอบ SAT ชุด VSA4 มี ข้อสอบ ๔๐ ข้อ เทียบมาตราไปสู่แบบสอบชุด XSA2 ซึ่งมี ๔๕ ข้อ คนสอบทุกคนได้ทำแบบสอบรวม ๔๐ ข้อ เหมือนกัน ผลการประมาณค่าพหุคูณข้อสอบ และความสามารถได้จากผู้สอบ ๒๒๕ คน เมื่อปี ๑๙๗๑ และ ๒๒๔๖ คน เมื่อปี ๑๙๗๕ การวิเคราะห์ความคลาดเคลื่อนได้บอกถึง การเทียบมาตรา IRT ไม่เป็นเส้นตรงเมื่อคะแนนอยู่นอกพิสัยจาก ๕๐ ถึง ๒๕๐ คะแนน ค่าความ คลาดเคลื่อนเพิ่มขึ้นขณะที่คะแนนจริงน้อยลงจนต่ำสุดที่คะแนนจริง - ๕.๕ แต่ด้านตรงกันข้าม ค่า ความคลาดเคลื่อนน้อยลงเมื่อคะแนนมาตราสูงสุดของพิสัย เหตุผลที่อธิบาย คือ คะแนนในส่วนนี้ เข้าใกล้คะแนนเต็ม (perfect score) ความคลาดเคลื่อนมีแนวโน้มใกล้ศูนย์ สำหรับคะแนน ในด้านค่าที่ต่ำกว่าระดับการเดาเฉลี่ย ใน IRT ไม่มีคำอธิบายในส่วนนี้ การหาความคลาดเคลื่อน จึงไม่คำนวณเพราะไม่มีความหมาย

๒. ความลำเอียงของการเทียบมาตรา (Bias error)

ความลำเอียง เป็นอีกองค์ประกอบหนึ่งของความคลาดเคลื่อนของการเทียบมาตรา ซึ่งจะคงพิจารณากันอย่างรอบคอบ เมื่อต้องการใช้ผลของการเทียบมาตรา เมื่อมีการสุ่มตัวอย่าง อย่างง่าย สักส่วนตัวอย่างย่อมเป็นตัวอย่างประมาณที่ไม่เอนเอียงของสัดส่วนประชากรที่สมนัยกัน ในกรณี เช่นนี้เมื่อไปใช้กับการเทียบมาตราของแบบสอบคู่ขนานสองชุด ด้วยการใช้รูปแบบการเทียบที่ตำแหน่ง เปอร์เซนต์ไคล์ (equipercentile method) ย่อมจะไม่มี ความลำเอียง แต่ถ้าในกรณีที่การ แจกแจงของแบบสอบสองชุดที่นำมาเทียบกันมีการแจกแจงคนละรูปร่างแตกต่างกันตั้งแต่ไม่แมนที่สอง ขึ้นไป การใช้ฟังก์ชันเทียบมาตราเชิงเส้นตรงย่อมจะทำให้เกิดความลำเอียงได้ (Jaeger 1981: 25)



บุญ และ ฮอดแอนท์ ได้เสนอทวิเคราะห์ไว้ดังนี้ (๑๙๕๖: ๓๖ - ๓๘) ความ  
 สำคัญเชิงสถิติ มีความหมาย ๒ ประการ คือ

๑. ความแตกต่างของค่าเฉลี่ยของค่าตัวประมาณตลอดการทำการสุ่มซ้ำ ๆ จาก  
 ประชากรเดียวกัน กับค่าของประชากรที่ถูกประมาณ

๒. ความแตกต่างระหว่างฟังก์ชันของการเทียบมาตราประมาณ (estimated  
 equating function) และค่าที่แท้จริงของฟังก์ชันการเทียบมาตรา (equating  
 function)

แหล่งของความสำคัญของของการเทียบมาตราที่สำคัญมี ๒ แหล่ง คือ ความผันแปร  
 ของประชากร (population variability) และความคลาดเคลื่อนของรูปแบบ (model  
 errors) ซึ่งมีคำอธิบาย ดังนี้

### ๒.๑ ความแปรเปลี่ยนของประชากร (population variability)

เนื่องจากการเทียบมาตราไม่ได้จำกัดกระทำกับประชากรเพียงกลุ่มเดียวโดย  
 ตลอด ฟังก์ชันของการเทียบมาตราที่สร้างขึ้นเพื่อใช้กับประชากรหนึ่ง อาจนำไปใช้กับอีกประชากรหนึ่ง  
 ไม่ได้ เช่น ถ้า  $X_P^*(y)$  เทียบมาตรา  $y$  ไปยังมาตรา  $x$  ในประชากร  $P$  และ  $Y_Q^*(z)$  เทียบ  $z$   
 ไปยัง  $y$  ในประชากร  $Q$  แล้ว จะรับความสรุปว่า ฟังก์ชันเชิงซ้อน  $X_P^*(Y_Q^*(z))$  เทียบมาตราจาก  
 $z$  ไปยัง  $x$  ในประชากร  $P$  หรือ  $Q$  ย่อมไม่ได้ ในแบบแผนการเทียบมาตราที่ใช้แบบสอบรวมได้แก่ไข  
 ไขเหล่านี้ โดยการสร้างประชากรสังเคราะห์ (synthesis population) ซึ่งเป็นผลรวมของ  
 ประชากร  $P$  และ  $Q$  ตามสัดส่วน แต่ก็ยังไม่รับประกันเลยว่า จะไม่มีความสำคัญในทุก ๆ ประชากร  
 ที่แปรเปลี่ยนไป

### ๒.๒ ความคลาดเคลื่อนรูปแบบ (model errors)

ความคลาดเคลื่อนนี้ เกิดขึ้นจากการจำกัดกระทำข้อสมมุติฐานที่ว่าด้วยรูปการ  
 แจกแจงให้่ายขึ้นอย่างฉีกพาด ซึ่งแบ่งได้เป็น ๒ พวก คือ

- (๑) ความคลาดเคลื่อนรูปแบบที่ทดสอบได้ (testable model errors) เช่น ความเป็นเส้นตรงของฟังก์ชันการถดถอย สามารถทดสอบได้จากข้อมูลจำนวนมาก ซึ่งจะเห็นได้ว่า มีความเพียงพอหรือไม่
- (๒) ความคลาดเคลื่อนรูปแบบที่ทดสอบไม่ได้ (Nontestable model errors) เป็นกรณีที่ไม่สามารถหาข้อมูลมาทดสอบเพื่อพิสูจน์ว่า สมมุติฐานที่ว่าไว้เพียงพอหรือไม่ ตัวอย่างที่พบ เช่น ความคลาดเคลื่อนที่เกิดจากการประมาณความถี่ในวิธีการเทียบมาตราของทักเกอร์ ซึ่งกล่าวว่า การแจกแจงความถี่ของ X อย่างมีเงื่อนไขเมื่อกำหนด V ใน P เหมือนกับการแจกแจงความถี่ของ X อย่างมีเงื่อนไขเมื่อกำหนด V ใน Q ความคลาดเคลื่อนของการกำหนดเป็นข้อสมมุติฐานนี้ หากเกิดขึ้นจะมีผลก่อให้เกิดความลำเอียงในการประมาณของฟังก์ชันการเทียบมาตราได้อย่างมาก

สำหรับเรื่องของการวิเคราะห์แหล่งความคลาดเคลื่อนทั้งที่เป็นความแปรปรวนของการสุ่มและความลำเอียง ยังต้องมีการวิจัยอีกมากเพื่อสร้างวิธีปฏิบัติที่เหมาะสมกับปัญหาในแง่คุณค่าในแง่คุณค่าต่าง ๆ (Braun and Holland 1982: 32) ทั้งนี้ การวิจัยครั้งนี้ผู้วิจัยจะประเมินรูปแบบการเทียบมาตราต่าง ๆ ด้วยการใช้มาตรการของความคลาดเคลื่อนมาตรฐานของการเทียบมาตราที่แนะนำโดยออร์คเท่านั้น จะไม่พิจารณาในประเด็นของความลำเอียงของการเทียบมาตรา

การประเมินความเพียงพอของการเทียบมาตรา

วิธีการเทียบมาตราจะแนบแต่ละวิธี ประกอบขึ้นด้วยรูปแบบของการเทียบมาตรา (models) ซึ่งมีข้อกำหนดที่ว่าด้วยสมมุติฐานเบื้องต้น (assumptions) ของแต่ละรูปแบบ และประกอบด้วยการออกแบบ (designs) เพื่อจัดเก็บข้อมูลให้เป็นไปตามรชคกลางต่าง ๆ ถ้าหากทุกอย่างเป็นไปตามเงื่อนไข เชื่อได้ว่าผลของการเทียบมาตราจะมีความถูกต้อง (accurate) และความคมชัด (precise) ตามทฤษฎี แต่ความเป็นจริงของการทดสอบมักไม่ได้เป็นไปตามอุดมการณ์ เพราะมีหลายสิ่งหลายอย่างที่อยู่นอกเหนือการควบคุม เช่น โปรแกรม

