

Comparison of some robust methods in the presence of problems of multicollinearity and high leverage points

المقارنة بين بعض الطرائق الحصينة في ظل وجود مشكلتي تعدد العلاقة الخطية ونقاط الانعطاف العالية

ا.م. غفران اسماعيل كمال / كلية الادارة والاقتصاد / جامعة بغداد

ghufuran62@gmail.com

الباحث / سيف الامام سعدي خزعل / كلية الادارة والاقتصاد / جامعة بغداد saif.alemams@gmail.com

25
19

OPEN ACCESS



P - ISSN 2518 - 5764
E - ISSN 2227 - 703X

Received: 27/11/2018

Accepted: 26/12/2018

مستخلص البحث

يعد نموذج الانحدار الخطي المتعدد من نماذج الانحدار المهمة والمستعملة في تحليل البيانات لمختلف مجالات العلم وعلى نطاق واسع مثل الاعمال والاقتصاد والطب والعلوم الاجتماعية، ان تعدد العلاقة الخطية مشكلة كبيرة في الانحدار الخطي المتعدد اذ تؤدي في ابسط حالتها الى ابتعاد معاملات النموذج المقدر على خصائصها العلمية وغالباً ما تعطي استنتاجات مظللة، ايضاً هناك مشكلة هامة في تحليل الانحدار هو وجود نقاط الانعطاف العالية في البيانات مما تؤدي الى تأثيرات غير مرغوب بها على نتائج التحليل . نستعرض في هذا البحث بعض الطرائق الحصينة في نموذج الانحدار الخطي المتعدد ومن هذه الطرائق طريقتي انحدار الحرف لمقدر ال جاكنايف (Jackknife Ridge Regression) بالاعتماد على مقدر (MM-estimator) (MM) ومقدر (GM2) (Modified Generalized M-estimator)، ومن خلال استعمال المحاكاة بأسلوب مونت كارلو تمت اجراء المقارنة بين هاتين الطريقتين وفق معيار المقارنة متوسط مربعات الخطأ (MSE) ولحجوم عينات (n=20, n=50, n=100) ونسب تلوث مختلفة (5%, 15%, 50%)، واتضح من خلال المقارنة ان طريقة (RJGM2) هي الافضل في تقدير معاملات نموذج الانحدار الخطي المتعدد و يمتلك اقل قيمة لمتوسط مربعات خطأ (MSE) مقارنة مع بقية المقدرات الأخرى.

المصطلحات الرئيسية للبحث/ الانحدار الخطي المتعدد، تعدد العلاقة الخطية، نقاط الانعطاف العالية، مقدر MM ، مقدر GM2، انحدار الحرف لل جاكنايف .





المقارنة بين بعض الطرائق الحصينة في ظل وجود مشكلتي تعدد العلاقة الخطية ونقاط الانعطاف العالية

1- المقدمة Introduction

يستعمل نموذج الانحدار الخطي المتعدد (MLR) (Multiple Linear Regression) في العديد من مجالات الدراسة مثل الاعمال والاقتصاد والطب والعلوم الاجتماعية، فعند دراسة أي ظاهرة يجب تحديد المتغيرات المؤثرة في تلك الظاهرة وصياغة العلاقة بين تلك المتغيرات على هيئة أنموذج، ان العلاقة الخطية بين عدة متغيرات احدهما متغير معتمد (Dependent Variable) والباقي متغيرات توضيحية (Explanatory variables) يطلق عليها أنموذج الانحدار الخطي المتعدد (MLR)، أن مشكلة تعدد العلاقة الخطية أصبحت معروفة لدى العديد من الباحثين الإحصائيين و تمثل حالة انعدام الاستقلالية بين المتغيرات التوضيحية مما يسبب وجودها تأثيرات على تقديرات وتباينات المعاملات عند تطبيق طريقة المربعات الصغرى الاعتيادية (OLS) وبالتالي تظهر النتائج غير دقيقة، مشكلة مهمة أخرى في تحليل الانحدار هي وجود نقاط الانعطاف العالية (HLPs) (High Leverage Points) والتي تمثل (القياسات غير طبيعية في قيم المتغير x) والتي تسبب في ميل خط الانحدار نحوها بعيداً عن موقعة الحقيقي، وهي على نوعين الجيدة التي تتفق مع خط الانحدار وتساهم في كفاءة التقدير ومنها السيئة التي تؤثر على القيم المحسوبة لمختلف التقديرات . [1] (Midi & Mohammed, 2015:P147)

