



ระเบียบวิธีวิจัย

การวิจัยเรื่องนี้มีมุ่งที่จะศึกษาหาสาเหตุหรือปัจจัยที่มีผลต่อสัมฤทธิ์ผลทางการศึกษาในระดับอุดมศึกษาของสถาบันอุดมศึกษาเอกชน โดยเก็บรวบรวมข้อมูลจากนักศึกษาที่เข้าศึกษาตั้งแต่ปี 2523 ของวิทยาลัยเอกชน 7 แห่ง คือ วิทยาลัยกรุงเทพ วิทยาลัยหอการค้า วิทยาลัยธุรกิจบัณฑิต วิทยาลัยเทคนิคสยาม วิทยาลัยเอเชียอาคเนย์ วิทยาลัยเกริก วิทยาลัยอัลสลัมชนูปบริหารธุรกิจ ซึ่งข้อมูลที่ได้นี้แบ่งออกเป็น 2 ส่วนคือ ข้อมูลเกี่ยวกับสถานการณ์ส่วนตัว ซึ่งได้แก่ เพศ สำขาริชาที่เรียน สำขาที่สำเร็จในระดับมัธยมศึกษาตอนปลาย สำขาที่สำเร็จในระดับอุดมศึกษา เกรดเฉลี่ยสะสมที่ได้จากมัธยมศึกษาตอนปลาย เกรดเฉลี่ยสะสมที่ได้ในระดับอุดมศึกษา อีกส่วนหนึ่งของข้อมูลคือ ข้อมูลที่เกี่ยวกับปัญหาส่วนตัวหรือสภาพแวดล้อมของนักศึกษาซึ่งได้แก่ การแบ่งเวลาในการเรียน ความสนใจในการเรียน ทัศนคติต่อวิชาและคณะที่เรียน ความรู้พื้นฐานและสติปัญญา การร่วมกิจกรรมนักศึกษา การปรับตัวในการเรียนด้านเศรษฐกิจ ด้านครอบครัว ด้านสุขภาพและทัศนคติเกี่ยวกับครูผู้สอน

2.1 ประชากรที่ศึกษา

ประชากรที่ใช้ในการศึกษาคือ นักศึกษาที่เข้าศึกษาในวิทยาลัยเอกชนทั้ง 7 แห่งดังกล่าวที่เข้าศึกษาในปี พ.ศ. 2523 ซึ่งแบ่งเป็น 2 กลุ่มคือ กลุ่มที่สำเร็จตามหลักสูตรและกลุ่มที่ไม่สำเร็จการศึกษารายละเอียดเกี่ยวกับจำนวนนักศึกษาทั้งหมดจำแนกตามสำขาริชาและสถาบันการศึกษาแสดงไว้ในตารางที่ 1

ตารางที่ 1 จำนวนนักศึกษาจำแนกตามสาขาวิชาการศึกษาและสถาบันการศึกษาได้ดังนี้

วิทยาลัยที่ศึกษา	วิทยาลัย กรุงเทพ		วิทยาลัย หอการค้าไทย		วิทยาลัย ธุรกิจบัณฑิต		วิทยาลัย เกริก		วิทยาลัย อัสสัมชัญ		วิทยาลัย เอเชียอาคเนย์		วิทยาลัย เทคนิคสยาม		รวมทุกวิทยาลัย	
	สำเร็จ	ไม่สำเร็จ	สำเร็จ	ไม่สำเร็จ	สำเร็จ	ไม่สำเร็จ	สำเร็จ	ไม่สำเร็จ	สำเร็จ	ไม่สำเร็จ	สำเร็จ	ไม่สำเร็จ	สำเร็จ	ไม่สำเร็จ	สำเร็จ	ไม่สำเร็จ
สาขาบริหารธุรกิจ	217	173	360	304	166	151	134	68	353	172	148	53	156	58	1534	979
สาขาบัญชี	82	162	152	174	126	94	148	43	360	146	86	42	32	27	987	688
สาขานิเทศศาสตร์	134	168	-	-	-	-	-	-	-	-	-	-	-	-	134	168
รวมทุกสาขาวิชา	433	503	513	478	292	245	282	111	713	318	234	95	188	85	2655	1835

ศูนย์วิทยทรัพยากร
จุฬาลงกรณ์มหาวิทยาลัย

2.2 การเก็บรวบรวมข้อมูล

การเก็บรวบรวมข้อมูลมาใช้ในการวิจัยครั้งนี้ใช้วิธีการสำรวจจากตัวอย่าง ซึ่งมีขั้นตอนที่สำคัญในการรวบรวมข้อมูลดังนี้

1. จากจำนวนนักศึกษา 4,490 คน ซึ่งเข้าศึกษาในปีการศึกษา 2523 ของสถาบันการศึกษาเอกชนทั้ง 7 แห่ง ซึ่งประกอบด้วยนักศึกษาที่สำเร็จการศึกษา 2655 คน และไม่สำเร็จการศึกษา 1835 คน ผู้วิจัยได้คำนวณขนาดตัวอย่างที่จะต้องใช้ในการวิจัยโดยใช้สูตรดังนี้

$$n = \frac{z_{\frac{\alpha}{2}}^2 s^2}{d}$$

เมื่อ $z_{\frac{\alpha}{2}}$ คือ ค่าจากตารางการแจกแจงแบบปกติมาตรฐานที่ระดับนัยสำคัญ α

s คือ ส่วนเบี่ยงเบนมาตรฐานของสัมฤทธิ์ผลทางการศึกษาเมื่อสำเร็จการศึกษาหรือพ้นสภาพนักศึกษา

d คือ ความผิดพลาดของสัมฤทธิ์ผลทางการศึกษาที่ยอมรับให้เกิดขึ้นได้

เนื่องจากการวิจัยครั้งนี้มีตัวแปรทั้งหมด 62 ตัว ดังนั้นขนาดตัวอย่างที่ใช้จึงอาจมีได้ทั้งหมดถึง 62 ค่า ซึ่งแตกต่างกัน แต่ผู้วิจัยจะใช้ขนาดตัวอย่างที่คำนวณจากตัวแปรที่มีอิทธิพลต่อสัมฤทธิ์ผลทางการศึกษามากที่สุด จากการศึกษาว่าตัวแปรใดที่ควรมีอิทธิพลกับสัมฤทธิ์ผลทางการศึกษาในระดับอุดมศึกษามากที่สุด ปรากฏว่าตัวแปรที่น่าจะมีอิทธิพลกับสัมฤทธิ์ผลทางการศึกษาในระดับอุดมศึกษามากที่สุดคือเกรดเฉลี่ยสะสมระดับมัธยมศึกษาตอนปลาย ผู้วิจัยยังศึกษาต่ออีกว่าความแปรปรวนของเกรดเฉลี่ยสะสมเมื่อแยกตามสาขาและแยกตามวิทยาลัยของนักศึกษาที่สำเร็จการศึกษากับนักศึกษาที่ไม่สำเร็จการศึกษา จะมีความแปรปรวนมากกว่าหรือน้อยกว่าความแปรปรวนของเกรดเฉลี่ยสะสมที่รวมทุกสาขาและทุกวิทยาลัยของนักศึกษาที่สำเร็จการศึกษากับนักศึกษาที่ไม่สำเร็จการศึกษา ปรากฏว่า ความแปรปรวนของเกรดเฉลี่ยสะสมเมื่อคิดรวมทุกสาขาและทุกวิทยาลัยมีมากกว่าความแปรปรวนของเกรดเฉลี่ยสะสมเมื่อแยกตามสาขาและแยกตามวิทยาลัย ดังนั้นผู้วิจัยจึงคำนวณหาขนาดตัวอย่างของนักศึกษาที่สำเร็จการศึกษาและไม่สำเร็จการศึกษาโดยใช้ความแปรปรวนของเกรดเฉลี่ยสะสมเมื่อรวมทุกสาขาและทุกวิทยาลัย ซึ่งเท่ากับ 0.549 และ 0.507 ตามลำดับ จากสูตรการหาขนาดตัวอย่างข้างต้น ถ้ากำหนดให้ระดับนัยสำคัญเท่ากับ 0.05 และความผิดพลาดของสัมฤทธิ์ผลทางการศึกษาไม่เกินร้อยละ 5 จะได้ขนาดตัวอย่างของนักศึกษาที่

สำเร็จการศึกษา

$$n = \left(\frac{1.96 \times 0.549^2}{0.05} \right) = 463 \text{ คน}$$

หรือคิดเป็นร้อยละ $\frac{463 \times 100}{2655} = 17$ ของนักศึกษาซึ่งเข้าศึกษาในปีการศึกษา 2523

ที่สำเร็จการศึกษา

สำหรับขนาดตัวอย่างของนักศึกษาที่ไม่สำเร็จการศึกษาสามารถหาได้โดยวิธีเดียวกัน

กล่าวคือ

$$n = \left(\frac{1.96 \times 0.507^2}{0.05} \right) = 395 \text{ คน}$$

หรือคิดเป็นร้อยละ $\frac{395 \times 100}{1835} = 21.5$ ของนักศึกษาซึ่งเข้าศึกษาในปีการศึกษา 2523

ที่ไม่สำเร็จการศึกษา

เพื่อให้ผลการวิจัยมีความถูกต้องเชื่อถือได้มากขึ้น ผู้วิจัยได้ลุ่มนักศึกษาที่สำเร็จการศึกษาแต่ละสาขาวิชามาร้อยละ 20 คิดเป็นจำนวนนักศึกษาทั้งหมด 530 คน ลุ่มนักศึกษาที่ไม่สำเร็จการศึกษาแต่ละสาขาของทุกวิทยาลัยมาร้อยละ 25 คิดเป็นจำนวนนักศึกษาทั้งหมด 450 คน

สำหรับสาขานิติศาสตร์ ซึ่งเปิดสอนในมหาวิทยาลัยกรุงเทพเพียงแห่งเดียวและมีผู้ไม่สำเร็จการศึกษามากกว่าสาขาอื่น ๆ ได้คำนวณหาขนาดตัวอย่างตามวิธีข้างต้น ได้ตัวอย่างนักศึกษาที่ไม่สำเร็จการศึกษา 100 คน หรือประมาณร้อยละ 60 ของจำนวนนักศึกษาที่ไม่สำเร็จการศึกษาทั้งหมดและได้นักศึกษาที่สำเร็จการศึกษา 65 คน หรือประมาณร้อยละ 40 ของจำนวนนักศึกษาที่สำเร็จการศึกษาทั้งหมด

2. การลุ่มตัวอย่างใช้วิธีการลุ่มตัวอย่างแบบมีระบบโดยแบ่งชั้นภูมิ (Stratified Systematic Random Sampling) กล่าวคือ แบ่งนักศึกษาที่เข้าศึกษาในปีการศึกษา 2523 ออกเป็น 2 กลุ่มได้แก่ กลุ่มนักศึกษาที่สำเร็จการศึกษาและไม่สำเร็จการศึกษาแล้ว ในแต่ละกลุ่มของนักศึกษาที่สำเร็จและไม่สำเร็จการศึกษาแบ่งนักศึกษาออกตามวิทยาลัยและแต่ละวิทยาลัยแบ่งนักศึกษาออกตามสาขาวิชาที่ศึกษาจากนั้นจึงลุ่มตัวอย่างนักศึกษาจากแต่ละสาขาวิชาของวิทยาลัยต่าง ๆ ให้เป็นสัดส่วนกับจำนวนนักศึกษาทั้งหมดที่มีอยู่ในสาขาวิชาของคนนั้น ๆ ในอัตราร้อยละ 20 สำหรับผู้ที่สำเร็จการศึกษาและร้อยละ 25 สำหรับผู้ที่ไม่สำเร็จการศึกษา

3. ส่งแบบสอบถามที่เกี่ยวกับความคิดเห็นของนักศึกษาทางด้านต่าง ๆ ที่อาจมีผลกระทบต่อผลการศึกษาในระดับอุดมศึกษา ไปยังนักศึกษาตัวอย่างทางไปรษณีย์ เมื่อวันที่ 1 ธันวาคม 2527 โดยกำหนดให้นักศึกษาตัวอย่างส่งแบบสอบถามคืนภายในวันที่ 15 มกราคม 2528 สำหรับแบบสอบถามที่ไม่ได้รับคืนภายในกำหนด ผู้วิจัยได้เลือกตัวอย่างนักศึกษาเพิ่มเติมโดยให้นักศึกษา

ที่เลือกเพิ่มเติมนี้มีผลสัมฤทธิ์ทางการเรียนเมื่อสำเร็จการศึกษาหรือพ้นสภาพการศึกษาใกล้เคียงกับนักศึกษาที่ไม่ตอบแบบสอบถามภายในกำหนดมากที่สุด โดยเริ่มส่งแบบสอบถามรอบที่ส่งเมื่อวันที่ 22 มกราคม 2528 และกำหนดให้ส่งคืนภายในวันที่ 28 กุมภาพันธ์ 2528 สำหรับแบบสอบถามที่ไม่ได้รับกลับคืนมาในรอบส่งนี้ ก็ทำเช่นเดียวกับรอบที่ส่งคือ เลือกตัวอย่างเพิ่มเติมแล้วส่งแบบสอบถามออกไปอีกในรอบที่สามระหว่างวันที่ 7 มีนาคม และกำหนดส่งคืนภายในวันที่ 15 เมษายน 2528 สำหรับแบบสอบถามที่ไม่ได้รับกลับคืนมาในรอบที่สามนี้ก็ทำเช่นเดียวกับรอบที่ส่งคือเลือกตัวอย่างเพิ่มเติมแล้วส่งแบบสอบถามออกไปอีกในรอบที่สี่ระหว่างวันที่ 20 เมษายนและกำหนดให้ส่งคืนภายในวันที่ 31 พฤษภาคม 2528 จำนวนตัวอย่างทั้งหมดที่เก็บรวบรวมได้จากการส่งแบบสอบถามทั้ง 4 รอบเท่ากับ 703 ราย คิดเป็นร้อยละ 71.73 ของจำนวนตัวอย่างทั้งหมด ในจำนวนนี้เป็นแบบสอบถามจากกลุ่มที่สำเร็จการศึกษา 350 ราย และกลุ่มที่ไม่สำเร็จการศึกษา 353 ราย รายละเอียดเกี่ยวกับจำนวนตัวอย่างนักศึกษาคำนวณตามสาขาวิชาและสถาบันการศึกษาแสดงไว้ในตารางที่ 2



ตารางที่ 2 จำนวนตัวอย่างนักศึกษาจำแนกตามสาขาวิชาการศึกษาและสถาบันการศึกษาได้ดังนี้