การทดสอบระดับชาติ ซึ่งจะมีทั้งนโยบาย และกฎหมายที่อยู่นอกเหนือข้อตกลงเชิงทฤษฎี ตัวแบบสอบเองก็จำเป็นต้องมีการเปลี่ยนแปลงไป ดังนั้น ข้อมูลที่นักจิตวิทยานำมาใช้ จึงไม่สามารถระบุว่าเป็นตัวอย่างประชากรได้อย่างชัดเจน และโดยความเป็นจริง เป็นการจัดกระทำการกับค่าประชากรมากกว่า (population quantities) ไม่ใช่ค่าประมาณ (sample estimates) (Braun and Holland 1982: 10)

ด้วยสภาพความเป็นจริงดังกล่าว ทำให้การเทียบมาตราที่ของจักรกระทำอยู่น้อยในภาวะที่มีเงื่อนไขน้อยกว่าความพอดีตามข้อตกลงในแต่ละรูปแบบ ดังนั้น การเทียบมาตราที่ได้พัฒนาขึ้น จึงจำเป็นต้องมีการตรวจสอบความเพียงพอของรูปแบบ (the adequacy of equating models) (Petersen, Marco and Stewart 1982: 71) วิธีการประเมินความเพียงพอมีข้อเสนอแนวความคิดและวิถีปฏิบัติไว้ดังนี้

• ขั้นที่ตรวจสอบความเพียงพอของเจเกอร์ (Jaeger)

เจเกอร์ (Jaeger 1981: 26) ได้เสนอขั้นที่ ๕ ทวี เพื่อตรวจสอบความเพียงพอของการใช้รูปแบบเชิงเส้นตรงว่า เทคนิควิธีที่ได้นำมาใช้ในการเทียบเพียงพอกับการปรับความแตกต่างระหว่างการแจกแจงของคะแนนจากแบบสอบหรือไม่ หรือจำเป็นต้องมีรูปแบบอื่นที่มีความสัมพันธ์มากขึ้น ขั้นที่ ๕ ทวี มีดังนี้ คือ

๑.๑ ขั้นที่ความคล้ายคลึงของการแจกแจงคะแนนสะสมของแบบสอบเก่าและชุดใหม่ โดยการปรับความแตกต่างระหว่างค่าเฉลี่ยและส่วนเบี่ยงเบนมาตรฐาน การทดสอบความเหมือนของการแจกแจงใช้การทดสอบด้วย the Kolmogorov-Smirnov two-sample test (Smirnov 1948 cited by Jaeger 1981: 27)

๑.๒ รูปแบบของการแจกแจงคะแนนดิบกับคะแนนแปลง (Shape of the raw score transformation) เหตุผลของการใช้ขั้นที่ตัวนี้ คือ ถ้าการเทียบมาตราด้วยรูปแบบเชิงเส้นตรงสามารถอธิบายความแตกต่างในการแจกแจงของคะแนนดิบทั้งสองชุดอย่างเพียงพอ ก็เป็นเหตุสุดอย่างเพียงพอ เช่นกันที่จะยอมรับว่าการแปลงคะแนนดิบจากแบบสอบชุดใหม่ไปยังแบบสอบชุดเก่าเป็นเส้นตรงอย่างแน่นอน

๑.๓ ความคงเส้นคงวาของผลลัพธ์ของการเทียบมาตรฐานรูปแบบเชิงเส้นตรงกับการเทียบที่ตำแหน่งเปอร์เซนไทล์ (Consistency of linear and equipercentile equating results) การวิเคราะห์ที่อาศัยข้อตกลงที่เป็นสมมุติฐานเบื้องต้นที่ว่า ถ้ารูปแบบเชิงเส้นตรงมีความเพียงพอแล้ว พังจน์ของรูปแบบการเทียบที่ตำแหน่งเปอร์เซนไทล์จะแปรผันไปโดยสุ่มรอบ ๆ พังจน์ของรูปแบบเชิงเส้นตรงที่สมนัยกัน

๑.๔ ความคล้ายคลึงของการแจกแจงความยาวของข้อสอบ (Similarity of item difficulty distributions) โดยอาศัยหลักการที่ว่า การใช้รูปแบบเชิงเส้นตรงมีความเพียงพออย่างแท้จริงกับแบบสอบที่มีคุณสมบัติเป็นคู่ขนานกัน ถ้ามีการเบี่ยงเบนจากความเป็นคู่ขนานมากเท่าใด แสดงว่าข้อถาวรรูปแบบการเทียบมาตรฐานที่ซับซ้อนขึ้น เพราะการแจกแจงของแบบสอบที่ไม่ใช่คู่ขนานจะมีความแตกต่างเกิดขึ้นในระกัมโมเมนต์ที่สูงขึ้น

๑.๕ ความคล้ายคลึงของค่าอำนาจจำแนกของข้อสอบ (Similarity of item discrimination distributions) มีเหตุผลทำนองเดียวกับที่ข้อที่ ๔

## ๒ ดัชนีความแตกต่าง (Discrepancy Indices)

ปีเตอร์สัน มาร์โค และ สตีเวอท์ (Petersen, Marco and Stewart 1982: 91) ได้เสนอวิธีประเมินความเพียงพอของรูปแบบการเทียบมาตรฐานโดยการหาความแตกต่างระหว่างคะแนนแปลง (an estimated criterion score,  $t'$ ) ซึ่งเป็นผลจากการเทียบมาตรฐานกับคะแนนเกณฑ์ (criterion score,  $t$ ) ที่สมนัยกัน ถ้าความแตกต่างมีค่าน้อย มีความหมายว่า ความคลาดเคลื่อนที่เกิดขึ้นจากการใช้รูปแบบการเทียบมาตรฐานนั้น ๆ น้อยกว่า รูปแบบการเทียบมาตรฐานดังกล่าวย่อมมีความเหมาะสมกับสถานการณ์ที่ใช้

ดัชนีความแตกต่างที่นำมาเปรียบเทียบระหว่างรูปแบบและสถานการณ์ที่ต่าง ๆ กัน คือ กำลังสองของค่าเฉลี่ยของความแตกต่างที่ถ่วงน้ำหนักด้วยความแปรปรวนของคะแนนเกณฑ์ (the weighted mean-square difference) เป็นค่ามาตรฐาน ค่าที่คำนวณออกมาเรียกว่า ความคลาดเคลื่อนรวม (the total error) ซึ่งมีสูตรดังนี้



$$\text{total error} = \frac{\sum_j f_j d_j^2}{n s_t^2} \text{ --- (๕๑)}$$

เมื่อ  $d_j = (t - t')$   
 $n =$  จำนวนคะแนนที่ใช้  
 $s_t^2 =$  ความแปรปรวนของคะแนน  $t$

๓ ดัชนีเปรียบเทียบเปอร์เซนไทล์ (the percentile comparison index)

เป็นมาตรการวัดความไม่สอดคล้องระหว่างการแจกแจงของคะแนนในแบบสอบชุด  $x$  กับแบบสอบชุด  $y$  ที่ได้แปลงไปอยู่ในมาตราคะแนนของ  $x$  แล้วตามวิธีที่ได้พัฒนาขึ้น ดัชนีเปรียบเทียบเปอร์เซนไทล์ คือ ค่าเฉลี่ยกำลังสองของความแตกต่าง (the mean-squared difference) ที่ได้จากการแจกแจงของคะแนนต่าง ๆ ของเกณฑ์  $x$  กับคะแนนแปลง  $x^*$  ที่แปลงมาจาก  $y$  ด้วย วิธีการเทียบมาตราที่ระบุไว้ ณ ที่ตำแหน่งเปอร์เซนไทล์เดียวกัน ดัชนีตัวนี้ได้เสนอโดย โคลเลน (Kolen 1982) ได้แนะนำให้ใช้ข้อมูลจากกลุ่มตัวอย่างสอบทานผลซึ่งได้สุ่มมาจากประชากรเดียวกันกับกลุ่มตัวอย่างที่ใช้พัฒนาการวางคะแนนแปลง สูตรการคำนวณมีดังนี้ คือ

$$c = \frac{\sum_i (x_i - x_i^*)^2}{nk} \text{ --- (๕๒)}$$

เมื่อ  $n$  คือ จำนวนของคะแนนดิบของกลุ่มสอบทานผล  
 $k$  คือ จำนวนข้อสอบในแบบสอบรวมที่ใช้

ค่า  $c$  ที่ได้ ถ้ามีค่าน้อยจะให้ความหมายว่า รูปแบบการเทียบมาตราที่นำมาสร้างการวางคะแนนแปลงนั้นมีความเหมาะสมและเพียงพอที่จะให้ผลของการแปลงคะแนนอย่างคงเส้นคงวา

วิธีการประเมินความเพียงพอที่ได้กล่าวมานี้ อาจจำแนกเป็นสองพวก พวกแรก เป็นการประเมินก่อนดำเนินการ เช่น ดัชนีความคล้ายคลึงของการแจกแจงคะแนนสะสมของแบบสอบสองชุด

