

บทที่ 1

บทนำ

ความเป็นมาและความสำคัญของปัญหา

แบบทดสอบเป็นเครื่องมือหลักที่ใช้ในการวัดและประเมินผลการศึกษา เพื่อตรวจสอบว่าผู้เข้าสอบมีคุณลักษณะแห่งหรือความสามารถใดอยู่ในระดับใด ดังนั้น การสร้างและการตรวจสอบคุณภาพของแบบทดสอบจะต้องคำนึงถึงความตรง (Validity) เป็นสำคัญ ทั้งนี้เพื่อระบุว่าความตรงเป็นคุณสมบัติของแบบทดสอบที่แสดงว่า คะแนนจากแบบทดสอบสามารถสรุปอ้างอิงไปยังสิ่งที่วัดได้อย่างเหมาะสม มีความหมาย และเป็นประโยชน์ (เสรี ชัดแจ้ง, 2544, หน้า 137) ถ้าผลการวัดได้ค่าที่ใกล้เคียงกับคุณลักษณะที่แท้จริงมากเพียงใด ก็ถือว่าการวัดมีความตรงมากขึ้นเพียงนั้น (ศรีชัย กาญจนวงศ์, 2545, หน้า 103) นักวัดผลตรวจสอบหลักฐานที่แสดงความตรงของแบบทดสอบได้ 3 ด้าน คือ (1) ความตรงเชิงเนื้อหา (Content Validity) (2) ความตรงเชิงเกณฑ์สัมพันธ์ (Criterion Related Validity) และ (3) ความตรงเชิงโครงสร้าง (Construct Validity) ส่วนการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ (Differential Item Functioning) เป็นการตรวจสอบคุณภาพของแบบทดสอบด้านความตรงเชิงโครงสร้าง (Mazor, Clauser, & Hambleton, 1992; Kim, Kim, & Cohen, 1994 อ้างถึงในวีลีนาศ แซ่จัง, 2543, หน้า 1)

การทำหน้าที่ต่างกันของข้อสอบ (Differential Item Functioning: DIF) มีผู้ให้นิยามไว้หลายความหมายด้วยกัน แต่ที่ยอมรับกันอย่างกว้างขวางก็คือ ข้อสอบทำหน้าที่ต่างกันภายใต้เงื่อนไขผู้เข้าสอบมีความสามารถเท่ากัน แต่มาจากการกลุ่มผู้เข้าสอบย่อยที่มีลักษณะต่างกัน มีความน่าจะเป็นในการตอบข้อสอบข้อนั้นไม่เท่ากัน (เสรี ชัดแจ้ง, 2540, หน้า 42)

หลักการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบใช้การเปรียบเทียบผลการตอบข้อสอบระหว่างกลุ่มผู้เข้าสอบย่อย 2 กลุ่ม ที่มีความสามารถระดับเดียวกัน โดยผู้เข้าสอบกลุ่มนั้นเป็นตัวแทนกลุ่มหลักในประชากรเรียกว่า “กลุ่มอ้างอิง” (Reference Group: R) ซึ่งเป็นกลุ่มพื้นฐาน ส่วนผู้เข้าสอบอีกกลุ่มนั้นเป็นตัวแทนกลุ่มรองในประชากรเรียกว่า “กลุ่มเปรียบเทียบ” (Focal Group: F) ซึ่งตามปกติแล้วเป็นกลุ่มผู้เข้าสอบที่สนใจจะศึกษาการทำหน้าที่ต่างกันของข้อสอบ (Angoff, 1993) ถ้าข้อสอบทำหน้าที่ต่างกันแล้วโอกาสตอบข้อสอบถูกของผู้เข้าสอบหั้งสองกลุ่มนั้นจะไม่เท่ากัน โดยคาดว่าผู้เข้าสอบกลุ่มอ้างอิงจะได้เปรียบในการตอบข้อสอบ ส่วนผู้เข้าสอบกลุ่มเปรียบเทียบจะเสียเปรียบในการตอบข้อสอบ

เงื่อนไขสำคัญในการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบระหว่างกลุ่มอ้างอิงกับกลุ่มเปรียบเทียบ ใช้เกณฑ์การจับคู่ (Matching Criteria) ที่นิยมใช้กันมี 2 แนวทางใหญ่ คือ แนวทางแรก เป็นการตรวจสอบโดยใช้เกณฑ์ภายนอก (External Criteria) โดยการนำคะแนนจากการตอบข้อสอบไปสัมพันธ์กับเกณฑ์ที่ได้จากภายนอกแบบทดสอบ เช่น เกรดเฉลี่ยสะสมหรือคะแนนจากชิ้นงานที่ทำเป็นต้น แต่จุดอ่อนของแนวทางนี้คือ เป็นการยกที่จะหาเกณฑ์ภายนอกที่มีความตรงและยุติธรรมซึ่งถ้าเกณฑ์การจับคู่ขาดคุณสมบัติเรื่องนี้จะทำให้ผลการตรวจสอบเกิดความคลาดเคลื่อนได้ ส่วนแนวทางที่สอง เป็นการตรวจสอบโดยใช้เกณฑ์ภายใน (Internal Criteria) โดยการนำวิธีการทางสถิติมาตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ เป็นการพิจารณาจากโครงสร้างภายในของแบบทดสอบเป็นหลัก (Rudner, Getson & Knight, 1980) ซึ่งในปัจจุบันเป็นแนวทางที่มีผู้นิยมศึกษา กันมาก

วิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบมีหลายวิธี แรมเบิลตัน และคณะ (Hambleton et al., 1993 อ้างอิงใน เสรี ชัด เช้ม, 2539, หน้า 1-4) ได้จำแนกออกเป็น 3 กลุ่มใหญ่ ๆ ดังนี้

