

บทที่ 4

การประยุกต์ใช้ NB2 กับข้อมูลจริง

สำหรับการประยุกต์ใช้ NB2 และ Robust standard error ที่เสนอในวิทยานิพนธ์ฉบับนี้เราใช้ข้อมูลซึ่งได้มีการเสนอในผลงานวิจัยที่เกี่ยวข้องกับตัวแบบถดถอยปัวซอง จำนวน 2 ชุด ได้แก่ Fabric fault data จาก Hinde and Demètrio (1998) กับ Quine data ซึ่งได้กล่าวถึงใน Aitkin et al. (1989)

4.1 Fabric fault data

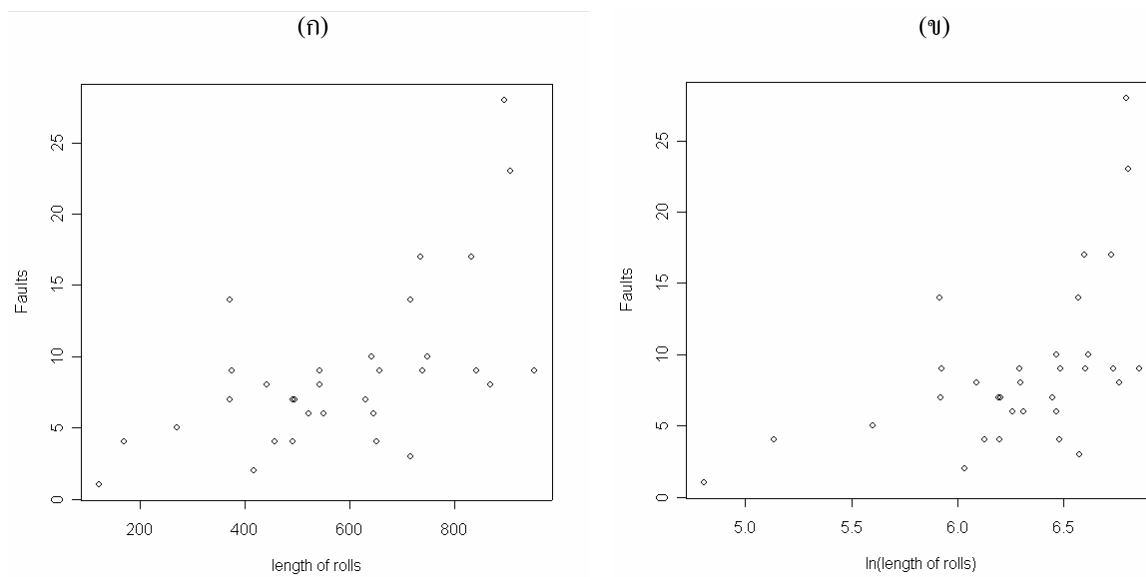
Fabric fault data เป็นข้อมูลจำนวนรอยตำหนิบนผ้า ซึ่งผลิตโดยโรงงานหนึ่ง ในการศึกษาผู้วิจัยได้สุ่มผ้าที่ผลิตจากโรงงานนี้ มาจำนวน 32 ม้วน แต่ละม้วนมีความยาว (หน่วย : ฟุต) แตกต่างกัน ข้อมูลดังกล่าวได้เสนอ ดังตารางที่ 4.1

ตารางที่ 4.1 Fabric fault data

Length of roll	Fault	Length of roll	Fault	Length of roll	Fault	Length of roll	Fault
551	6	543	8	651	4	842	9
832	17	905	23	375	9	542	9
715	14	522	6	868	8	122	1
271	5	657	9	630	7	170	4
491	7	738	9	372	7	371	14
645	6	735	17	441	8	749	10
895	28	495	7	458	4	716	3
642	10	952	9	492	4	417	2

