

บทที่ 5

สรุปผลการวิจัย อภิปรายผล และข้อเสนอแนะ

การวิจัยครั้งนี้มีวัตถุประสงค์ 2 ประการ ประการแรก เพื่อพัฒนาเกณฑ์ตัดสินคุณภาพการเปรียบเทียบคะแนนตามทฤษฎีการตอบสนองข้อสอบ ด้วยข้อมูลจำลองตามเงื่อนไขการเปรียบเทียบคะแนน โมเดลการตอบสนองข้อสอบ แบบแผนการเก็บรวบรวมข้อมูล และการหาคุณภาพของการเปรียบเทียบ ประการที่สอง เพื่อตรวจสอบคุณภาพของเกณฑ์ที่พัฒนาขึ้น ในด้านความตรงเชิงเกณฑ์สัมพันธ์ โดยตรวจสอบความสอดคล้องของผลการตัดสินด้วยเกณฑ์ที่พัฒนาขึ้นกับเกณฑ์ของปีเตอร์เซนและคณะ และความสอดคล้องของผลการตัดสินด้วยเกณฑ์ที่พัฒนาขึ้นกับเกณฑ์ความเสมอภาคของลอร์ด

ข้อมูลที่ใช้ในการวิจัยเป็นข้อมูลจำลองจากโปรแกรม IRTDATA เมื่อจำลองข้อมูลตามเงื่อนไขการเปรียบเทียบคะแนนที่อาจเกิดขึ้นได้เกือบทั้งหมดดังนี้คือ โมเดล 1 และ 3 พารามิเตอร์ แบบแผนการเก็บรวบรวมข้อมูลแบบใช้ข้อสอบร่วมและใช้กลุ่มสมมูล การหาคุณภาพการเปรียบเทียบคะแนนด้วยการเปรียบเทียบกลับสู่แบบสอบเดิมและการใช้กลุ่มสอบแทนผล โดยกำหนดเงื่อนไขในการจำลองข้อมูลเป็นจำนวนข้อสอบ 10, 20, 30, ..., 200 ข้อ และกลุ่มผู้สอบจำนวน 100, 200, 300, ..., 3,000 คน

วิธีดำเนินการวิจัยมีขั้นตอนดังนี้ คือ 1) ตรวจสอบความถูกต้องโปรแกรม IRTDATA โดยพิจารณาจากการนำข้อมูลที่เป็นคะแนนตอบรายข้อที่จำลองจากโปรแกรม IRTDATA ไปวิเคราะห์ด้วยโปรแกรม BILOG แล้วนำค่าพารามิเตอร์ผู้ตอบและค่าพารามิเตอร์ข้อสอบไปเทียบกับผลจากโปรแกรม IRTDATA ว่าแตกต่างกันหรือไม่ และตรวจสอบความเป็นเอกมิติ (Unidimensionality) ของข้อมูลด้วยการวิเคราะห์ตัวประกอบ (Factor Analysis) ด้วยวิธีการวิเคราะห์ตัวประกอบสำคัญ (Principal Component Analysis) และหมุนแกนด้วยวิธี Varimax 2) จำลองข้อมูลตามแบบแผนที่กำหนด 3) เปรียบเทียบพารามิเตอร์ข้อสอบจากแบบสอบฉบับที่ 1 ผ่านแบบสอบฉบับที่ 2 และแบบสอบฉบับที่ 3 แล้วเปรียบเทียบกลับสู่แบบสอบฉบับที่ 1 เดิม สำหรับการเปรียบเทียบกลับสู่แบบสอบเดิม และเปรียบเทียบพารามิเตอร์ข้อสอบจากแบบสอบฉบับที่ 1 ไปสู่แบบสอบฉบับที่ 2 สำหรับกลุ่มสอบแทนผล 4) คำนวณคะแนนจริง (True Score) จากทั้งพารามิเตอร์ข้อสอบจากแบบสอบฉบับที่ 1 ที่ไม่ได้เปรียบเทียบกับพารามิเตอร์ข้อสอบที่เปรียบเทียบแล้ว เมื่อใช้พารามิเตอร์ผู้สอบจากกลุ่มผู้สอบแบบสอบฉบับที่ 1 สำหรับการเปรียบเทียบกลับสู่แบบสอบเดิม และใช้พารามิเตอร์ผู้ตอบจากกลุ่มสอบแทนผลสำหรับการใช้กลุ่มสอบแทนผล 5) ตรวจสอบความแตกต่างระหว่างค่าเฉลี่ยคะแนนจริงที่ไม่ได้เปรียบเทียบกับค่าเฉลี่ยของคะแนนจริงที่เปรียบเทียบแล้ว ด้วยสถิติ t-test แบบ Two Dependent Samples เพื่อหาค่าวิกฤต t-test ที่นัยสำคัญ

