

บทที่ 1

บทนำ



1.1 ความเป็นมาและความสำคัญของปัญหา

ปัญหาประการหนึ่งในการวิเคราะห์ข้อมูลในทางสถิติคือ การที่ข้อมูลหรือค่าสังเกตที่เก็บรวบรวมได้ไม่เป็นไปตามสภาวะการณที่กำหนดหรือควบคุมไว้ เช่น การทดลองทางการแพทย์ ด้านชีววิทยา เป็นต้น ทำให้ข้อมูลหรือค่าสังเกตบางค่าสูงหรือต่ำมาก หรือเป็นค่าสังเกตที่ไม่ได้มาจากประชากรเดียวกับค่าสังเกตส่วนใหญ่ โดยกรณีหลังนี้จะพบมากในข้อมูลที่ได้จากการวางแผนการทดลองผิดพลาด ข้อมูลที่มีลักษณะดังกล่าวจะเรียกว่าเป็น "ค่าผิดปกติ" (outliers)

ในปี ค.ศ. 1960 อังโคมบ์ (Ancombe) กล่าวว่าสาเหตุของข้อมูลผิดปกติมี 3 ประการ ได้แก่

1. ความผันแปรที่มีอยู่ในประชากรที่ศึกษา (inherent variability) เป็นความผันแปรที่ไม่สามารถหลีกเลี่ยงได้ แม้จะมีการควบคุมหรือการวัดหรือการปฏิบัติเป็นอย่างดีก็ยังมีอยู่ และแก้ไขไม่ได้ นอกจากจะเปลี่ยนประชากรหรือวัตถุประสงค์ในการศึกษา
2. ความคลาดเคลื่อนที่เกิดจากการวัด (measurement error) ความผิดพลาดชนิดนี้เกิดจากการบันทึกข้อมูลหรือเครื่องมือเครื่องใช้ในการวัดมีคุณภาพต่ำ ความคลาดเคลื่อนนี้อาจแก้ไขให้หมดไปได้
3. ความคลาดเคลื่อนที่เกิดจากการปฏิบัติการ (execution error) เช่น การลงรหัส การเจาะบัตร เป็นต้น ความผิดพลาดชนิดนี้สามารถลดลงได้ด้วยการระมัดระวัง ป้องกันไว้ก่อน

ในปี ค.ศ. 1980 ฮอกกินส์ (Hawkins) ให้คำจำกัดความของค่าผิดปกติว่า "เป็นค่าสังเกตที่เบี่ยงเบนไปจากค่าสังเกตอื่นอย่างมากจนทำให้สงสัยว่าค่าสังเกตนั้นได้มาจากวิธีการ

(mechanism) อื่น" และในปี 1983 เบกแมนและคูก (Beckman and Cook) ได้สรุปความหมายของ คำผิดปรกติออกเป็น 2 ลักษณะคือ

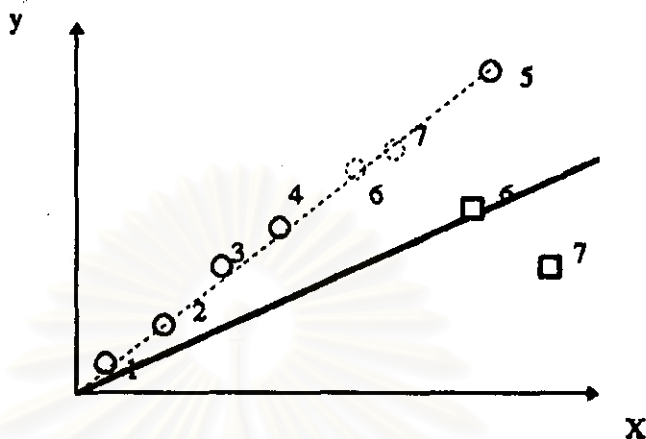
ก) ค่าสังเกตที่มีค่าสูงหรือต่ำมาก (extreme) หรือเป็นค่าที่เบี่ยงเบนไปจากกลุ่มที่ศึกษา เรียกคำผิดปรกตินี้ว่า discordant observation

ข) ค่าสังเกตที่มีลักษณะการแจกแจงแตกต่างจากลักษณะการแจกแจงของประชากรที่สนใจศึกษาและเรียกคำผิดปรกตินี้ว่า contaminate observation

การวิเคราะห์ความถดถอยเป็นวิธีการวิเคราะห์ข้อมูลทางสถิติที่มีการนำไปใช้ในงานวิจัยทางด้านสังคมศาสตร์และวิทยาศาสตร์ โดยการนำข้อมูลที่สังเกตหรือเก็บรวบรวมได้มาใช้ประมาณหรือพยากรณ์สิ่งที่สนใจ ข้อมูลที่ใช้ในการวิเคราะห์ความถดถอยบางครั้งอาจมีคำผิดปรกติปนอยู่ ทำให้ผลการวิเคราะห์ความถดถอยเกิดความผิดพลาดขึ้นได้ ทั้งนี้เพราะคำผิดปรกติดังกล่าวมีผลต่อการประมาณค่าสัมประสิทธิ์การถดถอย (β) ซึ่งทำให้สมการถดถอยที่ได้เบี่ยงเบนไปในทิศทางหรือตำแหน่งของคำผิดปรกติ การเบี่ยงเบนไปในทิศทางหรือตำแหน่งของคำผิดปรกติดังกล่าวทำให้เกิดปัญหาในการตรวจสอบคำผิดปรกติเกิดขึ้น 2 ชนิด คือ การตรวจสอบคำผิดปรกตินั้นไม่พบ และการตรวจสอบพบคำผิดปรกติแต่คำผิดปรกติที่ตรวจสอบไม่ใช่คำผิดปรกติจริงในข้อมูล กล่าวคือ หลังจากที่ทำกรวิเคราะห์ความถดถอยด้วยวิธีกำลังสองน้อยสุดแล้วค่าตัวแปรตามที่พยากรณ์ได้จะมีความคลาดเคลื่อนมากน้อยแตกต่างกัน คำผิดปรกติในการวิเคราะห์ความถดถอย หมายถึง ค่าสังเกตที่มีความคลาดเคลื่อนสูง แต่เนื่องจากคำผิดปรกติที่มีในข้อมูลดึงสมการการถดถอยเบี่ยงเบนไปในทิศทางหรือตำแหน่งของมัน ก็จะส่งผลให้ความคลาดเคลื่อนของมันต่ำ ทำให้ตรวจสอบไม่พบคำผิดปรกติดังกล่าวหรืออาจส่งผลให้ค่าสังเกตที่ดีหรือสะอาด* (good or clean observation) มีความคลาดเคลื่อนสูงจนกลายเป็นคำผิดปรกติได้เมื่อตรวจสอบด้วยวิธีพื้นฐาน