ทัศนคติความคล้ายคลึงของการแจกแจงของค่าความยากของข้อสอบ เป็นต้น พวกหลัง เป็นการประเมินผลของการเทียบมาตรา ซึ่งอาศัยคะแนนเกณฑ์ที่เลือกแล้วเป็นหลักในการเทียบหาความแตกต่างสำหรับทัศนคติที่เสนอโดย ปีเตอร์สันและคณะนั้น คะแนนเกณฑ์ที่ใช้ คือ ผลการแปลงคะแนนด้วยรูปแบบอิงพหุฎีกาการตอบข้อสอบ (๑๙๘๖: ๘๖) ค่าทัศนคติคำนวณมีลักษณะเป็นหน่วยมาตรฐานแล้ว สามารถนำค่าเหล่านี้ที่ได้จากการใช้รูปแบบที่ต่างกันตลอดจนสถานการณ์ใช้ข้อมูลที่คำนวณมาเปรียบเทียบกันโดยตรงได้ ส่วนทัศนคติที่เสนอโดย โคลเลน (Kolom 1982) ได้ใช้ข้อมูลคะแนนจากผู้สอบเองเป็นเกณฑ์ในการหาความแตกต่าง ข้อมูลเหล่านี้ได้มาจากการออกแบบด้วยการใช้กลุ่มตัวอย่างสอบทานผลซึ่งผู้สอบในกลุ่มตัวอย่างนี้ได้รับการทดสอบด้วยแบบสอบทั้งสองชุด ดังนั้น การใช้คะแนนของตนเองเป็นเกณฑ์จึงมีความอิสระไม่ข้องขึ้นกับกระบวนการแปลงคะแนนอื่น ๆ เช่นวิธีที่เสนอโดย ปีเตอร์สันและคณะด้วยเหตุผลดังกล่าว ผู้วิจัยจึงได้เลือกวิธีการประเมินความเพียงพอกจากการวิเคราะห์ผลในกลุ่มตัวอย่างสอบทานผล แต่การหาค่าทัศนคติได้ก็แตกต่างจากสูตรของโคลเลน ไปใช้ตามแนวความคิดของปีเตอร์สันและคณะ คือ ใช้ค่าความแปรปรวนเป็นตัวถ่วงน้ำหนักเพื่อให้ค่าที่ได้มีหน่วยเป็นมาตรฐาน

สูตรที่ก้กแปลงและนำมาใช้ในการศึกษารังนี้ คือ

$$C = \frac{\sum (x_i - x_i^*)^2}{n S_x^2} \text{ ----- (๘๓)}$$

- เมื่อ  $x_i$  คือ คะแนนเกณฑ์หรือคะแนนจากการสอบชุด  $x$  ของคนที่  $i$
- $x_i^*$  คือ คะแนนที่ได้เทียบด้วยการวางแปลงคะแนนที่สมนัยกันของคนี่  $i$
- $n$  คือ จำนวนคนในกลุ่มตัวอย่างสอบทานผลที่นำมาวิเคราะห์
- $S_x^2$  คือ ค่าความแปรปรวนของคะแนน  $x$

งานวิจัยที่เกี่ยวข้องกับการ เปรียบเทียบรูปแบบการเทียบมาตรา

งานวิจัยที่มุ่งตรวจสอบผลจากวิธีการเทียบมาตราแบบต่าง ๆ ในสถานการณ์ที่ได้ออกแบบจัดกระทำข้อมูลต่าง ๆ กันในปัจจุบันนี้ นับว่ามีมากพอที่จะให้วิเคราะห์ในเชิงวรรณคดี แต่การลงข้อสรุปที่ว่า วิธีการใดให้ผลที่แม่นยำที่สุดนั้นยังไม่อาจทำได้ ในการวิเคราะห์จากวรรณคดีพบประเด็นที่เป็น

เงื่อนไขก่อนผลการสรุปที่กองพิจารณาต่าง ๆ กัน ซึ่งจะกล่าวในรายละเอียดต่อไป

สไลด์ และ ลิน (Slindt and Linn 1977, 1978, 1979) ได้ทำการตรวจสอบถึงปัญหาของการเทียบมาตราในแนวตั้ง (vertical equating) ของแบบทดสอบสองชุดที่สร้างขึ้นเพื่อใช้กับประชากรที่มีระดับความสามารถต่างกันด้วยวิธีอื่น จากการศึกษาได้ให้ข้อเสนอแนะว่า การใช้รูปแบบการเทียบมาตราสามรูปแบบ คือ เริงเส้นตรง อีควิเปอร์เพนโทล และ IRT ชนิดหนึ่งพหุพารามิเตอร์รูปแบบโลจิสติกส์หรือจำกัดในกระบวนการเทียบมาตราในแนวตั้ง โดยเฉพาะอย่างยิ่ง เมื่อทำการเทียบมาตราโดยใช้แบบทดสอบสองชุดที่มีความแตกต่างกันมาก และกลุ่มตัวอย่างสองกลุ่มที่มีระดับความสามารถต่างกันมาก จากการศึกษาให้ข้อสังเกตว่า ถ้าใช้รูปแบบ IRT ชนิดสามพารามิเตอร์ของโลจิสติกส์ให้ผลการเทียบมาตราที่ดีกว่าในสถานการณ์เช่นนี้

ปีเตอร์สัน มาร์โค และ สตีเวอร์ท (Petersen, Marco and Stewart 1982; 71-135) เป็นกลุ่มนักวิจัยที่ทำการศึกษาระเบียงประจักษ์ถึงเทคนิคการเทียบมาตราต่าง ๆ อย่างกว้างขวางที่สุด การเทียบมาตราได้จัดกระทำโดยเทียบแบบสอบชุด SAT (Scholastic Aptitude Test) ที่จัดสอบเมื่อเดือนเมษายน ๑๙๗๕ กับแบบสอบชุด TSWE (Test of Standard Written English) ที่จัดสอบเมื่อพฤศจิกายน ๑๙๗๕ ข้อมูลที่นำมาใช้เป็นบันทึกผลการสอบรายชื่อของผู้เข้าสอบทั้งหมด กลุ่มตัวอย่างได้สร้างขึ้นโดยการกำหนดลักษณะเฉพาะ และเลือกคะแนนตามจุดมุ่งหมาย มีทั้งคะแนนสอบที่ต้องการเทียบ และแบบสอบรวม ค่ามาตรฐานสถิติของข้อสอบรายชื่อ และข้อมูลอื่น ๆ ที่ต้องการใช้ มีกลุ่มตัวอย่างพื้นฐาน ๓ กลุ่ม จากการศึกษาเมื่อเดือนเมษายน ๑๙๗๕ และอีก ๓ กลุ่ม จากการศึกษาเมื่อเดือนพฤศจิกายน ๑๙๗๕ กลุ่มตัวอย่างไม่มีคนสอบซ้ำประกอบด้วยกลุ่มละ ๔๗๗ คน แล้วจึงเป็นกลุ่มย่อยโดยออกแบบการสุ่มให้ลักษณะแตกต่างกัน ๓ กลุ่มย่อย คือ กลุ่มสุ่มอย่างง่าย กลุ่มความสามารถคล้ายคลึงกัน (ใช้ค่าเฉลี่ยความสามารถทางภาษาเป็นเกณฑ์) และกลุ่มไม่คล้ายคลึงกัน การออกแบบเทียบมาตราได้ใช้แบบสอบ SAT ส่วนที่เป็นภาษามาแยกเป็นแบบสอบ ๒ ชุด แล้วนำมาเทียบมาตรากันโดยใช้แบบสอบรวมลักษณะต่าง ๆ กัน คือ ต่างกันในประเด็นของตำแหน่งแบบสอบรวม (ภายนอกและภายใน) ประเด็นของเนื้อหา (คล้ายคลึงกันและไม่คล้ายคลึงกัน) ประเด็นของความยาก (คล้ายคลึงกันและไม่คล้ายคลึงกัน) ส่วนในเรื่องของกลุ่มตัวอย่างที่นำเข้มาทั้งหมด ๓ ลักษณะดังกล่าวแล้ว นอกจากนี้ได้ทำการเทียบมาตราแบบสอบ SAT กับแบบสอบที่สร้าง

ขึ้นจากแบบสอบ TSWE โดยใช้แบบสอบร่วมชนิดภายใน ลักษณะของแบบสอบร่วมที่ใช้กับแบบสอบที่  
 เที่ยงตัวเอง (SAT) ได้ศึกษาทั้งที่มีเนื้อหาคลาย และไม่คล้ายกับฉบับที่กองการเทียบ การประเมิน  
 ความเพียงพอของวิธีการเทียบมาตราในสถานการณ์ต่าง ๆ ที่ทำการศึกษาศักยภาพความแตกต่าง  
 (Discrepancy index) ก็ชนนี้ คือ ค่าที่ถ่วงน้ำหนักแล้วของค่าเฉลี่ยกำลังสองของความแตกต่าง  
 ระหว่างค่าประมาณ หรือคะแนนแปลงกับคะแนนเกณฑ์ (The weighted mean-square  
 difference) คะแนนเกณฑ์ที่ใช้ คือ คะแนนแปลงจากคะแนนดิบชุดเดียวกัน ภายรูปแบบการ  
 เทียบมาตรา IRT ค่านี้ทำให้เป็นมาตรฐานด้วยการทำให้เป็นสัดส่วนของความเบี่ยงเบนมาตรฐาน  
 ของคะแนนเกณฑ์ ซึ่งนำไปใช้เปรียบเทียบผลที่ได้จากการวิเคราะห์ทางสถานการณ์และทางรูปแบบที่  
 ใช้ ค่าก็ชนนี้ให้ความหมายตามหลักการ ดังนี้ ค่าก็ชนนี้ความแตกต่าง หรือที่เรียกว่า ค่าความคลาด  
 เคลื่อนรวม (TE) ซึ่งคำนวณจากสูตรที่ (๕๑) เมื่อส่วนเบี่ยงเบนมาตรฐานของคะแนนเกณฑ์ที่ใช้มีค่า  
 เท่ากับ ๑๐๐ ปริมาณของ TE มีความหมายในเชิงคุณภาพของวิธีการเทียบมาตรา ดังนี้

