

บทที่ 2

เอกสารและงานวิจัยที่เกี่ยวข้อง

เอกสารและงานวิจัยที่เกี่ยวกับการทำหน้าที่ต่างกันของข้อสอบ วิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบโดยวิธีชิปเพลส์ ผลสัมฤทธิ์ทางการเรียน และค่าความเที่ยงของแบบทดสอบ ผู้วิจัยได้เสนอเป็น 5 ตอน ดังนี้

ตอนที่ 1 การทำหน้าที่ต่างกันของข้อสอบ

ตอนที่ 2 วิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบโดยวิธีชิปเพลส์

ตอนที่ 3 งานวิจัยที่เกี่ยวข้องกับการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบโดยวิธีชิปเพลส์

ตอนที่ 4 ผลสัมฤทธิ์ทางการเรียน

ตอนที่ 5 ค่าความเที่ยงของแบบทดสอบ

ตอนที่ 1 การทำหน้าที่ต่างกันของข้อสอบ

ความเป็นมาของการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ

การศึกษาความลำเอียงหรือการทำหน้าที่ต่างกันของข้อสอบหรือแบบทดสอบนั้น เริ่มมาจากการสังเกตปรากฏการณ์ของผลการสอบคัดเลือก เช่น การสอบคัดเลือกเพื่อเข้าศึกษาต่อ เพื่อบรรจุเข้าทำงาน หรือการสอบเพื่อเดินตำแหน่ง หรือในการสอบอื่น ๆ ซึ่งพบว่าไม่เป็นไปตามสัดส่วนระดับสูงปัญญาของผู้เข้าสอบหรือโครงการสร้างของประชากร จึงเป็นเหตุให้เกิดความลำเอียง หรือความไม่ยุติธรรมในการสอบและส่งผลให้เกิดการศึกษาความลำเอียงของแบบทดสอบ (Test Bias) และข้อสอบ (Item Bias) ในระยะต่อ ๆ มาแนวคิดในการวิเคราะห์ความลำเอียงและการตรวจสอบความยุติธรรมของการใช้ข้อสอบ ได้เริ่มปรากฏและได้รับความสนใจตั้งแต่ปี ก.ศ. 1910 โดยอลเฟรด มิเนต์และไซมอน ได้ทดสอบเด็กที่มีฐานะทางเศรษฐกิจและสังคมที่แตกต่างกัน พบว่า มีข้อสอบบางข้อที่น่าจะวัดในสิ่งที่เป็นการกระทบกับการฝึกฝนทางด้านวัฒนธรรม (เช่นที่บ้านหรือที่โรงเรียน) มากกว่าที่จะวัดความสามารถทางสมอง หลังจากศึกษารั้งนี้ มิเนต์และไซมอน ได้มีการพิจารณาอย่างละเอียดและได้อาชีวศึกษาเหล่านี้ออกจากแบบทดสอบ เช่น ถ้าเป็นเด็กอายุ 7 ปี ข้อคำถามที่เกี่ยวกับน้ำมือ การคัดลอกประโยชน์และการเรียงตัวเรียงอู จะถูกคัดออกจากแบบทดสอบทั่วไป โดยในการทดสอบจะพิจารณาเฉพาะข้อคำถามที่ไม่เข้มแข็งกับการฝึกฝนของทางบ้าน (Home Training) ความดึงใจ (Attention) ภาษา (Language) นิสัยในการพูด (Habit of

Looking at Picture) และแบบฝึกหัดเชิงวิชาการ (Scholastic Exercise) ในปี ค.ศ. 1912

วิลเลียม สเตอร์น (William Stern) ผู้ที่ได้รับในนิยาม I.Q. ในรูปของอัตราส่วน ระหว่างอายุ สมองกับอายุปัญพิทิน ได้ศึกษาความแตกต่างระหว่างชนชั้นต่าง ๆ (Class Difference) ในประเทศเยอรมัน ผลที่ได้มีความคล้ายคลึงกับที่บีเนต์พบ คือ ความแตกต่างระหว่างชนชั้นมีนัยสำคัญ เสตอร์นพยายามที่จะค้นหาและทำความเข้าใจเกี่ยวกับแหล่งที่ทำให้เกิดความแตกต่างอย่างจริงจัง เพื่อที่จะเสนอแนะว่า แบบทดสอบต่าง ๆ ที่ได้พิสูจน์มาแล้วนั้นเข้าข้างชนชั้นหนึ่งมากกว่าชนชั้นอื่น ๆ อย่างชัดเจน

ต่อมาในปี ค.ศ. 1951 ได้มีการเริ่มค้นศึกษาเกี่ยวกับความลำเอียงของข้อสอบแนวใหม่ โดยเอลลส์ (Ells) และคณะสนใจความแตกต่างของการวัดบางอย่าง ซึ่งไม่มีความแม่นยำพอที่จะสะท้อนให้เห็นความสามารถได้ เช่น ความแตกต่างทางด้านสติปัญญาของนักเรียน อาจขึ้นอยู่กับเนื้อหาเฉพาะของข้อสอบแต่ละข้อในแบบทดสอบ และไม่ได้สะท้อนความสามารถที่แท้จริงของนักเรียน จุดมุ่งหมายหลักเพื่อที่จะอธิบายความแตกต่างระหว่างกลุ่มนักเรียนโดยใช้สติปัญญา โดยมุ่งไปที่ข้อคำถามแต่ละข้อเช่นเดียวกับบีเนต์และสเตอร์น คือตัวข้อคำถามใดมีความลำเอียง หรือไม่ แล้วมีคราฟท์ด้วยวิธีหากาคาวิเคราะห์ความยากของข้อสอบระหว่างกลุ่มที่แตกต่างกัน ผลการศึกษาพบว่าขึ้นอยู่กับความแตกต่างระหว่างกลุ่มนักเรียนที่มีระดับชั้นทางสังคม และภูมิหลังของเชื้อชาติและวัฒนธรรมด้วยกันในการตอบแบบทดสอบ เอลลส์และคณะได้พยายามพัฒนาวิธีการปฏิบัติและแก้ไขอิทธิพลของระดับชั้นทางสังคมและภูมิหลังของเชื้อชาติ ที่มีผลต่อการตอบข้อสอบ เพื่อให้ข้อสอบที่มีความลำเอียงในแบบทดสอบต่าง ๆ โดยปรับเปลี่ยนกลุ่มของเชื้อชาติที่อยู่ในชั้นทางสังคม

ผลจากการศึกษาของเอลลส์ได้นำไปสู่การศึกษาอย่างกว้างขวางในช่วง 30 กว่าปีที่ผ่านมา เมื่อรู้ว่ามีความมุ่งหวังที่จะทำให้เกิดความยุติธรรมในการสอนเข้าศึกษาและโอกาสในการมีงานทำ ดังนั้นในช่วงปลายปี ค.ศ. 1960 และ ค.ศ. 1970 ได้มีความตื่นตัวเกี่ยวกับความลำเอียงของข้อสอบ โดยมุ่งไปที่แบบทดสอบวัดสติปัญญา เพราะเป็นแบบทดสอบที่ถูกนำมาใช้ในทางด้านการศึกษา การคัดเลือก (Selection) และจัดนักเรียนเข้าสู่ตำแหน่ง (Placement) อย่างแพร่หลายจนกระทั่งปี ค.ศ. 1971 บริษัทกริกกัส วี ดักซ์ เพาเวอร์ จำกัด (Griggs V. Duke Power) ได้วินิจฉัยแบบทดสอบต่าง ๆ แล้วว่าให้เห็นว่าผลการทดสอบเหล่านั้นใช้ในการคัดเลือกอาจไม่เหมาะสม อาจพิจารณาตามเชื้อชาติของผู้สมัครงาน ยกเว้นกรณีที่ผลของการปฏิบัติงานมีความสอดคล้องกับการทดสอบ

ในช่วงเวลาเดียวกันมีการโต้แย้งในเรื่องความลำเอียงของข้อสอบมากขึ้น เมื่อเจนเซน (Jensen, 1969) ได้เสนอบทความลงในสารทางการศึกษาของ沙尔瓦ร์ด ฉบับที่ 39 ปี ค.ศ. 1969 โดยกล่าวว่า “โอลิวเป็นผลมาจากการถ่ายทอดทางพันธุกรรมมากกว่าสิ่งแวดล้อม ความแตกต่างของโอลิว

ระหว่างคนผิวขาวกับคนผิวผ่านน้ำไม่ได้มาจากการสิงแผลล้อมเพียงอย่างเดียว และพันธุกรรมอีกนายความแปรปรวนผลการปฏิบัติงานของแต่ละคน ได้ประมาณร้อยละ 80 จากบทความของเจนเซ่นนี้ เองทำให้สมาคมนักจิตวิทยาคนผิวผ่านน้ำเสนอว่า การทดสอบคนผิวผ่านน้ำทุกคนควรทำให้เท่าเทียมกันมากที่สุดเท่าที่จะเป็นไปได้ ยังมีบทความและวิธีการทางสถิติในการวิเคราะห์ความถูกต้องของข้อสอบเกิดขึ้นมากมาย ซึ่งแสดงให้เห็นว่าปัญหาเรื่องความถูกต้องของข้อสอบมีความสำคัญ เพื่อจะสร้างแบบทดสอบที่มีความยุติธรรมทางวัฒนธรรม และให้ผลการวิเคราะห์มีความถูกต้องมากที่สุด (Camilli & Shepard, 1994, pp. 4-7) นอกจากนี้ยังมีเหตุการณ์ที่ทำให้มีการตั้งตัวเป็นอันมากในเรื่องความถูกต้องของข้อสอบ คือ กรณีที่มาร์โค เดฟูนิส (Marco Defunis) และคณะซึ่งถูกปฏิเสธจากโรงเรียนกฎหมายของมหาวิทยาลัยแห่งวอชิงตัน ได้ฟ้องร้องว่าเขาได้คะแนนการสอบสูงกว่าผู้ได้รับการคัดเลือกของคนที่มหาวิทยาลัยเข้าศึกษา และได้ยื่นฎีกาฟ้อง ชาร์ล โอเดการ์ด (Charles Odsgard) อธิการบดีมหาวิทยาลัยดังกล่าวและคณะ เพื่อให้พิจารณาทบทวนการคัดเลือกนักศึกษาใหม่ หลังจากนั้นการพิจารณาตรวจสอบความถูกต้องของข้อสอบระหว่างผู้สอบกลุ่มย่อยที่มีลักษณะแตกต่างกันในเรื่องเพศ เชื้อชาติ ศาสนาหรือวัฒนธรรม และได้ปฏิบัติกันมาจนปัจจุบัน ซึ่งเป็นเสมือนส่วนหนึ่งของการพัฒนาแบบทดสอบ โดยทั่วๆ ไป โดยเฉพาะการพัฒนาแบบทดสอบมาตรฐาน เช่น ศูนย์บริการการทดสอบทางการศึกษา (Education Testing Service) ของประเทศไทย ได้มีการตรวจสอบความถูกต้องของข้อสอบ ทั้งก่อนนำไปทดลองใช้ ใช้การตรวจสอบโดยผู้เชี่ยวชาญ ซึ่งประกอบด้วยบุคคลหลายฝ่ายตรวจสอบความไว (Sensitivity) ของแบบทดสอบโดยพิจารณาถึงรูปแบบของข้อสอบ เนื้อหา ภาพประกอบ คำที่ใช้ และอื่นๆ เพื่อไม่ให้มีความถูกต้องหรือเกิดการได้เปรียบเสียเปรียบระหว่างผู้สอบกลุ่มย่อยที่มีลักษณะแตกต่างกัน ด้วยการตรวจสอบความถูกต้องของข้อสอบหลังจากการทดลองใช้นั้นจะใช้วิธีการทางสถิติ

“ความตรง” เป็นหัวใจสำคัญของคุณภาพแบบสอบ ใน การสร้างและการตรวจสอบคุณภาพของแบบสอบจะต้องคำนึงถึงคุณภาพด้านความตรงเป็นสำคัญ ทั้งนี้ เพราะว่าความตรงเป็นคุณสมบัติของแบบสอบที่แสดงถึงความสามารถในการวัดได้ถูกต้องแม่นยำ ถ้าผลการวัดได้ถูกต้อง ก็จะสอดคล้องกับความจริงของข้อสอบ แต่ถ้าผลการวัดไม่ถูกต้อง ก็จะแสดงถึงความไม่ถูกต้องของแบบสอบ ดังนั้นจะต้องคำนึงถูกต้องในกระบวนการประเมินคุณภาพแบบสอบ ที่สำคัญที่สุดคือ “ความตรง” ที่ต้องคำนึงถึงคุณภาพด้านความตรงเป็นสำคัญ ทั้งนี้ เพราะว่าความตรงเป็นคุณสมบัติของแบบสอบที่แสดงถึงความสามารถในการวัดได้ถูกต้องแม่นยำ ถ้าผลการวัดได้ถูกต้อง ก็จะสอดคล้องกับความจริงของข้อสอบ แต่ถ้าผลการวัดไม่ถูกต้อง ก็จะแสดงถึงความไม่ถูกต้องของแบบสอบ ดังนั้นจะต้องคำนึงถูกต้องในกระบวนการประเมินคุณภาพแบบสอบ (Item and Test Unfairness)

ความอยุติธรรมของข้อสอบ หรือแบบสอบเกิดขึ้นในกรณีที่ผู้สอบกลุ่มย่อยต่างกัน และมีลักษณะเฉพาะบางอย่างแตกต่างกัน มีความได้เปรียบหรือเสียเปรียบกันทั้งที่มีความสามารถจริงเท่าเทียมกัน การศึกษาคุณภาพด้านความอยุติธรรมของข้อสอบหรือแบบสอบระหว่างผู้สอบกลุ่ม

ต่าง ๆ เริ่มศึกษาภัยอย่างจริงจังในช่วงปลายทศวรรษของปี พ.ศ. 1960 มีการเสนอวิธีการต่าง ๆ เพื่อตรวจสอบความลำเอียงของข้อสอบ (Item Bias) ความลำเอียงของแบบสอบถาม (Test Bias) และความลำเอียงในการคัดเลือก (Selection Bias) โดยนิยามความลำเอียงว่าเป็น ความคาดเดาล่วงหน้า อย่างเป็นระบบ (Systematic Error) ที่เกิดขึ้นจากการวัด ความพยายามในการตรวจสอบความลำเอียงดังกล่าวดำเนินไปเพื่อจำแนกข้อสอบที่ทำหน้าที่ไม่เหมาะสมหรือไม่สอดคล้องกับความต้องการของผู้ทดสอบ ที่มีลักษณะบางอย่างแตกต่างกัน เช่น เชื้อชาติ ศาสนา วัฒนธรรม ภูมิลำเนา สังคม เพศ ภาษา อายุ ประสบการณ์ เป็นต้น เพื่อพัฒนาแบบสอบถามให้มีคุณภาพเหมาะสมสำหรับนำไปใช้ทดสอบต่อไป

ในเวลาต่อมาผู้วัดผลการศึกษาได้ทำการศึกษาความลำเอียงของข้อสอบกันอย่างกว้างขวาง ทำให้เกิดความสับสนของการใช้คำและความหมาย มีประเด็นโต้แย้งกันว่าความลำเอียงของข้อสอบ เป็นผลการตัดสินว่าข้อสอบมีความยุติธรรมหรือไม่ อันส่งผลต่อการบรรลุวัตถุประสงค์ของการใช้แบบสอบถาม หรือความลำเอียงของข้อสอบเป็นสารสนเทศทางสถิติที่ได้จากข้อสอบ เกี่ยวกับความสัมพันธ์ระหว่างคุณลักษณะที่ข้อสอบมุ่งวัดกับประสบการณ์ของผู้สอบกลุ่มต่าง ๆ ที่ทำการสอบ เมื่อกลุ่มผู้สอบต่างกลุ่มกันตอบข้อสอบข้อเดียวกัน ความแตกต่างที่เกิดขึ้นอาจมาจากการไม่เหมาะสมของข้อคำถาม ซึ่งสามารถเกิดขึ้นได้หลายลักษณะ หรือประสบการณ์ของผู้สอบซึ่งอาจมีลักษณะพื้นฐานเดิมแตกต่างกัน ในหลายสถานการณ์จะไม่เหมาะสมที่จะใช้คำว่า ข้อสอบลำเอียง (Biased Item) เมื่อจากเป็นภาษาที่มีความหมายในเชิงลบ ประกอบกับเกณฑ์ที่ใช้สำหรับตัดสินความลำเอียงยังมีความคลุมเครือและค่อนข้างสับสน ดังนั้นจึงควรเปลี่ยนมาใช้คำว่า การทำหน้าที่ต่างกันของข้อสอบ ซึ่งเป็นคำที่มีความเป็นกลางและเหมาะสมกว่า

การทำหน้าที่ต่างกันของข้อสอบกับความลำเอียงของข้อสอบมีแนวคิดที่แตกต่างกัน สำหรับการทำหน้าที่ต่างกันของข้อสอบ เป็นกระบวนการที่เน้นการใช้วิธีการทางสถิติสำหรับตรวจสอบ เพื่อให้ได้สารสนเทศเกี่ยวกับการทำหน้าที่ของข้อสอบสำหรับกลุ่มผู้สอบกลุ่มข้อยกเว้นที่มีลักษณะเฉพาะบางอย่างแตกต่างกัน ส่วนความลำเอียงของข้อสอบเป็นกระบวนการตัดสินความยุติธรรมของข้อสอบ โดยนำสารสนเทศการทำหน้าที่ต่างกันของข้อสอบมาวิเคราะห์เชิงตรรกะ (Logical Analysis) โดยผู้เชี่ยวชาญพิจารณาถึงการเขียนข้อสอบ เนื้อหาสาระของข้อสอบและคุณลักษณะของการวัด เพื่อรับว่าข้อสอบข้อนั้นลำเอียงเข้าข้างกลุ่มใดหรือไม่ เพราะเหตุใดจึงเป็นการตัดสินความลำเอียงของข้อสอบ (Camilli & Shapard, 1994)

ความหมายของความสำเร็จและการทำหน้าที่ต่างกันของข้อสอบ

นักวิจัยทางการวัดผลพยายามท่านได้ให้ความหมายของความสำเร็จของข้อสอบและการทำหน้าที่ต่างกันของข้อสอบไว้ดังนี้

ความสำเร็จของข้อสอบ หมายถึง สัดส่วนของผู้สอบที่ตอบข้อสอบได้ถูกต้องไม่เท่ากันในแต่ละกลุ่มประชากรที่ใช้ในการศึกษา เมื่อกลุ่มผู้สอบมีคะแนนเท่ากันและข้อสอบมีความเป็นเอกพันธ์ (Scheuneman, 1979)

ความสำเร็จของข้อสอบ หมายถึง ข้อสอบที่มีค่าความยากลำบากพอดำรงรับสมาร์กของผู้สอบกลุ่มนี้มากกว่าสมาชิกของผู้สอบอีกกลุ่มหนึ่ง (Rudner, Getson, & Knight, 1980)

ข้อสอบที่มีความสำเร็จด้านวัฒนธรรมว่า หมายถึง ข้อสอบที่ทำให้ผู้สอบซึ่งคุ้นเคยกับวิชาเฉพาะหรือกระบวนการเฉพาะมีโอกาสตอบถูกมากกว่าผู้สอบคนอื่น ๆ ที่มีลักษณะของกลุ่มที่แตกต่างกัน (Eells, 1951 citing Jensen, 1980)

ความสำเร็จของข้อสอบ หมายถึง ความโน้มเอียงของข้อสอบที่เมื่อใช้คะแนนจากข้อสอบนั้นแล้วทำให้การตัดสินใจผลเป็นไปอย่างไม่ยุติธรรม (Popham, 1981)

ความสำเร็จของข้อสอบ หมายถึง โอกาสในการตอบข้อสอบได้ถูกต้องแตกต่างกันสำหรับการวัดความสามารถ หรือโอกาสในการตอบข้อสอบในทางบวกแตกต่างกันสำหรับการวัดเจตคติ เมื่อผู้สอบที่มีคุณลักษณะของการวัดในปริมาณเท่ากัน แต่มาจากการกลุ่มประชากรย่อที่แตกต่างกัน (Hulin, Drasgow, & Parson, 1983)