* ข้อมูลดีหรือสะอาด หมายถึง ข้อมูลนั้นเป็นข้อมูลปรกติ

ปัญหาที่เกิดขึ้น เรียกว่า มาซคกิงเอฟเฟ็ค (masking effect) และ ซวอมพิงเอฟเฟ็ค (swamping effect) ตามลำดับ ลักษณะการเกิดมาซคกิงและซวอมพิงเอฟเฟ็คแสดงในรูปที่ 1.1



รูปที่ 1.1 : แสดงลักษณะของการเกิดมาซคกิงเอฟเฟ็ค ณ ตำแหน่งที่ 6 และ ซวอมพิงเอฟเฟ็ค ณ ตำแหน่งที่ 5

จากรูปที่ 1.1 ข้อมูลที่เรานำมาวิเคราะห์มีความสัมพันธ์เชิงเส้นตรงด้วยเส้นประ แต่ถ้าค่าที่ 6 และ 7 เบี่ยงเบนจากข้อมูลโดยส่วนใหญ่เราจะได้ความสัมพันธ์เชิงเส้นตรงด้วยเส้นทึบ^{***} เมื่อเราทำการตรวจสอบค่าผิดปกติด้วยวิธีพื้นฐาน คือ พิจารณาค่าที่มีความคลาดเคลื่อนสูงเป็นค่าผิดปกติและใช้ทุกค่าสังเกตในการวิเคราะห์ เราจะพบว่าค่าที่ 6 จะมีความคลาดเคลื่อนต่ำทำให้ตรวจสอบค่าผิดปกติดังกล่าวไม่พบ แต่ค่าที่ 5 จะมีความคลาดเคลื่อนสูงทำให้ค่าที่ 5 กลายเป็นค่าผิดปกติที่ถูกตรวจพบ ดังนั้นการตรวจสอบค่าผิดปกติจำเป็นต้องหาวิธีที่มีความแกร่ง (robust) พอที่จะไม่ทำให้เกิดปัญหาทั้ง 2 ชนิด หรือยอมให้เกิดได้น้อยที่สุด เพื่อจะได้ทำให้พบค่าผิดปกติจริงๆ และจะได้แก้ปัญหของข้อมูลมีค่าผิดปกติก่อนที่จะนำไปใช้ในการวิเคราะห์ความถดถอย เพื่อให้ได้ผลการพยากรณ์ทางสถิติซึ่งถูกต้องต่อไป

- * เหตุการณ์ที่ค่าผิดปกติค่าหนึ่งมีผลต่อค่าผิดปกติอีกค่าหนึ่งทำให้ไม่สามารถตรวจพบค่าผิดปกตินั้น
- ** เหตุการณ์ที่ค่าผิดปกติมีผลต่อค่าสังเกตอื่นๆ ที่ไม่ใช่ค่าผิดปกติ ทำให้ค่านั้นกลายเป็นค่าผิดปกติซึ่งถูกตรวจพบ
- *** การเบี่ยงเบนของค่าที่ 6 และ 7 อาจทำให้เรามองความสัมพันธ์ระหว่าง x กับ y ในรูปโพลีโนเมียลได้ แต่การใช้ความถดถอยโพลีโนเมียลให้ผลการประมาณค่าสัมประสิทธิ์ความถดถอยไม่แม่นยำเท่าการถดถอยเชิงเส้นตรงเนื่องจากแต่ละพจน์ของตัวแปรอิสระมีความสัมพันธ์กันเอง และข้อมูลที่เราศึกษาก็มีความสัมพันธ์เชิงเส้นระหว่าง x กับ y จึงไม่เหมาะสมที่จะใช้การถดถอยโพลีโนเมียล

