

## Unit 9

มี การพูดถึงหัวข้อหลักที่สำคัญเกี่ยวกับการวิเคราะห์ข้อมูลที่มีลักษณะ "collinear"

หรืออาจจะ ความสัมพันธ์เชิงเส้นที่ดูในกราฟหรืออะไร ในกรณีการถดถอย (Regression Analysis)

- ความหมายของ "Collinearity" หมายถึง Predictor Variables ใน multiple regression มีความสัมพันธ์เชิงเส้นที่

ตัวแปรบางตัวมีผลในค่าของตัวแปรอิสระบางตัวในแบบจำลองไม่สมบูรณ์แยกออกจากกันอย่างชัดเจน ที่ตัวแปรบางตัว

สัมพันธ์กันมากเกินไป จะทำให้การตีความค่าความสัมพันธ์ของสมการถดถอยทำได้ยาก เพราะการเปลี่ยนแปลงในตัวแปรหนึ่ง

อาจทำให้ได้การเปลี่ยนแปลงในค่าแปรอื่นโดยไม่สามารถบอกได้ว่าตัวแปรแต่ละตัวนั้น เป็นอิสระต่อกันได้ ในการตีความค่า

ตัวแปรที่มีความสัมพันธ์กันสูงๆ ตัวแปรบางตัวอาจกล่าวได้ว่า "Orthogonal" ซึ่งหมายถึง ไม่มีความสัมพันธ์

เส้นกันเลย แต่ในสถานการณ์ส่วนใหญ่ในชีวิตจริง การที่ตัวแปรอิสระเป็นอิสระจากกันสูงๆ นั้น ทำได้ยาก

- ผลกระทบของ Collinearity

- ค่าสัมประสิทธิ์ถดถอยจะไม่เสถียร ค่าสัมประสิทธิ์การถดถอยอาจเปลี่ยนแปลงอย่างคาดไม่ถึงในการเปลี่ยนแปลง

ตัวแปรใดๆ ออกจากแบบจำลอง แม้ว่าตัวแปรที่สัมพันธ์กันจะไม่มีความสัมพันธ์มากนักก็ตาม

- ข้อสรุป จากแบบจำลองที่ใช้ค่าสัมประสิทธิ์ที่ได้จากข้อมูลที่มี collinearity อาจทำให้การพยากรณ์ผิดได้

- การตรวจหา Collinearity

ในทางทฤษฎีการวิเคราะห์ข้อมูล การวิเคราะห์ข้อมูลต่างๆ เช่น : - ค่าสัมประสิทธิ์ความสัมพันธ์ : หากพบว่า ค่าสัมประสิทธิ์

ระหว่างตัวแปรอิสระคู่ใดคู่หนึ่ง แสดงว่าสมมติฐานที่จะเกิด collinearity

- Variance Inflation Factor (VIF) : ถ้าค่า VIF สูงกว่า 10 จะมีความน่าเชื่อถือต่ำของสมการที่สร้างขึ้นจาก collinearity

- Eigenvalues : การตรวจหาค่าที่มีค่าเหลือน้อยกว่า 0 จะช่วยชี้ให้เห็น collinearity

• วิธีแก้ปัญหาคอลลิเนียร์ Collinearity มีหลายวิธีแก้ปัญหาคอลลิเนียร์ ดังนี้

- ① "เพิ่มข้อมูล" การเพิ่มข้อมูลเพิ่มเติมอาจช่วยลดปัญหานี้ได้ แต่ก็ไม่ดีถ้าหากทำได้ลำบาก หรือจากข้อจำกัดด้านงบ
- ② "ลดจำนวนตัวแปร" ศึกษาความสัมพันธ์ของตัวแปรจากแบบจำลอง อาจช่วยแก้ปัญหานี้ได้ในบางกรณี
- ③ "การใช้วิธีทางสถิติอื่นๆ" เช่น PCA วิธีนี้จะช่วยแก้ปัญหาคอลลิเนียร์ โดยการปรับแต่งในเสถียร
- ④ "Ridge Regression" เป็นอีกวิธีหนึ่งที่จะช่วยลดปัญหาคอลลิเนียร์ โดยการเพิ่มเงื่อนไขลงในแบบจำลองการถดถอย ซึ่งทำให้อุปกรณ์การคำนวณมีเสถียรภาพต่อการคำนวณด้วย

★ ผู้เขียนเห็นว่าปัญหาคอลลิเนียร์ ไม่ใช่ข้อผิดพลาดในการสร้างแบบจำลอง แต่เป็นข้อจำกัดของข้อมูลที่นำมาใช้ ดังนั้นวิธีแก้ที่ดีที่สุดคือควรพิจารณาว่ามีการเพิ่มข้อมูลเพิ่มเติมหรือไม่