วิทยาลัยที่ศึกษา	วิทยาลัย กรุงเทพ		วิทยาลัย หอการค้าไทย		วิทยาลัย ธุรกิจบัณฑิตย์		วิทยาลัย เกริก		วิทยาลัย อีสต์สัสม์ซิว		วิทยาลัย เอเชียอาคเนย์		วิทยาลัย เทคนิคสยาม		รวมทุกวิทยาลัย	
	สำเร็จ	ไม่สำเร็จ	สำเร็จ	ไม่สำเร็จ	สำเร็จ	ไม่สำเร็จ	สำเร็จ	ไม่สำเร็จ	สำเร็จ	ไม่สำเร็จ	สำเร็จ	ไม่สำเร็จ	สำเร็จ	ไม่สำเร็จ	สำเร็จ	ไม่สำเร็จ
บริหารธุรกิจ	26	34	31	31	24	29	16	23	43	17	33	23	40	19	213	176
บัญชี	17	23	14	18	18	14	13	14	14	14	7	9	3	2	86	94
นิเทศศาสตร์	51	83	-	-	-	-	-	-	-	-	-	-	-	-	51	83
รวมทุกสาขาวิชา	94	140	45	49	42	43	29	37	57	31	40	32	43	21	350	353

จุฬาลงกรณ์มหาวิทยาลัย

2.3 การวิเคราะห์ข้อมูล

ในการวิเคราะห์ข้อมูล ผู้วิจัยได้นำข้อมูลที่ได้รับรวบรวมได้มาทำการวิเคราะห์ตามระเบียบวิธีสถิติโดยใช้โปรแกรมสำเร็จรูป SPSS (Statistical Package for the Social Science) และแยกวิเคราะห์ดังนี้

2.3.1 การวิเคราะห์เกี่ยวกับตัวแปรด้านสถานภาพส่วนตัวของนักศึกษา

ศึกษาว่าการสำเร็จการศึกษาขึ้นอยู่กับเพศ ลำหยาที่เรียน ผลสัมฤทธิ์ทางการเรียน ระดับมัธยมปลาย ผลสัมฤทธิ์ทางการเรียนระดับอุดมศึกษา ลำดับที่การสอบคัดเลือกเข้าในคณะวิชานั้น ๆ วุฒิส่งสุดเดิมก่อนเข้าศึกษาที่วิทยาลัย การศึกษาสูงสุดของบิดามารดา รายได้ของบิดา-มารดาหรือไม่มี และศึกษาว่าตัวแปรใดใน 61 ตัวที่มีอิทธิพลกับเกรดเฉลี่ยสะสมระดับอุดมศึกษา (X_4) ซึ่งมีตัวแปร 61 ตัว

ดังกล่าวคือ

$$X_1 = \text{เพศ}$$

$$X_2 = \text{วิทยาลัยที่ศึกษา}$$

$$X_3 = \text{ลำหยาวิชา}$$

$$X_4 = \text{ระดับคะแนนเฉลี่ยสะสมในระดับอุดมศึกษา}$$

$$X_5 = \text{ลำดับที่การสอบคัดเลือกเข้าในคณะวิชานั้น ๆ}$$

$$X_6 = \text{วุฒิส่งสุดเดิมก่อนเข้าศึกษาในอุดมศึกษา}$$

$$X_7 = \text{ระดับคะแนนเฉลี่ยสะสมในระดับมัธยมปลาย}$$

$$X_8 = \text{ลำเหตุที่ออกกลางคัน}$$

$$X_9 = \text{ผู้อุปการะการศึกษา}$$

$$X_{10} = \text{พักอาศัยอยู่กับใคร}$$

$$X_{11} = \text{สภาพความเป็นอยู่ของบิดามารดา}$$

$$X_{12} = \text{ผู้ที่ เป็นหัวหน้าครอบครัวที่อาศัยอยู่}$$

$$X_{13} = \text{อาชีพบิดามารดาหรือผู้อุปการะ}$$

$$X_{14} = \text{รายได้ของบิดาและมารดาหรือผู้อุปการะ (บาท/เดือน)}$$

$$X_{15} = \text{การศึกษาสูงสุดของบิดามารดาหรือผู้อุปการะ}$$

$$X_{16} = \text{การเข้าชั้นเรียนลำย}$$

$$X_{17} = \text{ออกจากชั้นเรียนก่อนเวลาเลิก}$$

- X_{18} = ขาดเรียนบ่อย ๆ
 X_{19} = ขาดเรียนบางวิชาเพื่อดูหนังสือสอบวิชาต่อไป
 X_{20} = ไม่สนใจเนื้อหาที่เรียน
 X_{21} = ไม่เตรียมตัวก่อนเข้าชั้นเรียน
 X_{22} = ไม่ทบทวนหลังการเรียน
 X_{23} = ไม่ส่งงานตามเวลาผู้สอนกำหนด
 X_{24} = ไม่ศึกษาค้นคว้าจากหนังสืออื่น ๆ ที่เกี่ยวข้อง
 X_{25} = ไม่เตรียมตัวในการสอบเท่าที่ควร
 X_{26} = มีความเบื่อหน่ายต่อการเรียนทุกวิชา
 X_{27} = ชอบหลีกเลี่ยงการทำงานที่ได้รับมอบหมาย
 X_{28} = รู้สึกว่าจำเป็นต้องเรียนบางวิชาที่ไม่ชอบเลย
 X_{29} = วิชาบังคับไม่น่าสนใจเท่าที่ควร
 X_{30} = วิชาเลือกไม่น่าสนใจเท่าที่ควร
 X_{31} = คิดว่าความรู้ที่คณะวิชาให้มันเป็นเรื่องแคบ ๆ
 X_{32} = คิดว่าการศึกษาในคณะวิชานี้ให้ประโยชน์น้อย
 X_{33} = รู้สึกว่าคณะวิชาที่ศึกษาด้อยกว่าคณะวิชาอื่น ๆ
 X_{34} = รู้สึกว่าวิชาที่ศึกษายากเกินไปสำหรับความรู้ของท่าน
 X_{35} = ข้อสอบยาก
 X_{36} = วุฒิส่งสุดท้ายที่จบมา (ก่อนเข้ามาศึกษาในวิทยาลัย) ไม่เอื้ออำนวยสำหรับการ
 ศึกษาในคณะวิชานี้
 X_{37} = ความรู้ความถนัดของท่านไม่ตรงกับคณะวิชาที่เรียน
 X_{38} = ในขณะที่ศึกษาอยู่ท่านให้ความสนใจกับการทำกิจกรรมมากกว่าการเรียน
 X_{39} = กิจกรรมเป็นสาเหตุสำคัญอันหนึ่งสำหรับท่านที่ทำให้ผลการเรียนไม่ดีเท่าที่ควร
 X_{40} = ไม่สามารถแบ่งเวลาเรียนกับเวลาทำกิจกรรมแยกออกจากกัน
 X_{41} = ท่านไม่ชอบทำกิจกรรม
 X_{42} = ในขณะที่ศึกษาอยู่ท่านมักจดคำบรรยายไม่ทัน
 X_{43} = ท่านไม่กล้าถามอาจารย์ผู้สอนเมื่อมีปัญหา

- X_{44} = ไม่ค่อยเข้าใจที่อาจารย์บรรยาย
 X_{45} = ไม่กล้าปรึกษาที่อาจารย์ที่ปรึกษาเมื่อมีปัญหา
 X_{46} = มีปัญหาเกี่ยวกับการทำแบบฝึกหัด
 X_{47} = ค่าใช้จ่ายในการเรียนคณะวิชาณี
 X_{48} = ต้องหาค่าใช้จ่ายในการเรียนด้วยตนเอง
 X_{49} = ค่าใช้จ่ายสำหรับการศึกษาไม่เพียงพอ
 X_{50} = ครอบครัวมีรายได้น้อยจึงไม่สามารถส่งเรียนจนจบหลักสูตรได้
 X_{51} = ปัญหาในครอบครัวของท่านมีส่วนกระทบต่อผลการเรียน
 X_{52} = บิดามารดาหรือผู้ปกครองไม่สนับสนุนการศึกษา
 X_{53} = ต้องช่วยทางบ้านในการประกอบอาชีพ
 X_{54} = มีภาระในการเลี้ยงครอบครัว
 X_{55} = เกิดปัญหาส่วนตัวซึ่งมีผลกระทบต่อ การเรียน
 X_{56} = มีปัญหาสุขภาพทางกายหรือทางจิตที่เป็นอุปสรรคต่อการเรียน
 X_{57} = ไม่สบายหนักจนทำให้ผลการเรียนต่ำลงมากเกินไป
 X_{58} = การคัดครู เข้าชั้นเรียนไม่เหมาะสมกับความรู้อและ ความถนัดในการสอน
 X_{59} = ผู้สอนมีประสบการณ์ในการสอนน้อย
 X_{60} = ผู้สอนไม่พยายามทำความเข้าใจปัญหาหรือความรู้สึกของนักศึกษา
 X_{61} = ผู้สอนมีทัศนคติไม่ดีต่อผู้ที่เรียนอ่อน
 X_{62} = ผู้สอนไม่มีเวลาที่จะให้นักศึกษา เข้าพบเพื่อซักถามปัญหา

การทดสอบดังกล่าวข้างต้นนี้จะใช้การทดสอบไคสแควร์ ซึ่งมีสมมติฐานของการทดสอบระหว่างเกรดเฉลี่ยสะสมระดับอุดมศึกษา กับปัจจัยตัวแปรที่ i ดังนี้

H_0 : เกรดเฉลี่ยสะสมระดับอุดมศึกษา เป็นอิสระกันกับปัจจัยหรือตัวแปรที่ i .

H_a : เกรดเฉลี่ยสะสมระดับอุดมศึกษา ไม่เป็นอิสระกันกับปัจจัยสำหรับตัวแปรที่ i

ตัวสถิติที่ใช้ในการทดสอบคือ

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

- เมื่อ $\chi^2 =$ ค่าสถิติไคสแควร์ ที่ขึ้นความเป็นอิสระ
 $r =$ จำนวนแถว
 $c =$ จำนวนลัดมภ์
 $O_{ij} =$ ค่าความถี่ที่สังเกตได้ ในแถวที่ i และลัดมภ์ที่
 $E_{ij} =$ ค่าความถี่ที่คาดหวัง ในแถวที่ j และลัดมภ์ที่

นำค่า χ^2 ที่คำนวณได้มาเปรียบเทียบกับค่า $\chi^2_{((r-1)(c-1), \alpha)}$ เมื่อ r และ c แทนจำนวนแถวและจำนวนลัดมภ์ตามลำดับ ถ้า χ^2 คำนวณมากกว่า $\chi^2_{((r-1)(c-1), \alpha)}$ จะปฏิเสธ H_0 หรือยอมรับ H_a สำหรับการเปรียบเทียบว่าปัจจัยใดมีอิทธิพลต่อการสำเร็จการศึกษามากกว่ากัน ใช้ตัวสถิติคราแมร์วี ซึ่งมีสูตรในการคำนวณต่อไปนี้

$$V^2 = \frac{\chi^2}{n \cdot \text{Min}(r-1, c-1)}$$

เมื่อ χ^2 คือ ค่าสถิติที่ได้จากการทดสอบความสัมพันธ์ระหว่างปัจจัยต่าง ๆ แต่ละปัจจัยกับสภาวะการสำเร็จการศึกษา

n คือ จำนวนนักศึกษาตัวอย่างที่นำมาทดสอบ

$\text{Min}(r-1, c-1)$ คือ ค่าที่น้อยระหว่าง $r-1$ และ $c-1$ โดยที่ r และ c คือจำนวนแถวและจำนวนลัดมภ์ของตารางการจรณ์ (Contingency table) ตามลำดับ

2.3.2 การวิเคราะห์ตัวแปรเกี่ยวกับปัญหาส่วนตัว

เปรียบเทียบความแตกต่างระหว่างกลุ่มที่สำเร็จการศึกษาและกลุ่มที่ไม่สำเร็จการศึกษาจำแนกตามหมวดปัญหา โดยใช้การทดสอบแบบที (t-test) ซึ่งมีสมมติฐานในการทดสอบสำหรับหมวดปัญหาที่ i ดังนี้

H_0 : คะแนนเฉลี่ยของหมวดปัญหาที่ i ของกลุ่มนักศึกษาที่สำเร็จและไม่สำเร็จ ไม่แตกต่างกัน

H_a : คะแนนเฉลี่ยของหมวดปัญหาที่ i ของกลุ่มนักศึกษาที่สำเร็จน้อยกว่ากลุ่มนักศึกษาที่ไม่สำเร็จ

ตัวสถิติที่ใช้ในการทดสอบคือ $t = \frac{(\bar{X}_{1i} - \bar{X}_{2i})}{\sqrt{\frac{s_{1i}^2}{n_{1i}} + \frac{s_{2i}^2}{n_{2i}}}}$

ผู้วิจัยเลือกใช้สูตรนี้เนื่องจากไม่ทราบว่าคุณค่าความแปรปรวนของคะแนนที่แท้จริงของทั้งสองกลุ่มเท่ากันหรือไม่

เมื่อ μ_{1i} และ μ_{2i} เป็นคะแนนเฉลี่ยจากประชากรของหมวดปัญหาที่ i ของนักศึกษา
กลุ่มที่สำเร็จและไม่สำเร็จตามลำดับ

\bar{X}_{1i} และ \bar{X}_{2i} เป็นคะแนนเฉลี่ยจากตัวอย่างของหมวดปัญหาที่ i ของนักศึกษา
กลุ่มที่สำเร็จและไม่สำเร็จตามลำดับ

S_{1i}^2 และ S_{2i}^2 เป็นค่าความแปรปรวนจากตัวอย่างของคะแนนของหมวดปัญหาที่
 i ของนักศึกษากลุ่มที่สำเร็จและไม่สำเร็จตามลำดับ

n_{1i} และ n_{2i} เป็นขนาดของนักศึกษาตัวอย่างของหมวดปัญหาที่ i ของนักศึกษา
กลุ่มที่สำเร็จและไม่สำเร็จตามลำดับ

นำค่า t ที่คำนวณได้มาเปรียบเทียบกับค่า t ที่เปิดจากตารางด้วยระดับนัยสำคัญที่ต้องการ
ซึ่งมีขึ้นของความเป็นอิสระเท่ากับ

$$\frac{\left[\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2} \right]^2}{\left[\frac{(S_1^2/n_1)^2}{n_1 - 1} + \frac{(S_2^2/n_2)^2}{n_2 - 1} \right]}$$

ถ้าค่า t ที่คำนวณได้มากกว่าค่า t ที่เปิดจากตาราง ที่ขึ้นค่าความเป็นอิสระข้างต้นแล้วระดับ
นัยสำคัญที่ต้องการจะปฏิเสธ H_0 หรือยอมรับ H_a