วิธีการตรวจสอบค่าผิดปกติในการวิเคราะห์ความถดถอยเชิงเส้นได้มีผู้ศึกษาไว้ดังนี้ ในปี ค.ศ. 1967 มิกคีย์ (Mickey) คัน (Dun) และ คลาสก (Clask) ได้ใช้วิธีการวิเคราะห์ความถดถอยแบบขั้นตอน (Stepwise Regression) และการเพิ่มตัวแปรหุ่น (Dummy Variable) เข้าไปในสมการการถดถอยเพื่อทำการแยกค่าผิดปกติ แต่วิธีนี้ไม่เหมาะสมในกรณีที่มีค่าผิดปกติมากกว่า 1 ค่า ในปี ค.ศ. 1985 เมอร์วิน จี มาราสิง (Merveyn G. Marasinghe) ได้เสนอวิธีการวิเคราะห์แบบหลายขั้นตอน (multistage procedure) และ ตัวสถิติ F_{α}^* เพื่อใช้ตรวจสอบค่าผิดปกติหลายค่าในการวิเคราะห์ความถดถอยเชิงเส้น ในปี ค.ศ. 1990 รูซโซ (Rousseue) และ แวน โซเมอร์เรน (van Zomeren) ได้นำตัวประมาณที่ได้จากทรงรีที่มีปริมาตรต่ำสุด (minimum volume ellipsoid (MVE)) มาหาระยะทางที่แกร่ง (robust distance) เพื่อใช้ในการตรวจสอบค่าผิดปกติแบบหลายตัวแปร (multivariate outliers) แต่การหาตัวประมาณ MVE นั้นจะต้องหาเซตย่อยที่มีจำนวนสมาชิก(ค่าสังเกต) อย่างน้อยครึ่งหนึ่งของข้อมูล ถ้าให้จำนวนสมาชิกอย่างน้อยครึ่งหนึ่งของข้อมูลเท่ากับ n และ ขนาดตัวอย่างเท่ากับ n จะต้องทำการคำนวณหาปริมาตรของทรงรีเท่ากับ $\binom{n}{k}$ ปริมาตรแล้วจึงจะเลือกทรงรีที่มีปริมาตรต่ำที่สุด ดังนั้นจึงทำให้เสียค่าใช้จ่ายและเวลาในการคำนวณสูงและบางครั้งการคำนวณอาจเป็นไปไม่ได้ ต่อมา ในปี ค.ศ. 1992 ฮาดี (Hadi) เสนอวิธีการหา MVE ซึ่งจะทำให้ได้ MVE เพียงค่าเดียว (unique) และยังใช้ได้กับเมทริกซ์ความแปรปรวนร่วมจะเป็นเมทริกซ์เอกฐาน (singular matrix) ก็ตาม ฮาดีได้แก้ปัญหาที่เมทริกซ์ความแปรปรวนร่วมเป็นเมทริกซ์เอกฐาน โดยการถ่วงน้ำหนักเวกเตอร์เจาะจง (eigenvectors) ด้วยค่าเจาะจงสูงสุด (maximum eigenvalues) ของเมทริกซ์ความแปรปรวนร่วม และในปี ค.ศ. 1994 ฮาดี (Hadi) ก็ได้แก้ปัญหาดังกล่าวอีกครั้ง โดยการรวมค่าสังเกตที่อยู่ดัด ๆ ไปเข้าไปในเซตย่อยเท่าที่จำเป็น จนกว่าเมทริกซ์ความแปรปรวนร่วมจะเป็นเมทริกซ์ไม่เอกฐาน (non-singular) ในปี ค.ศ. 1993 ฮาดี (Hadi) และไซมันอฟฟ์ (Simonoff) ได้เสนอวิธีการตรวจสอบค่าผิดปกติหลายค่า (multiple outliers) ในตัวแบบเชิงเส้นโดยอาศัยแนวคิดของฮาดี ในปี ค.ศ. 1992 โคนิฟาร์ด (Kianifard) และชวอดโล (Swallow) ได้เสนอวิธีการตรวจสอบค่าผิดปกติโดยใช้เกณฑ์ความคลาดเคลื่อนเวียนเกิด (recursive residual) และได้เสนอตัวสถิติทดสอบ 2 วิธี ได้แก่ วิธีเวียนเกิดโดยลำดับ (sequential recursive method) และวิธีเวียนเกิดดัดแปร (modified recursive method) ในปี พ.ศ. 2532 บุญสม ธรรมศิริพจน์ ได้ศึกษาวิธีตรวจสอบค่าผิดปกติในสมการการถดถอยเชิงพหุ โดยศึกษาเปรียบเทียบวิธีของ เคนนิส คูก (Dennis Cook, 1977), วิธีของแอนดรูและเพรตจีบอน (Andrew and Pregibon, 1978) และวิธีของ จี แบร์รี (G. Barre Weterill, 1986) ซึ่งศึกษาในกรณีที่การแจกแจงของความคลาดเคลื่อนมี 2 แบบคือ การแจกแจงปกติปลอมปนในตำแหน่งและการแจกแจงปกติปลอมปนในสเกลผลการศึกษานี้ปรากฏว่าวิธีของ จี แบร์รี และ วิธีของแอนดรูและเพรตจีบอน มีความสามารถในการควบคุมความผิดพลาดประเภทที่ 1 ได้ดี ในปี

พ.ศ. 2533 สมชาย รัตนเดิพนุตรณ์ ได้ศึกษาวิธีการตรวจสอบค่าผิดปกติในการวิเคราะห์ความถดถอยเชิงเส้น โดยการเปรียบเทียบวิธีของ จี บาร์รี (G. Barre Weterill, 1986) วิธีของทิตเจน, มัวร์และเบกแมน (Tietjen, Moore and Beckman, 1973) และวิธีของ เมอร์วิน จี มาราสิง (Mervyn G. Marasinghe, 1985) โดยศึกษาในกรณีที่มีการแจกแจงของความคลาดเคลื่อนมี 2 ลักษณะ คือ การแจกแจงแบบหางยาวกว่าปกติ ได้แก่ การแจกแจงปกติปลอมปนในสเกล การแจกแจงปกติปลอมปนในตำแหน่ง และการแจกแจงที่ ส่วนการแจกแจงแบบเบ้ขวา ได้แก่ การแจกแจงลอกนอร์มอล แกมมา และไวบูลต์ พบว่าวิธีของ จี บาร์รี และวิธีของเมอร์วิน จี มาราสิง ควบคุมความคลาดผิดพลาดประเภทที่ 1 ได้ดีใกล้เคียงกัน ส่วนวิธีของทิตเจน, มัวร์และเบกแมนควบคุมได้น้อย