TE ≤ ๒๕	คือ	น่าพอใจอย่างมาก	(Zone A)
๒๕ < TE ≤ ๑๐๐	คือ	น่าพอใจ	(Zone B)
๑๐๐ < TE ≤ ๒๒๕	คือ	ปานกลาง	(Zone C)
๒๒๕ < TE ≤ ๔๐๐	คือ	ไม่น่าพอใจ	(Zone D)
๔๐๐ < TE	คือ	ไม่น่าพอใจอย่างมาก	(Zone E)

ระดับต่าง ๆ ที่กำหนดนี้มาจากหลักเหตุผลที่ว่า ถ้าความคลาดเคลื่อนรวม (TE) เกิดขึ้น  
 เนื่องจากความลำเอียงเพียงสาเหตุเดียวแล้ว ย่อมเป็นเหตุผลที่ให้นักจิตวิทยายอมรับได้ว่า ผลการ  
 เทียบมาตราอยู่ในระดับน่าพอใจอย่างมากเมื่อค่าเฉลี่ยความคลาดเคลื่อนต่างจากความจริง (truth)  
 ไม่เกินไปกว่าร้อยละ ๕ ของส่วนเบี่ยงเบนมาตรฐานของคะแนนเกณฑ์ และจะไม่เป็นที่ยอมรับอย่าง  
 แน่นนอนถ้าผลการเทียบมีค่าเฉลี่ยต่างจากความจริงถึงร้อยละ ๒๐ ของส่วนเบี่ยงเบนมาตรฐานของ  
 คะแนนเกณฑ์ จุดตัดที่กำหนดขึ้นนี้ถึงแม้จะเป็นการกำหนดโดยคณะผู้วิจัยเอง แต่เป็นวิธีการที่มีเหตุผล  
 และสามารถทำให้ศึกษาเปรียบเทียบผลการเทียบมาตราที่เกิดจากการใช้วิธีการ และสถานการณ์ที่ต่าง  
 กันได้ ผลการวิจัยสรุปไว้ ดังนี้



(๑) รูปแบบการเทียบเชิงเส้นตรงโดยใช้เทคนิคของพอททอฟ (Potthoff) ให้นผลไม่น่าพอใจเกือบทุกสถานการณ์ที่ศึกษา

(๒) ในกรณีที่ไขกลุ่มตัวอย่างขนาดเท่ากัน รูปแบบอควิเปอร์ เชนไคโลให้ค่าความคลาดเคลื่อนมากกว่าการเทียบด้วยวิธีการเชิงเส้นตรง

(๓) รูปแบบเชิงเส้นตรงโดยเทคนิคของทักเกอร์ ๑ ให้ค่าความคลาดเคลื่อนขนาดใหญ่กว่าเชิงเส้นตรงในกรณีที่ไขกลุ่มตัวอย่างแตกต่างกัน

(๔) การใช้รูปแบบเชิงเส้นตรงให้ผลไม่เหมาะสม เมื่อเกิดความสัมพันธ์เชิงเส้นโค้ง ซึ่งเนื่องจากความยากของแบบสองสองซุกต่างกัน

(๕) รูปแบบอควิเปอร์ เชนไคโลให้ผลการเทียบดีกว่ารูปแบบเชิงเส้นตรง เมื่อแบบสองสองซุกมีความสัมพันธ์เชิงเส้นโค้ง ซึ่งเกิดจากความยากที่ต่างกันของแบบสองสองซุก

(๖) ถ้าความยากของแบบสองรวมต่างกับแบบสองฉบับที่เทียบตัวเองมาก ความคลาดเคลื่อนจะเกิดมากขึ้นด้วย

(๗) แบบสองที่ทำการเทียบมาตราสองซุก ถ้ามีความยากต่างกันมาก จะเกิดความคลาดเคลื่อนมากด้วย

(๘) ถ้าแบบสองรวมได้สร้างให้มีลักษณะคล้ายกับฉบับที่เทียบมามาก หรืออาจเรียกว่าเป็นฉบับย่อยส่วนแล้ว การเทียบมาตราจะให้ผลดีที่สุด

จุฬาลงกรณ์มหาวิทยาลัย

ผู้วิจัยได้ให้ข้อสังเกตว่า การสรุปผลการวิจัยเป็นการทั่วไปนั้นต้องระมัดระวังมาก เพราะ การศึกษานี้ใช้เฉพาะแบบสอบส่วนที่เป็นภาษาเท่านั้น การศึกษายังไม่ได้แยกผลของความแตกต่างใน ส่วนที่เกี่ยวกับเนื้อหา และความยากอย่างอิสระ และสุดท้าย คือ เกณฑ์การเทียบอาจให้ความลำเอียง ในสถานการณ์บางอย่าง อย่างไรก็ตาม อย่างไรก็ดีผลการวิจัยอาจให้ประโยชน์ต่อการนำไปใช้โดยที่ผู้ใช้ต้องพิจารณา ความเหมาะสมของรูปแบบกับสถานการณ์ที่กำหนด และลักษณะเฉพาะของกลุ่มตัวอย่าง นอกจากนี้ ผู้วิจัยมีความเห็นว่า ควรมีการพัฒนาในส่วนที่เป็นกฎการตัดสินใจ เพื่อให้เห็นว่า ควรเลือกใช้รูปแบบ ใดกับลักษณะแบบสอบและกลุ่มตัวอย่างคนที่แปรเปลี่ยนไป

ปีเตอร์สัน กูก และ สตอกกิง (Petersen, Cook and Stocking 1981) ทำการศึกษาเปรียบเทียบความคงเส้นคงวาของมาตราที่เกิดจากการเทียบด้วยรูปแบบ IRT กับรูปแบบ กึ่งเกมสองรูปแบบ คือ อีควิเปอร์เชนโกล์ และรูปแบบเชิงเส้นตรง รูปแบบเชิงเส้นตรงที่นำมาศึกษา มี ๓ วิธี คือ Tucker model, Levine Equally Reliable model และ Levine unequal Reliable model ส่วน IRT ใช้ ๓ วิธี คือ วิธีปัจจุบัน วิธีกิ่งวิเคราะห์ค่า โดยกำหนดค่าความยากคงที่ และวิธีกิ่งวิเคราะห์ค่าโดยกำหนดค่าความยากที่เทียบแล้ว การศึกษา ครั้งนี้ของการประเมินราคาของการ เบี่ยงเบนของมาตราในกระบวนการเทียบมาตราที่เป็นลูกโซ่ การเบี่ยงเบนของมาตรานี้ คือ ความแตกต่างของคะแนนแปลงในแบบสอบ (โดยได้ผ่านการแปลง หลายครั้งเหมือนลูกโซ่) กับคะแนนเก็บบ่อนแปลง เช่น คะแนน v4 เทียบกับ v4 ที่ได้จากการ ใช้วิธีการเทียบมา ๔ ครั้ง กลุ่มตัวอย่างสุ่มจากการทดสอบ ๔ ครั้ง ในรอบ ๒ ปี แยกเป็นกลุ่มของ ฉับบภาษา และกลุ่มของฉับบคณิตศาสตร์ แต่ละฉับบมีกลุ่มตัวอย่างย่อย ๑๒ กลุ่ม กลุ่มละ ๒๒๒๐ คน โดยประมาณ แบบสอบรวมมีเนื้อหากำหนดแน่นอนไม่มีการเปลี่ยนแปลง การเปรียบเทียบความคลาด-เคลื่อนใช้การพิจารณาจากภาพกราฟ และดัชนีความแตกต่าง (a discrepancy index) ซึ่งเป็นค่าถ่วงน้ำหนักแล้วของค่าเฉลี่ยความแตกต่างกำลังสอง (The weighted mean square difference) ความแตกต่างคำนวณได้จากผลต่างของคะแนนมาตราเกณฑ์,  $t$  (คะแนนกึ่งเกม ของ v4) กับคะแนนมาตราที่แปลงด้วยวิธีการต่าง ๆ ที่ศึกษา,  $t'$  ดัชนีนี้เป็นดัชนีสรุปช่วยในการ ประเมินประสิทธิภาพของรูปแบบการเทียบมาตราต่าง ๆ ผลการวิจัยพบว่า ในการใช้รูปแบบเชิงเส้น-ตรงกับการเทียบมาตราสายแบบสอบภาษา และสายแบบสอบคณิตศาสตร์ของ SAT มีผลที่คล้ายคลึงกัน