2.3.3 การคำนวณหาค่าสัมประสิทธิ์สหสัมพันธ์ของปัจจัย

ระหว่างผลสัมฤทธิ์ทางการเรียนระดับมัธยมศึกษาตอนปลายและระดับปริญญาตรี

จำแนกตามสาขาวิชา จำแนกตามกลุ่มการสำเร็จการศึกษา และรวมทุก ๆ สาขาวิชา คำนวณ
โดยใช้สูตร

$$r = \frac{n \sum X_i Y_i - \sum X_i \sum Y_i}{\sqrt{[n \sum X_i^2 - (\sum X_i)^2] [n \sum Y_i^2 - (\sum Y_i)^2]}}$$

เมื่อ r = ค่าสัมประสิทธิ์สหสัมพันธ์แบบเพียร์สัน

n = จำนวนนักศึกษาที่สำเร็จการศึกษาจากระดับอุดมศึกษาจากวิทยาลัยเอกชน

X_i = คะแนนเฉลี่ยสะสมระดับมัธยมศึกษา

Y_i = คะแนนเฉลี่ยสะสมระดับอุดมศึกษา

2.3.4 ทดสอบสมมติฐานเกี่ยวกับสัมประสิทธิ์สหสัมพันธ์

ทดสอบสมมติฐานเกี่ยวกับผลสัมฤทธิ์ทางการเรียนของการศึกษาระดับมัธยมศึกษาตอนปลายและระดับปริญญาตรีว่ามีความสัมพันธ์กันหรือไม่ จำแนกตามสาขาวิชา และรวมทุกสาขา โดยทดสอบความมีนัยสำคัญของสัมประสิทธิ์สหสัมพันธ์แบบเพียร์สัน ซึ่งมีสมมติฐานเพื่อการทดสอบเป็นดังนี้

H_0 : ผลสัมฤทธิ์ทางการเรียนระหว่างการศึกษาในระดับมัธยมศึกษาและระดับปริญญาตรี
ไม่มีความสัมพันธ์กัน

H_a : ผลสัมฤทธิ์ทางการเรียนระหว่างการศึกษาในระดับมัธยมศึกษาและระดับปริญญาตรี
มีความสัมพันธ์กัน

หรือ $H_0 : \rho = 0$

$H_a : \rho \neq 0$

ตัวสถิติที่ใช้ในการทดสอบคือ

$$t_{n-1} = \frac{r}{s_r}$$

$$\text{เมื่อ } r = \frac{n\sum x_i y_i - \sum x_i \sum y_i}{\sqrt{[n\sum x_i^2 - (\sum x_i)^2][n\sum y_i^2 - (\sum y_i)^2]}}$$

$$s_r = \frac{\sqrt{1-r^2}}{\sqrt{n-2}}$$

นำค่า t_{n-2} ที่คำนวณได้มาเปรียบเทียบกับค่า $t_{(n-2, \frac{\alpha}{2})}$ ถ้าค่า t ที่คำนวณมากกว่า $t_{(n-2, \frac{\alpha}{2})}$ จะปฏิเสธ H_0 หรือยอมรับ H_a

2.3.5 การวิเคราะห์ห้ำสหสัมพันธ์คาโนนิกอล (Canonical Correlation Analysis)

ตามปกติเมื่อมีตัวแปรลุ่มอยู่ 2 ตัวคือ X และ Y จะสามารถหาความสัมพันธ์ของตัวแปรทั้งสองได้ โดยพิจารณาจากสัมประสิทธิ์สหสัมพันธ์ (Correlation Coefficient) หรือ ρ เมื่อ

$$\rho = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}}$$

ถ้าค่า Y ขึ้นอยู่กับค่า X หลาย ๆ ตัว วิธีการที่จะสามารถหาความสัมพันธ์ระหว่าง Y และ X เหล่านั้นจะใช้วิธีวิเคราะห์หสัมพันธ์พหุคูณ (Multiple Correlation)

โดยให้ X เป็นปัจจัยหรือตัวแปรที่ต้องการศึกษาซึ่งมีอยู่ P ตัว

$$\text{ดังนั้น } X' = [X_1 \ X_2 \ \dots \ X_p]$$

และ \hat{Y} เป็นค่าประมาณของผลรวมเชิงเส้น

$$\hat{Y}_j = \hat{\beta}_1 X_{1j} + \hat{\beta}_2 X_{2j} + \dots + \hat{\beta}_p X_{pj}, \quad j = 1, 2, \dots, n$$

เมื่อ $\hat{\beta}_s$ ($s = 1, 2, \dots, p$) เป็นค่าสัมประสิทธิ์การถดถอยพหุคูณ (Multiple Regression Coefficient) ซึ่งหาค่าได้จากการประมาณโดยวิธี least-Square อันจะทำให้สามารถหาค่าสัมประสิทธิ์หสัมพันธ์โดยการหาค่าความสัมพันธ์อย่างง่าย (Simple Correlation) ระหว่าง Y และ \hat{Y} นั้นเอง ยังมีอีกกรณีหนึ่ง ซึ่งนอกจากจะมีตัวแปร X หลาย ๆ ตัวแล้วยังมี Y หลาย ๆ ตัวคือ

$$Y' = [Y_1 \ Y_2 \ \dots \ Y_q]$$

วิธีการที่จะหาความสัมพันธ์ระหว่าง X หลาย ๆ ตัว และ Y หลาย ๆ ตัวก็คือ

วิธีวิเคราะห์หสัมพันธ์คาโนนิกอล จะได้ค่าความสัมพันธ์คาโนนิกอล (Canonical Correlation) ซึ่งหาได้จากความสัมพันธ์ระหว่างผลรวมเชิงเส้นของกลุ่มตัวแปร X และผลรวมเชิงเส้นของกลุ่มตัวแปร Y

เมื่อกำหนดค่า

$$\begin{aligned} X_j^* &= \alpha_1 X_{1j} + \alpha_2 X_{2j} + \dots + \alpha_p X_{pj} \\ Y_j^* &= \beta_1 Y_{1j} + \beta_2 Y_{2j} + \dots + \beta_q Y_{qj} \end{aligned} \quad \dots \dots \dots (1)$$

สมการนี้ยังไม่ทราบค่าของ α_i และ β_i ที่จะทำให้ความสัมพันธ์ระหว่าง Y^* และ X^* มีค่ามากที่สุดที่จะเป็นไปได้

ถ้าให้ ρ_c = Canonical Correlation Coefficient ระหว่าง X และ Y

$$\text{ซึ่ง } \rho_c = \frac{\text{Cov}(X^*, Y^*)}{\sqrt{\text{Var}(X^*) \text{Var}(Y^*)}} \dots\dots\dots(2)$$

และ $\rho_c^{(1)}$ เป็น Canonical Correlation Coefficient ที่มีค่ามากที่สุด
เรียกว่า First Canonical Correlation ซึ่งเกิดจาก X_1^* และ Y_1^*

$\rho_c^{(2)}$ เป็น Canonical Correlation Coefficient ที่มีค่ามากที่สุด
เป็นอันดับสองเรียกว่า Second Canonical Correlation ซึ่งเกิดจาก X_2^* และ Y_2^*

$\rho_c^{(3)}$ เป็น Canonical Correlation Coefficient ที่มีค่ามากที่สุด
เป็นอันดับสามเรียกว่า Third Canonical Correlation ซึ่งเกิดจาก X_3^* และ Y_3^*

- \vdots
- \vdots
- \vdots
- \vdots
- X_1^* และ Y_1^* เรียกว่า first canonical variables
- X_2^* และ Y_2^* เรียกว่า second canonical variables
- X_3^* และ Y_3^* เรียกว่า third canonical variables
- \vdots
- \vdots
- \vdots
- \vdots

คู่ของ canonical variable ต่าง ๆ นั้น จะมีหลักการเหมือนกัน
กรณี $q < p$ ก็จะมี Canonical Correlation ทั้งหมด q ค่า
และ Canonical Variable ทั้งหมด q คู่
ผลลัพธ์ต่าง ๆ เหล่านี้ได้มาจากหลักการต่อไปนี้คือ

ถ้า X มีค่าเฉลี่ยเป็น μ_1 และความแปรปรวนเท่ากับ Σ_{11}
Y มีค่าเฉลี่ยเป็น μ_2 และความแปรปรวนเท่ากับ Σ_{22}

$$\text{ให้ } \Sigma = \begin{bmatrix} \Sigma_{11} & \vdots & \Sigma_{12} \\ \vdots & \ddots & \vdots \\ \Sigma_{21} & \vdots & \Sigma_{22} \end{bmatrix}_{(p+q) \times (p+q)}$$



โดยที่ Σ_{11} เป็นเมตริกซ์ของความแปรปรวน (Variance matrix) ของ X
ซึ่งมีขนาด $p \times p$

Σ_{22} เป็นเมตริกซ์ของความแปรปรวน (Variance matrix) ของ Y
ซึ่งมีขนาด $q \times q$

Σ_{12} เป็นเมตริกซ์ของความแปรปรวนร่วม (Covariance matrix) ของ X
และ Y ซึ่งมีขนาด $p \times q$

$$\Sigma_{21} = \Sigma_{12}'$$

สามารถเขียนรายละเอียดได้เป็น

$$\Sigma = \begin{bmatrix} \sigma_{x_1x_1} & \sigma_{x_1x_2} & \dots & \sigma_{x_1x_p} & \vdots & \sigma_{x_1y_1} & \sigma_{x_1y_2} & \dots & \sigma_{x_1y_q} \\ \sigma_{x_2x_1} & \sigma_{x_2x_2} & \dots & \sigma_{x_2x_p} & \vdots & \sigma_{x_2y_1} & \sigma_{x_2y_2} & \dots & \sigma_{x_2y_q} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \sigma_{x_px_1} & \sigma_{x_px_2} & \dots & \sigma_{x_px_p} & \vdots & \sigma_{x_py_1} & \sigma_{x_py_2} & \dots & \sigma_{x_py_q} \\ \hline \sigma_{x_1y_1} & \sigma_{x_2y_1} & \dots & \sigma_{x_py_1} & \vdots & \sigma_{y_1y_1} & \sigma_{y_1y_2} & \dots & \sigma_{y_1y_q} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \sigma_{x_1y_q} & \sigma_{x_2y_q} & \dots & \sigma_{x_py_q} & \vdots & \sigma_{y_qy_1} & \sigma_{y_qy_2} & \dots & \sigma_{y_qy_q} \end{bmatrix}$$

Canonical variables ในรูปของ matrix คือ

$$X^* = X\alpha$$

$$Y^* = Y\beta$$

เมื่อ X^* และ Y^* เป็นเวกเตอร์ของ Canonical variable

ขนาด $p \times 1$ และ $q \times 1$ ตามลำดับ

X และ Y เป็นเมตริกซ์ของตัวแปรสุ่มแรกเริ่ม (Original random variables)

ขนาด $p \times p$ และ $q \times q$ ตามลำดับ

α เป็น Coefficient vectors ขนาด $p \times 1$

β เป็น Coefficient vectors ขนาด $q \times 1$

ดังนั้นจากสมการ (2) สามารถเขียน Canonical Correlation Coefficient

ได้เป็น

$$\rho_c = \frac{\text{Cov}(X\alpha, Y\beta)}{\sqrt{\text{Var}(X\alpha) \text{Var}(Y\beta)}}$$

หรือ
$$\rho_c = \frac{\alpha' \Sigma_{12} \beta}{\sqrt{(\alpha' \Sigma_{11} \alpha) (\beta' \Sigma_{22} \beta)}} \dots \dots \dots (2a)$$

หาค่า ρ_c จากสมการ (2a) โดยที่ต้องการได้ ρ_c ที่มีค่าสูงสุดดังนั้นจะ Maximize ρ_c เทียบกับ α และ β ซึ่งวิธีดังกล่าวนี้จะง่ายขึ้น เมื่อ Canonical Variable แต่ละตัวมีคุณสมบัติที่ว่าความแปรปรวนเป็น 1 นั่นคือ

$$\text{Var}(X_j^*) = \alpha' \Sigma_{11} \alpha = 1 \dots \dots \dots (3)$$

$$\text{Var}(Y_j^*) = \beta' \Sigma_{22} \beta = 1 \dots \dots \dots (4)$$

ซึ่งเงื่อนไขนี้จะทำให้ได้ $\rho_c = \alpha' \Sigma_{12} \beta$ เมื่อ Maximize ρ_c โดยวิธี Least Square ซึ่งจะได้ดังนี้

$$\Sigma_{12} \beta - \lambda \Sigma_{11} \alpha = 0 \dots \dots \dots (5)$$

$$\Sigma_{21} \alpha - \lambda \Sigma_{22} \beta = 0 \dots \dots \dots (6)$$

เมื่อ λ เป็น Lagrange multiplier

คูณสมการ (5) ด้วย λ ได้

$$\Sigma_{12} \lambda \beta = \lambda^2 \Sigma_{11} \alpha \dots \dots \dots (7)$$

คูณสมการ (6) ด้วย Σ_{22}^{-1}

$$\Sigma_{22}^{-1} \Sigma_{21} \alpha = \lambda \beta \dots \dots \dots (8)$$

นำสมการ (8) แทนในสมการ (7)

$$\Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \alpha = \lambda^2 \Sigma_{11} \alpha$$

คูณทั้งสองข้างด้วย Σ_{11}^{-1} และจัดเทอมใหม่ จะได้

$$(\Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} - \lambda^2 I) \alpha = 0 \dots \dots \dots (9)$$

ทำนองเดียวกันจะได้

$$(\Sigma_{22}^{-1}\Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12} - \lambda^2 I) \beta = 0 \dots\dots\dots (10)$$

ทั้งสมการ (9) และ (10) เป็น homogeneous equations ซึ่งสามารถหาคำตอบได้ แต่จากการที่ได้กำหนดไว้ก่อนแล้วว่า $q \leq p$ ดังนั้น ถ้าหากทราบค่า λ^2 และ β แล้วก็จะสามารถหา α ได้จากสมการ (5)

เนื่องจาก

$$\alpha = \frac{\Sigma_{11}^{-1}\Sigma_{12}\beta}{\lambda} \dots\dots\dots (11)$$

จากสมการ (5) เมื่อคูณด้วย α' จะได้

$$\alpha'\Sigma_{12}\beta = \lambda\alpha'\Sigma_{11}\alpha$$

จากสมการ (6) เมื่อคูณด้วย β' จะได้

$$\beta'\Sigma_{21}\alpha = \lambda\beta'\Sigma_{22}\beta$$

$$\text{เนื่องจาก } \alpha'\Sigma_{11}\alpha = \beta'\Sigma_{22}\beta = 1$$

$$\text{ดังนั้น } \alpha'\Sigma_{12}\beta = \beta'\Sigma_{21}\alpha = \lambda \dots\dots\dots (12)$$

นั่นคือ จากสมการ (2a) จะได้ว่า $\alpha'\Sigma_{12}\beta = \lambda$ เป็น Canonical Correlation Coefficient ภายใต้เงื่อนไขของสมการ (3) และสมการ (4)

$$\rho_c = \lambda = \alpha'\Sigma_{12}\beta$$

ถ้าหากถอดรากที่สองของ characteristic root coefficient ที่มีค่ามากที่สุดของสมการ (9) หรือ (10) จะได้ first canonical correlation