เนื่องจากตัวสถิติทดสอบแต่ละตัวมีความสามารถในการตรวจสอบค่าผิดปกติได้ต่างกัน บางวิธีอาจทำให้เกิดมาซคคิงและซวอมฟิงเอฟเฟ็คได้เป็นผลให้ตรวจสอบค่าผิดปกติจริงไม่พบหรือพบไม่ครบทุกค่า ผู้วิจัยจึงสนใจที่จะเปรียบเทียบวิธีการตรวจสอบค่าผิดปกติในการวิเคราะห์ความถดถอย เชิงเส้นโดยเปรียบเทียบวิธีของ เมอร์วิน จี มาราสิง (Mervyn G. Marasinghe, 1985) วิธีของไคนิฟาร์ดและสวอลโล (Kianifard and Swallow, 1990) และวิธีของฮาโคและไซมันนอฟฟ์ (Hadi and Simonoff, 1993)

1.2 วัตถุประสงค์

เพื่อศึกษาเปรียบเทียบตัวสถิติที่ใช้ตรวจสอบข้อมูลผิดปกติ ในการวิเคราะห์ความถดถอยเชิงเส้น ของตัวสถิติทดสอบ 4 ตัว คือ

1. ตัวสถิติทดสอบของเมอร์วิน จี มาราสิง (MV)
2. ตัวสถิติทดสอบของฮาโคและไซมันนอฟฟ์ (HS)
3. ตัวสถิติทดสอบของไคนิฟาร์ดและสวอลโล ซึ่งมี 2 วิธี คือ
 - 3.1 วิธีเวียนเกิดโดยลำดับ (sequential recursive method(SRM))
 - 3.2 วิธีเวียนเกิดดัดแปร (modified recursive method(MRM))

ในการเปรียบเทียบตัวสถิติทั้ง 4 ตัว ผู้วิจัยจะพิจารณาจากเกณฑ์ความสามารถในการควบคุมความน่าจะเป็นของความผิดพลาดประเภทที่ 1, ความน่าจะเป็นซึ่งค่าผิดปกติที่ถูกตรวจพบเป็นค่าผิดปกติจริงทุกค่า (p_1), ความน่าจะเป็นซึ่งทำให้เกิดมาซคคิงเอฟเฟ็ค (p_2) และ ความน่าจะเป็นซึ่งทำให้เกิดซวอมฟิงเอฟเฟ็ค (p_3)

1.3 สมมติฐาน

ตัวสถิติทดสอบ HS มีประสิทธิภาพดีที่สุด (มีค่าของ p_1 สูงกว่าตัวอื่น หรือมีค่า p_2 ต่ำกว่าตัวอื่น หรือมีค่า p_3 ต่ำกว่าตัวอื่น) ในการตรวจพบค่าผิดปกติหลายค่า เพราะว่าการหาค่าผิดปกติของวิธีการนี้ มีการคัดเลือกข้อมูลที่ผิดปกติออกจากค่าคืออย่างเคร่งครัดก่อนจะทำการทดสอบ และการทดสอบค่าผิดปกติจะทดสอบจากค่าปลายน้อยสุดไปหาค่าปลายมากที่สุด (least to most extreme value) ซึ่งโอกาสที่เราจะพบค่าผิดปกติจริงมีมาก

1.4 ข้อตกลงเบื้องต้น

1. ความคลาดเคลื่อน (error) มีการแจกแจงแบบเดียวกัน (ยกเว้นกรณีค่าผิดปกติค่าความคลาดเคลื่อนอาจจะมีการปลอมปนของการแจกแจงอื่นได้และอิสระกัน)
2. การสร้างค่าผิดปกติเราจะกำหนดตำแหน่งของค่าผิดปกติเพื่อหาค่าความน่าจะเป็นของความผิดพลาดประเภทที่ 1 ค่า p_1 , p_2 และ p_3
3. การประมาณค่าสัมประสิทธิ์ความถดถอย ($\hat{\beta}$) จะใช้วิธีกำลังสองน้อยสุด (Least Square method)

1.5 ขอบเขตการวิจัย

1. ศึกษาวิธีการตรวจสอบค่าผิดปกติในการวิเคราะห์ความถดถอยเชิงเส้น โดยกำหนดจำนวนตัวแปรอิสระเท่ากับ 1,3 และ 5 ตัว
2. ศึกษากรณีขนาดตัวอย่างเท่ากับ 20, 50 และ 100
3. ศึกษาโดยกำหนดจำนวนค่าผิดปกติเท่ากับ 0, 1, 2 และ 3 ค่า
4. กำหนดระดับนัยสำคัญ (α) เท่ากับ 0.01 และ 0.05
5. การแจกแจงของความคลาดเคลื่อนที่ศึกษามีดังนี้

5.1 การแจกแจงหางยาวกว่าการแจกแจงปกติ ได้แก่

5.1.1 การแจกแจงปรกติปลอมปนในสเกล (Scale - contaminated normal distribution) เราจะศึกษาในกรณีที่ค่าพารามิเตอร์แสดงสเกล (c) เท่ากับ 3, 5 และ 10 ตามลำดับ และมีร้อยละการปลอมปนเท่ากับร้อยละของค่าผิดปรกติที่กำหนดในแต่ละขนาดตัวอย่าง ดังตารางที่ 1.1

5.1.2 การแจกแจงแบบปรกติปลอมปนในตำแหน่ง (location contaminated normal distribution) เราจะศึกษาในกรณีที่ค่าพารามิเตอร์แสดงตำแหน่ง (a) เท่ากับ 3, 5 และ 10 ตามลำดับ และมีร้อยละการปลอมปนเท่ากับร้อยละของค่าผิดปรกติที่กำหนดในแต่ละขนาดตัวอย่าง ดังตารางที่ 1.1