ทางสถิติระดับ .01 6) กำหนดจุดตัดโดยคำนวณค่าดัชนีความแตกต่างเมื่อทดสอบความแตกต่างค่าเฉลี่ยของคะแนนจริงทั้งสองด้วยสถิติ t-test แล้วพบค่าวิกฤตที่นัยสำคัญทางสถิติระดับ .01 จุดตัดนี้จะแบ่งกลุ่มค่าดัชนีความแตกต่างออกเป็น 2 กลุ่ม คือ กลุ่มดัชนีที่แสดงคุณภาพการปรับเทียบระดับสูง และกลุ่มดัชนีที่แสดงคุณภาพการปรับเทียบระดับต่ำ จำแนกตามเงื่อนไขการปรับเทียบคะแนนเป็น 7 เงื่อนไข ดังนี้ คือ รวมทุกเงื่อนไข โมเดล 1 พารามิเตอร์ โมเดล 3 พารามิเตอร์ การใช้กลุ่มสมมูล การใช้ข้อสอบร่วม การปรับเทียบกลับสู่แบบสอบเดิม และการใช้กลุ่มสอบทานผล 7) กำหนดจุดตัดซึ่งได้จากการคำนวณค่าดัชนีความแตกต่าง เมื่อพบค่าวิกฤต t-test ที่นัยสำคัญทางสถิติ .02 จากการทดสอบความแตกต่างค่าเฉลี่ยของคะแนนจริงทั้งสองด้วยสถิติ t-test จุดตัดนี้จะแบ่งกลุ่มดัชนีที่แสดงคุณภาพการปรับเทียบระดับสูงออกเป็น 2 กลุ่มย่อย และจุดตัดที่คำนวณค่าดัชนีความแตกต่างเมื่อพบค่าวิกฤตที่นัยสำคัญทางสถิติระดับ .001 เป็นจุดตัดที่แบ่งกลุ่มดัชนีที่แสดงคุณภาพการปรับเทียบระดับต่ำออกเป็น 2 กลุ่มย่อย 8) ให้ความหมายแต่ละกลุ่มดัชนีความแตกต่างกำหนดเป็นเกณฑ์ 9) ตรวจสอบความสอดคล้องของผลการตัดสินคุณภาพการปรับเทียบคะแนนด้วยเกณฑ์ที่พัฒนาขึ้นกับเกณฑ์ของปีเตอร์เซนและคณะ และความสอดคล้องของผลการตัดสินคุณภาพการปรับเทียบคะแนนด้วยเกณฑ์ที่พัฒนาขึ้นกับเกณฑ์ความเสมอภาคของลอร์ด

สรุปผลการวิจัย

ผลการวิจัยในการศึกษาครั้งนี้ สรุปได้ดังนี้

1. ผลการตรวจสอบความถูกต้องโปรแกรม IRTDATA พบว่า

1.1 ค่าพารามิเตอร์ผู้ตอบ ค่าพารามิเตอร์ข้อสอบจากการจำลองข้อมูลด้วยโปรแกรม IRTDATA และจากการวิเคราะห์หัดด้วยโปรแกรม BILOG เมื่อใช้ข้อมูลคะแนนการตอบรายข้อชุดเดียวกันไม่แตกต่างกัน โดยมีการจำลองข้อมูลทั้งโมเดลการตอบสนองข้อสอบ 1 พารามิเตอร์ และโมเดล 3 พารามิเตอร์ แสดงให้เห็นว่าข้อมูลที่จำลองจากโปรแกรม IRTDATA มีความเชื่อถือได้

1.2 ค่าไอเกน (Eigenvalue) สูงสุดของแบบสอบ คือ 11.712 และคิดเป็นร้อยละ 29.28 ของความแปรปรวนทั้งหมดสำหรับข้อมูลจำลองจากโปรแกรม IRTDATA ตามโมเดลการตอบสนองข้อสอบ 1 พารามิเตอร์ และค่าไอเกนสูงสุดของแบบสอบ คือ 8.734 คิดเป็นร้อยละ 21.834 ของความแปรปรวนทั้งหมด สำหรับข้อมูลจำลองตามโมเดลการตอบสนองข้อสอบ 1 พารามิเตอร์ พิจารณาแล้วเห็นว่ามีความพอที่จะสรุปว่ามีความเป็นเอกมิติประเภทที่มีตัวประกอบหลักเด่นกว่าตัวประกอบอื่น (Essential Unidimensionality) สรุปได้ว่าข้อมูลที่ได้จากการใช้โปรแกรม IRTDATA มีความเป็นเอกมิติเป็นไปตามข้อตกลงของทฤษฎีการตอบสนองข้อสอบ

2. เกณฑ์ตัดสินคุณภาพการเปรียบเทียบคะแนนตามทฤษฎีการตอบสนองข้อสอบ
ที่พัฒนาขึ้น มี 7 เกณฑ์ดังนี้

2.1 เกณฑ์ตัดสินคุณภาพการเปรียบเทียบคะแนนตามทฤษฎีการตอบสนองข้อสอบ
เกณฑ์รวมเงื่อนไขทั้งหมด

ระดับคุณภาพการเปรียบเทียบ	ดัชนีความแตกต่าง
น่าพอใจอย่างยิ่ง	AMD < 0.000183
	MAD < 0.002391
	RMS < 0.002907
น่าพอใจ	0.000183 ≤ AMD < 0.000244
	0.002391 ≤ MAD < 0.002481
	0.002907 ≤ RMS < 0.003032
ไม่น่าพอใจ	0.000244 ≤ AMD < 0.000352
	0.002481 ≤ MAD < 0.002621
	0.003032 ≤ RMS < 0.003181
ไม่น่าพอใจอย่างยิ่ง	0.000352 ≤ AMD
	0.002621 ≤ MAD
	0.003181 ≤ RMS

2.2 เกณฑ์ตัดสินคุณภาพการเปรียบเทียบคะแนนสำหรับโมเดล 1 พารามิเตอร์

ระดับคุณภาพการเปรียบเทียบ	ดัชนีความแตกต่าง
นำพอใจอย่างยิ่ง	AMD < 0.000182
	MAD < 0.002450
	RMS < 0.002924
นำพอใจ	0.000182 ≤ AMD < 0.000261
	0.002450 ≤ MAD < 0.002508
	0.002924 ≤ RMS < 0.003045
ไม่นำพอใจ	0.000261 ≤ AMD < 0.000355
	0.002508 ≤ MAD < 0.002574
	0.003045 ≤ RMS < 0.003073
ไม่นำพอใจอย่างยิ่ง	0.000355 ≤ AMD
	0.002574 ≤ MAD
	0.003073 ≤ RMS

2.3 เกณฑ์ตัดสินคุณภาพการเปรียบเทียบคะแนนตามทฤษฎีการตอบสนองข้อสอบ
เกณฑ์สำหรับโมเดล 3 พารามิเตอร์

ระดับคุณภาพการเปรียบเทียบ	ดัชนีความแตกต่าง
นำพอใจอย่างยิ่ง	AMD < 0.000183
	MAD < 0.002402
	RMS < 0.002924
นำพอใจ	0.000183 ≤ AMD < 0.000242
	0.002402 ≤ MAD < 0.002468
	0.002924 ≤ RMS < 0.003122
ไม่นำพอใจ	0.000242 ≤ AMD < 0.000355
	0.002468 ≤ MAD < 0.002613
	0.003122 ≤ RMS < 0.003398
ไม่นำพอใจอย่างยิ่ง	0.000355 ≤ AMD
	0.002613 ≤ MAD
	0.003398 ≤ RMS