Tucker model ให้ผลความคลาดเคลื่อนมากที่สุด ส่วน Levine equally Reliable model ให้ผลความคลาดเคลื่อนน้อยที่สุด รูปแบบอิกวิเปอร์เซนโกล์ให้ผลโดยทั่วไปไม่น่าพอใจเมื่อเทียบกับรูปแบบเชิงเส้นตรง เหตุผลที่น่าจะอธิบายได้ คือ รูปแบบอิกวิเปอร์เซนโกล์ไม่สามารถให้สารสนเทศได้ตลอดพิสัยของคะแนนของแบบสอบ ถ้ากลุ่มตัวอย่างที่ใช้ไม่มีการแจกแจงโดยตลอด รูปแบบ IRT ซึ่งทำการศึกษา ๓ วิธีนี้พบว่า ลักษณะผลการเทียบโดยทั่วไปที่ทำได้กับแบบสอบภาษา และแบบสอบคณิตศาสตร์ มีความแตกต่างกัน วิธีที่ใช้ค่า  $b$  ที่เทียบแล้วให้น่าพอใจที่สุดกับแบบสอบภาษา แต่ให้ผลตรงกันข้ามกับแบบสอบคณิตศาสตร์ สำหรับวิธี IRT ปัจจุบัน ให้ผลดีกว่าวิธีที่ใช้ค่า  $b$  เทียบแล้วเพียงเล็กน้อยในแบบสอบภาษา แต่เมื่อเปลี่ยนเป็นแบบสอบคณิตศาสตร์ วิธี IRT ปัจจุบันให้ผลดีกว่ามาก เปรียบเทียบผลของรูปแบบ IRT กับรูปแบบทั้งเดิม พบว่า ความสัมพันธ์ที่เกิดขึ้นกับแบบสอบภาษาต่างกับแบบสอบคณิตศาสตร์ สำหรับการเทียบในสายของภาษารูปแบบ IRT ดีกว่าแต่ในสายคณิตศาสตร์รูปแบบ Levine และ IRT ปัจจุบันให้ผลคล้ายคลึงกันมาก ผลการวิจัยได้สรุปเป็นข้อเสนอว่า ถ้าหากจำเป็นต้องเลือกรูปแบบการเทียบเพียงรูปแบบเดียวในการเทียบกับวิชาทางภาษาและคณิตศาสตร์แล้ว รูปแบบ IRT ปัจจุบันเป็นวิธีที่เหมาะสมที่สุดใน ๘ วิธีที่ได้ศึกษามา แต่ไม่ได้หมายความว่า รูปแบบ IRT ปัจจุบันจะให้ผลดีกว่าเสมอไปเมื่อข้อมูลที่ใช้ต่างไปจาก SAT

โคเลน (Kolen 1981) เปรียบเทียบผลการเทียบมาคราระหว่างรูปแบบทั้งเดิมสองวิธี คือ รูปแบบเชิงเส้นตรง และรูปแบบอิกวิเปอร์เซนโกล์ กับรูปแบบ IRT ๘ วิธี โดยใช้ข้อมูลการทดสอบนักเรียนในรัฐไอโอวา ปี ๑๙๗๕ ในโครงการ the Iowa Tests of Educational Development (ITED) ข้อสอบที่ทำการศึกษา คือ แบบสอบวัดสัมฤทธิ์ผล ที่ใช้กับระดับที่ ๑ คือ ระดับการศึกษาเกรด ๔ และ ๑๐ ระดับที่ ๒ คือ ที่ใช้กับเกรด ๑๑ และ ๑๒ ในแต่ละระดับแบบสอบได้พิมพ์ไว้ในปีต่าง ๆ สำหรับฉบับที่นำไปศึกษา คือ ฉบับที่พิมพ์ครั้งที่ ๘ เทียบไปหาครั้งที่ ๖ โดยให้  $X_7$  และ  $Y_7$  ใด ๆ เป็นแบบสอบสองชุดที่ต้องการเทียบกัน เป็นแบบสอบที่ไม่เป็นคู่ขนานกัน การเทียบมาครานั้น ต่างเทียบ  $X_7$  และ  $Y_7$  ไปสู่มาคราร่วมที่กำหนดขึ้น คือ  $X_6$  โคเลน ได้ออกแบบการรวบรวมข้อมูลโดยสุ่มโรงเรียน ๓๔ แห่ง กำหนดให้นักเรียนระดับ ๑ และ ๒ ได้รับการทดสอบในทำนองเดียวกัน คือ ใช้วิธีกำหนดให้นักเรียนได้รับการทดสอบเพียงชุดเดียว  $X_6$   $X_7$  หรือ  $Y_7$  โดยให้สอบชุดละ ๑ ใน ๓ ของคนสอบทั้งหมดโดยการสุ่ม สำหรับข้อมูลในกลุ่มสอบทานผล คือ นักเรียน



ทุกคนที่ ๓ ของแต่ละชุดแต่ละระดับของแบบสอบที่ใช้ ขั้นตอนการดำเนินการเทียบรูปแบบอิกวิเปอร์-  
 เทนไคด์ และรูปแบบเชิงเส้นตรงได้จักรกระทำตามวิธีของแองกอฟ (Angoff 1984) ส่วนรูปแบบ  
 IRT ที่ศึกษามี ๓ วิธี คือ ใช้ ๑, ๒ หรือ ๓ พารามิเตอร์ แต่ละวิธีเทียบด้วยการประมาณคะแนน  
 จริง (Estimated true score equating) และเทียบด้วยการประมาณค่าสังเกต  
 (Estimated observed score equating) วิธีสุดท้าย คือ Rasch model เกณฑ์  
 ใช้ในการประเมินผลการเทียบมาตราวิธีต่าง ๆ มีหลายวิธี สำหรับโคเสน ได้เลือกวิธีการตรวจสอบ  
 ผลสุดท้ายของการเทียบมาตราในแต่ละวิธีจากกลุ่มตัวอย่างสอบทานผล ซึ่งเป็นกลุ่มตัวอย่างอิสระ  
 กันที่ไรเปรียบเทียบ คือ ค่าเฉลี่ยกำลังสองของความแตกต่าง (mean squared difference)  
 ระหว่างคะแนนจากแบบสอบชุดเก่ากับชุดใหม่ที่เปลี่ยนไปสู่มาตราชุดเก่าด้วยวิธีต่าง ๆ กันในแต่ละ  
 ตำแหน่งเปอร์เซนต์ไคด์เดียวกัน ค่าสถิติจากการวิเคราะห์กลุ่มสอบทานผลนำมาทดสอบความแตกต่าง  
 อย่างมีนัยสำคัญทางสถิติด้วยการทดสอบฟรายด์แมน (Friedman test) ผลการวิจัยพบว่า ผล  
 การเทียบมาตรา ๓ วิธี แตกต่างอย่างมีนัยสำคัญที่ระดับ .๐๑ และ .๐๕ เมื่อใช้กับแบบทดสอบระดับ  
 ที่ ๑ และ ๒ ตามลำดับ เมื่อแยกพิจารณาแบบสอบระดับที่ ๑ ซึ่งมีความยากต่างกันระหว่างแบบสอบ  
 ชุดเก่าและชุดใหม่ พบว่า วิธี IRT สามพารามิเตอร์ประมาณค่าการสังเกตให้ผลได้ค่อนข้างแน่นอนมาก  
 ที่สุด ส่วนวิธีที่ให้ผลไม่แน่นอนที่สุด คือ วิธีรูปแบบเชิงเส้นตรง และถัดขึ้นมา คือ วิธี IRT พารา-  
 มิเตอร์เดียว สำหรับแบบสอบของระดับที่ ๒ วิธีที่ให้ผลแน่นอนที่สุด คือ วิธี IRT สามพารามิเตอร์  
 ประมาณคะแนนจริง ส่วนวิธีที่ให้ผลแน่นอนน้อยที่สุด คือ IRT พารามิเตอร์เดียว และถัดขึ้นมา คือ  
 วิธีรูปแบบเชิงเส้นตรง โคเสน ได้เสนอการอภิปรายสรุปโดยส่วนรวมว่า วิธีที่ให้ผลการเทียบมาตรา  
 ที่ไม่เพียงพอที่สุด คือ วิธี IRT พารามิเตอร์เดียว ทั้งนี้มีข้อน่าสังเกตปัจจัยที่อาจเป็นสาเหตุ คือ  
 การเดาซึ่งโดยรูปแบบการประมาณค่าพารามิเตอร์ของข้อสอบไม่ได้ระบุ แต่ถือว่าเป็นส่วนหนึ่งของ  
 ความสามารถ ทั้งนี้ ลักษณะข้อมูลที่ได้มาจากรายการข้อสอบที่มีความยากต่างกัน การเทียบมาตราอาจได้รับผล  
 กระทบ และทำให้ประสิทธิภาพต่ำลง ข้อสังเกตนี้เป็นที่น่าองเดียวกับผลการวิจัยของ สไลด์ และ ลินท์  
 (Slinde and Linn 1979) ในกรณีการเทียบโดยใช้วิธี IRT สามพารามิเตอร์ประมาณ  
 คะแนนจริง ซึ่งให้ผลการเทียบที่น่าพอใจที่สุดใน ๓ วิธีที่ศึกษา ยังมีปัญหาที่ต้องศึกษา คือ ผลของการ  
 เทียบมาตราที่ต่ำกว่าระดับการเดา ซึ่งต้องอาศัยวิธีการประมาณอื่น เช่น วิธีประมาณเชิงเส้นตรงที่นำ  
 มาใช้ในครั้งนี้ วิธีอิกวิเปอร์ เทนไคด์ให้ผลการเทียบมาตราอยู่ในระดับเพียงพอปานกลางสำหรับการ



ศึกษาในครั้งนี้ เป็นวิธีที่ผู้นำไปใช้ไม่สามารถคาดหวังที่จะได้ผลสมบูรณ์ แต่เป็นวิธีสำคัญในการเทียบแบบสอบที่มีความยากต่างกันในรอบเขตที่คล้ายการศึกษานี้ สำหรับรูปแบบการเทียบเชิงเส้นตรงให้ผลการเทียบที่ไม่น่าพอใจอย่างชัดเจนกับแบบสอบที่มีความยากไม่เท่ากัน