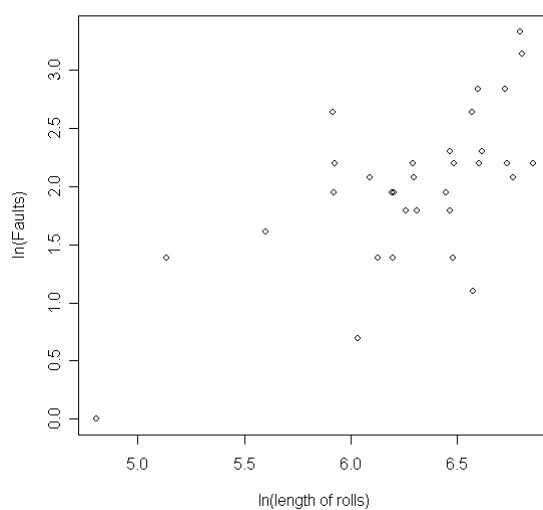
ในเบื้องต้นอาจกล่าวได้ว่าจำนวนรอยตำหนิบนผ้าแต่ละม้วนมีการแจกแจงปัวซอง แผนภาพการกระจายระหว่างความยาวผ้า (Length of roll) และจำนวนรอยตำหนิ (Fault) ซึ่งในกรณีนี้ใช้ข้อมูล

เดิม และแปลงความยาวผ้าเป็นลอการิทึมฐานธรรมชาติ $\ln(\text{length of roll})$ ดังภาพที่ 4.1 (ก) และ (ข) ตามลำดับ



ภาพที่ 4.1 แผนภาพการกระจายระหว่าง Length of roll กับ Faults

แสดงให้เห็นว่าความสัมพันธ์ระหว่างตัวแปรทั้งสองมีแนวโน้มเป็นเส้นโค้งการเติบโตแบบเลขชี้กำลัง (Exponential growth curve) และเมื่อแปลงจำนวนรอยตำหนิให้เป็นมาตรฐานลอการิทึมฐานธรรมชาติ ($\ln(\text{Faults})$) แผนภาพการกระจายระหว่าง $\ln(\text{Length of roll})$ กับ $\ln(\text{Faults})$ ดังภาพที่ 4.2



ภาพที่ 4.2 แผนภาพการกระจายระหว่าง $\ln(\text{Length of roll})$ กับ $\ln(\text{Faults})$

แสดงให้เห็นว่าความสัมพันธ์ระหว่างตัวแปรทั้งสองมีแนวโน้มเป็นเส้นตรง อย่างไรก็ตามในวิทยานิพนธ์ฉบับนี้จะได้ศึกษาการใช้การแจกแจงปกติประมาณการแจกแจงปัวซอง แต่จะใช้การแจกแจงปัวซองโดยตรง

เนื่องจากจำนวนรอยตำหนิซึ่งแทนด้วย Faults มีการแจกแจงปัวซองที่มี Link function คือ $\ln(\mu_i)$ เมื่อ μ_i เป็นค่าเฉลี่ยของจำนวนรอยตำหนิจากผ้าม้วนที่ i , $i = 1, 2, \dots, 32$ ดังนั้นตัวแบบถดถอยปัวซองของจำนวนรอยตำหนิ ซึ่งมี $ll = \ln(\text{Length of roll})$ เป็นตัวแปรอิสระ คือ

$$\ln(\mu_i) = \beta_0 + \beta_1 ll_i \quad (4.1)$$

และตัวแบบถดถอยปัวซองประมาณของ (4.1) คือ

$$\ln(\hat{\mu}_i) = \hat{\beta}_0 + \hat{\beta}_1 ll_i \quad (4.2)$$

ค่าประมาณความควรจะเป็นสูงสุด $\hat{\beta}_0$ และ $\hat{\beta}_1$ ค่าคลาดเคลื่อนมาตรฐาน (Standard error) และสถิติที่จำเป็นได้แสดงในตารางที่ 4.2 แถวที่ 1 – 2

ตารางที่ 4.2 ค่าสถิติต่าง ๆ สำหรับ Fabric fault data

ตัวแบบ ถดถอย	linear pred.coeff			Overdispersion		Score test NB2		Residual deviance	AIC	df
	$\hat{\beta}_0$	$\hat{\beta}_1$	$\hat{\alpha}$	Dean (1992)	Wang-Shu Lu (1997)					
ปัวซอง	-4.171 (1.135)	0.997 (0.176)	0 -	5.635	5.643	64.55	191.85	30		
NB2	-3.794 (1.419) (1.456) ⁺	0.937 (0.222) (0.230) ⁺	0.115 (0.055) (0.048) ⁺			30.67	181.39	29		

