

تحليل الانحدار اللوجستي Logistic Regression Model

1- الغرض من التحليل:

يعتبر نموذج الانحدار اللوجستي أحد النماذج الإحصائية التي تستخدم في كثير من المجالات التطبيقية كالمجالات الاقتصادية، والزراعية، والطبية، وغيرها من المجالات العلمية الأخرى، إذ يهتم هذا النموذج بدراسة وتحليل أثر عدة متغيرات مستقلة (لا يشترط نوعها) على متغير تابع نوعي مكون من مجموعتين أو أكثر من المجموعات المتنافية. ففي مثل هذه الحالات يصعب استخدام نموذج الانحدار الخطي، لما تتطلبه من افتراضات غير محققة في الجانب التطبيقي.

2- الهيكل النظري لنموذج الانحدار اللوجستي الثاني

بفرض أن المتغير التابع متغير نوعي ذو مجموعتين متنافيتين (استحالة وقوعهما في آن واحد)، ويعبر عنه بمتغير ثنائي $y(0,1)$ ، فإن التوزيع الاحتمالي لهذا المتغير هو التوزيع البيرنولي Bernoulli بمعلمة $\pi = pr(y=1)$ والتي تمثل احتمال النجاح، وبفرض أن (y_1, y_2, \dots, y_n) تمثل المشاهدات لـ n من المتغيرات العشوائية المستقلة، فإن نموذج الانحدار اللوجستي يعرف كالتالي:

$$\begin{aligned} y_i &\sim \text{Bernoulli}(\pi_i), \\ \pi_i &= pr(y_i = 1 | X_{i1}, X_{i2}, \dots, X_{ip}) \\ &= \frac{e^{\lambda_i}}{1 + e^{\lambda_i}}, \quad i = 1, 2, \dots, n \end{aligned} \quad (1)$$

حيث تمثل $\beta = (\beta_0, \beta_1, \dots, \beta_p)'$ متجه معاملات الانحدار من الدرجة $((p+1) \times 1)$ ، وتمثل (X_{i1}, \dots, X_{ik}) المشاهدات المفسرة للمفردة رقم i ، ويفترض أنه محدد Fixed ومعطاه، وتعبّر λ_i عن دالة اللوجت وتأخذ الصورة التالية.

$$\lambda_i = Ln(\pi_i / 1 - \pi_i) = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_p X_{ip} \quad (1)$$

ويقوم برنامج (SPSS statistics, version 22,2013) بإيجاد تقدير لهذه المعاملات ويرمز لها

$$\hat{\beta}_{mle} = (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k)'$$
 بالرمز

3- قاعدة التصنيف

تعرف النسبة $\pi_i / (1 - \pi_i)$ بنسبة الرجحان $Odds Ratio (OR)$ لحالة النجاح $(y_i = 1)$ ، وتقدير هذه النسبة هو :

$$OR_i = \frac{\hat{\pi}_i}{1 - \hat{\pi}_i} = e^{\hat{\lambda}_i} \quad (8)$$

ويلاحظ أن $OR_i > 0$ ، وعندما يكون احتمال النجاح $\hat{\pi}_i$ أكبر من احتمال الفشل $(1 - \hat{\pi}_i)$ تكون نسبة الرجحان $OR_i > 1$ ، والعكس صحيح ، ومن ثم يمكن تصنيف المفردة لأحد مجموعتي المتغير التابع وفقا لقاعدة التصنيف التالية.

$$\begin{aligned} & \text{The individual number } i \text{ must be classified to success group } (y_i = 1) \text{ if } OR_i > 1 \\ & \text{The individual number } i \text{ must be classified to failure group } (y_i = 0) \text{ if } OR_i < 1 \end{aligned} \quad (9)$$

وتدل القاعدة أعلاه أن المفردة رقم i يرجح تصنيفها لمجموعة النجاح $(y_i = 1)$ إذا كانت نسبة الرجحان أكبر من الواحد $(OR_i > 1)$ ، بينما يرجح تصنيفها لمجموعة الفشل $(y_i = 0)$ إذا كانت نسبة الرجحان أقل من الواحد $(OR_i < 1)$.