คุก ดันบาร์ และ ไอเนอร์ (Cook, Dunbar and Egnor 1981) ศึกษาวิธีการเทียบมาตราด้วยรูปแบบ IRT ที่จะนำไปสู่การแก้ปัญหาการทดสอบในการปฏิบัติ โดยการเปรียบเทียบผลการเทียบที่ใช้รูปแบบดั้งเดิม วิธีทั้งหมดที่ทำการศึกษามี ๔ วิธี คือ วิธีเชิงเส้นตรง วิธีอีควิเปอร์-เซนไทล์ธรรมดา วิธีอีควิเปอร์เซนไทล์ที่ใช้ประมาณความถี่ และวิธี IRT ประมาณค่าคะแนนจริง

เกณฑ์การเปรียบเทียบความแตกต่างของการเทียบมาตราได้อาศัยการพิจารณา ๒ ทาง ทางแรก คือ ความสอดคล้องสัมพันธ์โดยพิจารณาจากภาพการแจกแจงของคะแนนแปลง ที่ใช้วิธีเทียบแบบดั้งเดิมวิธีหนึ่ง กับการแจกแจงของคะแนนแปลงเมื่อใช้วิธีเทียบด้วยรูปแบบ IRT เกณฑ์การเปรียบเทียบอีกทางหนึ่ง คือ การเปรียบเทียบเชิงตัวเลข ซึ่งเป็นค่าดัชนีความแตกต่าง ซึ่งมาร์โก เคยใช้มาแล้ว (cited in Marco et al 1979) เป็นค่าที่ถ่วงน้ำหนักแล้วของค่าเฉลี่ยกำลังสองของความแตกต่างระหว่างค่าประมาณจากการเทียบมาตราวิธีที่กำหนดกับคะแนนเกณฑ์ เกณฑ์ที่ใช้ในการศึกษา คือ คะแนนแปลงที่ผ่านกระบวนการเทียบมาตราตามรูปแบบ IRT ข้อมูลที่ใช้ในการศึกษาเป็นผลการสอบรายบุคคลจากการสอบสองครั้งที่ได้จัดทำดำเนินการไปแล้ว แบบสอบที่ใช้มี ๒ ชุด ซึ่งสร้างขึ้นโดย the College Board Admissions Testing Program เป็นแบบสอบหลายตัวเลือก มีความยาวและความยากต่างกัน รูปแบบการเทียบมาตรา IRT ที่เลือกไว้ศึกษา คือ การเทียบโดยใช้การประมาณคะแนนจริงแบบสามพารามิเตอร์โลจิส การประมาณค่าพารามิเตอร์ข้อสอบ และความสามารถผู้สอบดำเนินการ เช่น เกี่ยวกับการวิจัยของผู้อื่น คือ ใช้โปรแกรมคอมพิวเตอร์ LOGIST ของวุก และ ลอร์ค ความน่าจะเป็นของการตอบข้อสอบข้อใด ๆ ได้ถูกต้องในรูปแบบ IRT มีพื้นฐาน ดังนี้

จุฬาลงกรณ์มหาวิทยาลัย

$$P_i(\theta) = c_i + (1 - c_i) \frac{a_i(\theta - b_i)}{1 + e^{1.7a_i(\theta - b_i)}} \quad \text{--- (๘๘)}$$

เมื่อ  $a_i$   $b_i$  และ  $c_i$  เป็นพารามิเตอร์ของข้อสอบที่บรรยายข้อสอบข้อที่  $i$  และ  $\theta$  คือระดับความสามารถของผู้สอบ คะแนนจริงของผู้สอบในแต่ละระดับความสามารถประมาณจากสมการ

$$\hat{\xi} = \frac{\sum_{i=1}^n \hat{P}_i(\theta) - \left[ \sum_{i=1}^n \hat{Q}_i(\theta) \right] / A - 1}{A - 1} \quad \text{--- (๘๙)}$$

$$\hat{\eta} = \frac{\sum_{j=1}^m \hat{P}_j(\theta) - \left[ \sum_{j=1}^m \hat{Q}_j(\theta) \right] / A - 1}{A - 1} \quad \text{--- (๙๐)}$$

เมื่อ  $A$  คือ จำนวนตัวเลือกในแต่ละข้อ  $\hat{\xi}$  และ  $\hat{\eta}$  เป็นคะแนนจริงจากการประมาณแบบสอบคนละชุด สำหรับคะแนนจริงที่อยู่ต่ำกว่าระดับการแก้ไขวิธีการ linear interpolation ผลการวิจัยพบวิธีการเทียบมาตราแบบดั้งเดิม ๓ วิธี ให้ผลสอดคล้องอย่างมากกับวิธี IRT ตรงช่วงคะแนนส่วนใหญ่ของมาตรา ความแตกต่างที่เบี่ยงเบนออกจากผลของ IRT เกิดขึ้นเพราะช่วงปลายของการแจกแจง ผลที่ปรากฏนี้จะสรุปว่าเป็นเช่นนั้นเลย ยังอาจมีข้อโต้แย้งได้ เพราะเกณฑ์ที่ใช้ประเมินนั้นเป็นเกณฑ์ยึดการเทียบมาตรา IRT เป็นหลัก ความเหมาะสมของเกณฑ์ถึงแม้จะได้รับการสนับสนุนจากผลการวิจัยก่อนก็ตาม ทั้งนี้ การอภิปรายจึงไม่สามารถสรุปได้หนักแน่น แต่ให้ความจริงจากการค้นพบในการศึกษาครั้งนี้ได้ว่า สภาพการทดสอบที่นำมาศึกษาประกอบด้วยแบบสอบที่มีความยากต่างกัน กลุ่มตัวอย่างไม่ได้เป็นกลุ่มสมมูลโดยสุ่ม แต่อย่างไรก็ตามแบบสอบ ๒ ชุดนี้ก็ได้ทำให้เกิดการแจกแจงคะแนนในแต่ละชุดของกลุ่มตัวอย่างมีลักษณะเหมือนกัน จึงน่าจะเป็นเหตุผลที่ทำให้รูปแบบเชิงเส้นตรงและเส้นโค้งมีผลใกล้เคียงกันมาก

ความแตกต่างที่ปรากฏตรงปลายของการแจกแจงคะแนนเมื่อใช้รูปแบบทั้งที่เป็นเส้นตรง และเส้นโค้ง ส่วนมากเกิดขึ้นเพราะข้อมูลในส่วนนี้มีน้อย หายากมาก ในทางทฤษฎีจึงแนะนำให้ใช้ รูปแบบการเทียบ IRT แบบสามพารามิเตอร์ เพราะไม่มีผลกระทบจากการขาดแคลนข้อมูล แต่อย่างไรก็ดีในส่วนที่อยู่ทางปลายด้านคะแนนค่า IRT เองไม่สามารถประมาณความสัมพันธ์ได้ ท้องอาศัยวิธีอื่น ดังนั้น ในประเด็นของการเปรียบเทียบผลการเทียบคะแนนที่อยู่ส่วนปลายจึงยังคงไม่สามารถให้ข้อสรุปได้

โคเลน และ วิทนี (Kolen and Whitney 1982) ทำการศึกษาเปรียบเทียบ วิธีการเพื่อเทียบมาตรา ๔ วิธี คือ อีควิเปอร์เซนไทล์ เชิงเส้นตรง IRT พารามิเตอร์เดียว และ IRT สามพารามิเตอร์ ว่าวิธีใดให้ความเที่ยงพอในการเทียบมาตราแบบสอบ General Educational Development (GED) มากกว่ากัน แบบสอบ GED เป็นแบบสอบวัดผลสัมฤทธิ์ เพื่อใช้ตัดสินให้ประกาศนียบัตรแก่ผู้ที่ต้องการเทียบความรู้ระดับเตรียมอุดมศึกษาทั่วประเทศ แบบสอบ ที่พัฒนาขึ้นนี้มี ๑๒ ชุด แต่ละชุดประกอบด้วย ๕ ฉบับย่อย การออกแบบรวบรวมข้อมูลใช้กลุ่มตัวอย่าง ร่วม คือ กำหนดให้ผู้สอบแต่ละคนสอบ ๒ ชุด ในจำนวน ๑๑ ชุด ส่วนชุดที่ ๑๒ เป็นแบบสอบร่วมที่ ทุกคนในกลุ่มตัวอย่างต้องทำเหมือนกันหมด การจัดให้ผู้สอบได้รับการทดสอบ ๒ ชุดใด ๆ ใช้การสุ่ม และให้เกิดความสมดุลในการจัดการ กลุ่มตัวอย่างผู้สอบมาจากการสอบในปี ๑๙๘๐ จำนวนมากกว่า ๘๐๐,๐๐๐ คน ทำการสุ่มอย่างแบ่งชั้นภูมิหลายชั้นจากประเภทโรงเรียนเขตทางภูมิศาสตร์ สถานภาพ ทางสังคม สุกท้ายได้ตัวอย่างผู้สอบรายคนในแต่ละโรงเรียน ๆ ละ ๒๒ คน จำนวนคนในแต่ละชุดของ แบบสอบที่ทำการศึกษาเท่ากับ ๒๐๐ คนโดยประมาณ ซึ่งเป็นขนาดกลุ่มตัวอย่างที่ค่อนข้างเล็ก ในการศึกษาเกี่ยวกับการเทียบมาตรา แต่ผู้วิจัยได้ให้เหตุผลว่า ถึงอย่างไรก็เป็นขนาดที่เป็นปกติของการทดสอบที่พบโดยทั่วไป นอกจากนี้การออกแบบโดยใช้กลุ่มตัวอย่างร่วมเป็นวิธีที่ทำให้เกิดความ แม่นยำได้ตามการสนับสนุนของแองกอฟ (๑๙๘๔) และลอว์ค (๑๙๘๑) และมีข้อสังเกตเพิ่มเติม ว่า หากข้อค้นพบใดที่เป็นข้อที่ตีจากการวิจัยนี้ ย่อมจะต้องคิดว่าเมื่อกลุ่มตัวอย่างมีขนาดใหญ่ขึ้น การ- ทำเนิการเทียบมาตราแต่ละชุดทั้ง ๑๑ ชุด ได้เทียบไปสู่มาตราคะแนนสังเกตของแบบสอบร่วม รูปแบบ เชิงเส้นตรง และรูปแบบอีควิเปอร์เซนไทล์ ใช้วิธีที่แองกอฟได้บรรยาย คือ IA - 1 และ IA - 2 ตามลำดับ ส่วนรูปแบบ IRT สามพารามิเตอร์ และพารามิเตอร์เดียว เริ่มด้วยการประมาณค่า



