

สำนักหอสมุด มหาวิทยาลัยบูรพา
๑.แผนสข อ.เมือง จ.ชลบุรี 20131

การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3
ด้วยวิธีการทดสอบอัตราส่วนไลค์ลิชฎ์ วิชชิปเพลท์ และวิธีแมนเทล-แฮนส์เซล

อรุณี แปลงกาย

b00254796

- ๙ ม.ค. 2562

381345

วิทยานิพนธ์นี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรวิทยาศาสตรมหาบัณฑิต
สาขาวิชาการวิจัยและสถิติทางวิทยาการปัญญา
วิทยาลัยวิทยาการวิจัยและวิทยาการปัญญา มหาวิทยาลัยบูรพา

กรกฎาคม 2561
ลิขสิทธิ์เป็นของมหาวิทยาลัยบูรพา

คณะกรรมการควบคุมวิทยานิพนธ์และคณะกรรมการสอบวิทยานิพนธ์ได้พิจารณา
วิทยานิพนธ์ของ อรุณี แปลงกาย ฉบับนี้แล้ว เห็นสมควรรับเป็นส่วนหนึ่งของการศึกษาตามหลักสูตร
วิทยาศาสตรมหาบัณฑิต สาขาวิชาการวิจัยและสถิติทางวิทยาการปัญญา ของมหาวิทยาลัยบูรพาได้

คณะกรรมการควบคุมวิทยานิพนธ์

.....
.....
(ดร. ปิยะทิพย์ ประดุจพรม)

คณะกรรมการสอบวิทยานิพนธ์

.....
.....
(รองศาสตราจารย์ ดร. เสรี ชัดแข็ม)
.....
(ดร. ปิยะทิพย์ ประดุจพรม)

.....
.....
(ผู้ช่วยศาสตราจารย์ ดร. วัชราวดี มากมี)
.....
(ดร. กนก พานทอง)

วิทยาลัยวิทยาการวิจัยและวิทยาการปัญญา อนุมัติให้รับวิทยานิพนธ์ฉบับนี้เป็นส่วนหนึ่ง
ของการศึกษาตามหลักสูตรวิทยาศาสตรมหาบัณฑิต สาขาวิชาการวิจัยและสถิติทางวิทยาการปัญญา
ของมหาวิทยาลัยบูรพา

.....
.....
(ผู้ช่วยศาสตราจารย์ ดร. สุชาดา กรเพชรปานี) และวิทยาการปัญญา
วันที่ 14 เดือน มกราคม พ.ศ. 2561

ประกาศคุณูปการ

วิทยานิพนธ์ เรื่อง การเปรียบเทียบผลการตรวจสอบการทำงานที่ต่างกันของข้อสอบ NT ขั้นประถมศึกษาปีที่ 3 ระหว่างวิธีการทดสอบอัตราส่วนไลค์ลิขิต วิธีชี้ปเศษท์ และวิธีแม่นเหล็กน์เซลลดับบันน์ สำเร็จลุล่วงได้ด้วยความกรุณาจาก ดร.ปิยะทิพย์ ประดุจพร อาจารย์ที่ปรึกษาหลัก ซึ่งเฝ้าอบรมบ่มเพาะ ถ่ายทอดความรู้ ประสิทธิ์ประสาทวิชา และให้คำแนะนำเป็นอย่างดี ทั้งในส่วนของเนื้อหาวิทยานิพนธ์และวิธีการปรับตัว ตลอดระยะเวลาที่ดำเนินการวิจัย ทำให้ผู้วิจัยสามารถพัฒนาอุปสรรคต่าง ๆ มาได้ด้วยดี และความเอาใจใส่ ช่วยเหลือในทุกขั้นตอนของการทำวิทยานิพนธ์ ผู้วิจัยรู้สึกซาบซึ้งและระลึกถึงพระคุณอันหาที่เบรียบมีได้ในครั้งนี้เสมอ จึงขอกราบขอบพระคุณเป็นอย่างสูงไว้ ณ โอกาสนี้

ขอกราบขอบพระคุณผู้ช่วยศาสตราจารย์ ดร.สุชาดา arpichapane คณบดีวิทยาลัย วิทยาการวิจัยและวิทยาการปัญญา สำหรับการสนับสนุน ส่งเสริมกิจกรรมเชิงวิชาการ และการให้แนวคิดในการทำงานอย่างเป็นระบบ ทำให้ผู้วิจัยมีความมุ่งมั่นในการทำงานจนประสบความสำเร็จ รวมถึงเจ้าหน้าที่ของวิทยาลัยวิทยาการวิจัยและวิทยาการปัญญาทุกคนที่ให้บริการ และการประสานงานที่ดีเยี่ยมอยู่เสมอ

ขอกราบขอบพระคุณบุคลากรทุกท่านและหน่วยงานที่เกี่ยวข้องทุกแห่งสำหรับการอำนวย ความสะดวกตลอดการดำเนินการวิจัย สำหรับการเสียสละเวลาอันมีค่าเพื่อให้ข้อมูลสำคัญสำหรับ การศึกษาในครั้งนี้ นอกจากนี้ ขอขอบคุณมิตรภาพของเพื่อนทุกกลุ่ม สำหรับการสนับสนุน และ กระตุ้นให้ผู้วิจัยมีแรงกาย แรงใจมากเพียงพอที่จะดำเนินการวิจัยจนผ่านพ้นไปได้ด้วยดี ผู้วิจัย ขอขอบคุณในมิตรภาพที่มีอปปะเสมอมา

สุดท้ายนี้ ผู้วิจัยกราบขอบพระคุณคณาจารย์ทุกท่าน ที่ได้ประสิทธิ์ประสาทความรู้ และขอบคุณครอบครัวที่คอยให้ความช่วยเหลือเป็นกำลังใจตลอดมา ประโยชน์ของวิทยานิพนธ์ฉบับนี้ ผู้วิจัยขอมอบเป็นกตัญญูกตเวทิตาแด่บุพการี บุรพาจารย์ และผู้มีพระคุณทุกท่านทั้งในอดีต และ ปัจจุบันที่ทำให้ข้าพเจ้าเป็นผู้มีการศึกษา และประสบความสำเร็จมาจนทราบเท่าทุกวันนี้

อรุณี แปลงกาญ

56910403: สาขาวิชา: การวิจัยและสถิติทางวิทยาการปัญญา;

วท.ม. (การวิจัยและสถิติทางวิทยาการปัญญา)

คำสำคัญ: การทำหน้าที่ต่างกันของข้อสอบ/ ข้อสอบ NT/ วิธีการทดสอบอัตราส่วนไลค์ลิชญด/ วิธีซิปเทสท์/ วิธีแมนเทล-แฮนส์เซล

อรุณ แปลงกาย: การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้วยวิธีการทดสอบอัตราส่วนไลค์ลิชญด วิธีซิปเทสท์ และวิธีแมนเทล-แฮนส์เซล (A COMPARISON OF THE IRT LIKELIHOOD RATIO TEST, SIBTEST, AND THE MANTEL-HAENSZEL METHOD FOR DETECTING DIFFERENTIAL ITEM FUNCTIONING USING RESULTS FROM THE GRADE 3 NATIONAL TEST) คณะกรรมการควบคุมวิทยานิพนธ์: ปิยะทิพย์ ประดุจพร, Ph.D., 256 หน้า. ปี พ.ศ. 2561.

การวิจัยนี้มีวัตถุประสงค์เพื่อวิเคราะห์คุณภาพของข้อสอบของแบบทดสอบระดับชาติ (NT) และตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน และด้านเหตุผล ภายใต้เงื่อนไขกลุ่มตัวอย่างขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน) ด้วยวิธีการทดสอบอัตราส่วนไลค์ลิชญด วิธีซิปเทสท์ และวิธีแมนเทล-แฮนส์เซล การดำเนินการวิจัยแบ่งเป็น 3 ระยะ ดังนี้ 1) วิเคราะห์คุณภาพของแบบทดสอบระดับชาติ ทั้ง 3 ด้าน 2) ตรวจสอบการทำหน้าที่ต่างกันของข้อสอบทั้ง 3 วิธี และ 3) เปรียบเทียบผลการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบด้วยวิธีการตรวจสอบ 3 วิธี ข้อมูลที่นำมาใช้วิเคราะห์เป็นข้อมูลทุติยภูมิ ซึ่งเป็นผลการตอบแบบทดสอบระดับชาติของนักเรียน ชั้นประถมศึกษาปีที่ 3 ปีการศึกษา 2556 จำนวนทั้งหมด 706,372 คน

ผลการวิจัยปรากฏว่า

1. แบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 มีค่าความยากของข้อสอบ (b) อยู่ในระดับค่อนข้างยาก มีค่าอำนาจจำแนกของข้อสอบ (a) อยู่ในระดับที่สามารถจำแนกผู้สอบได้ดีมาก และมีค่าโอกาสในการเดาของข้อสอบ (c) ไม่เกิน .30

2. การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบทั้ง 3 ด้าน พบร่วมกับกลุ่มตัวอย่างขนาดใหญ่ ส่งผลให้เกิดการทำหน้าที่ต่างกันของข้อสอบ จะได้เปรียบในการตอบข้อสอบ ด้านภาษา ด้านคำนวน และด้านเหตุผล มากกว่าขนาดเล็ก และขนาดกลาง โดยวิธีแมนเทล-แฮนส์เซล ตรวจพบข้อสอบทำหน้าที่ต่างกัน จำนวนมากที่สุด คิดเป็นร้อยละ 34 ของข้อสอบทั้งหมด รองลงมาคือ วิธีซิปเทสท์ และวิธีการทดสอบอัตราส่วนไลค์ลิชญด คิดเป็นร้อยละ 14 ของข้อสอบทั้งหมด

3. การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ พบร่วมกับ วิธีแมนเทล-แฮนส์เซล ตรวจพบ DIF มากกว่าวิธีการทดสอบอัตราส่วนไลค์ลิชญดในด้านภาษา ด้านคำนวน และด้านเหตุผล คิดเป็นร้อยละ 80, 13 และ 80 และวิธีแมนเทล-แฮนส์เซล ตรวจพบ DIF มากกว่า วิธีซิปเทสท์ ทั้ง 3 ด้าน คิดเป็นร้อยละ 87, 13 และ 73 วิธีการทดสอบอัตราส่วนไลค์ลิชญด ตรวจพบ DIF มากกว่าวิธีซิปเทสท์ในด้านภาษา คิดเป็นร้อยละ 7 และวิธีการทดสอบอัตราส่วนไลค์ลิชญด ตรวจพบ DIF เท่ากันกับวิธีซิปเทสท์ ในด้านคำนวน ส่วนวิธีซิปเทสท์ ตรวจพบ DIF มากกว่า วิธีการทดสอบอัตราส่วนไลค์ลิชญด ด้านเหตุผล คิดเป็นร้อยละ 7 ($p < .05$)

56910403: MAJOR: RESEARCH AND STATISTICS IN COGNITIVE SCIENCE;
M.Sc. (RESEARCH AND STATISTICS IN COGNITIVE SCIENCE)

KEYWORDS: DIFFERENTIAL ITEM FUNCTIONING / NATIONAL TESTS: NT/ IRT-LR
METHOD / SIBTEST METHOD / MANTEL-HAENSZEL METHOD

ARUNEE PLANGKAY: A COMPARISON OF THE IRT LIKELIHOOD RATIO TEST,
SIBTEST, AND THE MANTEL-HAENSZEL METHOD FOR DETECTING DIFFERENTIAL ITEM
FUNCTIONING USING RESULTS FROM THE GRADE 3 NATIONAL TEST. ADVISORY
COMMITTEE: PIYATHIP PRADUJUPROM, Ph.D. 256 P. 2018.

The objectives of this research were to analyze the quality of test items from National Tests (NT), and to determine if differential item functioning (DIF) was present. Three subjects were involved: Literacy, Numeracy, and Reasoning. Different sample sizes were used, from small ($n=300$), medium ($n=1,000$), to large ($n=2,000$). Three DIF methods were compared: IRT-LR, SIBTEST, and Mantel-Haenszel. The research procedures were divided into three phases: 1) Analyzing the item quality of NT for the three subjects; 2) Testing DIF detection of the items in NT using IRT-LR, SIBTEST, and Mantel-Haenszel methods; 3) Comparing the results of the three methods of DIF using secondary data from 706,372 NT Grade three students in the academic year 2013.

Results were as follows:

1. Items on the NT had relatively high difficult levels, very good discrimination levels, and guessing parameters which did not exceed .30.
2. The examination of results in the three subjects revealed that the larger samples had higher scores on Literacy, Numeracy, and Reasoning. The Mantel-Haenszel method detected the greatest number of DIF items across all three subjects, finding DIF on 34% of the test items; SIBTEST and IRT-LR both found DIF in 14% of the items.
3. Comparison of the DIF test results revealed that the Mantel-Haenszel method outperformed the IRT-LR method in terms of DIF detection, namely 80% for Literacy, 13% for Numeracy, and 80% for Reasoning. The Mantel-Haenszel method also outperformed the SIBTEST method in terms of DIF detection, namely 87% for Literacy, 13% for Numeracy, and 73% for Reasoning subjects. The IRT-LR method outperformed the SIBTEST method in terms of DIF detection, namely 7% for Literacy, and same DIF as for Numeracy subjects. Also, the SIBTEST method outperformed the IRT-LR method in terms of DIF detection on the Reasoning subject (7%) ($p < .05$).

สารบัญ

	หน้า
บทคัดย่อภาษาไทย.....	๑
บทคัดย่อภาษาอังกฤษ.....	๒
สารบัญ.....	๓
สารบัญตาราง.....	๔
สารบัญภาพ.....	๗
บทที่	
1 บทนำ.....	1
ความเป็นมาและความสำคัญของปัญหา.....	1
วัตถุประสงค์ของการวิจัย.....	4
กรอบแนวคิดการวิจัย.....	5
สมมติฐานของการวิจัย.....	7
ประโยชน์ที่คาดว่าจะได้รับจากการวิจัย.....	7
ขอบเขตของการวิจัย.....	8
นิยามศัพท์เฉพาะ.....	8
2 เอกสารและงานวิจัยที่เกี่ยวข้อง.....	11
ตอนที่ 1 การทดสอบความสามารถพื้นฐานของผู้เรียนระดับชาติ (National Testing: NT) และงานวิจัยที่เกี่ยวข้อง.....	11
ตอนที่ 2 ทฤษฎีการตอบสนองข้อสอบ (Item Response Theory: IRT) และงานวิจัยที่เกี่ยวข้อง.....	18
ตอนที่ 3 การทำหน้าที่ต่างกันของข้อสอบ (Differential item Functioning: DIF) และงานวิจัยที่เกี่ยวข้อง.....	34
ตอนที่ 4 การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบด้วยวิธีอัตราส่วนไลค์ลิลูด และงานวิจัยที่เกี่ยวข้อง.....	47
ตอนที่ 5 การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบด้วยวิธีชี้ปีเทส์ และงานวิจัยที่เกี่ยวข้อง.....	57
ตอนที่ 6 การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบด้วยวิธีเมนเทล-แ xen's เชล และงานวิจัยที่เกี่ยวข้อง.....	64
3 วิธีดำเนินการวิจัย.....	73
ระยะที่ 1 การวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 โดยใช้ หลักการของทฤษฎีการตอบสนองข้อสอบ แบบ 3 พารามิเตอร์ ภายใต้ เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน).....	75

สารบัญ (ต่อ)

บทที่	หน้า
ระยะที่ 2 การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้วยวิธีการทดสอบอัตราส่วนไลค์ลิชูด วิธีชิปเทสท์ และวิธีแมนเทล-แ昏ส์เซล ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน).....	89
ระยะที่ 3 การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ ในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้วยวิธีการทดสอบ อัตราส่วนไลค์ลิชูด วิธีชิปเทสท์ และวิธีแมนเทล-แ昏ส์เซล ภายใต้เงื่อนไข ¹ ขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน).....	108
4 ผลการวิจัย.....	110
ตอนที่ 1 ผลการวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 โดยใช้หลักการของทฤษฎีการตอบสนองข้อสอบ แบบ 3 พารามิเตอร์ ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน).....	111
ตอนที่ 2 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้วยวิธีการทดสอบอัตราส่วนไลค์ลิชูด วิธีชิปเทสท์ และวิธีแมนเทล-แ昏ส์เซล ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน).....	126
ตอนที่ 3 ผลการเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ ในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้วยวิธีการทดสอบ อัตราส่วนไลค์ลิชูด วิธีชิปเทสท์ และวิธีแมนเทล-แ昏ส์เซล ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน).....	142
5 สรุปและอภิปรายผล.....	159
สรุปผลการวิจัย.....	159
อภิปรายผล.....	161
ข้อเสนอแนะในการนำผลการวิจัยไปใช้.....	163
ข้อเสนอแนะสำหรับการวิจัยต่อไป.....	163
บรรณานุกรม.....	164

สารบัญ (ต่อ)

บทที่	หน้า
ภาคผนวก.....	169
ภาคผนวก ก หนังสือรับรองผลการพิจารณาจริยธรรมการวิจัย.....	170
ภาคผนวก ข หนังสือขอความอนุเคราะห์ข้อมูลเพื่อการวิจัย.....	172
ภาคผนวก ค ตัวอย่างแสดงข้อมูลผลการตอบข้อสอบ NT ปีการศึกษา 2556 ชั้นประถมศึกษาปีที่ 3.....	174
ภาคผนวก จ ตัวอย่าง Print Out การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ ด้วยวิธี IRT-LR.....	194
ภาคผนวก ฉ ตัวอย่าง Print Out การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ ด้วยวิธี SIBTEST.....	209
ภาคผนวก ช ตัวอย่าง Print Out การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ ด้วยวิธี MH.....	241
ประวัติย่อของผู้วิจัย.....	256

สารบัญตาราง

ตารางที่	หน้า
2-1 วิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบที่มีการตรวจให้ค่าคะแนนแบบทวิภาค (Dichotomous DIF) และพหุวิภาค (Polytomous DIF).....	41
2-2 ผลการตอบข้อสอบข้อหนึ่งระหว่างผู้สอบกลุ่มอ้างอิงและกลุ่มเปรียบเทียบที่มีค่าแనนรวม j.....	64
2-3 สัดส่วนของผลการตอบข้อสอบข้อหนึ่งระหว่างผู้สอบกลุ่มอ้างอิงและกลุ่มเปรียบเทียบที่มีค่าแnan j.....	65
4-1 ผลการวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านภาษา จำนวน 30 ข้อ โดยใช้ทฤษฎีการตอบสนองข้อสอบ (IRT) แบบ 3 พารามิเตอร์ กลุ่มตัวอย่างขนาดเล็ก (300 คน)	111
4-2 ผลการวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านคำนวณ จำนวน 30 ข้อ โดยใช้ทฤษฎีการตอบสนองข้อสอบ (IRT) แบบ 3 พารามิเตอร์ กลุ่มตัวอย่างขนาดเล็ก (300 คน)	113
4-3 ผลการวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านเหตุผล จำนวน 30 ข้อ โดยใช้ทฤษฎีการตอบสนองข้อสอบ (IRT) แบบ 3 พารามิเตอร์ กลุ่มตัวอย่างขนาดเล็ก (300 คน)	114
4-4 ผลการวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านภาษา จำนวน 30 ข้อ โดยใช้ทฤษฎีการตอบสนองข้อสอบ (IRT) แบบ 3 พารามิเตอร์ กลุ่มตัวอย่างขนาดกลาง (1,000 คน)	115
4-5 ผลการวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านคำนวณ จำนวน 30 ข้อ โดยใช้ทฤษฎีการตอบสนองข้อสอบ (IRT) แบบ 3 พารามิเตอร์ กลุ่มตัวอย่างขนาดกลาง (1,000 คน).....	117
4-6 ผลการวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านเหตุผล จำนวน 30 ข้อ โดยใช้ทฤษฎีการตอบสนองข้อสอบ (IRT) แบบ 3 พารามิเตอร์ กลุ่มตัวอย่างขนาดกลาง (1,000 คน).....	118
4-7 ผลการวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านภาษา จำนวน 30 ข้อ โดยใช้ทฤษฎีการตอบสนองข้อสอบ (IRT) แบบ 3 พารามิเตอร์ กลุ่มตัวอย่างขนาดใหญ่ (2,000 คน).....	119
4-8 ผลการวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านคำนวณ จำนวน 30 ข้อ โดยใช้ทฤษฎีการตอบสนองข้อสอบ (IRT) แบบ 3 พารามิเตอร์ กลุ่มตัวอย่างขนาดใหญ่ (2,000 คน).....	121

สารบัญตาราง (ต่อ)

ตารางที่	หน้า
4-9 ผลการวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านเหตุผล จำนวน 30 ข้อ โดยใช้ทฤษฎีการตอบสนองข้อสอบ (IRT) แบบ 3 พารามิเตอร์ กลุ่มตัวอย่างขนาดใหญ่ (2,000 คน).....	122
4-10 สรุปผลการวิเคราะห์คุณภาพของแบบทดสอบรายข้อ และหั้งฉบับชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน และด้านเหตุผล.....	125
4-11 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ระดับชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี IRT-LR ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างขนาดเล็ก (300 คน).....	126
4-12 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ระดับชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี IRT-LR ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างขนาดกลาง (1,000 คน)	128
4-13 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ระดับชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี IRT-LR ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างขนาดใหญ่ (2,000 คน).....	129
4-14 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ระดับชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี SIBTEST ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างขนาดเล็ก (300 คน).....	131
4-15 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ระดับชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี SIBTEST ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างขนาดกลาง (1,000 คน).....	132
4-16 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ระดับชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี SIBTEST ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างขนาดใหญ่ (2,000 คน).....	134
4-17 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ระดับชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี MH ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างขนาดเล็ก (300 คน).....	136
4-18 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ระดับชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี MH ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างขนาดกลาง (1,000 คน).....	137
4-19 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ระดับชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี MH ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างขนาดใหญ่ (2,000 คน).....	139

สารบัญตาราง (ต่อ)

ตารางที่	หน้า
4-20 ผลการวิเคราะห์คุณภาพของแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน และด้านเหตุผล.....	141
4-21 การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบ ระดับชาติ ด้านภาษา ระหว่างวิธี IRT-LR กับวิธี SIBTEST.....	142
4-22 การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบ ระดับชาติ ด้านคำนวน ระหว่างวิธี IRT-LR กับ วิธี SIBTEST.....	144
4-23 การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบ ระดับชาติ ด้านเหตุผล ระหว่างวิธี IRT-LR กับวิธี SIBTEST.....	145
4-24 การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบ ระดับชาติ ด้านภาษา ระหว่างวิธี IRT-LR กับวิธี MH	147
4-25 การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบ ระดับชาติ ด้านคำนวน ระหว่างวิธี IRT-LR กับวิธี MH	148
4-26 การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบ ระดับชาติ ด้านเหตุผล ระหว่างวิธี IRT-LR กับวิธี MH	150
4-27 การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบ ระดับชาติ ด้านภาษา ระหว่างวิธี SIBTEST กับวิธี MH.....	151
4-28 การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบ ระดับชาติ ด้านคำนวน ระหว่างวิธี SIBTEST กับวิธี MH.....	153
4-29 การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบ ระดับชาติ ด้านเหตุผล ระหว่างวิธี SIBTEST กับวิธี MH	154
4-30 การเปรียบเทียบการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบด้วย วิธีการทดสอบ อัตราส่วนໄลค์ลิญด์ วิธีซิปเพสท์ และวิธีแมนเนล-แยนส์เซล ภายใต้เงื่อนไขขนาดกลุ่ม ตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน).....	156
ค-1 แสดงข้อมูลผลการตอบข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านภาษา จำนวน 300 คน.....	175
ค-2 แสดงข้อมูลผลการตอบข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านคำนวน จำนวน 300 คน.....	176
ค-3 แสดงข้อมูลผลการตอบข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านเหตุผล จำนวน 300 คน.....	177
ค-4 แสดงข้อมูลผลการตอบข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านภาษา จำนวน 1,000 คน.....	178

สารบัญตาราง (ต่อ)

ตารางที่	หน้า
ค-5 แสดงข้อมูลผลการตอบข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านคำนวณ จำนวน 1,000 คน.....	179
ค-6 แสดงข้อมูลผลการตอบข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านเหตุผล จำนวน 1,000 คน.....	180
ค-7 แสดงข้อมูลผลการตอบข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านภาษา จำนวน 2,000 คน.....	181
ค-8 แสดงข้อมูลผลการตอบข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านคำนวณ จำนวน 2,000 คน.....	182
ค-9 แสดงข้อมูลผลการตอบข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านเหตุผล จำนวน 2,000 คน.....	183
ง-10 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวณ ด้านเหตุผล ด้วยวิธี IRT-LR ภายใต้เงื่อนไขขนาด กลุ่มตัวอย่างขนาดใหญ่ (300 คน).....	185
ง-11 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวณ ด้านเหตุผล ด้วยวิธี IRT-LR ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง ขนาดใหญ่ (1,000 คน).....	186
ง-12 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวณ ด้านเหตุผล ด้วยวิธี IRT-LR ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง ขนาดใหญ่ (2,000 คน).....	187
ง-13 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวณ ด้านเหตุผล ด้วยวิธี SIBTEST ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง ขนาดใหญ่ (300 คน).....	188
ง-14 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวณ ด้านเหตุผล ด้วยวิธี SIBTEST ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง ขนาดใหญ่ (1,000 คน).....	189
ง-15 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวณ ด้านเหตุผล ด้วยวิธี SIBTEST ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง ขนาดใหญ่ (2,000 คน).....	190
ง-16 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวณ ด้านเหตุผล ด้วยวิธี MH ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง ขนาดใหญ่ (300 คน).....	191

สารบัญตาราง (ต่อ)

ตารางที่	หน้า
ง-17 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวณ ด้านเหตุผล ด้วยวิธี MH ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง ขนาดใหญ่ (1,000 คน).....	192
ง-18 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวณ ด้านเหตุผล ด้วยวิธี MH ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง ขนาดใหญ่ (2,000 คน).....	193

สารบัญภาพ

ภาพที่	หน้า
1-1 กรอบแนวคิดการวิจัย.....	6
2-1 ข้อสอบทำหน้าที่ต่างกันแบบเอกสาร (Uniform DIF).....	36
2-2 ข้อสอบทำหน้าที่ต่างกันแบบไม่มีทิศทาง (Non-unidirectional DIF).....	37
2-3 ข้อสอบทำหน้าที่ต่างกันแบบมีทิศทางเดียวกัน (Unidirectional DIF).....	37
2-4 การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบและแบบสอบถามโดยใช้เทคนิควิเคราะห์สมการรถถอย.....	38
3-1 ขั้นตอนการดำเนินงานวิจัย.....	74
3-2 ขั้นตอนการการวิเคราะห์คุณภาพของข้อสอบ NT ระดับชั้นประถมศึกษาปีที่ 3 โดยใช้หลักการของทฤษฎีการตอบสนองข้อสอบ แบบ 3 พารามิเตอร์.....	75
3-3 ตัวอย่างเปิดไฟล์ Lertap5.xls.....	76
3-4 ตัวอย่างข้อมูลที่ต้องการจะวิเคราะห์ไฟล์ Lertap5.xls.....	77
3-5 ตัวอย่างตัวอย่างข้อมูลที่บันทึกไว้ใน Excel ที่จะวิเคราะห์ไฟล์	77
3-6 ตัวอย่างข้อมูลแผ่นงาน CCs ที่มีรายคำตอบ.....	78
3-7 ตัวอย่างเลือกเมนู Interpret	78
3-8 ตัวอย่างเลือกเมนู Ok	79
3-9 ตัวอย่างแสดงผลลัพธ์ Interpret	79
3-10 ตัวอย่างเลือกเมนู Elmillion และคลิก Ok	80
3-11 ตัวอย่างแสดงภาพหน้าจอหลังจากเลือกเมนู OK	80
3-12 ตัวอย่างแสดงภาพหน้าจอหลังจากเลือกเมนู More เลือก Item Scores and Correlation	81
3-13 ตัวอย่างแสดงภาพหน้าจอหลังจากเลือกเมนู More เลือก Item Scores and Correlation และคลิก OK.....	81
3-14 ตัวอย่างแสดงภาพหน้าจอหลังจากคลิก OK และคลิก OK	82
3-15 ตัวอย่างแสดงภาพหน้าจอ Output ของการ Run เมนู Item Scores and Correlation.....	82
3-16 ตัวอย่างโปรแกรม Xcalibre.....	83
3-17 ตัวอย่างการเลือกไฟล์ Data.txt สำหรับวิเคราะห์โปรแกรม Xcalibre.....	84
3-18 ตัวอย่างไฟล์ที่ต้องการเลือกวิเคราะห์โปรแกรม Xcalibre ใน Data Matrix File.....	84
3-19 ตัวอย่างการเลือกไฟล์ ICF.txt วิเคราะห์โปรแกรม Xcalibre.....	85
3-20 ตัวอย่างไฟล์ที่ต้องการเลือกวิเคราะห์โปรแกรม Xcalibre ใน Item Control File.....	85
3-21 ตัวอย่างระบุไฟล์ที่ต้องการจะเก็บช่อง Output File.....	86
3-22 ตัวอย่างไฟล์ที่ต้องการเลือกวิเคราะห์โปรแกรม Xcalibre ในช่อง Output File.....	86

สารบัญภาพ (ต่อ)

ภาพที่	หน้า
3-23 ตัวอย่างไฟล์ที่ต้องการเลือกวิเคราะห์โปรแกรม Xcalibre.....	87
3-24 ตัวอย่างไฟล์ที่ต้องการเลือกวิเคราะห์โปรแกรม Xcalibre เลือกที่ Yes.....	87
3-25 แสดงตัวอย่างไฟล์ ผลการวิเคราะห์ด้วยโปรแกรม Xcalibre.....	88
3-26 ขั้นตอนการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ระดับชั้นประถมศึกษาปีที่ 3.....	89
3-27 ตัวอย่างการจัดไฟล์ข้อมูลสำหรับวิเคราะห์ด้วยวิธี IRT-LR ในรูปแบบไฟล์ .ssig.....	90
3-28 ตัวอย่างการเรียกไฟล์ข้อมูลสำหรับวิเคราะห์ด้วยวิธี IRT-LR ในรูปแบบไฟล์ .ssig.....	91
3-29 ตัวอย่างการเปิดไฟล์ข้อมูลสำหรับวิเคราะห์ด้วยวิธี IRT-LR ในรูปแบบไฟล์ .ssig.....	92
3-30 ตัวอย่างการเลือก Analysis การวิเคราะห์ข้อมูล.....	92
3-31 ตัวอย่างการเลือก Yes เพื่อการวิเคราะห์ข้อมูล.....	93
3-32 ตัวอย่างการเลือก Unidimensional Analysis เพื่อการวิเคราะห์ข้อมูล.....	93
3-33 ตัวอย่างการเลือกใช้ชี้ไฟล์เพื่อการวิเคราะห์ข้อมูล.....	94
3-34 ตัวอย่างการเลือก Group และเลือก Gender เพื่อการวิเคราะห์ข้อมูล.....	94
3-35 ตัวอย่างการเลือก Item และเลือกจำนวนข้อสอบทั้งหมดเพื่อการวิเคราะห์ข้อมูล.....	95
3-36 ตัวอย่างการจัดไฟล์ข้อมูลสำหรับวิเคราะห์ด้วยวิธี SIBTEST ในรูปแบบไฟล์ .dat.....	96
3-37 ตัวอย่างการบันทึกไฟล์ข้อมูลสำหรับวิเคราะห์ด้วยวิธี SIBTEST ในรูปแบบไฟล์ .dat	96
3-38 ตัวอย่างการจัดไฟล์ข้อมูลสำหรับวิเคราะห์ด้วยวิธี SIBTEST ในรูปแบบไฟล์ .dat ในรูปแบบ NotePad.....	97
3-39 ตัวอย่างการจัดไฟล์ข้อมูลที่ลับจุลภาค (,) ออกและขิดของข้อมูลสำหรับวิเคราะห์ด้วย วิธี SIBTEST ในรูปแบบไฟล์ .dat.....	98
3-40 ตัวอย่างการวิเคราะห์ด้วยวิธี SIBTEST ในรูปแบบไฟล์ .dat เรียกไฟล์ข้อมูลวิเคราะห์.....	99
3-41 ตัวอย่างการวิเคราะห์ด้วยวิธี SIBTEST ในรูปแบบไฟล์ .dat เมื่อไฟล์ข้อมูลเข้าใน โปรแกรม.....	100
3-42 ตัวอย่างการวิเคราะห์ด้วยวิธี SIBTEST ในรูปแบบไฟล์ .dat.....	101
3-43 ตัวอย่างการเลือกข้อมูลวิเคราะห์ด้วยวิธี SIBTEST ในรูปแบบไฟล์ .dat	102
3-44 ตัวอย่างการจัดไฟล์ข้อมูลสำหรับวิเคราะห์ด้วยวิธี Mantel-Haenszel ตารางการณ์จร แบบ 2x2.....	103
3-45 ตัวอย่างคำสั่งที่ Analysis การวิเคราะห์ด้วยวิธี Mantel-Haenszel โดยโปรแกรม SPSS.	104
3-46 ตัวอย่างคำสั่งที่ Crosstabs การวิเคราะห์ด้วยวิธี Mantel-Haenszel โดยโปรแกรม SPSS.....	104
3-47 ตัวอย่างคำสั่งเรียกไฟล์ที่ต้องการจะวิเคราะห์ด้วยวิธี Mantel-Haenszel โดยโปรแกรม SPSS.....	105

สารบัญภาพ (ต่อ)

ภาพที่	หน้า
3-48 ตัวอย่างคำสั่งเลือก Chi-Square การวิเคราะห์ด้วยวิธี Mantel-Haenszel โดยโปรแกรม SPSS.....	106
3-49 ตัวอย่างผลการวิเคราะห์ Chi-Square วิธี Mantel-Haenszel.....	106
3-50 ขั้นตอนการเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ.....	108

บทที่ 1

บทนำ

ความเป็นมาและความสำคัญของปัญหา

การวัดประเมินผลการศึกษามีความสำคัญกับการศึกษา ซึ่งเป็นตัวบ่งชี้ที่สำคัญของความสำเร็จของการจัดการศึกษาของประเทศไทยที่พัฒนาแล้วทำให้องค์กรระหว่างประเทศของโลกได้ให้ความสำคัญกับการประเมินผลการศึกษา โดยกำหนดการทดสอบ PISA (Programme for International Student Assessment: PISA) ขึ้นมาเพื่อเป็นการประเมินผลการศึกษาของแต่ละประเทศว่าประเทศใดสามารถที่จะพัฒนาผู้เรียนให้รู้จักคิดวิเคราะห์ มีเหตุผล และรู้จักแก้ปัญหาอยู่ในระดับใด (ปนัดดา หัสสรา, 2557, หน้า 1) สำหรับประเทศไทยจึงได้กำหนดหลักสูตรแกนกลางการศึกษาขั้นพื้นฐาน พุทธศักราช 2551 ขึ้นมาเพื่อการวัดผลและประเมินผลการเรียนรู้ โดยแบ่งเป็น 4 ระดับ คือ 1) ระดับชั้นเรียน 2) ระดับสถานศึกษา 3) ระดับเขตพื้นที่การศึกษา และ 4) ระดับชาติ ซึ่งการประเมินระดับชาติจะประเมินในระดับชั้นประถมศึกษาปีที่ 3 ชั้นประถมศึกษาปีที่ 6 ชั้นมัธยมศึกษาปีที่ 3 และชั้นมัธยมศึกษาปีที่ 6 (สำนักทดสอบทางการศึกษา สำนักงานคณะกรรมการการศึกษาขั้นพื้นฐาน, 2557, หน้า 2-4)

การวัดผลประเมินผลทางการศึกษาเพื่อการตรวจสอบความรู้ความสามารถของนักเรียน หรือเป็นการวัดคุณลักษณะที่ต้องการวัดอยู่ในระดับใด (พิรญา สูงเนิน เสรี ชัดแข็ง และสมโภชน์ อนงกสุข 2552, หน้า 51) และการทดสอบระดับชาติ (National Testing: NT) เป็นการประเมินคุณภาพทางการศึกษาที่มีจุดมุ่งหมาย เพื่อประเมินคุณภาพการศึกษาขั้นพื้นฐาน โดยจะนำข้อมูลเพื่อจัดทำแผนการพัฒนาให้นักเรียนได้มีความสามารถอ่านออกเขียนได้ รู้จักคิดวิเคราะห์ สังเคราะห์ และเป็นการวัดผลสัมฤทธิ์ทางการเรียนของผู้เรียนแต่ละช่วงชั้น จะมีการทดสอบที่ชั้นเรียนสูงสุด ของช่วงชั้น คือ ระดับชั้นประถมศึกษาปีที่ 3 ซึ่งการทดสอบจะส่งผลสะท้อนให้เห็นถึงการประสบความสำเร็จของการพัฒนาคุณภาพการศึกษาที่สอดคล้องกับพระราชบัญญัติการศึกษาแห่งชาติ โดยมีสำนักงานคณะกรรมการการศึกษาขั้นพื้นฐานเป็นหน่วยงานที่รับผิดชอบในการจัดการศึกษา ซึ่งจะจัดสอบโดยสำนักงานคณะกรรมการการศึกษาขั้นพื้นฐาน ในปีการศึกษา 2555 โดยสำนักทดสอบทางการศึกษาได้กำหนดนโยบายจัดสอบระดับชาติ (NT) ในระดับชั้นประถมศึกษาปีที่ 3 มีทั้งหมด 3 ด้าน คือ ความสามารถด้านภาษา (Literacy Ability) ความสามารถด้านคำนวณ (Numeracy Ability) และความสามารถด้านเหตุผล (Reason Ability) การวัดทั้ง 3 ด้านนี้เพื่อให้สอดคล้องกับจุดเน้นการพัฒนาคุณภาพผู้เรียน คือ ชั้นประถมศึกษาปีที่ 1-3 นักเรียนมีทักษะความสามารถในการอ่านออกเขียนได้ คิดเลขเป็น และมีทักษะการคิดขั้นพื้นฐาน ซึ่งลักษณะของแบบทดสอบจะเป็นแบบปรนัยเลือกตอบ แบบ 4 ตัวเลือก ข้อสอบด้านภาษา มีจำนวนข้อสอบ 30 ข้อ ใช้เวลาในการทดสอบ 50 นาที ข้อสอบด้านคำนวณ มีจำนวนข้อสอบ 30 ข้อ ใช้เวลาในการทดสอบ 50 นาที (สำนักทดสอบทางการศึกษา สำนักงานคณะกรรมการการศึกษาขั้นพื้นฐาน, 2557, หน้า 1-4)

แบบทดสอบจึงเป็นเครื่องมือที่ใช้การวัดผลประเมินผลทางการศึกษา โดยจะเป็นการตรวจสอบว่าผู้สอบบัน្តจะมีคุณลักษณะแห่งหรือความสามารถอยู่ในระดับดี การสร้างเครื่องมือในการทดสอบและการตรวจสอบคุณภาพของแบบทดสอบจะต้องคำนึงถึงความตรงตามเนื้อหา (Content Validity) ความตรงตามเกณฑ์ (Criterion Validity) และความตรงเชิงโครงสร้าง (Construct Validity) เป็นสำคัญ (พรญา สูงเนิน เสรี ชัชแฉม และสมโภชน์ อนกสุข, 2552) ความตรงเป็นคุณสมบัติของข้อสอบที่แสดงว่า คะแนนจากข้อสอบสามารถสรุปอ้างอิงไปยังสิ่งที่วัดได้อย่างเหมาะสม (ปิยะพิพิญ ตินาว ม.ร.ว.สมพร สุหัศนีย์ และเสรี ชัชแฉม, 2550) ซึ่งแบบทดสอบในบางฉบับยังขาดความตรงหรือแสดงถึงการทำหน้าที่ต่างกันของข้อสอบ และเมื่อนำไปทดสอบกับกลุ่มย่อยตั้งแต่ 2 กลุ่มขึ้นไป พบว่าผลการวิเคราะห์การตอบของผู้สอบที่มีความสามารถที่เท่ากัน แต่โอกาสในการตอบได้ถูกต้องไม่เท่ากัน คือ แบบทดสอบเกิดปัจจัยที่เอื้อต่อผู้สอบบางกลุ่มเท่านั้น และนักการศึกษาพยายามที่จะตัดข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบออกเพื่อให้เกิดความยุติธรรม กับผู้สอบก็ยังเกิดปัญหาอีกเช่นเดิม และในต่างประเทศมีการศึกษาการทำหน้าที่ต่างกันของข้อสอบโดยใช้ข้อสอบ PISA 2003 และใช้วิธีการวิเคราะห์องค์ประกอบจำกัด (Restricted Factor Analysis: RFA) วิธีแม่นเทล-แyenส์เซล และทฤษฎีการตอบสนองข้อสอบในการวิเคราะห์อัตราส่วนไลค์ลิขิต พบว่า มีข้อสอบที่ทำหน้าที่ต่างกัน (ชัยวัฒน์ ฤทธิพันธ์, 2558)

จากการวิจัยของ Li, Hunter, and Oshima (2013) ได้ศึกษาการตรวจสอบการทำหน้าที่ต่างกันในการทดสอบ การอ่านและเหตุผล ที่เป็นไปได้ระหว่างเพศด้วยวิธีอัตราส่วนไลค์ลิขิต และวิธีแม่นเทล-แyenส์เซล พบว่า เพศมีผลต่อการทำแบบทดสอบในการอ่านและเหตุผล ที่เป็นไปได้ วิธีอัตราส่วนไลค์ลิขิต สามารถตรวจสอบการทำหน้าที่ต่างกันของข้อสอบได้ดีในแบบทดสอบที่มีความยาวตั้งแต่ 20 ข้อขึ้นไป และ Kabasakal, Arsan, Gok, and Kelecioglu (2014) ได้ศึกษาแบบจำลองการเปรียบเทียบอัตราความคลาดเคลื่อนประเทศไทยที่ 1 และอัตราความถูกต้องโดยใช้วิธีอัตราส่วนไลค์ลิขิต วิธีชิปเทสท์ และวิธีแม่นเทล-แyenส์เซล ภายใต้เงื่อนไข ขนาดกลุ่มตัวอย่าง ความยาวแบบทดสอบที่ทำหน้าที่แตกต่างกันของข้อสอบ พบว่า วิธีชิปเทสท์ มีอัตราความคลาดเคลื่อนประเทศไทยที่ 1 สูงสุดแบบทดสอบที่ทำหน้าที่ต่างกันของข้อสอบแบบเอกสารูปແຕวิธีแม่นเทล-แyenส์เซล มีอัตราความถูกต้องสูงภายใต้เงื่อนไขของขนาดกลุ่มตัวอย่างความยาวแบบทดสอบ นอกจากนี้ค่าร้อยละของการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบแบบเอกสารูปภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง ความยาวแบบทดสอบจะมีอิทธิพลต่ออัตราความคลาดเคลื่อนประเทศไทยที่ 1 ของวิธีอัตราส่วนไลค์ลิขิต และค่าร้อยละของการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ มีความแตกต่างความสามารถในการทดสอบ ระยะเวลาในการทดสอบการทำหน้าที่ต่างกันของข้อสอบความยาวแบบทดสอบสั่งผลให้ วิธีชิปเทสท์มีความคลาดเคลื่อนประเทศไทยที่ 1 ส่วนวิธีแม่นเทล-แyenส์เซล ปัจจัยที่มีผลต่อความคลาดเคลื่อนประเทศไทยที่ 1 คือ ขนาดกลุ่มตัวอย่าง ระยะเวลาในการทดสอบของการทำหน้าที่ต่างกันของข้อสอบที่มีความแตกต่างความสามารถ และไม่มีปัจจัยใดที่ส่งผลให้เกิดอัตราความถูกต้อง วิธีชิปเทสท์ และวิธีแม่นเทล-แyenส์เซล แต่จะมีผลกระทบต่อการประเมินอัตราความถูกต้องต่อวิธีอัตราส่วนไลค์ลิขิต และการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบพหุมิติ ระหว่างวิธีการวิเคราะห์องค์ประกอบจำกัด (RFA) กับวิธีทดสอบโดยโลจิสติก (LR)

ภายใต้เงื่อนไข 18 เงื่อนไข โดยกำหนดขนาดกลุ่มตัวอย่างขนาดใหญ่ จำนวน 2,000 คน กลุ่มตัวอย่างขนาดกลาง จำนวน 1,000 คน และกลุ่มตัวอย่างขนาดเล็ก จำนวน 300 คน พบว่า วิธีดัดแปลงจิตติก มีประสิทธิภาพในการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบดีกว่าเกือบทุกเงื่อนไขโดยพิจารณาจากอัตราความถูกต้องประเภทที่ 1 และอัตราความคลาดเคลื่อนประเภทที่ 1 (ปิยะพิพัฒน์ ตินวรา ม.ร.ว.สมพร สุทธานนท์ และเสรี ชัดแข็ม, 2550) การศึกษาขนาดกลุ่มตัวอย่างของ Stout, Li, Nandakumar, and Bolt (1997) ได้เสนอการวิเคราะห์โดยใช้โปรแกรม SIBTEST กลุ่มตัวอย่างขนาดเล็กที่ควรใช้ คือ ขนาดกลุ่ม จำนวน 100 คน จากงานวิจัยของ Narayanan and Swaminathan (1994) ได้ศึกษาการใช้กลุ่มตัวอย่างขนาดเล็ก จำนวน 300 คน ก็พอที่จะตรวจสอบการทำหน้าที่ต่างกันของข้อสอบได้อย่างมีประสิทธิภาพ และ (พรญา สูงเนิน เสรี ชัดแข็ม และสมโภชน์ อเนกสุข, 2552) ได้ศึกษาเพื่อพัฒนาเกณฑ์ตัดสินข้อสอบลำเอียงทางเพศ พบว่า วิธีชิปเทส์ และวิธีแมนเทล-แyenส์เซล ควรใช้กลุ่มตัวอย่างจำนวน 600 คนขึ้นไป ได้เปรียบเทียบประสิทธิภาพ วิธีแมนเทล-แyenส์เซล และวิธีชิปเทส์ พบร้า เมื่อกลุ่มตัวอย่างมีขนาดจำนวน 200 คน และจำนวน 600 คน วิธีทั้งสองสามารถตรวจสอบการทำหน้าที่ต่างกันของข้อสอบได้ถูกต้องร้อยละ 50 แต่ถ้ากลุ่มตัวอย่างขนาดจำนวน 1,000 คน สามารถตรวจสอบได้ถูกต้อง ร้อยละ 100

การทำหน้าที่ต่างกันของข้อสอบ (Differential Item Functioning: DIF) มีขนาดที่ศึกษา การการทำหน้าที่ต่างกันเปลี่ยนไปตามความสามารถของผู้เข้าสอบที่แตกต่างกัน แบ่งออกเป็น 2 ประเภท คือ 1) ข้อสอบทำหน้าที่ต่างกันแบบเอกรูป (Uniform DIF) โอกาสของผู้เข้าสอบที่คะแนนสูงกว่าจะไม่ปฏิสัมพันธ์กับตัวอย่าง 2) ข้อสอบทำหน้าที่ต่างกันแบบอนุรูป (Nonuniform DIF) โอกาสของผู้เข้าสอบที่คะแนนสูงกว่าจะปฏิสัมพันธ์กับตัวอย่าง ความสามารถของผู้เข้าสอบที่มีคะแนนต่ำกว่าตามทฤษฎีการตอบสนองข้อสอบ (IRT) โดยสามารถพิจารณาปฏิสัมพันธ์จากความแตกต่างของพารามิเตอร์อำนาจจำแนกของข้อสอบระหว่างผู้เข้าสอบ กลุ่มย่อย 2 กลุ่ม และถ้าหากข้อสอบระหว่างผู้เข้าสอบ 2 กลุ่มย่อยมีค่าพารามิเตอร์อำนาจจำแนกของข้อสอบระหว่างผู้เข้าสอบกลุ่มย่อย 2 กลุ่ม ซึ่งมีค่าพารามิเตอร์อำนาจจำแนกเท่ากันแล้วนั้น ໂດຍคุณลักษณะของข้อสอบ (Item Characteristic Curves: ICC) ของผู้เข้าสอบ 2 กลุ่มจะขนาดกันนั้นก็แสดงว่าข้อสอบที่ทำหน้าที่ต่างกันแบบสม่ำเสมอถ้าข้อสอบระหว่างกลุ่มผู้เข้าสอบ 2 กลุ่มย่อยมีค่าพารามิเตอร์อำนาจจำแนกไม่เท่ากันแล้ว ໂດຍคุณลักษณะข้อสอบ (Item Characteristic Curves: ICC) ของผู้เข้าสอบ 2 กลุ่มจะไม่ขนาดกัน นั่นก็แสดงว่าข้อสอบทำหน้าที่ต่างกันแบบไม่สม่ำเสมอ การการทำหน้าที่ต่างกันของข้อสอบมาวิเคราะห์เชิงตรรกะ (Logical Analysis) นอกจากนี้การทำหน้าที่ต่างกันของข้อสอบเป็นการเปรียบเทียบผลการตอบข้อสอบระหว่างกลุ่มผู้สอบอย่างน้อย 2 กลุ่มขึ้นไป ประกอบด้วยกลุ่มแรกที่เรียกว่ากลุ่มเปรียบเทียบ (Focal Group: F) เป็นกลุ่มที่สนใจศึกษาและคาดว่าจะเป็นกลุ่มที่เสียเปรียบในการตอบข้อสอบ และกลุ่มที่สองเรียกว่า กลุ่มอ้างอิง (Reference Group: R) เป็นกลุ่มที่คาดว่าจะได้เปรียบในการตอบข้อสอบได้ถูกต้อง ซึ่งในการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบ (ศรีชัย กาญจนวนวاسي, 2555, หน้า 115-118)

การเปรียบเทียบผลการตอบข้อสอบในระหว่างกลุ่มอ้างอิงกับกลุ่มเปรียบเทียบต้องจับคู่ (Matching) ผู้สอบตามความสามารถ โดยเป็นเงื่อนไขของการตรวจสอบการทำหน้าที่ต่างกัน

ของข้อสอบเกณฑ์การจับคู่ (Matching Criteria) ที่นิยมใช้กันมี 2 วิธี ดังนี้ (ศิริชัย กาญจนวาสี, 2555, หน้า 120-122) เกณฑ์ภายนอก (External Criterion) การวิเคราะห์การทำหน้าที่ต่างกัน จะใช้เกณฑ์ภายนอกนี้สามารถใช้ได้ทั้งข้อสอบรายข้อและข้อสอบทั้งฉบับ โดยการใช้คะแนนจาก ข้อสอบอื่นเป็นเกณฑ์ภายนอกแล้วใช้เทคนิคการวิเคราะห์การทดสอบ (Regression Analysis) เพื่อเป็นการเปรียบเทียบเส้นกราฟความสัมพันธ์ระหว่างตัวแปรเกณฑ์ กับตัวแปรทำนายระหว่าง กลุ่มอ้างอิงและกลุ่มเปรียบเทียบและเกณฑ์ภายใน (Internal Criterion) การวิเคราะห์การทำหน้าที่ ต่างกัน โดยใช้เกณฑ์ภายในเป็นการนำวิธีการทางสถิติมาตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ หรือแบบสอบจะเป็นการเน้นการพิจารณาจากโครงสร้างภายในของแบบสอบเป็นหลัก และ การวิเคราะห์ผลจากการตอบข้อสอบและความสามารถหรือคะแนนจริงของผู้เข้าสอบได้จาก แบบทดสอบฉบับนั้น เพื่อนำมาเปรียบเทียบระหว่างผู้เข้าสอบจากกลุ่มอ้างอิงและกลุ่มเปรียบเทียบ ที่มีความสามารถหรือคะแนนจริงเท่ากัน ว่าจะมีผลการตอบหรือโอกาสในการตอบข้อสอบได้ถูกต้อง แตกต่างกันหรือไม่เพื่อบ่งชี้การทำนายของข้อสอบ

ผู้วิจัยจึงมีความสนใจที่จะศึกษาการเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกัน ของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ระหว่างวิธีอัตราส่วนไลคลิชูด วิธีซิปเทสท์ และวิธีแมนเทล-ແ xen สเซล ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก จำนวน 300 คน ขนาดกลาง จำนวน 1,000 คน และขนาดใหญ่ จำนวน 2,000 คน และข้อค้นพบ ที่ได้จะเป็นแนวทางในการดำเนินการจัดทำแบบทดสอบระดับชาติในชั้นประถมศึกษาปีที่ 3 ในครั้ง ต่อไปว่าการที่จะนำแบบทดสอบที่จะนำมาทดสอบจะต้องมีความเป็นกลางไม่เกิดการทำหน้าที่ต่างกัน ของข้อสอบต่อไป

วัตถุประสงค์ของการวิจัย

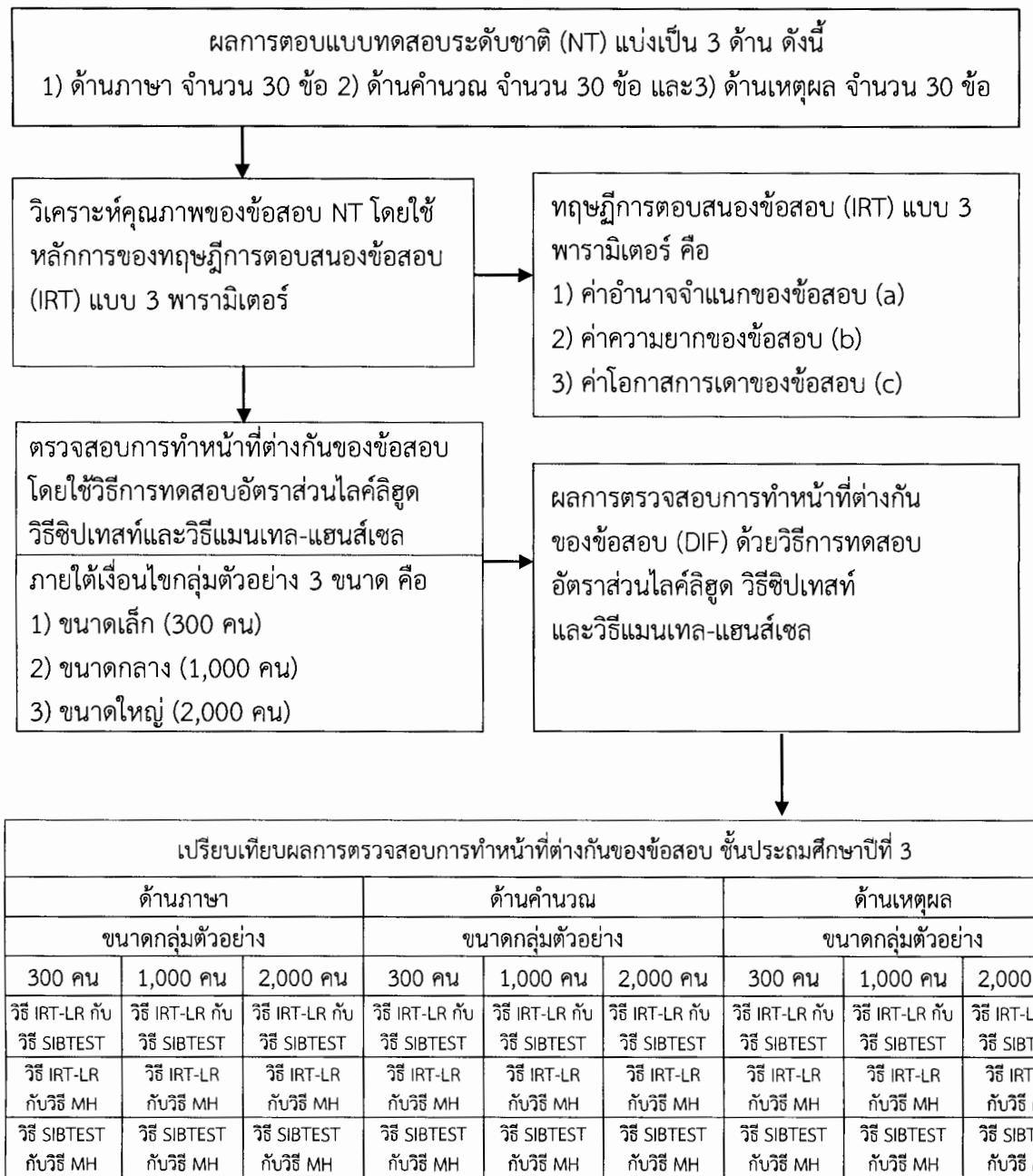
- เพื่อวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 โดยใช้หลักการของทฤษฎี การตอบสนองของข้อสอบ แบบ 3 พารามิเตอร์ ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน)

- เพื่อตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้วยวิธีการทดสอบอัตราส่วนไลคลิชูด วิธีซิปเทสท์ และวิธีแมนเทล-ແ xen สเซล ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน)

- เพื่อเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบ ระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้วยวิธีการทดสอบอัตราส่วนไลคลิชูด วิธีซิปเทสท์ และวิธีแมนเทล-ແ xen สเซล ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน)

กรอบแนวคิดการวิจัย

จากการศึกษาแนวคิด ทฤษฎี และงานวิจัยที่เกี่ยวข้อง ผู้วิจัยได้วิเคราะห์และสังเคราะห์ การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบด้วยวิธีการทดสอบอัตราส่วนไลคลิชูด วิธีซิปเทส์ และวิธีแมนเทล-แยนส์เซล เป็นวิเคราะห์ข้อสอบที่อยู่บนพื้นฐานทฤษฎีตอบสนองข้อสอบทั้ง 3 วิธี (Li, Hunter, & Oshima, 2013) และได้เปรียบเทียบประสิทธิภาพในการตรวจสอบการทำหน้าที่ ต่างกันของข้อสอบระหว่างวิธีการทดสอบอัตราส่วนไลคลิชูด วิธีซิปเทส์ และวิธีแมนเทล-แยนส์เซล โดย Kabasakal, Arsan, Gok, and Kelecioglu (2014) ได้ศึกษาประสิทธิภาพในการตรวจสอบ DIF ด้วยวิธีแมนเทล-แยนส์เซล วิธีซิปเทส์ และวิธีการทดสอบอัตราส่วนไลคลิชูด พบว่า วิธีการทดสอบ อัตราส่วนไลคลิชูด มีประสิทธิภาพในการตรวจสอบ DIF ได้ดี และงานวิจัยของ Li, Hunter, and Oshima (2013) ได้ศึกษาการตรวจสอบการทำหน้าที่ต่างกันในแบบทดสอบด้านการอ่าน และ ด้านเหตุผล ด้วยวิธีการทดสอบอัตราส่วนไลคลิชูด และวิธีแมนเทล-แยนส์เซล พบว่า วิธีการทดสอบ อัตราส่วนไลคลิชูด มีประสิทธิภาพในการตรวจสอบ DIF ได้ดีในแบบทดสอบที่มีความยาวของข้อสอบ ตั้งแต่ 20 ข้อขึ้นไป มากำหนดเป็นกรอบแนวคิดการวิจัย ดังภาพที่ 1-1



ภาพที่ 1-1 กรอบแนวคิดการวิจัย

สมมติฐานของการวิจัย

1. เมื่อกลุ่มตัวอย่างขนาดเล็ก (300 คน) ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ขั้นประเพณีศึกษาปีที่ 3 ด้านภาษา โดยใช้วิธีการทดสอบอัตราส่วนไลร์คลิยูด วิธีซิปเหลท์ และวิธีแมนเทล-แซนส์เซล แตกต่างกัน

2. เมื่อกลุ่มตัวอย่างขนาดเล็ก (300 คน) ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ขั้นประเพณีศึกษาปีที่ 3 ด้านคำวณ โดยใช้วิธีการทดสอบอัตราส่วนไลค์ลิชญด์ วิธีซิปเทส์ และวิธีแมนเกล-แยนส์เซล แตกต่างกัน

3. เมื่อกลุ่มตัวอย่างขนาดเล็ก (300 คน) ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ขั้นประถมศึกษาปีที่ 3 ด้านเหตุผล โดยใช้วิธีการทดสอบอัตราส่วนไลค์ลิขิต วิธีซีปเทส์ และวิธีแม่นเทล-แชนส์เซล แตกต่างกัน

4. เมื่อกลุ่มตัวอย่างขนาดกลาง (1,000 คน) ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านภาษา โดยใช้วิธีการทดสอบอัตราส่วนไลค์ลิขิต วิธีซิปเพลท์ และวิธีแมนเทล-แฮนส์เซล แตกต่างกัน

5. เมื่อกลุ่มตัวอย่างขนาดกลาง (1,000 คน) ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านความสนใจใช้วิธีการทดสอบอัตราส่วนไลค์ลิชญด วิธีซีปเทสท์ และวิธีแมนเกล-แยนส์เซล แตกต่างกัน

6. เมื่อกลุ่มตัวอย่างขนาดกลาง (1,000 คน) ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านเหตุผล โดยใช้วิธีการทดสอบอัตราส่วนไลค์ลิขิต วิธีซิปเพสท์ และวิธีแม่นเทล-แซนส์เซล แตกต่างกัน

7. เมื่อกลุ่มตัวอย่างขนาดใหญ่ (2,000 คน) ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านภาษา โดยใช้วิธีการทดสอบอัตราส่วนไลค์ลิขิต วิธีซิปเพลท และวิธีแมเนเกล-แคนส์เซล แตกต่างกัน

8. เมื่อกลุ่มตัวอย่างขนาดใหญ่ (2,000 คน) ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านคำนวณ โดยใช้วิธีการทดสอบอัตราส่วนไลค์ลิลย์ด วิรชิปเทสท์ และวิธีแมนเกล-แยนส์เซล แตกต่างกัน

9. เมื่อกลุ่มตัวอย่างขนาดใหญ่ (2,000 คน) ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านเหตุผล โดยใช้วิธีการทดสอบอัตราส่วนไลค์ลิขิต วิธีซีปเพลท์ และวิธีแมนเทล-แซนส์เซล แตกต่างกัน

ประโยชน์ที่คาดว่าจะได้รับจากการวิจัย

1. ได้วิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบแบบหล่ายั่วเลือก
วิธีใดเหมาะสมกับการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบเมื่อกลุ่มตัวอย่างที่ต่างกัน

2. ได้รับการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ (DIF) ทั้ง 3 ด้าน คือ ด้านภาษา ด้านคำนวณ และด้านเหตุผล ตรวจพบ DIF ได้ดีในกลุ่มตัวอย่างขนาดใหญ่ (2,000 คน)

3. ได้วิธีตรวจสอบการทำหน้าที่ต่างกันของข้อสอบเมื่อกลุ่มตัวอย่างขนาดต่างกัน 3 ขนาด คือ ขนาดใหญ่ (2,000 คน) ขนาดกลาง (1,000 คน) และขนาดเล็ก (300 คน) ว่าวิธีแม่นเทล-แชนส์เซล ตรวจพบ DIF ได้มากที่สุด รองลงมาคือ วิธีซิบเทส์ และวิธีการทดสอบอัตราส่วนไลร์คลิชูด

ขอบเขตของการวิจัย

การวิจัยครั้งนี้ใช้ข้อมูลทุติยภูมิ ที่เป็นผลการตอบแบบทดสอบระดับชาติ (NT) ขั้นประถมศึกษาปีที่ 3 ทั้ง 3 ด้าน คือ ด้านภาษา ด้านคำนวณ และด้านเหตุผล ของนักเรียนสังกัด สำนักงานคณะกรรมการการศึกษาขั้นพื้นฐานโดยสำนักทดสอบทางการศึกษา จำนวนทั้งหมด 706,372 คน

ตัวแปรที่ศึกษา

1. ตัวแปรต้น มี 2 ตัว ได้แก่

1.1 วิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ จำนวน 3 วิธี ได้แก่

1.1.1 วิธีการทดสอบอัตราส่วนไลร์คลิชูด (วิธี IRT-LR)

1.1.2 วิธีซิบเทส์ (วิธี SIBTEST)

1.1.3 วิธีแม่นเทล-แชนส์เซล (วิธี MH)

1.2 ขนาดกลุ่มตัวอย่าง จำนวน 3 ขนาด ได้แก่

1.2.1 ขนาดเล็ก (300 คน)

1.2.2 ขนาดกลาง (1,000 คน)

1.2.3 ขนาดใหญ่ (2,000 คน)

2. ตัวแปรตาม คือ ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ (DIF) โดยพิจารณา จากความถี่ (จำนวนข้อ) และร้อยละที่ตรวจพบการทำหน้าที่ต่างกันของข้อสอบ

นิยามคัพท์เฉพาะ

การทดสอบระดับชาติ (National Testing: NT) หมายถึง การประเมินคุณภาพผู้เรียน ระดับชาติของผู้เรียน ขั้นประถมศึกษาปีที่ 3 ตามมาตรฐานและตัวชี้วัดตามหลักสูตรแกนกลาง การศึกษาขั้นพื้นฐาน พุทธศักราช 2551 ซึ่งดำเนินการจัดสอบโดยสำนักทดสอบทางการศึกษา สำนักงานคณะกรรมการการศึกษาขั้นพื้นฐาน ความสามารถด้านภาษา (Literacy Ability) ความสามารถด้านคำนวณ (Numeracy Ability) และความสามารถด้านเหตุผล (Reason Ability)

แบบทดสอบระดับชาติ (National Test: NT) หมายถึง แบบทดสอบประเมินคุณภาพ ผู้เรียนระดับชาติของผู้เรียน ขั้นประถมศึกษาปีที่ 3 ตามมาตรฐานและตัวชี้วัดตามหลักสูตรแกนกลาง การศึกษาขั้นพื้นฐาน พุทธศักราช 2551 ซึ่งดำเนินการจัดสอบโดยสำนักทดสอบทางการศึกษา สำนักงานคณะกรรมการการศึกษาขั้นพื้นฐาน ความสามารถด้านภาษา (Literacy Ability) ความสามารถด้านคำนวณ (Numeracy Ability) และความสามารถด้านเหตุผล (Reason Ability)

การทำหน้าที่ต่างกันของข้อสอบ (Differential Item Functioning: DIF) หมายถึง ข้อสอบที่ผู้ตอบข้อสอบ ซึ่งมีความสามารถหรือคุณลักษณะที่ต้องการวัดเท่ากัน มีโอกาสตอบข้อสอบ

ข้อนี้ได้ถูกต้องไม่เท่ากัน เนื่องจากอยู่ในกลุ่มผู้เข้าสอบย่อยที่มีลักษณะต่างกันในที่นี่ คือ กลุ่มผู้ตอบข้อสอบนักเรียนชาย กับกลุ่มผู้ตอบนักเรียนหญิง

การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกัน หมายถึง การนำผลการตรวจพบข้อสอบที่ก่อให้เกิดการทำหน้าที่ต่างกันของข้อสอบ ระหว่างวิธีการตรวจสอบของวิธี โดยทำการเปรียบเทียบ จำนวนข้อ (ร้อยละ) ที่พบรากอนหน้าที่ต่างกันของข้อสอบ

ทฤษฎีการตอบสนองข้อสอบ (Item Response Theory: IRT) หมายถึง ทฤษฎีการวัดที่อธิบายความสัมพันธ์ระหว่างความสามารถที่อยู่ภายในตัวของบุคคล กับผลการตอบข้อคำถามหรือแบบทดสอบโดยใช้คิ้งคุณลักษณะข้อสอบ แบบ 3 พารามิเตอร์ คือ ค่าอำนาจจำแนกของข้อสอบ (a) ค่าความยากของข้อสอบ (b) และค่าโอกาสในการเดาของข้อสอบ (c)

วิธีการทดสอบอัตราส่วนไลค์ลิhood (Item Response Theory-Likelihood Ratio: IRT-LR) หมายถึง การใช้หลักของอัลกอริทึมในการประมาณค่าความเป็นไปได้ของค่าพารามิเตอร์ ตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบเอกสาร และออนไลน์ โดยใช้โปรแกรมสำเร็จรูป IRTPRO ตามทฤษฎีการตอบสนองข้อสอบ IRT

วิธีซิบเทส์ (SIBTEST) หมายถึง วิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบที่อยู่ภายในทฤษฎีการตอบสนองข้อสอบ (Item Response Theory: IRT) ที่พัฒนาโดย Shealy and Stout (1993) ซึ่งในการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบจะคำนวณดัชนีจากสัดส่วนของการตอบข้อสอบถูกของผู้เข้าสอบกลุ่มอ้างอิงกับกลุ่มเปรียบเทียบในข้อสอบ NT ซึ่งจะใช้คะแนนที่ได้มาจับคู่เปรียบเทียบระหว่างกลุ่มผู้เข้าสอบที่ข้อสอบมีความตรง

วิธีเมนเทล-เอนส์เซล (Mantel-Haenszel: MH) หมายถึง วิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบโดยทำการทดสอบอัตราส่วนเปรียบเทียบด้วยไค-สแควร์ (Chi-Square Statistic) จะนำผลการตอบข้อสอบของกลุ่มผู้เข้าสอบระหว่าง 2 กลุ่ม กลุ่มหนึ่งเรียกว่า “กลุ่มเปรียบเทียบ” (Focal Group) และอีกกลุ่มหนึ่งเรียกว่า “กลุ่มอ้างอิง” (Reference Group) เอามาพิจารณาเปรียบเทียบสัดส่วนการตอบข้อสอบว่า 2 กลุ่ม นี้ที่มีความสามารถในระดับเดียวกัน

ขนาดกลุ่มตัวอย่าง (Sample Size) หมายถึง จำนวนของผู้เข้าสอบตอบข้อคำถาม ทั้งกลุ่มอ้างอิงและกลุ่มเปรียบเทียบที่ซึ่งผู้วิจัยกำหนดขึ้นเพื่อการวิจัยครั้งนี้ โดยกำหนดไว้ 3 ขนาด คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน)

กลุ่มเปรียบเทียบ (Focal Group: F) หมายถึง กลุ่มผู้เข้าสอบที่คาดว่าจะเสียเปรียบจากการตอบข้อสอบทำหน้าที่ต่างกันของข้อสอบ เป็นกลุ่มที่แนวโน้มที่ตอบข้อสอบได้ถูกต้องมากกว่าผู้เข้าสอบอีกกลุ่มหนึ่งที่มีความสามารถในระดับเดียวกัน

กลุ่มอ้างอิง (Reference Group: R) หมายถึง กลุ่มผู้เข้าสอบที่คาดว่าจะได้เปรียบจากการตอบข้อสอบทำหน้าที่ต่างกันของข้อสอบ เป็นกลุ่มที่แนวโน้มที่ตอบข้อสอบได้ถูกต้องสูงกว่าผู้เข้าสอบอีกกลุ่มหนึ่งที่มีความสามารถในระดับเดียวกัน

ค่าอำนาจจำแนกของข้อสอบ (a-parameter) หมายถึง ค่าแสดงถึงข้อสอบที่มีความสามารถในการจำแนกคุณลักษณะของข้อสอบ ค่าอำนาจจำแนกข้อสอบมีค่าเข้าใกล้ 1 แสดงถึงข้อสอบมีอำนาจจำแนกได้ดี ค่าอำนาจจำแนกเป็นลบ แสดงว่าข้อสอบใช้ไม่ได้ ค่าอำนาจจำแนกเป็นศูนย์แสดงว่าข้อสอบจำแนกไม่ได้ และค่าอำนาจจำแนกของข้อสอบอยู่ระหว่าง

0.50 ถึง 2.50

ค่าความยากของข้อสอบ (b-parameter) หมายถึง ค่าความยากข้อสอบที่เข้าใกล้ 2 แสดงว่าข้อสอบมีความยากค่อนข้างดี ถ้าข้อสอบมีค่าความยากของข้อสอบที่ติดลบ แสดงว่าข้อสอบค่อนข้างง่าย และค่าความยากของข้อสอบอยู่ระหว่าง -2.50 ถึง 2.50

ค่าโอกาสการเดาของข้อสอบ (c-parameter) หมายถึง ค่าแสดงการเดาของข้อสอบ มีโอกาสการเดาข้อสอบได้ถูกต้อง เป็น 0 และค่าการเดาข้อสอบมากกว่า 0.3 แสดงว่าค่อนข้างใช้ไม่ได้ ถ้าค่าโอกาสการเดาข้อสอบเป็น 0 แสดงว่าข้อสอบใช้ได้ ค่าโอกาสการเดาข้อสอบอยู่ระหว่าง 0 ถึง 1

บทที่ 2

เอกสารและงานวิจัยที่เกี่ยวข้อง

ผู้วิจัยได้ศึกษาเอกสารและงานวิจัยที่เกี่ยวข้องกับการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ด้วยวิธีการทดสอบอัตราส่วนไลค์ลิชญด์ วิธีซิปเทส์ และวิธีแมนเทล-แยนส์เซล โดยนำเสนอเป็น 6 ตอน ดังนี้

ตอนที่ 1 การทดสอบความสามารถพื้นฐานของผู้เรียนระดับชาติ (National Testing: NT) และงานวิจัยที่เกี่ยวข้อง

ตอนที่ 2 ทฤษฎีการตอบสนองข้อสอบ (Item Response Theory: IRT) และงานวิจัยที่เกี่ยวข้อง

ตอนที่ 3 การทำหน้าที่ต่างกันของข้อสอบ (Differential Item Functioning: DIF) และงานวิจัยที่เกี่ยวข้อง

ตอนที่ 4 การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบด้วยวิธีการทดสอบอัตราส่วนไลค์ลิชญด์ และงานวิจัยที่เกี่ยวข้อง

ตอนที่ 5 การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบด้วยวิธีซิปเทส์ และงานวิจัยที่เกี่ยวข้อง

ตอนที่ 6 การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบด้วยวิธีแมนเทล-แยนส์เซล และงานวิจัยที่เกี่ยวข้อง

ตอนที่ 1 การทดสอบความสามารถพื้นฐานของผู้เรียนระดับชาติ (National Testing: NT) และงานวิจัยที่เกี่ยวข้อง

การทดสอบความสามารถพื้นฐานของผู้เรียนระดับชาติ (National Testing: NT) คือ การสอบประเมินคุณภาพการศึกษาระดับชาติด้านพื้นฐาน เพื่อการประกันคุณภาพผู้เรียน ตรวจสอบ กำกับดูแล และพัฒนาคุณภาพการศึกษาของโรงเรียน ของสำนักงานคณะกรรมการการศึกษา ขั้นพื้นฐาน (สพฐ.) จัดสอบโดย สำนักงานคณะกรรมการการศึกษาขั้นพื้นฐาน ให้กับนักเรียน ในช่วงชั้นที่ 1 หรือ ชั้นประถมศึกษาปีที่ 3 โดยพระราชบัญญัติการศึกษาแห่งชาติ พ.ศ. 2542 มาตรา 47 ได้กำหนดให้มีระบบการประกันคุณภาพการศึกษา เพื่อให้มีการพัฒนาคุณภาพและ มาตรฐานการศึกษาในทุกระดับ และมาตรา 48 โดยให้น่วยงานต้นสังกัดและสถานศึกษา จัดให้มีระบบการประกันคุณภาพการศึกษาภายในสถานศึกษา และการประกันคุณภาพภายใน ยังเป็น ส่วนหนึ่งของกระบวนการบริหารการศึกษาที่ต้องดำเนินการอย่างต่อเนื่อง มีการจัดทำรายงาน ประจำปีเสนอต่อหน่วยงานต้นสังกัด หน่วยงานที่เกี่ยวข้องเปิดเผยต่อสาธารณะ และนำไปสู่ การพัฒนาคุณภาพมาตรฐานการศึกษา เพื่อรับการประกันคุณภาพภายนอก การประเมินคุณภาพ การศึกษาขั้นพื้นฐานจึงเป็นกระบวนการ เพื่อให้ข้อมูลที่จะเป็นตัวปัจฉีกสำคัญในการจัดการศึกษา ซึ่งเป็นส่วนประกอบสำคัญส่วนหนึ่งในการประกันคุณภาพภายใน แนวทางการวัดและประเมินผล การเรียนรู้ เพื่อให้ได้ข้อมูลสารสนเทศที่แสดงพัฒนาการความก้าวหน้า และความสำเร็จทางการเรียน

ของผู้เรียน โดยสถานศึกษาต้องจัดให้มีการประเมินผลการเรียนให้เป็นมาตรฐานเดียวกัน ทั้งในระดับชั้นเรียน ระดับสถานศึกษา ระดับเขตพื้นที่การศึกษา และระดับชาติ ข้อมูลที่ได้จากการประเมินจะนำไปใช้ในการพัฒนาคุณภาพของผู้เรียน และคุณภาพการจัดการศึกษาของสถานศึกษา และเพื่อเป็นสารสนเทศองรับปริบทางการประเมินภายนอก

ดังนั้นการประเมินคุณภาพทางการศึกษา เป็นกลไกสำคัญประการหนึ่ง ที่สะท้อนให้เห็นถึงความสำเร็จของการพัฒนาคุณภาพการศึกษาตามเจตนาของพระราชบัญญัติการศึกษาแห่งชาติรวมทั้งเป็นข้อมูลพื้นฐานประกอบการวางแผนและพัฒนาคุณภาพของหน่วยงานต่าง ๆ ที่เกี่ยวข้องกับการจัดการศึกษา สำนักงานคณะกรรมการการศึกษาขั้นพื้นฐานเป็นหน่วยงานที่รับผิดชอบในการจัดการศึกษาตั้งแต่ระดับปฐมวัยจนถึงระดับมัธยมศึกษา จำเป็นต้องใช้ประโยชน์จากสารสนเทศผลการประเมินคุณภาพการศึกษาของสถานศึกษาในสังกัด ไปใช้ในการวางแผนเพื่อพัฒนาและปรับปรุงคุณภาพการศึกษาอย่างเป็นระบบและต่อเนื่อง เพื่อให้เกิดการพัฒนาคุณภาพผู้เรียน อีกทั้งการพัฒนาหลักสูตร กระบวนการจัดการศึกษาที่ส่งเสริมให้ผู้เรียนสามารถพัฒนาตามธรรมชาติ และเต็มตามศักยภาพของตนของต่อไป

ในปีการศึกษา 2555 สำนักงานคณะกรรมการการศึกษาขั้นพื้นฐาน โดยสำนักทดสอบทางการศึกษาได้กำหนดนโยบายจัดการทดสอบความสามารถพื้นฐานของผู้เรียนระดับชาติ (National Testing: NT) ในระดับชั้นประถมศึกษาปีที่ 3 ทั้ง 3 ด้าน คือ ด้านภาษา (Literacy Ability) ด้านคำนวณ (Numeracy Ability) และด้านเหตุผล (Reason Ability) ให้สอดคล้องกับจุดเน้นการพัฒนาคุณภาพผู้เรียน คือชั้นประถมศึกษาปีที่ 1-3 นักเรียนมีทักษะความสามารถในการอ่านออกเขียนได้ คิดเลขเป็น มีทักษะการคิดขั้นพื้นฐาน (สำนักทดสอบทางการศึกษา สำนักงานคณะกรรมการการศึกษาขั้นพื้นฐาน กระทรวงศึกษาธิการ, 2555, หน้า 1) โดยมีกรอบการประเมินผลสัมฤทธิ์ทางการเรียนตามหลักสูตร ปีการศึกษา 2555 ซึ่งมีตัวชี้วัดด้านภาษา 6 ตัวชี้วัด ตัวชี้วัดด้านคำนวณ 5 ตัวชี้วัด และตัวชี้วัดด้านเหตุผล 4 ตัวชี้วัด รวม 15 ตัวชี้วัด และใช้เครื่องมือที่ใช้ในการประเมินผลสัมฤทธิ์ทางการเรียนของผู้เรียน คือ แบบทดสอบ โดยลักษณะของเครื่องมือเป็นปรนัยแบบเลือกตอบ มีจำนวนแบบทดสอบที่ใช้ในการสอบ แบ่งเป็น 3 ด้าน คือ ด้านภาษา มีแบบทดสอบ 30 ข้อ เวลา 50 นาที ด้านคำนวณ มีแบบทดสอบ 30 ข้อ เวลา 50 นาที และด้านเหตุผล มีแบบทดสอบ 30 ข้อ เวลา 50 นาที ผลการประเมินผลสัมฤทธิ์ทางการเรียนของผู้เรียนระดับชั้นประถมศึกษาปีที่ 3 ปีการศึกษา 2555 พบว่า ด้านเหตุผลมีคะแนนเฉลี่ยสูงสุด เท่ากับ 45.92 ด้านภาษา มีคะแนนเฉลี่ยเท่ากับ 42.94 และด้านคำนวณมีคะแนนเฉลี่ย เท่ากับ 37.45 ซึ่งการประเมินคุณภาพผู้เรียนในระดับชาติ ระดับชั้นประถมศึกษาปีที่ 3 จำแนกออกเป็น 4 ภูมิภาค ได้แก่ ภาคเหนือ ภาคกลาง ภาคตะวันออกเฉียงเหนือ และภาคใต้ พบร่วมกันว่า ภาคตะวันออกเฉียงเหนือ เป็นภูมิภาคที่มีคะแนนค่าเฉลี่ยสูงกว่าคะแนนเฉลี่ยระดับประเทศ โดยมีคะแนนสูงสุดเท่ากับ 43.21 ภาคเหนือมีคะแนนค่าเฉลี่ยเท่ากับ 41.57 ภาคกลางมีคะแนนค่าเฉลี่ย เท่ากับ 41.17 และภาคใต้มีคะแนนค่าเฉลี่ยเท่ากับ 40.89 โดยคะแนนเฉลี่ยระดับประเทศ เท่ากับ 42.10 (สำนักทดสอบทางการศึกษา สำนักงานคณะกรรมการการศึกษาขั้นพื้นฐาน กระทรวงศึกษาธิการ, 2557, หน้า 4-6)

ผลการประเมินคุณภาพการศึกษาขั้นพื้นฐาน เพื่อการประกันคุณภาพผู้เรียน (Nation Test: NT) ชั้นประถมศึกษาปีที่ 3 พบร่วมกันว่า ด้านภาษา ด้านคำนวณ และด้านเหตุผล คะแนนเฉลี่ยต่อ

กว่าร้อยละ 50 ทั้ง 3 ด้าน เนื่องจากปีการศึกษา 2555 สำนักงานคณะกรรมการการศึกษาขั้นพื้นฐานได้มีการปรับกรอบการประเมินไปเป็น 3 ด้าน คือ ด้านภาษา (Literacy Ability) ด้านคำนวณ (Numeracy Ability) และด้านเหตุผล (Reason Ability) ซึ่งเดิมจะเป็นการประเมิน 3 กลุ่มสาระ การเรียนรู้ คือ กลุ่มสาระการเรียนรู้ภาษาไทย กลุ่มสาระการเรียนรู้คณิตศาสตร์ และกลุ่มสาระ การเรียนรู้วิทยาศาสตร์ ดังนั้นจึงทำให้นักเรียนครูผู้สอนมีการนำข้อสอบแบบเดิมมาให้นักเรียนทดลองทำ เพราะยังไม่มีตัวอย่างข้อสอบแบบปรับใหม่อ่อนมาให้ทดลองทำ ส่งผลให้นักเรียนไม่คุ้นเคยกับการประเมินผลสัมฤทธิ์ในครั้งนี้ การประเมินผลสัมฤทธิ์ทางการเรียนของผู้เรียนจำแนกตามกลุ่มสถานศึกษาที่มีบริบทแตกต่างกัน

1. กลุ่มสถานศึกษาที่เป็นต้นแบบส่งเสริมสนับสนุนในการพัฒนาคุณภาพในภาพรวม ได้แก่ โรงเรียนดีศรีตำบล โรงเรียนในฝัน รุ่น 1-3 โรงเรียนมาตรฐานสากล โรงเรียนจุฬาภรณ์ โรงเรียนอัตรา การแข่งขันสูง โรงเรียนเฉลิมพระเกียรติ และโรงเรียนไทยรัฐวิทยา สรุปได้ดังนี้

1.1 โรงเรียนดีศรีตำบล ผลการประเมินผลสัมฤทธิ์ทางการเรียนชั้นประถมศึกษาปีที่ 3 พบว่าคะแนนเฉลี่ยรวมต่ำกว่าระดับประเทศ เมื่อพิจารณาความสามารถ พบว่าทุกความสามารถมีคะแนนเฉลี่ยต่ำกว่าคะแนนเฉลี่ยระดับประเทศทุกด้าน เนื่องจากกระบวนการคัดเลือกโรงเรียนดีศรีตำบลในบางเขตพื้นที่ ยังมีข้อบกพร่อง เช่น ไม่ได้คัดเลือกจากโรงเรียนที่มีความพร้อมในแต่ละตำบล จริงๆ ระยะเวลาในจัดสรรงบประมาณโดยมุ่งเน้นทางด้านครุภัณฑ์ในการก่อสร้างอาคารสถานที่ จัดซื้อคอมพิวเตอร์ และเน้นด้านคุณธรรมจริยธรรมมากกว่า แต่ไม่เน้นไปที่การขาดแคลนครุภัณฑ์ สาขาวิชา และไม่เน้นทางวิชาการจึงทำให้คะแนนเฉลี่ยต่ำกว่าคะแนนเฉลี่ยระดับประเทศทุกกลุ่มสาระ การเรียนรู้

1.2 โรงเรียนในฝัน รุ่นที่ 1-3 ผลการประเมินผลสัมฤทธิ์ทางการเรียนระดับชั้นประถมศึกษาปีที่ 3 ของกลุ่มโรงเรียนในฝัน พบว่า โรงเรียนในฝันรุ่นที่ 2 คะแนนเฉลี่ยรวมสูงกว่า คะแนนระดับประเทศ และโรงเรียนในรุ่นที่ 1,3 คะแนนเฉลี่ยรวมต่ำกว่าคะแนนระดับประเทศ เนื่องจากโรงเรียนในรุ่นที่ 3 ได้ผ่านการประเมินรับรองเป็นต้นแบบโรงเรียนในฝันแล้ว จึงทำให้ไม่ต่อเนื่องในการส่งเสริมการจัดการศึกษาร่วมไปถึงการนิเทศที่รับผิดชอบโรงเรียนในฝันยังทำหน้าที่รับผิดชอบโรงเรียนดีศรีตำบลอีกด้วยจึงทำให้การบริหารจัดการไม่ดีเท่าที่ควรทำให้โรงเรียนในฝันรุ่นที่ 3 มีผลคะแนนเฉลี่ยรวมต่ำกว่าระดับประเทศ

1.3 โรงเรียนมาตรฐานสากล ผลการประเมินผลสัมฤทธิ์ทางการเรียนระดับชั้นประถมศึกษาปีที่ 3 มีคะแนนเฉลี่ยรวมสูงกว่าระดับประเทศ เท่ากับ 47.18 พบว่า มีคะแนนเฉลี่ยสูงกว่าคะแนนเฉลี่ยระดับประเทศทุกด้าน เนื่องจากการคัดเลือกโรงเรียนเพื่อพัฒนาและยกระดับ เป็นโรงเรียนมาตรฐานสากล โรงเรียนจะมีการคัดเลือกนักเรียนที่สอดคล้องความพร้อมเข้าสู่โรงเรียน อีกทั้งผู้บริหาร ครูผู้สอนได้รับการพัฒนาเพื่อเพิ่มสมรรถนะ และมีภาคีเครือข่ายรวมไปถึงโรงเรียนร่วมพัฒนาทั้งในประเทศและต่างประเทศ จึงส่งผลให้มีคะแนนเฉลี่ยสูงกว่าระดับประเทศทุกกลุ่มสาระการเรียนรู้

1.4 โรงเรียนอัตราการแข่งขันสูง ผลการประเมินผลสัมฤทธิ์ทางการเรียนระดับชั้นประถมศึกษาปีที่ 3 พบว่า ความสามารถทั้ง 3 ด้าน มีคะแนนเฉลี่ยรวมสูงกว่าระดับประเทศ เนื่องจากโรงเรียนมีกระบวนการคัดเลือกนักเรียนที่เข้าศึกษามีการแข่งขันสูงและเป็นนักเรียน

ที่มีผลการเรียนที่ดีถึงดีมาก และโรงเรียนก็เน้นการพัฒนาทางด้านวิชาการ จึงส่งให้คะแนนเฉลี่ยสูงกว่าระดับประเทศ

1.5 โรงเรียนเฉลี่ยมพระเกียรติ ผลการประเมินผลสัมฤทธิ์ทางการเรียนในระดับชั้นประถมศึกษาปีที่ 3 มีคะแนนเฉลี่ยรวมสูงกว่าระดับประเทศ เท่ากับ 53.13 พบว่า ด้านภาษาและด้านเหตุผลมีค่าเฉลี่ย เท่ากับ 46.98 และ 46.07 สูงกว่าคะแนนเฉลี่ยระดับประเทศ และด้านคำนวณ มีคะแนนเฉลี่ย เท่ากับ 37.07 ใกล้เคียงกับคะแนนเฉลี่ยระดับประเทศ

1.6 โรงเรียนไทยรัฐวิทยา ผลการประเมินผลสัมฤทธิ์ทางการเรียนในระดับชั้นประถมศึกษาปีที่ 3 มีคะแนนค่าเฉลี่ยรวมต่ำกว่าคะแนนเฉลี่ยระดับประเทศ เท่ากับ 40.96 พบว่า มีคะแนนเฉลี่ยต่ำกว่าคะแนนระดับประเทศในทุกด้าน เนื่องจากผลสัมฤทธิ์ในภาพรวมในรายกลุ่ม สาระการเรียนรู้ต่ำกว่าระดับประเทศในทุกรายการ ซึ่งโรงเรียนไทยรัฐวิทยามีจุดเด่นด้านคุณลักษณะของผู้เรียน ความซื่อสัตย์ เพื่อสร้างพลเมืองโลกที่ดี อาจทำให้โครงการยกระดับผลสัมฤทธิ์ของโรงเรียนไม่ได้เท่าที่ควร และบางส่วนเป็นโรงเรียนขนาดเล็ก ส่งผลให้ขาดความพร้อมหลายด้าน ครูผู้สอนไม่ตรงตามเอกวิชาที่สอน

2. กลุ่มสถานศึกษาต้นแบบส่งเสริมสนับสนุนในการพัฒนาคุณภาพการศึกษาเฉพาะด้านระดับชั้นประถมศึกษาปีที่ 3 คือ โรงเรียนส่งเสริมนิสัยรักการอ่าน (โรงเรียนต้นแบบ, โรงเรียนแกนนำ และโรงเรียนเครือข่าย) โรงเรียนยกระดับผลสัมฤทธิ์ทางการเรียนในวิชาวิทยาศาสตร์และคณิตศาสตร์ และโครงการยกระดับคุณภาพ โรงเรียนขนาดเล็กผลการทดสอบ ดังนี้

2.1 โรงเรียนส่งเสริมนิสัยรักการอ่าน (โรงเรียนต้นแบบ) ระดับชั้นประถมศึกษาปีที่ 3 พบว่า มีคะแนนเฉลี่ยสูงกว่าคะแนนเฉลี่ยระดับประเทศ ซึ่งมีค่าเฉลี่ยใกล้เคียงกับค่าเฉลี่ยระดับประเทศ เท่ากับ 42.10 และพบว่ามี 2 ด้าน ที่มีคะแนนค่าเฉลี่ยสูงกว่าระดับประเทศ คือ ด้านภาษา เท่ากับ 42.10 และด้านเหตุผล เท่ากับ 42.10

2.2 โรงเรียนส่งเสริมนิสัยรักการอ่าน (โรงเรียนแกนนำ) ระดับชั้นประถมศึกษาปีที่ 3 พบว่า มีคะแนนเฉลี่ยสูงกว่าคะแนนเฉลี่ยระดับประเทศ เท่ากับ 42.57 ซึ่งมีค่าใกล้เคียงกับค่าเฉลี่ยระดับประเทศ เท่ากับ 42.10 และมี 2 ด้านที่มีคะแนนเฉลี่ยสูงกว่าระดับประเทศ คือ ด้านภาษา เท่ากับ 43.61 และด้านคำนวณ เท่ากับ 37.62 โรงเรียนที่มีวัตถุประสงค์และการบริหารจัดการเฉพาะด้านซึ่งเน้นด้านการอ่านส่งผลให้ผลสัมฤทธิ์ทางการเรียนสูงกว่าระดับชาติ

2.3 โรงเรียนส่งเสริมนิสัยรักการอ่าน (โรงเรียนเครือข่าย) ระดับชั้นประถมศึกษาปีที่ 3 พบว่า มีค่าเฉลี่ยต่ำกว่าคะแนนเฉลี่ยระดับประเทศ เท่ากับ 42.03 และมี 1 ด้านที่มีค่าเฉลี่ยสูงกว่าระดับประเทศ คือ ด้านเหตุผล เท่ากับ 46.15 ซึ่งอาจเป็นเพราะโรงเรียนที่เริ่มเข้าร่วมโครงการ การบริหารจัดการ การสนับสนุนด้านสื่อ วัสดุ อุปกรณ์ บุคลากร และงบประมาณ และการสนับสนุนอย่างเต็มที่

2.4 โรงเรียนห้องสมุดมีชีวิต (โรงเรียนต้นแบบ) พบว่ามีคะแนนเฉลี่ยต่ำกว่าคะแนนเฉลี่ยระดับประเทศ เท่ากับ 41.87 เนื่องจากการนิเทศกำกับ ติดตามอย่างเป็นระบบ จะมี 1 ด้านเท่านั้นที่มีค่าเฉลี่ยสูงกว่าระดับประเทศ คือ ความสามารถด้านภาษา เท่ากับ 43.07 ซึ่งโรงเรียนอาจมีวัตถุประสงค์เฉพาะด้าน และเน้นกิจกรรม ส่งเสริมการอ่าน การใช้ภาษา อาจส่งผลให้มีคะแนนค่าเฉลี่ยเฉพาะด้านสูงกว่าระดับประเทศ

2.5 โรงเรียนห้องสมุดมีชีวิต (โรงเรียนแก่นนำ) พบว่า มีคะแนนเฉลี่ยสูงกว่าคะแนนเฉลี่ยระดับประเทศ เท่ากับ 42.91 ในทุกด้าน อาจเป็นเพราะเป็นโรงเรียนที่มีวัตถุประสงค์เฉพาะด้านเน้นกิจกรรมต่างๆ ส่งเสริมการอ่าน การใช้ภาษา และระยะเวลาการประเมินความก้าวหน้าอย่างต่อเนื่อง เพื่อเลื่อนระดับเป็นโรงเรียนต้นแบบ จึงทำให้ส่งผลให้ค่าเฉลี่ยทุกด้านสูงกว่าระดับประเทศ

2.6 โรงเรียนห้องสมุดมีชีวิต (โรงเรียนเครือข่าย) พบว่า มีคะแนนเฉลี่ยต่ำกว่าคะแนนเฉลี่ยระดับประเทศ เท่ากับ 41.94 จะมีเพียง 1 ด้านที่มีค่าเฉลี่ยสูงกว่าระดับประเทศ คือ ความสามารถด้านเหตุผล เท่ากับ 46.02 อาจเป็นเพราะเป็นที่มีวัตถุประสงค์เฉพาะด้าน เน้นกิจกรรมส่งเสริมการอ่าน การใช้ภาษา และช่วงระยะเวลาของการประเมินความก้าวหน้าอย่างต่อเนื่อง เพื่อเลื่อนระดับเป็นโรงเรียนแก่นนำ ทำให้ส่งผลให้ค่าเฉลี่ยด้านดังกล่าวสูงกว่าระดับประเทศ

2.7 โครงการยกระดับคุณภาพโรงเรียนขนาดเล็ก เป็นสถานศึกษาต้นแบบส่งเสริมสนับสนุนในการพัฒนาคุณภาพการศึกษาเฉพาะด้าน พบว่า มีคะแนนเฉลี่ยสูงกว่าคะแนนเฉลี่ยระดับประเทศ เท่ากับ 43.78 ค่าเฉลี่ยใกล้เคียงค่าเฉลี่ยระดับประเทศ เท่ากับ 42.10 เนื่องจาก การสนับสนุน สื่อ วัสดุ อุปกรณ์ อาคารสถานที่ งบประมาณ การสร้างเครือข่าย และโรงเรียนที่ผ่านกระบวนการคัดเลือกและฝึกอบรมเพื่อพัฒนาบุคลากร ทำให้ส่งผลกระทบโดยตรงต่อความสามารถด้านภาษา ความสามารถด้านคำนวณ และความสามารถด้านเหตุผลสูงกว่าระดับประเทศ เนื่องจากการสนับสนุนในหลายด้าน โดยเฉพาะครูผู้สอนที่จบการศึกษาตรงตามสาขา วิชาเอก ทำให้มีความรู้ความสามารถเฉพาะด้าน ทั้งสนับสนุนด้านวิชาการ วัสดุ อุปกรณ์ที่ทันสมัย งบประมาณเพียงพอ

3. โรงเรียนที่มีขนาดของสถานศึกษาแตกต่างกัน ได้แก่ โรงเรียนที่มีขนาดเล็ก (มีจำนวนนักเรียน 0-120 คน) โรงเรียนที่มีขนาดกลาง (มีจำนวนนักเรียนอยู่ระหว่าง 121-300 คน) โรงเรียนที่มีขนาดใหญ่ (มีจำนวนนักเรียนอยู่ระหว่าง 301-500 คน) และโรงเรียนที่มีขนาดใหญ่พิเศษ (มีจำนวนนักเรียนตั้งแต่ 501) ผลสัมฤทธิ์ทางการเรียนระดับประถมศึกษาปีที่ 3 ของโรงเรียนที่มีขนาดต่างกัน โรงเรียนที่มีขนาดใหญ่พิเศษมีคะแนนสูงสุดในทุกด้าน รองลงมาเป็นโรงเรียนขนาดเล็ก โรงเรียนขนาดกลาง เนื่องจากโรงเรียนขนาดใหญ่พิเศษ โดยส่วนใหญ่เป็นโรงเรียนประจำจังหวัด จึงทำให้มีความพร้อมในทุกด้าน ส่วนโรงเรียนขนาดเล็ก มีนักเรียนน้อยทำให้ครูผู้สอนดูแลได้ทั่วถึง มีการสนับสนุนเป็นรูปธรรม และโรงเรียนขนาดกลาง มักประสบปัญหาบุคลากรมีจำนวนไม่เพียงพอ นักเรียนมีจำนวนมากมาก จึงทำให้ส่งเสริมการจัดการเรียนรู้ค่อนข้างยาก (สำนักทดสอบทางการศึกษา สำนักงานคณะกรรมการการศึกษาขั้นพื้นฐาน กระทรวงศึกษาธิการ, 2557, หน้า 53-81)