ความสำเร็จของข้อสอบ หมายถึง โอกาสในการตอบข้อสอบได้ถูกต้องของผู้สอบกลุ่มนี้ นิค่าคำกว่าหรือสูงกว่าผู้สอบอีกกลุ่มหนึ่งที่มีระดับความสามารถเดียวกัน (Dorans & Kulick, 1986)

ความสำเร็จของข้อสอบ หมายถึง คะแนนข้อสอบของกลุ่มผู้สอบที่มีความสามารถเท่ากันแต่มาจากต่างกลุ่มกัน มีความแตกต่างกันอย่างเป็นระบบ (Kederman, 1990)

การทำหน้าที่ต่างกันของข้อสอบ หมายถึง สารสนเทศทางสถิติของข้อสอบที่ได้จากผลการตอบของผู้สอบต่างกลุ่มกันและมีความสามารถเท่ากัน แต่มีโอกาสในการตอบข้อสอบได้ถูกต้องแตกต่างกัน (Holland & Wainer, 1993)

การทำหน้าที่ต่างกันของข้อสอบ หมายถึง ความเป็นพหุนิธิในการวัดของข้อสอบ ซึ่งแสดงได้จากการแยกแยะความสามารถหลัก (Primary Ability) ของผู้สอบตัวตัว 2 กลุ่มขึ้นไป มีความสามารถเท่ากัน และมีการแยกแยะความสามารถรอง (Secondary Ability) แตกต่างกัน (Camilli & Shepard, 1994)

การทำหน้าที่ต่างกันของข้อสอบ หมายถึง ฟังก์ชันการตอบสนองข้อสอบซึ่งคำนวณจากกลุ่มผู้สอบข้อสอบที่ต่างกัน มีค่าไม่เท่ากัน (Naryanan & Swaminathan, 1996)

ข้อสอบทำหน้าที่ต่างกัน เป็นข้อสอบที่วัดความสามารถหรือคุณลักษณะทางจิตวิทยาของแค่กลุ่มไม่ตรงกัน (Allen & Yen, 1979 อ้างถึงใน อารี วัชรสอดดิกุล, 2543, หน้า 18)

การทำหน้าที่ต่างกันของข้อสอบ หมายถึง ความแตกต่างในการทำหน้าที่ของข้อสอบ หลังจากกลุ่มของผู้สอบได้ถูกจับคู่ตามความสามารถหรือคุณลักษณะที่ข้อสอบนั้นวัด (Millsap & Everson, 1993 อ้างถึงใน อารี วัชรสอดดิกุล, 2543, หน้า 8)

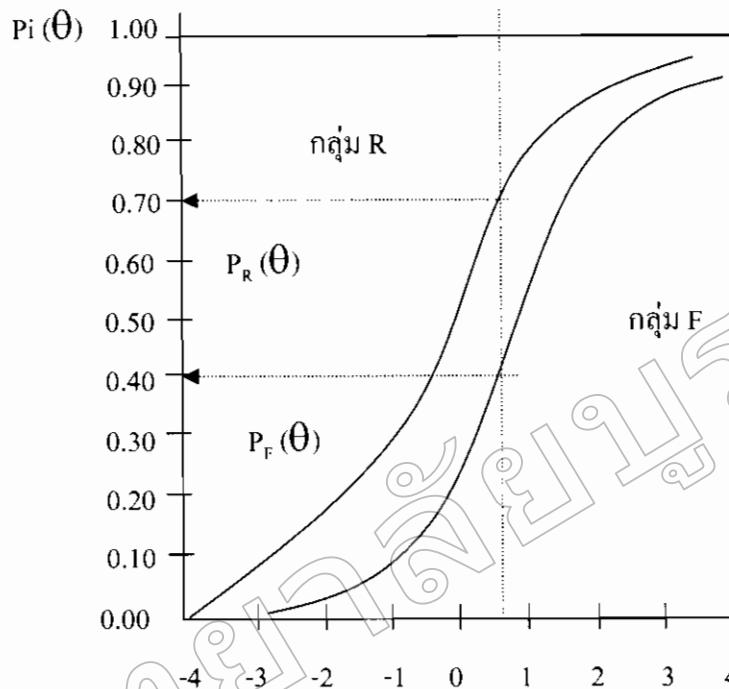
ดังนั้นจึงสรุปได้ว่า การทำหน้าที่ต่างกันของข้อสอบ หมายถึง การที่ข้อสอบทำให้ผู้สอบจากต่างกลุ่มกันที่มีความสามารถหรือคุณลักษณะที่มุ่งวัดเท่ากัน มีโอกาสในการตอบข้อสอบได้ถูกต้องแตกต่างกัน หรือมีฟังก์ชันการตอบสนองข้อสอบแตกต่างกัน การทำหน้าที่ต่างกันของข้อสอบเกิดขึ้นเมื่อนำข้อสอบไปทดสอบกับผู้สอบกลุ่มข้อสอบต่างกัน ที่มีความสามารถหลัก (Primary Ability) ระดับเดียวกันหรือมีคุณลักษณะแฝง (Latent Trait) ที่ต้องการวัดเท่ากัน แต่มีความสามารถรอง (Secondary Ability) แตกต่างกัน ทำให้ผู้สอบต่างกลุ่มที่นำมาจับคู่เบรย์นเทิร์บมีโอกาสตอบข้อสอบถูกต้องต่างกัน

ประเภทของการทำหน้าที่ต่างกันของข้อสอบ

การทำหน้าที่ต่างกันของข้อสอบ เป็นการเบรย์นเทิร์บผลการตอบข้อสอบระหว่างกลุ่มผู้สอบอย่างน้อย 2 กลุ่มขึ้นไป ปกตินิยมทำการเบรย์นเทิร์บ 2 กลุ่ม ประกอบด้วยกลุ่มแรก เรียกว่า กลุ่มเบรย์นเทิร์บ (Focal Group หรือกลุ่ม F) เป็นกลุ่มที่สนใจศึกษาและคาดว่าจะเป็นกลุ่มที่เสียเบรย์นในการตอบข้อสอบ และกลุ่มที่สอง เรียกว่า กลุ่มอ้างอิง (Reference Group หรือกลุ่ม R) เป็นกลุ่มที่คาดว่าจะได้เบรย์นในการตอบข้อสอบได้ถูกต้อง

ในการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบ พบว่า ข้อสอบสามารถทำหน้าที่แตกต่างกันได้ 2 ประเภท (Mellenbergh, 1982) ได้แก่ การทำหน้าที่ต่างกันของข้อสอบแบบเอกรูป (Uniform) และแบบอนุรูป (Nonuniform)

1. ข้อสอบทำหน้าที่ต่างกันแบบเอกรูป (Uniform DIF) หมายถึง ข้อสอบที่ทำให้ผู้สอบกลุ่มหนึ่ง มีโอกาสในการตอบข้อสอบถูกมากกว่าผู้สอบอีกกลุ่มหนึ่งสม่ำเสมอ กัน ในทุกระดับความสามารถ เมื่อพิจารณาโถึงคุณลักษณะข้อสอบของผู้สอบ 2 กลุ่ม จะพบว่าไม่มีปฏิสัมพันธ์ระหว่างระดับความสามารถของผู้สอบกับการเป็นสมาชิกกลุ่ม (Group Membership) ดังภาพที่ 1



ภาพที่ 1 ข้อสอบทำหน้าที่ต่างกันแบบเอกรูป (Uniform DIF)

2. ข้อสอบทำหน้าที่ต่างกันแบบอเนกรูป (Nonuniform DIF) หมายถึง ข้อสอบที่ทำให้โอกาสในการตอบข้อสอบถูกของผู้สอบระหว่างกลุ่มแตกต่างกันอย่างไม่สม่ำเสมอ กันในทุกระดับความสามารถ เมื่อพิจารณาโดยคุณลักษณะข้อสอบของผู้สอบ 2 กลุ่ม พบร่วมกับ “ปฏิสัมพันธ์” ระหว่างระดับความสามารถของผู้สอบกับการเป็นสมาชิกกลุ่ม เช่น ที่ระดับความสามารถอีกระดับหนึ่ง กลุ่มผู้สอบกลุ่ม R มีโอกาสในการตอบข้อสอบถูกมากกว่ากลุ่มผู้สอบกลุ่ม F แต่ที่ระดับความสามารถอีกระดับหนึ่งกลุ่มผู้สอบกลุ่ม F มีโอกาสในการตอบข้อสอบถูกมากกว่ากลุ่มผู้สอบกลุ่ม R

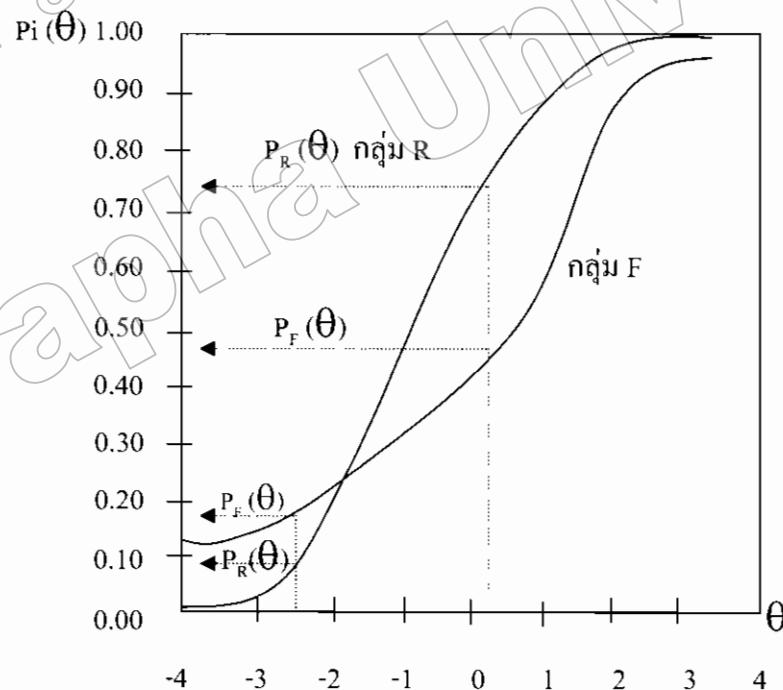
ตามทฤษฎีการตอบสนองข้อสอบ (Item Response Theory: IRT) สามารถพิจารณา “ปฏิสัมพันธ์” ดังกล่าวได้จากความแตกต่างของค่าพารามิเตอร์ “อำนาจจำแนกของข้อสอบ” ระหว่างผู้สอบกลุ่มย่อยสองกลุ่ม กล่าวคือ ถ้าข้อสอบทำหน้าที่ต่างกันแบบเอกรูปแล้ว ให้ลักษณะข้อสอบ (Item Characteristic Curves: ICCs) ระหว่างผู้สอบกลุ่มย่อยสองกลุ่มจะนานกัน หรือมีฟังก์ชันการตอบสนองข้อสอบ (Item Response Theory: IRT) เหมือนกัน แต่ถ้าข้อสอบทำหน้าที่ต่างกันแบบอเนกรูปแล้ว ให้ลักษณะข้อสอบระหว่างผู้สอบกลุ่มย่อยสองกลุ่มจะไม่นานกัน หรือมีฟังก์ชันการตอบสนองข้อสอบต่างกัน ดังนั้นความแตกต่างระหว่างคุณลักษณะข้อสอบทั้งสอง

แบบจะปั่งบอกถึงขนาดและทิศทางของข้อสอบทำหน้าที่ต่างกัน ซึ่งสามารถคำนวณได้โดยใช้สูตร การคำนวณพื้นที่ของ Raju (1990)

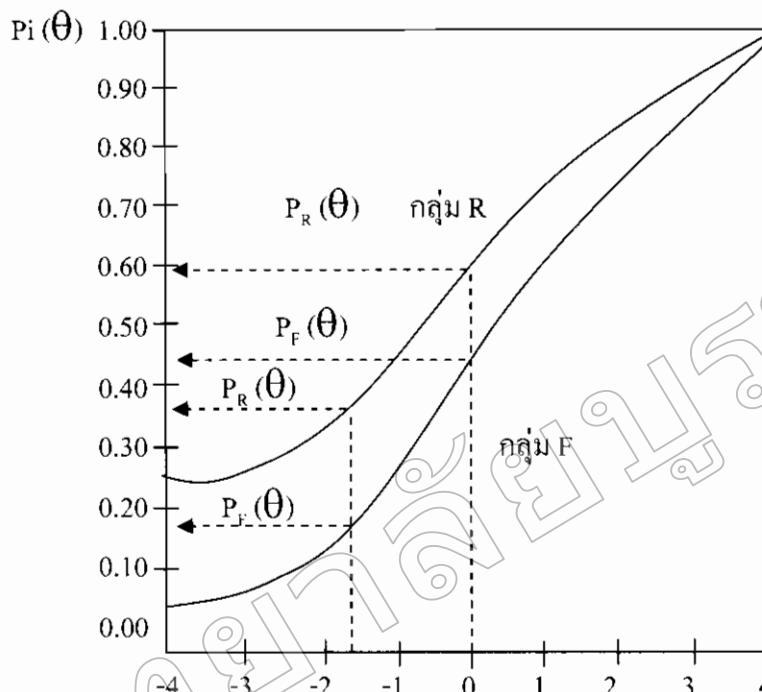
ข้อสอบทำหน้าที่ต่างกันแบบอเนกกรุปสามารถจำแนกได้เป็น 2 ลักษณะ (Swaminathan & Rogers, 1990) ดังนี้

1. ข้อสอบทำหน้าที่ต่างกันแบบอเนกกรุปโดยมีปฏิสัมพันธ์ไม่เป็นลำดับ (Disordinal Interaction) เป็นการทำหน้าที่ต่างกันสำหรับกลุ่มผู้สอบซึ่งเกิดขึ้น เมื่อโครงสร้างของข้อสอบตัดกัน ระหว่างช่วงความสามารถของผู้สอบหรือเรียกว่าข้อสอบทำหน้าที่ต่างกันแบบไม่มีทิศทาง (Non-Directional DIF) ดังภาพที่ 2

2. ข้อสอบทำหน้าที่ต่างกันแบบอเนกกรุปโดยมีปฏิสัมพันธ์เป็นลำดับ (Ordinal Interaction) เป็นการทำหน้าที่ต่างกันสำหรับกลุ่มผู้สอบซึ่งเกิดขึ้น เมื่อโครงสร้างของข้อสอบคู่ตัดกัน อย่างไม่สม่ำเสมอแต่ไม่ตัดกัน หรืออาจตัดกันน้อยกว่าช่วงความสามารถของผู้สอบตรงปลายสุดของ ช่วงความสามารถต่ำหรือสูง อาจเรียกข้อสอบลักษณะนี้ว่า ข้อสอบทำหน้าที่ต่างกันแบบมีทิศทาง เดียว (Unidirectional DIF) ดังภาพที่ 3



ภาพที่ 2 ข้อสอบทำหน้าที่ต่างกันแบบไม่มีทิศทาง (Non-Directional DIF)



ภาพที่ 3 ข้อสอบทำหน้าที่ต่างกันแบบมีทิศทางเดียว (Unidirectional DIF)

หลักการและวิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ

1. หลักการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ

การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ เป็นการเปรียบเทียบผลการตอบข้อสอบ เป็นรายชื่อระหว่างกลุ่มผู้สอบอย่างน้อย 2 กลุ่ม ที่มีความสามารถหลัก (Primary Ability) หรือ ความสามารถปัจจัยที่ต้องการวัดเท่ากัน แต่คาดว่าจะมีความได้เปรียบหรือเสียเปรียบกัน โดย กลุ่มหนึ่งถือเป็นกลุ่มอ้างอิงซึ่งตอบข้อสอบได้ถูกต้องมากกว่า ส่วนอีกกลุ่มคือกลุ่มเปรียบเทียบ ซึ่ง เป็นกลุ่มที่สนใจศึกษาและคาดว่าจะเป็นกลุ่มที่เสียเปรียบ

ในการเปรียบเทียบผลการตอบข้อสอบระหว่างกลุ่มอ้างอิงและกลุ่มเปรียบเทียบ จำเป็นต้องจับคู่ (Matching) ผู้สอบตามความสามารถ ซึ่งเป็นเงื่อนไขสำคัญของการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ เกณฑ์การจับคู่ (Matching Criteria) ที่นิยมใช้กันมี 2 วิธี ดังนี้

1.1 เกณฑ์ภายนอก (External Criterion)

การวิเคราะห์การทำหน้าที่ต่างกัน โดยใช้เกณฑ์ภายนอกนี้ สามารถนำไปใช้ได้ทั้ง ข้อสอบรายชื่อและแบบสอบทั้งฉบับ โดยการใช้คะแนนจากแบบสอบสอนอื่นเป็นเกณฑ์ภายนอก แล้วใช้เทคนิคการวิเคราะห์การถดถอย (Regression Analysis) เพื่อทำการเปรียบเทียบเส้นกราฟ ความสัมพันธ์ระหว่างตัวแปรเกณฑ์กับตัวแปรทำนาย ระหว่างกลุ่มอ้างอิงและกลุ่มเปรียบเทียบ

หลักการนี้มีจุดมุ่งหมายเพื่อสร้างสมการทำนายตัวแปรเกณฑ์ ซึ่งเป็นคะแนนของแบบสอนอื่นจากตัวแปรทำนายซึ่งเป็นคะแนนรายข้อ หรือคะแนนสอบระหว่างกลุ่มอ้างอิงและกลุ่มเปรียบเทียบ ในการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบจะใช้คะแนนรายข้อเป็นตัวแปรทำนาย แต่ถ้าการวิเคราะห์การทำหน้าที่ต่างกันของแบบสอน จะใช้คะแนนรวมของแบบสอนทั้งฉบับเป็นตัวแปรทำนาย สำหรับตัวแปรเกณฑ์ที่ใช้เป็นเกณฑ์ภายนอก อาจใช้คะแนนรวมทั้งฉบับ หรือเกรดเฉลี่ย หรือผลสัมฤทธิ์ในงานที่เกี่ยวข้องของผู้สอบ (Cronbach & Furby, 1970) สมการทำนายสำหรับกลุ่มอ้างอิงและกลุ่มเปรียบเทียบแสดงได้ดังนี้

กลุ่มอ้างอิง (R)

$$R_i = A_R + B_R X_i$$

กลุ่มเปรียบเทียบ (F)

$$Y_i = A_F + B_F X_i$$

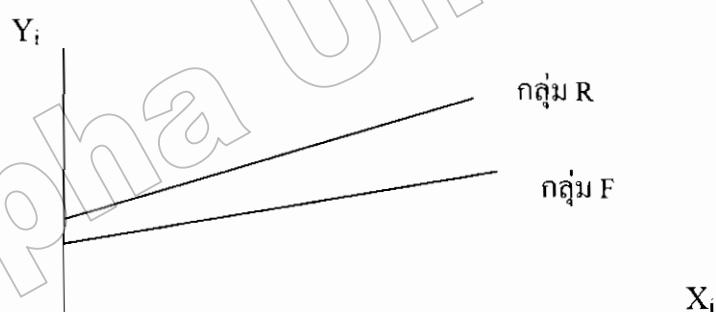
เมื่อ

Y_i = คะแนนของตัวแปรเกณฑ์ภายนอก

X_i = คะแนนของตัวแปรทำนาย

A = ค่าคงที่หรือค่าตัดแกน (Intercept)

B = ค่าความชัน (Slope)



ภาพที่ 4 การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบและแบบสอนโดยใช้วิธีการวิเคราะห์สมการลด削