พารามิเตอร์ของข้อสอบ และความสามารถของผู้สอบด้วยโปรแกรมสำเร็จรูป LOGIST ของ Wood, Wingersky และ Lord เฉพาะแบบสอบรวมก่อน จากนั้นใช้ค่าพารามิเตอร์ความสามารถของผู้สอบกำหนดให้เป็นค่าคงที่ในการวิเคราะห์แบบสอบชุดอื่น ๆ ๑๑ ชุด เพื่อประมาณค่าพารามิเตอร์ของข้อสอบอีก ๑๑ ชุด ให้ค่าประมาณต่าง ๆ มีมาตราเดียวกัน ขึ้นต่อไปสร้างคะแนนสมมูลระหว่างแบบสอบรวม กับแบบสอบอีก ๑๑ ชุดนั้น โดยใช้การประมาณคะแนนจริงของผู้สอบซึ่งคำนวณจากผลบวกของค่าความน่าจะเป็นของการตอบถูกต้องตลอดทุกข้อในแบบสอบ สำหรับรูปแบบสามพารามิเตอร์ใช้วิธีการ linear interpolation หากค่าคะแนนสมมูลของส่วนที่อยู่ต่ำกว่าคะแนนเดา ผลการเทียบมาตราทั้ง ๔ วิธี ให้นำไปประเมินความเที่ยงพอโดยการวิเคราะห์จากกลุ่มตัวอย่างสอบทานผล ซึ่งเป็นกลุ่มอิสระมีขนาดร้อยละ ๑๐ ของกลุ่มตัวอย่างเทียบมาตรา ทุก ๆ คนในกลุ่มนี้ได้รับการทดสอบแบบสอบทั้ง ๒ ชุด วิธีที่ใช้ คือ percentile comparison index เป็นมาตรการที่ใช้วัดความแตกต่างระหว่างการแจกแจงของคะแนนในแบบสอบรวม และคะแนนแปลงของอีกชุดหนึ่ง วิธีนี้คือค่าเฉลี่ยกำลังสองของความแตกต่างของผู้สอบทุกคนในกลุ่มตัวอย่างสอบทานผลระหว่างตัวเลขคะแนนในแบบสอบรวมกับคะแนนแปลงจากแบบสอบอีกชุดหนึ่ง ณ ที่ตำแหน่งเปอร์เซ็นต์เดียวกัน วิธีนี้โคเลน (Kolen 1981) ได้พัฒนาขึ้นใช้ ผลการวิจัยได้ข้อค้นพบดังนี้ จากการประเมินผลสุดท้ายของการเทียบมาตราวิธีต่าง ๆ กับกลุ่มสอบทานผล หากการทดสอบความมีนัยสำคัญของค่าอันดับจากการแปลงตามปริมาณของค่าดัชนีด้วยการทดสอบนัยพาราเมตริก Friedman Test สรุปได้ว่า รูปแบบเชิงเส้นตรงให้ผลการเทียบที่มีความเที่ยงพอมากกว่ารูปแบบอิกวิเปอร์เซนต์ และรูปแบบ IRT สามพารามิเตอร์ ขณะเดียวกันรูปแบบ IRT พารามิเตอร์เดียว ให้ผลการเทียบที่เที่ยงพอกว่ารูปแบบ IRT สามพารามิเตอร์ ข้อสรุปเกี่ยวกับผลการวิจัยครั้งนี้ ผู้วิจัยได้ตั้งข้อสังเกตไว้หลายประการ ดังนี้ คือ การประมาณค่าพารามิเตอร์ความถ่วงน้ำหนักของสามพารามิเตอร์ ทำให้มีค่าที่ประมาณสูงมาก ซึ่งเป็นธรรมชาติของพหุคูณที่สามารถควบคุมความยากของแบบสอบ เมื่อมาใช้ในกรณีของกลุ่มตัวอย่างขนาดเล็กเกินไป ทำให้เกิดความผันแปรเชิงสุ่มจนทำให้ผลการเทียบมาตราขาดความแน่นอนไป แต่จุดอ่อนนี้สามารถแก้ไขได้โดยเพิ่มตัวอย่างให้มีขนาดใหญ่ขึ้นตามคำแนะนำของลอร์ด (Lord 1980: 209-210) รูปแบบอิกวิเปอร์เซนต์ให้ผลการเทียบมาตราดีกว่ารูปแบบเชิงเส้นตรง เนื่องจากตัวอย่างมีขนาดเล็กเกินไป ทำให้คะแนนแปลงมีลักษณะผิดปกติมาก ซึ่งสันนิษฐานว่ามาจากความคลาดเคลื่อนเชิงสุ่ม การเกลาเส้นให้เรียบในส่วนที่มีลักษณะผิดปกติได้ใช้วิธีเกลาค้วยมือ ทำให้ความคลาด



เคลื่อนที่ปรากฏในกลุ่มสอบทานผลน้อยลง แต่ก็ยังไม่สามารถปรับปรุงได้มากกว่านั้น จึงได้ลงความเห็นว่า ตามสถานการณ์ที่ศึกษาในกรณีของแบบสอบ QED นี้ เมื่อออกแบบด้วยการใช้กลุ่มตัวอย่างร่วม หรืออาจจะเป็นกลุ่มตัวอย่างสุ่มสมบูรณ์แล้ว รูปแบบที่เหมาะสมที่สุด คือ วิธีการเชิงเส้นตรงเป็นวิธีที่ง่าย และให้ผลการเทียบที่แน่นอนมากกว่า

สรุปการวิเคราะห์วรรณคดีที่เกี่ยวข้องกับการ เปรียบเทียบรูปแบบการเทียบมาตรา แยกได้เป็น ๒ ประเด็น คือ ประเด็นของผลการศึกษา และประเด็นของวิธีการทางทัศนการ เปรียบเทียบประเด็นของผลการศึกษา จากการวิจัยที่เริ่มมาตั้งแต่ปี ๑๙๗๓ จนถึง ๑๙๘๒ สรุปได้ว่า ไม่มีรูปแบบการเทียบมาตราใดที่ให้ผลในการประยุกต์ใช้ดีที่สุดกับทุกสถานการณ์ที่กำหนด แต่จากผลการวิจัยที่ดำเนินการมานั้นได้ให้ข้อความจริงจากการค้นพบในเชิงประจักษ์ที่พอจะเป็นแนวทาง เพื่อเลือกใช้ให้ใกล้เคียง เหมาะสมกับสภาพที่ท้องถื่นของแต่ละกรณี การพิจารณาทางเลือกที่สำคัญก่อนตัดสินใจเลือกรูปแบบที่เป็นหลักในการกำหนดคะแนนสมบูรณ์ คือ ลักษณะของแบบสอบ ๒ ชุดที่ท้องถื่นเทียบ หรือมากกว่า ๒ ชุด หรือการเทียบอยู่ในลักษณะหลายชุดเทียบกันเป็นสายโซ่ ลักษณะของกลุ่มตัวอย่างผู้สอบ ซึ่งเป็นเรื่องของการออกแบบรวบรวมข้อมูลคะแนนเพื่อสร้างคะแนนสมบูรณ์ นอกจากนี้ยังมีตัวแปรปลี่ยนย่อยซึ่งเป็นข้อจำกัดเฉพาะกรณี ข้อสอบอย่างกว้าง ๆ มีดังนี้ (๑) เมื่อแบบสอบที่นำมาเทียบมีความยากแตกต่างกัน กลุ่มตัวอย่างไม่ได้เป็นกลุ่มสมบูรณ์ ซึ่งส่วนมากเป็นการเทียบมาตราในแนวตั้ง (vertical equating) รูปแบบ IRT สามพารามิเตอร์ให้ผลน่าพอใจกว่ารูปแบบอื่น (Slindt and Linn 1979; Petersen, Cook and Stocking 1981; Kolen 1981; Cook, Dunbar, Eignor 1981) (๒) เมื่อแบบสอบมีความคล้ายคลึงในระดับความยาก และกลุ่มตัวอย่างเป็นกลุ่มตัวอย่างร่วม หรือกลุ่มสุ่มสมบูรณ์ ควรใช้วิธีการของรูปแบบเชิงเส้นตรง ซึ่งสะดวกและให้ความแน่นอนกว่า (Petersen, Cook and Stocking 1981; Kolen, Whitney 1982) (๓) เมื่อแบบสอบมีความยากต่างกัน แต่กลุ่มตัวอย่างเป็นกลุ่มสุ่มสมบูรณ์แล้ว วิธีการของรูปแบบอิคิวเปอร์เซนโตล ซึ่งให้ความสัมพันธ์ของการแปลงที่ไม่เป็นเส้นตรงให้ความเหมาะสมต่อการแปลงได้ แต่ถ้าไม่เป็นกลุ่มสมบูรณ์บางครั้งอาจเกิดการแจกแจงคะแนนที่แตกต่างมาก จนทำให้เกิดความคลาดเคลื่อนที่ทำให้ผลการเทียบขาดความเพียงพอของการสร้างคะแนนสมบูรณ์ (Cook, Dunbar, Eignor 1981; Kolen 1982) (๔) ในสภาพข้อมูลที่ค่อนข้างน้อย และแบบสอบโดยส่วนใหญ่