ตารางที่ 1.1 แสดงร้อยละการปลอมปนของการแจกแจงปรกติปลอมปนในตำแหน่ง และการแจกแจงปรกติปลอมปนในสเกล

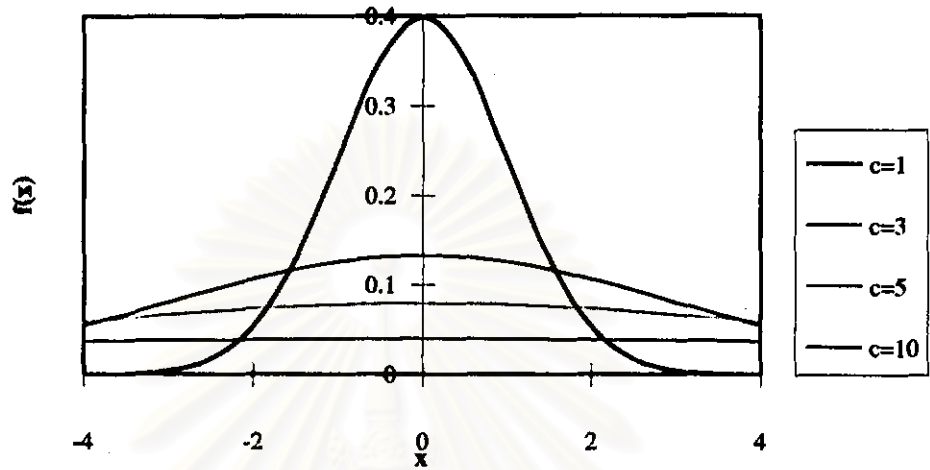
จำนวนค่าผิดปรกติ	ขนาดตัวอย่าง		
	20	50	100
1	5%	2%	1%
2	10%	4%	2%
3	15%	6%	3%

5.1.3 การแจกแจงที (t -distribution) เราจะศึกษากรณีที่ขนาดตัวอย่างเท่ากับ 20 ณ ระดับนัยความเสรี (d.f.) เท่ากับ 14, 16 และ 18 ของจำนวนตัวแปรอิสระเท่ากับ 1, 3 และ 5 ตามลำดับ เนื่องจากเมื่อขนาดตัวอย่างมากๆ การแจกแจงทีจะเข้าสู่การแจกแจงปรกติ และทำให้ความคลาดเคลื่อนต่ำกว่าความเป็นจริง

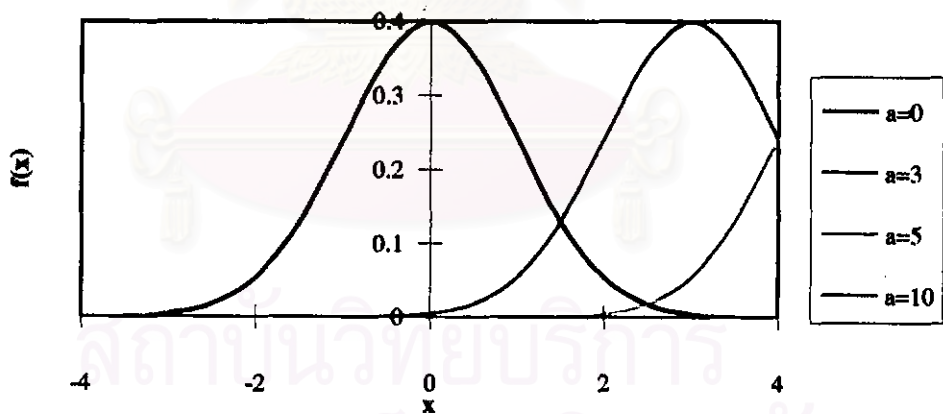
* ผู้วิจัยทดลองศึกษาเมื่อ $\mu = 0$ ด้วยค่าพารามิเตอร์แสดงสเกล(c)ต่างๆ พบว่าลักษณะโค้งการแจกแจงเปลี่ยนแปลงไปตามค่า c ดังรูปที่ 1.2 ซึ่งแสดงให้เห็นว่าเมื่อค่า c มากๆ โอกาสที่จะเกิดค่าผิดปรกติมีมากขึ้น

** ผู้วิจัยทดลองศึกษาเมื่อ $\sigma^2 = 1$ ด้วยค่าพารามิเตอร์แสดงตำแหน่ง(a)ต่างๆ พบว่าลักษณะโค้งการแจกแจงเปลี่ยนแปลงไปตามค่า a ดังรูปที่ 1.3 ซึ่งแสดงให้เห็นว่าเมื่อค่า a มากๆ โอกาสที่จะเกิดค่าผิดปรกติมีมากขึ้น

รูปที่ 1.2 แสดงการแจกแจงปกติปอดมปนในสเกล



รูปที่ 1.3 แสดงการแจกแจงปกติปอดมปนในตำแหน่ง



5.2 การแจกแจงแบบเบ้* (skewed distribution) เราจะศึกษาการแจกแจงดังนี้

5.2.1 การแจกแจงลอการิทึมปกติ** (lognormal distribution) ผู้วิจัยจะศึกษาในกรณีที่มีค่า $\mu = 0$, $\sigma^2 = 0.1, 0.2$ และ 0.7 การแจกแจงลอการิทึมปกติมีฟังก์ชันความหนาแน่นดังนี้

$$f(x) = \frac{1}{x\sigma\sqrt{2\pi}} \exp\left[-\frac{1}{2} \frac{(\log(x) - \mu)^2}{\sigma^2}\right]; \quad -x > 0, \sigma > 0, -\infty < \mu < \infty$$