สรุป คือ จากผลการประเมินคุณภาพการศึกษาขั้นพื้นฐาน เพื่อการประกันคุณภาพผู้เรียน ระดับชั้นประถมศึกษาปีที่ 3 พบร้า ความสามารถด้านภาษา ความสามารถด้านคำนวณ และความสามารถด้านเหตุผล มีคะแนนเฉลี่ยต่ำกว่าร้อยละ 50 ทั้ง 3 ด้าน เนื่องจากในปีการศึกษา 2555 สำนักงานคณะกรรมการการศึกษาขั้นพื้นฐานได้ปรับเปลี่ยนกรอบการประเมินผล จากเดิม

ที่ประเมินผลสัมฤทธิ์ทางการเรียน 3 กลุ่มสาระการเรียนรู้ คือ ภาษาไทย คณิตศาสตร์ และ วิทยาศาสตร์ เปเลี่ยนเป็นการประเมินความสามารถด้านภาษา (Literacy Ability) ความสามารถด้านคำนวณ (Numeracy Ability) และความสามารถด้านเหตุผล (Reason Ability) จึงทำให้มีตัวอย่างข้อสอบเพื่อให้ครูผู้สอนนำไปเทียบเคียง ในการสร้างเครื่องมือวัดและประเมินผลผู้เรียน

งานวิจัยที่เกี่ยวข้องกับการทดสอบความสามารถพื้นฐานของผู้เรียนระดับชาติ (NT) มิติหนึ่ง

วรรณ ศรีกัล่า อุน เจริญวงศ์ระยับ และนิคม นาคอ้าย (2559) ศึกษาปัจจัยพหุระดับที่ส่งผลต่อคะแนนการสอบประเมินคุณภาพการศึกษาระดับชาติ ด้านความสามารถทางภาษา การศึกษาของโรงเรียนที่มีผล NT ต่ำในจังหวัดพิษณุโลก มีวัตถุประสงค์เพื่อศึกษาปัจจัยพหุระดับได้แก่ ปัจจัยพหุระดับนักเรียนและระดับห้องเรียนที่ส่งผลกระทบ NT ด้านความสามารถทางภาษาโรงเรียนที่มีคะแนน NT ต่ำในจังหวัดพิษณุโลก กลุ่มตัวอย่าง ได้แก่ นักเรียนจำนวน 1,260 คน และครูจำนวน 68 ห้องเรียน ข้อมูลที่เก็บรวมข้อมูลโดยใช้แบบทดสอบ 2 ฉบับ (แบบสอบตามระดับนักเรียนและระดับห้องเรียน) แบบสอบตามระดับนักเรียน ประกอบด้วย 5 ตัวแปร แบบสอบตามวัดความรู้พื้นฐานเดิม แบบสอบตามวัดแรงจูงใจเพื่อสัมฤทธิ์ในการทำข้อสอบ NT แบบสอบตามวัดเจตคติต่อการเรียนภาษาไทยแบบสอบตามวัดสภาพทางบ้าน และแบบสอบตามวัดความอาใจใส่ผู้ปกครองในการส่งเสริมการเรียน ซึ่งมีค่าความเชื่อมั่นเท่ากับ 0.92 แบบสอบตามระดับห้องเรียน ประกอบด้วย 2 ตัวแปร แบบสอบวัดคุณภาพการสอนครุภาษาไทยและแบบสอบตามวัดบรรยายกาศในชั้นเรียน มีค่าความเชื่อมั่นเท่ากับ 0.87 ทำการวิเคราะห์ข้อมูล โดยการวิเคราะห์พหุระดับ (Multilevel Analysis) และผลการวิเคราะห์ที่สำคัญ พบว่า โรงเรียน ที่มีคะแนน NT ต่ำ 1) ในระดับนักเรียน ตัวแปรความรู้พื้นฐานเดิมส่งผลต่อคะแนนการสอบ NT ด้านความสามารถทางภาษาอย่างมีนัยสำคัญทางสถิติ 2) ในระดับห้องเรียนไม่มีตัวแปรอิสระใดที่ส่งผลต่อคะแนน NT ด้านความสามารถทางภาษา ตัวแปรความรู้พื้นฐานเดิมสามารถอธิบายความแปรปรวนของคะแนนการสอบ NT ด้านความสามารถทางภาษา ได้ร้อยละ 12.85

เอกลักษณ์ คล้ายสุบรรณ สั่ງรณ์ งัดกระโภก และนลินี ณ นคร (2559) ศึกษามโนเดล การวัดมูลค่าเพิ่มทางการศึกษาสำหรับวัดคุณภาพสถานศึกษาด้วยการใช้ผลรวมของผลสัมฤทธิ์ทางการเรียนและการประเมินและรับรองคุณภาพของโรงเรียน โดยมีวัตถุประสงค์เพื่อ 1) พัฒนาวิธีการวัดมูลค่าเพิ่มทางการศึกษาเพื่อใช้ประเมินคุณภาพสถานศึกษาด้วยการวัดมูลค่าเพิ่มจากผลสัมฤทธิ์ทางการเรียนและการประเมินคุณภาพของโรงเรียน 2) ศึกษา ความสอดคล้องของผลการประเมินคุณภาพสถานศึกษาด้วยการวัดมูลค่าเพิ่มที่พัฒนาขึ้น ซึ่งมี การกำหนดน้ำหนักของการรวมคะแนนแตกต่างกัน และ 3) เปรียบเทียบมูลค่าเพิ่มทางการศึกษาของสถานศึกษาที่มีบริบทต่างกัน การวิจัยนี้ใช้ข้อมูลทุติยภูมิ กลุ่มตัวอย่างที่ใช้ในการวิจัย คือ โรงเรียนประถมศึกษาในสังกัดสำนักงานคณะกรรมการการศึกษาขั้นพื้นฐาน จำนวน 96 โรงเรียน และมีนักเรียนขั้นประถมศึกษาปีที่ 3 ปีการศึกษา 2557 จำนวน 7,988 คน ซึ่งได้จากการสุ่มแบบหลายขั้นตอน ตัวแปรที่ศึกษาประกอบด้วย ตัวแปรระดับนักเรียน ได้แก่ เพศ การศึกษา ของผู้ปกครองความสัมพันธ์ในครอบครัว ผลการประเมินคุณภาพการศึกษาระดับชาติขั้นพื้นฐาน

ตัวแปรระดับสถานศึกษา ได้แก่ ผลการประเมินและรับรองคุณภาพของโรงเรียน ที่ตั้ง และขนาดของสถานศึกษา เครื่องมือที่ใช้ คือ แบบบันทึกข้อมูลพื้นฐานของโรงเรียน ผลการประเมินคุณภาพการศึกษาระดับชาติขั้นพื้นฐานและ ผลการประเมินคุณภาพสถานศึกษา รอบ 3 ผู้วิจัย วิเคราะห์ข้อมูลค่านิยมคุณภาพการศึกษาที่เป็นผลรวมของคะแนนผลสัมฤทธิ์ทางการเรียนกับผลการประเมินและรับรองคุณภาพของโรงเรียนโดยมีการถ่วงน้ำหนักต่างกันสามโมเดล คือ 40:60, 50:50, และ 60:40 สถิติที่ใช้ในการวิเคราะห์ข้อมูล คือ การวิเคราะห์พหุระดับผลการวิจัยสรุปได้ว่า (1) โมเดลการวัดมูลค่าเพิ่มทางการศึกษาที่มีความกลมกลืนกับข้อมูลมากที่สุด คือ โมเดลที่ 1 (40:60) รองลงมาคือโมเดลที่ 2 (50:50) และ โมเดลที่ 3 (60:40) (2) ความสอดคล้องของผลการประเมินคุณภาพสถานศึกษาด้วยการวัดมูลค่าเพิ่มระหว่างโมเดลที่ 1 (40:60) กับโมเดลที่ 2 (50:50) มีความสอดคล้องกันมากที่สุด มีความสอดคล้อง 91.67% และรองลงมา ระหว่างโมเดลที่ 2 (50:50) กับโมเดลที่ 3 (60:40) มีความสอดคล้อง 88.54% และระหว่างโมเดลที่ 1 (40:60) กับโมเดลที่ 3 (60:40) มีความสอดคล้อง 80.21% ตามลำดับ และ(3) สถานศึกษาที่มีที่ตั้ง และขนาดต่างกัน มีคะแนนมูลค่าเพิ่มทางการศึกษาไม่แตกต่างกัน

สุราทิพย์ ตรีสิน และปิยะทิพย์ ประดุจพร (2560) ได้ศึกษาการเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ด้านภาษา ด้านคำนวน และด้านเหตุผล ขั้นประถมศึกษาปีที่ 3 ด้วยวิธี HGLM วิธี MIMIC และวิธี IRT-LR โดยมีวัตถุประสงค์เพื่อวิเคราะห์คุณภาพของแบบทดสอบระดับชาติ (NT) และตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติขั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน และด้านเหตุผล ด้วยวิธี HGLM วิธี MIMIC และวิธี IRT-LR 1) วิเคราะห์คุณภาพของแบบทดสอบระดับชาติ ทั้ง 3 ด้าน 2) ตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ด้วยวิธี HGLM วิธี MIMIC และวิธี IRT-LR และ 3) เปรียบเทียบผลการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบ ด้วยวิธีการตรวจสอบ 3 วิธี ข้อมูลวิเคราะห์เป็นข้อมูลทุติยภูมิ จากผลการตอบแบบทดสอบระดับชาติของนักเรียนขั้นประถมศึกษาปีที่ 3 ปีการศึกษา 2556 จำนวน 9,600 คน ผลการวิจัยพบว่า แบบทดสอบระดับชาติขั้นประถมศึกษาปีที่ 3 มีค่าความยากของข้อสอบ (b) อยู่ในระดับค่อนข้างยากมีค่าอำนาจจำแนกของข้อสอบ (a) อยู่ในระดับที่สามารถจำแนกผู้สอบได้ดี และมีค่าโอกาสในการเดาของข้อสอบ (c) ไม่เกิน 0.3 และการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบทั้ง 3 ด้าน จะเห็นได้ว่าเพียงส่งผลให้เกิดการทำหน้าที่ต่างกันของข้อสอบ โดยเพียงอยู่จะได้เปรียบในการตอบข้อสอบ ด้านภาษา และด้านเหตุผล

Kjellstrom and Pettersson (2005) ได้ศึกษาความรู้ความเข้าใจของหลักสูตรการศึกษานำไปทดสอบระดับชาติในวิชาคณิตศาสตร์ในประเทศไทย จะเกี่ยวกับระบบการทดสอบของประเทศไทย โดยมีการทดสอบระดับชาติ ที่เป็นการมุ่งเป้าหมายองค์ความรู้ ระบบการทดสอบที่สำคัญในหลักสูตรของการเรียนการสอนรายวิชาคณิตศาสตร์ ที่ได้อธิบายอิทธิพลต่อการทดสอบระดับชาติ ซึ่งได้เปลี่ยนแปลงการประเมินที่แตกต่างไปจากเดิม เพื่อไปสู่เป้าหมาย เกณฑ์ในการอ้างอิง การให้คะแนนและการแก้ปัญหาของนักเรียน ของกระบวนการเรียนการสอนรายวิชาคณิตศาสตร์ มีผลต่อคะแนนระดับชาติ ดังนั้น สถานศึกษาจึงควรปรับปรุงหลักสูตรการศึกษาให้มีประสิทธิภาพยิ่งขึ้น

Brown, De Four-Babb, Bristol and Conrad (2014) ได้ศึกษาการทดสอบระดับชาติเกี่ยวกับคำพูดของครูในตรินเดดและโตเบโก ซึ่งได้มีการกล่าวถึงข้อเสนอแนะของการทดสอบในตรินเดดและโตเบโก รวมไปถึงขอบเขตที่ใช้ในการตัดสินใจของหลักสูตรที่ส่งผลกระทบต่อการเรียนรู้ของนักเรียน โดยกลุ่มตัวอย่างโดยจะประกอบไปด้วยครูประถมศึกษา จำนวน 133 คน แบ่งเป็น 79 คน จากโรงเรียนที่มีประสิทธิภาพต่ำ ส่วน 54 คน จากโรงเรียนประสิทธิภาพสูง และผู้บริหาร 10 คน ซึ่งผลการวิจัยเชิงปริมาณและข้อมูลเชิงคุณภาพ พบว่า มีครูจำนวนมากที่รู้สึกไม่สบายใจกับการตีความข้อมูลที่นำเสนอด้วยงานที่เกี่ยวข้องกับการทดสอบระดับชาติ ซึ่งครูในโรงเรียนที่มีประสิทธิภาพสูง ที่ผ่านการทำางานร่วมกันเพื่อเป็นข้อมูลในการตัดสินใจเรื่องการเรียนการสอนในหลักสูตร ดังนั้นจึงมีความจำเป็นที่จะให้มีการฝึกอบรมครูในการใช้ข้อมูลการประเมินของปัญหา อื่นๆ ที่จะเกิดขึ้นจากข้อมูลและหัวข้อที่จะเป็นไปได้สำหรับการวิจัยเพิ่มเติม เพื่อการแสดงถึงสัญลักษณ์ของโรงเรียนที่มีประสิทธิภาพต่ำ และโรงเรียนที่ไม่มีประสิทธิภาพในการทดสอบระดับชาติ

สรุปจากการศึกษางานวิจัยที่เกี่ยวข้องการทดสอบระดับชาตินั้นมีความสำคัญกับการศึกษา เป็นอย่างมาก เนื่องจากการทดสอบระดับชาติมีเป้าหมายในการพัฒนาผู้เรียนเป็นสำคัญ ซึ่งจะมี การนำผลของการทดสอบระดับชาติไปวิเคราะห์ผลเพื่อให้มีการปรับปรุงหลักสูตรการเรียนการสอน ให้สอดคล้องกับผู้เรียนรวมไปถึงการฝึกอบรมครูผู้สอนให้มีกระบวนการจัดการเรียนการสอน ที่มีประสิทธิภาพมากขึ้น เพื่อให้ส่งผลให้กับผลสัมฤทธิ์ทางเรียนดีขึ้นอย่างมีประสิทธิภาพ

ตอนที่ 2 ทฤษฎีการตอบสนองข้อสอบ (Item Response Theory: IRT) และงานวิจัยที่เกี่ยวข้อง

ทฤษฎีการตอบสนองข้อสอบ หรือ IRT สามารถจำแนกได้เป็น 2 ประเภท ได้แก่ ทฤษฎีการตอบสนองข้อสอบแบบตัวเลือกที่คัดแนน 2 ค่า (Binary or Dichotomous IRT) ซึ่งเป็นโมเดลการตอบสนองข้อสอบที่ใช้กับการตรวจคัดแนนรายข้อแบบ 2 ค่า เช่น ข้อสอบ หรือข้อคําถามที่ต้องให้คัดแนนแบบ 0,1 (ตอบผิดได้ 0, ตอบถูกได้ 1) แบบถูก/ผิด ใช/ไม่ใช และทฤษฎีการตอบสนองข้อสอบแบบตัวเลือกที่คัดแนนมากกว่า 2 ค่า (Polytomous IRT) ซึ่งเป็นโมเดลการตอบสนองข้อสอบที่ใช้กับการตรวจคัดแนนรายข้อแบบมากกว่า 2 ค่า เช่น ข้อสอบ หรือข้อคํามาตรประมาณค่า (Rating Scale) การตรวจข้อสอบแบบให้คัดแนนความรู้บางส่วน (Partial Credit)

แนวคิดพื้นฐานของ IRT

โมเดลการวัด (Measurement Model) โมเดลการวัดในที่นี้เป็นโมเดลเชิงคณิตศาสตร์ (Mathematical Model) ซึ่งเป็นระบบความสัมพันธ์ระหว่างตัวแปรอิสระที่รวมกันอย่างเหมาะสม ในการทำนายตัวแปรตาม

1) โมเดลการวัดแบบดั้งเดิม (CTT Model)

$$X_p = T_p + E_p \quad (1)$$

เมื่อ X_p เป็นคะแนนที่สังเกตได้ของผู้สอบ (คะแนนดิบรวม) เกิดจากการรวมเชิงเส้นตรงของผลบวกของตัวแปรแต่ละตัว 2 ตัว ได้แก่ คะแนนจริงของผู้สอบ (T_p) คะแนนความคลาดเคลื่อน (E_p) โดยมีข้อตกลงเบื้องต้นของคะแนนความคลาดเคลื่อนว่า (1) ค่าเฉลี่ยของ E ของกลุ่มผู้สอบมีค่าเป็นศูนย์ (2) ค่า E ไม่มีความสัมพันธ์กับ T และ E อื่นๆ จากโมเดลดังกล่าวมีข้อจำกัดอย่างด้านที่สำคัญได้แก่ (1) ค่า X เป็นค่าเฉพาะที่ได้จากการวัดแต่ละครั้ง ซึ่งได้จากข้อสอบที่มีค่าความยากและอำนาจจำแนกเฉพาะของฉบับที่ใช้ดังนั้น การเปรียบเทียบคะแนน (X หรือ T) ระหว่างแบบสอบที่ใช้วัดคุณลักษณะเดียวกัน จึงต้องอยู่บนพื้นฐานข้อตกลงเบื้องต้นของแบบสอบที่มีข้อสอบคู่ขนานกัน (Parallelism) (2) คะแนนรายข้อไม่ได้ เชื่อมโยงกับพฤติกรรมการตอบข้อสอบและคะแนนจริงของผู้สอบ ปัจจัยสำคัญเกี่ยวกับคุณลักษณะข้อสอบจึงถูกละเอียไปจากโมเดลการวัดค่าพารามิเตอร์ของข้อสอบ เช่น ค่าความยาก ค่าอำนาจจำแนก เป็นต้น จึงแปรผันไปตามกลุ่มผู้สอบ

2) โมเดลการตอบสนองข้อสอบ (IRT Models)

โมเดลการวัดเป็นระบบความสัมพันธ์ระหว่างตัวแปรอิสระที่รวมกันสำหรับทำนายตัวแปรตามสำหรับโมเดลการตอบสนองของข้อสอบ ตัวแปรอิสระประกอบด้วย ตัวแปรແ Pang คือความสามารถที่แท้จริงของผู้สอบ (θ) และคุณลักษณะของข้อสอบ (B) หรือค่าพารามิเตอร์ของข้อสอบ (a, b, c) ส่วนตัวแปรอิสระเป็นตัวแปรที่สังเกตได้ คือ โอกาสการตอบข้อสอบได้ถูกต้อง

ทฤษฎีการตอบสนองข้อสอบเป็นทฤษฎีการวัดที่อธิบายความสัมพันธ์ระหว่างคุณลักษณะภายในหรือความสามารถที่มีอยู่ภายในตัวบุคคลกับพฤติกรรมการตอบสนองข้อสอบของบุคคลนั้นโดยจะใช้โค้งลักษณะข้อสอบ (Item Characteristic Curve: ICC) มีการกำหนดความยาก (b) อำนาจจำแนก (a) และโอกาสการเดาข้อสอบถูก (c) IRT จึงอยู่บนพื้นฐานความคิดที่สำคัญ 2 ประการ คือ 1) ผลการตอบข้อสอบหรือข้อคำถามของผู้ตอบ สามารถอธิบายได้ด้วยความสามารถที่มีอยู่ภายในของผู้ตอบ และ 2) ความสัมพันธ์ระหว่างผลการตอบข้อสอบกับความสามารถที่มีอยู่ภายในสามารถอธิบายได้ด้วยฟังก์ชันลักษณะข้อสอบ หรือโค้งลักษณะข้อสอบ (ICC) อันมีลักษณะเป็นฟังก์ชันทางคณิตศาสตร์ เรียกว่าฟังก์ชันโลจิส (Logistic Function) หรือ ไอล์คิยองก์บันฟังก์ชันปกติสะสม (Normal Ogive Function)

ทฤษฎีการตอบสนองข้อสอบ พยายามอธิบายความสัมพันธ์ระหว่างคุณลักษณะระหว่างคุณลักษณะภายในหรือความสามารถที่มีอยู่ภายในตัวบุคคล กับพฤติกรรมการตอบสนองข้อสอบของบุคคลนั้นว่ามีโอกาสตอบข้อสอบถูกมากน้อยเพียงไร ทฤษฎีนี้มีพื้นฐานความเชื่อว่า พฤติกรรมการตอบสนองต่อข้อสอบของผู้สอบ ซึ่งเป็นสิ่งที่สังเกตได้โดยตรงว่าถูกหรือผิดจะถูกกำหนดโดยคุณลักษณะภายในหรือความสามารถที่อยู่ภายในตัวบุคคล ซึ่งเป็นสิ่งที่ไม่สามารถสังเกตได้โดยตรง ทฤษฎีนี้ได้โดยอธิบายความสามารถสัมพันธ์ดังกล่าวในรูปของฟังก์ชันคณิตศาสตร์ หรือโมเดลที่แสดงความสัมพันธ์ระหว่างระดับความสามารถ คุณลักษณะ ของข้อสอบและโอกาสของการตอบข้อสอบได้ถูก ที่เรียกว่า ฟังก์ชันการตอบสนองข้อสอบ ซึ่งมีลักษณะความสัมพันธ์เป็นแบบฟังก์ชันโลจิสหรือฟังก์ชันปกติสะสม

ฟังก์ชันการตอบสนองข้อสอบสามารถนำมาใช้ศึกษาความสัมพันธ์ระหว่างความนำ้จะเป็นในการตอบข้อสอบแต่ละข้อได้ถูกต้อง [$P_i(\theta)$] กับระดับความสามารถของผู้สอบที่วัดได้โดยแบบสอบฉบับนั้น (θ) เมื่อนำมาเขียนเป็นกราฟได้โค้งลักษณะข้อสอบ (Item Characteristic Curve: ICC)

โค้งลักษณะข้อสอบมีได้หลายลักษณะขึ้นอยู่กับโมเดล (Model) หรือแบบจำลองที่ใช้อธิบายความสัมพันธ์ตั้งกล่าว โมเดลที่นิยมใช้กันคือ โมเดลแบบหนึ่งพารามิเตอร์ (One-Parameter Model) โมเดลแบบสองพารามิเตอร์ (Two- Parameter Model) และโมเดลแบบสามพารามิเตอร์ (Three-Parameter Model)

โมเดลการตอบสนองข้อสอบ (Item Response Model)

คุณลักษณะของโมเดลการตอบสนองข้อสอบเป็นระบบความสัมพันธ์ระหว่างโอกาสตอบข้อสอบถูก (P_i) กับความสามารถที่มีอยู่ภายในผู้ตอบ (θ) ในรูปของโค้งลักษณะข้อสอบ (ICC) ซึ่งมีลักษณะเป็นฟังก์ชันโลจิส (Logistic Function) หรือฟังก์ชันปกติสะสม (Normal Ogive Function) บางครั้ง อาจเรียกว่า โมเดลโลจิส หรือ โมเดลปกติสะสม

โมเดลปกติสะสม ใช้ฟังก์ชันปกติสะสมแสดงความสัมพันธ์ระหว่างผลการตอบข้อสอบกับความสามารถของผู้สอบ ส่วนโมเดลโลจิสฟังก์ชันโลจิสแสดงความสัมพันธ์ระหว่างผลการตอบกับความสามารถดังกล่าว ซึ่งฟังก์ชันทั้งสองให้ผลลัพธ์ของการประมาณค่าใกล้เคียงกันมาก แต่ฟังก์ชันโลจิสมีลักษณะของสูตรทางคณิตศาสตร์ และวิธีคำนวณง่ายและสะดวกกว่า โมเดลโลจิสยังมีความสามารถทนทานต่อความคาดเคลื่อนที่เกิดขึ้นกับผู้สอบที่มีความสามารถสูงจะตอบข้อสอบได้ดีกว่าจึงทำให้โมเดลโลจิสเป็นที่นิยมกันมากในการนำไปใช้จริง

พารามิเตอร์ในทฤษฎีการตอบสนองข้อสอบ

พารามิเตอร์ในทฤษฎีการตอบสนองข้อสอบ แบ่งออกเป็น 2 ชนิด คือ พารามิเตอร์ข้อสอบ (Item Parameter) ได้แก่ ค่าความยากของข้อสอบ (b) ค่าอำนาจจำแนกของข้อสอบ (a) ค่าโอกาสการเดาของข้อสอบ (c) และความรอบคอบ (γ) ส่วนพารามิเตอร์ของผู้สอบ (Person Parameter) ได้แก่ ระดับความสามารถ หรือ คุณลักษณะของผู้สอบ (θ) ซึ่งค่าพารามิเตอร์ต่าง ๆ มีลักษณะและการแปลความหมาย ดังนี้ (ศิริชัยกาญจนวารี, 2550, หน้า 53-55)

1. พารามิเตอร์ข้อสอบ (Item Parameter)

พารามิเตอร์ข้อสอบประกอบด้วย ค่าความยากของข้อสอบ (b) ค่าอำนาจจำแนกของข้อสอบ (a) ค่าโอกาสการเดาของข้อสอบ (c) และความรอบคอบ ซึ่งพารามิเตอร์แต่ละชนิดมีรายละเอียด ดังนี้

1.1 ค่าความยากของข้อสอบ (b)

ค่าความยากของข้อสอบได้มาจากการค่าความสามารถที่ตรงจุดเปลี่ยนโค้ง (Inflexion Point) ซึ่งเป็นจุดที่โค้งมีความชันมากที่สุด หรือมีความหมายอีกนัยหนึ่งคือ ผู้สอบที่มีความสามารถถึงระดับ ณ จุดเปลี่ยนโค้งข้อสอบนั้น จะมีโอกาสตอบข้อสอบถูกอยู่ 0.5 หรือในทางปฏิบัติจากจุดบนแกน y ที่แสดงถึงตำแหน่งโอกาสในการตอบข้อสอบขึ้นถูกมีอยู่ 0.5 ถ้าหากเส้นตั้งขานากับแกน x 伸展กับเส้นโค้ง ซึ่งจะเป็นจุดเปลี่ยนโค้งด้วยนั้น ในทางตรงกันข้าม เมื่อหากเส้นตั้งขานากับจุดดังกล่าวให้มารดแกน x ค่าที่วัดได้ในแกน x คือ ค่าความยากของข้อสอบข้อนั้น ๆ ซึ่งข้อสอบทั้งหมดในแบบทดสอบฉบับหนึ่ง ๆ ที่นำมาวิเคราะห์จะมีค่าความยากของข้อสอบกระจายอยู่ในแกน x จากค่า $-\infty$ ถึง ∞ แต่ในทางปฏิบัติ นิยมใช้ช่วง -3 ถึง $+3$ และแบบทดสอบทั่วไป มักจะมีค่า b อยู่ระหว่าง -2.5 ถึง $+2.5$ ถ้าค่า b เข้าใกล้ -2.5 แสดงว่าข้อสอบง่าย ตรงกันข้ามถ้าค่า b อยู่ใกล้ $+2.5$ แสดงว่าข้อสอบยาก ในกรณีที่เป็นรูปแบบที่มีค่าพารามิเตอร์ 3 ตัว หรือคำนึงถึงโอกาส

ในการเดาค่าตอบ (c) ค่าความยากซึ่งเริ่มต้นจากจุดบนแกน y นั้น จะใช้จุดตั้งต้นตรงที่ค่าโอกาสในการตอบข้อสอบถูก กล่าวคือ ค่าความยากมีค่าเริ่มต้นที่ $\frac{1-c}{2}$

1.2 ค่าอำนาจจำแนกของข้อสอบ (a)

ค่าอำนาจจำแนกของข้อสอบ (a) เป็นสัดส่วนกับค่าความชัน (Slope) ของ $p_i(\theta)$ ที่จุดเปลี่ยนโถง หรือที่จุด $\theta = b$ ซึ่งทฤษฎีค่าอำนาจจำแนกของข้อสอบ จะมีค่าอยู่ในช่วง $-\infty \text{ ถึง } \infty$ แต่ในทางการนำเสนอใช้ประโยชน์นั้น ข้อสอบข้อใดที่มีค่า a ติดลบ แสดงว่าข้อสอบข้อนั้นไม่ดี และควรจะต้องถูกตัดออกไป ส่วนข้อสอบที่มีค่า a สูงขึ้น แสดงว่า ความน่าจะเป็นของการตอบข้อสอบข้อนั้นเพิ่มขึ้นเมื่อระดับความสามารถของผู้สอบสูงขึ้น ตามปกติ ค่า a มีค่าไม่เกิน +2.5 ในทางปฏิบัติ นิยมใช้ข้อสอบที่มีค่า a อยู่ระหว่าง +0.5 ถึง +2.5

1.3 ค่าโอกาสการเดาของข้อสอบ (c)

ค่าโอกาสการเดาของข้อสอบ (c) เป็นค่าที่อยู่ปลายนอกด้านตัว (Lower Asymptote) ของข้อสอบ ค่านี้เป็นค่าแทนความน่าจะเป็น หรือโอกาสที่คนซึ่งมีความสามารถตัว แต่สามารถตอบข้อสอบข้อนั้นได้ถูกต้องโดยการเดา ในทางทฤษฎีพารามิเตอร์การเดา มีค่าระหว่าง 0.00 ถึง 1.00 โดยทั่วไปนิยมใช้ข้อสอบที่มีค่าโอกาสการเดาของข้อสอบไม่เกิน 0.30

1.4 ความรอบคอบ (γ)

McDonald (1967 อ้างถึงใน Hambleton & Swaminathan, 1985) ได้เสนอ พารามิเตอร์ที่แสดงถึงความรอบคอบของผู้สอบ เป็นค่าพารามิเตอร์ที่บ่งชี้ว่าผู้สอบ ที่มีความสามารถสูงจะตอบข้อสอบได้ไม่ถูกต้องเสมอไป ซึ่งอาจเกิดความไม่รอบคอบในการพิจารณาคำตอบ หรือผู้สอบอาจมีสารสนเทศอื่น ๆ เกี่ยวกับผู้ออกข้อสอบทำให้เลือกตอบในตัวเลือกที่ไม่ใช่คำตอบที่ถูกต้อง โดย Barton and Lord (1981 อ้างถึงใน Hambleton & Swaminathan, 1985) กล่าวว่า พารามิเตอร์ตัวนี้จะเหมาะสมในการศึกษาทางทฤษฎีเท่านั้น ซึ่งในทางปฏิบัติแล้วไม่สามารถพบพารามิเตอร์นี้ได้

2. พารามิเตอร์ผู้สอบ

พารามิเตอร์ผู้สอบ (θ) เป็นระดับความสามารถของผู้สอบ (θ) ที่ประมาณได้จากโมเดลตามทฤษฎีการตอบสนองข้อสอบ นิยมปรับให้เป็นคะแนนมาตรฐานที่มีค่าเฉลี่ยเป็น 0 และส่วนเบี่ยงเบนมาตรฐานเป็น 1 ซึ่งพารามิเตอร์ผู้สอบมีค่าระหว่าง $-\infty \text{ ถึง } \infty$ แต่ส่วนใหญ่จะมีค่าอยู่ในช่วง -3.0 ถึง +3.0 ค่าที่เป็นลบ แสดงว่าผู้สอบมีความสามารถต่ำ และค่าที่เป็นบวก แสดงว่าผู้สอบมีความสามารถสูง สรุปได้ว่า พารามิเตอร์ตามทฤษฎีการตอบสนองข้อสอบจำแนกเป็น 2 ชนิด คือ พารามิเตอร์ข้อสอบ และพารามิเตอร์ผู้สอบ ซึ่งพารามิเตอร์ข้อสอบประกอบด้วยพารามิเตอร์ ความยากของข้อสอบ พารามิเตอร์อำนาจจำแนกของข้อสอบ พารามิเตอร์โอกาสการเดาของข้อสอบ และความรอบคอบ ส่วนพารามิเตอร์ผู้สอบเป็นพารามิเตอร์ที่แสดงระดับความสามารถของผู้สอบ ซึ่งข้อตกลงเบื้องต้น และพารามิเตอร์ ที่กล่าวมาเนี้ย มีความหมายเด่นชัดในกรณีที่ข้อสอบนั้นให้คะแนนแบบสองค่าในการประยุกต์ทฤษฎีเพื่อใช้กับข้อสอบที่ให้คะแนนแบบมากกว่าสองค่า ข้อตกลงเบื้องต้น ทั้งหมดก็เทียบเคียงในทำนองเดียวกันแต่ต่างกันเพียงรายละเอียดปลีกย่อยเกี่ยวกับเงื่อนไข

เฉพาะของแต่ละโมเดลเท่านั้นโน้มเดลการตอบสนองข้อสอบแบบตรวจให้คะแนนสองค่า Dichotomous IRT Models) เนื่องจากข้อมูลที่ได้จากการสอบมีหลายลักษณะ ได้แก่ ข้อมูลแบบมีสองค่า (Dichotomous) และข้อมูลแบบมีมากกว่าสองค่า (Polytomous) ดังนั้น จึงมี การพัฒนารูปแบบเพื่อให้สอดคล้องกับลักษณะของข้อมูลตั้งแต่ขั้นมากmany แต่สำหรับข้อมูลที่เป็นแบบมี 2 ค่ารูปแบบที่นิยมใช้เป็นรูปแบบโลจิสติก (Logistic Model) ซึ่งแตกต่างไปตามจำนวนพารามิเตอร์ที่ใช้ในแต่ละรูปแบบ มีรายละเอียด ดังนี้

โน้มเดลการตอบสนองข้อสอบแบบ 1 พารามิเตอร์ (One-Parameter Model)
รูปแบบนี้บางครั้งเรียกว่า รูปแบบราสช์ (Rasch Model) คือ โน้มเดลที่มีการแปรเปลี่ยนค่าพารามิเตอร์เพียงพารามิเตอร์ค่าความยากของข้อสอบ (b) เพียงอย่างเดียวโค้งคุณลักษณะข้อสอบสามารถเขียนสมการดังนี้ (Hambleton et al., 1991, pp. 12-14)

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

เมื่อ $P_i(\theta)$ แทน ความน่าจะเป็นที่ผู้ตอบซึ่งมีความสามารถจะตอบข้อสอบข้อที่ i ได้ถูกต้อง

b_i แทน ค่าพารามิเตอร์ความยากของข้อสอบข้อที่ i

e แทน ค่าคงที่มีค่าเท่ากับ 2.718

n แทน ลำดับข้อสอบข้อที่ i

ถึงแม้ว่ารูปแบบนี้จะเป็นกรณีเฉพาะของรูปแบบ 2 พารามิเตอร์ และ 3 พารามิเตอร์ แต่ก็ยังมีคุณสมบัติพิเศษที่ทำให้นิยมใช้กัน คือ ประการแรก เนื่องจากรูปแบบนี้มี จำนวนพารามิเตอร์ ไม่มากจึงสะดวกต่อการใช้งาน ประการที่สอง ปัญหาที่เกิดจากการประมาณค่าพารามิเตอร์มีน้อยกว่า การประมาณค่าพารามิเตอร์สำหรับรูปแบบที่มีพารามิเตอร์หลาย ๆ ตัว

โน้มเดลการตอบสนองข้อสอบแบบ 2 พารามิเตอร์ (Two-Parameter Model)

โน้มเดลแบบ 2 พารามิเตอร์ เป็นโค้งคุณลักษณะข้อสอบและเป็นฟังก์ชันของการแจกแจงที่มี 2 พารามิเตอร์ คือ ค่าความยากของข้อสอบ (b) ค่าอำนาจจำแนกของข้อสอบ (a) และโน้มเดล 2PL มีความเหมาะสมสมสำหรับการวัดคุณลักษณะ fenced ที่แต่ละคนมีไม่เท่ากัน สามารถเขียนเป็นสมการได้ดังนี้ (Hambleton et al., 1991, pp. 14-17)

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

เมื่อ $P_i(\theta)$ แทน เป็นค่าความน่าจะเป็นของผู้สอบที่มีความสามารถ θ

สามารถตอบข้อสอบข้อที่ i ได้ถูกต้อง

b_i แทน ค่าพารามิเตอร์ความยากของข้อสอบข้อที่ i

สำนักหอสมุด มหาวิทยาลัยบูรพา
๑.แผนสข อ.เมือง จ.ชลบุรี 2013।

23

a_i แทน ค่าพารามิเตอร์อำนาจจำแนกของข้อสอบข้อที่ i

D แทน ค่าสเกลองค์ประกอบ มีค่าเท่ากับ 1.7

$P_i(\theta)$ จากโครงความถี่สะสมกับโค้งโลจิสติก จะมีค่าที่ต่างกันน้อยกว่า 0.01

สำหรับทุกค่าของ θ จากรูปแบบนี้อยู่บนข้อตกลงที่ว่าการเดาคำตอบจะไม่เกิดขึ้น ซึ่งถ้าจะเป็นเช่นนี้ได้ก็ต่อเมื่อค่าพารามิเตอร์ $a_i > 0$ (ข้อสอบที่มีความสัมพันธ์ด้านบวกระหว่างคะแนนจากการสอบกับความสามารถของผู้สอบที่วัดโดยแบบทดสอบนั้น) และค่าความน่าจะเป็นในการตอบข้อสอบได้ถูกจะลดลงถึงศูนย์เมื่อความสามารถลดลง

โมเดลการตอบสนองข้อสอบแบบ 3 พารามิเตอร์ (Three-Parameter Model)

โมเดลแบบ 3 พารามิเตอร์ เป็นการปรับปรุงมาจาก 2 พารามิเตอร์โดยการเพิ่มพารามิเตอร์ที่ 3 คือ พารามิเตอร์ของการเดาของข้อสอบ หรือพารามิเตอร์ (c) เข้าไปในรูปแบบนี้ ดังในข้อสอบแบบหลายตัวเลือก ความน่าจะเป็นของการตอบถูกมากกว่า 0 เมื่อผู้สอบจะมีความสามารถต่ำซึ่งสามารถเขียนในรูปแบบสมการได้ดังนี้ (Hambleton et al., 1991, pp. 17-18)

$$P_i(\theta) = c_i + (1 - c_i) \frac{e^{D a_i (\theta - b_i)}}{1 + e^{D a_i (\theta - b_i)}} \quad (i = 1, 2, 3, \dots, k) \quad (4)$$

เมื่อ $P_i(\theta)$ แทน ความน่าจะเป็นของผู้สอบที่มีความสามารถ

ตอบข้อสอบข้อที่ i ได้ถูกต้อง

b_i แทน พารามิเตอร์ความยากของข้อสอบข้อที่ i

a_i แทน พารามิเตอร์อำนาจจำแนกของข้อสอบข้อที่ i

c_i แทน พารามิเตอร์ของการเดาข้อสอบได้ถูกของข้อสอบข้อที่ i

D แทน ค่าสเกลองค์ประกอบ มีค่าเท่ากับ 1.7

พารามิเตอร์ c_i เป็นจุดต่ำสุดที่โค้งคุณลักษณะข้อสอบ ซึ่งพารามิเตอร์นี้จะใช้เมื่อคิดว่า การเดาเป็นองค์ประกอบในการตอบข้อสอบ บางครั้งเรียกพารามิเตอร์นี้ว่าโอกาสที่จะตอบข้อสอบได้ถูกต้อง สำหรับคนที่มีความสามารถต่ำในการปรับรูปแบบสามพารามิเตอร์ให้เป็นรูปแบบ 2 พารามิเตอร์ ต้องอยู่บนข้อตกลงที่ว่า $c_i = 0$

1. ข้อตกลงเบื้องต้น (Assumptions) และการตรวจสอบ

ฟังก์ชันทางคณิตศาสตร์ที่ใช้ใน IRT ได้กำหนดความน่าจะเป็นในการตอบข้อสอบที่ตอบถูกนั้นจะขึ้นอยู่กับความสามารถของผู้เข้าสอบและคุณลักษณะของข้อสอบ ดังนั้นโมเดลการตอบสนองข้อสอบจึงอยู่บนพื้นฐานของความเชื่อหรือข้อตกลงเบื้องต้นหลายประการด้วยกัน ซึ่งลักษณะข้อมูลที่จะทำให้โมเดลสามารถนำไปใช้ได้อย่างเหมาะสม ข้อตกลงเบื้องต้นบางประการก็ไม่สามารถที่จะตรวจสอบไปได้โดยตรง แต่สามารถเก็บรวมข้อมูลได้โดยทางอ้อมเพื่อนำมาสนับสนุนในข้อตกลงเบื้องต้นที่สำคัญของ IRT คือดังนี้

2. ความเป็นเอกมิตร (Unidimensionality)

ข้อตกลงเบื้องต้นที่ใช้กันทั่วไปสำหรับ IRT คือ ข้อคำถาม/ ข้อสอบทุกข้อในเครื่องมือ/ แบบสอบถามมุ่งวัดเพียงคุณลักษณะเดียว หรือ ความสามารถเดียว (One Ability) ซึ่งเรียกว่า ความเป็นเอกมิตร (Unidimensionality)

การตรวจสอบความเป็นเอกมิตรของเครื่องมือ หรือแบบสอบถาม สามารถทำได้โดยใช้ เทคนิคทางสถิติ ได้แก่ การวิเคราะห์ตัวประกอบ (Factor Analysis) เพื่อคำนวณค่า ไอกน (Eigen value) สำหรับศึกษาอัตราส่วนระหว่างค่าไอกนของตัวประกอบแรกกับตัวประกอบถัดไป ถ้ามี อัตราส่วนที่สูงแสดงถึงเครื่องมือ หรือข้อสอบวัดคุณลักษณะเด่นเดียว (Single Dominant Factor) หรือทำการวิเคราะห์ให้เกิดความมั่นใจยิ่งขึ้น ด้วยการวิเคราะห์ตัวประกอบเชิงยืนยัน (Confirmatory Factor Analysis) เพื่อตรวจสอบยืนยันว่า เครื่องมือหรือแบบสอบถามมุ่งวัดเพียงคุณลักษณะเดียวหรือ ความสามารถเดียว

3. ความเป็นอิสระ (Local Independence)

แนวคิดเกี่ยวกับ "ความเป็นอิสระระหว่างข้อสอบและผู้สอบ" มีความเกี่ยวข้อง และ เชื่อมโยงมาจาก "ความเป็นเอกมิตรของแบบสอบถาม" ความเป็นอิสระระหว่างข้อสอบและผู้เข้าสอบ หมายถึง เมื่อมีการควบคุมความสามารถ (θ) ที่ส่งผลต่อการตอบข้อสอบ หรือให้ θ คงที่แล้ว ผลการตอบข้อสอบแต่ละข้อจะต้องเป็นอิสระจากกัน หรือเมื่อควบคุมอิทธิพลของ θ แล้ว ผลการตอบสนองข้อสอบรายข้อ ไม่มีความสัมพันธ์กัน นั่นคือ ไม่เดลการตอบสนองข้อสอบมีเพียง θ ปัจจัยเดียวเท่านั้นที่มีอิทธิพลต่อการตอบรายข้อ ความเป็นอิสระสามารถจำแนกพิจารณาเป็น ความอิสระระหว่างข้อสอบและความเป็นอิสระระหว่างผู้สอบดังนี้

1) ความเป็นอิสระระหว่างข้อสอบ

เมื่อสุ่มผู้เข้าสอบ ซึ่งมีความสามารถ θ ขึ้นมา 1 คน ในการข้อสอบ k ข้อ ให้ U_j เป็นผล การตอบหรือคะแนนข้อสอบข้อที่ j หลังจากควบคุม θ ของผู้สอบแล้ว คะแนนผลการตอบของผู้สอบ นั้นในแต่ละข้อไม่สัมพันธ์กัน

ถ้าผลการตอบข้อสอบรายข้อของผู้สอบคนเดียวกันเป็นอิสระจากกันความน่าจะเป็นของ แผนการตอบข้อสอบ k ข้อของผู้สอบที่มีความสามารถ θ จะเท่ากับผลคูณระหว่างความน่าจะเป็นของ ผลการตอบข้อสอบแต่ละข้อ

2) ความเป็นอิสระระหว่างผู้สอบ

เมื่อสุ่มข้อสอบขึ้นมา 1 ข้อ ในการตอบข้อสอบของผู้สอบ g คน ให้ U_i เป็นผลการตอบ หรือคะแนนข้อสอบของผู้สอบคนที่ i หลังจากควบคุม θ ของผู้สอบแต่ละคนแล้ว คะแนนผลการตอบ ข้อนั้นของผู้สอบแต่ละคนไม่สัมพันธ์กัน

ถ้าผลการตอบข้อสอบข้อเดียวกันของผู้สอบแต่ละคนเป็นอิสระจากกันความน่าจะเป็น ของแบบแผนการตอบข้อสอบของผู้สอบ g คน จะเท่ากับ ผลคูณระหว่างความน่าจะเป็นของ ผลการตอบข้อนั้นของผู้สอบแต่ละคน

การตรวจสอบความเป็นอิสระระหว่างข้อสอบและผู้สอบสามารถกระทำได้โดยการ พิจารณาเมตริกซ์ ความแปรปรวนและความแปรปรวนร่วม (Variance - Covariance Matrix) หรือเมตริกซ์สหสัมพันธ์ (Correlation Matrix) ของคะแนนคำตอบรายข้อ สำหรับกลุ่มผู้สอบ

ที่มีช่วงคะแนนความสามารถเท่ากัน โดยค่านอกแนวภาพແยงมุมครัวมีค่าต่ำหรือเข้าใกล้ 0

4. โมเดลการตอบสนองข้อสอบ (Item Response Models)

IRT อยู่บนพื้นฐานความเชื่อว่า พิنج์ชันลักษณะข้อสอบ หรือโค้งลักษณะข้อสอบ (ICC) สามารถสะท้อนความสัมพันธ์จริงระหว่างความสามารถของผู้สอบกับลักษณะของข้อสอบและผลการตอบข้อสอบ โมเดลการตอบสนองข้อสอบเสนอ ICC ซึ่งเป็นพิنج์ชันโลจิสติก ด้วยรูปลักษณ์ที่แตกต่างกัน ตามจำนวนพารามิเตอร์ที่ใช้บรรยายลักษณะของข้อสอบ โมเดลการตอบสนองข้อสอบแบบตรวจให้คะแนน 2 ค่า (Dichotomous Item Models) ที่ใช้กันแพร่หลายได้แก่ โมเดลการตอบสนองข้อสอบแบบ 1 พารามิเตอร์ โมเดลการตอบสนองข้อสอบแบบ 2 พารามิเตอร์ โมเดลการตอบสนองข้อสอบแบบ 3 พารามิเตอร์ การเลือกใช้จึงขึ้นกับจุดมุ่งหมายของงานและธรรมชาติของข้อมูล

โมเดลการตอบสนองข้อสอบแบบ 1 พารามิเตอร์ มีข้อตกลงเบื้องต้นว่า ข้อสอบแต่ละข้อ มีพารามิเตอร์ $c=0$ และพารามิเตอร์ a เท่ากัน แต่มีความแตกต่างกันเฉพาะพารามิเตอร์ b เท่านั้น โมเดลนี้จึงเหมาะสมสำหรับใช้กับข้อสอบบิ๊กเกนท์ที่ไม่สลับซับซ้อน ข้อสอบที่ค่อนข้างเรียบง่ายสำหรับพัฒนาคลังข้อสอบที่มีความเป็นเอกพันธ์

โมเดลการตอบสนองข้อสอบแบบ 2 พารามิเตอร์ มีข้อตกลงเบื้องต้นว่า ข้อสอบแต่ละข้อ มีพารามิเตอร์ $c=0$ มีความแตกต่างกันของพารามิเตอร์ a และ b โมเดลนี้จึงเหมาะสมสำหรับใช้กับข้อสอบที่ต้องเติมคำตอบ หรือข้อสอบแบบเลือกตอบที่ไม่ยากมากนักและกลุ่มผู้สอบมีความพร้อมในการตอบ

โมเดลการตอบสนองข้อสอบแบบ 3 พารามิเตอร์ มีข้อตกลงเบื้องต้นว่า ข้อสอบแต่ละข้อ มีความแตกต่างกันได้ทั้ง พารามิเตอร์ a , b , และ c โมเดลนี้จึงเหมาะสมสำหรับใช้กับข้อสอบแบบเลือกตอบแบบทั่วไปข้อสอบแบบหลายตัวเลือก เนื่องจากผู้สอบสามารถเดาคำตอบได้

การตรวจสอบความเหมาะสมของโมเดลการตอบสนองข้อสอบกับข้อมูล (Model Data Fit) ไม่ว่าจะเลือกใช้โมเดลใดก็ตาม โมเดลจะต้องสอดคล้องกับข้อมูล เพื่อให้ผลการวิเคราะห์มีความถูกต้องและน่าเชื่อถือ การตรวจสอบความสอดคล้องควรพิจารณาทั้งความไม่แปรเปลี่ยนของค่าประมาณความสามารถ และความไม่แปรเปลี่ยนของค่าประมาณพารามิเตอร์ของข้อสอบ

ความไม่แปรเปลี่ยนของค่าประมาณความสามารถ ตรวจสอบได้โดยการเปรียบเทียบค่าประมาณความสามารถของผู้สอบที่ได้จากการกลุ่มข้อสอบที่แตกต่างกัน เช่น กลุ่มข้อสอบง่าย หรือกลุ่มข้อสอบจากคลังข้อสอบเดียวกัน แต่มีความครอบคลุมเนื้หาแตกต่างกัน เป็นต้น ค่าประมาณความสามารถจะถือว่าไม่แปรเปลี่ยน เมื่อความแตกต่างเกิดขึ้นไม่เกินความคลาดเคลื่อนมาตรฐานของการประมาณค่า (ศิริชัย กาญจนวานิช, 2555, หน้า 77)

ความไม่แปรเปลี่ยนของค่าประมาณพารามิเตอร์ของข้อสอบ ตรวจสอบได้โดยเปรียบเทียบค่าประมาณพารามิเตอร์แต่ละตัวของข้อสอบที่ได้จากการกลุ่มตัวอย่างประชากรผู้สอบหลายกลุ่ม เช่น กลุ่มผู้สอบชาย/หญิง กลุ่มผู้สอบจำแนกตามภูมิภาค เป็นต้น ค่าประมาณพารามิเตอร์ของข้อสอบจะถือว่าไม่แปรเปลี่ยน เมื่อผลการพล็อตกราฟออกแบบเป็นเส้นตรง โดยมีการกระจายไม่แตกต่างจากผลที่ได้จากการกลุ่มตัวอย่าง 2 กลุ่ม ซึ่งเป็นกลุ่มสุ่มที่ทัดเทียมกัน

5. การสอบที่ไม่แข่งขันด้านเวลา (Nonspeeded Test Administration)

IRT ถือว่าความสามารถ (θ) เป็นปัจจัยสำคัญต่อผลการสอบ ความเร็วในการตอบ จะต้องไม่มีอิทธิพลต่อผลการตอบ การจัดการสอบจึงไม่อยู่ในสถานการณ์ที่สอบแข่งขันด้วยเวลา การสอบจะต้องอยู่ในลักษณะที่ผู้สอบซึ่งมีความสามารถมีเวลาเพียงพอในการทำข้อสอบ (Power Test Administration)

การตรวจสอบถึงความเหมาะสม ของมิติด้านเวลา สำหรับการดำเนินการสอบสามารถ พิจารณาได้จาก สัดส่วนหรือร้อยละของจำนวนผู้สอบที่ทำข้อสอบได้ครบทุกข้อ โดยผู้สอบส่วนใหญ่ (เช่น ร้อยละ 80 เป็นต้น) สามารถตอบข้อสอบได้ครบหรือเกือบครบทุกข้อ นอกจากนี้ควรพิจารณา เปรียบเทียบระหว่างความแปรปรวนของจำนวนข้อที่เว้น กับความแปรปรวนของจำนวนข้อที่ตอบผิด ถ้าอัตราส่วนของความแปรปรวนเข้าใกล้ 0 แสดงว่าการจัดการสอบเป็นไปตามข้อตกลงเบื้องต้นข้อนี้ (ศิริชัย ภานุจนาวาสี, 2555, หน้า 75-78)

คุณสมบัติของความไม่แปรเปลี่ยนของค่าพารามิเตอร์ (Invariance)

เมื่อ IRT Model สอดคล้อง (Fit) กับข้อมูลจะทำให้เกิดความไม่แปรเปลี่ยนของ พารามิเตอร์ของข้อสอบ (Item Parameter) และพารามิเตอร์ความสามารถของผู้สอบ (Ability Parameter) ซึ่งเป็นคุณสมบัติของ IRT

1) ความไม่แปรเปลี่ยนของพารามิเตอร์ข้อสอบ (Item Parameter) ประกอบด้วย ค่าพารามิเตอร์ของข้อสอบไม่แปรเปลี่ยนไปตามกลุ่มผู้สอบ และ ICC ลักษณะเดียวกัน (a, b, c) สำหรับทุกกลุ่มความสามารถของผู้สอบ แสดงว่า ICC มีความคงที่ข้ามกลุ่มผู้สอบ

2) ความไม่แปรเปลี่ยนของพารามิเตอร์ความสามารถของผู้สอบ (Ability Invariance) ประกอบด้วย ค่าพารามิเตอร์ของผู้สอบไม่แปรเปลี่ยนไปตามชุดของข้อสอบ และเมื่อนำข้อสอบ ต่างชุด (ทุกข้อมุ่งวัดเดียวกัน) เข้า ข้อสอบค่อนข้างง่าย กับข้อสอบชุดค่อนข้างยากมาสอบวัด ผู้สอบกลุ่มเดียวกัน ค่า θ ที่ประมาณได้จากข้อสอบทั้ง 2 ชุด มีความแตกต่างกันไม่เกิน SEE แสดงว่า การประมาณค่าความสามารถมีความคงที่ข้ามชุดของข้อสอบ

สรุป คือ ทฤษฎีการตอบสนองข้อสอบเป็นทฤษฎีการวัดที่อธิบายความสัมพันธ์ระหว่าง คุณลักษณะภายใต้ความสามารถที่มีอยู่ภายในตัวบุคคลกับพฤติกรรมการตอบสนองข้อสอบของ บุคคลนั้น โดยจะใช้โค้งลักษณะข้อสอบ (Item Characteristic Curve: ICC) มีการกำหนดความยาก (b) อำนาจจำแนก (a) และโอกาสการเดาข้อสอบถูก (c) IRT จึงอยู่บนฐานความคิดที่สำคัญ 2 ประการ คือ 1) ผลการตอบข้อสอบหรือข้อคำถาของผู้ตอบ สามารถอธิบายได้ด้วยความสามารถ ที่มีอยู่ภายในของผู้ตอบ และ 2) ความสัมพันธ์ระหว่างผลการตอบข้อสอบกับความสามารถที่มีอยู่ ภายใน สามารถอธิบายได้ด้วยฟังก์ชันลักษณะข้อสอบ หรือโค้งลักษณะข้อสอบ (ICC) ขึ้นมีลักษณะ เป็นฟังก์ชันทางคณิตศาสตร์ เรียกว่าฟังก์ชันโลจิส (Logistic Function) หรือใกล้เคียงกับฟังก์ชันปกติ สะสม (Normal Ogive Function)

งานวิจัยที่เกี่ยวข้องกับทฤษฎีการตอบสนองข้อสอบ (IRT) มีดังนี้

ชนะศัก นิชานนท์ (2554) ได้ศึกษาประสิทธิภาพของการประมาณค่าพารามิเตอร์แบบเบส์ โดยใช้การสรุปอ้างอิงความน่าเชื่อถือของโมเดลการตอบสนองข้อสอบมีวัตถุประสงค์เพื่อเปรียบเทียบ ประสิทธิภาพของวิธีการประมาณค่าพารามิเตอร์ของวิธีการสรุปอ้างอิงความน่าเชื่อถือของโมเดลการ

ตอบสนองข้อสอบ (Generalizability in Item Response) 4 รูปแบบ ได้แก่ รูปแบบที่ 1 Original GIRM พัฒนาโดย Brigg and Wilson (2007) รูปแบบที่ 2 AGIRM A รูปแบบที่ 3 AGIRM B และ รูปแบบที่ 4 Numerical Bayesian GIRM ผู้วิจัยเป็นผู้พัฒนาขึ้นนอกจากรากนี้ยังศึกษาถึงอิทธิพลของ ขนาดกลุ่มตัวอย่างและจำนวนข้อสอบ รวมทั้งยังศึกษาความไว (Sensitivity) ของรูปแบบต่าง ๆ ต่อการกำหนดลักษณะการแจกแจงเริ่มแรกของค่าพารามิเตอร์ของข้อสอบและผู้สอบที่ส่งผลต่อ ประสิทธิภาพของวิธีการประมาณ ซึ่งวัดได้จากดัชนี 3 ประเภท ได้แก่ ความล้าเฉียง ในการประมาณค่า (Biased estimator) คำนวณจากการวิเคราะห์ค่าความคลาดเคลื่อนเฉลี่ย (Mean Average Deviation) ความไม่แน่นอนในการประมาณค่า (Uncertainty estimator) คำนวณจากการวิเคราะห์ ค่าส่วนเบี่ยงเบนมาตรฐาน (Standard Deviation) และประสิทธิภาพในการประมาณค่าของค่าประกอบ ความแปรปรวนซึ่งคำนวณจากการวิเคราะห์ระยะทางยุคลิด (Euclidean Distance) ข้อมูลที่ใช้ใน การศึกษารั้งนี้ใช้ข้อมูลจำลอง (Simulation) จากโปรแกรม R และทำการประมาณผลภายใต้จาก การเขียนคำสั่งการประมาณด้วยโปรแกรม WinBUGS ด้วย Package R2 WinBUGS พบว่า เมื่อ เปรียบเทียบประสิทธิภาพในการประมาณค่าของรูปแบบของวิธีการสรุปอ้างอิงความน่าเชื่อถือของผล การวัดด้วยโมเดลการตอบสนองข้อสอบ (GIRM) พบว่า ความล้าเฉียงในการประมาณค่า รูปแบบที่ 1 กับรูปแบบที่ 4 ให้ค่าประสิทธิภาพสูงที่สุด โดยรูปแบบที่ 4 สามารถประมาณค่าพารามิเตอร์ได้เฉพาะ ลักษณะการแจกแจงเริ่มแรกของผู้สอบและข้อสอบแบบปกติ สำหรับความไม่แน่นอนในการประมาณ ค่า พบว่า รูปแบบที่ 4 ให้ค่าประสิทธิภาพสูงที่สุด สำหรับลักษณะการแจกแจงเริ่มแรกของผู้สอบ และข้อสอบแบบปกติส่วนลักษณะการแจกแจงเริ่มแรกของค่าพารามิเตอร์ตัวใดตัวหนึ่งของผู้สอบ หรือข้อสอบที่ไม่มีลักษณะการแจกแจงเริ่มแรกแบบปกติ พบว่า รูปแบบที่ 1 ให้ค่าประสิทธิภาพสูง ที่สุด และเมื่อพิจารณาในด้านประสิทธิภาพขององค์ประกอบความแปรปรวนยุคลิด พบว่า รูปแบบ ที่ 2 ให้ค่าประสิทธิภาพสูงที่สุด การศึกษาอิทธิพลของขนาดกลุ่มตัวอย่างและจำนวนข้อสอบ พบว่า ขนาดกลุ่มตัวอย่างไม่ส่งผลต่อประสิทธิภาพในการประมาณค่า สำหรับความยาวแบบสอบ พบว่า ส่งผลต่อ การวัดประสิทธิภาพในความล้าเฉียงในการประมาณค่า และการวิเคราะห์ประสิทธิภาพองค์ประกอบ ความแปรปรวนยุคลิด ในทุกรูปแบบ ส่วนความไม่แน่นอนในการประมาณค่า พบว่า การแจกแจง เริ่มแรกของผู้เข้าสอบไม่ส่งผลต่อประสิทธิภาพด้านความล้าเฉียง ส่วนความไม่แน่นอน ในการ ประมาณค่าและประสิทธิภาพขององค์ประกอบความแปรปรวนยุคลิด พบว่า การแจกแจงเริ่มแรกของ ข้อสอบ พบว่า ส่งผลต่อการวัดประสิทธิภาพความล้าเฉียงในการประมาณค่าความไม่แน่นอนใน การประมาณค่าทุกรูปแบบ และส่งผลการวิเคราะห์ประสิทธิภาพองค์ประกอบความแปรปรวนยุคลิด เฉพาะในกรณีที่การแจกแจงเริ่มแรกของผู้สอบเป็นแบบแกรมม่าเท่านั้น

ศัตรา แสนปัญญา (2555) ได้ศึกษาพัฒนาการความสามารถทางวิทยาศาสตร์ของนักเรียน ชั้นมัธยมศึกษาปีที่ 3 ไปยังชั้นมัธยมศึกษาปีที่ 4 โดยการศึกษาภาคตัดขวาง (Cross-Sectional Study) เมื่อปรับเทียบคะแนนแนวตั้ง (Vertical Equating) ตามทฤษฎีการตอบสนองข้อสอบ (IRT) แบบโลจิสติก 3 พารามิเตอร์ ร่วมกับวิธีค่าเฉลี่ยและซิกมา (Mean and Sigma Method) และสมการ ทดถอยโดยวิธีกำลังสองน้อยที่สุด (OLS) ผลที่ได้นำไปสร้างสมการโดยวิธีการทดถอย (Regression

Method) เพื่อท่านายความสามารถทางวิทยาศาสตร์ของนักเรียนชั้นมัธยมศึกษาปีที่ 3 ในอนาคต หลังจากมีพัฒนาการแล้ว เครื่องมือที่ใช้ในการวิจัยครั้งนี้เป็นแบบวัดความสามารถทางวิทยาศาสตร์ ของมัธยมศึกษาปีที่ 3 และมัธยมศึกษาปีที่ 4 ระดับชั้นละ 1 ฉบับ ฉบับละ 30 ข้อ มีความเชื่อมั่นของ แบบทดสอบเท่ากัน 0.70 และ 0.71 ตามลำดับ และมีข้อสอบร่วมภายใน 6 ข้อ ผลการวิจัย พบว่า พัฒนาการทางวิทยาศาสตร์ (Growth Rate) ของนักเรียนชั้นมัธยมศึกษาปีที่ 3 เมื่อขึ้นชั้นมัธยมศึกษา ปีที่ 4 กลุ่ม สูง ปานกลาง และต่ำ มีพัฒนาการเรียงลำดับจากมากไปน้อยตามลำดับ และสมการ มีสัมประสิทธิ์การทำนาย (r^2) ของคะแนนสอบ (x) ความสามารถแฝง (θ) และคะแนนจริง (T) เท่ากับ 0.997 1.000 และ 0.979 ตามลำดับ โดยค่าสัมประสิทธิ์การทำนายของสมการทั้งสามสมการ มีค่านัยสำคัญทางสถิติที่ระดับ 0.05

อนุชิต กลินคำเนต อรจิรา สิทธิ์ศักดิ์ และทศนารรณ ศุนย์กลาง (2555) ได้ศึกษา ผลการประเมินระบบบริหารจัดการการเรียนรู้แบบปรับเปลี่ยน กรณีศึกษาเรื่อง องค์ประกอบ ของระบบสารสนเทศ มีวัตถุประสงค์เพื่อพัฒนาและหาประสิทธิภาพของระบบบริหารจัดการเรียนรู้ แบบปรับเปลี่ยนรูปของเว็บแอปพลิเคชัน (Web Application) นักเรียนเข้าเรียนเนื้อหา โดยแบ่งนักเรียนออกเป็น 3 ระดับ คือ กลุ่มเก่ง กลุ่มปานกลาง และกลุ่มอ่อน ทำการวัดค่า ความสามารถของนักเรียนก่อนเรียนและหลังเรียนผ่านระบบบริหารจัดการเรียนรู้แบบปรับเปลี่ยน ในการวัดค่าความสามารถของนักเรียนใช้ทฤษฎีการตอบสนองข้อสอบ (Item Response Theory: IRT) และการประมาณค่าความสามารถ (θ) โดยใช้กลวิธีของเบย์ (Bayesian Updating) สำหรับ ครุภู่สอนระบบได้จัดเตรียมเครื่องมือสำหรับช่วยให้ผู้สอนเตรียมเนื้อหาที่เรียนให้เหมาะสมสำหรับ นักเรียนแต่ละระดับโดยใช้เทคนิคสื่อหلامมิติแบบปรับตัว และจัดเตรียมแบบทดสอบโดยระบุ ค่าพารามิเตอร์ของข้อสอบเพื่อนำไปใช้ในการทำแบบทดสอบแบบปรับเปลี่ยน ได้ ผลการวิจัยระบบ บริหารจัดการเรียนรู้แบบปรับเปลี่ยนพบร่วมกับค่าความสามารถทางการเรียนก่อนเรียนและหลังเรียน แตกต่างกันอย่างมีนัยสำคัญทางสถิติที่ระดับ .05 โดยค่าความสามารถทางการเรียนหลังเรียนสูงกว่า ก่อนเรียน และผลการประเมินความพึงพอใจของครุภู่สอนที่มีต่อระบบพบว่าครุภู่สอน มีความพึงพอใจ ต่อระบบในภาพรวมอยู่ในระดับความพึงพอใจมากที่สุด ($X = 4.68$, $S.D. = 0.47$)

นุภาพรรณ ปลื้มใจ ปิยะพิพิญ ตินوار และโ戍ส สุขานนท์สวัสดิ์ (2558) ได้พัฒนาโปรแกรม การทดสอบแบบปรับเปลี่ยนด้วยคอมพิวเตอร์สำหรับการจัดสอบ O-NET ระดับชั้นมัธยมศึกษาปีที่ 6 โดยมีวัตถุประสงค์เพื่อวิเคราะห์คุณภาพข้อสอบ O-NET จัดทำคลังข้อสอบ O-NET และพัฒนา โปรแกรมการทดสอบแบบปรับเปลี่ยนด้วยคอมพิวเตอร์ สำหรับการจัดสอบ O-NET ระดับชั้น มัธยมศึกษาปีที่ 6 จำนวน 8 กลุ่มสาระการเรียนรู้ ได้พัฒนาโปรแกรมการทดสอบแบบปรับเปลี่ยนด้วย คอมพิวเตอร์ในรูปแบบของ Web Application การดำเนินการวิจัยมี 4 ขั้นตอน ดังนี้ 1) วิเคราะห์ คุณภาพของข้อสอบ O-NET ตามทฤษฎีการตอบสนองข้อสอบ (Item Response Theory: IRT) แบบ 3 พารามิเตอร์จำนวน 8 กลุ่มสาระการเรียนรู้ 2) จัดทำคลังข้อสอบ O-NET จำนวน 8 กลุ่มสาระการเรียนรู้ 3) พัฒนาโปรแกรมการทดสอบแบบปรับเปลี่ยนด้วยคอมพิวเตอร์ และ 4) ประเมินความคิดเห็นของผู้ทดสอบใช้โปรแกรมการทดสอบแบบปรับเปลี่ยนด้วยคอมพิวเตอร์โดย ทำการทดสอบบนเว็บไซต์ (www.onetcat.net/onem6) กับนักเรียนระดับชั้นมัธยมศึกษาปีที่ 6 ที่กำลังศึกษาอยู่ในภาคเรียนที่ 2 ปีการศึกษา 2557 จำนวน 61 คน ผลการวิจัยปรากฏว่า 1. ข้อสอบ

O-NET ระดับชั้นมัธยมศึกษาปีที่ 6 จำนวน 8 กลุ่มสาระการเรียนรู้ที่ผ่านเกณฑ์การคัดเลือก จำนวน 1,197 ข้อ มีค่าอำนาจจำแนกของข้อสอบเฉลี่ย เท่ากับ 1.3693 ค่าความยากของข้อสอบเฉลี่ย เท่ากับ 0.8624 และค่าการเดาของข้อสอบเฉลี่ย เท่ากับ 0.2024 แสดงให้เห็นว่า ข้อสอบที่อยู่ในคลังข้อสอบ O-NET ค่อนข้างยาก 2. คลังข้อสอบ O-NET ระดับชั้นมัธยมศึกษาปีที่ 6 สามารถบรรจุข้อสอบแบบหลายตัวเลือก ชนิด 4 ตัวเลือก ได้ไม่จำกัดขึ้นอยู่กับขนาดของ Server 3. โปรแกรมการทดสอบแบบปรับเหมาะสมด้วยคอมพิวเตอร์ สำหรับการจัดสอบ O-NET ระดับชั้นมัธยมศึกษาปีที่ 6 มีความถูกต้องและปลอดภัยในการใช้งานเป็นที่ยอมรับของผู้เชี่ยวชาญ และนักเรียนที่ทดลองใช้โปรแกรมประเมินว่ามีความสะดวกในการนำไปใช้งาน

จากรุจิตร สิทธิปูรุ ปิยะทิพย์ ตินوار และโสฬส สุขานันท์ (2559) ได้พัฒนาโปรแกรมการทดสอบแบบปรับเหมาะสมด้วยคอมพิวเตอร์สำหรับการจัดสอบ O-NET ระดับชั้นมัธยมศึกษาปีที่ 3 โดยมีวัตถุประสงค์ เพื่อวิเคราะห์คุณภาพข้อสอบ O-NET จัดทำคลังข้อสอบ O-NET และพัฒนาโปรแกรมการทดสอบแบบปรับเหมาะสมด้วยคอมพิวเตอร์ สำหรับการจัดสอบ O-NET ระดับชั้นมัธยมศึกษาปีที่ 3 จำนวน 8 กลุ่มสาระการเรียนรู้ ซึ่งผู้วิจัยพัฒนาโปรแกรมการทดสอบ แบบปรับเหมาะสมด้วยคอมพิวเตอร์ในรูปแบบของ Web Application วิเคราะห์คุณภาพข้อสอบโดยใช้ทฤษฎี การตอบสนองข้อสอบ (Item Response Theory: IRT) แบบ 3 พารามิเตอร์ โดยใช้โปรแกรมสำเร็จรูป Xcalibre Version 4.1.7 ซึ่งแบ่งวิธีดำเนินการวิจัยออกเป็น 4 ขั้นตอน ดังนี้
 1) การวิเคราะห์คุณภาพข้อสอบ O-NET จำนวน 8 กลุ่มสาระการเรียนรู้ 2) การจัดทำคลังข้อสอบ O-NET จำนวน 8 กลุ่มสาระการเรียนรู้ ระหว่างปี พ.ศ. 2551–2553 3) การพัฒนาโปรแกรมการทดสอบแบบปรับเหมาะสมด้วยคอมพิวเตอร์ และ 4) การประเมินประสิทธิภาพของโปรแกรมการทดสอบแบบปรับเหมาะสมด้วยคอมพิวเตอร์ โดยสรุปผลกระทบต่อบรรจุผลคิดเห็นในการทดลองใช้งานโปรแกรม ซึ่งกลุ่มตัวอย่างเป็นนักเรียนระดับชั้นมัธยมศึกษาปีที่ 3 โรงเรียนอ่างศิลาพิทยาคม จังหวัดชลบุรี จำนวน 30 คนเครื่องมือที่ใช้ในการวิจัย ได้แก่ 1) ข้อสอบ O-NET จำนวน 8 กลุ่มสาระการเรียนรู้ และ 2) โปรแกรมการทดสอบแบบปรับเหมาะสมด้วยคอมพิวเตอร์ โดยทำการทดสอบใน website: www.onetcat.net/onetcatm3 และวิเคราะห์ระดับความคิดเห็นของนักเรียนที่ทดลองใช้โปรแกรม ด้วยสถิติพื้นฐาน ได้แก่ ค่าเฉลี่ย และส่วนเบี่ยงเบนมาตรฐาน ผลการวิจัยปรากฏว่า
 1. การวิเคราะห์คุณภาพข้อสอบ O-NET 8 กลุ่มสาระการเรียนรู้ แสดงให้เห็นว่า ข้อสอบ O-NET ระดับชั้นมัธยมศึกษาปีที่ 3 มีค่าความยากของข้อสอบ (b) เฉลี่ยค่อนข้างยาก 2. คลังข้อสอบ O-NET สามารถบรรจุข้อสอบแบบหลายตัวเลือก (Multiple Choice) ชนิด 4 ตัวเลือกได้ไม่จำกัด ทั้งนี้ขึ้นอยู่กับขนาดของ Server โดยได้บรรจุข้อสอบ O-NET ที่วิเคราะห์ตามทฤษฎีการตอบสนองข้อสอบแบบ 3 พารามิเตอร์ และผ่านเกณฑ์การคัดเลือกข้อสอบ ทั้ง 8 กลุ่มสาระการเรียนรู้ ในระดับชั้นมัธยมศึกษาปีที่ 3 จำนวนทั้งหมด 469 ข้อ

สุชาดา กรเพชรปานี ปิยะทิพย์ ตินوار และโสฬส สุขานันท์สวัสดิ์ (2559) ได้พัฒนาโปรแกรมการทดสอบแบบปรับเหมาะสมด้วยคอมพิวเตอร์สำหรับการจัดสอบ O-NET โดยมีวัตถุประสงค์ เพื่อจัดทำคลังข้อสอบ O-NET และ พัฒนาโปรแกรมการทดสอบแบบปรับเหมาะสมด้วยคอมพิวเตอร์ สำหรับการจัดสอบ O-NET ระดับชั้นประถมศึกษาปีที่ 6 ชั้นมัธยมศึกษาปีที่ 3 และ ชั้นมัธยมศึกษา

ปีที่ 6 ระดับชั้นละ 8 กลุ่มสาระการเรียนรู้ การจัดทำคลังข้อสอบใช้ฐานข้อมูล MySQL และคัดเลือกข้อสอบ O-NET ของสถาบันทดสอบทางการศึกษาแห่งชาติ (องค์กรมหาชน) ระหว่างปี พ.ศ. 2551-2553 ที่ผ่านการวิเคราะห์ข้อสอบตามทฤษฎีการตอบสนองข้อสอบ แบบ 3 พารามิเตอร์ ในระดับชั้นประถมศึกษาปีที่ 6 จำนวน 258 ข้อ ชั้นมัธยมศึกษาปีที่ 3 จำนวน 469 ข้อ และชั้นมัธยมศึกษาปีที่ 6 จำนวน 1,197 ข้อโดยข้อสอบในคลังข้อสอบ O-NET อยู่ในระดับค่อนข้างยากพัฒนาโปรแกรมในรูปแบบของ Web Application โดยใช้ภาษา PHP พัฒนาโปรแกรมการทดสอบแบบปรับเท่ากับความพิเศษตามแนวคิดของ Thompson and Weiss (2011) ผู้ใช้สามารถเข้าใช้โปรแกรมคอมพิวเตอร์สำหรับการจัดสอบ O-NET <http://www.onetcat.net> การประเมินความเหมาะสมใน การนำไปใช้ของโปรแกรมการทดสอบแบบปรับเท่ากับความพิเศษ ในเดือนมิถุนายน พ.ศ. 2557 โดยนักเรียนระดับชั้นประถมศึกษาปีที่ 6 จำนวน 224 คน ชั้นมัธยมศึกษาปีที่ 3 จำนวน 432 คน และชั้นมัธยมศึกษาปีที่ 6 จำนวน 435 คน ปรากฏว่าการใช้งานของโปรแกรมการทดสอบแบบปรับเท่ากับความพิเศษสำหรับการจัดสอบ O-NET อยู่ในเกณฑ์ดี เป็นที่พึงพอใจของนักเรียน

สุชาดา สกลกิจรุ่งโรจน์ เสรี ชัดแข็ง และสมพร สุทัศนีย์ (2559) ได้ศึกษาโปรแกรมการทดสอบแบบปรับเท่ากับความพิเศษสำหรับมาตรฐานขั้นตามแนวคิดทฤษฎีการตอบสนองของข้อสอบร่วมกับการทดสอบแบบปรับเท่ากับความพิเศษ สำหรับวัดความอยู่ดีมีสุขเชิงอัตโนมัติ การวิจัยนี้มีวัตถุประสงค์เพื่อตรวจสอบประสิทธิภาพของโปรแกรมในด้านความประยุกต์และความรวดเร็วในการทดสอบ และศึกษาความสัมพันธ์ของผลการประมาณค่าความสุขที่ได้จากการทดสอบด้วยโปรแกรมการทดสอบแบบปรับเท่ากับความพิเศษ สำหรับมาตรฐานขั้นตามแนวคิดทฤษฎีการตอบสนองของข้อสอบร่วมกับการทดสอบด้วยข้อคำถามทั้งหมดในคลังข้อคำถาม เครื่องมือที่ใช้ในการวิจัย คือ โปรแกรมการทดสอบแบบปรับเท่ากับความพิเศษ สำหรับมาตรฐานความสุขของคนไทยที่พัฒนาขึ้นในงานวิจัยก่อนหน้านี้ ในรูปแบบเว็บแอปพลิเคชัน ซึ่งสามารถเข้าใช้งานผ่านทางเว็บไซต์ <http://www.thscat.com/test/> โดยสามารถเลือกทดสอบได้ 2 ลักษณะ ได้แก่ 1) การทดสอบแบบปรับเท่ากับความพิเศษ (THS-CAT) และ 2) การทดสอบด้วยข้อคำถามทั้งหมดในคลังข้อคำถาม (THS-ทุกข้อ) กลุ่มตัวอย่างที่ใช้ในการวิจัยมีจำนวนทั้งสิ้น 30 คน กลุ่มตัวอย่างจะได้รับการทดสอบทั้งสองรูปแบบ คือ การทดสอบด้วย THS-CAT และ THS-ทุกข้อ วิเคราะห์ข้อมูลด้วยสถิติบรรยาย และวิเคราะห์ความสัมพันธ์ของผลการประมาณค่าความสุขที่ได้จากการทดสอบทั้งสองรูปแบบ ด้วยการหาค่าสัมประสิทธิ์สหสัมพันธ์ เพียร์สัน ผลการวิจัย ปรากฏว่า โปรแกรมการทดสอบแบบปรับเท่ากับความพิเศษ สำหรับ มาตรวัดความสุขของคนไทยสามารถดำเนินการสอบโดยใช้จำนวนข้อคำถามและระยะเวลาในการทดสอบน้อยกว่าการทดสอบด้วยข้อคำถามทั้งหมด ในคลังข้อคำถาม อีกทั้งมีความแม่นยำในการประมาณค่าสูงผลการประมาณค่าที่ได้จากการทดสอบทั้งสองรูปแบบมีความสัมพันธ์กันทางบวก

สุทธิวรรณา พิรศักดิ์สกุล ปิยพงษ์ คล้ายคลึง และสมกิจ กิจพุนวงศ์ (2560) ได้ศึกษาคุณภาพแบบทดสอบความถนัดทางการเรียน SWUSAT ปีการศึกษา 2553-2555 โดยมีวัตถุประสงค์เพื่อ 1) วิเคราะห์คุณภาพของแบบทดสอบความถนัดทางการเรียน SWUSAT ตามทฤษฎีการทดสอบแบบมาตรฐานเดิมและวิธีทฤษฎีการตอบข้อสอบ 2) หาความสัมพันธ์ระหว่างค่าพารามิเตอร์ข้อสอบที่

วิเคราะห์โดยทฤษฎีการทดสอบแบบมาตรฐานเดิมและวิธีทฤษฎีการตอบข้อสอบ และ 3) เปรียบเทียบจำนวนข้อสอบที่ผ่านเกณฑ์ใช้ได้ตามวิธีทฤษฎีแบบมาตรฐานเดิมและวิธีทฤษฎีการตอบข้อสอบข้อมูลที่นำมาวิเคราะห์คุณภาพได้จากการผลการตอบแบบวัดความถนัดทางการเรียน SWUSAT ของนิสิตระดับปริญญาตรี ชั้นปีที่ 1 ระหว่างปีการศึกษา 2553-2555 จำนวน 27 ฉบับ การวิเคราะห์ข้อมูลใช้การวิเคราะห์ข้อสอบเป็นรายข้อตามแนวทางทฤษฎีการทดสอบแบบมาตรฐานเดิมและทฤษฎีการตอบข้อสอบและค่าสัมประสิทธิ์สหสัมพันธ์ระหว่างพารามิเตอร์ของข้อสอบ ผลการวิจัยสรุปได้ดังนี้ ผลการวิเคราะห์ข้อสอบความถนัดทางการเรียน SWUSAT53-55 ด้วยวิธีการทดสอบแบบมาตรฐานเดิมพบว่า แบบทดสอบฉบับที่มีค่าความยากอยู่ในเกณฑ์ใช้ได้ทุกข้อ คือ N53-2, R54-1, S54-2, R55-2, S55-2 และ S55-3 แบบทดสอบฉบับที่มีค่าอำนาจจำแนกอยู่ในเกณฑ์ใช้ได้ทุกข้อ คือ S53-1, S53-2, S54-3, S54-4, S55-2 และ S55-3 นอกนั้นไม่พบแบบทดสอบฉบับใดที่มีค่าความยากหรือค่าอำนาจจำแนกอยู่ในเกณฑ์ใช้ได้ทุกข้อผลการวิเคราะห์ข้อสอบความถนัดทางการเรียน SWUSAT53-55 ด้วยวิธีทฤษฎีการตอบข้อสอบ พบร่วม แบบทดสอบที่มีค่าความยากอยู่ในเกณฑ์ใช้ได้ทุกข้อมีเพียงฉบับเดียว คือ S55-3 แบบทดสอบที่มีค่าความยากเฉลี่ยต่ำสุดคือ S54-3 และความยากเฉลี่ยสูงสุดคือ V55-2 แบบทดสอบที่มีค่าอำนาจจำแนกอยู่ในเกณฑ์ใช้ได้ทุกข้อ คือ V53-1, V53-2, N53-1, R53-1, S53-1, V54-1, N54-1, R54-2, S54-2, V55-1, V55-2, N55-2, R55-2, S55-1, S55-2 และ S55-3 แบบทดสอบที่มีค่าอำนาจจำแนกเฉลี่ยต่ำสุดคือ V53-2 และที่มีค่าอำนาจจำแนกเฉลี่ยสูงสุดคือ N55-1 ค่าสัมประสิทธิ์สหสัมพันธ์ระหว่างค่าความยากที่วิเคราะห์โดยวิธีทฤษฎีการทดสอบแบบมาตรฐานเดิมและวิธีทฤษฎีการตอบข้อสอบ พบร่วมทุกฉบับมีความสัมพันธ์กันทางลบอย่างมีนัยสำคัญทางสถิติ แต่ค่าสัมประสิทธิ์สหสัมพันธ์ระหว่างค่าอำนาจจำแนกที่วิเคราะห์โดยวิธีทฤษฎีการทดสอบแบบมาตรฐานเดิมและวิธีทฤษฎีการตอบข้อสอบ พบร่วม มีความสัมพันธ์กันบางฉบับ คือ V53-2, R53-1, S53-1, S53-2, V54-2, S54-2, N54-2, R54-2, S54-1, V55-1, N55-2, R55-1 และ R55-2 นอกนั้นไม่พบว่ามีความสัมพันธ์กัน 3) การเปรียบเทียบจำนวนข้อที่ใช้ได้ระหว่างการวิเคราะห์ด้วยวิธีการทดสอบ แบบมาตรฐานเดิมและวิธีทฤษฎีการตอบข้อสอบ พบร่วม ฉบับที่มีความแตกต่างกันอย่างมีนัยสำคัญทางสถิติ คือ V53-2, S54-3 และ V54-1 นอกนั้นไม่พบความแตกต่างกัน

สรุบทิพย์ ตรีสิน และปิยะทิพย์ ประคุจพร (2560) ได้เปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบ ระดับชาติ ด้านภาษา ด้านคำนวน และด้านเหตุผล ชั้นประถมศึกษาปีที่ 3 ด้วยวิธี HGLM วิธี MIMIC และวิธี IRT-LR โดยมีวัตถุประสงค์เพื่อวิเคราะห์คุณภาพของแบบทดสอบระดับชาติ (NT) และตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน และด้านเหตุผล ด้วยวิธี HGLM วิธี MIMIC และวิธี IRT-LR ดังนี้ 1) วิเคราะห์คุณภาพของแบบทดสอบระดับชาติ ทั้ง 3 ด้าน 2) ตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ด้วยวิธี HGLM วิธี MIMIC และวิธี IRT-LR และ 3) เปรียบเทียบผลการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบด้วยวิธีการตรวจสอบ 3 วิธีข้อมูลที่นำมาใช้ วิเคราะห์เป็นข้อมูลทุติยภูมิ จากผลการตอบแบบทดสอบระดับชาติของนักเรียนชั้นประถมศึกษาปีที่ 3 ปีการศึกษา 2556 จำนวน 9,600 คน ผลการวิจัยปรากฏว่า 1) แบบทดสอบระดับชาติชั้นประถมศึกษาปีที่ 3 มีค่าความยากของข้อสอบ (b) อยู่ในระดับค่อนข้างยากมีค่าอำนาจจำแนกของข้อสอบ (a) อยู่ในระดับที่สามารถจำแนกผู้สอบได้ดี และมี

ค่าโอกาสในการเดาของข้อสอบ (c) ไม่เกิน 0.3 2) การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบทั้ง 3 ด้าน ซึ่งให้เห็นว่า เพศส่งผลให้เกิดการทำหน้าที่ต่าง กันของข้อสอบ โดยเพศหญิงจะได้เปรียบในการตอบข้อสอบด้านภาษา และด้านเหตุผล ในขณะที่เพศชาย จะได้เปรียบในการตอบข้อสอบด้านคำนวณ โดยวิธี HGLM ตรวจพบข้อสอบทำหน้าที่ต่างกัน จำนวนมากที่สุด คิดเป็นร้อยละ 69 ของข้อสอบทั้งฉบับ รองลงมาคือ วิธี IRT-LR ร้อยละ 54 และวิธี MIMIC ร้อยละ 16 ตามลำดับ 3) การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ พบว่า วิธี HGLM ตรวจพบ DIF มากกว่า วิธี MIMIC ในด้านภาษา ด้านคำนวณ และด้านเหตุผล คิดเป็นร้อยละ 70, 36 และ 53 ตามลำดับ และวิธี HGLM ตรวจ พบ DIF มากกว่าวิธี IRT-LR ด้านภาษา และด้านคำนวณ คิดเป็นร้อยละ 37 และ 13 และวิธี IRT-LR ตรวจพบ DIF มากกว่าวิธี MIMIC ทั้ง 3 ด้าน คิดเป็นร้อยละ 33, 43 และ 40 ตามลำดับ ส่วนวิธี HGLM ตรวจพบ DIF น้อยกว่า วิธี IRT-LR ด้านคำนวณคิดเป็นร้อยละ 7 ($p<.05$)

Yildirim and Berberoglu (2009) ได้วิเคราะห์ข้อมูลการตัดสินใจและทางสถิติของรายการวิชาคณิตศาสตร์ PISA-2003 การเปรียบเทียบลักษณะของนุյงย์ในกลุ่มภาษาและวัฒนธรรม ที่แตกต่างกันกลยุทธ์เป็นสิ่งที่สำคัญมากขึ้นในแนวทางการประเมินผลการศึกษา ความสนใจที่เพิ่มขึ้นในการศึกษาเปรียบเทียบระหว่างประเทศผลลัพธ์ในด้านภาษาและกลุ่มทางวัฒนธรรมที่แตกต่างกัน ทำให้ความสนใจของนักวิจัย ในการศึกษาครั้งนี้มีสมมติฐานที่สร้างขึ้นเกี่ยวกับการทำงานของข้อสอบ ที่แตกต่างกัน (DIF) จะได้รับการศึกษาโดยใช้ชุดข้อมูลสำหรับการประเมินนักศึกษาต่างชาติ PISA 2003 ข้อสอบที่พบ DIF ผ่านการวิเคราะห์ของค์ประกอบ Mantel-Haenszel (MH) และการวิเคราะห์อัตราส่วนไลคลิสต์ความเป็นไปได้ในการตอบสนองข้อสอบ (IRT-LR) ได้รับการพิจารณาพร้อม ๆ กัน สำหรับการประเมินเนื้อหาที่ขยายผล กลุ่มผู้ทดสอบ ด้านทักษะทางปัญญาที่วัดได้จากข้อบกพร่อง ด้านข้อผิดพลาดในการแปลและด้านการใช้คำเชิงปริมาณเป็นข้อผิดพลาดสำคัญสามข้อที่อาจก่อให้เกิด DIF ในการศึกษาของ PISA 2003 ผลการวิจัยพบว่าศึกษาการวิเคราะห์ข้อมูลการตัดสินใจและทางสถิติของวิชาคณิตศาสตร์ PISA-2003 พบว่า วิธี IRT-LR มีประสิทธิภาพในการตรวจสอบ DIF ได้ดีในด้านความสามารถทางคณิตศาสตร์ของโครงการประเมินผลนักเรียนนานาชาติ (PISA, 2003)

Muninsakorn, Tinnaworn, and Sukhanonsawat (2016) ได้พัฒนาโปรแกรมทดสอบแบบปรับเหมาะสมด้วยคอมพิวเตอร์สำหรับการจัดสอบ O-NET ระดับชั้นประถมศึกษาปีที่ 6 โดยมีวัตถุประสงค์เพื่อวิเคราะห์คุณภาพของข้อสอบ O-NET จัดทำคลังข้อสอบ O-NET และพัฒนาโปรแกรมการทดสอบแบบปรับเหมาะสมด้วยคอมพิวเตอร์ สำหรับการจัดสอบ O-NET ระดับชั้นประถมศึกษาปีที่ 6 จำนวน 8 กลุ่มสาระการเรียนรู้ ที่ดำเนินการวิจัยแบ่งเป็น 4 ขั้นตอน ดังนี้ 1) วิเคราะห์คุณภาพของข้อสอบ O-NET จำนวน 8 กลุ่มสาระการเรียนรู้ 2) จัดทำคลังข้อสอบ O-NET จำนวน 8 กลุ่มสาระการเรียนรู้ 3) พัฒนาโปรแกรมการทดสอบแบบปรับเหมาะสมด้วยคอมพิวเตอร์ และ 4) ประเมินความคิดเห็นของผู้ทดลองใช้โปรแกรมการทดสอบแบบปรับเหมาะสมด้วยคอมพิวเตอร์ ซึ่งเป็นนักเรียนระดับชั้นประถมศึกษาปีที่ 6 จังหวัดชลบุรี จำนวน 30 คน ผลการวิจัยปรากฏว่า 1. ข้อสอบ O-NET ระดับชั้นประถมศึกษาปีที่ 6 มีค่าความยากของข้อสอบ (b) ในระดับค่อนข้างยาก 2. คลังข้อสอบ O-NET ระดับชั้นประถมศึกษาปีที่ 6 บรรจุข้อสอบแบบหลายตัวเลือก (Multiple Choice) ชนิด 4 ตัวเลือก ได้โดยไม่จำกัด ขึ้นอยู่กับขนาดของ Server ซึ่งมีข้อสอบจำนวน 258 ข้อ

ที่ผ่านเกณฑ์การวิเคราะห์คุณภาพของข้อสอบตามโมเดลโลจิส แบบ 3 พารามิเตอร์ใน 8 กลุ่มสาระ การเรียนรู้ 3. โปรแกรมการทดสอบแบบปรับเหมาะสมด้วยคอมพิวเตอร์ สำหรับการจัดสอบ O-NET ระดับชั้นประถมศึกษาปีที่ 6 มีความเหมาะสมของโปรแกรมในระดับมากที่สุด ไม่มีปัญหาด้านการนำไปใช้และเป็นที่ยอมรับของที่ผู้ทดลองใช้โปรแกรม

Moghadamzadeh, Salehi, and Khodaei (2011) ได้เปรียบเทียบฟังก์ชันสารสนเทศ ของข้อสอบ และฟังก์ชันสารสนเทศของแบบสอบในโมเดลการตอบสนองข้อสอบแบบ 1 พารามิเตอร์ 2 พารามิเตอร์ และ 3 พารามิเตอร์ โดยเริ่มต้นจากการสร้างฟังก์ชันสารสนเทศภายใต้ทฤษฎี การตอบสนองข้อสอบ จากนั้นจึงใช้ข้อมูลการตอบข้อสอบที่ได้จากคลังข้อสอบซึ่งเป็นข้อสอบแบบ หลายตัวเลือก ซึ่งได้มาจากการสอบของนักศึกษามหาวิทยาลัยอิหร่านที่ทำการทดสอบ ระดับชาติที่มีชื่อว่า “National Organization for Educational Testing (NOET)” ซึ่งเป็น การทดสอบคัดเลือกเข้าศึกษาต่อในระดับมหาวิทยาลัย ในสาขาวิชาด้านคณิตศาสตร์และพิสิกส์ ในปี ค.ศ. 2009 ผู้วิจัยเลือกข้อมูลการทดสอบ 2,000 ชุด ด้วยวิธีการเลือกอย่างเป็นระบบ จากนั้น จัดวิเคราะห์ข้อมูลโดยใช้โปรแกรม SPSS และ BILOG ผลการศึกษาปรากฏว่าค่าสารสนเทศของ ข้อสอบ และค่าสารสนเทศของแบบสอบที่ได้จากการประมาณค่าด้วยโมเดลการตอบสนองข้อสอบ แบบ 2 พารามิเตอร์ มีค่าสูงกว่าโมเดลการตอบสนองข้อสอบแบบ 1 พารามิเตอร์ และแบบ 3 พารามิเตอร์ นอกจากนี้ ผลการวิจัยยังชี้ให้เห็นว่าค่าฟังก์ชันสารสนเทศของข้อสอบ และ ค่าฟังก์ชันสารสนเทศของแบบสอบที่ได้จากโมเดลการตอบสนองข้อสอบแบบ 1 พารามิเตอร์ ไม่แตกต่างจากค่าที่ได้จากโมเดลการตอบสนองข้อสอบแบบ 3 พารามิเตอร์ ซึ่งหมายความว่าโมเดล การตอบสนองข้อสอบแบบ 2 พารามิเตอร์ ทำให้การทดสอบมีความแม่นยำมากกว่าโมเดล การตอบสนองข้อสอบแบบหนึ่งพารามิเตอร์และแบบสามพารามิเตอร์ โดยเหตุผลที่โมเดล การตอบสนองข้อสอบแบบสามพารามิเตอร์มีค่าสารสนเทศต่ำกว่าอาจเนื่องมาจากพารามิเตอร์ การเดาที่เพิ่มเข้ามาในโมเดลซึ่งทำให้เกิดความคลาดเคลื่อนในการประมาณค่าความสามารถ ได้มากกว่า เมื่อมีความคลาดเคลื่อนในการประมาณค่ามากกว่า จึงทำให้สารสนเทศของแบบสอบ มีค่าต่ำกว่าด้วย

Kose and Demirtasli (2012) ได้เปรียบเทียบโมเดลการตอบสนองข้อสอบแบบเอก มิติและโมเดลการตอบสนองข้อสอบแบบพหุมิติภายในตัวอย่างที่มีความยาวของแบบสอบ และขนาดกลุ่ม ตัวอย่างที่แตกต่างกัน โดยใช้ข้อมูลการทดสอบภาษาตุรกี โดยขนาดของกลุ่มตัวอย่าง และความยาว ของแบบสอบเป็นตัวแปรต้นที่ถูกจัดกระทำในการศึกษาครั้งนี้ ข้อมูลการทดสอบได้จากนักเรียนเกรด 8 จำนวน 1,516 คน ผลการศึกษาปรากฏว่าการประมาณค่าพารามิเตอร์ข้อสอบ และพารามิเตอร์ ความสามารถของบุคคลโดยใช้โมเดลการตอบสนองข้อสอบแบบพหุมิติมีค่าความคลาดเคลื่อนใน การประมาณค่าต่ำกว่า ขณะที่มีความแม่นยำในการวัดสูงกว่า นอกจากนี้พบว่าข้อมูลการทดสอบมี ความสอดคล้องกับโมเดลการตอบสนองข้อสอบแบบพหุมิติเป็นอย่างดี และยังพบด้วยว่าความยาว ของแบบสอบ และขนาดของกลุ่มตัวอย่างไม่มีใด ๆ ต่อความสามารถคล้องของข้อมูลการทดสอบกับ โมเดลการตอบสนองข้อสอบแบบเอกมิติ แต่ขนาดกลุ่มตัวอย่างที่ใหญ่ขึ้น และแบบสอบที่มีความยาว มากขึ้นสามารถช่วยลดความคลาดเคลื่อนในการทดสอบได้และมีผลต่อความสามารถคล้องของข้อสอบ

กับโมเดลการตอบสนองข้อสอบแบบพหุมิติ

สรุปจากการศึกษางานวิจัยที่เกี่ยวข้อง ทฤษฎีการตอบสนองข้อสอบ จะมีความสัมพันธ์ระหว่างคุณลักษณะระหว่างคุณลักษณะภายในหรือความสามารถที่มีอยู่ภายในตัวบุคคลกับพฤติกรรมการตอบสนองข้อสอบของบุคคลนั้นว่ามีโอกาสตอบข้อสอบถูกมากน้อยเพียงไร ทฤษฎีนี้มีพื้นฐานความเชื่อว่าพฤติกรรมการตอบสนองต่อ ข้อสอบของผู้สอบ ซึ่งเป็นสิ่งที่สังเกตได้โดยตรงว่าถูกหรือผิด จะถูกกำหนดโดยคุณลักษณะภายในหรือความสามารถที่อยู่ภายในตัวบุคคลซึ่งเป็นสิ่งที่ไม่สามารถสังเกตได้โดยตรง ทฤษฎีนี้ได้โดยอิสายความสัมพันธ์ดังกล่าวในรูปของฟังก์ชันคณิตศาสตร์หรือโมเดลที่แสดงความสัมพันธ์ระหว่างระดับความสามารถคุณลักษณะของข้อสอบและโอกาสของการตอบข้อสอบได้ถูก ที่เรียกว่า ฟังก์ชันการตอบสนองข้อสอบ ซึ่งมีลักษณะความสัมพันธ์เป็นแบบฟังก์ชันโลจิสหรือฟังก์ชันปกติสะสม