1. กลุ่มวิธีที่ใช้ทฤษฎีการทดสอบแบบดั้งเดิม (Methods Using Classical Test Theory: CTT) วิธีในกลุ่มนี้พัฒนามาจากหลักการของทฤษฎีการทดสอบแบบดั้งเดิม วิธีการในกลุ่มนี้ได้แก่ การวิเคราะห์ความแปรปรวน (Analysis of Variance) วิธีสหสัมพันธ์ (Correlational Methods) (Green & Draper, 1972 cited in Scheuneman & Bleistein, 1989) วิธีแปลงค่าความยากของข้อสอบ (Transformed Item Difficulty Method, TID) ฯลฯ

2. กลุ่มวิธีที่ใช้ทฤษฎีการตอบสนองข้อสอบ (Methods Using Item Response Theory: IRT) วิธีในกลุ่มนี้ตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ ตามกรอบแนวคิดของทฤษฎีการตอบสนองข้อสอบ วิธีการในกลุ่มนี้ได้แก่ วิธี Analysis of fit (Durovic, 1975, cited in Hambleton & Others, 1993) วิธี Difficulty shift (Wright, Mead & Draba, 1976, cited in Hambleton & Others, 1993) ฯลฯ

3. กลุ่มวิธีที่ใช้วิธีไค-สแควร์ (Methods Using Chi-Square Methods) วิธีการในกลุ่มนี้ได้แก่ วิธีแมนเทล-แฮนส์เซล (Mantel-Haenszel Method: MH) (Holland & Thayer, 1986, 1988) และวิธีลดด้อยโลจิสติก (Logistic Regression Methods: LR) (Swaminathan & Rogers, 1990) ฯลฯ

วิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ ได้พัฒนาขึ้นเรื่อย ๆ ปัจจุบันวิธีการวิเคราะห์องค์ประกอบน้ำหนัก (Restricted Factor Analysis Methods: RFA) ซึ่งได้รับการพัฒนาขึ้นโดย ออร์ท (Oort, 1998) เป็นวิธีการใหม่ล่าสุดในการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ ซึ่งอาจขัดอยู่ในกลุ่มวิธีไค-สแควร์ และออร์ท (Oort, 1998, pp. 107-124) ได้ศึกษาเปรียบเทียบ

ประสิทธิภาพในการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบระหว่างวิธีการตอบสนองข้อสอบแบบ 1 พารามิเตอร์ กับวิธีการวิเคราะห์องค์ประกอบฉบับจำกัด โดยใช้สถานการณ์จำลองสร้างข้อสอบ 3 แบบ คือ ข้อสอบทำหน้าที่ต่างกันมาก ข้อสอบทำหน้าที่ต่างกันน้อย และข้อสอบทำหน้าที่ไม่ต่างกัน และกำหนดขนาดกลุ่มตัวอย่าง 2 ขนาด คือ กลุ่มตัวอย่างขนาดใหญ่ (2,000 คน) และขนาดเล็ก (200 คน) ผลปรากฏว่า วิธีการวิเคราะห์องค์ประกอบฉบับจำกัด มีประสิทธิภาพในการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบได้ใกล้เคียงกับวิธีการตอบสนองข้อสอบแบบ 1 พารามิเตอร์ แต่ในกรณีกลุ่มตัวอย่างขนาดเล็ก วิธีการวิเคราะห์องค์ประกอบฉบับจำกัด มีประสิทธิภาพดีกว่าวิธีการตอบสนองข้อสอบแบบ 1 พารามิเตอร์ ในทุกเงื่อนไขของการทดสอบ นอกจากนี้ วิธีการวิเคราะห์องค์ประกอบฉบับจำกัดขึ้นประหัดเวลาและค่าใช้จ่ายในการวิเคราะห์ ส่วนนิคม กิริโวรงกูร (2542) ได้ศึกษาผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบด้วยวิธีการวิเคราะห์องค์ประกอบฉบับจำกัด วิธีแม่นเทล-แฮนส์เซล และวิธีการตอบสนองข้อสอบแบบ 2 พารามิเตอร์ พบว่า โดยภาพรวม วิธีการวิเคราะห์องค์ประกอบฉบับจำกัด มีประสิทธิภาพในการตรวจสอบสูงที่สุด รองลงมาวิธีแม่นเทล-แฮนส์เซล และวิธีการตอบสนองข้อสอบแบบ 2 พารามิเตอร์ และวิธีการตอบสนองข้อสอบแบบ 2 พารามิเตอร์ มีอัตราความคลาดเคลื่อนประเภทที่ 1 สูงกว่าวิธีแม่นเทล-แฮนส์เซล และวิธีการวิเคราะห์องค์ประกอบฉบับจำกัด