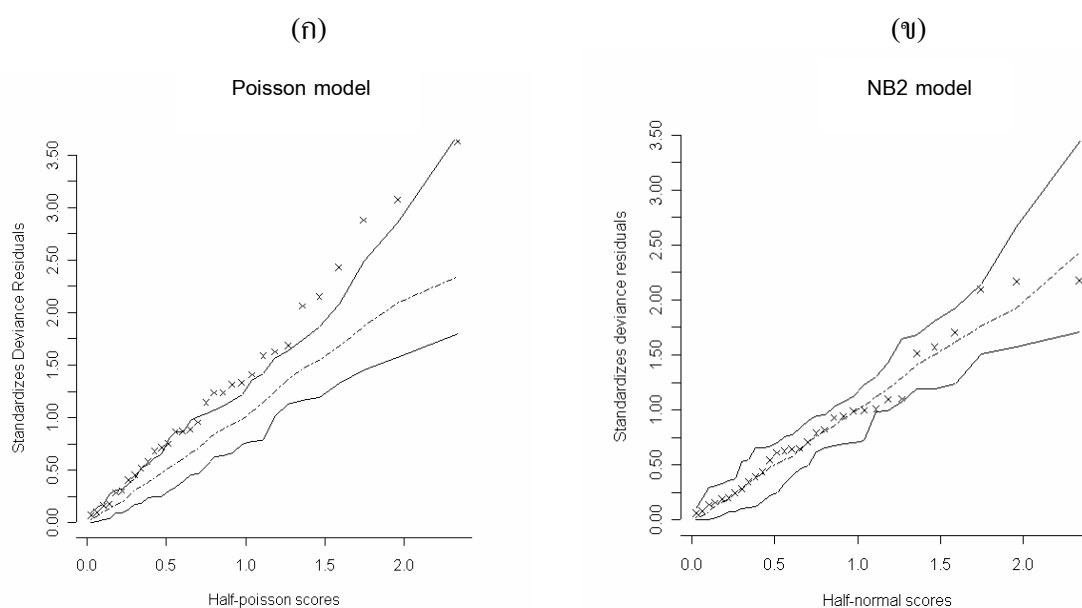
หมายเหตุ + คือ Robust standard error

ตัวแบบถดถอยปัวซอง

$$\ln(\mu_i) = -4.171 + 0.997 ll_i \quad (4.3)$$

(1.135) (0.176)

ตัวแบบถดถอยปัวซองประมาณของ (4.2) ให้ค่า Residual deviance เท่ากับ 64.55 และ df เท่ากับ 30 จะเห็นว่า ค่า Residual deviance มีค่ามากเป็นสองเท่าของ df เป็นดัชนีที่บ่งบอกว่าตัวแบบถดถอยปัวซองไม่เหมาะสมกับข้อมูล เนื่องจากเกิด Overdispersion นอกจากนี้ ค่า Adjusted score test ของ Dean (1992) และ Wang – Shu Lu (1997) ซึ่งแสดงในตารางที่ 4.2 ในสดมภ์ 5 และ 6 ตามลำดับ รวมทั้ง Half normal plot with simulated envelope ภาพที่ 4.3 (ก) ก็เป็นสถิติที่สนับสนุนข้อสรุปดังกล่าว ด้วยเหตุนี้เราจึงทำการประมาณค่า Overdispersion โดยการ fit ตัวแบบถดถอย NB2 กับข้อมูล พร้อมกับคำนวณค่า Robust standard error เพื่อเปรียบเทียบกับ Asymptotic standard error ของค่าประมาณที่ได้จากกระบวนการสุดท้ายของกระบวนการวนซ้ำนิวตัน – ราฟสัน ที่ได้กล่าวไปแล้วในบทที่ 3 โดยสถิติที่กล่าวถึงนี้ ได้แสดงในตารางที่ 4.2 แถวที่ 3 – 5 และตัวแบบการถดถอย NB2 ที่ได้คือ