ตอนที่ 3 การทำหน้าที่ต่างกันของข้อสอบ (Differential Item Functioning: DIF) และงานวิจัยที่เกี่ยวข้อง

การตรวจสอบคุณภาพของข้อที่เป็นผลมาจากการตัดสินว่าข้อสอบนั้นมีความยุติธรรมหรือไม่อย่างไร ซึ่งจะเป็นข้อมูลสารสนเทศทางสถิติและจะมีความสัมพันธ์ที่ส่งผลไปยังคุณลักษณะที่ข้อสอบมุ่งที่จะวัด กับมวลประสบการณ์ของผู้สอบในกลุ่มต่าง ๆ ที่มีบางอย่างที่แตกต่างกัน เช่น เขื้อชาติ ศาสนา วัฒนธรรม ภูมิลำเนา สังคม เพศ ภาษา อายุ และประสบการณ์เป็นต้น โดยข้อสอบที่จะมุ่งวัดนั้นอาจทำให้เกิดการได้เปรียบเสียเปรียบทั้งที่ความสามารถของผู้สอบมีความเท่าเทียมกันทำให้เกิดความแตกต่างกับความสามารถของข้อสอบจึงทำใช้ คำที่ว่า ข้อสอบลำเอียง (Biased item) ซึ่งคำที่ว่านี้จะมีความหมายในเชิงลบ แล้วเกณฑ์ที่ใช้สำหรับในการตัดสินว่าข้อสอบนั้นมีความลำเอียงยังไง ชัดเจนเท่าที่ควรนัก จึงได้เปลี่ยนมาใช้คำว่า การทำหน้าที่ต่างกันของข้อสอบ (Differential Item Functioning: DIF) จะเห็นได้ว่าความหมายจะเป็นกลางสำหรับผู้สอบ (Holland & Thayer, 1988; Honlland & Wainer, 1993; ศิริชัย กาญจนวารี, 2555, หน้า 115-116) และมีนักวัดผลหลายท่านที่ได้ให้ความหมายของการทำหน้าที่ต่างกันของข้อสอบไว้ดังนี้

ความลำเอียงของข้อสอบ หมายถึง ผู้ที่เข้าสอบจะมีสัดส่วนในการตอบคำถามข้อสอบได้ถูกต้องไม่เท่ากันของแต่ละกลุ่มประชากรที่จะศึกษา และเมื่อกลุ่มผู้เข้าสอบได้คะแนนที่เท่ากันนั้น ข้อสอบก็เป็นเอกพันธ์ (Scheuneman, 1979; ศิริชัย กาญจนวารี, 2555, หน้า 116)

ความลำเอียงของข้อสอบ หมายถึง ข้อสอบที่มีค่าความยากจะมีความสัมพันธ์ต่อเนื่องกับกลุ่มผู้เข้าสอบกลุ่มนั้นมากกว่ากลุ่มผู้เข้าสอบอีกกลุ่ม (Rudner, Getson, & Knight, 1980)

ความลำเอียงของข้อสอบ หมายถึง แนวโน้มความลำเอียงที่ใช้คะแนนจากผลการสอบนำมาตัดสินผลเป็นจะทำให้เกิดความไม่ยุติธรรม (Popham, 1981) ความลำเอียงของข้อสอบหมายถึง ข้อสอบที่วัดความสามารถในการที่เลือกคำตอบมีโอกาสในการตอบที่ถูกได้แตกต่างกัน และโอกาสการตอบในเชิงบวกจะแตกต่างกันในการวัดเจตคติ เนื่องจากผู้เข้าสอบมีคุณลักษณะของการประเมินที่เท่ากัน แต่จะมาจากกลุ่มประชากรย่อยที่แตกต่างกัน (Hulin, Drasgow, & Parson, 1983)

ความลำเอียงของข้อสอบ หมายถึง ในการเลือกตอบข้อสอบนั้นจะมีโอกาสในการตอบถูกต้องของผู้เข้าสอบกลุ่มนั้นมากกว่าค่าคาดคะเนต่างกันหรือคะแนนอาจสูงกว่าอีกกลุ่มนั้นของผู้เข้าสอบโดยที่มีความสามารถในระดับเดียวกัน (Lederman, 1990)

การทำหน้าที่ต่างกันของข้อสอบ หมายถึง จะเป็นแบบพหุมิติของข้อสอบที่จะวัดจะได้จากข้อมูลการแยกแยะความสามารถหลัก (Primary Ability) ของผู้เข้าสอบตั้งแต่ 2 กลุ่ม ขึ้นไปที่มีความสามารถเท่ากันโดยมีการแยกแยะความสามารถของลงมา (Secondary) จะแตกต่างจากกัน (Camilli & Shepard, 1994)

การทำหน้าที่แตกต่างกันของข้อสอบ หมายถึง จะเป็นการตอบสนองที่มีพังก์ชันคำนวนจากกลุ่มผู้เข้าสอบที่เป็นกลุ่มย่อยที่ต่างกัน และก็มีค่าที่ไม่เท่ากัน (Narayanan & Swaminathan, 1996)

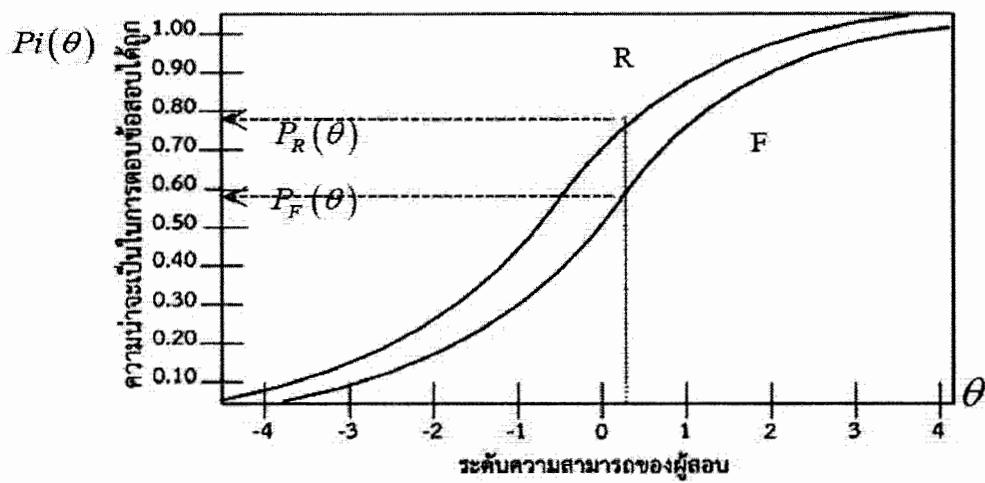
การทำหน้าที่ต่างกันของข้อสอบ ตามที่ได้กล่าวมานี้สามารถที่สรุปได้ว่า การทำหน้าที่ต่างกันของข้อสอบ หมายถึง ผู้เข้าสอบข้อสอบที่ต่างกลุ่มกันที่มีความสามารถหรือคุณลักษณ์ที่มุ่งวัดเท่ากัน ซึ่งทำให้มีโอกาสในการตอบข้อสอบที่ถูกได้แตกต่างกันออกไป หรือมีพังก์ชันการตอบสนองข้อสอบที่แตกต่างกัน และการทำหน้าที่ต่างกันของข้อสอบจะเกิดขึ้นเมื่อนำข้อสอบไปทดสอบกับผู้เข้าสอบกลุ่มย่อยที่ต่างกัน โดยจะมีความสามารถหลัก (Primary Ability) ในระดับเดียวกันหรือมีคุณลักษณะแฝง (Latent Trait) ที่จะวัดเท่ากัน และจะมีความสามารถของลงมา (Secondary Ability) ที่แตกต่างซึ่งจะทำให้ผู้เข้าสอบที่ต่างกลุ่มกันเมื่อนำมาเปรียบเทียบกันจะทำให้โอกาสในการตอบข้อสอบนั้นแตกต่างกัน (ศิริชัย กาญจนวารี, 2555, หน้า 116-117)

ประเภทของการทำหน้าที่ต่างกันของข้อสอบ

การเปรียบเทียบผลการตอบของข้อสอบ ที่ทำหน้าที่ต่างกันของข้อสอบจะเป็นการเปรียบเทียบกลุ่มผู้เข้าสอบ 2 กลุ่มขึ้นไป โดยจะประกอบด้วยกลุ่มแรกคือ กลุ่มเปรียบเทียบ (Focal Group หรือกลุ่ม F) เป็นกลุ่มที่นำเสนอด้วยศึกษาและคาดว่าจะเป็นกลุ่มที่จะเสียเปรียบในตอบข้อสอบ ส่วนกลุ่มที่ 2 เป็นกลุ่มอ้างอิง (Reference Group หรือ กลุ่ม R) เป็นกลุ่มที่คาดว่าจะได้เปรียบในการตอบข้อสอบได้ถูกต้องมากกว่ากลุ่มแรก

การวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบ ซึ่งข้อสอบจะทำหน้าที่ต่างกันของข้อสอบแตกต่างกันได้ 2 ประเภท (Mellenbergh, 1982) คือ การทำหน้าที่ต่างกันของข้อสอบแบบเอกรูป (Unifrom DIF) และแบบอนุรูป (Nonunifrom DIF)

1. ข้อสอบทำหน้าที่ต่างกันแบบเอกรูป (Unifrom DIF) หมายถึง ข้อสอบที่กลุ่มผู้เข้าสอบกลุ่มนั้นมีโอกาสในการตอบข้อสอบได้มากกว่าอีกกลุ่มผู้เข้าสอบอย่างสม่ำเสมอ กัน ในทุกระดับความสามารถ และพิจารณาโคงลักษณะข้อสอบของผู้สอบทั้ง 2 กลุ่ม จะเห็นได้ว่าไม่มีปฏิสัมพันธ์กันระหว่างระดับความสามารถของผู้สอบกับการเป็นสมาชิกของกลุ่ม (Group Membership) ดังภาพที่ 2-1



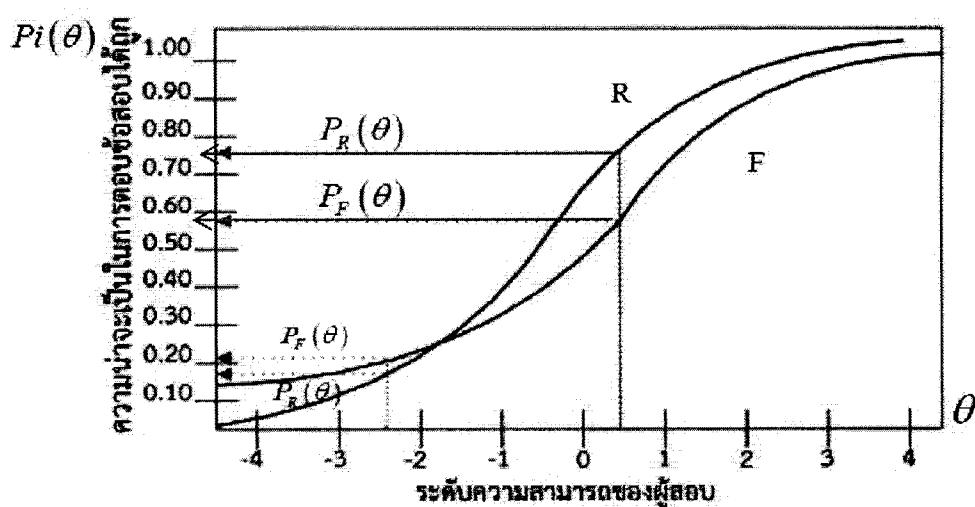
ภาพที่ 2-1 ข้อสอบทำหน้าที่ต่างกันแบบเอกรูป (Uniform DIF) (ศิริชัย กาญจนวاسي, 2555,
หน้า 118)

2. ข้อสอบทำหน้าที่ต่างกันแบบเอนกรูป (Nonuniform DIF) หมายถึง ข้อสอบที่ผู้เข้าสอบ มีโอกาสในการตอบถูกในระหว่างผู้เข้าสอบทั้ง 2 กลุ่มมีความสามารถต่างกันไม่สม่ำเสมอ กันในทุกระดับความสามารถ และเมื่อพิจารณาโดยคุณลักษณะของข้อสอบทั้งกลุ่ม 2 กลุ่ม จะเห็นได้ว่าจะมีปฏิสัมพันธ์กันในระหว่างระดับความสามารถของผู้เข้าสอบ กับการเป็นสมาชิกกลุ่ม เช่น ระดับความสามารถหนึ่งกลุ่มผู้เข้าสอบกลุ่มอ้างอิง (R) จะมีโอกาสในการตอบข้อสอบถูกได้มากกว่ากลุ่มผู้เข้าสอบเปรียบเทียบ (F) และที่ระดับความสามารถอีกระดับหนึ่งกลุ่มผู้เข้าสอบกลุ่มเปรียบเทียบ (F) มีโอกาสตอบข้อสอบได้น้อยกว่ากลุ่มอ้างอิง (R)

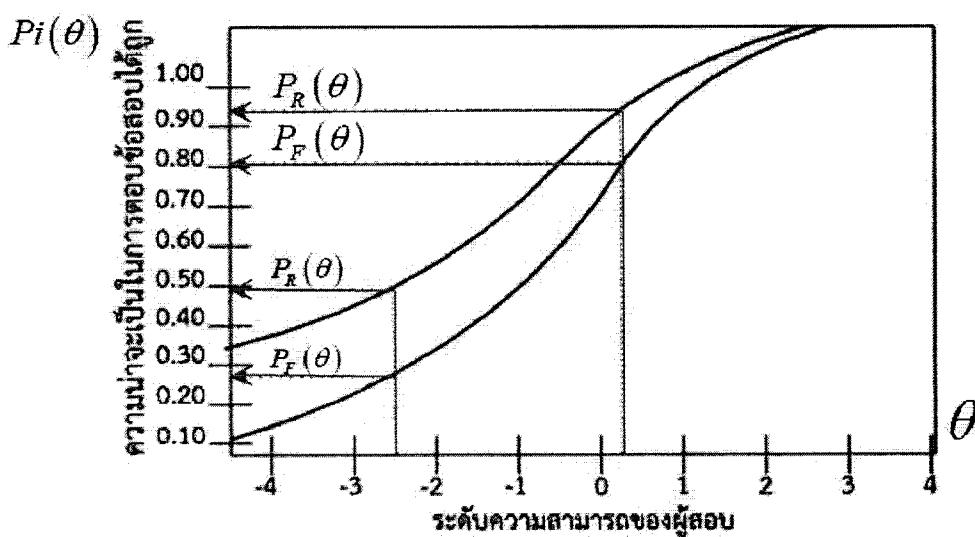
ตามทฤษฎีการตอบสนองข้อสอบ (Item Response Theory: IRT) สามารถพิจารณา “ปฏิสัมพันธ์” ที่ได้จากการความแตกต่างของค่าพารามิเตอร์อำนาจจำแนกของข้อสอบในระหว่างผู้เข้าสอบกลุ่มย่อยสองกลุ่ม ถ้าข้อสอบทำหน้าที่ต่างกันแบบเอกรูป โดยคุณลักษณะของข้อสอบ (Item Characteristic Curves: ICCs) ในกลุ่มผู้เข้าสอบทั้ง 2 กลุ่มจะเป็นเส้นที่ขนานกัน หรือฟังก์ชันการตอบสนองข้อสอบ (Item Response Functions: IRFs) เมื่อเทียบกัน ส่วนข้อสอบทำหน้าที่ต่างกันแบบเอนกรูป โดยคุณลักษณะข้อสอบระหว่างผู้เข้าสอบกลุ่มย่อยทั้ง 2 กลุ่มไม่ขนาน หรือ มีฟังก์ชันการตอบสนองข้อสอบต่างกัน และความแตกต่างระหว่างโดยคุณลักษณะข้อสอบทั้งสองแบบ จะแสดงถึงขนาดและทิศทางของข้อสอบที่ทำหน้าที่ต่างกัน

ข้อสอบที่ทำหน้าที่ต่างกันแบบเอนกรูป จำแนกได้เป็น 2 ลักษณะ (Swaminathan & Roger, 1990)

1. ข้อสอบทำหน้าที่ต่างแบบเอนกรูปที่มีปฏิสัมพันธ์ไม่เป็นลำดับ (Disordinal interaction) เป็นการทำหน้าที่ต่างกันสำหรับกลุ่มผู้เข้าสอบที่เกิดขึ้น เมื่อโดยคุณลักษณะของข้อสอบจะตัดกันระหว่างช่วงความสามารถของผู้เข้าสอบเรียกว่า ข้อสอบทำหน้าที่ต่างกันแบบไม่มีทิศทาง (Non-Unidirectional DIF) ดังภาพที่ 2-2



ภาพที่ 2-2 ข้อสอบทำหน้าที่ต่างกันแบบไม่มีทิศทาง (Non-unidirectional DIF)
(ศิริชัย กาญจนวاسي, 2555, หน้า 119)



ภาพที่ 2-3 ข้อสอบทำหน้าที่ต่างกันแบบมีทิศทางเดียวกัน (Unidirectional DIF)
(ศิริชัย กาญจนวاسي, 2555, หน้า 120)

หลักการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ
การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ (DIF Detection) เป็นการเปรียบเทียบ
ผลการตอบข้อสอบเป็นรายข้อระหว่างกลุ่มผู้สอบอย่างน้อย 2 กลุ่ม ที่มีความสามารถหลัก (Primary Ability) ที่มุ่งวัดเท่ากัน แต่คาดว่าจะมีความได้เปรียบในการตอบข้อสอบข้อนั้น หรือมีโอกาสตอบ

ข้อสอบได้ถูกต้องมากกว่า ส่วนอีกกลุ่มคือ กลุ่มเปรียบเทียบ (Focal Group) ซึ่งเป็นกลุ่มที่สนใจศึกษาและคาดว่าจะเป็นกลุ่มที่เสียเปรียบ

ในการเปรียบเทียบผลการตอบข้อสอบระหว่างกลุ่มอ้างอิงและกลุ่มเปรียบเทียบจำเป็นต้องจับคู่ (Matching) ผู้สอบตามความสามารถ ซึ่งเป็นเงื่อนไขสำคัญของการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ เกณฑ์การจับคู่ (Matching Criteria) ที่นิยมใช้กันมี 2 วิธี ดังนี้

1. เกณฑ์ภายนอก (External Criterion)

การวิเคราะห์การทำหน้าที่ต่างกัน โดยใช้คะแนนจากแบบสอบถามเป็นเกณฑ์ภายนอก แล้วใช้เทคนิคการวิเคราะห์การทดถอย (Regression Analysis) เพื่อทำการเปรียบเทียบเส้นกราฟความล้มพั้นธ์ระหว่างตัวแปรเกณฑ์ กับตัวแปรทำนายระหว่างกลุ่มอ้างอิงและกลุ่มเปรียบเทียบ

หลักการนี้มีจุดมุ่งหมาย เพื่อสร้างสมการทำนายตัวแปรเกณฑ์ ซึ่งเป็นคะแนนของ แบบสอบถามอื่นจากตัวแปรทำนายซึ่งเป็นคะแนนรายข้อหรือคะแนนแบบสอบถามระหว่างกลุ่มอ้างอิง และกลุ่มเปรียบเทียบในการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบ จะใช้คะแนนรวมข้อเป็นตัวแปรทำนาย แต่ถ้าเป็นการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบ จะใช้คะแนนรวมของแบบสอบถาม ทั้งฉบับ เป็นตัวแปรทำนายสำหรับตัวแปรเกณฑ์ที่ใช้เป็นเกณฑ์ภายนอก อาจใช้คะแนนรวมทั้งฉบับ หรือเกรดเฉลี่ย หรือผลสัมฤทธิ์ในงานที่เกี่ยวข้องของผู้สอบ (Cronbach, 1970) สมการทำนายสำหรับกลุ่มอ้างอิงและกลุ่มเปรียบเทียบแสดงได้ดังนี้

กลุ่มอ้างอิง (R)

$$Y_i = A_R + B_R X_i$$

กลุ่มเปรียบเทียบ(F)

$$Y_i = A_F + B_F X_i$$

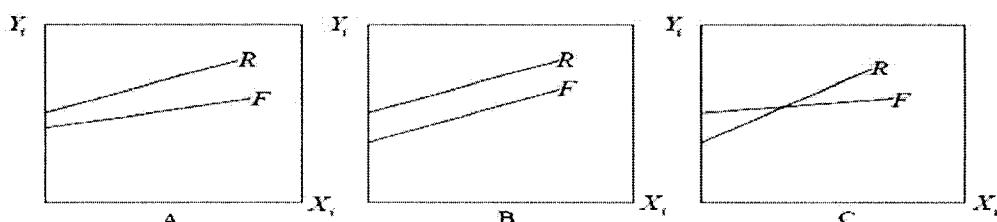
เมื่อ

Y_i = คะแนนของตัวแปรเกณฑ์ภายนอก

X_i = คะแนนของตัวแปรทำนาย

A = ค่าคงที่หรือค่าตัดแกน (Intercept)

B = ค่าความชัน (Slope)



ภาพที่ 2-4 การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบและแบบสอบถามโดยใช้เทคนิคการวิเคราะห์สมการทดถอย (ศิริชัย กาญจนวารี, 2555, หน้า 121)

จากฟังก์ชันการทำนายทั้ง 2 สมการ สามารถเปรียบเทียบค่าตัดแกน (A) และค่าความชัน (B) ของเส้นกราฟระหว่างกลุ่มอ้างอิง (R) และกลุ่มเปรียบเทียบ (F) ได้ ถ้าเส้นกราฟดังกล่าวมีค่าความชันหรือค่าตัดแกน แตกต่างกันสำหรับข้อสอบใดหรือแบบสอบถามใดแสดงว่าข้อสอบ

หรือแบบสอบนั้นมีการทำหน้าที่ต่างกัน โดยลำเอียงเข้าข้างกลุ่มผู้สอบที่มีค่าตัดแgnสูงกว่าหรือค่าความซันที่สูงกว่า

การใช้เกณฑ์ภายนอกมีข้อดี คือเกณฑ์ที่ใช้มีความเป็นอิสระจากข้อสอบ และแบบสอบที่ต้องการตรวจสอบ แต่มีจุดอ่อนตรงที่ความเหมาะสมของเกณฑ์ที่จะนำมาใช้ ในทางปฏิบัติ เป็นการยากที่จะหาตัวแปรเกณฑ์ภายนอกจากแบบสอบฉบับอื่นที่มีความตรงเชิงทำนาย และมีความยุติธรรมสำหรับกลุ่มอ้างอิงและกลุ่มเปรียบเทียบ ถ้าตัวแปรเกณฑ์ภายนอกขาดคุณสมบัติดังกล่าวจะทำให้ผลการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบหรือแบบสอบขาดความแม่นยำและความสมบูรณ์

2. เกณฑ์ภายใน (Internal Criterion)

การวิเคราะห์การทำหน้าที่ต่างกัน โดยใช้เกณฑ์ภายในเป็นการนำวิธีการทางสถิติมาตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ หรือแบบสอบ โดยเน้นการพิจารณาโครงสร้างภายในของแบบสอบเป็นหลัก ด้วยการวิเคราะห์ผลการจากตอบข้อสอบและความสามารถหรือคะแนนจริงของผู้สอบที่ได้จากการสอบฉบับนั้น เพื่อนำมาเปรียบเทียบระหว่างผู้สอบจากกลุ่มอ้างอิงและกลุ่มเปรียบเทียบ ที่มีความสามารถหรือคะแนนจริงเท่ากัน ว่าจะมีผลการตอบหรือโอกาสในการตอบข้อสอบได้ถูกต้องแตกต่างกันหรือไม่ เพื่อบ่งชี้ถึงการทำหน้าที่ต่างกันของข้อสอบ การวิเคราะห์ในลักษณะนี้นิยมใช้ค่าสถิติต่าง ๆ เป็นตัวชี้ถึงการทำหน้าที่ต่างกันของข้อสอบ ค่าสถิติทดสอบที่นิยมนำมาใช้พื้นฐานได้ดังนี้

2.1 การทดสอบปฎิสัมพันธ์ (Interaction)

ในระยะเริ่มแรกของการศึกษาความลำเอียงของข้อสอบ มีการใช้สถิติทดสอบ เอฟ (F-test) จากการวิเคราะห์ความแปรปรวน (ANOVA) เพื่อทดสอบปฎิสัมพันธ์ระหว่างกลุ่มผู้สอบ กับข้อสอบ ถ้าการทดสอบมีนัยสำคัญเป็นสัญญาณของการทำหน้าที่ต่างกันของข้อสอบ (Cleary & Hilton 1968; Jensen, 1974) จากนั้นจึงทำการวิเคราะห์ต่อด้วยวิธีการ Post Hoc เพื่อรับข้อสอบที่มีผลต่อการเกิดปฎิสัมพันธ์ ซึ่งเป็นข้อที่ทำหน้าที่ต่างกัน

วิธีการนี้มีข้อดีที่สามารถศึกษาผู้สอบหลายกลุ่มได้สะดวก แต่มีจุดอ่อนในเรื่องการควบคุมกลุ่มต่างๆ ให้มีความสามารถทัดเทียมกัน ขนาดกลุ่มตัวอย่างของกลุ่มต่าง ๆ และอัตราความคลาดเคลื่อนประเพณีที่ 1 จะสูงขึ้น ถ้าจำนวนข้อสอบเพิ่มขึ้น

2.2 การวัดความเบี่ยงเบนสัมพันธ์ (Reletive Deviation)

การคำนวณค่าความยากของข้อสอบ เช่น p, b เป็นต้น เมื่อคำนวณแยกระหว่างกลุ่มและแปลงให้เป็นค่ายากมาตรฐาน สามารถนำมาเปรียบเทียบเป็นรายข้อ ถ้าข้อใดเบี่ยงเบนไปมากหลักที่คาดหมาย หรือเบี่ยงเบนเกินจากความคลาดเคลื่อนมาตรฐานของค่าความยากที่กำหนด ย่อมแสดงถึงหน้าที่ต่างกันของข้อสอบ (Cleary & Kilton, 1968; Angoff & Ford, 1973) รวมทั้งสามารถคำนวณค่าสหพันธ์เข้าใกล้ 1.00 แสดงว่าความยากสัมพันธ์ของข้อสอบมีค่าใกล้เคียงกันระหว่างกลุ่ม ดังนั้นแบบสอบวัดคุณลักษณะคล้ายกันระหว่างกลุ่ม

วิธีการนี้มีข้อดีและข้อเสียคล้ายการทดสอบปฎิสัมพันธ์ นอกจานนี้ค่าความยากของข้อสอบ (P) มิใช่ตัวแทนของค่าความยากจริงของข้อสอบ และได้รับอิทธิพลจากค่าแทรกซ้อนอื่นได้แก่ ค่าอำนาจจำแนก และความสามารถของผู้สอบ

การเปรียบเทียบน้ำหนักตัวประกอบ (Factor Loading)

การวิเคราะห์ตัวประกอบ (Factor Analysis) เป็นเทคนิคทางสถิติที่นิยมใช้ในการตรวจสอบความตรงเชิงทฤษฎีหรือโครงสร้าง (Construct Validity) เมื่อนำการวิเคราะห์ตัวประกอบมาใช้ในการวิเคราะห์โครงสร้างของแบบสอบถามแยกตามกลุ่มผู้สอบ ความไม่สอดคล้องกันระหว่างน้ำหนักตัวประกอบคุณลักษณะสำคัญที่มุ่งวัด หรือ ความแตกต่างของค่าเฉลี่ยคะแนนตัวประกอบ (Factor Scores) ระหว่างกลุ่มผู้สอบ ย่อมสะท้อนการทำหน้าที่ต่างกันของข้อสอบและแบบสอบถาม

การใช้เทคนิคการวิเคราะห์ตัวประกอบเชิงสำรวจ (Exploratory Factor Analysis: EFA) สำหรับศึกษาการทำหน้าที่ต่างกัน จะมีจุดอ่อนในเรื่องความสอดคล้องระหว่างน้ำหนักตัวประกอบอาจเกิดความแตกต่างของความสามารถระหว่างกลุ่มก็ได้ แนวทางที่เหมาะสมจึงควรใช้เทคนิคการวิเคราะห์ตัวประกอบเชิงยืนยัน (Confirmatory Factor Analysis: CFA) นอกจากนี้ยังสามารถใช้ CFA สำหรับตรวจสอบความแตกต่างระหว่างกลุ่ม ในด้านคุณลักษณะหรือความสามารถหลักและความสามารถรองได้อีกด้วย (Camilli & Shepard, 1994)

การเปรียบเทียบโอกาสตอบข้อสอบถูก

การวิเคราะห์โอกาสตอบข้อสอบถูกของผู้สอบจากกลุ่มอ้างอิงและกลุ่มเปรียบเทียบที่มีความสามารถเท่ากัน เป็นแนวทางสำคัญที่นิยมใช้กันและเป็นที่ยอมรับในปัจจุบัน สำหรับปัจจัยการทำหน้าที่ต่างกันของข้อสอบ มีการคำนวณค่าสถิติ 2 แนวทาง ดังนี้

1) เปรียบเทียบค่าสัดส่วนหรือความนำžeเป็นการตอบข้อสอบถูกของผู้สอบต่างกลุ่มที่มีความสามารถเท่ากัน เช่น วิธีแมนเทล-แ xen สเซล เป็นต้น

2) เปรียบเทียบค่าฟังก์ชันการตอบสนองข้อสอบ หรือคोэฟฟิคิลล์สักขณะข้อสอบระหว่างกลุ่มที่มีระดับความสามารถเท่ากัน เป็นวิธีที่อยู่บนพื้นฐานของทฤษฎี IRT เช่น วิธีวัดความแตกต่างของพื้นที่ วิธีวัดความแตกต่างของค่าพารามิเตอร์ความยาก วิธีการทดสอบไป-แสควร์ของ ลอร์ด (Lord's χ^2 -test) เป็นต้น

วิธีการนี้มีข้อดีที่สำคัญ ได้แก่ การคำนวณค่าสถิติของข้อสอบมีความนำžeเชื่อถือมีกิลไก ควบคุมความสามารถของผู้สอบโดยการจับคู่กลุ่มความสามารถ เพื่อทำการเปรียบเทียบ ณ ตำแหน่งต่าง ๆ ที่มีความสามารถเท่ากัน จึงเป็นวิธีการที่ยอมกันทั่วไป แต่มีข้อจำกัดในด้านความสัลับซับซ้อนของแนวคิดพื้นฐาน และการวิเคราะห์มีความจำเป็นต้องใช้โปรแกรมคอมพิวเตอร์โดยเฉพาะ

วิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ (DIF Detection) จำแนก ตามลักษณะการตรวจให้คะแนนได้เป็น 2 ประเภท คือ ข้อสอบที่มีการให้คะแนนแบบทวิภาค หรือสองค่า (Dichotomous Scoring) และข้อสอบที่มีการให้คะแนนแบบพหุวิภาค หรือหลายค่า (Polytomous Scoring) วิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบแต่ละประเภท ยังสามารถจำแนกได้อีก 2 มิติ ได้แก่ มิติลักษณะของตัวแปรเกณฑ์ ซึ่งแบ่งเป็นกลุ่มวิธีที่ใช้คะแนนสังเกตได้ (Observed Score) และกลุ่มวิธีที่ใช้คะแนนสังเกตไม่ได้หรือคะแนนของตัวแปรแฝง (Latent variable) และมิติลักษณะของสถิติวิเคราะห์ ซึ่งแบ่งเป็นกลุ่มวิธีที่ใช้สถิติราพามตริก (Parametric Approach) รายชื่อวิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบที่สำคัญ ๆ ดังตารางที่ 2-1

ตารางที่ 2-1 วิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบที่มีการตรวจให้คะแนนแบบทวิภาค (Dichotomous DIF) และพหุวิภาค (Polytomous DIF) (ศิริชัย กาญจนวงศ์, 2555, หน้า 124)

ประเภทและตัวแปรเกณฑ์	พารามเมตริก	นันพารามเมตริก
1. DIF แบบทวิภาค		
1.1 คะแนนที่สังเกตได้ (observed score)	ANOVA Logistic Regression	TID MH STND
1.2 คุณลักษณะ/ตัวแปรแฝง (latent variable)	IRT-D ² Lord's χ^2 General IRTLR Loglinear IRTLR	SIBTEST
2. DIF แบบพหุวิภาค		
2.1 คะแนนที่สังเกตได้ (observed score)	ANOVA Polytomous Logistic Regression	Polytomous STND
2.2 คุณลักษณะ/ตัวแปรแฝง (latent variable)	General IRTLR PCM	GMH Polytomous SIBTEST GPCM

1). วิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบที่ให้คะแนนแบบทวิภาค

1.1 กลุ่มวิธีที่ใช้คะแนนที่สังเกตได้

วิธีการตรวจสอบคะแนนที่สังเกตได้แบบทวิภาคจะวิเคราะห์ตามทฤษฎีการทดสอบแบบดั้งเดิม (CTT) และกลุ่มที่ไม่ใช้ทฤษฎีการตอบสนองข้อสอบ (NON-IRT approach) โดยใช้คะแนนรวมของผู้เข้าสอบเป็นเกณฑ์การจับคู่ของกลุ่มผู้สอบ วิธีการตรวจสอบที่สำคัญในกลุ่มนี้ คือ

การวิเคราะห์ความแปรปรวน (ANOVA) (Cleary & Hilton, 1968)

วิธีการวิเคราะห์การทดสอบโดยโลจิสติก (Logistic Regression: LR)
(Swaminathan & Rogers, 1990)

วิธีแปลงค่าความยากของข้อสอบ (Transformed Item Difficulty: TID)
(Cleary & Hilton, 1968; Angoff & Ford, 1973)

วิธีแมนเทล-แฮนส์เซล (Mantel-Haenszel: MH)
(Holland & Thayer, 1988: 1989)

วิธีดัชนีมาตรฐาน (Standardization: STND) การปรับให้เป็นมาตรฐานด้วยน้ำหนักตัวประกอบ (Dorans & Kulick, 1986)

1.2 กลุ่มวิธีที่ใช้คุณลักษณะแฝง

วิธีการตรวจสอบคุณลักษณะ/ ตัวแปรแฝงแบบทวิภาค ซึ่งเป็นการวิเคราะห์อยู่บนพื้นฐานของทฤษฎีการตอบสนองข้อสอบ (IRT) สำหรับใช้เป็นเกณฑ์การจับคู่กลุ่มสอบ วิธีการตรวจสอบที่สำคัญในกลุ่มนี้ คือ

วิธีวัดพื้นที่ความแตกต่างระหว่างโค้งการตอบสนองข้อสอบ (IRT-D²) (Linn et al., 1981; Shepard et al., 1984; Raju, 1990; Kim & Cohen, 1991)

วิธีเค-สแควร์ของลอร์ด (Lord's χ^2) (Lord, 1980)

วิธีอัตราส่วนไลค์ลิขิตทั่วไป (General IRT Likelihood Ratio) (Thissen, Steinberg & Wainer, 1993)

วิธีอัตราส่วนไลค์ลิขิต โลกลินเนียร์ (Loglinear Likelihood Ratio) (Thissen, Steinberg & Wainer, 1993)

วิธีชิปเทสท์ (SIBTEST) (Shealy & Stout, 1993)

2.) วิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบที่ให้คะแนนแบบพหุวิภาค

2.1 กลุ่มวิธีที่ใช้คะแนนที่สังเกตได้

วิธีการวิเคราะห์ความแปรปรวน (ANOVA) (Cleary & Hilton, 1968)

วิธีการวิเคราะห์การคาดถอยโลจิสติกพหุวิภาค (Polytomous Logistic Regression) (Swaminathan & Rogers, 1990)

วิธีดัชนีมาตรฐานพหุวิภาค (Polytomous Standardization) (Dorans & Thayer, 1988, 1989)

วิธีแมนเทล-แฮนស์เซลทั่วไป (General Mantel-Haenzel; GMH) (Holland & Thayer, 1988, 1989)

2.2 กลุ่มวิธีที่ใช้คุณลักษณะแฝง

วิธีอัตราส่วนไลค์ลิขิตในรูปทั่วไป (General Mantel Likelihood Ratio) (Thissen, Steinberg & Wainer, 1993)

วิธีการให้คะแนนบางส่วน (Partial Credit Model: PCM) (Master, 1982)

วิธีชิปเทสท์พหุวิภาค (Polytomous SIBTEST) (Shealy & Stout, 1993)

วิธีการให้คะแนนบางส่วนทั่วไป (Generalized Partial Credit Model: GPCM) (Muraki, 1992, 1993) (ศรีษะ กาญจนวاسي, 2555, หน้า 116-126)

สรุป คือ การทำหน้าที่ต่างกันของข้อสอบ ตามที่ได้กล่าวมาแล้วสามารถที่สรุปได้ว่า การทำหน้าที่ต่างกันของข้อสอบ หมายถึง ผู้เข้าสอบข้อสอบที่ต่างกลุ่มกันที่มีความสามารถหรือ คุณลักษณะที่มุ่งวัดเท่ากัน ซึ่งทำให้มีโอกาสในการตอบข้อสอบที่ถูกได้แตกต่างกันออกไป หรือ มีฟังก์ชันการตอบสนองข้อสอบที่แตกต่างกัน และการทำหน้าที่ต่างกันของข้อสอบจะเกิดขึ้น เมื่อนำข้อสอบไปทดสอบกับผู้เข้าสอบกลุ่มย่อยที่ต่างกัน โดยจะมีความสามารถหลัก (Primary Ability) ในระดับเดียวกันหรือมีคุณลักษณะแฝง (Latent Trait) ที่จะวัดเท่ากัน และจะมี

ความสามารถรองลงมา (Secondary Ability) ที่แตกต่างซึ่งจะทำให้ผู้เข้าสอบที่ต่างกันมีอ่อนไหวเบริญเทียบกันจะทำให้โอกาสในการตอบข้อสอบนั้นแตกต่างกัน โดยการเบริญเทียบผลการตอบของข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบจะเป็นการเบริญเทียบกลุ่มผู้เข้าสอบ 2 กลุ่มขึ้นไปโดยประกอบด้วยกลุ่มแรกคือ กลุ่มเบริญเทียน (Focal Group หรือกลุ่ม F) เป็นกลุ่มที่นาสนใจที่จะศึกษาและคาดว่าจะเป็นกลุ่มที่จะเสียเบริญในตอบข้อสอบ ส่วนกลุ่มที่ 2 เป็นกลุ่มอ้างอิง (Reference Group หรือ กลุ่ม R) เป็นกลุ่มที่คาดว่าจะได้เบริญในการตอบข้อสอบได้ถูกต้องมากกว่ากลุ่มแรก

งานวิจัยที่เกี่ยวข้องกับการทำหน้าที่ต่างกันของข้อสอบ (DIF) มีดังนี้

พรญา สูงเนิน เสรี ชัดเช้ม และสมโภชน์ อนงสุข (2552) ได้เบริญเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบแบบพหุมิติ ระหว่างข้อสอบรายข้อกับหมวดข้อสอบ โดยวิธีซิปเทส์ ภายใต้เงื่อนไขขนาดของกลุ่มตัวอย่างที่ต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน) กลุ่มตัวอย่างเป็นนักเรียนชั้นประถมศึกษาปีที่ 6 สังกัดเขตพื้นที่การศึกษานครศรีธรรมราช ปีการศึกษา 2546 ที่เข้าสอบวัดผลสัมฤทธิ์ทางการเรียนระดับชาติ จำนวน 2,000 คน ข้อมูลทุติยภูมิที่ใช้เป็นคะแนนจากแบบทดสอบวิชาภาษาไทย ขั้นประถมศึกษาปีที่ 6 จำนวน 40 ข้อ จำแนกเป็น 2 หมวดข้อสอบ คือ หมวดที่ 1 วัดด้านโครงสร้างความรู้ จำนวน 15 ข้อ และหมวดที่ 2 วัดด้านกระบวนการ จำนวน 25 ข้อ วิเคราะห์ค่าสถิติพื้นฐานโดยใช้โปรแกรม SPSS ตรวจสอบความตรงเชิงโครงสร้างด้วยการวิเคราะห์องค์ประกอบเชิงยืนยัน อันดับสอง โดยใช้โปรแกรม LISREL 8.50 และตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ โดยใช้โปรแกรม SIBTEST ผลการวิจัยปรากฏว่า ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก ขนาดกลาง และขนาดใหญ่ ระหว่างการตรวจสอบการทำหน้าที่ต่างกันเป็น รายข้อกับรายหมวด ข้อสอบ พบที่ต้องตอบหน้าที่ต่างกันแตกต่างกันการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบเป็นรายข้อ เมื่อกลุ่มตัวอย่างขนาดเล็ก พบที่ต้องตอบหน้าที่ต่างกันจำนวน 4 ข้อ คิดเป็นร้อยละ 10 ขนาดกลางพบที่ต้องตอบหน้าที่ต่างกันจำนวน 13 ข้อ คิดเป็นร้อยละ 32.5 และขนาดใหญ่พบที่ต้องตอบหน้าที่ต่างกันจำนวน 15 ข้อ คิดเป็นร้อยละ 37.5 การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบเป็นรายหมวดข้อสอบ เมื่อกลุ่มตัวอย่างขนาดเล็ก พบที่ต้องตอบหน้าที่ต่างกันจำนวน 4 ข้อ คิดเป็นร้อยละ 10 ขนาดกลางพบที่ต้องตอบหน้าที่ต่างกันจำนวน 8 ข้อ คิดเป็นร้อยละ 20 และขนาดใหญ่พบที่ต้องตอบหน้าที่ต่างกันจำนวน 16 ข้อ คิดเป็นร้อยละ 40 การตรวจสอบการทำหน้าที่ต่างกันของหมวดข้อสอบ ชี้ให้เห็นว่า หมวดที่ 2 ภายใต้เงื่อนไขกลุ่มตัวอย่างขนาดกลาง มีนัยสำคัญทางสถิติที่ระดับ .05

ณรงค์ จันทร์มหา (2554) ได้เบริญเทียบค่าความเที่ยงของแบบทดสอบผลสัมฤทธิ์ทางการเรียนที่มีจำนวนข้อสอบทำหน้าที่ต่างกันแตกต่างกัน 7 เงื่อนไข คือ 0% 5% 10% 15% 20% 25% และ 30% กลุ่มตัวอย่างเป็นนักเรียนชั้นประถมศึกษาปีที่ 6 สังกัดสำนักงานคณะกรรมการการศึกษาขั้นพื้นฐาน ปีการศึกษา 2546 จำนวน 2,000 คน ที่ทำแบบทดสอบวัดผลสัมฤทธิ์ทางการเรียนระดับชาติ วิชาภาษาไทย จำนวน 40 ข้อ วิเคราะห์ค่าสถิติพื้นฐานโดยใช้โปรแกรม SPSS ตรวจสอบการทำหน้าที่ต่างกันของข้อสอบโดยใช้โปรแกรมซิปเทส์ และทดสอบความแตกต่างของค่าความเที่ยงของแบบทดสอบ โดยใช้การทดสอบสถิติชี ผลการวิจัย พบว่า ค่าความเที่ยงของแบบทดสอบที่มีจำนวนข้อสอบที่ต่างกัน 0% 5% 10% 15% 20% 25% และ 30% ไม่แตกต่างกัน

ศิริรัตน์ สุคันธพฤกษ์ โขติกา ภารี และศิริชัย กัญจนวานิช (2554) ได้วิเคราะห์ข้อคำถามในแบบวัดความวิตกกังวลในการสอบคณิตศาสตร์การเปรียบเทียบระหว่างไทรัชิกอลินเนียร์โนมเดล พาเชียลเครดิตโมเดล และเกรตเดรสพอนโนมเดล ซึ่งมีการตรวจสอบการทำหน้าที่ต่างกันในแบบวัดความวิตกกังวล ในการสอบคณิตศาสตร์ โดยเปรียบเทียบระหว่างไทรัชิกอลินเนียร์โนมเดล พาเชียล เครดิตโมเดลและเกรตเดรสพอนโนมเดลโดยกลุ่มตัวอย่างที่ใช้ในการวิจัยเป็นนักเรียนมัธยมศึกษาปีที่ 6 สายวิทย์-คณิต ปีการศึกษา 2552 จำนวน 1,715 คน จาก 29 โรงเรียนในสังกัดสำนักงานเขตพื้นที่การศึกษาพระนครศรีอยุธยาเขต 1 และเขต 2 สำนักงานเขตพื้นที่การศึกษาอ่างทองและสำนักงานเขตพื้นที่การศึกษานนทบุรี ซึ่งได้มาจากการสุ่มตัวอย่างแบบแบยกัน เครื่องมือที่ใช้ในการวิจัย คือ แบบวัดความวิตกกังวลในการสอบคณิตศาสตร์ การวิเคราะห์ข้อมูล 3 ขั้นตอน คือ วิเคราะห์การทำหน้าที่ต่างกันของข้อสอบด้วย Hierarchical Linear Model (HLM) โดยใช้โปรแกรม HLM การวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบด้วย Partial Credit Model (PCM) และ Graded Response Model: GRM ด้วยโปรแกรม PRASCALE และเปรียบเทียบผลการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบด้วย HLM โดยใช้โปรแกรม HLM กับ PCM กับ GRM ด้วยโปรแกรม PRASCALE ผลการวิจัย พบว่า ข้อคำถามที่ทำหน้าที่ต่างกันของข้อร่วมระหว่าง HLM, PCM และ GRM มี 6 ข้อ จาก 39 ข้อ คิดเป็นร้อยละ 15.38 ข้อคำถามที่ทำหน้าที่ต่างกันของข้อร่วมระหว่าง HLM กับ PCM มี 7 ข้อ จาก 39 ข้อ คิดเป็นร้อยละ 17.94 และข้อคำถามที่ทำหน้าที่ต่างกันของร่วมระหว่าง HLM กับ GRM มี 9 ข้อ จาก 39 ข้อ คิดเป็นร้อยละ 23.07

ชัยวัฒน์ หาทัยพันธ์ (2558) ได้พัฒนาวิธีการสำหรับการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบโดยผู้เชี่ยวชาญ และประการที่สองเพื่อเปรียบเทียบประสิทธิภาพการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในด้านอัตราความถูกต้อง และอัตราความคลาดเคลื่อนของผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบกับข้อสอบที่คัดสรรมาสำหรับ การทดลองเมื่อใช้วิธีการวิเคราะห์ด้วยแบบวินิจฉัยโดยผู้เชี่ยวชาญ วิธีการประยุกต์ใช้เทคนิคการประชุม แบบเดลฟายจากกลุ่มผู้เชี่ยวชาญและวิธีการประยุกต์ใช้เทคนิคໂປຣໂຕຄລະລາວົດ ตัวอย่างที่ใช้ในการวิจัย คือ ผู้เชี่ยวชาญจำนวน 21 คนและนักเรียนระดับมัธยมศึกษาปีที่ 6 ปีการศึกษา 2556 จำนวน 139 คน ซึ่งได้จาก การเลือกตัวอย่างแบบเจาะจง เครื่องมือที่ใช้ในการวิจัยประกอบด้วยแบบวินิจฉัยการทำหน้าที่ต่างกันของข้อสอบจากผู้เชี่ยวชาญ แบบยืนยันการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ โดยประยุกต์เทคนิคการประชุมแบบเดลฟาย แบบสอบถามสำหรับการตรวจสอบความลำเอียงของข้อสอบสำหรับนักเรียน ชุดข้อสอบสาระการเรียนรู้สุขศึกษาและพลศึกษาสำหรับการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบสำหรับผู้เชี่ยวชาญ ข้อสอบที่คัดสรรมาได้นำมาผ่านการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบโดยวิธีແນ່ເລ-ແຍນສ-ເຊລ ด้วยโปรแกรม DDFS 1.0 และโปรแกรม DIFAS 5.0 พบว่า วิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบโดยการตัดสินของผู้เชี่ยวชาญที่สำคัญมี 3 วิธี ได้แก่ วิธีที่ 1 การวินิจฉัยการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบโดยผู้เชี่ยวชาญ วิธีที่ 2 การประยุกต์ใช้เทคนิคการประชุมแบบเดลฟายจากกลุ่มผู้เชี่ยวชาญ และวิธีที่ 3 การประยุกต์ใช้เทคนิคໂປຣໂຕຄລະລາວົດ 2) ข้อสอบที่ทำหน้าที่ต่างกันด้านเพศของแบบสอบถามสาระการเรียนรู้สุขศึกษาและพลศึกษา จากผลการวิเคราะห์เปรียบเทียบอัตราความถูกต้องระหว่างวิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบด้วยวิธีทางสถิติกับการตรวจสอบโดยผู้เชี่ยวชาญ พบว่า วิธีที่ 1 การตรวจสอบด้วยแบบ

วินิจฉัยโดยผู้เชี่ยวชาญมีอัตราความถูกต้องโดยเฉลี่ยคิดเป็น ร้อยละ 50 และมีอัตราความคลาดเคลื่อนของการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบโดยเฉลี่ยคิดเป็นร้อยละ 50 วิธีที่ 2 การประยุกต์ใช้เทคนิคการประชุมแบบเดลฟายจากกลุ่มผู้เชี่ยวชาญ มีอัตราความถูกต้องตามฉบับตามติกกลุ่มเชี่ยวชาญ คิดเป็นร้อยละ 0 และมีอัตราความคลาดเคลื่อนของการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ คิดเป็นร้อยละ 100 วิธีที่ 3 การประยุกต์ใช้เทคนิคໂປຣໂຕຄອລະລາວด้วยมีอัตราความถูกต้องโดยเฉลี่ยคิดเป็นร้อยละ 25 และมีอัตราความคลาดเคลื่อนของการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบโดยเฉลี่ยคิดเป็นร้อยละ 75

สุราทิพย์ ตรีสิน และปิยะทิพย์ ประดุจพร (2560) ได้เปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบ ระดับชาติ ด้านภาษา ด้านคำนวน และด้านเหตุผล ชั้นประถมศึกษาปีที่ 3 ด้วยวิธี HGLM วิธี MIMIC และวิธี IRT-LR โดยมีวัตถุประสงค์ เพื่อวิเคราะห์คุณภาพของแบบทดสอบระดับชาติ (NT) และตรวจสอบการทำหน้าที่ต่างกัน ของข้อสอบในแบบทดสอบระดับชาติชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน และด้านเหตุผล ด้วยวิธี HGLM วิธี MIMIC และวิธี IRT-LR ดังนี้ 1) วิเคราะห์คุณภาพของแบบทดสอบระดับชาติ ทั้ง 3 ด้าน 2) ตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ด้วยวิธี HGLM วิธี MIMIC และวิธี IRT-LR และ 3) เปรียบเทียบผลการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบ ด้วยวิธีการตรวจสอบ 3 วิธีข้อมูลที่นำมาใช้ วิเคราะห์เป็นข้อมูลทุติยภูมิ จากผลการตอบแบบทดสอบ ระดับชาติของนักเรียนชั้นประถมศึกษาปีที่ 3 ปีการศึกษา 2556 จำนวน 9,600 คน ผลการวิจัย ปรากฏว่า 1) แบบทดสอบระดับชาติชั้นประถมศึกษาปีที่ 3 มีค่าความยากของข้อสอบ (b) อยู่ใน ระดับค่อนข้างยาก มีค่าอำนาจจำแนกของข้อสอบ (a) อยู่ในระดับที่สามารถจำแนกผู้สอบได้ดี และ มีค่าโอกาสในการเดาของข้อสอบ (c) ไม่เกิน 0.3 2) การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ ทั้ง 3 ด้าน ซึ่งให้เห็นว่า เพศส่งผลให้เกิดการทำหน้าที่ต่าง กันของข้อสอบ โดยเพศหญิงจะได้เปรียบ ในการตอบข้อสอบด้านภาษา และด้านเหตุผล ในขณะที่เพศชาย จะได้เปรียบในการตอบข้อสอบ ด้านคำนวน โดยวิธี HGLM ตรวจพบข้อสอบทำหน้าที่ต่างกัน จำนวนมากที่สุด คิดเป็นร้อยละ 69 ของข้อสอบทั้งฉบับ รองลงมาคือ วิธี IRT-LR ร้อยละ 54 และวิธี MIMIC ร้อยละ 16 ตามลำดับ 3) การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ พบร่วม วิธี HGLM ตรวจพบ DIF มากกว่า วิธี MIMIC ในด้านภาษา ด้านคำนวน และด้านเหตุผล คิดเป็นร้อยละ 70, 36 และ 53 ตามลำดับ และวิธี HGLM ตรวจ พบร่วม DIF มากกว่าวิธี IRT-LR ด้านภาษา และด้านคำนวน คิดเป็นร้อยละ 37 และ 13 และวิธี IRT-LR ตรวจพบ DIF มากกว่าวิธี MIMIC ทั้ง 3 ด้าน คิดเป็นร้อยละ 33, 43 และ 40 ตามลำดับ ส่วนวิธี HGLM ตรวจพบ DIF น้อยกว่า วิธี IRT-LR ด้านคำนวน คิดเป็นร้อยละ 7 ($p < .05$)

Mendes-Barnett and Ercikan (2006) ได้ศึกษาการทำหน้าที่ต่างกันของข้อสอบ ในการสอบวิชาคณิตศาสตร์ โดยใช้วิธีชิบเพรสท์ โดยใช้ตัวแปรเพศ ผลการวิจัยพบว่า เพศชาย มีความสามารถในการแก้ปัญหา และวิธีการทางปัญญาที่สูงเป็นอย่างมากที่จัดไว้ในข้อสอบ ขณะที่เพศหญิงมีความสามารถในด้านการคำนวนสมการ ซึ่งการคำนวนไม่ได้ถูกจัดให้อยู่ในข้อสอบ จึงสรุปได้ว่า ข้อสอบวิชาคณิตศาสตร์นี้เกิดการทำหน้าที่ต่างกัน โดยลำเอียงเข้าทางเพศชายมากกว่า เพศหญิง

Breland and Lee (2007) ได้ศึกษาการทำหน้าที่ต่างกันของข้อสอบของการทดสอบปรับเหมาะสมด้วยคอมพิวเตอร์ที่มีภาษาอังกฤษและภาษาต่างประเทศ ในรายการข้อสอบของ TOEFL-CBT จากกรณีตัวอย่างจำนวน 5,660 ซึ่งใช้วิเคราะห์สมการคัดถ่ายโลจิสติก สำหรับข้อสอบที่มีการให้คะแนนแบบหลายค่า โดยวิเคราะห์จากตัวแปรที่แยกตามประเภทของข้อสอบคือ ข้อสอบที่ทำหน้าที่ต่างกันแบบเอกสาร (Uniform DIF) กับข้อสอบที่ทำหน้าที่ต่างกันแบบอนุกรูป (Non-uniform DIF) และตัวแปรเพศ จากการวิจัยพบว่าข้อสอบทำหน้าที่ต่างกันอย่างอนุกรูป มีการทำหน้าที่ต่างกันของข้อสอบการเรียงความ (Writing) ที่เพศหญิงจะได้เปรียบ จึงสามารถสรุป ได้ว่าข้อสอบ TOEFL-CBT มีการทำหน้าที่ต่างกันของข้อสอบด้านตัวแปรเพศยังมีน้อย

Barnes and Wells (2009) ได้ศึกษาการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบของโครงการ National Doctoral Program Survey (NDPS) ระหว่างตัวแปรเพศกับเชื้อชาติ ทางด้านความคิด โดยผลการวิจัยพบว่า ข้อคำถามจำนวน 29 จาก 48 ข้อ พน DIF โดยที่เพศหญิง ที่มีผิวสี มีโอกาสมากหรือน้อยที่จะมีความคิดตรงกับเพื่อนเพศชายชาวผิวขาวของเข้า และควรให้มีการระมัดระวังในการใช้ข้อคำถามสำหรับกลุ่มผู้เรียนที่มีความหลากหลาย เช่น เพศ และเชื้อชาติ จึงสามารถสรุปได้ว่า เพศ และเชื้อชาติมีผลต่อความคิด และการตอบคำถามของผู้เรียนที่ศึกษาในสถานศึกษาที่มีลักษณะที่แตกต่างกัน

Le (2009) ได้ศึกษาการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบระหว่างเพศ การอยู่นอกพื้นที่และการทดสอบภาษาในการทดสอบด้านวิทยาศาสตร์ ใน PISA โดยใช้ข้อมูล การทดลองภาคสนามของ PISA รอบ 3 เพื่อตรวจสอบความสัมพันธ์ระหว่าง การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบระหว่างเพศ การอยู่นอกพื้นที่และการทดสอบภาษาในการทดสอบด้านวิทยาศาสตร์และรูปแบบอื่น ๆ ที่กำหนดไว้ใน PISA โดยมุ่งเน้นที่บริบทความสามารถและ ความรู้ทางวิทยาศาสตร์ข้อมูลที่ใช้ได้รับการรวบรวมจาก 60 กลุ่มภาษาทดสอบโดย 50 ประเทศ ที่เข้ารวมด้วยรวมประมาณ 83,000 คน ที่เป็นนักเรียนอายุ 15 ปี งานวิจัยนี้ใช้วิธี IRT เพื่อตรวจสอบการทำหน้าที่ต่างกันของข้อสอบระหว่างเพศสำหรับแต่ละกลุ่มภาษาและทั่วโลก ข้อสอบแบบเลือกตอบและแบบตอบกลับ แบบปิดมีแนวโน้มที่จะเข้าข้างเพศชาย การศึกษายังแสดงให้เห็นถึง ผลกระทบของประเทศไทยและการทดสอบภาษา กับเพศ การทำหน้าที่ต่างกันของข้อสอบในข้อมูล ระหว่างประเทศ ผลการวิจัย ปรากฏว่ามีคุณค่าการมีส่วนร่วมในการพัฒนาการทดสอบเพื่อการใช้งาน ระหว่างประเทศ

Taylor and Lee (2012) ได้ศึกษาการทำหน้าที่ต่างกันของข้อสอบ (DIF) สำหรับการอ่าน ระดับชั้นประถมศึกษาปีที่ 4 มัธยมศึกษาปีที่ 1 ของข้อสอบคณิตศาสตร์จากการทดสอบตามเกณฑ์ ของรัฐ การทดสอบประกอบด้วย ข้อสอบที่มีตัวเลือกหลายตัวเลือกและสร้างการตอบสนอง เพื่อตรวจสอบ DIF เมื่อจำแนกตามเพศ โดยใช้ Poly-SIBTEST และขั้นตอน Rasch โดยขั้นตอน Rasch ถูกตั้งค่าสถานะไว้เพิ่มเติมสำหรับข้อสอบที่ DIF มากกว่าที่ทำในข้อสอบพร้อมกันโดยเฉพาะ 20 ข้อสอบที่มีหลายตัวเลือก สำหรับการอ่านและการทดสอบทางคณิตศาสตร์ทั้งสองแบบ ข้อสอบที่เพศชาย ได้เปรียบในขณะที่รายการตอบสนองการสร้างที่เข้าข้างเพศหญิง การวิเคราะห์เนื้อหา แสดงให้เห็นว่ามีการอ่านค่า อ่านหนังสือที่ถูกตั้งค่าสถานะไว้ การตีความข้อความหรือความหมาย โดยนัย เพศชายมีแนวโน้มที่จะได้เปรียบจากสิ่งต่าง ๆ จึงมีการระบุการตีความและการวิเคราะห์

ข้อความที่ให้ข้อมูลที่เหมาะสม ข้อสอบที่เข้าข้างเพศหญิงขอให้นักเรียนทำการตีความของตัวเองและวิเคราะห์ทั้งข้อความวรรณกรรมและข้อมูลโดยได้รับการสนับสนุนจากหลักฐานจากข้อความ การวิเคราะห์เนื้อหาของรายการคณิตศาสตร์แสดงให้เห็นว่าข้อสอบที่เข้าข้างเพศชาย คือ เรขาคณิต ความน่าจะเป็นและพีชคณิต ข้อสอบคณิตศาสตร์ที่เข้าข้างเพศหญิง การตีความทางสถิติ การแก้ปัญหาหลายขั้นตอนและการให้เหตุผลเชิงคณิตศาสตร์

สรุปจากการศึกษางานวิจัยที่เกี่ยวข้องการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ เป็นผลมาจากการตัดสินว่าข้อสอบนั้นมีความยุติธรรมหรือไม่อย่างไร ซึ่งเป็นข้อมูลสารสนเทศทางสถิติ มีความสัมพันธ์ที่ส่งผลไปยังคุณลักษณะที่ข้อสอบมุ่งที่จะวัด กับมวลประสบการณ์ของผู้สอบในกลุ่มต่างๆ ที่มีบางอย่างที่แตกต่างกัน เช่น เชื้อชาติ ศาสนา วัฒนธรรม ภูมิลำเนา สังคม เพศ ภาษา อายุ และประสบการณ์เป็นต้น การทำหน้าที่ต่างกันของข้อสอบของผู้ที่เข้าสอบจะมีสัดส่วนในการตอบคำถามข้อสอบได้ถูกต้องไม่เท่ากันของแต่ละกลุ่มประชากรที่จะศึกษา และเมื่อกลุ่มผู้เข้าสอบได้คะแนนที่เท่ากันนั้นข้อสอบก็เป็นเอกพันธ์ ซึ่งการทำหน้าที่ต่างกันของข้อสอบ จะเป็นแบบพหุมิติ ของข้อสอบที่จะวัด

ตอนที่ 4 การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบด้วยวิธีการทดสอบอัตราส่วน ไลค์ลิขิต และงานวิจัยที่เกี่ยวข้อง

ทฤษฎีการตอบสนองข้อสอบ IRT นี้ สามารถจำแนกได้เป็น 2 ประเภท ได้แก่ ทฤษฎีการตอบสนองข้อสอบแบบตรูจิให้ค่า 2 ค่า (Binary or Dichotomous IRT) ซึ่งเป็นโมเดลการตอบสนองข้อสอบที่ใช้กับการตรวจสอบคุณภาพรายข้อแบบ 2 ค่า เช่น ข้อสอบหรือข้อคำถามที่ตรวจให้ค่าแบบ 0,1 (ตอบผิดได้ 0, ตอบถูกได้ 1) แบบถูก/ผิด ใช่/ไม่ใช่ เป็นต้น และทฤษฎี การตอบสนองข้อสอบแบบตรูจิให้ค่ามากกว่า 2 ค่า (Polytomous IRT) ซึ่งเป็นโมเดลการตอบสอบ ข้อสอบที่ใช้กับการตรวจคุณภาพรายข้อแบบมากกว่า 2 ค่า เช่น ข้อสอบหรือข้อคำถามมาตรฐานมาตราประมาณค่า (Rating Scale) การตรวจข้อสอบแบบให้ค่าความรู้บางส่วน (Partial Credit) เป็นต้น (ศิริชัย กาญจนวารี, 2555, หน้า 51)

โมเดลการตอบสนองข้อสอบ (IRT Models) เดิมมีแนวคิดพื้นฐานมาจากทฤษฎีการวัดแบบดั้งเดิม (Classical Test Theory) โดยโมเดลการตอบสนองข้อสอบจะวัดความสัมพันธ์ระหว่างตัวแปรอิสระ ที่ประกอบด้วยตัวแปรคง ดังนี้ ความสามารถของผู้สอบ (θ) หรือค่าพารามิเตอร์ของผู้สอบ (a, b, c) และตัวแปรอิสระที่เป็นตัวแปรที่สังเกตได้นั่น คือ โอกาสในการตอบข้อสอบได้ถูกต้องของผู้สอบซึ่งทฤษฎีการตอบสนองข้อสอบ (IRT) เป็นการอธิบายความสัมพันธ์ระหว่างความสามารถภายในของตัวบุคคล (Latent Trait or Ability) กับผลของการตอบข้อสอบถูกหรือค้างคุณลักษณะข้อสอบ (ItemCharacteristic Curve: ICC) ที่มีการกำหนดลักษณะข้อสอบด้วยพารามิเตอร์ของข้อสอบ คือค่าอำนาจจำแนกของข้อสอบ (a) ค่าความยากของข้อสอบ (b) และค่าโอกาสในการเดาของข้อสอบได้ถูก (c) เพราะเหตุนี้ IRT จึงอยู่บนพื้นฐานของความคิด 2 ประการ คือ 1) ผลการตอบข้อสอบได้ถูกต้องของผู้สอบ ที่อธิบายความสามารถที่มีอยู่ภายในตัวผู้สอบได้ และ 2) ความสัมพันธ์ระหว่างผลการตอบข้อสอบกับความสามารถที่อยู่ภายในตัวของผู้สอบ ซึ่งอธิบาย

ได้ด้วยโคงคุณลักษณะข้อสอบ (ICC) และ พิงค์ชันของลักษณะข้อสอบ (Embretson & Resie, 2000, p. 8)

แนวคิดทฤษฎีการตอบสนองข้อสอบ

ทฤษฎีการทดสอบเริ่มแรกมีการพัฒนาจากผลงาน และความพยายามของนักจิตวิทยา ทั้งในยุโรปและอเมริกา นักจิตวิทยาได้หัวเรื่องการแก้ปัญหาการวัดต่าง ๆ เพื่อพัฒนาศาสตร์แห่งการวัด และตรวจสอบจนมีความมั่นคง ในศตวรรษที่ 12 เมื่อมีเริ่มผลิตกระดาษขึ้นใช้แทนการสอบปากเปล่า และปี พ.ศ. 2293 มหาวิทยาลัยเคมบริดจ์ ได้เริ่มใช้การสอบด้วยกระดาษอย่างเป็นทางการเพื่อวัดผล การเรียนของนักศึกษา มหาวิทยาลัยออกซ์ฟอร์ดได้ใช้ข้อสอบข้อเขียนสำหรับวัดจำแนกความสามารถ ของนักศึกษาเพื่อวัดผลกระทบดับผ่านหรือระดับเกียรตินิยม

แนวคิดพื้นฐานของทฤษฎีการตอบสนองข้อสอบ

ทฤษฎีนี้เกิดขึ้นโดยมีข้อจำกัดของทฤษฎีการทดสอบแบบดั้งเดิมหลายประการ คือ (Hambleton, Swaminathan, & Rogers, 1991, pp. 7-12)

1. ค่าสถิติของข้อสอบ เช่น ความยากจะขึ้นอยู่กับลักษณะของกลุ่มผู้สอบ กล่าวคือ ถ้าผู้สอบมีความสามารถสูง ข้อสอบจะกล้ายเป็นข้อสอบที่ง่าย แต่ถ้าผู้สอบมีความสามารถต่ำข้อสอบ ดังกล่าวจะกล้ายเป็นข้อสอบที่ยาก ส่วนค่าอำนาจจำแนกของข้อสอบขึ้นอยู่กับความเป็นเอก พันธ์ของความสามารถของผู้สอบ ถ้าผู้สอบมีความสามารถแตกต่างกันมากข้อสอบก็จะมีค่าอำนาจจำแนกของข้อสอบสูงซึ่งมีผลทำให้ความเที่ยงของแบบทดสอบมีค่าสูงตามไปด้วย เนื่องจากความเที่ยง ของแบบทดสอบมีความสัมพันธ์ทางบวกกับค่าความแปรปรวนของคะแนนจากการทดสอบ

2. การเปรียบเทียบความสามารถของผู้สอบบัน្ត จะต้องใช้แบบทดสอบฉบับเดียวกันหรือแบบทดสอบคู่ขนาน ปัญหาที่เกิดขึ้น คือ แบบทดสอบวัดผลสัมฤทธิ์ และแบบทดสอบวัดความถนัดนั้น ส่วนใหญ่แล้วจะเหมาะสมกับผู้ที่มีความสามารถปานกลาง ดังนั้น ความถูกต้อง แม่นยำของ การวัด ผู้สอบที่มีความสามารถสูงและผู้สอบที่มีความสามารถต่ำจึงลดลง

3. ค่าความเที่ยงของแบบทดสอบถูกนิยามในรูปของผลที่ได้จากการใช้แบบทดสอบคู่ขนาน ซึ่งในทางปฏิบัติจริงนั้นนับว่าเป็นเรื่องยากที่จะให้การสอบ 2 ครั้ง มีสภาพที่เหมือนกัน ถึงแม้ว่าแบบทดสอบคู่ขนานจะนานกันจริง แต่ผู้สอบอาจจะมีลักษณะที่แตกต่างกันไปจากการสอบครั้งแรกเกี่ยวกับแรงจูงใจ ความกังวล การลีม หรือการพัฒนาตนเองในบางทักษะ เป็นต้น

4. ทฤษฎีการทดสอบแบบดั้งเดิมไม่สามารถบอกได้ว่าผู้สอบจะตอบข้อสอบอย่างไร ยกเว้นแต่ว่าจะได้ใช้ข้อสอบข้อนั้นกับผู้สอบที่มีลักษณะคล้ายคลึงกันมาแล้ว

5. ทฤษฎีการทดสอบแบบดั้งเดิมใช้ค่าความแปรปรวนของความคลาดเคลื่อนในการวัด (Variance of Error of Measurement) ให้มีน้อยกับผู้สอบทุกคน ซึ่งตามความเป็นจริงแล้ว ผู้สอบที่มีความสามารถสูงและต่ำ จะมีค่าความแปรปรวนของความคลาดเคลื่อนในการวัดต่างจาก ผู้สอบที่มีความสามารถปานกลาง

หลักการของทฤษฎีการตอบสนองข้อสอบ

ทฤษฎีการตอบสนองข้อสอบ (Item Response Theory) มีความเชื่อเกี่ยวกับ ค่าพารามิเตอร์ของข้อสอบ (Item Parameter) คือ ค่าความยากของข้อสอบ (b) ค่าอำนาจจำแนก

ของข้อสอบ (a) ค่าโอกาสในการเดาข้อสอบ (c) ของข้อสอบแต่ละข้อว่าเป็นคุณลักษณะที่คิดที่ในตัวข้อสอบนั้น เพราะฉะนั้นค่าพารามิเตอร์เหล่านี้จึงไม่ควรเปลี่ยนไปตามกลุ่มผู้สอบ และในทำนองเดียวกัน ค่าพารามิเตอร์ของผู้สอบ (Person Parameter) หรือความสามารถที่แท้จริงของผู้สอบก็เป็นคุณลักษณะที่มีอยู่ภายในตัวผู้สอบจึงไม่ควรเปลี่ยนไปตามชุดข้อสอบที่เลือกใช้แต่เนื่องจากความสามารถของผู้สอบเป็นคุณลักษณะแฟงไม่สามารถสังเกต หรือวัดได้โดยตรง (Unobservable) จึงจำเป็นต้องใช้การทำนาย (Predict) หรืออธิบาย (Explain) คุณลักษณะดังกล่าวโดยอาศัยผลที่ได้จากการตอบแบบทดสอบ (Test Performance) หรือคะแนน (Score) ซึ่งเป็นสิ่งที่สามารถสังเกตและวัดได้ (Observable) นักวัดผลจึงได้พยายามหาความสัมพันธ์ระหว่างผลที่ได้จากการตอบแบบทดสอบหรือคะแนน (Test Performance or Score) กับระดับความสามารถ (Ability) ของผู้ตอบแต่ละคน เพื่อเขียนเป็นโมเดลทางคณิตศาสตร์ (Mathematical Model) ความสัมพันธ์ระหว่างผลที่ได้จากการตอบแบบทดสอบกับระดับความสามารถของผู้สอบ สามารถเขียนในรูปของความสัมพันธ์ ได้ดังนี้

$$P=f(U_i/\theta_1, \theta_2, \theta_3, \dots, \theta_k; \beta_k) \quad (5)$$

เมื่อ P แทน ผลการตอบแบบทดสอบ (Test Performance)

f แทน พังก์ชัน (Function)

U_i แทน ผลการตอบแบบทดสอบข้อที่ i (ตอบถูก = 1, ตอบผิด = 0)

$\theta_1, \theta_2, \theta_3, \dots, \theta_k$ แทน ระดับความสามารถ (Ability) ที่ 1, 2, 3, ..., k

β_k แทน ค่าพารามิเตอร์ของข้อสอบข้อที่ j

เนื่องจากความสัมพันธ์ดังกล่าวเป็นเพียงพังก์ชันความสัมพันธ์ในลักษณะทั่วไป นักวัดผลการศึกษาจึงต้องหาโมเดลทางคณิตศาสตร์ที่เหมาะสม เพื่อใช้แทนพังก์ชันความสัมพันธ์ ดังกล่าว โดยอาศัยข้อตกลงเบื้องต้นของทฤษฎีการตอบสนองข้อสอบ

ข้อตกลงเบื้องต้นของทฤษฎีการตอบสนองข้อสอบ

ข้อตกลงเบื้องต้นของทฤษฎีการตอบสนองข้อสอบที่สำคัญ คือ (Hambleton et al.,

1991, pp. 9-12, ศิริชัย กาญจนวารี, 2555, หน้า 75-76)

1. ความเป็นเอกมิตร (Unidimensional) การวัดตามแนวทฤษฎีการตอบสนองข้อสอบได้ระบุความเป็นเอกมิตรของคุณลักษณะในข้อตกลงเบื้องต้นว่า ข้อคำถาม หรือข้อสอบทุกข้อในแบบทดสอบนั้นจะต้องมุ่งวัดความสามารถ หรือคุณลักษณะเดียวเท่านั้น โดยที่ร่วงแล้วข้อตกลงข้อนี้ เป็นไปได้ค่อนข้างยากเนื่องจากมีปัจจัยที่มีผลต่อคะแนนสอบ เช่น ปัจจัยด้านความรู้ความเข้าใจ (Cognitive) บุคลิกภาพและปัจจัยเกี่ยวกับการจัดการสอบ ปัจจัยเหล่านี้อาจรวมถึงแรงจูงใจ ความวิตกกังวลในการสอบความสามารถในการทำงานได้รวดเร็ว ความรู้สึกภัยกับการใช้กระดาษคำตอบเมื่อเป็นเช่นนี้สิ่งที่ทำให้ข้อตกลงนี้เป็นไปได้ คือ การพิจารณาว่าแบบทดสอบฉบับนั้นมีองค์ประกอบใดหรือปัจจัยใดที่เด่นที่สุด ก็ถือว่าแบบทดสอบได้วัดในสิ่งนั้นการตรวจสอบความเป็นเอกมิตรของคุณลักษณะที่ใช้ในการทดสอบมีวิธีการตรวจสอบได้หลายวิธี ดังนี้

1.1 การหาค่าความสัมพันธ์ระหว่างค่าน้ำหนักองค์ประกอบรายข้อ (Factor Loading) ขององค์ประกอบที่หนึ่งกับค่าสหสัมพันธ์แบบบิเซเรียล (Biserial Correlation Coefficient) ของข้อสอบรายข้อกับคะแนนรวม ถ้ามีค่าสัมประสิทธิ์สหสัมพันธ์มากกว่า .80 ทำให้สามารถสรุปได้ว่า ข้อสอบหรือแบบสอบถามฉบับนั้นมีความเป็นเอกมิตรของคุณลักษณะที่ใช้ในการทดสอบ

1.2 การวิเคราะห์องค์ประกอบ (Factor Analysis) ของข้อสอบทั้งฉบับ พิจารณา ได้จากค่าไอegen (Eigen Value) โดยผลการวิเคราะห์องค์ประกอบใหม่มีค่าไอegenในองค์ประกอบใด องค์ประกอบหนึ่งสูงกว่าค่าอื่นอย่างชัดเจน สามารถสรุปได้ว่าข้อสอบหรือแบบสอบถามฉบับนั้น มีความเป็นเอกมิตรของคุณลักษณะในการทดสอบ

1.3 การใช้โปรแกรม TESTFACT ในการพิจารณาความเป็นมิติของแบบสอบถาม จากการวิเคราะห์องค์ประกอบเชิงยืนยัน (Confirmatory Factor Analysis) ซึ่งพัฒนาขึ้น โดย (Wilson & Hoskens, 1991 อ้างถึงใน ชนะศึก นิชานนท์, 2553) โดยวิเคราะห์ข้อสอบแล้ว ทดสอบความตรงของโครงสร้าง ด้วย χ^2 สำหรับ Likelihood Ratio G^2 ในการตรวจสอบมิติ ของแบบสอบถามดังนี้ที่ใช้นี้ทดสอบด้วยการกำหนดจำนวน องค์ประกอบของชุดข้อมูลไว้ล่วงหน้าแล้ว ทดสอบด้วย χ^2 ที่ประมาณค่าด้วยวิธี G^2 เพื่อทดสอบความเหมาะสมของโมเดล เมื่อค่า G^2 ไม่มีนัยสำคัญแสดงว่า ข้อมูลมีจำนวนขององค์ประกอบเท่าที่กำหนดในการทดสอบ

1.4 การวิเคราะห์ การจัดกลุ่มระดับขั้น (Hierarchical Cluster Analysis) เป็นเทคนิค สำหรับการทดสอบความเป็นพหุมิติของชุดแบบสอบถาม โดยการพิจารณาการแบ่งกลุ่มของตัวแปร โดยกระบวนการแบ่งกลุ่มนี้เป็นการแบ่งกลุ่มจำนวนข้อสอบที่มีลักษณะคล้ายคลึงกันให้อยู่ในกลุ่มเดียวกัน นอกจากนี้ยังใช้วิธีในการหมุนซ้ำ (Iteration) จนสำเร็จหรืออยู่ในระดับที่น่าพอใจซึ่งทำให้ได้ ความเป็นไปได้ของผลลัพธ์ (Outcome) โดยสามารถช่วยอธิบายให้ข้อมูลมีความถูกต้องมากยิ่งขึ้น ซึ่งการวิเคราะห์ด้วยวิธีการนี้สามารถใช้โปรแกรมคอมพิวเตอร์สำเร็จรูป CCPROX และ HCA ใน การวิเคราะห์ได้

1.5 การใช้โปรแกรม DETECT เป็นการตรวจสอบมิติแฟรงเชิงยืนยัน แบบ Nonparametric ซึ่งจะใช้ในการประมาณค่าจำนวนของมิติแฟรงที่มีคุณลักษณะเด่นในชุดของข้อมูล และสามารถตรวจสอบความเป็นเอกมิตรของแบบสอบถาม โดยระบุคุณลักษณะเด่นของมิติแฟรงในแต่ละข้อ ซึ่งผู้ใช้โปรแกรมสามารถระบุจำนวนมิติแฟรงสูงสุดที่ต้องการศึกษาได้ เนื่องจากการจัดกลุ่มชุดของ ข้อสอบแต่กระบวนการการดังกล่าวยังมีลักษณะแบบไม่เป็นทางการเท่าไนก็ เนื่องจากการระบุการจัด กลุ่มเพื่อจำแนกความแตกต่างของมิติจะอาศัยกระบวนการในการระบุความเป็นหนึ่งเดียว

1.6 การใช้โปรแกรม DIMTEST เป็นกระบวนการ ตรวจสอบสมมติฐานของแบบสอบถาม ด้วย Nonparametric Statistical โดยมีลักษณะคล้ายคลึงกับการตรวจสอบด้วยโปรแกรม DETECT โดยการตรวจสอบจะตรวจสอบความสัมพันธ์ระหว่างชุดข้อสอบบ่อยภายในตัวเองและกับความแปรปรวน ร่วมของข้อสอบ ซึ่งแตกต่างจากการใช้โปรแกรม DETECT ที่มีลักษณะคล้ายการวิเคราะห์ องค์ประกอบเชิงยืนยัน

2. ความเป็นอิสระในการตอบข้อสอบ (Local Independent) ความเป็นอิสระใน การตอบข้อสอบ หมายถึง ความน่าจะเป็นในการตอบข้อสอบแต่ละข้อได้ถูกต้องเป็นอิสระจากกัน คือ การตอบข้อสอบข้อใดข้อหนึ่งถูกหรือผิดจะไม่มีผลกระทบต่อการตอบข้ออื่นด้วย หรืออาจจะกล่าวใน

เชิงคณิตศาสตร์ได้ว่า ความเป็นอิสระในการตอบข้อสอบ หมายถึง ความน่าจะเป็นในการตอบข้อสอบถูกทั้งหมดมีค่าเท่ากับ ผลคูณของความน่าจะเป็นในการตอบข้อสอบถูกเป็นรายข้อ คือ ผู้สอบที่มีความสามารถ (θ) จะมีความน่าจะเป็นที่จะตอบข้อสอบทั้งข้อ 1 และข้อ 2 ถูกเท่ากัน ซึ่งได้มาจากความน่าจะเป็นในการตอบข้อสอบข้อที่ 1 ถูก และความน่าจะเป็นในการตอบข้อสอบข้อที่ 2 ถูก คือ ถ้าผู้สอบมีความสามารถ (θ) เท่ากับ 1.5 มีความน่าจะเป็นในการตอบข้อสอบข้อที่ 1 ถูกเท่ากับ 0.5 และมีความน่าจะเป็นในการตอบข้อสอบข้อที่ 2 ถูก เท่ากับ 0.6 ดังนั้นผู้สอบที่มีความสามารถเท่ากับ 1.5 มีความน่าจะเป็นในการตอบข้อสอบทั้งสองข้อถูกภายใต้เงื่อนไขความเป็นอิสระ มีค่าเท่ากับ 0.3 และ Hambleton and Swaminathan (1985) กล่าวว่า ถ้าแบบทดสอบมีความเป็นเอกมิตรอยู่แล้ว ความเป็นอิสระในการตอบข้อสอบก็จะเกิดขึ้นตามไปด้วย

3. โค้งคุณลักษณะข้อสอบ (Item Characteristic Curve) โค้งคุณลักษณะข้อสอบ เป็นฟังก์ชันทางคณิตศาสตร์ สามารถใช้อธิบายความสัมพันธ์ระหว่างความน่าจะเป็น หรือโอกาส ที่ผู้สอบจะตอบข้อสอบถูกกับระดับความสามารถที่วัดได้โดยใช้ชุดของข้อสอบ หรือแบบทดสอบฉบับนั้น ทั้งนี้ ความน่าจะเป็น หรือโอกาสในการตอบข้อสอบถูกจะขึ้นอยู่กับโค้งลักษณะข้อสอบในแต่ละโมเดลที่เลือกใช้ โดยที่รูปร่าง (Shape) ของโค้งคุณลักษณะข้อสอบในแต่ละข้อมีคุณสมบัติไม่แปรเปลี่ยน (Invariant) ไปตามกลุ่มตัวอย่างที่ใช้ ดังนั้น จึงทำให้ความน่าจะเป็น หรือโอกาสในการตอบข้อสอบถูกในแต่ละข้อไม่แปรเปลี่ยนด้วยคุณสมบัตินี้ถือเป็นลักษณะเด่นของโมเดลต่าง ๆ ในทฤษฎีการตอบสนองข้อสอบโค้งคุณลักษณะข้อสอบมีหลายรูปแบบขึ้นอยู่กับว่าเลือกใช้พารามิเตอร์ ของข้อสอบกี่พารามิเตอร์

4. ข้อสอบที่ใช้ต้องไม่เป็นข้อสอบประเภทความเร็ว (Speediness) ผู้สอบทุกคน ควรมีโอกาสในการทำข้อสอบทุกข้อ เพื่อให้คะแนนรวมจากการสอบเป็นค่าความสามารถที่แท้จริง ของผู้สอบไม่มีข้อจำกัดเกี่ยวกับเวลาในการสอบสรุปได้ว่าข้อตกลงเบื้องต้นของทฤษฎีการตอบสนอง ข้อสอบ คือ ข้อสอบต้องมีความเป็นเอกมิตร กล่าวคือ วัดความสามารถหรือคุณลักษณะเดียว ความน่าจะเป็นในการตอบข้อสอบแต่ละข้อได้ถูกต้อง จะต้องเป็นอิสระต่อกัน นอกเหนือนั้น โค้งคุณลักษณะข้อสอบเป็นฟังก์ชันทางคณิตศาสตร์ที่อธิบายความสัมพันธ์ระหว่าง ความน่าจะเป็น ในการตอบข้อสอบถูกกับระดับความสามารถของผู้สอบ และข้อสอบที่ใช้ต้องไม่เป็นข้อสอบประเภท ความเร็ว

วิธีการทดสอบอัตราส่วนไลค์ลิชต์

วิธีการทดสอบอัตราส่วนไลค์ลิชต์ (IRT-LR) ใช้หลักของอัลกอริทึมในการประมาณค่า ความเป็นไปได้สูงสุดในการประมาณค่าพารามิเตอร์ ในข้อมูลที่จำกัดวิธีการทดสอบอัตราส่วนไลค์ลิชต์ (IRT-LR) ใช้การประมาณการกำลังสองน้อยมากสำหรับแบบจำลองการตอบสนองข้อสอบของข้อสอบแบบปกติ และยังใช้ข้อมูลที่ต่ำกว่าเกณฑ์ของการตอบสนองรูปแบบการจำแนกผู้ตอบข้อสอบ (Marie Wiberg, 2007) โดยที่วิธี IRT-LR จะประเมินความสำคัญระหว่างกลุ่มอ้างอิงและกลุ่มเปรียบเทียบ ระหว่างความแตกต่างของค่าพารามิเตอร์ (Kim & Cohen, 1998) และ Thissen, David, Lynne Steinberg, and Howard Wainer (1988) เสนอไว้วิธีการทดสอบ IRT-LR เป็นที่นิยม เพราะมี การเปรียบเทียบของพารามิเตอร์และวัดพื้นที่ที่ต้องมีการประมาณการที่ถูกต้อง ในวิธี IRT-LR ข้อสอบบางข้ออาจจะเข้าข้างผู้สอบทั้งกลุ่มอ้างอิงและกลุ่มเปรียบเทียบ ซึ่งจะเรียกข้อสอบเหล่านั้นว่า DIF-

free โดยมีข้อจำกัดอยู่ระหว่างสองกลุ่มวิธี IRT-LR สามารถตรวจสอบการทำหน้าที่ต่างกันของข้อสอบได้ทั้งแบบเอกสารและonenewซึ่งมีการเปรียบเทียบกันกับแบบจำลองที่มีข้อจำกัดเท่ากันทั้งสองกลุ่มสถิติที่ใช้ในการทดสอบ คือ ความแตกต่างระหว่างค่าความเป็นไปได้ของ $-2\log$ สำหรับกลุ่ม L_C และกลุ่ม L_A ตามค่าที่ถูกกำหนด Acar (2010) โดยสามารถเขียนสมการการวิเคราะห์ได้ดังนี้

$$G^2(d.f.) = -2\log L_c - (-2\log L_A) \quad (6)$$

เมื่อ G^2 แสดงการกระจายของค่าไค-สแควร์ χ^2 และค่าพารามิเตอร์รายข้อทั้งยังสามารถพิจารณาจากค่าของ p -value ที่มีนัยสำคัญทางสถิติที่ .05 เมื่อข้อสอบมีการทำหน้าที่ต่างกัน

David, Lynne, and Meg (1986) ได้เสนอวิธีการทดสอบอัตราส่วนความเป็นไปได้ที่ใช้ในการตรวจสอบการทำหน้าที่เบี่ยงเบนของข้อสอบ โดยเป็นการทดสอบผลการตอบข้อสอบระหว่างผู้สอบกลุ่มย่อย 2 กลุ่ม ซึ่งวิธีทดสอบอัตราส่วนความเป็นไปได้ แบ่งเป็น 3 วิธี คือ 1) วิธีการทดสอบอัตราส่วนความเป็นไปได้ในทฤษฎีการตอบข้อสอบ ในรูปทั่วไป (General LR) ซึ่งเป็นวิธีการตรวจสอบการทำหน้าที่เบี่ยงเบนของข้อสอบซึ่งใช้การประมาณค่าพารามิเตอร์ในไมเดลการตอบข้อสอบวิธีความเป็นไปได้สูงสุดแบบมาร์จินอล (Marginal Maximum likelihood: MML) 2) วิธีการทดสอบอัตราส่วนความเป็นไปได้ในทฤษฎีการตอบข้อสอบแบบโลกลิнейร์ (Loglinear: LR) ซึ่งเป็นการประมาณค่าพารามิเตอร์แบบวิธีความเป็นไปได้สูงสุด (Maximum likelihood: ML) และ 3) วิธีการทดสอบอัตราส่วนความเป็นไปได้ของทฤษฎีการตอบข้อสอบสารสนเทศที่มีขอบเขตแบบจำกัด (Limited Information: LR) วิธีนี้เป็นการประมาณค่าพารามิเตอร์ในไมเดลการตอบข้อสอบแบบ Normal Ogive เป็นวิธีกำลังสองน้อยที่สุดในรูปทั่วไป (Generalized least Squares: GLS) ซึ่ง 3 วิธีดังกล่าวจะใช้ในการทดสอบอัตราส่วนความเป็นไปได้เพื่อทดสอบหากค่านัยสำคัญของการทำหน้าที่เบี่ยงเบนของข้อสอบ ซึ่งความหมายของโมเดลของฟังก์ชันในทฤษฎีการตอบข้อสอบ (IRT) ซึ่งค่าความหมายของฟังก์ชันเป็นดัชนีไมเดลความหมายสมกับข้อมูลกับการประมาณค่าความเป็นไปได้สูงสุดของการประมาณค่าพารามิเตอร์ข้อสอบ

Embretson and Resie (2000) ไมเดลการตอบสนองข้อสอบ (IRT Models) เดิมมีแนวคิดพื้นฐานมาจากทฤษฎีการวัดแบบดั้งเดิม (Classical Test Theory) โดยไมเดลการตอบสนองข้อสอบเป็นไมเดลการวัดความสัมพันธ์ระหว่างตัวแปรอิสระ ที่จะประกอบไปด้วยตัวแปรແ Pang คือ ความสามารถของผู้สอบ (θ) หรือค่าพารามิเตอร์ของผู้สอบ (a, b, c) และตัวแปรอิสระที่สามารถสังเกตได้ คือ โอกาสในการตอบข้อสอบได้ถูกต้องของผู้สอบ ทฤษฎีการตอบสนองข้อสอบ (IRT) จะเป็นการอธิบายความสัมพันธ์ความสามารถภายในของตัวบุคคล (Latent Trait or Ability) กับผลของการตอบข้อสอบถูกหรือโค้งคุณลักษณะข้อสอบ (ItemCharacteristic Curve: ICC) โดยมีการกำหนดลักษณะข้อสอบด้วยพารามิเตอร์ของข้อสอบ คือค่าอำนาจจำแนกของข้อสอบ (a) ค่าความยากของข้อสอบ (b) และค่าโอกาสในการเดาของข้อสอบได้ถูก (c) ดังนั้น ทฤษฎีการตอบสนองข้อสอบ (IRT) จึงอยู่บนพื้นฐานของความคิด 2 ประการ คือ 1) ผลการตอบข้อสอบได้ถูกต้องของผู้สอบ ที่จะอธิบายความสามารถที่มีอยู่ภายในตัวผู้สอบได้ และ 2) ความสัมพันธ์ระหว่างผล

การตอบข้อสอบกับความสามารถที่อยู่ภายในตัวของผู้สอบ โดยจะอธิบายได้ด้วยโครงคุณลักษณะข้อสอบ (ICC) และ ฟังก์ชันของลักษณะข้อสอบ

สรุป คือ ทฤษฎีการตอบสนองข้อสอบ หรือ IRT นี้ สามารถจำแนกได้เป็น 2 ประเภท ได้แก่ ทฤษฎีการตอบสนองข้อสอบแบบตรวจให้คะแนน 2 ค่า ซึ่งเป็นโมเดลการตอบสนองข้อสอบ ที่ใช้กับการตรวจสอบค่าคะแนนรายข้อแบบ 2 ค่า เช่น ข้อสอบหรือข้อคำถามที่ตรวจให้คะแนนแบบ 0, 1 (ตอบผิดได้ 0, ตอบถูกได้ 1) แบบถูก/ผิด ใช่/ไม่ใช่ ซึ่งโมเดลการตอบสนองข้อสอบ (IRT Models) เดิมมีแนวคิดพื้นฐานมาจากทฤษฎีการวัดแบบดั้งเดิม โดยโมเดลการตอบสนองข้อสอบเป็นโมเดล การวัดความสัมพันธ์ระหว่างตัวแปรอิสระ ที่จะประกอบไปด้วยตัวแปรແ Pang คือ ความสามารถของผู้สอบ (θ) หรือค่าพารามิเตอร์ของผู้สอบ (a, b และ c) และตัวแปรอิสระที่สามารถสังเกตได้ คือ โอกาสในการตอบข้อสอบได้ถูกต้องของผู้สอบ ทฤษฎีการตอบสนองข้อสอบ (IRT) จะเป็นการอธิบาย ความสัมพันธ์ความสามารถภายนอกของตัวบุคคล กับผลของการตอบข้อสอบถูกหรือโดยคุณลักษณะ ข้อสอบ

งานวิจัยที่เกี่ยวข้องกับการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบด้วยวิธีการทดสอบอัตราส่วนไลค์ลิขิต มีดังนี้

รุ่งภา แสนอำนวย ประกฤติยา ทักษิณ และชนะศึก นิชานนท์ (2555) ได้ศึกษา ประสิทธิภาพของแบบทดสอบวัดผลสัมฤทธิ์ทางการเรียนรูปแบบผสมการประยุกต์ ใช้ทฤษฎี การตอบสนองข้อสอบแบบตรวจให้คะแนนความรู้บางส่วน และทฤษฎีการตอบสนองข้อสอบ แบบตรวจให้คะแนนความรู้บางส่วนแบบทั่วไป มีวัตถุประสงค์ เพื่อศึกษาประสิทธิภาพของแบบทดสอบรูปแบบผสม และเพื่อศึกษาปฏิสัมพันธ์ระหว่างโมเดลการตรวจให้คะแนนสัดส่วนของข้อสอบ ที่ตรวจให้คะแนนแบบสองค่าและมากกว่าสองค่า และความยาวของแบบทดสอบ และเปรียบเทียบ ประสิทธิภาพของแบบทดสอบรูปแบบผสม เงื่อนไขที่ทำการศึกษามี 18 เงื่อนไข ประกอบด้วยโมเดล การตรวจให้คะแนน 2 โมเดล คือโมเดลโลจิสติก 1 พารามิเตอร์กับโมเดลการตรวจให้คะแนนความรู้ บางส่วน (PCM) และโมเดลโลจิสติก 3 พารามิเตอร์กับโมเดลการตรวจให้คะแนนความรู้บางส่วน (GPCM) สัดส่วนข้อสอบที่ตรวจให้คะแนนแบบสองค่าและมากกว่าสองค่า 3 สัดส่วน คือ 20:80 50:50 และ 80:20 และความยาวของแบบทดสอบ 3 เงื่อนไข คือ 10 30 และ 50 ข้อ การประเมินประสิทธิภาพของแบบทดสอบรูปแบบ ผสมพิจารณาจากดัชนีความคลาดเคลื่อน มาตรฐานในการประเมินค่า $SE(\theta)$ ค่าความล้าเอียง (Bias) พร้อมทั้งวิเคราะห์ความแปรปรวนแบบ พหุจำแนก 3 ทาง (Three-way MANOVA) เพื่อเปรียบเทียบค่าเฉลี่ยของดัชนี $SE(\theta)$ และค่าความล้าเอียง (Bias) พบว่า โมเดลโลจิสติก 1 พารามิเตอร์กับ PCM และโมเดลโลจิสติก 3 พารามิเตอร์กับ GPCM มีค่า $SE(\theta)$ และค่าความล้าเอียง (Bias) ต่ำสุดที่สัดส่วนข้อสอบที่ตรวจให้คะแนนสองค่าและมากกว่าสองค่า คือ 20:80 และความยาว ของแบบทดสอบ 50 มีปฏิสัมพันธ์ระหว่างโมเดลการตรวจให้คะแนน สัดส่วนของข้อสอบที่ตรวจให้คะแนนสองค่า และมากกว่าสองค่าและความยาวของแบบทดสอบที่ส่งผลต่อค่า $SE(\theta)$ และค่าความล้าเอียง (Bias) อย่างมีนัยสำคัญทางสถิติที่ระดับ .05 ส่วนปฏิสัมพันธ์รายคู่ พบร่วมกับ ปฏิสัมพันธ์ระหว่างโมเดลการตรวจให้คะแนนกับสัดส่วนของข้อสอบ ที่ตรวจให้คะแนน สองค่าและมากกว่าสองค่า ระหว่างโมเดลการตรวจให้คะแนน กับความยาวของแบบทดสอบ และระหว่างสัดส่วน ของข้อสอบที่ตรวจให้คะแนนสองค่าและมากกว่าสองค่ากับ

ความยากของแบบทดสอบ ส่งผลต่อค่า $SE(\theta)$ และค่าความลำเอียง (Bias) อย่างมีนัยสำคัญทางสถิติ ที่ระดับ .05 นอกจากนี้พบว่าไม่เดลการตรวจให้คะแนน สัดส่วนของข้อสอบที่ตรวจให้คะแนนสองค่า และมากกว่าสองค่า และความยากของแบบทดสอบที่แตกต่างกันส่งผลต่อค่า $SE(\theta)$ และค่าความลำเอียง (Bias) ที่ต่างกันอย่างมีนัยสำคัญทางสถิติที่ระดับ .05