โดยที่ μ และ σ^2 คือค่าเฉลี่ยและความแปรปรวนของ y เมื่อ $y = \log(x)$ มีการแจกแจงปกติ ค่าเฉลี่ยและความแปรปรวนของ x มีดังนี้

$$E(x) = \exp\left(\mu + \frac{1}{2}\sigma^2\right)$$

$$Var(x) = \exp(2\mu + \sigma^2) * [\exp(\sigma^2) - 1]$$

5.2.2 การแจกแจงแกมมา*** (gamma distribution) ผู้วิจัยจะศึกษากรณีที่ค่าพารามิเตอร์แสดงตำแหน่ง $\beta = 1$ และค่าพารามิเตอร์แสดงรูปร่าง $\alpha = 1, 2, 3$ และ 10 การแจกแจงแกมมามีฟังก์ชันความหนาแน่นดังนี้

$$f(x) = \frac{x^{\alpha-1} \exp(-x/\beta)}{\beta^\alpha \Gamma(\alpha)}; \quad x > 0, \alpha > 0, \beta > 0$$

โดยที่ β เป็นพารามิเตอร์แสดงตำแหน่ง

และ α เป็นพารามิเตอร์แสดงรูปร่าง

ค่าเฉลี่ย และความแปรปรวนของตัวแปร x คือ

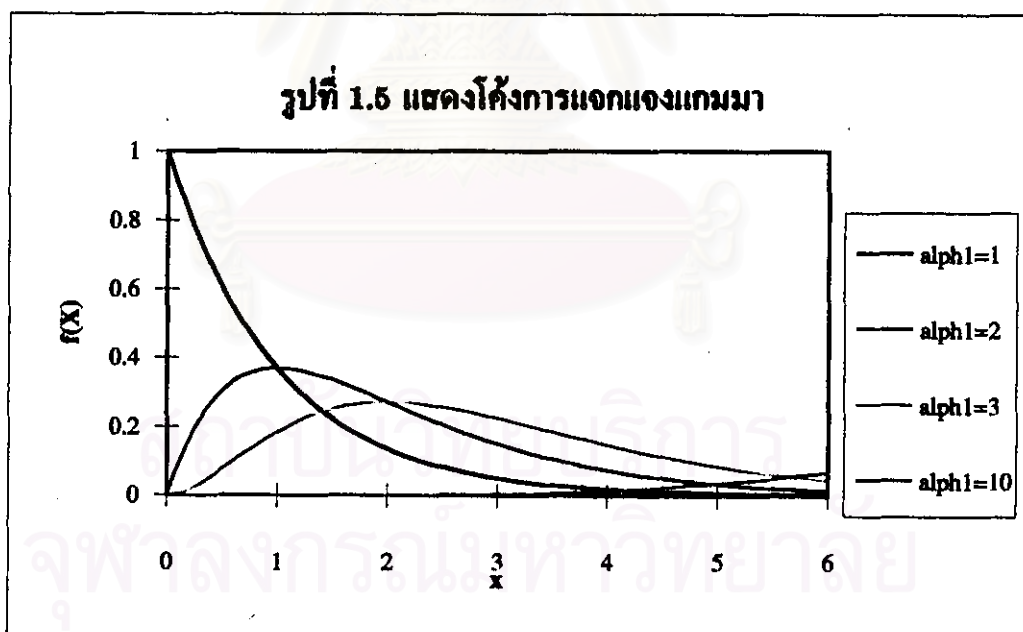
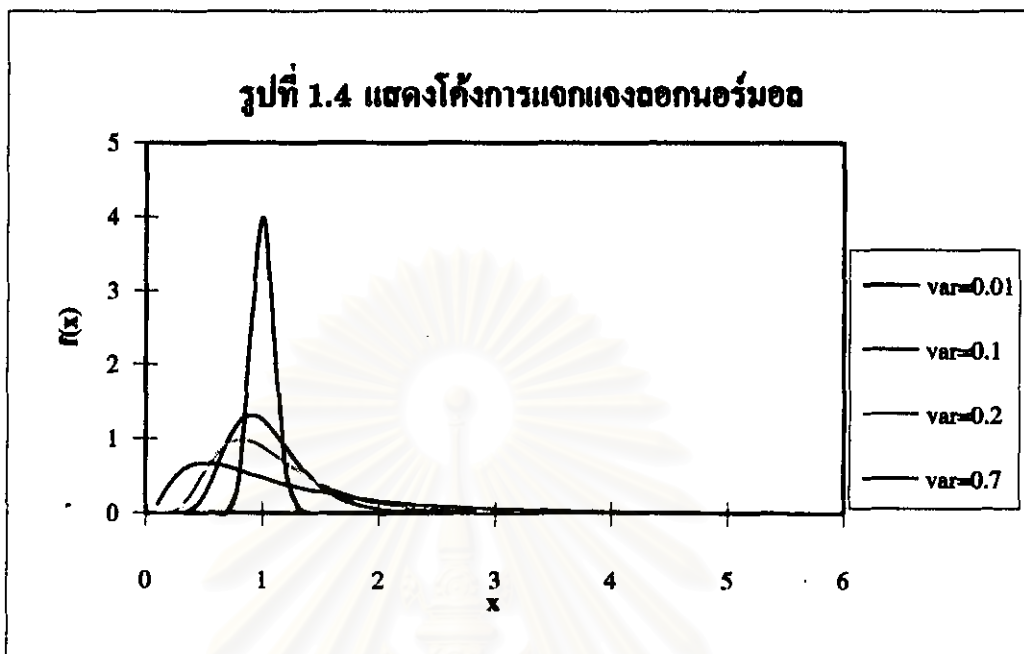
$$E(x) = \alpha\beta$$

$$Var(x) = \alpha\beta^2$$

* ศึกษาเฉพาะการแจกแจงเบ้ขวา เนื่องจากการแจกแจงเบ้ซ้ายเป็นส่วนกลับของการแจกแจงเบ้ขวา