ภาพที่ 4.3 แผนภาพ Half normal plot with simulated envelope ของตัวแบบถดถอยปัวซอง และ NB2

$$\ln(\mu_i) = -3.794 + 0.9371I_i; \quad \hat{\alpha} = 0.115 \quad (4.4)$$

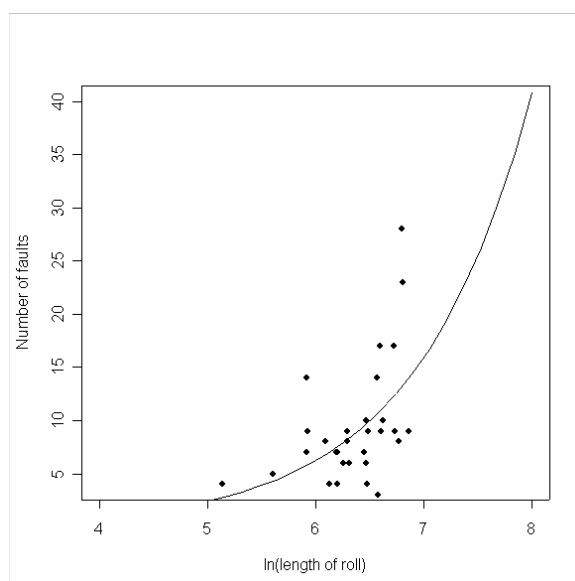
(1.419) (0.222) (0.055)

นอกจากนี้ค่า Robust standard error มีค่าใกล้เคียงกับค่า Asymptotic standard error แสดงว่าการแจกแจงของจำนวนรอยตำหนิที่ผลิตโดยโรงงานนี้มีการแจกแจงแบบ NB2 และ Half normal plot with simulated envelope ภาพที่ 4.3 (ข) ยืนยันว่าตัวแบบ (4.4) สอดคล้องกับข้อสมมติสำหรับ NB2

ดังนั้นตัวแบบที่แปลงให้อยู่ในมาตราชี้กำลัง สำหรับอัตราการเกิดรอยตำหนิที่ความยาวผ้า (ที่แปลงให้เป็นมาตราลอการิทึมฐานธรรมชาติ) ที่ศึกษา คือ

$$\hat{\mu}_i = \exp[-3.794 + 0.937\ln_i] \quad (4.5)$$

และเส้นโค้งการเติบโตแบบเลขชี้กำลังที่ได้จากการแทนค่า Length of roll ในสมการ (4.5) ได้ แผนภาพการกระจาย ดังภาพที่ 4.6



ภาพที่ 4.4 แผนภาพแสดงเส้นโค้งการเติบโตแบบเลขชี้กำลังที่ได้จากการแทนค่า $\ln(\text{Length of roll})$ ในสมการ (4.5)

4.2 Quine data

Quine data ซึ่งได้กล่าวถึงใน Aitkin et al. (1989) เป็นข้อมูลเกี่ยวกับจำนวนวันการขาดเรียนของนักเรียน ซึ่งได้นำข้อมูลมาจากการศึกษาด้านสังคมศาสตร์ในประเทศออสเตรเลียของนักเรียนจำนวน 146 คน โดยตัวแปรที่รวบรวมมี ดังนี้

จำนวนวันการขาดเรียน แทนด้วย Days

อายุ แทนด้วย A มี 4 ระดับ คือ

1: final grade in primary schools; 2: first form in secondary schools

3: second form in secondary schools; 4: third form in secondary schools

เพศ แทนด้วย S มี 2 ระดับ คือ 1: ชาย และ 2: หญิง

ชนชาติ (Cultural group) แทนด้วย C มี 2 ระดับ คือ 1: ชนพื้นเมือง; 2: คนผิวขาว

ระดับการเรียนรู้ (Learning rate) แทนด้วย L มี 2 ระดับ คือ 1: ช้า; 2: ปานกลาง
ข้อมูลดังกล่าวได้เสนอดังตารางที่ 4.3

