

บทที่ ๑

บทนำ

๑.๑ คำนำ

ในการแก้ปัญหาทางสถิติทั่วไปนั้น ไม่ว่าจะเป็นการแก้ปัญหาด้าน การประมาณค่า (estimation problems) หรือในการทดสอบสมมติฐาน (hypothesis testing problems) ของค่าพารามิเตอร์ (parameter) ความที่ค่าทางสถิติทั้งชั้นต้นและชั้นกลางได้กล่าวถึงนั้น จะกล่าวถึงแต่การใช้ ข้อมูลตัวอย่างที่เป็นค่าสัมบูรณ์ (absolute value) เท่านั้น แต่ความจริง นั้นปัญหาทางสถิติเหล่านี้ยังสามารถจะแก้ได้ด้วยข้อมูลประเภทอื่น ๆ ได้ลึก วิทยาลัยพหุวัฒนธรรมนี้มีวัตถุประสงค์ที่จะแสดงให้เห็นถึงความสำคัญทั้งทางด้านทฤษฎี และทางด้านการใช้ประโยชน์จากการใช้ข้อมูลสถิติอีกประเภทหนึ่งที่เน้นถึง ความสำคัญของตำแหน่งของข้อมูล โดยมีการจัดเรียงลำดับค่าของข้อมูล ตามขนาด (arranged in order of magnitude) ที่เรียกกันว่า "Order Statistics" สถิติที่จะกล่าวถึงนี้ไม่เพียงแต่จะสามารถใช้แก้ ปัญหาทางสถิติได้ เช่น วิธีการทางสถิติทั่วไป (typical statistics) ที่ใช้ค่าสัมบูรณ์ แต่ยังมีคุณสมบัติอื่นที่ดีกว่าอีกหลายประการ อาทิเช่น ง่ายใน การทำความเข้าใจ ใช้งานได้สะดวก (simplicity) มีประสิทธิภาพ (efficiency) ที่ดี เป็นต้น

*ค่าสัมบูรณ์ หมายถึงค่าทางพีชคณิต (algebraic value) กล่าวคือ เป็น ค่าตัวเลขที่แท้จริง ซึ่งอาจเป็นได้ทั้งค่าบวกหรือค่าลบ แต่เราไม่คำนึงถึง เครื่องหมาย เช่น -๗ และ $+๗$ เราก็คือว่าเป็น ๗ เท่านั้น

ด้วยควาสำคัญของ order statistics ดังกล่าวมาแล้วนั้น วิทยานิพนธ์นี้มีเจตจำนงที่จะเผยแพร่ให้เป็นที่รู้จักทั่วไป เพราะเท่าที่เป็นอยู่ในปัจจุบันนี้ ตำราทางสถิติภาคภาษาไทยทั้งในขั้นต้นและขั้นกลางยังมีได้เป็นการกล่าวถึงสถิติที่ได้จากข้อมูลประเภทนี้เลย หรือแม้แต่ในภาษาอังกฤษซึ่งมีอยู่ไม่มากนักก็กล่าวถึงสถิติประเภทนี้ก็ยังเป็นที่ยากยิ่งต้องการที่จะทำการศึกษาค้นคว้า จะมีอยู่มากก็ เป็นบทความที่เขียนลงไว้ในวารสารภาษาอังกฤษด้านวิชาสถิติและคณิตศาสตร์หลายฉบับกระจายอยู่ทั่วไป วิทยานิพนธ์นี้ได้ทำการรวบรวมทฤษฎีจากตำราและวารสารดังกล่าว โดยได้มีการเรียบเรียงเสียใหม่ให้ต่อเนื่อง เพื่อความสะดวกในการศึกษาและทำความเข้าใจ พร้อมทั้งนี้ได้นำตัวอย่างที่เป็นตัวเลขมาเสนอไว้ประกอบในการทำความเข้าใจด้วย ตัวอย่างเหล่านี้จะเป็นตัวเลขที่ได้จากการใช้ประโยชน์ของ order statistics จริง ๆ ในทางปฏิบัติ ด้วยเหตุนี้จึงหวังว่าวิทยานิพนธ์นี้จะพอเป็นประโยชน์แก่ผู้ที่ทำการศึกษาด้านวิชาสถิติและคณิตศาสตร์หรือสาขาอื่นใดที่มีส่วนสัมพันธ์ต่อทฤษฎีเหล่านี้ รวมทั้งผู้ที่สนใจในสถิติประเภทนี้เพื่ออาศัยเป็นบรรทัดฐานในการค้นคว้าศึกษาเพิ่มเติมต่อไป