จากสมการ (10) หาค่า λ^2 และ β ได้ และจาก (11) หาค่า α ได้

ให้ $\lambda_1^2 > \lambda_2^2 > \dots > \lambda_q^2$ เป็น characteristic root ของสมการ (10)

$\beta_1, \beta_2, \dots, \beta_q$ เป็น characteristic vectors ที่สัมพันธ์กับ

characteristic root λ_i^2 ดังกล่าว

$\alpha_1, \alpha_2, \dots, \alpha_q$ เป็น characteristic vectors ซึ่งคำนวณได้จากอาร์แทนค่า λ_i และ β_i ลงใน (11)

ดังนั้น Canonical Correlation Coefficient ที่ i คือ λ_i

และเซตของ canonical variables ก็คือ $X_i^* = X\alpha_i$
 $Y_i^* = Y\beta_i$

ค่าสัมประสิทธิ์สหสัมพันธ์คาโนนิคอลล (Canonical Correlation Coefficient) ที่ได้มีคุณสมบัติเช่นเดียวกับค่าสัมประสิทธิ์สหสัมพันธ์อย่างง่าย (Simple Correlation Coefficient) กล่าวคือค่าสัมบูรณ์ของค่าสัมประสิทธิ์สหสัมพันธ์คาโนนิคอลลจะมีค่าอยู่ระหว่าง 0 และ 1 และจะไม่เปลี่ยนค่าแม้ว่าจะเปลี่ยนหน่วยการวัดก็ตาม

การคำนวณ เนื่องจากคุณสมบัติดังกล่าวข้างต้นของ Canonical Correlation Coefficient ดังนั้นเมื่อไม่ทราบค่า Σ เราจะประมาณ Σ ด้วย Covariance matrix (X) แทนโดยที่

$$S = \begin{bmatrix} S_{11} & \cdots & S_{12} \\ \cdots & \cdots & \cdots \\ S_{21} & \cdots & S_{22} \end{bmatrix}$$

เมื่อ S เป็น Covariance Matrix ของตัวอย่างหรืออาจจะคำนวณได้จาก Correlation matrix (R) โดยที่

$$R = \begin{bmatrix} R_{11} & \cdots & R_{12} \\ \cdots & \cdots & \cdots \\ R_{21} & \cdots & R_{22} \end{bmatrix}$$

เมื่อ R เป็น Correlation matrix ของตัวอย่าง

R_{11} เป็น matrix ขนาด $p \times p$ ซึ่งแสดงความสัมพันธ์ระหว่างตัวแปร X ซึ่งมี p ตัว

R_{22} เป็น matrix ขนาด $q \times q$ ซึ่งแสดงความสัมพันธ์ระหว่างตัวแปร Y ซึ่งมี q ตัว

R_{12} เป็น matrix ขนาด $p \times q$ ซึ่งแสดงความสัมพันธ์ระหว่างตัวแปร X และตัวแปร Y

$$R_{21} = R'_{12}$$

ค่าของ Canonical Correlation Coefficients ที่คำนวณจาก S หรือ R นั้น
 จะได้ผลเช่นเดียวกัน แต่ค่า α_i และ β_i นั้นจะได้ไม่เหมือนกัน กล่าวคือ

ถ้าหากคำนวณ α_i และ β_i จาก S matrix แล้ว ค่าที่ได้จะเป็นค่าตัวแปรแรกเริ่ม
 (Original Variables) ถ้าหากคำนวณจาก R matrix ค่าที่ได้จะเปลี่ยนอยู่ในรูปตัวแปรแรกเริ่ม
 มาตรฐาน (Standardized Original Variables) อย่างไรก็ตามก็สามารถที่จะเปลี่ยนค่า
 สัมประสิทธิ์ที่คำนวณได้จากการใช้ R นั้นไปเป็นค่าสัมประสิทธิ์ที่ได้จาก S ได้ ดังจะเห็นได้จาก

$$\begin{bmatrix} \frac{1}{\sqrt{s_{11}}} & 0 \\ 0 & \frac{1}{\sqrt{s_{22}}} \end{bmatrix} \begin{bmatrix} s_{11} & s_{12} \\ s_{21} & s_{22} \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{s_{11}}} & 0 \\ 0 & \frac{1}{\sqrt{s_{22}}} \end{bmatrix} = \begin{bmatrix} 1 & R_{12} \\ R_{21} & 1 \end{bmatrix}$$

ดังนั้น เมื่อนำ diagonal matrix ของส่วนกลับของค่าส่วนเบี่ยงเบนมาตรฐาน
 มาคูณทั้งหน้าและหลังของ Covariance Matrix แล้ว ก็จะสามารถหาค่าของ Correlation
 matrix α, β ได้ กล่าวคือ

$$\text{ถ้า } \alpha_{i1}^*, \alpha_{i2}^*, \dots, \alpha_{ip}^* \text{ และ } \beta_{i1}^*, \beta_{i2}^*, \dots, \beta_{iq}^* ;$$

$$(i = 1, 2, \dots, q)$$

เป็นค่าที่คำนวณได้จาก matrix R แล้วจะสามารถหาค่า α, β ซึ่งคำนวณได้จาก S ได้โดย
 การหารแต่ละค่าสัมประสิทธิ์เหล่านั้นด้วยค่าส่วนเบี่ยงเบนมาตรฐานของตัวแปรแรกเริ่ม (Original
 Variable) นั่นคือ

$$\alpha_{i1}^*/s_{x_{11}}, \alpha_{i2}^*/s_{x_{22}}, \alpha_{i3}^*/s_{x_{33}}, \dots, \alpha_{ip}^*/s_{x_{pp}}$$

และ

$$\beta_{i1}^*/s_{y_{11}}, \beta_{i2}^*/s_{y_{22}}, \dots, \beta_{iq}^*/s_{y_{qq}}$$

2.3.5.1 การทดสอบนัยสำคัญของ การวิเคราะห์หลักส่วนประกอบ

เมื่อทำการวิเคราะห์หลักส่วนประกอบแล้ว สิ่งที่จะได้คือ

- (1) p_c
- (2) α และ β

ค่าของ α และ β ซึ่งสอดคล้องกับ λ^2 ของเมตริกซ์ $\Sigma_{11}^{-1}\Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}$ โดยที่ λ^2 ก็คือความแปรปรวนที่ร่วมกันของตัวแปรทั้งสองชุดนั่นเอง ด้วยเหตุนี้ก็จะนิยาม α และ β อยู่หลายชุดด้วยกัน ดังนั้นสิ่งที่จะต้องคำนึงถึงคือ Canonical Variables ทั้ง q คู่ นั้นมีความสัมพันธ์กันหรือไม่ คือต้องมีการทดสอบสมมติฐานที่ว่า ตัวแปรทั้งสองชุดนั้นมีความสัมพันธ์กันหรือไม่ นั่นคือทดสอบว่า Σ_{12} หรือ $\Sigma_{21} = 0$ หรือไม่

$$H_0 : \Sigma_{12} = 0 \quad \text{หรือ} \quad H_0 : \Sigma_{21} = 0$$

$$H_a : \Sigma_{12} \neq 0 \quad \text{หรือ} \quad H_a : \Sigma_{21} \neq 0$$

ใช้วิธีการทดสอบของ Bartlett ดังนี้

วิธีการทดสอบของ Bartlett

ใช้ตัวสถิติ V_0 โดยที่

$$V_0 = -\left[n-1-\frac{1}{2}(p+q+1)\right] \ln \Lambda_0$$

เมื่อ

$$\Lambda_0 = \frac{1}{\prod_{i=1}^q \left[1 + \frac{\lambda_i^2}{(1 - \lambda_i^2)} \right]}$$

$$= \frac{1}{\prod_{i=1}^q \left[\frac{1}{(1 - \lambda_i^2)} \right]}$$

$$= \prod_{i=1}^q (1 - \lambda_i^2)$$

V_0 จะมีการกระจายแบบไคลสแควร์ที่องศาความเป็นอิสระเท่ากับ pq ที่ระดับนัยสำคัญที่ต้องการ ถ้าค่า V_0 ที่คำนวณได้มากกว่าค่าไคลสแควร์จากตาราง จะสรุปได้ว่า $\rho_c^{(1)}$ ซึ่งเป็น first canonical correlation coefficient นั้นไม่เท่ากับ 0 และจะทดสอบความสัมพันธ์ของ ρ_c ของคู่ที่เหลือต่อไปคือ

$$\Lambda_1 = \prod_{i=2}^q (1 - \lambda_i^2)$$

ซึ่งจะนำไปสู่ Barlett V_1 ดังนี้

$$V_1 = -\left[n - 2 - \frac{1}{2}(p+q+1)\right] \ln \Lambda_1$$

V_1 มีการกระจายแบบไคลสแควร์ที่องศาความเป็นอิสระเท่ากับ $(p-1)(q-1)$ การสรุปผลจะทำเหมือน V_0 ดังนั้น ถ้าหากมี k canonical correlation coefficients ที่ทดสอบแล้วปรากฏว่ามีนัยสำคัญส่วนที่เหลือจะถูกทดสอบว่า $(q-k)$ Canonical Correlation Coefficient เหล่านั้นจะเท่ากับ 0 หรือไม่ กล่าวคือ

$$\Lambda_k = \prod_{i=(k+1)}^q (1 - \lambda_i^2)$$

ซึ่งจะได้
$$V_k = -\left[n - 1 - k - \frac{1}{2}(p+q+1)\right] \ln \Lambda_k$$

V_k มีการกระจายแบบไคลสแควร์ที่องศาความเป็นอิสระเท่ากับ $(p-k)(q-k)$

ในการประมวลผลด้วยโปรแกรมสำเร็จรูป SPSS นั้น ผลลัพธ์ที่ได้จะมีนัยสำคัญหรือค่าของความน่าจะเป็นที่จะไม่ยอมรับ H_0 เมื่อ H_0 จริง ในกรณีที่ทดสอบสมมติฐาน ณ ระดับนัยสำคัญ $\alpha = 0.05$ นั้น อาจใช้ค่าของความน่าจะเป็นที่จะไม่ยอมรับ H_0 เมื่อ H_0 จริง มาเปรียบเทียบกับ $\alpha = 0.05$ เพื่อใช้เป็นเกณฑ์ในการตัดสินว่าจะยอมรับสมมติฐาน H_0 หรือไม่โดยที่

ถ้าค่าความน่าจะเป็นที่จะไม่ยอมรับ H_0 เมื่อ H_0 จริง < 0.05 จะไม่ยอมรับ H_0

ถ้าค่าความน่าจะเป็นที่จะไม่ยอมรับ H_0 เมื่อ H_0 จริง > 0.05 จะยอมรับ H_0

2.3.6 การวิเคราะห์จำแนกประเภท (Discriminant Analysis)

เป็นวิธีวิเคราะห์ทางสถิติโดยมีจุดมุ่งหมายที่จะคัดเลือกตัวแปรชุดหนึ่ง ซึ่งนักวิจัยคิดว่าตัวแปรชุดนี้มีความสัมพันธ์กับสิ่งที่ต้องการศึกษา จนถึงขั้นที่ตัวแปรชุดนี้เป็นตัวแบ่งแยกประชากรออกเป็นกลุ่มต่าง ๆ ได้อย่างชัดเจน

2.3.6.1 ขั้นตอนในการวิเคราะห์จำแนกประเภท มี 2 ขั้นตอนคือ

ขั้นที่ 1 การคัดเลือกตัวแปรชุดหนึ่ง เพื่อสร้างสมการที่ใช้ในการจำแนกประชากรออกเป็นกลุ่มต่าง ๆ กัน ได้อย่างชัดเจน สมการนี้คือ สมการจำแนกประเภท (Discriminant equation)

ขั้นที่ 2 การจำแนกตัวอย่างที่ได้ศึกษามานั้น เข้า เป็นสมาชิกของประชากร แต่ละกลุ่มโดยอาศัยสมการจำแนกประเภท

ในกรณีที่มีประชากร 2 กลุ่ม จะใช้วิธีการของ Fisher มาใช้ในการจำแนกประเภท กล่าวคือ ถ้ากำหนดให้

π_1 เป็นประชากรกลุ่มที่ 1

π_2 เป็นประชากรกลุ่มที่ 2

ตัวแปรที่ศึกษา คือ ตัวแปร X ซึ่งมีทั้งหมด p ตัวคือ

$$X' = [x_1 \ x_2 \ \dots \ x_p]$$

วิธีการของ Fisher นั้น จะแปลงตัวแปร X เหล่านี้ไปเป็นค่าของตัวแปรเพียงตัวเดียว คือ Y โดยที่ Y_1 และ Y_2 เป็นค่าสังเกตที่ได้จากประชากร π_1 และ π_2 ตามลำดับ

$$\left. \begin{aligned} \text{ถ้า } \mu_{1Y} &= \text{ค่าเฉลี่ยของค่า } Y \text{ ซึ่งได้มาจากค่า } X \text{ ของประชากร } \pi_1 \\ \mu_{2Y} &= \text{ค่าเฉลี่ยของค่า } Y \text{ ซึ่งได้มาจากค่า } X \text{ ของประชากร } \pi_2 \\ \mu_1 &= E(X/\pi_1) = \text{ค่าคาดหวังของตัวแปร } X \text{ ที่มาจากประชากร } \pi_1 \\ \mu_2 &= E(X/\pi_2) = \text{ค่าคาดหวังของตัวแปร } X \text{ ที่มาจากประชากร } \pi_2 \end{aligned} \right\} \dots (1)$$

และโควาเรียนซ์เมตริกซ์ คือ

$$\Sigma = E(X_i - \mu_i)(X_i - \mu_i)' \quad ; \quad i = 1, 2, \dots \dots \dots (2)$$

พิจารณาผลรวมเชิงเส้น

$$\begin{matrix} Y & = & \beta' X & \dots\dots\dots (3) \\ (1 \times 1) & & (1 \times p) (p \times 1) \end{matrix}$$

เมื่อ Y คือ Fisher's Linear Discriminant Function (ผลรวมเชิงเส้น)

β คือ ค่าที่แสดงถึงความสำคัญของตัวแปร X_1, X_2, \dots, X_p

ซึ่ง $\beta' = [\beta_1 \quad \beta_2 \quad \dots \quad \beta_p]$

X คือตัวแปรที่ต้องการศึกษา ซึ่งมีทั้งหมด p ตัว

และ $X' = [X_1 \quad X_2 \quad \dots \quad X_p]$

ดังนั้น จะได้ $\mu_{1Y} = E(Y/\pi_1) = E(\beta'X/\pi_1) = \beta'\mu_1 \dots\dots\dots (4)$

$$\mu_{2Y} = E(Y/\pi_2) = E(\beta'X/\pi_2) = \beta'\mu_2$$

และค่าความแปรปรวนของ Y จากทั้ง 2 ประชากรนั้นคือ

$$\begin{aligned} \sigma_Y^2 &= \text{Var}(\beta'X) = \beta' \text{Cov}(X) \beta \\ &= \beta' \Sigma \beta \dots\dots\dots (5) \end{aligned}$$