ตารางที่ 4.3 Quine data

Days	C	S	A	L	Days	C	S	A	L	Days	C	S	A	L
2	1	1	1	1	17	1	2	2	2	10	2	1	4	2
5	1	1	1	2	14	1	2	3	1	41	2	1	4	2
6	1	1	2	1	60	1	2	3	1	11	2	2	1	2
14	1	1	2	2	2	1	2	4	2	1	2	2	2	1
57	1	1	3	1	14	1	2	4	2	5	2	2	2	1
17	1	1	3	2	6	2	1	1	1	11	2	2	2	1
8	1	1	4	2	0	2	1	1	2	6	2	2	2	2
34	1	1	4	2	12	2	1	1	2	0	2	2	3	1
5	1	2	1	2	5	2	1	2	1	5	2	2	3	1
5	1	2	2	1	3	2	1	2	2	14	2	2	3	1
13	1	2	2	1	36	2	1	3	1	3	2	2	4	2
53	1	2	2	1	7	2	1	3	2	18	2	2	4	2
11	1	2	2	2	0	2	1	4	2	5	1	1	1	2
13	1	2	3	1	30	2	1	4	2	22	1	1	1	2
48	1	2	3	1	10	2	2	1	2	7	1	1	2	2
0	2	1	1	2	0	2	2	2	1	53	1	1	3	1
10	1	2	4	2	5	2	2	2	1	16	1	1	3	2
40	1	2	4	2	7	2	2	2	1	46	1	1	3	2
0	2	1	1	2	6	2	2	2	2	28	1	1	4	2
11	2	1	1	2	28	2	2	2	2	3	1	2	1	1
5	2	1	2	1	3	2	2	3	1	45	1	2	1	2
17	2	1	2	1	12	2	2	3	1	9	1	2	2	1
30	2	1	3	1	3	2	2	4	2	32	1	2	2	1

ตารางที่ 4.3 (ต่อ)

Days	C	S	A	L	Days	C	S	A	L	Days	C	S	A	L
5	2	1	3	2	15	2	2	4	2	5	1	2	2	2
27	2	1	3	2	37	2	2	4	2	8	1	2	3	1
27	2	1	4	2	14	1	1	1	1	47	1	2	3	1
25	2	2	1	1	20	1	1	1	2	2	1	2	3	2
33	2	2	1	2	15	1	1	2	1	5	1	2	4	2
5	2	2	2	1	32	1	1	3	1	36	1	2	4	2
7	2	2	2	1	16	1	1	3	2	67	2	1	1	1
5	2	2	2	2	43	1	1	3	2	7	2	1	1	2
14	2	2	2	2	23	1	1	4	2	0	2	1	2	1
2	2	2	3	1	38	1	1	4	2	11	2	1	2	1
10	2	2	3	1	24	1	2	1	2	22	2	1	3	1
1	2	2	4	2	6	1	2	2	1	1	2	1	3	2
9	2	2	4	2	25	1	2	2	1	16	2	1	3	2
22	2	2	4	2	5	1	2	2	1	14	2	1	4	2
11	1	1	1	1	19	1	2	2	2	69	2	1	4	2
13	1	1	1	2	20	1	2	3	1	20	2	2	1	2
6	1	1	2	1	81	1	2	3	1	5	2	2	2	1
6	1	1	3	1	3	1	2	4	2	5	2	2	2	1
14	1	1	3	2	21	1	2	4	2	15	2	2	2	1
40	1	1	3	2	17	2	1	1	1	7	2	2	2	2
23	1	1	4	2	2	2	1	1	2	2	2	2	3	1
36	1	1	4	2	0	2	1	2	1	8	2	2	3	1
11	1	2	1	2	5	2	1	2	1	1	2	2	3	2
6	1	2	2	1	4	2	1	2	2	5	2	2	4	2
23	1	2	2	1	0	2	1	3	2	22	2	2	4	2
54	1	2	2	1	8	2	1	3	2					