สำหรับการศึกษาเกี่ยวกับปัจจัยที่ส่งผลต่อประสิทธิภาพการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ ปรากฏว่า ข้อค้นพบริ่ำมีความสอดคล้องกัน กล่าวคือ เมื่อขนาดกลุ่มตัวอย่างใหญ่ขึ้น อัตราการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบจะเพิ่มขึ้น (Roger & Swaminathan, 1993) แต่ในเรื่องขนาดของกลุ่มตัวอย่างยังไม่ชัดเจนว่า ต้องใช้กลุ่มตัวอย่างขนาดเท่าใดจึงจะเพียงพอ เมเชอร์ และຄณา (Mazor et al., 1992) กล่าวว่า ขนาดของกลุ่มตัวอย่างที่เหมาะสมสำหรับ วิธีแม่นเทล-แฮนส์เซล ควรใช้ระหว่าง 100 และ 300 คน สำหรับกลุ่มไดกฤุณหนึ่ง หรือห้องสอบกลุ่ม (กลุ่มอ้างอิง และกลุ่มเปรียบเทียบ) นารายานัน และสวามินาทาน (Narayanan & Swaminathan, 1994) ได้เสนอแนะว่า โดยทั่วไปใช้กลุ่มตัวอย่างขนาด กลุ่มละ 300 คน ที่เพียงพอที่จะตรวจสอบการทำหน้าที่ต่างกันของข้อสอบได้อย่างมีประสิทธิภาพ และพบว่าอัตราการตรวจสอบข้อสอบการทำหน้าที่ต่างกันระหว่างวิธีแม่นเทล-แฮนส์เซล กับวิธีซิปเกสท์ จะได้รับผลกระทบจากกลุ่มเปรียบเทียบที่มีขนาดเล็กมากกว่ากลุ่มอ้างอิงที่มีขนาดเล็ก นอกจากรูปแบบเดียวกันนี้ แม้เมบลิตัน และຄณา (Hambleton et al., 1993) ยังเสนอแนะว่า กลุ่มตัวอย่างที่ใช้ในการวิเคราะห์ด้วยวิธีแม่นเทล-แฮนส์เซล ควรอยู่ระหว่าง 200 ถึง 1,000 คน แต่ในการวิเคราะห์บางเงื่อนไข การใช้กลุ่มตัวอย่าง 200 คน ในกลุ่มไดกฤุณหนึ่งอาจจะไม่เพียงพอ ถ้าใช้กลุ่มตัวอย่างขนาดใหญ่จะทำให้ได้ผล

การตรวจสอบที่ดีกว่า จากที่กล่าวมาอาจกล่าวได้ว่าขนาดกลุ่มตัวอย่างน่าจะมีผลกระทบต่อประสิทธิภาพการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ

การศึกษาของ คิมและโคเคน (Kim & Cohen, 1994) และ อัททาโรและมิลล์เชพ (Uttaro & Millsap, 1994) พบว่า ความยาวของแบบทดสอบมีผลกระทบต่ออัตราการตรวจพบข้อสอบทำหน้าที่ต่างกัน และจากงานวิจัยของ นาราيانัน และสวามินาทาน (Narayanan & Swaminathan, 1996) พบว่า แบบทดสอบขนาด 40 ข้อ เป็นแบบทดสอบที่มีความยาวเพียงพอสำหรับการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ แต่ โรเจอร์ และสวามินาทาน (Roger & Swaminathan, 1993) พบว่า ความยาวของแบบทดสอบไม่มีผลกระทบต่ออัตราการตรวจพบข้อสอบทำหน้าที่ต่างกัน ดังนั้นเพื่อเป็นการยืนยันข้อค้นพบให้มีความชัดเจนมากยิ่งขึ้น จึงควรศึกษาความยาวของแบบทดสอบว่ามีผลต่อการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบหรือไม่

นอกจากนี้ยังมีการศึกษาปัจจัยอื่น ๆ ที่มีผลกระทบต่อประสิทธิภาพการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ คือ เกณฑ์การจับคู่ (Matching Criteria) เป็นการหาเกณฑ์การจับคู่เรียบเทียบที่เหมาะสมสำหรับตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ ที่ผ่านมา มีการศึกษาวิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบที่ใช้ทั้งเกณฑ์ภายนอก และเกณฑ์ภายใน เช่น นุพนาก พิพัฒนสูข (2541) ได้เปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบระหว่างวิธีแม่นเทล-แ昏ส์เซล กับวิธีถดถอยโลจิสติก เมื่อใช้เกณฑ์ภายในเป็นเกณฑ์การจับคู่พบว่า เมื่อใช้เกณฑ์การจับคู่ค่าคะแนนรวมทั้งฉบับ วิธีแม่นเทล-แ昏ส์เซล มีประสิทธิภาพมากกว่าวิธีถดถอยโลจิสติก และเมื่อใช้เกณฑ์การจับคู่ค่าคะแนนแบบทดสอบย่อย วิธีทั้งสองมีประสิทธิภาพไม่แตกต่างกัน สำหรับวิธีถดถอยโลจิสติกมีประสิทธิภาพมากกว่าวิธีแม่นเทล-แ昏ส์เซล เมื่อใช้เกณฑ์การจับคู่ค่าคะแนนจากแบบทดสอบย่อยหลายฉบับ ซึ่งสอดคล้องกับผลการศึกษาของเคลาเซอร์ และคอลล์ (Clauser et al., 1996, pp. 202-214) ที่ศึกษาเกณฑ์การจับคู่ในการตรวจสอบการทำหน้าที่ต่างกันโดยใช้เกณฑ์การจับคู่ 2 เกณฑ์ ได้แก่ ค่าคะแนนรวมทั้งฉบับ ค่าคะแนนแบบทดสอบย่อย และค่าคะแนนจากแบบทดสอบย่อยหลายฉบับ ระหว่างวิธีแม่นเทล-แ昏ส์เซลกับวิธีถดถอยโลจิสติก ทั้งในข้อมูลจริง และข้อมูลจำลอง พบว่า ในแบบทดสอบพหุมิติ การใช้ค่าคะแนนรวมทั้งฉบับเป็นเกณฑ์การจับคู่ไม่เหมาะสม แต่เมื่อใช้ค่าคะแนนจากแบบทดสอบย่อยหลายฉบับเป็นเกณฑ์การจับคู่ จะมีประสิทธิภาพในการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ ได้ดีกว่า โดยพิจารณาจากจำนวนข้อสอบทำหน้าที่ต่างกันที่ตรวจพบน้อยที่สุด และเมื่อเปรียบเทียบกับเกณฑ์การจับคู่ค่าคะแนนรวมทั้งฉบับ และค่าคะแนนแบบทดสอบย่อย การนำค่าคะแนนรวมทั้งฉบับมาใช้เป็นเกณฑ์การจับคู่เปรียบเทียบนั้นไม่เหมาะสม เพราะจะทำให้อัตราความคลาดเคลื่อนประ踉ที่ 1 สูงขึ้น