วิธีการของ Fisher คือพยายามหาผลรวมเชิงเส้นของค่า X ซึ่งจะทำให้ระยะทางระหว่าง μ_{1Y} และ μ_{2Y} มีค่ามากที่สุดที่จะเป็นไปได้

ค่าผลรวมเชิงเส้นซึ่งดีที่สุดจะต้องมีคุณสมบัติในการแบ่งแยกประชากรทั้ง 2 กลุ่มออกจากกันได้มากที่สุด ซึ่งวิธีการที่จะได้ผลรวมเชิงเส้นที่ดีนั้นก็โดยการหาค่า β ที่ทำให้อัตราส่วน

(ระยะทางระหว่างค่าเฉลี่ยของ Y)² มีค่ามากที่สุด
ความแปรปรวนของ Y

$$\begin{aligned} \text{ให้ } \ell &= \frac{(\text{ระยะทางระหว่างค่าเฉลี่ยของ } Y)^2}{\text{ความแปรปรวนของ } Y} \\ &= \frac{(\mu_{1Y} - \mu_{2Y})^2}{\sigma_Y^2} \end{aligned}$$

$$\begin{aligned}
 &= \frac{(\beta' \mu_1 - \beta' \mu_2)^2}{\beta' \Sigma \beta} \\
 &= \frac{\beta' (\mu_1 - \mu_2) (\mu_1 - \mu_2)' \beta}{\beta' \Sigma \beta} \dots\dots\dots (6)
 \end{aligned}$$

ค่า β ที่ได้จะมีค่ามากที่สุดเมื่อ

$$\frac{\partial L}{\partial \beta} = 0$$

$$\begin{aligned}
 \text{นั่นคือ } C(\mu_1 - \mu_2) - \Sigma \beta &= 0 \\
 \beta &= C \Sigma^{-1} (\mu_1 - \mu_2)
 \end{aligned}$$

กำหนดให้ $C = 1$

จะได้ผลรวมเชิงเส้นเป็น

$$Y = \beta' X = (\mu_1 - \mu_2)' \Sigma^{-1} X \dots\dots\dots (7)$$

สมการ $Y = (\mu_1 - \mu_2)' \Sigma^{-1} X$ สามารถใช้เป็นตัวจำแนกค่าสังเกตที่ได้มานั้นว่าจะอยู่ในประชากรกลุ่ม π_1 หรือ π_2 ได้โดยให้

$$Y_0 = (\mu_1 - \mu_2)' \Sigma^{-1} X_0$$

เป็นค่าของ discriminant function ของค่าสังเกต X_0 และให้

$$\begin{aligned}
 m &= \frac{1}{2}(\mu_{1Y} + \mu_{2Y}) \\
 &= \frac{1}{2}(\beta' \mu_1 + \beta' \mu_2) \\
 &= \frac{1}{2} (\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 + \mu_2) \dots\dots\dots (8)
 \end{aligned}$$

เป็นจุดกึ่งกลาง (centriod) ระหว่างค่าเฉลี่ยของค่า Y จากทั้ง 2 ประชากร

สามารถเขียนกฎการจำแนกประเภทได้ดังนี้คือ

$$\text{จัด } X_0 \text{ ให้อยู่ใน } \pi_1 \text{ ถ้าหากว่า } Y_0 = (\mu_1 - \mu_2)' \Sigma^{-1} X_0 < m \dots\dots\dots (9)$$

$$\text{จัด } X_0 \text{ ให้อยู่ใน } \pi_2 \text{ ถ้าหากว่า } Y_0 = (\mu_1 - \mu_2)' \Sigma^{-1} X_0 > m$$

ในทางปฏิบัติจะไม่ทราบค่า μ_1 , μ_2 และ Σ ดังนั้นถ้าสุ่มตัวอย่างจาก Π_1 และ Π_2 มาเป็นจำนวน n_1 และ n_2 แล้วโดยที่วัดค่าสังเกต

$$\begin{aligned} \tilde{x}' &= [x_1 \quad x_2 \quad \dots \quad x_p] \quad \text{จะได้ว่า} \\ \tilde{x}_{1j} &= [x_{11} \quad x_{12} \quad \dots \quad x_{1n_1}] \\ & \text{(pxn}_1\text{)} \\ \tilde{x}_{2j} &= [x_{21} \quad x_{22} \quad \dots \quad x_{2n_2}] \\ & \text{(pxn}_2\text{)} \end{aligned}$$

จะได้ค่าประมาณของ μ_1 , μ_2 และ Σ เป็น \bar{x}_1 , \bar{x}_2 และ S_{pooled}^{-1} ตามลำดับโดยที่

$$\bar{x}_1 = \frac{1}{n_1} \sum_{j=1}^{n_1} x_{1j}; \quad \bar{x}_2 = \frac{1}{n_2} \sum_{j=1}^{n_2} x_{2j}$$

(px1) (px1)

$$\begin{aligned} S &= \left[\frac{n_1-1}{(n_1-1)+(n_2-1)} \right] S_1 + \left[\frac{n_2-1}{(n_1-1)+(n_2-1)} \right] S_2 \\ &= \frac{(n_1-1)S_1 + (n_2-1)S_2}{n_1+n_2-2} \dots \dots \dots (12) \end{aligned}$$

เมื่อ

$$S_1 = \frac{1}{n_1-1} \sum_{j=1}^{n_1} (x_{1j} - \bar{x}_1)(x_{1j} - \bar{x}_1)'$$

(pxp)

$$S_2 = \frac{1}{n_2-1} \sum_{j=1}^{n_2} (x_{2j} - \bar{x}_2)(x_{2j} - \bar{x}_2)'$$

(pxp)

ดังนั้น เมื่อแทนค่า \bar{x}_1 , \bar{x}_2 และ S_{pooled}^{-1} ลงใน (7) จะได้ Fisher's Sample Linear Discriminant Function เป็น

$$\begin{aligned} X &= \hat{\beta}' X \\ &= (\bar{x}_1 - \bar{x}_2)' S_{\text{pooled}}^{-1} X \dots \dots \dots (13) \end{aligned}$$

และค่าจุดกึ่งกลางระหว่างตัวแปร $\bar{Y}_1 = \hat{\beta}'\bar{X}_1$ และ $\bar{Y}_2 = \hat{\beta}'\bar{X}_2$ จะเขียนได้เป็น

$$\begin{aligned} \hat{m} &= \frac{1}{2} (\bar{X}_1 + \bar{X}_2) \\ &= \frac{1}{2} (\bar{X}_1 - \bar{X}_2)' S_{pooled}^{-1} (\bar{X}_1 + \bar{X}_2) \dots\dots\dots(14) \end{aligned}$$

กฎการจำแนกประเภทจะเขียนดังนี้

- จัดให้ X_0 อยู่ใน π_1 ถ้า $Y_0 < \hat{m}$
- จัดให้ X_0 อยู่ใน π_2 ถ้า $Y_0 \geq \hat{m}$
- เมื่อ $Y_0 = (\bar{X}_1 - \bar{X}_2)' S_{pooled}^{-1} X_0$



จากสมการสามารถเขียนได้เป็น

$$Y = \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p$$

ค่า $\beta_1, \beta_2, \dots, \beta_p$ จะเป็นค่าแสดงถึงความสำคัญของตัวแปร X_1, X_2, \dots, X_p ซึ่งโดยปกติ X_i แต่ละค่ามันจะมีหน่วยไม่เหมือนกัน การเปรียบเทียบค่า β_i เหล่านี้จะทำได้ก็โดยการปรับค่า β_i เหล่านี้ให้เป็นมาตรฐาน (Standardized) เสียก่อน

การแปลงค่า β_i เหล่านี้สามารถทำได้โดยเอาสมาชิกในแนวเส้นทแยงมุม (diagonal) ของเมตริกซ์ W มาถอดรากที่สอง แล้วนำไปคูณกับค่า β_i ทุก ๆ ค่าตามสูตรต่อไปนี้

$$\beta_i^* = (\sqrt{W_{ii}}) (\beta_i)$$

เมื่อ W คือ $(X_i - \bar{X}_{ig0})(X_i - \bar{X}_{ig0})'$

$$X_i = \begin{bmatrix} X_{111} & X_{121} & \dots & X_{1n1} & \vdots & X_{112} & X_{122} & \dots & X_{1n2} \\ X_{211} & X_{221} & \dots & X_{2n1} & \vdots & X_{212} & X_{222} & \dots & X_{2n2} \\ \vdots & \vdots & & \vdots & & \vdots & \vdots & & \vdots \\ X_{p11} & X_{p21} & \dots & X_{pn1} & \vdots & X_{p12} & X_{p22} & \dots & X_{pn2} \end{bmatrix}$$

$$\bar{x}_{g0} = \begin{bmatrix} \bar{x}_{11} \dots \bar{x}_{11} & \vdots & \bar{x}_{12} \dots \bar{x}_{12} \\ \bar{x}_{21} \dots \bar{x}_{21} & \vdots & \bar{x}_{22} \dots \bar{x}_{22} \\ \vdots & \vdots & \vdots \\ \bar{x}_{p1} \dots \bar{x}_{p1} & \vdots & \bar{x}_{p2} \dots \bar{x}_{p2} \end{bmatrix}$$

β_i^* คือ Standardized discriminant weight ของ discriminant function

β_i คือ Row discriminant weight ของ discriminant function

w_{ii} คือ Diagonal element ของ W เมตริกซ์ โดยที่ $i = 1, 2, \dots, p$

สมการจำแนกประเภทที่ได้มานั้นจะมีความสามารถในการแบ่งแยกกลุ่มได้ดีหรือไม่นั้น สามารถดูได้จากค่า miss classification error ซึ่งทราบได้จากการที่เปรียบเทียบกับสิ่งที่เราศึกษาอยู่นั้นเป็นสมาชิกของประชากรกลุ่มหนึ่งและเมื่อใช้สมการจำแนกประเภทมาตัดสินการจำแนกแล้วกลับไปเป็นสมาชิกของอีกกลุ่มหนึ่งซึ่งทำให้ได้ผลที่ผิดจากความเป็นจริงไป

2.3.6.2 วิธีการคัดเลือกตัวแปรเข้าไปในสมการจำแนกประเภท (Selection of Variables)

จากการที่กล่าวมาแล้วข้างต้นนั้นเป็นหลักการโดยทั่ว ๆ ไปของวิธีการสร้างสมการจำแนกประเภทในกรณีนี้หากว่าในการวิจัยนั้น มีตัวแปรเป็นจำนวนมาก การคัดเลือกตัวแปรให้เหลือจำนวนน้อยที่สุดแต่มีความสามารถในการใช้เป็นตัวจำแนกมากที่สุดนั้น เราสามารถทำได้โดยใช้การคัดเลือกตัวแปรทีละตัว โดยที่จะหาตัวแปรที่ดีที่สุดตัวแรกและตัวแปรที่ดีที่สุดตัวที่สองที่จะทำให้สมการจำแนกประเภทมีความสามารถจำแนกประเภทได้ดีที่สุด จากนั้นก็จะเลือกตัวที่สามและตัวต่อไป ที่จะช่วยการจำแนกให้ดีขึ้นตามลำดับ ในแต่ละขั้นตอนตัวแปรที่ได้รับการคัดเลือกมาก่อนแล้วนั้น อาจถูกตัดทิ้งไปหากพบว่าเมื่อนำมารวมกับตัวแปรตัวอื่น ๆ แล้วไม่ช่วยให้สมการจำแนกประเภทได้ดีขึ้น วิธีการนี้เรียกว่าวิธีการสร้างสมการจำแนกประเภทแบบขั้นตอน (Stepwise discriminant analysis)

วิธีการคัดเลือกตัวแปรแบบขั้นตอนนี้เป็นวิธีการคัดเลือกตัวแปรอิสระทีละตัวตามอำนาจการจำแนกของตัวแปรนั้น โดยที่ตัวแปรตัวแรกที่ถูกเลือกเข้าไปนั้นจะมีอำนาจในการจำแนกสูงสุดเมื่อเทียบกับตัวแปรอื่น ๆ ตัวแปรที่เหลือก็จะถูกคัดเลือกเช่นเดียวกัน ทำเช่นนี้จนกระทั่งถึงเกณฑ์หนึ่ง ๆ ซึ่งเราได้กำหนดไว้ว่าตัวแปรที่เหลือนั้นจะไม่สามารถเข้าไปในสมการจำแนกได้อีกแล้ว วิธีการนี้จะคล้ายคลึงกับวิธีการวิเคราะห์หัตถดถอยพหุแบบขั้นตอน (Stepwise multiple regression analysis) ซึ่งโดยปกติเกณฑ์ที่ใช้คือ

$$\text{กำหนดให้ } \hat{\ell}_p = (\bar{X}_1 - \bar{X}_2)' S_{\text{pooled}}^{-1} (\bar{X}_1 - \bar{X}_2)$$

เป็นค่าที่เกิดจากตัวแปร p ตัว

$$\hat{\ell}_{p+q} = (\bar{X}_1 - \bar{X}_2)' S_{\text{pooled}}^{-1} (\bar{X}_1 - \bar{X}_2)$$

เป็นค่าที่เกิดจากตัวแปร $p + q$ ตัว

จะกล่าวได้ว่าตัวแปร q ตัวที่เพิ่มขึ้นมานั้น จะไม่มีอำนาจในการจำแนกนั้นคือไม่ได้ช่วยให้สมการจำแนกเดิมที่ใช้ตัวแปรอยู่ p ตัวแล้วนั้นมีประสิทธิภาพดีขึ้น ถ้าหากว่า

$$F = \frac{n_1 + n_2 - p - q - 1}{q} \times \frac{m(\hat{\ell}_{p+q} - \hat{\ell}_p)}{1 + m\hat{\ell}_p} < F_{\alpha, q, n_1 + n_2 - p - q - 1}$$

โดยที่ n_1 = จำนวนตัวอย่างจากประชากร π_1

n_2 = จำนวนตัวอย่างจากประชากร π_2

$$m = \frac{n_1 n_2}{(n_1 + n_2)(n_1 + n_2 - 2)}$$

q = ค่าตัวแปรที่จะเข้าไปในสมการจำแนกในทีละขั้นตอน q จึงมีค่าเป็น 1 เสมอ

p = จำนวนตัวแปรที่มีอยู่เต็มแล้วในสมการจำแนก

α = ระดับนัยสำคัญที่กำหนดไว้

จากวิธีการดังกล่าวข้างต้น จะทำให้ได้สมการจำแนกประเภทและค่า eigenvalue ซึ่งเป็นค่าที่ได้จากขั้นตอนการหาสมการวิเคราะห์จำแนกประเภทโดยเป็นค่าความแปรปรวนของค่า Y ซึ่งได้จากการแปลงรูป (Transform) จากค่า X ต่าง ๆ (X_1, X_2, \dots, X_p)