2.4 เกณฑ์สำหรับแบบแผนการเก็บรวบรวมข้อมูลใช้กลุ่มสมมูล

ระดับคุณภาพการเปรียบเทียบ	ดัชนีความแตกต่าง
น่าพอใจอย่างยิ่ง	AMD < 0.000197
	MAD < 0.002353
	RMS < 0.002924
น่าพอใจ	0.000197 ≤ AMD < 0.000237
	0.002353 ≤ MAD < 0.002375
	0.002924 ≤ RMS < 0.002930
ไม่น่าพอใจ	0.000237 ≤ AMD < 0.000374
	0.002375 ≤ MAD < 0.002510
	0.002930 ≤ RMS < 0.003022
ไม่น่าพอใจอย่างยิ่ง	0.000374 ≤ AMD
	0.002510 ≤ MAD
	0.003022 ≤ RMS

2.5 เกณฑ์สำหรับแบบแผนการเก็บรวบรวมข้อมูลใช้ข้อสอบร่วม

ระดับคุณภาพการเปรียบเทียบ	ดัชนีความแตกต่าง
น่าพอใจอย่างยิ่ง	AMD < 0.000176
	MAD < 0.002410
	RMS < 0.002899
น่าพอใจ	0.000176 ≤ AMD < 0.000247
	0.002410 ≤ MAD < 0.002535
	0.002899 ≤ RMS < 0.003083
ไม่น่าพอใจ	0.000247 ≤ AMD < 0.000341
	0.002535 ≤ MAD < 0.002677
	0.003083 ≤ RMS < 0.003261
ไม่น่าพอใจอย่างยิ่ง	0.000341 ≤ AMD
	0.002677 ≤ MAD
	0.003261 ≤ RMS

2.6 เกณฑ์สำหรับการเปรียบเทียบกลับสู่แบบสอบถาม

ระดับคุณภาพการเปรียบเทียบ	ดัชนีความแตกต่าง
น่าพอใจอย่างยิ่ง	AMD < 0.000169
	MAD < 0.002370
	RMS < 0.002875
น่าพอใจ	0.000169 ≤ AMD < 0.000261
	0.002370 ≤ MAD < 0.002508
	0.002875 ≤ RMS < 0.003045
ไม่น่าพอใจ	0.000261 ≤ AMD < 0.000336
	0.002508 ≤ MAD < 0.002637
	0.003045 ≤ RMS < 0.003124
ไม่น่าพอใจอย่างยิ่ง	0.000336 ≤ AMD
	0.002637 ≤ MAD
	0.003124 ≤ RMS

2.7 เกณฑ์สำหรับการเปรียบเทียบกลับสู่แบบสอบถาม

ระดับคุณภาพการเปรียบเทียบ	ดัชนีความแตกต่าง
น่าพอใจอย่างยิ่ง	AMD < 0.000189
	MAD < 0.002402
	RMS < 0.002924
น่าพอใจ	0.000189 ≤ AMD < 0.000235
	0.002402 ≤ MAD < 0.002468
	0.002924 ≤ RMS < 0.003026
ไม่น่าพอใจ	0.000235 ≤ AMD < 0.000361
	0.002468 ≤ MAD < 0.002613
	0.003026 ≤ RMS < 0.003210
ไม่น่าพอใจอย่างยิ่ง	0.000361 ≤ AMD
	0.002613 ≤ MAD
	0.003210 ≤ RMS

3. ผลการตรวจสอบคุณภาพของเกณฑ์ที่พัฒนาขึ้น มีดังต่อไปนี้

3.1 ผลการตัดสินคุณภาพการเปรียบเทียบคะแนนเมื่อใช้เกณฑ์ที่พัฒนาขึ้นทั้ง 7 เกณฑ์ เปรียบเทียบกับผลการใช้ตามเกณฑ์ของปีเตอร์เซนและคณะ พบว่าไม่สอดคล้องกัน โดยผลการตัดสินคุณภาพการเปรียบเทียบทั้ง 4 ระดับจากเกณฑ์ที่พัฒนาขึ้น คือระดับคุณภาพการปรับเทียบน่าพอใจอย่างยิ่ง คุณภาพการปรับเทียบน่าพอใจ คุณภาพการปรับเทียบไม่น่าพอใจ และคุณภาพการปรับเทียบไม่น่าพอใจอย่างยิ่ง อยู่ในระดับน่าพอใจอย่างยิ่งเมื่อเทียบกับเกณฑ์ของปีเตอร์เซนและคณะ แสดงว่าเกณฑ์ที่พัฒนาขึ้นให้ความคลาดเคลื่อนของการปรับเทียบคะแนนน้อยกว่าเกณฑ์ของปีเตอร์เซนและคณะ

3.2 ผลการตัดสินคุณภาพการปรับเทียบคะแนนเมื่อใช้เกณฑ์ที่พัฒนาขึ้นทั้ง 7 เกณฑ์เปรียบเทียบกับผลการใช้เกณฑ์ความเสมอภาคของลอร์ด พบว่ามีความสอดคล้องกัน โดยผลการตัดสินคุณภาพการปรับเทียบระดับน่าพอใจอย่างยิ่งหรือระดับน่าพอใจและระดับไม่น่าพอใจหรือระดับไม่น่าพอใจอย่างยิ่งจากการใช้เกณฑ์ที่พัฒนาขึ้น อยู่ในระดับคุณภาพน่าพอใจและระดับคุณภาพไม่น่าพอใจตามลำดับ เมื่อใช้เกณฑ์ความเสมอภาคของลอร์ด แสดงว่าเกณฑ์ทั้งสองให้ผลการตัดสินคุณภาพการปรับเทียบคะแนนเป็นไปในทิศทางเดียวกัน