จากการศึกษางานวิจัยที่ผ่านมา มีประเด็นปัญหาที่ผู้วิจัยสนใจ ดังนี้

1. ยังไม่มีการเปรียบเทียบประสิทธิภาพในการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบระหว่างวิธีการวิเคราะห์องค์ประกอบจำกัด กับวิธีทดสอบโดยโลจิสติก ว่าวิธีใดมีประสิทธิภาพดีกว่ากัน
2. ปัจจัยที่ส่งผลกระทบต่อประสิทธิภาพการตรวจพิจารณาข้อสอบทำหน้าที่ต่างกัน คือ ขนาดของกลุ่มตัวอย่าง และความขาวของแบบทดสอบ ดังนั้น น่าจะศึกษาว่าขนาดของกลุ่มตัวอย่าง และความขาวของแบบทดสอบ มีผลต่อการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบอย่างไร
3. ภาษาได้เงื่อนไขการใช้เกณฑ์ภาษาในเป็นเกณฑ์การจับคู่ ได้แก่ คะแนนรวมทั้งฉบับ และคะแนนแบบทดสอบย่อย วิธีการวิเคราะห์องค์ประกอบจำกัดกับวิธีทดสอบโดยโลจิสติก มีประสิทธิภาพในการตรวจสอบต่างกันหรือไม่

การวิจัยครั้งนี้ผู้วิจัยตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบวัดผลสัมฤทธิ์ทางการเรียนระดับชาติ วิชาภาษาไทย ชั้นประถมศึกษาปีที่ 6 ซึ่งวัดองค์ประกอบเดียว โครงสร้างความรู้ และค้านกระบวนการ (กระทรวงศึกษาธิการ, 2546) จึงเป็นแบบทดสอบพหุมิติ

สำหรับตัวแปรจำแนกกลุ่มประชากรที่เลือกศึกษา ได้แก่ ตัวแปรเพศ เนื่องจากพิจารณาเห็นว่าเป็นตัวแปรที่สามารถแบ่งได้โดยปราศจากความคลาดเคลื่อน (Millsap & Everson, 1993; Zieky, 1993 อ้างอิงใน เสรี ชั้ดเชน, 2539, หน้า 12) จึงหมายความที่จะนำมาใช้ศึกษาในเชิงวิธีการกับการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบด้วยข้อมูลจริง ในงานวิจัยนี้เพศหญิงเป็นกลุ่มอ้างอิง เพราะแบบทดสอบพหุมิติวิชาภาษาไทยส่วนใหญ่ จะดำเนินการเข้าข้างเพศหญิง (สุพัฒน์ สุกมลลัตน์, 2534; กาญจนा วัชันสุนทร, 2537)

จากประเด็นปัญหาดังกล่าวข้างต้น ผู้วิจัยจึงสนใจที่จะเปรียบเทียบประสิทธิภาพการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบพหุมิติ ระหว่างวิธีการวิเคราะห์องค์ประกอบจำกัดกับวิธีทดสอบโดยโลจิสติก ภายใต้เงื่อนไขที่จะใช้เป็นแนวทางในการเลือกวิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบพหุมิติ ว่าควรเลือกวิธีใด ใช้กลุ่มตัวอย่างขนาดเท่าใด แบบทดสอบยาวเท่าใด และเกณฑ์การจับคู่แบบใด จึงจะมีประสิทธิภาพในการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบพหุมิติ

วัตถุประสงค์ของการวิจัย

1. เพื่อตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบพหุมิติ ระหว่างวิธีการวิเคราะห์องค์ประกอบจำกัดกับวิธีทดสอบโดยโลจิสติก ภายใต้เงื่อนไข 18 เงื่อนไข ($3 \times 3 \times 2$) คือ ขนาดของกลุ่มตัวอย่าง 3 ขนาด (2,000 คน 1,000 คน และ 300 คน) ความขาวของแบบทดสอบ

3 ขนาด (40 ข้อ 30 ข้อ และ 20 ข้อ) และเกณฑ์การจับคู่ 2 เกณฑ์ (คะแนนรวมทั้งฉบับ และคะแนนแบบทดสอบย่อย)

2. เพื่อเปรียบเทียบประสิทธิภาพของวิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบพหุมิติ ระหว่างวิธีการวิเคราะห์องค์ประกอบจำกัดกับวิธีทดสอบโดยโลจิสติก ภายใต้เงื่อนไข 18 เมื่อนำไป ($3 \times 3 \times 2$) คือ ขนาดของกลุ่มตัวอย่าง 3 ขนาด (2,000 คน 1,000 คน และ 300 คน) ความยาวของแบบทดสอบ 3 ขนาด (40 ข้อ 30 ข้อ และ 20 ข้อ) และเกณฑ์การจับคู่ 2 เกณฑ์ (คะแนนรวมทั้งฉบับ และคะแนนแบบทดสอบย่อย) เมื่อใช้ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ โดยวิธีชิปเพลทที่เป็นเกณฑ์สำหรับการเปรียบเทียบในเรื่อง อัตราความถูกต้อง และอัตราความคลาดเคลื่อน