จากฟังก์션การทำนายทั้ง 2 สมการ สามารถเปรียบเทียบค่าตัดแกน (A) และค่าความชัน (B) ของเส้นกราฟระหว่างกลุ่มอ้างอิงและกลุ่มเปรียบเทียบได้ ถ้าเส้นกราฟตั้งกล่ำมีค่าความชันหรือค่าตัดแกนแตกต่างกันสำหรับข้อสอบใดหรือแบบสอนใด แสดงว่าข้อสอบหรือแบบสอนนั้นมีการทำหน้าที่ต่างกัน โดยเข้าข้างกลุ่มผู้สอบที่มีค่าตัดแกนหรือค่าความชันที่สูงกว่า

การใช้เกณฑ์ภายนอกมีข้อดี คือ เกณฑ์ที่ใช้มีความเป็นอิสระจากข้อสอบและแบบสอบที่ต้องการตรวจสอบ แต่มีจุดอ่อนตรงที่ความหมายของเกณฑ์ที่จะนำมาใช้ ในทางปฏิบัติเป็นการยากที่จะหาได้แบบเกณฑ์ภายนอกจากแบบสอบฉบับอื่นที่มีความตรงเชิงทำนาย ซึ่งจะทำให้ผลการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบหรือแบบสอบ ขาดความแม่นยำและความสมบูรณ์

1.2 เกณฑ์ภายใน (Internal Criterion)

การวิเคราะห์การทำหน้าที่ต่างกันโดยใช้เกณฑ์ภายใน เป็นการนำวิธีการทางสถิติมาตรวจสอบการทำหน้าที่ต่างกันของข้อสอบหรือแบบสอบ โดยเน้นการพิจารณาจากโครงสร้างภายในของแบบสอบเป็นหลัก ด้วยการวิเคราะห์ผลจากการตอบข้อสอบ และความสามารถหรือคะแนนจริงของผู้สอบที่ได้จากแบบสอบฉบับนั้น เพื่อนำมาเปรียบเทียบระหว่างผู้สอบจากกลุ่มอ้างอิงและกลุ่มเปรียบเทียบที่มีความสามารถมาตรฐานหรือคะแนนจริงเท่ากัน ว่าจะมีผลการตอบหรือโอกาสในการตอบข้อสอบได้ถูกต้องมากต่างกันหรือไม่ เพื่อบ่งชี้ถึงการทำหน้าที่ต่างกันของข้อสอบ การวิเคราะห์ในลักษณะนี้นิยมใช้ค่าสถิติต่าง ๆ เป็นตัวบ่งชี้ถึงการทำหน้าที่ต่างกันของข้อสอบ ค่าสถิติทดสอบที่นิยมนำมาใช้พิสูจน์ได้ดังนี้

1.2.1 การทดสอบปฎิสัมพันธ์ (Interaction)

ในระยะเริ่มแรกของการศึกษาความลำเอียงของข้อสอบ มีการใช้สถิติทดสอบเอฟ (*F*- Test) ในการวิเคราะห์ความแปรปรวน เพื่อทดสอบปฎิสัมพันธ์ระหว่างกลุ่มผู้สอบกับข้อสอบ ถ้าการทดสอบมีนัยสำคัญ เป็นสัญญาณของการทำหน้าที่ต่างกันของข้อสอบ (Cleary & Hilton, 1968; Jensen, 1974) จากนั้นจึงทำการวิเคราะห์ต่อด้วยวิธีการ Post Hoc เพื่อระบุข้อสอบที่มีผลต่อการเกิดปฎิสัมพันธ์ซึ่งเป็นข้อที่ทำหน้าที่ต่างกัน

วิธีการนี้มีข้อดีที่สามารถศึกษาผู้สอบหลาย ๆ กลุ่ม ได้สะดวก แต่มีจุดอ่อนในเรื่องการควบคุมกลุ่มต่าง ๆ ให้มีความสามารถที่ตัดเทียบกัน ขนาดกลุ่มตัวอย่างของกลุ่มต่าง ๆ และอัตราความคลาดเคลื่อนประเภทที่ 1 จะสูงขึ้น ถ้าจำนวนข้อสอบเพิ่มมากขึ้น

1.2.2 การวัดความเบี่ยงเบนสัมพันธ์ (Relative Deviation)

การคำนวณค่าความยากของข้อสอบ เช่น *p*, *b* เป็นต้น เมื่อคำนวณแยกระหว่างกลุ่ม และแปลงให้เป็นค่าความยากมาตรฐาน (Δ) สามารถนำไปพื้นอัตรากับค่าเดียวกัน หรือเบี่ยงเบนเกินความคลาดเคลื่อนมาตรฐานของค่าความยากที่กำหนด ย้อมแสดงถึงการกระทำหน้าที่ต่างกันของข้อสอบ (Cleary & Hilton, 1968; Angoff & Ford, 1973) รวมทั้งสามารถคำนวณค่าสหสัมพันธ์ระหว่างค่าความยากรายข้อระหว่างกลุ่ม เพื่อแสดงถึงการกระทำหน้าที่ต่างกันของข้อสอบ ถ้าสหสัมพันธ์เข้าใกล้ 1.00 แสดงว่าค่าความ

หากสัมพัทธ์ของข้อสอบมีค่าไกล์เดียงกันระหว่างกลุ่ม ดังนั้นแบบสอบถามวัดคุณลักษณะคล้ายกันระหว่างกลุ่ม

วิธีการนี้มีข้อดีและข้อเสียคือการทดสอบปฎิสัมพันธ์ นอกจากนี้ค่าความยากของข้อสอบ มิใช่ค่าแทนของค่าความยากจริงของข้อสอบ และได้รับอิทธิพลจากค่าแทรกซ้อนอื่นๆ ได้แก่ ค่าอำนาจจำแนกและความสามารถของผู้สอบ

1.2.3 การเปรียบเทียบน้ำหนักตัวประกอบ (Factor Loading)

การวิเคราะห์ตัวประกอบ (Factor Analysis) เป็นเทคนิคทางสถิติกิที่นิยมใช้ในการตรวจสอบความตรงเชิงทฤษฎีหรือโครงสร้าง (Construct Validity) เมื่อมีการวิเคราะห์ตัวประกอบมาใช้ในการวิเคราะห์โครงสร้างของแบบสอบถามแยกตามกลุ่มผู้สอบ ความไม่สอดคล้องกันระหว่างน้ำหนักตัวประกอบบนคุณลักษณะสำคัญที่มุ่งวัด หรือความแตกต่างของค่าเฉลี่ยคะแนนตัวประกอบ (Factor Scores) ระหว่างกลุ่มผู้สอบ ย้อมสะท้อนการทำหน้าที่ต่างกันของข้อสอบและแบบสอบถาม

การใช้เทคนิคการวิเคราะห์ตัวประกอบเชิงสำรวจ (Exploratory Factor Analysis: EFA) สำหรับศึกษาการทำหน้าที่ต่างกันของข้อสอบ จะมีจุดอ่อนในเรื่องความไม่สอดคล้องระหว่างน้ำหนักตัวประกอบ อาจเกิดจากความต้องการของความสามารถระหว่างกลุ่มก็ได้ แนวทางที่เหมาะสมจะควรใช้เทคนิคการวิเคราะห์ตัวประกอบเชิงยืนยัน (Confirmatory Factor Analysis: CFA) นอกจากนี้ยังสามารถใช้ CFA สำหรับตรวจสอบความแยกต่างระหว่างกลุ่ม ในด้านคุณลักษณะหรือความสามารถหลักและความสามารถรองได้อีกด้วย (Camilli & Shepard, 1994)

1.2.4 การเปรียบเทียบโอกาสตอบข้อสอบถูก

การวิเคราะห์โอกาสตอบข้อสอบถูกของผู้สอบจากกลุ่มองอ้างอิงและกลุ่มเปรียบเทียบที่มีความสามารถเท่ากัน เป็นแนวทางสำคัญที่นิยมใช้กันและเป็นที่ยอมรับในปัจจุบัน สำหรับมีชีวิตรากฐานที่มีความสามารถเท่ากัน เช่น วิธีแมเนเกล - แชนส์เซล เป็นต้น

เปรียบเทียบค่าสัดส่วนหรือความน่าจะเป็นในการตอบข้อสอบถูกของผู้สอบต่างกันที่มีความสามารถเท่ากัน เช่น วิธีแมเนเกล - แชนส์เซล เป็นต้น

เปรียบเทียบค่าพังก์ชันการตอบสนองของข้อสอบ หรือโดยลักษณะข้อสอบระหว่างกลุ่มที่มีระดับความสามารถเท่ากัน เป็นวิธีที่อยู่บนพื้นฐานของทฤษฎี IRT เช่น วิธีวัดความแตกต่างของพื้นที่ วิธีวัดความแยกต่างของค่าพารามิเตอร์ความยาก วิธีการทดสอบไอ-สแควร์ของคอร์ด เป็นต้น

วิธีการนี้มีข้อดีที่สำคัญ ได้แก่ การคำนวณค่าสถิติของข้อสอบมีความน่าเชื่อถือ มีกลไกควบคุมความสามารถของผู้สอบโดยการจับคู่กับความสามารถ เพื่อทำการเปรียบเทียบ

ณ ตำแหน่งต่าง ๆ ที่มีความสามารถเท่ากัน จึงเป็นวิธีการที่ยอมรับกันทั่วไปแต่มีข้อจำกัดในด้าน ความสลับซับซ้อนของแนวคิดพื้นฐาน และการวิเคราะห์มีความจำเป็นต้องใช้โปรแกรมคอมพิวเตอร์โดยเฉพาะ

2. วิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ

วิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ จำแนกตามลักษณะการตรวจให้คะแนนได้เป็น 2 ประเภท คือ ข้อสอบที่มีการตรวจให้คะแนนแบบทวิภาคหรือ 2 ค่า และข้อสอบที่มีการให้คะแนนแบบพหุวิภาคหรือหลายค่า วิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ แต่ละประเภทสามารถจำแนกได้อีก 2 มิติ ได้แก่ มิติลักษณะของตัวประเมณท์ ซึ่งแบ่งเป็นกลุ่มวิธีที่ใช้คะแนนสังเกตได้ และกลุ่มวิธีที่ใช้คะแนนสังเกตไม่ได้หรือคะแนนของตัวประเมณ และมิติลักษณะของสถิติวิเคราะห์ ซึ่งแบ่งเป็นกลุ่มวิธีที่ใช้สถิติพารามิตริก และกลุ่มวิธีที่ใช้สถิตินักพารามิตริก

2.1 วิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบที่ให้คะแนนแบบทวิภาค

2.1.1 กลุ่มวิธีที่ใช้คะแนนสังเกตได้

วิธีในการคุณนี้มักวิเคราะห์ความทุณภูมิการทดสอบแบบดั้งเดิม หรือกลุ่มนี้ไม่ใช้ทฤษฎีการตอบสนองข้อสอบ โดยใช้คะแนนรวมของผู้สอบเป็นเกณฑ์การจับคู่ของกลุ่มข้อสอบ วิธีการตรวจสอบที่สำคัญในกลุ่มนี้ ได้แก่

2.1.1.1 วิธีวิเคราะห์ความแปรปรวน (Cleary & Hilton, 1968)

2.1.1.2 วิธีวิเคราะห์การลดด้อยโลจิสติก (Swaminathan & Rogers, 1990)

2.1.1.3 วิธีการแปลงค่าความยากของข้อสอบ (Clear & Hilton, 1968; Angoff & Ford, 1973)

2.1.1.4 วิธีแมนเทล – แ昏ส์เซล (Holland & Thayer, 1989)

2.1.1.5 วิธีดัชนีมาตรฐาน การปรับให้เป็นมาตรฐานด้วยน้ำหนักตัวประกอบ (Dorans & Kulick, 1986)

2.1.2 กลุ่มที่ใช้คุณลักษณะแฝง

วิธีในการนี้ใช้คุณลักษณะหรือตัวแปรแฝง ซึ่งวิเคราะห์บนพื้นฐานของทฤษฎีการตอบสนองข้อสอบ สำหรับใช้เป็นเกณฑ์การจัดอันดับผู้สอบ วิธีการตรวจสอบที่สำคัญในกลุ่มนี้ ได้แก่

2.1.2.1 วิธีวัดพื้นที่ความแตกต่างระหว่างโ้างการตอบสนองข้อสอบ

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

2.1.2.2 วิธีไอ – สแคร์ของลอร์ด (Lord, 1980)

2.1.2.3 วิธีอัตราส่วนไลค์ลิขุดทั่วไป (Thissen, Steinberg, & Wainer, 1993)

2.1.2.4 วิธีอัตราส่วนไลค์ลิขุด ลอกลินีเยอร์ (Loglinear IRT Likelihood Ratio)

(Thissen, Steinberg, & wainer, 1993)

2.1.2.6 วิธีชิปเทสท์ (Shealy & Stout, 1993)

2.2 วิธีการตรวจสอบการทำหน้าที่แตกต่างกันของข้อสอบที่ใช้คะแนนแบบพหุวิภาค

2.2.1 กลุ่มวิธีที่ใช้คะแนนสังเกตได้

2.2.1.1 วิธีการวิเคราะห์ความแปรปรวน (Clear & Hilton, 1968)

2.2.1.2 วิธีการวิเคราะห์การทดสอบโลจิสติกพหุวิภาค (Swaminathan &

Rogers, 1990)

2.2.1.3 วิธีดัชนีมาตรฐานพหุวิภาค (Dorans & Kulick, 1968)

2.2.1.4 วิธีแม่นเทล - แฮนส์เซลทั่วไป (Holland & Thayer, 1988, 1989)

2.2.2 กลุ่มที่ใช้คะแนนลักษณะแห่ง

2.2.2.1 วิธีอัตราส่วนไลค์ลิขุดในรูปทั่วไป (Thissen, Steinberg, & Wainer,

1993)

2.2.2.2 วิธีการให้คะแนนบางส่วน (Master, 1982)

2.2.2.3 วิธีชิปเทสท์พหุวิภาค (Shealy & Stout, 1993)

2.2.2.4 วิธีการใช้คะแนนบางส่วนทั่วไป (Muraki, 1992, 1993)

การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบตามแนววิถีการตอบสนอง

ข้อสอบ (Item Response Theory)

วิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบวิธีนี้ จะพิจารณาเปรียบเทียบฟังก์ชันการตอบข้อสอบ (Item Response Function) หรือเส้นคุณลักษณะข้อสอบ (Item Characteristic Curves) ระหว่างกลุ่มตัวอย่าง 2 กลุ่มขึ้นไปว่ามีความแตกต่างกันหรือไม่

การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบตามแนววิถีการตอบข้อสอบ มีจุดเด่นกว่าวิธีตามแนววิถีการทดสอบแบบเดิม ดังนี้

1. พฤติกรรมการตอบข้อสอบสามารถประมาณค่าพารามิเตอร์ของข้อสอบ เช่น ค่าความยาก (b) ค่าอำนาจจำแนก (a) และค่าการเค้า (c) ทำให้ลดความปนกันของพารามิเตอร์ข้อสอบลงได้

2. คุณสมบัติทางสถิติของข้อสอบ (Statistical Properties of Item) สามารถอธิบายได้ละเอียดและชัดเจนกว่า เมื่อข้อสอบทำหน้าที่ต่างกันระหว่าง 2 กลุ่ม

3. คุณสมบัติทางสถิติของข้อสอบ สามารถแสดงค่วยแพนภาพที่เป็นกราฟ ทำให้เข้าใจง่ายและนำไปใช้ได้อย่างกว้างขวาง

การวิเคราะห์การทำหน้าที่ต่างกันของแบบสอบถามและข้อสอบ

1. แนวคิดพื้นฐานของการวิเคราะห์การทำหน้าที่ต่างกันของแบบสอบถาม

ถ้าให้ $P_i(\theta_s)$ เป็นความน่าจะเป็นในการตอบข้อสอบข้อที่ i ได้ถูกต้องของผู้สอบคนที่ S ซึ่งมีความสามารถ θ ในกรณีที่เรามีผู้สอบร่วมกัน 2 กลุ่ม ได้แก่ กลุ่มอ้างอิง (R) และกลุ่มเปรียบเทียบ (F) แต่ละกลุ่มทำแบบสอบถามฉบับเดียวกันซึ่งมีจำนวน k ข้อ เมื่อทำการแยกวิเคราะห์ตามกลุ่ม แต่ละกลุ่มต่างก็มีชุดของค่าพารามิเตอร์ของข้อสอบ (a , b และ c) ซึ่งสามารถคำนวณหาค่าความน่าจะเป็นในการตอบถูกได้ ดังนี้

$P_{iR}(\theta_s) =$ ความน่าจะเป็นในการตอบข้อสอบข้อที่ i ได้ถูกต้อง ของผู้สอบคนที่ S ซึ่งเป็นสมาชิกของกลุ่มอ้างอิงและมีความสามารถ θ

$P_{iF}(\theta_s) =$ ความน่าจะเป็นในการตอบข้อสอบข้อที่ i ได้ถูกต้อง ของผู้สอบคนที่ R ซึ่งเป็นสมาชิกของกลุ่มเปรียบเทียบและมีความสามารถ θ

ตามทฤษฎี IRT เราสามารถคำนวณค่าคะแนนจริง (True Scores) ของผู้สอบคนที่ S ซึ่งมีความสามารถ θ ได้ดังนี้

$$T_S = \sum_{i=1}^k P_i(\theta_s)$$

ผู้สอบแต่ละคนไม่ว่าจะเป็นสมาชิกของกลุ่มใด มี T_S ได้ 2 ค่า ในฐานะสมาชิกของกลุ่ม R และ F จึงคำนวณค่าได้ทั้ง T_{SR} และ T_{SF}

ถ้า $T_{SR} = T_{SF}$ แสดงว่า ผู้สอบ S เป็นอิสระจากกลุ่มสมาชิก ดังนั้นแบบสอบถามจึงทำหน้าที่ไม่แตกต่างกัน

$T_{SR} \neq T_{SF}$ แสดงว่า การเป็นสมาชิกกลุ่มของผู้สอบคนที่ S มีผลทำให้ T_S แตกต่างกันจึงเป็นสัญญาณของการทำหน้าที่ต่างกันของแบบสอบถาม

2. การวิเคราะห์การทำหน้าที่ต่างกันของแบบสอบถามและการทำหน้าที่ต่างกันของข้อสอบ

ราจูและคณะ (Raju et al., 1995) ได้เสนอกระบวนการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบและแบบสอบถาม (DIF and DTF Procedures) บนพื้นฐานของทฤษฎี IRT มีการคำนวณที่ซับซ้อนแต่ให้ผลที่น่าเชื่อถือ โดยนำวิธีทางสถิติมาใช้ในการตรวจสอบโครงสร้างแบบสอบถาม ด้วยเกณฑ์ภายใน กระบวนการตรวจสอบประกอบด้วยการคำนวณค่านี้การทำหน้าที่ต่างกัน 2 ระดับ

ได้แก่ ดัชนีการทำหน้าที่ต่างกันระดับแบบสอบ (DTF Index) และดัชนีการทำหน้าที่ต่างกันระดับข้อสอบ (DIF Index)

2.1 ดัชนีการทำหน้าที่ต่างกันระดับข้อสอบ (DIF Index)

$$\text{ถ้าให้ } d_{is} = P_{iR}(\theta_s) - P_{iF}(\theta_s)$$

$$\text{และ } D_S = T_{SR} - T_{SF}$$

ดังนั้นดัชนีการทำหน้าที่ต่างกันของข้อสอบแบบชดเชย (Compensatory DIF Index:

CDIF) เป็นดังนี้

$$\begin{aligned} \text{CDIF}_i &= E(d_i, D) \\ &= Cov(d_i, D) + \mu_{di}\mu_D \end{aligned}$$

เมื่อ E = ค่าความคาดหมาย (Expectation)

$Cov(d_i, D)$ = ความแปรปรวนระหว่างความแตกต่างของความน่าจะเป็นในการตอบข้อที่ i ได้ถูก กับความแตกต่างของคะแนนจริงระหว่างกลุ่มอ้างอิงกับกลุ่มเปรียบเทียบที่มีค่า θ เท่ากัน