อภิปรายผลการวิจัย

จากผลการวิจัยที่กล่าวมาข้างต้นสามารถอภิปรายได้ใน 3 ประเด็น ได้แก่ ประเด็นที่หนึ่ง เกณฑ์ตัดสินคุณภาพการปรับเทียบคะแนนตามทฤษฎีการตอบสนองข้อสอบที่พัฒนาขึ้น ประเด็นที่สอง ผลการตัดสินคุณภาพการปรับเทียบคะแนนตามทฤษฎีการตอบสนองข้อสอบเมื่อใช้เกณฑ์ที่พัฒนาขึ้นกับเกณฑ์ของปีเตอร์เซนและคณะ และสุดท้ายประเด็นที่สาม ผลการตัดสินคุณภาพการปรับเทียบคะแนนตามทฤษฎีการตอบสนองข้อสอบ เมื่อใช้เกณฑ์ที่พัฒนาขึ้นกับเกณฑ์ความเสมอภาคของลอร์ด โดยมีรายละเอียดในแต่ละประเด็น ดังนี้

ประเด็นที่ 1 เกณฑ์ตัดสินคุณภาพการปรับเทียบคะแนนตามทฤษฎีการตอบสนองข้อสอบที่พัฒนาขึ้น

เกณฑ์ที่พัฒนาขึ้นในครั้งนี้มี 7 เกณฑ์ มีความหลากหลายเพื่อให้ผู้ดำเนินการปรับเทียบคะแนนได้เลือกใช้ให้เหมาะสมกับสถานการณ์การปรับเทียบคะแนน ดังเช่น แฮร์ริส และ คราวส์ (Harris and Crouse, 1993) กล่าวว่าไม่มีเกณฑ์ที่สมบูรณ์ที่สุดที่จะใช้ตัดสินผลการปรับเทียบคะแนนได้ทุกเงื่อนไขการปรับเทียบคะแนน และเมื่อใช้เกณฑ์ที่ต่างกันแล้วจะทำให้ผลการตัดสิน

คุณภาพการเปรียบเทียบต่างกันด้วย (Skaggs, 1990b; Livingston and others cited in Harris and Crouse, 1993) ฉะนั้นผู้ใช้ผลการเปรียบเทียบคะแนนต้องพิจารณาเลือกใช้เกณฑ์ให้สอดคล้องกับกระบวนการเปรียบเทียบคะแนน ทำให้ผลตัดสินคุณภาพการเปรียบเทียบเป็นที่ยอมรับได้ (Harris and Crouse, 1993)

เกณฑ์ที่พัฒนาขึ้นทั้ง 7 เกณฑ์ เป็นกลุ่มค่าดัชนีความแตกต่างเรียงลำดับจากน้อยไปหามาก 4 กลุ่ม กลุ่มค่าดัชนีเหล่านี้แสดงถึงระดับคุณภาพการเปรียบเทียบคะแนน 4 ระดับ โดยกลุ่มแรกแสดงถึงระดับคุณภาพการเปรียบเทียบน่าพอใจอย่างยิ่ง ค่าดัชนีในกลุ่มที่ 2 แสดงถึงระดับคุณภาพการเปรียบเทียบน่าพอใจ ค่าดัชนีในกลุ่มที่ 3 แสดงถึงระดับคุณภาพการเปรียบเทียบไม่น่าพอใจ และกลุ่มดัชนีแสดงถึงระดับคุณภาพการเปรียบเทียบไม่น่าพอใจอย่างยิ่ง โดยจุดตัดที่ใช้แบ่งกลุ่มดัชนี คือ ค่าดัชนีความแตกต่างระหว่างคะแนนจริงที่เปรียบเทียบแล้วกับคะแนนจริงที่ยังไม่ได้เปรียบเทียบ เมื่อพบค่าวิกฤตจากการทดสอบความแตกต่างค่าเฉลี่ยของคะแนนจริงทั้งสองด้วยสถิติ t-test แบบ Two Dependent Sample Test ซึ่งจุดตัดมีความสำคัญกับการกำหนดช่วงของเกณฑ์ (Glaser and Nitko, 1971) การกำหนดจุดตัดที่มีหลักการและมีความเชื่อถือได้ จะทำให้เกณฑ์มีคุณภาพด้วยเพราะจุดตัดจะแบ่งช่วงของเกณฑ์ได้อย่างถูกต้องและเหมาะสม

เกณฑ์ที่พัฒนาขึ้นเป็นช่วงของค่าดัชนี AMD (Absolute Mean Difference), MAD (Mean Absolute Difference) และ RMS (Root Mean Square) ที่บอกระดับคุณภาพการเปรียบเทียบคะแนนตามทฤษฎีการตอบสนองข้อสอบ ซึ่งเป็นดัชนีที่นิยมใช้ในกระบวนการเปรียบเทียบคะแนน (Harris and Crouse, 1993) เพราะเป็นดัชนีที่คำนวณได้ง่ายไม่ซับซ้อน แปลผลชัดเจน โดยค่าดัชนีที่มีค่าน้อยจะบอกถึงคุณภาพการเปรียบเทียบที่ดี เมื่อใช้ข้อมูลชุดเดียวกัน ดัชนี RMS จะให้ค่ามากที่สุด รองลงมาคือดัชนี MAD และดัชนี AMD ตามลำดับ

ค่าดัชนีความแตกต่างที่กำหนดเป็นเกณฑ์มีค่าน้อย เช่น ดัชนี AMD มีค่าตั้งแต่ 0.000056 ถึง 0.270909 ดัชนี MAD มีค่าตั้งแต่ 0.002223 ถึง 0.270909 และดัชนี RMS มีค่าตั้งแต่ 0.002704 ถึง 0.280773 เป็นเพราะค่าดัชนีเหล่านี้คำนวณจากความแตกต่างระหว่างคะแนนจริงที่เปรียบเทียบแล้วกับคะแนนจริงที่ยังไม่ได้เปรียบเทียบ จากผู้สอบกลุ่มเดียวกัน เนื่องจากดัชนีความแตกต่างมีค่าน้อยจะแสดงถึงการเกิดความคลาดเคลื่อนจากการเปรียบเทียบคะแนนมีน้อย สอดคล้องกับคำกล่าวของ แฮร์ริส และ เคิร์ส (Harris and Crouse, 1993) ที่ว่าในการประเมินผลการเปรียบเทียบคะแนนเมื่อใช้ดัชนีความแตกต่าง เช่น ดัชนี RMS ผู้ใช้ผลการเปรียบเทียบต้องการที่จะให้ได้ค่าดัชนีความแตกต่างที่มีค่าน้อย ๆ