وبأخذ لوغاريثم نسبة الرجحان $Log_e(OR_i)$ يمكن التوصل إلى نموذج يأخذ شكل معادلة خطية في معاملات الانحدار اللوجستية يسمى بنموذج اللوجت Logit، ويأخذ الصورة التالية:

$$\begin{aligned} \text{Logit}(\hat{\pi}_i) &= \text{Log}_e(OR_i) \\ \hat{\lambda}_i &= \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_k x_{ik} \end{aligned} \quad (10)$$

ويلاحظ من المعادلة أعلاه أن $\text{Logit}(\pi_i) > 0$ إذا كانت نسبة الرجحان $(OR_i > 1)$ ، وتكون $\text{Logit}(\pi_i) < 0$ إذا كانت نسبة الرجحان $(OR_i < 1)$ ومن ثم يمكن إعادة صياغة قاعدة التصنيف (9) من خلال نموذج اللوجت $\hat{\lambda}_i$ كما يلي:

$$\begin{aligned} & \text{The individual number } i \text{ must be classified to success group } (y_i = 1) \text{ if } \text{Logit}(\pi_i) > 0 \\ & \text{The individual number } i \text{ must be classified to failure group } (y_i = 0) \text{ if } \text{Logit}(\pi_i) < 0 \end{aligned} \quad (11)$$

ومن ثم إذا كان المعامل $\beta_j > 0$ يحدث زيادة نسبية في نسبة الرجحان للنجاح $(y_i = 1)$ ، وإذا كان المعامل

$\beta_j < 0$ يحدث انخفاض نسبي في نسبة الرجحان للنجاح ($y_i = 1$).

4- اختبارات الفروض.

بمجرد الحصول على تقديرات الإمكانية العظمى لمعاملات الانحدار اللوجستية يمكن إجراء العديد من الاختبارات الإحصائية التي تسعى لتحقيق الهدف من تطبيق نموذج الانحدار اللوجستي ومنها ما يلي:

اختبار مناسبة النموذج اللوجستي للتنبؤ.

يهدف هذا الاختبار إلى معرفة ما إذا كان نموذج الانحدار اللوجستي المقترح مناسباً للتنبؤ بمجموعة المتغير التابع التي تنتمي لها المفردة ذات المشاهدات المفسرة (x_1, x_2, \dots, x_k) . ومن ثم يأخذ فرض العدم H_0 والفرض البديل H_1 الصورة التالية:

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_k = 0 \quad \text{vs} \quad H_1 : \text{at least one of } \beta_j \neq 0 \quad (13)$$

ويستخدم اختبار نسبة الإمكان (LR) *Likelihood ratio* لإجراء هذا الاختبار، حيث يعبر عن إحصائية الاختبار بالمعادلة التالية:

$$\chi^2 = -2 \log_e(LR) = -2[\text{LL}(\hat{\beta}_{mle} | H_0) - \text{LL}(\hat{\beta}_{mle} | H_1)] \quad (14)$$

حيث أن $\text{LL}(\hat{\beta}_{mle} | H_1)$ ، $\text{LL}(\hat{\beta}_{mle} | H_0)$ تعبر عن القيمة المحسوبة لـ $\text{LL}(\beta)$ في حالة صحة الفرض العدم $H_0 : \beta_1 = \beta_2 = \dots = \beta_k = 0$ ، وفي حالة عدم صحته $H_1 : \text{at least one of } \beta_j \neq 0$. وتحت صحة الفرض العدم H_0 يقترب توزيع إحصائية الاختبار χ^2 من توزيع مربع كاي بدرجات حرية k أي أن $\chi^2 \approx \chi^2_{(k)}$.

اختبار معنوية تأثير المتغير المفسر.