μ_{di} = ค่าเฉลี่ยของความแตกต่างของความน่าจะเป็นในการตอบข้อสอบที่ i ได้ถูกต้อง ระหว่างกลุ่มอ้างอิงและกลุ่มเปรียบเทียบที่มีค่า θ เท่ากัน

μ_D = ค่าเฉลี่ยความแตกต่างของคะแนนจริงระหว่างกลุ่มอ้างอิงกับกลุ่มเปรียบเทียบที่มี θ เท่ากัน

2.2 ดัชนีการทำหน้าที่ต่างกันระดับแบบสอบ (DTF Index)

$$\begin{aligned} \text{DTF} &= E(T_{SR} - T_{SF})^2 \\ &= E(D_S)^2 \\ &= E(D_S)^2 - \mu_D^2 + \mu_D^2 \\ &= \sigma_D^2 + (\mu_{TR} - \mu_{TF})^2 \\ &= \sigma_D^2 + \mu_D^2 \end{aligned}$$

$$\begin{aligned} \text{เมื่อ } D &= T_{SF} - T_{SR} = \sum_{i=1}^K P_{iR}(\theta_S) - \sum_{i=1}^K P_{iF}(\theta_S) \\ &= \sum_{i=1}^K P_{iR}(\theta_S) - \sum_{i=1}^K P_{iF}(\theta_S) \end{aligned}$$

$$\begin{aligned}
 &= \sum_{i=1}^K d_i \\
 \therefore DTF &= E_D D^2 \\
 &= E_D \left(\sum_{i=1}^K d_i \right)^2 \\
 &= \sum_{i=1}^K [\text{cov}(d_i, D) + \mu_{d_i} \mu_D]
 \end{aligned}$$

สาเหตุของการทำหน้าที่ต่างกันของข้อสอบ

สาเหตุที่น่าจะก่อให้เกิดการทำหน้าที่ต่างกันของข้อสอบมีหลายสาเหตุด้วยกัน ซึ่งสามารถสรุปสาเหตุที่ก่อให้เกิดการทำหน้าที่ต่างกันของข้อสอบมากที่สุด ได้ดังนี้

1. การเดา (Guessing) อาจจะเกิดจากข้อสอบนั้นยากเกินไปหรือเวลาไม่เพียงพอ จะก่อให้เกิดความไม่เท่าเทียมกันในโอกาสการตอบข้อสอบลูกของผู้สอนแต่ละคน
2. ความเร็ว (Speed) หรือเวลาในการตอบ จะทำให้เกิดการคาดหรือในกรณีเวลาเหลืออย่างจำกัดข้อสอบไม่ทัน ซึ่งจะมีผลกับข้อสอบข้อหลัง ๆ โดยเฉพาะในการศึกษาความจำอธิบายของข้อสอบวัดความถนัด
3. ความกำหนดหรือความไม่ชัดเจน (Unclear) ของข้อคำถาม นั่นคือ ข้อคำถามขาดความเป็นปัจจัย การใช้ภาษาท้องถิ่นหรือใช้คำที่ไม่เป็นภาษากลางในการสื่อสารความหมาย ซึ่งจะก่อให้เกิดความจำอธิบายกับกลุ่มภาษาใดภาษาหนึ่งขึ้นได้
4. ลำดับขั้นของคำถาม (Series) อาจจะเป็นสิ่งที่ก่อให้เกิดความสับสน หรือซีเรียส กำหนดของข้อสอบบางข้อได้
5. สถานภาพทางสังคมหรือความเกี่ยวข้องทางสังคม (Social Implication) ก็เป็นสิ่งที่ก่อให้เกิดความแตกต่างระหว่างกลุ่มผู้สอนแต่ละกลุ่มได้
6. ประสบการณ์หรือการฝึกฝนของแต่ละกลุ่มย่อย เป็นสิ่งที่ก่อให้เกิดการได้เปรียบเสียเปรียบของแต่ละกลุ่มค่อนข้างชัดเจน
7. องค์ประกอบทางวัฒนธรรม ความเป็นอยู่ ขนบธรรมเนียมประเพณี เชื้อชาติ ศาสนา ก็จะเอื้อให้กับบางกลุ่มย่อย จึงทำให้เกิดการได้เปรียบเสียเปรียบในบางเนื้อหาวิชาได้ นอกจากนี้การทำหน้าที่ต่างกันของข้อสอบอาจจะเกิดจากสาเหตุหรือแหล่งที่สำคัญ 2 แหล่งคือ

1. การเลือกเนื้อหา (Bias in Selection) คือผู้สร้างข้อสอบเลือกเนื้อหาเฉพาะส่วนใดส่วนหนึ่งมาสร้างข้อสอบ ทำให้ได้ข้อสอบที่มีเนื้อหาไม่ครอบคลุมและไม่ได้สัดส่วนที่สมดุลกัน

2. การสร้างข้อสอบ (Bias in Construction) กือการใช้ภาษาหรือข้อความบางอย่างที่เอื้อให้เกิดประโยชน์กับผู้สอบกลุ่มใดกลุ่มหนึ่ง

ตอนที่ 2 วิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบโดยวิธีชิปเทสท์

วิธีชิปเทสท์ พัฒนาโดยเชียลีและสโตต์ (Shealy & Stout, 1993) เพื่อใช้ตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ การทำหน้าที่ต่างกันของแบบสอบ และการทำหน้าที่ต่างกันของกลุ่มข้อสอบ (Differential Bundle Functioning: DBF) วิธีนี้สามารถคุณิเคราะห์ได้ทั้งในแบบสอบเอกมิติ และแบบสอบพหุมิติ (Stout, Li, & Nandakumar, 1997) วิธีชิปเทสท์ใช้สถิติ ทดสอบแบบนันพารา เมตริก ซึ่งพัฒนาบนพื้นฐานของทฤษฎี IRT ชนิดพหุมิติ แต่ไม่ต้องใช้ฟังก์ชันการตอบสนอง ข้อสอบหรือการประมาณค่าความสามารถแฝง วิธีชิปเทสท์ถูกออกแบบมาสำหรับการตรวจสอบ การทำหน้าที่ต่างกันของข้อสอบแบบมีทิศทางเดียว ตั้งนั้นจึงไม่มีความไวในการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบแบบไม่มีทิศทาง (Li & Stout, 1996) จุดเด่นของวิธีชิปเทสท์ คือ สามารถคำนวณได้ง่ายไม่ซับซ้อน ประยุกต์ค่าใช้จ่าย และไม่จำเป็นต้องใช้กลุ่มตัวอย่างที่มีขนาดใหญ่ ทั้งยังใช้สถิติทดสอบนัยสำคัญ (Narayanan & Swaminathan, 1996) นอกจากนี้ยังสามารถนำไปประยุกต์ใช้กับการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ ที่มีการให้คะแนนแบบพหุวิภาค (Polytomous DIF) (Chang, Mazzeo, & Roussos, 1995)

ในการศึกษาการทำหน้าที่ต่างกันของข้อสอบด้วยวิธีชิปเทสท์ของแบบสอบเอกมิติถือว่า ข้อสอบในแบบสอบจะต้องมุ่งวัดคุณลักษณะ หรือความสามารถแฝงเพียงลักษณะเดียว ความสามารถแฝงสามารถจำแนกเป็นความสามารถเป้าหมายที่ต้องการวัด (θ) กับความสามารถแทรกซ้อนที่ไม่ใช้เป้าหมายของการวัด (η) ตัวอย่าง เช่น แบบสอบคำศัพท์ในวิชาภาษาอังกฤษ ข้อสอบบางข้ออาจถูกประเมินตามความรู้สำหรับผู้ชายเป็นพิเศษ เช่น ความรู้เกี่ยวกับกีฬา ในขณะที่ข้อสอบบางข้ออาจถูกประเมินตามความรู้เกี่ยวกับผู้หญิง โดยเฉพาะ เช่น ความรู้เกี่ยวกับงานในบ้าน จากสถานการณ์ดังกล่าวทักษะความรู้เกี่ยวกับคำศัพท์ในวิชาภาษาอังกฤษ เป็นความสามารถเป้าหมายที่ต้องการวัด (θ) ส่วนความสามารถทางด้านกีฬาและงานในบ้านเป็นความสามารถแทรกซ้อน ที่ไม่ใช่เป้าหมายของ การวัด (η_1 และ η_2) ข้อสอบทุกข้อในแบบสอบจะวัดความสามารถเป้าหมาย ส่วนข้อสอบบางข้อทำหน้าที่ต่างกันจะวัดทั้งความสามารถเป้าหมาย และความสามารถแทรกซ้อน (Nandakumar, 1993)

ถ้าให้ฟังก์ชันการตอบสนองข้อสอบ (IRF) ข้อที่ i ซึ่งขึ้นอยู่กับความสามารถ θ เพียงอย่างเดียวแทนด้วย $P_i(\theta)$ ส่วน IRF ข้อที่ j ที่ขึ้นอยู่กับความสามารถทั้ง θ และ η แทนด้วย

$P_i(\theta, \eta)$ ฟังก์ชันการตอบสนองข้อสอบของข้อสอบดังกล่าวแบบ 3 พารามิเตอร์จะเป็นดังนี้
(Shealy & Stout, 1993)

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

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

ดังนี้ฟังก์ชันความน่าจะเป็นอย่างมีเงื่อนไขของแบบแผนการตอบข้อสอบทั้งฉบับเป็นดังนี้

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

เชียลล์และสโตต (Shealy & Stout, 1993) ได้ใช้ Marginal IRFs อธิบายการทำหน้าที่ต่างกันของข้อสอบ ดังนี้

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

เมื่อ $M_{ig}(\theta)$ = Marginal IRF สำหรับความสามารถเป้าหมาย θ ที่ต้องการวัดของผู้สอบกลุ่มอ้างอิงหรือกลุ่มเปรียบเทียบ

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

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

การเปรียบเทียบ Marginal IRF ระหว่างกลุ่มอ้างอิง (R) กับกลุ่มเปรียบเทียบ (F) จะทำให้ทราบถึงทิศทางของการได้เปรียบหรือเสียเปรียบ กล่าวคือ ถ้า $M_{if}(\theta) < M_{ir}(\theta)$ ทุกค่าของ θ แสดงว่า ข้อสอบทำหน้าที่ต่างกันแบบมีทิศทางเดียว โดยข้อสอบจะเข้าข้างผู้สอบกลุ่มอ้างอิงและถ้า $M_{if}(\theta) > M_{ir}(\theta)$ ทุกค่าของ θ แสดงว่าข้อสอบทำหน้าที่ต่างกันแบบมีทิศทางเดียว โดยข้อสอบจะเข้าข้างผู้สอบกลุ่มเปรียบเทียบ การทำหน้าที่ต่างกันของข้อสอบแบบมีทิศทางเดียวอาจเรียกอีกอย่างหนึ่งว่า “การทำหน้าที่ต่างกันแบบไม่ตัดกัน” (Noncrossing DIF)

ในการเปรียบเทียบผลการตอบข้อสอบระหว่างกลุ่มอ้างอิงและกลุ่มเปรียบเทียบตามวิธีซิปเทสท์ จะแบ่งข้อสอบออกเป็น 2 ชุดย่อย (Subtests) คือ (1) ชุดแบบสอบที่มีความตรง (Valid Subtests) หรือชุดแบบสอบที่ใช้ในการจับคู่เปรียบเทียบ (Matching Subtests) แบบสอบชุดนี้ประกอบด้วยข้อสอบที่ไม่ทำหน้าที่ต่างกัน และ (2) ชุดแบบสอบที่ต้องการศึกษา (Studied Subtests)

ประกอบด้วยข้อสอบที่สงสัยว่าทำหน้าที่ต่างกัน ถ้าแบบสอบชุดแรกมีจำนวน n ข้อ (ข้อที่ 1 ถึง n) แล้วแบบสอบชุดที่ 2 จะมีจำนวน $N-n$ ข้อ (ข้อที่ $n+1$ ถึง N) เมื่อ N เป็นจำนวนข้อสอบทั้งหมด

พึงก็ชันการตอบสนองข้อสอบของแบบสอบที่ต้องการศึกษา จากผู้สอบกลุ่มอ้างอิงและกลุ่มเปรียบเทียบ กำหนดในรูปพึงก์ชัน Marginal ดังนี้

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

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

เมื่อ $M_{SR}(\theta)$ = ผลรวมของ Marginal IRFs ของข้อสอบที่ต้องการศึกษาจากผู้สอบกลุ่มอ้างอิง ณ ระดับความสามารถ θ

$M_{SF}(\theta)$ = ผลรวมของ Marginal IRFs ของข้อสอบที่ต้องการศึกษาจากผู้สอบกลุ่มเปรียบเทียบ ณ ระดับความสามารถ θ

ขนาดของความแตกต่างระหว่าง $M_{SR}(\theta)$ และ $M_{SF}(\theta)$ แสดงถึงปริมาณความเข้มของ การทดสอบการทำหน้าที่ต่างกันของข้อทดสอบแบบมีทิศทางเดียว หรือการทำหน้าที่แบบไม่ตัดกันจากชุดแบบสอบตามที่ต้องการศึกษา ณ ระดับความสามารถ θ ซึ่งสามารถคำนวณโดยใช้ช่วงความสามารถ θ ของผู้สอบ ด้วยการอินทิเกรท ดังนี้

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

เมื่อ β_{un} = ดัชนีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบแบบมีทิศทางเดียว

$f_p(\theta)$ = พึงก์ชันความหนาแน่นของโอกาสการแยกแจงความสามารถ θ ของผู้สอบทั้ง 2 กลุ่ม

ดัชนี β_{un} ที่คำนวณได้จากสูตรดังกล่าว นำมาทดสอบสมมติฐานของการทำหน้าที่ต่างกันของข้อสอบแบบมีทิศทางเดียว ดังนี้

$$H_0 : \beta_{un} = 0$$

$$H_1 : \beta_{un} > 0$$

สมมติฐานอื่น (H_1) มีลักษณะทิศทางเดียวกัน ซึ่งใช้ทดสอบการทำหน้าที่ต่างกันของข้อสอบที่เข้าข้างผู้สอบกลุ่มอ้างอิง สำหรับค่าประมาณของ β_{un} คำนวณได้จากคะแนนรวมของชุดแบบสอบถามที่มีความตรงและชุดแบบสอบถามที่ต้องการศึกษา ซึ่งกำหนดด้วยสัญลักษณ์ดังนี้

$$X = \sum_{i=1}^n U_i$$

$$Y = \sum_{i=n+1}^N U_i$$

เมื่อ X = คะแนนรวมของชุดแบบสอบถามที่มีความตรง

Y = คะแนนรวมของชุดแบบสอบถามที่ต้องการศึกษา

U_i = ผลการตอบข้อสอบข้อที่ i (ตอบถูกได้ 1 คะแนน และตอบผิดได้ 0 คะแนน)

นำคะแนนรวมของชุดแบบสอบถามที่มีความตรง (X) เป็นเกณฑ์ในการจับคู่ที่มีความสามารถระดับเดียวกัน แล้วคำนวณคะแนนเฉลี่ยรายข้อจากผลการตอบข้อสอบชุดแบบสอบถามที่ต้องการศึกษา ระหว่างผู้สอบกลุ่มอ้างอิงและกลุ่มเปรียบเทียบที่มีความสามารถระดับเดียวกันมาจับคู่เปรียบเทียบกัน ณ $X = k$ โดยเขียนในรูปสัญลักษณ์ได้ ดังนี้

$$\bar{Y}_{Rk} - \bar{Y}_{Ik}; \quad k = 0, 1, 2, \dots, n$$

เมื่อ \bar{Y}_{Rk} = ค่าเฉลี่ยของคะแนนรายข้อจากชุดแบบสอบถามที่ต้องการศึกษาของผู้สอบกลุ่มอ้างอิง ซึ่งได้คะแนน $X = k$

\bar{Y}_{Ik} = ค่าเฉลี่ยของคะแนนรายข้อจากชุดแบบสอบถามที่ต้องการศึกษาของผู้สอบกลุ่มเปรียบเทียบ ซึ่งได้คะแนน $X = k$

k = คะแนนรวมจากชุดแบบสอบถามที่มีความตรง

ค่า $\bar{Y}_{Rk} - \bar{Y}_{Ik}$ ดังกล่าวเป็นความแตกต่างของผลการตอบข้อสอบในชุดแบบสอบถามที่ต้องการศึกษาระหว่างผู้สอบกลุ่มอ้างอิง และกลุ่มเปรียบเทียบที่มีความสามารถระดับเดียวกัน ถ้า $\bar{Y}_{Rk} - \bar{Y}_{Ik} = 0$ ทุกคะแนน k แสดงว่า ข้อสอบที่ต้องการศึกษาทำหน้าที่ไม่ต่างกัน แต่ถ้า $\bar{Y}_{Rk} - \bar{Y}_{Ik} > 0$ ทุกคะแนน k แสดงว่า ข้อสอบทำหน้าที่ต่างกันแบบมีทิศทางเดียว โดยจะเข้าข้างผู้สอบกลุ่มอ้างอิง ค่าความแตกต่างของผลการตอบข้อสอบสามารถประมาณค่าในรูป β_{un} ได้ดังนี้

$$\hat{\beta}_{un} = \sum_{k=0}^n \hat{P}_k (\bar{Y}_{Rk} - \bar{Y}_{Ik})$$

$$\text{โดยที่ } \hat{P}_k = \frac{(J_{Rk} + J_{Fk})}{\sum_{k=0}^n (J_{Rk} + J_{Fk})}$$

เมื่อ P_k = สัดส่วนของจำนวนผู้สอบกู้มเปรียบเทียบและกลุ่มอ้างอิงที่ได้ค่าคะแนนรวม $X = k$ จากจำนวนผู้สอบทั้งหมด
 J_{Fk} = จำนวนผู้สอบกู้มเปรียบเทียบที่ได้ค่าคะแนนรวม $X = k$
 J_{Rk} = แทนจำนวนผู้สอบกู้มอ้างอิงที่ได้ค่าคะแนนรวม $X = k$
 สำหรับการทดสอบสมมติฐานศูนย์ของ กอ DIF ใช้สถิติ β_{um} ดังนี้

$$\beta_{um} = \frac{\hat{\beta}_{um}}{\sigma(\hat{\beta}_{um})}$$

$$\text{โดยที่ } \hat{\sigma}(\hat{\beta}_{um}) = \sqrt{\sum_{k=0}^n \hat{P}_k^2 \left[\frac{1}{J_{Rk}} \hat{\sigma}^2(Y|k, R) + \frac{1}{J_{Fk}} \hat{\sigma}^2(Y|k, F) \right]}$$

เมื่อ $\sigma(\hat{\beta}_{um})$ = ค่าประมาณความคลาดเคลื่อนมาตรฐานของ β_{um}
 $\hat{\sigma}^2(Y|k, g)$ = ค่าประมาณความแปรปรวนของคะแนนจากชุดแบบสอบที่ต้องการศึกษา สำหรับผู้สอบกู้ม g (R หรือ F) ซึ่งมีคะแนนรวมเท่ากับ k
 J_{gk} = จำนวนผู้สอบกู้ม g (R หรือ F) ซึ่งตอบชุดแบบสอบที่มีความตรงແຕ່ວໄດ้ค่าคะแนนรวม $X = k$

β_{um} มีการแจกแจงในลักษณะปกติมาตรฐาน $[N(0,1)]$ จึงสามารถทดสอบนัยสำคัญด้วยสถิติทดสอบ Z ถ้าผลการทดสอบปรากฏว่า $\beta_{um} > Z_\alpha$ แสดงว่าการทดสอบมีนัยสำคัญจึงปฏิเสธ H_0 นั่นคือ ข้อสอบที่นำมาตรวจสอบทำหน้าที่ต่างกัน โดยจะเข้าข้างผู้สอบกู้มอ้างอิง เมื่อ β_{um} มีค่าเป็นบวก และจะเข้าข้างผู้สอบกู้มเปรียบเทียบเมื่อ β_{um} มีค่าเป็นลบ