ประเด็นที่ 2 ผลการตัดสินคุณภาพการเปรียบเทียบคะแนนตามทฤษฎีการตอบสนอง ข้อสอบ เมื่อใช้เกณฑ์ที่พัฒนาขึ้นกับเกณฑ์ของปีเตอร์เซนและคณะ

เกณฑ์ที่พัฒนาขึ้นเป็นช่วงของค่าดัชนีความแตกต่างที่บอกระดับคุณภาพการเปรียบเทียบคะแนนตามทฤษฎีการตอบสนองข้อสอบ 4 ระดับ คือ ระดับคุณภาพเปรียบเทียบน่าพอใจอย่างยิ่ง คุณภาพการเปรียบเทียบน่าพอใจ คุณภาพการเปรียบเทียบไม่น่าพอใจ และคุณภาพการเปรียบเทียบไม่น่าพอใจอย่างยิ่ง ส่วนเกณฑ์ของปีเตอร์เซนและคณะเป็นช่วงของกำลังสองของผลคูณระหว่างค่าคงที่กับความแปรปรวนของคะแนน มี 5 ช่วง แต่ละช่วงบอกระดับความพึงพอใจในการเปรียบเทียบคะแนน คือ ระดับน่าพอใจอย่างมาก ระดับน่าพอใจ ระดับปานกลาง ระดับไม่น่าพอใจ และระดับไม่น่าพอใจอย่างยิ่ง เมื่อต้องการทราบคุณภาพการเปรียบเทียบคะแนน สามารถคำนวณค่าดัชนีความแตกต่างแล้วพิจารณาว่าดัชนีตกอยู่ในช่วงใด ซึ่งจะบอกระดับคุณภาพการเปรียบเทียบคะแนนตามต้องการได้ จากการจำลองข้อมูลตามเงื่อนไขการปรับเทียบคะแนน แล้ววิเคราะห์หาคุณภาพการปรับเทียบคะแนนโดยใช้เกณฑ์ที่พัฒนาขึ้นทั้ง 7 เกณฑ์ พบว่าค่าดัชนีความแตกต่างที่แสดงถึงคุณภาพการปรับเทียบน่าพอใจอย่างยิ่ง คุณภาพการปรับเทียบน่าพอใจ คุณภาพการปรับเทียบไม่น่าพอใจ และคุณภาพการปรับเทียบไม่น่าพอใจอย่างยิ่ง ทั้ง 4 ระดับ จะตกอยู่ในระดับน่าพอใจอย่างยิ่งเมื่อใช้เกณฑ์ของปีเตอร์เซนและคณะ จากแนวคิดที่ แฮร์ริสและเคร้าส์ (Harris and Crouse, 1993) ได้เสนอแนะว่าในกระบวนการปรับเทียบคะแนนต้องการให้มีค่าดัชนีความแตกต่างที่น้อย แสดงว่าเกณฑ์ที่พัฒนาขึ้นให้ผลการตัดสินคุณภาพการปรับเทียบคะแนนได้ดีกว่าเกณฑ์ของปีเตอร์เซนและคณะ เพราะให้ค่าดัชนีความแตกต่างที่น้อยกว่าและค่าดัชนีความแตกต่างที่กำหนดเป็นจุดตัดแบ่งกลุ่มดัชนีที่แสดงถึงคุณภาพการปรับเทียบคะแนนสำหรับเกณฑ์ที่พัฒนาขึ้นมีค่าน้อย เพราะเป็นค่าที่ได้จากการทดสอบความแตกต่างค่าเฉลี่ยของคะแนนแล้วพบค่าวิกฤตที่มีนัยสำคัญทางสถิติ ส่วนจุดตัดในเกณฑ์ของปีเตอร์เซนและคณะ ได้จากการกำหนดของผู้เชี่ยวชาญ

เกณฑ์ของปีเตอร์เซนและคณะสามารถนำไปใช้ได้ทุกสถานการณ์การปรับเทียบคะแนนเมื่อนำไปใช้ในสถานการณ์เฉพาะ เช่น การปรับเทียบโมเดลการตอบสนองข้อสอบ 3 พารามิเตอร์ เกณฑ์ที่เหมาะสมกว่าควรเป็นเกณฑ์สำหรับโมเดล 3 พารามิเตอร์ ดังเช่น แฮร์ริส และ เคร้าส์ (Harris and Crouse, 1993) กล่าวว่าไม่มีเกณฑ์ที่สมบูรณ์ที่สุดที่จะใช้ตัดสินผลการปรับเทียบคะแนนได้ทุกเงื่อนไขการปรับเทียบคะแนน และเมื่อใช้เกณฑ์ที่ต่างกันแล้วจะทำให้ผลการตัดสินคุณภาพการปรับเทียบต่างกันด้วย (Skaggs, 1990b; Livingston and others cited in Harris and Crouse, 1993)