เป็นค่าที่ใช้วัดความสำคัญเชิงเปรียบเทียบของสมการที่ได้ โดยที่สมการวิเคราะห์จำแนกประเภทที่ได้มานั้นได้ตามลำดับความสำคัญของค่า eigenvalue ต่าง ๆ จากมากที่สุดและรองลงไปตามลำดับ ค่าผลรวมของ eigenvalue ทั้งหมดนั้นเป็นค่าของความแปรปรวน (Total variance) ทั้งหมดของตัวแปรจำแนกประเภท (ตัวแปร X) ค่า eigenvalue แต่ละค่าจึงคิดเป็นอัตราส่วนร้อยละของค่ารวมของ eigenvalue ทั้งหมด ทำให้สามารถใช้ค่านี้ไปอ้างอิงถึงความสำคัญเชิงเปรียบเทียบของสมการจำแนกประเภทได้

2.3.6.3 การทดสอบนัยสำคัญของสมการจำแนกประเภท

เมื่อได้สมการจำแนกประเภทมาแล้วและอยากทราบว่าสมการนี้สามารถมีอำนาจในการจำแนกกลุ่มได้อย่างมีนัยสำคัญทางสถิติหรือไม่ เราสามารถทดสอบได้จาก

$$V_m = [N-1-(p+k)/2] \ln(1+\lambda_m)$$

โดยที่ V_m คือ Bartlett V Statistic ซึ่งมีการกระจายเป็นแบบ Chi-Square

ที่มี d.f. = $p+k-2m$

N คือ จำนวนตัวอย่างทั้งหมดจากประชากรทั้ง 2 กลุ่ม

p คือ จำนวนตัวแปรทั้งหมด

k คือ จำนวนกลุ่ม

λ_m คือ Eigenvalues ต่าง ๆ เช่น ค่าตัวที่ 1, 2, 3,, m

ค่า λ_m นี้คำนวณได้จาก

$$|A - \lambda_m I| = 0$$

เมื่อ A คือ โควาเรียนซ์ เมตริกซ์ของ X

ค่า λ_m นี้จะเลือกเอาเฉพาะค่าที่ไม่เป็นศูนย์

ถ้าค่า V_m ที่คำนวณได้มากกว่าค่า Chi-Square ที่ระดับนัยสำคัญที่กำหนดไว้แล้วแสดงว่าสมการนี้สามารถใช้จำแนกกลุ่มได้อย่างมีนัยสำคัญทางสถิติ

การคำนวณค่าอำนาจการจำแนกของกลุ่มของตัวแปร ((Total Discriminatory Power)

$$\hat{W}^2 = 1 - \frac{N}{(m-k)(1+\lambda_1) + (1+\lambda_2) + \dots + (1+\lambda_m) + 1}$$

- เมื่อ $\hat{w}^2 =$ ค่าอำนาจในการแยกตัวแปรได้จากค่าประมาณค่า w^2
 $N =$ จำนวนตัวอย่างทั้งหมดจากประชากรทั้ง 2 กลุ่ม
 $k =$ จำนวนกลุ่ม
 $\lambda_m =$ Eigenvalues ต่าง ๆ
 $m =$ จำนวนค่า Eigenvalues

2.3.7 การวิเคราะห์ความถดถอยพหุคูณ (Multiple Regression Analysis)

ตัวแบบทางคณิตศาสตร์ที่จะนำมาใช้อธิบายวิธีหนึ่งที่จะอธิบายความสัมพันธ์ระหว่าง

สัมฤทธิ์ผลทางการเรียนกับสถานการณ์ส่วนตัวและหมวดปัญหาต่าง ๆ ที่เกี่ยวข้องกับตัวนักศึกษา

ได้แก่ ตัวแบบเชิงเส้นของการถดถอยซึ่งเป็นที่นิยมใช้กันอย่างกว้างขวางในการวิเคราะห์เชิงเส้น

ทั้งนี้เพราะส่วนใหญ่แล้วจะอยู่ในรูปของความสัมพันธ์ระหว่างตัวแปร โดยเฉพาะอย่างยิ่ง

ถ้าหากว่าความสัมพันธ์อยู่ในรูปของความสัมพัทธ์ระหว่างตัวแปรหลาย ๆ ตัว ดังเช่น

ที่เราได้ทดลอง χ^2 แล้วว่า เกรดเฉลี่ยสะสมในระดับอุดมศึกษานั้นขึ้นอยู่กับ เพศ เกรดเฉลี่ยสะสม

ในระดับมัธยมศึกษาตอนปลายและปัจจัยอื่น ๆ อีก ในการวิเคราะห์เช่นนี้วิธีการถดถอยพหุคูณก็มี

ความเหมาะสมเพราะเราสามารถใช้ในการคำนวณเพื่อหาลักษณะทำนายใช้ในการสำรวจหาตัวแปร

อิสระอีกครั้งหนึ่งหลังจากที่เราได้ตัดตัวแปรโดยใช้ χ^2 แล้ว และนำตัวแปรที่ได้นี้มาตัดอีกทีโดยใช้

วิธีการในการเลือกกลุ่มการถดถอยที่ดีที่สุด โดยการถดถอยพหุคูณนี้มีตัวแบบดังต่อไปนี้คือ

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + \epsilon$$

โดยที่ β_i เป็นค่าสัมประสิทธิ์การถดถอยพหุคูณ ไม่ทราบค่าของตัวแบบ; $i = 0, 1, 2, \dots, k$

และถ้าหากว่าเรามีจำนวนข้อมูล n ตัว ต่อตัวแปรแต่ละตัวแล้วจะได้รูปของข้อมูลเป็น

ดังนี้คือ

ข้อมูล	ตัวแปรตาม	ตัวแปรอิสระที่ 1	ที่ 2	ที่ k
1	Y_1	X_{11}	X_{21}	X_{k1}
2	Y_2	X_{12}	X_{22}	X_{k2}
3	Y_3	X_{13}	X_{23}	X_{k3}
\vdots	\vdots	\vdots	\vdots	\vdots
n	Y_n	X_{1n}	X_{2n}	X_{kn}

จากข้อมูลข้างต้นนี้ ถ้าหากเราจะนำมาเขียนในรูปของสมการเมตริกซ์จะได้เป็น

$$\begin{bmatrix} Y_1 \\ Y_2 \\ Y_3 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} 1 & X_{11} & X_{21} & \dots & X_{k1} \\ 1 & X_{12} & X_{22} & \dots & X_{k2} \\ 1 & X_{13} & X_{23} & \dots & X_{k3} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{1n} & X_{2n} & \dots & X_{kn} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

หรือ $Y = X\beta + \varepsilon$

โดยที่ Y คือข้อมูลของตัวแปรตาม โดยมีเวกเตอร์อันดับที่ $n \times 1$

X คือข้อมูลของตัวแปรอิสระ โดยมีเมตริกซ์อันดับที่ $n \times (k+1)$

β คือสัมประสิทธิ์ของการถดถอย โดยมีเวกเตอร์อันดับที่ $(k+1) \times 1$

ε คือค่าความคลาดเคลื่อนของค่าประมาณจากค่าจริง โดยมีเวกเตอร์อันดับที่ $n \times 1$

2.3.7.1 ข้อสมมติเบื้องต้น

ในการที่เราจะศึกษาถึงวิธีการประมาณและคุณสมบัติของรูปแบบซึ่งเราไม่ทราบค่า นั้น เราจะต้องตั้งข้อสมมติ ซึ่งจำเป็นสำหรับการประมาณเสียก่อนและข้อสมมติที่จะกล่าวถึงต่อไปก็มีความสำคัญมาก เพราะคุณสมบัติของตัวประมาณโดยวิธีกำลังสองน้อยที่สุด (The Least Squares Method) ซึ่งเป็นวิธีที่เรานำมาใช้จะขึ้นอยู่กับข้อสมมติเหล่านั้น ซึ่งจะต้องนำมาพิจารณาประกอบก่อนลงมือทำการวิเคราะห์คือ

$$1. Y = X\beta + \varepsilon$$

นั่นคือ Y_i เป็นฟังก์ชันเชิงเส้นของ X_{ij} กับ ε_i โดยที่ $i = 1, 2, \dots, n$

$$j = 0, 1, \dots, k \text{ และ } X_{0i} = 1$$

$$2. E(\varepsilon) = 0$$

นั่นคือ $E(\varepsilon_i) = 0$ สำหรับทุกค่า i

$$3. E(\varepsilon\varepsilon') = \sigma^2 I$$

นั่นคือ $E(\varepsilon_i^2) = \sigma^2$ สำหรับทุกค่า i นั่นคือ ε_i มีความแปรปรวนคงที่

$E(\varepsilon_i \varepsilon_j) = 0$ เมื่อ $i \neq j$ แสดงว่าแต่ละคู่ของ ε ไม่เกี่ยวข้องกัน
(pairwise uncorrelated)

4. ค่าความคลาดเคลื่อนของค่าประมาณจากค่าจริง ε_i เป็นตัวแปรเชิงสุ่มที่มีการกระจายแบบปกติ $N(0, \sigma^2)$ ซึ่งข้อนี้จำเป็นต้องใช้ในกรณีที่ต้องการสร้างช่วงความเชื่อมั่นของ $\beta_0, \beta_1, \beta_2, \dots, \beta_k$ หรือเมื่อต้องการทดสอบสมมติฐานใด ๆ เกี่ยวกับตัวประมาณ

5. X เป็นเมตริกซ์อันดับที่ $n \times (k+1)$ ซึ่งมีค่าคงที่ สิ่งทำให้ X และ ε เป็นตัวแปรอิสระกัน เราจึงได้ว่า

$$E(\varepsilon|X) = E(\varepsilon) = 0$$

$$E(X'\varepsilon) = X'E(\varepsilon) = 0$$

$$E(\varepsilon\varepsilon'|X) = E(\varepsilon\varepsilon') = \sigma^2 I$$

6. X มี rank $k+1 < n$ หมายความว่าจำนวนค่าสังเกต n จะมากกว่าจำนวนพารามิเตอร์ $(k+1)$ ที่จะต้องประมาณค่าออกมา

จากข้อสมมติดังกล่าวมาแล้วข้างต้นเราจะเห็นได้ว่า

$$1. E(Y|X) = X\beta + E(\varepsilon|X) = X\beta$$

$$\text{เช่น } E(Y_i | X_{1i}, X_{2i}, \dots, X_{ki}) = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_k X_{ki}$$

$$2. E[(Y - E(Y))(Y - E(Y))' | X] = E(\varepsilon\varepsilon' | X) = \sigma^2 I$$

$$\text{เช่น } E[(Y_i - E(Y_i))(Y_i - E(Y_i))' | X_{1i}, X_{2i}, \dots, X_{ki}] = \sigma^2 \text{ สำหรับทุกค่า } i$$

2.3.7.2 การประมาณค่าพารามิเตอร์

วิธีการในการประมาณค่าพารามิเตอร์การถดถอยนั้น กระทำได้ 2

วิธีคือ

1. วิธีการกำลังสองน้อยที่สุด (Least Squares Method)

2. วิธีน่าจะเป็นสูงสุด (Maximum Likelihood Method)

ทั้ง 2 วิธีดังกล่าวเป็นวิธีในการคำนวณหาค่า เส้นถดถอยจากข้อมูลที่กำหนดให้หรือ
กล่าวอีกนัยหนึ่งก็คือ ถ้ากำหนดสมการ

$$\begin{aligned} Y &= X\beta + \epsilon && \text{แล้วเทคนิคทั้งหมดข้างต้นก็จะเป็นเครื่องมือในการคำนวณ} \\ \hat{Y} &= X\hat{\beta} && \text{นั่นเอง หากแต่ว่าแต่ละวิธีก็มีคุณสมบัติของตัวเอง}$$

แต่ในที่นี้เราจะนำมากล่าวถึงเพียงวิธีเดียวเท่านั้นคือวิธีที่ 1 ซึ่งถือว่าเป็นวิธีที่สำคัญ
และนิยมใช้กันมาก โดยมีหลักการของการประมาณแบบกำลังสองน้อยที่สุด คือการทำให้ผลรวม
กำลังสองของส่วนเบี่ยงเบนของค่าสังเกตต่างจากค่าเฉลี่ยมีค่าน้อยที่สุดและจะทำให้ตัวประมาณที่
ไม่เอนเอียง

$$\begin{aligned} \text{จากรูปแบบ} \quad Y &= X\beta + \epsilon \\ \text{และ} \quad E(Y) &= X\beta && \because E(\epsilon) = 0 \\ \text{ซึ่ง} \quad E(Y) &&& \text{คือค่าเฉลี่ยของ } Y \\ \text{ดังนั้น} \quad \epsilon &= Y - E(Y) \\ &= Y - X\beta\end{aligned}$$

โดยวิธีกำลังสองน้อยที่สุด เราจะต้องหาหาค่ากำลังสองของค่าความคลาดเคลื่อนเท่ากับ

$$\begin{aligned}\epsilon'\epsilon &= (Y - X\beta)'(Y - X\beta) \\ &= Y'Y - \beta'X'Y - Y'X\beta + \beta'X'X\beta \\ &= Y'Y - 2\beta'X'Y + \beta'X'X\beta\end{aligned}$$

[เนื่องจาก $Y'X\beta$ เป็นสเกลลาร์ ดังนั้น ทรานสโพสของมันจึงมีค่าเท่าเดิม จึงได้ว่า

$$(Y'X\beta)' = \beta'X'Y]$$

จากนั้นจึงเริ่มหาตัวประมาณโดยทำให้กำลังสองของค่าความคลาดเคลื่อนที่ได้มานั้นมีค่า
น้อยที่สุดโดยใช้ การดิฟเฟอเรนเชียล

$$\begin{aligned}\frac{\partial}{\partial \beta} (\epsilon'\epsilon) &= \frac{\partial}{\partial \beta} (Y'Y - 2\beta'X'Y + \beta'X'X\beta) = 0 \\ &= -2X'Y + 2X'X\beta = 0 \\ &X'X &= X'X\beta\end{aligned}$$

