

การเปรียบเทียบวิธีการประมาณค่า สู่หาบในการวิเคราะห์ตัวแปรพหุ



นางสาวพรศิริ หมื่นไวยงค์

ศูนย์วิทยทรัพยากร

วิทยานิพนธ์นี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตรปริญญาตรีศึกษาศาสตร์มหาบัณฑิต

ภาควิชาสถิติ

บัณฑิตวิทยาลัย จุฬาลงกรณ์มหาวิทยาลัย


พ.ศ. 2529

ISBN 974-566-475-8

013539

I16683040

A COMPARISON OF MISSING VALUES ESTIMATION METHODS IN
MULTIVARIATE ANALYSIS



Miss Pornsiri Muenchaisri

A Thesis Submitted in Partial Fulfillment of the Requirements

for the Degree of Master of Science

Department of Statistics

Graduate School

Chulalongkorn University

1986

ISBN 974-566-475-8

หัวข้อวิทยานิพนธ์ การเปรียบเทียบวิธีการประมาณค่าสู่เหย้าในการวิเคราะห์ตัวแปรพหุ
โดย นางสาว พรศิริ หมื่นไชยศรี
ภาควิชา สถิติ
อาจารย์ที่ปรึกษา รองศาสตราจารย์ ดร. สรชัย พิศาลบุตร



บัณฑิตวิทยาลัย จุฬาลงกรณ์มหาวิทยาลัย อนุมัติให้บัณฑิตวิทยาลัยฉบับนี้เป็นส่วนหนึ่งของ
การศึกษาตามหลักสูตรปริญญามหาบัณฑิต

.....
(รองศาสตราจารย์ ดร. สรชัย พิศาลบุตร)
รักษาการในตำแหน่งรองคณบดีฝ่ายวิชาการ
ปฏิบัติราชการแทนรักษาการในตำแหน่ง คณบดีบัณฑิตวิทยาลัย

คณะกรรมการสอบวิทยานิพนธ์

..... ประธานกรรมการ
(รองศาสตราจารย์ ดร. สุชาติ ธีระนันท์)

..... กรรมการ
(รองศาสตราจารย์ ดร. สรชัย พิศาลบุตร)

..... กรรมการ
(รองศาสตราจารย์ กัลยา ครองแก้ว)

..... กรรมการ
(อาจารย์ ดร. สุพล ตูรงค์วัฒนา)

ลิขสิทธิ์ของบัณฑิตวิทยาลัย จุฬาลงกรณ์มหาวิทยาลัย

หัวข้อวิทยานิพนธ์ การเปรียบเทียบวิธีการประมาณค่าสูญเสียในการวิเคราะห์ตัวแปรพหุ

ชื่อผลิต นางสาว พรศิริ หมื่นไวยยศร์

อาจารย์ที่ปรึกษา รองศาสตราจารย์ ดร. ลระชัย พิศาลบุตร

ภาควิชา สถิติ

ปีการศึกษา 2528



บทคัดย่อ

ในการวิเคราะห์ตัวแปรพหุ หากมีปัญหาเกี่ยวกับข้อมูลบางตัวสูญเสียจะทำให้ไม่สามารถวิเคราะห์ข้อมูลได้ วิธีการแก้ปัญหอย่างหนึ่งก็คือ ตัดค่าสังเกตขุดทิ้งทั้งไป แต่การแก้ปัญหโดยวิธีนี้จะมีผลทำให้จำนวนค่าสังเกตน้อยลง และสูญเสียรายละเอียดของข้อมูลบางตัวไป วิธีการแก้ปัญหอีกวิธีหนึ่งก็คือ ต้องทำการประมาณค่าสูญหายนั้น แต่เนื่องจากวิธีการประมาณค่าสูญหายหลายวิธีซึ่งแต่ละวิธีต่างก็มีข้อดีและข้อเสียแตกต่างกันไป ดังนั้นการวิจัยนี้จึงสนใจเปรียบเทียบวิธีการประมาณค่าสูญหายที่นิยมใช้กันทั่วไป 4 วิธีคือ วิธีใช้ค่าเฉลี่ย วิธีวิเคราะห์ความถดถอยพหุเชิงเส้น วิธีวิเคราะห์ความถดถอยพหุเชิงเส้นตัดแปลง และวิธีวิเคราะห์ส่วนประกอบหลัก โดยใช้ค่าความคลาดเคลื่อนเฉลี่ยเป็นเกณฑ์ในการเปรียบเทียบสถานการณ์ต่าง ๆ ซึ่งจำลองการทดลองขึ้นโดยใช้เทคนิคมอนติคาร์โล แต่ละสถานการณ์ต่างกันขึ้นอยู่กับขนาดตัวอย่าง $n = 30, 50, 70, 100, 200$ จำนวนตัวแปร $p = 3, 5, 7, 10$ และขนาดความสัมพันธ์ระหว่างตัวแปร $\rho = 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9$ แต่เนื่องจากบางช่วงของการทดลอง ถ้า ρ มีค่าแตกต่างกัน ผลการเรียงลำดับของค่าความคลาดเคลื่อนเฉลี่ยเหมือนกัน ดังนั้นผู้วิจัยจึงไม่ได้นำสถานการณ์เหล่านั้นมาเล่นในวิทยานิพนธ์ แต่จะเล่นสถานการณ์ที่แตกต่างกันเพียง 106 สถานการณ์เท่านั้น