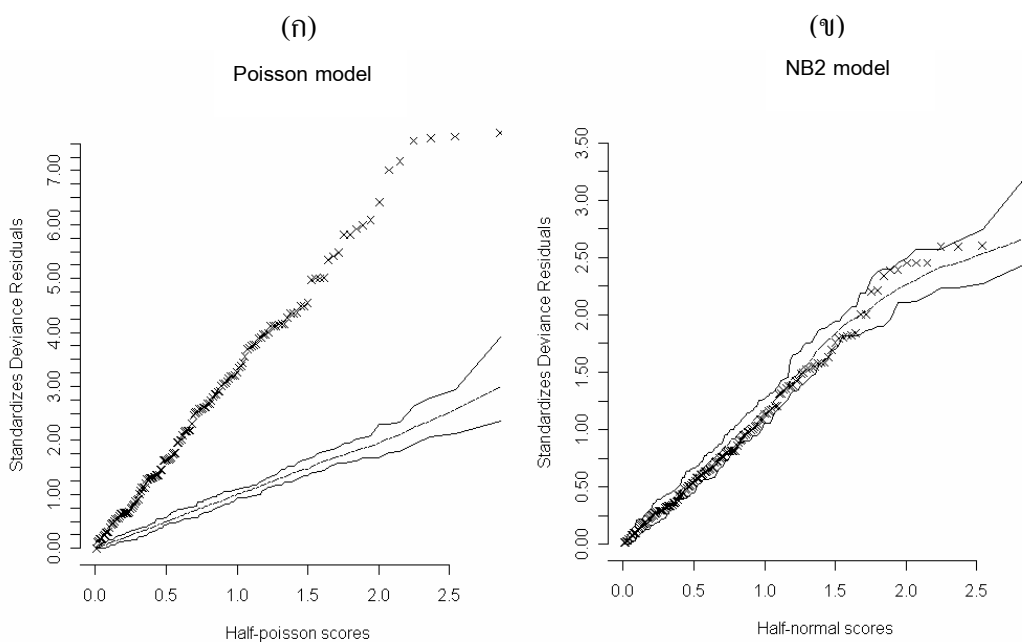
ในเบื้องต้นอาจกล่าวได้ว่าจำนวนวันขาดเรียน มีการแจกแจงปัวซองที่มี Link function คือ $\ln(\mu)$ เมื่อ μ เป็นเวกเตอร์ค่าเฉลี่ยของ Days ที่มีตัวแปรอิสระคือ A , S , C และ L ตัวแบบถดถอยปัวซองของ Days ที่สำคัญ พร้อมค่าสถิติ แสดงในตารางที่ 4.4 ดังนี้

ตารางที่ 4.4 Residual deviance, df, AIC และค่า Score Test ของตัวแบบถดถอยปัวซอง

Model $\ln(\mu)$	Residual deviance	df	AIC	Score Test	
				Dean(1992)	Wang-Shu Lu(1997)
$C*S*A*L$	1166.3	118	1810.8	65.409	65.680
$C*S*A$	1434.1	130	2054.6	80.640	80.614
$C*S*L$	1746.1	138	2350.6	103.070	103.061
$C*A*L$	1455.0	132	2071.4	80.565	80.671
$S*A*L$	1631.5	132	2248.0	94.659	94.811
$C+S+A+L$	1674.9	139	2277.4	95.249	95.210
$(C+S)*A+C*S*L$	1301.0	129	1923.5	71.769	71.666

จากตารางที่ 4.4 จะเห็นว่า Full interaction model ($C*S*A*L$) มีค่า AIC น้อยที่สุด เท่ากับ 1810.8 ดังนั้น Full interaction model น่าเป็นตัวแทนที่เหมาะสมที่สุด แต่ Full interaction model มีค่า residual deviance เท่ากับ 1166.3 มีค่ามากกว่า df ของมันเกือบ 10 เท่า อีกทั้งค่า Score test ของ Dean(1992) และ Wang - Shu Lu(1997) มีค่ามากกว่า $Z_{0.05} = 1.65$ ก็แสดงให้เห็นว่า ตัวแบบถดถอยปัวซองไม่สอดคล้องกับข้อมูลเนื่องจากเกิด Overdispersion นอกจากนี้ Half normal plot with simulated envelope ของตัวแบบถดถอยปัวซอง รูปที่ 4.5 (ก) แสดงว่าข้อมูลจำนวนวันขาดเรียนเป็น Overdispersed Poisson counts ด้วยเหตุนี้เราจึงทำการประมาณค่า Overdispersion โดยการ fit ตัวแบบถดถอย NB2 โดยพิจารณาตัวแทนที่เป็นไปได้ทั้งหมด และสรุปสำหรับตัวแทนที่สำคัญโดยพิจารณาจากค่า AIC ดังแสดงในตารางที่ 4.5