สุราทิพย์ ตรีสิน และปิยะทิพย์ ประดุจพร (2560) ได้ศึกษาการเปรียบเทียบผล การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ด้านภาษา ด้านคำนวน และด้านเหตุผล ขั้นประถมศึกษาปีที่ 3 ด้วยวิธี HGLM วิธี MIMIC และวิธี IRT-LR มีวัตถุประสงค์ เพื่อวิเคราะห์คุณภาพของแบบทดสอบระดับชาติ (NT) และตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ ในแบบทดสอบระดับชาติขั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน และด้านเหตุผล ด้วยวิธี HGLM วิธี MIMIC และวิธี IRT-LR ปรากฏว่า แบบทดสอบระดับชาติขั้นประถมศึกษาปีที่ 3 มีค่าความยากของข้อสอบ (b) อยู่ในระดับค่อนข้างยาก มีค่าอำนาจจำแนกของข้อสอบ (a) อยู่ในระดับที่สามารถจำแนกผู้สอบได้ดี และมีค่าโอกาสในการเดาของข้อสอบ (c) ไม่เกิน 0.3 ในการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบทั้ง 3 ด้าน ซึ่งให้เห็นว่า เพศส่งผลให้เกิดการทำหน้าที่ ต่างกันของข้อสอบ โดยเพศหญิงจะได้เปรียบในการตอบข้อสอบด้านภาษา และด้านเหตุผล ในขณะที่ เพศชาย จะได้เปรียบในการตอบข้อสอบด้านคำนวน โดยวิธี HGLM ตรวจพบข้อสอบทำหน้าที่ต่างกัน จำนวนมากที่สุด คิดเป็นร้อยละ 69 ของข้อสอบทั้งฉบับ รองลงมาคือ วิธี IRT-LR ร้อยละ 54 และ วิธี MIMIC ร้อยละ 16 ตามลำดับ และการเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกัน ของข้อสอบ พบร่วมกับ วิธี HGLM ตรวจพบ DIF มากกว่า วิธี MIMIC ในด้านภาษา ด้านคำนวน และ ด้านเหตุผลคิดเป็นร้อยละ 70, 36 และ 53 ตามลำดับ และวิธี HGLM ตรวจพบ DIF มากกว่า วิธี IRT-LR ด้านภาษา และด้านคำนวน คิดเป็นร้อยละ 37 และ 13 และวิธี IRT-LR ตรวจพบ DIF มากกว่า วิธี MIMIC ทั้ง 3 ด้าน คิดเป็นร้อยละ 33, 43 และ 40 ตามลำดับ ส่วนวิธี HGLM ตรวจพบ DIF น้อยกว่า วิธี IRT-LR ด้านคำนวน คิดเป็นร้อยละ 7 ($p < .05$)

ธีระวัฒน์ สุขีสาร ดุษฎี โยเหลา เสกสรรค์ ทองคำบรรจง และนิษดา จิตต์จรัส (2555) ได้ศึกษาความเที่ยงตรงของการประมาณค่าในการวิเคราะห์ไม่เดลสมการโครงสร้างพหุระดับภายใต้ เงื่อนไขวิธีการประมาณค่าและขนาดตัวอย่างที่ต่างกัน มีวัตถุประสงค์เพื่อศึกษาความเที่ยงตรง ของวิธีการประมาณค่าพารามิเตอร์สำหรับงานวิจัยไม่เดลสมการโครงสร้างพหุระดับภายใต้เงื่อนไข ขนาดตัวอย่างของแต่ละระดับการวิเคราะห์ที่แตกต่างกัน ซึ่งข้อมูลที่ใช้ในการวิเคราะห์ เป็นข้อมูล ทุติยภูมิ โดยคัดเลือกข้อมูลจากโครงการศึกษาโอกาสในการเรียนรู้ทางด้านคณิตศาสตร์และ วิทยาศาสตร์ของนักเรียนชั้นมัธยมศึกษาปีที่ 3 ของสถาบันส่งเสริมการสอนวิทยาศาสตร์และ เทคโนโลยี (สวท.) โดยมีการสุ่มข้อมูลตามเงื่อนไขขนาดตัวอย่างระดับที่สอง 50 โรงเรียน 30 โรงเรียน และ 15 โรงเรียน ขนาดตัวอย่างระดับที่หนึ่ง โรงเรียนละ 10 คน โรงเรียนละ 5 คน และขนาด เท่ากับแต่ละโรงเรียนที่เก็บจริง (10-40) วิเคราะห์ข้อมูลแต่ละเงื่อนไขด้วยวิธีการประมาณค่า 2 วิธี คือ Full Information Maximum Likelihood (FIML) และ Robust Maximum Likelihood (RML) และก็ตรวจสอบความเที่ยงตรงโดยการพิจารณาจากค่าการประมาณค่าน้ำหนักองค์ประกอบ ค่าสัมประสิทธิ์เส้นทาง และค่าสัมประสิทธิ์อิทธิพลข้ามระดับโดยการเปรียบเทียบค่าสถิติกับ ค่าพารามิเตอร์ของประชากร ด้วยการวิเคราะห์สถิติทดสอบที่ พบร่วมกับ การทดลองความเที่ยงตรง

ของการประมาณค่าพารามิเตอร์ระหว่าง วิธีการประมาณค่าพารามิเตอร์แบบ FIML และ RML พบว่า เมื่อกลุ่มตัวอย่างระดับที่สอง 50 โรงเรียน และ 30 โรงเรียน ส่วนกลุ่มตัวอย่างระดับที่หนึ่ง โรงเรียน ละ 5 คน โรงเรียนละ 10 คน และจำนวนนักเรียนที่เก็บจริงแต่ละโรงเรียน (10-40) ซึ่งทั้งสองวิธี ประมาณค่าน้ำหนักองค์ประกอบ ค่าสัมประสิทธิ์เส้นทางและค่าสัมประสิทธิ์อิทธิพลข้ามระดับ ไม่ต่างกัน คือการประมาณค่าทั้งสองวิธีให้ผลการประมาณค่าที่มีความเที่ยงตรงไม่ต่างกันส่วนกรณี ขนาดตัวอย่างระดับที่สอง 15 โรงเรียน ตัวอย่างระดับที่หนึ่งเก็บจริงแต่ละโรงเรียน (10-40 คน) วิธีการประมาณค่าแบบ RML ประมาณค่าต่าง ๆ ครบ ส่วนวิธีการประมาณค่าพารามิเตอร์ แบบ FIML ไม่สามารถประมาณค่าต่าง ๆ ได้ครบ

Teresi (2007) ได้ศึกษาการประเมินผลการวัดความเท่าเทียมโดยการใช้ทฤษฎี การตอบสนองข้อสอบแบบอัตราส่วนไลค์ลิลี่ (IRT-LR) สำหรับการตรวจสอบการทำหน้าที่ต่างกัน ของข้อสอบ (DIF) ในแบบวัดความสามารถและความทุกข์ในการทำงาน ซึ่งมีเงื่อนไขในการตรวจสอบ การทำหน้าที่ต่างกันของข้อสอบ คือ เพศ อายุ โดยมีกลุ่มตัวอย่างเป็นผู้ป่วยมะเร็ง จำนวน 1,714 คน จะใช้เครื่องมือในการวิจัยเป็นแบบวัดความสามารถทางกายภาพ จำนวน 23 ข้อ และแบบวัด ทางอารมณ์ จำนวน 15 ข้อ จากการศึกษาพบว่า อายุ และความสามารถ มีผลต่อการที่จะได้รับ การปฏิบัติที่แตกต่างกันออกไป เป็นปัจจัยที่อาจจะส่งผลต่อผู้ป่วยที่จะเลือกเข้ารับการรักษาของ สถานรักษาพยาบาล

Acar and Kelecioglu (2010) ได้ศึกษาเปรียบเทียบการตรวจสอบการทำหน้าที่ต่างกัน ของข้อสอบ ระหว่างวิธี HGLM วิธี LR และวิธี IRT-LR มีวัตถุประสงค์ของการวิจัยครั้งนี้คือ การตรวจสอบความสามารถคล้องระหว่างวิธีที่ใช้ตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ 3 วิธี มีเงื่อนไขในการตรวจสอบ คือ เพศ โดยมีกลุ่มตัวอย่างเป็นนักเรียนในประเทศไทย ซึ่งจะใช้ เครื่องมือในการวิจัยคือ แบบทดสอบของวิชาสังคมศาสตร์และวิทยาศาสตร์ จากการศึกษาพบว่า วิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบทั้ง 3 วิธี คือ วิธี HGLM วิธี LR และวิธี IRT-LR มีการตรวจพับข้อสอบที่ทำหน้าที่ต่างกันในปริมาณที่ใกล้เคียงกัน แต่ไม่ตรงกัน โดยวิธี LR และ วิธี IRT-LR ตรวจพับข้อสอบที่ทำหน้าที่ต่างกันได้น้อยทั้งแบบทดสอบสังคมศาสตร์และวิทยาศาสตร์ ส่วนวิธี HGLM ตรวจสอบพับข้อสอบที่ทำหน้าที่ต่างกันในแบบทดสอบทั้ง 2 แบบทดสอบในปริมาณ ที่มากที่สุด

Pae (2012) ได้ศึกษาการทำงานของข้อสอบที่แตกต่างระหว่างเพศ ในการทดสอบย่อ ภาษาอังกฤษของการทดสอบความถนัดทางวิชาการของนักวิชาการวิทยาลัยเกาหลี (Korean College Scholastic Aptitude Test: KCSAT) ในช่วงเวลา 9 ปี ชุดข้อมูล 3 ชุด โดยใช้ทั้งอัตราส่วน Mantel-Haenszel (MH) และทฤษฎีการตอบสนองข้อสอบแบบอัตราส่วนไลค์ลิลี่ (IRT-LR) ขั้นตอน ต่าง ๆ การศึกษายังระบุถึงปัจจัย 2 ประการ ได้แก่ กลยุทธ์การอ่านและการรับรู้ความสนใจซึ่งอธิบาย ถึงความแปรปรวนของขนาดการทำหน้าที่ต่างกันที่มีต่อเพศ โดยใช้การวิเคราะห์การทดสอบเชิงเส้น หลายแบบ จากการศึกษา พบว่า ปฏิสัมพันธ์การทำหน้าที่ต่างกันของข้อสอบระหว่างข้อสอบรายข้อ กับเพศ และความสัมพันธ์ที่สำคัญระหว่างความแตกต่างของเพศ ความสนใจผลของการทดสอบและ ขนาดทดสอบในการทำหน้าที่ต่างกันของข้อสอบต่อเพศ ซึ่งการรับรู้และการอ่านมีผลต่อการทำหน้าที่ ต่างกันของข้อสอบต่อเพศ

Li, Hunter, and Oshima (2013) ได้ศึกษาการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในการทดสอบการอ่านและเหตุผลที่เป็นไปได้ระหว่างเพศ จากการศึกษาอย่างละเอียดพบว่า ข้อสอบ 1,210 ข้อ จาก 18 บทความได้รับรวมไว้ในการวิเคราะห์ขั้นสุดท้าย พบร่วม 23.3% ของข้อสอบ แสดงให้เห็นถึงการทำหน้าที่ต่างกันของข้อสอบในด้านแพร่เพศ และมีการเปลี่ยนแปลงค่าเปอร์เซ็นต์ของข้อสอบที่ได้กำหนดไว้ว่าแสดงการทำหน้าที่ต่างกันของข้อสอบระหว่างการศึกษาตั้งแต่ 0 ถึง 77% ของข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบที่ครึ่งหนึ่งลำเอียงเข้าข้างเพศชายและอีกครึ่งหนึ่งที่ลำเอียงเข้าข้างเพศหญิง โดยรูปแบบที่เป็นจริง สำหรับการศึกษาโดยใช้การทดสอบการตอบสนองต่อข้อสอบ (IRT-LR) และสำหรับผู้ที่ใช้ วิธี Mantel-Haenszel (MH) นอกจากนี้ แบบทดสอบที่สันก่อว่ามีแนวโน้มที่จะได้รับการพิจารณาทำหน้าที่ต่างกันของข้อสอบมากกว่าที่แบบทดสอบที่ยกกว่า รูปแบบการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบอื่น ๆ จะขึ้นอยู่กับวิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ และแบบทดสอบ

Kabasakal, Arsan, Gok, and Kelecioglu (2014) ได้ศึกษาแบบจำลองเปรียบเทียบผลตราความคลาดเคลื่อนประเพณีที่ 1 และประสิทธิภาพในการตรวจสอบการทำหน้าที่ต่างกันของวิธี Mantel-Haenszel (MH) วิธี SIBTEST และวิธีการทดสอบอัตราส่วนไลค์ลิชต์ (IRT-LR) โดยมีภัยได้เงื่อนไขบางประการ ปัจจัยการจัดการคือขนาดตัวอย่างความแตกต่างระหว่างกลุ่มที่มีความสามารถความพยายามแบบทดสอบร้อยละของการทำหน้าที่ต่างกันของข้อสอบ (DIF) และรูปแบบพื้นฐานที่ใช้ในการสร้างข้อมูลผลการวิจัยแสดงให้เห็นว่า SIBTEST มีข้อผิดพลาดประเพณีที่สูงที่สุดในการตรวจสอบแบบข้อสอบทำหน้าที่ต่างกันแบบเอกสาร ส่วนวิธี MH มีอัตราความถูกต้องสูงภายใต้เงื่อนไขทั้งหมด นอกจากนี้ร้อยละของการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ (DIF) และโมเดลพื้นฐานยังมีผลที่อัตราความคลาดเคลื่อนประเพณีที่ 1 ของ IRT-LR ความแตกต่างของความสามารถระหว่างกลุ่มความพยายามแบบทดสอบร้อยละของ DIF รูปแบบ และความสัมพันธ์ระหว่างความสามารถแตกต่างของความสามารถร้อยละของ DIF ความแตกต่างของความสามารถความพยายามแบบทดสอบมีผลต่อ วิธี SIBTEST อัตราความคลาดเคลื่อนประเพณีที่ 1 ในขั้นตอน MH ปัจจัยที่มีผลต่ออัตราความคลาดเคลื่อนประเพณีที่ 1 ได้แก่ ขนาดตัวอย่างระยะเวลาในการทดสอบร้อยละ DIF ความแตกต่างของความสามารถร้อยละของ DIF ความแตกต่างของความสามารถรูปแบบความสามารถและความแตกต่างของ DIF ไม่มีผลต่อประสิทธิภาพของ วิธี SIBTEST และ วิธี MH แต่รูปแบบพื้นฐานมีผลต่ออัตราการใช้ วิธี IRT-LR

Lopez (2012) ได้ศึกษาเรื่องตรวจสอบและจัดหมวดหมู่ของประเพณี DIF ใช้ Parametric และ Nonparametric โดยการเปรียบเทียบวิธีการทดสอบอัตราส่วนไลค์ลิชต์ (IRT-LR) วิธี SIBTEST และวิธีการทดสอบโลจิสติก มีวัตถุประสงค์ของการตรวจสอบนี้มีวัตถุประสงค์เพื่อเปรียบเทียบประสิทธิภาพของสามวิธีการที่สำหรับการตรวจสอบรายการที่แตกต่างกันทำงาน (DIF) ประสิทธิภาพการทำงานของข้ามทดสอบคุณภาพการพร้อมกัน ความน่าจะเป็นทฤษฎีการตอบสนองข้อสอบ วิธีการทดสอบอัตราส่วนไลค์ลิชต์ (IRT-LR) และการทดสอบโลจิสติกถูกตรวจสอบในช่วงของเงื่อนไขการทดลองรวมทั้งความสามารถทดสอบที่แตกต่างกันขนาดตัวอย่าง DIF และการทดสอบการทำงานที่แตกต่างกัน (DTF) ขนาดและหมายถึงความสามารถแตกต่างในการกระจายลักษณะพื้นฐานของกลุ่ม เปรียบเทียบในที่นี้จะเรียกว่าเป็นข้อมูลอ้างอิงและกลุ่มฟอกส์ นอกจากนี้ในแต่ละขั้นตอนได้ดำเนินการ

โดยใช้เทียบทั้งหมดอีน ๆ วิธีการซึ่งในรูปแบบพื้นฐานวิธีการทดสอบอัตราส่วนไลค์ลิขิต (IRT-LR) วิธี CSIBEST การจับคุณภาพทดสอบบ่อยและลักษณะ LOGREG ประมาณการอยู่บนพื้นฐานของ การทดสอบทุกรายการยกเว้นหนึ่งภายใต้การศึกษาและวิธีการที่ผู้ประกาศข่าวอย่างต่อเนื่องซึ่งใน รูปแบบพื้นฐานที่ตรงกับการทดสอบบ่อย และประมาณการลักษณะอยู่บนพื้นฐานของเขตที่กำหนดไว้ ล่วงหน้าระยะเวลอิสระของรายการ DIF ข้อมูลการตอบสนองสำหรับการอ้างอิงและกลุ่มไฟกัสที่ถูก สร้างขึ้นโดยใช้พารามิเตอร์ที่รู้จักกันในรายการที่อยู่บนพื้นฐานของสามพารามิเตอร์โลจิสติกรูปแบบ ทฤษฎีการตอบสนองข้อสอบ (3 PLM) ประเภทต่าง ๆ ของ DIF ถูกจำลองโดยขับพารามิเตอร์ การสร้างรายการของรายการที่เลือกเพื่อให้บรรลุ DIF ต้องการและขนาด DTF ขึ้นอยู่กับพื้นที่ระหว่าง รายการที่กลุ่มฟังก์ชันการตอบสนอง อัตราความถูกต้อง ความผิดพลาดประเภทและประเภทที่สาม อัตราความผิดพลาดถูกคำนวณสำหรับแต่ละสภาพขั้นอยู่กับการทดลองซ้ำ 100 และผลการวิเคราะห์ ผ่านการวิเคราะห์ความแปรปรวน ผลการศึกษาพบว่าวิธีการที่แตกต่างกันในการรับรู้ความสามารถ มี LOGREG เมื่อดำเนินการโดยใช้วิธีการอื่น ๆ ทั้งหมดให้สมดุลที่ดีที่สุดของอำนาจจำแนกและ อัตราความผิดพลาด การจำลองครั้งนี้ แต่ยังไม่มีวิธีการที่มีประสิทธิภาพในการระบุชนิดของ DIF

สรุปจากการศึกษางานวิจัยที่เกี่ยวข้อง วิธี IRT-LR สามารถตรวจสอบข้อสอบการทำหน้าที่ ต่างกันของเพศได้มากกว่าวิธี SIBTEST และสามารถควบคุมอัตราความคลาดเคลื่อนประเภทที่ 1 ได้ดี ในแบบทดสอบที่มีความยาวไม่เกิน 30 ข้อ และตรวจพบข้อสอบการทำหน้าที่ต่างกันได้น้อยกว่า วิธี HGLM และควบคุมอัตราความคลาดเคลื่อนประเภทที่ 1 ได้น้อย เมื่อแบบทดสอบมีความยาวกว่า 30 ข้อและขนาดของกลุ่มตัวอย่างมีขนาดใหญ่

ตอนที่ 5 การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบด้วยวิธีซิปเทส์ และงานวิจัย ที่เกี่ยวข้อง

วิธีซิปเทส์ (Simultaneous Item Bias Test: SIBTEST) ได้พัฒนาโดยเชียเลียและ สเตาท์ (Shealy & Stout, 1993) เพื่อนำมาใช้การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ (Differential Item Functioning: DIF) การทำหน้าที่ต่างกันของแบบทดสอบ (Differential Test Functioning: DTF) และการทำหน้าที่ต่างกันของกลุ่มข้อสอบ (Differential Bundle Functioning: DBF) สามารถใช้ได้กับแบบทดสอบเอกมิติ (Unidimensional Test) และแบบทดสอบพหุมิติ (Multidimensional Test) (Stout, Li & Nandakumar, 1997) วิธีซิปเทส์เป็นสถิติ แบบนั้นพารามิตริก (Nonparametric) ซึ่งไม่ต้องใช้ฟังก์ชันการตอบสนองข้อสอบหรือการประมาณ ค่าความสามารถแฟง วิธีซิปเทส์ใช้กับข้อสอบแบบเอกรูป (Uniform DIF) จึงทำให้มีความໄ้ใน การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบแบบอนุรูป (Nonuniform DIF) (Li & Stout, 1996) ข้อดีของวิธีซิปเทส์ คือ คำนวณได้ง่าย ไม่ซับซ้อน ประยัดค่าใช้จ่าย ซึ่งไม่จำเป็นที่จะต้องใช้ กับกลุ่มตัวอย่างที่มีขนาดใหญ่ และเป็นสถิติทดสอบนัยสำคัญ (Narayanan & Swaminathan, 1996) และนำไปประยุกต์ใช้กับการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบที่มีการให้คะแนน แบบพหุวิภาค (Polytomous DIF)

วิธีซิปเทส์ใช้ในการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบเอกมิติจะต้อง มีข้อตกลงในแบบทดสอบ คือ ข้อสอบในแบบทดสอบจะต้องมุ่งวัดคุณลักษณะ หรือความสามารถแฟง

ลักษณะเดียวเท่านั้น ซึ่งความสามารถแฝงประเภทที่หนึ่ง เรียกว่า ความสามารถเป้าหมายที่ต้องการวัด (Target Ability: θ) แต่จะมีความสามารถแฝงอีกประเภทหนึ่งที่มีอิทธิพลต่อการตอบข้อสอบ เรียกว่า ความสามารถแทรกซ้อนที่ไม่ต้องการวัด (Nuisance Ability: η) เช่น แบบทดสอบคำศัพท์ในรายวิชาภาษาต่างประเทศ ซึ่งในข้อสอบบางข้ออาจมีคำถกความรู้สำหรับผู้ชายเกี่ยวกับ อิเล็กทรอนิกส์ และบางข้ออาจถกความรู้สำหรับผู้หญิงเกี่ยวกับการเย็บปักถักร้อย เป็นต้น จากแบบทดสอบดังกล่าวเป็นทักษะความรู้เกี่ยวกับคำศัพท์วิชาภาษาอังกฤษ ซึ่งเป็นเป้าหมายที่ต้องการวัด (θ) ส่วนความรู้เกี่ยวกับอิเล็กทรอนิกส์และการเย็บปักถักร้อย จัดเป็นความสามารถแทรกซ้อนที่ไม่ต้องที่จะวัด (η_1 และ η_2) ตามลำดับ ข้อสอบทุกข้อในแบบทดสอบจะวัดความสามารถเป้าหมาย แต่ข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบจะวัดทั้งความสามารถเป้าหมาย และความสามารถแทรกซ้อน

ถ้าให้ฟังก์ชันการตอบสนองข้อสอบ (IRF) ข้อที่ i ซึ่งขึ้นอยู่กับความสามารถ θ เพียงอย่างเดียวแทนด้วย $P_i(\theta)$ ส่วน IRF ข้อที่ i ที่ขึ้นอยู่กับความสามารถทั้ง θ และ η แทนด้วย $P_i(\theta, \eta)$ ฟังก์ชันการตอบสนองข้อสอบของข้อสอบดังกล่าวแบบ 3 พารามิเตอร์เป็นดังนี้ (Shealy & Stout, 1993)

$$P_i(\theta) = c_i \frac{(1-c)}{1 + \exp[-1.7(a_{10}(\theta - b_{10})]} , i = 1, \dots, N \quad (7)$$

$$P_i(\theta, \eta) = c_i + \frac{(1-c)}{1 + \exp[-1.7(a_{10}(\theta - b_{10}) + a_{in}(\eta - b_{in}))]} , i = 1, \dots, N \quad (8)$$

ดังนั้นฟังก์ชันความน่าจะเป็นอย่างมีเงื่อนไขของแบบแผนการตอบข้อสอบทั้งฉบับเป็นดังนี้

$$P[U/(\Theta=\theta, \eta)] = \prod_{i=1}^N p_i(\theta, \eta)^{U_i} [1 - P_i(\theta, \eta)]^{1-U_i} \quad (9)$$

จีเลียและสเตาท์ (Shealy & Stout, 1993) ได้ใช้ marginal IRFs อธิบายการทำหน้าที่ต่างกันของข้อสอบ

$$M_{ig}(\theta) = \int_n P_i(\theta, \eta) f_g(\eta | \theta) d\eta \quad (10)$$

เมื่อ $M_{ig}(\theta)$ = marginal IRF สำหรับความสามารถเป้าหมายที่ต้องการวัด θ ของผู้เข้าสอบกลุ่มอ้างอิงหรือกลุ่มเปรียบเทียบ

$P_i(\theta, \eta)$ = IRF ของข้อสอบข้อที่ i

$f_g(\eta | \theta)$ = การแจกแจงแบบมีเงื่อนไขของกลุ่มผู้เข้าสอบ

การเปรียบเทียบ marginal IRF ระหว่างกลุ่มอ้างอิง (R) กับกลุ่มเปรียบเทียบ (F) จะทำให้ทราบถึงทิศทางของการได้เปรียบหรือเสียเปรียบ กล่าวคือ ถ้า $M_{iF}(\theta) < M_{iR}(\theta)$ ทุกค่าของ θ แสดงว่า ข้อสอบทำหน้าที่ต่างกันแบบมีทิศทางเดียว โดยข้อสอบจะเข้าข้างผู้สอบกลุ่มอ้างอิง และ ถ้า $M_{iF}(\theta) >$

$M_{iR}(\theta)$ ทุกค่าของ θ แสดงว่าข้อสอบทำหน้าที่ต่างกันแบบมีทิศทางเดียว โดยข้อสอบจะเข้าข้างผู้สอบกลุ่มเปรียบเทียบการทำหน้าที่ต่างกันของข้อสอบแบบมีทิศทางเดียวอาจเรียก อีกอย่างหนึ่งว่า “การทำหน้าที่ต่างกันแบบไม่ตัดกัน” (Noncrossing DIF)

วิธีชิปเทสท์มีหลักการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบแบบเอกสาร โดยการเปรียบเทียบผลการตอบข้อสอบระหว่างกลุ่มอ้างอิงกับกลุ่มเปรียบเทียบ จะแบ่งแบบทดสอบออกเป็น 2 ชุด ย่อย (Subtest) คือ (1) ชุดแบบทดสอบที่มีความตรง (Valid Subtest) คือ แบบทดสอบที่ใช้เป็นเกณฑ์ในการเปรียบเทียบ (Matching Subtest) โดยแบบทดสอบประกอบด้วย ข้อสอบที่ทำหน้าที่ไม่ต่างกัน (2) ชุดทดสอบที่ต้องการศึกษา (Studied Subtest) ประกอบด้วย ข้อสอบที่สงสัยว่าทำหน้าที่ต่างกัน และข้อสอบชุดแรกมีจำนวน n ข้อ (ข้อที่ 1 ถึง n) และ แบบทดสอบชุดที่สองจะมีจำนวน $N-n$ ข้อ (ข้อที่ $n+1$ ถึง N) เมื่อ N เป็นจำนวนข้อสอบทั้งหมด พังก์ชันการตอบสนองข้อสอบของแบบทดสอบที่ต้องการจะศึกษาจากผู้เข้าสอบกลุ่มอ้างอิง และกลุ่มเปรียบเทียบ กำหนดในรูปฟังก์ชัน marginal ดังนี้

$$M_{SR}(\theta) = \sum_{i=n+1}^N M_{iR}(\theta) \quad (11)$$

$$M_{SF}(\theta) = \sum_{i=n+1}^N M_{iF}(\theta) \quad (12)$$

เมื่อ $M_{SR}(\theta)$ แทน ผลรวมของ Marginal IRFs ของข้อสอบที่ต้องการศึกษา จากผู้เข้าสอบกลุ่มอ้างอิง ณ ระดับความสามารถ θ

เมื่อ $M_{SF}(\theta)$ แทน ผลรวม Marginal IRFs ของข้อสอบที่ต้องการศึกษา จาก ผู้เข้าสอบกลุ่มเปรียบเทียบ ณ ระดับความสามารถ θ

ขนาดของความสามารถแตกต่างระหว่าง $M_{SR}(\theta)$ กับ $M_{SF}(\theta)$ แสดงถึงปริมาณของ การทดสอบของการทำหน้าที่ต่างกันของข้อสอบแบบเอกสาร หรือการทำหน้าที่ต่างกันแบบไม่ตัดกัน จากชุดแบบทดสอบที่ต้องการศึกษา ณ ระดับความสามารถ θ สามารถคำนวณในรูปการอินทิเกรท ดังนี้

$$\beta_{uni} = \int_{\theta} [M_{SR}(\theta) - M_{SF}(\theta)] f_p(\theta) d\theta \quad (13)$$

เมื่อ β_{uni} แทน ดัชนีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบแบบเอกสาร

$f_p(\theta)$ แทน พังก์ชันความหนาแน่นของโอกาสการแจกแจงความสามารถ ของผู้สอบ θ ทั้ง 2 กลุ่ม

ดัชนี β_{uni} ที่คำนวณได้จากสูตรดังกล่าวข้างต้น นำมาทดสอบสมมติฐาน ของการทำหน้าที่ต่างกันของข้อสอบแบบเอกสาร ดังนี้

$$H_0 : \beta_{uni} = 0 \quad (14)$$

$$H_1 : \beta_{uni} > 0 \quad (15)$$

สมมติฐานอื่น (H_1) ใช้ทดสอบการทำหน้าที่ต่างกันของข้อสอบที่มีลักษณะทิศทางเดียวกับที่เข้าข้างผู้เข้าสอบกลุ่มเปรียบเทียบ และค่าประมาณของ β_{uni} คำนวณได้จากการรวมของชุดแบบทดสอบที่มีความตรงและชุดแบบทดสอบที่ต้องการศึกษา ซึ่งกำหนดด้วยสัญลักษณ์ดังนี้

$$X = \sum_{i=1}^n U_i \quad (16)$$

$$Y = \sum_{i=n+1}^n U_i \quad (17)$$

เมื่อ X แทน คะแนนรวมของชุดแบบทดสอบที่มีความตรงใช้เป็นเกณฑ์ในการเปรียบเทียบ

Y แทน คะแนนรวมของชุดแบบทดสอบที่ต้องการศึกษา

U_i แทน ผลการตอบข้อสอบข้อที่ i (ตอบถูกต้องได้ 1 คะแนน และตอบผิดได้ 0 คะแนน)

นำผลคะแนนเฉลี่ยของการตอบข้อสอบที่ต้องการศึกษาระหว่างผู้เข้าสอบกลุ่มอ้างอิงกับกลุ่มเปรียบเทียบที่มีความสามารถระดับเดียวกันมาจับคู่เปรียบเทียบกัน จะพิจารณาจากคะแนนรวมที่เท่ากันของชุดแบบทดสอบที่มีความตรง ($X=k$) ด้วยสัญลักษณ์ดังนี้

$$\bar{Y}_{RK} - \bar{Y}_{FK} \quad ; K = 0, 1, 2, \dots, n \quad (18)$$

เมื่อ \bar{Y}_{RK} = ค่าเฉลี่ยของคะแนนรายข้อ จากชุดแบบทดสอบที่ต้องการของผู้เข้าสอบกลุ่มอ้างอิง ซึ่งได้คะแนน $X=k$

\bar{Y}_{FK} = ค่าเฉลี่ยของคะแนนรายข้อ จากชุดแบบทดสอบที่ต้องการศึกษาของผู้เข้าสอบกลุ่มเปรียบเทียบ ซึ่งได้คะแนน $X=k$

K = คะแนนรวมจากชุดแบบทดสอบที่มีความตรง

ค่า $\bar{Y}_{RK} - \bar{Y}_{FK}$ เป็นผลแตกต่างของผลการตอบข้อสอบในชุดแบบทดสอบที่ต้องการศึกษาระหว่างผู้เข้าสอบกลุ่มอ้างอิงกับกลุ่มเปรียบเทียบที่มีความสามารถระดับเดียวกัน

ถ้า $\bar{Y}_{RK} - \bar{Y}_{FK} = 0$ ทุกคะแนน K หมายถึงว่าข้อสอบที่ต้องการศึกษาทำหน้าที่ไม่ต่างกัน

และ ถ้า $\bar{Y}_{RK} - \bar{Y}_{FK} > 0$ ทุกคะแนน K หมายถึงข้อสอบทำหน้าที่ต่างกันแบบเอกสารโดยจะลำเอียงเข้าข้างผู้เข้าสอบกลุ่มอ้างอิง ค่าความแตกต่างของผลการตอบข้อสอบสามารถประมาณค่าในรูป β_{uni} ดังนี้

$$\hat{\beta}_{uni} = \sum_{K=0}^n \hat{P}_K (\bar{Y}_{RK} - \bar{Y}_{FK}) \quad (19)$$

$$\hat{P}_k = \frac{(J_{RK} + J_{FK})}{\sum_{k=0}^n (J_{RK} + J_{FK})} \quad (20)$$

เมื่อ P_k = สัดส่วนของจำนวนผู้เข้าสอบทั้งหมด(กลุ่มอ้างอิงและกลุ่มเปรียบเทียบ) ซึ่งตอบชุดแบบทดสอบที่มีความตรงแล้วได้คะแนนรวม $X = k$ จากจำนวนผู้สอบทั้งหมด

J_{FK} = จำนวนผู้เข้าสอบกลุ่มเปรียบเทียบซึ่งตอบชุดแบบทดสอบที่มีความตรงแล้วได้คะแนนรวม $X=k$

J_{RK} = จำนวนผู้เข้าสอบกลุ่มอ้างอิงซึ่งตอบชุดแบบทดสอบที่มีความตรงแล้วได้คะแนนรวม $X=k$ (ศิริชัย กาญจนวารี, 2555, หน้า 139-142)

สรุป คือ วิธีชิปเทสท์ใช้ในการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบเอกสาร มีข้อตกลงในแบบทดสอบ คือ ข้อสอบในแบบทดสอบจะต้องมุ่งวัดคุณลักษณะ หรือ ความสามารถแห่งลักษณะเดียวเท่านั้น ซึ่งความสามารถแห่งประเภทที่หนึ่ง เรียกว่า ความสามารถ เป้าหมายที่ต้องการวัด (Target Ability: Θ) แต่จะมีความสามารถแห่งอีกประเภทหนึ่งที่มืออธิผล ต่อการตอบข้อสอบเรียกว่า ความสามารถแห่งรากซ้อนที่ไม่ต้องการวัด และวิธีชิปเทสท์มีหลักการ ตรวจสอบการทำหน้าที่ต่างกันของข้อสอบแบบเอกสารโดยการเปรียบเทียบผลการตอบข้อสอบ ระหว่างกลุ่มอ้างอิงกับกลุ่มเปรียบเทียบ จะแบ่งแบบทดสอบออกเป็น 2 ชุด ย่อย (Subtest) คือ (1) ชุดแบบทดสอบที่มีความตรง (Valid Subtest) คือแบบทดสอบที่ใช้เป็นเกณฑ์ในการเปรียบเทียบ (Matching Subtest) โดยแบบทดสอบประกอบด้วยข้อสอบที่ทำหน้าที่ไม่ต่างกัน (2) ชุดทดสอบ ที่ต้องการศึกษา (Studied Subtest) ประกอบด้วยข้อสอบที่สงสัยว่าทำหน้าที่ต่างกัน และข้อสอบ ชุดแรกมีจำนวน n ข้อ (ข้อที่ 1 ถึง n) และแบบทดสอบชุดที่สองจะมีจำนวน $N-n$ ข้อ (ข้อที่ $n+1$ ถึง N) เมื่อ N เป็นจำนวนข้อสอบทั้งหมด

งานวิจัยที่เกี่ยวข้องกับการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบด้วยวิธีชิปเทสท์ มีดังนี้

พรัญญา สูงเนิน เสรี ชัดแฉม และสมโภชน์ วนกสุข (2552) ได้ศึกษาการทำหน้าที่ต่างกัน ของข้อสอบในแบบทดสอบพหุมิติเป็นการเปรียบเทียบระหว่างข้อสอบรายข้อกับหมวดข้อสอบ โดยใช้วิธีชิปเทสท์ ซึ่งมีวัตถุประสงค์ของศึกษาวิจัยเพื่อเปรียบเทียบผลการตรวจสอบการทำหน้าที่ ต่างกันของข้อสอบในแบบทดสอบพหุมิติ ระหว่างการตรวจสอบเป็นรายข้อกับรายหมวดข้อสอบ ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และ ขนาดใหญ่ (2,000 คน) กลุ่มประชากรที่ศึกษาเป็นนักเรียนชั้นประถมศึกษาปีที่ 6 สังกัดเขตพื้นที่ การศึกษานครศรีธรรมราช ปีการศึกษา 2546 ที่เข้าสอบวัดผลสัมฤทธิ์ระดับชาติรายวิชาภาษาไทย จำนวนที่เข้าสอบ 2,000 คน และจำนวนข้อสอบ 40 ข้อ จำแนกเป็น 2 หมวดข้อสอบ คือ หมวดที่ 1 วัดด้านโครงสร้างความรู้ จำนวน 15 ข้อ และหมวดที่ 2 วัดด้านกระบวนการ จำนวน 25 ข้อ จะใช้วิธีการ

สุ่มแบบแบ่งชั้น (Stratified Random Sampling) แบบจัดสรรเท่าเทียมกัน (Equal Allocation) โดยแบ่งชั้นตามระดับความสามารถเป็น 3 ระดับ คือ ดี พอดี และปรับปรุง ใช้นักเรียนเป็นหน่วย การสุ่ม สุ่มมาจำนวน 2,000 คน แบ่งเป็นเพศชาย จำนวน 1,000 คน และเพศหญิง จำนวน 1,000 คน การวิเคราะห์ข้อมูลค่าสถิติพื้นฐานจะใช้โปรแกรม SPSS ตรวจสอบความตรงเชิงโครงสร้าง ด้วยการวิเคราะห์องค์ประกอบเชิงยืนยันอันดับสอง โดยใช้โปรแกรม LISREL 8.50 และชั้นตอนต่อไป ก็จะเป็นการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบและหมวดข้อสอบโดยใช้โปรแกรม SIBTEST ผลการวิจัยพบว่า (1) ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก ขนาดกลาง และ ขนาดใหญ่ ระหว่างการตรวจสอบการทำหน้าที่ต่างกันเป็นรายข้อกับรายหมวดข้อสอบ พบร่วมกัน พบว่าข้อสอบทำหน้าที่ต่างกันแตกต่าง (2) การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบรายข้อ เมื่อกลุ่มตัวอย่างขนาดเล็ก พบร่วมกันของข้อสอบทำหน้าที่ต่างกันจำนวน 4 ข้อ คิดเป็นร้อยละ 10 ขนาดกลางพบร่วมกัน 13 ข้อ คิดเป็นร้อยละ 32.5 และขนาดใหญ่ พบร่วมกัน 15 ข้อ คิดเป็นร้อยละ 37.5 (3) การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบเป็นรายหมวดข้อสอบ เมื่อกลุ่มตัวอย่างขนาดเล็ก พบร่วมกันของข้อสอบทำหน้าที่ต่างกันจำนวน 4 ข้อ คิดเป็นร้อยละ 10 ขนาดกลางพบร่วมกัน 8 ข้อ คิดเป็นร้อยละ 20 ขนาดใหญ่ พบร่วมกัน 16 ข้อ คิดเป็นร้อยละ 40 และ (4) การตรวจสอบการทำหน้าที่ต่างกันของหมวดข้อสอบ ซึ่งให้เห็นว่า หมวดที่ 2 ภายใต้เงื่อนไขกลุ่มตัวอย่างขนาดกลาง มีนัยสำคัญทางสถิติ ที่ระดับ .50

ณรงค์ จันทร์มหา (2554) ได้ศึกษาเปรียบเทียบค่าความเที่ยงของแบบทดสอบผลสัมฤทธิ์ทางการเรียนที่มีจำนวนข้อสอบทำหน้าที่ต่างกันแตกต่างกันโดยมีวัตถุประสงค์เพื่อเปรียบเทียบค่าความเที่ยงของแบบทดสอบผลสัมฤทธิ์ทางการเรียนที่มีจำนวนข้อสอบทำหน้าที่ต่างกันแตกต่างกัน 7 เงื่อน คือ 0% 5% 10% 15% 20% 25% และ 30% กลุ่มตัวอย่างเป็นนักเรียนชั้นประถมศึกษาปีที่ 6 สังกัดสำนักงานคณะกรรมการการศึกษาขั้นพื้นฐาน ปีการศึกษา 2546 จำนวน 2,000 คน ที่ทำแบบทดสอบวัดผลสัมฤทธิ์ทางการเรียนระดับชาติ วิชาภาษาไทย จำนวน 40 ข้อ วิเคราะห์ค่าสถิติพื้นฐานโดยใช้โปรแกรม SPSS ตรวจสอบการทำหน้าที่ต่างกันของข้อสอบโดยใช้โปรแกรมซิปเหลท์ และทดสอบความแตกต่างของค่าความเที่ยงของแบบทดสอบ โดยใช้การทดสอบสถิติซี ผลการวิจัยปรากฏว่า ค่าความเที่ยงของแบบทดสอบที่มีจำนวนข้อสอบทำหน้าที่ต่างกัน 0% 5% 10% 15% 20% 25% และ 30% ไม่แตกต่างกัน

Apinyapibal, Lawthong, and Kanjanawasee (2015) ได้ศึกษาเรื่องการวิเคราะห์เปรียบเทียบประสิทธิภาพของรายการที่แตกต่างกัน การตรวจสอบการทำงานสำหรับ Dichotomously คะแนนรายการหมู่ วิธีโลจิสติกการทดสอบโดย วิธี SIBTEST และวิธีการ Raschtree ที่แตกต่างกันส่วนใหญ่รายการที่ทำงาน (DIF) วิธีการตรวจสอบที่ใช้ในรูปแบบ Rasch จะขึ้นอยู่กับการเปรียบเทียบของพารามิเตอร์ของการทดสอบการประเมินระหว่างสองกลุ่มหรือมากกว่า กลุ่มที่มีการกำหนดไว้ล่วงหน้า เช่นกลุ่มของผู้ชาย พบร่วมกัน ความหลากหลายของวิธีการทำงานสำหรับ ที่สามารถใช้ได้สำหรับการตรวจสอบรายการที่แตกต่างกันการทำงาน (DIF) ในรูปแบบ Rasch ส่วนใหญ่เหล่านี้วิธีการประกอบด้วยสองวิธี โดยวิธีแรกจะขึ้นอยู่กับการเปรียบเทียบของพารามิเตอร์รายการประมาณการกลุ่มที่กำหนดไว้ล่วงหน้าของอาสาสมัคร วิธีที่สองจะขึ้นอยู่กับการเปรียบเทียบระหว่างกลุ่มที่เป็นไปได้ทั้งหมดของอาสาสมัครโดยไม่คำนึงถึงตัวแปรคน วัตถุประสงค์ของงานวิจัยนี้

คือการเปรียบเทียบประสิทธิภาพของรายการที่แตกต่างกันในการทำงาน (DIF) การตรวจสอบระหว่าง 3 วิธี คือ (1) วิธีการทดสอบโดยโลจิสติกตามทฤษฎีการทดสอบคลาสสิก (2) วิธี SIBTEST บนพื้นฐานของทฤษฎีการตอบสนองข้อสอบและ (3) วิธี Raschtree ขั้นอยู่กับแบบที่ใช้การแบ่งส่วน การจำลองการตรวจหา DIF ในช่วงแรกของการศึกษาพบข้อดีและข้อเสียของทั้งสามวิธีการจากการบททวนวรรณกรรม ขั้นที่สองของการวิจัยเปรียบเทียบประสิทธิภาพของการตรวจสอบ DIF โดยใช้การจำลองข้อมูลประเภทของการการจำแนน DIF ความยาวทดสอบและขนาดของกลุ่มตัวอย่าง ซึ่งการวิจัยครั้งนี้ การวิจัยครั้งนี้เมื่อเทียบกับประสิทธิภาพของสามวิธีของการตรวจสอบ DIF ของรายการคะแนน Dichotomously วิธีการโลจิสติกทดสอบโดย วิธี SIBTEST และ วิธี Raschtree โดยตรวจสอบเอกสารงานวิจัยที่เกี่ยวข้องผลการวิจัยพบว่า (1) โลจิสติกการทดสอบที่มีประสิทธิภาพในการตรวจสอบหัวข้อสอบทำหน้าที่ต่างกันแบบเอกสาร (Uniform) และข้อสอบที่ทำหน้าที่ต่างกันแบบอนุกรูป (Non-Uniform DIF) และสามารถนำมาใช้ในการตรวจสอบคะแนน Polytomously DIF รูปแบบ การวิเคราะห์มีความยืดหยุ่นและง่ายต่อการใช้งาน ก็ยังสามารถนำมาใช้ในการตรวจสอบ DIF กับหลายกลุ่มสอบและสามารถใช้สถิติในการทดสอบความสำคัญของสมมติฐาน แต่มีข้อจำกัดที่มีความมั่นคงสถิติที่แตกต่างกันตามกลุ่มตัวอย่างและคะแนนรวมเป็นเกณฑ์ในการจับคู่ (2) วิธี SIBTEST สามารถนำมาใช้ในการตรวจสอบหัวข้อในเครื่องแบบและนอกเครื่องแบบ DIF และการตรวจสอบการทดสอบที่แตกต่างกันหัวที่ต่างกัน (DTF) และทำหน้าที่ต่างกัน (DBF) นอกจากนี้ยังสามารถนำมาใช้ในการตรวจสอบ Polytomously ฝ่าย DIF และการวิเคราะห์การทดสอบหัวข้อในมิติเดียวและหลายมิติ นั้นสามารถคำนวณได้ง่ายไม่ซับซ้อนประยัดและค่าใช้จ่ายและไม่ต้องใช้กลุ่มตัวอย่างที่มีขนาดใหญ่ นอกจากนี้จะใช้การทดสอบนัยสำคัญทางสถิติสมมติฐานแต่มีข้อจำกัด เกี่ยวกับความไวในการตรวจสอบที่ไม่ได้ทิศทางเดียว DIF (3) จุดที่โดดเด่นของวิธี Raschtree ที่สามารถตรวจจับหัวข้อสองกลุ่มที่กำหนดไว้ล่วงหน้าและແຜງ ผลการวิเคราะห์สามารถแสดงในรูปแบบของแผนภูมิที่ชัดเจน อย่างถึงกลุ่มรายการ และสามารถเข้าใจได้ง่ายโดยสอบถาม แต่มีข้อจำกัด คล้ายกับวิธีการทดสอบโดยโลจิสติกการศึกษาต่อไปจะแสดงการเปรียบเทียบประสิทธิภาพของการตรวจสอบ DIF จากข้อมูล จำลองโดยใช้การจำลองข้อมูลภายใต้หนึ่งพารามิเตอร์รูปแบบทฤษฎีการตอบสนองรายการข้อมูล จะถูกจัดให้สอดคล้องกับปัจจัยสี่ตัวแปร: สามประเภทของการตัวเลขสองของ DIF สองระดับความยาวของการทดสอบและการสาระดับของขนาดตัวอย่าง จะมีหัวหมด 36 เสื่อนไข ต้องมีการศึกษา (3 ประเภท \times 2 ขนาด \times 2 ขนาด \times 3 ขนาด) เพื่อทำให้การเปรียบเทียบที่ชัดเจนของ การรับรู้ความสามารถของสามวิธีการ

สรุปจากการศึกษางานวิจัยที่เกี่ยวข้อง วิธีซิปเทส์สามารถตรวจสอบการทำงานที่ต่างกันของข้อสอบในขนาดกลุ่มตัวอย่างที่มีขนาดเล็ก ขนาดกลาง และขนาดใหญ่ นั้น สามารถตรวจสอบพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบในกลุ่มตัวอย่างที่มีขนาดใหญ่ก็จะตรวจพบข้อสอบได้มากกว่ากลุ่มขนาดตัวอย่างขนาดเล็ก

ตอนที่ 6 การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบด้วยวิธีแมนเทล-แ昏ส์เซล และงานวิจัยที่เกี่ยวข้อง

วิธีแมนเทล-แ昏ส์เซล พัฒนาโดย (Camilli & Shepard, 1994; cited Mantalel-Haenszel, 1959) เดิมเป็นสถิติสำหรับเปรียบเทียบอัตราส่วนแต้มต่อร่วม และการทดสอบอัตราส่วนเปรียบเทียบด้วยโคลสแควร์ และชอนแลน (Holland, 1985; cited Holland & Thayer, 1988) ได้นำวิธีแมนเทล-แ昏ส์เซล ไปประยุกต์ในการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ และหน่วยงานบริการทดสอบทางการศึกษาของประเทศไทย วิธีแมนเทล-แ昏ส์เซล เป็นที่ยอมรับจากนักวิจัย เป็นวิธีที่ใช้ง่าย สะดวก และประหยัด โดยใช้หลักการของตารางการณ์จร แบบทฤษฎีการตอบข้อสอบ และไม่มีการคำนวนซ้ำ สามารถนำไปใช้ตรวจสอบการทำหน้าที่ต่างกันของข้อสอบได้ง่าย มีขั้นตอนการคำนวนที่ไม่ слับซับซ้อน เป็นการทดสอบทางสถิติแบบอนพารามetric ไม่จำเป็นต้องใช้โมเดลประมาณค่า

หลักการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบโดยวิธีแมนเทล-แ昏ส์เซล คือ จะนำเอาผลการตอบของข้อสอบระหว่างผู้สอบสองกลุ่มมาเปรียบเทียบกัน นั้น คือ กลุ่มเปรียบเทียบ กับกลุ่มอ้างอิง ซึ่งจะเปรียบเทียบทุกระดับความสามารถของผู้เข้าสอบกลุ่มย่อยสองกลุ่มที่มีระดับความสามารถเท่าเทียมกัน ในทางการปฏิบัติจะใช้คะแนนรวมของแบบทดสอบเป็นเกณฑ์การจับคู่ กลุ่มผู้สอบ

การวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบเป็นรายข้อ เมื่อจับคู่กลุ่มผู้สอบแล้ว จะนำข้อมูลผลการตอบข้อสอบนั้น ทั้งสองกลุ่มย่อยมาจัดลงตารางการณ์จรแบบ 2×2 (กลุ่มผู้สอบ 2 กลุ่ม x ผลการตอบ 2 แบบ) ตารางการณ์จร 1 ตารางแทนคะแนนรวม 1 ระดับ และถ้ามีคะแนนรวมของกลุ่มผู้เข้าสอบทั้งสิ้น k ระดับ จะต้องมีตารางการณ์ 2×2 ทั้งหมด k ตารางสำหรับตารางการณ์จรแบบ 2×2 ของข้อสอบแต่ละข้อที่มีคะแนนรวมระดับ j ดังตารางที่ 2-2

ตารางที่ 2-2 ผลการตอบข้อสอบข้อหนึ่งระหว่างผู้สอบกลุ่มอ้างอิงและกลุ่มเปรียบเทียบที่มีคะแนนรวม j

กลุ่ม	คะแนนผลของการตอบข้อสอบที่ศึกษา		รวม
	ตอบถูก (1)	ตอบผิด (0)	
R	A_j	B_j	nR_j
F	C_j	D_j	nF_j
รวม	m_{1j}	m_{0j}	T_j

- เมื่อ A_j = จำนวนผู้เข้าสอบกลุ่มอ้างอิงที่มีคะแนนรวม j ซึ่งตอบข้อสอบถูก
- B_j = จำนวนผู้เข้าสอบกลุ่มอ้างอิงที่มีคะแนนรวม j ซึ่งตอบข้อสอบผิด
- C_j = จำนวนผู้เข้าสอบกลุ่มเปรียบเทียบที่มีคะแนนรวม j ซึ่งตอบข้อสอบถูก
- D_j = จำนวนผู้เข้าสอบกลุ่มเปรียบเทียบที่มีคะแนนรวม j ซึ่งตอบข้อสอบผิด
- m_{1j} = จำนวนผู้เข้าสอบทั้งหมดที่มีคะแนนรวม j ซึ่งตอบข้อสอบถูก

- m_{0j} = จำนวนผู้เข้าสอบทั้งหมดที่มีคะแนนรวม j ซึ่งตอบข้อสอบผิด
 n_{Rj} = จำนวนผู้เข้าสอบกลุ่มอ้างอิงที่มีคะแนนรวม j
 n_{Fj} = จำนวนผู้เข้าสอบกลุ่มเปรียบเทียบที่มีคะแนนรวม j
 T_j = จำนวนผู้เข้าสอบทั้งหมดที่มีคะแนนรวม j

จากตารางที่ 2-2 นำมาคำนวณสัดส่วนของผลการตอบข้อสอบถูกและผิด ดังตารางที่ 2-3

ตารางที่ 2-3 สัดส่วนของผลการตอบข้อสอบข้อหนึ่งระหว่างผู้สอบกลุ่มอ้างอิงและกลุ่มเปรียบเทียบที่มีคะแนน j

กลุ่ม	คะแนนผลของการตอบข้อสอบที่ศึกษา		รวม
	ตอบถูก (1)	ตอบผิด (0)	
R	pR_j	qR_j	1.00
F	pF_j	qF_j	1.00

- เมื่อ pR_j = สัดส่วนของผู้สอบกลุ่มอ้างอิงที่มีคะแนนรวม j ซึ่งตอบข้อสอบถูก
 qR_j = สัดส่วนของผู้สอบกลุ่มอ้างอิงที่มีคะแนนรวม j ซึ่งตอบข้อสอบผิด
 pF_j = สัดส่วนของผู้สอบกลุ่มเปรียบเทียบที่มีคะแนนรวม j ซึ่งตอบข้อสอบถูก
 qF_j = สัดส่วนของผู้สอบกลุ่มเปรียบเทียบที่มีคะแนนรวม j ซึ่งตอบข้อสอบผิด

ในการทดสอบสมมติฐานของการทำหน้าที่ต่างกันของข้อสอบจะต้องทดสอบสมมติฐานศูนย์ (H_0) และสมมติฐานอื่น (H_1) ดังนี้

$$H_0: \frac{pR_j}{qR_j} = \frac{pF_j}{qF_j} \quad ; j = 1, 2, 3, \dots, k \quad (21)$$

$$H_1: \frac{pR_j}{qR_j} = \alpha \frac{pF_j}{qF_j} \quad ; j = 1, 2, 3, \dots, k \text{ เมื่อ } \alpha \neq 1 \quad (22)$$

สมมติฐานศูนย์ตั้งขึ้นเป็นอิสระอย่างมีเงื่อนไขระหว่างกลุ่มผู้สอบและคะแนนผลการตอบข้อสอบที่ศึกษา ดังนั้นสัดส่วนคะแนนที่ได้จากการที่ 3.2 ภายใต้สมมติฐานศูนย์ที่ตั้งไว้สามารถคำนวณเป็นค่าคาดหมาย (Expected values) ในแต่ละเซลล์ได้ดังนี้

$$E(A_{ij}) = \frac{nR_j m_{0j}}{T_j} \quad (23)$$

$$E(B_{ij}) = \frac{nR_j m_{0j}}{T_j} \quad (24)$$

$$E(C_j) = \frac{nF_j m_{1j}}{T_j} \quad (25)$$

$$E(D_j) = \frac{nF_j m_{0j}}{T_j} \quad (26)$$

สำหรับพารามิเตอร์ α ภายใต้สมมติฐานอื่นเรียกว่า อัตราส่วนแต้มต่อร่วม (Common oddsratio) ซึ่งคำนวณได้จาก

$$\alpha = \frac{\frac{pR_j}{qR_j}}{\frac{pF_j}{qF_j}} = \frac{pR_j qF_j}{qR_j pF_j} \quad (27)$$

ถ้า $\alpha = 1$ แสดงโอกาสของการตอบข้อสอบระหว่างผู้สอบทั้งสองกลุ่ม มีค่าเท่ากัน

ถ้า $\alpha > 1$ แสดงว่าผู้สอบกลุ่มอ้างอิงมีโอกาสตอบข้อสอบถูกมากกว่าผู้สอบกลุ่ม

เปรียบเทียบ

ถ้า $\alpha < 1$ แสดงว่าผู้สอบกลุ่มเปรียบเทียบมีโอกาสตอบข้อสอบถูกมากกว่าผู้สอบกลุ่มอ้างอิง

แมนเทล-แฮนเซล ได้เสนอวิธีประมาณค่า α จากตารางแบบ 2×2 จำนวน k ระดับ ดังนี้

$$\hat{\alpha}_{MH} = \frac{\sum_{j=1}^k A_j D_j / T_j}{\sum_{j=1}^k B_j C_j / T_j} \quad (28)$$

$\hat{\alpha}_{MH}$ เป็นค่าประมาณขนาดอิทธิพลของการทำหน้าที่ต่างกันของข้อสอบ (DIF effect size) ซึ่งมีค่าอยู่ระหว่าง 0 ถึง ∞ ยอลแลนด์และ雷耶อร์ (1985 cited in Holland ; & Thayer, 1988) ได้เสนอให้แบ่ง $\hat{\alpha}_{MH}$ ค่าเป็นสเกลมาตราฐานเดลต้า (Delta scale ; $MH_{D - DIF}$ หรือ MH_{DIF}) ดังนี้

$$MH_{DIF} = -2.35 \ln(\alpha_{MH}) \quad (29)$$

ค่า MH_{DIF} ดักกล่าวสามารถนำไปพิจารณาค่าความยากของข้อสอบ กล่าวคือ ถ้า MH_{DIF} มีค่าเป็นศูนย์ แสดงว่าข้อสอบนั้นมีความยากสำหรับแต่ละกลุ่มเท่ากัน ถ้า MH_{DIF} มีค่าเป็นลบ แสดงว่าข้อสอบยากสำหรับผู้สอบกลุ่มเปรียบเทียบมากกว่ากลุ่มอ้างอิง และถ้า MH_{DIF} มีค่าเป็นบวก แสดงว่าแสดงว่าข้อสอบยากสำหรับผู้สอบกลุ่มอ้างอิงมากกว่ากลุ่มเปรียบเทียบ ส่วนการประมาณค่าความคลาดเคลื่อนมาตราฐานของ MH_{DIF} สามารถคำนวณได้จากสูตร ดังนี้

$$SE(MH_{DIF}) = 2.35 \sqrt{\text{var}[\ln(\hat{\alpha}_{MH})]} \quad (30)$$

$$\text{โดยที่ } \text{Var}[\ln(\hat{\alpha}_{MH})] = \frac{\sum_{j=1}^k U_j V_j / T_j^2}{2 \left[\sum_{j=1}^k A_j D_j / T_j \right]^2} \quad (31)$$

$$\text{ขณะที่ } U_i = A_j D_j + \alpha_{MH} (B_j C_j)$$

$$\text{และ } V_j = (A_j + D_j) + \alpha_{MH} (B_j + C_j)$$

ในการทดสอบนัยสำคัญของสมมติฐาน จะต้องนำค่า α_{MH} หรือค่า MHDF ไปทดสอบ กับสถิติเมเนเทล-แ昏ส์เซลล์-สแควร์ (χ^2_{MH}) ที่ระดับชั้นความเป็นอิสระเท่ากับ 1 ($df = 1$) สถิติ χ^2_{MH} มีสูตรในการคำนวณดังนี้

$$\chi^2_{HM} = \frac{\left[\left| \sum_{j=1}^k A_j - E(A_j) \right| - 0.5 \right]^2}{\sum_{j=1}^k \text{Var}(A_j)} \quad (32)$$

$$\text{โดยที่ } E(A_j) = \frac{nR_j m_{1j}}{T_j} \quad (33)$$

$$\text{Var}(A_j) = \frac{nR_j nF_j m_{1j} m_{0j}}{T_j^2 (T_j - 1)} \quad (34)$$

เมื่อ $E(A_j)$ = ค่าคาดหมายของจำนวนผู้สอบกลุ่มอ้างอิงที่ระดับคะแนน j ซึ่งตอบ ข้อสอบถูก

$\text{Var}(A_j)$ = ค่าความแปรปรวนของจำนวนผู้สอบกลุ่มอ้างอิงที่ระดับคะแนน j ซึ่งตอบข้อสอบถูก

สำหรับเกณฑ์ในการตัดสินข้อสอบที่ทำหน้าที่ต่างกัน คือ ข้อสอบที่มีค่าแตกต่าง จาก 1 อย่างมีนัยสำคัญทางสถิติหรือค่าแตกต่างจาก 0 อย่างมีนัยสำคัญทางสถิติ โดยมีเกณฑ์ ในการพิจารณา ดังนี้

สำหรับค่าแปลผลการทดสอบได้ดังนี้

1) $c = 1$ หรือไม่แตกต่างจาก 1 อย่างมีนัยสำคัญ และแสดงว่าข้อสอบนั้นทำหน้าที่ ไม่แตกต่างกันระหว่างกลุ่ม (No DIF)

2) $c > 1$ และแสดงว่าข้อสอบนั้นทำหน้าที่แตกต่างกันระหว่างกลุ่มโดยจะเข้าข้างกลุ่ม อ้างอิง

3) ค่า < 1 แสดงว่าข้อสอบนั้นทำหน้าที่แตกต่างกันระหว่างกลุ่มโดยจะเข้าข้างกลุ่มเปรียบเทียบ

ส่วนค่า MH DIF แปลผลการทดสอบได้ดังนี้

1) ค่า MH DIF เท่ากับ 0 หรือไม่แตกต่างจาก 0 อย่างมีนัยสำคัญทางสถิติ แสดงว่าข้อสอบนั้นทำหน้าที่ไม่แตกต่างกันระหว่างกลุ่ม (No DIF)

2) ค่า MH DIF แตกต่างจาก 0 อย่างมีนัยสำคัญทางสถิติ และมีค่าเป็นบวก (positive) แสดงว่าข้อสอบนั้นทำหน้าที่แตกต่างกันระหว่างกลุ่ม โดยจะเข้าข้างกลุ่มเปรียบเทียบ (F)

3) ค่า MH DIF แตกต่างจาก 0 อย่างมีนัยสำคัญทางสถิติ และมีค่าเป็นลบ (negative) แสดงว่าข้อสอบนั้นทำหน้าที่ไม่แตกต่างกันระหว่างกลุ่ม โดยจะเข้าข้างกลุ่มอ้างอิง (R) นอกจากนี้ขนาดของ $|MH DIF|$ สามารถนำไปใช้แปลผลถึงระดับของการทำหน้าที่แตกต่างกันของข้อสอบได้ ถ้า $0 < |MH DIF| < 1.00$ แสดงว่า ข้อสอบทำหน้าที่แตกต่างกันเล็กน้อย ถ้า $1.00 \leq |MH DIF| \leq 1.50$ แสดงว่า ข้อสอบทำหน้าที่แตกต่างกันปานกลาง แต่ถ้า $|MH DIF| > 1.50$ แสดงว่าข้อสอบทำหน้าที่แตกต่างกันมาก (Zieky, 1993)

การวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบโดยวิธี MH นี้ใช้คะแนนรวมของแบบทดสอบเป็นเกณฑ์การจับคู่ ซึ่งมีจุดอ่อนในด้านความไม่เป็นอิสระของคะแนนรวมกับคะแนนรายข้อที่ทำการศึกษาอยลalienและเทเยอร์ (Holland & Thayer, 1988) ได้เสนอวิธีแก้จุดอ่อนดังกล่าว เพื่อทำให้เกณฑ์การจับคู่ผู้สอบบรรห่วงกลุ่มมีความบริสุทธิ์ยิ่งขึ้น โดยใช้วิธีการ 2 ขั้นตอน ดังนี้

ขั้นตอนแรก นำคะแนนรวมของแบบสอบทั้งฉบับเป็นเกณฑ์การจับคู่ผู้สอบบรรห่วงกลุ่มย่อย 2 กลุ่ม และวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบเป็นรายข้อ เมื่อพบว่าข้อสอบข้อใดทำหน้าที่ต่างกันให้นำคะแนนของข้อสอบข้อนั้นออกจากคะแนนรวมของผู้สอบแต่ละคน

ขั้นตอนที่สอง ใช้คะแนนรวมของแบบสอบที่นำเอาคะแนนข้อสอบที่ทำหน้าที่ต่างกันซึ่งตรวจพบในขั้นตอนแรกออกไป เพื่อใช้เป็นเกณฑ์การจับคู่แล้ววิเคราะห์การทำหน้าที่ต่างกันของข้อสอบข้าวีกครั้งหนึ่ง สำหรับนำไปใช้สรุปผลการตรวจสอบ

ที่ผ่านมา มีผู้นำวิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ วิธี MH มาศึกษา เช่น คลอสเซอร์และคณะ (Clauzer et al., 1991) ใช้วิธี MH วิเคราะห์การทำหน้าที่ต่างกันของข้อสอบ อย่างสมำเสมอ พบร่วมกับวิธี MH จะใช้ได้ผลดีในกรณีที่ข้อสอบมีค่าอำนาจจำแนกสูง แต่จะไม่สามารถวิเคราะห์ข้อสอบที่มีค่าความยากมากได้ ซึ่งสอดคล้องกับค้นพบทองเมเชอร์และคณะ (Mazor et al., 1991) ชุคเว็คส์ และทอลแมน (Sudweeks & Tolman, 1990) ที่พบว่าข้อสอบที่ทำหน้าที่ต่างกันมักเป็นข้อสอบที่ยาก บากิและเพอรารา (Baghi & Ferrara, 1990) พบร่วม เมื่อใช้กลุ่มตัวอย่างขนาด 750 คนขึ้นไป วิธี MH ใช้แทนวิธีทฤษฎีการตอบข้อสอบ 3 พารามิเตอร์ สามมินารานและโรเจอร์ (Swaminathan & Rogers, 1990) ได้ศึกษาเอาข้อมูลจำลองพบร่วมวิธี MH วิเคราะห์ได้ดีกว่าการทดลองแบบโลจิสติกน้อยโดยตรวจค้นได้ถูกต้องร้อยละ 75 กรณีใช้กลุ่มตัวอย่างจำนวน 250 คน และตรวจค้นได้ถูกต้องร้อยละ 100 กรณีกลุ่มตัวอย่างจำนวน 500 คน กรณีที่การทำหน้าที่ต่างกันของข้อสอบอย่างสมำเสมอ และกรณีที่ไม่สมำเสมอที่ติดกันปลายข้างใดข้างหนึ่ง แต่วิธี MH มีค่าใช้จ่ายย่ำຍน้อยกว่าวิธีโลจิสติกประมาณ 3-4 เท่า แย่และเบลตันและคณะ (Hambleton et al., 1986) พบร่วม MH ให้ค่าใกล้เคียงกับทฤษฎีการตอบข้อสอบ ทั้งที่ใช้ค่าความแตกต่างของค่าเฉลี่ยกำลัง

สองและการตรวจสอบความแตกต่างของพื้นที่รวมใต้โค้ง แต่ไวริ MH มีค่าใช้จ่ายต่ำกว่าและใช้เวลาน้อยกว่าลินเครอร์ (Linacre, 1988) พบว่า ใช้ได้ดีในสถานการณ์จำลองทุกสถานการณ์พอ ๆ กับโปรแกรม PROX ของราสช์ (Rasch) แต่ควรใช้เกณฑ์ในการคำนวณการทำหน้าที่ต่างกันของข้อสอบและความคลาดเคลื่อนหลอย ๆ เกณฑ์ และใช้ตัวประมาณค่าความคลาดเคลื่อนมาตรฐานมากกว่า 1 ตัว เพิร์ลแมนและคณะ (Perlman et al., 1988) พบว่า ไวริ MH มีปัญหาด้านความเชื่อมั่น เมื่อจำนวนกลุ่มตัวอย่างน้อยกว่า 660 คน ริสเซน สไตน์เบอร์ก และ瓦因เบอร์ (Thissen, Steinberg & Wainer, 1988) พบว่า ไวริ MH ให้ผลการวิเคราะห์คล้ายกับวิธีทฤษฎีการตอบข้อสอบแบบการทดสอบ เชิงเส้น และอาจใช้วิเคราะห์การทำหน้าที่ต่างกันของข้อสอบก่อนใช้วิธีทฤษฎีการตอบข้อสอบแบบการทดสอบเชิงเส้น (ศิริชัย กาญจนวารี, 2555, หน้า 126-131)

สรุป คือ การทำหน้าที่ต่างกันของข้อสอบโดยวิธีแมนเทล-แบรนส์เซล คือจะนำเอาผลการตอบของข้อสอบระหว่างผู้สอบสองกลุ่มมาเปรียบเทียบกัน นั้น คือ กลุ่มเปรียบเทียบกับกลุ่มอ้างอิง ซึ่งจะเปรียบเทียบทุกรอบตับความสามารถของผู้เข้าสอบกลุ่มย่อยสองกลุ่มที่มีระดับความสามารถเท่าเทียมกัน ซึ่งจะใช้คะแนนรวมของแบบทดสอบเป็นเกณฑ์การจับคู่กลุ่มผู้สอบ

งานวิจัยที่เกี่ยวข้องการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบด้วยวิธีแมนเทล-แบรนส์เซล มีดังนี้

ยกเวียรติกมล ทองออก โซติกา ภานี และศิริชัย กาญจนวารี (2556) ได้ศึกษาประสิทธิภาพการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบภายใต้วิธีทดสอบโดยโลจิสติกสำหรับข้อสอบที่ตรวจให้คะแนนแบบทวิภาค เปรียบเทียบขนาดอิทธิพลสองเกณฑ์ ที่จะเป็นการสอบเพื่อสรุปบุคคลจากคะแนนที่ได้จากการเรียนรู้ที่ต้องดำเนินด้วยความยุติธรรม ซึ่งมีวัตถุประสงค์ที่จะศึกษา คือ 1) เพื่อเปรียบเทียบอัตราความถูกต้องและอัตราความคลาดเคลื่อนประเภทที่ 1 ในวิธีทดสอบโดยโลจิสติกระหว่างการวัดขนาดอิทธิพลตามเกณฑ์ Jodoin and Gierl และการวัดขนาดอิทธิพลตามเกณฑ์ Zumbo and Thomas ภายใต้เงื่อนไขเดียวกันของปัจจัยที่แปรเปลี่ยน 4 ปัจจัยและปฏิสัมพันธ์สองทางระหว่างปัจจัยที่แปรเปลี่ยน 4 ปัจจัย 2) เพื่อเปรียบเทียบอัตราความถูกต้องและอัตราความคลาดเคลื่อนประเภทที่ 1 ในวิธีทดสอบโดยโลจิสติก ด้วยการวัดขนาดอิทธิพลตามเกณฑ์ Jodoin and Gierl และการวัดขนาดอิทธิพลตามเกณฑ์ Zumbo and Thomas ภายใต้เงื่อนไขต่างกันของปัจจัยที่แปรเปลี่ยน 4 ปัจจัยและปฏิสัมพันธ์สองทางระหว่างปัจจัยที่แปรเปลี่ยน 4 ปัจจัย 3) เพื่อเปรียบเทียบอัตราความถูกต้องและมีวิธีการวิจัย ดังนี้ คือ ตัวแปรที่ศึกษาตัวแปรอิสระ (1) เกณฑ์ขนาดอิทธิพลเกณฑ์ Jodoin, Gierl, Zumbo, and Thomas (2) รูปแบบของข้อสอบที่ทำหน้าที่ต่างกันแบบเอกสารและแบบออนไลน์ (3) ขนาดของการทำหน้าที่ต่างกันขนาด 0.1 0.2 และ 0.4 (4) จำนวนข้อสอบที่ทำหน้าที่ต่างกันทั้งฉบับคิดเป็นร้อยละ 10 และ 20 และ (5) ความยาวของแบบสอบถามทั้งฉบับ 40 และ 50 ข้อ ส่วนตัวแปรตาม (1) อัตราความถูกต้องคำนวณจากจำนวนของข้อสอบที่ตรวจสอบได้ถูกต้อง ว่าทำหน้าที่ต่างกันต่อจำนวนข้อสอบที่ทำหน้าที่ต่างกันทั้งหมดในแบบสอบถามคำนวณเป็นค่าร้อยละ (2) อัตราความคลาดเคลื่อนประเภทที่ 1 เป็นการระบุผิดพลาดว่าข้อสอบทำหน้าที่ต่างกัน ซึ่งในความจริงข้อสอบทำหน้าที่ไม่ต่างกันคำนวณจากสัดส่วนของจำนวนข้อสอบที่ตรวจสอบผิดพลาดว่าทำหน้าที่ต่างกันทั้งที่ในความเป็นจริงข้อสอบไม่ได้ทำหน้าที่ต่างกันคำนวณเป็นค่าร้อยละและใช้วิเคราะห์ข้อมูลจะเป็นการวิเคราะห์ค่าพารามิเตอร์ซึ่งไม่เดลที่ใช้ทฤษฎีการตอบสนองของข้อสอบด้วยโปรแกรม

MULTILOG-MG พัฒนาขึ้นโดย DAVID Thissen ที่มีคุณสมบัติในการประมาณค่าพารามิเตอร์ของข้อสอบและผู้สอบ ความเหมาะสมสำหรับการตอบสนองที่มีการให้คะแนนแบบทวิภาค วิเคราะห์คุณภาพข้อสอบด้วยโปรแกรม SPSS และวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบด้วยวิธีถดถอยโลจิสติก โดยมีการเขียนคำสั่งประมวลผลด้วยโปรแกรม R จากนั้นก็นำมาเปรียบเทียบอัตราความถูกต้องและอัตราความคลาดเคลื่อนประเภทที่ 1 ที่ภายใต้ตัวแปรอิสระ โดยสถิติการวิเคราะห์ความแปรปรวนหลายตัวแปร (MANOVA) ที่ระดับนัยสำคัญ .01 จะกำหนดการวิเคราะห์ให้มีปฏิสัมพันธ์ระหว่างตัวแปรอิสระไม่เกินอันดับสอง

ผลการวิจัยปรากฏว่า 1) ประสิทธิภาพการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบภายใต้วิธีถดถอยโลจิสติก ระหว่างการวัดขนาดอิทธิพลตามเกณฑ์ Jodoin, Gierl, Zumbo และ Thomas ภายใต้การจำลองตามปัจจัยที่แปรเปลี่ยนและศึกษาปฏิสัมพันธ์สองทางระหว่างปัจจัยที่แปรเปลี่ยน 4 ปัจจัย ในการวัดขนาดอิทธิพลตามเกณฑ์ Jodoin และ Gierl มีความถูกต้องสูงกว่าเกณฑ์ Zumbo และ Thomas เกือบทุกเงื่อนไข คือ (1.1) รูปแบบของข้อสอบที่ทำหน้าที่ต่างกัน และขนาดของการทำหน้าที่ต่างกัน มีผลต่ออัตราความถูกต้องของการวัดขนาดอิทธิพลตามเกณฑ์ทั้งสองเกณฑ์ (1.2) ความยาวของแบบสอบถามทั้งฉบับ ไม่มีผลต่ออัตราความถูกต้องและ TE เนื่องจากกำหนดความยาวของแบบสอบถามใกล้เคียงกันมากระหว่าง 40 กับ 50 ข้อทำให้มีความแตกต่างกัน และ (1.3) ปฏิสัมพันธ์สองทางของปัจจัยที่แปรเปลี่ยนระหว่างรูปแบบของการทำหน้าที่ต่างกันกับขนาดของการทำหน้าที่ต่างกันปฏิสัมพันธ์สองทาง ระหว่างรูปแบบของข้อสอบที่ทำหน้าที่ต่างกันกับจำนวนข้อสอบที่ทำหน้าที่ต่างกัน และปฏิสัมพันธ์สองทางระหว่างจำนวนข้อสอบที่ทำหน้าที่ต่างกันกับขนาดของการทำหน้าที่ต่างกัน มีผลต่ออัตราความถูกต้องของการวัดขนาดอิทธิพลตามเกณฑ์ทั้งสองเกณฑ์ 2) ประสิทธิภาพการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบของวิธีโลจิสติก โดยการวัดขนาดอิทธิพลของการทำหน้าที่ต่างกัน ปัจจัยที่ศึกษาพบว่า (2.1) รูปแบบของข้อสอบที่ทำหน้าที่ต่างกันแบบองค์กรูปมีอัตราความถูกต้องสูงกว่าแบบเอกสารูป รูปแบบองค์กรูปตามเกณฑ์ Jodoin and Gierl ให้ค่า TE ต่ำกว่าเกณฑ์ Zumbo and Thomas (2.2) ขนาดของการทำหน้าที่ต่างกัน มีผลต่อประสิทธิภาพการตรวจสอบระหว่างการวัดขนาดอิทธิพลตามเกณฑ์ Jodoin and Gierl และเกณฑ์ Zumbo and Thomas และ (2.3) จำนวนข้อสอบที่ทำหน้าที่ต่างกันและความยาวของแบบสอบถามทั้งฉบับมีผลต่อประสิทธิภาพการตรวจสอบขนาดอิทธิพลทั้งสองเกณฑ์ อย่างมีนัยสำคัญทางสถิติที่ระดับ .001

สุภารณ์ แดงเพ็ง ชิดชนก เชิงเชาว์ และบุญญิสา แซ่หล่อ (2554) ได้ศึกษาการเปรียบเทียบความลำเอียงของแบบทดสอบคณิตศาสตร์ในการประเมินภาพการศึกษา ระดับห้องเรียนขั้นประถมศึกษาปีที่ 5 จังหวัดปัตตานี ระหว่างวิธีแมนเทล-เ xen'szel และวิธีโค้งลักษณะข้อสอบ 3 พารามิเตอร์ โดยมีวัตถุประสงค์ของ การวิจัย 1) เปรียบเทียบจำนวนข้อของแบบทดสอบที่ตรวจพบความลำเอียงจากการวิเคราะห์ระหว่างวิธีแมนเทล-เ xen'szel และวิธีโค้งลักษณะข้อสอบ 3 พารามิเตอร์ 2) ศึกษาความสอดคล้องของดัชนีความลำเอียงจากการวิเคราะห์ระหว่างวิธีแมนเทล- xen'szel และวิธีโค้งลักษณะข้อสอบ 3 พารามิเตอร์ และ 3) เปรียบเทียบความซื่อมั่นของแบบทดสอบหลังจากคัดเลือกข้อสอบข้อที่มีความลำเอียงอย่างระหว่างวิธีแมนเทล- xen'szel และวิธีโค้งลักษณะข้อสอบ 3 พารามิเตอร์ พぶว่า จำนวนข้อสอบที่มีความลำเอียงในแบบทดสอบใน

กลุ่มสาระการเรียนรู้คณิตศาสตร์ เมื่อวิเคราะห์ความลำเอียงด้วยวิธี 2 วิธี คือ วิธีแมนเทล-แ昏ส์เซล และวิธีโคิงลักษณะข้อสอบ 3 พารามิเตอร์ ซึ่งวิธีแมนเทล-แ昏ส์เซลตรวจพบข้อสอบที่ทำหน้าที่ต่างกันจำนวน 15 ข้อ มีความลำเอียงเข้าหากลุ่มอ้างอิงจำนวน 7 ข้อ และลำเอียงเข้าหากลุ่มเปรียบเทียบจำนวน 8 ข้อ ส่วนวิธีโคิงลักษณะข้อสอบ 3 พารามิเตอร์ จะวัดความแตกต่างค่าพารามิเตอร์ความยาก พบว่า มีข้อสอบที่มีความลำเอียงจำนวน 7 ข้อ ลำเอียงเข้าหากลุ่มอ้างอิงจำนวน 4 ข้อ และลำเอียงเข้าหากลุ่มเปรียบเทียบจำนวน 3 ข้อ จำนวนข้อสอบที่มีความลำเอียงจากการวิเคราะห์ ด้วยวิธีแมนเทล-แ昏ส์เซลและวิธีโคิงลักษณะข้อสอบ มีความแตกต่างกันอย่างมีนัยสำคัญทางสถิติที่ระดับ.05 และค่าสัมประสิทธิ์สหสัมพันธ์ระหว่างดัชนีความลำเอียงของข้อสอบจากการวิเคราะห์ ด้วยวิธีแมนเทล-แ昏ส์เซล และวิธีโคิงลักษณะข้อสอบ 3 พารามิเตอร์ มีค่าพารามิเตอร์ มีค่าเท่ากับ 0.846 ซึ่งอยู่ในระดับสูง และเมื่อทดสอบนัยสำคัญทางสถิติ พบร่วมดัชนีความลำเอียงจากการวิเคราะห์ด้วยทั้ง 2 วิธี มีความสอดคล้องกันสูงอย่างมีนัยสำคัญทางสถิติ ที่ระดับ .05 ส่วนค่าความเชื่อมั่นของแบบทดสอบก่อนการคัดเลือกข้อสอบที่มีความลำเอียงออก มีค่าความเชื่อมั่นเท่า 0.646 ค่าความยากเฉลี่ย เท่ากับ 0.318 และค่าอำนาจจำแนกเฉลี่ยเท่ากับ 0.164 และหลังจากการคัดเลือกข้อสอบที่มีความลำเอียงออก ด้วยวิธีแมนเทล-แ昏ส์เซล เป็นจำนวน 15 ข้อ เมื่อนำไปทดสอบด้วยทฤษฎีการทดสอบดังเดิม พบว่า ข้อสอบที่เหลือ จำนวน 25 ข้อ มีค่าความเชื่อมั่นเท่ากับ 0.462 ค่าความเชื่อมั่นลดลง .184 ถือว่ายังอยู่ในระดับปานกลางมีค่าความยากเฉลี่ยเท่ากับ 0.318 ค่าอำนาจจำแนกเฉลี่ยเท่ากับ 0.121 และเมื่อวิเคราะห์ความลำเอียงข้อสอบด้วยวิธีโคิงลักษณะข้อสอบ 3 พารามิเตอร์ แล้วเลือกข้อสอบที่มีความลำเอียงจำนวน 7 ข้อ ออกเมื่อประเมินคุณภาพของแบบทดสอบโดยทฤษฎีการทดสอบแบบดังเดิม พบว่า ข้อสอบที่เหลือ จำนวน 33 ข้อ ซึ่งค่าความเชื่อมั่นอยู่ในระดับปานกลางเท่ากับ 0.561 ซึ่งลดลงเพียงเล็กน้อย คือ ลดลง .085 มีค่าความยากเพิ่มขึ้น .008 คือ เท่ากับ 0.310 และค่าอำนาจจำแนกเฉลี่ย 0.147 ข้อสอบที่ลำเอียงทั้ง 2 วิธี จำนวน 7 ข้อ คือ ข้อที่ 7, 16, 19, และ 28 จะลำเอียงเข้าหากลุ่มผู้สอบที่ใช้ภาษาไทยในชีวิตประจำวันเป็นข้อสอบเกี่ยวกับตัวประกอบของจำนวน คุณสมบัติของรูปเรขาคณิต ทรงเรขาคณิต และโจทย์ปัญหาการค้าขาย และข้อที่ 21, 34 และ 40 ลำเอียงเข้าหากลุ่มผู้สอบที่ใช้ภาษาลາວ/กຳມົງກອນ ปัจจุบัน เป็นข้อสอบที่เกี่ยวกับการเรียนลำดับเหตุการณ์ การใช้แผนภูมิ ร้อยละ ใช้สัญลักษณ์และเน้นตัวเลข

Awuor (2008) ได้ศึกษาเรื่องผลของขนาดตัวอย่างไม่เท่ากันในพลังของการตรวจสอบ DIF กับ วิธี IRT-Based ศึกษา Monte Carlo กับวิธี SIBTEST และวิธี Mantel-昏斯เซล พบร่วมการจำลองการศึกษานี้มุ่งเน้นไปที่การกำหนดผลของขนาดตัวอย่างที่ไม่เท่ากันในอัตราความถูกต้องทางสถิติของวิธี SIBTEST และวิธี Mantel-Haenzsel ใน การตรวจหา DIF ของขนาดปานกลางและขนาดใหญ่ ข้อสอบเป็นแบบ พารามิเตอร์โดยสใช้แบบ 2 PLM โดยใช้โปรแกรม WinGen2 (Han, 2006) ถูกใช้ในการจำลองการประมาณความสามารถและการสร้างข้อมูลการตอบสนองที่ถูกวิเคราะห์โดย วิธี SIBTEST ขั้นตอนวิธี SIBTEST กับกระบวนการการการทดสอบโดยการตรวจความถูกต้องใช้ใน การคำนวณสถิติ DIF คือขนาดผล DIF และสถิติ ความสำคัญของการมีลำเอียง วิธี SIBTEST เป็นวิธีที่ถูกนำมาใช้ในการคำนวณสถิติ DIF สำหรับ ขั้นตอนวิธี Mantel-Haenzsel โดยมีเงื่อนไขในสิ่งแวดล้อมที่พารามิเตอร์สามารถถูกจำลองข้อมูลการตอบสนองและการดำเนินการสร้าง DIF

วิเคราะห์ ข้อสอบที่ใช้ทดสอบเพื่อตรวจสอบพบร่วมกับรายการที่แสดงให้เห็นถึงเบื้องต้นของการทำหน้าที่ต่างกันของข้อสอบ ผลการศึกษาชี้ให้เห็นว่ามีตัวอย่างที่ขนาดต่างกันในอัตราส่วนได้ ၇ วิธี Mantel-Haenzsel มีข้อผิดพลาดประเภทที่ 1 ดีขึ้นกว่า วิธี SIBTEST ผลยังชี้ให้เห็นว่าไม่เพียงแต่สัดส่วนแต่ยังขนาดตัวอย่างและขนาดของการทำหน้าที่ต่างกันของข้อสอบที่มีอิทธิพลต่อรูปแบบของวิธี SIBTEST และ วิธี Mantel-Haenzsel เกี่ยวกับการพุฒารูปแบบอัตราความผิดพลาด

Acar (2011) ได้ศึกษาขนาดของกลุ่มตัวอย่างของการทำหน้าที่ต่างกันของข้อสอบ โดยใช้วิธี Hierarchical Linear Modeling (HGLM) เนื่องมาจากวิธี Hierarchical Linear Modeling (HGLM) วิธี Mantel-Haenszel วิธี Logistics Regression มีลักษณะที่มีความคล้ายคลึงกันจึงเลือกใช้วิธี HGLM ในการหาความแตกต่างกันของกลุ่มตัวอย่างประชากรที่ทำการศึกษา คือ นักเรียนที่ทำการสอบ OKS ในปี 2006 จำนวน 798,307 คน โดยสูงตัวอย่างจำนวน 10,727 คน ทำการทดสอบในสาขาวิชาภาษาตุรกี วิชาภาษาศาสตร์และวิชาสังคม วิชาละ 25 ข้อ พบร่วมตัวอย่าง จำนวน 2,681 คน เริ่มส่งผลให้เกิดการทำหน้าที่ต่างกันของข้อสอบโดยเมื่อจำนวนกลุ่มตัวอย่างที่เพิ่มมากขึ้นส่งผลให้การทำหน้าที่ต่างกันของข้อสอบเพิ่มสูงขึ้น

Steinmayr, Bergold, and Stiksrud (2015) ได้ศึกษาเรื่อง เพศที่แตกต่างในการทดสอบความรู้ที่นำไปของ การทำงานที่แตกต่างกันของข้อสอบ พบร่วม เพศที่แตกต่างในการทดสอบความรู้ความนิยมในหมู่คุณที่มีความแตกต่างทางเพศมีสถิติรากฐานมากที่สุดที่พบในมาตรฐานการความสามารถทางปัญญาแม้ว่าหลายคนพยายามที่จะทำให้อธิบายการค้นพบนี้การศึกษาส่วนใหญ่ยังไม่ได้พิจารณาอย่างพอเพียง ด้านระเบียบวิธีการที่แตกต่างกัน เช่นรายการที่ทำงาน (DIF) การศึกษาครั้นนี้ตรวจสอบว่ามีการทดสอบความรู้ที่นำไปของความรู้ที่มีความแตกต่างทางเพศและไม่ว่าจะเป็นเพศที่แตกต่างเหล่านี้อาจจะอธิบายได้ด้วย DIF ด้วยเหตุนี้เราบริหารงานการทดสอบความรู้ เพื่อเป็นตัวอย่างของ $N = 977$ เยอรมันนักเรียนมัธยมปลายเรารสังเกตเห็นความแตกต่างทางเพศขนาดใหญ่ในคะแนนรวมของที่นำไปของความรู้ ($|d| = 0.78$) บนพื้นฐานของวิธีการตรวจสอบ nonparametric DIF พบร่วม 40 จาก 84 รายการที่แสดงให้เห็นจำนวนมากของ DIF จัดรายการเหล่านี้จากคะแนนการทดสอบความรู้โดยรวมลดความแตกต่างทางเพศที่จะสังเกต $|d| = 0.32$ ผลการค้นหาจะมีการหารือเกี่ยวกับการแตกต่างทางเพศในการทดสอบความรู้ในการพิจารณาที่นำไปและระเบียบวิธีการ

สรุปจากการศึกษางานวิจัยที่เกี่ยวข้อง การทำหน้าที่ต่างกันของข้อสอบโดยวิธีเมนเทล-แยนเซล กลุ่มตัวอย่างของการทำหน้าที่ต่างกันของข้อสอบ โดยเมื่อจำนวนกลุ่มตัวอย่างที่เพิ่มมากขึ้นจะส่งผลให้การทำหน้าที่ต่างกันของข้อสอบเพิ่มสูงขึ้น

บทที่ 3

วิธีดำเนินการวิจัย

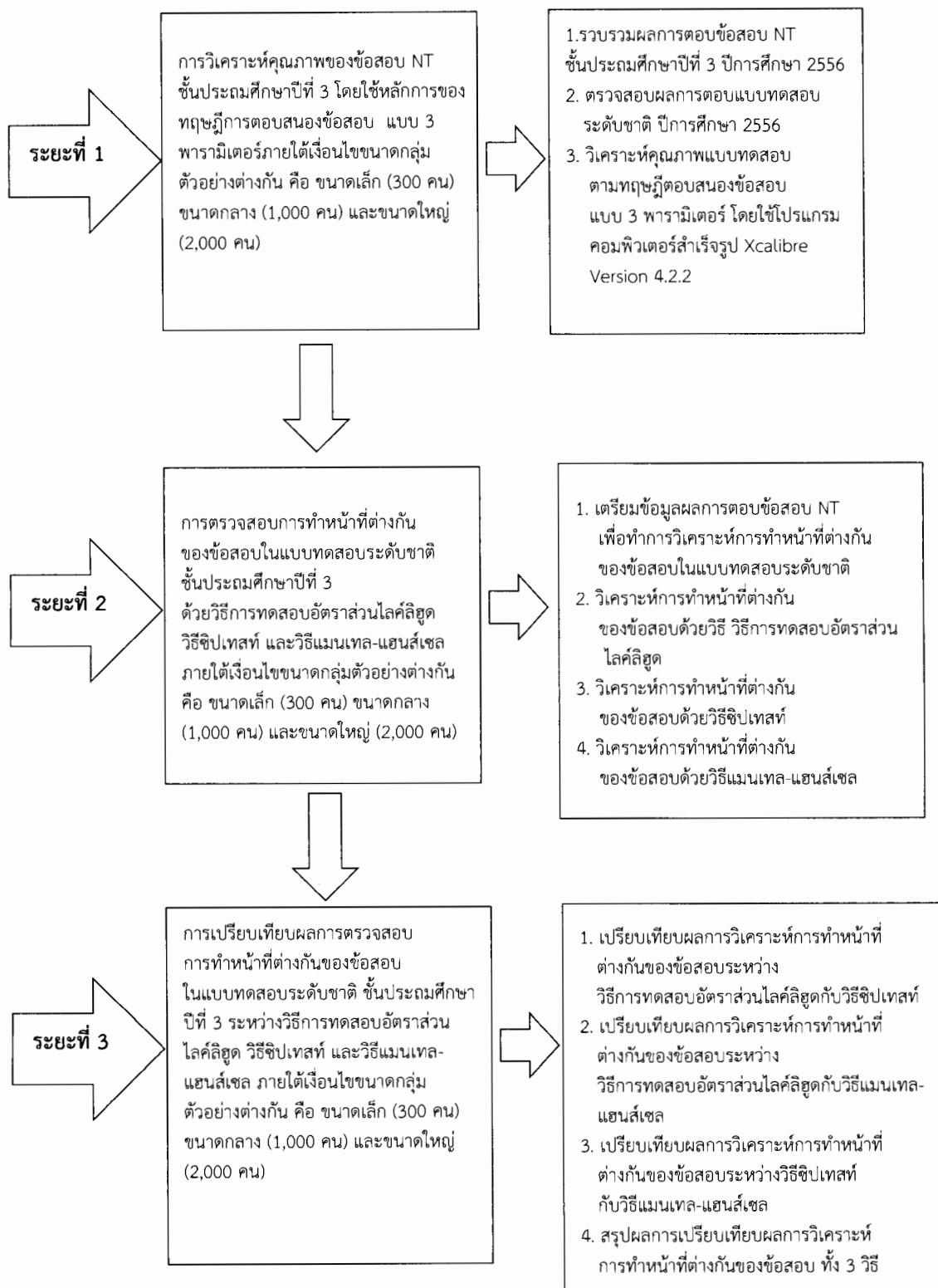
การวิจัยครั้งนี้มีวัตถุประสงค์เพื่อวิเคราะห์คุณภาพของข้อสอบ NT ตรวจสอบ การทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ และเปรียบเทียบผลการตรวจสอบ การทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ของนักเรียนชั้นประถมศึกษาปีที่ 3 จำนวน 3 ด้าน คือ ด้านภาษา ด้านคำนวณ และด้านเหตุผล ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน) ด้วยวิธีการทดสอบ อัตราส่วนไลค์ลิขิต วิธีซิปเทสท์ และวิธีแมนเทล-แ昏ส์เซล ผู้วิจัยเสนอวิธีการดำเนินการวิจัย แบ่งออกเป็น 3 ระยะ ดังนี้

ระยะที่ 1 การวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 โดยใช้หลักการ ของทฤษฎีการตอบสนองข้อสอบ แบบ 3 พารามิเตอร์ ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน)

ระยะที่ 2 การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้วยวิธีการทดสอบอัตราส่วนไลค์ลิขิต วิธีซิปเทสท์ และวิธีแมนเทล-แ昏ส์เซล ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน)

ระยะที่ 3 การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบใน แบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้วยวิธีการทดสอบอัตราส่วนไลค์ลิขิต วิธีซิปเทสท์ และวิธีแมนเทล-แ昏ส์เซล ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน)

โดยมีวิธีดำเนินการวิจัยแบ่งเป็น 3 ระยะ ดังภาพที่ 3-1



ภาพที่ 3-1 ขั้นตอนการดำเนินงานวิจัย

ระยะที่ 1 การวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 โดยใช้หลักการของทฤษฎีการตอบสนองข้อสอบ แบบ 3 พารามิเตอร์ ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน)

การวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน) ดังภาพที่ 3-2



ภาพที่ 3-2 ขั้นตอนการวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 โดยใช้หลักการของทฤษฎีการตอบสนองข้อสอบ แบบ 3 พารามิเตอร์

จากภาพที่ 3-2 แสดงขั้นตอนการวิเคราะห์คุณภาพของ NT ชั้นมัธยมศึกษาปีที่ 3 โดยใช้หลักการของทฤษฎีการตอบสนองข้อสอบ แบบ 3 พารามิเตอร์ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน) มีขั้นตอนดังนี้

- ผู้จัดได้ดำเนินการขอหนังสือขอความอนุเคราะห์ข้อมูลผลการตอบแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ปีการศึกษา 2556 จากวิทยาลัยวิทยาการวิจัยและวิทยาการปัญญา

ถึงสำนักทดสอบทางการศึกษา สำนักงานคณะกรรมการการศึกษาขั้นพื้นฐาน (สพฐ.)

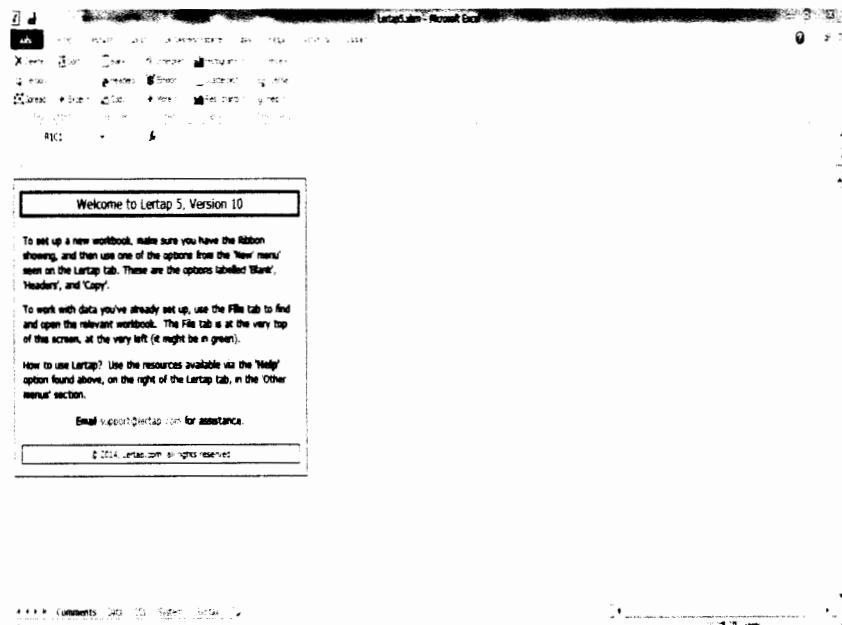
2. ตรวจสอบผลการตอบแบบทดสอบระดับชาติ ปีการศึกษา 2556 ทั้งโจทย์คdam
ตัวเลือกคำตอบ เฉลยคำตอบ และคำตอบของผู้เข้าสอบมีครบถ้วนในการตอบแบบทดสอบ ทั้ง 3 ด้าน^{คือ ด้านภาษา ด้านคำนวณ และด้านเหตุผล}

3. วิเคราะห์คุณภาพของข้อสอบ NT ตามทฤษฎีตอบสนองข้อสอบ แบบ 3 พารามิเตอร์
โดยใช้โปรแกรม Lertap5.xls มีขั้นตอน ดังนี้

3.1 เริ่มต้นด้วยເວຼຣ່າໜັກ 5 ໄຟ໌ Excel xls จะเรียกว່າ "งาน" ສມັດຈານຄື່ອງຊຸດຂອງແຜ່ນ
ງານຈະເຫັນຕ້ວຍຢ່າງດັ່ງກ່າວເມື່ອເປີດໄຟ໌ Lertap5.xls ດັ່ງການທີ່ 3-3

3.2 ພຶບແນບເມນູແລະແນບເຄື່ອງມືສອງແບບຈະປາກູ້ທີ່ທ້ານບນຂອງໜ້າຈອ ແບມໍນ
ເຮັມຈາກຕ້ວເລືອກທີ່ຂອງ File, Edit ແລະ View

3.3 ໄຟ໌ Lertap5.xls ມີ 5 ແຜ່ນງານທີ່ມີອານຸໃດໆ ໂດຍມີຊື່ຈະແສດງອູ່ດ້ານບນເປັນ-
Comments, Data, CCs, System ແລະ Syntax ໃນແຕ່ບັນດາທີ່ເກີ່ວຂ້ອງສິ່ງປາກູ້ທີ່ດ້ານລ່າງຂອງ
ໜ້າຈອ



ກາພທີ່ 3-3 ຕ້ວຍຢ່າງເປີດໄຟ໌ Lertap5.xls

3.4 ຄລິກ Bank ດັ່ງການທີ່ 3-4

A screenshot of a Microsoft Excel spreadsheet window. The title bar says 'Data' and 'Microsoft Excel'. The menu bar includes 'File', 'Edit', 'View', 'Insert', 'Format', 'Tools', 'Data', 'Window', and 'Help'. The ribbon tabs are 'Home', 'Insert', 'Page Layout', 'Formulas', 'Data', 'Page Break Preview', and 'Sort & Filter'. The sheet tab at the bottom is labeled 'Data'. The grid shows rows from 1 to 24 and columns from 1 to 24. The first few rows contain numbers: Row 1 has 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24. Row 2 has 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24. Row 3 has 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24.