สถิติที่ใช้ในการทดสอบสำหรับสรุปอ้างอิงการทำหน้าที่ต่างกันของข้อสอบดังกล่าว นักจะมีปัญหาในกรณีที่มีความแตกต่างของการแจกแจงความสามารถเป้าหมาย ระหว่างกลุ่มผู้สอบกล่าวคือ ถ้าผู้สอบกู้มอ้างอิงมีความสามารถสามารถเป้าหมายสูงกว่าผู้สอบกู้มเปรียบเทียบ จะเกิดผลกระทบทำให้สถิติ β_{um} มีค่าเพิ่ม (Inflate) หรือนมีค่าสูงผิดปกติ ถึงแม้ว่าในความเป็นจริงข้อสอบทำหน้าที่ไม่ต่างกัน จึงควรแก้ไขความแตกต่างของการแจกแจงความสามารถเป้าหมาย ด้วยการ

ปรับแก้ค่าความถดถอย เพื่อขัดผลผลกระทบดังกล่าว โดยปรับแก้ค่า \bar{Y}_{Rk} และ \bar{Y}_{Hk} เป็นรายคู่ ก่อนการคำนวณ β_{mn}

ตอนที่ 3 งานวิจัยที่เกี่ยวข้องกับการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ โดยวิธีชิปเทสท์

วิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบมีวิัฒนาการมาเป็นเวลากว่า โดยพยากรณ์แก้ไขข้อพกพร่องของวิธีการที่มีมาก่อน ประกอบกับความก้าวหน้าของเทคโนโลยีคอมพิวเตอร์ ทำให้มีการพัฒนาวิธีการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบหลายวิธี และได้มีผลการศึกษาค้นคว้าวิจัยโดยใช้ข้อมูลจำลองจำนวนมาก เพื่อหลักเลี่ยงข้อจำกัดบางประการของข้อมูลเชิงประจักษ์ งานวิจัยที่เกี่ยวข้องกับการทำหน้าที่ต่างกันของข้อสอบมีดังนี้

กาญจน วัฒนสุนทร (2538) "ได้พัฒนาเกณฑ์ในการตัดสินข้อสอบลำเอียงทางเพศ โดยใช้ข้อมูลเชิงประจักษ์" จากการสอบคัดเลือกเข้าระดับอุดมศึกษาปีการศึกษา 2535 ใช้วิธีการตรวจสอบ 3 วิธี คือ วิธีความแనวคิตทฤษฎีการตอบข้อสอบ วิธีแม่นแทลแลรอนส์เซล และวิธีชิปเทสท์ ดัชนีที่พัฒนาเพื่อเป็นเกณฑ์ในการตัดสินข้อสอบลำเอียงคือ SA, UA, α_{MH} , β_{SIB} ตามลำดับซึ่งในการวิจัยครั้งนี้ได้จากการวิเคราะห์ค่าเฉลี่ยของดัชนีแต่ละตัว ตัวแปรที่ใช้ในการศึกษา ได้แก่ ความยาวของแบบทดสอบขนาด 20, 30, 40 ข้อ สำหรับวิชาคณิตศาสตร์ และ 50, 60, 70 และ 80 ข้อ สำหรับวิชาภาษาอังกฤษ ขนาดกลุ่มตัวอย่างผู้ตอบข้อสอบ 100, 200, 400, 600, 800 และ 1,000 คน ผลการวิจัย พบว่า ขนาดผู้ตอบข้อสอบมีอิทธิพลต่อค่าเฉลี่ยของดัชนีทุกตัว ความยาวของแบบทดสอบมีผลต่อค่าเฉลี่ยของดัชนี SA และ UA แต่ไม่อิทธิพลต่อค่าเฉลี่ยของดัชนี α_{MH} และ β_{SIB} เกณฑ์ที่พัฒนาขึ้นเพื่อใช้ตัดสินความลำเอียง ระหว่างผู้ตอบข้อสอบเพศหญิงและเพศชาย มีดังนี้

1. $|SA| > .80$ และ $UA > .50$ เมื่อความยาวแบบทดสอบน้อยกว่า 50 ข้อ
2. $|SA| > .40$ และ $UA > 1.20$ เมื่อความยาวแบบทดสอบ 50 ข้อ ขึ้นไป
3. $\alpha_{MH} > .60$ และ $\alpha_{MH} < 1.40$ สำหรับทุกขนาดของผู้ตอบข้อสอบและความยาวแบบทดสอบ

4. $|\beta_{SIB}| > .06$ สำหรับทุกขนาดผู้ตอบข้อสอบ และความยาวแบบทดสอบที่นี้ในการใช้ดัชนี SA หรือ UA ควรใช้ผู้ตอบข้อสอบขนาด 800 คนขึ้นไป ส่วนดัชนี α_{MH} และ β_{SIB} ควรใช้ผู้ตอบข้อสอบอย่างน้อย 600 คน

จิตima วรรณศรี (2539) ได้เปรียบเทียบประสิทธิภาพการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบระหว่างวิธีแมลงเหล� - แ昏ส์เซลล์และวิธีชิปเพสท์ โดยศึกษาจากข้อมูลจำลอง ด้วยการที่ศึกษาได้แก่ ความยาวของแบบทดสอบ 3 ขนาด คือ 30, 60 และ 90 ข้อ ขนาดกลุ่มตัวอย่าง 3 ขนาด คือ 200, 600 และ 1,000 คน โดยแต่ละขนาดมีอัตราส่วนระหว่างผู้เข้าสอบกลุ่มอ้างอิงกับกลุ่มเปรียบเทียบต่างกัน คือ 1:1, 1:0.9, 1:0.75 และ 1:0.5 ผลการศึกษาปรากฏว่า วิธีแมลงเหล� - แ昏ส์เซลล์และวิธีชิปเพสท์มีประสิทธิภาพเท่าเทียมกัน ใน การตรวจสอบการทำหน้าที่ต่างกันของข้อสอบในทุกขนาดกลุ่มตัวอย่างและทุกอัตราส่วนภายใต้ความยาวของแบบทดสอบเดียวกัน และเมื่อใช้แบบทดสอบที่มีความยาวปานกลาง (60 ข้อ) ห้องสอบสามารถครองสอบได้อย่างมีประสิทธิภาพที่สุด นอกจากนี้เมื่อใช้ขนาดกลุ่มตัวอย่างมากขึ้นจะสามารถลดเวลาตรวจสอบข้อสอบทำหน้าที่ต่างกันได้ถูกต้องมากขึ้น โดยส่วนมากวิธีชิปเพสท์มีอัตราความคลาดเคลื่อนประมาณที่ 1 มากกว่าวิธีแมลงเหล�-แ昏ส์เซลล์เล็กน้อย

พรรณิ จินตมาส (2540) ได้ศึกษาผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบกับ เพศ จากแบบทดสอบคณิตศาสตร์โจทย์ปัญหาที่ผู้วิจัยสร้างขึ้นเอง โดยใช้วิเคราะห์การทำหน้าที่ต่างกันของข้อสอบ 3 วิธี คือ วิธีแปลงค่าความยาก วิธีแมลงเหล� - แ昏ส์เซลล์และวิธีชิปเพสท์ ในขนาดกลุ่มผู้เข้าสอบ 500 คน และขนาดกลุ่มผู้เข้าสอบ 1,000 คน โดยเปรียบเทียบจำนวนข้อสอบทำหน้าที่ต่างกัน และเปรียบเทียบค่าความเที่ยงแบ่งครึ่งบนของแบบทดสอบหลังคัดเลือกข้อสอบทำหน้าที่ต่างกันออกແล้า กลุ่มตัวอย่างเป็นนักเรียนชั้นมัธยมศึกษาปีที่ 1 ภาคเรียนที่ 2 ปีการศึกษา 2539 ของโรงเรียนสังกัดกรมสามัญศึกษาส่วนกลางจำนวน 2,200 คน ซึ่งได้มารายงานการสุ่มแบบชั้น มีขนาดโรงเรียนเป็นชั้นและโรงเรียนเป็นหน่วยการสุ่ม ผลการศึกษาปรากฏว่าเมื่อวิเคราะห์จากกลุ่มผู้เข้าสอบขนาด 500 คน วิธีชิปเพสท์พบข้อสอบทำหน้าที่ต่างกันมากที่สุด และวิธีแปลงค่าความยากพบข้อสอบทำหน้าที่ต่างกันน้อยที่สุด โดยจำนวนข้อสอบทำหน้าที่ต่างกันแตกต่างกันอย่างไม่มีนัยสำคัญทางสถิติกวิเคราะห์ และเมื่อวิเคราะห์จากกลุ่มผู้เข้าสอบ 1,000 คน วิธีแมลงเหล� - แ昏ส์เซลล์ พบรหัสผลการทำหน้าที่ต่างกันมากที่สุด วิธีแปลงค่าความยากไม่พบรหัสผลการทำหน้าที่ต่างกัน โดยจำนวนข้อสอบทำหน้าที่ต่างกัน จากการวิเคราะห์ด้วยวิธีแปลงค่าความยาก กับวิธีแมลงเหล� - แ昏ส์เซลล์ และวิธีแปลงค่าความยากกับวิธีชิปเพสท์แตกต่างกันอย่างนัยสำคัญทางสถิติที่ระดับ .05 นอกจากนี้แตกต่างกันอย่างไม่มีนัยสำคัญทางสถิติ สำหรับจำนวนข้อสอบทำหน้าที่ต่างกัน จากการวิเคราะห์ด้วยวิธีแปลงค่าความยาก ระหว่างกลุ่มผู้เข้าสอบ 500 คนและขนาดกลุ่มผู้เข้าสอบ 1,000 คน จะมีจำนวนข้อสอบทำหน้าที่ต่างกันแตกต่างกันอย่างมีนัยสำคัญทางสถิติที่ระดับ .05 นอกจากนี้มีค่าแตกต่างกันอย่างไม่มีนัยสำคัญทางสถิติ

วีโนน่า แซ่อ่อง (2543) เปรียบเทียบอำนาจการทดสอบและอัตราความคลาดเคลื่อน

ประเภทที่ 1 ในการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบonenกรุป ระหว่างวิธีชิบเทสท์ปรับใหม่ และวิธีการทดสอบโดยโลจิสติก โดยการศึกษาจำลองด้วยโมเดลโลจิสติกชนิด 3 พารามิเตอร์ ในเงื่อนไขค่า g 324 เงื่อนไข การศึกษาพบว่าอำนาจการทดสอบของวิธีชิบเทสท์ปรับใหม่และวิธีการทดสอบโดยโลจิสติกมีค่าเท่าเทียมกันทุกเงื่อนไข และทั้งสองวิธีดังกล่าวมีอำนาจการทดสอบสูงกว่าวิธีชิบเทสท์และเมนแทล - แซนส์เซลเกือบทุกเงื่อนไข อัตราความคลาดเคลื่อนประมาณ 10 % ที่ตรวจสอบด้วยวิธีการทั้ง 4 ข้างต้น มีค่าอยู่ภายในเกณฑ์ของอัตราความคลาดเคลื่อน ที่ระดับ 10 % เกือบทุกเงื่อนไข

นารายานาน และสวามินิธาน (Narayanan & Swaminathan, 1994, pp. 315- 328 อ้างถึงใน อารี วัชร โสดคิกุล, 2543, หน้า 28) ได้เปรียบเทียบผลการวิเคราะห์การทำหน้าที่ต่างกันของข้อสอบแบบเอกสาร ด้วยวิธีเมนแทล - แซนส์เซลและวิธีชิบเทสท์ โดยใช้ข้อมูลที่จำลองภายใต้เงื่อนไข 1,296 เงื่อนไข ($9 \times 3 \times 2 \times 4 \times 6$) และวัดจำนวนทั้ง 100 ครั้ง แต่ละเงื่อนไขประกอบด้วย

1. ขนาดกลุ่มตัวอย่างที่แตกต่างกัน 3 ขนาด คือกลุ่มเปรียบเทียบใช้ขนาด 100, 200 และ 300 คน กลุ่มอ้างอิงใช้ขนาด 300, 500 และ 1,000 คน วงไชวัkin ได้กลุ่มตัวอย่างทั้งหมด 9 กลุ่ม

2. ความแตกต่างในการแยกแจงความสามารถ 3 ระดับ ระดับที่ 1 ให้ทั้ง 2 กลุ่ม มีค่าเฉลี่ยการแยกแจงความสามารถที่ 0.0 ระดับที่ 2 ให้ค่าเฉลี่ยของความสามารถในกลุ่มอ้างอิงและกลุ่มเปรียบเทียบเป็น 0.0 และ 0.5 ตามลำดับ และระดับที่ 3 ค่าเฉลี่ยของความสามารถในกลุ่มอ้างอิงและเปรียบเทียบเป็น 0.0 และ -1.0 ตามลำดับ ส่วนเบี่ยงเบนมาตรฐานของทั้ง 3 ระดับ เป็น 1.0

3. ร้อยละของข้อสอบทำหน้าที่ต่างกันมี 2 ขนาด คือ ร้อยละ 10 และ ร้อยละ 20 จากแบบทดสอบจำนวน 40 ข้อ

4. ขนาดอิทธิพลของข้อสอบทำหน้าที่ต่างกัน 4 ระดับ คือ 0.4, 0.6, 0.8 และ 1.0
5. ลักษณะของข้อสอบ 6 ระดับ (ลักษณะ b ต่ำกับ a ปานกลาง, b ต่ำกับ a สูง, b ปานกลาง กับ a ต่ำ, b ปานกลางกับ a สูง, b สูงกับ a ต่ำ, b สูงกับ a ปานกลาง)

ผลการศึกษาปรากฏว่า วิธีเมนแทล-แซนส์เซล และวิธีชิบเทสท์สามารถตรวจพบข้อสอบทำหน้าที่ต่างกันได้ใกล้เคียงกัน โดยตรวจพบได้ดีในขนาดกลุ่มตัวอย่าง 600 คน (กลุ่มอ้างอิงและกลุ่มเปรียบเทียบกลุ่มละ 300 คน) ขนาดของข้อสอบทำหน้าที่ต่างกันทั้ง 4 ขนาดมีผลต่อการตรวจพบการทำหน้าที่ต่างกันของข้อสอบทั้งสองวิธี อย่างมีนัยสำคัญทางสถิติ ความแตกต่างในการแยกแจงความสามารถไม่มีผลต่อวิเคราะห์การทำหน้าที่ต่างกันโดยวิธีชิบเทสท์ แต่มีผลต่อวิธีเมนแทล-แซนส์เซล ที่เป็นเช่นนี้อาจเป็นเพราะวิธีชิบเทสท์ได้มีการปรับแก้การทดสอบขนาดของแบบทดสอบและร้อยละของข้อสอบทำหน้าที่ต่างกันไม่มีผลต่อวิธีทั้งสอง สำหรับการ

เกิดความคลาดเคลื่อนแบบที่ 1 วิธีแม่นเทล-แชนส์เซล จะมีอัตราการเกิดอยู่ในระดับปกติ แต่วิธีชิปเทสท์จะเพิ่มเล็กน้อยในกลุ่มที่มีความสามารถเท่ากัน แต่ในกลุ่มที่มีความสามารถไม่เท่ากันการเกิดความคลาดเคลื่อนแบบที่ 1 จะมีอัตราเกิดสูงขึ้นทั้งสองวิธี ซึ่งการเพิ่มความยาวของแบบทดสอบจะทำให้เกิดความเที่ยงสูงขึ้น มีผลให้ความคลาดเคลื่อนประเภทที่ 1 น่าจะมีอัตราการเกิดลดลง

นารายานและสวามินาทาน (Narayanan & Swaminathan, 1996, หน้า 257-274 ถ้าถูกใน บุญชริน ในโพธิ์, 2544, หน้า 46) ได้เปรียบเทียบการทำหน้าที่ต่างกันของข้อสอบแบบอนกรูป ด้วยวิธีแม่นเทล - แชนส์เซล วิธีดังดอยโดยจิสติกและวิธีชิปเทสท์ โดยศึกษาจำนวนการตรวจสอบและการจำแนกพิเศษโดยการจำลองข้อมูล ตัวแปรที่ศึกษา ได้แก่ ขนาดกลุ่มตัวอย่าง ความแตกต่างของการกระจายความสามารถระหว่างกลุ่มตัวอย่างกับกลุ่มเปรียบเทียบ สัดส่วนของข้อสอบทำหน้าที่ต่างกันที่มีภายในแบบทดสอบ ขนาดพื้นที่ระหว่างโถงลักษณะข้อสอบของผู้เข้าสอบสองกลุ่ม ค่าความยากและค่าอำนาจจำแนกของแบบทดสอบ ผลการศึกษาปรากฏว่าวิธีชิปเทสท์และวิธีดอยโดยจิสติก มีประสิทธิภาพในการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบแบบอนกรูป ได้เท่าที่ยอมกันในทุกเงื่อนไขที่ศึกษา ส่วนวิธีแม่นเทล - แชนส์เซลไม่สามารถตรวจสอบทำหน้าที่ต่างกันแบบอนกรูป และการจำแนกพิเศษวิธีแม่นเทล - แชนส์เซลจะสูงกว่าวิธีดอยโดยจิสติกและวิธีชิปเทสท์

บุญชริน ในโพธิ์ (2544) เปรียบเทียบความสอดคล้องของผลการตรวจสอบข้อสอบทำหน้าที่ต่างกันระหว่างวิธี ไอ - สแคร์ของลอร์ด วิธีแม่นเทล - แชนส์เซล และวิธีชิปเทสท์ โดยใช้กลุ่มตัวอย่างเป็นนักเรียนชั้นมัธยมศึกษาปีที่ 1 ปีการศึกษา 2543 ในเขตกรุงเทพมหานคร จำนวน 790 คน เครื่องมือที่ใช้วัดเป็นแบบทดสอบวัดความถนัดทางคณิตศาสตร์ จำนวน 60 ข้อ ผลการวิจัยพบว่า

1. ข้อสอบทำหน้าที่ต่างกันจากการตรวจสอบด้วย วิธี ไอ - สแคร์ของลอร์ด, วิธีแม่นเทล - แชนส์เซล และวิธีชิปเทสท์ มีจำนวนข้อสอบที่ตรวจพบแตกต่างกัน เมื่อขนาดกลุ่มตัวอย่างและจำนวนข้อสอบต่างกัน ได้แก่ กรณีแบบทดสอบ 30 ข้อ ที่ใช้กลุ่มตัวอย่างที่มีขนาด 300 คน และ 500 คน แบบทดสอบ 40 ข้อ ใช้กลุ่มตัวอย่าง 500 คน และ 700 คน

2. ความสอดคล้องของผลการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ เมื่อพิจารณาจากข้อสอบจำนวน 20 ข้อ และ 30 ข้อ ที่ตรวจสอบโดยวิธี ไอ - สแคร์ของลอร์ด วิธีแม่นเทล - แชนส์เซลและวิธีชิปเทสท์ พบรากลุ่มตัวอย่างขนาด 500 คน และ 700 คน มีจำนวนข้อสอบที่ตรวจพบสอดคล้องกันสูงสุด ส่วนข้อสอบจำนวน 40 ข้อ ซึ่งตรวจด้วยวิธี ไอ - สแคร์ของลอร์ด พบรากลุ่มตัวอย่างขนาด 300 คน และ 500 คน มีจำนวนข้อสอบที่ตรวจพบสอดคล้องกันสูงสุด

แต่เมื่อตรวจสอบด้วยวิธีแบบเกล - แอนส์เชลและวิธีชิปเพสท์ พบรากคุณด้วยย่างขนาด 500 และ 700 คน มีจำนวนข้อสอบที่ตรวจพบสอดคล้องกันสูงสุด