๑.๒ Order Statistics คืออะไร

ถ้าเรามีตัวอย่างขนาด n หน่วยย่อยชุดหนึ่ง สมมติว่าเป็น $x_1, x_2, x_3, \dots, x_n$. โดยที่ตัวอย่างชุดนี้ถูกเลือกโดยการสุ่มจากประชากรที่มีฟังก์ชันของการกระจายของความน่าจะเป็น $f(x)$ จากนั้นถ้าเราจะจัดเรียงข้อมูลเสียใหม่ตามลำดับความสำคัญ เช่น เรียงจากค่าต่ำสุดไปยังค่าสูงสุดดังปี คือ $x_{(1)}; x_{(2)}; \dots; x_{(n)}$. ด้วยคุณสมบัติที่ว่า

$$x_{(1)} \leq x_{(2)} \leq x_{(3)} \leq \dots \leq x_{(n)},$$

แล้ว เราเรียกข้อมูลที่ได้รับการจัดเรียงตามขนาดของความสำคัญนี้ว่า

Order Statistics และค่า $X_{(n)}$ กับค่า $X_{(1)}$ เรียกว่า "ค่าสูงสุด" กับ "ค่าต่ำสุด" หรือโดยทั่ว ๆ ไปเราเรียก $X_{(i)}$ ว่าเป็นค่าของตัวที่มีตำแหน่งที่ i ของ order statistic ความแตกต่างระหว่างค่าสูงสุดและค่าต่ำสุด $X_{(n)} - X_{(1)} = R$ เรียกว่า "พิสัย" ของข้อมูลตัวอย่าง

ในกรณีที่ฟังก์ชันของการกระจายของความน่าจะเป็นเป็น continuous distribution แล้ว order statistics อาจเขียนได้ดังนี้

$$X_{(1)} < X_{(2)} < X_{(3)} < \dots < X_{(n)} .$$

ตัวอย่างที่ ๑.๒.๑ จากหนังสือรายงานเศรษฐกิจประจำเดือน มกราคม พ.ศ.๒๕๑๑ ของธนาคารแห่งประเทศไทยได้รายงานไว้ว่า ราคาเฉลี่ยของข้าวชนิด ๑๐๐ % ต่อเบตริกตันในช่วงระยะเวลาจาก พ.ศ.๒๕๐๓ - ๒๕๑๐ เป็นตามตารางที่ ๑.๒.๑ ดังนี้

ตารางที่ ๑.๒.๑

แสดงการเรียงลำดับข้อมูลตามระยะเวลาที่เกิดเหตุการณ์

ปี	ราคา (บาท)
๒๕๐๓	๑,๖๔๑
๒๕๐๔	๑,๗๓๑
๒๕๐๕	๑,๘๘๖
๒๕๐๖	๑,๗๘๕
๒๕๐๗	๑,๖๘๐
๒๕๐๘	๑,๖๔๕
๒๕๐๙	๒,๑๘๕
๒๕๑๐	๒,๕๓๒

ถ้าเราจะได้จัดเรียงข้อมูลชุดนี้เสียใหม่เป็น "ราคาเฉลี่ยของข้าวชนิด ๑๐๐ % คอเมกริกตันใน ช่วงระยะเวลาจาก พ.ศ.๒๕๐๓ - ๒๕๑๐ จาก ราคาสูงสุดไปยังราคาต่ำสุด" จะได้อันนี้

ตารางที่ ๑.๒.๒

แสดงการจัดเรียงข้อมูลตามลำดับขนาดของความสำคัญที่เราสนใจ

ปี	ราคา (บาท)
๒๕๑๐	๒,๕๓๒
๒๕๐๙	๒,๑๘๕
๒๕๐๕	๑,๙๙๒
๒๕๐๖	๑,๙๕๕
๒๕๐๔	๑,๙๓๑
๒๕๐๗	๑,๖๘๐
๒๕๐๘	๑,๖๔๕
๒๕๐๓	๑,๖๔๑



ข้อมูลที่ได้รับการจัดเรียงใหม่นี้เองที่เรียกว่า "Order Statistics" ตัวอย่างของ "Order Statistics" ที่เรากำลังใช้อยู่อาจนำมาเสนอได้อันนี้

กึ่งพิสัย (Mid-range) หมายถึงมัธยิมที่ได้จากค่าเฉลี่ยของคะแนนสูงสุดกับคะแนนต่ำสุด เช่น จากตัวอย่างที่ ๑.๒.๑ จะได้อีกึ่งพิสัยเป็น $1/2 (๒๕๓๒/๑๖๔๑)$ หรือเท่ากับ ๒๐๘๖.๕