สมมติฐานของการวิจัย

โรเจอร์ และสวามินาธาน (Roger & Swaminathan, 1993, pp. 105-116) ได้ศึกษาพบว่า วิธีคะแนนแทนสีเซลล์ มีประสิทธิภาพในการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบแบบเอกสารเป็นที่ยอมรับวิธีทดสอบโดยโลจิสติก และวิธีทดสอบโดยโลจิสติก มีประสิทธิภาพในการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบแบบนี้มากกว่าวิธีคะแนนแทนสีเซลล์ ออร์ท (Oort, 1998, pp. 107-124) ได้เปรียบเทียบประสิทธิภาพในการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบระหว่างวิธีการตอบสนองของข้อตอบ แบบ 1 พารามิเตอร์ กับวิธีการวิเคราะห์องค์ประกอบจำกัด ในข้อสอบแบบให้คะแนนเป็น 2 ค่า (ตอบถูกได้ 1 และตอบผิดได้ 0) เมื่อกำหนดเงื่อนไขให้ข้อสอบทำหน้าที่ต่างกันมาก ปานกลาง และน้อย กลุ่มตัวอย่างมีขนาด 2 ขนาด คือ ขนาดใหญ่ 2,000 คน และขนาดเล็ก 200 คน ใช้แบบทดสอบเดียวกัน(40 ข้อ) พบว่า โดยภาพรวม การตรวจสอบด้วยวิธีการวิเคราะห์องค์ประกอบจำกัด ให้ผลการตรวจสอบที่ดีกว่าวิธีการตอบสนองของข้อตอบ แบบ 1 พารามิเตอร์ สำหรับกลุ่มตัวอย่างขนาดเล็ก (200 คน) การตรวจสอบด้วยวิธีการวิเคราะห์องค์ประกอบจำกัด ให้ผลการตรวจสอบดีกว่า แต่เมื่อกลุ่มตัวอย่างมีขนาดใหญ่ ทั้งสองวิธีจะให้ผลลัพธ์ที่ต่างกันน้อย นักงานศึกษาของนิคม กิรติวงศ์ (2542) ยังสนับสนุนว่า โดยภาพรวมวิธีการวิเคราะห์องค์ประกอบจำกัด มีประสิทธิภาพในการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบดีกว่าวิธีคะแนนแทนสีเซลล์ และวิธีการตอบสนองของข้อสอบ แบบ 2 พารามิเตอร์ และวิธีการตอบสนองของข้อสอบแบบ 2 พารามิเตอร์ มีอัตราความคลาดเคลื่อนประเภทที่ 1 สูงกว่าวิธีคะแนนแทนสีเซลล์ และวิธีการวิเคราะห์องค์ประกอบจำกัด จึงทำให้ผู้วิจัยตั้งสมมติฐานของการวิจัยดังนี้

1. เมื่อถูกนั่งตัวอย่างขนาดใหญ่ (2,000 คน) แบบทดสอบ (40 ข้อ 30 ข้อ และ 20 ข้อ) และเกณฑ์การจับคู่คะแนนรวมทั้งฉบับ วิธีการวิเคราะห์องค์ประกอบจำภาคมีประสิทธิภาพดีกว่าวิธีทดลองโดยโลจิสติก
2. เมื่อถูกนั่งตัวอย่างขนาดใหญ่ (2,000 คน) แบบทดสอบ (40 ข้อ 30 ข้อ และ 20 ข้อ) และเกณฑ์การจับคู่คะแนนแบบทดสอบย่อ วิธีการวิเคราะห์องค์ประกอบจำภาคมีประสิทธิภาพดีกว่าวิธีทดลองโดยโลจิสติก
3. เมื่อถูกนั่งตัวอย่างขนาดกลาง (1,000 คน) แบบทดสอบ (40 ข้อ 30 ข้อ และ 20 ข้อ) และเกณฑ์การจับคู่คะแนนรวมทั้งฉบับ วิธีการวิเคราะห์องค์ประกอบจำภาคมีประสิทธิภาพดีกว่าวิธีทดลองโดยโลจิสติก
4. เมื่อถูกนั่งตัวอย่างขนาดกลาง (1,000 คน) แบบทดสอบ (40 ข้อ 30 ข้อ และ 20 ข้อ) และเกณฑ์การจับคู่คะแนนแบบทดสอบย่อ วิธีการวิเคราะห์องค์ประกอบจำภาคมีประสิทธิภาพดีกว่าวิธีทดลองโดยโลจิสติก
5. เมื่อถูกนั่งตัวอย่างขนาดเล็ก (300 คน) แบบทดสอบ (40 ข้อ 30 ข้อ และ 20 ข้อ) และเกณฑ์การจับคู่คะแนนรวมทั้งฉบับ วิธีการวิเคราะห์องค์ประกอบจำภาคมีประสิทธิภาพดีกว่าวิธีทดลองโดยโลจิสติก
6. เมื่อถูกนั่งตัวอย่างขนาดเล็ก (300 คน) แบบทดสอบ (40 ข้อ 30 ข้อ และ 20 ข้อ) และเกณฑ์การจับคู่คะแนนแบบทดสอบย่อ วิธีการวิเคราะห์องค์ประกอบจำภาคมีประสิทธิภาพดีกว่าวิธีทดลองโดยโลจิสติก

ประโยชน์ที่คาดว่าจะได้รับจากการวิจัย

งานวิจัยครั้งนี้ มุ่งเปรียบเทียบประสิทธิภาพการตรวจสอบการดำเนินการที่ต่างกันของข้อสอบในแบบทดสอบพหุมิติ ระหว่างวิธีการวิเคราะห์องค์ประกอบจำภาคกับวิธีทดลองโดยโลจิสติก ภายใต้เงื่อนไข 18 เงื่อนไข ($3 \times 3 \times 2$) คือ ขนาดของกลุ่มตัวอย่าง 3 ขนาด (2,000 คน 1,000 คน และ 300 คน) ความขาวของแบบทดสอบ 3 ขนาด (40 ข้อ 30 ข้อ และ 20 ข้อ) และเกณฑ์การจับคู่ 2 เกณฑ์ (คะแนนรวมทั้งฉบับ และคะแนนแบบทดสอบย่อ) โดยศึกษากับข้อมูลจริง ผู้วิจัยเชื่อคาดว่า จะเป็นประโยชน์ ดังนี้