** เลือกค่า σ^2 ตามลักษณะโค้งการแจกแจงที่เปลี่ยนไปดังรูปที่ 1.4

*** เลือกค่า α ตามลักษณะโค้งการแจกแจงที่เปลี่ยนไปดังรูปที่ 1.5



5.2.3 การแจกแจงไวบูลล์ (Weibull distribution) ผู้วิจัยจะศึกษากรณีที่ค่าพารามิเตอร์แสดงตำแหน่ง $\beta = 1$ และค่าพารามิเตอร์แสดงรูปร่าง $\alpha = 1, 2, 3$ และ 10 การแจกแจงไวบูลล์มีฟังก์ชันความหนาแน่นดังนี้

$$f(x) = \frac{\alpha x^{\alpha-1} \exp[-(x/\beta)^\alpha]}{\beta^\alpha} ; x >, \alpha > 0, \beta > 0$$

โดยที่ β เป็นพารามิเตอร์แสดงตำแหน่ง

และ α เป็นพารามิเตอร์แสดงรูปร่าง

ค่าเฉลี่ย และความแปรปรวนของตัวแปร x คือ

$$E(x) = \beta \Gamma\left(1 + \frac{1}{\alpha}\right)$$

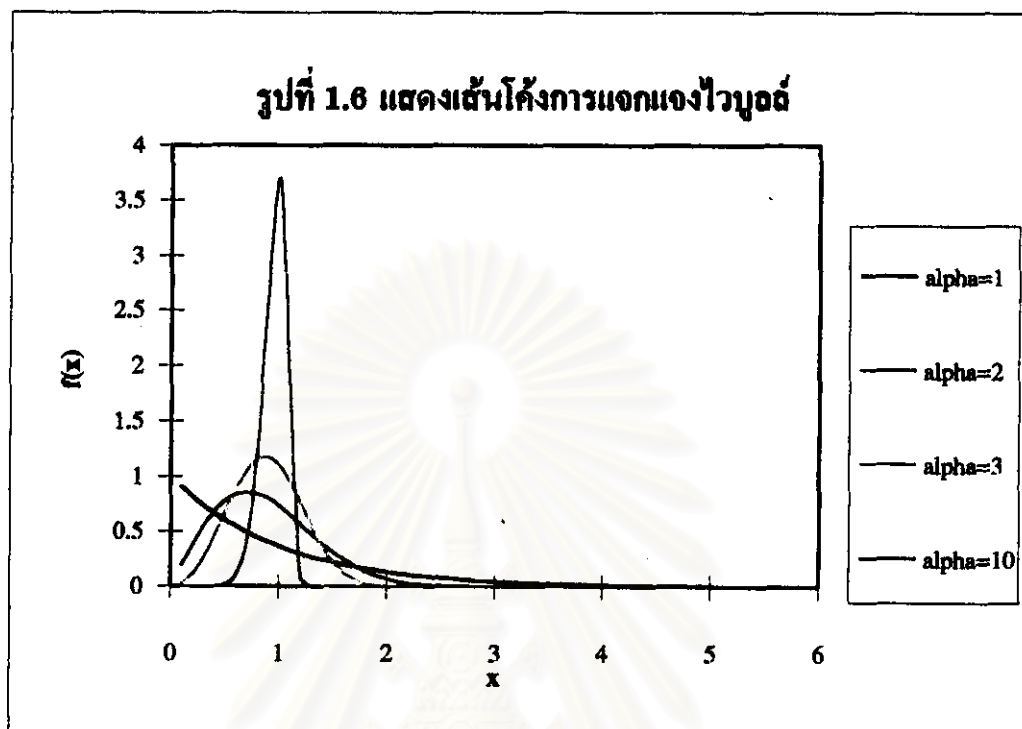
$$Var(x) = \beta^2 \left[\Gamma\left(1 + \frac{2}{\alpha}\right) - \Gamma^2\left(1 + \frac{1}{\alpha}\right) \right]$$

6. ทำการสร้างแบบจำลองข้อมูลให้มีสถานการณ์ตามต้องการ โดยใช้วิธีมอนติคาร์โลด้วยเครื่อง IBM/3031 ณ สถาบันบริการคอมพิวเตอร์ จุฬาลงกรณ์มหาวิทยาลัย และเขียนโปรแกรมด้วยภาษา FORTRAN

1.6 คำจำกัดความ

1. **ค่าผิดปกติ (Outliers)** หมายถึงค่าสังเกตที่มีค่ามากกว่าหรือน้อยกว่าค่าสังเกตอื่นๆ หรือค่าสังเกตที่ไม่ได้มาจากประชากรเดียวกัน
2. **มาซคิงเอฟเฟกต์ (Masking effect)** หมายถึง เหตุการณ์ที่ค่าผิดปกติค่าหนึ่งมีผลต่อค่าผิดปกติอีกค่าหนึ่งทำให้ไม่สามารถตรวจพบค่าผิดปกตินั้น
3. **ซวอมมิงเอฟเฟกต์ (Swamping effect)** หมายถึง เหตุการณ์ที่ค่าผิดปกติมีผลต่อค่าสังเกตอื่นๆ ที่ไม่ใช่ค่าผิดปกติ ทำให้นั้นกลายเป็นค่าผิดปกติซึ่งถูกตรวจพบ