สิริรัตน์ วิภาสศิลป์ (2545) ได้เปรียบเทียบการตรวจสอบการทำหน้าที่ต่างกันของ ข้อสอบ หมวดข้อสอบและแบบทดสอบ จากข้อมูลการตอบข้อสอบที่ใช้ความสามารถทางสมอง ระหว่างวิธีชิปเพสท์กับวิธีดีอฟไอย์ ภายใต้เงื่อนไขความยาวของแบบทดสอบ 30, 40 และ 50 ข้อ กว่าคุณด้วยขนาด 50, 100, 200, 500, และ 1,000 คน กลุ่มด้วยย่างการศึกษาได้จากการสุ่มแบบ ใส่คืนจากประชาชนเทียม ซึ่งกำหนดจากนักเรียนชายและนักเรียนหญิงชั้นมัธยมศึกษาปีที่ 1 จังหวัดคุณฑบุรี แต่ละขนาดสุ่มกลุ่มด้วยย่าง 50 ครั้ง เครื่องมือที่ใช้ในการวิจัยเป็นแบบทดสอบ วิชาคณิตศาสตร์ที่ผู้วุฒิสร้างขึ้น ซึ่งเป็นข้อสอบแบบหลักๆ ตัวเลือก 5 ตัวเลือก จำนวน 50 ข้อ มี ข้อสอบที่ผู้เชี่ยวชาญพิจารณาว่าเป็นข้อสอบทำหน้าที่ต่างกันต่อเพศชายจำนวน 16 ข้อ หลังจาก เก็บรวบรวมข้อมูลแล้วคัดเลือกข้อสอบตามสัดส่วนในตารางสัดส่วนกำหนด จัดเป็นแบบทดสอบที่ มีความยาว 30 และ 40 ข้อ แล้วตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ หมวดข้อสอบ และ แบบทดสอบในเงื่อนไขดังๆ โดยโปรแกรมชิปเพสท์และดีอฟไอย์ จากนั้นนำผลที่ได้ไป

เปรียบเทียบความถูกต้อง และการระบุผิดพลาดในการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบ ด้วยวิธีเดียวกันและต่างกัน โดยการวิเคราะห์ความแปรปรวนแบบด้วยแพรพหุ คำนวณความ สอดคล้องในการตรวจสอบการทำหน้าที่ต่างกันของหมวดข้อสอบ และแบบทดสอบด้วยวิธี ชิปเพสท์และวิธีดีอฟไอย์ โดยใช้สถิติ Z-test ผลการศึกษาปรากฏว่า

- เมื่อแบบทดสอบประกอบด้วยข้อสอบ 30, 40 และ 50 ข้อ กว่าคุณด้วยขนาด 50, 100 และ 200 คน ส่งผลต่อความถูกต้องในการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบด้วยวิธี ชิปเพสท์ไม่แตกต่างกัน กลุ่มด้วยย่างขนาด 500 และ 1,000 คน ส่งผลต่อความถูกต้องในการ ตรวจสอบการทำหน้าที่ต่างกันของข้อสอบด้วยวิธีชิปเพสท์สูงกว่ากลุ่มด้วยย่างขนาด 50, 100 และ 200 คน แต่การระบุผิดพลาดในการตรวจสอบสูงกว่าด้วย เมื่อตรวจสอบด้วยวิธีดีอฟไอย์ พบราก กลุ่มด้วยย่างขนาด 50, 100, 200, 500 และ 1,000 คน ส่งผลต่อความถูกต้องในการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบไม่แตกต่างกัน

- ทุกเงื่อนไขความยาวแบบทดสอบและกลุ่มด้วยย่างขนาดแตกต่างกัน วิธีชิปเพสท์มี ความถูกต้องในการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบน้อยกว่าวิธีดีอฟไอย์ และพบว่า ความสอดคล้องในการตรวจสอบการทำหน้าที่ต่างกันของข้อสอบด้วยวิธีทั้งสองมีค่าต่ำกว่าร้อยละ 1

- วิธีชิปเพสท์มีความถูกต้องในการตรวจสอบการทำหน้าที่ต่างกันของหมวดข้อสอบ มากกว่าดีอฟไอย์ เมื่อแบบทดสอบมีข้อสอบ 30 ข้อ กลุ่มด้วยย่าง 1,000 คน และเมื่อแบบทดสอบ แบบมี 40 ข้อ กลุ่มด้วยย่างขนาด 500 คน

4. วิชีซิปเพสท์มีความถูกต้องในการตรวจสอบการทำหน้าที่ด่างกันของแบบทดสอบมากกว่าวิชีดีอีฟ์ไอที เมื่อแบบทดสอบมีข้อสอบ 50 ข้อ กลุ่มตัวอย่างขนาด 100, 200 และ 1,000 คน

ราจู แวนเดอร์ ลินเดน และฟลีร์ (Raju, Vander, Linden, & Fleer, 1995, pp. 353 – 368 อ้างถึงใน สิริรัตน์ วิภาสศิลป์, 2545, หน้า 40) เสนอวิธีการตรวจสอบการทำหน้าที่เบี่ยงเบนของข้อสอบและแบบทดสอบ (Internal Measures of Differential Functioning of Items and Tests: DFIT) โดยอาศัยกรอบความคิดของทฤษฎีการตอบข้อสอบ ศึกษาโดยใช้ข้อมูลจำลองที่สร้างขึ้นมาใช้แบบจำลองโลจิสติก ชนิด 2 พารามิเตอร์ ด้วยโปรแกรม RANDGEN (Fleer, Kiley, & Raju, 1991) สร้างแบบทดสอบที่มีความยาว 40 ข้อ ประกอบด้วยข้อสอบที่แสดง DIF 0 %, 5 %, 10 % และ 20 % ขนาดของกลุ่มตัวอย่างกลุ่มละ 500 คน สำหรับเงื่อนไขกลุ่มตัวอย่างขนาดเล็ก และกลุ่มละ 1000 คน สำหรับกลุ่มตัวอย่างขนาดใหญ่ ตรวจสอบการทำหน้าที่เบี่ยงเบนในเงื่อนไขต่างๆ ด้วยวิธีการของไค – สแคร์บของลอร์ด การวัดพื้นที่แบบคิดเครื่องหมาย การวัดพื้นที่โดยไม่คิดเครื่องหมายและวิชีดีอีฟ์ไอที ผลที่ได้พบว่าวิธีการดีอีฟ์ไอทีสามารถตรวจสอบการทำหน้าที่เบี่ยงเบนของข้อสอบทั้งแบบที่มีการขาดเชยและแบบที่ไม่มีการขาดเชย (CDIF และ NCDIF) ได้อย่างถูกต้อง

ดักล้าส รูสโซ และ สเตาต์ (Douglas, Roussos, & Stout, 1996, pp. 465 – 484 อ้างถึงใน สิริรัตน์ วิภาสศิลป์, 2545, หน้า 40) ได้ตรวจสอบการทำหน้าที่เบี่ยงเบนของข้อสอบ โดยพิจารณาข้อสอบพร้อม ๆ กันครั้งละหลายข้อ โดยเสนอวิธีการตรวจสอบหมวดข้อสอบที่คาดว่าจะแสดงการทำหน้าที่เบี่ยงเบนของข้อสอบเหล่านี้ตามมาพิจารณาพร้อมกัน 2 วิธี คือ ใช้เฉพาะความคิดเห็นของผู้เชี่ยวชาญ และใช้วิธีทางสถิติและตรวจสอบข้อความเห็นของผู้เชี่ยวชาญ พร้อมทั้งแสดงตัวอย่างตรวจสอบการทำหน้าที่เบี่ยงเบน 3 ตัวอย่าง ดังนี้ ตัวอย่างที่ 1 เลือกหมวดข้อสอบโดยใช้ความเห็นของผู้เชี่ยวชาญอย่างเดียว ให้ผู้เชี่ยวชาญชาย 3 คน และหญิง 1 คน พิจารณาแบบทดสอบย่อยใช้เหตุผลเชิงตรรกศาสตร์ (Logical Reasoning Subtest) ที่ได้จากการดำเนินการสอบในเดือนธันวาคม ค.ศ. 1991 (December 1991 Administation of the Low School Admission Test: LSAT) จำนวน 49 ข้อ ในการตรวจให้เพศชายเป็นกลุ่มอ้างอิงและเพศหญิงเป็นกลุ่มสนใจ ผู้เชี่ยวชาญจัดหมวดข้อสอบออกเป็น 8 หมวด และพิจารณาว่าหมวดใดน่าจะให้ประโยชน์แก่เพศชายหรือเพศหญิง หลังจากนั้นวิเคราะห์ด้วยโปรแกรมคอมพิวเตอร์สำเร็จรูปซิปเพสท์ ใช้ก่อตัวตัวอย่างเพศชาย 3,000 คน พบว่าความคิดเห็นของผู้เชี่ยวชาญและผลวิเคราะห์ด้วยโปรแกรมซิปเพสท์ สอดคล้องกัน จึงทำการวิเคราะห์ปริมาณ DIF ของแต่ละข้อ พบว่าเดลล์คะแนน DIF ด้วยปริมาณที่น้อยมาก จึงไม่มีการคัดเลือกข้อสอบออกจากแบบทดสอบ ตัวอย่างที่ 2 ใช้วิธีการทางสถิติในการตรวจสอบหมวดข้อสอบ คือ วิธี HCA (Agglomerative Hierachical Cluster Analysis;

Jian & Dubes, 1988) และ DIMTEST (Nandacumar & Stout, 1993; Stout, 1987) แล้วจึงตรวจสอบด้วยโปรแกรมซิปเทสท์ ในการตรวจสอบใช้ข้อสอบ (National Assessment of Educational Progress (NEAP) จำนวน 36 ข้อ ใช้กู้มตัวอย่างจำนวน 500 คน ผลการตรวจสอบพบว่าการตรวจสอบการทำหน้าที่เบี่ยงเบนของหมวดข้อสอบ ประสบความสำเร็จเป็นอย่างดี ตัวอย่างที่ 3 ใช้แบบทดสอบความเข้าใจในการอ่าน ที่สอบในเดือนธันวาคม ค.ศ. 1991 ซึ่งเป็นส่วนหนึ่งของการทดสอบ LSAT แบ่งเป็นบทความที่กำหนดให้อ่านเป็น 4 บทความ แต่ละบทความมีข้อสอบให้อ่าน 5-8 ข้อ แล้วทำการตรวจสอบหมวดข้อสอบตามโดยใช้วิธี เอชซีเอ และ ดิมเทสท์ หลังจากนั้นจึงวิเคราะห์ด้วยโปรแกรมซิปเทสท์ ใช้กู้มตัวอย่างกลุ่มละ 1,000 คน พบว่าทั้ง 4 บทความ แสดงการทำหน้าที่เบี่ยงเบนของหมวดข้อสอบ ข้อสอบ 2 หมวดให้ประโยชน์แก่เพศชาย และ 2 หมวด ให้ประโยชน์แก่เพศหญิง

สมາลี แก้วทันงค์ (2547) ได้ศึกษาการทำหน้าที่ต่างกันของข้อสอบด้านเพศ ภาษาพูด และเชื้อชาติ ของแบบทดสอบสาระการเรียนรู้ภาษาไทย และสาระการเรียนรู้สังคมศึกษา ศาสนา และวัฒนธรรม และตรวจสอบสาเหตุของการทำหน้าที่ต่างกันของข้อสอบด้านเพศ ภาษาพูด และ เชื้อชาติ ของแบบทดสอบกลุ่มสาระการเรียนรู้ภาษาไทย และกลุ่มสาระการเรียนรู้สังคมศึกษา ศาสนา และวัฒนธรรม กลุ่มตัวอย่างที่ใช้ในการวิจัย คือ นักเรียนระดับชั้นมัธยมศึกษาปีที่ 1 ปีการศึกษา 2546 สำนักงานเขตพื้นที่การศึกษาสงขลา พัทลุง ตรัง ยะลา จำนวน 1,320 คน เครื่องมือที่ใช้ในการวิจัยประกอบด้วยแบบสอบถามเพื่อสำรวจความสนใจในการอ่านของผู้เรียน แบบสอบถามในกลุ่มสาระการเรียนรู้ภาษาไทยและกลุ่มสาระการเรียนรู้สังคมศึกษา ศาสนา และวัฒนธรรม และแบบวินิจฉัยสาเหตุการทำหน้าที่ต่างกันของข้อสอบ ในกลุ่มสาระการเรียนรู้ภาษาไทยและสังคมศึกษา ศาสนา และวัฒนธรรม วิเคราะห์การทำหน้าที่ต่างกันของข้อสอบโดยวิธีแม่นเทล - แนนส์เซล และโปรแกรมซิปเทสท์ และตรวจสอบสาเหตุการทำหน้าที่ต่างกันโดยผู้เชี่ยวชาญ ผลการวิจัยพบว่า

1. แบบสอบถามกลุ่มสาระการเรียนรู้ภาษาไทยมี 3 ฉบับ ประกอบด้วยข้อสอบจำนวน 120 ข้อ มีข้อสอบที่ทำหน้าที่ต่างกันด้านเพศ 9 ข้อ ภาษาพูด 15 ข้อ และเชื้อชาติ 28 ข้อ และกลุ่มสาระการเรียนรู้สังคมศึกษา ศาสนา และวัฒนธรรมมีข้อสอบที่ทำหน้าที่ต่างกันด้านเพศ 22 ข้อ ด้านภาษาพูด 52 ข้อ และเชื้อชาติ 20 ข้อ

2. ข้อสอบที่ทำหน้าที่ต่างกันด้านเพศของแบบสอบถามสาระการเรียนรู้ภาษาไทย ส่วนใหญ่มีสาเหตุมาจากการเนื้อร่องที่น่าสนใจและภาษาที่ใช้ในแบบสอบถาม โดยมีลักษณะของข้อสอบที่ออกตามความสนใจและการใช้ภาษาในแบบสอบถามที่เป็นสำนวนและคำศัพท์ และข้อสอบที่ทำหน้าที่ต่างกันด้านเพศ ของแบบสอบถามสาระการเรียนรู้สังคมศึกษา ศาสนา และวัฒนธรรม มีสาเหตุมาจาก

เนื้อเรื่องที่สนใจและเนื้อเรื่องเกี่ยวกับวัฒนธรรม โดยมีลักษณะของข้อสอบที่มีเนื้อเรื่องที่สนใจ และมีเนื้อเรื่องเกี่ยวกับวัฒนธรรมของผู้สอน

3. ข้อสอบที่ทำหน้าที่ต่างกันค้านภาษาพูดของแบบสอบสาระการเรียนรู้ภาษาไทย ส่วนใหญ่มีสาเหตุมาจาก เนื้อเรื่องที่สนใจ และภาษาที่ใช้ในแบบสอบ โดยมีลักษณะข้อสอบที่ออกแบบ ความสนใจและมีการใช้ภาษาที่เป็นคำศัพท์เฉพาะ เช่น ราชศัพท์ และข้อสอบที่ทำหน้าที่ต่างกัน ค้านภาษาพูดของแบบสอบสาระการเรียนรู้สังคมศึกษา ศาสนา และวัฒนธรรม มีสาเหตุมาจาก เนื้อเรื่องที่สนใจ และเนื้อเรื่องเกี่ยวกับวัฒนธรรม โดยมีลักษณะของข้อสอบที่ออกแบบความสนใจ เช่น ภูมิศาสตร์ และในเรื่องเกี่ยวกับวัฒนธรรมไทย

4. ข้อสอบที่ทำหน้าที่ต่างกันในด้านเชื้อชาติของแบบสอบสาระการเรียนรู้ภาษาไทย ส่วนใหญ่มีสาเหตุมาจากเนื้อเรื่องที่สนใจ และภาษาที่ใช้ในแบบสอบตาม โดยมีลักษณะข้อสอบที่ออกแบบความสนใจ และมีการใช้คำราชศัพท์หรือคำศัพท์ยากมาออกข้อสอบ และข้อสอบที่ทำหน้าที่ต่างกันในด้านเชื้อชาติของแบบสอบสาระการเรียนรู้สังคมศึกษา ศาสนา และวัฒนธรรม มีสาเหตุจากเนื้อเรื่องที่สนใจและเนื้อเรื่องเกี่ยวกับวัฒนธรรมและประเพณี

ศุภวัฒน์ มะลิเพ็อก (2548) ได้ศึกษาอิทธิพลของการทำหน้าที่ต่างกันของข้อสอบที่ส่งผลต่อคุณภาพของแบบทดสอบวัดผลสัมฤทธิ์ทางการเรียนระดับชาติวิชาภาษาไทย ชั้นประถมศึกษา ปีที่ 6 ปีการศึกษา 2546 ด้วยวิธีซิปเทสท์กับวิธีอุดถอยโลจิสติก โดยการเปรียบเทียบค่าความเที่ยง ความตรงเชิง โครงสร้าง ความคงที่ของโครงสร้างองค์ประกอบ และค่าสัมประสิทธิ์สหสัมพันธ์ อันดับของผู้สอน ระหว่างแบบทดสอบฉบับก่อนกับฉบับหลังตัดข้อสอบทำหน้าที่ต่างกันออก ผลปรากฏว่า พนข้อสอบที่ต่างกันเมื่อจำแนกกลุ่มผู้สอนตามด้วยประเภทจำนวน 12 ข้อ กิตเป็นร้อยละ 30 แบบทดสอบฉบับก่อนกับฉบับหลังตัดข้อสอบทำหน้าที่ต่างกันมีค่าความเที่ยงแตกต่างกัน โครงสร้างองค์ประกอบคงที่และค่าสัมประสิทธิ์สหสัมพันธ์อันดับของผู้สอน มีความสัมพันธ์กันอย่างมีนัยสำคัญทางสถิติที่ระดับ 0.5

รูสโซ และสโตท (Roussos & Stout, 1996, pp. 215-230 อ้างถึงใน ปิยะพิพิญ ตินวร, 2549, หน้า 34) ได้ศึกษาสถานการณ์จำลองของผลกระทบของกลุ่มตัวอย่างขนาดเล็ก และค่าพารามิเตอร์ของข้อสอบที่มีต่ออัตราความคลาดเคลื่อนประเภทที่ 1 ด้วยวิธีซิปเทสท์และวิธีแมนเทล-แซนส์เซล โดยจำลองข้อมูล 2 ครั้ง ครั้งแรกศึกษากลุ่มตัวอย่างขนาดเล็ก ใช้กลุ่มตัวอย่าง 4 ขนาด (100, 200, 500, และ 1,000) และความแตกต่างของการแจกแจงค่าความสามารถระหว่างผู้เข้าสอบกลุ่มอ้างอิงกับกลุ่มเปรียบเทียบ 3 ขนาด (0, 0.5 และ 1) ใช้แบบทดสอบ 25 ข้อ ครั้งที่ 2 ศึกษาค่าพารามิเตอร์ของข้อสอบ ใช้กลุ่มตัวอย่าง 3 ขนาด (500, 1,000 และ 3,000) และความแตกต่างของการแจกแจงค่าความสามารถระหว่างผู้เข้าสอบกลุ่มอ้างอิงกับกลุ่มเปรียบเทียบ 2 ขนาด

(0 และ 1) เลือกค่าพารามิเตอร์อ่านจากจำแนก 3 ค่า (0.4, 1.0 และ 2.5) พารามิเตอร์ความยาก 5 ค่า (-1.5, -0.5, 0.0.5 และ 1.5) และค่าพารามิเตอร์การเดา 1 ค่า (.20) ผลการศึกษาพบว่า ครั้งที่ 1 เมื่อศึกษากลุ่มตัวอย่างขนาดเล็ก พบว่าอัตราความคลาดเคลื่อนประเพณีที่ 1 ระหว่างวิธีซิปเหล็กและวิธีแม่นเหล็ก-แ昏ส์เซล มีค่าไม่แตกต่างกัน ครั้งที่ 2 เมื่อศึกษาค่าพารามิเตอร์ของข้อสอบ พบว่าเมื่อความแตกต่างของการแจกแจงค่าความสามารถ ระหว่างผู้เข้าสอบกลุ่มอ้างอิงกับกลุ่มเปรียบเทียบ ไม่แตกต่างกัน และอัตราความคลาดเคลื่อนประเพณีที่ 1 ของวิธีซิปเหล็กมีค่าต่ำกว่าวิธีแม่นเหล็ก-แ昏ส์เซลในทุกเงื่อนไข