1. เป็นแนวทางในการเลือกวิธีการตรวจสอบการดำเนินการที่ต่างกันของข้อสอบในแบบทดสอบพหุมิติ ที่เหมาะสมต่อการปฏิบัติ โดยคำนึงถึงประสิทธิภาพในการตรวจสอบการดำเนินการที่ต่างกันของข้อสอบ

2. เป็นแนวทางในการเลือกใช้ขนาดของกลุ่มตัวอย่าง ความขาวของแบบทดสอบ และเกณฑ์การจับคู่ ว่าควรใช้กลุ่มตัวอย่างเท่าใด ความขาวของแบบทดสอบขนาดเท่าใด และเกณฑ์การจับคู่แบบใด จึงจะทำให้มีประสิทธิภาพในการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ ในแบบทดสอบพหุมิติที่สุด

3. เป็นแนวทางในการศึกษาเกี่ยวกับการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ ในแบบทดสอบพหุมิติต่อไป

ขอบเขตของการวิจัย

1. ประชากรและกลุ่มตัวอย่าง

1.1 ประชากร เป็นนักเรียนชั้นประถมศึกษาปีที่ 6 สังกัดสำนักงานคณะกรรมการการศึกษาขั้นพื้นฐาน กระทรวงศึกษาธิการ ปีการศึกษา 2546 ทั้งประเทศ ที่เข้าสอบวัดผลสัมฤทธิ์ทางการเรียนเพื่อประเมินคุณภาพการศึกษาระดับชาติ จำนวน 29,204 โรงเรียน และมีนักเรียนที่เข้าสอบ วิชาภาษาไทย จำนวน 750,978 คน

1.2 กลุ่มตัวอย่าง เป็นนักเรียนชั้นประถมศึกษาปีที่ 6 สังกัดสำนักงานคณะกรรมการการศึกษาขั้นพื้นฐาน ได้มาโดยการสุ่มแบบแบ่งชั้น ไม่กำหนดคัดส่วน โดยใช้ระดับความสามารถ (คะแนนรวม) แบ่งเป็น 3 ระดับ คือ ดี พอดี และปรับปรุง ให้นักเรียนเป็นหน่วยการสุ่ม สูงมา จำนวน 2,000 คน แบ่งเป็น นักเรียนชาย 1,000 คน และนักเรียนหญิง 1,000 คน

2. ตัวแปรที่ศึกษา

2.1 ตัวแปรต้น มี 4 ตัว ได้แก่

2.1.1 วิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ 2 วิธี ได้แก่

2.1.1.1 วิธีการวิเคราะห์องค์ประกอบประจำกลุ่ม

2.1.1.2 วิธีคอมโบทโลจิสติก

2.1.2 ขนาดของกลุ่มตัวอย่าง 3 ขนาด ได้แก่

2.1.2.1 กลุ่มตัวอย่างขนาดใหญ่ (2,000 คน)

2.1.2.2 กลุ่มตัวอย่างขนาดกลาง (1,000 คน)

2.1.2.3 กลุ่มตัวอย่างขนาดเล็ก (300 คน)

2.1.3 ความขาวของแบบทดสอบ 3 ขนาด ได้แก่

2.1.3.1 แบบทดสอบ 40 ข้อ

2.1.3.2 แบบทดสอบ 30 ข้อ

2.1.3.3 แบบทดสอบ 20 ข้อ

2.1.4 เกณฑ์การจับคู่ 2 เกณฑ์ ได้แก่

2.1.4.1 เกณฑ์การจับคู่คะแนนรวมทั้งฉบับ

2.1.4.2 เกณฑ์การจับคู่คะแนนแบบทดสอบย่อย

2.2 ตัวแปรตาม คือ ประสิทธิภาพในการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในเรื่อง

2.2.1 อัตราความถูกต้อง

2.2.2 อัตราความคลาดเคลื่อน

3. เกณฑ์ที่ใช้แสดงว่าข้อสอบทำหน้าที่ต่างกัน มีดังนี้

3.1 เกณฑ์สำหรับวิธีซิปเทสท์ ได้แก่ ค่า $\beta_n > 0$ และการทดสอบค่าสถิติ Z

ที่ระดับนัยสำคัญ .05

3.2 เกณฑ์สำหรับวิธีการวิเคราะห์องค์ประกอบจำกัด ได้แก่ ค่าดัชนีตัดแบ่งไมเดล (MI) และค่าพารามิเตอร์ที่คาดหวัง (EPC) แตกต่างจาก 0 หรือไม่ ที่ระดับนัยสำคัญ .05

3.3 เกณฑ์สำหรับวิธีทดสอบโดยโลจิสติก ได้แก่ ค่าอิทธิพลจากปฏิสัมพันธ์ระหว่างกลุ่มผู้เข้าสอบกับความสามารถของผู้เข้าสอบหรือค่าอิทธิพลจากกลุ่มผู้เข้าสอบ และการทดสอบค่าสถิติ χ^2 ที่ระดับนัยสำคัญ .05

ข้อตกลงเบื้องต้น

วิธีซิปเทสท์ เป็นการนำวิธีการทางสถิติมาใช้ตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ โดยพิจารณาโครงสร้างภายในของแบบทดสอบพหุมิติ (Multidimensional Test) ใช้หลักการของทฤษฎีการตอบสนองข้อสอบแบบพหุมิติ จึงเป็นวิธีตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบพหุมิติ ที่ให้ผลการตรวจสอบที่เชื่อถือได้ (Shealy & Stout, 1993) ดังนั้น ผู้วิจัยจึงเลือกวิธีซิปเทสท์เป็นเกณฑ์สำหรับเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบระหว่างวิธีการวิเคราะห์องค์ประกอบจำกัดกับวิธีทดสอบโดยโลจิสติก ในศึกษาระดับนี้ผู้วิจัย วิเคราะห์การทำหน้าที่ต่างกันของข้อสอบโดยใช้โปรแกรม SIBTEST ที่สั่งซื้อมาจาก Assessment Systems Corporation ในประเทศสหรัฐอเมริกา

