

تحليل الانحدار الخطي المتعدد Multiple Linear Regression Analysis

المحاضرة الرابعة

مثال (02): اعتمادا على المثال السابق (01) أحسب مصفوفة التباين / التغاير لمعاملات نموذج وزن الطفل، ومن ثم حدد تقديرات الأخطاء (الانحرافات) المعيارية لهذه المعاملات.

الحل:

أولا نقوم بحساب التباين (S_e^2):

$$S_e^2 = \frac{Y^T \cdot Y - b^T \cdot X^T \cdot Y}{n - (k+1)}$$

$$(Y^T Y) = (11.5 \quad 16 \quad . \quad . \quad . \quad 1.4) \cdot \begin{pmatrix} 11.5 \\ 16 \\ . \\ . \\ . \\ 1.4 \end{pmatrix} = 6653.3853$$

و عليه نجد:

$$[Y^T \cdot Y - b^T \cdot X^T \cdot Y] = \left[(6653.3853) - (-2.1819 \quad 1.2008 \quad 0.12457) \cdot \begin{pmatrix} 507.7 \\ 1739.2 \\ 45081.8 \end{pmatrix} \right]$$

$$[Y^T \cdot Y - b^T \cdot X^T \cdot Y] = [(6653.3853) - (6596.6948)] = 56.6904$$

وبالتالي فإن:

$$S_e^2 = \frac{56.6904}{50-(2+1)} = 1.20618$$

بضرب التباين (S_e^2) في المصفوفة $(X^T X)^{-1}$ التي حسبنا في المثال السابق نتحصل على مصفوفة التباين / التباين الخاصة بمقدرات معالم نموذج الانحدار (b) كما يلي:

$$\begin{aligned} S_b^2 &= \widehat{\text{var}}(b) = S_e^2 \cdot (X^T X)^{-1} \\ &= 1.20618 \cdot \begin{pmatrix} 0.789948 & 0.144239 & -0.014509 \\ 0.144239 & 0.036936 & -0.003023 \\ -0.014509 & -0.003023 & 0.000283 \end{pmatrix} \\ &= \begin{pmatrix} \mathbf{0.95282} & 0.173978 & -0.017501 \\ 0.173978 & \mathbf{0.044551} & -0.003646 \\ -0.017501 & -0.003646 & \mathbf{0.000341} \end{pmatrix} \end{aligned}$$

وبالتالي من خلال هذه المصفوفة نجد:

$$\widehat{\text{var}}(b_0) = 0.95282$$

$$\widehat{\text{var}}(b_1) = 0.044551$$

$$\widehat{\text{var}}(b_2) = 0.000341$$

ومنه فإن الانحرافات المعيارية لمعاملات نموذج الانحدار هي:

$$\hat{\sigma}(b_0) = \sqrt{0.95282} = \mathbf{0.976125}$$

$$\hat{\sigma}(b_1) = \sqrt{0.044551} = \mathbf{0.21107}$$

$$\hat{\sigma}(b_2) = \sqrt{0.000341} = 0.018466$$

6-II. خصائص البواقي أو الأخطاء

البواقي هي القيم المقدرة لحد الخطأ العشوائي (ε_i) وهي عبارة عن الفروق بين القيم الفعلية والقيم المقدرة للمتغير التابع، وتتميز البواقي بعدة خصائص أهمها:

1- القيمة المتوقعة لأي عنصر من عناصر متجه البواقي تساوي الصفر، أي:

$$\mathbf{E}(\mathbf{e}) = \mathbf{0} \quad (14)$$

2- استقلال البواقي عن المتغيرات المستقلة أي:

$$\mathbf{X}^T \cdot \mathbf{e} = \mathbf{0} \quad (15)$$

3- استقلال البواقي عن القيم المقدرة للمتغير التابع أي:

$$\hat{\mathbf{y}}^T \cdot \mathbf{e} = \mathbf{0} \quad (16)$$

7-II. معامل التحديد المتعدد

معامل التحديد في الانحدار الخطي المتعدد له نفس التفسير كما رأينا في الانحدار الخطي البسيط. فهو يقيس نسبة التغيرات أو الاختلافات في قيم المتغير التابع التي تفسر بواسطة المتغيرات المستقلة في معادلة انحدار المربعات الصغرى.

وكما رأينا في الانحدار الخطي البسيط فإن معامل التحديد في هذه الحالة يحسب بنفس الطريقة أي:

$$R^2 = \frac{SSR}{SST} = \frac{SST - SSE}{SST} \quad (17)$$

وعند استخدام رموز المصفوفات يمكن كتابة: SST ، SSR ، SSE على النحو التالي:

$$SST = \sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum Y_i^2 - n\bar{Y}^2 = \mathbf{Y}^T \cdot \mathbf{Y} - n\bar{Y}^2 \quad (18)$$

$$SSR = \sum_{i=1}^n (\hat{y}_i - \bar{Y})^2 = SST - SSE = \mathbf{b}^T \cdot \mathbf{X}^T \cdot \mathbf{Y} - n\bar{Y}^2 \quad (19)$$

$$SSE = \sum_{i=1}^n (Y_i - \hat{y}_i)^2 = \mathbf{e}^T \cdot \mathbf{e} = \mathbf{Y}^T \cdot \mathbf{Y} - \mathbf{b}^T \cdot \mathbf{X}^T \cdot \mathbf{Y} \quad (20)$$