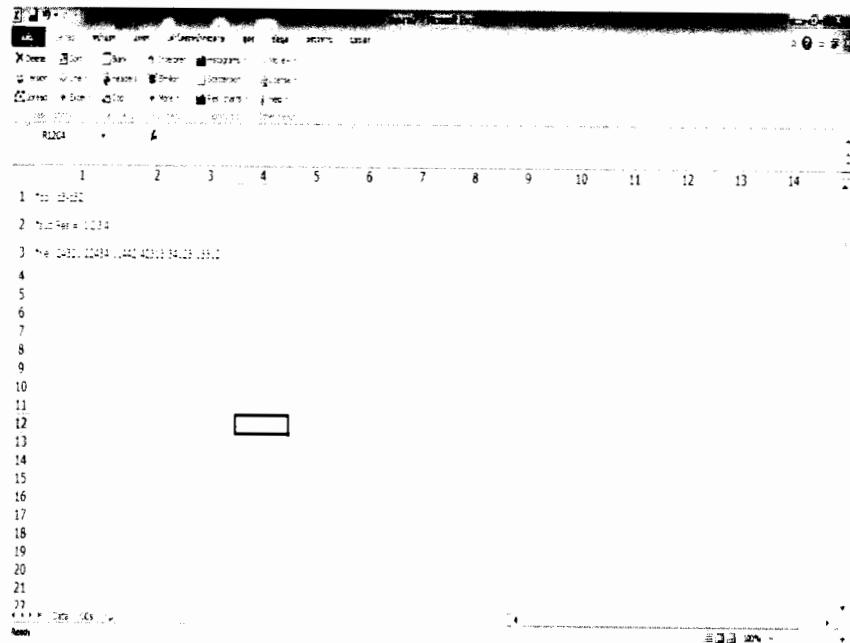
ภาพที่ 3-4 ตัวอย่างข้อมูลที่ต้องการจะวิเคราะห์ไฟล์ Lertap5.xls

3.5 เปิดไฟล์ที่บันทึกไว้ใน Excel และ copy ข้อมูลมาวางไว้ใน Data ดังภาพที่ 3-5

A screenshot of a Microsoft Excel spreadsheet window. The title bar says 'Data' and 'Microsoft Excel'. The menu bar includes 'File', 'Edit', 'View', 'Insert', 'Format', 'Tools', 'Data', 'Page Break Preview', and 'Sort & Filter'. The ribbon tabs are 'Home', 'Insert', 'Page Layout', 'Formulas', 'Data', 'Page Break Preview', and 'Sort & Filter'. The sheet tab at the bottom is labeled 'Data'. The grid shows rows from 1 to 16 and columns from 1 to 24. The first few rows contain numbers: Row 1 has 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24. Row 2 has 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24. Row 3 has 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24. Row 4 has 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24. Row 5 has 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24. Row 6 has 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24. Row 7 has 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24. Row 8 has 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24. Row 9 has 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24. Row 10 has 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24. Row 11 has 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24. Row 12 has 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24. Row 13 has 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24. Row 14 has 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24. Row 15 has 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24. Row 16 has 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24. A cell at row 8, column 8 is highlighted with a yellow border.

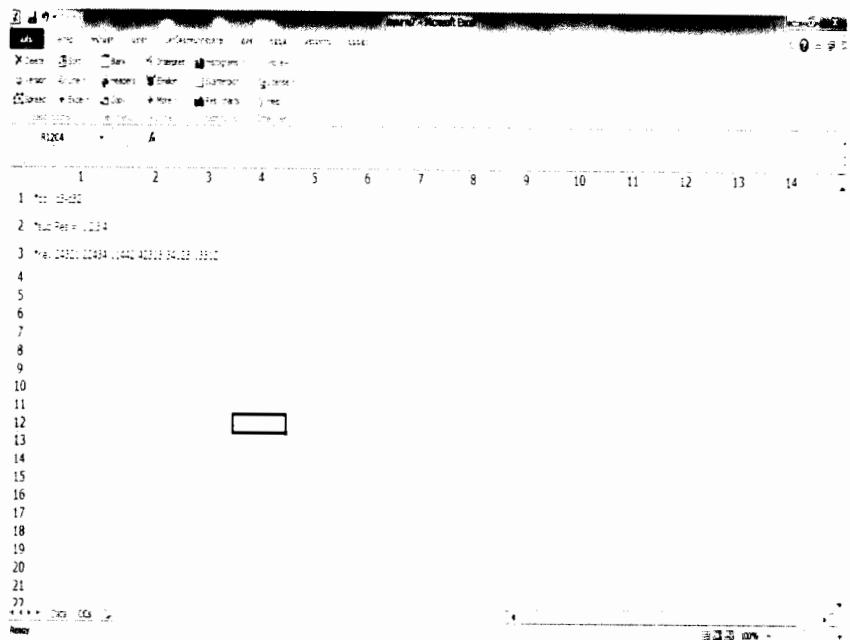
ภาพที่ 3-5 ตัวอย่างข้อมูลที่บันทึกไว้ใน Excel ที่จะวิเคราะห์ไฟล์

3.6 เปิดไฟล์ที่บันทึกไว้ใน Excel และ copy เฉลยแบบทดสอบมาวางไว้ใน CCs ดังภาพที่ 3-6



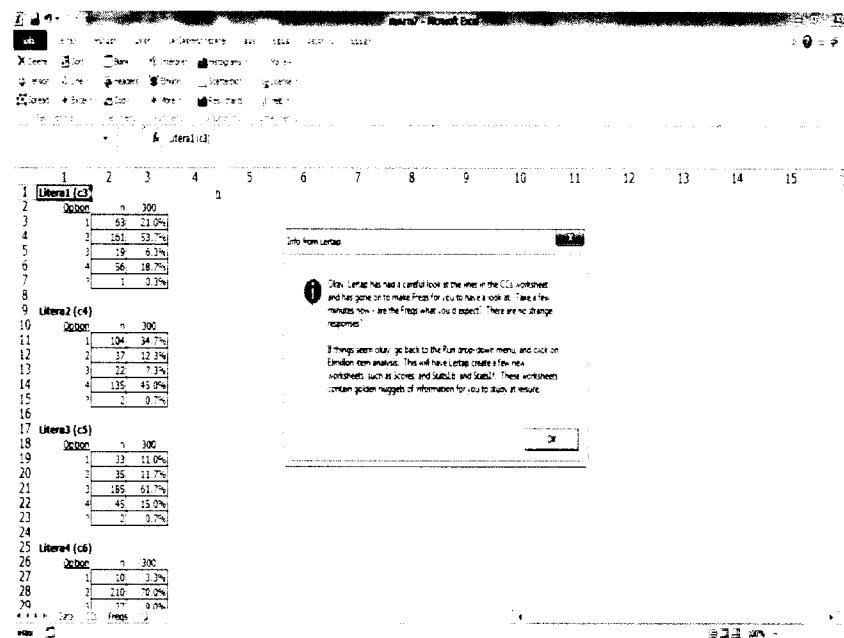
ภาพที่ 3-6 ตัวอย่างข้อมูลแผ่นงาน CCs ที่มีเฉลยคำตอบ

3.7 เลือก Interpret ดังภาพที่ 3-7



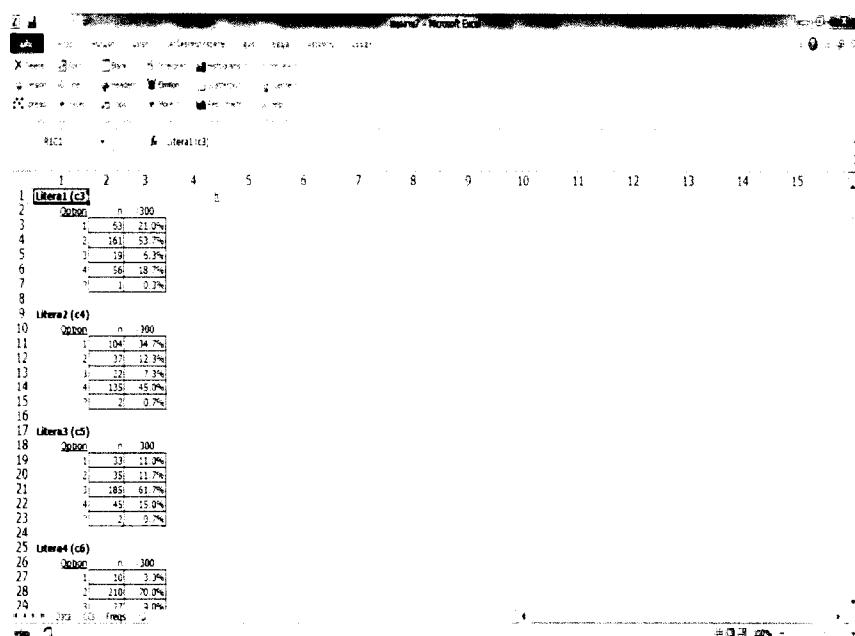
ภาพที่ 3-7 ตัวอย่างเลือกเมนู Interpret

3.8 แผ่นงานจะแสดงหน้าต่างขึ้นมา แล้วคลิก Ok ดังภาพที่ 3-8



ภาพที่ 3-8 ตัวอย่างเลือกเมนู Ok

3.9 เมื่อคลิก Ok จะแสดงผล ดังภาพที่ 3-9



ภาพที่ 3-9 ตัวอย่างแสดงผลลัพธ์ Interpret

The screenshot shows a software window titled 'Item scores and correlations for Test1 created 10/8/2011'. The main area displays a table of item statistics. A tooltip message is overlaid on the right side of the table.

Item scores and correlations for Test1 created 10/8/2011

Item	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	
1	Lectora's own item stats for 'Test1' created 10/8/2011																		
2	Options >	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
3	Item1	21%	54%	6%	19%	7%	0%	0%	0.54										
4	Item2	35%	12%	7%	45%	7%	0%	0%	0.45										
5	Item3	11%	12%	62%	15%	7%	0%	0%	0.62										
6	Item4	3%	70%	3%	17%	7%	0%	0%	0.70										
7	Item5	51%	2%	19%	27%	7%	0%	0%	0.51										
8	Item6	17%	70%	3%	10%	3%	0%	0%	0.70										
9	Item7	8%	70%	3%	13%	3%	0%	0%	0.70										
10	Item8	15%	8%	29%	48%	3%	0%	0%	0.48										
11	Item9	5%	14%	78%	3%	0%	0%	0%	0.78										
12	Item10	18%	17%	24%	42%	3%	0%	0%	0.42										
13	Item11	70%	18%	5%	6%	7%	0%	0%	0.70	0.16									
14	Item12	58%	3%	11%	22%	3%	0%	0%	0.58	0.29									
15	Item13	32%	12%	7%	48%	3%	0%	0%	0.48	0.18									
16	Item14	13%	12%	27%	48%	3%	0%	0%	0.48	0.26									

Info from Lectora

1 Data score: Lectora has had its Sorenson program produce a few new worksheets for you to study.
You should see these names below. Just click 'Select' and go to them. Click on these to see results. Then when you like:
3 You are going to edit things, be aware that each 'item' option is different. You can change where pages of items by using 'item' codes. Page 1 starts 'Item1'.
Don't lose the results! You can go back to the 'Edit' menu, make changes and then go through the full drop down menu options again. If you do this, Lectora will replace the contents of the scores and start worksheets with new results. The contents of the present scores and their sheets will be deleted.

ภาพที่ 3-10 ตัวอย่างเลือกเมนู Elmillion และคลิก OK

3.10 เมื่อคลิก OK จะแสดงผล แล้วคลิก OK อีกครั้ง ดังภาพที่ 3-11

The screenshot shows a software window titled 'Item scores and correlations for Test1 created 10/8/2011'. The main area displays a table of item statistics. A tooltip message is overlaid on the right side of the table.

Item scores and correlations for Test1 created 10/8/2011

Item	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	
1	Lectora's own item stats for 'Test1' created 10/8/2011																		
2	Options >	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
3	Item1	21%	54%	6%	19%	7%	0%	0%	0.54	0.16									
4	Item2	35%	12%	7%	45%	7%	0%	0%	0.45	0.21									
5	Item3	11%	12%	62%	15%	7%	0%	0%	0.62	0.27									
6	Item4	3%	70%	3%	17%	7%	0%	0%	0.70	0.24									
7	Item5	51%	2%	19%	27%	7%	0%	0%	0.51	0.26									
8	Item6	17%	70%	3%	10%	3%	0%	0%	0.70	0.11									
9	Item7	8%	70%	3%	13%	3%	0%	0%	0.70	0.14									
10	Item8	15%	8%	29%	48%	3%	0%	0%	0.48	0.18									
11	Item9	5%	14%	78%	3%	0%	0%	0%	0.78	0.19									
12	Item10	18%	17%	24%	42%	3%	0%	0%	0.42	0.12									
13	Item11	70%	18%	5%	6%	7%	0%	0%	0.70	0.16									
14	Item12	58%	3%	11%	22%	3%	0%	0%	0.58	0.29									
15	Item13	32%	12%	7%	48%	3%	0%	0%	0.48	0.18									
16	Item14	13%	12%	27%	48%	3%	0%	0%	0.48	0.26									

Info from Lectora

1 Data score: Lectora has had its Sorenson program produce a few new worksheets for you to study.
You should see these names below. Just click 'Select' and go to them. Click on these to see results. Then when you like:
3 You are going to edit things, be aware that each 'item' option is different. You can change where pages of items by using 'item' codes. Page 1 starts 'Item1'.
Don't lose the results! You can go back to the 'Edit' menu, make changes and then go through the full drop down menu options again. If you do this, Lectora will replace the contents of the scores and start worksheets with new results. The contents of the present scores and their sheets will be deleted.

ภาพที่ 3-11 ตัวอย่างแสดงภาพหน้าจอหลังจากเลือกเมนู OK

3.11 เลือก More ดังภาพที่ 3-11

3.12 เลือก Item scores and correlation ดังภาพที่ 3-12

Item	Item stats for "Test1" created: 10/8/2561																			
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18		
1	Item stats for "Test1" created: 10/8/2561	Item	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
2	Options >	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	
3	Item1	21%	54%	6%	19%	6%	0.54	0.18												
4	Item2	35%	12%	7%	45%	2%	0.45	0.21												
5	Item3	11%	12%	62%	15%	2%	0.62	0.27												
6	Item4	3%	20%	9%	17%	2%	0.70	0.24												
7	Item5	51%	2%	19%	27%	2%	0.51	0.26												
8	Item6	17%	20%	10%	3%		0.70	0.11												
9	Item7	8%	20%	9%	13%		0.70	0.34												
10	Item8	15%	8%	29%	48%		0.48	0.18												
11	Item9	5%	14%	23%	3%		0.78	0.19												
12	Item10	18%	17%	24%	42%		0.42	0.12												
13	Item11	20%	18%	5%	5%	1%	0.70	0.16												
14	Item12	58%	8%	11%	23%		0.58	0.29												
15	Item13	32%	12%	7%	48%		0.48	0.18												
16	Item14	13%	12%	27%	48%		0.48	0.26												

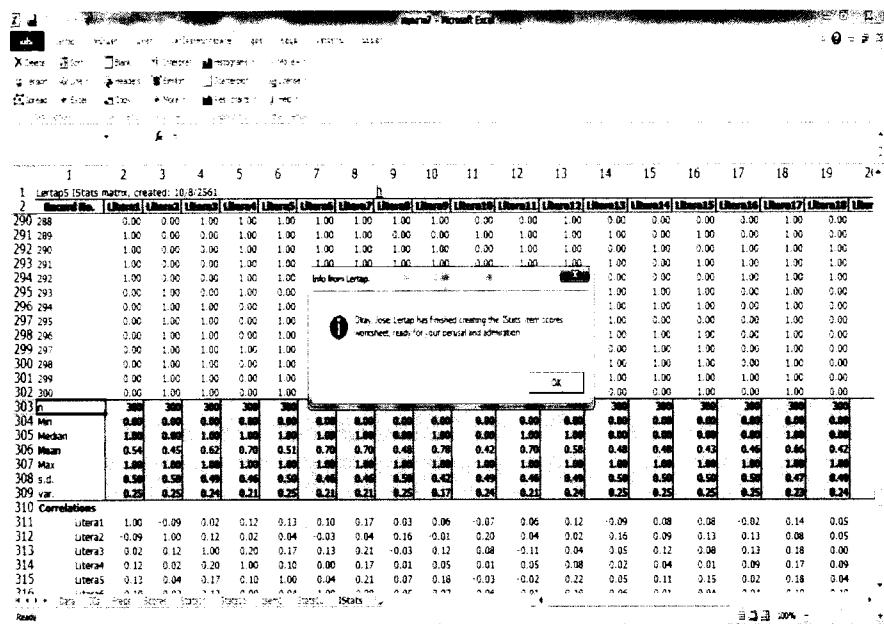
ภาพที่ 3-12 ตัวอย่างแสดงภาพหน้าจอหลังจากเลือกเมนู More เลือก Item Scores and Correlation

3.13 คลิก OK ดังภาพที่ 3-13

Item	Item stats for "Test1" created: 10/8/2561																			
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18		
1	Item stats for "Test1" created: 10/8/2561	Item	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
2	Options >	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	
3	Item1	21%	54%	6%	19%	6%	0.54													
4	Item2	35%	12%	7%	45%	2%	0.45													
5	Item3	11%	12%	62%	15%	2%	0.62													
6	Item4	3%	20%	9%	17%	2%	0.70													
7	Item5	51%	2%	19%	27%	2%	0.51													
8	Item6	17%	20%	10%	3%		0.70													
9	Item7	8%	20%	9%	13%		0.70													
10	Item8	15%	8%	29%	48%		0.48													
11	Item9	5%	14%	23%	3%		0.78													
12	Item10	18%	17%	24%	42%		0.42													
13	Item11	20%	18%	5%	5%	1%	0.70	0.16												
14	Item12	58%	8%	11%	23%		0.58	0.29												
15	Item13	32%	12%	7%	48%		0.48	0.18												
16	Item14	13%	12%	27%	48%		0.48	0.26												

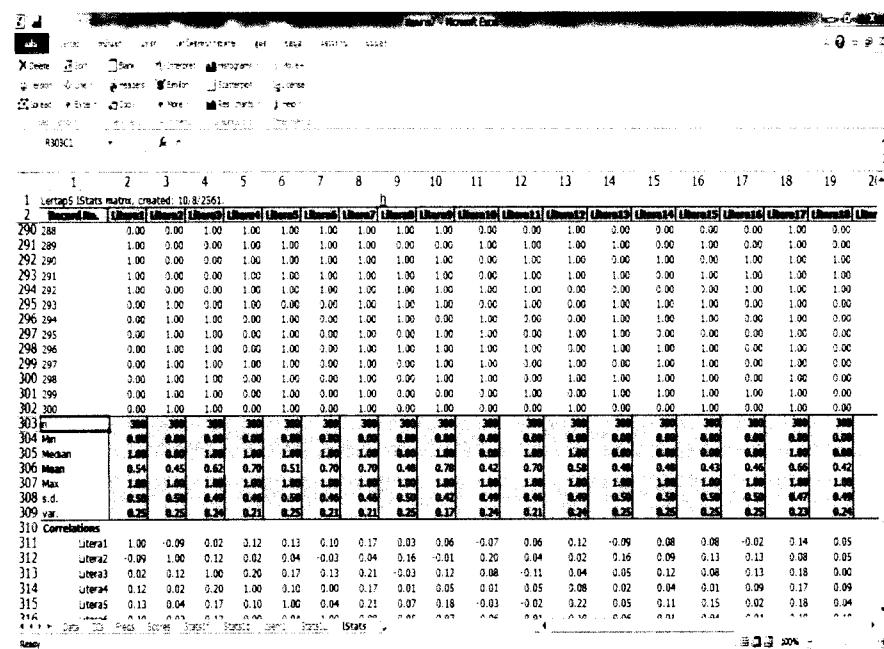
ภาพที่ 3-13 ตัวอย่างแสดงภาพหน้าจอหลังจากเลือกเมนู More เลือก Item Scores and Correlation และคลิก OK

3.14 คลิก OK ดังภาพที่ 3-14



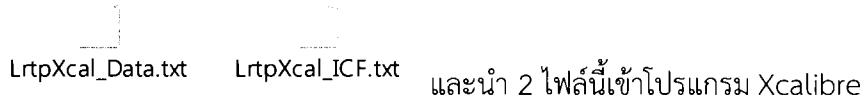
ภาพที่ 3-14 ตัวอย่างแสดงภาพหน้าจอหลังจากคลิก OK และคลิก OK

3.15 คลิก OK ดังภาพที่ 3-15

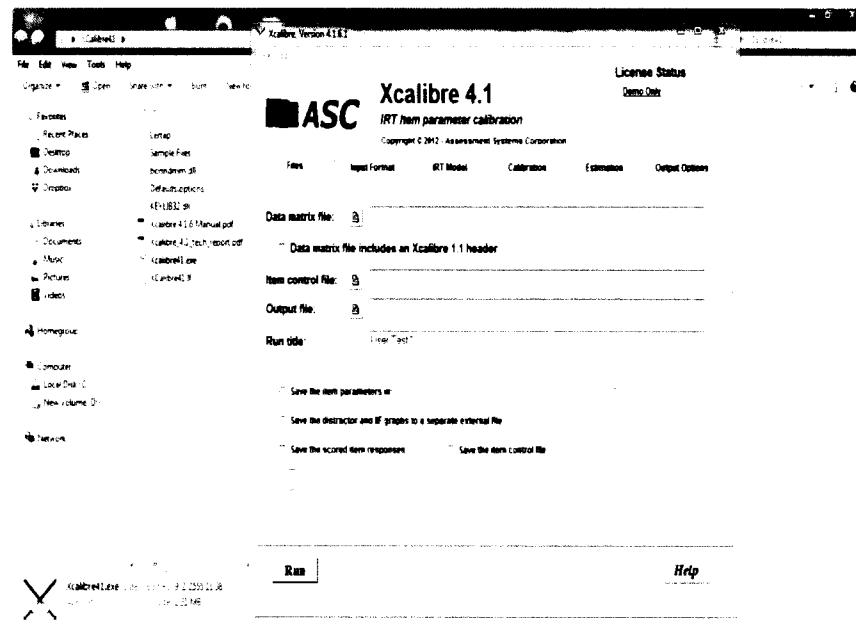


ภาพที่ 3-15 ตัวอย่างแสดงภาพหน้าจอ Output ของการ Run เมนู Item Scores and Correlation

3.16 เมื่อวิเคราะห์ Run แล้วจะได้ไฟล์ ดังนี้เป็น 2 ไฟล์ดังนี้



4. การวิเคราะห์คุณภาพของข้อสอบ NT ตามทฤษฎีตอบสนองข้อสอบแบบ 3 พารามิเตอร์ โดยใช้โปรแกรม Xcalibre มีขั้นตอนดังนี้

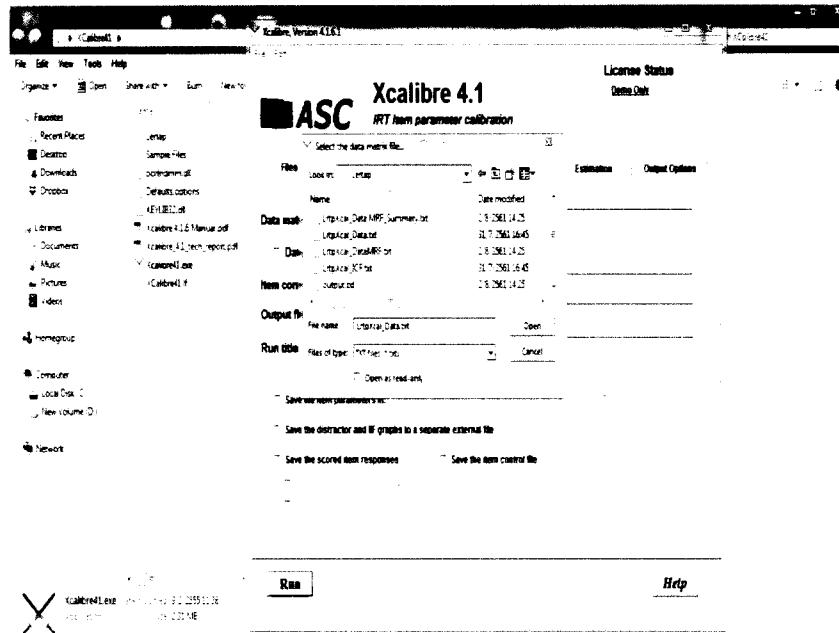


ภาพที่ 3-16 ตัวอย่างโปรแกรม Xcalibre

4.1 โปรแกรม Xcalibre ต้องมี 2 ไฟล์ ได้แก่ Data Matrix File และ Item Control File

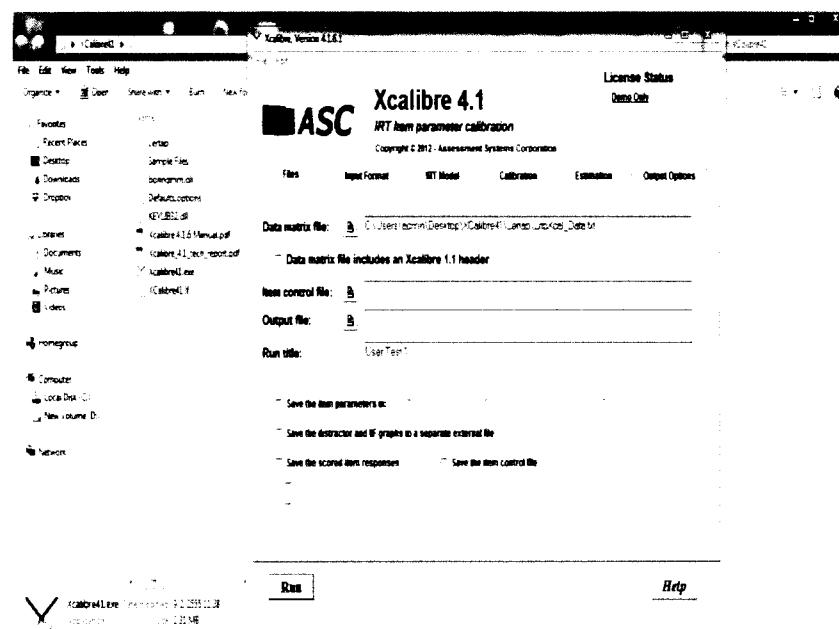
4.2 ระบุไฟล์ในแท็บ Files

4.3 แต่ละไฟล์ ซึ่งจะเปิดใช้หน้าต่างมาตรฐานเพื่อรับ ชื่อของแต่ละไฟล์

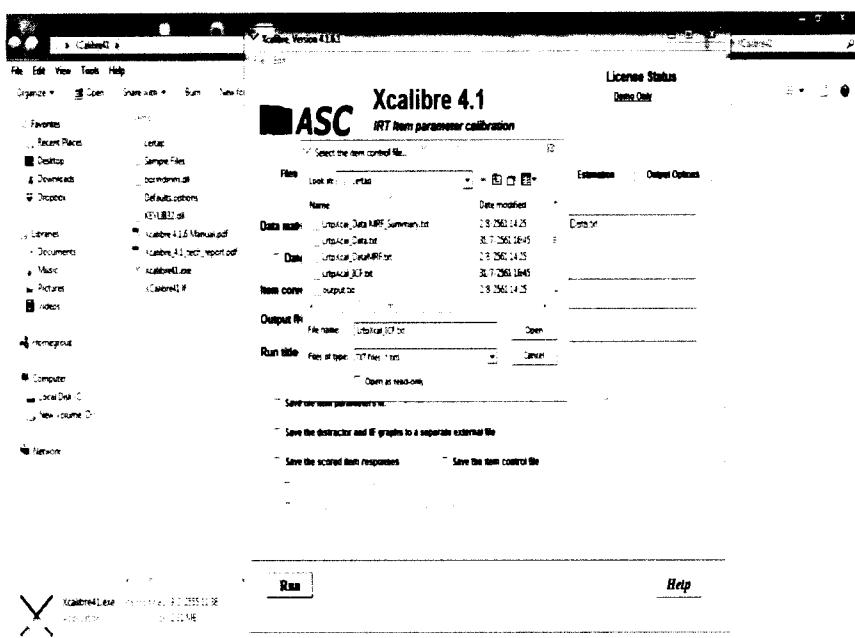


ภาพที่ 3-17 ตัวอย่างการเลือกไฟล์ Data.txt สำหรับวิเคราะห์โปรแกรม Xcalibre

4.4 ไฟล์ที่เรียกเข้า Data Matrix File คือ LtpXcal_Data.txt ดังภาพที่ 3-18

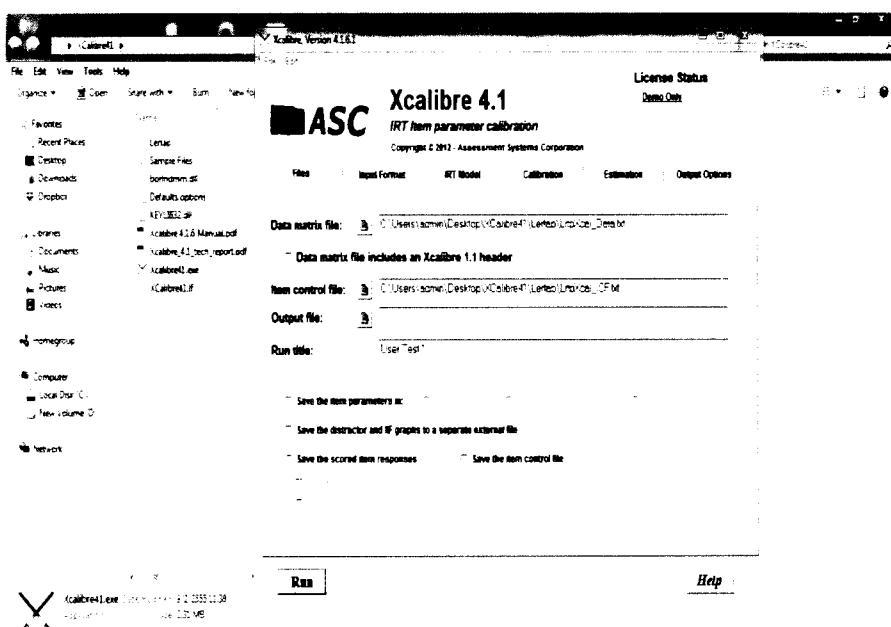


ภาพที่ 3-18 ตัวอย่างไฟล์ที่ต้องการเลือกวิเคราะห์โปรแกรม Xcalibre ใน Data Matrix File



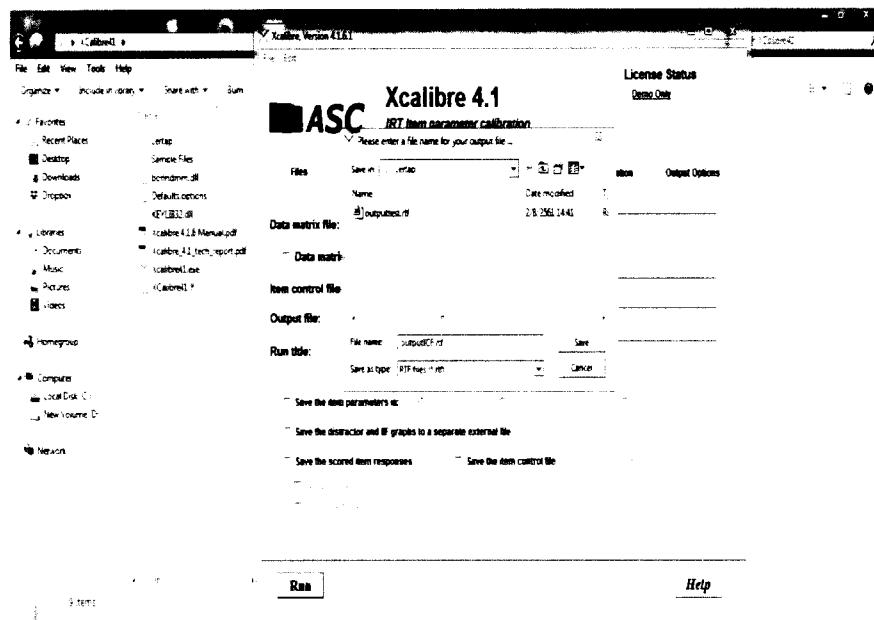
ภาพที่ 3-19 ตัวอย่างการเลือกไฟล์ ICF.txt วิเคราะห์โปรแกรม Xcalibre

4.5 ไฟล์ที่เรียกเข้า Item Control File คือ LrtpXcal_ICF.txt ดังภาพที่ 3-20



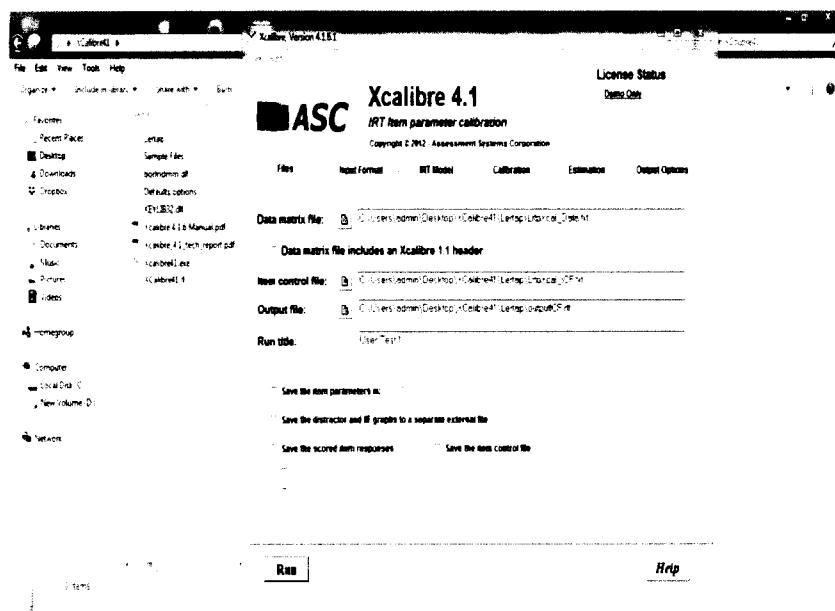
ภาพที่ 3-20 ตัวอย่างไฟล์ที่ต้องการเลือกวิเคราะห์โปรแกรม Xcalibre ใน Item Control File

4.6 ระบุชื่อไฟล์ที่ต้องการจะเก็บไว้ Output File คือ ดังภาพที่ 3-21

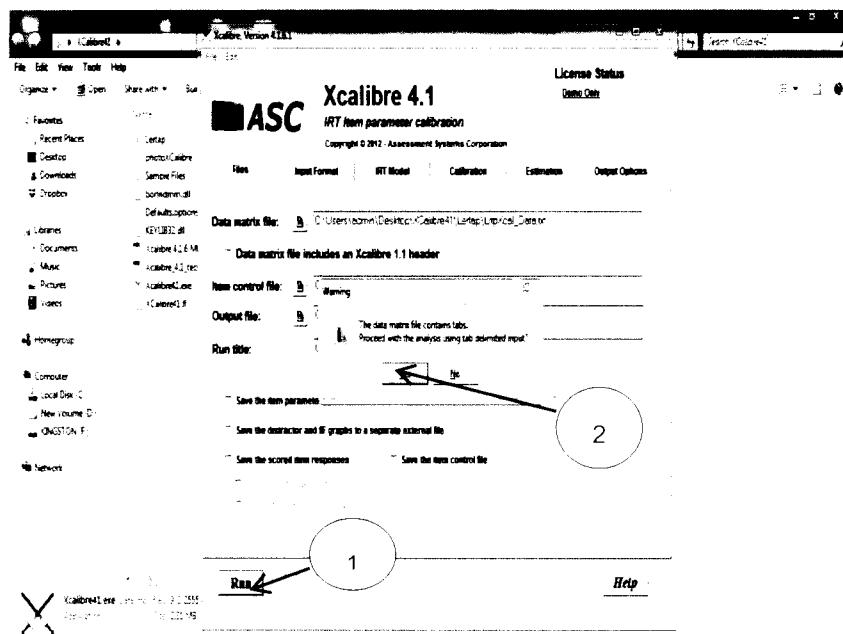


ภาพที่ 3-21 ตัวอย่างระบบไฟล์ที่ต้องการจะเก็บช่อง Output File

4.7 ไฟล์ที่ระบุต้องการจะเก็บไว้จะแสดงผลที่ช่อง Output File คือ ดังภาพที่ 3-22



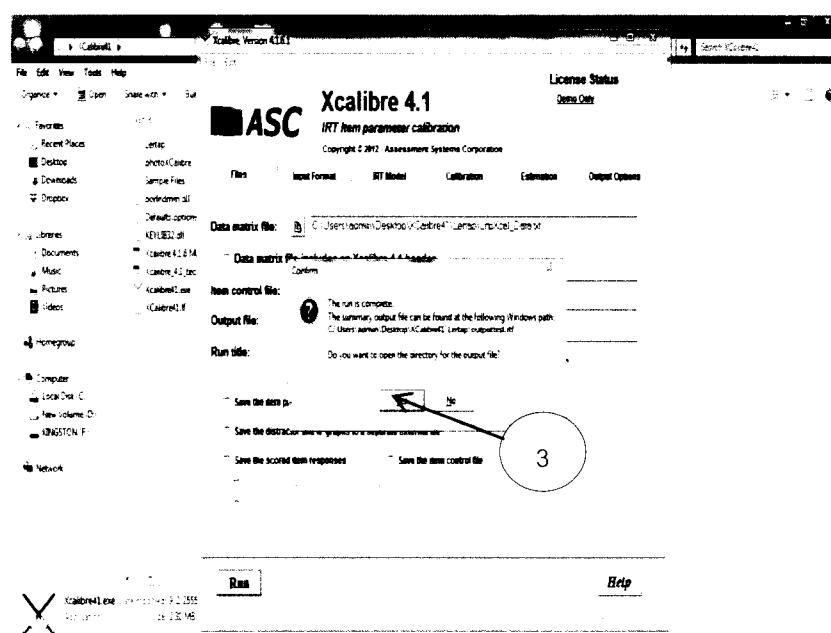
ภาพที่ 3-22 ตัวอย่างไฟล์ที่ต้องการเลือกวิเคราะห์โปรแกรม Xcalibre ในช่อง Output File



ภาพที่ 3-23 ตัวอย่างไฟล์ที่ต้องการเลือกวิเคราะห์โปรแกรม Xcalibre

4.8 คลิกที่ หมายเลข 1 เพื่อ Run ข้อมูลการวิเคราะห์ข้อมูล คือ ดังภาพที่ 3-23

4.9 คลิกที่ หมายเลข 2 เพื่อ Yes ดังภาพที่ 3-23



ภาพที่ 3-24 ตัวอย่างไฟล์ที่ต้องการเลือกวิเคราะห์โปรแกรม Xcalibre เลือกที่ Yes

4.10 คลิกที่ หมายเลข 3 เพื่อ Yes ดังภาพที่ 3-24



***IRT Item Parameter
Calibration Report***

User Test 1

Report created on 3/8/2018

Xcalibre 4.1: IRT Item Parameter Estimation Software
Copyright © 2012 - Assessment Systems Corporation



ภาพที่ 3-25 แสดงตัวอย่างไฟล์ ผลการวิเคราะห์ด้วยโปรแกรม Xcalibre

ระยะที่ 2 การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้วยวิธีการทดสอบอัตราส่วนไลค์ลิชูด วิธีซิปเทส์ และ วิธีแมนเทล-แ昏ส์เซล ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน)

การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้วยวิธีการทดสอบอัตราส่วนไลค์ลิชูด วิธีซิปเทส์ และวิธีแมนเทล-แ昏ส์เซล ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน) ดังภาพที่ 3-26



ภาพที่ 3-26 ขั้นตอนการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3

จากภาพที่ 3-26 แสดงขั้นตอนการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ขั้นประถมศึกษาปีที่ 3 ด้วยวิธีการทดสอบอัตราส่วนไลค์ลิขิต วิชีซิปเทสท์ และวิชีแมนเทล-แฮนส์เซล ดังนี้

1. การวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ด้วยวิธี IRT-LR โดยใช้โปรแกรม IRTPRO

การวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบด้วยวิธี IRT-LR โดยใช้โปรแกรม IRTPRO เป็นการวิเคราะห์ด้วยสถิติทางคณิตศาสตร์ โดยการเปรียบเทียบความแตกต่างของ ค่าพารามิเตอร์ ในกลุ่มอ้างอิงและกลุ่มเปรียบเทียบ ผู้วิจัยได้เตรียมการวิเคราะห์ กำหนดค่าตัวแปร ด้านภาษา ด้านคำนวณ และด้านเหตุผล โดยกำหนดค่าตัวแปรเพศ ให้ตัวแปรเพศชายเป็น 1 และตัวแปรเพศหญิงเป็น 2 ดังนี้

ขั้นตอนที่ 1 เตรียมไฟล์ข้อมูลสำหรับวิเคราะห์

การวิเคราะห์ข้อมูลด้วยวิธี IRT-LR โดยใช้โปรแกรม IRTPRO เตรียมข้อมูล ในรูปแบบของไฟล์ .ssig วิเคราะห์การทำหน้าที่ต่างกันของข้อสอบด้วยโปรแกรม IRTPRO ประกอบด้วยข้อมูล คือ (Student) นักเรียนที่เข้าสอบ (Gender) เพศชาย เพศหญิง และผลการตอบแบบทดสอบของผู้สอบ (Response) ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน) ดังภาพที่ 3-27

	Student	Gender	Literacy1	Literacy2	Literacy3	Literacy4	Literacy5	Literacy6	Literacy7	Literacy8	Literacy9	Literacy10	Literacy11	Literacy12	Gender
1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1
2	2	1	0	1	1	1	1	1	1	1	1	1	1	1	1
3	3	1	0	1	1	1	1	1	1	1	1	1	1	1	1
4	4	1	0	1	1	1	1	0	1	1	1	1	1	1	1
5	5	1	0	1	1	1	1	0	1	1	1	1	1	1	1
6	6	1	0	0	1	0	1	0	0	0	0	1	0	0	1
7	7	1	1	0	1	1	0	1	1	0	0	0	1	0	1
8	8	1	1	0	1	1	0	0	1	0	1	0	0	1	0
9	9	1	0	0	0	1	0	0	1	0	1	0	1	0	0
10	10	1	0	1	1	1	0	1	1	1	1	1	0	1	0
11	11	1	1	0	1	0	1	0	1	0	1	1	1	1	1
12	12	1	1	0	0	0	0	1	0	0	0	1	0	0	0
13	13	1	0	1	0	0	0	1	1	1	1	1	0	1	0
14	14	1	1	0	1	1	1	1	1	1	1	1	0	1	0
15	15	1	0	1	1	1	0	0	1	1	1	1	1	1	1
16	16	1	0	1	1	0	0	0	0	0	0	0	0	0	0
17	17	1	0	0	1	1	1	1	1	0	1	1	0	1	1
18	18	1	1	0	0	1	0	1	1	1	1	1	0	0	0
19	19	1	1	0	1	1	1	1	1	0	1	0	1	0	0

ภาพที่ 3-27 ตัวอย่างการจัดไฟล์ข้อมูลสำหรับวิเคราะห์ด้วยวิธี IRT-LR ในรูปแบบไฟล์ .ssig

ขั้นตอนที่ 2 การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ (DIF)

การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ (DIF) ด้วยวิธี IRT-LR โดยใช้โปรแกรม IRTPRO จะประเมินความสำคัญระหว่างกลุ่มอ้างอิงและกลุ่มเปรียบเทียบในแบบทดสอบที่แตกต่างของแต่ละพารามิเตอร์ ดังนี้

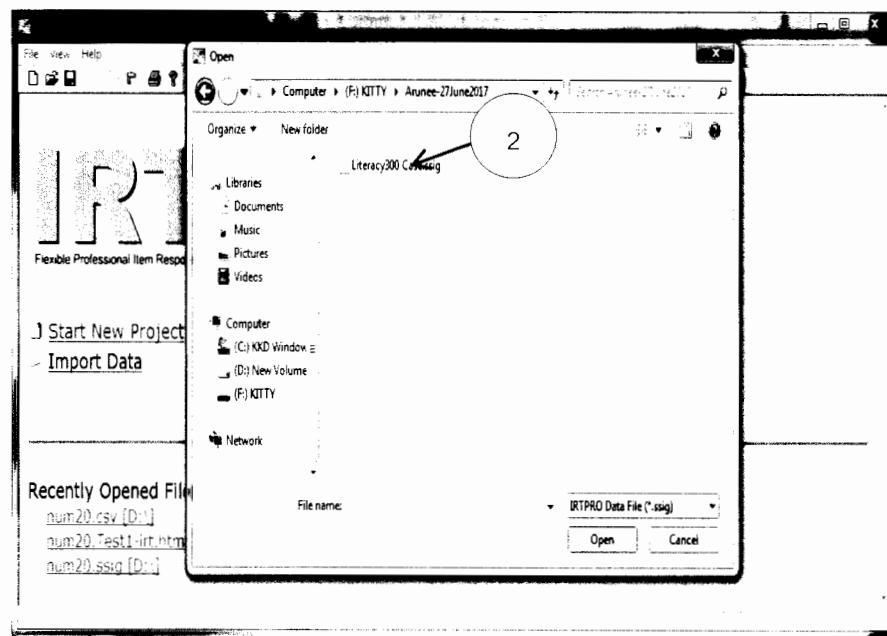
1. เลือกวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบด้วยข้อสอบตามทฤษฎีการตอบทฤษฎี การตอบสนองข้อสอบแบบเอกมิติ (Unidimensional IRT)
2. กำหนดตัวแปร ในด้านภาษา และด้านเหตุผล ให้เพศหญิงเป็นกลุ่มอ้างอิง (Referent Group: R) และเพศชายเป็นกลุ่มเปรียบเทียบ (Focal Group: F) ส่วนด้านคำนวณ ให้เพศชายเป็นกลุ่มอ้างอิง (R) และเพศหญิงเป็นกลุ่มเปรียบเทียบ (F)
3. เลือกโมเดลสำหรับการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ เพื่อวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบ ทำการวิเคราะห์แบบทดสอบทุกข้อจากนั้นพิจารณาว่าข้อใดมี DIF โดยดูจากค่า p-value ที่มีนัยสำคัญทางสถิติที่ระดับ .05 ดังภาพที่ 3-28

1. เปิดโปรแกรม IRTPRO ดังภาพที่ 3-28



ภาพที่ 3-28 ตัวอย่างการเรียกไฟล์ข้อมูลสำหรับวิเคราะห์ด้วยวิธี IRT-LR ในรูปแบบไฟล์ .ssig

2. เลือกที่หมายเลข 1 คลิกที่ File และเลือก Open ดังภาพที่ 3-29



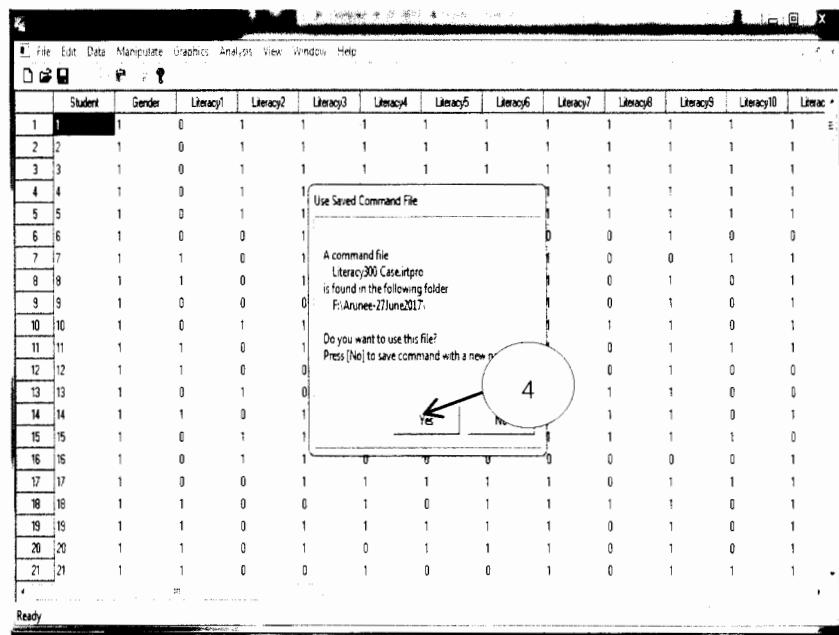
ภาพที่ 3-29 ตัวอย่างการเปิดไฟล์ข้อมูลสำหรับวิเคราะห์ด้วยวิธี IRT-LR ในรูปแบบไฟล์ .ssig

3. เลือกที่หมายเลข 3 เพื่อเลือก Analysis ดังภาพที่ 3-30

Student	Gender	Literacy1	Literacy2	Literacy3	Literacy4	Literacy5	Literacy6	Literacy7	Literacy8	Literacy9	Literacy10	Literacy11
1	1	0	1	1	1	1	1	1	1	1	1	1
2	1	0	1	1	1	1	1	1	1	1	1	1
3	1	0	1	1	1	1	1	1	1	1	1	1
4	1	0	1	1	1	1	0	1	1	1	1	1
5	1	0	1	1	1	1	0	1	1	1	1	1
6	1	0	0	1	0	1	0	0	0	1	0	0
7	1	1	0	1	1	0	1	1	0	0	1	1
8	1	1	0	1	1	0	0	1	0	1	0	1
9	1	0	0	0	1	0	0	1	0	1	0	1
10	1	0	1	1	1	0	1	1	1	1	0	1
11	1	1	0	1	0	1	0	1	0	1	1	1
12	1	1	0	0	0	0	1	0	0	1	0	0
13	1	0	1	0	0	0	1	1	1	0	0	0
14	1	1	0	1	1	1	1	1	1	0	1	1
15	1	0	1	1	1	0	0	1	1	1	1	0
16	1	0	1	1	0	0	0	0	0	0	1	1
17	1	0	0	1	1	1	1	0	1	1	1	1
18	1	1	0	0	1	0	1	1	1	1	0	1
19	1	1	0	1	1	1	1	0	1	0	1	1
20	1	1	0	1	0	1	1	1	0	1	0	1
21	1	1	0	0	1	0	0	1	0	1	1	1

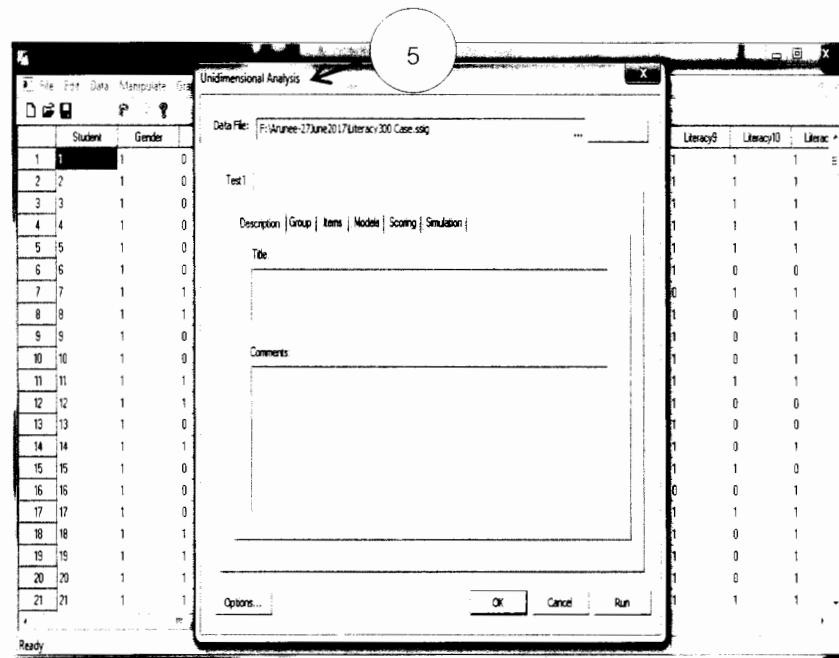
ภาพที่ 3-30 ตัวอย่างการเลือก Analysis การวิเคราะห์ข้อมูล

4. เลือกที่หมายเลข 4 เพื่อเลือก Yes ดังภาพที่ 3-31



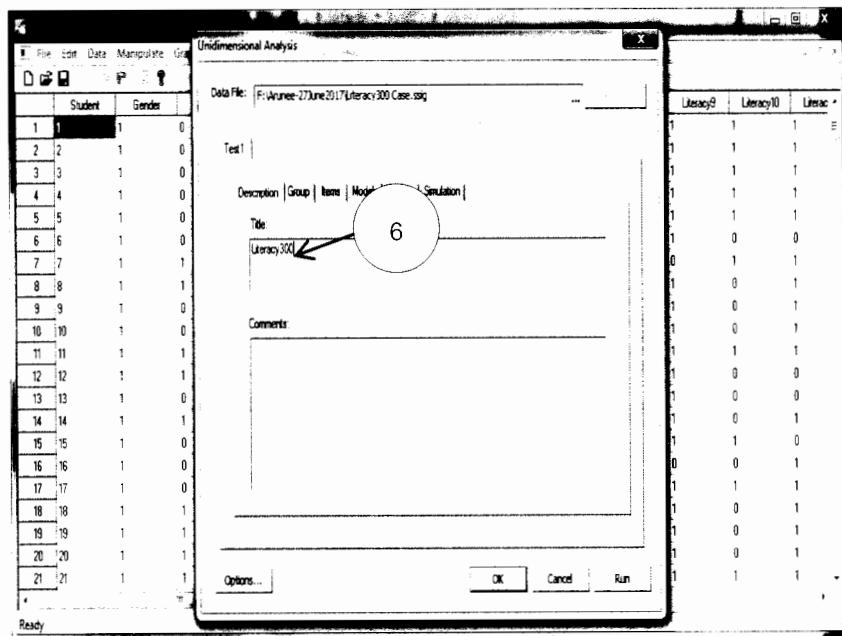
ภาพที่ 3-31 ตัวอย่างการเลือก Yes เพื่อการวิเคราะห์ข้อมูล

5. เลือกที่หมายเลข 5 เพื่อเลือก Unidimensional Analysis ดังภาพที่ 3-32



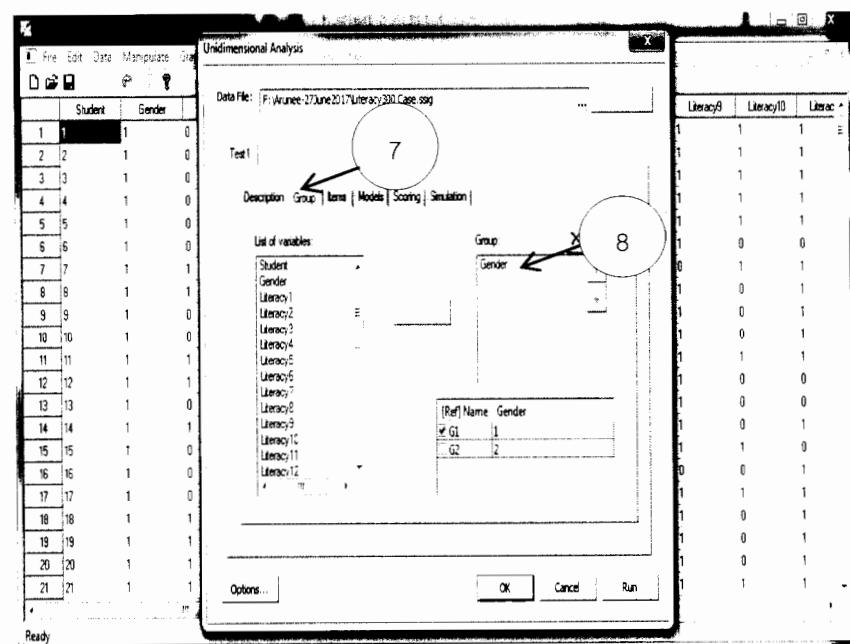
ภาพที่ 3-32 ตัวอย่างการเลือก Unidimensional Analysis เพื่อการวิเคราะห์ข้อมูล

6. เลือกที่หมายเลข 6 เพื่อเลือก ใส่ชื่อเรื่องไฟล์ ดังภาพที่ 3-33



ภาพที่ 3-33 ตัวอย่างการเลือกไฟล์เพื่อการวิเคราะห์ข้อมูล

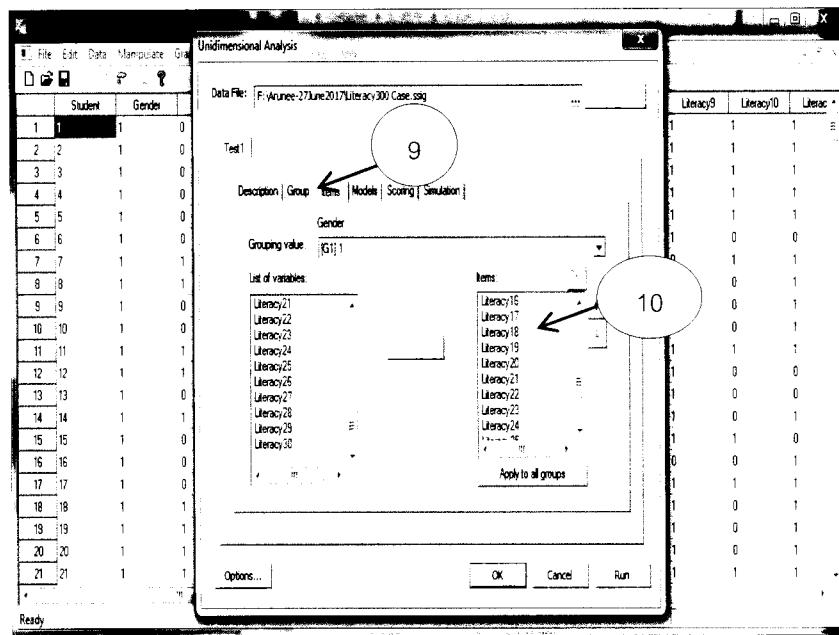
7. เลือกที่หมายเลข 7 เพื่อเลือก Group ดังภาพที่ 3-34
8. เลือกที่หมายเลข 8 เพื่อเลือก Gender ดังภาพที่ 3-34



ภาพที่ 3-34 ตัวอย่างการเลือก Group และเลือก Gender เพื่อการวิเคราะห์ข้อมูล

9. เลือกที่หมายเลข 9 เพื่อเลือก Item ดังภาพที่ 3-35

10. เลือกที่หมายเลข 10 เพื่อเลือก จำนวนข้อสอบทั้งหมด ดังภาพที่ 3-35



ภาพที่ 3-35 ตัวอย่างการเลือก Item และเลือกจำนวนข้อสอบทั้งหมดเพื่อการวิเคราะห์ข้อมูล

2. การวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ด้วยวิธี SIBTEST

การวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบด้วยวิธี SIBTEST โดยใช้โปรแกรม SIBTEST ในการเปรียบเทียบความแตกต่างของ ค่าพารามิเตอร์ในกลุ่มอ้างอิงและกลุ่มเปรียบเทียบ ผู้วิจัยได้เตรียมการวิเคราะห์ กำหนดค่าตัวแปร ด้านภาษา ด้านคำนวณ และด้านเหตุผล โดยกำหนดค่าตัวแปร เพศ ให้ตัวแปรเพศชายเป็น 1 และตัวแปรเพศหญิงเป็น 2 ดังนี้

ขั้นตอนที่ 1 เตรียมไฟล์ข้อมูลสำหรับวิเคราะห์

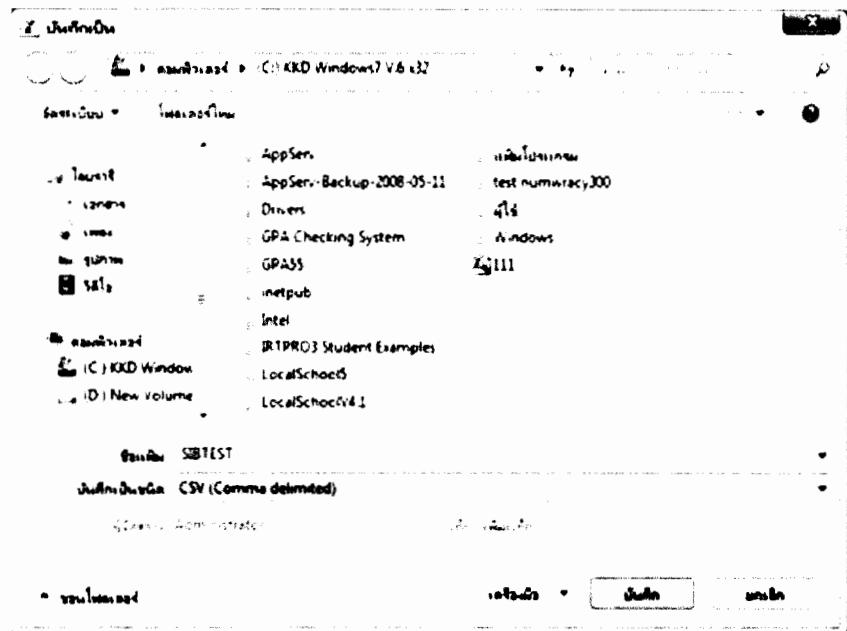
การวิเคราะห์ข้อมูลด้วยวิธี SIBTEST ผู้วิจัยได้เตรียมข้อมูล วิเคราะห์การทำหน้าที่ต่างกัน ของข้อสอบด้วยโปรแกรม SIBTEST ประกอบด้วยข้อมูล คือ นักเรียน (Student) ที่เข้าสอบ (Gender) เพศชาย เพศหญิง และผลการตอบแบบทดสอบของผู้สอบ (Response) ภายใต้เงื่อนไข ขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน) ดังภาพที่ 3-36

The screenshot shows a Microsoft Excel spreadsheet with 32 rows and 18 columns. The columns are labeled 1 through 18. The first few rows contain binary data (0s and 1s) representing the SIBTEST dataset. Row 1 has a header 'SIBTEST' at the top. The data starts from row 2, column 1.

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
1	1	1	1	1	1	1	0	0	1	1	1	1	0	0	1	1	1
2	2	1	1	1	1	1	0	0	1	1	1	1	0	0	1	1	1
3	3	1	1	1	1	1	0	0	1	1	1	1	0	0	1	1	1
4	4	1	1	1	1	0	1	0	0	1	0	0	0	0	0	1	1
5	5	1	0	0	0	0	0	0	0	1	0	0	0	0	1	1	0
6	6	1	0	0	1	0	0	0	0	0	0	1	0	1	1	0	0
7	7	1	0	0	1	0	0	0	0	0	0	0	1	0	0	1	1
8	8	1	0	0	0	0	0	0	0	1	1	0	1	0	0	1	1
9	9	1	0	0	1	0	1	0	1	1	1	0	0	0	0	0	0
10	10	1	0	0	0	1	1	0	1	0	1	1	1	1	0	0	1
11	11	1	0	1	0	1	0	1	0	0	1	0	0	0	0	0	1
12	12	1	0	0	0	0	1	1	1	0	1	0	0	0	0	0	1
13	13	1	0	0	0	0	0	0	0	0	1	0	1	1	0	0	0
14	14	1	0	1	0	0	0	0	0	1	0	0	1	0	0	0	0
15	15	1	0	0	1	0	0	0	1	0	0	0	0	1	1	0	0
16	16	1	0	1	1	0	1	0	0	0	0	0	0	0	1	0	0
17	17	1	0	1	1	0	1	0	1	1	0	1	0	0	1	0	1
18	18	1	0	0	0	0	0	0	0	0	1	0	0	1	0	1	0
19	19	1	0	1	0	0	0	0	0	0	0	0	0	0	0	1	0
20	20	1	0	1	1	0	0	0	0	1	1	0	0	0	0	0	1
21	21	1	0	1	1	1	0	0	0	1	1	1	0	0	1	0	0
22	22	1	0	1	0	0	0	0	0	1	1	1	0	0	1	0	1
23	23	1	0	1	1	0	0	0	1	1	1	1	0	1	0	1	1
24	24	1	0	0	0	0	0	1	0	0	1	0	1	0	0	0	0
25	25	1	0	0	0	0	0	0	0	0	1	1	0	0	0	0	0
26	26	1	0	1	1	1	0	0	0	1	1	1	0	0	0	0	0
27	27	1	0	0	0	0	0	0	1	0	1	0	0	0	1	0	0

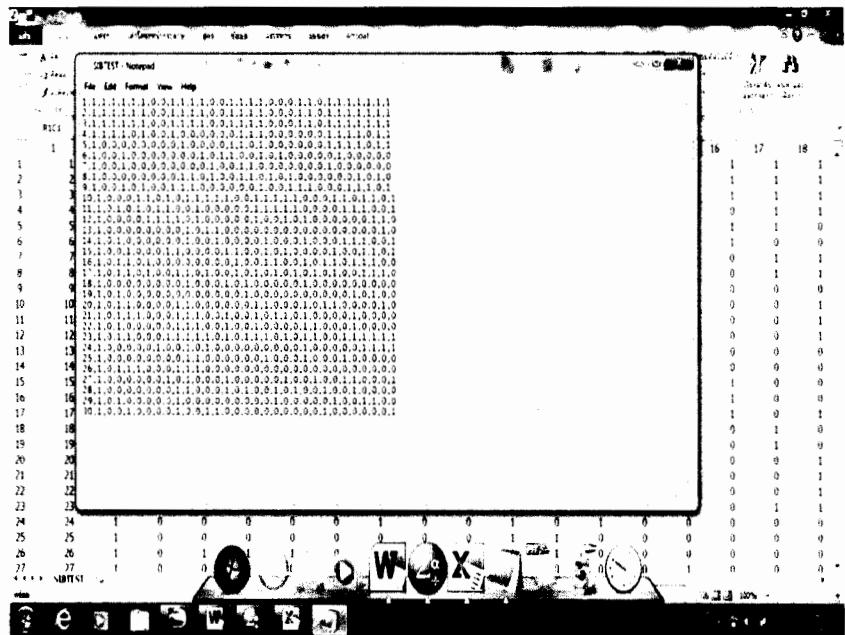
ภาพที่ 3-36 ตัวอย่างการจัดไฟล์ข้อมูลสำหรับวิเคราะห์ด้วยวิธี SIBTEST ในรูปแบบไฟล์ .dat

1. เปิดโปรแกรม Excel



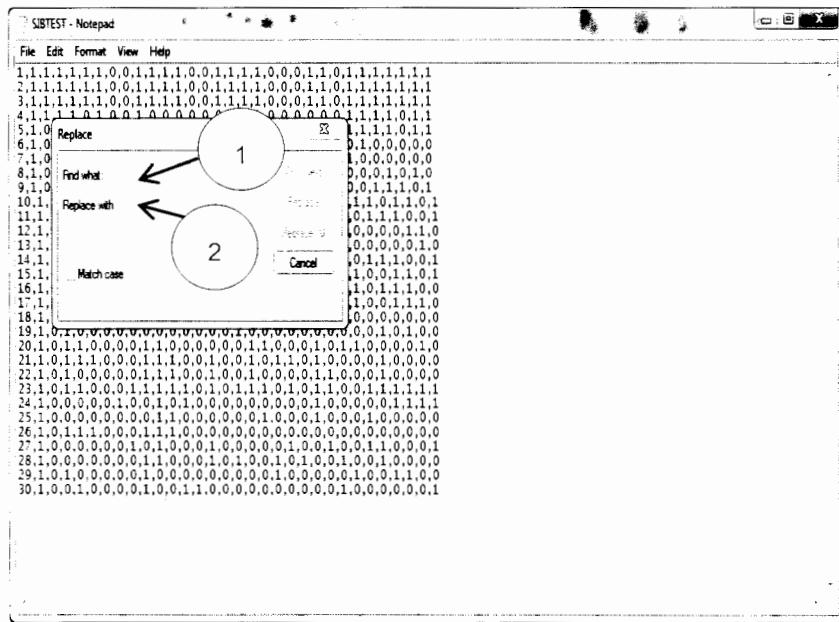
ภาพที่ 3-37 ตัวอย่างการบันทึกไฟล์ข้อมูลสำหรับวิเคราะห์ด้วยวิธี SIBTEST ในรูปแบบไฟล์ .dat

2. เลือก แฟ้ม แล้วเลือก บันทึกเป็น
3. ใส่ชื่อแฟ้ม
4. เลือก CSV (Comma delimited)
5. เลือกบันทึก เลือก Ok เลือก Yes



ภาพที่ 3-38 ตัวอย่างการจัดไฟล์ข้อมูลสำหรับวิเคราะห์ด้วยวิธี SIBTEST ในรูปแบบไฟล์ .dat
ในรูปแบบ NotePad

6. เปิดโปรแกรม NotePad

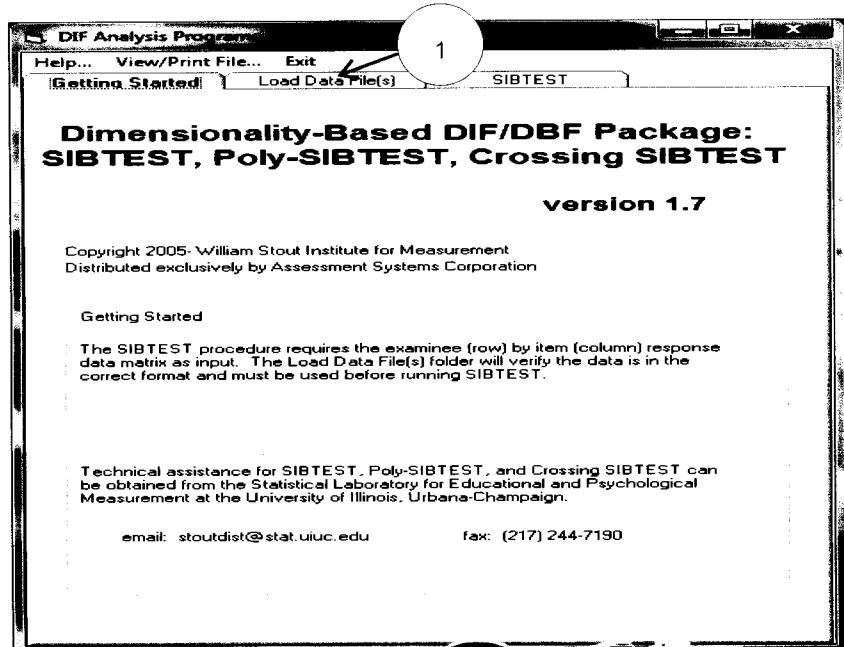


ภาพที่ 3-39 ตัวอย่างการจัดไฟล์ข้อมูลที่ลบจุลภาค (,) ออก และซิดของข้อมูลสำหรับวิเคราะห์ด้วยวิธี SIBTEST ในรูปแบบไฟล์ .dat

7. เลือกที่ Edit เลือก หมายเลข 1 ให้ใส่เครื่องหมาย (,) ช่องที่ 2 ไม่ต้องใส่อะไรแล้วไปเลือกที่ Replace เพื่อเอาจุลภาค (,) ออก และให้ข้อมูลซิดกัน
8. เลือกที่ File เลือก Save As
9. เลือกที่ All File แล้วไปตั้งชื่อ File ตามด้วย .dat และบันทึก

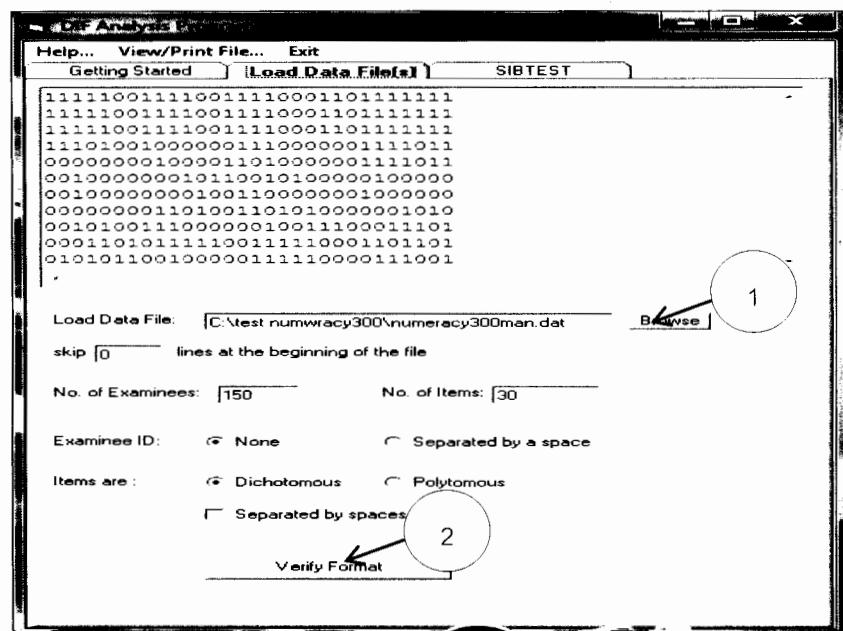
ขั้นตอนที่ 2 การวิเคราะห์การทำหน้าที่กันของข้อสอบด้วยวิธี SIBTEST มีดังนี้

1. เปิดโปรแกรม SIBTEST



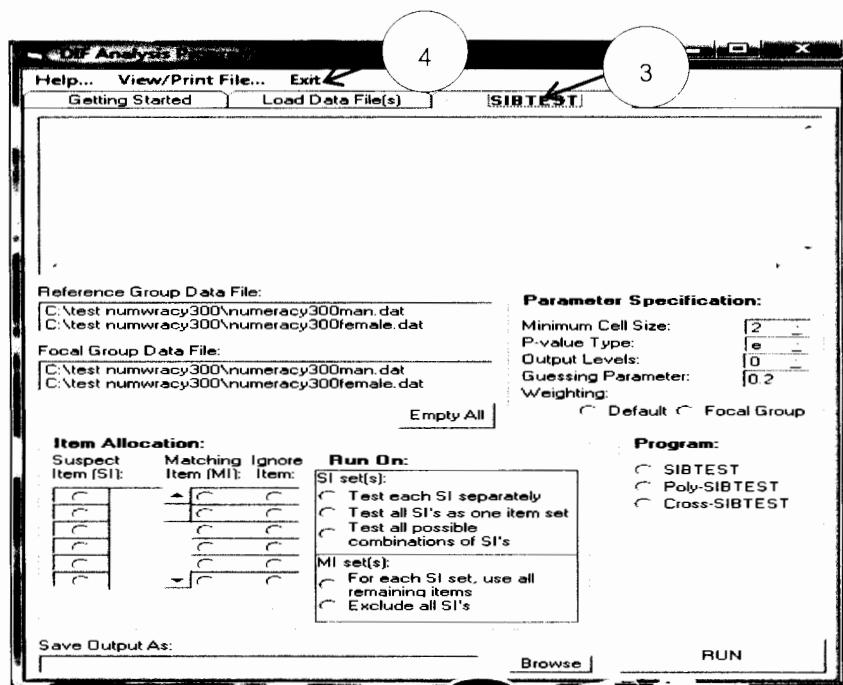
ภาพที่ 3-40 ตัวอย่างการวิเคราะห์ด้วยวิธี SIBTEST ในรูปแบบไฟล์ .dat เรียกไฟล์ข้อมูลวิเคราะห์

2. เลือกที่ Load Data File หมายเลข 1



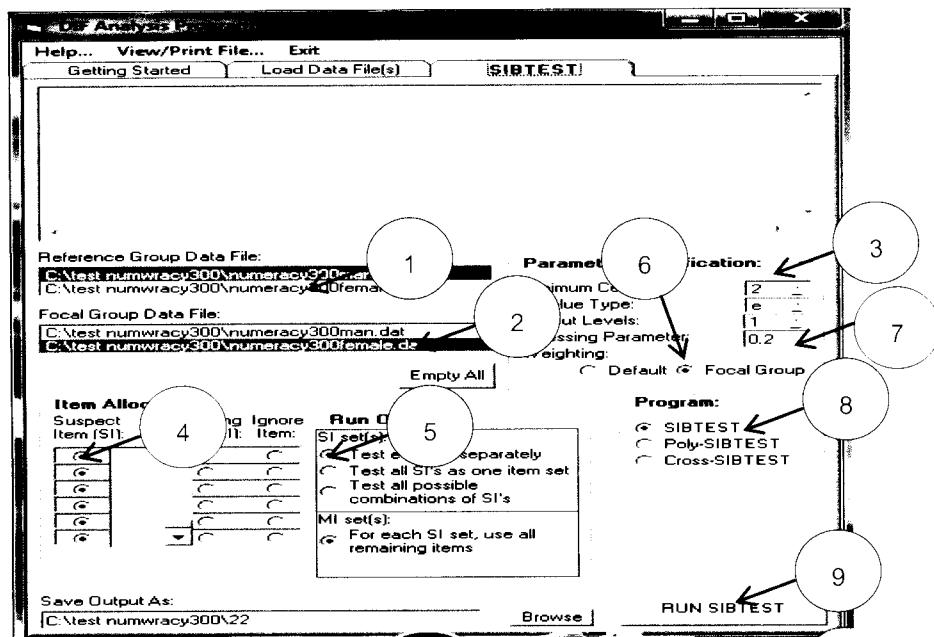
ภาพที่ 3-41 ตัวอย่างการวิเคราะห์ด้วยวิธี SIBTEST ในรูปแบบไฟล์ .dat เมื่อไฟล์ข้อมูลเข้าในโปรแกรม

3. เลือกที่ Browse หมายเลข 1
4. เลือกไฟล์ที่ Save ไว้ และตรวจเช็คข้อมูลให้เรียบร้อย
5. เลือกที่ Verify Fomat หมายเลข 2



ภาพที่ 3-42 ตัวอย่างการวิเคราะห์ด้วยวิธี SIBTEST ในรูปแบบไฟล์ .dat

6. เลือก SIBTEST เพื่อเตรียมไฟล์อีกไฟล์เรียกเข้ามา หมายเลข 3
7. เลือก Load Data File หมายเลข 4 และทำตาม ข้อ 3-5



ภาพที่ 3-43 ตัวอย่างการเลือกข้อมูลวิเคราะห์ด้วยวิธี SIBTEST ในรูปแบบไฟล์ .dat

8. เลือก Reference Group Data File หมายเลข 1 และหมายเลข 2
9. เลือก คูตรง หมายเลข 3 เพื่อตรวจเช็คความเรียบร้อย
10. เลือกตระง หมายเลข 4 เพื่อคลิกจำนวนข้อสอบให้ครบจำนวนข้อ
11. เลือกตระง หมายเลข 5 เพื่อ Test จำนวนข้อสอบ
12. เลือกตระง หมายเลข 6 เพื่อ เลือก Group
13. เลือกตระง หมายเลข 7 เพื่อ Save File Output
14. เลือกตระง หมายเลข 8 เพื่อเลือกโปรแกรม SIBTEST
15. เลือกตระง หมายเลข 9 เพื่อ Run โปรแกรม
16. พิจารณาว่าข้อใดมี DIF โดยดูจากค่า p-value ที่มีนัยสำคัญทางสถิติที่ระดับ .05

3. การวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติด้วยวิธี Mantel-Haenszel

การวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบด้วยวิธี Mantel-Haenszel โดยใช้โปรแกรม SPSS เป็นการทดสอบอัตราส่วนเบรียบเทียบด้วยไค-สแควร์ (Chi-square Statistic) ในการเปรียบเทียบความแตกต่างของ ค่าพารามิเตอร์ในกลุ่มอ้างอิงและกลุ่มเปรียบเทียบ ผู้วิจัยได้เตรียมการวิเคราะห์ กำหนดค่าตัวแปร ด้านภาษา ด้านคำนวน และด้านเหตุผล โดยกำหนดค่าตัวแปร เพศ ให้ตัวแปรเพศชายเป็น 1 และตัวแปรเพศหญิงเป็น 2 ดังนี้

ขั้นตอนที่ 1 เตรียมไฟล์ข้อมูลสำหรับวิเคราะห์

การวิเคราะห์ข้อมูลด้วยวิธี Mantel-Haenszel ผู้วิจัยได้เตรียมข้อมูลผลการตอบแบบทดสอบ ระหว่างกลุ่มอ้างอิง (R) และกลุ่มเปรียบเทียบ (F) มาจัดลงในตารางการณ์จาร์แบบ 2×2 (กลุ่มผู้สอบ 2 กลุ่ม \times ผลการตอบ 2 แบบ) ประกอบด้วยข้อมูล คือ (ID) จำนวนข้อสอบ (Gender) เพศชาย เพศหญิง (Answer) คำตอบถูกเป็น 1 คำตอบผิดเป็น 0 และ (Count) จำนวนผลการตอบแบบทดสอบของผู้สอบ ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน) ดังภาพที่ 3-44

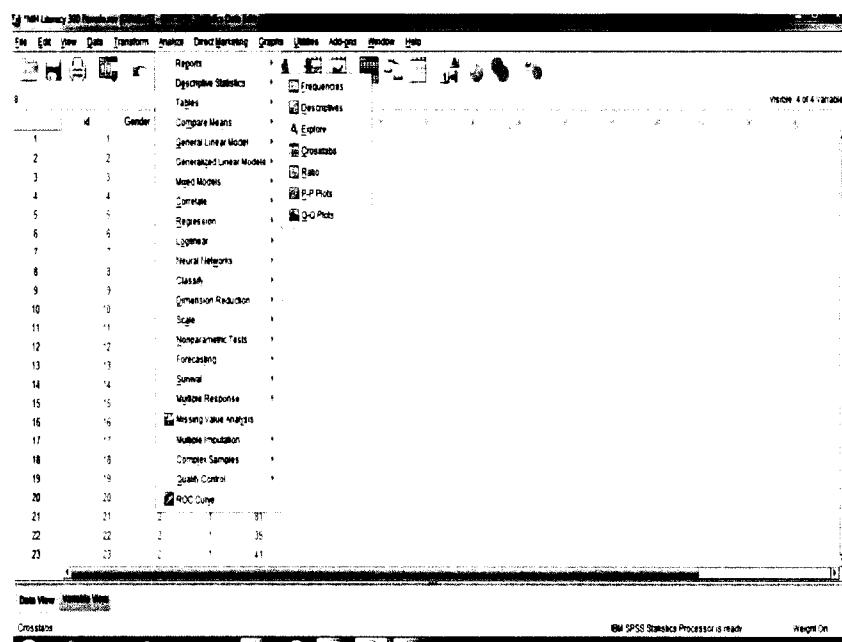
				Count	.0%	.0%	.0%	.0%	.0%	.0%	.0%
					1R	2R	3R	4R	5R	6R	7R
				1							
1		1	1	84							
2		2	1	55							
3		3	1	83							
4		4	1	96							
5		5	1	74							
6		6	1	104							
7		7	1	97							
8		8	1	73							
9		9	1	115							
10		10	1	59							
11		11	1	99							
12		12	1	88							
13		13	1	70							
14		14	1	68							
15		15	1	63							
16		16	1	61							
17		17	1	85							
18		18	1	62							
19		19	1	100							
20		20	1	93							
21		21	1	81							
22		22	1	34							
23		23	1	37							

ภาพที่ 3-44 ตัวอย่างการจัดไฟล์ข้อมูลสำหรับวิเคราะห์ด้วยวิธี Mantel-Haenszel ตารางการณ์จาร์แบบ 2×2

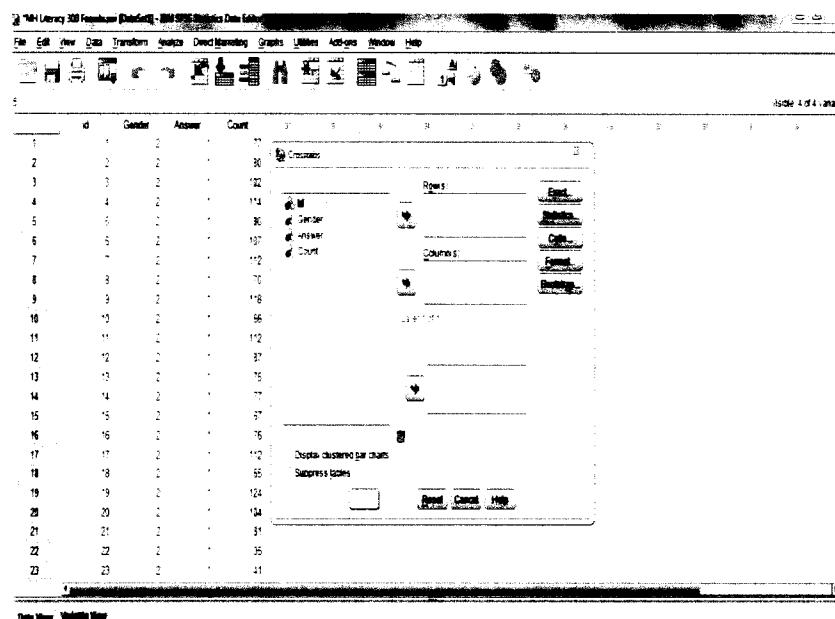
ขั้นตอนที่ 2 การตรวจสอบการทำหน้าที่กันของข้อสอบ

1. กำหนดตัวแปร ในด้านภาษา และด้านเหตุผล ให้เพศหญิงเป็นกลุ่มอ้างอิง (Referent Group: R) และเพศชายเป็นกลุ่มเปรียบเทียบ (Focal Group: F) ส่วนด้านคำนวน ให้เพศชายเป็นกลุ่มอ้างอิง (R) และเพศหญิงเป็นกลุ่มเปรียบเทียบ (F)

2. การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ (DIF) ด้วยวิธี Mantel-Haenszel เป็นการเปรียบเทียบผลการตอบข้อสอบระหว่างกลุ่มอ้างอิงและกลุ่มเปรียบเทียบ ในความแตกต่างของแต่ละพารามิเตอร์ โดยโปรแกรม SPSS คำสั่ง Analyze \rightarrow Descriptive Statistics \rightarrow Crosstabs ดังภาพที่ 3-43 ภาพที่ 3-44 และภาพที่ 3-45

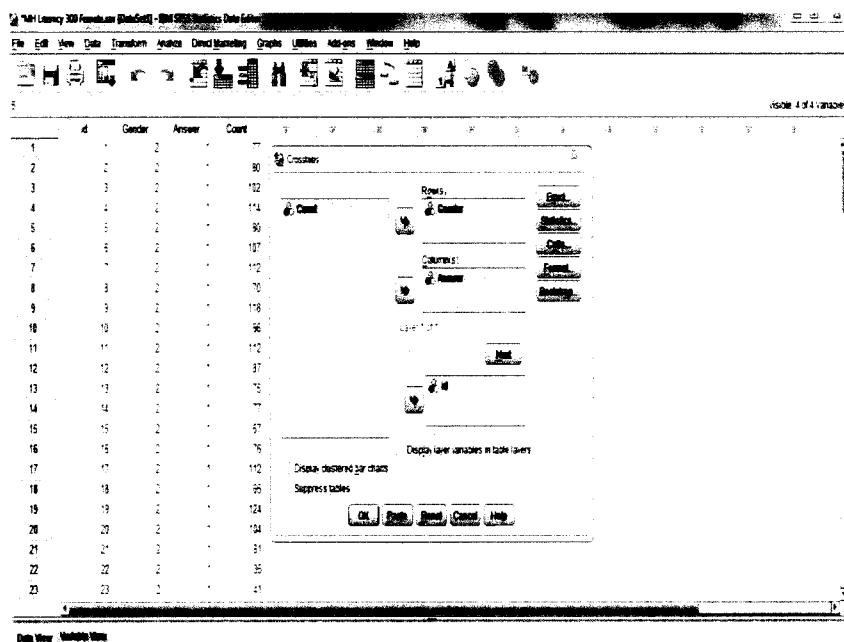


ภาพที่ 3-45 ตัวอย่างคำสั่งที่ Analysis การวิเคราะห์ด้วยวิธี Mantel-Haenszel โดยโปรแกรม SPSS



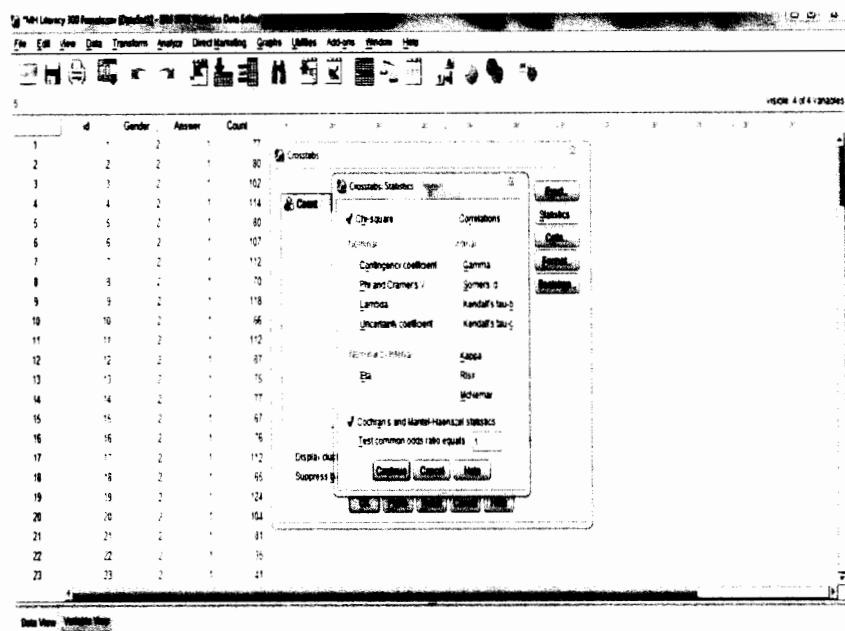
ภาพที่ 3-46 ตัวอย่างคำสั่งที่ Crosstabs การวิเคราะห์ด้วยวิธี Mantel-Haenszel โดยโปรแกรม SPSS

3. การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ (DIF) ด้วยวิธี Mantel-Haenszel
เป็นการเปรียบเทียบผลการตอบข้อสอบระหว่างกลุ่มอ้างอิงและกลุ่มเปรียบเทียบ ในความแตกต่าง
ของแต่ละพารามิเตอร์ โดยโปรแกรม SPSS โดยเลือก Gender ไปที่ช่อง Row(s) > Answer เลือก
ไปที่ช่อง Column(s) > ID เลือกไปที่ช่อง Layer 1 of 1 ดังภาพที่ 3-47



ภาพที่ 3-47 ตัวอย่างคำสั่งเรียกไฟล์ที่ต้องการจะวิเคราะห์ด้วยวิธี Mantel-Haenszel โดยโปรแกรม SPSS

4. การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ (DIF) ด้วยวิธี Mantel-Haenszel
เป็นการเปรียบเทียบผลการตอบข้อสอบระหว่างกลุ่มอ้างอิงและกลุ่มเปรียบเทียบ ในความแตกต่าง
ของแต่ละพารามิเตอร์ โดยโปรแกรม SPSS โดยเลือก Statistics > เลือกที่ Chi-square >
Cochram's and Mantel-Haenszel Statistics ดังภาพที่ 3-48



ภาพที่ 3-48 ตัวอย่างคำสั่งเลือก Chi-Square การวิเคราะห์ด้วยวิธี Mantel-Haenszel โดยโปรแกรม SPSS

Case Processing Summary						
	Cases					
	Valid		Missing		Total	
	N	Percent	N	Percent	N	Percent
Sample * Response	270	100.0%	0	0.0%	270	100.0%
Subject * Response	270	100.0%	0	0.0%	270	100.0%
Methods * Response	270	100.0%	0	0.0%	270	100.0%

Sample * Response						
Crosstab						
Count	Response					
	NO-DIF		DIF		Total	
	Sample 300 persons	242	28		270	
Total	242	28		270		

Chi-Square Tests		Value
Pearson Chi-Square		a
N of Valid Cases		270

a. No statistics are computed because Sample is a constant.