นิยามศัพท์เฉพาะ

การทำหน้าที่ต่างกันของข้อสอบ (Differential Item Functioning: DIF) หมายถึง ข้อสอบที่ผู้เข้าสอบซึ่งมีความสามารถเท่ากันในเนื้อหาที่ต้องการวัด มีโอกาสตอบข้อสอบข้อนี้ได้ถูกต้อง

ไม่เท่ากัน เนื่องจากอยู่ในกลุ่มผู้เข้าสอบบ่อยที่มีลักษณะต่างกันในที่นี้คือ กลุ่มผู้เข้าสอบเพศหญิง กับกลุ่มผู้เข้าสอบเพศชาย

แบบทดสอบพหุมิติ (Multidimensional Test) หมายถึง แบบทดสอบที่วัดคุณลักษณะเด่นตั้งแต่ 2 คุณลักษณะขึ้นไป ในที่นี้คือแบบทดสอบวัดผลสัมฤทธิ์ทางการเรียนระดับชาติ วิชาภาษาไทย ชั้นประถมศึกษาปีที่ 6 ซึ่งแบบทดสอบของคุณลักษณะนี้ประกอบด้วย 2 ด้าน ได้แก่ ด้านโครงสร้างความรู้ และด้านกระบวนการ มีลักษณะเป็นข้อสอบเลือกตอบ ชนิด 4 ตัวเลือก

วิธีซิบเทสท์ (SIBTEST) หมายถึง วิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในกลุ่มวิธี IRT ที่พัฒนาโดย เชียลี และสโตท (Shealy & Stout, 1993) ซึ่งคำนวณดัชนีการทำหน้าที่ ต่างกันของข้อสอบจากค่าเฉลี่ยสัดส่วนการตอบข้อสอบถูกระหว่างผู้เข้าสอบกลุ่มอ้างอิงกับกลุ่มเปรียบเทียบในชุดแบบทดสอบที่ศึกษา (Studied Subtest) โดยใช้คะแนนการจับคู่เปรียบเทียบระหว่างกลุ่มผู้เข้าสอบจากชุดแบบทดสอบที่มีความตรง (Valid Subtest) แล้วทดสอบนัยสำคัญด้วยค่าสถิติ Z ในการศึกษาครั้งนี้ผู้วิจัยวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบโดยใช้โปรแกรม SIBTEST

วิธีการวิเคราะห์องค์ประกอบจำกัด (Restricted Factor Analysis: RFA) หมายถึง วิธีตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในกลุ่มวิธี ไอ-แสควร์ ที่พัฒนาโดย ออร์ท (Oort, 1998) เป็นวิธีที่ใช้คะแนนรวมของการตอบแทนความสามารถของผู้เข้าสอบ เกณฑ์การพิจารณาการทำหน้าที่ ต่างกันของข้อสอบ ได้แก่ ค่าดัชนีคัดแปลงโมเดล (Modification Indices: MI) และค่าพารามิเตอร์ที่คาดหวัง (Expected Parameter Change: EPC) แตกต่างจาก 0 หรือไม่ ที่ระดับนัยสำคัญ .05 ใน การศึกษาครั้งนี้ผู้วิจัยวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบโดยใช้โปรแกรม LISREL 8.50

วิธีจดดอยโลจิสติก (Logistic Regression: LR) หมายถึง วิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในกลุ่มวิธี ไอ-แสควร์ ที่พัฒนาโดย สวามินาธาน และโรเจอร์ (Swaminathan & Rogers, 1990) ซึ่งคำนวณดัชนีการทำหน้าที่ต่างกันของข้อสอบจากผลการตอบข้อสอบถูกระหว่างผู้เข้าสอบกลุ่มอ้างอิงกับกลุ่มเปรียบเทียบ เกณฑ์การพิจารณาการทำหน้าที่ต่างกันของข้อสอบ ได้แก่ ค่าอิทธิพลจากปฏิสัมพันธ์ระหว่างกลุ่มผู้เข้าสอบกับความสามารถของผู้เข้าสอบ หรือค่าอิทธิพลจากกลุ่มผู้เข้าสอบ และการทดสอบค่าสถิติ χ^2 ที่ระดับนัยสำคัญ .05 ใน การศึกษาครั้งนี้ผู้วิจัยวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบโดยใช้โปรแกรม SPSS

กลุ่มอ้างอิง (Reference Group: R) หมายถึง กลุ่มผู้เข้าสอบที่คาดว่าจะได้ประโยชน์จากการตอบข้อสอบทำหน้าที่ต่างกัน เป็นกลุ่มที่มีความน่าจะเป็นในการตอบข้อสอบได้ถูกต้องสูงกว่า ผู้เข้าสอบอีกกลุ่มนั้นทั้ง ๆ ที่มีความสามารถเท่ากัน ใน การวิจัยนี้กลุ่มอ้างอิง คือ เพศหญิง

กลุ่มเปรียบเทียบ (Focal Group: F) หมายถึง กลุ่มผู้เข้าสอบที่คาดว่าจะเสียประโยชน์จากการตอบข้อสอบทำหน้าที่ต่างกัน เป็นกลุ่มที่มีความน่าจะเป็นในการตอบข้อสอบได้ถูกต้องมากกว่าผู้เข้าสอบอีกกลุ่มนั่นเองทั้งๆ ที่มีความสามารถเท่ากัน ในการวิจัยนี้กลุ่มเปรียบเทียบ คือ เพศชาย

ขนาดกลุ่มตัวอย่าง (Sample Size) หมายถึง จำนวนผู้เข้าสอบในกลุ่มอ้างอิงและกลุ่มเปรียบเทียบที่ใช้ในการศึกษา มี 3 ขนาด คือ ขนาดใหญ่ (2,000 คน) ขนาดกลาง (1,000 คน) และขนาดเล็ก (300 คน)