ตอนที่ 4 ผลสัมฤทธิ์ทางการเรียน

ผลสัมฤทธิ์ทางการเรียน หมายถึงคุณลักษณะรวมถึงความรู้ความสามารถของบุคคลอันเป็นผลมาจากการเรียนการสอน หรือมวลประสบการณ์ทั้งปวงที่บุคคลได้รับจากการเรียนการสอน ทำให้บุคคลเกิดการเปลี่ยนแปลงพฤติกรรมในด้านต่างๆ ของสมรรถภาพสมอง แสดงว่าเป็นการตรวจสอบพฤติกรรมของผู้เรียนในด้านพุทธิพิสัยนั้นเอง โดยการวัด 2 องค์ประกอบนั้น คือ

1. การวัดด้านปฏิบัติ โดยให้ผู้เรียนได้ลงมือปฏิบัติจริงให้เป็นผลงานปรากฏออกมາ ทำให้สังเกตและวัดได้ เช่น ศิลปศึกษา พลศึกษา เป็นต้น การวัดแบบนี้จึงค้องใช้ข้อสอบภาคปฏิบัติ ซึ่งเป็นการประเมินผลและพิจารณาที่วิธีปฏิบัติและผลงานที่ปฏิบัติ

2. การวัดด้านเนื้อหาเป็นการตรวจสอบความรู้ความสามารถเกี่ยวกับเนื้อหาวิชา รวมถึง พฤติกรรมความสามารถในด้านต่างๆ อันเป็นผลมาจากการทำกิจกรรมการเรียนการสอน มีวิธีสอบ วัดได้ 2 ลักษณะ คือ การสอบปากเปล่า และการสอบแบบให้เขียนตอบ ในการสอบแบบเขียนตอบ นี้จะใช้กันแพร่หลาย เครื่องมือที่ใช้ในการสอบวัด เรียกว่า “ข้อสอบวัดผลสัมฤทธิ์” หรือ “แบบทดสอบวัดผลสัมฤทธิ์” เพื่อให้การวัดผลสัมฤทธิ์ทางการเรียนเกิดประสิทธิภาพสามารถแยก พฤติกรรมด้านต่างๆ ได้ชัดเจน จึงต้องมีการสร้างแบบทดสอบที่เหมาะสมกับจุดมุ่งหมายของ เนื้อหาขึ้น

จะเห็นว่าผลสัมฤทธิ์ทางการเรียน เป็นผลที่เกิดขึ้นจากการเรียนการสอนหรือมวลประสบการณ์ที่ได้รับจากสิ่งแวดล้อม ซึ่งทำให้บุคคลเกิดการเปลี่ยนแปลงพฤติกรรมในด้านต่างๆ ของสมรรถภาพทางสมอง ซึ่งในการศึกษาครั้งนี้ผู้รายงานใช้แบบทดสอบในลักษณะเขียนตอบ โดยเรียกว่าแบบทดสอบวัดผลสัมฤทธิ์ทางการเรียน เพราะสามารถทราบสมรรถภาพทางสมองและ สามารถวัดพฤติกรรมต่างๆ ได้ชัดเจนครบถ้วนตามจุดมุ่งหมายที่ต้องการวัด

จุดมุ่งหมายของการวัดผลสัมฤทธิ์ทางการเรียน

การวัดผลสัมฤทธิ์ทางการเรียนมีจุดมุ่งหมายเพื่อ

1. ตรวจสอบความสามารถด้านพุทธิพิสัย เพื่อเปรียบเทียบการเปลี่ยนแปลงในการเรียนรู้ โดยวัดทางด้านเนื้อหา

2. ปรับปรุงแก้ไขกระบวนการเรียนการสอน โดยใช้วัดผลเป็นองค์ประกอบหนึ่งของการเรียนการสอน

3. ประเมินผลผู้เรียนว่าบรรลุเป้าหมายในสิ่งที่สอนไปหรือไม่ โดยการประเมินผลผู้เรียน มีประโยชน์ต่อการจัดการศึกษาดัง

3.1 ช่วยสำรวจความต้องการด้านการศึกษาของแต่ละบุคคล

3.2 ช่วยตรวจสอบประสิทธิผลของวิธีการจัดการเรียนการสอน

3.3 ช่วยให้สารสนเทศแก่ผู้เรียน อาจารย์ที่ปรึกษาและนักแนะแนว

3.4 ช่วยให้สารสนเทศแก่ผู้บริหารการศึกษา เพื่อนำไปใช้ในการจัดทรัพยากรของโรงเรียน

3.5 ช่วยให้สารสนเทศแก่ผู้บริหารการศึกษา เพื่อนำไปใช้จัดการบริหารงานบุคคล
3.6 เพื่อเลือกวิธีหรือกระบวนการสอนที่เหมาะสม ที่จะเป็นแรงกระตุ้นให้ผู้เรียนมีความสนใจในการเรียนยิ่งขึ้น

วิธีการทดสอบวัดผลสัมฤทธิ์ทางการเรียน

วิธีการทดสอบวัดผลสัมฤทธิ์ทางการเรียนมี 2 วิธี

1. การทดสอบแบบอิงคู่ เป็นการทดสอบที่เปลี่ยนความหมายของคะแนนโดยการนำผลการปฏิบัติงาน ไปเปรียบเทียบกับผลการปฏิบัติงานของคนอื่น ๆ ภายในกลุ่ม การรายงานผลการทดสอบใช้คะแนนมาตรฐานในการบ่งบอก เช่น เปรอร์เซ็นต์ไทย เกรดเทียบเท่า (Grade Equivalent)

2. การทดสอบแบบอิงเกณฑ์ เป็นการทดสอบที่เปลี่ยนความหมายของคะแนนโดยการนำผลการปฏิบัติงาน ไปเปรียบเทียบกับมาตรฐานที่แท้จริง (Absolute Standard) ซึ่งเป็นเกณฑ์ภายนอกกลุ่มที่กำหนดไว้อย่างรอบคอบโดยไม่เปรียบเทียบกับผลงานคนอื่นภายในกลุ่ม ดังนั้นผลงานของนักเรียนจะอยู่ในระดับมาตรฐานหรือไม่ ต้องพิจารณาและเปรียบเทียบกับมาตรฐานที่แท้จริงเท่านั้น การรายงานผลการทดสอบเสนอในพจน์ของจำนวน หรือเปอร์เซ็นต์การตอบถูกของแต่ละบุคคล แบบทดสอบแบบอิงเกณฑ์จำแนกเป็น 2 ประเภท ดังนี้

2.1 แบบทดสอบอิงจุดประสงค์ เป็นแบบทดสอบอิงเกณฑ์ที่สร้างขึ้นโดยใช้ชุดประสงค์รายวิชา ส่วนมากจะมีการกำหนดเกณฑ์เพื่อบ่งชี้ระดับความรอบรู้ของผู้สอบ ซึ่งมักจะใช้คะแนนจุดตัดของแบบทดสอบ

2.2 แบบทดสอบอิงความรู้ เป็นแบบทดสอบอิงเกณฑ์ที่สร้างขึ้นโดยยึดหลัก

เฉพาะของมวลความรู้ ใช้การประเมินความสามารถของผู้สอบที่สามารถตอบข้อสอบถูกในประชากรข้อสอบ

แบบทดสอบวัดผลสัมฤทธิ์ทางการเรียน (Achievement Test)

แนวคิดของแบบทดสอบวัดผลสัมฤทธิ์ทางการเรียน (เยาวศิริ วิญญูลย์ศรี, 2545,
หน้า 14-27)

1. แบบทดสอบวัดผลสัมฤทธิ์ทางการเรียน (Achievement Test) เป็นเครื่องมือสำหรับช่วยให้ครุศาสตร์ตัดสินผลสัมฤทธิ์ทางการเรียนได้อย่างมีประสิทธิภาพ เพราะเป็นวิธีการประเมินพฤติกรรมของนักเรียนที่มีความเป็นอิสระได้มากกว่าวิธีอื่น ๆ แบบทดสอบวัดผลสัมฤทธิ์ทางการเรียนที่ใช้ในโรงเรียนมุ่งวัดความรู้ในแต่ละวิชา และทักษะต่าง ๆ โดยมีวัตถุประสงค์พื้นฐานที่สำคัญ 2 ประการ คือ ประการแรก เพื่อเป็นเครื่องมือในการวัดผลสัมฤทธิ์ทางการเรียน อันเป็นข้อมูลที่ได้รับการประเมินผลการเรียนการสอนเป็นรายบุคคล ประการที่สอง เพื่อเป็นการตรวจสอบความสามารถของนักเรียนแต่ละคนซึ่งแตกต่างกันโดยธรรมชาติ เมื่อความสามารถของแต่ละคนมีความแตกต่างกันทั้งเด็กและผู้ใหญ่ ดังนั้นการที่จะให้นักเรียนต่าง ๆ ได้รับการพัฒนาความสามารถเฉพาะตนที่มีอยู่อย่างเหมาะสม จึงก่อให้เกิดการพัฒนาแบบทดสอบวัดผลสัมฤทธิ์ที่มีประสิทธิภาพขึ้น เพื่อใช้เป็นเครื่องมือในการวัดระดับความสามารถของบุคคล หรือเพื่อจำแนกความสามารถของบุคคลที่แตกต่างกัน ทำให้เราสามารถศึกษาให้สอดคล้องกับระดับความสามารถของบุคคลต่าง ๆ ได้

2. ประเภทของแบบทดสอบวัดผลสัมฤทธิ์ทางการเรียน เราสามารถจำแนกตามมิติ ดังนี้ ได้หลายมิติ ดังนี้

2.1 จำแนกตามขอบข่ายเนื้อหาวิชาของแบบทดสอบวัดผลสัมฤทธิ์ทางการเรียน อาจกำหนดให้กวดง่าย ให้หรือกำหนดให้ยาก ตามปกติแล้วจะไม่มีมาตรฐานอ้างอิงสากลที่จะนำไปใช้ในการกำหนดเนื้อหาวิชาสำหรับแบบทดสอบวัดผลสัมฤทธิ์ ผู้ใช้แบบทดสอบเท่านั้นจะต้องกำหนดเนื้อหาวิชาขึ้นเอง โดยให้สอดคล้องกับวัตถุประสงค์ของการสอน ผู้สร้างแบบทดสอบสามารถที่จะพัฒนาแบบทดสอบให้มีเนื้อหาได้ตามขอบข่ายที่ต้องการ สำหรับแบบทดสอบมาตรฐานวัดผลสัมฤทธิ์ ซึ่งจัดพิมพ์ในต่างประเทศจำนวนมากมักจะรวมแบบทดสอบต่าง ๆ ไว้เป็นชุด (Batteries) แต่ละชุดจะครอบคลุมเนื้อหาวิชาสำหรับระดับชั้นเรียนที่ต่าง ๆ กัน ทั้งนี้ก็เพื่อให้โรงเรียนทั้งหลายสามารถวัดผลสัมฤทธิ์ของนักเรียนได้อย่างต่อเนื่อง ตามความเจริญของงานทางวิชาการในแต่ละยุคแต่ละสมัย ตัวอย่างเช่น แบบทดสอบวัดผลสัมฤทธิ์ทั่วไป (General Achievement Test Batteries) หรือเรียกย่อ ๆ ว่า GATB ซึ่งเป็นแบบทดสอบวัดผลสัมฤทธิ์มาตรฐานที่นิยมใช้กันในระดับประเทศศึกษาและมัธยมศึกษาเป็นส่วนใหญ่ ทั้งนี้เพราะเนื้อหาที่

หมายความกับช่วงชั้นดังกล่าวเป็นอย่างดี เนื่องในระดับประถมศึกษานั้นลักษณะของแบบทดสอบจะประกอบด้วยเนื้อหาซึ่งสัมพันธ์กับทักษะสามประการ (3R) คือ การอ่าน (Reading) การเขียน (Writing) และการคำนวณ (Arithmetic) ส่วนในระดับมัธยมศึกษาแบบทดสอบ GATB มีความมุ่งหมายเพื่อวัดผลสัมฤทธิ์ทั่วไป รวมทั้งเพื่อวัดพื้นฐานความรู้ของแต่ละวิชา เช่น พื้นฐานความรู้ทางสังคมศาสตร์ พื้นฐานความรู้ของวิชาวิทยาศาสตร์ธรรมชาติ เป็นต้น นอกจากนั้นก็จะเน้นการวัดความสามารถด้านการอ่าน ความเข้าใจในสาขาวิชา รวมทั้งความคิดเกี่ยวกับจัํนวนและคำศัพท์ ทั่วไป ซึ่งเป็นการวัดในด้านการพัฒนาความสามารถเชิงวิชาการ ดังนั้นแบบทดสอบ GATB ในระดับนี้ จึงมักจะใช้ในการสอบคัดเลือกนักศึกษาต่อในสถาบันอุดมศึกษา หรือในการสอบคัดเลือกเข้าแข่งขันชิงทุนประเภทต่าง ๆ เป็นต้น

2.2 จำแนกตามลักษณะหน้าที่ทั่วไปของแบบทดสอบ ซึ่งแบ่งออกเป็น 3 ลักษณะดังนี้

2.2.1 แบบทดสอบเพื่อการสำรวจผลสัมฤทธิ์ (Survey Tests) เป็นแบบทดสอบวัดผลสัมฤทธิ์ที่ทำหน้าที่ในการสำรวจความสามารถทั่ว ๆ ไปของนักเรียน โดยประเมินความรู้ในเนื้อหาวิชาหรือทักษะต่าง ๆ เพื่อแสดงระดับความสามารถของนักเรียน ดังนั้นแบบทดสอบวัดผลสัมฤทธิ์จะมักจะครอบคลุมเนื้อหาทั้งในระดับกว้างและระดับทั่วไป โดยถือคะแนนรวมที่ได้จากแบบทดสอบเป็นตัวบ่งชี้ถึงระดับความสามารถทั่วคุ้นได้

2.2.2 แบบทดสอบเพื่อวินิจฉัยผลสัมฤทธิ์ (Diagnostic Test) เป็นแบบทดสอบวัดผลสัมฤทธิ์ที่ทำหน้าที่วินิจฉัยเกี่ยวกับปัจจุบันและปัจจุบốiขององค์ประกอบสำคัญทางด้านทักษะต่าง ๆ ของนักเรียน ตัวอย่างเช่น แบบทดสอบเพื่อวินิจฉัยผลสัมฤทธิ์ทางด้านภาษา อาจจะรวมแบบทดสอบย่อยหลายชุด ซึ่งครอบคลุมเนื้อหาเกี่ยวกับการรู้จักใช้คำ ความเข้าใจเกี่ยวกับถ้อยคำ ต่าง ๆ รวมทั้งคำศัพท์ ตลอดจนอัตราการอ่าน ความเข้าใจในเรื่องที่อ่าน การจำแนกเสียงและการจำแนกพยางค์ ฯลฯ คะแนนที่ได้จากแต่ละองค์ประกอบของแบบทดสอบวินิจฉัยดังกล่าว จะช่วยให้นักจิตวิทยาหรือครูสามารถคัดสินใจได้ว่า อะไรคือจุดضعفของผู้สอบ ซึ่งจะช่วยให้สามารถสอนเสริมในส่วนของเนื้อหาวิชาหรือทักษะที่ยังขาดอยู่ได้อย่างมีประสิทธิภาพ

2.2.3 แบบทดสอบเพื่อวัดความพร้อม (Readiness Test) เป็นแบบทดสอบวัดผลสัมฤทธิ์ ซึ่งทำหน้าที่ในการวัดทักษะที่จำเป็นสำหรับการเรียนในชั้นที่สูงขึ้น แบบทดสอบที่วัดความพร้อม ใช้สำหรับนำways การกระทำในอนาคต ซึ่งทำหน้าที่เป็นเครื่องมือในการวัดความตันด ไปในด้านด้วย

2.3 จำแนกตามคำตอบที่ใช้ โดยทั่วไปแล้วแบบทดสอบวัดผลสัมฤทธิ์ส่วนใหญ่ที่ใช้กันมักจะเป็นประเภทข้อเขียน และที่ใช้กันค่อนข้างมาก คือแบบทดสอบภาคปฏิบัติ (Performance

Test) เป็นแบบทดสอบที่ต้องการให้นักเรียนหรือผู้เข้าสอบได้สังเกตทักษะของเขาวง เช่น ให้แสดงทักษะในการแก้ไขเครื่องยนต์กลไกที่ไม่ทำงาน หรือแสดงทักษะในการเล่นดนตรี เป็นต้น

สำหรับแบบทดสอบประเภทข้อเขียนนั้น ยังแยกออกได้อีก 2 ระดับ คือ (1) ระดับการเลือกคำตอบจากที่กำหนดไว้แล้ว (Recognition) และ (2) ระดับของการเขียนคำตอบจากความรู้หรือความทรงจำที่มีอยู่เดิม (Recall) ในแบบทดสอบระดับที่ 1 แต่ละข้อจะมีคำตอบที่ถูกตัว และจะประกอบด้วยคำลือกหลาย ๆ ตัวที่เป็นไปได้รวมอยู่ในคำตอบที่เกี่ยวข้อง ผู้เข้าสอบจะต้องตัดสินใจเลือกคำตอบอย่างรอบคอบและถูกต้องให้สอดคล้องกับชนิดของคำถามที่ระบุไว้ ตัวอย่างของข้อสอบระดับนี้ เช่น แบบทดสอบหาค่าตัวลือก (Multiple Choice) แบบทดสอบประเภทถูก – ผิด (True - False) และแบบทดสอบประเภทจับคู่ (Matching)

ส่วนแบบทดสอบระดับที่ 2 ซึ่งต้องใช้ความรู้และความทรงจำที่มีอยู่เดิมมาเขียนตอนนั้นลักษณะของคำตอบอาจจะไม่ถูกตัว ขึ้นอยู่กับเหตุผลและความถูกต้องในเชิงวิชาการ ผสมผสานกับความคิดหรือเริ่มสร้างสรรค์ของผู้เข้าสอบเป็นสำคัญ แบบทดสอบระดับนี้ ได้แก่ แบบทดสอบประเภทเติมคำหรือข้อความในช่องว่าง (Completion) แบบทดสอบประเภทตอบสั้น (Short Answer) และแบบทดสอบประเภทความเรียง

3. แบบทดสอบวัดผลสัมฤทธิ์มาตรฐาน

เป็นแบบทดสอบที่สร้างขึ้นโดยกลุ่มผู้เชี่ยวชาญมากกว่าที่จะสร้างขึ้นโดยบุคคลใดบุคคลหนึ่งเพียงบุคคลเดียวเท่านั้น ตามปกติแล้วผู้สร้างแบบทดสอบวัดผลสัมฤทธิ์มาตรฐานมักจะประกอบด้วยผู้เชี่ยวชาญทางด้านการวัดและประเมินผล รวมทั้งผู้เชี่ยวชาญในสาขาวิชานั้น ๆ ตลอดจนครูในโรงเรียนต่าง ๆ ซึ่งมีบทบาทในการกำหนดขอบข่ายเนื้อหาที่ต้องการทดสอบให้เหมาะสม แบบทดสอบวัดผลสัมฤทธิ์มาตรฐาน “ไม่จำเป็นต้องครอบคลุมเนื้อหาและทักษะที่มีในหลักสูตร เนื้อหาและทักษะของแบบทดสอบวัดผลสัมฤทธิ์มาตรฐานส่วนมากมักจะได้จำกัดไว้เรียนและความคิดเห็นของผู้เชี่ยวชาญทางด้านหลักสูตร เนื้อหาโดยทั่วไปจะเป็นความรู้และทักษะในระดับที่กว้าง ๆ เพื่อให้สามารถนำไปใช้กับนักเรียนโรงเรียนต่าง ๆ ได้ สำหรับขั้นตอนในการสร้างแบบทดสอบวัดผลสัมฤทธิ์มาตรฐาน จะต้องมีการวางแผนสร้างอย่างเป็นระบบ คือ มีการระบุหลักเกณฑ์และเหตุผลของการสร้างแบบทดสอบ มีการกำหนดวัตถุประสงค์ของการสร้างที่ชัดเจน มีการทดลองใช้แบบทดสอบที่สร้างขึ้น เพื่อตรวจสอบความเป็นมาตรฐาน โดยการวิเคราะห์ระดับความยากง่าย และอำนาจจำแนกของข้อสอบ มีการหาค่าความตรง และความเที่ยงของแบบทดสอบ พิริมทั้งพื้นฐานตารางปักติวิสัย (Norm Table) เพื่อใช้ในการเปรียบเทียบ มีการกำหนดเวลาของการทดสอบและวิธีดำเนินการสอบ ตลอดจนมีคู่มือประกอบการใช้แบบทดสอบ