จากการวิจัยที่ระดับนัยสำคัญ 0.05 พบว่าวิธีการประมาณค่าสูญหายทั้ง 4 วิธีให้ค่าความคลาดเคลื่อนเฉลี่ยไม่แตกต่างกันอย่างมีนัยสำคัญ ดังนั้นจึงกล่าวได้ว่า ไม่ว่าจะ เป็นสถานการณ์ใดก็ตามที่กำหนดเหล่านี้ ถ้าหากมีข้อมูลสูญหายเกิดขึ้น สามารถเลือกวิธีการประมาณค่าสูญหายวิธีใดก็ได้ใน 4 วิธีนี้ แต่วิธีการประมาณค่าสูญหายที่ง่ายที่สุดและใช้เวลาในการประมวล

ผลน้อยที่สุดคือวิธีตัวเฉลี่ย ซึ่งเป็นวิธีการประมาณค่าสูญหายที่จะทำให้ได้ค่าความคลาดเคลื่อนเฉลี่ย ไม่แตกต่างไปจากการใช้วิธีการประมาณอีก 3 วิธีที่เหลือ แต่อย่างไรก็ตาม ถ้าพิจารณาให้ละเอียด ในแต่ละสถานการณ์ เมื่อเปรียบเทียบค่าความคลาดเคลื่อนเฉลี่ยแล้วพบว่ามีความแตกต่างกัน แม้ว่าจะ ไม่แตกต่างกันอย่างมีนัยสำคัญก็ตาม แต่ในการประมาณต่าง ๆ ผู้วิจัยต้องพยายามทำให้ค่าความ คลาดเคลื่อนเฉลี่ยมีค่าน้อยที่สุด อาจกล่าวได้ว่า ถ้า $p = 3$ วิธีที่ใช้ค่าเฉลี่ยจะดีที่สุดเมื่อ $\rho = 0.1$ วิธีวิเคราะห์ความถดถอยพหุเชิงเส้นตัดแปลง จะดีที่สุดเมื่อ $\rho = 0.2 - 0.7$ วิธีวิเคราะห์ ส่วนประกอบหลัก จะดีที่สุดเมื่อ $\rho = 0.9$ ถ้า $p = 5$ วิธีที่ใช้ค่าเฉลี่ย จะดีที่สุดเมื่อ $\rho = 0.1 - 0.2$ วิธีวิเคราะห์ความถดถอยพหุเชิงเส้นตัดแปลง จะดีที่สุดเมื่อ $\rho = 0.3$ วิธีวิเคราะห์ส่วนประกอบหลักจะดีที่สุดเมื่อ $\rho = 0.5 - 0.9$ ถ้า $p = 7$ วิธีที่ใช้ค่าเฉลี่ยจะ ดีที่สุดเมื่อ $\rho = 0.1 - 0.2$ วิธีวิเคราะห์ความถดถอยพหุเชิงเส้นตัดแปลงจะดีที่สุดเมื่อ $\rho = 0.3 - 0.4$ วิธีวิเคราะห์ส่วนประกอบหลักจะดีที่สุดเมื่อ $\rho = 0.5 - 0.8$ ถ้า $p = 10$ วิธีวิเคราะห์ส่วนประกอบหลักจะดีที่สุดเมื่อ $\rho = 0.2 - 0.5$

ศูนย์วิทยทรัพยากร
จุฬาลงกรณ์มหาวิทยาลัย

Thesis Title A comparison of missing values estimation methods in
multivariate analysis.