ان نقاط الانعطاف العالية (HLPs) لها تأثيرات كبيرة على الأنموذج الخطي منها التسبب في تضليل البيانات وفشل التحليل كونها تعتبر مصدر اخر لمصادر تعدد العلاقة الخطية وظهرت نتائج دراسة المحاكاة لكل من الباحثان (Imon & Kamruzzaman) عام (2002 م) ان نقاط الانعطاف العالية يمكن ان تولد تعدداً خطياً كبيراً وأن درجته قد تزداد مع زيادة نسبتها في البيانات . [3] (Kamruzzaman & Imon, 2002:P445)

مما يستوجب استعمال طرائق بديلة عن الطرائق التقليدية لمعالجة تلك المشاكل تكون اكثر كفاءة في التقدير والتي تدعى (بالطرائق التقدير الحصينة) (Robust estimation methods) حيث تتصف انها قليلة الحساسية تجاه نقاط الانعطاف العالية (HLPs) اذ يتم الحصول عليها من خلال مقدرات حصينة تمتلك كفاءة عالية .

2- مشكلة البحث The problem of the research

أن أنموذج الانحدار الخطي المتعدد (MLR) يتضمن عدد من المتغيرات التوضيحية التي تكون ذات أداء ضعيف للمقدرات لوجود علاقة خطية بين المتغيرات التوضيحية مما يؤدي لعواقب سيئة مثل زيادة في تباينات المقدرات وانخفاض معامل الدقة وغالباً ما تكون النتائج مربكة وتعطي استنتاجات مظلمة، بالإضافة الى ذلك وجود نقاط الانعطاف العالية (HLPs) في البيانات التي ستضاعف من المشكلة .

3- هدف البحث The Aim of the Research

الهدف من البحث هو الحصول على افضل تقدير لمعاملات أنموذج الانحدار الخطي المتعدد (MLR) عند وجود مشكلة تعدد العلاقة الخطية ونقاط الانعطاف العالية (HLPs) من خلال استعمال معيار المقارنة متوسط مربعات الخطأ (MSE) للحصول على أفضل مقدر بين المقدرات الاخرى .

4- الجانب النظري

يعطي تحليل الانحدار الاجابات على الاسئلة المتعلقة بالعلاقة الوظيفية للمتغير المعتمد (Y) مع واحد او اكثر من المتغيرات التوضيحية (X's)، يضاف مصطلح الخطأ العشوائي الى الأنموذج الاحصائي لحساب الفروق الفردية، أن الدالة الخطية بين Y و X تسمى أنموذج الانحدار الخطي المتعدد (MLR) ويمكن تعريفها على أنها :

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + \varepsilon_i \quad \dots (1)$$

$i = 1, 2, \dots, n$ ، $j = 0, 1, \dots, K$

الأنموذج اعلاه يمكن كتابته بشكل مختصر وكالاتي :-



المقارنة بين بعض الطرائق الحصينة في ظل وجود مشكلتي تعدد العلاقة الخطية ونقاط الانعطاف العالية

$$Y_i = \sum_{j=0}^k B_j X_{ij} + \varepsilon_i \quad \dots (2)$$

وبدلالة المصفوفات يمكن صياغة نموذج الخطي العام كما في الشكل الاتي :

$$\underline{y} = X\underline{B} + \underline{\varepsilon} \quad \dots (3)$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_i \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & X_{11} & X_{12} & \dots & X_{1j} & \dots & X_{1k} \\ 1 & X_{21} & X_{22} & \dots & X_{2j} & \dots & X_{2k} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & X_{i1} & X_{i2} & \dots & X_{ij} & \dots & X_{ik} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & X_{n1} & X_{n2} & \dots & X_{nj} & \dots & X_{nk} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_j \\ \vdots \\ \beta_k \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_i \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