จากการใช้เกณฑ์ที่พัฒนาขึ้นกับเกณฑ์ของปีเตอร์เซนและคณะ ทำให้ได้ระดับคุณภาพการเปรียบเทียบคะแนนแตกต่างกัน อาจเป็นเพราะเกณฑ์ที่พัฒนาขึ้นเป็นเกณฑ์สัมบูรณ์ ค่าดัชนีความแตกต่างที่บอกระดับคุณภาพการเปรียบเทียบคะแนนจะไม่แปรเปลี่ยนไปตามเงื่อนไขการเปรียบเทียบคะแนน และเป็นเกณฑ์ที่ได้จากเงื่อนไขที่อาจเป็นไปได้เกือบจะทั้งหมดของการเปรียบเทียบคะแนน ส่วนเกณฑ์ของปีเตอร์เซนและคณะจะแปรเปลี่ยนไปตามส่วนเบี่ยงเบนมาตรฐานของคะแนน ซึ่งการเปรียบเทียบคะแนนในแต่ละครั้งจะใช้เกณฑ์ที่แตกต่างกัน และในการวิจัยครั้งนี้ได้นำเกณฑ์ของปีเตอร์เซนและคณะ มาใช้กับข้อมูลจำลองหลากหลายสถานการณ์การเปรียบเทียบคะแนนและพบว่า การเปรียบเทียบคะแนนทุกสถานการณ์มีคุณภาพในระดับน่าพอใจอย่างยิ่ง แต่เมื่อใช้เกณฑ์ที่พัฒนาขึ้นและเกณฑ์ความเสมอภาคของลอร์ด ปรากฏว่าคุณภาพการเปรียบเทียบมีหลายระดับ แสดงว่าเกณฑ์ของปีเตอร์เซนและคณะให้ผลการตัดสินคุณภาพการเปรียบเทียบไม่คงที่ อาจเป็นเพราะการพัฒนาของปีเตอร์เซนและคณะมีการกำหนดช่วงของจุดตัดในเกณฑ์ที่มีความกว้างเกินไป ทำให้ผลการตัดสินคุณภาพการเปรียบเทียบคะแนนมีความคลาดเคลื่อน

การใช้เกณฑ์ของปีเตอร์เซนและคณะในการตัดสินคุณภาพการเปรียบเทียบคะแนน มีขั้นตอนวิเคราะห์ค่อนข้างจะยุ่งยากซับซ้อน เพราะต้องนำข้อมูลไปคำนวณค่าจุดตัด 4 จุด ที่แบ่งระดับคุณภาพการเปรียบเทียบคะแนนออกเป็น 5 ระดับ และนำข้อมูลชุดเดิมไปคำนวณค่าดัชนีความแตกต่าง แล้วนำมาเทียบกับเกณฑ์ที่ทำไว้ว่าคุณภาพการเปรียบเทียบอยู่ในระดับใด ส่วนเกณฑ์ที่พัฒนาขึ้นจะคำนวณค่าดัชนีความแตกต่าง แล้วนำมาเทียบกับเกณฑ์ว่าคุณภาพการเปรียบเทียบอยู่ในระดับใด

ประเด็นที่ 3 ผลการตัดสินคุณภาพการเปรียบเทียบคะแนนตามทฤษฎีการตอบสนองข้อสอบ เมื่อใช้เกณฑ์ที่พัฒนาขึ้นกับเกณฑ์ความเสมอภาคของลอร์ด

เกณฑ์ความเสมอภาคของลอร์ดเป็นเกณฑ์ที่กำหนดจากคุณสมบัติที่ดีของการเปรียบเทียบคะแนน คือ ความเสมอภาคของลอร์ด (Lord's Equity) ที่ว่าการเปรียบเทียบคะแนนระหว่างแบบสอบฉบับที่ 1 และฉบับที่ 2 จะถือว่ามีความเสมอภาคสำหรับผู้สอบก็ต่อเมื่อ ไม่ว่าผู้สอบมีความสามารถระดับใดก็ตาม เมื่อสอบแบบสอบฉบับที่ 1 หรือฉบับที่ 2 แล้วจะไม่ก่อให้เกิดผลที่แตกต่างกับผู้ใช้สอบ นั่นคือคะแนนที่ผ่านการเปรียบเทียบแล้วจะต้องมีการแจกแจงเหมือนกัน ดังนั้นเกณฑ์ความเสมอภาคของลอร์ดจึงแบ่งออกเป็น 2 ระดับ คือระดับคุณภาพการเปรียบเทียบน่าพอใจ เมื่อมีการตรวจสอบการแจกแจงของคะแนนจากแบบสอบ 2 ฉบับที่เปรียบเทียบแล้วปรากฏว่าไม่แตกต่างกัน ส่วนอีกระดับหนึ่งแสดงถึงคุณภาพการเปรียบเทียบไม่น่าพอใจเมื่อมีการตรวจสอบการแจกแจงของคะแนนจากแบบสอบทั้ง 2 ฉบับแล้วแตกต่างกัน โดยการตรวจสอบด้วยสถิติ Wilcoxon signed-rank test จากการจำลองตามเงื่อนไขการเปรียบเทียบคะแนนทั้งหมด แล้ว

ผลการตัดสินคุณภาพการเปรียบเทียบคะแนนเมื่อใช้เกณฑ์ที่พัฒนาขึ้นทั้ง 7 เกณฑ์เปรียบเทียบกับผลการใช้เกณฑ์ความเสมอภาคของลอร์ด พบว่ามีความสอดคล้องกัน โดยผลการตัดสินคุณภาพการเปรียบเทียบระดับนำพอใจอย่างยิ่งหรือระดับนำพอใจและระดับไม่นำพอใจหรือระดับไม่นำพอใจอย่างยิ่งจากการใช้เกณฑ์ที่พัฒนาขึ้น อยู่ในระดับคุณภาพนำพอใจและระดับคุณภาพไม่นำพอใจตามลำดับ เมื่อใช้เกณฑ์ความเสมอภาคของลอร์ด เป็นเพราะทั้งสองเกณฑ์ได้ใช้กระบวนการทางสถิติกำหนดจุดตัดแบ่งกลุ่มดัชนีที่แสดงคุณภาพการเปรียบเทียบคะแนน โดยเกณฑ์ที่พัฒนาขึ้นใช้สถิติ t-test แบบ Two Dependent Sample Test ส่วนเกณฑ์ความเสมอภาคของลอร์ดใช้สถิติ Wilcoxon signed-rank test ทดสอบ ทั้งสองสถิติมีข้อตกลงเบื้องต้นและกระบวนการวิเคราะห์ที่แตกต่างกัน แต่อยู่ภายใต้วัตถุประสงค์เดียวกันคือทดสอบค่าความคลาดเคลื่อนจากการเปรียบเทียบคะแนน