Name Miss Pornsiri Muenchaisri

Thesis Advisor Associate Professor Sorachai Bhisalbutra , Ph.D.

Department Statistics

Academic year 1985



ABSTRACT

The purpose of this study is to investigate the four well known missing value estimation methods in multivariate analysis namely,
1) Mean 2) Multiple Linear Regression 3) Modified Multiple Linear Regression 4) Principal Component, using mean square errors as means of comparison.

The data for each experiment were obtained through simulation using the Monte Carlo technique. The computer program was designed to calculate the mean square error for each methods in different situations with varying sample size $n = 30 \ 50 \ 70 \ 100 \ 200$ number of variables $p = 3 \ 5 \ 7 \ 10$ and correlation coefficient $\rho = 0.1 \ 0.2 \ 0.3 \ 0.4 \ 0.5 \ 0.6 \ 0.7 \ 0.8 \ 0.9$ However, some intervals of ρ provide the same ranking results of mean square error. These situations are omitted and thus only 106 different situations are presented in the thesis.

The result of this study shows that, at 5% level of significance, mean square errors of the four methods are not significantly different. So if there exists the missing value problem, any one of these methods can be used to estimate the missing value. Nevertheless, the easiest method which also uses least processing time is the first method, Mean.

In addition, attempt to obtain the smallest mean square error is made by considering for each situation which method has the smallest mean square error. The results are up to n , p and ρ . Conclusively, if the number of variables are three, Mean is the best when $\rho = 0.1$, Modified multiple linear regression is the best when $\rho = 0.2 - 0.7$ and Principal component is the best when $\rho = 0.9$. If the number of variables are five, Mean is the best when $\rho = 0.1 - 0.2$ Modified multiple linear regression is the best when $\rho = 0.3$ and Principal component is the best when $\rho = 0.5 - 0.9$. If the number of variables are seven Mean is the best when $\rho = 0.1 - 0.2$, Modified multiple linear regression is the best when $\rho = 0.3 - 0.4$ and Principal component is the best when $\rho = 0.5 - 0.8$. If the number of variables are ten, Principal component is the best when $\rho = 0.2 - 0.5$.



ศูนย์วิทยทรัพยากร
จุฬาลงกรณ์มหาวิทยาลัย



กิตติกรรมประกาศ

วิทยานิพนธ์ฉบับนี้สำเร็จลุล่วงได้ด้วยความกรุณาของรองศาสตราจารย์ ดร. สรชัย พิศาบุตร คณะพาณิชย์ศาสตร์และการบัญชี จุฬาลงกรณ์มหาวิทยาลัย ที่ให้คำแนะนำ ปรึกษา ตลอดจนแก้ไขข้อบกพร่องต่าง ๆ เป็นอย่างดียิ่ง ผู้วิจัยขอกราบขอบพระคุณเป็นอย่างสูงไว้ ณ โอกาสนี้

ขอกราบขอบพระคุณท่านอาจารย์ทุก ๆ ท่านที่ได้ประสิทธิ์ประสาทวิชาความรู้ให้แก่ผู้วิจัย มาโดยตลอด

ขอขอบคุณเจ้าหน้าที่สำนักงานบริการคอมพิวเตอร์ จุฬาลงกรณ์มหาวิทยาลัยทุกท่านที่ช่วยประมวลผลงานวิจัยตั้งแต่ต้นจนแล้วเสร็จ

ขอขอบคุณ เพื่อน ๆ และพี่ ๆ ทุกคนที่ให้ความช่วยเหลือ และคำแนะนำต่าง ๆ

สุดท้ายนี้ ขอกราบขอบพระคุณ คุณพ่อ คุณแม่ และพี่สาวที่ช่วยส่งเสริมและสนับสนุนการเรียบเรียงของผู้นี้ตลอดมา

พรศิริ หมั่นไชยศิริ

ศูนย์วิทยทรัพยากร
จุฬาลงกรณ์มหาวิทยาลัย

สารบัญ



หน้า

บทคัดย่อภาษาไทย	ง
บทคัดย่อภาษาอังกฤษ	จ
กิตติกรรมประกาศ	ข
สารบัญตาราง	ฅ
บทที่	
1. บทนำ	1
2. ทฤษฎีที่ใช้ในการวิจัย	6
3. ระเบียบวิธีวิจัย	21
4. ผลการวิจัย	30
5. สรุปผลการวิจัยและข้อเสนอแนะ	75
บรรณานุกรม	79
ภาคผนวก	81
ประวัติ	115