Subject * Response						
Crosstab						
Count	Response					
	NO-DIF		DIF		Total	
	Subject Literacy	80	10		90	
Subject Numeracy	82	8		90		
Subject Reasoning	80	10		90		
Total	242	28		270		

ภาพที่ 3-49 ตัวอย่างผลการวิเคราะห์ Chi-Square วิธี Mantel-Haenszel

5. การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ เพื่อวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบ ทำการวิเคราะห์แบบทดสอบทุกข้อจากนั้นพิจารณาว่าข้อใดมี DIF โดยดูจากค่าไค-สแควร์ (Chi-square Statistic) ที่มีนัยสำคัญทางสถิติที่ระดับ .05

จากการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบ ผู้วิจัยได้ดำเนินการรวมผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ ของวิธี IRT- LR วิธี SIBTEST และวิธี Mantel-Haenszel ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน) และขั้นตอนการเปรียบเทียบผลการสอบการทำหน้าที่ต่างกันของข้อสอบ มีดังนี้

ระยะที่ 3 การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้วยวิธีการทดสอบอัตราส่วนไลค์ลิลูด วิธีซิปเพลท์ และวิธีแมนเทล-แ昏ส์เซล ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน)

การเปรียบผลการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ มีขั้นตอนดำเนินการ ดังนี้

เริ่มต้น



ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ทั้ง 3 ด้าน ได้แก่ ด้านภาษา ด้านคำนวณ และด้านเหตุผล ด้วยวิธี IRT- LR วิธี SIBTEST และวิธี MH ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน)



เปรียบเทียบผลการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบ

ด้านภาษา	1) วิธี IRT- LR กับ วิธี SIBTEST ขนาดกลุ่มตัวอย่าง 3 ขนาด คือ ^{300 คน} ^{1,000 คน} ^{2,000 คน}	2) วิธี IRT- LR กับ วิธี MH ขนาดกลุ่มตัวอย่าง 3 ขนาด คือ ^{300 คน} ^{1,000 คน} ^{2,000 คน}	3) วิธี SIBTEST กับ วิธี MH ขนาดกลุ่มตัวอย่าง 3 ขนาด คือ ^{300 คน} ^{1,000 คน} ^{2,000 คน}
	1) วิธี IRT- LR กับ วิธี SIBTEST ขนาดกลุ่มตัวอย่าง 3 ขนาด คือ ^{300 คน} ^{1,000 คน} ^{2,000 คน}	2) วิธี IRT- LR กับ วิธี MH ขนาดกลุ่มตัวอย่าง 3 ขนาด คือ ^{300 คน} ^{1,000 คน} ^{2,000 คน}	3) วิธี SIBTEST กับ วิธี MH ขนาดกลุ่มตัวอย่าง 3 ขนาด คือ ^{300 คน} ^{1,000 คน} ^{2,000 คน}
	1) วิธี IRT- LR กับ วิธี SIBTEST ขนาดกลุ่มตัวอย่าง 3 ขนาด คือ ^{300 คน} ^{1,000 คน} ^{2,000 คน}	2) วิธี IRT- LR กับ วิธี MH ขนาดกลุ่มตัวอย่าง 3 ขนาด คือ ^{300 คน} ^{1,000 คน} ^{2,000 คน}	3) วิธี SIBTEST กับ วิธี MH ขนาดกลุ่มตัวอย่าง 3 ขนาด คือ ^{300 คน} ^{1,000 คน} ^{2,000 คน}



สรุปผลการเปรียบเทียบจำนวนข้อสอบที่ตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ ทั้ง 3 ด้าน



สิ้นสุด

ภาพที่ 3-50 ขั้นตอนการเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ

ผู้วิจัยได้เปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้วยวิธี IRT-LR วิธี SIBTEST และ วิธี MH ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน) ตามสมมติฐานการวิจัย ดังนี้

1. เปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ขั้นประถมศึกษาปีที่ 3 ด้านภาษา กลุ่มตัวอย่างขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และ ขนาดใหญ่ (2,000 คน) ระหว่างวิธี IRT-LR กับวิธี SIBTEST

2. เปรียบเทียบผลการตรวจสอบการทําหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติชั้นประถมศึกษาปีที่ 3 ด้านภาษา กลุ่มตัวอย่างขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน) ระหว่างวิธี IRT-LR กับ วิธี MH

3. เปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้านภาษา กลุ่มตัวอย่างขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และ ขนาดใหญ่ (2,000 คน) ระหว่างวิธี SIBTEST กับวิธี MH

4. เปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติชั้นประถมศึกษาปีที่ 3 ด้านความน่าสนใจ (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน) ระหว่างวิธี IRT-LR กับวิธี SIBTEST

5. เปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้านคำนวน กลุ่มตัวอย่างขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และ ขนาดใหญ่ (2,000 คน) ระหว่างวิธี IRT-LR กับ วิธี MH

6. เปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้านความคุ้มครอง กลุ่มตัวอย่างขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และ ขนาดใหญ่ (2,000 คน) ระหว่างวิธี SIBTEST กับวิธี MH

7. เปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้านเหตุผล กลุ่มตัวอย่างขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และ ขนาดใหญ่ (2,000 คน) ระหว่างวิธี IRT-LR กับวิธี SIBTEST

8. เปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ขั้นประถมศึกษาปีที่ 3 ด้านเหตุผล กลุ่มตัวอย่างขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และ ขนาดใหญ่ (2,000 คน) ระหว่างวิธี IRT-LR กับ วิธี MH

9. เปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ขั้นปฐมศึกษาปีที่ 3 ด้านเหตุผล กลุ่มตัวอย่างขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และ ขนาดใหญ่ (2,000 คน) ระหว่างวิธี SIBTEST กับวิธี MH

บทที่ 4 ผลการวิจัย

การวิจัยครั้งนี้มีวัตถุประสงค์เพื่อวิเคราะห์คุณภาพของข้อสอบ ตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ และเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ระหว่างวิธีการทดสอบอัตราส่วนไลค์ลิชุด วิธีซิปเหล็ฟ และวิธีแมนเทล-แยนส์เซล ภายใต้เงื่อนไข ขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน) ผู้วิจัยนำเสนอผลการวิจัยเป็น 3 ตอน ดังนี้

ตอนที่ 1 ผลการวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 โดยใช้หลักการของทฤษฎีการตอบสนองข้อสอบ แบบ 3 พารามิเตอร์ ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน)

ตอนที่ 2 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้วยวิธีการทดสอบอัตราส่วนไลค์ลิชุด วิธีซิปเหล็ฟ และวิธีแมนเทล-แยนส์เซล ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน)

ตอนที่ 3 ผลการเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้วยวิธีการทดสอบอัตราส่วนไลค์ลิชุด วิธีซิปเหล็ฟ และวิธีแมนเทล-แยนส์เซล ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน)

สัญลักษณ์และอักษรย่อที่ใช้ในการเสนอผลการวิเคราะห์ข้อมูล มีดังนี้

- a หมายถึง จำนวนตัวอย่าง
- b หมายถึง ค่าอำนาจจำแนกของข้อสอบ
- c หมายถึง ค่าความยากของข้อสอบ
- d หมายถึง ค่าโอกาสการเดาของข้อสอบ

ตอนที่ 1 ผลการวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 โดยใช้หลักการของทฤษฎีการตอบสนองข้อสอบ แบบ 3 พารามิเตอร์ ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน)

การวิเคราะห์คุณภาพข้อสอบ NT ระดับชั้นประถมศึกษาปีที่ 3 เป็นการวิเคราะห์โดยใช้ทฤษฎีการตอบสนองข้อสอบแบบ 3 พารามิเตอร์ คือ ค่าอำนาจจำแนกของข้อสอบ (a) ค่าความยากของข้อสอบ (b) และค่าโอกาสการเดาของข้อสอบ (c) โดยใช้โปรแกรมสำเร็จรูป Xcalibre Version 4.2.2 ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน) การประมาณค่าพารามิเตอร์ของข้อสอบจะมีเกณฑ์ในการคัดเลือกข้อสอบดังนี้

เกณฑ์การคัดเลือกข้อสอบ (Urry, 1977)

1. ค่าอำนาจจำแนกของข้อสอบ (a) มีค่าตั้งแต่ 0.50 ถึง 2.50
2. ค่าความยากของข้อสอบ (b) มีค่าตั้งแต่ -2.50 ถึง 2.50
3. ค่าการเดาของข้อสอบ (c) มีค่าไม่เกิน 0.30

เกณฑ์การประเมินค่าความยากของข้อสอบ (b) ระดับชั้นประถมศึกษาปีที่ 3 ดังนี้

ค่าความยากเฉลี่ยมากกว่า	2.0000	หมายถึง ข้อสอบยากมาก
ค่าความยากเฉลี่ยตั้งแต่	1.0001 ถึง 2.0000	หมายถึง ข้อสอบยาก
ค่าความยากเฉลี่ยตั้งแต่	0.5001 ถึง 1.0000	หมายถึง ข้อสอบค่อนข้างยาก
ค่าความยากเฉลี่ยตั้งแต่	-0.4999 ถึง 0.5000	หมายถึง ข้อสอบปานกลาง
ค่าความยากเฉลี่ยตั้งแต่	-0.9999 ถึง -0.5000	หมายถึง ข้อสอบค่อนข้างง่าย
ค่าความยากเฉลี่ยตั้งแต่	-1.9999 ถึง -1.0000	หมายถึง ข้อสอบง่าย
ค่าความยากเฉลี่ยน้อยกว่า	-2.0000	หมายถึง ข้อสอบง่ายมาก

ตารางที่ 4-1 ผลการวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านภาษา

จำนวน 30 ข้อ โดยใช้ทฤษฎีการตอบสนองข้อสอบ (IRT) แบบ 3 พารามิเตอร์ กลุ่มตัวอย่างขนาดเล็ก (300 คน)

แบบทดสอบ NT	ข้อที่	ค่าพารามิเตอร์		
		a	b	c
ด้านภาษา	1	0.314	1.066	0.264
	2	0.377	1.709	0.255
	3	0.574	0.499	0.364
	4	0.549	-0.488	0.254
	5	0.646	0.736	0.250
	6	0.423	-0.538	0.259

ตารางที่ 4-1 (ต่อ)

แบบทดสอบ NT	ข้อที่	ค่าพารามิเตอร์		
		a	b	c
ด้านภาษา	7	0.680	-0.438	0.247
	8	0.558	1.209	0.258
	9	0.516	-1.120	0.252
	10	0.582	1.790	0.263
	11	0.472	-0.581	0.255
	12	0.628	0.280	0.248
	13	0.553	1.166	0.258
	14	0.612	1.007	0.251
	15	0.602	1.348	0.246
	16	0.579	1.257	0.252
	17	0.782	-0.189	0.248
	18	0.743	1.158	0.238
	19	0.839	-0.660	0.252
	20	0.654	-0.173	0.251
	21	0.478	0.802	0.258
	22	0.931	2.316	0.214
	23	0.695	2.829	0.231
	24	0.676	1.385	0.250
	25	0.652	2.153	0.278
	26	0.647	1.212	0.262
	27	0.656	1.289	0.251
	28	0.617	0.245	0.256
	29	0.611	0.910	0.244
	30	0.788	0.629	0.238
ค่าความเที่ยงทั้งฉบับ (Reliability)		0.712		

จากตารางที่ 4-1 ข้อสอบด้านภาษา จากจำนวนข้อสอบ 30 ข้อ มีข้อสอบที่ผ่านเกณฑ์การวิเคราะห์คุณภาพข้อสอบ จำนวน 30 ข้อ เป็นข้อสอบที่มีค่าอำนาจจำแนกของข้อสอบ (a) ตั้งแต่ 0.314 ถึง 0.931 ค่าความยากของข้อสอบ (b) ตั้งแต่ -1.120 ถึง 2.829 และค่าโอกาสการเดาของข้อสอบ (c) ตั้งแต่ 0.214 ถึง 0.364 และมีค่าความเที่ยงทั้งฉบับ 0.712

ตารางที่ 4-2 ผลการวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านคำนวน
จำนวน 30 ข้อ โดยใช้ทฤษฎีการตอบสนองข้อสอบ (IRT) แบบ 3 พารามิเตอร์
กลุ่มตัวอย่างขนาดเล็ก (300 คน)

แบบทดสอบ NT	ข้อที่	ค่าพารามิเตอร์		
		a	b	c
ด้านคำนวน	1	0.731	1.430	0.079
	2	0.624	1.351	0.081
	3	0.563	2.395	0.248
	4	0.940	1.362	0.339
	5	0.901	1.315	0.225
	6	0.895	2.137	0.216
	7	0.795	2.408	0.229
	8	0.767	0.701	0.247
	9	0.678	0.735	0.254
	10	0.791	1.365	0.241
	11	0.699	0.893	0.260
	12	0.715	2.537	0.256
	13	0.791	1.281	0.250
	14	0.983	1.250	0.230
	15	0.738	1.732	0.251
	16	0.860	0.905	0.248
	17	0.540	0.271	0.256
	18	0.676	1.457	0.239
	19	0.648	1.636	0.246
	20	0.527	0.728	0.257
	21	0.807	1.772	0.254
	22	0.699	2.112	0.238
	23	0.921	3.503	0.222
	24	0.706	1.867	0.255
	25	0.733	1.718	0.239
	26	0.915	1.592	0.235
	27	0.807	0.662	0.241

ตารางที่ 4-2 (ต่อ)

แบบทดสอบ NT	ข้อที่	ค่าพารามิเตอร์		
		a	b	c
ด้านคำนวณ	28	1.024	1.004	0.242
	29	0.746	0.989	0.249
	30	1.073	1.180	0.222
ค่าความเที่ยงทั้งฉบับ (Reliability)		0.763		

จากตารางที่ 4-2 ข้อสอบด้านคำนวณ จากจำนวนข้อสอบ 30 ข้อ มีข้อสอบที่ผ่านเกณฑ์การวิเคราะห์คุณภาพข้อสอบ จำนวน 30 ข้อ เป็นข้อสอบที่มีค่าอำนาจจำแนกของข้อสอบ (a) ตั้งแต่ 0.527 ถึง 1.073 ค่าความยากของข้อสอบ (b) ตั้งแต่ 0.271 ถึง 3.503 และค่าโอกาสการเดาของข้อสอบ (c) ตั้งแต่ 0.079 ถึง 0.339 และมีค่าความเที่ยงทั้งฉบับ 0.763

ตารางที่ 4-3 ผลการวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านเหตุผล จำนวน 30 ข้อ โดยใช้ทฤษฎีการตอบสนองข้อสอบ (IRT) แบบ 3 พารามิเตอร์ กลุ่มตัวอย่างขนาดเล็ก (300 คน)

แบบทดสอบ NT	ข้อที่	ค่าพารามิเตอร์		
		a	b	c
ด้านเหตุผล	1	0.629	0.628	0.286
	2	0.806	0.400	0.275
	3	0.712	0.945	0.379
	4	0.787	-0.432	0.248
	5	0.455	0.390	0.251
	6	0.710	0.098	0.248
	7	0.733	0.327	0.245
	8	0.752	1.234	0.270
	9	0.679	0.789	0.251
	10	0.798	1.432	0.246
	11	0.609	0.497	0.252
	12	0.707	0.476	0.251
	13	0.552	1.053	0.254
	14	0.899	1.596	0.241

ตารางที่ 4-3 (ต่อ)

แบบทดสอบ NT	ข้อที่	ค่าพารามิเตอร์		
		a	b	c
ด้านเหตุผล	15	0.701	1.292	0.256
	16	0.954	2.022	0.254
	17	0.781	2.820	0.247
	18	0.821	0.850	0.243
	19	0.668	1.189	0.249
	20	0.772	2.384	0.270
	21	0.791	0.403	0.244
	22	0.787	0.597	0.248
	23	0.781	0.723	0.239
	24	0.821	0.331	0.247
	25	0.610	1.169	0.258
	26	0.786	0.637	0.251
	27	0.802	1.811	0.234
	28	0.737	2.031	0.236
	29	0.933	1.555	0.228
	30	1.289	0.886	0.220
ค่าความเที่ยงทั้งฉบับ (Reliability)		0.766		

จากตารางที่ 4-3 ข้อสอบด้านเหตุผล จากจำนวนข้อสอบ 30 ข้อ มีข้อสอบที่ผ่านเกณฑ์การวิเคราะห์คุณภาพข้อสอบ จำนวน 30 ข้อ เป็นข้อสอบที่มีค่าอำนาจจำแนกของข้อสอบ (a) ตั้งแต่ 0.455 ถึง 1.289 ค่าความยากของข้อสอบ (b) ตั้งแต่ -0.432 ถึง 2.820 และค่าโอกาสการเดาของข้อสอบ (c) ตั้งแต่ 0.220 ถึง 0.379 และมีค่าความเที่ยงทั้งฉบับ 0.766

ตารางที่ 4-4 ผลการวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านภาษา
จำนวน 30 ข้อ โดยใช้ทฤษฎีการตอบสนองข้อสอบ (IRT) แบบ 3 พารามิเตอร์
กลุ่มตัวอย่างขนาดกลาง (1,000 คน)

แบบทดสอบ NT	ข้อที่	ค่าพารามิเตอร์		
		a	b	c
ด้านภาษา	1	0.337	0.483	0.188
	2	0.323	3.475	0.209
	3	0.644	0.443	0.328

ตารางที่ 4-4 (ต่อ)

แบบทดสอบ NT ด้านภาษา	ข้อที่	ค่าพารามิเตอร์		
		a	b	c
	4	0.503	-0.950	0.253
	5	0.575	0.087	0.242
	6	0.571	-0.749	0.254
	7	0.849	-0.710	0.244
	8	0.609	0.873	0.247
	9	0.587	-1.207	0.251
	10	0.645	1.972	0.283
	11	0.469	-0.572	0.262
	12	0.760	0.285	0.246
	13	0.523	1.859	0.269
	14	0.786	0.796	0.244
	15	0.481	1.965	0.270
	16	0.704	1.120	0.245
	17	0.783	-0.400	0.246
	18	0.648	1.100	0.241
	19	1.010	-0.903	0.246
	20	0.742	-0.408	0.243
	21	0.679	0.627	0.252
	22	0.993	2.133	0.200
	23	0.556	2.297	0.246
	24	0.618	0.882	0.243
	25	0.702	2.695	0.303
	26	0.440	1.510	0.280
	27	0.759	1.070	0.242
	28	0.685	0.516	0.262
	29	0.587	1.396	0.236
	30	0.766	0.302	0.236
ค่าความเที่ยงทั้งฉบับ (Reliability)		0.738		

จากตารางที่ 4-4 ข้อสอบด้านภาษา จากจำนวนข้อสอบ 30 ข้อ มีข้อสอบที่ผ่านเกณฑ์การวิเคราะห์คุณภาพข้อสอบ จำนวน 30 ข้อ เป็นข้อสอบที่มีค่าอำนาจจำแนกของข้อสอบ (a) ตั้งแต่ 0.323 ถึง 1.010 ค่าความยากของข้อสอบ (b) ตั้งแต่ -1.207 ถึง 3.475 และค่าโอกาสการเดาของข้อสอบ (c) ตั้งแต่ 0.188 ถึง 0.328 และมีค่าความเที่ยงทั้งฉบับ 0.738

ตารางที่ 4-5 ผลการวิเคราะห์คุณภาพของข้อสอบ NT ขั้นประถมศึกษาปีที่ 3 ด้านคำนวน จำนวน 30 ข้อ โดยใช้ทฤษฎีการตอบสนองข้อสอบ (IRT) แบบ 3 พารามิเตอร์ กลุ่มตัวอย่างขนาดกลาง (1,000 คน)

แบบทดสอบ NT	ข้อที่	ค่าพารามิเตอร์		
		a	b	c
ด้านคำนวน	1	0.619	1.672	0.080
	2	0.827	0.733	0.077
	3	0.796	2.426	0.277
	4	1.188	0.783	0.307
	5	0.843	1.444	0.214
	6	0.959	1.484	0.202
	7	1.044	1.771	0.197
	8	1.057	0.144	0.244
	9	0.958	0.358	0.248
	10	1.008	1.043	0.259
	11	0.804	1.019	0.271
	12	0.566	2.370	0.252
	13	0.593	0.748	0.250
	14	0.786	1.835	0.237
	15	0.777	1.720	0.262
	16	1.163	0.710	0.243
	17	0.750	0.016	0.253
	18	0.803	1.040	0.228
	19	0.745	1.005	0.241
	20	0.524	0.815	0.254
	21	0.654	1.817	0.238
	22	1.147	2.452	0.216
	23	1.165	3.688	0.203
	24	0.918	1.069	0.258

ตารางที่ 4-5 (ต่อ)

แบบทดสอบ NT	ข้อที่	ค่าพารามิเตอร์		
		a	b	c
ด้านคำนวน	25	0.943	0.935	0.224
	26	0.822	1.529	0.250
	27	0.761	0.567	0.236
	28	0.994	1.384	0.229
	29	0.823	1.257	0.261
	30	0.978	1.127	0.212
ค่าความเที่ยงทั้งฉบับ (Reliability)		0.791		

จากตารางที่ 4-5 ข้อสอบด้านคำนวน จากจำนวนข้อสอบ 30 ข้อ มีข้อสอบที่ผ่านเกณฑ์การวิเคราะห์คุณภาพข้อสอบ จำนวน 30 ข้อ เป็นข้อสอบที่มีค่าอำนาจจำแนกของข้อสอบ (a) ตั้งแต่ 0.524 ถึง 1.188 ค่าความยากของข้อสอบ (b) ตั้งแต่ 0.016 ถึง 3.688 และค่าโอกาสการเดาของข้อสอบ (c) ตั้งแต่ 0.077 ถึง 0.307 และมีค่าความเที่ยงทั้งฉบับ 0.791

ตารางที่ 4-6 ผลการวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษานิปีที่ 3 ด้านเหตุผล จำนวน 30 ข้อ โดยใช้ทฤษฎีการตอบสนองข้อสอบ (IRT) แบบ 3 พารามิเตอร์ กลุ่มตัวอย่างขนาดกลาง (1,000 คน)

แบบทดสอบ NT	ข้อที่	ค่าพารามิเตอร์		
		a	b	c
ด้านเหตุผล	1	0.686	0.177	0.247
	2	0.796	0.020	0.236
	3	0.655	0.317	0.336
	4	0.761	-0.533	0.246
	5	0.537	-0.088	0.243
	6	0.581	-0.067	0.250
	7	0.754	0.235	0.241
	8	0.806	0.801	0.273
	9	0.672	0.232	0.252
	10	0.833	1.119	0.237
	11	0.435	0.705	0.253
	12	0.633	0.728	0.238
	13	0.496	0.774	0.251

ตารางที่ 4-6 (ต่อ)

แบบทดสอบ NT	ข้อที่	ค่าพารามิเตอร์		
		a	b	c
ด้านเหตุผล	14	0.661	1.315	0.240
	15	0.646	1.548	0.256
	16	0.908	1.940	0.266
	17	0.999	2.908	0.231
	18	0.694	0.903	0.226
	19	0.530	0.965	0.235
	20	0.925	3.056	0.270
	21	0.806	0.675	0.254
	22	0.964	0.398	0.262
	23	0.934	0.142	0.227
	24	1.113	0.113	0.248
	25	0.516	1.744	0.271
	26	0.776	0.409	0.236
	27	0.901	1.144	0.217
	28	0.830	1.224	0.202
	29	0.884	1.677	0.208
	30	1.115	1.090	0.225
ค่าความเที่ยงทั้งฉบับ (Reliability)		0.784		

จากตารางที่ 4-6 ข้อสอบด้านเหตุผล จากจำนวนข้อสอบ 30 ข้อ มีข้อสอบที่ผ่านเกณฑ์การวิเคราะห์คุณภาพข้อสอบ จำนวน 30 ข้อ เป็นข้อสอบที่มีค่าอำนาจจำแนกของข้อสอบ (a) ตั้งแต่ 0.435 ถึง 1.115 ค่าความยากของข้อสอบ (b) ตั้งแต่ -0.533 ถึง 3.056 และค่าโอกาสการเดาของข้อสอบ (c) ตั้งแต่ 0.202 ถึง 0.336 และมีค่าความเที่ยงทั้งฉบับ 0.784

ตารางที่ 4-7 ผลการวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านภาษา จำนวน 30 ข้อ โดยใช้ทฤษฎีการตอบสนองข้อสอบ (IRT) แบบ 3 พารามิเตอร์ กลุ่มตัวอย่างขนาดใหญ่ (2,000 คน)

แบบทดสอบ NT	ข้อที่	ค่าพารามิเตอร์		
		a	b	c
ด้านภาษา	1	0.480	0.453	0.211
	2	0.585	2.064	0.216

ตารางที่ 4-7 (ต่อ)

แบบทดสอบ NT	ข้อที่	ค่าพารามิเตอร์		
		a	b	c
ด้านภาษา	3	0.545	-0.030	0.208
	4	0.625	-0.931	0.214
	5	0.539	0.153	0.210
	6	0.708	-0.772	0.215
	7	1.089	-0.432	0.211
	8	0.670	0.614	0.208
	9	0.627	-1.140	0.215
	10	0.503	1.508	0.231
	11	0.588	-0.582	0.222
	12	0.928	0.051	0.201
	13	0.415	1.520	0.216
	14	0.822	0.575	0.199
	15	0.496	1.260	0.218
	16	0.698	0.525	0.201
	17	1.001	-0.398	0.213
	18	0.733	0.733	0.203
	19	0.949	-0.902	0.211
	20	0.737	-0.384	0.213
	21	0.594	0.546	0.206
	22	0.984	1.606	0.179
	23	0.546	1.904	0.207
	24	0.559	0.698	0.206
	25	0.438	2.113	0.232
	26	0.318	1.339	0.220
	27	0.722	1.021	0.222
	28	0.689	0.193	0.220
	29	0.555	0.877	0.194
	30	0.741	0.102	0.203
ค่าความเที่ยงทั้งฉบับ (Reliability)			0.729	

จากตารางที่ 4-7 ข้อสอบด้านภาษา จำนวนข้อสอบ 30 ข้อ มีข้อสอบที่ผ่านเกณฑ์การวิเคราะห์คุณภาพข้อสอบ จำนวน 30 ข้อ เป็นข้อสอบที่มีค่าอำนาจจำแนกของข้อสอบ (a) ตั้งแต่ 0.318 ถึง 1.089 ค่าความยากของข้อสอบ (b) ตั้งแต่ -1.140 ถึง 2.113 และค่าโอกาสการเดาของข้อสอบ (c) ตั้งแต่ 0.179 ถึง 0.232 และมีค่าความเที่ยงทั้งฉบับ 0.729

ตารางที่ 4-8 ผลการวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านคำนวน จำนวน 30 ข้อ โดยใช้ทฤษฎีการตอบสนองข้อสอบ (IRT) แบบ 3 พารามิเตอร์ กลุ่มตัวอย่างขนาดใหญ่ (2,000 คน)

แบบทดสอบ NT ด้านคำนวน	ข้อที่	ค่าพารามิเตอร์		
		a	b	c
	1	0.695	1.620	0.080
	2	1.007	1.267	0.171
	3	0.693	2.681	0.276
	4	1.061	0.767	0.241
	5	0.855	1.335	0.206
	6	0.854	1.748	0.200
	7	1.037	1.640	0.203
	8	0.969	0.286	0.257
	9	0.839	0.449	0.243
	10	0.947	0.989	0.239
	11	0.746	1.023	0.255
	12	0.480	2.713	0.262
	13	0.638	0.590	0.234
	14	0.589	2.214	0.251
	15	0.711	1.607	0.251
	16	1.002	0.986	0.243
	17	0.744	0.061	0.255
	18	0.829	1.290	0.220
	19	0.903	1.340	0.234
	20	0.483	1.018	0.256
	21	0.605	2.024	0.244

ตารางที่ 4-8 (ต่อ)

แบบทดสอบ NT	ข้อที่	ค่าพารามิเตอร์		
		a	b	c
ด้านคำนวน	22	1.125	2.886	0.213
	23	1.208	4.000	0.193
	24	0.960	1.341	0.271
	25	0.997	1.068	0.233
	26	0.921	1.445	0.239
	27	0.773	0.796	0.234
	28	0.714	1.755	0.214
	29	0.727	1.651	0.240
	30	0.973	1.345	0.192
ค่าความเที่ยงทั้งฉบับ (Reliability)			0.775	

จากตารางที่ 4-8 ข้อสอบด้านคำนวน จากจำนวนข้อสอบ 30 ข้อ มีข้อสอบที่ผ่านเกณฑ์การวิเคราะห์คุณภาพข้อสอบ จำนวน 30 ข้อ เป็นข้อสอบที่มีค่าอำนาจจำแนกของข้อสอบ (a) ตั้งแต่ 0.480 ถึง 1.208 ค่าความยากของข้อสอบ (b) ตั้งแต่ 0.061 ถึง 4.000 และค่าโอกาสการเดาของข้อสอบ (c) ตั้งแต่ 0.080 ถึง 0.276 และมีค่าความเที่ยงทั้งฉบับ 0.775

ตารางที่ 4-9 ผลการวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านเหตุผล จำนวน 30 ข้อ โดยใช้ทฤษฎีการตอบสนองข้อสอบ (IRT) แบบ 3 พารามิเตอร์ กลุ่มตัวอย่างขนาดใหญ่ (2,000 คน)

แบบทดสอบ NT	ข้อที่	ค่าพารามิเตอร์		
		a	b	c
ด้านเหตุผล	1	0.681	0.413	0.257
	2	0.925	0.214	0.248
	3	0.746	0.309	0.315
	4	0.667	-0.442	0.241
	5	0.542	-0.107	0.237
	6	0.551	-0.100	0.252
	7	0.703	0.363	0.230
	8	0.822	0.758	0.264
	9	0.613	0.311	0.255

ตารางที่ 4-9 (ต่อ)

แบบทดสอบ NT	ข้อที่	ค่าพารามิเตอร์		
		a	b	c
ด้านเหตุผล	10	0.737	1.113	0.225
	11	0.503	0.455	0.237
	12	0.540	0.549	0.238
	13	0.584	0.900	0.241
	14	0.794	1.084	0.225
	15	0.646	1.604	0.252
	16	1.004	1.738	0.275
	17	0.898	2.323	0.227
	18	0.730	0.756	0.222
	19	0.628	0.702	0.221
	20	0.955	2.968	0.255
	21	0.793	0.751	0.244
	22	0.792	0.571	0.251
	23	0.833	0.286	0.216
	24	1.068	0.200	0.240
	25	0.508	1.613	0.260
	26	0.697	0.460	0.214
	27	0.843	1.188	0.210
	28	0.727	1.451	0.196
	29	0.822	1.802	0.205
	30	0.914	1.311	0.210
ค่าความเที่ยงทั้งฉบับ (Reliability)		0.791		

จากตารางที่ 4-9 ข้อสอบด้านเหตุผล จากจำนวนข้อสอบ 30 ข้อ มีข้อสอบที่ผ่านเกณฑ์การวิเคราะห์คุณภาพข้อสอบ จำนวน 30 ข้อ เป็นข้อสอบที่มีค่าอำนาจจำแนกของข้อสอบ (a) ตั้งแต่ 0.503 ถึง 1.068 ค่าความยากของข้อสอบ (b) ตั้งแต่ -0.442 ถึง 2.968 และค่าโอกาสการเดาของข้อสอบ (c) ตั้งแต่ 0.196 ถึง 0.315 และมีค่าความเที่ยงทั้งฉบับ 0.791

จากการวิเคราะห์คุณภาพของแบบทดสอบระดับชาติ โดยการวิเคราะห์ค่าพารามิเตอร์ ตามทฤษฎีการตอบสนองข้อสอบ พบร่วม

แบบทดสอบระดับชาติ ด้านภาษา ขนาดกลุ่มตัวอย่าง 300 คน มีค่าอำนาจจำแนกของข้อสอบ (a) มีค่าเฉลี่ยเท่ากับ 0.614 ค่าความยากของข้อสอบ (b) มีค่าเฉลี่ยเท่ากับ 0.760 และ

ขนาดกลุ่มตัวอย่าง 2,000 คน มีอำนาจจำแนกข้อสอบอยู่ในระดับดีมาก มีค่าความยากของข้อสอบอยู่ในระดับข้อสอบยาก และค่าโอกาสการเดาของข้อสอบ ไม่เกิน 0.3

แบบทดสอบระดับชาติ ด้านเหตุผล ขนาดกลุ่มตัวอย่าง 2,000 คน มีค่าอำนาจจำแนกของข้อสอบ (a) มีค่าเฉลี่ยเท่ากับ 0.742 ค่าความยากของข้อสอบ (b) มีค่าเฉลี่ยเท่ากับ 0.851 และค่าโอกาสการเดาของข้อสอบ (c) มีค่าเฉลี่ยเท่ากับ 0.239 สรุปได้ว่าแบบทดสอบด้านเหตุผลขนาดกลุ่มตัวอย่าง 2,000 คน มีอำนาจจำแนกข้อสอบอยู่ในระดับดีมาก มีค่าความยากของข้อสอบอยู่ในระดับข้อสอบค่อนข้างยาก และค่าโอกาสการเดาของข้อสอบ ไม่เกิน 0.3

ตารางที่ 4-10 สรุปผลการวิเคราะห์คุณภาพของแบบทดสอบรายข้อ และทั้งฉบับ ชั้นประถมศึกษา ปีที่ 3 ด้านภาษา ด้านคำนวน และด้านเหตุผล

แบบทดสอบ	กลุ่มตัวอย่าง (คน)	ค่าอำนาจจำแนก ของข้อสอบ (a)	ค่าความยาก ของข้อสอบ (b)	ค่าโอกาสการเดา ของข้อสอบ (c)	ค่าความเที่ยง ทั้งฉบับ
ด้านภาษา	300	0.314 - 0.931	-1.120 - 2.829	0.214 - 0.364	0.712
	1,000	0.323 - 1.010	-1.207 - 3.475	0.188 - 0.328	0.738
	2,000	0.318 - 1.089	-1.140 - 2.113	0.179 - 0.232	0.729
ด้านคำนวน	300	0.527 - 1.073	0.271 - 3.503	0.079 - 0.339	0.763
	1,000	0.524 - 1.188	0.016 - 3.688	0.077 - 0.307	0.791
	2,000	0.480 - 1.208	0.061 - 4.000	0.080 - 0.276	0.775
ด้านเหตุผล	300	0.455 - 1.289	-0.432 - 2.820	0.220 - 0.379	0.766
	1,000	0.435 - 1.115	-0.533 - 3.056	0.202 - 0.336	0.784
	2,000	0.503 - 1.068	-0.442 - 2.968	0.196 - 0.315	0.791

จากตารางที่ 4-10 ปรากฏว่า แบบทดสอบด้านภาษา เมื่อกลุ่มตัวอย่างขนาดเล็ก (300 คน) มีค่าความเที่ยงทั้งฉบับเท่ากับ 0.712 กลุ่มตัวอย่างขนาดกลาง (1,000 คน) มีค่าความเที่ยงทั้งฉบับเท่ากับ 0.738 และกลุ่มตัวอย่างขนาดใหญ่ (2,000 คน) มีค่าความเที่ยงทั้งฉบับเท่ากับ 0.729 แบบทดสอบด้านคำนวน เมื่อกลุ่มตัวอย่างขนาดเล็ก (300 คน) มีค่าความเที่ยงทั้งฉบับเท่ากับ 0.763 กลุ่มตัวอย่างขนาดกลาง (1,000 คน) มีค่าความเที่ยงทั้งฉบับเท่ากับ 0.791 และกลุ่มตัวอย่างขนาดใหญ่ (2,000 คน) มีค่าความเที่ยงทั้งฉบับเท่ากับ 0.775 ส่วนแบบทดสอบด้านเหตุผล เมื่อกลุ่มตัวอย่างขนาดเล็ก (300 คน) มีค่าความเที่ยงทั้งฉบับเท่ากับ 0.766 กลุ่มตัวอย่างขนาดกลาง (1,000 คน) มีค่าความเที่ยงทั้งฉบับเท่ากับ 0.784 และกลุ่มตัวอย่างขนาดใหญ่ (2,000 คน) มีค่าความเที่ยงทั้งฉบับเท่ากับ 0.791

ตอนที่ 2 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้วยวิธีการทดสอบอัตราส่วนไลค์ลิชญด วิธีซิบเทส์ และวิธีเมน เหล-ແ xen ส์เซล ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน)

การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษา ปีที่ 3 แบบทดสอบแบ่งออกเป็น 3 ด้าน คือ ด้านภาษา 30 ข้อ ด้านคำนวณ 30 ข้อ ด้านเหตุผล 30 ข้อ ด้วยวิธี IRT-LR วิธี SIBTEST และวิธี MH ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน)

ตารางที่ 4-11 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ระดับชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวณ ด้านเหตุผล ด้วยวิธี IRT-LR ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างขนาดเล็ก (300 คน)

ข้อที่	ผลการตรวจสอบ DIF ด้วยวิธี IRT-LR		
	กลุ่มตัวอย่าง (300 คน)		
	ด้านภาษา	ด้านคำนวณ	ด้านเหตุผล
1	NO-DIF	NO-DIF	NO-DIF
2	<u>DIF</u>	NO-DIF	NO-DIF
3	NO-DIF	NO-DIF	<u>DIF</u>
4	NO-DIF	NO-DIF	NO-DIF
5	NO-DIF	NO-DIF	NO-DIF
6	NO-DIF	NO-DIF	NO-DIF
7	NO-DIF	NO-DIF	NO-DIF
8	NO-DIF	NO-DIF	NO-DIF
9	NO-DIF	NO-DIF	NO-DIF
10	NO-DIF	NO-DIF	NO-DIF
11	NO-DIF	<u>DIF</u>	NO-DIF
12	NO-DIF	<u>DIF</u>	<u>DIF</u>
13	NO-DIF	NO-DIF	NO-DIF
14	NO-DIF	NO-DIF	<u>DIF</u>
15	NO-DIF	NO-DIF	NO-DIF
16	NO-DIF	NO-DIF	NO-DIF
17	<u>DIF</u>	NO-DIF	NO-DIF
18	NO-DIF	NO-DIF	NO-DIF

ตารางที่ 4-11 (ต่อ)

ข้อที่	ผลการตรวจสอบ DIF ด้วยวิธี IRT-LR		
	กลุ่มตัวอย่าง (300 คน)		
	ด้านภาษา	ด้านคำนวน	ด้านเหตุผล
19	NO-DIF	NO-DIF	NO-DIF
20	NO-DIF	NO-DIF	NO-DIF
21	NO-DIF	DIF	NO-DIF
22	NO-DIF	NO-DIF	NO-DIF
23	NO-DIF	NO-DIF	NO-DIF
24	NO-DIF	NO-DIF	NO-DIF
25	NO-DIF	NO-DIF	NO-DIF
26	NO-DIF	NO-DIF	NO-DIF
27	NO-DIF	NO-DIF	NO-DIF
28	NO-DIF	NO-DIF	NO-DIF
29	NO-DIF	NO-DIF	NO-DIF
30	DIF	NO-DIF	NO-DIF
จำนวนข้อที่พบ DIF พบว่า	(3 ข้อ) 10%	(3 ข้อ) 10%	(3 ข้อ) 10%
หมายเหตุ	DIF	หมายถึง ข้อสอบที่ทำหน้าที่ต่างกัน	
	NO-DIF	หมายถึง ข้อสอบที่ไม่พบว่าทำหน้าที่ต่างกัน	

จากตารางที่ 4-11 แสดงผลการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี IRT-LR ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง จำนวน 300 คน พบว่า

ด้านภาษา ตรวจพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ จำนวน 3 ข้อ ได้แก่ ข้อที่ 2, 17, 30 คิดเป็นร้อยละ 10 ของข้อสอบทั้งหมด

ด้านคำนวน ตรวจพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ จำนวน 3 ข้อ ได้แก่ ข้อที่ 11, 12, 21 คิดเป็นร้อยละ 10 ของข้อสอบทั้งหมด

ด้านเหตุผล ตรวจพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ จำนวน 3 ข้อ ได้แก่ ข้อที่ 3, 12, 14 คิดเป็นร้อยละ 10 ของข้อสอบทั้งหมด

ตารางที่ 4-12 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ
ระดับชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี IRT-LR
ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างขนาดกลาง (1,000 คน)

ข้อที่	ผลการตรวจสอบ DIF ด้วยวิธี IRT-LR		
	กลุ่มตัวอย่าง (1,000 คน)		
	ด้านภาษา	ด้านคำนวน	ด้านเหตุผล
1	NO-DIF	NO-DIF	<u>DIF</u>
2	NO-DIF	NO-DIF	NO-DIF
3	NO-DIF	NO-DIF	<u>DIF</u>
4	NO-DIF	NO-DIF	NO-DIF
5	<u>DIF</u>	NO-DIF	NO-DIF
6	NO-DIF	NO-DIF	NO-DIF
7	NO-DIF	NO-DIF	NO-DIF
8	NO-DIF	NO-DIF	NO-DIF
9	NO-DIF	NO-DIF	NO-DIF
10	NO-DIF	<u>DIF</u>	NO-DIF
11	NO-DIF	NO-DIF	<u>DIF</u>
12	NO-DIF	<u>DIF</u>	<u>DIF</u>
13	NO-DIF	NO-DIF	NO-DIF
14	NO-DIF	NO-DIF	NO-DIF
15	<u>DIF</u>	NO-DIF	NO-DIF
16	NO-DIF	<u>DIF</u>	NO-DIF
17	NO-DIF	NO-DIF	NO-DIF
18	<u>DIF</u>	NO-DIF	NO-DIF
19	NO-DIF	<u>DIF</u>	NO-DIF
20	NO-DIF	NO-DIF	NO-DIF
21	NO-DIF	NO-DIF	NO-DIF
22	NO-DIF	NO-DIF	NO-DIF
23	NO-DIF	NO-DIF	NO-DIF
24	NO-DIF	NO-DIF	<u>DIF</u>
25	NO-DIF	NO-DIF	NO-DIF
26	NO-DIF	NO-DIF	NO-DIF
27	NO-DIF	NO-DIF	NO-DIF
28	NO-DIF	NO-DIF	NO-DIF

ตารางที่ 4-12 (ต่อ)

ข้อที่	ผลการตรวจสอบ DIF ด้วยวิธี IRT-LR		
	กลุ่มตัวอย่าง (1,000 คน)		
	ด้านภาษา	ด้านคำนวน	ด้านเหตุผล
29	NO-DIF	NO-DIF	NO-DIF
30	DIF	NO-DIF	NO-DIF
จำนวนข้อที่พบ DIF พบว่า	(4 ข้อ) 13%	(4 ข้อ) 13%	(5 ข้อ) 16%
หมายเหตุ	DIF หมายถึง ข้อสอบที่ทำหน้าที่ต่างกัน		
	NO-DIF หมายถึง ข้อสอบที่ไม่พบว่าทำหน้าที่ต่างกัน		

จากตารางที่ 4-12 แสดงผลการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้านทางภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี IRT-LR ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง จำนวน 1,000 คน พบว่า

ด้านภาษา ตรวจพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ จำนวน 4 ข้อ ได้แก่ ข้อที่ 5, 15, 18, 30 คิดเป็นร้อยละ 13 ของข้อสอบทั้งหมด

ด้านคำนวน ตรวจพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ จำนวน 4 ข้อ ได้แก่ ข้อที่ 10, 12, 16, 19 คิดเป็นร้อยละ 13 ของข้อสอบทั้งหมด

ด้านเหตุผล ตรวจพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ จำนวน 5 ข้อ ได้แก่ ข้อที่ 1, 3, 11, 12, 24 คิดเป็นร้อยละ 16 ของข้อสอบทั้งหมด

ตารางที่ 4-13 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติระดับชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี IRT-LR ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างขนาดใหญ่ (2,000 คน)

ข้อที่	ผลการตรวจสอบ DIF ด้วยวิธี IRT-LR		
	กลุ่มตัวอย่าง (2,000 คน)		
	ด้านภาษา	ด้านคำนวน	ด้านเหตุผล
1	NO-DIF	NO-DIF	DIF
2	NO-DIF	NO-DIF	NO-DIF
3	NO-DIF	NO-DIF	DIF
4	NO-DIF	DIF	NO-DIF
5	DIF	NO-DIF	NO-DIF
6	NO-DIF	NO-DIF	NO-DIF
7	NO-DIF	NO-DIF	NO-DIF
8	NO-DIF	NO-DIF	NO-DIF

ตารางที่ 4-13 (ต่อ)

ข้อที่	ผลการตรวจสอบ DIF ด้วยวิธี IRT-LR		
	กลุ่มตัวอย่าง (2,000 คน)		
	ด้านภาษา	ด้านคำนวน	ด้านเหตุผล
9	NO-DIF	NO-DIF	DIF
10	NO-DIF	NO-DIF	NO-DIF
11	NO-DIF	NO-DIF	DIF
12	NO-DIF	DIF	DIF
13	NO-DIF	NO-DIF	NO-DIF
14	NO-DIF	NO-DIF	DIF
15	NO-DIF	NO-DIF	NO-DIF
16	NO-DIF	DIF	NO-DIF
17	NO-DIF	NO-DIF	NO-DIF
18	NO-DIF	NO-DIF	NO-DIF
19	NO-DIF	NO-DIF	DIF
20	NO-DIF	NO-DIF	NO-DIF
21	NO-DIF	NO-DIF	DIF
22	NO-DIF	NO-DIF	NO-DIF
23	NO-DIF	NO-DIF	NO-DIF
24	NO-DIF	NO-DIF	NO-DIF
25	NO-DIF	NO-DIF	NO-DIF
26	NO-DIF	NO-DIF	NO-DIF
27	NO-DIF	NO-DIF	NO-DIF
28	NO-DIF	NO-DIF	NO-DIF
29	DIF	NO-DIF	DIF
30	DIF	NO-DIF	DIF
จำนวนข้อที่พบ DIF พบร่วมกัน	(3 ข้อ) 10%	(3 ข้อ) 10%	(10 ข้อ) 33%
หมายเหตุ	DIF	หมายถึง ข้อสอบที่ทำหน้าที่ต่างกัน	
	NO-DIF	หมายถึง ข้อสอบที่ไม่พบร่วมกัน	

จากตารางที่ 4-13 แสดงผลการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี IRT-LR ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง จำนวน 2,000 คน พบร่วมกัน

ด้านภาษา ตรวจพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ จำนวน 3 ข้อ ได้แก่ ข้อที่ 5, 29, 30 คิดเป็นร้อยละ 10% ของข้อสอบทั้งหมด

ด้านคำนวน ตรวจพบทข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ จำนวน 3 ข้อ ได้แก่ ข้อที่ 4, 12, 16 คิดเป็นร้อยละ 10% ของข้อสอบทั้งหมด

ด้านเหตุผล ตรวจพบทข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ จำนวน 10 ข้อ ได้แก่ ข้อที่ 1, 3, 9, 11, 12, 14, 19, 21, 29, 30 คิดเป็นร้อยละ 33% ของข้อสอบทั้งหมด

ตารางที่ 4-14 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ

ระดับชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี SIBTEST
ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างขนาดเล็ก (300 คน)

ข้อที่	ผลการตรวจสอบ DIF ด้วยวิธี SIBTEST		
	กลุ่มตัวอย่าง (300 คน)		
	ด้านภาษา	ด้านคำนวน	ด้านเหตุผล
1	NO-DIF	NO-DIF	NO-DIF
2	NO-DIF	NO-DIF	NO-DIF
3	NO-DIF	NO-DIF	NO-DIF
4	NO-DIF	NO-DIF	NO-DIF
5	NO-DIF	NO-DIF	NO-DIF
6	NO-DIF	NO-DIF	NO-DIF
7	NO-DIF	NO-DIF	NO-DIF
8	NO-DIF	NO-DIF	NO-DIF
9	NO-DIF	NO-DIF	NO-DIF
10	NO-DIF	NO-DIF	NO-DIF
11	NO-DIF	NO-DIF	NO-DIF
12	NO-DIF	DIF	NO-DIF
13	NO-DIF	NO-DIF	NO-DIF
14	NO-DIF	NO-DIF	NO-DIF
15	NO-DIF	NO-DIF	NO-DIF
16	NO-DIF	NO-DIF	NO-DIF
17	NO-DIF	NO-DIF	NO-DIF
18	NO-DIF	NO-DIF	DIF
19	DIF	NO-DIF	NO-DIF
20	NO-DIF	NO-DIF	DIF
21	NO-DIF	DIF	NO-DIF
22	NO-DIF	NO-DIF	NO-DIF
23	NO-DIF	NO-DIF	NO-DIF
24	NO-DIF	NO-DIF	NO-DIF

ตารางที่ 4-14 (ต่อ)

ข้อที่	ผลการตรวจสอบ DIF ด้วยวิธี SIBTEST		
	กลุ่มตัวอย่าง (300 คน)		
	ด้านภาษา	ด้านคำนวน	ด้านเหตุผล
25	NO-DIF	NO-DIF	NO-DIF
26	NO-DIF	<u>DIF</u>	NO-DIF
27	NO-DIF	NO-DIF	NO-DIF
28	NO-DIF	NO-DIF	NO-DIF
29	NO-DIF	NO-DIF	NO-DIF
30	<u>DIF</u>	NO-DIF	NO-DIF
จำนวนข้อที่พบ DIF พบร่วม	(2 ข้อ) 7%	(3 ข้อ) 10%	(2 ข้อ) 7%
หมายเหตุ	DIF	หมายถึง ข้อสอบที่ทำหน้าที่ต่างกัน	
	NO-DIF	หมายถึง ข้อสอบที่ไม่พบว่าทำหน้าที่ต่างกัน	

จากตารางที่ 4-14 แสดงผลการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี SIBTEST ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง จำนวน 300 คน พบร่วม

ด้านภาษา ตรวจพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ จำนวน 2 ข้อ ได้แก่ ข้อที่ 19, 30 คิดเป็นร้อยละ 7% ของข้อสอบทั้งหมด

ด้านคำนวน ตรวจพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ จำนวน 3 ข้อ ได้แก่ ข้อที่ 12, 21, 26 คิดเป็นร้อยละ 10% ของข้อสอบทั้งหมด

ด้านเหตุผล ตรวจพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ จำนวน 2 ข้อ ได้แก่ ข้อที่ 18, 20 คิดเป็นร้อยละ 7% ของข้อสอบทั้งหมด

ตารางที่ 4-15 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ระดับชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี SIBTEST ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างขนาดกลาง (1,000 คน)

ข้อที่	ผลการตรวจสอบ DIF ด้วยวิธี SIBTEST		
	กลุ่มตัวอย่าง (1,000 คน)		
	ด้านภาษา	ด้านคำนวน	ด้านเหตุผล
1	NO-DIF	NO-DIF	<u>DIF</u>
2	NO-DIF	NO-DIF	NO-DIF
3	NO-DIF	NO-DIF	NO-DIF
4	NO-DIF	NO-DIF	NO-DIF

ตารางที่ 4-15 (ต่อ)

ข้อที่	ผลการตรวจสอบ DIF ด้วยวิธี SIBTEST		
	กลุ่มตัวอย่าง (1,000 คน)		
	ด้านภาษา	ด้านคำนวน	ด้านเหตุผล
5	DIF	NO-DIF	NO-DIF
6	NO-DIF	NO-DIF	NO-DIF
7	NO-DIF	NO-DIF	NO-DIF
8	NO-DIF	NO-DIF	NO-DIF
9	NO-DIF	NO-DIF	DIF
10	DIF	NO-DIF	NO-DIF
11	NO-DIF	NO-DIF	DIF
12	NO-DIF	NO-DIF	DIF
13	NO-DIF	NO-DIF	NO-DIF
14	NO-DIF	NO-DIF	NO-DIF
15	NO-DIF	NO-DIF	NO-DIF
16	NO-DIF	NO-DIF	NO-DIF
17	NO-DIF	NO-DIF	NO-DIF
18	NO-DIF	NO-DIF	NO-DIF
19	NO-DIF	DIF	DIF
20	NO-DIF	NO-DIF	DIF
21	NO-DIF	NO-DIF	NO-DIF
22	NO-DIF	NO-DIF	NO-DIF
23	NO-DIF	NO-DIF	NO-DIF
24	NO-DIF	NO-DIF	NO-DIF
25	NO-DIF	NO-DIF	NO-DIF
26	NO-DIF	NO-DIF	NO-DIF
27	NO-DIF	NO-DIF	NO-DIF
28	NO-DIF	NO-DIF	NO-DIF
29	NO-DIF	NO-DIF	NO-DIF
30	NO-DIF	NO-DIF	NO-DIF
จำนวนข้อที่พบ DIF พบว่า		(2 ข้อ) 7%	(1 ข้อ) 3%
หมายเหตุ DIF หมายถึง ข้อสอบที่ทำหน้าที่ต่างกัน		(6 ข้อ) 20%	
NO-DIF หมายถึง ข้อสอบที่ไม่พบว่าทำหน้าที่ต่างกัน			

จากตารางที่ 4-15 แสดงผลการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี SIBTEST ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง จำนวน 1,000 คน พบว่า

ด้านภาษา ตรวจพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ จำนวน 2 ข้อ ได้แก่ ข้อที่ 5, 10 คิดเป็นร้อยละ 7% ของข้อสอบทั้งหมด

ด้านคำนวน ตรวจพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ จำนวน 1 ข้อ ได้แก่ ข้อที่ 19 คิดเป็นร้อยละ 3% ของข้อสอบทั้งหมด

ด้านเหตุผล ตรวจพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ จำนวน 6 ข้อ ได้แก่ ข้อที่ 1, 9, 11, 12, 19, 20 คิดเป็นร้อยละ 20% ของข้อสอบทั้งหมด

ตารางที่ 4-16 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ระดับชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี SIBTEST ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง ขนาดใหญ่ (2,000 คน)

ข้อที่	ผลการตรวจสอบ DIF ด้วยวิธี SIBTEST		
	กลุ่มตัวอย่าง (2,000 คน)		
	ด้านภาษา	ด้านคำนวน	ด้านเหตุผล
1	NO-DIF	NO-DIF	<u>DIF</u>
2	NO-DIF	NO-DIF	NO-DIF
3	NO-DIF	NO-DIF	NO-DIF
4	NO-DIF	<u>DIF</u>	NO-DIF
5	<u>DIF</u>	NO-DIF	NO-DIF
6	NO-DIF	NO-DIF	NO-DIF
7	NO-DIF	NO-DIF	NO-DIF
8	NO-DIF	NO-DIF	NO-DIF
9	NO-DIF	<u>DIF</u>	<u>DIF</u>
10	NO-DIF	NO-DIF	NO-DIF
11	NO-DIF	NO-DIF	<u>DIF</u>
12	NO-DIF	NO-DIF	<u>DIF</u>
13	NO-DIF	<u>DIF</u>	NO-DIF
14	NO-DIF	NO-DIF	<u>DIF</u>
15	NO-DIF	NO-DIF	NO-DIF
16	NO-DIF	<u>DIF</u>	NO-DIF
17	NO-DIF	NO-DIF	NO-DIF
18	NO-DIF	NO-DIF	NO-DIF
19	NO-DIF	<u>DIF</u>	<u>DIF</u>

ตารางที่ 4-16 (ต่อ)

ข้อที่	ผลการตรวจสอบ DIF ด้วยวิธี SIBTEST		
	กลุ่มตัวอย่าง (2,000 คน)		
	ด้านภาษา	ด้านคำนวน	ด้านเหตุผล
20	NO-DIF	NO-DIF	<u>DIF</u>
21	NO-DIF	NO-DIF	<u>DIF</u>
22	NO-DIF	NO-DIF	NO-DIF
23	NO-DIF	NO-DIF	NO-DIF
24	<u>DIF</u>	NO-DIF	NO-DIF
25	NO-DIF	NO-DIF	<u>DIF</u>
26	NO-DIF	NO-DIF	<u>DIF</u>
27	NO-DIF	NO-DIF	NO-DIF
28	NO-DIF	<u>DIF</u>	NO-DIF
29	<u>DIF</u>	NO-DIF	<u>DIF</u>
30	<u>DIF</u>	NO-DIF	<u>DIF</u>
จำนวนข้อที่พบ DIF พบร่วม	(4 ข้อ) 13%	(6 ข้อ) 20%	(12 ข้อ) 40%
หมายเหตุ	DIF	หมายถึง ข้อสอบที่ทำหน้าที่ต่างกัน	
	NO-DIF	หมายถึง ข้อสอบที่ไม่พบว่าทำหน้าที่ต่างกัน	

จากตารางที่ 4-16 แสดงผลการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี SIBTEST ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง จำนวน 2,000 คน พบร่วม

ด้านภาษา ตรวจพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ จำนวน 4 ข้อ ได้แก่ ข้อที่ 5, 24, 29, 30 คิดเป็นร้อยละ 13% ของข้อสอบทั้งหมด

ด้านคำนวน ตรวจพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ จำนวน 6 ข้อ ได้แก่ ข้อที่ 4, 9, 13, 16, 19, 28 คิดเป็นร้อยละ 20% ของข้อสอบทั้งหมด

ด้านเหตุผล ตรวจพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ จำนวน 12 ข้อ ได้แก่ ข้อที่ 1, 9, 11, 12, 14, 19, 20, 21, 25, 26, 29, 30 คิดเป็นร้อยละ 40% ของข้อสอบทั้งหมด

ตารางที่ 4-17 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ระดับชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวณ ด้านเหตุผล ด้วยวิธี MH ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างขนาดเล็ก (300 คน)

ข้อที่	ผลการตรวจสอบ DIF ด้วยวิธี MH		
	กลุ่มตัวอย่าง (300 คน)		
	ด้านภาษา	ด้านคำนวณ	ด้านเหตุผล
1	NO-DIF	NO-DIF	NO-DIF
2	<u>DIF</u>	NO-DIF	NO-DIF
3	<u>DIF</u>	NO-DIF	NO-DIF
4	<u>DIF</u>	NO-DIF	<u>DIF</u>
5	NO-DIF	NO-DIF	NO-DIF
6	NO-DIF	NO-DIF	NO-DIF
7	NO-DIF	NO-DIF	NO-DIF
8	NO-DIF	NO-DIF	NO-DIF
9	NO-DIF	NO-DIF	NO-DIF
10	NO-DIF	NO-DIF	NO-DIF
11	NO-DIF	NO-DIF	<u>DIF</u>
12	NO-DIF	<u>DIF</u>	<u>DIF</u>
13	NO-DIF	NO-DIF	<u>DIF</u>
14	NO-DIF	NO-DIF	NO-DIF
15	NO-DIF	NO-DIF	NO-DIF
16	NO-DIF	NO-DIF	NO-DIF
17	<u>DIF</u>	NO-DIF	NO-DIF
18	NO-DIF	NO-DIF	NO-DIF
19	<u>DIF</u>	NO-DIF	NO-DIF
20	NO-DIF	NO-DIF	NO-DIF
21	NO-DIF	<u>DIF</u>	NO-DIF
22	NO-DIF	NO-DIF	NO-DIF
23	NO-DIF	NO-DIF	NO-DIF
24	NO-DIF	NO-DIF	NO-DIF
25	NO-DIF	NO-DIF	NO-DIF
26	NO-DIF	NO-DIF	<u>DIF</u>
27	NO-DIF	NO-DIF	NO-DIF

ตารางที่ 4-17 (ต่อ)

ข้อที่	ผลการตรวจสอบ DIF ด้วยวิธี MH		
	กลุ่มตัวอย่าง (300 คน)		
	ด้านภาษา	ด้านคำนวน	ด้านเหตุผล
28	NO-DIF	NO-DIF	NO-DIF
29	NO-DIF	NO-DIF	NO-DIF
30	NO-DIF	NO-DIF	NO-DIF
จำนวนข้อที่พบ DIF พบว่า	(5 ข้อ) 17%	(2 ข้อ) 7%	(5 ข้อ) 17%
หมายเหตุ	DIF	หมายถึง ข้อสอบที่ทำหน้าที่ต่างกัน	
	NO-DIF	หมายถึง ข้อสอบที่ไม่เพนกว่าทำหน้าที่ต่างกัน	

จากตารางที่ 4-17 แสดงผลการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี MH ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง จำนวน 300 คน พบว่า

ด้านภาษา ตรวจพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ จำนวน 5 ข้อ ได้แก่ ข้อที่ 2, 3, 4, 17, 19 คิดเป็นร้อยละ 17% ของข้อสอบทั้งหมด

ด้านคำนวน ตรวจพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ จำนวน 2 ข้อ ได้แก่ ข้อที่ 12, 21 คิดเป็นร้อยละ 7% ของข้อสอบทั้งหมด

ด้านเหตุผล ตรวจพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ จำนวน 5 ข้อ ได้แก่ ข้อที่ 4, 11, 12, 13, 26 คิดเป็นร้อยละ 17% ของข้อสอบทั้งหมด

ตารางที่ 4-18 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ระดับชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี MH ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างขนาดกลาง (1,000 คน)

ข้อที่	ผลการตรวจสอบ DIF ด้วยวิธี MH		
	กลุ่มตัวอย่าง (1,000 คน)		
	ด้านภาษา	ด้านคำนวน	ด้านเหตุผล
1	NO-DIF	NO-DIF	NO-DIF
2	NO-DIF	NO-DIF	NO-DIF
3	<u>DIF</u>	NO-DIF	NO-DIF
4	<u>DIF</u>	NO-DIF	<u>DIF</u>
5	NO-DIF	NO-DIF	<u>DIF</u>
6	NO-DIF	NO-DIF	<u>DIF</u>
7	<u>DIF</u>	NO-DIF	NO-DIF

ตารางที่ 4-18 (ต่อ)

ข้อที่	ผลการตรวจสอบ DIF ด้วยวิธี MH		
	กลุ่มตัวอย่าง (1,000 คน)		
	ด้านภาษา	ด้านคำนวน	ด้านเหตุผล
8	NO-DIF	NO-DIF	NO-DIF
9	<u>DIF</u>	NO-DIF	NO-DIF
10	NO-DIF	NO-DIF	<u>DIF</u>
11	NO-DIF	NO-DIF	<u>DIF</u>
12	NO-DIF	NO-DIF	<u>DIF</u>
13	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>
14	NO-DIF	NO-DIF	<u>DIF</u>
15	NO-DIF	NO-DIF	NO-DIF
16	<u>DIF</u>	<u>DIF</u>	NO-DIF
17	<u>DIF</u>	<u>DIF</u>	NO-DIF
18	<u>DIF</u>	NO-DIF	<u>DIF</u>
19	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>
20	NO-DIF	NO-DIF	NO-DIF
21	NO-DIF	NO-DIF	NO-DIF
22	NO-DIF	NO-DIF	<u>DIF</u>
23	NO-DIF	NO-DIF	<u>DIF</u>
24	<u>DIF</u>	NO-DIF	<u>DIF</u>
25	NO-DIF	NO-DIF	NO-DIF
26	NO-DIF	NO-DIF	<u>DIF</u>
27	NO-DIF	NO-DIF	NO-DIF
28	NO-DIF	NO-DIF	<u>DIF</u>
29	<u>DIF</u>	NO-DIF	NO-DIF
30	NO-DIF	NO-DIF	NO-DIF
จำนวนข้อที่พบ DIF พบร่วมกัน		(11 ข้อ) 37%	(4 ข้อ) 13% (15 ข้อ) 50%
หมายเหตุ	DIF	หมายถึง ข้อสอบที่ทำหน้าที่ต่างกัน	
	NO-DIF	หมายถึง ข้อสอบที่ไม่พบว่าทำหน้าที่ต่างกัน	

จากตารางที่ 4-18 แสดงผลการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี MH ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง จำนวน 1,000 คน พบว่า

ด้านภาษา ตรวจพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ จำนวน 11 ข้อ ได้แก่ ข้อที่ 3, 4, 7, 9, 13, 16, 17, 18, 19, 24, 29 คิดเป็นร้อยละ 37% ของข้อสอบทั้งหมด

ด้านคำนวน ตรวจพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ จำนวน 4 ข้อ ได้แก่ ข้อที่ 13, 16, 17, 19 คิดเป็นร้อยละ 13% ของข้อสอบทั้งหมด

ด้านเหตุผล ตรวจพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ จำนวน 15 ข้อ ได้แก่ ข้อที่ 4, 5, 6, 10, 11, 12, 13, 14, 18, 19, 22, 23, 24, 26, 28 คิดเป็นร้อยละ 50% ของข้อสอบทั้งหมด

ตารางที่ 4-19 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติระดับ

ชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี MH

ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างขนาดใหญ่ (2,000 คน)

ข้อที่	ผลการตรวจสอบ DIF ด้วยวิธี MH		
	กลุ่มตัวอย่าง (2,000 คน)		
	ด้านภาษา	ด้านคำนวน	ด้านเหตุผล
1	NO-DIF	NO-DIF	NO-DIF
2	NO-DIF	NO-DIF	<u>DIF</u>
3	<u>DIF</u>	NO-DIF	<u>DIF</u>
4	<u>DIF</u>	NO-DIF	<u>DIF</u>
5	<u>DIF</u>	NO-DIF	<u>DIF</u>
6	NO-DIF	NO-DIF	<u>DIF</u>
7	<u>DIF</u>	NO-DIF	<u>DIF</u>
8	NO-DIF	<u>DIF</u>	<u>DIF</u>
9	<u>DIF</u>	<u>DIF</u>	NO-DIF
10	NO-DIF	NO-DIF	<u>DIF</u>
11	<u>DIF</u>	NO-DIF	<u>DIF</u>
12	<u>DIF</u>	NO-DIF	<u>DIF</u>
13	NO-DIF	<u>DIF</u>	<u>DIF</u>
14	<u>DIF</u>	NO-DIF	<u>DIF</u>
15	<u>DIF</u>	NO-DIF	<u>DIF</u>
16	<u>DIF</u>	<u>DIF</u>	NO-DIF
17	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>
18	<u>DIF</u>	NO-DIF	<u>DIF</u>
19	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>

ตารางที่ 4-19 (ต่อ)

ข้อที่	ผลการตรวจสอบ DIF ด้วยวิธี MH		
	กลุ่มตัวอย่าง (2,000 คน)		
	ด้านภาษา	ด้านคำนวน	ด้านเหตุผล
20	<u>DIF</u>	<u>DIF</u>	NO-DIF
21	<u>DIF</u>	NO-DIF	NO-DIF
22	NO-DIF	NO-DIF	<u>DIF</u>
23	NO-DIF	NO-DIF	<u>DIF</u>
24	<u>DIF</u>	NO-DIF	<u>DIF</u>
25	NO-DIF	NO-DIF	NO-DIF
26	NO-DIF	NO-DIF	<u>DIF</u>
27	<u>DIF</u>	NO-DIF	<u>DIF</u>
28	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>
29	<u>DIF</u>	NO-DIF	NO-DIF
30	NO-DIF	NO-DIF	NO-DIF
จำนวนข้อที่พบร่วมกัน	(19 ข้อ) 63%	(8 ข้อ) 27%	(22 ข้อ) 73%
หมายเหตุ	DIF	หมายถึง ข้อสอบที่ทำหน้าที่ต่างกัน	
	NO-DIF	หมายถึง ข้อสอบที่ไม่พบร่วมกัน	

จากตารางที่ 4-19 แสดงผลการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี MH ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง จำนวน 2,000 คน พบร่วมกัน

ด้านภาษา ตรวจพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ จำนวน 19 ข้อ ได้แก่ ข้อที่ 3, 4, 5, 7, 9, 11, 12, 14, 15, 16, 17, 18, 19, 21, 24, 27, 28, 29 คิดเป็นร้อยละ 63 ของข้อสอบทั้งหมด

ด้านคำนวน ตรวจพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ จำนวน 8 ข้อ ได้แก่ ข้อที่ 8, 9, 13, 16, 17, 19, 20, 28, คิดเป็นร้อยละ 27 ของข้อสอบทั้งหมด

ด้านเหตุผล ตรวจพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ จำนวน 22 ข้อ ได้แก่ ข้อที่ 2, 3, 4, 5, 6, 7, 8, 10, 11, 12, 13, 14, 15, 17, 18, 19, 22, 23, 24, 26, 27, 28 คิดเป็นร้อยละ 73 ของข้อสอบทั้งหมด

ตารางที่ 4-20 ผลการวิเคราะห์คุณภาพของแบบทดสอบระดับชาติชั้นประถมศึกษาปีที่ 3 ด้านภาษา
ด้านคำนวน และด้านเหตุผล

แบบ ทดสอบ	กลุ่ม ตัวอย่าง (คน)	จำนวน ข้อสอบ (ข้อ)	วิธี IRT-LR		วิธี SIBTEST		วิธี MHI	
			จำนวน ข้อที่พบร DIF	ร้อยละ	จำนวน ข้อที่พบร DIF	ร้อยละ	จำนวน ข้อที่พบร DIF	ร้อยละ
ด้าน ภาษา	300	30	3	10	2	7	5	17
	1,000	30	4	13	2	7	11	37
	2,000	30	3	10	4	13	19	63
ด้าน คำนวน	300	30	3	10	3	10	2	7
	1,000	30	4	13	1	3	4	13
	2,000	30	3	10	6	20	8	27
ด้าน เหตุผล	300	30	3	10	2	7	5	17
	1,000	30	5	16	6	20	15	50
	2,000	30	10	33	12	40	22	73
รวมทั้ง 3 ด้าน		270	38	14	38	14	91	34

จากตารางที่ 4-20 แสดงผลการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบ
ระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน และด้านเหตุผล ปรากฏว่า วิธี Mantel-
Haenszel ตรวจพบ DIF จำนวน 91 ข้อ คิดเป็นร้อยละ 34 ของข้อสอบทั้งฉบับ รองลงมาคือ
วิธี SIBTEST และวิธี IRT-LR ตรวจพบ DIF จำนวน 38 ข้อ เท่ากัน คิดเป็นร้อยละ 14

ตอนที่ 3 ผลการเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้วยวิธีการทดสอบอัตราส่วนໄลค์ลิลลูด วิธีซิปเทสท์ และวิธีแมนเทล-ແ xen ส์เซล ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน)

การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 คือ ด้านภาษา จำนวน 30 ข้อ ด้านคำวณ จำนวน 30 ข้อ และด้านเหตุผล จำนวน 30 ข้อ โดยเปรียบเทียบระหว่างวิธี 2 วิธี วิธีใดตรวจพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบจำนวนมากกว่ากัน ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน)

ตารางที่ 4-21 การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ด้านภาษา ระหว่างวิธี IRT-LR กับวิธี SIBTEST

ข้อที่	ผลการเปรียบเทียบการทำหน้าที่ต่างกันของข้อสอบ					
	ด้านภาษา					
	วิธี IRT-LR			วิธี SIBTEST		
	300 คน	1,000 คน	2,000 คน	300 คน	1,000 คน	2,000 คน
1	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
2	<u>DIF</u>	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
3	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
4	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
5	NO-DIF	<u>DIF</u>	<u>DIF</u>	NO-DIF	<u>DIF</u>	<u>DIF</u>
6	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
7	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
8	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
9	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
10	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	NO-DIF
11	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
12	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
13	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
14	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
15	NO-DIF	<u>DIF</u>	NO-DIF	NO-DIF	NO-DIF	NO-DIF
16	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
17	<u>DIF</u>	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
18	NO-DIF	<u>DIF</u>	NO-DIF	NO-DIF	NO-DIF	NO-DIF

ตารางที่ 4-21 (ต่อ)

ข้อที่	ผลการเปรียบเทียบการทำหน้าที่ต่างกันของข้อสอบ					
	ด้านภาษา					
	วิธี IRT-LR			วิธี SIBTEST		
	300 คน	1,000 คน	2,000 คน	300 คน	1,000 คน	2,000 คน
19	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	NO-DIF	NO-DIF
20	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
21	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
22	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
23	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
24	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>
25	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
26	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
27	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
28	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
29	NO-DIF	NO-DIF	<u>DIF</u>	NO-DIF	NO-DIF	<u>DIF</u>
30	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>	NO-DIF	<u>DIF</u>
จำนวนข้อที่พบDIF	(3 ข้อ)	(4 ข้อ)	(3 ข้อ)	(2 ข้อ)	(2 ข้อ)	(4 ข้อ)
	10%	13%	10%	7%	7%	13%

หมายเหตุ DIF หมายถึง ข้อสอบที่ทำหน้าที่ต่างกัน
NO-DIF หมายถึง ข้อสอบที่ไม่พบว่าทำหน้าที่ต่างกัน

จากตารางที่ 4-21 ผลการเปรียบเทียบการทำหน้าที่ต่างกันของข้อสอบ ด้านภาษา ระหว่างวิธี IRT-LR กับวิธี SIBTEST พบว่า เมื่อกรุ่มตัวอย่างขนาดเล็ก (300 คน) วิธี IRT-LR พบร DIF มากกว่า วิธี SIBTEST จำนวน 1 ข้อ คิดเป็นร้อยละ 3 ขนาดกลาง (1,000 คน) วิธี IRT-LR พบร DIF มากกว่า วิธี SIBTEST จำนวน 2 ข้อ คิดเป็นร้อยละ 7 ขนาดใหญ่ (2,000 คน) วิธี IRT-LR พบร DIF น้อยกว่า วิธี SIBTEST จำนวน 1 ข้อ คิดเป็นร้อยละ 3 ของข้อสอบทั้งหมด

ตารางที่ 4-22 การเปรียบเทียบผลการตรวจสอบการทำงานที่ต่างกันของข้อสอบในแบบทดสอบ ระดับชาติ ด้านค่านวน ระหว่างวิชี IRT-LR กับ วิชี SIBTEST

ตารางที่ 4-22 (ต่อ)

ข้อที่	ผลการเปรียบเทียบการทำหน้าที่ต่างกันของข้อสอบ						
	ด้านคำนวณ						วิธี SIBTEST
	วิธี IRT-LR			วิธี SIBTEST			
	300 คน	1,000 คน	2,000 คน	300 คน	1,000 คน	2,000 คน	
28	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>
29	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
30	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
จำนวนข้อที่พับDIF	(3 ข้อ)	(4 ข้อ)	(3 ข้อ)	(3 ข้อ)	(1 ข้อ)	(6 ข้อ)	
	10%	13%	10%	10%	3%	20%	
หมายเหตุ	DIF	หมายถึง ข้อสอบที่ทำหน้าที่ต่างกัน					
	NO-DIF	หมายถึง ข้อสอบที่ไม่พบว่าทำหน้าที่ต่างกัน					

จากตารางที่ 4-22 ผลการเปรียบเทียบการทำหน้าที่ต่างกันของข้อสอบ ด้านคำนวณ ระหว่างวิธี IRT-LR กับวิธี SIBTEST พบร่วม เมื่อคู่มุ่งตัวอย่างขนาดเล็ก (300 คน) วิธี IRT-LR พบรหัส DIF เท่ากัน กับวิธี SIBTEST จำนวน 0 ข้อ คิดเป็นร้อยละ 0 ขนาดกลาง (1,000 คน) วิธี IRT-LR พบรหัส DIF มากกว่า วิธี SIBTEST จำนวน 3 ข้อ คิดเป็นร้อยละ 10 ขนาดใหญ่ (2,000 คน) วิธี IRT-LR พบรหัส DIF น้อยกว่า วิธี SIBTEST จำนวน 3 ข้อ คิดเป็นร้อยละ 10 ของข้อสอบทั้งหมด

ตารางที่ 4-23 การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบ ระดับชาติ ด้านเหตุผล ระหว่างวิธี IRT-LR กับวิธี SIBTEST

ข้อที่	ผลการเปรียบเทียบการทำหน้าที่ต่างกันของข้อสอบ						
	ด้านเหตุผล						วิธี SIBTEST
	วิธี IRT-LR			วิธี SIBTEST			
	300 คน	1,000 คน	2,000 คน	300 คน	1,000 คน	2,000 คน	
1	NO-DIF	<u>DIF</u>	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>	
2	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
3	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>	NO-DIF	NO-DIF	NO-DIF	NO-DIF
4	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
5	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
6	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
7	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
8	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
9	NO-DIF	NO-DIF	<u>DIF</u>	NO-DIF	<u>DIF</u>	<u>DIF</u>	

ตารางที่ 4-23 (ต่อ)

ข้อที่	ผลการเปรียบเทียบการทำหน้าที่ต่างกันของข้อสอบ					
	ด้านเหตุผล					
	วิธี IRT-LR			วิธี SIBTEST		
	300 คน	1,000 คน	2,000 คน	300 คน	1,000 คน	2,000 คน
10	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
11	NO-DIF	<u>DIF</u>	<u>DIF</u>	NO-DIF	<u>DIF</u>	<u>DIF</u>
12	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>	NO-DIF	<u>DIF</u>	<u>DIF</u>
13	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
14	<u>DIF</u>	NO-DIF	<u>DIF</u>	NO-DIF	NO-DIF	<u>DIF</u>
15	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
16	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
17	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
18	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	NO-DIF	NO-DIF
19	NO-DIF	NO-DIF	<u>DIF</u>	NO-DIF	<u>DIF</u>	<u>DIF</u>
20	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>
21	NO-DIF	NO-DIF	<u>DIF</u>	NO-DIF	NO-DIF	<u>DIF</u>
22	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
23	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
24	NO-DIF	<u>DIF</u>	NO-DIF	NO-DIF	NO-DIF	NO-DIF
25	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>
26	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>
27	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
28	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
29	NO-DIF	NO-DIF	<u>DIF</u>	NO-DIF	NO-DIF	<u>DIF</u>
30	NO-DIF	NO-DIF	<u>DIF</u>	NO-DIF	NO-DIF	<u>DIF</u>
จำนวนข้อที่พบDIF	(3 ข้อ)	(5 ข้อ)	(10 ข้อ)	(2 ข้อ)	(6 ข้อ)	(12 ข้อ)
	10%	16%	33%	7%	20%	40%

หมายเหตุ DIF หมายถึง ข้อสอบที่ทำหน้าที่ต่างกัน

NO-DIF หมายถึง ข้อสอบที่ไม่พบว่าทำหน้าที่ต่างกัน

จากตารางที่ 4-23 ผลการเปรียบเทียบการทำหน้าที่ต่างกันของข้อสอบ ด้านเหตุผล ระหว่างวิธี IRT-LR กับวิธี SIBTEST พบว่า เมื่อกรุ่มตัวอย่างขนาดเล็ก (300 คน) วิธี IRT-LR พบร DIF มากกว่าวิธี SIBTEST จำนวน 1 ข้อ คิดเป็นร้อยละ 3 ขนาดกลาง (1,000 คน) วิธี IRT-LR พบร DIF

น้อยกว่า วิธี SIBTEST จำนวน 1 ข้อ คิดเป็นร้อยละ 3 ขนาดใหญ่ (2,000 คน) วิธี IRT-LR พบ DIF น้อยกว่า วิธี SIBTEST จำนวน 2 ข้อ คิดเป็นร้อยละ 7 ของข้อสอบทั้งหมด

ตารางที่ 4-24 การเปรียบเทียบผลการตรวจสอบการทำงานที่ต่างกันของข้อสอบในแบบทดสอบ ระดับชาติ ด้านภาษา ระหว่างวิธี IRT-LR กับวิธี MH

ตารางที่ 4-24 (ต่อ)

ข้อที่	ผลการเปรียบเทียบการทำหน้าที่ต่างกันของข้อสอบ						
	ด้านภาษา						วิธี MH
	วิธี IRT-LR			300 คน 1,000 คน 2,000 คน			
27	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>
28	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>
29	NO-DIF	NO-DIF	<u>DIF</u>		NO-DIF	<u>DIF</u>	<u>DIF</u>
30	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>		NO-DIF	NO-DIF	NO-DIF
จำนวนข้อที่พับDIF	(3 ข้อ)	(4 ข้อ)	(3 ข้อ)		(5 ข้อ)	(11 ข้อ)	(18 ข้อ)
	10%	13%	10%		17%	37%	60%
หมายเหตุ DIF	หมายถึง ข้อสอบที่ทำหน้าที่ต่างกัน						
NO-DIF	หมายถึง ข้อสอบที่ไม่พบว่าทำหน้าที่ต่างกัน						

จากตารางที่ 4-24 ผลการเปรียบเทียบการทำหน้าที่ต่างกันของข้อสอบ ด้านภาษา ระหว่างวิธี IRT-LR กับวิธี MH พบร่วม เมื่อกลุ่มตัวอย่างขนาดเล็ก (300 คน) วิธี IRT-LR พบรหัส DIF น้อยกว่าวิธี MH จำนวน 2 ข้อ คิดเป็นร้อยละ 7 ขนาดกลาง (1,000 คน) วิธี IRT-LR พบรหัส DIF น้อยกว่าวิธี MH จำนวน 7 ข้อ คิดเป็นร้อยละ 23 ขนาดใหญ่ (2,000 คน) วิธี IRT-LR พบรหัส DIF น้อยกว่าวิธี MH จำนวน 15 ข้อ คิดเป็นร้อยละ 50 ของข้อสอบทั้งหมด

ตารางที่ 4-25 การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบ ระดับชาติ ด้านคำนวน ระหว่างวิธี IRT-LR กับวิธี MH

ข้อที่	ผลการเปรียบเทียบการทำหน้าที่ต่างกันของข้อสอบ						
	ด้านคำนวน						วิธี MH
	วิธี IRT-LR			300 คน 1,000 คน 2,000 คน			
1	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
2	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
3	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
4	NO-DIF	NO-DIF	<u>DIF</u>		NO-DIF	NO-DIF	NO-DIF
5	NO-DIF	NO-DIF	NO-DIF		NO-DIF	NO-DIF	NO-DIF
6	NO-DIF	NO-DIF	NO-DIF		NO-DIF	NO-DIF	NO-DIF
7	NO-DIF	NO-DIF	NO-DIF		NO-DIF	NO-DIF	NO-DIF
8	NO-DIF	NO-DIF	NO-DIF		NO-DIF	NO-DIF	<u>DIF</u>

ตารางที่ 4-25 (ต่อ)

ข้อที่	ผลการเปรียบเทียบการทำหน้าที่ต่างกันของข้อสอบ						
	ด้านคำนวณ						วิธี MH
	วิธี IRT-LR			300 คน 1,000 คน 2,000 คน			
	300 คน	1,000 คน	2,000 คน	300 คน	1,000 คน	2,000 คน	
9	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>
10	NO-DIF	<u>DIF</u>	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
11	<u>DIF</u>	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
12	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>		<u>DIF</u>	NO-DIF	NO-DIF
13	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>
14	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
15	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
16	NO-DIF	<u>DIF</u>	<u>DIF</u>		NO-DIF	<u>DIF</u>	<u>DIF</u>
17	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>
18	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
19	NO-DIF	<u>DIF</u>	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>
20	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>
21	<u>DIF</u>	NO-DIF	NO-DIF		<u>DIF</u>	NO-DIF	NO-DIF
22	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
23	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
24	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
25	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
26	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
27	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
28	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>
29	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
30	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
จำนวนข้อที่พับDIF	(3 ข้อ)	(4 ข้อ)	(3 ข้อ)	(2 ข้อ)	(4 ข้อ)	(8 ข้อ)	
	10%	13%	10%	7%	13%	27%	

หมายเหตุ DIF หมายถึง ข้อสอบที่ทำหน้าที่ต่างกัน

NO-DIF หมายถึง ข้อสอบที่ไม่พับว่าทำหน้าที่ต่างกัน

จากตารางที่ 4-25 ผลการเปรียบเทียบการทำหน้าที่ต่างกันของข้อสอบ ด้านคำนวณ ระหว่างวิธี IRT-LR กับวิธี MH พบร้า เมื่อคุณตัวอย่างขนาดเล็ก (300 คน) วิธี IRT-LR พับ DIF มากกว่าวิธี MH จำนวน 1 ข้อ คิดเป็นร้อยละ 3 ขนาดกลาง (1,000 คน) วิธี IRT-LR พับ DIF

เท่ากัน กับวิธี MH จำนวน 0 ข้อ คิดเป็นร้อยละ 0 ขนาดใหญ่ (2,000 คน) วิธี IRT-LR พบ DIF น้อยกว่า วิธี MH จำนวน 5 ข้อ คิดเป็นร้อยละ 17 ของข้อสอบทั้งหมด

ตารางที่ 4-26 การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบ
ระดับชาติ ด้านเหตุผล ระหว่างวิธี IRT-LR กับวิธี MH

ข้อที่	ผลการเปรียบเทียบการทำหน้าที่ต่างกันของข้อสอบ						
	ด้านเหตุผล						
	วิธี IRT-LR			วิธี MH			
	300 คน	1,000 คน	2,000 คน	300 คน	1,000 คน	2,000 คน	
1	NO-DIF	<u>DIF</u>	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
2	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>
3	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>
4	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>	
5	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>
6	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>
7	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>
8	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>
9	NO-DIF	NO-DIF	<u>DIF</u>	NO-DIF	NO-DIF	NO-DIF	NO-DIF
10	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>	
11	NO-DIF	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>	
12	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>	
13	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>	
14	<u>DIF</u>	NO-DIF	<u>DIF</u>	NO-DIF	<u>DIF</u>	<u>DIF</u>	
15	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>
16	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
17	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>
18	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>	
19	NO-DIF	NO-DIF	<u>DIF</u>	NO-DIF	<u>DIF</u>	<u>DIF</u>	
20	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
21	NO-DIF	NO-DIF	<u>DIF</u>	NO-DIF	NO-DIF	NO-DIF	NO-DIF
22	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>	
23	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>	
24	NO-DIF	<u>DIF</u>	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>	
25	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
26	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>	

ตารางที่ 4-26 (ต่อ)

ข้อที่	ผลการเปรียบเทียบการทำหน้าที่ต่างกันของข้อสอบ						
	ด้านเหตุผล				วิธี MH		
	วิธี IRT-LR			วิธี MH			
	300 คน	1,000 คน	2,000 คน	300 คน	1,000 คน	2,000 คน	
27	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	DIF	DIF
28	NO-DIF	NO-DIF	NO-DIF	NO-DIF	DIF	DIF	
29	NO-DIF	NO-DIF	DIF	NO-DIF	NO-DIF	NO-DIF	
30	NO-DIF	NO-DIF	DIF	NO-DIF	NO-DIF	NO-DIF	
จำนวนข้อที่พนDIF	(3 ข้อ)	(5 ข้อ)	(10 ข้อ)	(5 ข้อ)	(15 ข้อ)	(22 ข้อ)	
	10%	16%	33%	17%	50%	73%	

หมายเหตุ	DIF	หมายถึง ข้อสอบที่ทำหน้าที่ต่างกัน
	NO-DIF	หมายถึง ข้อสอบที่ไม่พบร่วมทำหน้าที่ต่างกัน

จากตารางที่ 4-26 ผลการเปรียบเทียบการทำหน้าที่ต่างกันของข้อสอบ ด้านเหตุผลระหว่างวิธี IRT-LR กับวิธี MH พบร่วม เมื่อกลุ่มตัวอย่างขนาดเล็ก (300 คน) วิธี IRT-LR พบร DIF น้อยกว่าวิธี MH จำนวน 2 ข้อ คิดเป็นร้อยละ 7 ขนาดกลาง (1,000 คน) วิธี IRT-LR พบร DIF น้อยกว่าวิธี MH จำนวน 10 ข้อ คิดเป็นร้อยละ 33 ขนาดใหญ่ (2,000 คน) วิธี IRT-LR พบร DIF น้อยกว่าวิธี MH จำนวน 12 ข้อ คิดเป็นร้อยละ 40 ของข้อสอบทั้งหมด

ตารางที่ 4-27 การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบ ระดับชาติต้านทานภาษา ระหว่างวิธี SIBTEST กับวิธี MH

ตารางที่ 4-27 (ต่อ)

ข้อที่	ผลการเปรียบเทียบการทำหน้าที่ต่างกันของข้อสอบ					
	ด้านภาษา					
	วิธี SIBTEST			วิธี MH		
	300 คน	1,000 คน	2,000 คน	300 คน	1,000 คน	2,000 คน
9	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>
10	NO-DIF	<u>DIF</u>	NO-DIF	NO-DIF	NO-DIF	NO-DIF
11	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>
12	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>
13	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	NO-DIF
14	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>
15	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>
16	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>
17	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>
18	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>
19	<u>DIF</u>	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>
20	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
21	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>
22	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
23	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
24	NO-DIF	NO-DIF	<u>DIF</u>	NO-DIF	<u>DIF</u>	<u>DIF</u>
25	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
26	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
27	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>
28	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>
29	NO-DIF	NO-DIF	<u>DIF</u>	NO-DIF	<u>DIF</u>	<u>DIF</u>
30	<u>DIF</u>	NO-DIF	<u>DIF</u>	NO-DIF	NO-DIF	NO-DIF
จำนวนข้อที่พบDIF	(2 ข้อ)	(2 ข้อ)	(4 ข้อ)	(5 ข้อ)	(11 ข้อ)	(18 ข้อ)
	7%	7%	13%	17%	37%	60%

หมายเหตุ DIF หมายถึง ข้อสอบที่ทำหน้าที่ต่างกัน

NO-DIF หมายถึง ข้อสอบที่ไม่พบว่าทำหน้าที่ต่างกัน

จากตารางที่ 4-27 ผลการเปรียบเทียบการทำหน้าที่ต่างกันของข้อสอบ ด้านภาษา ระหว่างวิธี SIBTEST กับวิธี MH พบร่วม เมื่อกลุ่มตัวอย่างขนาดเล็ก (300 คน) วิธี SIBTEST พบร DIF น้อยกว่าวิธี MH จำนวน 3 ข้อ คิดเป็นร้อยละ 10 ขนาดกลาง (1,000 คน) วิธี SIBTEST พบร DIF

น้อยกว่า วิธี MH จำนวน 9 ข้อ คิดเป็นร้อยละ 30 ขนาดใหญ่ (2,000 คน) วิธี SIBTEST พบ DIF น้อยกว่า วิธี MH จำนวน 14 ข้อ คิดเป็นร้อยละ 47 ของข้อสอบทั้งหมด

ตารางที่ 4-28 การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบ
ระดับชาติ ด้านคำนวณ ระหว่างวิธี SIBTEST กับวิธี MH

ข้อที่	ผลการเปรียบเทียบการทำหน้าที่ต่างกันของข้อสอบ					
	ด้านคำนวน					
	วิธี SIBTEST			วิธี MH		
	300 คน	1,000 คน	2,000 คน	300 คน	1,000 คน	2,000 คน
1	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
2	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
3	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
4	NO-DIF	NO-DIF	<u>DIF</u>	NO-DIF	NO-DIF	NO-DIF
5	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
6	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
7	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
8	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>
9	NO-DIF	NO-DIF	<u>DIF</u>	NO-DIF	NO-DIF	<u>DIF</u>
10	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
11	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
12	<u>DIF</u>	NO-DIF	NO-DIF	<u>DIF</u>	NO-DIF	NO-DIF
13	NO-DIF	NO-DIF	<u>DIF</u>	NO-DIF	<u>DIF</u>	<u>DIF</u>
14	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
15	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
16	NO-DIF	NO-DIF	<u>DIF</u>	NO-DIF	<u>DIF</u>	<u>DIF</u>
17	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>
18	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
19	NO-DIF	<u>DIF</u>	<u>DIF</u>	NO-DIF	<u>DIF</u>	<u>DIF</u>
20	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>
21	<u>DIF</u>	NO-DIF	NO-DIF	<u>DIF</u>	NO-DIF	NO-DIF
22	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
23	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
24	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
25	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
26	<u>DIF</u>	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF

ตารางที่ 4-28 (ต่อ)

ข้อที่	ผลการเปรียบเทียบการทำหน้าที่ต่างกันของข้อสอบ ด้านความน่าจะเป็น					
	วิธี SIBTEST			วิธี MH		
	300 คน	1,000 คน	2,000 คน	300 คน	1,000 คน	2,000 คน
27	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
28	NO-DIF	NO-DIF	DIF	NO-DIF	NO-DIF	DIF
29	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
30	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
จำนวนข้อที่พับDIF	(3 ข้อ)	(1 ข้อ)	(6 ข้อ)	(2 ข้อ)	(4 ข้อ)	(8 ข้อ)
	10%	3%	20%	7%	13%	27%

หมายเหตุ DIF	หมายถึง ข้อสอบที่ทำหน้าที่ต่างกัน
NO-DIF	หมายถึง ข้อสอบที่ไม่พบร่วมทำหน้าที่ต่างกัน

จากตารางที่ 4-28 ผลการเปรียบเทียบการทำหน้าที่ต่างกันของข้อสอบ ด้านคำนวณระหว่างวิธี SIBTEST กับวิธี MH พบร่วม เมื่อกลุ่มตัวอย่างขนาดเล็ก (300 คน) วิธี SIBTEST พบรหัสกว่า วิธี MH จำนวน 1 ข้อ คิดเป็นร้อยละ 3 ขนาดกลาง (1,000 คน) วิธี SIBTEST พบรหัสกว่า วิธี MH จำนวน 3 ข้อ คิดเป็นร้อยละ 10 ขนาดใหญ่ (2,000 คน) วิธี SIBTEST พบรหัสกว่า วิธี MH จำนวน 2 ข้อ คิดเป็นร้อยละ 7 ของข้อสอบทั้งหมด

ตารางที่ 4-29 การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบ
ระดับชาติ ด้านเหตุผล ระหว่างวิธี SIBTEST กับวิธี MH

ข้อที่	ผลการเปรียบเทียบการทำหน้าที่ต่างกันของข้อสอบ						
	ด้านเหตุผล						
	วิธี SIBTEST			วิธี MH			
	300 คน	1,000 คน	2,000 คน	300 คน	1,000 คน	2,000 คน	
1	NO-DIF	<u>DIF</u>	<u>DIF</u>	NO-DIF	NO-DIF	NO-DIF	
2	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	
3	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	
4	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>	
5	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>	
6	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>	
7	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	
8	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	

ตารางที่ 4-29 (ต่อ)

ข้อที่	ผลการเปรียบเทียบการทำหน้าที่ต่างกันของข้อสอบ					
	ด้านเหตุผล					
	วิธี SIBTEST			วิธี MH		
	300 คน	1,000 คน	2,000 คน	300 คน	1,000 คน	2,000 คน
9	NO-DIF	<u>DIF</u>	<u>DIF</u>	NO-DIF	NO-DIF	NO-DIF
10	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>
11	NO-DIF	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>
12	NO-DIF	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>
13	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>
14	NO-DIF	NO-DIF	<u>DIF</u>	NO-DIF	<u>DIF</u>	<u>DIF</u>
15	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>
16	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF
17	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>
18	<u>DIF</u>	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>
19	NO-DIF	<u>DIF</u>	<u>DIF</u>	NO-DIF	<u>DIF</u>	<u>DIF</u>
20	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>	NO-DIF	NO-DIF	NO-DIF
21	NO-DIF	NO-DIF	<u>DIF</u>	NO-DIF	NO-DIF	NO-DIF
22	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>
23	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>
24	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>
25	NO-DIF	NO-DIF	<u>DIF</u>	NO-DIF	NO-DIF	NO-DIF
26	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>	<u>DIF</u>
27	NO-DIF	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>
28	NO-DIF	NO-DIF	NO-DIF	NO-DIF	<u>DIF</u>	<u>DIF</u>
29	NO-DIF	NO-DIF	<u>DIF</u>	NO-DIF	NO-DIF	NO-DIF
30	NO-DIF	NO-DIF	<u>DIF</u>	NO-DIF	NO-DIF	NO-DIF
จำนวนข้อที่พบDIF	(2 ข้อ)	(6 ข้อ)	(12 ข้อ)	(5 ข้อ)	(15 ข้อ)	(22 ข้อ)
	7%	20%	40%	17%	50%	73%

หมายเหตุ DIF หมายถึง ข้อสอบที่ทำหน้าที่ต่างกัน

NO-DIF หมายถึง ข้อสอบที่ไม่พบว่าทำหน้าที่ต่างกัน

จากตารางที่ 4-29 ผลการเปรียบเทียบการทำหน้าที่ต่างกันของข้อสอบ ด้านเหตุผล ระหว่างวิธี SIBTEST กับวิธี MH พบร่วม เมื่อกลุ่มตัวอย่างขนาดเล็ก (300 คน) วิธี SIBTEST พบร DIF น้อยกว่าวิธี MH จำนวน 3 ข้อ คิดเป็นร้อยละ 10 ขนาดกลาง (1,000 คน) วิธี SIBTEST พบร DIF

น้อยกว่า วิธี MH จำนวน 9 ข้อ คิดเป็นร้อยละ 30 ขนาดใหญ่ (2,000 คน) วิธี SIBTEST พบ DIF น้อยกว่า วิธี MH จำนวน 10 ข้อ คิดเป็นร้อยละ 33 ของข้อสอบทั้งหมด

ตารางที่ 4-30 การเปรียบเทียบการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบด้วย วิธีการทดสอบอัตราส่วนไลค์ลิชุด วิธีซิปเทสท์ และวิธีเมนเทล-แยนส์เซล ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน)

เปรียบเทียบร้อยละของการตรวจพบ DIF					
แบบทดสอบ NT	กลุ่มตัวอย่าง (คน)	วิธี IRT-LR กับ วิธี SIBTEST	วิธี IRT-LR กับ วิธี MH	วิธี SIBTEST กับ วิธี MH	
ด้านภาษา	300	วิธี IRT-LR > วิธี SIBTEST* (1 ข้อ ร้อยละ 3)	วิธี IRT-LR < วิธี MH* (2 ข้อ ร้อยละ 7)	วิธี SIBTEST < วิธี MH* (3 ข้อ ร้อยละ 10)	
	1,000	วิธี IRT-LR > วิธี SIBTEST* (2 ข้อ ร้อยละ 7)	วิธี IRT-LR < วิธี MH* (7 ข้อ ร้อยละ 23)	วิธี SIBTEST < วิธี MH* (9 ข้อ ร้อยละ 30)	
	2,000	วิธี IRT-LR < วิธี SIBTEST* (1 ข้อ ร้อยละ 3)	วิธี IRT-LR < วิธี MH* (16 ข้อ ร้อยละ 53)	วิธี SIBTEST < วิธี MH* (15 ข้อ ร้อยละ 50)	
ด้านคำนวณ	300	วิธี IRT-LR = วิธี SIBTEST * (0 ข้อ ร้อยละ 0)	วิธี IRT-LR > วิธี MH* (1 ข้อ ร้อยละ 3)	วิธี SIBTEST > วิธี MH* (1 ข้อ ร้อยละ 3)	
	1,000	วิธี IRT-LR > วิธี SIBTEST* (3 ข้อ ร้อยละ 10)	วิธี IRT-LR = วิธี MH* (0 ข้อ ร้อยละ 0)	วิธี SIBTEST < วิธี MH* (3 ข้อ ร้อยละ 10)	
	2,000	วิธี IRT-LR < วิธี SIBTEST* (3 ข้อ ร้อยละ 10)	วิธี IRT-LR < วิธี MH* (5 ข้อ ร้อยละ 17)	วิธี SIBTEST < วิธี MH* (2 ข้อ ร้อยละ 7)	
ด้านเหตุผล	300	วิธี IRT-LR > วิธี SIBTEST* (1 ข้อ ร้อยละ 3)	วิธี IRT-LR < วิธี MH* (2 ข้อ ร้อยละ 7)	วิธี SIBTEST < วิธี MH* (3 ข้อ ร้อยละ 10)	
	1,000	วิธี IRT-LR < วิธี SIBTEST* (1 ข้อ ร้อยละ 3)	วิธี IRT-LR < วิธี MH* (10 ข้อ ร้อยละ 33)	วิธี SIBTEST < วิธี MH* (9 ข้อ ร้อยละ 30)	
	2,000	วิธี IRT-LR < วิธี SIBTEST* (2 ข้อ ร้อยละ 7)	วิธี IRT-LR < วิธี MH* (12 ข้อ ร้อยละ 40)	วิธี SIBTEST < วิธี MH* (10 ข้อ ร้อยละ 33)	

* หมายเหตุ $p < .05$

วิธี IRT-LR > วิธี SIBTEST หมายถึง วิธี IRT-LR ตรวจพบ DIF มากกว่า วิธี SIBTEST

วิธี IRT-LR < วิธี SIBTEST หมายถึง วิธี IRT-LR ตรวจพบ DIF น้อยกว่า วิธี SIBTEST

วิธี IRT-LR = วิธี SIBTEST หมายถึง วิธี IRT-LR ตรวจพบ DIF เท่ากับ วิธี SIBTEST

วิธี IRT-LR > วิธี MH หมายถึง วิธี IRT-LR ตรวจพบ DIF มากกว่า วิธี MH

วิธี IRT-LR < วิธี MH หมายถึง วิธี IRT-LR ตรวจพบ DIF น้อยกว่า วิธี MH

วิธี IRT-LR = วิธี MH หมายถึง วิธี IRT-LR ตรวจพบ DIF เท่ากับ วิธี MH

วิธี SIBTEST > วิธี MH หมายถึง วิธี SIBTEST ตรวจพบ DIF มากกว่า วิธี MH

วิธี SIBTEST < วิธี MH หมายถึง วิธี SIBTEST ตรวจพบ DIF น้อยกว่า วิธี MH

จากตารางที่ 4-30 ผลการเปรียบเทียบการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบด้วย วิธีอัตราส่วนไลค์ลิชุด วิธีซิปเทสท์ และวิธีเมนเทล-แยนส์เซล ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน) พบว่า

สรุปผลการเปรียบเทียบการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวณ และด้านเหตุผล ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน) ด้วยวิธีการทดสอบอัตราส่วนไลคอสิฎ्ट วิธีซิปเทส์ และวิธีแมนเทล-แยนส์เซล พบร่วมด้านภาษา วิธีการทดสอบอัตราส่วนไลคอสิฎ्ट พบ DIF มากกว่า วิธีซิปเทส์ ร้อยละ 7 ของข้อสอบทั้งหมด วิธีการทดสอบอัตราส่วนไลคอสิฎ्ट พบ DIF น้อยกว่า วิธีแมนเทล- แยนส์เซล ร้อยละ 80 ของข้อสอบทั้งหมด และ วิธีซิปเทส์ พบ DIF น้อยกว่า วิธีแมนเทล-แยนส์เซล ร้อยละ 87 ของข้อสอบทั้งหมด ส่วนด้านคำนวณ วิธีอัตราส่วนไลคอสิฎ्ट พบ DIF เท่ากันกับ วิธีซิปเทส์ ร้อยละ 0 ของข้อสอบทั้งหมด วิธีอัตราส่วนไลคอสิฎ्ट พบ DIF น้อยกว่าวิธีแมนเทล-แยนส์เซล ร้อยละ 13 ของข้อสอบทั้งหมด และ วิธีซิปเทส์ พบ DIF น้อยกว่าวิธีแมนเทล-แยนส์เซล ร้อยละ 13 ของข้อสอบทั้งหมด และส่วนด้านเหตุผล วิธีอัตราส่วนไลคอสิฎ्ट พบ DIF น้อยกว่าวิธีซิปเทส์ ร้อยละ 7 ของข้อสอบทั้งหมด วิธีอัตราส่วนไลคอสิฎ्ट พบ DIF น้อยกว่าวิธีแมนเทล- แยนส์เซล ร้อยละ 80 ของข้อสอบทั้งหมด และ วิธีซิปเทส์ พบ DIF น้อยกว่าวิธีแมนเทล-แยนส์เซล ร้อยละ 73 ของข้อสอบทั้งหมดอย่างมีนัยสำคัญทางสถิติที่ .05