* เลือกค่า α ตามลักษณะใส่การแจกแจงที่เปลี่ยนไปดังรูปที่ 1.6



สถาบันวิทยบริการ
จุฬาลงกรณ์มหาวิทยาลัย

4. **ความผิดพลาดประเภทที่ 1 (type I error)** หมายถึง ความคลาดเคลื่อนที่เกิดจากการปฏิเสธสมมติฐานว่าง (null hypothesis) เมื่อสมมติฐานว่างเป็นจริง
5. **ความผิดพลาดประเภทที่ 2 (type II error)** หมายถึง ความคลาดเคลื่อนที่เกิดจากการยอมรับสมมติฐานว่างเมื่อสมมติฐานว่างไม่จริง
6. **อำนาจการทดสอบ (power of the test)** หมายถึง ความน่าจะเป็นที่จะปฏิเสธสมมติฐานว่างเมื่อสมมติฐานว่างไม่จริง
7. **ความน่าจะเป็นซึ่งค่าสถิติที่ตรวจพบเป็นค่าสถิติจริงทุกค่า (p_1)** หมายถึง ความน่าจะเป็นซึ่งตัวสถิติทดสอบจะตรวจพบค่าสถิติโดยที่ค่าสถิตินั้นเป็นค่าสถิติจริงที่เรากำหนดทุกค่า
8. **ความน่าจะเป็นซึ่งทำให้เกิดมาทคิงเอฟเฟกต์ (p_2)** หมายถึง ความน่าจะเป็นที่จะไม่พบค่าสถิติทั้งที่ในข้อมูลที่ทำการทดสอบนั้นมีค่าสถิติอยู่
9. **ความน่าจะเป็นซึ่งทำให้เกิดขอมทิงเอฟเฟกต์ (p_3)** หมายถึงความน่าจะเป็นที่จะตรวจพบค่าสถิติทั้งที่ในข้อมูลนั้นไม่มีค่าสถิติ หรือความน่าจะเป็นที่จะตรวจพบค่าสถิติโดยที่ค่าสถิตินั้นเป็นค่าสถิติไม่จริง (เป็นค่าปรกติ)
10. **ความแกร่งของการทดสอบ (Robustness)** หมายถึง คุณสมบัติของการทดสอบที่ไม่ไวต่อการเปลี่ยนแปลงของปัจจัยอื่น ที่ไม่ใช่ปัจจัยที่ต้องการทดสอบ เช่น การฝ่าฝืนข้อตกลงเบื้องต้นของการทดสอบนั้น ซึ่งสิ่งที่ใช้พิจารณาความแกร่งของการทดสอบ คือ ความผิดพลาดประเภทที่ 1

1.7 เกณฑ์การตัดสินใจ

เกณฑ์การตัดสินใจว่าวิธีการตรวจสอบค่าสถิติวิธีใดมีประสิทธิภาพสูงสุดจะพิจารณาภายใต้สมมติฐานว่าง

H : ไม่มีค่าผิดปกติ

เทียบกับ K : มีค่าผิดปกติอย่างน้อย 1 ค่า

* เกณฑ์ที่ใช้วัดดังกล่าวในข้อ 8, 9 และ 10 เดวิดและพอลสัน (David and Paulson) ได้เสนอไว้ในปีค.ศ. 1965 ต่อมาในปี ค.ศ.1990 คีนิฟาร์ดและชวอลโล (Kianifard and Swollow) และ ในปี ค.ศ. 1993 ฮาดี (Hadi) ได้นำมาคิดแปลงใช้ในการวัดประสิทธิภาพของวิธีการตรวจสอบค่าผิดปกติ

ผู้วิจัยจะทำการเปรียบเทียบความน่าจะเป็นของตัวสถิติทดสอบ 2 ลักษณะ ดังนี้

1. ความสามารถในการควบคุมความน่าจะเป็นของความผิดพลาดประเภทที่ 1 ของการทดสอบแต่ละสถานการณ์ โดยใช้เกณฑ์ของคอกเรน (Cochran) และบราวเดย์ (Bradley) ซึ่งเราจะพิจารณาดังนี้

ก) เกณฑ์ของคอกเรน (Cochran) เราจะพิจารณาว่าตัวสถิติทดสอบใดมีค่าความน่าจะเป็นของความผิดพลาดประเภทที่ 1 ซึ่งได้จากการทดลอง อยู่ในช่วง $[0.007, 0.015]$ และ $[0.04, 0.06]$ ณ ระดับนัยสำคัญ 0.01 และ 0.05 ตามลำดับ เราจะถือว่าตัวสถิติทดสอบนั้นควบคุมความน่าจะเป็นของความผิดพลาดประเภทที่ 1 ได้

ข) เกณฑ์ของบราวเดย์ (Bradley) เราจะพิจารณาว่าตัวสถิติทดสอบใดมีค่าความน่าจะเป็นของความผิดพลาดประเภทที่ 1 ซึ่งได้จากการคำนวณ อยู่ในช่วง $[0.005, 0.015]$ และ $[0.025, 0.075]$ ณ ระดับนัยสำคัญ 0.01 และ 0.05 ตามลำดับ เราจะถือว่าวิธีการตรวจสอบนั้นควบคุมความน่าจะเป็นของความผิดพลาดประเภทที่ 1 ได้

2. พิจารณาความน่าจะเป็นซึ่งค่าผิดปรกติที่ถูกตรวจพบเป็นค่าผิดปรกติจริงทุกค่า (p_1), ความน่าจะเป็นซึ่งทำให้เกิดมาซคคิงเอฟเฟ็ค (p_2) และความน่าจะเป็นซึ่งทำให้เกิดชวอมฟิงเอฟเฟ็ค (p_3) เราจะทำการเปรียบเทียบเฉพาะตัวสถิติทดสอบที่สามารถควบคุมของความน่าจะเป็นของความผิดพลาดประเภทที่ 1 ได้เท่านั้น

1.8 ประโยชน์ที่คาดว่าจะได้รับ

ผลการศึกษาทำให้ผู้ใช้สามารถเลือกตัวสถิติทดสอบที่เหมาะสมในกรณีที่ข้อมูลมีค่าผิดปรกติหลายค่าในสมการการถดถอยเชิงเส้น