มัธยฐาน (Median) หมายถึงมัธยิมที่แสดงให้ทราบว่า ครึ่งหนึ่งของจำนวนคะแนนมีค่าสูงกว่ามัธยิมนี้ และอีกครึ่งหนึ่งมีค่าต่ำกว่า หรืออาจกล่าวว่าเป็นค่าของคะแนนที่อยู่ในตำแหน่งที่ ๕๐ ของคะแนนทั้งหมดที่แบ่งได้เป็น ๑๐๐ ตำแหน่ง

จากตัวอย่างที่ ๑.๒.๑ บัญชีฐานก็คือ $๑,๓๕๕ \neq ๑,๓๓๑$ หรือกล่าว
ว่าบัญชีฐานของราคาเฉลี่ยของชาวชนิต ๑๐๐ % คอเมตริกตันจากช่วงระยะ
เวลาจาก พ.ศ.๒๕๐๓ - ๒๕๑๐ เป็น ๑,๓๖๕ บาท

ควอไทล์, เดไซล์, และเปอร์เซ็นต์ไทล์ (Quartile, Decile, Percentile) เหล่านี้มีใช้บ่อยๆ แต่เป็นตัวที่แสดงให้เห็นให้ทราบว่านิโคเนนที่
ต่ำกว่าตัวมันเองก็ตัวและที่สูงกว่าอีกก็ตัว หรือกล่าวว่าเป็นการแสดงตำแหน่ง
ของคะแนนเพื่อเทียบกับคะแนนทั้งหมด อธิบายได้ว่า

ควอไทล์ คือคะแนนที่แสดงให้เห็นว่า จำนวนคะแนนอื่น ๆ ที่ต่ำกว่า
คะแนนนี้เป็น $๑/๔$ ของจำนวนคะแนนทั้งหมดเรียกว่าควอไทล์ที่หนึ่ง หรือ
จำนวนคะแนนอื่น ๆ ที่ต่ำกว่านี้เป็น $๓/๔$ ของคะแนนทั้งหมดก็เรียกว่า
ควอไทล์ที่สาม

เดไซล์ คือ คะแนนที่แสดงให้เห็นว่า จำนวนคะแนนอื่น ๆ ที่
ต่ำกว่าคะแนนนี้เป็น $๑/๑๐, ๒/๑๐, ๓/๑๐, \dots, ๙/๑๐$.

เปอร์เซ็นต์ไทล์ คือ คะแนนที่แสดงให้เห็นว่า จำนวนคะแนนอื่น ๆ
ที่ต่ำกว่าคะแนนนี้เป็นกี่เปอร์เซ็นต์ ซึ่งก็คล้ายกับของควอไทล์และเดไซล์
เป็นแต่เพียงว่าเราแบ่งคะแนนทั้งหมดให้อยู่ในตำแหน่ง ๑๐๐ ตำแหน่งด้วยกัน
ทั้งหมด

พิสัย (range) เป็นการวัดการกระจายของข้อมูลวิธีหนึ่งโดยดู
จากความแตกต่างของคะแนนที่อยู่ในตำแหน่งสูงสุดกับตำแหน่งต่ำสุด เช่น
จากตัวอย่างที่ ๑.๒.๑ พิสัยของราคาเฉลี่ยของชาวชนิต ๑๐๐ % คอเมตริกตัน
จากช่วงระยะเวลาจาก พ.ศ.๒๕๐๓ - ๒๕๑๐ เป็น $๒๕๓๒ - ๑๖๔๑ = ๘๙๑$ บาท

ในการวัดการกระจายของข้อมูลโดยพิสัยนี้ เราวินิจฉัยว่า ข้อมูล
ชุดที่มีพิสัยมากกว่าย่อมมีการกระจายของข้อมูลมากกว่า ฉะนั้นจะเห็นว่า พิสัย
เป็นการวัดการกระจายที่คร่าว ๆ เพราะชี้เฉพาะเพียงคะแนนตำแหน่งสูงสุดกับ

ต่ำสุดเท่านั้น โดยต้องการความรวดเร็วเป็นสำคัญ

ส่วนเบี่ยงเบนควอไทล์ (Quartile Deviation) เป็นการวัดการกระจายของข้อมูลวิธีหนึ่ง โดยดูจากส่วนต่างของควอไทล์ที่สามกับที่หนึ่ง ($Q_3 - Q_1$) ถ้าข้อมูลใดมีส่วนเบี่ยงเบนควอไทล์มากก็ถือว่าข้อมูลมีการกระจายมาก ถ้าส่วนเบี่ยงเบนควอไทล์น้อยก็ถือว่าข้อมูลมีการกระจายน้อย