จากนี้ จะหาตัวประมาณของ β คือ $\hat{\beta}$

ซึ่งเป็นชุดสมการปกติ (normal equations) และถ้าหากเราจะนำมาเขียนในรูป

สเกลล่าจะได้ว่า

$$\begin{aligned} \Sigma Y_i &= b_0 n + b_1 \Sigma X_{1i} + b_2 \Sigma X_{2i} + \dots + b_k \Sigma X_{ki} \\ \Sigma X_{1i} Y_i &= b_0 \Sigma X_{1i} + b_1 \Sigma X_{1i}^2 + b_2 \Sigma X_{1i} X_{2i} + \dots + b_k \Sigma X_{1i} X_{ki} \\ \Sigma X_{2i} Y_i &= b_0 \Sigma X_{2i} + b_1 \Sigma X_{1i} X_{2i} + b_2 \Sigma X_{2i}^2 + \dots + b_k \Sigma X_{2i} X_{ki} \\ \vdots \\ \Sigma X_{ki} Y_i &= b_0 \Sigma X_{ki} + b_1 \Sigma X_{1i} X_{ki} + b_2 \Sigma X_{2i} X_{ki} + \dots + b_k \Sigma X_{ki}^2 \end{aligned}$$

การจะหา β ได้ต้องมีเงื่อนไขว่าชุดสมการปกติทั้งหมด $k+1$ สมการนี้ต้องเป็นอิสระซึ่งกันและกันโดยที่ $X'X$ จะต้องเป็นเมตริกซ์ที่ไม่เอกเทศ เพราะจะทำให้สามารถหา $(X'X)^{-1}$ ได้

$$\beta = (X'X)^{-1} X'Y$$

และเมื่อเราทำการดิฟเฟอเรนเชียลอันดับที่ 2 ผลปรากฏว่า

$$\frac{\partial}{\partial \beta^2} (\epsilon'\epsilon) = 2X'X$$

ซึ่งเป็นเทอริททางบวกเมื่อ X มี rank เต็มตามจำนวนคอสมันคือ $k+1$

เพราะฉะนั้น $\beta = (X'X)^{-1} X'Y$ จึงเป็นตัวประมาณที่ได้จากการทำให้ค่าฟังก์ชันของความคลาดเคลื่อนมีค่าน้อยที่สุด

$$\text{จึงได้ว่า } \hat{Y} = X\beta$$

ต่อไปเราจะพิจารณาค่าความแปรปรวนของ β

$$\begin{aligned} \text{จาก } \beta &= (X'X)^{-1} X'Y \\ &= (X'X)^{-1} X'(X\beta + \epsilon) \\ &= \beta + (X'X)^{-1} X'\epsilon \end{aligned}$$

$$\beta - \beta = (X'X)^{-1} X'\epsilon$$

$$\begin{aligned} E(\beta - \beta)(\beta - \beta)' &= E[(X'X)^{-1} X'\epsilon] [(X'X)^{-1} X'\epsilon]' \\ &= E[(X'X)^{-1} X'\epsilon \epsilon' X (X'X)^{-1}] \\ &= (X'X)^{-1} X' E(\epsilon \epsilon') X (X'X)^{-1} \\ &= \sigma^2 (X'X)^{-1} \quad (\because E(\epsilon \epsilon') = \sigma^2 I) \end{aligned}$$

$$\text{เราจึงได้ค่า } \text{var}(\beta) = \sigma^2 (X'X)^{-1}$$

จากข้างต้นจึงได้ว่า $E(\hat{b} - \beta) (\hat{b} - \beta)'$ เป็น VAR - COV ของ \hat{b} แต่
เนื่องจากไม่ทราบค่าของ σ^2 จึงต้องประมาณค่าของ σ^2 ก่อน

$$\text{จาก } Y = X\beta + \epsilon$$

$$\text{และ } \hat{Y} = X\hat{b}$$

$$\text{ให้ } \epsilon = Y - \hat{Y}$$

$$= Y - X\hat{b}$$

$$= Y - X(X'X)^{-1}X'Y$$

$$= (I - X(X'X)^{-1}X')Y$$

$$= MY \quad (\text{ให้ } M = I - X(X'X)^{-1}X')$$

$$= M(X\beta + \epsilon)$$

$$= MX\beta + M\epsilon$$

$$= M\epsilon$$

$$\text{ทั้งนี้เนื่องจาก } MX = 0$$

$$\text{โดย } MX = (I - X(X'X)^{-1}X')X$$

$$= X - X(X'X)^{-1}X'X$$

$$= 0$$

$$\text{หา } e'e = \epsilon'MM\epsilon$$

$$= \epsilon'M\epsilon$$

$$\text{ทั้งนี้เนื่องจาก } MM = (I - X(X'X)^{-1}X')(I - X(X'X)^{-1}X')$$

$$= I - X(X'X)^{-1}X' - X(X'X)^{-1}X' + X(X'X)^{-1}X'X(X'X)^{-1}X'$$

$$= I - X(X'X)^{-1}X'$$

$$= M$$

$$\therefore e'e = \text{tr} \epsilon'M\epsilon$$

ทั้งนี้เนื่องจาก M เป็นสเกลาร์จึงเท่ากับ Trace ของตัวมันเอง

$$e'e = \text{tr} M\epsilon\epsilon'$$

$$\therefore \text{tr}(AB) = \text{tr}(BA)$$

$$\therefore E(e'e) = E(\text{tr} M\epsilon\epsilon')$$

$$= \text{tr} M \cdot E(\epsilon\epsilon')$$



เพราะว่า trace เป็นฟังก์ชันเชิงเส้น

$$\begin{aligned} E(e'e) &= \sigma^2 \text{tr } M & (E(\xi\xi') &= \sigma^2 I) \\ &= \sigma^2 (n-k-1) \end{aligned}$$

$$\begin{aligned} \text{ทั้งนี้เนื่องจาก } \text{tr } M &= \text{tr}(I - X(X'X)^{-1}X') \\ &= \text{tr}(I_n) - \text{tr}(X(X'X)^{-1}X') \\ &= n - \text{tr}X'X(X'X)^{-1} & (\text{tr}(AB) = \text{tr}(BA)) \\ &= n - \text{tr}(I_{k+1}) \\ &= n-k-1 \end{aligned}$$

$$\begin{aligned} \text{ดังนั้นจึงได้ว่า } S^2 &= \frac{e'e}{n-k-1} \\ &= \frac{Y'MY}{n-k-1} & \text{ซึ่งเป็นตัวประมาณที่ไม่เอนเอียงของ } \sigma^2 \end{aligned}$$

และ $S^2(X'X)^{-1}$ ก็เป็นตัวประมาณที่ไม่เอนเอียงของ $\sigma^2(X'X)^{-1}$

2.3.7.3 คุณสมบัติของตัวประมาณโดยวิธีกำลังสองน้อยที่สุด

จากข้อสมมติที่กำหนดไว้ในหัวข้อข้างต้น ถ้าหากเป็นจริงแล้วเราจะ

ได้ว่าการหาค่าตัวประมาณโดยวิธีกำลังสองน้อยที่สุดจะมีคุณสมบัติดังนี้คือ

$$1. \quad b = (X'X)^{-1}X'Y \quad \text{เป็นตัวประมาณที่ไม่เอนเอียงของ } \beta$$

พิสูจน์

$$\begin{aligned} b &= (X'X)^{-1}X'Y \\ &= (X'X)^{-1}X'(X\beta + \xi) \\ &= (X'X)^{-1}X'X\beta + (X'X)^{-1}X'\xi \\ &= \beta + (X'X)^{-1}X'\xi \end{aligned}$$

$$\begin{aligned} E(b) &= E(\beta) + (X'X)^{-1}X'E(\xi) \\ &= \beta \end{aligned}$$

นั่นคือ จะได้ว่า $b = (X'X)^{-1}X'Y$ เป็นตัวประมาณที่ไม่เอนเอียงของ β

2. $\hat{\beta} = (X'X)^{-1}X'Y$ จะเป็นตัวประมาณค่าไม่เอนเอียงเชิงเส้นที่ดีที่สุด (Best Linear Unbiased Estimator) ในความหมายคือตัวประมาณอื่น ๆ ของ β ซึ่งเป็นเชิงเส้นกับ Y และไม่เอนเอียงเช่นกัน จะมีค่าความแปรปรวน-โควาเรียนซ์เมตริก ซึ่งมีค่าเกินกว่าค่าความแปรปรวน-โควาเรียนซ์ของตัวประมาณ $\hat{\beta}$ หรืออาจกล่าวได้อย่างง่าย ๆ ก็คือในบรรดาตัวประมาณในจำพวกของตัวประมาณค่าไม่เอนเอียงเชิงเส้นด้วยกันแล้วตัวประมาณโดยวิธีกำลังสองน้อยที่สุดนั้นจะมีค่าความแปรปรวนต่ำสุด โดยทฤษฎีของกอลด์ส-มาคอฟฟ์และเป็นที่ยอมรับกันอย่างแพร่หลาย

จากการที่ตัวประมาณของ β เป็นเชิงเส้นกับ Y เราจึงสามารถเขียนตัวประมาณนั้นได้เป็น A^*Y โดย A^* เป็นเมตริกซ์ $(k+1) \times n$ ซึ่งไม่ขึ้นอยู่กับค่าของ Y

$$\text{ให้ } A^* = (X'X)^{-1}X' + A$$

$$\begin{aligned} \therefore A^*Y &= [(X'X)^{-1}X' + A] Y \\ &= [(X'X)^{-1}X' + A] (X\beta + \epsilon) \end{aligned}$$

$$A^*Y = [I + AX]\beta + [(X'X)^{-1}X' + A]\epsilon$$

$$\begin{aligned} E(A^*Y) &= (I + AX)\beta + 0 \quad (\text{จากข้อสมมติที่ 2}) \\ &= \beta + AX\beta \end{aligned}$$

ตัวประมาณ A^*Y จะเป็นตัวประมาณค่าไม่เอนเอียง ถ้า $AX = 0$

$$\therefore \text{จาก } A^*Y = \beta + [(X'X)^{-1}X' + A]\epsilon$$

$$A^*Y - \beta = [(X'X)^{-1}X' + A]\epsilon$$

$$(A^*Y - \beta)(A^*Y - \beta)' = [(X'X)^{-1}X' + A]\epsilon\epsilon'[(X'X)^{-1}X' + A]'$$

$$\begin{aligned} E(A^*Y - \beta)(A^*Y - \beta)' &= [(X'X)^{-1}X' + A] E(\epsilon\epsilon') [X(X'X)^{-1} + A'] \\ &= \sigma^2 [(X'X)^{-1}X' + A] [X(X'X)^{-1} + A'] \\ &= \sigma^2 [(X'X)^{-1} + AX(X'X)^{-1} + (X'X)^{-1}X'A' + AA'] \\ &= \sigma^2 [(X'X)^{-1} + AA'] \end{aligned}$$

เนื่องจากเรากำหนดให้ $AX = 0 \rightarrow X'A' = 0$ ด้วย

$$\begin{aligned} \text{จึงได้ว่า } \text{Var}(A^*Y) &= \sigma^2 (X'X)^{-1} + \sigma^2 AA' \\ &= \text{Var}(\hat{\beta}) + \sigma^2 AA' \end{aligned}$$

นั่นคือ ค่าความแปรปรวน-โควาเรียนซ์ของตัวประมาณจะมีค่าเกินกว่าค่าความแปรปรวน-โควาเรียนซ์ของตัวประมาณ ซึ่งประมาณโดยวิธีกำลังสองน้อยที่สุด

$\rightarrow \hat{\beta} = (X'X)^{-1}X'Y$ จะเป็นตัวประมาณที่ไม่เอนเอียงเชิงเส้นที่ดีที่สุด

2.3.7.4 วิธีการทางสถิติที่ใช้ในการเลือกสมการถดถอยที่ดีที่สุด¹

ในการเลือกสมการถดถอยเพื่อหาตัวแบบนั้น มีหลักเกณฑ์การเลือกที่เป็นหลักหรือข้อคิดที่สำคัญสำหรับการพิจารณาของนักวิจัยคือ ประการแรกต้องการให้ตัวแบบของสมการที่เลือกหรือที่สร้างขึ้นมาั้นมีประสิทธิภาพสูงและเป็นประโยชน์ในการพยากรณ์มากที่สุด ซึ่งในสมการนั้นจะประกอบด้วยตัวแปรตาม ซึ่งขึ้นอยู่กับตัวแปรอิสระมากที่สุด ดังตัวแบบเช่น

$$Y = f(X_1, X_2, \dots, X_k)$$

เพื่อว่าค่าของตัวแปรตามที่คำนวณได้จะเป็นค่าพยากรณ์ที่ใกล้เคียงและอีกประการหนึ่งเนื่องจากนักวิจัยทั้งหลายมักประสบกับปัญหาเกี่ยวกับเวลาและค่าใช้จ่ายในการวิเคราะห์ข้อมูล ยิ่งถ้าหากว่าการวิจัยนั้นใช้ตัวแปรอิสระมากเท่าไร (เพื่อที่จะให้มีประสิทธิภาพสูง) ก็ยิ่งจะทำให้ต้องสิ้นเปลืองค่าใช้จ่ายและเวลาในการคำนวณมากยิ่งขึ้นเท่านั้น

จากเหตุผลดังกล่าวมาแล้วข้างต้นจะเห็นได้ว่ามีความขัดแย้งในตัวเอง เพื่อที่จะเป็นการหาตัวแบบสมการถดถอยที่ดีที่สุด คือให้ได้ตัวแปรอิสระที่มีความเหมาะสมจำนวนน้อยแต่ให้สมการนั้นมีประสิทธิภาพสูงมีความเชื่อถือได้ จึงได้นำเอาวิธีการทางสถิติมาช่วยในการเลือกสมการถดถอยซึ่งกระทำได้หลายวิธี ในแต่ละวิธีนั้นก็จะมีเหตุผลในการเลือกตามความต้องการของแต่ละบุคคล และทั้งนี้ก็ยังขึ้นอยู่กับงบประมาณและความสะดวกด้วย ในที่นี้จะขอกล่าวถึงวิธีการทางสถิติในการเลือกสมการถดถอย ซึ่งมีอยู่หลายวิธีด้วยกันคือ

1. All Possible Regression
2. The Backward Elimination Procedure
3. The Forward Selection Procedure
4. The Stepwise Regression Procedure
5. Two Variations on the Four Previous Methods
6. The Stagewise Regression Procedure

¹ ดูรายละเอียดใน N.R. Draper and H. Smith, Applied Regression Analysis (New York: John Wiley & Sons, Inc., 1966), p 163-177.