บทที่ 5

สรุปและอภิปรายผล

การวิจัยครั้งนี้มีวัตถุประสงค์เพื่อวิเคราะห์คุณภาพของแบบทดสอบระดับชาติ ตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ และเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ ในแบบทดสอบระดับชาติ ของนักเรียนชั้นประถมศึกษาปีที่ 3 จำนวน 3 ตัวน คือ ต้านภาษา ด้านคำนวน และด้านเหตุผล โดยมีวิธีการดำเนินการวิจัยเป็น 3 ระยะ ดังนี้ ระยะที่ 1 การวิเคราะห์คุณภาพของแบบทดสอบระดับชาติ โดยใช้ทฤษฎีการตอบสนองข้อสอบ (IRT) แบบ 3 พารามิเตอร์ ระยะที่ 2 การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ ด้วยวิธีการตรวจสอบ 3 วิธี คือ วิธีการทดสอบอัตราส่วนไลค์ลิสต์ วิเคราะห์ด้วยโปรแกรมสำเร็จรูป IRTPRO วิธีซิปเหลท์ วิเคราะห์ด้วยโปรแกรมสำเร็จรูป SIBTEST และวิธีแมนเทล-แyenส์เซล วิเคราะห์ด้วยโปรแกรมสำเร็จรูป SPSS และระยะที่ 3 การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติตัวอย่างทั่วไป 3 วิธี ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน)

สรุปผลการวิจัย

1. ผลการวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน และด้านเหตุผล ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และ ขนาดใหญ่ (2,000 คน) โดยการวิเคราะห์ตามทฤษฎีการตอบสนองข้อสอบแบบ 3 พารามิเตอร์ พบว่า กลุ่มตัวอย่างขนาดเล็ก (300 คน) ด้านภาษา มีค่าอำนาจจำแนกของข้อสอบ (a) เฉลี่ยเท่ากับ 0.614 ค่าความยากของข้อสอบ (b) เฉลี่ยเท่ากับ 0.760 และค่าการเดาของข้อสอบ (c) เฉลี่ยเท่ากับ 0.255 สรุปได้ว่าแบบทดสอบด้านภาษา กลุ่มตัวอย่างขนาดเล็ก (300 คน) มีค่าอำนาจจำแนกข้อสอบอยู่ในระดับดีมาก มีค่าความยากของข้อสอบอยู่ในระดับข้อสอบค่อนข้างยาก และมีค่าการเดาของข้อสอบไม่เกิน 0.3 สำหรับ ด้านคำนวน มีค่าอำนาจจำแนกของข้อสอบ (a) เฉลี่ยเท่ากับ 0.776 ค่าความยากของข้อสอบ (b) เฉลี่ยเท่ากับ 1.476 และมีค่าการเดาของข้อสอบ (c) เฉลี่ยเท่ากับ 0.235 สรุปได้ว่าแบบทดสอบด้านคำนวน กลุ่มตัวอย่างขนาดเล็ก (300 คน) มีค่าอำนาจจำแนกข้อสอบอยู่ในระดับดีมาก มีค่าความยากของข้อสอบอยู่ในระดับข้อสอบยาก และมีค่าการเดาของข้อสอบไม่เกิน 0.3 และด้านเหตุผล มีค่าอำนาจจำแนกของข้อสอบ (a) เฉลี่ยเท่ากับ 0.762 ค่าความยากของข้อสอบ (b) เฉลี่ยเท่ากับ 1.004 และมีค่าการเดาของข้อสอบ (c) เฉลี่ยเท่ากับ 0.254 สรุปได้ว่าแบบทดสอบด้านเหตุผล กลุ่มตัวอย่างขนาดเล็ก (300 คน) มีค่าอำนาจจำแนกข้อสอบอยู่ในระดับดีมาก มีค่าความยากของข้อสอบอยู่ในระดับข้อสอบยาก และมีค่าการเดาของข้อสอบไม่เกิน 0.3

กลุ่มตัวอย่างขนาดกลาง (1,000 คน) ด้านภาษา มีค่าอำนาจจำแนกของข้อสอบ (a) เฉลี่ยเท่ากับ 0.644 ค่าความยากของข้อสอบ (b) เฉลี่ยเท่ากับ 0.733 และมีค่าการเดาของข้อสอบ (c) เฉลี่ยเท่ากับ 0.250 สรุปได้ว่าแบบทดสอบด้านภาษา กลุ่มตัวอย่างขนาดกลาง (1,000 คน)

มีค่าอำนาจจำแนกข้อสอบอยู่ในระดับต่ำมาก มีค่าความยากของข้อสอบอยู่ในระดับข้อสอบค่อนข้างยาก และ มีค่าการเดาของข้อสอบไม่เกิน 0.3 สำหรับด้านคำนวน มีค่าอำนาจจำแนกของข้อสอบ (a) เฉลี่ยเท่ากับ 0.867 ค่าความยากของข้อสอบ (b) เฉลี่ยเท่ากับ 1.297 และมีค่าการเดาของข้อสอบ (c) เฉลี่ยเท่ากับ 0.231 สรุปได้ว่าแบบทดสอบด้านคำนวน กลุ่มตัวอย่างขนาดกลาง (1,000 คน) มีค่าอำนาจจำแนกข้อสอบอยู่ในระดับต่ำมาก มีค่าความยากของข้อสอบอยู่ในระดับข้อสอบยาก และ มีค่าการเดาของข้อสอบไม่เกิน 0.3 และด้านเหตุผล มีค่าอำนาจจำแนกของข้อสอบ (a) เฉลี่ยเท่ากับ 0.762 ค่าความยากของข้อสอบ (b) เฉลี่ยเท่ากับ 0.856 และมีค่าการเดาของข้อสอบ (c) เฉลี่ยเท่ากับ 0.246 สรุปได้ว่าแบบทดสอบ ด้านเหตุผล กลุ่มตัวอย่างขนาดกลาง (1,000 คน) มีค่าอำนาจจำแนกข้อสอบอยู่ในระดับต่ำมาก มีค่าความยากของข้อสอบอยู่ในระดับข้อสอบค่อนข้างยาก และมีค่าการเดาของข้อสอบไม่เกิน 0.3

กลุ่มตัวอย่างขนาดใหญ่ (2,000 คน) ด้านภาษา มีค่าอำนาจจำแนกของข้อสอบ (a) เฉลี่ยเท่ากับ 0.643 ค่าความยากของข้อสอบ (b) เฉลี่ยเท่ากับ 0.476 และมีค่าการเดาของข้อสอบ (c) เฉลี่ยเท่ากับ 0.211 สรุปได้ว่าแบบทดสอบ ด้านภาษา กลุ่มตัวอย่างขนาดใหญ่ (2,000 คน) มีอำนาจจำแนกข้อสอบอยู่ในระดับต่ำมาก มีค่าความยากของข้อสอบอยู่ในระดับข้อสอบปานกลาง และมีค่าการเดาของข้อสอบไม่เกิน 0.3 สำหรับ ด้านคำนวน มีค่าอำนาจจำแนกของข้อสอบ (a) เฉลี่ยเท่ากับ 0.836 ค่าความยากของข้อสอบ (b) มีค่าเฉลี่ยเท่ากับ 1.465 และมีค่าการเดาของข้อสอบ (c) มีค่าเฉลี่ยเท่ากับ 0.228 สรุปได้ว่าแบบทดสอบ ด้านคำนวน กลุ่มตัวอย่างขนาดใหญ่ (2,000 คน) มีอำนาจจำแนกข้อสอบอยู่ในระดับต่ำมาก มีค่าความยากของข้อสอบอยู่ในระดับข้อสอบยาก และมีค่าการเดาของข้อสอบไม่เกิน 0.3 และ ด้านเหตุผล มีค่าอำนาจจำแนกของข้อสอบ (a) เฉลี่ยเท่ากับ 0.742 ค่าความยากของข้อสอบ (b) เฉลี่ยเท่ากับ 0.851 และมีค่าการเดาของข้อสอบ (c) มีค่าเฉลี่ยเท่ากับ 0.239 สรุปได้ว่าแบบทดสอบ ด้านเหตุผล กลุ่มตัวอย่างขนาดใหญ่ (2,000 คน) มีค่าอำนาจจำแนกข้อสอบอยู่ในระดับต่ำมาก มีค่าความยากของข้อสอบอยู่ในระดับข้อสอบค่อนข้างยาก และมีค่าการเดาของข้อสอบไม่เกิน 0.3

2. ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ขั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน และ ด้านเหตุผล ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง ต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน) พบร่วมกัน วิธีการทดสอบอัตราส่วนไลคริสตี้ ตรวจพบข้อสอบที่ทำหน้าที่ต่างกัน (DIF) กลุ่มตัวอย่างขนาดเล็ก (300 คน) ด้านภาษา จำนวน 3 ข้อ คิดเป็นร้อยละ 10 สำหรับ ด้านคำนวน พบร DIF จำนวน 3 ข้อ คิดเป็นร้อยละ 10 และด้านเหตุผล พบร DIF จำนวน 3 ข้อ คิดเป็นร้อยละ 10 ส่วนกลุ่มตัวอย่างขนาดกลาง (1,000 คน) ด้านภาษา จำนวน 4 ข้อ คิดเป็นร้อยละ 13 สำหรับ ด้านคำนวน พบร DIF จำนวน 4 ข้อ คิดเป็นร้อยละ 13 และ ด้านเหตุผล พบร DIF จำนวน 5 ข้อ คิดเป็นร้อยละ 16 และ กลุ่มตัวอย่างขนาดใหญ่ (2,000 คน) ด้านภาษา จำนวน 3 ข้อ คิดเป็นร้อยละ 10 สำหรับ ด้านคำนวน พบร DIF จำนวน 3 ข้อ คิดเป็นร้อยละ 10 และ ด้านเหตุผล พบร DIF จำนวน 10 ข้อ คิดเป็นร้อยละ 33 ส่วนวิธีซิปเพ斯ท์ ตรวจพบข้อสอบที่ทำหน้าที่ต่างกัน (DIF) กลุ่มตัวอย่างขนาดเล็ก (300 คน) ด้านภาษา จำนวน 2 ข้อ คิดเป็นร้อยละ 7 ด้านคำนวน พบร DIF จำนวน 3 ข้อ คิดเป็นร้อยละ 10 และ ด้านเหตุผล พบร DIF จำนวน 2 ข้อ คิดเป็นร้อยละ 7 กลุ่มตัวอย่างขนาดกลาง (1,000 คน)

ด้านภาษา จำนวน 2 ข้อ คิดเป็นร้อยละ 7 ด้านคำนวน พบ DIF จำนวน 1 ข้อ คิดเป็นร้อยละ 3 และด้านเหตุผล พบ DIF จำนวน 6 ข้อ คิดเป็นร้อยละ 20 และกลุ่มตัวอย่างขนาดใหญ่ (2,000 คน) ด้านภาษาพบว่า จำนวน 4 ข้อ คิดเป็นร้อยละ 13 ด้านคำนวน พบ DIF จำนวน 6 ข้อ คิดเป็นร้อยละ 20 และด้านเหตุผล พบ DIF จำนวน 12 ข้อ คิดเป็นร้อยละ 40 และวิธีแมนเทล-แยนส์เซล ตรวจพบข้อสอบที่ทำหน้าที่ต่างกัน (DIF) กลุ่มตัวอย่างขนาดเล็ก (300 คน) ด้านภาษา จำนวน 5 ข้อ คิดเป็นร้อยละ 17 สำหรับด้านคำนวน พบ DIF จำนวน 2 ข้อ คิดเป็นร้อยละ 7 และด้านเหตุผลพบ DIF จำนวน 5 ข้อ คิดเป็นร้อยละ 17 ส่วนกลุ่มตัวอย่างขนาดกลาง (1,000 คน) ด้านภาษา จำนวน 11 ข้อ คิดเป็นร้อยละ 37 สำหรับด้านคำนวน พบ DIF จำนวน 4 ข้อ คิดเป็นร้อยละ 13 และด้านเหตุผล พบ DIF จำนวน 15 ข้อ คิดเป็นร้อยละ 50 และกลุ่มตัวอย่างขนาดใหญ่ (2,000 คน) ด้านภาษา จำนวน 19 ข้อ คิดเป็นร้อยละ 63 ด้านคำนวน พบ DIF จำนวน 8 ข้อ คิดเป็นร้อยละ 27 และด้านเหตุผลพบ DIF จำนวน 22 ข้อ คิดเป็นร้อยละ 73 อย่างมีนัยสำคัญทางสถิติที่ .05

3. ผลการเปรียบเทียบการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน และด้านเหตุผล ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน) ด้วยวิธีการทดสอบอัตราส่วนไลค์ลิขุด วิธีซิปเทสท์ และวิธีแมนเทล-แยนส์เซล พบว่าด้านภาษา วิธีการทดสอบอัตราส่วนไลค์ลิขุด พบ DIF มากกว่า วิธีซิปเทสท์ ร้อยละ 7 ของข้อสอบทั้งหมด วิธีการทดสอบอัตราส่วนไลค์ลิขุด พบ DIF น้อยกว่า วิธีแมนเทล- แยนส์เซล ร้อยละ 80 ของข้อสอบทั้งหมด และ วิธีซิปเทสท์ พบ DIF น้อยกว่า วิธีแมนเทล- แยนส์เซล ร้อยละ 87 ของข้อสอบทั้งหมด ส่วนด้านคำนวน วิธีอัตราส่วนไลค์ลิขุด พบ DIF เท่ากันกับ วิธีซิปเทสท์ ร้อยละ 0 ของข้อสอบทั้งหมด วิธีอัตราส่วนไลค์ลิขุด พบ DIF น้อยกว่า วิธีแมนเทล- แยนส์เซล ร้อยละ 13 ของข้อสอบทั้งหมด และ วิธีซิปเทสท์ พบ DIF น้อยกว่า วิธีแมนเทล- แยนส์เซล ร้อยละ 13 ของข้อสอบทั้งหมด และ ส่วนด้านเหตุผล วิธีอัตราส่วนไลค์ลิขุด พบ DIF น้อยกว่า วิธีซิปเทสท์ ร้อยละ 7 ของข้อสอบทั้งหมด วิธีอัตราส่วนไลค์ลิขุด พบ DIF น้อยกว่า วิธีแมนเทล- แยนส์เซล ร้อยละ 80 ของข้อสอบทั้งหมด และ วิธีซิปเทสท์ พบ DIF น้อยกว่า วิธีแมนเทล- แยนส์เซล ร้อยละ 73 ของข้อสอบทั้งหมด อย่างมีนัยสำคัญทางสถิติที่ .05

อภิปรายผล

ผลการวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน และด้านเหตุผล โดยใช้ทฤษฎีการตอบสนองข้อสอบแบบ 3 พารามิเตอร์ การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน) และการเปรียบเทียบผลการตรวจพบข้อสอบทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ทั้ง 3 วิธี มีประเด็นที่ควรอภิปราย ดังนี้

1. การวิเคราะห์คุณภาพของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ด้านภาษา ด้านคำนวน และด้านเหตุผล การวิเคราะห์คุณภาพของแบบทดสอบโดยใช้หลักการทฤษฎีการตอบสนองข้อสอบ (IRT) แบบ 3 พารามิเตอร์ ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก (300 คน) ขนาดกลาง (1,000 คน) และขนาดใหญ่ (2,000 คน) ประกอบด้วยค่าอำนาจจำแนกของข้อสอบ (a)

ค่าความยากของข้อสอบ (b) และค่าโอกาสการเดาของข้อสอบ (c) แบบทดสอบระดับชาติ ด้านภาษา มีค่าอำนาจจำแนกของข้อสอบ (a) ทั้งฉบับอยู่ในระดับดีมาก มีค่าความยากของข้อสอบ (b) ทั้งฉบับอยู่ในระดับปานกลาง ถึงค่อนข้างยาก และมีค่าโอกาสการเดาของข้อสอบ (c) ทั้งฉบับไม่เกิน 0.3 สำหรับด้านคำนวน มีค่าอำนาจจำแนกของข้อสอบ (a) ทั้งฉบับอยู่ในระดับดีมาก มีค่าความยากของข้อสอบ (b) ทั้งฉบับอยู่ในระดับยาก และมีค่าโอกาสการเดาของข้อสอบ (c) ทั้งฉบับไม่เกิน 0.3 และ ส่วนด้านเหตุผล มีค่าอำนาจจำแนกของข้อสอบ (a) ทั้งฉบับอยู่ในระดับดีมาก มีค่าความยากของข้อสอบ (b) ทั้งฉบับอยู่ในระดับค่อนข้างยาก และมีค่าโอกาสการเดาของข้อสอบ (c) ทั้งฉบับไม่เกิน 0.3

2. การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ ด้วยวิธีการตรวจสอบ 3 วิธี คือ วิธี IRT-LR วิธี SIBTEST และ วิธี Mantel-Haenszel ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก ขนาดกลาง และขนาดใหญ่ พบว่า วิธี Mantel-Haenszel สามารถตรวจพบข้อสอบที่ทำหน้าที่ต่างกัน (DIF) ได้มากที่สุด รองลงมา คือ วิธี IRT-LR และวิธี SIBTEST ตามลำดับ ซึ่ง สอดคล้องกับงานวิจัยของ สุราทิพย์ ตรีสิน และปิยะพิพิร ประดุจพร (2560) พบว่า วิธี IRT-LR จะมีประสิทธิภาพในการตรวจสอบ DIF ได้ดี นอกจากนี้ยังสอดคล้องกับงานวิจัยของ Awuor (2008) ที่ได้ศึกษาขนาดกลุ่มตัวอย่างไม่เท่ากันในการตรวจสอบ DIF ด้วยวิธี IRT-Based วิธี SIBTEST และวิธี Mantel-Haenszel พบว่า ขนาดตัวอย่างต่างกันจะส่งผลต่อการตรวจสอบ DIF โดยวิธี SIBTEST และ วิธี Mantel-Haenszel มีประสิทธิภาพในการตรวจสอบ DIF ได้ดีเมื่อกลุ่มตัวอย่างขนาดกลาง และ ขนาดใหญ่ และสอดคล้องงานวิจัยของ พีรญา สูงเนิน เสรี ชัดแจ่ม และ สมโภชน์ อเนกสุข (2552) เปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบแบบพหุมิติระหว่าง ข้อสอบรายข้อกับหมวดข้อสอบ ด้วยวิธี SIBTEST ผลการศึกษาพบว่า ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก ขนาดกลาง และขนาดใหญ่ ระหว่างการตรวจสอบการทำหน้าที่ต่างกัน เป็นรายข้อกับรายหมวดข้อสอบ พบที่ต่างกันแต่ต่างกัน มีนัยสำคัญทางสถิติที่ระดับ .05 และขนาดของกลุ่มตัวอย่างต่างกันมีผลต่อการตรวจพบข้อสอบทำหน้าที่ต่างกัน เมื่อขนาดของกลุ่มตัวอย่างใหญ่ขึ้นจะทำให้สามารถตรวจพบข้อสอบทำหน้าที่ต่างกันได้ดีกว่ากลุ่มตัวอย่างขนาดเล็ก ซึ่งงานวิจัยของ Kabasakal, Arsan, Gok, and Kelecioglu (2014) ศึกษาประสิทธิภาพในการตรวจสอบ DIF ด้วยวิธี Mantel-Haenszel วิธี SIBTEST และวิธี IRT-LR พบว่า ประสิทธิภาพในการตรวจสอบ DIF ได้ดี และ Li, Hunter and Oshima (2013) ได้ศึกษาการตรวจสอบการทำหน้าที่ต่างกันในแบบทดสอบด้านการอ่าน และด้านเหตุผล ด้วยวิธีการทดสอบอัตราส่วนไลค์ลิชฎ และวิธี แมนเทล-แ昏ส์เซล พบว่า มีประสิทธิภาพในการตรวจสอบ DIF ได้ดีในแบบทดสอบที่มีความยาวของข้อสอบตั้งแต่ 20 ข้อขึ้นไป

3. การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ ทั้ง 3 วิธี ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่างต่างกัน คือ ขนาดเล็ก ขนาดกลาง และขนาดใหญ่ พบว่า วิธี Mantel-Haenszel ตรวจพบ DIF มากกว่าวิธี SIBTEST และวิธี IRT-LR ทั้ง 3 ด้าน คือ ด้านภาษา ด้านคำนวน และด้านเหตุผล ส่วนวิธี Mantel-Haenszel ตรวจพบ DIF มากกว่าวิธี IRT-LR ในด้านภาษา และด้านเหตุผล ซึ่งสอดคล้องกับงานวิจัยของ Awuor (2008) ที่ได้ศึกษาขนาดกลุ่มตัวอย่างไม่เท่ากันในการตรวจสอบ DIF ด้วยวิธี IRT-Based วิธี SIBTEST และวิธี Mantel-Haenszel

พบว่าขนาดตัวอย่างที่แตกต่างกัน วิธี SIBTEST และวิธี Mantel-Haenszel ตรวจพบ DIF ได้ดีในกลุ่มตัวอย่างขนาดกลางและขนาดใหญ่ และตรวจพดได้ใกล้เคียงกัน และสอดคล้องกับงานวิจัยของ Kabasakal, Arsan, Gok, and Kelecioglu (2014) ได้ศึกษาประสิทธิภาพในการตรวจสอบ DIF ด้วยวิธี Mantel-Haenszel วิธี SIBTEST และวิธี IRT-LR ผลการศึกษาพบว่า มีประสิทธิภาพในการตรวจสอบ DIF ได้ดี ส่วนงานวิจัยของ Li, Hunter and Oshima (2013) ได้ศึกษาการตรวจสอบการทำหน้าที่ต่างกันในแบบทดสอบด้านการอ่าน และด้านเหตุผล ด้วยวิธีการทดสอบอัตราส่วนไลค์ลิขิต และวิธีแมนเทล-แยนส์เซล พบรวม มีประสิทธิภาพในการตรวจสอบ DIF ได้ดีในแบบทดสอบที่มีความยาวของข้อสอบตั้งแต่ 20 ข้อขึ้นไป และงานวิจัยของ Yildirim and Berberoglu (2009) ได้วิเคราะห์ข้อมูลการตัดสินใจและทางสถิติของวิชาคณิตศาสตร์ PISA-2003 พบรวม วิธี IRT-LR มีประสิทธิภาพในการตรวจสอบ DIF ได้ดีในด้านความสามารถทางคณิตศาสตร์ของโครงการประเมินผลนักเรียนนานาชาติ (PISA, 2003)

ข้อเสนอแนะในการนำผลการวิจัยไปใช้

จากผลการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ และการเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบระดับชาติ มีข้อเสนอแนะในการนำผลการวิจัยไปใช้ ดังนี้

- สำนักทดสอบทางการศึกษาสามารถนำผลการวิเคราะห์คุณภาพของแบบทดสอบระดับชาติที่ผ่านเกณฑ์การวิเคราะห์คุณภาพโดยใช้หลักการของทฤษฎีการตอบสนองข้อสอบแบบ 3 พารามิเตอร์ ไปใช้สอบในครั้งต่อไป เพื่อใช้สำหรับวัดความสามารถของนักเรียน ขั้นปฐมศึกษาปีที่ 3 ของสำนักทดสอบทางการศึกษาสำนักงานคณะกรรมการการศึกษาขั้นพื้นฐาน (สพฐ.)
- นักวิจัยและนักวัดผลการศึกษาที่สนใจเกี่ยวกับการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ ด้วยวิธีการตรวจสอบที่อยู่บนพื้นฐานของทฤษฎีการตอบสนองข้อสอบ (IRT) ทั้ง 3 วิธี คือ วิธีอัตราส่วนไลค์ลิขิต วิธีชิปเพลท์ และวิธีแมนเทล-แยนส์เซล ขนาดของกลุ่มตัวอย่างต้องมีขนาดใหญ่ จะทำให้สามารถตรวจพบข้อสอบทำหน้าที่ต่างกันได้ดีกว่ากลุ่มตัวอย่างขนาดเล็ก

ข้อเสนอแนะสำหรับการวิจัยต่อไป

- วิธีอัตราส่วนไลค์ลิขิต วิธีชิปเพลท์ และวิธีแมนเทล-แยนส์เซล มีประสิทธิภาพในการตรวจพบข้อสอบที่ทำหน้าที่ต่างกันของข้อสอบ เมื่อกลุ่มตัวอย่างมีขนาดใหญ่ จะทำให้สามารถตรวจสอบการทำหน้าที่ต่างกันของข้อสอบได้ดีกว่ากลุ่มตัวอย่างขนาดเล็ก จึงควรมีการเปรียบเทียบเพิ่มเติมกับวิธีการตรวจสอบอื่น ๆ และศึกษาตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบที่มีการตรวจให้คะแนนแบบมากกว่า 2 ค่า
- ควรมีการศึกษาวิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบที่มีการตรวจให้คะแนนแบบมากกว่า 2 ค่า ด้วยวิธี Standard MIMIC (M-ST) วิธี MIMIC with Scale Purification (M-SP) และวิธี MIMIC with Pure Anchor (M-PA) ว่า วิธีใดมีประสิทธิภาพในการตรวจสอบ DIF มากกว่ากัน

บรรณานุกรม

- จารุจิตร สิทธิปรุ ปิยะทิพย์ ตินวร และโสพส สุขานนท์สวัสดิ์. (2559). การพัฒนาโปรแกรมการทดสอบ แบบปรับเหมาะสมด้วยคอมพิวเตอร์ สำหรับการจัดสอบ O-NET ระดับชั้นมัธยมศึกษาปีที่ 3. วารสารการวัดผลการศึกษา มหาวิทยาลัยมหาสารคาม, 22(2), 47-62.
- ชนะศึก นิชานนท์. (2553). ประสิทธิภาพของการประมาณค่าพารามิเตอร์แบบเบส์โดยใช้การสรุปอ้างอิงความน่าเชื่อถือของโมเดลการตอบสนองข้อสอบ. วารสารวิจัย มสด., 11(2), 61-75.
- ชัยวัฒน์ ฤทธิพันธ์. (2558). การพัฒนาวิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบโดยการตัดสินของผู้เชี่ยวชาญ. วารสารครุศาสตร์, 43(1), 1-18.
- ณรงค์ จันทร์. (2554). การเปรียบเทียบค่าความเที่ยงของแบบทดสอบผลสัมฤทธิ์ทางการเรียนที่มีจำนวนข้อสอบทำหน้าที่ต่างกันแตกต่างกัน. วิทยาการวิจัยและวิทยาการปัญญา, 8(2), 58-71.
- ธเกียรติกมล ทองอก, โขติกา ภาษี และศิริชัย ภานจนวาสี. (2556). ประสิทธิภาพการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบภาษาไทยวิจิตด้วยโลจิสติกสำหรับข้อสอบที่ตรวจให้คะแนนแบบทวิภาค: เปรียบเทียบขนาดอิทธิพลสองเกณฑ์. วารสารวิจัย มสด., 31-49.
- ธีรวัฒน์ สุขีสาร, ดุษฎี โยเหลา, เสกสรรค์ ทองคำบรรจง และนิยมดา จิตต์จรัสร. (2555). การศึกษาความเที่ยงตรงของการประมาณค่าในการวิเคราะห์โมเดลสมการโครงสร้างพหุระดับภายใต้เงื่อนไขวิธีการประมาณค่าและขนาดตัวอย่างที่ต่างกัน. วารสารการวัดผลการศึกษา, 17(1), 95-106.
- นุภาพรณ ปลื้มใจ ปิยะทิพย์ตินวร และโสพส สุขานนท์สวัสดิ์. (2558). การพัฒนาโปรแกรมการทดสอบแบบปรับเหมาะสมด้วยคอมพิวเตอร์สำหรับการจัดสอบ O-NET ระดับชั้นมัธยมศึกษาปีที่ 6. วิทยาการวิจัยและวิทยาการปัญญา, 13(2), 109-125.
- ปนัดดา หัส平原. (2557). แนวทางการนำผลการทดสอบทางการศึกษาไปใช้ในการพัฒนาคุณภาพผู้เรียน. สถาบันทดสอบทางการศึกษาแห่งชาติ (องค์กรมหาชน). วันที่ค้นข้อมูล 1 สิงหาคม 2561, เข้าถึงได้จาก www.niets.or.th/th/content/download/262.
- ปิยะทิพย์ ตินวร, ม.ร.ว.สมพร สุทัศน์ และเสรี ชัดแข็ม. (2550). การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบพหุมิติ: การเปรียบเทียบประสิทธิภาพระหว่างวิธีการวิเคราะห์องค์ประกอบจำกัดกับวิจิตด้วยโลจิสติก. วารสารวิจัยและวัดผลการศึกษา, 5(1), 63-80.
- พีรญา สูงเนิน, เสรี ชัดแข็ม และสมโภชน์ อเนกสุข (2552). การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในแบบทดสอบพหุมิติ: การเปรียบเทียบระหว่างข้อสอบรายข้อกับหมวดข้อสอบ โดยใช้วิธีซีปเทสท์. วิทยาการวิจัยและวิทยาการปัญญา, 6(2), 49-62.
- รุ่งนภา แสนอำนวย, ประกฤติยา ทักษิโน และชนะศึก นิชานนท์ (2555). ประสิทธิภาพของแบบทดสอบ วัดผลสัมฤทธิ์ทางการเรียนรูปแบบผสม: การประยุกต์ใช้ทฤษฎีการตอบสนอง ข้อสอบแบบ ตรวจให้คะแนนความรู้บางส่วน และทฤษฎีการตอบสนองข้อสอบแบบตรวจให้คะแนนความรู้บางส่วนแบบทั่วไป. วารสารศึกษาศาสตร์ มหาวิทยาลัยขอนแก่น, 35(1), 58-66.

- วรรณ ศรีกัลា, อุ่น เจริญวงศ์ระยับ และนิคม นาคอ้าย (2559). ปัจจัยพหุระดับที่ส่งผลต่อ
คะแนนการสอบประเมินคุณภาพการศึกษาระดับชาติ ด้านความสามารถทางภาษา:
การศึกษาของโรงเรียน ที่มีผล NT ต่ำ ในจังหวัดพิษณุโลก. วารสารราชภัฏสุราษฎร์ธานี,
3(2), 81-98.
- ศาสตรา แสนปัญญา. (2555). การศึกษาพัฒนาการความสามารถทางวิทยาศาสตร์โดยใช้วิธีปรับแนวตั้ง¹
ตามทฤษฎีการตอบสนองข้อสอบ. ใน การประชุมวิชาการแห่งชาติ ครั้งที่ 9 วันที่ 6-7
ธันวาคม 2555 (หน้า 1493 - 1499) นครปฐม: มหาวิทยาลัยเกษตรศาสตร์ วิทยาเขต
กำแพงแสน.
- ศรีชัย กาญจนาวาสี. (2555). ทฤษฎีการทดสอบแนวใหม่ (พิมพ์ครั้งที่ 4). กรุงเทพฯ: โรงพิมพ์
จุฬาลงกรณ์มหาวิทยาลัย.
- สำนักทดสอบทางการศึกษา. (2557). ผลการประเมินคุณภาพผู้เรียนระดับชาติ ปีการศึกษา 2555
บทสรุปและข้อเสนอแนะเชิงนโยบาย. กรุงเทพฯ: โรงพิมพ์จุฬาลงกรณ์การเกษตร
แห่งประเทศไทย.
- สุชาดา มงคลกิจรุ่งโรจน์ เสรี ชัดแข็ง และสมพร สุทัคนีย์. (2559). ประสิทธิภาพของโปรแกรม
การทดสอบแบบปรับเหมาะสมด้วยคอมพิวเตอร์สำหรับมาตรฐานความสูงของคนไทย.
วารสารการวัดผลการศึกษา มหาวิทยาลัยมหาสารคาม, 22(2), 321-331.
- สุทธิวรรณ พิรศักดิ์สกาน, ปิยพงษ์ คล้ายคลึง และสมกิจ กิจพุนวงศ์. (2560). การศึกษาคุณภาพ
แบบทดสอบความถนัดทางการเรียน SWUSAT ปีการศึกษา 2553 -2555 ตามทฤษฎี
การทดสอบแบบมาตรฐานเดิมและทฤษฎีการตอบข้อสอบ. วารสารศึกษาศาสตร์
มหาวิทยาลัยหกชั้น, 17(1), 65-73.
- สุราทิพย์ ตรีสิน และปิยะพิพิญ ประดุจพร. (2560). การเปรียบเทียบการตรวจสอบการทำหน้าที่
ต่างกันของข้อสอบในแบบทดสอบระดับชาติ ด้านภาษา ด้านคำนวณ และด้านเหตุผล
ชั้นประถมศึกษาปีที่ 3 ด้วยวิธี HGLM วิธี MIMIC และวิธี IRT-LR. วิทยาการวิจัยและ
วิทยาการปัญญา, 15(2), 109-119.
- สุภาภรณ์ แดงเพ็ง, ชิดชนก เชิงเข้าว์ และบุญญิสิสา แซ่หล่อ. (2554). การเปรียบเทียบความจำเอียงของ
แบบทดสอบคณิตศาสตร์ในการประเมินคุณภาพการศึกษาระดับท้องถิ่นชั้นประถมศึกษา²
ปีที่ 5 จังหวัดปัตตานี ระหว่างวิธีแมนเทล-แணส์เซลและโค้งลักษณะข้อสอบ 3 พารามิเตอร์.
วารสารสังคมครiminology, 17(2), 135-326.
- อนุชิต กลืนกำเนิด, อรจิรา สิทธิศักดิ์ และทศนวรรตน์ ศุนย์กลาง. (2555). ผลการประเมินระบบบริหาร
จัดการการเรียนรู้แบบปรับเหมาะสม กรณีศึกษาเรื่อง องค์ประกอบของระบบสารสนเทศ.
วารสารวิชาการทางเทคโนโลยีคอมพิวเตอร์และระบบสารสนเทศประยุกต์, 1(2), 1-7.
- เอกลักษณ์ คล้ายสุบรรณ, สั่งวนิษ์ จัดกระไก กะ และนันี ณ นคร. (2559). ไมเดลการวัดมูลค่าเพิ่ม³
ทางการศึกษาสำหรับคุณภาพสถานศึกษาด้วยการใช้ผลรวมของผลสัมฤทธิ์ทางการเรียน
และผลการประเมินและรับรองคุณภาพของโรงเรียน. Veridian E-Journal, 9(1),
1041-1052.

- Acar, T., & Kelecioglu, H. (2010). Comparison of Differential Item Functioning Determination Techniques: HGLM, LR and IRT-LR. *Educational Sciences: Theory and Practice*, 10(2), 639-649.
- Acar.T (2011) Sample Size in Differential Item Functioning: An Application of Hierarchical Linear Modeling. *Educational Sciences: Theory & Practice*, 11(1), 284-288.
- Apinyapibal, S., Lawthong, N., Kanjanawasee, S. (2015). A Comparative Analysis of the Efficacy of Differential Item Functioning Detection For Dichotomously Scored Items among Logistic Regression SIBTEST and Raschtree Methods. *Procedia-Social and Behavioral Sciences*, 191, 21-25.
- Atar, B. (2007). *Differential Item Functioning Analyses for Mixed Response data using IRT Likelihood-Ratio Test, Logistic Regression, and GLLMM Procedures*. Doctor of Philosophy, Department of Educational Psychology and Learning Systems.
- Awuor, R. A. (2008). *Effect of Unequal Sample Sizes on the Power of DIF Detection: An IRT-Based Monte Carlo Study with SIBTEST and Mantel-Haenszel Procedures*. Doctor of Philosophy, Educational Research and Evaluation, Virginia Polytechnic Institute and State University, 1-111.
- Barnes, B. J., & Wells, C. S. (2009). Differential Item Functional Analysis by Gender and Race of the National Doctoral Program Survey. *International Journal of Doctoral Studies*, 4, 77-96.
- Breland, H., Lee, Y. (2007). Investigating uniform and non-uniform gender DIF in computer-based ESL writing assessment. *Applied Measurement in Education*, 20(4), 377-403.
- Brown, L. I., Bristol, L., De Four-Babb, J., & Conrad, D. A. (2014). National Tests and Diagnostic Feedback: What Say Teachers in Trinidad and Tobago?. *Journal of Educational Research*, 107(3), 241-251.
- Embretson, S. E., & Reise, S. P. (2000). *Item response theory for psychologists*. Mahwah, New Jersey: Lawrence Erlbaum Associates, Publishers.
- González-Betanzos, F., Abad, F. J., Barrada, J. R. (2014). Fixed item parameter calibration for assessing differential item functioning in computerized adaptive tests. *Psicológica*, 35, 331-359.
- Kabasakal, K. A., Arsan, N., Gok, B., Kelecioglu, H. (2014). Comparing Performances (Type I error and Power)of IRT Likelihood Ratio SIBTEST and Mantel-Haenszel Methods in the Determination of Differential Item Functioning. *Educational Sciences: Theory & Practice*, 14(6), 2186-2193.

- Kjellstrom, K., & Pettersson, A. (2005). The curriculum's view of knowledge transferred to national tests in mathematics in Sweden. *ZDM*, 37(4), 308-316.
- Kose, I. A., & Demirtasli, N. C. (2012). Comparison of unidimensional and multidimensional models based on item response theory in terms of both variables of test length and sample size. *Procedia - Social and Behavioral Sciences*, 46, 135–140.
- Le, L. T. (2009). Investigating gender differential item functioning across countries and test languages for PISA science items. *International Journal of Testing*, 9(2), 122-133.
- Li, H., Hunter, C. V., & Oshima, T. C. (2013). *Gender DIF in Reading Tests: A Synthesis of Research*. In New Developments in Quantitative Psychology (pp. 489-506). Springer, New York: Springer Science Business Media.
- Mendes-Barnett, S., & Ercikan, K. (2006). Examining sources of gender DIF in mathematics assessments using a confirmatory multidimensional model approach. *Applied Measurement in Education*, 19(4), 289-304.
- Moghadamzadeh, A., Salehi, K., & Khodaie, E. (2011). A Comparison the Information Functions of the Item and Test in One, Two and Three Parametric Model of the Item Response Theory (IRT). *Procedia - Social and Behavioral Sciences*, 29, 1359–1367.
- Muninsakorn, Y., Tinnaworn, P., & Sukhanonsawat, S. (2015). Development of the Computerized Adaptive Testing Program for O-NET at the Grade 6 Level. In *Burapha University Internationnal Conference 2015: Moving Forward to a Prosperous and Sustainable Community, July 10-12, 2015 Bangsaen Heritafe Hotel Chonburi*. Thailand: Burapha University.
- Pae, T. I. (2012). Causes of gender DIF on an EFL language test: A multiple-data analysis over nine years. *Language Testing*, 29(4), 533-554.
- Rivas, G. E. L. (2012). *Detection and Classification of DIF Types Using Parametric and Nonparametric Methods: A comparison of the IRT-Likelihood Ratio Test, Crossing-SIBTEST, and Logistic Regression Procedures*. University of South Florida, Department of Psychology College of Arts and Sciences.
- Sterinmayr, R., Bergold, S., Stiksrud, J. M., & Freund, P. A. (2015). Gender differences on general knowledge tests: Are they due to Differential Item Functioning. *Intelligence*, 50, 164-174.
- Taylor, C. S., & Lee, Y. (2012). Gender DIF in reading and mathematics tests with mixed item formats. *Applied Measurement in Education*, 25(3), 246-280.

- Teresi, J. A., Ocepek-Welikson, K., Kleinman, M., Cook, K. F., Crane, P. K., Gibbons, L. E., & Cellia, D. (2007). Evaluating measurement equivalence using the item response theory log-likelihood ratio (IRTLR) method to assess differential item functioning (DIF): applications (with illustrations) to measures of physical functioning ability and general distress. *Quality of Life Research*, 16(1), 43-68.
- Thissen, D., Steinberg, L., & Gerrard, M. (1986). Beyond group-mean differences: The concept of item bias. *Psychological Bulletin*, 99(1), 118-128.
- Urry, V. (1977). Tailored testing: A successful application of latent trait theory. *Journal of Education Measurement*, 14, 181-196.
- Yildirim, H. H., & Berberoglu G. (2009). Judgmental and statistical DIF analyses of the PISA-2003 mathematics literacy items. *International Journal of Testing*, 9, 108-121.

ภาคผนวก

ภาคผนวก ก
หนังสือรับรองผลการพิจารณาจริยธรรมการวิจัย



**แบบรายงานผลการพิจารณาจuryธรรมการวิจัยในคน
วิทยาลัยวิทยาการวิจัยและวิทยาการปัญญา
มหาวิทยาลัยบูรพา**

๑. ชื่อเรื่องวิทยานิพนธ์

ชื่อเรื่องวิทยานิพนธ์ (ภาษาไทย) การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ระดับชั้นประถมศึกษาปีที่ ๓: ระหว่างวิธีการทดสอบอัตราส่วนไลค์ลิคิਊต วิธีซิบเพสท์ และวิธีเมนเทล-เอนล์เซล ชื่อเรื่องวิทยานิพนธ์ (ภาษาอังกฤษ) COMPARISON OF THE DIFFERENTIAL ITEM FUNCTIONING OF THE NATIONAL TEST ITEM AT GRADE 3 LEVEL BETWEEN IRT LIKELIHOOD RATIO, SIBTEST AND MENTEL-HAENSZEL METHODS

๒. ชื่อนิสิต (นาย, นาง, นางสาว): อรุณี แปลงกาย

หลักสูตรวิทยาศาสตรมหาบัณฑิต (M.Sc.) สาขาวิชาการวิจัยและสถิติทางวิทยาการปัญญา

ภาคปกติ ภาคพิเศษ

รหัสประจำตัว ๕๖๔๑๐๔๓ คณะ/วิทยาลัยวิทยาการวิจัยและวิทยาการปัญญา

๓. หน่วยงานที่สังกัด: วิทยาลัยวิทยาการวิจัยและวิทยาการปัญญา

๔. ผลการพิจารณาของคณะกรรมการพิจารณาจuryธรรมการวิจัยในคน:

คณะกรรมการพิจารณาจuryธรรมการวิจัยในคน ได้พิจารณารายละเอียดวิทยานิพนธ์เรื่องดังกล่าว ข้างต้นแล้ว ในประเด็นที่เกี่ยวกับ

- (๑) การเคารพในสักครึ่ง และสิทธิของคนที่ใช้เป็นตัวอย่างการวิจัย
- (๒) วิธีการอย่างเหมาะสมในการได้รับความยินยอมจากกลุ่มตัวอย่างก่อนเข้าร่วมโครงการวิจัย (Informed consent) รวมทั้งการป้องกันสิทธิประโยชน์ และรักษาความลับกลุ่มตัวอย่างในการวิจัย
- (๓) การดำเนินการวิจัยอย่างเหมาะสม เพื่อไม่ก่อความเสียหายต่อสิ่งที่ศึกษาวิจัย ไม่ว่าจะเป็นสิ่งที่มีชีวิต หรือไม่มีชีวิต

คณะกรรมการพิจารณาจuryธรรมการวิจัยในคน มีมติเห็นชอบ ดังนี้

(✓) รับรองโครงการวิจัย

() ไม่รับรอง

๕. วันที่ให้การรับรอง: ๒๘ เดือน มีนาคม พ.ศ. ๒๕๕๙

ลงนาม.....
อรุณี แปลงกาย

(ผู้ช่วยศาสตราจารย์ ดร.สุชาดา กรเพชรปานี)
ประธานกรรมการพิจารณาจuryธรรมการวิจัยในคน
คณบดีวิทยาลัยวิทยาการวิจัยและวิทยาการปัญญา
วันที่ ๒๘ มีนาคม พ.ศ. ๒๕๕๙

ภาคผนวก ข
หนังสือขอความอนุเคราะห์ข้อมูลเพื่อการวิจัย



ที่ ศธ ๖๙๒๔/๐๗๒๙

วิทยาลัยวิทยาการวิจัยและวิทยาการปัญญา
มหาวิทยาลัยบูรพา
ต.แสนสุข อ.เมือง จ.ชลบุรี ๒๐๑๓๑

๑๗๔ พฤษภาคม ๒๕๕๘

เรื่อง ขอความอนุเคราะห์ขอข้อมูลเพื่อการวิจัย

เรียน ผู้อำนวยการสำนักทดสอบทางการศึกษา

สังกัดสำนักงานเขตฯ ๑. เค้าโครงวิทยานิพนธ์ฉบับย่อ จำนวน ๑ ชุด

๒. ผลการวิเคราะห์คุณภาพข้อสอบด้านความสามารถด้านภาษา ตัวน้ำคำนวน และ

ด้านเหตุผล ขั้นประถมศึกษาปีที่ ๓ ปีการศึกษา ๒๕๕๗ จำนวน ๑ ชุด

๓. แบบรายงานผลการพิจารณาการวิจัยธรรมการวิจัยในคน จำนวน ๑ ชุด

ด้วย นางสาวอรุณี แปลงกาย รหัสประจำตัว ๕๖๓๑๕๐๓ นิสิตหลักสูตรวิทยาศาสตร์
มหาบัณฑิต สาขาวิชาการวิจัยและสถิติทางวิทยาการปัญญา ได้รับอนุมัติให้ทำวิทยานิพนธ์เรื่อง^๑
“การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ระดับชั้นประถมศึกษาปีที่ ๓
ระหว่างวิธีการทดสอบอัตราร่วมในคลิกซิท วิชชิปเทสท์ และวิชเมเน่�템-แอนส์ເຊລ” ซึ่งอยู่ในความควบคุม
ดูแลของ ดร.ปิยะพิพัฒ ประดุจพร อาจารย์ที่ปรึกษาหลัก ในกรณี ผู้วิจัยมีความประสงค์ขอความอนุเคราะห์
ขอข้อสอบ NT (National Test) ระดับชั้นประถมศึกษาปีที่ ๓ พร้อมedly รหัสโรงเรียน เพศ และผล
การตอบข้อสอบ ในปีการศึกษา ๒๕๕๗ จำนวน ๓ ตัวนักภาษา ตัวนักการคิดคำนวน และตัวนักเหตุผล

จึงเรียนมาเพื่อโปรดพิจารณา วิทยาลัยวิทยาการวิจัยและวิทยาการปัญญา หวังเป็นอย่างยิ่ง^๒
ว่าคงจะได้รับความอนุเคราะห์จากท่านด้วยดี และขอขอบคุณมา ณ โอกาสสืบ

ขอแสดงความนับถือ

ธีระ คงวิช

(ผู้ช่วยศาสตราจารย์ ดร.สุชาดา กรเพชรปานี)

คณบดีวิทยาลัยวิทยาการวิจัยและวิทยาการปัญญา

วิทยาลัยวิทยาการวิจัยและวิทยาการปัญญา

โทร. ๐ ๓๔๑๐ ๒๐๗๗-๔

โทร/ โทรสาร ๐ ๓๔๓๘ ๓๔๕๔

<http://www.mcs.buu.ac.th>

ภาคผนวก ค
ตัวอย่างแสดงข้อมูลผลการตอบข้อสอบ NT
ปีการศึกษา 2556 ชั้นประถมศึกษาปีที่ 3

ตารางที่ ค-1 แสดงข้อมูลผลการตอบข้อสอบ NT ชั้นประถมศึกษาปีที่ 3 ต่างภาษา จำนวน 300 คน

#ID	Student	Gender	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	รวม
1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1	0	0	0	1	1	1	0	1	22			
2	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1	0	0	0	1	1	0	1	1	22			
3	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1	0	0	1	1	1	0	1	22			
4	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	0	1	1	1	0	1	18			
5	1	0	1	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	0	1	1	1	0	1	22			
6	1	0	0	1	0	1	0	0	1	0	0	1	1	1	1	1	1	1	1	1	1	0	0	1	0	1	1	0	0	9			
7	1	1	0	1	1	0	1	0	1	0	0	1	0	0	1	0	0	1	1	1	1	0	0	1	1	1	1	1	1	19			
8	1	1	0	1	0	0	1	0	1	0	1	0	0	1	0	1	0	1	0	1	0	0	0	0	0	0	0	0	0	9			
9	1	0	0	0	0	0	0	1	0	1	0	1	0	0	0	0	1	1	0	0	0	0	0	1	0	0	0	0	0	9			
10	1	0	1	1	1	0	1	0	1	1	0	1	1	1	0	1	1	1	1	1	0	0	1	0	0	1	1	1	1	19			
11	1	1	0	1	0	1	0	1	0	1	1	1	1	1	0	1	1	1	1	1	1	0	0	0	0	0	1	1	1	19			
12	1	1	0	0	0	0	1	0	0	1	0	0	0	0	0	0	1	0	0	1	1	0	0	0	1	1	0	1	0	8			
13	1	0	1	0	0	0	0	1	1	1	0	0	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	12			
14	1	1	0	1	1	1	1	1	1	1	0	1	1	1	1	1	1	1	1	1	1	0	1	1	1	1	0	0	0	20			
15	1	0	1	1	1	0	1	1	1	1	0	1	1	1	1	1	1	1	1	1	1	0	0	1	1	1	0	1	0	16			
16	2	1	1	1	1	0	1	1	1	1	0	1	1	1	1	1	1	1	1	1	1	0	0	1	1	1	1	0	20				
17	2	1	1	0	0	1	1	1	1	0	1	0	1	1	1	1	1	1	1	1	1	0	0	1	1	1	1	0	19				
18	2	0	1	1	1	1	0	1	0	1	0	1	1	1	1	1	1	1	1	1	1	0	0	1	1	1	1	0	1	16			
19	2	0	1	0	1	1	1	1	0	0	0	1	0	1	0	0	0	1	1	1	1	0	1	1	1	1	0	1	1	15			
20	2	0	1	1	1	1	0	0	0	1	0	1	1	1	1	0	0	1	1	1	1	0	0	1	1	1	1	0	1	17			
21	2	0	1	1	0	1	1	0	1	1	0	1	1	1	1	1	0	1	1	1	1	0	0	1	1	1	1	0	1	21			
22	2	0	1	0	1	0	0	1	0	1	0	0	1	1	1	1	1	1	1	1	1	0	0	1	1	1	1	0	0	9			
23	2	1	0	1	0	1	0	1	0	0	1	1	0	0	0	1	1	1	1	1	1	0	1	1	1	1	1	1	1	19			
24	2	1	0	1	1	0	0	1	0	1	0	1	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	9				
25	2	1	1	1	1	0	1	0	1	0	1	1	1	1	1	1	1	1	1	1	1	0	0	1	0	0	0	1	0	18			
26	2	0	1	1	0	1	1	1	1	0	1	1	1	1	1	1	0	0	1	1	1	1	0	0	1	0	1	1	1	20			
27	2	1	0	1	1	0	1	1	0	0	1	1	0	0	1	0	0	1	1	1	1	0	0	1	1	1	1	0	1	18			
28	2	1	0	1	1	0	1	0	1	0	0	1	1	0	0	1	0	1	1	1	1	0	0	1	1	1	1	1	1	18			
29	2	1	0	1	1	0	1	0	1	0	0	1	1	0	1	0	1	1	1	1	1	0	0	1	1	1	1	0	1	18			
30	2	0	0	0	1	0	1	1	0	1	1	0	1	1	1	1	1	1	1	1	1	0	0	1	0	1	1	1	1	19			
300	2	0	1	1	0	1	0	1	0	1	0	0	1	0	1	0	1	1	1	1	0	0	0	0	0	1	1	0	13				
รวม	161	135	185	210	154	211	175	195	211	125	233	143	209	143	211	130	137	197	127	224	197	162	69	78	127	120	140	131	179	151			

หมายเหตุ 1 คือ เพศชาย 2 คือ เพศหญิง

ตารางที่ ค-2 แสดงข้อมูลผลการทดสอบ NT ชั้นประถมศึกษาปีที่ 3 ดำเนินงานจำนวน 300 คน

ลำดับ	Student	Gender	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	รวม
1	1	1	1	1	1	0	0	1	1	1	1	1	1	1	1	0	0	0	1	1	1	1	1	1	1	1	1	1	1	1	1	23	
2	1	1	1	1	1	1	0	0	1	1	1	0	0	1	1	1	1	0	0	1	1	1	1	1	1	1	1	1	1	1	1	23	
3	1	1	1	1	1	1	0	0	1	1	1	1	0	0	1	1	1	0	0	1	1	1	1	1	1	1	1	1	1	1	1	23	
4	1	1	1	1	0	1	0	0	1	0	0	0	0	0	1	1	1	0	0	0	0	0	1	1	1	1	1	1	1	1	1	15	
5	1	0	0	0	0	0	0	0	0	1	0	0	0	0	1	1	0	1	0	0	0	0	1	1	1	1	0	1	1	1	1	11	
6	1	0	0	1	0	0	0	0	0	0	1	0	0	0	1	0	1	0	0	0	0	0	1	0	0	0	0	0	0	0	0	8	
7	1	0	0	1	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	6	
8	1	0	0	0	0	0	0	0	0	0	1	0	0	0	1	0	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0	10	
9	1	0	0	0	1	0	0	0	1	1	0	0	0	0	0	0	0	1	0	0	1	1	0	0	0	1	1	1	1	1	1	14	
10	1	0	0	0	1	1	0	1	0	1	1	0	0	1	1	1	1	1	1	1	0	0	0	1	1	1	1	1	1	1	1	19	
11	1	0	1	0	1	0	1	1	0	0	0	0	0	0	1	1	1	1	1	1	0	0	0	1	1	1	0	0	0	0	0	15	
12	1	0	0	0	0	0	1	1	1	0	1	0	0	0	0	0	1	0	0	1	0	0	0	0	0	0	1	1	0	0	11		
13	1	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	5		
14	1	0	1	0	0	0	0	0	1	0	0	0	0	0	0	1	0	0	0	1	0	0	0	1	1	1	1	1	1	1	1	10	
15	1	1	1	1	1	1	0	0	1	1	0	1	1	0	0	1	1	1	1	1	0	0	1	1	1	1	1	1	1	1	1	23	
16	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	0	0	0	0	0	1	0	0	0	0	0	0	0	0	8	
17	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	1	0	0	0	0	0	0	0	0	0	0	0	12	
18	2	1	0	0	0	0	1	0	0	0	0	1	0	0	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	5	
19	2	1	1	1	1	1	0	0	1	1	1	0	0	1	1	1	1	0	0	0	1	0	0	0	1	1	1	1	1	1	1	20	
20	2	1	0	0	0	0	0	0	0	1	1	1	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	9	
21	2	0	1	1	0	0	0	0	0	1	1	1	0	0	1	1	1	0	0	0	0	0	0	1	1	0	1	0	1	0	1	15	
22	2	0	0	1	0	0	0	0	0	0	0	1	0	0	1	1	0	0	1	0	0	0	0	1	0	0	0	0	0	0	1	12	
23	2	0	0	1	0	0	0	0	0	0	0	0	0	0	1	1	1	0	0	0	0	0	0	1	0	0	0	0	0	0	0	7	
24	2	0	0	0	0	0	0	0	0	0	0	1	0	0	1	1	0	1	0	0	0	0	0	1	0	0	0	1	0	0	0	15	
25	2	0	0	0	0	0	0	0	0	0	0	1	0	0	1	1	1	1	0	0	0	0	0	1	1	0	1	1	0	1	1	14	
26	2	0	0	0	0	0	0	0	0	0	0	1	1	0	1	0	0	1	1	0	0	0	0	1	0	0	0	1	0	0	0	12	
27	2	0	0	1	0	0	0	0	0	1	1	0	0	0	0	1	0	0	0	1	0	0	1	1	0	0	0	0	0	0	0	12	
28	2	0	0	0	1	0	1	1	0	1	0	0	0	1	1	0	1	1	0	0	1	0	0	1	0	0	0	0	0	0	0	15	
29	2	0	1	0	0	0	1	0	1	1	0	1	0	0	1	1	0	0	0	1	0	0	0	1	0	0	0	0	0	0	0	11	
30	2	0	1	0	1	0	0	1	0	1	1	0	0	0	1	1	0	0	0	1	1	0	0	1	0	0	1	0	1	0	1	17	
300	2	1	1	1	1	0	0	0	1	1	0	1	0	0	1	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	14	
รวม	83	93	107	130	110	119	151	157	153	157	119	97	128	114	114	142	180	119	118	162	111	97	113	108	104	150	131	141	111	111			

หมายเหตุ 1 คือ เพศชาย 2 คือ เพศหญิง

ตารางที่ ค-3 แสดงข้อมูลผลการตอบรับชื่อสอบ NT บนประมวลศึกษาปีที่ 3 ดำเนินการ จำนวน 300 คน

#	Student	Gender	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	Total
1	1	1	1	1	1	1	1	1	0	0	1	0	0	0	1	0	0	1	1	1	1	1	0	0	0	0	0	0	0	0	0	16	
2	1	1	1	1	1	1	1	1	1	0	0	1	0	0	0	1	0	0	1	1	1	1	1	0	0	0	0	0	0	0	0	1	21
3	1	1	1	1	1	1	1	1	1	0	0	1	0	0	1	0	0	1	1	1	1	1	1	0	0	0	0	0	0	0	0	1	21
4	1	1	1	1	1	1	1	1	1	0	0	1	0	0	1	0	0	1	1	0	1	1	0	1	1	0	0	0	0	0	0	1	20
5	1	1	1	1	1	1	1	1	1	0	0	1	0	0	1	0	0	1	1	0	1	1	0	1	1	0	0	0	0	0	0	1	20
6	1	0	0	1	1	0	0	1	0	0	1	1	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	9
7	1	0	1	0	1	1	1	1	1	0	0	1	0	0	0	0	0	0	1	0	0	0	1	1	1	0	1	0	1	0	1	15	
8	1	0	1	0	1	1	1	1	1	0	0	1	0	0	0	1	0	0	0	0	0	1	0	0	1	1	1	0	1	0	1	14	
9	1	1	1	0	0	0	0	0	0	0	0	1	1	0	0	0	0	0	0	0	0	1	1	0	0	0	0	0	0	0	1	10	
10	1	1	0	0	1	0	1	1	0	1	1	0	0	0	1	0	0	0	0	0	1	0	0	1	1	0	0	0	0	0	0	0	12
11	1	1	1	0	1	1	0	0	1	0	1	1	0	0	0	1	0	0	1	1	1	0	1	1	0	0	0	0	0	0	0	0	17
12	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	1	1	0	0	0	0	0	0	0	9
13	1	1	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	5	
14	1	1	1	0	1	0	1	1	0	1	1	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	15
15	1	1	0	0	1	1	1	1	0	1	1	0	0	0	1	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	17
16	2	0	0	1	1	1	1	1	0	0	1	1	0	0	0	1	0	0	1	1	1	1	1	0	1	1	0	1	0	1	0	1	13
17	2	0	0	1	0	1	1	1	0	0	1	1	0	0	0	1	1	0	0	0	1	0	0	1	0	0	1	0	1	0	1	0	13
18	2	1	1	0	1	1	0	1	1	1	1	0	0	0	1	1	0	0	0	1	0	0	1	1	0	0	1	0	1	0	1	0	19
19	2	0	1	1	1	1	1	1	0	0	1	1	0	0	0	1	0	0	0	1	0	0	0	1	0	0	0	0	0	0	0	1	14
20	2	1	0	1	1	0	0	1	0	0	1	1	0	0	0	1	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	11
21	2	1	0	1	1	1	1	1	0	0	1	1	1	0	0	0	1	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	14
22	2	0	0	1	0	1	0	0	1	0	1	1	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	9
23	2	0	1	0	1	1	1	1	1	0	1	1	0	0	0	1	0	0	0	0	1	0	0	1	0	0	0	0	0	0	0	0	15
24	2	0	1	0	1	1	1	1	0	1	0	0	0	1	0	0	0	0	0	0	0	0	0	1	1	0	1	0	1	0	1	14	
25	2	1	0	1	1	1	1	1	0	0	1	1	0	0	0	0	0	0	0	1	0	0	1	0	0	1	0	0	1	0	1	14	
26	2	1	0	1	1	1	1	1	0	1	1	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	13
27	2	0	1	0	1	1	1	1	0	0	1	0	0	0	0	1	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	1	14
28	2	0	1	0	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	1	1	0	0	1	0	1	1	13	
29	2	0	1	0	1	1	1	1	0	0	1	0	0	0	0	1	0	0	0	0	1	0	0	1	1	0	0	1	0	1	1	18	
30	2	0	1	1	0	0	1	1	0	0	1	1	0	0	1	1	1	0	0	0	1	0	0	1	1	0	0	1	0	1	1	0	14
300	2	0	1	1	1	1	0	0	0	0	1	1	0	0	0	1	1	0	0	0	1	0	0	1	1	0	0	1	0	1	0	0	14
รวม	166	171	169	211	174	183	170	140	152	118	167	166	147	106	132	99	86	141	134	105	164	157	146	170	141	156	98	95	99	120			

หมายเหตุ 1 คือ เพศชาย 2 คือ เพศหญิง

ตารางที่ ค-4 แสดงข้อมูลผลการตอบข้อสอบ NT บนประมาณศึกษาปี 3 ต้นภาษา จำนวน 1,000 คน

ลำดับ	เพศ	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	รวม
Student	Gender																															
1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	0	1	0	0	0	1	1	0	1	
2	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	1	1	0	1	22		
3	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	1	1	0	1	22		
4	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	0	1	0	1	1	0	1	22		
5	1	0	1	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	0	1	0	1	1	0	1	22		
6	1	0	0	1	0	1	0	0	0	1	0	0	1	1	0	0	0	0	1	0	0	0	1	1	0	0	0	0	0	9		
7	1	1	0	1	0	1	1	0	1	0	1	0	0	1	0	0	1	1	1	1	1	0	1	0	0	1	1	1	1	19		
8	1	1	0	1	1	0	0	1	0	1	0	1	0	1	0	1	0	1	0	1	0	0	0	0	0	0	0	0	0	9		
9	1	0	0	1	0	1	0	0	1	0	1	0	1	1	0	0	0	0	1	0	1	0	0	0	0	0	0	0	0	9		
10	1	0	1	1	0	1	1	1	1	0	1	1	1	1	1	0	1	1	1	1	1	0	0	0	0	1	1	1	1	19		
11	1	1	0	1	0	1	0	1	1	1	1	1	1	1	1	0	1	1	1	1	1	1	0	0	0	1	1	1	0	8		
12	1	1	0	0	0	0	1	0	0	0	1	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0			
13	1	0	1	0	0	0	1	1	1	0	1	1	1	1	1	0	1	1	1	1	1	0	0	0	1	0	0	0	0	12		
14	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	0	1	1	1	1	1	19		
15	1	0	1	1	0	0	1	1	1	1	0	1	1	1	1	0	1	1	1	1	1	1	0	0	0	1	1	1	1	19		
16	2	1	1	1	1	0	1	1	1	1	0	1	1	1	1	0	1	1	1	1	1	0	0	0	1	1	1	1	0	20		
17	2	1	1	0	1	1	1	0	1	1	1	0	1	1	1	1	1	1	1	1	1	0	0	0	1	1	1	1	0	19		
18	2	0	1	1	1	0	0	1	0	1	0	1	0	1	1	0	1	1	1	1	1	0	0	0	1	1	0	1	1	16		
19	2	0	1	0	1	1	1	1	0	0	1	0	0	1	0	0	0	0	1	1	1	0	0	1	1	0	1	1	1	15		
20	2	0	1	1	0	0	1	1	0	0	1	0	0	0	1	1	1	1	1	1	1	0	0	0	1	1	1	1	0	17		
21	2	0	1	1	0	1	1	1	0	1	1	0	1	0	1	1	1	1	1	1	1	0	0	1	0	1	1	1	1	21		
22	2	0	0	1	0	1	0	0	0	1	0	1	0	1	1	0	0	0	1	1	1	0	0	0	1	0	0	0	0	9		
23	2	1	0	1	1	0	1	1	0	1	0	1	0	1	0	1	1	1	1	1	1	0	0	1	1	1	1	1	1	19		
24	2	1	0	1	1	0	0	1	0	1	0	1	0	1	1	1	1	1	1	1	1	0	0	0	1	1	1	1	0	9		
25	2	1	1	1	1	0	1	0	1	0	1	1	1	1	1	0	1	1	1	1	1	0	0	1	1	1	1	1	0	18		
26	2	0	1	1	1	0	1	1	1	1	0	1	1	1	1	0	0	1	1	1	1	0	0	1	1	1	1	1	0	20		
27	2	1	0	1	1	0	1	1	0	0	1	1	0	0	1	0	1	1	1	1	1	0	1	0	1	1	1	1	1	18		
28	2	1	0	1	1	0	0	1	0	0	1	1	0	0	1	0	1	1	1	1	1	0	1	1	1	1	1	1	1	18		
29	2	1	0	1	1	0	1	1	0	0	1	1	0	0	1	1	1	1	1	1	1	0	0	1	1	1	1	1	1	19		
30	2	0	0	1	1	0	1	1	1	0	1	1	0	1	1	0	1	1	1	1	1	0	0	1	1	1	1	1	1	19		
1,000	2	0	0	0	1	1	0	1	1	1	1	1	1	1	1	0	0	0	0	0	1	0	1	1	1	1	1	1	1	12		
รวม		542	323	601	750	606	736	495	749	401	705	576	419	487	423	448	633	453	791	689	530	247	349	483	367	444	554	420	565			

หมายเหตุ 1 คือ เพศชาย 2 คือ เพศหญิง

หมายเหตุ 1 คือ เพศชาย 2 คือ เพศหญิง

ตารางที่ ค-5 แสดงข้อมูลผลการตอบทุกสอบ NT ชั้นประถมศึกษาปีที่ 3 ตามคำนวณ จำนวน 1,000 คน

#	Student	Gender	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	รวม
1	1	1	1	1	1	1	1	0	0	1	1	1	1	0	0	0	1	1	0	1	1	1	1	1	1	1	1	1	1	1	23		
2	1	1	1	1	1	1	1	0	0	1	1	1	0	0	0	0	0	1	1	0	1	1	1	1	1	1	1	1	1	1	23		
3	1	1	1	1	1	1	1	0	0	1	1	1	0	0	1	1	1	0	0	1	1	1	1	1	1	1	1	1	1	1	23		
4	1	1	1	1	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	15		
5	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	11		
6	1	0	0	1	0	0	0	0	0	0	0	0	0	0	0	1	0	0	1	0	0	0	0	0	0	0	0	0	0	0	8		
7	1	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	6		
8	1	0	0	0	0	0	0	0	0	0	1	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	10		
9	1	0	0	1	0	0	1	0	0	1	1	0	0	0	0	0	1	0	0	1	1	0	0	0	1	1	0	1	0	1	14		
10	1	0	0	0	1	0	1	0	1	1	1	1	0	0	1	1	1	1	1	1	0	0	0	1	1	0	1	1	0	1	19		
11	1	0	1	0	1	0	1	1	0	0	0	0	0	0	0	1	1	1	1	1	1	0	0	0	1	1	0	0	0	1	15		
12	1	0	0	0	0	0	1	1	1	0	1	0	0	0	0	0	1	0	0	1	0	0	0	0	0	1	0	0	0	0	11		
13	1	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	5		
14	1	0	1	0	0	0	0	0	0	0	0	1	0	0	0	1	0	0	0	0	0	1	1	0	0	0	1	0	0	1	10		
15	1	1	1	1	1	0	0	0	1	0	0	1	1	0	0	1	1	0	0	1	1	0	1	1	1	1	1	1	1	1	23		
16	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	0	0	1	0	0	0	1	0	0	0	0	0	0	8		
17	2	0	0	0	0	0	0	0	0	0	1	1	0	0	1	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0	12		
18	2	1	0	0	0	1	0	0	0	0	1	1	0	0	1	1	1	0	0	1	0	0	0	0	0	0	0	0	0	1	12		
19	2	1	1	1	1	1	0	0	1	1	1	0	0	0	1	1	1	0	0	1	1	0	0	0	1	1	0	1	0	1	20		
20	2	1	0	0	0	0	0	0	0	1	1	1	0	0	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	1	9		
21	2	0	1	0	0	0	0	0	0	1	1	1	0	0	1	1	1	0	0	0	0	0	0	0	0	1	0	1	0	1	15		
22	2	0	0	1	0	0	0	0	0	0	1	0	1	0	0	1	0	0	1	0	0	0	0	0	0	0	0	0	0	9			
23	2	0	0	1	0	0	0	0	0	0	0	1	0	0	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	7			
24	2	0	0	0	0	0	0	0	0	0	0	1	1	0	0	0	1	0	1	0	0	0	0	0	0	0	0	0	0	11			
25	2	0	0	0	0	1	0	0	0	0	1	1	0	0	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	1	15		
26	2	0	0	0	0	0	0	0	0	0	1	1	1	0	0	1	1	1	0	1	0	0	0	0	0	0	0	0	0	1	14		
27	2	0	0	1	0	0	0	0	1	1	0	0	0	0	1	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	12		
28	2	0	0	0	1	0	1	0	1	1	0	1	1	0	0	1	1	0	0	1	0	0	0	0	0	0	0	0	0	0	15		
29	2	0	0	1	0	0	0	0	0	1	1	0	1	0	0	0	1	0	0	1	0	0	0	0	0	0	0	0	0	0	11		
30	2	0	1	0	1	0	0	0	1	1	0	0	0	0	1	1	0	0	0	1	0	0	0	1	0	0	1	0	0	1	17		
1,000	2	1	1	1	1	0	0	0	1	1	0	1	0	0	1	0	1	1	1	0	0	0	0	0	0	0	0	0	0	0	14		
รวม	263	383	344	496	350	276	595	556	477	350	517	302	385	482	425	451	519	364	240	192	406	432	386	517	366	433	386	14					

ตารางที่ ค-7 แสดงข้อมูลผลการตอบชี้ตอบ NT ชั้นประถมศึกษาปีที่ 3 ต่างภาษา จำนวน 2,000 คน

#ID	Student	Gender	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	รวม
1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1	0	0	0	1	1	0	1	1	0	1	22		
2	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1	0	0	0	1	1	0	1	1	0	1	22		
3	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1	0	0	0	1	1	0	1	1	0	1	22		
4	1	0	1	1	1	1	1	1	0	1	1	1	1	1	1	1	1	0	1	1	0	0	1	0	1	1	0	1	1	0	1	22	
5	1	0	1	1	1	1	1	0	1	1	1	1	1	1	1	1	1	0	1	1	0	0	1	0	1	1	0	1	1	0	1	22	
6	1	0	0	1	0	1	0	0	0	0	1	0	0	0	0	0	0	0	1	0	0	0	1	0	0	0	0	0	0	0	0	9	
7	1	1	0	1	1	0	1	1	0	0	0	1	1	0	0	0	1	0	1	1	1	0	0	1	1	1	1	1	1	1	1	19	
8	1	1	0	1	1	0	1	0	0	1	0	0	1	0	0	1	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	9	
9	1	0	0	0	1	0	0	1	0	1	0	1	1	0	0	0	0	0	1	0	0	0	0	0	0	1	0	0	0	0	9		
10	1	0	1	1	0	1	1	1	0	1	1	1	1	0	0	1	0	1	1	1	1	0	0	1	0	1	1	1	1	1	1	19	
11	1	1	0	1	0	1	0	1	0	1	1	1	1	0	1	0	1	1	1	1	1	0	0	0	0	1	1	1	1	1	19		
12	1	1	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	1	0	1	0	1	0	1	8		
13	1	0	1	0	0	0	1	1	0	0	1	1	0	1	0	1	1	0	0	0	0	0	1	0	0	0	0	0	0	0	0	12	
14	1	1	0	1	1	1	1	1	0	1	1	0	0	0	1	1	1	1	0	0	1	1	0	0	1	1	0	0	0	0	0	20	
15	1	0	1	1	0	1	1	0	0	1	1	0	1	0	0	1	0	1	1	1	0	0	1	1	0	1	0	1	0	1	16		
16	2	1	1	1	0	1	1	1	0	1	1	0	1	0	0	1	1	0	1	1	0	0	1	1	1	1	0	1	1	0	20		
17	2	1	0	1	0	1	1	1	0	1	1	0	1	1	0	1	1	0	0	0	0	1	0	1	1	0	1	1	0	1	19		
18	2	0	1	1	1	0	0	1	0	1	0	1	0	0	1	0	1	1	0	0	0	1	1	0	1	1	0	1	1	0	9		
19	2	0	1	0	1	1	1	1	0	0	1	0	0	1	0	0	0	1	1	1	0	1	1	0	1	1	0	1	1	1	17		
20	2	0	1	1	1	0	0	0	1	1	0	0	1	0	0	1	1	1	1	1	0	0	1	1	1	1	0	1	1	1	17		
21	2	0	1	1	0	1	1	1	0	1	1	0	1	1	0	1	1	1	1	1	1	0	1	1	1	1	1	1	1	1	21		
22	2	0	0	1	0	1	0	0	0	1	0	1	1	0	0	0	0	1	0	0	1	1	0	0	0	0	0	0	0	0	9		
23	2	1	0	1	1	0	1	0	0	0	1	1	0	0	1	1	1	1	1	1	1	0	1	1	1	1	1	1	1	1	19		
24	2	1	0	1	1	0	0	1	0	1	0	1	0	0	1	0	0	1	0	0	1	0	0	0	0	0	0	0	0	0	9		
25	2	1	1	1	1	0	1	0	1	1	1	1	0	1	1	1	0	1	1	1	1	0	0	1	0	0	1	0	0	1	18		
26	2	0	1	1	1	0	1	1	0	1	1	0	0	1	1	1	1	1	1	1	1	0	1	0	1	1	1	1	1	1	20		
27	2	2	1	0	1	1	0	1	0	0	1	1	0	0	1	1	1	1	1	1	1	0	0	1	1	1	1	1	1	1	18		
28	2	2	1	0	1	1	0	1	1	0	0	1	1	0	0	1	0	1	1	1	1	0	0	1	1	1	1	1	1	1	18		
29	2	1	0	1	1	0	1	1	0	0	1	1	0	0	1	1	1	1	1	1	1	0	1	0	1	1	1	1	1	1	18		
30	2	0	0	1	1	0	1	1	0	1	1	0	1	1	1	1	1	1	1	1	1	0	0	1	0	1	1	1	1	1	19		
2,000	2	1	0	0	0	0	1	0	0	0	0	1	0	0	0	1	0	0	0	0	1	1	0	0	1	0	1	0	0	0	10		
รวม	1070	642	1200	1511	1147	1467	1503	949	1577	822	1401	1159	845	937	860	948	1380	903	1580	1346	1004	527	748	966	674	842	1133	695	1145				

หมายเหตุ 1 คือ เพศชาย 2 คือ เพศหญิง

ภาคผนวก ง
ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ
ด้วยวิธี IRT-LR วิธี SIBTEST และวิธี MH

ตารางที่ ง-10 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3
**ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี IRT-LR ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง
 ขนาดใหญ่ (300 คน)**

ข้อที่	ผลการตรวจสอบ DIF ด้วยวิธี IRT-LR		
	กลุ่มตัวอย่าง (300 คน)		
	ด้านภาษา	ด้านคำนวน	ด้านเหตุผล
1	0.22(NO-DIF)	0.91(NO-DIF)	0.10(NO-DIF)
2	0.01(DIF)	0.84(NO-DIF)	0.89(NO-DIF)
3	0.30(NO-DIF)	0.79(NO-DIF)	0.01(DIF)
4	0.33(NO-DIF)	0.76(NO-DIF)	0.20(NO-DIF)
5	0.73(NO-DIF)	0.89(NO-DIF)	0.12(NO-DIF)
6	0.10(NO-DIF)	0.76(NO-DIF)	0.56(NO-DIF)
7	0.20(NO-DIF)	0.27(NO-DIF)	0.87(NO-DIF)
8	0.53(NO-DIF)	0.20(NO-DIF)	0.95(NO-DIF)
9	0.76(NO-DIF)	0.72(NO-DIF)	0.28(NO-DIF)
10	0.35(NO-DIF)	0.61(NO-DIF)	0.50(NO-DIF)
11	0.42(NO-DIF)	0.02(DIF)	0.17(NO-DIF)
12	0.34(NO-DIF)	0.04(DIF)	0.05(DIF)
13	0.54(NO-DIF)	0.92(NO-DIF)	0.11(NO-DIF)
14	0.88(NO-DIF)	0.79(NO-DIF)	0.05(DIF)
15	0.18(NO-DIF)	0.39(NO-DIF)	0.72(NO-DIF)
16	0.55(NO-DIF)	0.25(NO-DIF)	0.26(NO-DIF)
17	0.03(DIF)	0.30(NO-DIF)	0.72(NO-DIF)
18	0.21(NO-DIF)	0.44(NO-DIF)	0.06(NO-DIF)
19	0.07(NO-DIF)	0.30(NO-DIF)	0.56(NO-DIF)
20	0.77(NO-DIF)	0.13(NO-DIF)	0.08(NO-DIF)
21	0.76(NO-DIF)	0.00(DIF)	0.06(NO-DIF)
22	0.40(NO-DIF)	0.92(NO-DIF)	0.88(NO-DIF)
23	0.72(NO-DIF)	0.39(NO-DIF)	0.70(NO-DIF)
24	0.19(NO-DIF)	0.58(NO-DIF)	0.86(NO-DIF)
25	0.96(NO-DIF)	0.59(NO-DIF)	0.28(NO-DIF)
26	0.75(NO-DIF)	0.30(NO-DIF)	0.30(NO-DIF)
27	0.63(NO-DIF)	0.87(NO-DIF)	0.94(NO-DIF)
28	0.13(NO-DIF)	0.28(NO-DIF)	0.18(NO-DIF)
29	0.64(NO-DIF)	0.69(NO-DIF)	0.36(NO-DIF)
30	0.02(DIF)	0.28(NO-DIF)	0.13(NO-DIF)
จำนวนข้อที่พน DIF		(3 ข้อ) 10%	(3 ข้อ) 10%

ตารางที่ ง-11 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3
**ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี IRT-LR ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง
 ขนาดใหญ่ (1,000 คน)**

ข้อที่	ผลการตรวจสอบ DIF ด้วยวิธี IRT-LR		
	กลุ่มตัวอย่าง (1,000 คน)		
	ด้านภาษา	ด้านคำนวน	ด้านเหตุผล
1	0.97(NO-DIF)	0.22(NO-DIF)	0.00(DIF)
2	0.12(NO-DIF)	0.99(NO-DIF)	0.63(NO-DIF)
3	0.76(NO-DIF)	0.44(NO-DIF)	0.01(DIF)
4	0.42(NO-DIF)	0.56(NO-DIF)	0.90(NO-DIF)
5	0.00(DIF)	0.34(NO-DIF)	0.59(NO-DIF)
6	0.55(NO-DIF)	0.58(NO-DIF)	0.45(NO-DIF)
7	0.88(NO-DIF)	0.41(NO-DIF)	0.11(NO-DIF)
8	0.99(NO-DIF)	0.33(NO-DIF)	0.83(NO-DIF)
9	0.78(NO-DIF)	0.56(NO-DIF)	0.07(NO-DIF)
10	0.28(NO-DIF)	0.03(DIF)	0.29(NO-DIF)
11	0.78(NO-DIF)	0.96(NO-DIF)	0.00(DIF)
12	0.23(NO-DIF)	0.03(DIF)	0.00(DIF)
13	0.42(NO-DIF)	0.20(NO-DIF)	0.11(NO-DIF)
14	0.75(NO-DIF)	0.57(NO-DIF)	0.09(NO-DIF)
15	0.01(DIF)	0.96(NO-DIF)	0.74(NO-DIF)
16	0.37(NO-DIF)	0.05(DIF)	0.71(NO-DIF)
17	0.24(NO-DIF)	0.10(NO-DIF)	0.89(NO-DIF)
18	0.01(DIF)	0.72(NO-DIF)	0.23(NO-DIF)
19	0.50(NO-DIF)	0.04(DIF)	0.15(NO-DIF)
20	0.23(NO-DIF)	0.36(NO-DIF)	0.15(NO-DIF)
21	0.14(NO-DIF)	0.48(NO-DIF)	0.15(NO-DIF)
22	0.24(NO-DIF)	0.22(NO-DIF)	0.93(NO-DIF)
23	0.70(NO-DIF)	0.36(NO-DIF)	0.30(NO-DIF)
24	0.44(NO-DIF)	0.06(NO-DIF)	0.05(DIF)
25	0.97(NO-DIF)	0.48(NO-DIF)	0.20(NO-DIF)
26	0.60(NO-DIF)	0.27(NO-DIF)	0.76(NO-DIF)
27	0.99(NO-DIF)	0.97(NO-DIF)	0.85(NO-DIF)
28	0.79(NO-DIF)	0.57(NO-DIF)	0.98(NO-DIF)
29	0.07(NO-DIF)	0.90(NO-DIF)	0.86(NO-DIF)
30	0.01(DIF)	0.08(NO-DIF)	0.08(NO-DIF)
จำนวนข้อที่พิบ DIF	(4 ข้อ) 13%	(4 ข้อ) 13%	(5 ข้อ) 17%

ตารางที่ ง-12 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3
ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี IRT-LR ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง
ขนาดใหญ่ (2,000 คน)

ข้อที่	ผลการตรวจสอบ DIF ด้วยวิธี IRT-LR		
	กลุ่มตัวอย่าง (2,000 คน)		
	ด้านภาษา	ด้านคำนวน	ด้านเหตุผล
1	0.11(NO-DIF)	0.22(NO-DIF)	0.00(DIF)
2	0.07(NO-DIF)	0.82(NO-DIF)	0.15(NO-DIF)
3	0.21(NO-DIF)	0.18(NO-DIF)	0.01(DIF)
4	0.39(NO-DIF)	0.00(DIF)	0.39(NO-DIF)
5	0.00(DIF)	0.14(NO-DIF)	0.10(NO-DIF)
6	0.45(NO-DIF)	0.91(NO-DIF)	0.20(NO-DIF)
7	0.49(NO-DIF)	0.20(NO-DIF)	0.11(NO-DIF)
8	0.70(NO-DIF)	0.10(NO-DIF)	0.15(NO-DIF)
9	0.37(NO-DIF)	0.26(NO-DIF)	0.00(DIF)
10	0.95(NO-DIF)	0.06(NO-DIF)	0.29(NO-DIF)
11	0.68(NO-DIF)	0.98(NO-DIF)	0.00(DIF)
12	0.39(NO-DIF)	0.02(DIF)	0.00(DIF)
13	0.59(NO-DIF)	0.15(NO-DIF)	0.19(NO-DIF)
14	0.68(NO-DIF)	0.12(NO-DIF)	0.00(DIF)
15	0.06(NO-DIF)	0.54(NO-DIF)	0.79(NO-DIF)
16	0.20(NO-DIF)	0.04(DIF)	0.11(NO-DIF)
17	0.22(NO-DIF)	0.06(NO-DIF)	0.13(NO-DIF)
18	0.69(NO-DIF)	0.95(NO-DIF)	0.92(NO-DIF)
19	0.14(NO-DIF)	0.06(NO-DIF)	0.00(DIF)
20	0.24(NO-DIF)	0.52(NO-DIF)	0.24(NO-DIF)
21	0.79(NO-DIF)	0.06(NO-DIF)	0.01(DIF)
22	0.41(NO-DIF)	0.58(NO-DIF)	0.31(NO-DIF)
23	0.79(NO-DIF)	0.36(NO-DIF)	0.67(NO-DIF)
24	0.07(NO-DIF)	0.08(NO-DIF)	0.15(NO-DIF)
25	0.92(NO-DIF)	0.79(NO-DIF)	0.08(NO-DIF)
26	0.17(NO-DIF)	0.37(NO-DIF)	0.12(NO-DIF)
27	0.17(NO-DIF)	0.29(NO-DIF)	0.22(NO-DIF)
28	0.45(NO-DIF)	0.055(NO-DIF)	0.52(NO-DIF)
29	0.03(DIF)	0.68(NO-DIF)	0.03(DIF)
30	0.01(DIF)	0.37(NO-DIF)	0.00(DIF)
จำนวนข้อที่พ布 DIF	(3 ข้อ) 10%	(3 ข้อ) 10%	(10 ข้อ) 33%

ตารางที่ ง-13 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3
ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี SIBTEST ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง
ขนาดใหญ่ (300 คน)

ข้อที่	ผลการตรวจสอบ DIF ด้วยวิธี SIBTEST		
	กลุ่มตัวอย่าง (300 คน)		
	ด้านภาษา	ด้านคำนวน	ด้านเหตุผล
1	0.11(NO-DIF)	0.81(NO-DIF)	0.58(NO-DIF)
2	0.09(NO-DIF)	0.77(NO-DIF)	0.61(NO-DIF)
3	0.42(NO-DIF)	0.53(NO-DIF)	0.07(NO-DIF)
4	0.42(NO-DIF)	0.49(NO-DIF)	0.19(NO-DIF)
5	0.94(NO-DIF)	0.90(NO-DIF)	0.35(NO-DIF)
6	0.54(NO-DIF)	0.53(NO-DIF)	0.54(NO-DIF)
7	0.67(NO-DIF)	0.60(NO-DIF)	0.49(NO-DIF)
8	0.45(NO-DIF)	0.08(NO-DIF)	0.37(NO-DIF)
9	0.50(NO-DIF)	0.60(NO-DIF)	0.69(NO-DIF)
10	0.96(NO-DIF)	0.17(NO-DIF)	0.61(NO-DIF)
11	0.26(NO-DIF)	0.67(NO-DIF)	0.14(NO-DIF)
12	0.16(NO-DIF)	0.03(DIF)	0.10(NO-DIF)
13	0.86(NO-DIF)	0.18(NO-DIF)	0.14(NO-DIF)
14	0.93(NO-DIF)	0.97(NO-DIF)	0.09(NO-DIF)
15	0.71(NO-DIF)	0.45(NO-DIF)	0.90(NO-DIF)
16	0.64(NO-DIF)	0.10(NO-DIF)	0.19(NO-DIF)
17	0.10(NO-DIF)	0.88(NO-DIF)	0.95(NO-DIF)
18	0.52(NO-DIF)	0.79(NO-DIF)	0.05(DIF)
19	0.01(DIF)	0.36(NO-DIF)	0.11(NO-DIF)
20	0.94(NO-DIF)	0.77(NO-DIF)	0.05(DIF)
21	0.15(NO-DIF)	0.02(DIF)	0.07(NO-DIF)
22	0.77(NO-DIF)	0.86(NO-DIF)	0.43(NO-DIF)
23	0.57(NO-DIF)	0.36(NO-DIF)	0.75(NO-DIF)
24	0.62(NO-DIF)	0.80(NO-DIF)	0.46(NO-DIF)
25	0.35(NO-DIF)	0.30(NO-DIF)	0.30(NO-DIF)
26	0.46(NO-DIF)	0.04(DIF)	0.77(NO-DIF)
27	0.42(NO-DIF)	0.77(NO-DIF)	0.92(NO-DIF)
28	0.10(NO-DIF)	0.13(NO-DIF)	0.32(NO-DIF)
29	0.36(NO-DIF)	0.61(NO-DIF)	0.20(NO-DIF)
30	0.00(DIF)	0.16(NO-DIF)	0.18(NO-DIF)
จำนวนข้อที่พบ DIF	(2 ข้อ) 7%	(3 ข้อ) 10%	(2 ข้อ) 7%

ตารางที่ ง-14 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3
ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี SIBTEST ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง
ขนาดใหญ่ (1,000 คน)

ข้อที่	ผลการตรวจสอบ DIF ด้วยวิธี SIBTEST		
	กลุ่มตัวอย่าง (1,000 คน)		
	ด้านภาษา	ด้านคำนวน	ด้านเหตุผล
1	0.88(NO-DIF)	0.71(NO-DIF)	0.00(DIF)
2	0.85(NO-DIF)	0.76(NO-DIF)	0.31(NO-DIF)
3	0.53(NO-DIF)	0.27(NO-DIF)	0.72(NO-DIF)
4	0.26(NO-DIF)	0.28(NO-DIF)	0.71(NO-DIF)
5	0.00(DIF)	0.47(NO-DIF)	0.72(NO-DIF)
6	0.46(NO-DIF)	0.74(NO-DIF)	0.56(NO-DIF)
7	0.88(NO-DIF)	0.91(NO-DIF)	0.21(NO-DIF)
8	0.99(NO-DIF)	0.49(NO-DIF)	0.83(NO-DIF)
9	0.95(NO-DIF)	0.19(NO-DIF)	0.02(DIF)
10	0.43(NO-DIF)	0.08(NO-DIF)	0.76(NO-DIF)
11	0.90(NO-DIF)	0.93(NO-DIF)	0.00(DIF)
12	0.09(NO-DIF)	0.33(NO-DIF)	0.00(DIF)
13	0.27(NO-DIF)	0.12(NO-DIF)	0.24(NO-DIF)
14	0.61(NO-DIF)	0.50(NO-DIF)	0.59(NO-DIF)
15	0.15(NO-DIF)	0.87(NO-DIF)	0.86(NO-DIF)
16	0.07(NO-DIF)	0.07(NO-DIF)	0.43(NO-DIF)
17	0.11(NO-DIF)	0.13(NO-DIF)	0.70(NO-DIF)
18	0.30(NO-DIF)	0.93(NO-DIF)	0.30(NO-DIF)
19	0.48(NO-DIF)	0.04(DIF)	0.01(DIF)
20	0.25(NO-DIF)	0.10(NO-DIF)	0.05(DIF)
21	0.33(NO-DIF)	0.71(NO-DIF)	0.17(NO-DIF)
22	0.09(NO-DIF)	0.12(NO-DIF)	0.52(NO-DIF)
23	0.62(NO-DIF)	0.76(NO-DIF)	0.45(NO-DIF)
24	0.23(NO-DIF)	0.65(NO-DIF)	0.58(NO-DIF)
25	0.80(NO-DIF)	0.30(NO-DIF)	0.06(NO-DIF)
26	0.29(NO-DIF)	0.28(NO-DIF)	0.25(NO-DIF)
27	0.91(NO-DIF)	0.52(NO-DIF)	0.53(NO-DIF)
28	0.98(NO-DIF)	0.24(NO-DIF)	0.93(NO-DIF)
29	0.07(NO-DIF)	0.22(NO-DIF)	0.88(NO-DIF)
30	0.00(DIF)	0.64(NO-DIF)	0.08(NO-DIF)
จำนวนข้อที่พบร DIF	(2 ข้อ) 7%	(1 ข้อ) 3%	(6 ข้อ) 20%

ตารางที่ ง-15 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3
ด้านภาษา ด้านคำนวณ ด้านเหตุผล ด้วยวิธี SIBTEST ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง
ขนาดใหญ่ (2,000 คน)

ข้อที่	ผลการตรวจสอบ DIF ด้วยวิธี SIBTEST		
	กลุ่มตัวอย่าง (2,000 คน)		
	ด้านภาษา	ด้านคำนวณ	ด้านเหตุผล
1	0.14(NO-DIF)	0.35(NO-DIF)	0.00(DIF)
2	0.35(NO-DIF)	0.87(NO-DIF)	0.21(NO-DIF)
3	0.10(NO-DIF)	0.61(NO-DIF)	0.26(NO-DIF)
4	0.63(NO-DIF)	0.04(DIF)	0.25(NO-DIF)
5	0.00(DIF)	0.45(NO-DIF)	0.08(NO-DIF)
6	0.28(NO-DIF)	0.84(NO-DIF)	0.34(NO-DIF)
7	0.68(NO-DIF)	0.23(NO-DIF)	0.14(NO-DIF)
8	0.46(NO-DIF)	0.24(NO-DIF)	0.45(NO-DIF)
9	0.10(NO-DIF)	0.00(DIF)	0.00(DIF)
10	0.78(NO-DIF)	0.42(NO-DIF)	0.29(NO-DIF)
11	0.86(NO-DIF)	0.85(NO-DIF)	0.00(DIF)
12	0.13(NO-DIF)	0.33(NO-DIF)	0.00(DIF)
13	0.87(NO-DIF)	0.02(DIF)	0.08(NO-DIF)
14	0.46(NO-DIF)	0.07(NO-DIF)	0.00(DIF)
15	0.14(NO-DIF)	0.26(NO-DIF)	0.41(NO-DIF)
16	0.07(NO-DIF)	0.04(DIF)	0.38(NO-DIF)
17	0.31(NO-DIF)	0.37(NO-DIF)	0.34(NO-DIF)
18	0.98(NO-DIF)	0.75(NO-DIF)	0.54(NO-DIF)
19	0.47(NO-DIF)	0.00(DIF)	0.00(DIF)
20	0.29(NO-DIF)	0.15(NO-DIF)	0.05(DIF)
21	0.64(NO-DIF)	0.41(NO-DIF)	0.02(DIF)
22	0.17(NO-DIF)	0.67(NO-DIF)	0.67(NO-DIF)
23	0.89(NO-DIF)	0.86(NO-DIF)	0.26(NO-DIF)
24	0.03(DIF)	0.96(NO-DIF)	0.74(NO-DIF)
25	0.58(NO-DIF)	0.78(NO-DIF)	0.02(DIF)
26	0.33(NO-DIF)	0.43(NO-DIF)	0.03(DIF)
27	0.38(NO-DIF)	0.41(NO-DIF)	0.29(NO-DIF)
28	0.51(NO-DIF)	0.02(DIF)	0.82(NO-DIF)
29	0.02(DIF)	0.95(NO-DIF)	0.02(DIF)
30	0.00(DIF)	0.19(NO-DIF)	0.00(DIF)
จำนวนข้อที่พิบ DIF	(4 ข้อ) 13%	(6 ข้อ) 20%	(12 ข้อ) 40%

ตารางที่ ง-16 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3
ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี MH ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง
ขนาดใหญ่ (300 คน)

ข้อที่	ผลการตรวจสอบ DIF ด้วยวิธี MH		
	กลุ่มตัวอย่าง (300 คน)		
	ด้านภาษา	ด้านคำนวน	ด้านเหตุผล
1	0.42(NO-DIF)	0.69(NO-DIF)	0.35(NO-DIF)
2	0.00(DIF)	0.53(NO-DIF)	0.41(NO-DIF)
3	0.02(DIF)	0.72(NO-DIF)	0.08(NO-DIF)
4	0.02(DIF)	0.82(NO-DIF)	0.01(DIF)
5	0.49(NO-DIF)	1.00(NO-DIF)	0.06(NO-DIF)
6	0.71(NO-DIF)	0.90(NO-DIF)	0.19(NO-DIF)
7	0.06(NO-DIF)	0.15(NO-DIF)	0.24(NO-DIF)
8	0.73(NO-DIF)	0.13(NO-DIF)	0.49(NO-DIF)
9	0.68(NO-DIF)	0.56(NO-DIF)	0.82(NO-DIF)
10	0.41(NO-DIF)	0.56(NO-DIF)	0.16(NO-DIF)
11	0.10(NO-DIF)	0.91(NO-DIF)	0.02(DIF)
12	0.91(NO-DIF)	0.04(DIF)	0.01(DIF)
13	0.56(NO-DIF)	0.82(NO-DIF)	0.02(DIF)
14	0.30(NO-DIF)	0.81(NO-DIF)	0.09(NO-DIF)
15	0.64(NO-DIF)	0.48(NO-DIF)	0.82(NO-DIF)
16	0.08(NO-DIF)	0.17(NO-DIF)	0.18(NO-DIF)
17	0.00(DIF)	0.35(NO-DIF)	1.00(NO-DIF)
18	0.73(NO-DIF)	0.41(NO-DIF)	0.20(NO-DIF)
19	0.00(DIF)	0.16(NO-DIF)	0.10(NO-DIF)
20	0.18(NO-DIF)	0.11(NO-DIF)	0.07(NO-DIF)
21	1.00(NO-DIF)	0.00(DIF)	0.64(NO-DIF)
22	0.89(NO-DIF)	0.90(NO-DIF)	0.30(NO-DIF)
23	0.60(NO-DIF)	0.20(NO-DIF)	0.36(NO-DIF)
24	0.60(NO-DIF)	0.55(NO-DIF)	0.16(NO-DIF)
25	1.00(NO-DIF)	0.47(NO-DIF)	0.30(NO-DIF)
26	0.25(NO-DIF)	0.15(NO-DIF)	0.02(DIF)
27	0.91(NO-DIF)	0.82(NO-DIF)	0.62(NO-DIF)
28	0.41(NO-DIF)	0.13(NO-DIF)	0.06(NO-DIF)
29	0.73(NO-DIF)	0.91(NO-DIF)	0.07(NO-DIF)
30	0.13(NO-DIF)	0.12(NO-DIF)	0.48(NO-DIF)
จำนวนข้อที่พบ DIF	(5 ข้อ) 17%	(2 ข้อ) 7%	(5 ข้อ) 17%

ตารางที่ ง-17 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3
ด้านภาษา ด้านคำนวน ด้านเหตุผล ด้วยวิธี MH ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง
ขนาดใหญ่ (1,000 คน)