๑.๓ ความหมายของคำที่ใช้

เพื่อความสะดวกในการทำความเข้าใจกับวิทยานิพนธ์ฉบับนี้จึงได้อธิบายความหมายของคำบางคำที่จะต้องใช้บ่อย ๆ ดังต่อไปนี้

ข้อมูลทางสถิติ (Statistical data) อาจจำแนกได้สามวิธีการเก็บรวบรวมเป็นส่วนประเภทดังนี้คือ

ก. Sample data เป็นข้อมูลที่ได้จากการสุ่มตัวอย่างจากประชากรที่เราสนใจ ในการรวบรวมข้อมูลเราสนใจในลักษณะของตัวอย่างเท่านั้น เช่น การสำรวจรายได้ การสำรวจการใช้จ่าย การสำรวจทัศนคติ เป็นต้น แต่ในการนี้เราไม่สนใจถึงสาเหตุหรือปัจจัยที่ทำให้เกิดลักษณะดังกล่าว กล่าวคือเราไม่สนใจว่าในการสำรวจรายได้นั้นอะไรเป็นสาเหตุทำให้คนส่วนหนึ่งมีรายได้สูงบางคนอีกส่วนหนึ่งมีรายได้ต่ำ

ข. Experimental data เป็นข้อมูลที่ได้จากการทดลอง คือการแบ่งตัวอย่างที่ถูกเลือกมานั้นออกเป็นกลุ่ม ๆ ให้อยู่ภายใต้สิ่งแวดล้อมต่าง ๆ กัน เช่น การใช้ปุ๋ยต่างชนิดกันในพืชชนิดเดียวกัน ในอุณหภูมิที่แตกต่างกัน ข้อมูลประเภทนี้ใช้วัดถึงสาเหตุและปัจจัยที่ทำให้เกิดลักษณะประจำแก่ประชากร

๓. Observational data เป็นข้อมูลที่ได้จากการทดลองเช่นเดียวกับ Experimental data แตกต่างที่ถ้าใน Experimental data นั้นเราสามารถที่จะกำหนดขนาดของตัวอย่างได้ เพราะเราสามารถที่จะหาตัวอย่างได้ตามขนาดที่ต้องการ แต่ในกรณีของ Observational data เราไม่สามารถที่จะมีหรือกำหนดขนาดของตัวอย่างได้เสมอไป ที่เป็นเช่นนั้นก็เพราะลักษณะพิเศษของประชากรนั้น ดังนั้นเราจึงใช้ตัวอย่างจากสิ่งที่มีอยู่ เช่น การทดสอบว่า การสูบบุหรี่เป็นสาเหตุของการเป็นมะเร็ง เราก็จะต้องศึกษาจากผู้ที่เป็นหรืออยู่แล้ว เราไม่สามารถที่จะเกณฑ์ให้คนที่ไม่สูบบุหรี่ให้มาสูบบุหรี่เพื่อเป็นเครื่องทดลองของเราได้ อีกทั้งเราก็ไม่สามารถที่จะกำหนดขนาดของตัวอย่างทดลองมากกว่าจำนวนผู้สูบบุหรี่ในขณะนั้นไปได้

พารามิเตอร์ (parameter) คือ "ลักษณะ" ของประชากร หรือ "ลักษณะ" ของการกระจายของประชากร เช่น อัตราส่วนประชากรชายต่อหญิงของจังหวัดชลบุรีเป็น ๑.๑ จำนวนครัวเรือนการเกษตรของประเทศไทยเป็น ๗๓.๕ % รายได้เฉลี่ยต่อบุคคลของประเทศไทยเป็น ๒,๔๐๐ บาทต่อปี

Statistics คือ "ลักษณะ" ของตัวอย่าง หรือ "ลักษณะ" ของการกระจายของตัวอย่าง โดยปกติจะเป็นฟังก์ชันของ observable random variables หน้าที่อันสำคัญยิ่งของ Statistics ในทางสถิติคือใช้เป็นตัวประมาณพารามิเตอร์ เช่น $\bar{X} = \frac{1}{n} \sum_{i=1}^n x_i$ เป็น Statistic ที่ใช้ประมาณค่าของ $\mu = \frac{1}{N} \sum_{i=1}^N x_i$ ซึ่งเป็นพารามิเตอร์ที่เรายังไม่ทราบค่าที่แท้จริง

นอกจากค่าเหล่านี้แล้ว ค่าอื่น ๆ ที่มีความหมายในทางวิชาการ เมื่อถูกนำมาใช้ในชีวิตประจำวันก็จะได้มีคำอธิบายเพิ่มเติมเฉพาะคำไป在本นั้น ๆ