วิธีการของการเลือกสมการถดถอยเพื่อจะดูถึงตัวแปรอิสระนี้เป็นเทคนิคทางสถิติวิธีหนึ่งที่ใช้ในการวิเคราะห์ความสัมพันธ์ระหว่างตัวแปรตาม (dependent Variable Y) กับกลุ่มของตัวแปรอิสระ (independent variables X_1, X_2, \dots, X_k) เพื่อที่จะเลือกว่าตัวแปรอิสระใดมีความสำคัญในขอบเขตที่มีนัยสำคัญ

ในการวิจัยครั้งนี้ เนื่องจากผู้วิจัยต้องการที่จะพิจารณหาตัวแบบของสมการถดถอยที่เหมาะสม เพื่อที่จะดูว่าตัวแปรอิสระอะไรบ้างที่มีความสัมพันธ์ที่จะอธิบายถึงสัมฤทธิ์ผลทางการเรียนของนักศึกษา จากการศึกษาเบื้องต้นได้ว่าตัวแบบอิสระ 62 ตัว ที่มีความสัมพันธ์กับสัมฤทธิ์ผลทางการเรียนโดยใช้ χ^2 ทำการทดสอบ ได้แก่ตัวแปรดังนี้

X_1	คือ	เพศ
X_2	คือ	วิทยาลัยที่ศึกษาอยู่
X_4	คือ	G.P.A. ระดับอุดมศึกษา
X_7	คือ	เกรดเฉลี่ยสะสมระดับมัธยมปลาย
X_{20}	คือ	ไม่สนใจเนื้อหาที่เรียน
X_{24}	คือ	ไม่ศึกษาค้นคว้าจากหนังสืออื่น ๆ ที่เกี่ยวข้อง
X_{33}	คือ	รู้สึกว่าคณะวิชาที่ศึกษาต่อยกกว่าคณะวิชาอื่น ๆ
X_{34}	คือ	รู้สึกว่าวิชาที่ศึกษายากเกินไปสำหรับความรู้ของทำงาน
X_{35}	คือ	ข้อสอบยาก
X_{41}	คือ	ท่านไม่ชอบทำกิจกรรม
X_{42}	คือ	ในขณะที่ศึกษาอยู่ท่านมักจดจำบรรยายไม่ทัน
X_{53}	คือ	ต้องช่วยหางบ้านในการประกอบอาชีพ
X_{55}	คือ	เกิดปัญหาส่วนตัวซึ่งมีผลกระทบต่อ การเรียน
X_{58}	คือ	การตัดครุเข้าชั้นเรียน ไม่เหมาะสมกับความรู้และความถนัดในการสอน

วิธีการถดถอยแบบเป็นขั้นตอน (The Stepwise Regression Procedure) เป็นวิธี ซึ่งผู้วิจัยได้ใช้ SPSS มาใช้ในการเลือกตัวแปรอีกครั้งหนึ่ง เพื่อจะได้ตัวแบบของสมการถดถอยครั้งนี้ เนื่องจากเป็นวิธีที่รัดกุมและแก้ไขจุดบกพร่องของวิธีอื่น กล่าวคือในการตั้งตัวแปรใหม่เข้ามาอยู่ในสมการถดถอยแต่ละครั้งนั้น จะมีการตรวจสอบดูว่าตัวแปรที่ตั้งเข้ามาใหม่นี้จะมีอิทธิพลต่อตัวแปรที่อยู่ในสมการก่อนแล้วหรือไม่เพราะตัวแปรบางตัวที่เราเลือกเข้ามาอยู่ในตัวแบบในตอนแรกนั้นอาจไม่จำเป็น ถ้า

เราดูความสัมพันธ์ของมันกับตัวแปรที่เราดึงเข้ามาใหม่ ในการตรวจสอบดังกล่าวนี้จะใช้หลักของการทดสอบค่า F เพียงบางส่วน (Partial F-test) โดยถือว่าตัวแปรแต่ละตัวนั้นถูกใส่เข้าไปในตัวแบบเป็นตัวสุดท้าย ตัวแปรใดที่ไม่มีนัยสำคัญทางสถิติ (non-significance) ให้ตัดออกจากตัวแบบ คงไว้เฉพาะตัวที่มีนัยสำคัญ (Significance) วิธีการเช่นนี้จะหยุดต่อเมื่อเราไม่สามารถเพิ่มหรือลดตัวแปรใด ๆ ในตัวแบบได้อีก ซึ่งในที่นี้ผู้วิจัยจะขอกกล่าวถึงขั้นตอนการเลือกตัวแปรโดยวิธีการถดถอยแบบเป็นขั้นตอนเพื่อให้ได้ตัวแบบที่ดีที่สุดโดยละเอียดดังนี้คือ

ขั้นที่ 1 คำนวณหาค่าของสัมประสิทธิ์สหสัมพันธ์อย่างง่าย (Simple correlation coefficient) ระหว่างตัวแปรอิสระทั้งหมด X กับตัวแปรตาม Y โดยใช้สัญลักษณ์เป็น $r_{x,y}$ โดยที่

$$r_{x,y} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2 \sum_{i=1}^n (Y_i - \bar{Y})^2}}$$

$$= \frac{\sum_{i=1}^n X_i Y_i - \frac{(\sum_{i=1}^n X_i)(\sum_{i=1}^n Y_i)}{n}}{\sqrt{\left(\sum_{i=1}^n X_i^2 - \frac{(\sum_{i=1}^n X_i)^2}{n}\right) \left(\sum_{i=1}^n Y_i^2 - \frac{(\sum_{i=1}^n Y_i)^2}{n}\right)}}$$

สัมประสิทธิ์สหสัมพันธ์นี้เป็นตัวที่ใช้ในการพิจารณาว่าตัวแปร 2 ชุดนั้นมีความสัมพันธ์กันหรือไม่อย่างไร โดยที่ถ้าตัวแปรชุดแรกแทนด้วย random variable X และตัวแปรชุดที่สองแทนด้วย random variable Y ค่าของความสัมพันธ์นี้จะอยู่ระหว่าง -1 ถึง 1 นั่นคือ

$$-1 \leq r_{xy} \leq 1$$

ถ้า $-1 \leq r_{x,y} < 0$ จะได้ว่า X กับ Y มีความสัมพันธ์กันในทางตรงข้าม นั่นคือถ้าหากค่าของ X มากขึ้นแล้วค่าของ Y จะลดลง หรือถ้าหากค่าของ X ลดลงแล้ว ค่าของ Y จะมากขึ้น

ถ้า $r_{X,Y} = 0$ จะได้ว่า X กับ Y นั้นไม่มีความสัมพันธ์กันเลย

ถ้า $0 < r_{X,Y} \leq 1$ จะได้ว่า X กับ Y มีความสัมพันธ์กันในทางเดียวกันนั่นคือถ้าหากค่าของ X มากขึ้นแล้ว ค่าของ Y จะมากขึ้นด้วย หรือถ้าหากค่าของ X ลดลงแล้ว ค่าของ Y จะลดลงด้วย

ซึ่งในการคำนวณนี้เนื่องจากเรามีตัวแปรอิสระหลายตัว ดังนั้นถ้าหากจะดูในรูปของเมตริกซ์แล้วจะสะดวกกว่าและได้เป็น

$$\begin{bmatrix} 1 & r_{X_1 Y} & r_{X_2 Y} & \dots & r_{X_k Y} \\ r_{Y X_1} & 1 & r_{X_2 X_1} & \dots & r_{X_k X_1} \\ r_{Y X_2} & r_{X_1 X_2} & 1 & \dots & r_{X_k X_2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ r_{Y X_k} & r_{X_1 X_k} & r_{X_2 X_k} & \dots & 1 \end{bmatrix}$$

(เมื่อ k คือจำนวนตัวแปรอิสระที่เรานำมาพิจารณา)

จากนี้จะดูว่าตัวแปรอิสระใดที่ให้ค่าสัมประสิทธิ์สัมพันธ์กับตัวแปรตามสูงสุดก็จะได้รับเลือกเข้าสมการถดถอยก่อนในที่นี้สมมติว่า $r_{X_1 Y}$ มีค่าสูงสุด ดังนั้น X_1 จึงเป็นตัวแปรอิสระที่จะเลือกมาพิจารณาส่งสมการก่อนได้สมการถดถอยเป็น

$$\hat{Y} = b_0 + b_1 X_1$$

จากสมการนี้จะนำมาทดสอบดูว่าตัวแปรอิสระ X_1 ซึ่งนำเข้ามาอยู่ในสมการนั้นจะมีส่วนช่วยในการอธิบายค่าของตัวแปรตาม อย่างมีนัยสำคัญทางสถิติหรือไม่ โดยการใช้ F-test ดูค่า F เทียบกับ F จากตารางในระดับนัยสำคัญที่กำหนดไว้ | ในที่นี้เทียบกับ $F(1, n-1, 1-\alpha)$ ในตารางโดย n คือ จำนวนข้อมูล, α คือ ระดับนัยสำคัญที่กำหนด | ถ้าค่า $F \geq F$ จากตาราง แสดงว่า X_1 ยังคงอยู่ในสมการ แต่ถ้าหากว่า

Sequential $F \leq F$ จากตารางแสดงว่า Y ไม่มีความสัมพันธ์กับตัวแปรอิสระใดอย่างมีนัยสำคัญเลย ก็จะหยุดอยู่เพียงแค่นี้

ในการทดสอบนี้ จะเห็นได้ว่า ในการทดสอบเรากล่าวว่าทดสอบ ซึ่งอันที่จริงแล้วคือ การทดสอบ F ของ Regression Coefficient ของตัวแปรอิสระที่กำลังพิจารณาอยู่นั้นเอง อย่างไรก็ตามในการนี้สมมติว่า X_1 คือตัวแปรอิสระที่ถูกเลือกเข้าในสมการ ดังนั้นจึงทดสอบสมมติฐาน

$$H_0 : \beta_1 = 0$$

ขั้นที่ 2 คำนวณหาค่าของสัมประสิทธิ์สหสัมพันธ์เพียงบางส่วน (Partial Correlation Coefficient) ระหว่างตัวแปรอิสระที่เหลือ X_j ($j \neq 1$) กับตัวแปรตาม Y โดยกำหนดให้ตัวแปรอิสระ X_1 ซึ่งเราเลือกไว้ในสมการถดถอยแล้วคงที่และใช้สัญลักษณ์เป็น $r_{X_j Y \cdot X_1}$ โดยที่

$$r_{X_j Y \cdot X_1} = \frac{r_{X_j Y} - r_{X_j X_1} \cdot r_{Y X_1}}{\sqrt{(1 - r_{X_j X_1}^2)(1 - r_{Y X_1}^2)}}$$

ตัวแปรอิสระที่เหลือ X_j ตัวใดที่ให้ค่าสัมประสิทธิ์สหสัมพันธ์เพียงบางส่วนกับตัวแปรตาม Y สูงสุด จะได้รับเลือกเข้าสมการต่อไป สมมติว่าในที่นี้ $r_{X_2 Y \cdot X_1}$ ให้ค่าสูงสุด ดังนั้น X_2 จะเป็นตัวแปรอิสระที่เลือกมาพิจารณาและสมการถดถอยจะเป็น

$$\hat{Y} = b_0 + b_1 X_1 + b_2 X_2$$

จากสมการนี้ จะมีการทดสอบ F ของ X_2 แต่ถ้า X_2 ไม่สำคัญพอ ก็จะ

ตัด X_2 ออกเหลือ $\hat{Y} = b_0 + b_1 X_1$ เช่นเดิมและหยุดอยู่เพียงแค่นี้ แต่ถ้า X_2 สำคัญแล้ว แสดงว่าต้องคงสมการ $\hat{Y} = b_0 + b_1 X_1 + b_2 X_2$

ขั้นที่ 3 เราจะทดสอบ partial F ของตัวแปรอิสระที่เลือกเข้าสมการก่อนหน้านี้ ในที่นี้คือ X_1 นั่นคือ ดูว่าหลังจากที่เลือก X_2 เข้ามาเพิ่มในสมการแล้ว X_1 ซึ่งเป็นตัวแปรอิสระ

ซึ่งเลือกเข้ามาในสมการ ก่อนจะยังมีความสำคัญอยู่หรือไม่ โดยจะทดสอบ Partial F ของ X_1 เมื่อ X_2 อยู่ในสมการเรียบร้อยแล้วเทียบกับ F จากตาราง $[F(1/n-3, 1-\alpha)]$ ถ้าหากว่าค่า partial $F \leq F$ จากตารางแล้วจะได้ว่า X_1 ซึ่งอยู่ในสมการตอนแรกนั้นจะ ไม่มีความสำคัญอีกต่อไป ดังนั้นจะตัด X_1 ออกสมการที่ได้จะเปลี่ยนใหม่เป็น $\hat{Y} = b_0 + b_2 X_2$ แต่ถ้า partial $F > F$ จากตารางแล้วแสดงว่า สมการ $\hat{Y} = b_0 + b_1 X_1 + b_2 X_2$ จะยังคงใช้ได้ และเราจะเริ่มหาตัวแปรอิสระตัวใหม่ต่อไป

จากขั้นตอนนี้จะพบว่ามีการทดสอบค่า Partial F ซึ่งในที่นี้ก็คือการทดสอบดูว่า ตัวแปรอิสระที่เข้าไปอยู่ในสมการแล้ว ถ้าให้เข้าเป็นตัวสุดท้ายโดยให้ตัวแปรอิสระตัวอื่นที่เลือก แล้วอยู่ในสมการก่อนจะสามารถอธิบายค่า Y ได้อีกมากน้อยเพียงใด

ขั้นที่ 4 จากตอนนี้เราได้แล้วว่า X_1, X_2 เป็นตัวแปรอิสระที่เลือกอยู่ในสมการ แล้ววิธีการต่อไปคือ จะย้อนกลับไปทำเช่นเดียวกับขั้นที่ 2 อีกคือหาค่าของสัมประสิทธิ์สหสัมพันธ์ เพียงบางส่วนระหว่างตัวแปรอิสระที่เหลือ X_j ($j \neq 1, 2$) กับตัวแปรตาม Y โดยกำหนดให้ตัวแปรอิสระ X_1, X_2 ซึ่งเราเลือกไว้แล้วในสมการถดถอยนั้นคงที่และใช้สัญลักษณ์เป็น $r_{X_j Y \cdot X_1 X_2}$ โดยที่

$$r_{X_j Y \cdot X_1 X_2} = \frac{r_{X_j Y X_2} - r_{X_j X_1 \cdot X_2} r_{Y X_1 \cdot X_2}}{\sqrt{(1 - r_{X_j X_1 \cdot X_2}^2)(1 - r_{Y X_1 \cdot X_2}^2)}}$$

หรือ

$$= \frac{r_{X_j Y X_1} - r_{X_j X_2 \cdot X_1} r_{Y X_2 \cdot X_1}}{\sqrt{(1 - r_{X_j X_2 \cdot X_1}^2)(1 - r_{Y X_2 \cdot X_1}^2)}}$$

ตัวแปรอิสระใดที่หาค่าของสัมประสิทธิ์สหสัมพันธ์เพียงบางส่วนสูงสุดก็จะได้รับการเลือกเป็น ตัวแปรอิสระเข้าสมการตัวต่อไปและจะทำการทดสอบเช่นเดียวกับที่ผ่านมาจนกว่าจะไม่มีตัวแปรอิสระใด ที่เป็นที่ยอมรับคือไม่ผ่านการทดสอบ F test ก็จะได้สมการอันสุดท้ายเป็นอันที่ต้องการ