ومنه يمكن كتابة معامل التحديد كما يلي:

$$R^2 = \frac{\mathbf{b}^T \cdot \mathbf{X}^T \cdot \mathbf{Y} - n\bar{Y}^2}{\mathbf{Y}^T \cdot \mathbf{Y} - n\bar{Y}^2} \quad (21)$$

ويتصف معامل التحديد المتعدد كما هو الحال في نموذج الانحدار الخطي البسيط بالخاصية التالية:

$$0 \leq R^2 \leq 1$$

ملاحظة:

كلما كانت قيمة معامل التحديد صغيرة، كلما كان الجزء الأكبر من التباين في المتغير التابع راجع الى متغيرات لا توجد في نموذج الانحدار، أما إذا كانت قيمة معامل التحديد قريبة من الواحد الصحيح فان هذا يدل على أن الجزء الأكبر من التباين في المتغير التابع قد تم تفسيره بواسطة المتغيرات المستقلة الموجودة في نموذج الانحدار.

- الاستخدام المُعيب لـ R^2 :

في سياق الحديث عن تحليل الانحدار الخطي المتعدد عادة ما يُساء فهم (R^2) أو يُساء استخدامه. ومن الأهمية أن نعلم أن قيمة (R^2) لا يمكن أن تقل عندما تُضاف متغيرات مفسرة الى نموذج الانحدار، حتى ولو كانت هذه المتغيرات لا تساهم بمعلومات إضافية للتنبؤ بقيمة (Y) . وهذا صحيح لأن الاختلاف غير المفسر في العينة لقيم (Y) ، كما قيست بواسطة SSE تتناقص بوضوح عندما يكون هناك حد إضافي، قد أضيف لنموذج الانحدار بينما يظل مجموع المربعات الكلي SST ثابتا بصرف النظر عن عدد المكونات في النموذج (لأن SST تكون محددة كليا بواسطة قيم X). لهذا فان مجموع مربعات الانحدار SSR (الاختلاف المفسر) يجب أن يزيد على الأقل عندما تُضاف حدود جديدة (متغيرات مستقلة أو مفسرة) الى النموذج.

فإذا استخدمنا (R^2) لتحديد ما إذا كان يجب إضافة عناصر جديدة للنموذج أم لا. إذن السؤال لا يكون ما إذا كانت هناك زيادة في قيمة (R^2) عند إضافة متغيرات جديدة ولكن بكم تزيد (R^2) ؟

(R^2) الكبيرة لا تعني بالضرورة نموذج أفضل، في الحقيقة (R^2) الكبيرة بدرجة كافية يمكن تحقيقها ببساطة بإضافة متغيرات تفسيرية، البعض منها ربما يُساهم في تفسير القليل من التغيرات في قيم (Y) بالعينة. [بعض المحللين يُخطئون بضم عدد كبير من المتغيرات المفسرة في النموذج كأساس للحصول على قيمة عالية لقيمة (R^2)].

- معامل التحديد المعدل (The Adjusted Coefficient of Determination)

كما سبق وأن ذكرنا فان قيمة (R^2) تزيد بزيادة العناصر المضافة لنموذج الانحدار؛ فإضافة عناصر كافية (متغيرات مفسرة) للنموذج يمكن أن تقترب (R^2) من الواحد الصحيح. لهذا السبب فان هناك علاقة بديلة لقياس جودة التوفيق قد أُقترحت لتأخذ عناصر النموذج في الحسبان، وهذا المقياس الوصفي لجودة توفيق معادلة المربعات الصغرى يسمى "معامل التحديد المعدل أو المصحح" ويحسب وفق العلاقة التالية:

$$R_a^2 = 1 - \left(\frac{SSE/(n-k-1)}{SST/(n-1)} \right) = 1 - \left[\frac{n-1}{n-k-1} \right] \cdot \frac{SSE}{SST} \quad (22)$$

ونعلم أن:

$$R^2 = \frac{SSR}{SST} = \frac{SST-SSE}{SST} = 1 - \frac{SSE}{SST} \Rightarrow \frac{SSE}{SST} = 1 - R^2$$

وعليه فإن العلاقة رقم (22) تصبح كما يلي:

$$R_a^2 = 1 - \left[\frac{n-1}{n-k-1} \right] \cdot (1 - R^2) \quad (23)$$

والشيء الملاحظ على معامل التحديد المعدل هو:

- أنه يأخذ قيما أقل من قيم معامل التحديد (قيمة R_a^2 دائما أقل من قيمة R^2).
 - أنه يمكن أن يأخذ قيما سالبة، في حين نجد أن قيم R^2 تكون دائما موجبة.
 - قيمة (R_a^2) يمكن أن تتناقص عندما تُضاف متغيرات مفسرة غير مناسبة لنموذج الانحدار.
- لهذا فإنه في تحليل الانحدار الخطي المتعدد يُفضل استخدام معامل التحديد المعدل على معامل التحديد للمقارنة بين نماذج الانحدار المتنافسة.