ตัวอย่างที่ผู้เขียนใช้ในหนังสือเป็นการทำนายยอดขายของบริษัทต่างๆ โดยใช้ตัวแปร เช่น การผลิตในประเทศ (domestic production) การบริโภคในประเทศ และสินค้าคงคลัง ในช่วงปี 1949-1966 เมื่อตัวแปรเหล่านี้มี "collinearity" กันอย่างรุนแรง การคำนวณพหุคูณจะได้อัตราที่ผิดพลาดไปมากที่ควรจะเป็น แม้จะวัด  $R^2$  สูงถึง 0.99 แต่การพยากรณ์ก็ไม่น่าเชื่อถือ

หลักการใช้ Principal Component Analysis (PCA) และ

Ridge Regression ในการแก้ไขปัญหาความสัมพันธ์เชิงเส้น (Collinearity) รวมถึงข้อจำกัดและข้อควรระวังในการใช้วิธีเหล่านี้

## ★ ทฤษฎีองค์ประกอบหลัก (PCA)

PCA เป็นวิธีการที่ใช้ในการเปลี่ยนตัวแปรต้นที่มีลักษณะสัมพันธ์กัน (collinear) ให้กลายเป็นชุดของตัวแปรใหม่ที่ไม่มีความสัมพันธ์เชิงเส้นต่อกัน ซึ่งเราเรียกตัวแปรใหม่นี้ว่า "องค์ประกอบหลัก" (PC) โดยใช้ตัวแปรใหม่สี่ๆ เป็นกร รวบรวมให้เส้น ของ ตัวแปร ต้นฉบับ

ขั้นตอนการทำงานของ PCA :

- ① คำนวณ Eigenvalues และ Eigenvectors : จากเมทริกซ์ความสัมพันธ์ (correlation matrix) ของตัวแปรต้นฉบับ โดยตัวแปรในแต่ละตัว จะถูกสร้างจาก Eigenvectors ในค่านี้
- ② เลือกจำนวนองค์ประกอบหลัก (PC) : ตัวแปรใหม่เหล่านี้จะถูกจัดเรียงตามลำดับความสำคัญ โดยองค์ประกอบหลักแรกจะมีค่าแปรปรวนมากที่สุด ดังนั้นเราสามารถเลือกใช้เฉพาะองค์ประกอบที่สำคัญ เมื่อแทนที่ตัวแปรต้นฉบับที่มีปัญหาคอลลิเนอริตี้
- ③ ทำองค์ประกอบหลักมาวิเคราะห์ : นวัตกรรมที่ได้ดูขององค์ประกอบหลักที่ไม่มีความสัมพันธ์กัน แล้วจะนำข้อมูลตัวแปรนี้ไปใช้ในการสร้างแบบจำลองใหม่

## ★ ข้อดีของ PCA

- ช่วยแก้ปัญหา collinearity ได้อย่างมีประสิทธิภาพ
- ลดจำนวนตัวแปรในแบบจำลองลงได้ ทำให้แบบจำลองง่ายต่อการตีความ
- ปรับปรุงประสิทธิภาพของค่าความสัมพันธ์ในแบบจำลองการถดถอย

## ★ ข้อจำกัดของ PCA :

- การตีความที่ยากขึ้น
- ข้อมูลที่สูญหาย

★ Ridge Regression → เป็นอีกหนึ่งวิธีที่นิยมใช้ในการจัดการกับปัญหา collinearity โดยเพิ่มพารามิเตอร์ควบคุมในแบบจำลองการถดถอย (regression model) ที่ช่วยปรับปรุงประสิทธิภาพการคำนวณและลดผลกระทบจาก multicollinearity

หลักการการทำงานของ Ridge Regression :

- ① เพิ่มพารามิเตอร์ควบคุม
- ② เลือกค่าพารามิเตอร์  $k$  ที่เหมาะสม

การเลือกใช้วิธี Ridge Regression และ PCA ขึ้นอยู่กับลักษณะของข้อมูลและวัตถุประสงค์ของการวิเคราะห์ : การจัดการกับปัญหา Collinearity และต้องการรักษาตัวแปรอิสระไว้มากที่สุดในแบบจำลอง Ridge Regression จะเป็นทางเลือกที่ดีกว่า

: ทั้งสองการลดจำนวนตัวแปรและลดการวิเคราะห์ ที่ง่ายต่อการจัดการ \* PCA เน้นการ

## Summary

การแก้ปัญหา Collinearity เป็นส่วนสำคัญในการวิเคราะห์การถดถอย โดย PCA และ Ridge Regression เป็นวิธีที่ใช้ได้ผลดีในการจัดการกับปัญหาที่ขึ้นกับลักษณะข้อมูลและเป้าหมายในการวิเคราะห์ ผู้ใช้จึงควรเลือกใช้วิธีที่เหมาะสมเพื่อใช้ได้แบบจำลองที่มีความเสถียรและแม่นยำมากที่สุด