ศูนย์วิทยทรัพยากร
จุฬาลงกรณ์มหาวิทยาลัย

สารบัญตาราง

		หน้า
ตาราง 1	ผลการวิเคราะห์ความแปรปรวนกรณีที่มีตัวแปรตาม 4 ตัว	32
ตารางที่ 2	ผลการวิเคราะห์ความแปรปรวนซึ่งเป็นผลจากขนาดตัวอย่าง จำนวนตามวิธีที่ใช้ประมาณ	33
ตารางที่ 3	ผลการวิเคราะห์ความแปรปรวนซึ่งเป็นผลจากจำนวนตัวแปร จำนวนตามวิธีที่ใช้ประมาณ	33
ตารางที่ 4	ผลการวิเคราะห์ความแปรปรวนซึ่งเป็นผลจากขนาดความสัมพันธ์ ระหว่าง 2 ตัวแปรใด ๆ จำนวนตามวิธีที่ใช้ประมาณ	34
ตารางที่ 5	ผลการวิเคราะห์ความแปรปรวนซึ่งเป็นผลจากจำนวนตัวแปรและ ขนาดความสัมพันธ์ระหว่าง 2 ตัวแปรใด ๆ จำนวนตามวิธีที่ใช้ ประมาณ	35
ตารางที่ 6	ผลการวิเคราะห์ความแปรปรวนซึ่งเป็นผลจากจำนวนตัวแปรและ ขนาดตัวอย่างจำนวนตามวิธีที่ใช้ประมาณ	36
ตารางที่ 7	ผลการวิเคราะห์ความแปรปรวนซึ่งเป็นผลจากขนาดตัวอย่างและ ขนาดความสัมพันธ์ระหว่าง 2 ตัวแปรใด ๆ จำนวนตามวิธีที่ใช้ ประมาณ	36
ตารางที่ 8	ผลการวิเคราะห์ความแปรปรวนทดสอบความแตกต่างของวิธีการ ประมาณค่าสู่สุดท้ายทั้ง 4 วิธี	37
ตารางที่ 9	ผลของจำนวนครั้งของแต่ละวิธีที่ได้อันดับ 1 2 3 4 และ คะแนนรวมถ่วงน้ำหนัก ในกรณีที่ $p = 3$ $n = 30$ $\rho = 0.2$	38
ตารางที่ 10	ค่าความคลาดเคลื่อนเฉลี่ย ของวิธีการประมาณค่าสู่สุดท้ายแบบ ต่าง ๆ เมื่อ $n = 30$ $p = 3$ จำนวนตามค่า ρ	40
ตารางที่ 11	ค่าคะแนนรวมถ่วงน้ำหนัก ของวิธีการประมาณค่าสู่สุดท้ายแบบ ต่าง ๆ เมื่อ $n = 30$ $p = 3$ จำนวนตามค่า ρ	40

ตารางที่ 38	ค่าความคลาดเคลื่อนเฉลี่ยของวิธีการประมาณค่าสู่สูญแบบ ต่าง ๆ เมื่อ $n = 200$ $p = 7$ จำแนกตามค่า p	68
ตารางที่ 39	ค่าคะแนนความถ่วงน้ำหนักของวิธีการประมาณค่าสู่สูญแบบ ต่าง ๆ เมื่อ $n = 200$ $p = 7$ จำแนกตามค่า p	68
ตารางที่ 40	ค่าความคลาดเคลื่อนเฉลี่ยและคะแนนรวมถ่วงน้ำหนักของวิธี การประมาณค่าสู่สูญแบบต่าง ๆ เมื่อ $n = 70$ $p = 10$ จำแนกตามค่า p	70
ตารางที่ 41	ค่าความคลาดเคลื่อนเฉลี่ยของวิธีการประมาณค่าสู่สูญแบบ ต่าง ๆ เมื่อ $n = 100$ $p = 10$ จำแนกตามค่า p	72
ตารางที่ 42	ค่าคะแนนรวมถ่วงน้ำหนักของวิธีการประมาณค่าสู่สูญแบบ ต่าง ๆ เมื่อ $n = 100$ $p = 10$ จำแนกตามค่า p	73
ตารางที่ 43	ค่าความคลาดเคลื่อนเฉลี่ยและคะแนนรวมถ่วงน้ำหนักของวิธี การประมาณค่าสู่สูญแบบต่าง ๆ เมื่อ $n = 200$ $p = 10$ จำแนกตามค่า p	73