يهدف هذا الاختبار إلى معرفة مدى الدلالة الإحصائية لتأثير كل متغير من المتغيرات المفسرة قيد الدراسة على المتغير التابع، ويعبر عن فرض العدم H_0 والفرض البديل H_1 كالتالي:

$$H_0 : \beta_j = 0 \quad \text{vs} \quad H_1 : \beta_j \neq 0 \quad , \quad j = 1, 2, \dots, k \quad (15)$$

ويستخدم اختبار (Wald) لإجراء هذا الاختبار، ويعبر عن إحصائية الاختبار بالمعادلة التالية:

$$wald = (\hat{\beta}_j / S.E_{\hat{\beta}_j})^2 \quad (16)$$

حيث أن $\hat{\beta}_j$ يعبر عن تقدير الإمكانية العظمى لمعامل الانحدار β_j والذي يعكس تأثير المتغير المفسر x_j على المتغير التابع، $S.E_{\hat{\beta}_j}$ يمثل الخطأ المعياري للتقدير $\hat{\beta}_j$. وتحت صحة الفرض العدم H_0 يقترب توزيع إحصائية الاختبار $wald$ من توزيع مربع كاي بدرجة حرية واحدة، أي أن $wald \approx \chi^2_{(1)}$.

اختبار جودة المطابقة.

ويقصد بهذا الاختبار مدى مطابقة القيم المتوقعة للمتغير التابع للقيم المشاهدة، وفي هذا الاختبار يعبر عن فرض العدم H_0 والفرض البديل H_1 كالتالي:

$$\begin{aligned} H_0 : & \text{The expected values from using prediction model match the observed values.} \\ \text{vs } H_1 : & \text{The expected values from using prediction model does not match the observed values.} \end{aligned} \quad (17)$$

ويدل الفرض العدم H_0 على أن القيم المتوقعة من استخدام نموذج التنبؤ تطابق القيم المشاهدة، وقبول هذا الفرض يستدل منه على أن النموذج تمثيل جيد للبيانات وأن ليس به قصور. وتستخدم طريقة Hosmer and Lemeshow (2000) لإجراء هذا الإختبار، حيث تستند هذه الطريقة على فكرة تقسيم المشاهدات الكلية n إلى عدد G من المجموعات المتساوية هي: n_1, n_2, \dots, n_G ، ثم ترتب احتمالات النجاح المتنبأ بها $\hat{\pi}_i$'s تصاعديا ويحسب متوسط احتمالات النجاح لكل مجموعة ففي المجموعة g يكون المتوسط هو $\bar{\pi}_g = \sum_{i=1}^{n_g} \hat{\pi}_i / n_g$ ، ومن ثم تحسب إحصائية الاختبار HL بالمعادلة التالية.

$$HL = \sum_{g=1}^G \left[\frac{(O_g - n_g \bar{\pi}_g)^2}{n_g \bar{\pi}_g (1 - \bar{\pi}_g)} \right] \quad (18)$$

حيث أن O_g هو التكرار المشاهد لعدد حالات النجاح داخل المجموعة g ، $\sum_{g=1}^G n_g = n$ ، وتحت صحة الفرض العدم H_0 يقترب توزيع الإحصائية HL من توزيع مربع كاي بدرجات حرية $(G - 2)$ أي أن $HL \approx \chi^2_{(G-2)}$.

نسبة التصنيف الصحيح

يقصد بالتصنيف الصحيح للمفردة، أنها إذا كانت فعلا تنتمي للمجموعة k ($success$ ، $failure$) وتم تصنيفها وفقا للمعيار (9) أو (11) لنفس المجموعة k ($success$ ، $failure$)، وعلى هذا الأساس يمكن

حساب نسبة التصنيف الصحيح لكل مجموعة وكذلك نسبة التصنيف الصحيح العام. وكلما كانت نسبة التصنيف الصحيح عالية كان ذلك مؤشرا لجودة النموذج.

تطبيق

المتغير التابع: موقف العائلة من تفضيلها أو عدم تفضيلها شراء سلعة معينة، حيث أن $y = 1$: إذا كانت العائلة تفضل شراء السلعة. $y = 0$: إذا كانت العائلة لا تفضل شراء السلعة.