Regression Analysis by example จะแสดงให้เห็นการทำงานของข้อมูลที่มีปัญหา

## Unit 10

\* Collinearity \* และขบวนการเกี่ยวกับวิธีวิธีต่าง ๆ ในการจัดการปัญหานี้ โดยในบทที่ 10 เราจะศึกษาวิธีวิธีสองแบบที่นิยมใช้ ได้แก่ Principal Component Regression (PCR) และ Ridge Regression ซึ่งเป็นเทคนิคหลักในการจัดการกับปัญหา collinearity

การทำงานกับข้อมูลที่มี Collinearity (Unit 10) ผู้เขียนได้เสนอวิธีวิธีเพิ่มเติมสำหรับการทำงานกับข้อมูลที่มี

ปัญหา \* collinearity \* โดยแนะนำวิธีวิธี < PCR > และ Ridge Regression

ความสำคัญของการทำงานกับปัญหา Collinearity เป็นปัญหาสำคัญที่พบได้บ่อย ในแบบจำลองการถดถอยเชิงพหุ (Multiple Regression) เมื่อจากตัวแปรอิสระหลายตัวในข้อมูลมักมีความสัมพันธ์กัน ซึ่งพบใช้ค่าสัมประสิทธิ์ของตัวแปรเหล่านั้นไม่สอดคล้องกัน ทำให้ได้การตีความที่ผิดพลาด และการทำนายที่ไม่แม่นยำ การแก้ไขปัญหานี้จึงมีความสำคัญในการปรับปรุงความถูกต้องของแบบจำลองโดยวิธีวิธีต่าง ๆ เช่น Principal Component Regression และ Ridge Regression จะช่วยลดปัญหา Collinearity และเพิ่มความเสถียรของค่าสัมประสิทธิ์ในแบบจำลอง

**สรุป** การวิเคราะห์ข้อมูลที่มีปัญหา Collinearity เป็นกระบวนการที่ต้องใช้เครื่องมือที่เหมาะสมเพื่อให้ได้ผลลัพธ์ที่ถูกต้อง และเชื่อถือได้ (PCA) และ Ridge Regression เป็นสองวิธีที่สามารถจัดการกับปัญหานี้ได้อย่างมีประสิทธิภาพ การเลือกใช้วิธีใดวิธีหนึ่ง ควรขึ้นอยู่กับเป้าหมายของการวิเคราะห์ ความซับซ้อนของข้อมูล และความสามารถในการตีความผลลัพธ์

★ ทงเลือกอื่นในการจัดการกับ Collinearity นอกจาก Principal Component Regression (PCR) และ Ridge Regression และยังมีวิธีอื่นที่สามารถใช้ในการจัดการกับปัญหา Collinearity

เช่น ① Lasso Regression ② Elastic Net Regression

③ Partial Least Squares Regression

ส่วนจากหนังสือ PCR และ Ridge Regression เป็นสองวิธีที่นิยมใช้ในการแก้ปัญหา Collinearity

- PCR เหมาะสำหรับการแก้ที่ลดจำนวนตัวแปรและลดผลกระทบของ Collinearity

โดยใช้ตัวแปรใหม่ที่สร้างขึ้นจากองค์ประกอบหลัก

- Ridge Regression เหมาะสำหรับการแก้ที่ลดการบิดเบือนตัวแปรอิสระทั้งหมดไว้ในแบบจำลอง

และลดผลกระทบจาก Collinearity โดยใช้ตัวแปรใหม่ที่สร้างขึ้นจากองค์ประกอบหลัก

• Ridge Regression เหมาะสำหรับการแก้ที่ลดการบิดเบือนตัวแปรอิสระทั้งหมดไว้ในแบบจำลองและลดผลกระทบจาก Collinearity โดยเพิ่มพารามิเตอร์การปรับ

• นอกจากวิธียังมีวิธีการทางเลือกอื่นๆ เช่น Lasso Regression, Elastic Net และ PLS ที่สามารถใช้ในการแก้ปัญหา Collinearity ได้เช่นกัน

การประเมินผลของโมเดลมีความสำคัญอย่างมาก ซึ่งสามารถทำได้โดยการแบ่งชุดข้อมูลทดสอบที่ใช้ cross-validation และการวิเคราะห์ Residual

การใช้เทคนิคที่เบาะสลในการแก้ปัญหา Collinearity จะช่วยทำให้โมเดลมีความแม่นยำและเสถียรมากขึ้น รวมถึงการนำไปใช้งานในทางพยากรณ์หรือการวิเคราะห์ข้อมูลเชิงธุรกิจ