เมื่อ fit ตัวแบบ NB2 พบว่าตัวแทนที่เหมาะสมกับข้อมูลมากที่สุดคือ $\ln(\mu_{\text{DAYS}}) = (C+S)*A+C*S*L$ โดยมีค่า Residual deviance เท่ากับ 167.91 และ df ที่มีค่าเท่ากับ 128 ซึ่งเป็นตัวแทนที่มีค่า AIC ต่ำสุด เท่ากับ 1092.4 โดยตัวแทนดังกล่าวเขียนให้อยู่ในรูปฟังก์ชันของตัวแปรอิสระได้ ดังสมการ (4.5)



ภาพที่ 4.5 แผนภาพ Half normal plot with simulated envelope ของตัวแบบถดถอยปัวซอง และ NB2

ตารางที่ 4.5 ค่าสถิติต่าง ๆ ของตัวแบบถดถอย NB2 สำหรับ Quine data

Model $\ln(\mu)$	Residual deviance	df	AIC
$\underline{C}*\underline{S}*\underline{A}*\underline{L}$	167.27	117	1095.1
$\underline{C}*\underline{S}*\underline{A}$	167.76	129	1103.4
$\underline{C}*\underline{S}*\underline{L}$	167.84	137	1112.9
$\underline{C}*\underline{A}*\underline{L}$	167.87	131	1102.7
$\underline{S}*\underline{A}*\underline{L}$	167.21	131	1114.6
$\underline{C}+\underline{S}+\underline{A}+\underline{L}$	167.93	138	1107.4
$(\underline{C}+\underline{S})*\underline{A}+\underline{C}*\underline{S}*\underline{L}$	167.91	128	1092.4

$$\begin{aligned} \ln(\hat{\mu}) = & \hat{\beta}_0 + \hat{\beta}_1 C_2 + \hat{\beta}_2 S_2 + \hat{\beta}_3 A_2 + \hat{\beta}_4 A_3 + \hat{\beta}_5 A_4 + \hat{\beta}_6 L_2 + \hat{\beta}_7 C_2 : A_2 + \hat{\beta}_8 C_2 : A_3 \\ & + \hat{\beta}_9 C_2 : A_4 + \hat{\beta}_{10} S_2 : A_2 + \hat{\beta}_{11} S_2 : A_3 + \hat{\beta}_{12} S_2 : A_4 + \hat{\beta}_{13} C_2 : S_2 + \hat{\beta}_{14} C_2 : L_2 \\ & + \hat{\beta}_{15} S_2 : L_2 + \hat{\beta}_{16} C_2 : S_2 : L_2 \end{aligned} \quad (4.5)$$

เมื่อ A_2, A_3, A_4, C_2, L_2 และ S_2 เป็นตัวแปรดัมมี่ ที่

$$\begin{aligned} A_2 = & \begin{cases} 1: & \text{first form in secondary schools} \\ 0: & \text{others} \end{cases} \\ A_3 = & \begin{cases} 1: & \text{second form in secondary schools} \\ 0: & \text{others} \end{cases} \\ A_4 = & \begin{cases} 1: & \text{third form in secondary schools} \\ 0: & \text{others} \end{cases} \end{aligned}$$

$$C_2 = \begin{cases} 1: & \text{ชนผิวขาว} \\ 0: & \text{ชนพื้นเมือง} \end{cases} \quad L_2 = \begin{cases} 1: & \text{เรียนรู้ปานกลาง} \\ 0: & \text{เรียนรู้ช้า} \end{cases} \quad S_2 = \begin{cases} 1: & \text{หญิง} \\ 0: & \text{ชาย} \end{cases}$$