ความยาวของแบบทดสอบ (Test Length) หมายถึง จำนวนข้อสอบในแบบทดสอบที่ใช้ในการศึกษา มี 3 ขนาด คือ 40 ข้อ 30 ข้อ และ 20 ข้อ

เกณฑ์การจับคู่ (Matching Criteria) หมายถึง คะแนนที่ใช้แทนความสามารถที่แท้จริงของผู้เข้าสอบสองกลุ่ม ซึ่งใช้ในการจับคู่เปรียบเทียบกลุ่มผู้เข้าสอบย่อย เพื่อตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ มี 2 เกณฑ์ คือ คะแนนรวมทั้งฉบับ และคะแนนแบบทดสอบย่อย

เกณฑ์การจับคู่คะแนนรวมทั้งฉบับ (Total Matching Criteria) หมายถึง คะแนนที่ใช้แทนความสามารถที่แท้จริงของผู้เข้าสอบ ซึ่งเป็นผลรวมของคะแนนที่ได้จากการสอบแบบทดสอบวิชาภาษาไทยทั้งฉบับ ในการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบครั้งนี้ใช้คะแนนรวมของผลการสอบทั้งฉบับ เป็นเกณฑ์จับคู่ระหว่างกลุ่มผู้เข้าสอบเพศหญิง กับกลุ่มผู้เข้าสอบเพศชาย

เกณฑ์การจับคู่คะแนนแบบทดสอบย่อย (Subtest Matching Criteria) หมายถึง คะแนนที่ใช้แทนความสามารถที่แท้จริงของผู้เข้าสอบ ซึ่งเป็นผลรวมของคะแนนที่ได้จากการสอบแบบทดสอบวิชาภาษาไทย ของแต่ละองค์ประกอบของการวัดแบบทดสอบย่อย ในการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบครั้งนี้ ใช้คะแนนรวมของผลการตอบข้อสอบ 2 องค์ประกอบ คือ ด้านโครงสร้างความรู้ และด้านกระบวนการ เป็นเกณฑ์จับคู่ระหว่างกลุ่มผู้เข้าสอบเพศหญิง กับกลุ่มผู้เข้าสอบเพศชาย โดยแยกวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบที่ละองค์ประกอบ

ประสิทธิภาพของการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ (Efficiency of Detecting DIF Items) หมายถึง ความถูกต้องของการระบุข้อสอบทำหน้าที่ต่างกัน จากการตรวจสอบโดยวิเคราะห์องค์ประกอบจำกัด และวิธีลดด้อยโลจิสติก ซึ่งพิจารณาได้จากอัตราความถูกต้อง และอัตราความคลาดเคลื่อน เมื่อเปรียบเทียบกับผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบด้วยวิธีชิปเพลท

อัตราความถูกต้อง (Correct Rates) หมายถึง สัดส่วนของจำนวนข้อสอบที่ตรวจสอบพบว่าถูกต้องตรงกับวิธีชิปเพลท มี 2 ประเภท คือ อัตราความถูกต้องประเภทที่ 1 และอัตราความถูกต้องประเภทที่ 2

อัตราความถูกต้องประเภทที่ 1 (Type I Correct Rate) หมายถึง สัดส่วนของจำนวนข้อสอบที่ตรวจพบว่าทำหน้าที่ไม่ต่างกัน ได้ถูกต้องตรงกับวิธีชิปเกสท์ ซึ่งคำนวณจากสัดส่วนของจำนวนข้อสอบทำหน้าที่ไม่ต่างกันถูกต้องต่อจำนวนข้อสอบทำหน้าที่ไม่ต่างกันทั้งหมดในแบบทดสอบ

อัตราความถูกต้องประเภทที่ 2 (Type II Correct Rate) หมายถึง สัดส่วนของจำนวนข้อสอบที่ตรวจพบว่าทำหน้าที่ต่างกัน ได้ถูกต้องตรงกับวิธีชิปเกสท์ ซึ่งคำนวณจากสัดส่วนของจำนวนข้อสอบทำหน้าที่ต่างกันถูกต้องต่อจำนวนข้อสอบทำหน้าที่ต่างกันทั้งหมดในแบบทดสอบ

อัตราความคลาดเคลื่อน (Error Rates) หมายถึง สัดส่วนของจำนวนข้อสอบที่ระบุผิดพลาด มี 2 ประเภท คือ ความคลาดเคลื่อนประเภทที่ 1 และความคลาดเคลื่อนประเภทที่ 2

ความคลาดเคลื่อนประเภทที่ 1 (Type I Error) หมายถึง การระบุผิดพลาดว่าข้อสอบทำหน้าที่ต่างกัน (False Positive) ทั้ง ๆ ที่ในความเป็นจริงแล้วข้อสอบทำหน้าที่ไม่ต่างกัน ซึ่งในการวิจัยนี้คำนวณ ได้จากสัดส่วนของจำนวนข้อสอบที่ระบุผิดพลาดว่าทำหน้าที่ต่างกันต่อจำนวนข้อสอบที่ทำหน้าที่ไม่ต่างกันทั้งหมดในแบบทดสอบ

ความคลาดเคลื่อนประเภทที่ 2 (Type II Error) หมายถึง การระบุผิดพลาดว่าข้อสอบทำหน้าที่ไม่ต่างกัน (False Negative) ทั้ง ๆ ที่ในความเป็นจริงแล้วข้อสอบทำหน้าที่ต่างกัน ซึ่งในการวิจัยนี้คำนวณ ได้จากสัดส่วนของจำนวนข้อสอบที่ระบุผิดพลาดว่าทำหน้าที่ไม่ต่างกันต่อจำนวนข้อสอบที่ทำหน้าที่ต่างกันทั้งหมดในแบบทดสอบ