ข้อที่	ผลการตรวจสอบ DIF ด้วยวิธี MH		
	กลุ่มตัวอย่าง (1,000 คน)		
	ด้านภาษา	ด้านคำนวน	ด้านเหตุผล
1	0.25(NO-DIF)	0.94(NO-DIF)	0.30(NO-DIF)
2	0.74(NO-DIF)	0.56(NO-DIF)	0.15(NO-DIF)
3	0.03(DIF)	0.29(NO-DIF)	0.26(NO-DIF)
4	0.02(DIF)	0.70(NO-DIF)	0.01(DIF)
5	0.24(NO-DIF)	0.36(NO-DIF)	0.04(DIF)
6	0.57(NO-DIF)	0.28(NO-DIF)	0.03(DIF)
7	0.01(DIF)	0.57(NO-DIF)	0.41(NO-DIF)
8	0.15(NO-DIF)	0.27(NO-DIF)	0.16(NO-DIF)
9	0.03(DIF)	0.45(NO-DIF)	0.95(NO-DIF)
10	0.11(NO-DIF)	0.13(NO-DIF)	0.00(DIF)
11	0.24(NO-DIF)	0.66(NO-DIF)	0.00(DIF)
12	0.80(NO-DIF)	0.23(NO-DIF)	0.00(DIF)
13	0.05(DIF)	0.04(DIF)	0.01(DIF)
14	0.23(NO-DIF)	0.89(NO-DIF)	0.04(DIF)
15	0.06(NO-DIF)	0.75(NO-DIF)	0.48(NO-DIF)
16	0.01(DIF)	0.01(DIF)	0.60(NO-DIF)
17	0.00(DIF)	0.02(DIF)	0.46(NO-DIF)
18	0.01(DIF)	0.23(NO-DIF)	0.00(DIF)
19	0.00(DIF)	0.01(DIF)	0.00(DIF)
20	0.31(NO-DIF)	0.18(NO-DIF)	0.06(NO-DIF)
21	1.38(NO-DIF)	0.60(NO-DIF)	0.85(NO-DIF)
22	0.61(NO-DIF)	0.14(NO-DIF)	0.01(DIF)
23	0.74(NO-DIF)	0.42(NO-DIF)	0.00(DIF)
24	0.00(DIF)	0.61(NO-DIF)	0.00(DIF)
25	0.74(NO-DIF)	0.61(NO-DIF)	0.44(NO-DIF)
26	0.11(NO-DIF)	0.36(NO-DIF)	0.00(DIF)
27	0.13(NO-DIF)	0.66(NO-DIF)	0.08(NO-DIF)
28	0.08(NO-DIF)	0.29(NO-DIF)	0.04(DIF)
29	0.00(DIF)	0.48(NO-DIF)	0.34(NO-DIF)
30	0.34(NO-DIF)	0.36(NO-DIF)	0.90(NO-DIF)
จำนวนข้อที่พบ DIF	(11 ข้อ) 37%	(4ข้อ) 13%	(15 ข้อ) 50%

ตารางที่ ง-18 ผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ชั้นประถมศึกษาปีที่ 3
ด้านภาษา ด้านคำนวณ ด้านเหตุผล ด้วยวิธี MH ภายใต้เงื่อนไขขนาดกลุ่มตัวอย่าง
ขนาดใหญ่ (2,000 คน)

ข้อที่	ผลการตรวจสอบ DIF ด้วยวิธี MH		
	กลุ่มตัวอย่าง (2,000 คน)		
	ด้านภาษา	ด้านคำนวณ	ด้านเหตุผล
1	0.59(NO-DIF)	0.64(NO-DIF)	1.00(NO-DIF)
2	0.92(NO-DIF)	0.14(NO-DIF)	0.00(DIF)
3	0.00(DIF)	0.32(NO-DIF)	0.01(DIF)
4	0.08(DIF)	0.45(NO-DIF)	0.00(DIF)
5	0.01(DIF)	0.21(NO-DIF)	0.00(DIF)
6	0.11(NO-DIF)	0.09(NO-DIF)	0.00(DIF)
7	0.00(DIF)	0.85(NO-DIF)	0.01(DIF)
8	0.09(NO-DIF)	0.00(DIF)	0.01(DIF)
9	0.00(DIF)	0.00(DIF)	0.93(NO-DIF)
10	0.32(NO-DIF)	0.65(NO-DIF)	0.00(DIF)
11	0.01(DIF)	0.08(NO-DIF)	0.00(DIF)
12	0.03(DIF)	0.24(NO-DIF)	0.00(DIF)
13	0.14NO-DIF	0.00(DIF)	0.00(DIF)
14	0.02(DIF)	0.82(NO-DIF)	0.00(DIF)
15	0.01(DIF)	0.06(NO-DIF)	0.00(DIF)
16	0.00(DIF)	0.00(DIF)	0.20(NO-DIF)
17	0.00(DIF)	0.00(DIF)	0.01(DIF)
18	0.01(DIF)	0.05(NO-DIF)	0.00(DIF)
19	0.00(DIF)	0.00(DIF)	0.00(DIF)
20	0.01(DIF)	0.02(DIF)	0.19(NO-DIF)
21	0.04(DIF)	0.58(NO-DIF)	0.25(NO-DIF)
22	0.58(NO-DIF)	0.63(NO-DIF)	0.00(DIF)
23	0.30(NO-DIF)	0.25(NO-DIF)	0.00(DIF)
24	0.00(DIF)	0.32(NO-DIF)	0.00(DIF)
25	0.46(NO-DIF)	0.17(NO-DIF)	0.93(NO-DIF)
26	0.11(NO-DIF)	0.82(NO-DIF)	0.00(DIF)
27	0.00(DIF)	0.82(NO-DIF)	0.00(DIF)
28	0.00(DIF)	0.01(DIF)	0.00(DIF)
29	0.00(DIF)	0.62(NO-DIF)	0.85(NO-DIF)
30	0.89(NO-DIF)	0.60(NO-DIF)	0.78(NO-DIF)
จำนวนข้อที่พบ DIF	(19 ข้อ) 63%	(8 ข้อ) 27%	(22 ข้อ) 73%

ภาคผนวก จ
ตัวอย่าง Print Out การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ
ด้วยวิธี IRT-LR

IRTPRO Version 2.1**Output generated by IRTPRO estimation engine Version 4.54 (32-bit)**

Project:	Literacy 2000
Description:	
Date:	06 August 2017
Time:	02:15 PM

Table of Contents2PL Model Item Parameter Estimates for Group 1, logit: $a\theta + c$ or $a(\theta - b)$ Summed-Score Based Item Diagnostic Tables and χ^2 s for Group 12PL Model Item Parameter Estimates for Group 2, logit: $a\theta + c$ or $a(\theta - b)$ Summed-Score Based Item Diagnostic Tables and χ^2 s for Group 2

Group Parameter Estimates

DIF Statistics for Graded Items

Marginal fit (χ^2) and Standardized LD χ^2 Statistics for Group 1Marginal fit (χ^2) and Standardized LD χ^2 Statistics for Group 2Item Information Function Values for Group 1 at 15 Values of θ from -2.8 to 2.8Item Information Function Values for Group 2 at 15 Values of θ from -2.8 to 2.8

Likelihood-based Values and Goodness of Fit Statistics

Summary of the Data and Control Parameters

2PL Model Item Parameter Estimates for Group 1, logit: $a\theta + c$ or $a(\theta - b)$ (Back to TOC)

Item	Label	a	s.e.	c	s.e.	b	s.e.
1	Litera1	² 0.63	0.09	¹ 0.35	0.08	-0.57	0.11
2	Litera2	⁴ 0.44	0.08	³ -0.63	0.07	1.44	0.36
3	Litera3	⁶ 0.67	0.09	⁵ 0.49	0.08	-0.73	0.11
4	Litera4	⁸ 0.92	0.11	⁷ 1.50	0.11	-1.63	0.14
5	Litera5	¹⁰ 0.82	0.10	⁹ 0.77	0.09	-0.94	0.10
6	Litera6	¹² 0.93	0.11	¹¹ 1.50	0.11	-1.60	0.14
7	Litera7	¹⁴ 1.23	0.12	¹³ 1.28	0.11	-1.04	0.08
8	Litera8	¹⁶ 0.65	0.09	¹⁵ 0.09	0.07	-0.13	0.11
9	Litera9	¹⁸ 0.80	0.11	¹⁷ 1.56	0.11	-1.94	0.19
10	Litera10	²⁰ 0.33	0.08	¹⁹ -0.30	0.07	0.88	0.34
11	Litera11	²² 0.64	0.09	²¹ 1.03	0.09	-1.61	0.19
12	Litera12	²⁴ 1.07	0.11	²³ 0.67	0.09	-0.62	0.07
13	Litera13	²⁶ 0.38	0.08	²⁵ -0.26	0.07	0.67	0.27
14	Litera14	²⁸ 0.87	0.10	²⁷ 0.05	0.08	-0.06	0.09
15	Litera15	³⁰ 0.51	0.08	²⁹ -0.25	0.07	0.49	0.18
16	Litera16	³² 0.75	0.09	³¹ 0.02	0.07	-0.03	0.10
17	Litera17	³⁴ 1.09	0.11	³³ 1.11	0.10	-1.02	0.08
18	Litera18	³⁶ 0.64	0.09	³⁵ -0.10	0.07	0.16	0.12
19	Litera19	³⁸ 1.11	0.12	³⁷ 1.77	0.13	-1.59	0.12
20	Litera20	⁴⁰ 0.80	0.10	³⁹ 0.98	0.09	-1.22	0.12
21	Litera21	⁴² 0.61	0.09	⁴¹ 0.13	0.07	-0.22	0.11
22	Litera22	⁴⁴ 0.56	0.09	⁴³ -0.92	0.08	1.66	0.33

23	Litera23	⁴⁶	0.32	0.08	⁴⁵	-0.63	0.07	1.98	0.61
24	Litera24	⁴⁸	0.57	0.09	⁴⁷	-0.08	0.07	0.13	0.13
25	Litera25	⁵⁰	0.18	0.08	⁴⁹	-0.49	0.07	2.66	1.31
26	Litera26	⁵²	0.32	0.08	⁵¹	-0.06	0.07	0.17	0.23
27	Litera27	⁵⁴	0.45	0.08	⁵³	-0.32	0.07	0.72	0.24
28	Litera28	⁵⁶	0.65	0.09	⁵⁵	0.38	0.08	-0.58	0.11
29	Litera29	⁵⁸	0.60	0.09	⁵⁷	-0.24	0.07	0.41	0.15
30	Litera30	⁶⁰	0.97	0.10	⁵⁹	0.69	0.09	-0.72	0.08

Summed-Score Based Item Diagnostic Tables and χ^2 's for Group 1 (Back to TOC)

S- χ^2 Item Level Diagnostic Statistics

Item	Label	χ^2	d.f.	Probability
1	Litera1	15.68	20	0.7369
2	Litera2	21.82	20	0.3523
3	Litera3	38.89	20	0.0069
4	Litera4	20.84	19	0.3477
5	Litera5	21.51	19	0.3084
6	Litera6	22.38	19	0.2648
7	Litera7	14.79	18	0.6773
8	Litera8	27.45	21	0.1561
9	Litera9	17.87	19	0.5323
10	Litera10	21.61	22	0.4850
11	Litera11	25.20	20	0.1934
12	Litera12	23.25	18	0.1806
13	Litera13	20.01	21	0.5219
14	Litera14	23.04	19	0.2351
15	Litera15	26.91	20	0.1375
16	Litera16	22.88	19	0.2418
17	Litera17	31.80	18	0.0231
18	Litera18	28.36	20	0.1009
19	Litera19	32.64	18	0.0184
20	Litera20	27.96	20	0.1100
21	Litera21	23.68	21	0.3081
22	Litera22	26.71	20	0.1432
23	Litera23	27.35	20	0.1254
24	Litera24	17.33	20	0.6326
25	Litera25	24.59	21	0.2645
26	Litera26	24.21	21	0.2821
27	Litera27	22.29	21	0.3844
28	Litera28	11.40	20	0.9353
29	Litera29	22.95	20	0.2903
30	Litera30	23.41	20	0.2683

2PL Model Item Parameter Estimates for Group 2, logit: $a\theta + c$ or $a(\theta - b)$ (Back to TOC)

Item	Label	<i>a</i>	s.e.	<i>c</i>	s.e.	<i>b</i>	s.e.
1	Litera1	⁶² 0.45	0.08	⁶¹ 0.17	0.07	-0.39	0.16
2	Litera2	⁶⁴ 0.20	0.08	⁶³ -0.75	0.07	3.80	1.55
3	Litera3	⁶⁶ 0.65	0.09	⁶⁵ 0.66	0.07	-1.00	0.16
4	Litera4	⁶⁸ 0.71	0.10	⁶⁷ 1.40	0.09	-1.98	0.26
5	Litera5	⁷⁰ 0.56	0.08	⁶⁹ 0.19	0.07	-0.35	0.13
6	Litera6	⁷² 0.92	0.11	⁷¹ 1.34	0.09	-1.46	0.16
7	Litera7	⁷⁴ 1.38	0.14	⁷³ 1.47	0.11	-1.06	0.09
8	Litera8	⁷⁶ 0.68	0.09	⁷⁵ 0.01	0.07	-0.02	0.10
9	Litera9	⁷⁸ 0.85	0.12	⁷⁷ 1.75	0.10	-2.06	0.24
10	Litera10	⁸⁰ 0.34	0.08	⁷⁹ -0.32	0.07	0.95	0.28
11	Litera11	⁸² 0.76	0.10	⁸¹ 1.09	0.08	-1.45	0.18
12	Litera12	⁸⁴ 1.07	0.11	⁸³ 0.52	0.08	-0.48	0.08
13	Litera13	⁸⁶ 0.27	0.08	⁸⁵ -0.25	0.06	0.92	0.35
14	Litera14	⁸⁸ 0.78	0.09	⁸⁷ -0.03	0.07	0.03	0.09
15	Litera15	⁹⁰ 0.29	0.08	⁸⁹ -0.16	0.06	0.56	0.27
16	Litera16	⁹² 0.69	0.09	⁹¹ 0.18	0.07	-0.26	0.10
17	Litera17	⁹⁴ 1.29	0.13	⁹³ 1.37	0.10	-1.06	0.09
18	Litera18	⁹⁶ 0.75	0.09	⁹⁵ -0.09	0.07	0.13	0.09
19	Litera19	⁹⁸ 1.45	0.16	⁹⁷ 2.14	0.14	-1.48	0.11
20	Litera20	¹⁰⁰ 1.05	0.11	⁹⁹ 1.03	0.09	-0.99	0.10
21	Litera21	¹⁰² 0.67	0.09	¹⁰¹ 0.11	0.07	-0.16	0.10
22	Litera22	¹⁰⁴ 0.56	0.09	¹⁰³ -1.07	0.08	1.91	0.31
23	Litera23	¹⁰⁶ 0.39	0.08	¹⁰⁵ -0.65	0.07	1.67	0.37
24	Litera24	¹⁰⁸ 0.55	0.08	¹⁰⁷ 0.14	0.07	-0.25	0.12
25	Litera25	¹¹⁰ 0.23	0.08	¹⁰⁹ -0.49	0.07	2.13	0.76
26	Litera26	¹¹² 0.14	0.07	¹¹¹ -0.02	0.06	0.17	0.47
27	Litera27	¹¹⁴ 0.63	0.09	¹¹³ -0.20	0.07	0.31	0.11
28	Litera28	¹¹⁶ 0.80	0.10	¹¹⁵ 0.47	0.07	-0.59	0.10
29	Litera29	¹¹⁸ 0.57	0.08	¹¹⁷ 0.00	0.07	-0.00	0.12
30	Litera30	¹²⁰ 0.80	0.09	¹¹⁹ 0.34	0.07	-0.43	0.09

Summed-Score Based Item Diagnostic Tables and X's for Group 2 (Back to TOC)**S-X² Item Level Diagnostic Statistics**

Item	Label	X ²	d.f.	Probability
1	Litera1	23.75	20	0.2526
2	Litera2	30.56	20	0.0611
3	Litera3	18.53	19	0.4887
4	Litera4	18.60	18	0.4186
5	Litera5	20.78	20	0.4120
6	Litera6	14.29	18	0.7110
7	Litera7	21.18	18	0.2699
8	Litera8	23.44	20	0.2668
9	Litera9	19.97	19	0.3982
10	Litera10	25.42	20	0.1855
11	Litera11	15.38	19	0.6987
12	Litera12	13.58	18	0.7569
13	Litera13	21.96	20	0.3446

14	Litera14	26.26	19	0.1227
15	Litera15	18.50	20	0.5561
16	Litera16	26.63	19	0.1133
17	Litera17	26.12	18	0.0969
18	Litera18	18.41	19	0.4965
19	Litera19	20.65	17	0.2416
20	Litera20	19.75	18	0.3490
21	Litera21	16.95	19	0.5943
22	Litera22	26.40	19	0.1190
23	Litera23	18.91	19	0.4639
24	Litera24	22.33	20	0.3221
25	Litera25	25.45	20	0.1841
26	Litera26	17.25	20	0.6376
27	Litera27	36.60	19	0.0089
28	Litera28	13.37	19	0.8196
29	Litera29	14.74	20	0.7917
30	Litera30	22.12	19	0.2777

Group Parameter Estimates (Back to TOC)

Group	Label	μ	s.e.	σ^2	s.e.	σ	s.e.
1	G1	-0.36	----	0.98	----	0.99	----
2	G2	0.00	----	1.00	----	1.00	----

DIF Statistics for Graded Items (Back to TOC)

Item numbers in:

Group 1	Group 2	Total χ^2	d.f.	p	χ^2_a	d.f.	p	χ^2_{cla}	d.f.	p
1	1	4.4	2	0.1124	2.3	1	0.1324	2.1	1	0.1489
2	2	5.2	2	0.0729	4.2	1	0.0411	1.1	1	0.3028
3	3	3.1	2	0.2122	0.0	1	0.9159	3.1	1	0.0782
4	4	1.9	2	0.3870	1.9	1	0.1713	0.0	1	0.8791
5	5	28.4	2	0.0001	4.3	1	0.0372	24.1	1	0.0001
6	6	1.6	2	0.4463	0.0	1	0.9466	1.6	1	0.2050
7	7	1.4	2	0.4880	0.6	1	0.4415	0.8	1	0.3590
8	8	0.7	2	0.6976	0.1	1	0.7877	0.6	1	0.4212
9	9	2.0	2	0.3684	0.1	1	0.7802	1.9	1	0.1663
10	10	0.1	2	0.9508	0.0	1	0.9450	0.1	1	0.7567
11	11	0.8	2	0.6802	0.7	1	0.3874	0.0	1	0.8802
12	12	1.9	2	0.3864	0.0	1	0.9658	1.9	1	0.1685
13	13	1.1	2	0.5869	1.0	1	0.3146	0.1	1	0.8177
14	14	0.8	2	0.6818	0.4	1	0.5289	0.4	1	0.5435
15	15	5.5	2	0.0639	3.8	1	0.0517	1.7	1	0.1913
16	16	3.2	2	0.2027	0.2	1	0.6404	3.0	1	0.0851
17	17	3.1	2	0.2166	1.4	1	0.2453	1.7	1	0.1912
18	18	0.7	2	0.6883	0.7	1	0.3890	0.0	1	0.9485
19	19	3.9	2	0.1411	2.8	1	0.0936	1.1	1	0.2940
20	20	2.9	2	0.2353	2.6	1	0.1048	0.3	1	0.6123
21	21	0.5	2	0.7866	0.3	1	0.5778	0.2	1	0.6803
22	22	1.8	2	0.4109	0.0	1	0.9977	1.8	1	0.1827

23	23	0.5	2	0.7925	0.4	1	0.5338	0.1	1	0.7807
24	24	5.4	2	0.0688	0.0	1	0.8382	5.3	1	0.0212
25	25	0.2	2	0.9178	0.2	1	0.6824	0.0	1	0.9504
26	26	3.5	2	0.1747	3.0	1	0.0826	0.5	1	0.4836
27	27	3.6	2	0.1704	2.4	1	0.1191	1.1	1	0.2927
28	28	1.6	2	0.4534	1.3	1	0.2523	0.3	1	0.6036
29	29	6.8	2	0.0332	0.1	1	0.7780	6.7	1	0.0095
30	30	9.7	2	0.0078	1.3	1	0.2466	8.4	1	0.0038

Marginal fit (χ^2) and Standardized LD χ^2 Statistics for Group 1 [\(Back to TOC\)](#)

Item	Label	χ^2	Margin al									
			1	2	3	4	5	6	7	8	9	10
1	Litera1	0.0										
2	Litera2	0.0	-0.7									
3	Litera3	0.0	3.0	0.1								
4	Litera4	0.0	-0.7	-0.7	1.2							
5	Litera5	0.0	1.3	-0.5	-0.7	-0.6						
6	Litera6	0.0	-0.2	2.3	0.5	-0.6	2.1					
7	Litera7	0.0	0.4	-0.5	-0.3	-0.5	-0.1	-0.5				
8	Litera8	0.0	-0.7	-0.1	0.6	0.0	-0.4	0.7	-0.4			
9	Litera9	0.0	0.5	3.0	-0.7	-0.5	-0.6	2.6	-0.7	1.7		
10	Litera10	0.0	-0.2	-0.1	-0.6	-0.6	-0.7	-0.7	-0.5	-0.7	-0.7	
11	Litera11	0.0	-0.1	-0.5	-0.4	-0.6	-0.3	-0.7	0.8	0.8	-0.5	-0.5
12	Litera12	0.0	-0.2	-0.4	-0.3	-0.7	-0.7	-0.3	-0.4	-0.7	-0.6	-0.6
13	Litera13	0.0	-0.4	1.2	2.6	0.3	0.4	-0.4	-0.7	-0.4	-0.2	2.3
14	Litera14	0.0	-0.7	-0.1	1.6	-0.6	0.1	0.6	-0.2	0.6	-0.6	0.6
15	Litera15	0.0	-0.7	0.1	-0.6	-0.5	-0.6	-0.7	-0.7	0.3	-0.7	-0.7
16	Litera16	0.0	2.4	-0.5	-0.7	-0.6	1.2	-0.1	-0.1	1.9	-0.6	-0.6
17	Litera17	0.0	-0.5	-0.7	2.8	-0.7	-0.7	-0.5	0.5	-0.5	-0.6	-0.1
18	Litera18	0.0	1.4	-0.7	7.0	-0.5	-0.6	-0.6	-0.7	-0.7	0.0	-0.7
19	Litera19	0.0	1.8	2.2	0.6	0.7	-0.7	-0.4	1.7	-0.7	-0.7	0.5
20	Litera20	0.0	-0.6	1.9	-0.5	-0.4	-0.3	-0.7	-0.2	0.8	-0.6	-0.0
21	Litera21	0.0	1.7	1.5	-0.4	0.2	-0.5	-0.7	-0.7	0.0	-0.5	0.4
22	Litera22	0.0	0.3	1.4	-0.6	-0.1	-0.7	-0.6	0.2	-0.7	1.2	-0.6
23	Litera23	0.0	0.4	1.2	1.2	-0.6	0.2	1.0	-0.6	-0.7	-0.7	-0.5
24	Litera24	0.0	0.7	0.2	3.7	0.1	-0.7	-0.5	-0.7	-0.7	1.0	0.0
25	Litera25	0.0	4.0	1.0	1.2	-0.5	0.5	1.2	-0.7	-0.1	-0.7	-0.6

		5											
26	Litera26	0.0	0.4	0.4	1.8	-0.7	-0.1	-0.7	-0.2	-0.5	-0.5	-0.5	-0.6
27	Litera27	0.0	-0.7	1.6	0.1	0.9	1.1	-0.7	-0.2	-0.7	0.6	0.2	
28	Litera28	0.0	3.2	4.8	-0.0	0.3	2.2	1.3	-0.7	-0.4	-0.7	0.2	
29	Litera29	0.0	-0.6	0.5	-0.7	-0.7	0.6	-0.1	-0.2	1.4	-0.4	-0.7	
30	Litera30	0.0	1.7	-0.7	-0.3	-0.3	-0.7	-0.5	-0.1	0.5	-0.5	1.1	

Item	Label	χ^2	Marginal										
			11	12	13	14	15	16	17	18	19	20	
11	Litera11	0.0											
12	Litera12	0.0	0.1										
13	Litera13	0.0	0.1	-0.3									
14	Litera14	0.0	-0.6	2.9	5.1								
15	Litera15	0.0	-0.7	0.8	0.5	2.4							
16	Litera16	0.0	-0.3	0.5	-0.3	-0.6	-0.5						
17	Litera17	0.0	0.7	-0.7	-0.2	-0.7	2.2	-0.7					
18	Litera18	0.0	0.3	-0.6	-0.2	-0.7	-0.2	-0.5	2.9				
19	Litera19	0.0	0.6	1.6	-0.7	0.9	-0.5	-0.7	-0.6	-0.7			
20	Litera20	0.0	-0.6	-0.3	-0.4	0.1	-0.7	-0.2	-0.5	-0.7	9.5		
21	Litera21	0.0	1.3	0.8	-0.1	-0.6	0.4	0.7	-0.5	-0.3	-0.6	-0.6	
22	Litera22	0.0	-0.7	-0.5	-0.3	0.5	1.7	1.4	-0.0	-0.6	-0.1	-0.0	
23	Litera23	0.0	-0.7	-0.5	-0.7	-0.7	3.5	-0.4	0.3	-0.7	-0.7	0.1	
24	Litera24	0.0	-0.7	1.3	-0.4	-0.7	-0.7	-0.6	0.7	-0.6	-0.7	0.8	
25	Litera25	0.0	-0.7	1.2	-0.5	-0.5	2.6	0.2	-0.4	-0.4	-0.4	0.3	
26	Litera26	0.0	-0.2	0.8	-0.5	2.7	-0.0	-0.6	-0.6	-0.6	-0.5	-0.3	
27	Litera27	0.0	-0.6	-0.6	-0.5	1.7	-0.7	-0.7	-0.7	0.1	-0.5	0.2	
28	Litera28	0.0	0.3	-0.7	3.6	-0.7	-0.6	-0.7	-0.5	-0.5	-0.6	0.1	
29	Litera29	0.0	-0.6	0.2	-0.6	-0.4	-0.7	-0.7	3.0	0.1	-0.5	-0.2	
30	Litera30	0.0	-0.7	-0.7	2.7	-0.6	0.0	0.5	-0.7	-0.3	-0.4	-0.7	

Item	Label	χ^2	Marginal								
			21	22	23	24	25	26	27	28	29
21	Litera21	0.0									
22	Litera22	0.0	1.5								
23	Litera23	0.0	0.1	-0.7							
24	Litera24	0.0	1.1	-0.7	-0.7						
25	Litera25	0.0	3.0	0.7	0.0	-0.2					
26	Litera26	0.0	-0.3	-0.6	-0.6	4.9	6.4				
27	Litera27	0.0	1.9	-0.6	1.6	2.0	-0.5	4.1			

28	Litera2 8	0.0	-0.7	-0.6	-0.3	-0.7	3.8	0.0	0.4	
29	Litera2 9	0.0	-0.7	-0.5	1.5	3.8	1.8	0.0	0.7	8.4
30	Litera3 0	0.0	0.1	0.0	2.5	-0.6	-0.5	-0.3	-0.4	0.2

Marginal fit (χ^2) and Standardized LD χ^2 Statistics for Group 2 (Back to TOC)

Item	Label	χ^2	Marginal									
			1	2	3	4	5	6	7	8	9	10
1	Litera1	0.0										
2	Litera2	0.0	-0.6									
3	Litera3	0.0	3.1	-0.4								
4	Litera4	0.0	5.6	-0.6	0.2							
5	Litera5	0.0	-0.1	0.1	-0.2	-0.7						
6	Litera6	0.0	-0.7	0.7	-0.5	2.5	-0.6					
7	Litera7	0.0	0.0	0.7	-0.7	0.8	-0.4	1.4				
8	Litera8	0.0	-0.5	-0.3	3.7	-0.3	0.4	0.4	0.0			
9	Litera9	0.0	-0.1	-0.6	-0.6	-0.2	0.2	-0.6	-0.3	-0.6		
10	Litera1 0	0.0	0.1	2.1	-0.7	0.3	-0.7	-0.1	-0.4	-0.5	-0.7	
11	Litera1 1	0.0	0.9	-0.5	0.9	-0.6	-0.6	1.4	-0.6	-0.5	-0.3	0.1
12	Litera1 2	0.0	-0.7	-0.5	0.8	-0.7	-0.4	0.2	-0.3	-0.1	0.4	-0.2
13	Litera1 3	0.0	0.8	6.1	3.3	-0.1	1.4	-0.7	-0.7	-0.7	-0.7	-0.7
14	Litera1 4	0.0	-0.7	-0.6	-0.4	1.0	-0.2	-0.4	0.9	-0.7	-0.7	0.6
15	Litera1 5	0.0	1.3	1.7	-0.3	4.7	-0.7	3.9	-0.2	-0.5	-0.7	2.2
16	Litera1 6	0.0	-0.6	-0.6	-0.7	-0.7	-0.5	-0.3	-0.4	-0.6	-0.0	-0.4
17	Litera1 7	0.0	-0.7	0.2	-0.4	0.7	-0.3	-0.1	-0.6	-0.7	1.1	-0.7
18	Litera1 8	0.0	-0.3	-0.3	1.3	-0.6	-0.4	1.3	0.4	2.3	-0.3	-0.4
19	Litera1 9	0.0	0.2	-0.6	-0.6	-0.2	-0.6	1.5	-0.4	-0.0	-0.6	-0.6
20	Litera2 0	0.0	0.7	3.8	-0.5	-0.1	0.3	-0.7	-0.6	-0.7	-0.6	-0.6
21	Litera2 1	0.0	0.9	-0.7	-0.3	0.2	-0.7	0.2	-0.6	-0.7	-0.7	-0.2
22	Litera2 2	0.0	0.9	0.1	0.2	-0.6	2.2	-0.4	-0.4	0.8	-0.4	-0.6
23	Litera2 3	0.0	0.2	2.4	2.1	0.8	-0.6	-0.5	-0.5	1.0	-0.6	0.3
24	Litera2 4	0.0	-0.6	-0.4	1.3	-0.6	-0.1	-0.7	1.3	-0.7	-0.1	3.5
25	Litera2 5	0.0	-0.2	-0.5	2.1	2.4	6.2	-0.7	0.2	-0.4	-0.7	-0.5
26	Litera2 6	0.0	-0.7	2.0	6.6	0.4	-0.2	-0.5	-0.7	-0.7	-0.1	-0.4
27	Litera2	0.0	-0.6	-0.6	0.1	0.3	-0.6	-0.5	-0.7	-0.3	-0.5	-0.5

		7											
28	Litera2 8	0.0	2.8	-0.6	5.4	2.9	-0.5	3.2	0.1	0.2	-0.5	0.0	
29	Litera2 9	0.0	-0.1	2.4	1.2	0.1	1.4	0.1	-0.7	-0.5	0.1	-0.7	
30	Litera3 0	0.0	1.6	0.4	-0.5	-0.3	0.8	-0.4	-0.3	1.4	-0.4	-0.7	

Item	Label	χ^2	Marginal										
			11	12	13	14	15	16	17	18	19	20	
11	Litera1 1	0.0											
12	Litera1 2	0.0	-0.6										
13	Litera1 3	0.0	-0.6	-0.2									
14	Litera1 4	0.0	0.2	-0.5	0.5								
15	Litera1 5	0.0	-0.2	-0.0	-0.7	-0.7							
16	Litera1 6	0.0	-0.6	-0.5	0.5	0.1	-0.6						
17	Litera1 7	0.0	-0.2	-0.6	-0.7	-0.5	-0.6	1.2					
18	Litera1 8	0.0	-0.2	0.2	1.2	0.3	-0.7	1.4	3.7				
19	Litera1 9	0.0	-0.5	1.9	-0.6	-0.2	-0.7	0.0	-0.6	-0.6			
20	Litera2 0	0.0	1.8	-0.7	-0.7	1.0	-0.7	-0.2	-0.7	-0.6	0.4		
21	Litera2 1	0.0	-0.7	-0.7	-0.7	-0.2	-0.7	-0.7	-0.3	-0.5	-0.4	-0.1	
22	Litera2 2	0.0	-0.7	-0.2	1.7	0.8	0.7	-0.7	-0.5	-0.5	-0.3	-0.7	
23	Litera2 3	0.0	-0.5	-0.4	-0.2	-0.4	-0.7	-0.4	2.0	0.1	-0.3	-0.7	
24	Litera2 4	0.0	-0.1	-0.1	0.4	-0.2	0.1	-0.4	-0.7	-0.6	3.5	-0.7	
25	Litera2 5	0.0	-0.7	1.3	-0.5	1.7	2.8	0.4	-0.7	1.3	-0.6	-0.6	
26	Litera2 6	0.0	-0.1	-0.3	-0.4	-0.3	0.2	-0.1	-0.7	1.7	0.0	0.3	
27	Litera2 7	0.0	-0.6	-0.0	-0.3	-0.7	-0.5	0.1	0.6	-0.7	-0.5	-0.1	
28	Litera2 8	0.0	-0.4	-0.7	0.3	1.9	0.1	-0.7	-0.6	-0.1	1.2	-0.6	
29	Litera2 9	0.0	-0.6	-0.7	0.2	-0.5	-0.7	-0.7	0.1	1.3	0.2	-0.2	
30	Litera3 0	0.0	-0.7	-0.5	-0.4	-0.2	-0.7	2.3	-0.3	-0.3	-0.6	-0.7	

Marginal

Item	Label	χ^2	21	22	23	24	25	26	27	28	29
21	Litera2 1	0.0									
22	Litera2 2	0.0	4.2								
23	Litera2 3	0.0	-0.5	2.6							
24	Litera2 4	0.0	-0.5	-0.7	-0.1						
25	Litera2 5	0.0	1.2	-0.6	-0.5	-0.4					
26	Litera2 6	0.0	0.0	-0.1	-0.6	3.7	-0.3				
27	Litera2 7	0.0	0.2	-0.4	-0.4	-0.7	-0.5	-0.4			
28	Litera2 8	0.0	1.6	0.9	-0.7	-0.7	-0.7	1.2 -0.4			
29	Litera2 9	0.0	0.7	-0.7	1.5	-0.7	-0.6	-0.4 -0.4	8.8		
30	Litera3 0	0.0	-0.4	0.5	0.2	-0.3	0.4	0.3 -0.5	-0.7	-0.4	

Item Information Function Values for Group 1 at 15 Values of θ from -2.8 to 2.8 (Back to TOC)

$\theta:$																
Item	Label	-2.8	-2.4	-2.0	-1.6	-1.2	-0.8	-0.4	-0.0	0.4	0.8	1.2	1.6	2.0	2.4	2.8
1	Litera1	0.06	0.07	0.08	0.09	0.09	0.10	0.10	0.09	0.09	0.08	0.07	0.06	0.05	0.05	0.04
2	Litera2	0.02	0.03	0.03	0.03	0.03	0.04	0.04	0.04	0.05	0.05	0.05	0.05	0.05	0.05	0.04
3	Litera3	0.07	0.08	0.09	0.10	0.11	0.11	0.11	0.10	0.10	0.09	0.08	0.06	0.05	0.04	0.04
4	Litera4	0.16	0.19	0.20	0.21	0.20	0.18	0.16	0.13	0.10	0.07	0.05	0.04	0.03	0.02	0.01
5	Litera5	0.10	0.12	0.14	0.16	0.17	0.17	0.16	0.15	0.13	0.11	0.08	0.07	0.05	0.04	0.03
6	Litera6	0.16	0.19	0.21	0.22	0.21	0.19	0.16	0.13	0.10	0.08	0.06	0.04	0.03	0.02	0.01
7	Litera7	0.14	0.20	0.27	0.34	0.38	0.37	0.33	0.26	0.19	0.13	0.09	0.05	0.03	0.02	0.01
8	Litera8	0.05	0.06	0.07	0.08	0.09	0.10	0.10	0.11	0.10	0.10	0.09	0.08	0.07	0.06	0.05
9	Litera9	0.14	0.16	0.16	0.16	0.15	0.13	0.11	0.09	0.07	0.06	0.04	0.03	0.03	0.02	0.01
10	Litera1 0	0.02	0.02	0.02	0.02	0.02	0.03	0.03	0.03	0.03	0.03	0.03	0.03	0.03	0.03	0.03
11	Litera1 1	0.09	0.10	0.10	0.10	0.10	0.10	0.09	0.08	0.07	0.06	0.05	0.04	0.03	0.03	0.02
12	Litera1 2	0.09	0.13	0.17	0.22	0.26	0.28	0.28	0.26	0.21	0.17	0.12	0.09	0.06	0.04	0.03
13	Litera1 3	0.02	0.03	0.03	0.03	0.03	0.03	0.03	0.04	0.04	0.04	0.04	0.04	0.03	0.03	0.03
14	Litera1 4	0.06	0.08	0.10	0.12	0.15	0.17	0.18	0.19	0.18	0.16	0.14	0.12	0.09	0.07	0.05
15	Litera1 5	0.03	0.04	0.04	0.05	0.05	0.06	0.06	0.06	0.06	0.06	0.06	0.06	0.05	0.05	0.05
16	Litera1 6	0.06	0.07	0.09	0.10	0.12	0.13	0.14	0.14	0.14	0.13	0.12	0.10	0.08	0.07	0.05
17	Litera1 7	0.13	0.18	0.23	0.27	0.29	0.29	0.26	0.22	0.17	0.13	0.09	0.06	0.04	0.03	0.02
18	Litera1 8	0.05	0.06	0.07	0.08	0.08	0.09	0.10	0.10	0.10	0.09	0.08	0.07	0.06	0.05	0.05

19	Litera1 9	0.20	0.25	0.29	0.31	0.30	0.26	0.21	0.15	0.11	0.08	0.05	0.03	0.02	0.01	0.01
20	Litera2 0	0.11	0.13	0.15	0.16	0.16	0.16	0.14	0.13	0.11	0.09	0.07	0.05	0.04	0.03	0.02
21	Litera2 1	0.05	0.06	0.07	0.08	0.08	0.09	0.09	0.09	0.09	0.08	0.08	0.07	0.06	0.05	0.04
22	Litera2 2	0.02	0.03	0.03	0.04	0.04	0.05	0.06	0.06	0.07	0.07	0.08	0.08	0.08	0.07	0.07
23	Litera2 3	0.01	0.02	0.02	0.02	0.02	0.02	0.02	0.02	0.02	0.02	0.02	0.03	0.03	0.03	0.02
24	Litera2 4	0.04	0.05	0.06	0.06	0.07	0.08	0.08	0.08	0.08	0.08	0.08	0.07	0.06	0.06	0.05
25	Litera2 5	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01
26	Litera2 6	0.02	0.02	0.02	0.02	0.03	0.03	0.03	0.03	0.03	0.03	0.03	0.02	0.02	0.02	0.02
27	Litera2 7	0.03	0.03	0.04	0.04	0.04	0.04	0.05	0.05	0.05	0.05	0.05	0.05	0.05	0.04	0.04
28	Litera2 8	0.07	0.08	0.09	0.10	0.10	0.11	0.11	0.10	0.10	0.09	0.08	0.07	0.06	0.05	0.04
29	Litera2 9	0.04	0.05	0.06	0.06	0.07	0.08	0.08	0.09	0.09	0.09	0.08	0.08	0.07	0.06	0.06
30	Litera3 0	0.10	0.13	0.16	0.20	0.22	0.23	0.23	0.21	0.18	0.14	0.11	0.08	0.06	0.04	0.03

Test Information: 3.19 3.66 4.12 4.49 4.71 4.74 4.57 4.26 3.87 3.47 3.09 2.75 2.46 2.22 2.01

Expected s.e.: 0.56 0.52 0.49 0.47 0.46 0.46 0.47 0.48 0.51 0.54 0.57 0.60 0.64 0.67 0.71

Marginal Reliability for Response Pattern Scores: 0.75

Item Information Function Values for Group 2 at 15 Values of θ from -2.8 to 2.8 [\(Back to TOC\)](#)

		Lite																	
11	ra1	0.11	0.13	0.14	0.14	0.14	0.13	0.12	0.11	0.09	0.07	0.06	0.05	0.04	0.03	0.02			
12	ra1	0.08	0.12	0.16	0.21	0.25	0.28	0.29	0.27	0.23	0.19	0.14	0.10	0.07	0.05	0.03			
13	ra1	0.01	0.02	0.02	0.02	0.02	0.02	0.02	0.02	0.02	0.02	0.02	0.02	0.02	0.02	0.02	0.02	0.02	0.02
14	ra1	0.05	0.07	0.09	0.10	0.12	0.14	0.15	0.15	0.15	0.14	0.13	0.11	0.09	0.07	0.06			
15	ra1	0.02	0.02	0.02	0.02	0.02	0.02	0.02	0.02	0.02	0.02	0.02	0.02	0.02	0.02	0.02	0.02	0.02	0.02
16	ra1	0.06	0.07	0.09	0.10	0.11	0.12	0.12	0.12	0.11	0.11	0.09	0.08	0.07	0.06	0.05			
17	ra1	0.14	0.21	0.29	0.37	0.41	0.40	0.35	0.27	0.19	0.13	0.08	0.05	0.03	0.02	0.01			
18	ra1	0.05	0.06	0.08	0.09	0.11	0.12	0.13	0.14	0.14	0.13	0.12	0.10	0.09	0.07	0.06	0.06	0.05	0.05
19	ra1	0.23	0.35	0.46	0.52	0.51	0.42	0.30	0.20	0.12	0.07	0.04	0.02	0.01	0.01	0.00			
20	ra2	0.12	0.17	0.21	0.25	0.27	0.27	0.25	0.21	0.17	0.13	0.09	0.06	0.04	0.03	0.02			
21	ra2	0.06	0.07	0.08	0.09	0.10	0.11	0.11	0.11	0.11	0.10	0.09	0.08	0.07	0.06	0.05	0.05		
22	ra2	0.02	0.02	0.03	0.03	0.04	0.05	0.05	0.06	0.07	0.07	0.07	0.07	0.08	0.08	0.08	0.07		
23	ra2	0.02	0.02	0.02	0.03	0.03	0.03	0.03	0.03	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04		
24	ra2	0.05	0.05	0.06	0.07	0.07	0.07	0.08	0.08	0.07	0.07	0.06	0.06	0.05	0.05	0.04			
25	ra2	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01		
26	ra2	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00		
27	ra2	0.04	0.05	0.06	0.07	0.08	0.09	0.09	0.10	0.10	0.10	0.09	0.08	0.08	0.07	0.06			
28	ra2	0.08	0.10	0.12	0.14	0.15	0.16	0.16	0.15	0.14	0.12	0.10	0.08	0.06	0.05	0.04			
29	ra2	0.05	0.05	0.06	0.07	0.07	0.08	0.08	0.08	0.08	0.08	0.07	0.07	0.06	0.05	0.05	0.05		

	Lite	ra3	0.07	0.09	0.11	0.13	0.15	0.16	0.16	0.16	0.15	0.13	0.11	0.09	0.07	0.05	0.04
30	0																

	Test	n:	3.11	3.66	4.23	4.72	4.98	4.95	4.68	4.27	3.81	3.38	2.99	2.65	2.36	2.12	1.92
Information	Expected	s.e.:	0.57	0.52	0.49	0.46	0.45	0.45	0.46	0.48	0.51	0.54	0.58	0.61	0.65	0.69	0.72

Marginal Reliability for Response Pattern Scores: 0.74

Likelihood-based Values and Goodness of Fit Statistics (Back to TOC)

Statistics based on the loglikelihood

-2loglikelihood:	74708.63
Akaike Information Criterion (AIC):	74948.63
Bayesian Information Criterion (BIC):	75620.73

Statistics based on the full item x item x ... classification

The table is too sparse to compute the general multinomial goodness of fit statistics.

Statistics based on one- and two-way marginal tables

The M₂ statistics were not requested.

Summary of the Data and Control Parameters (Back to TOC)

Group:	Group 1	Group 2
Sample Size	1000	1000
Number of Items	30	30
Number of Dimensions	1	1

Group 1

Item	Label	Categories	Model
1	Litera1	2	2PL
2	Litera2	2	2PL
3	Litera3	2	2PL
4	Litera4	2	2PL
5	Litera5	2	2PL
6	Litera6	2	2PL
7	Litera7	2	2PL
8	Litera8	2	2PL
9	Litera9	2	2PL
10	Litera10	2	2PL
11	Litera11	2	2PL
12	Litera12	2	2PL
13	Litera13	2	2PL
14	Litera14	2	2PL
15	Litera15	2	2PL
16	Litera16	2	2PL
17	Litera17	2	2PL
18	Litera18	2	2PL
19	Litera19	2	2PL

20	Litera20	2	2PL
21	Litera21	2	2PL
22	Litera22	2	2PL
23	Litera23	2	2PL
24	Litera24	2	2PL
25	Litera25	2	2PL
26	Litera26	2	2PL
27	Litera27	2	2PL
28	Litera28	2	2PL
29	Litera29	2	2PL
30	Litera30	2	2PL

Group 2

Item	Label	Categories	Model
1	Litera1	2	2PL
2	Litera2	2	2PL
3	Litera3	2	2PL
4	Litera4	2	2PL
5	Litera5	2	2PL
6	Litera6	2	2PL
7	Litera7	2	2PL
8	Litera8	2	2PL
9	Litera9	2	2PL
10	Litera10	2	2PL
11	Litera11	2	2PL
12	Litera12	2	2PL
13	Litera13	2	2PL
14	Litera14	2	2PL
15	Litera15	2	2PL
16	Litera16	2	2PL
17	Litera17	2	2PL
18	Litera18	2	2PL
19	Litera19	2	2PL
20	Litera20	2	2PL
21	Litera21	2	2PL
22	Litera22	2	2PL
23	Litera23	2	2PL
24	Litera24	2	2PL
25	Litera25	2	2PL
26	Litera26	2	2PL
27	Litera27	2	2PL
28	Litera28	2	2PL
29	Litera29	2	2PL
30	Litera30	2	2PL

Parameter Estimation Control Values

Bock-Aitkin EM Algorithm

Maximum number of cycles:	500
Convergence criterion:	1.00e-005
Maximum number of M-step iterations:	50

Convergence criterion for iterative M-steps: 1.00e-006
 Number of rectangular quadrature points: 49
 Minimum, Maximum quadrature points: -6.00 6.00
 SEM algorithm tolerance: 1.00e-003
 Standard error computation algorithm: Supplemented EM

DIF Analysis

All items are evaluated for DIF
(Conditional on population distribution estimates obtained with all items constrained equal)
Contrasts among groups:

Group:	1	2
Contrast 1	1.000	-1.000

Miscellaneous Control Values

Print parameter numbers? Yes
 Z tolerance, max. abs. logit value: 50.00
 Number of processor cores used: 8
 Number of cycles completed: 34
 Maximum parameter change: 0.00e+000
 Number of free parameters: 120

Processing times (in seconds)

E-step computations: 0.07
 M-step computations: 0.05
 Standard error computations: 1.65
 Goodness-of-fit statistics: 0.12
 Total: 1.89

Output Files

HTML results and control parameters: F:\Run IRT 06-08-2560\Literacy\Literacy 2000\Literacy 2000.Test1-irt.htm

Convergence and Numerical Stability

Engine status: Normal termination
 SEM algorithm status: Normal
 First-order test: Convergence criteria satisfied
 Condition number of information matrix: 1.28e+001
 Second-order test: Solution is a possible local maximum

ภาคผนวก ฉ

ตัวอย่าง Print Out การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ
ด้วยวิธี SIBTEST

number of items on test = 30

name of file for Ref. grp. scores = d:\RMCS หม้อวน วิเคราะห์โปรแกรม\Run
 SIBTEST\Run SIBTEST 20-08-2560\Literacy\Literacy 2000\Literacy 2000 Female.dat
 name of file for Focal grp. scores = d:\RMCS หม้อวน วิเคราะห์โปรแกรม\Run
 SIBTEST\Run SIBTEST 20-08-2560\Literacy\Literacy 2000\Literacy 2000 Male.dat
 minimum no. of examinees per matching score cell = 2
 number of runs for this data set = 30
 number of examinees in Reference group = 1000
 number of examinees in Focal group = 1000

Examinee Test Score Summary Statistics

Reference Group: Mean = 16.78

Standard deviation = 4.66

Focal Group: Mean = 15.33

Standard deviation = 4.87

Standardized Score Difference = 0.30

Item Statistics

= item number

p = proportion right on item (Classical Test Theory p-value)

r = point biserial (item score-test score correlation)

#: 1 2 3 4 5 6 7 8 9 10

p: 0.535 0.321 0.600 0.756 0.573 0.743 0.701 0.484 0.789 0.411

r: 0.308 0.237 0.335 0.348 0.345 0.381 0.478 0.356 0.337 0.241

#: 11 12 13 14 15 16 17 18 19 20

p: 0.701 0.580 0.423 0.468 0.430 0.492 0.690 0.451 0.790 0.673

r: 0.344 0.445 0.244 0.402 0.274 0.373 0.454 0.364 0.428 0.387

#: 21 22 23 24 25 26 27 28 29 30

p: 0.502 0.264 0.337 0.483 0.374 0.476 0.421 0.567 0.447 0.572

r: 0.352 0.283 0.218 0.332 0.195 0.212 0.316 0.380 0.341 0.414

\$

OUTPUT FOR RUN NUMBER 1 OUTPUT FOR RUN NUMBER 1

Suspect subtest items:

1

Matching subtest items:

2 3 4 5 6 7 8 9 10 11
 12 13 14 15 16 17 18 19 20 21
 22 23 24 25 26 27 28 29 30

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.019

proportion of Focal grp. examinees eliminated = 0.036

KR-20 for Ref. grp. = 0.727

KR-20 for Foc. grp. = 0.757

Matching Subtest Summary Statistics

Reference Group: Mean = 16.24

Standard deviation = 4.54

Focal Group: Mean = 14.80

Standard deviation = 4.73

Standardized Score Difference = 0.31

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.

-0.034 0.023 0.139065

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$

OUTPUT FOR RUN NUMBER 2 OUTPUT FOR RUN NUMBER 2

Suspect subtest items:

2

Matching subtest items:

1 3 4 5 6 7 8 9 10 11
 12 13 14 15 16 17 18 19 20 21
 22 23 24 25 26 27 28 29 30

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.017

proportion of Focal grp. examinees eliminated = 0.032

KR-20 for Ref. grp. = 0.727

KR-20 for Foc. grp. = 0.756

Matching Subtest Summary Statistics

Reference Group: Mean = 16.45

Standard deviation = 4.59

Focal Group: Mean = 15.01

Standard deviation = 4.76

Standardized Score Difference = 0.31

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.

-0.020 0.021 0.348945

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$

OUTPUT FOR RUN NUMBER 3 OUTPUT FOR RUN NUMBER 3

Suspect subtest items:

3

Matching subtest items:

1 2 4 5 6 7 8 9 10 11

12 13 14 15 16 17 18 19 20 21

22 23 24 25 26 27 28 29 30

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.018

proportion of Focal grp. examinees eliminated = 0.036

KR-20 for Ref. grp. = 0.725

KR-20 for Foc. grp. = 0.758

Matching Subtest Summary Statistics

Reference Group: Mean = 16.13

Standard deviation = 4.53

Focal Group: Mean = 14.78

Standard deviation = 4.73

Standardized Score Difference = 0.29

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.

0.037	0.022	0.094296
-------	-------	----------

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$

OUTPUT FOR RUN NUMBER 4 OUTPUT FOR RUN NUMBER 4

Suspect subtest items:

4

Matching subtest items:

1	2	3	5	6	7	8	9	10	11
12	13	14	15	16	17	18	19	20	21
22	23	24	25	26	27	28	29	30	

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.020

proportion of Focal grp. examinees eliminated = 0.040

KR-20 for Ref. grp. = 0.725

KR-20 for Foc. grp. = 0.756

Matching Subtest Summary Statistics

Reference Group: Mean = 16.00

Standard deviation = 4.55

Focal Group: Mean = 14.60

Standard deviation = 4.73

Standardized Score Difference = 0.30

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.

0.010	0.020	0.629005
-------	-------	----------

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$

OUTPUT FOR RUN NUMBER 5 OUTPUT FOR RUN NUMBER 5

Suspect subtest items:

5

Matching subtest items:

1	2	3	4	6	7	8	9	10	11
12	13	14	15	16	17	18	19	20	21
22	23	24	25	26	27	28	29	30	

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.017

proportion of Focal grp. examinees eliminated = 0.037

KR-20 for Ref. grp. = 0.724

KR-20 for Foc. grp. = 0.755

Matching Subtest Summary Statistics

Reference Group: Mean = 16.23

 Standard deviation = 4.52

Focal Group: Mean = 14.73

 Standard deviation = 4.71

Standardized Score Difference = 0.33

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.

-0.120	0.022	0.000000
--------	-------	----------

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$

OUTPUT FOR RUN NUMBER 6 OUTPUT FOR RUN NUMBER 6

Suspect subtest items:

6

Matching subtest items:

1 2 3 4 5 7 8 9 10 11
12 13 14 15 16 17 18 19 20 21
22 23 24 25 26 27 28 29 30

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.020

proportion of Focal grp. examinees eliminated = 0.040

KR-20 for Ref. grp. = 0.721

KR-20 for Foc. grp. = 0.755

Matching Subtest Summary Statistics

Reference Group: Mean = 16.02

 Standard deviation = 4.52

Focal Group: Mean = 14.60

 Standard deviation = 4.72

Standardized Score Difference = 0.31

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.

-0.022 0.020 0.278894

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$

OUTPUT FOR RUN NUMBER 7

OUTPUT FOR RUN NUMBER 7

Suspect subtest items:

7

Matching subtest items:

1 2 3 4 5 6 8 9 10 11
12 13 14 15 16 17 18 19 20 21
22 23 24 25 26 27 28 29 30

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.018

proportion of Focal grp. examinees eliminated = 0.039

KR-20 for Ref. grp. = 0.716

KR-20 for Foc. grp. = 0.750

Matching Subtest Summary Statistics

Reference Group: Mean = 16.03

Standard deviation = 4.47

Focal Group: Mean = 14.68

Standard deviation = 4.67

Standardized Score Difference = 0.30

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.

0.008 0.021 0.679891

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$

OUTPUT FOR RUN NUMBER 8

OUTPUT FOR RUN NUMBER 8

Suspect subtest items:

8

Matching subtest items:

1 2 3 4 5 6 7 9 10 11
 12 13 14 15 16 17 18 19 20 21
 22 23 24 25 26 27 28 29 30

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.016

proportion of Focal grp. examinees eliminated = 0.035

KR-20 for Ref. grp. = 0.721

KR-20 for Foc. grp. = 0.756

Matching Subtest Summary Statistics

Reference Group: Mean = 16.27

Standard deviation = 4.50

Focal Group: Mean = 14.87

Standard deviation = 4.72

Standardized Score Difference = 0.31

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.

-0.017 0.022 0.460945

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$

OUTPUT FOR RUN NUMBER 9 OUTPUT FOR RUN NUMBER 9

Suspect subtest items:

9

Matching subtest items:

1 2 3 4 5 6 7 8 10 11
 12 13 14 15 16 17 18 19 20 21
 22 23 24 25 26 27 28 29 30

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.020

proportion of Focal grp. examinees eliminated = 0.041

KR-20 for Ref. grp. = 0.724

KR-20 for Foc. grp. = 0.758

Matching Subtest Summary Statistics

Reference Group: Mean = 15.95

Standard deviation = 4.55

Focal Group: Mean = 14.58

Standard deviation = 4.75

Standardized Score Difference = 0.30

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.

0.031 0.019 0.103676

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$

OUTPUT FOR RUN NUMBER 10 OUTPUT FOR RUN NUMBER 10

Suspect subtest items:

10

Matching subtest items:

1 2 3 4 5 6 7 8 9 11
 12 13 14 15 16 17 18 19 20 21

22 23 24 25 26 27 28 29 30

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.016

proportion of Focal grp. examinees eliminated = 0.040

KR-20 for Ref. grp. = 0.727

KR-20 for Foc. grp. = 0.761

Matching Subtest Summary Statistics

Reference Group: Mean = 16.35

Standard deviation = 4.56

Focal Group: Mean = 14.93

Standard deviation = 4.78

Standardized Score Difference = 0.30

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.

-0.006 0.023 0.781364

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$

OUTPUT FOR RUN NUMBER 11 OUTPUT FOR RUN NUMBER 11

Suspect subtest items:

11

Matching subtest items:

1 2 3 4 5 6 7 8 9 10

12 13 14 15 16 17 18 19 20 21

22 23 24 25 26 27 28 29 30

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.021
 proportion of Focal grp. examinees eliminated = 0.041
 KR-20 for Ref. grp. = 0.723
 KR-20 for Foc. grp. = 0.759

Matching Subtest Summary Statistics

Reference Group:	Mean = 16.05
	Standard deviation = 4.52
Focal Group:	Mean = 14.66
	Standard deviation = 4.74

Standardized Score Difference = 0.30

SIBTEST-Focal weighting Results

p-value

Beta estimate	standard error	for DIF against either grp.
0.004	0.021	0.856078

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$

OUTPUT FOR RUN NUMBER 12 OUTPUT FOR RUN NUMBER 12

Suspect subtest items:

12

Matching subtest items:

1	2	3	4	5	6	7	8	9	10
11	13	14	15	16	17	18	19	20	21
22	23	24	25	26	27	28	29	30	

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.017
 proportion of Focal grp. examinees eliminated = 0.032

KR-20 for Ref. grp. = 0.716

KR-20 for Foc. grp. = 0.751

Matching Subtest Summary Statistics

Reference Group: Mean = 16.17

Standard deviation = 4.46

Focal Group: Mean = 14.78

Standard deviation = 4.68

Standardized Score Difference = 0.31

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.

-0.033	0.022	0.131220
--------	-------	----------

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$

OUTPUT FOR RUN NUMBER 13 OUTPUT FOR RUN NUMBER 13

Suspect subtest items:

13

Matching subtest items:

1	2	3	4	5	6	7	8	9	10
11	12	14	15	16	17	18	19	20	21
22	23	24	25	26	27	28	29	30	

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.022

proportion of Focal grp. examinees eliminated = 0.040

KR-20 for Ref. grp. = 0.730

KR-20 for Foc. grp. = 0.759

Matching Subtest Summary Statistics

Reference Group: Mean = 16.34

Standard deviation = 4.58

Focal Group: Mean = 14.93

Standard deviation = 4.77

Standardized Score Difference = 0.30

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.

0.004	0.023	0.866235
-------	-------	----------

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$

OUTPUT FOR RUN NUMBER 14

OUTPUT FOR RUN NUMBER 14

Suspect subtest items:

14

Matching subtest items:

1 2 3 4 5 6 7 8 9 10

11 12 13 15 16 17 18 19 20 21

22 23 24 25 26 27 28 29 30

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.016

proportion of Focal grp. examinees eliminated = 0.032

KR-20 for Ref. grp. = 0.719

KR-20 for Foc. grp. = 0.751

Matching Subtest Summary Statistics

Reference Group: Mean = 16.28

Standard deviation = 4.49

Focal Group: Mean = 14.89

Standard deviation = 4.69

Standardized Score Difference = 0.30

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.

-0.016 0.022 0.457467

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$

OUTPUT FOR RUN NUMBER 15 OUTPUT FOR RUN NUMBER 15

Suspect subtest items:

15

Matching subtest items:

1 2 3 4 5 6 7 8 9 10
11 12 13 14 16 17 18 19 20 21
22 23 24 25 26 27 28 29 30

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.020

proportion of Focal grp. examinees eliminated = 0.032

KR-20 for Ref. grp. = 0.729

KR-20 for Foc. grp. = 0.757

Matching Subtest Summary Statistics

Reference Group: Mean = 16.32

Standard deviation = 4.57

Focal Group: Mean = 14.93

Standard deviation = 4.75

Standardized Score Difference = 0.30

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.

0.034	0.023	0.137777
-------	-------	----------

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$

OUTPUT FOR RUN NUMBER 16 OUTPUT FOR RUN NUMBER 16

Suspect subtest items:

16

Matching subtest items:

1	2	3	4	5	6	7	8	9	10
11	12	13	14	15	17	18	19	20	21
22	23	24	25	26	27	28	29	30	

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.015

proportion of Focal grp. examinees eliminated = 0.036

KR-20 for Ref. grp. = 0.722

KR-20 for Foc. grp. = 0.754

Matching Subtest Summary Statistics

Reference Group: Mean = 16.24

Standard deviation = 4.51

Focal Group: Mean = 14.89

Standard deviation = 4.71

Standardized Score Difference = 0.29

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.

0.040 0.022 0.071020

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$

OUTPUT FOR RUN NUMBER 17

OUTPUT FOR RUN NUMBER 17

Suspect subtest items:

17

Matching subtest items:

1 2 3 4 5 6 7 8 9 10
11 12 13 14 15 16 18 19 20 21
22 23 24 25 26 27 28 29 30

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.018

proportion of Focal grp. examinees eliminated = 0.041

KR-20 for Ref. grp. = 0.716

KR-20 for Foc. grp. = 0.752

Matching Subtest Summary Statistics

Reference Group: Mean = 16.04

Standard deviation = 4.48

Focal Group: Mean = 14.69

Standard deviation = 4.68

Standardized Score Difference = 0.29

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.
0.021 0.021 0.307216

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$
OUTPUT FOR RUN NUMBER 18 OUTPUT FOR RUN NUMBER 18

Suspect subtest items:

18

Matching subtest items:

1 2 3 4 5 6 7 8 9 10
11 12 13 14 15 16 17 19 20 21
22 23 24 25 26 27 28 29 30

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.017

proportion of Focal grp. examinees eliminated = 0.033

KR-20 for Ref. grp. = 0.719

KR-20 for Foc. grp. = 0.755

Matching Subtest Summary Statistics

Reference Group: Mean = 16.30

Standard deviation = 4.49

Focal Group: Mean = 14.91

Standard deviation = 4.73

Standardized Score Difference = 0.30

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.
 0.001 0.022 0.979689

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$
 OUTPUT FOR RUN NUMBER 19 OUTPUT FOR RUN NUMBER 19

Suspect subtest items:

19

Matching subtest items:

1 2 3 4 5 6 7 8 9 10
 11 12 13 14 15 16 17 18 20 21
 22 23 24 25 26 27 28 29 30

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.018

proportion of Focal grp. examinees eliminated = 0.038

KR-20 for Ref. grp. = 0.718

KR-20 for Foc. grp. = 0.754

Matching Subtest Summary Statistics

Reference Group: Mean = 15.95

 Standard deviation = 4.51

Focal Group: Mean = 14.58

 Standard deviation = 4.71

Standardized Score Difference = 0.30

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.

0.014	0.019	0.471486
-------	-------	----------

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$

OUTPUT FOR RUN NUMBER 20 OUTPUT FOR RUN NUMBER 20

Suspect subtest items:

20

Matching subtest items:

1	2	3	4	5	6	7	8	9	10
11	12	13	14	15	16	17	18	19	21
22	23	24	25	26	27	28	29	30	

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.021

proportion of Focal grp. examinees eliminated = 0.037

KR-20 for Ref. grp. = 0.719

KR-20 for Foc. grp. = 0.757

Matching Subtest Summary Statistics

Reference Group: Mean = 16.08

 Standard deviation = 4.49

Focal Group: Mean = 14.69

 Standard deviation = 4.72

Standardized Score Difference = 0.30

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.

-0.023	0.022	0.290775
--------	-------	----------

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$

OUTPUT FOR RUN NUMBER 21

OUTPUT FOR RUN NUMBER 21

Suspect subtest items:

21

Matching subtest items:

1 2 3 4 5 6 7 8 9 10
11 12 13 14 15 16 17 18 19 20
22 23 24 25 26 27 28 29 30

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.018

proportion of Focal grp. examinees eliminated = 0.036

KR-20 for Ref. grp. = 0.721

KR-20 for Foc. grp. = 0.756

Matching Subtest Summary Statistics

Reference Group: Mean = 16.25

Standard deviation = 4.51

Focal Group: Mean = 14.85

Standard deviation = 4.73

Standardized Score Difference = 0.30

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.

-0.011 0.022 0.640146

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$\$\$\$\$\$\$\$\$\$\$\$\$\$\$\$\$\$
 OUTPUT FOR RUN NUMBER 22 OUTPUT FOR RUN NUMBER 22

Suspect subtest items:

22

Matching subtest items:

1 2 3 4 5 6 7 8 9 10
 11 12 13 14 15 16 17 18 19 20
 21 23 24 25 26 27 28 29 30

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.018

proportion of Focal grp. examinees eliminated = 0.034

KR-20 for Ref. grp. = 0.721

KR-20 for Foc. grp. = 0.753

Matching Subtest Summary Statistics

Reference Group: Mean = 16.51

 Standard deviation = 4.56

Focal Group: Mean = 15.07

 Standard deviation = 4.76

Standardized Score Difference = 0.31

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.

-0.026 0.019 0.173378

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$\$\$\$\$\$\$\$\$\$\$\$\$\$\$\$\$\$
 OUTPUT FOR RUN NUMBER 23 OUTPUT FOR RUN NUMBER 23

Suspect subtest items:

23

Matching subtest items:

1 2 3 4 5 6 7 8 9 10
 11 12 13 14 15 16 17 18 19 20
 21 22 24 25 26 27 28 29 30

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.019

proportion of Focal grp. examinees eliminated = 0.037

KR-20 for Ref. grp. = 0.726

KR-20 for Foc. grp. = 0.761

Matching Subtest Summary Statistics

Reference Group: Mean = 16.43

Standard deviation = 4.57

Focal Group: Mean = 15.01

Standard deviation = 4.80

Standardized Score Difference = 0.30

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.

-0.003 0.021 0.885020

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$

OUTPUT FOR RUN NUMBER 24 OUTPUT FOR RUN NUMBER 24

Suspect subtest items:

24

Matching subtest items:

1 2 3 4 5 6 7 8 9 10
 11 12 13 14 15 16 17 18 19 20
 21 22 23 25 26 27 28 29 30

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.017

proportion of Focal grp. examinees eliminated = 0.035

KR-20 for Ref. grp. = 0.725

KR-20 for Foc. grp. = 0.756

Matching Subtest Summary Statistics

Reference Group: Mean = 16.25

Standard deviation = 4.53

Focal Group: Mean = 14.90

Standard deviation = 4.73

Standardized Score Difference = 0.29

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.

0.050 0.023 0.028911

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$

OUTPUT FOR RUN NUMBER 25 OUTPUT FOR RUN NUMBER 25

Suspect subtest items:

25

Matching subtest items:

1 2 3 4 5 6 7 8 9 10

11 12 13 14 15 16 17 18 19 20

21 22 23 24 26 27 28 29 30

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.019

proportion of Focal grp. examinees eliminated = 0.038

KR-20 for Ref. grp. = 0.729

KR-20 for Foc. grp. = 0.763

Matching Subtest Summary Statistics

Reference Group: Mean = 16.40

Standard deviation = 4.58

Focal Group: Mean = 14.97

Standard deviation = 4.81

Standardized Score Difference = 0.30

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.

-0.012 0.022 0.580680

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$

OUTPUT FOR RUN NUMBER 26 OUTPUT FOR RUN NUMBER 26

Suspect subtest items:

26

Matching subtest items:

1 2 3 4 5 6 7 8 9 10

11 12 13 14 15 16 17 18 19 20

21 22 23 24 25 27 28 29 30

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.018

proportion of Focal grp. examinees eliminated = 0.043

KR-20 for Ref. grp. = 0.734

KR-20 for Foc. grp. = 0.762

Matching Subtest Summary Statistics

Reference Group: Mean = 16.28

Standard deviation = 4.60

Focal Group: Mean = 14.87

Standard deviation = 4.78

Standardized Score Difference = 0.30

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.

0.022	0.023	0.334398
-------	-------	----------

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$

OUTPUT FOR RUN NUMBER 27 OUTPUT FOR RUN NUMBER 27

Suspect subtest items:

27

Matching subtest items:

1 2 3 4 5 6 7 8 9 10

11 12 13 14 15 16 17 18 19 20

21 22 23 24 25 26 28 29 30

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.017

proportion of Focal grp. examinees eliminated = 0.038

KR-20 for Ref. grp. = 0.721

KR-20 for Foc. grp. = 0.758

Matching Subtest Summary Statistics

Reference Group: Mean = 16.32

Standard deviation = 4.51

Focal Group: Mean = 14.94

Standard deviation = 4.76

Standardized Score Difference = 0.30

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.

0.019 0.022 0.378562

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$

OUTPUT FOR RUN NUMBER 28 OUTPUT FOR RUN NUMBER 28

Suspect subtest items:

28

Matching subtest items:

1 2 3 4 5 6 7 8 9 10

11 12 13 14 15 16 17 18 19 20

21 22 23 24 25 26 27 29 30

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.018

proportion of Focal grp. examinees eliminated = 0.036

KR-20 for Ref. grp. = 0.720

KR-20 for Foc. grp. = 0.755

Matching Subtest Summary Statistics

Reference Group:	Mean = 16.18
	Standard deviation = 4.49
Focal Group:	Mean = 14.80
	Standard deviation = 4.71

Standardized Score Difference = 0.30

SIBTEST-Focal weighting Results

p-value

Beta estimate	standard error	for DIF against either grp.
0.015	0.023	0.512855

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$

OUTPUT FOR RUN NUMBER 29 OUTPUT FOR RUN NUMBER 29

Suspect subtest items:

29

Matching subtest items:

1 2 3 4 5 6 7 8 9 10
 11 12 13 14 15 16 17 18 19 20
 21 22 23 24 25 26 27 28 30

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.015

proportion of Focal grp. examinees eliminated = 0.032

KR-20 for Ref. grp. = 0.723

KR-20 for Foc. grp. = 0.755

Matching Subtest Summary Statistics

Reference Group: Mean = 16.28

 Standard deviation = 4.52

Focal Group: Mean = 14.94

 Standard deviation = 4.73

Standardized Score Difference = 0.29

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.

0.053 0.022 0.018511

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$

OUTPUT FOR RUN NUMBER 30 OUTPUT FOR RUN NUMBER 30

Suspect subtest items:

30

Matching subtest items:

1 2 3 4 5 6 7 8 9 10

11 12 13 14 15 16 17 18 19 20

21 22 23 24 25 26 27 28 29

estimate of guessing on this matching subtest = 0.20

proportion of Ref. grp. examinees eliminated = 0.018

proportion of Focal grp. examinees eliminated = 0.034

KR-20 for Ref. grp. = 0.720

KR-20 for Foc. grp. = 0.751

Matching Subtest Summary Statistics

Reference Group: Mean = 16.20

Standard deviation = 4.49

Focal Group: Mean = 14.76

Standard deviation = 4.67

Standardized Score Difference = 0.31

SIBTEST-Focal weighting Results

p-value

Beta estimate standard error for DIF against either grp.

-0.073	0.022	0.000939
--------	-------	----------

SIBTEST error flag = 0

No errors. This SIBTEST run had a normal successful completion.

\$

SUMMARY OF THE RUNS

p-value notation:

R denotes p-value for test of DIF/DBF against Ref. group

F denotes p-value for test of DIF/DBF against Foc. group

E denotes p-value for test of DIF/DBF against either the
Ref. or Foc. group.

NOTES:

MS/SSD = Matching Subtest Standardized Score Difference.

Standardized difference in mean observed scores
between Reference group and Focal group on the
matching subtest.

p-elim = proportion of Reference (R) and Focal (F) groups
eliminated (not used) in SIBTEST calculations.

Positive Beta estimate indicates DIF/DBF favoring Ref. grp.

Negative Beta estimate indicates DIF/DBF favoring Foc. grp.

FLAG = error flag indicator. FLAG=0 indicates a normal successful completion of a SIBTEST run. All other values of FLAG come with short error messages.

SIBTEST-Focal weighting							F
Run	Suspect	Subtest	Beta	standard	p-elim	MS A	L
no.	Item	Numbers	estimate	error	p-value	R F SSD G	
1	1		-0.034	0.023	0.139 E	.02 .04	0.31 0
2	2		-0.020	0.021	0.349 E	.02 .03	0.31 0
3	3		0.037	0.022	0.094 E	.02 .04	0.29 0
4	4		0.010	0.020	0.629 E	.02 .04	0.30 0
5	5		-0.120	0.022	0.000 E	.02 .04	0.33 0
6	6		-0.022	0.020	0.279 E	.02 .04	0.31 0
7	7		0.008	0.021	0.680 E	.02 .04	0.30 0
8	8		-0.017	0.022	0.461 E	.02 .04	0.31 0
9	9		0.031	0.019	0.104 E	.02 .04	0.30 0
10	10		-0.006	0.023	0.781 E	.02 .04	0.30 0
11	11		0.004	0.021	0.856 E	.02 .04	0.30 0
12	12		-0.033	0.022	0.131 E	.02 .03	0.31 0
13	13		0.004	0.023	0.866 E	.02 .04	0.30 0
14	14		-0.016	0.022	0.457 E	.02 .03	0.30 0
15	15		0.034	0.023	0.138 E	.02 .03	0.30 0
16	16		0.040	0.022	0.071 E	.02 .04	0.29 0
17	17		0.021	0.021	0.307 E	.02 .04	0.29 0
18	18		0.001	0.022	0.980 E	.02 .03	0.30 0
19	19		0.014	0.019	0.471 E	.02 .04	0.30 0
20	20		-0.023	0.022	0.291 E	.02 .04	0.30 0
21	21		-0.011	0.022	0.640 E	.02 .04	0.30 0
22	22		-0.026	0.019	0.173 E	.02 .03	0.31 0
23	23		-0.003	0.021	0.885 E	.02 .04	0.30 0
24	24		0.050	0.023	0.029 E	.02 .04	0.29 0
25	25		-0.012	0.022	0.581 E	.02 .04	0.30 0
26	26		0.022	0.023	0.334 E	.02 .04	0.30 0
27	27		0.019	0.022	0.379 E	.02 .04	0.30 0

28 28	0.015	0.023	0.513 E .02 .04	0.30	0
29 29	0.053	0.022	0.019 E .02 .03	0.29	0
30 30	-0.073	0.022	0.001 E .02 .03	0.31	0

ภาคผนวก ช

**ตัวอย่าง Print Out การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ
ด้วยวิธี MH**

```

GET
FILE='C:\Users\KITTY\Desktop\MH Literacy 300 1.sav'.
DATASET NAME DataSet1 WINDOW=FRONT.

SAVE OUTFILE='C:\Users\KITTY\Desktop\MH Literacy 2000 Literacy
Female.sav'
/COMPRESSED.

CROSSTABS
/TABLES=Gender BY Answer BY id
/FORMAT=AVALUE TABLES
/STATISTICS=CHISQ RISK CMH(1)
/CELLS=COUNT
/COUNT ROUND CELL.

```

Crosstabs

[DataSet1] C:\Users\KITTY\Desktop\MH Literacy 2000 Literacy Female.sav

Case Processing Summary

	Cases					
	Valid		Missing		Total	
	N	Percent	N	Percent	N	Percent
Gender * Answer * id	60000	100.0%	0	0.0%	60000	100.0%

Gender * Answer * id Crosstabulation

Count

id	Gender	Answer		Total
		0	1	
1	1	471	529	1000
	2	459	541	1000
	Total	930	1070	2000
2	1	680	320	1000
	2	678	322	1000
	Total	1358	642	2000
3	1	445	555	1000
	2	355	645	1000
	Total	800	1200	2000
4	1	270	730	1000
	2	219	781	1000
	Total	489	1511	2000
5	1	398	602	1000
	2	455	545	1000

	Total		853	1147	2000
6	Gender	1	272	728	1000
		2	241	759	1000
7	Total		513	1487	2000
	Gender	1	345	655	1000
8		2	252	748	1000
	Total		597	1403	2000
9	Gender	1	534	466	1000
		2	497	503	1000
10	Total		1031	969	2000
	Gender	1	246	754	1000
11		2	177	823	1000
	Total		423	1577	2000
12	Gender	1	600	400	1000
		2	578	422	1000
13	Total		1178	822	2000
	Gender	1	325	675	1000
14		2	274	726	1000
	Total		599	1401	2000
15	Gender	1	444	556	1000
		2	397	603	1000
16	Total		841	1159	2000
	Gender	1	594	406	1000
17		2	561	439	1000
	Total		1155	845	2000
18	Gender	1	557	443	1000
		2	506	494	1000
19	Total		1063	937	2000
	Gender	1	601	399	1000
		2	539	461	1000
	Total		1140	860	2000
	Gender	1	556	444	1000
		2	460	540	1000
	Total		1016	984	2000
	Gender	1	359	641	1000
		2	261	739	1000
	Total		620	1380	2000
	Gender	1	576	424	1000
		2	521	479	1000
	Total		1097	903	2000
	Gender	1	249	751	1000

		2	171	829	1000
	Total		420	1580	2000
20	Gender	1	353	647	1000
		2	301	699	1000
	Total		654	1346	2000
21	Gender	1	520	480	1000
		2	476	524	1000
	Total		996	1004	2000
22	Gender	1	742	258	1000
		2	731	269	1000
	Total		1473	527	2000
23	Gender	1	674	326	1000
		2	652	348	1000
	Total		1326	674	2000
24	Gender	1	566	434	1000
		2	468	532	1000
	Total		1034	966	2000
25	Gender	1	634	366	1000
		2	618	382	1000
	Total		1252	748	2000
26	Gender	1	542	458	1000
		2	506	494	1000
	Total		1048	952	2000
27	Gender	1	613	387	1000
		2	545	455	1000
	Total		1158	842	2000
28	Gender	1	468	532	1000
		2	399	601	1000
	Total		867	1133	2000
29	Gender	1	605	395	1000
		2	500	500	1000
	Total		1105	895	2000
30	Gender	1	429	571	1000
		2	426	574	1000
	Total		855	1145	2000
Total	Gender	1	14668	15332	30000
		2	13223	16777	30000
	Total		27891	32109	60000

Chi-Square Tests

id	Value	df	Asymp. Sig.	Exact Sig. (2-sided)	Exact Sig. (1-sided)
			(2-sided)		
1	Pearson Chi-Square	.289 ^c	1	.591	
	Continuity Correction ^b	.243	1	.622	
	Likelihood Ratio	.289	1	.591	
	Fisher's Exact Test				.622 .311
	Linear-by-Linear Association	.289	1	.591	
	N of Valid Cases	2000			
2	Pearson Chi-Square	.009 ^d	1	.924	
	Continuity Correction ^b	.002	1	.962	
	Likelihood Ratio	.009	1	.924	
	Fisher's Exact Test				.962 .481
	Linear-by-Linear Association	.009	1	.924	
	N of Valid Cases	2000			
3	Pearson Chi-Square	16.875 ^e	1	.000	
	Continuity Correction ^b	16.502	1	.000	
	Likelihood Ratio	16.903	1	.000	
	Fisher's Exact Test				.000 .000
	Linear-by-Linear Association	16.867	1	.000	
	N of Valid Cases	2000			
4	Pearson Chi-Square	7.040 ^f	1	.008	
	Continuity Correction ^b	6.767	1	.009	
	Likelihood Ratio	7.050	1	.008	
	Fisher's Exact Test				.009 .005
	Linear-by-Linear Association	7.037	1	.008	
	N of Valid Cases	2000			
5	Pearson Chi-Square	6.642 ^g	1	.010	
	Continuity Correction ^b	6.411	1	.011	
	Likelihood Ratio	6.646	1	.010	
	Fisher's Exact Test				.011 .006
	Linear-by-Linear Association	6.638	1	.010	
	N of Valid Cases	2000			
6	Pearson Chi-Square	2.520 ^h	1	.112	
	Continuity Correction ^b	2.360	1	.125	
	Likelihood Ratio	2.521	1	.112	

	Fisher's Exact Test				.124	.062
	Linear-by-Linear Association	2.518	1	.113		
	N of Valid Cases	2000				
	Pearson Chi-Square	20.652 ⁱ	1	.000		
	Continuity Correction ^b	20.210	1	.000		
	Likelihood Ratio	20.716	1	.000		
7	Fisher's Exact Test				.000	.000
	Linear-by-Linear Association	20.642	1	.000		
	N of Valid Cases	2000				
	Pearson Chi-Square	2.741 ^j	1	.098		
	Continuity Correction ^b	2.594	1	.107		
	Likelihood Ratio	2.741	1	.098		
8	Fisher's Exact Test				.107	.054
	Linear-by-Linear Association	2.739	1	.098		
	N of Valid Cases	2000				
	Pearson Chi-Square	14.274 ^k	1	.000		
	Continuity Correction ^b	13.864	1	.000		
	Likelihood Ratio	14.326	1	.000		
9	Fisher's Exact Test				.000	.000
	Linear-by-Linear Association	14.267	1	.000		
	N of Valid Cases	2000				
	Pearson Chi-Square	1.000 ^l	1	.317		
	Continuity Correction ^b	.911	1	.340		
	Likelihood Ratio	1.000	1	.317		
10	Fisher's Exact Test				.340	.170
	Linear-by-Linear Association	.999	1	.318		
	N of Valid Cases	2000				
	Pearson Chi-Square	6.199 ^m	1	.013		
	Continuity Correction ^b	5.958	1	.015		
	Likelihood Ratio	6.204	1	.013		
11	Fisher's Exact Test				.015	.007
	Linear-by-Linear Association	6.196	1	.013		
	N of Valid Cases	2000				
12	Pearson Chi-Square	4.533 ⁿ	1	.033		
	Continuity Correction ^b	4.342	1	.037		

	Likelihood Ratio	4.534	1	.033		
	Fisher's Exact Test				.037	.019
	Linear-by-Linear					
	Association	4.530	1	.033		
	N of Valid Cases	2000				
	Pearson Chi-Square	2.232 ^o	1	.135		
	Continuity Correction ^b	2.098	1	.147		
	Likelihood Ratio	2.232	1	.135		
13	Fisher's Exact Test				.147	.074
	Linear-by-Linear					
	Association	2.230	1	.135		
	N of Valid Cases	2000				
	Pearson Chi-Square	5.223 ^p	1	.022		
	Continuity Correction ^b	5.020	1	.025		
	Likelihood Ratio	5.225	1	.022		
14	Fisher's Exact Test				.025	.013
	Linear-by-Linear					
	Association	5.220	1	.022		
	N of Valid Cases	2000				
	Pearson Chi-Square	7.842 ^q	1	.005		
	Continuity Correction ^b	7.591	1	.006		
	Likelihood Ratio	7.847	1	.005		
15	Fisher's Exact Test				.006	.003
	Linear-by-Linear					
	Association	7.838	1	.005		
	N of Valid Cases	2000				
	Pearson Chi-Square	18.437 ^r	1	.000		
	Continuity Correction ^b	18.055	1	.000		
	Likelihood Ratio	18.465	1	.000		
16	Fisher's Exact Test				.000	.000
	Linear-by-Linear					
	Association	18.428	1	.000		
	N of Valid Cases	2000				
	Pearson Chi-Square	22.450 ^s	1	.000		
	Continuity Correction ^b	21.994	1	.000		
	Likelihood Ratio	22.521	1	.000		
17	Fisher's Exact Test				.000	.000
	Linear-by-Linear					
	Association	22.439	1	.000		
	N of Valid Cases	2000				
18	Pearson Chi-Square	6.107 ^t	1	.013		

	Continuity Correction ^b	5.887	1	.015		
	Likelihood Ratio	6.111	1	.013		
	Fisher's Exact Test				.015	.008
	Linear-by-Linear					
	Association	6.104	1	.013		
	N of Valid Cases	2000				
	Pearson Chi-Square	18.336 ^u	1	.000		
	Continuity Correction ^b	17.869	1	.000		
	Likelihood Ratio	18.422	1	.000		
19	Fisher's Exact Test				.000	.000
	Linear-by-Linear					
	Association	18.327	1	.000		
	N of Valid Cases	2000				
	Pearson Chi-Square	6.143 ^v	1	.013		
	Continuity Correction ^b	5.909	1	.015		
	Likelihood Ratio	6.148	1	.013		
20	Fisher's Exact Test				.015	.008
	Linear-by-Linear					
	Association	6.140	1	.013		
	N of Valid Cases	2000				
	Pearson Chi-Square	3.872 ^w	1	.049		
	Continuity Correction ^b	3.698	1	.054		
	Likelihood Ratio	3.873	1	.049		
21	Fisher's Exact Test				.054	.027
	Linear-by-Linear					
	Association	3.870	1	.049		
	N of Valid Cases	2000				
	Pearson Chi-Square	.312 ^x	1	.577		
	Continuity Correction ^b	.258	1	.612		
	Likelihood Ratio	.312	1	.577		
22	Fisher's Exact Test				.612	.306
	Linear-by-Linear					
	Association	.312	1	.577		
	N of Valid Cases	2000				
	Pearson Chi-Square	1.083 ^y	1	.298		
	Continuity Correction ^b	.987	1	.321		
	Likelihood Ratio	1.083	1	.298		
23	Fisher's Exact Test				.321	.160
	Linear-by-Linear					
	Association	1.083	1	.298		
	N of Valid Cases	2000				

	Pearson Chi-Square	19.230 ^z	1	.000		
	Continuity Correction ^b	18.840	1	.000		
	Likelihood Ratio	19.261	1	.000		
24	Fisher's Exact Test				.000	.000
	Linear-by-Linear Association	19.221	1	.000		
	N of Valid Cases	2000				
	Pearson Chi-Square	.547 ^{aa}	1	.460		
	Continuity Correction ^b	.481	1	.488		
	Likelihood Ratio	.547	1	.460		
25	Fisher's Exact Test				.488	.244
	Linear-by-Linear Association	.546	1	.460		
	N of Valid Cases	2000				
	Pearson Chi-Square	2.598 ^{ab}	1	.107		
	Continuity Correction ^b	2.456	1	.117		
	Likelihood Ratio	2.599	1	.107		
26	Fisher's Exact Test				.117	.059
	Linear-by-Linear Association	2.597	1	.107		
	N of Valid Cases	2000				
	Pearson Chi-Square	9.485 ^{ac}	1	.002		
	Continuity Correction ^b	9.208	1	.002		
	Likelihood Ratio	9.493	1	.002		
27	Fisher's Exact Test				.002	.001
	Linear-by-Linear Association	9.480	1	.002		
	N of Valid Cases	2000				
	Pearson Chi-Square	9.693 ^{ad}	1	.002		
	Continuity Correction ^b	9.415	1	.002		
	Likelihood Ratio	9.702	1	.002		
28	Fisher's Exact Test				.002	.001
	Linear-by-Linear Association	9.689	1	.002		
	N of Valid Cases	2000				
	Pearson Chi-Square	22.296 ^{ae}	1	.000		
	Continuity Correction ^b	21.873	1	.000		
	Likelihood Ratio	22.339	1	.000		
29	Fisher's Exact Test				.000	.000
	Linear-by-Linear Association	22.285	1	.000		

	N of Valid Cases	2000				
	Pearson Chi-Square	.018 ^{af}	1	.892		
	Continuity Correction ^b	.008	1	.928		
	Likelihood Ratio	.018	1	.892		
30	Fisher's Exact Test				.928	.464
	Linear-by-Linear Association	.018	1	.892		
	N of Valid Cases	2000				
	Pearson Chi-Square	139.893 ^a	1	.000		
	Continuity Correction ^b	139.699	1	.000		
	Likelihood Ratio	139.949	1	.000		
Total	Fisher's Exact Test				.000	.000
	Linear-by-Linear Association	139.891	1	.000		
	N of Valid Cases	60000				

a. 0 cells (.0%) have expected count less than 5. The minimum expected count is 13945.50.

b. Computed only for a 2x2 table

c. 0 cells (.0%) have expected count less than 5. The minimum expected count is 465.00.

d. 0 cells (.0%) have expected count less than 5. The minimum expected count is 321.00.

e. 0 cells (.0%) have expected count less than 5. The minimum expected count is 400.00.

f. 0 cells (.0%) have expected count less than 5. The minimum expected count is 244.50.

g. 0 cells (.0%) have expected count less than 5. The minimum expected count is 426.50.

h. 0 cells (.0%) have expected count less than 5. The minimum expected count is 256.50.

i. 0 cells (.0%) have expected count less than 5. The minimum expected count is 298.50.

j. 0 cells (.0%) have expected count less than 5. The minimum expected count is 484.50.

k. 0 cells (.0%) have expected count less than 5. The minimum expected count is 211.50.

l. 0 cells (.0%) have expected count less than 5. The minimum expected count is 411.00.

m. 0 cells (.0%) have expected count less than 5. The minimum expected count is 299.50.

n. 0 cells (.0%) have expected count less than 5. The minimum expected count is 420.50.

o. 0 cells (.0%) have expected count less than 5. The minimum expected count is 422.50.

p. 0 cells (.0%) have expected count less than 5. The minimum expected count is 468.50.

q. 0 cells (.0%) have expected count less than 5. The minimum expected count is 430.00.

r. 0 cells (.0%) have expected count less than 5. The minimum expected count is 492.00.

s. 0 cells (.0%) have expected count less than 5. The minimum expected count is 310.00.

t. 0 cells (.0%) have expected count less than 5. The minimum expected count is 451.50.

u. 0 cells (.0%) have expected count less than 5. The minimum expected count is 210.00.

v. 0 cells (.0%) have expected count less than 5. The minimum expected count is 327.00.

w. 0 cells (.0%) have expected count less than 5. The minimum expected count is 498.00.

x. 0 cells (.0%) have expected count less than 5. The minimum expected count is 263.50.

y. 0 cells (.0%) have expected count less than 5. The minimum expected count is 337.00.

- z. 0 cells (.0%) have expected count less than 5. The minimum expected count is 483.00.
 aa. 0 cells (.0%) have expected count less than 5. The minimum expected count is 374.00.
 ab. 0 cells (.0%) have expected count less than 5. The minimum expected count is 476.00.
 ac. 0 cells (.0%) have expected count less than 5. The minimum expected count is 421.00.
 ad. 0 cells (.0%) have expected count less than 5. The minimum expected count is 433.50.
 ae. 0 cells (.0%) have expected count less than 5. The minimum expected count is 447.50.
 af. 0 cells (.0%) have expected count less than 5. The minimum expected count is 427.50.

Risk Estimate

id	Value	95% Confidence Interval		
		Lower	Upper	
	Odds Ratio for Gender (1 / 2)	1.049	.880	1.251
1	For cohort Answer = 0	1.026	.934	1.127
	For cohort Answer = 1	.978	.901	1.061
	N of Valid Cases	2000		
	Odds Ratio for Gender (1 / 2)	1.009	.836	1.218
2	For cohort Answer = 0	1.003	.944	1.065
	For cohort Answer = 1	.994	.875	1.129
	N of Valid Cases	2000		
	Odds Ratio for Gender (1 / 2)	1.457	1.217	1.744
3	For cohort Answer = 0	1.254	1.125	1.397
	For cohort Answer = 1	.860	.801	.925
	N of Valid Cases	2000		
	Odds Ratio for Gender (1 / 2)	1.319	1.075	1.619
4	For cohort Answer = 0	1.233	1.056	1.440
	For cohort Answer = 1	.935	.889	.983
	N of Valid Cases	2000		
	Odds Ratio for Gender (1 / 2)	.792	.663	.946
5	For cohort Answer = 0	.875	.790	.969
	For cohort Answer = 1	1.105	1.024	1.192
	N of Valid Cases	2000		
	Odds Ratio for Gender (1 / 2)	1.177	.962	1.439
6	For cohort Answer = 0	1.129	.972	1.311
	For cohort Answer = 1	.959	.911	1.010

	N of Valid Cases	2000		
	Odds Ratio for Gender (1 / 2)	1.563	1.289	1.897
7	For cohort Answer = 0	1.369	1.194	1.570
	For cohort Answer = 1	.876	.827	.928
	N of Valid Cases	2000		
	Odds Ratio for Gender (1 / 2)	1.160	.973	1.382
8	For cohort Answer = 0	1.074	.987	1.170
	For cohort Answer = 1	.926	.846	1.014
	N of Valid Cases	2000		
	Odds Ratio for Gender (1 / 2)	1.517	1.221	1.885
9	For cohort Answer = 0	1.390	1.170	1.651
	For cohort Answer = 1	.916	.875	.959
	N of Valid Cases	2000		
	Odds Ratio for Gender (1 / 2)	1.095	.916	1.309
10	For cohort Answer = 0	1.038	.965	1.117
	For cohort Answer = 1	.948	.853	1.053
	N of Valid Cases	2000		
	Odds Ratio for Gender (1 / 2)	1.276	1.053	1.546
11	For cohort Answer = 0	1.186	1.037	1.357
	For cohort Answer = 1	.930	.878	.985
	N of Valid Cases	2000		
	Odds Ratio for Gender (1 / 2)	1.213	1.015	1.449
12	For cohort Answer = 0	1.118	1.009	1.240
	For cohort Answer = 1	.922	.856	.994
	N of Valid Cases	2000		
	Odds Ratio for Gender (1 / 2)	1.145	.959	1.367
13	For cohort Answer = 0	1.059	.982	1.141
	For cohort Answer = 1	.925	.835	1.025
	N of Valid Cases	2000		
	Odds Ratio for Gender (1 / 2)	1.228	1.030	1.464
14	For cohort Answer = 0	1.101	1.014	1.195
	For cohort Answer = 1	.897	.817	.985
	N of Valid Cases	2000		

	Odds Ratio for Gender (1 / 2)	1.288	1.079	1.538
15	For cohort Answer = 0	1.115	1.033	1.204
	For cohort Answer = 1	.866	.782	.958
	N of Valid Cases	2000		
	Odds Ratio for Gender (1 / 2)	1.470	1.233	1.753
16	For cohort Answer = 0	1.209	1.108	1.319
	For cohort Answer = 1	.822	.752	.900
	N of Valid Cases	2000		
	Odds Ratio for Gender (1 / 2)	1.586	1.310	1.920
17	For cohort Answer = 0	1.375	1.204	1.571
	For cohort Answer = 1	.867	.818	.920
	N of Valid Cases	2000		
	Odds Ratio for Gender (1 / 2)	1.249	1.047	1.490
18	For cohort Answer = 0	1.106	1.021	1.197
	For cohort Answer = 1	.885	.803	.975
	N of Valid Cases	2000		
	Odds Ratio for Gender (1 / 2)	1.607	1.292	1.999
19	For cohort Answer = 0	1.456	1.224	1.733
	For cohort Answer = 1	.906	.866	.948
	N of Valid Cases	2000		
	Odds Ratio for Gender (1 / 2)	1.267	1.051	1.528
20	For cohort Answer = 0	1.173	1.034	1.331
	For cohort Answer = 1	.926	.871	.984
	N of Valid Cases	2000		
	Odds Ratio for Gender (1 / 2)	1.193	1.001	1.421
21	For cohort Answer = 0	1.092	1.000	1.193
	For cohort Answer = 1	.916	.839	1.000
	N of Valid Cases	2000		
	Odds Ratio for Gender (1 / 2)	1.058	.867	1.291
22	For cohort Answer = 0	1.015	.963	1.070
	For cohort Answer = 1	.959	.828	1.111
	N of Valid Cases	2000		

	Odds Ratio for Gender (1 / 2)	1.104	.917	1.328
23	For cohort Answer = 0	1.034	.971	1.100
	For cohort Answer = 1	.937	.828	1.059
	N of Valid Cases	2000		
	Odds Ratio for Gender (1 / 2)	1.482	1.243	1.768
24	For cohort Answer = 0	1.209	1.110	1.317
	For cohort Answer = 1	.816	.744	.894
	N of Valid Cases	2000		
	Odds Ratio for Gender (1 / 2)	1.071	.893	1.283
25	For cohort Answer = 0	1.026	.959	1.098
	For cohort Answer = 1	.958	.855	1.073
	N of Valid Cases	2000		
	Odds Ratio for Gender (1 / 2)	1.155	.969	1.377
26	For cohort Answer = 0	1.071	.985	1.165
	For cohort Answer = 1	.927	.846	1.017
	N of Valid Cases	2000		
	Odds Ratio for Gender (1 / 2)	1.322	1.107	1.580
27	For cohort Answer = 0	1.125	1.043	1.212
	For cohort Answer = 1	.851	.767	.943
	N of Valid Cases	2000		
	Odds Ratio for Gender (1 / 2)	1.325	1.110	1.582
28	For cohort Answer = 0	1.173	1.061	1.297
	For cohort Answer = 1	.885	.820	.956
	N of Valid Cases	2000		
	Odds Ratio for Gender (1 / 2)	1.532	1.283	1.829
29	For cohort Answer = 0	1.210	1.117	1.310
	For cohort Answer = 1	.790	.716	.872
	N of Valid Cases	2000		
	Odds Ratio for Gender (1 / 2)	1.012	.848	1.209
30	For cohort Answer = 0	1.007	.910	1.115
	For cohort Answer = 1	.995	.922	1.073
	N of Valid Cases	2000		

Odds Ratio for Gender (1 / 2)	1.214	1.175	1.253
Total For cohort Answer = 0	1.109	1.090	1.129
For cohort Answer = 1	.914	.900	.928
N of Valid Cases	60000		

Tests of Homogeneity of the Odds Ratio

	Chi-Squared	df	Asymp. Sig. (2-sided)
Breslow-Day	86.936	29	.000
Tarone's	86.936	29	.000

Tests of Conditional Independence

	Chi-Squared	df	Asymp. Sig. (2-sided)
Cochran's	152.058	1	.000
Mantel-Haenszel	151.772	1	.000

Under the conditional independence assumption, Cochran's statistic is asymptotically distributed as a 1 df chi-squared distribution, only if the number of strata is fixed, while the Mantel-Haenszel statistic is always asymptotically distributed as a 1 df chi-squared distribution. Note that the continuity correction is removed from the Mantel-Haenszel statistic when the sum of the differences between the observed and the expected is 0.

Mantel-Haenszel Common Odds Ratio Estimate

Estimate			1.234
In(Estimate)			.210
Std. Error of In(Estimate)			.017
Asymp. Sig. (2-sided)			.000
Asymp. 95% Confidence Interval	Common Odds Ratio	Lower Bound	1.194
		Upper Bound	1.276
	In(Common Odds Ratio)	Lower Bound	.177
		Upper Bound	.244

The Mantel-Haenszel common odds ratio estimate is asymptotically normally distributed under the common odds ratio of 1.000 assumption. So is the natural log of the estimate.

ผลงานวิจัย

อรุณี แบลลงกาย และปิยะทิพย์ ประดุจพรม. (2561). การเปรียบเทียบผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ NT ระดับชั้นประถมศึกษาปีที่ 3 ระหว่างวิธีการทดสอบอัตราส่วนไลค์ลิขิต วิธีซิปเกสท์ และวิธีแมนเกล-แyenส์เซล. ใน การประชุมวิชาการระดับชาติครั้งที่ 3 ประจำปี 2561 “นวัตกรรมการจัดการ; งการขับเคลื่อนธุรกิจชุมชนสู่เศรษฐกิจดิจิทัล”, วันที่ 1 มิถุนายน พ.ศ. 2561. (หน้า 720-730). ปทุมธานี: มหาวิทยาลัยราชภัฏวไลยอลงกรณ์ ในพระบรมราชูปถัมภ์.