เมื่อเปรียบเทียบจำนวนระดับการตัดสินคุณภาพของการเปรียบเทียบคะแนนเมื่อใช้เกณฑ์ความเสมอภาคของลอร์ดและเกณฑ์ที่พัฒนาขึ้น จะเห็นว่าเกณฑ์ที่พัฒนาขึ้นให้ประโยชน์มากกว่า เพราะเกณฑ์ความเสมอภาคของลอร์ดให้ระดับคุณภาพเปรียบเทียบน้อยเพียง 2 ระดับเท่านั้น คือระดับนำพอใจและระดับไม่นำพอใจ ซึ่งให้ทางเลือกในการนำเสนอสารสนเทศเพื่อการตัดสินใจเกี่ยวกับการเปรียบเทียบคะแนนไปใช้ประโยชน์ได้น้อย ส่วนเกณฑ์ที่พัฒนาขึ้นมี 4 ระดับ คือ ระดับคุณภาพการเปรียบเทียบนำพอใจอย่างยิ่ง คุณภาพการเปรียบเทียบนำพอใจ คุณภาพการเปรียบเทียบไม่นำพอใจและระดับคุณภาพการเปรียบเทียบไม่นำพอใจอย่างยิ่ง

การใช้เกณฑ์ความเสมอภาคของลอร์ดกับเกณฑ์ที่พัฒนาขึ้นสำหรับการเปรียบเทียบคะแนนมีกระบวนการคล้ายคลึงกัน เพราะเป็นการคำนวณค่าดัชนีจากการเปรียบเทียบคะแนนแล้วนำไปเทียบกับดัชนีในเกณฑ์ว่าตกในช่วงใดของเกณฑ์ ซึ่งได้คำนวณค่าดัชนีความแตกต่างแล้วพิจารณาว่าค่าดัชนีตกอยู่ในระดับคุณภาพใดสำหรับเกณฑ์ที่พัฒนาขึ้น และทดสอบความแตกต่างการแจกแจงคะแนนที่เปรียบเทียบแล้วกับคะแนนที่ยังไม่ได้เปรียบเทียบ ถ้าการแจกแจงของคะแนนมีความแตกต่างกันอย่างมีนัยสำคัญแสดงว่าคุณภาพการเปรียบเทียบไม่นำพอใจ ส่วนการแจกแจงคะแนนไม่มีความแตกต่างกันอย่างมีนัยสำคัญ แสดงว่ามีคุณภาพการเปรียบเทียบคะแนนในระดับนำพอใจ

เกณฑ์ที่พัฒนาขึ้นมี 7 เกณฑ์ จำแนกตามเงื่อนไขการเปรียบเทียบคะแนน เพื่อให้ผู้ใช้ผลการเปรียบเทียบได้มีทางเลือกใช้เกณฑ์ที่เหมาะสมกับเงื่อนไข ส่วนเกณฑ์ความเสมอภาคของลอร์ดเป็นเกณฑ์รวมใช้ได้ทุกเงื่อนไข แต่ถ้าผู้ใช้ผลการเปรียบเทียบคะแนนต้องการให้ความสำคัญกับเงื่อนไขการเปรียบเทียบเป็นพิเศษ ก็ควรใช้เกณฑ์ที่พัฒนาขึ้นมี 7 เกณฑ์ คือ เกณฑ์รวมทุกเงื่อนไข เกณฑ์สำหรับโมเดล 1 พารามิเตอร์ เกณฑ์สำหรับโมเดล 3 พารามิเตอร์ เกณฑ์สำหรับกลุ่มสมมูล เกณฑ์สำหรับการใช้ข้อสอบร่วม เกณฑ์สำหรับการเปรียบเทียบกลับสู่แบบสอบเดิมและเกณฑ์สำหรับกลุ่มสอบทานผล

ข้อเสนอแนะในการนำเกณฑ์ที่พัฒนาขึ้นไปใช้

1. ในการนำเกณฑ์ไปใช้ตัดสินผลการเปรียบเทียบคะแนน สำหรับแบบสอบต่างฉบับ ซึ่งแบบสอบแต่ละฉบับต้องวัดคุณลักษณะเดียวกัน เป็นแบบสอบที่สร้างมาจากจุดประสงค์หรือตารางวิเคราะห์หลักสูตรเดียวกัน

2. เกณฑ์นี้จะเป็นประโยชน์อย่างยิ่งเมื่อมีการวางแผนและดำเนินการประเมินผลการเปรียบเทียบคะแนน และใช้เกณฑ์นี้ตัดสินคุณภาพการเปรียบเทียบไปพร้อมกับการเปรียบเทียบคะแนน

3. เกณฑ์เหล่านี้เป็นเกณฑ์ที่พัฒนาขึ้นเพื่อตัดสินคุณภาพการเปรียบเทียบคะแนนจริง (True Score) ตามทฤษฎีการตอบสนองข้อสอบ แต่สามารถใช้เป็นเกณฑ์ตัดสินคุณภาพการเปรียบเทียบคะแนนดิบ (Raw Score) ได้ โดยใช้แบบแผนการเก็บรวบรวมข้อมูลแบบกลุ่มสมมูล ซึ่งเป็นการสุ่มผู้สอบเข้ากลุ่มแต่ละกลุ่ม หรือจัดผู้สอบต่างกลุ่มให้มีความสามารถใกล้เคียงกัน หากคุณภาพการเปรียบเทียบด้วยการใช้กลุ่มสอบทานผล เป็นกลุ่มผู้สอบอีกกลุ่มที่มีความสามารถคล้ายกับกลุ่มผู้สอบที่ทำแบบสอบฉบับที่ 1 และกลุ่มผู้สอบที่ทำแบบสอบฉบับที่ 2 กลุ่มสอบทานผลนี้จะทำแบบสอบทั้งสองฉบับ วิเคราะห์หาค่าพารามิเตอร์ผู้สอบ (θ) ของกลุ่มสอบทานผล วิเคราะห์แยกเป็น 2 ครั้ง โดยกำหนดให้ค่าเฉลี่ยและส่วนเบี่ยงเบนมาตรฐานของค่าพารามิเตอร์ผู้สอบ (θ) ให้เท่ากันก่อนการวิเคราะห์ จากนั้นนำคะแนนดิบจากแบบสอบฉบับที่ 1 และแบบสอบฉบับที่ 2 สำหรับกลุ่มสอบทานผลที่ระดับความสามารถเดียวกันไปคำนวณหาค่าดัชนีความแตกต่าง AMD, MAD และ RMS แล้วนำไปเทียบกับเกณฑ์ที่พัฒนาขึ้นว่าคุณภาพการเปรียบเทียบคะแนนอยู่ในระดับใด