ซึ่งจะระบุคุณลักษณะของแบบทดสอบ ประสิทธิภาพของแบบทดสอบรวมทั้งวิธีการใช้แบบทดสอบ และวิธีการตรวจสอบหรือวิธีการให้คะแนน พร้อมทั้งตารางปกติวิสัยของกลุ่ม

4. ความแตกต่างระหว่างแบบทดสอบวัดผลสัมฤทธิ์มาตรฐานกับแบบทดสอบวัดผลสัมฤทธิ์ที่ครูสร้างขึ้น สามารถจำแนกความแตกต่างที่ชัดเจนได้ 5 ประการ ดังนี้

4.1 การจำกัดเนื้อหาวิชาที่สอน แบบทดสอบวัดผลสัมฤทธิ์มาตรฐานจะสุมเนื้อหา สำหรับนำมาสอบในระดับที่กว้างและทั่วไป เพื่อใช้กับโรงเรียนค่าง ๆ คัดเลือกชนมีการกลั่นกรอง เนื้อหาในการสร้างข้อกระทง โดยผู้เชี่ยวชาญทางเนื้อหาและหลักสูตร สำหรับแบบทดสอบ วัดผลสัมฤทธิ์ที่ครูสร้างขึ้น จะเน้นเนื้อหาเฉพาะที่เกี่ยวกับการเรียนการสอนในห้องเรียน ครูจะทำหน้าที่เป็นผู้เชี่ยวชาญ ซึ่งอาจประกอบด้วยครุคนเดียวหรือคณะครุเป็นผู้กำหนดเนื้อหาที่เหมาะสม ในการสอบ

4.2 การทดลองใช้แบบทดสอบ แบบทดสอบวัดผลสัมฤทธิ์มาตรฐานเมื่อสร้างขึ้น แล้วจะต้องมีการทดลองใช้ เพื่อทำการวิเคราะห์ประสิทธิภาพของแบบทดสอบด้วยค่าสถิติต่าง ๆ ต่อจากนั้นก็จะรายงานในคู่มือการใช้แบบทดสอบ เช่น ค่าความตรง ความเที่ยง ระดับความยาก ง่ายและค่าอำนาจจำแนก ในทำนองตรงกันข้าม สำหรับแบบทดสอบวัดผลสัมฤทธิ์ที่ครูสร้างขึ้น โดยทั่วไปจะไม่มีการทดลองใช้เพื่อวิเคราะห์ค่าสถิติในการตรวจสอบประสิทธิภาพของ แบบทดสอบมาก่อน

4.3 วิธีดำเนินการสอบ แบบทดสอบวัดผลสัมฤทธิ์มาตรฐานโดยปกติจะต้องมีคู่มือ อธิบายวิธีดำเนินการสอบอย่างเป็นมาตรฐาน เช่น วิธีการตอบ เวลาในการสอบ ฯลฯ ผู้ใช้แบบทดสอบด้วยปฏิบัติคำสอนอย่างเคร่งครัด สำหรับแบบทดสอบวัดผลสัมฤทธิ์ที่ครูสร้างขึ้น โดยทั่วไปไม่มีคู่มือประกอบการใช้ เพราะตัวครุเองเป็นผู้กำหนดมาตรฐานในการปฏิบัติเกี่ยวกับ วิธีการสอบ

4.4 วิธีการให้คะแนน แบบทดสอบวัดผลสัมฤทธิ์มาตรฐานจะต้องมีคำแปลยสำหรับ การตรวจสอบให้คะแนนตามที่ระบุอยู่ในคู่มือการใช้แบบทดสอบ ส่วนแบบทดสอบวัดผลสัมฤทธิ์ที่ครูสร้างขึ้น ครูจะเป็นผู้ให้คะแนนตามมาตรฐานที่กำหนดขึ้นเอง

4.5 ตารางปกติวิสัยเพื่อการเปรียบเทียบ โดยปกติแล้วการสร้างแบบทดสอบ วัดผลสัมฤทธิ์มาตรฐาน จะมีการนำໄไปใช้กับกลุ่มอ้างอิงหรือที่เรียกว่า (Norm Group) เพื่อทำการ ปกติวิสัย (Norm Table) ไว้ในคู่มือของการใช้แบบทดสอบมาตรฐาน โดยมีจุดมุ่งหมายให้ผู้ใช้แบบทดสอบสามารถนำไปใช้ดีความสำหรับคะแนนสอบที่ได้รับ รวมทั้งใช้เป็นตารางเพื่อ เปรียบเทียบทองคะแนนดังกล่าวด้วย ส่วนแบบทดสอบวัดผลสัมฤทธิ์ที่ครูสร้างขึ้นจะมีเพียง คะแนนของกลุ่มผู้เข้าสอบด้วยกัน ซึ่งอาจใช้เปรียบเทียบได้เฉพาะภายในกลุ่มเท่านั้น

แบบทดสอบวัดผลสัมฤทธิ์ทางการเรียนระดับชาติ ชั้นประถมศึกษาปีที่ 6

ปีการศึกษา 2546 (สำนักทดสอบทางการศึกษา, 2546 อ้างถึงใน ศุภวัฒน์ มะลิเผือก, 2548, หน้า 53-55)

กระทรวงศึกษาธิการเห็นความสำคัญและความจำเป็นในการรักษาและยกระดับคุณภาพการศึกษาของสถานศึกษาต่าง ๆ ให้มีมาตรฐานทัดเทียมกัน จึงกำหนดให้มีการประเมินคุณภาพการศึกษาระดับชาติในทุก ๆ ปี โดยจะประเมินในปลายปีการศึกษา โดยมอบหมายให้สำนักทดสอบทางการศึกษา สำนักงานคณะกรรมการการศึกษาขั้นพื้นฐาน กระทรวงศึกษาธิการดำเนินการจัดทดสอบวัดผลสัมฤทธิ์ทางการเรียนของนักเรียนชั้นประถมศึกษาปีที่ 6 ให้ได้ผลการประเมินที่น่าเชื่อถือและสามารถบ่งชี้คุณภาพการศึกษาระดับชาติ ระดับเขตพื้นที่การศึกษา ระดับสถานศึกษา และระดับผู้เรียนเป็นรายบุคคล ได้อย่างสมเหตุสมผล

วัตถุประสงค์ของการวัดผลสัมฤทธิ์ทางการเรียน

1. เพื่อนำผลไปใช้กำกับ คูแล ติดตามและประเมินการจัดการศึกษาในระดับประเทศ สำหรับพัฒนาคุณภาพทางการศึกษา
2. เพื่อนักเรียนได้นำผลจากการสอบวัดไปพัฒนาความรู้ ความสามารถและการศึกษาต่อของนักเรียนในระดับที่สูงขึ้น

ขอบข่ายของการประเมิน

ดำเนินการประเมินนักเรียนทุกคนในสังกัดคณะกรรมการการศึกษาขั้นพื้นฐาน และสังกัดอื่นที่จัดการศึกษาตามระบบ ตามหลักสูตรประถมศึกษา พุทธศักราช 2521 (ฉบับปรับปรุง 2533) ทุกชั้นหัวด้านประเทศ ได้แก่ สำนักงานคณะกรรมการการศึกษาขั้นพื้นฐาน สำนักงานคณะกรรมการการศึกษาเอกชน สำนักงานคณะกรรมการอุตุกรรมการการอุดมศึกษา สำนักการศึกษากรุงเทพมหานคร สำนักบริหารการศึกษาท้องถิ่น กองบังคับการตำรวจนครบาล สำนักงานตำรวจนครบาล สำนักงานพัฒนาการศึกษาและนักงานการ

ลักษณะทั่วไปของแบบทดสอบวัดผลสัมฤทธิ์ทางการเรียน

เป็นแบบทดสอบที่สร้างและพัฒนาขึ้นจนเป็นแบบทดสอบมาตรฐาน เพื่อใช้ประเมินผลการเรียนรู้ของนักเรียนที่ได้สะสมมาตลอดการเรียน ตามหลักสูตรประถมศึกษา พุทธศักราช 2521 (ฉบับปรับปรุง 2533) ของกระทรวงศึกษาธิการ

โครงสร้างของแบบทดสอบ

1. เนื้อหาสาระของแบบทดสอบวัดผลสัมฤทธิ์ทางการเรียนชั้นประถมศึกษาปีที่ 6 ประกอบด้วย เนื้อหาตามหลักสูตรประถมศึกษา พุทธศักราช 2521 (ฉบับปรับปรุง 2533) วิชาที่สอบในครั้งนี้ ได้แก่ วิชาภาษาไทย คณิตศาสตร์ วิทยาศาสตร์ และภาษาอังกฤษ

2. แบบทดสอบแต่ละวิชานั้นวัดความรู้ความเข้าใจ ความคิดวิเคราะห์ และทักษะเชิงกระบวนการเฉพาะวิชา

การเตรียมการสอน

1. แต่งตั้งคณะกรรมการคุณสอบบ

ให้สำนักงานเขตพื้นที่การศึกษาแต่งตั้งคณะกรรมการคุณสอบลับโรงเรียนภายในเขตพื้นที่การศึกษานั้น ๆ โดยในการคุณสอบแต่ละห้องต้องมีคณะกรรมการคุณสอบ 2 คน

2. การเตรียมตัวนักเรียน

ก่อนดำเนินการสอนให้ผู้บริหารโรงเรียน และครู – อาจารย์ ชี้แจงทำความเข้าใจกับนักเรียนให้ครบถ้วนถึงประโยชน์และความสำคัญของการสอน ตั้งใจตอบแบบทดสอบให้เต็มความสามารถที่แท้จริง ทั้งนี้ข้อมูลที่ได้จากการประเมินจะเกิดประโยชน์ด่อนักเรียนและส่วนรวมดังนี้

2.1 ประโยชน์ด่อนักเรียน นักเรียนจะได้ประโยชน์จากการสอบวัดผลสมฤทธิ์ทางการเรียนของนักเรียนในระดับชาติ ทำให้เกิดความชำนาญในการคิด ได้เห็นรูปแบบของแบบทดสอบมาตรฐาน เป็นการเพิ่มพูนประสบการณ์ในการสอน และคะแนนที่ได้จะเป็นข้อมูลประกอบในการพิจารณาเข้าศึกษาต่อของนักเรียน รวมทั้งได้รับทราบความสามารถของตนว่ามีอุดมคิด ควรปรับปรุงและพัฒนาด้านใด เพราะในการสอบครั้งนี้เป็นการวัดความสามารถหลายด้าน คือ ภาษาไทย คณิตศาสตร์ วิทยาศาสตร์ และภาษาอังกฤษ ผลการสอบต้องมานับทั้งกลุ่มสาระหลักฐานการศึกษา

2.2 ประโยชน์ด้านสถานศึกษา คณะครู – อาจารย์ นำผลการประเมินไปใช้เป็นข้อมูลในการแนะนำทางการศึกษาให้กับผู้เรียน และพัฒนาการจัดการเรียนการสอนให้ได้คุณภาพมาตรฐานยั่งยืน

2.3 ประโยชน์ต่อประเทศ รัฐบาลและหน่วยงานที่เกี่ยวข้องกับการจัดการศึกษานำผลการประเมินไปใช้เป็นข้อมูลในการกำหนดนโยบายทางการศึกษา และส่งเสริมสนับสนุนให้สถานศึกษาและหน่วยงานที่เกี่ยวข้อง พัฒนาคุณภาพการศึกษาให้ได้มาตรฐานทั่วทั้งประเทศ

แนวปฏิบัติกรณีจำเป็นในการดำเนินการสอน

กระทรวงศึกษาธิการต้องการให้นักเรียนชั้นที่ประมิน ได้รับการประเมินทุกคนพร้อมกันทั่วประเทศ เพื่อควบคุมความเป็นมาตรฐานในการดำเนินงาน และผลการประเมินมีความตรงเป็นที่น่าเชื่อถือ แต่ในกรณีที่โรงเรียนไม่สามารถดำเนินการสอนในวันที่กำหนด อาจเนื่องมาจาก

เหตุจำเป็นต้องปฏิบัติ ถ้าไม่ดำเนินการจะทำให้ราชการเสียหาย หรือเหตุสุดวิสัยอื่น โรงเรียนอาจเลื่อนการสอบไปได้ โดยเสนอประธานคณะกรรมการดำเนินงานระดับเขตพื้นที่การศึกษาพิจารณาแล้วเจ้งสำนักงานคณะกรรมการการศึกษาขั้นพื้นฐานทราบ

การดำเนินการสอบให้เป็นไปตามระเบียบการสอบของกระทรวงศึกษาธิการอย่างเคร่งครัด กรณีที่นักเรียนขาดสอบให้โรงเรียนพิจารณาถึงสาเหตุและความจำเป็น แล้วจึงพิจารณาให้นักเรียนเข้าสอบชดเชย โดยเก็บรักษาข้อสอบเป็นความลับ และดำเนินการสอบตามระเบียบของกระทรวงศึกษาธิการ

การดำเนินการหลังการสอบ

การดำเนินการหลังการสอบ แบ่งเป็น 3 ส่วน ดังนี้

1. เมื่อกรรมการดำเนินการสอบเก็บและตรวจสอบความเรียบร้อยของระยะเวลาคำตوبแล้ว ให้นำส่งคณะกรรมการดำเนินการสอบระดับโรงเรียน จากนั้นโรงเรียนรวบรวมคณะกรรมการคำตوبทั้งหมดส่งคณะกรรมการที่สำนักงานเขตพื้นที่การศึกษากำหนดไว้ สำหรับส่วนกลางให้รวมรวมระยะเวลาคำตوبส่งที่สำนักทดสอบทางการศึกษา สำนักงานคณะกรรมการการศึกษาขั้นพื้นฐาน กรุงเทพมหานคร
2. โรงเรียนรวบรวมแบบทดสอบทั้งหมดส่งคืนคณะกรรมการที่สำนักงานเขตพื้นที่การศึกษาแต่ตัว เพื่อดำเนินการย่อย/เพาทำลายต่อไป สำหรับส่วนกลาง นำส่งที่โรงพิมพ์ครุสภากาดพร้าว
3. คณะกรรมการดำเนินงานระดับเขตพื้นที่การศึกษา รวบรวมระยะเวลาคำตوبทั้งหมด ส่งสำนักทดสอบทางการศึกษา สำนักงานคณะกรรมการการศึกษาขั้นพื้นฐาน กระทรวงศึกษาธิการ ทันที เนื่องจากสำนักทดสอบทางการศึกษาจะต้องตรวจสอบให้ทั้งหมด ซึ่งมีปริมาณมากให้เสร็จสิ้นโดยเร็ว เพื่อให้โรงเรียนได้ทำบันทึกลงในเอกสารหลักฐานการศึกษาได้ทันก่อนปิดภาคเรียน

ตอนที่ 5 ค่าความเที่ยงของแบบทดสอบ

ความเที่ยงของแบบทดสอบ หมายถึง ความคงที่ของคะแนนที่ได้จากการสอบนักเรียน คนเดียวกันหลายครั้ง ในแบบทดสอบชุดเดิม ถ้าคะแนนเปลี่ยนไปจากเดิม แสดงว่าแบบทดสอบขาดความเที่ยง ทำให้ผลการสอบมีความคลาดเคลื่อนไปจากความรู้จริงของนักเรียน ความคลาดเคลื่อนชนิดนี้ เรียกว่า ความคลาดเคลื่อนในการวัด (Error of Measurement) ในการวัดผลนั้น จะต้องสร้างเครื่องมือที่ต้องการนำไปวัดผลให้มีคุณภาพที่เชื่อมั่นได้ เพื่อผลการวัดที่ออกมายังไห เป็นคะแนนความรู้จริงของนักเรียนที่ปราศจากความคลาดเคลื่อนในการวัด และค่าเที่ยงจะมีค่าอยู่

ระหว่าง -1 ถึง +1 และจะพิจารณาเฉพาะค่าที่เป็นบวกเท่านั้น ซึ่งกรณีค่ามากกว่า 0.70 จึงจะเป็นแบบทดสอบที่มีความเชื่อถือได้

การคำนวณหาค่าความเที่ยงนั้นจะคำนวณในรูปของการประมาณค่า ในรูปของค่าสัมประสิทธิ์ของความเที่ยง มักใช้สัญลักษณ์ว่า r_u , r_{xx} หรือ r_{cc} โดยทั่วไปแล้วความเที่ยงแบ่งออกเป็น 2 ประเภท คือ ความเที่ยงของแบบทดสอบอิงกลุ่ม (Reliability of Norm-Referenced Test) กับความเที่ยงของแบบทดสอบอิงเกณฑ์ (Reliability of Criterion-Referenced Test) และในเดลต้าประเภทของความเที่ยงน้ำใจความเที่ยงนั้น จะมีสูตรที่ใช้ประมาณค่าความเที่ยงอยู่หลายสูตร ซึ่งวิธีคำนวณจำแนกตามลักษณะของแบบทดสอบอิงกลุ่มและอิงเกณฑ์ สำหรับค่าความเที่ยงที่เกี่ยวข้องกับงานวิจัยนี้ เป็นค่าความเที่ยงที่คำนวณโดยใช้ความสอดคล้องภายใน (Internal Consistency Reliability) เป็นการหาความเที่ยงที่ใช้แบบทดสอบฉบับเดียวทำการทดสอบเพียงครั้งเดียว วิธีหาค่าความเที่ยงที่เกี่ยวข้องกับงานวิจัยนี้ได้แก่ วิธีของ cronbach (Cronbach Alpha Procedure)

cronbach ได้พัฒนาสูตรหาความเที่ยงในรูปสัมประสิทธิ์แอลฟ่า (α -Coefficient) ในปี ค.ศ. 1951 โดยพัฒนามาจากสูตร KR-20 ทั้งนี้เป็น เพราะว่าจะได้ใช้หาความเที่ยงกับเครื่องมือที่ไม่ได้ตรวจสอบให้คะแนนเป็น 1 กับ 0 จะตรวจให้คะแนนลักษณะใดก็ได้ เช่น ถ้าทำถูกได้คะแนนเป็น 10, 8 หรือในลักษณะแบบสอบถามที่ให้คะแนนแต่ละข้อเป็น 3, 2, 1 หรือ 5, 4, 3, 2, 1 ก็ได้ สูตรที่ใช้คือ

$$\alpha = \frac{k}{k-1} \left\{ 1 - \frac{\sum \sigma_i^2}{\sigma^2} \right\}$$

เมื่อ	α	คือ สัมประสิทธิ์ความเที่ยง
	k	คือ จำนวนข้อสอบ
	σ_i^2	คือ คะแนนความแปรปรวนเป็นรายข้อ
	σ^2	คือ คะแนนความแปรปรวนของทั้งฉบับ

โดยที่

$$\sigma_i^2 = \frac{N \sum X_i^2 - (\sum X_i)^2}{N^2}$$

เมื่อ	$\sum X_i^2$	คือ ผลรวมกำลังสองของคะแนนในข้อที่ i
	$(\sum X_i)^2$	คือ ผลรวมทั้งหมดของคะแนนในข้อที่ i ยกกำลังสอง

N คือ จำนวนคนเข้าสอบ

แล้ว

$$\sigma^2 = \frac{N \sum X^2 - (\sum X)^2}{N^2}$$

เมื่อ $\sum X^2$ คือ ผลรวมกำลังสองของคะแนนของข้อสอบทั้งฉบับ
 $(\sum X)^2$ คือ ผลรวมทั้งหมดของคะแนนของข้อสอบทั้งฉบับยกกำลังสอง