ทำหน้าที่วัดในเรื่องเดียวกัน รูปแบบเชิงเส้นตรงและ IRT พารามิเตอร์เดียว จะให้ผลการเทียบที่เพียงพอว่ารูปแบบอควิปเปอร์เซนไทล์กับ IRT สามพารามิเตอร์ ประเด็นของวิธีการพิจารณาทัศนการเปรียบเทียบที่ใช้ในการศึกษาทั้งหมดเป็นเกณฑ์ภายนอก (external criteria) โดยใช้พิจารณาจากความคลาดเคลื่อน (discrepancy) ระหว่างคะแนนแปลงที่เกิดจากการใช้ผลการเทียบมาตราวิธีที่ของการศึกษากับคะแนนเกณฑ์ ค่าความคลาดเคลื่อนใช้เป็นดัชนีชี้ความแตกต่างเชิงปริมาณ การทัศนความเพียงพอของวิธีเทียบมาตราวิธีหลักของการเปรียบเทียบค่าดัชนีความคลาดเคลื่อนเป็นการทัศนด้วยความเพียงพอสัมพัทธ์ (relative adequacy) รายละเอียดของเทคนิคการวิเคราะห์แตกต่างกันในแต่ละการวิจัยดังนี้ (๑) ดัชนีความคลาดเคลื่อน (discrepancy index) เป็นค่าที่ถ่วงน้ำหนักแล้วของค่าเฉลี่ยกำลังสองของความแตกต่างระหว่างคะแนนที่ประมาณกับคะแนนเกณฑ์ คำนี้นำการถ่วงน้ำหนักให้เป็นมาตรฐานด้วยความเบี่ยงเบนมาตรฐานของการแจกแจงในคะแนนเกณฑ์ ซึ่งจะทำให้ค่าดัชนีที่ได้จากการเทียบมาตราต่างกลุ่มผู้สอบ หรือต่างชุดของแบบสอบสามารถเปรียบเทียบได้โดยตรง ค่าตัวเลขของดัชนีที่น้อยให้ความหมายว่า วิธีการเทียบมาตราที่ใช้กับสภาพการทดสอบนั้น ๆ มีความแน่นอนในการแปลงคะแนน (stability) สำหรับการวิจัยที่ใช้วิธีการเปรียบเทียบนี้ ไคแก่ ปีเตอร์สัน มาร์โค และ สตีเวอร์ท (Petersen Marco and Stewart 1982) ซึ่งใช้คะแนนเกณฑ์ที่เป็นคะแนนแปลงโดยรูปแบบ IRT กับข้อมูลคะแนนเดียวกันกับการแปลงด้วยวิธีอื่น ๆ ที่ของการศึกษา การใช้ผลการเทียบด้วย IRT เป็นเกณฑ์ ถึงแม้ผู้วิจัยจะได้อ้างอิงผลการวิจัยต่าง ๆ ให้เป็นเหตุผลที่หนักแน่นขึ้นก็ตาม แต่ผู้วิจัยเองก็ไม่กล้าให้คำยืนยันว่าเมื่อเปลี่ยนสถานการณ์ทดสอบออกไปจากที่ศึกษาแล้ว ความลำเอียงจะไม่เกิดขึ้น ปีเตอร์สัน คูก และ สตอกคิง (Petersen, Cook and Stocking 1981) ได้เปลี่ยนเกณฑ์เป็นคะแนนดิบของแบบสอบชุดแรกในกระบวนการเทียบมาตราแบบสายโซ่ แบบสอบชุดหลังแต่ละครั้งจะเทียบไปยังชุดก่อนหน้านั้น หลังจากผ่านกระบวนการเทียบไป ๒ ถึง ๓ ครั้ง จะนำคะแนนแปลงไปสู่มาตราของชุดแรกมาเทียบกับคะแนนดั้งเดิมของชุดแรกซึ่งใช้เป็นคะแนนเกณฑ์ ความแตกต่างที่เกิดขึ้นจะให้เป็นค่าดัชนีความแตกต่าง คำนีจะให้ความหมายของความคงเส้นคงวาของการแปลง (consistency) ในกระบวนการที่เป็นสายโซ่ ส่วนงานวิจัยของโคเลน (Kolen 1981) ได้ใช้เทคนิคการวิเคราะห์กลุ่มสอบทานผลเพื่อสร้างเกณฑ์ที่ปราศจากความลำเอียง การวิเคราะห์กลุ่มตัวอย่างสอบทานผล เป็นวิธีตรวจสอบความเพียงพอของรูปแบบการเทียบมาตราตอนปลายทาง โดยใช้กลุ่มตัวอย่าง

อิสระจากประชากรเกี่ยวกับกลุ่มเทียบมาตรา พิจารณาจากความแตกต่างระหว่างการแจกแจงของคะแนนในแบบสอบชุดแรกกับคะแนนแปลงของแบบสอบชุดหลังในมาตราของชุดแรก ความแตกต่างในแนวคิดนี้ คือ ค่าเฉลี่ยกำลังสองของความแตกต่างของผู้สอบทุกคนในกลุ่มตัวอย่างสอบทานผลระหว่างตัวเลขคะแนนของแบบสอบชุดแรกกับตัวเลขของคะแนนแปลงที่ตำแหน่งเปอร์เซ็นต์เดียวกัน โคลเลน (๑๘๘๑) เรียกดัชนีนี้ว่า percentile comparison index โคลเลน และ วิทนี (Kolen and Whitney 1982) ได้ปรับปรุงดัชนีนี้เพื่อให้ใช้เปรียบเทียบกับ การเทียบมาตราที่เกิดจากการใช้แบบสอบที่มีความยาวของแบบสอบต่างกัน โดยการถ่วงน้ำหนักด้วยจำนวนข้อสอบ ( $k$ ) ซึ่งช่วยให้ค่าดัชนีสามารถนำมาเปรียบเทียบ เมื่อเปลี่ยนชุดของแบบสอบ

จากการวิเคราะห์ห้วงเวลาที่เกี่ยวกับการวิจัยเปรียบเทียบวิธีการเทียบมาตราทั้งกล่าวมาแล้วนี้ ผู้วิจัยได้พัฒนาวิธีการที่จะใช้ประเมินเพื่อตัดสินวิธีการเทียบมาตรา ๓ รูปแบบ จากแนวคิดของการใช้วิเคราะห์กลุ่มสอบทานผล ดัชนีความแตกต่างที่ใช้ คือ ดัชนีความแตกต่างมาตรฐาน (index of standard discrepancy) ซึ่งได้จากค่าเฉลี่ยกำลังสองความแตกต่างของคะแนนแปลงกับคะแนนเกณฑ์ของบุคคลเกี่ยวกับทุกคนในกลุ่มตัวอย่างสอบทานผล และได้ถ่วงน้ำหนักให้เป็นมาตรฐานด้วยส่วนเบี่ยงเบนมาตรฐานของการแจกแจงของคะแนนเกณฑ์ สำหรับคะแนนเกณฑ์ที่ใช้ คือ คะแนนการสังเกตของแบบสอบชุด  $x$  ที่กลุ่มตัวอย่างสอบทานผลทำได้ ค่านี้ให้ความหมายของการเปรียบเทียบเชิงปริมาณไม่ว่าจะเป็นผลจากการทดสอบด้วยแบบสอบใด หรือรูปแบบการเทียบใด ค่าที่น้อยกว่าหมายถึง รูปแบบการทดสอบนั้น ๆ มีความเพียงพอในการแปลงคะแนนแบบสอบชุดหลังให้อยู่ในคะแนนมาตราของชุดแรก การลงข้อสรุปใช้หลักการเปรียบเทียบ และทดสอบว่าการทดสอบนั้นพาราเมทริก เพื่อทดสอบความมีนัยสำคัญทางสถิติของดัชนีความแตกต่างด้วยนันทพาราเมทริก

จุฬาลงกรณ์มหาวิทยาลัย