4. เกณฑ์ที่พัฒนาขึ้นนี้ไม่เหมาะที่จะนำไปใช้ในกรณีที่มีการเปรียบเทียบคะแนนตามทฤษฎีการวัดแบบดั้งเดิม (Classical Test Theory) แต่เมื่อต้องการทราบคุณภาพของการเปรียบเทียบคะแนนสามารถทำได้โดย นำคะแนนสอบรายข้อที่ได้จากการสอบไปตรวจสอบว่ามีการวัดมิติเดียว (Unidimensionality) หรือไม่ แล้วนำไปวิเคราะห์เปรียบเทียบคะแนนและหาคุณภาพการเปรียบเทียบตามทฤษฎีการตอบสนองข้อสอบ แล้วจึงใช้เกณฑ์ที่พัฒนาขึ้นนี้ตัดสินคุณภาพการเปรียบเทียบ

5. เกณฑ์ที่พัฒนาขึ้นมีทั้งหมด 7 เกณฑ์ คือ เกณฑ์รวมทุกเงื่อนไข เกณฑ์สำหรับโมเดล 1 พารามิเตอร์ เกณฑ์สำหรับโมเดล 3 พารามิเตอร์ เกณฑ์สำหรับแบบแผนการเก็บรวบรวมข้อมูลแบบกลุ่มสมมูล เกณฑ์สำหรับแบบแผนการเก็บรวบรวมข้อมูลแบบใช้ข้อสอบร่วม เกณฑ์สำหรับการหาคุณภาพการเปรียบเทียบด้วยการเปรียบเทียบกลับสู่แบบสอบเดิม และเกณฑ์สำหรับการหาคุณภาพการเปรียบเทียบด้วยการใช้กลุ่มสอบทานผล ซึ่งจะอยู่ในดุลพินิจของผู้ใช้ผลการเปรียบเทียบคะแนนจะเลือกใช้ ตามความสำคัญของเงื่อนไขที่กำหนดไว้หรือตามจุดประสงค์ของการเปรียบเทียบคะแนนในแต่ละครั้ง

6. เกณฑ์ที่พัฒนาขึ้นมีทั้งหมด 7 เกณฑ์ แต่ละเกณฑ์เป็นช่วงของค่าดัชนีความแตกต่างระหว่างคะแนนจริงที่ไม่ได้ปรับกับคะแนนจริงที่ปรับเทียบแล้ว ได้แก่ดัชนี AMD (Absolute Mean

Diference) ดัชนี MAD (Mean Absolute Difference) และดัชนี RMS (Root Mean Square) แต่ดัชนี RMS เป็นค่าความคลาดเคลื่อนในการปรับเทียบคะแนน ซึ่งคำนวณได้จากรากที่สองของค่าเฉลี่ยของกำลังสองความแตกต่างระหว่างคะแนนจริงที่ปรับเทียบแล้วกับคะแนนจริงที่ไม่ได้ปรับเทียบ การได้มาสำหรับดัชนี RMS นี้มีความคงเส้นคงวาและมีความเชื่อถือได้มากกว่าดัชนี AMD และดัชนี MAD เพราะค่าดัชนี AMD และดัชนี MAD มีโอกาสเท่ากันในกรณีที่คะแนนจริงที่ปรับเทียบแล้วทุกค่ามากกว่าคะแนนจริงที่ไม่ได้ปรับเทียบ ค่าดัชนี RMS มีค่ามากกว่าดัชนี AMD และดัชนี MAD พิสูจน์ในแต่ละช่วงของเกณฑ์มากกว่า ทำให้เห็นความแตกต่างได้ชัดเจนกว่าดัชนีทั้งสอง และในงานวิจัยโดยทั่วไปจะนิยมใช้ดัชนี RMS ฉะนั้นเมื่อใช้เกณฑ์ที่พัฒนาขึ้นและต้องการเลือกใช้ค่าดัชนีเพียงดัชนีเดียว ควรเลือกใช้ดัชนี RMS

ข้อเสนอแนะเพื่อทำวิจัยต่อไป

1. เกณฑ์ที่พัฒนาขึ้นนี้ใช้สำหรับการปรับเทียบคะแนนระหว่างแบบสอบที่มีการให้คะแนนแบบ 0-1 แต่ยังมีแบบสอบที่มีการให้คะแนนแบบพหุวิภาคหรือให้คะแนนหลายค่า เช่น 5, 4, 3, 2, 1 เพื่อขยายองค์ความรู้เรื่องเกณฑ์การปรับเทียบคะแนน และได้เกณฑ์ที่มีความเชื่อถือได้ จึงควรพัฒนาเกณฑ์สำหรับการปรับเทียบคะแนนตามทฤษฎีการตอบสนองข้อสอบสำหรับการปรับเทียบคะแนนที่มีการให้คะแนนแบบพหุวิภาค

2. เกณฑ์ที่พัฒนาขึ้นในครั้งนี้เป็นเกณฑ์ตัดสินคุณภาพการปรับเทียบคะแนนตามทฤษฎีการตอบสนองข้อสอบ ควรพัฒนาเกณฑ์ตัดสินคุณภาพการปรับเทียบคะแนนตามทฤษฎีการวัดแบบดั้งเดิม (Classical Test Theory) เพื่อให้ผู้นำผลการปรับเทียบคะแนนได้เลือกใช้เกณฑ์ให้เหมาะสมกับกระบวนการปรับเทียบคะแนนระหว่างแบบสอบ