المتغيرات المفسرة:

incom: الدخل بالألف ريال *Family*: حجم الأسرة

البيانات

No	التفضيل	الدخل	حجم الأسرة	No	التفضيل	الدخل	حجم الأسرة
1	0	2.1	1	16	1	4.9	3
2	0	2.6	1	17	1	5	1
3	0	3.1	3	18	1	5	3
4	0	4.6	2	19	1	5	3
5	0	3.3	2	20	1	5.2	4
6	0	3.3	2	21	1	5.4	4
7	0	3.3	2	22	1	5.5	4
8	0	5.6	2	23	1	3.8	5
9	0	3.8	4	24	1	5.8	5
10	0	4.2	4	25	1	6.7	5
11	0	4.2	4	26	1	6.7	6
12	0	4.4	3	27	1	7.6	6
13	1	4.2	4	28	1	7.8	6
14	1	4.6	3	29	1	8.6	6
15	1	3.2	4	30	1	9.7	8

النتائج وتفسيرها

Logistic Regression

Case Processing Summary

Unweighted Cases ^a		N	Percent
Selected Cases	Included in Analysis	30	100.0
	Missing Cases	0	.0
	Total	30	100.0
Unselected Cases		0	.0
Total		30	100.0

a. If weight is in effect, see classification table for the total number of cases.

Dependent Variable Encoding

Original Value	Internal Value
No	0
Yes	1

Omnibus Tests of Model Coefficients

		Chi-square	df	Sig.
Step 1	Step	19.028	2	.000
	Block	19.028	2	.000
	Model	19.028	2	.000

Model Summary

Step	-2 Log likelihood	Cox & Snell R Square	Nagelkerke R Square
1	21.352 ^a	.470	.635

a. Estimation terminated at iteration number 7 because parameter estimates changed by less than .001.

Hosmer and Lemeshow Test

Step	Chi-square	df	Sig.
1	9.186	8	.327

Contingency Table for Hosmer and Lemeshow Test

		التفضيل = No		التفضيل = Yes		Total
		Observed	Expected	Observed	Expected	
Step 1	1	2	1.990	0	.010	2
	2	3	2.836	0	.164	3
	3	1	2.370	2	.630	3

4	3	1.683	0	1.317	3
5	2	1.486	2	2.514	4
6	0	1.140	4	2.860	4
7	1	.420	2	2.580	3
8	0	.074	3	2.926	3
9	0	.002	3	2.998	3
10	0	.000	2	2.000	2

Classification Table^a

Observed		Predicted		
		التفضيل		Percentage Correct
		No	Yes	
Step 1	التفضيل No	9	3	75.0
	Yes	2	16	88.9
Overall Percentage				83.3

a. The cut value is .500

Variables in the Equation

		B	S.E.	Wald	df	Sig.	Exp(B)
Step 1 ^a	Incom	1.692	.835	4.110	1	.043	5.430
	Family	.963	.572	2.835	1	.092	2.620
	Constant	-10.359	4.694	4.871	1	.027	.000

Variables in the Equation

		95% C.I. for EXP(B)	
		Lower	Upper
Step 1 ^a	Incom	1.058	27.878
	Family	.854	8.037
	Constant		

a. Variable(s) entered on step 1: Incom, Family.

Step number: 1

Observed Groups and Predicted Probabilities

التفضيل	الدخل	حجم الأسرة	Predicted probability	Predicted group
0	2.1	1	0.00289	0
0	2.6	1	0.00671	0
0	3.1	3	0.09753	0
0	4.6	2	0.343	0
0	3.3	2	0.0547	0
0	3.3	2	0.0547	0
0	3.3	2	0.0547	0
0	5.6	2	0.73924	1
0	3.8	4	0.48061	0
0	4.2	4	0.64548	1
0	4.2	4	0.64548	1
0	4.4	3	0.49366	0
1	4.2	4	0.64548	1
1	4.6	3	0.57762	1
1	3.2	4	0.25109	0
1	4.9	3	0.69437	1
1	5	1	0.28167	0
1	5	3	0.72905	1
1	5	3	0.72905	1
1	5.2	4	0.90815	1
1	5.4	4	0.93274	1
1	5.5	4	0.94261	1
1	3.8	5	0.70794	1
1	5.8	5	0.9862	1
1	6.7	5	0.99696	1
1	6.7	6	0.99884	1
1	7.6	6	0.99975	1
1	7.8	6	0.99982	1
1	8.6	6	0.99995	1
1	9.7	8	1.0	1

التعليق على النتائج

- يقوم الطالب بالتعليق من خلال ما تم فهمه في المحاضرة