ตัวแบบ (4.5) มีตัวแปรอิสระน้อยกว่าตัวแบบถดถอยปัวซอง ซึ่งสนับสนุนข้อเท็จจริงที่ว่า ถ้าเราไม่นำ Overdispersion มาพิจารณาในการวิเคราะห์ข้อมูล ทำให้การนำตัวแปรอิสระเข้าหรือออกจากตัวแบบมีค่ามากเกินไป อาจทำให้เลือกตัวแบบที่มีความซับซ้อนมากเกินไป ค่าประมาณพารามิเตอร์ของตัวแบบถดถอย NB2 (4.5) พร้อมค่าประมาณ Asymptotic standard error และ Robust standard error ได้แสดงในตารางที่ 4.6 ตัวแบบนี้ให้ค่า Residual deviance เท่ากับ 167.91 ซึ่งมีค่าใกล้เคียงกับ Residual df ซึ่งมีค่าเท่ากับ 128 นอกจากนี้ค่า Robust standard error มีค่าใกล้เคียงกับค่า Asymptotic standard error แสดงว่า NB2 เป็นตัวแบบสำหรับ Quine Data และ Half normal plot with simulated envelope ภาพที่ 4.5 (ข) สนับสนุนว่า ตัวแบบ (4.5) เหมาะสมกับข้อมูล

ตารางที่ 4.6 ค่าสถิติต่าง ๆ ตัวแบบถดถอย NB2 สำหรับ Quine data

ตัวแบบ	Linear pred. coeff										
ถดถอย	$\hat{\beta}_0$	$\hat{\beta}_1$	$\hat{\beta}_2$	$\hat{\beta}_3$	$\hat{\beta}_4$	$\hat{\beta}_5$	$\hat{\beta}_6$	$\hat{\beta}_7$	$\hat{\beta}_8$	$\hat{\beta}_9$	
NB2	2.435 (0.343) (0.288) ⁺	0.603 (0.449) (0.448) ⁺	0.962 (0.470) (0.519) ⁺	-0.373 (0.403) (0.318) ⁺	1.126 (0.342) (0.292) ⁺	1.136 (0.380) (0.372) ⁺	-0.104 (0.335) (0.240) ⁺	-0.548 (0.432) (0.428) ⁺	-0.886 (0.426) (0.444) ⁺	0.119 (0.439) (0.452) ⁺	
ตัวแบบ	Linear pred. coeff							Overdispersion		AIC	df
ถดถอย	$\hat{\beta}_{10}$	$\hat{\beta}_{11}$	$\hat{\beta}_{12}$	$\hat{\beta}_{13}$	$\hat{\beta}_{14}$	$\hat{\beta}_{15}$	$\hat{\beta}_{16}$	$\hat{\alpha}$			
NB2	0.089 (0.454) (0.439) ⁺	-1.113 (0.450) (0.505) ⁺	-1.459 (0.442) (0.443) ⁺	-1.301 (0.445) (0.380) ⁺	-1.027 (0.455) (0.385) ⁺	-0.308 (0.464) (0.462) ⁺	1.666 (0.583) (0.524) ⁺	0.592 (0.080) (0.082) ⁺	1092.4	128	

หมายเหตุ + คือ Robust standard error

ในการทำนายค่าเฉลี่ยของจำนวนวันขาดเรียน เราได้ตัวแบบที่แปลงให้อยู่ในมาตราซึ่งกำลัง คือ

$$\begin{aligned}
 \mu_{\text{DAYS}} = & \exp[2.435 + 0.603C2 + 0.962S2 - 0.373A2 + 1.125A3 + 1.136A4 \\
 & - 0.104L2 - 0.548C2:A2 - 0.886C2:A3 + 0.119C2:A4 + 0.089S2:A2 \\
 & - 1.113S2:A3 - 1.459S2:A4 - 1.301C2:S2 - 1.027C2:L2 \\
 & - 0.308S2:L2 + 1.666C2:S2:L2] \quad (4.6)
 \end{aligned}$$