Y : موجه مشاهدات المتغير المعتمد من الدرجة $(n * 1)$.
 X : مصفوفة من الدرجة $(n * (k + 1))$ وتمثل مشاهدات المتغيرات التوضيحية علماً بأن العمود الاول من هذه المصفوفة يمثل الحد الثابت .
 β : موجه المعالم المطلوب تقديرها من الدرجة $1 * (k + 1)$.
 k : عدد المتغيرات التوضيحية .
 n : عدد المشاهدات .

ε : موجه الاخطاء العشوائية من الدرجة $(n * 1)$.
لتقدير معلمات نموذج الانحدار الخطي المتعدد (MLR) يتم الاعتماد على بعض المقدرات الحصينة وهي :
[3] (كاظم ومسلم، 2002 : 50)

5- بعض مقدرات الانحدار الحصينة Regression Estimates Some of Robust

5-1- MM- estimator MM مقدر

أن مقدر (MM) في التسمية يشير الى استعمال لأكثر من عملية تقدير لمقدر M، اقترح من الباحث (Yohai) عام (1987 م) أن هذا المقدر له العديد من الخصائص الجيدة منها أن لها كفاءة عالية في حالة التوزيع الطبيعي للأخطاء مع نقطة الانهيار العالية، تعد طريقة (MM) واحدة من أكثر المقدرات شيوعاً التي تستعمل في مجال الانحدار الحصين ويكون وفقاً للخطوات التالية :-

1- يتم تحديد مقدر أولي ذو نقطة انهيار عالية، ولكن ليس بالضرورة أن يكون كفوء، نرسم له بالرمز $(\hat{\beta}_s)$ وباستخدامه يتم حساب البواقي الأولية وفقاً للصيغة التالية :

$$r_i(\hat{\beta}_s) = y_i - x_i' \hat{\beta}_s, \quad 1 < i < n \quad \dots (4)$$

2- يتم حساب مقدر M القياس (σ_n) للبواقي الأولية $r_i(\hat{\beta}_s)$ وفق معادلة M التقديرية لمعلمة القياس وبالشكل الاتي :

$$\frac{1}{n} \sum_{i=1}^n p \left(\frac{r_i(\hat{\beta}_s)}{\sigma} \right) = 0,5 \quad \dots (5)$$



المقارنة بين بعض الطرائق الحصينة في ظل وجود مشكلتي تعدد العلاقة الخطية ونقاط الانعطاف العالية

3- مقدر MM يعرف كمقدر M لـ β باستعمال دالة (re-descending score) :

$$\psi_1(u) = \frac{\partial p_1(u)}{\partial u} \quad \dots (6)$$

عليه فإن مقدر MM والذي نرسم له بالرمز $\hat{\beta}_{MM}$ يكون الحل للمعادلة التالية :

$$\sum_{i=1}^n X_{ij} \psi_1 \left(\frac{y_i - x_i' \beta}{\sigma_n} \right) = 0, \quad j = 1, 2, \dots, k \quad \dots (7)$$

أن تقدير القياس σ_n يتم أيجاده في الخطوة (2) .
اذ أن :

$$r_i(\beta) = y_i - x_i' \beta \quad [5] \text{ (Lukman, et.al, 2015:P58)}$$

6- مقدر (GM2) GM2-estimator

اقترح من قبل الباحثة (Bagheri) عام (2011 م)، أن فكرة هذا المقدر تقوم على اساس التعديل لمقدر (GM) (Modified generalized M-estimator) للحصول على مقدر حصين ذو خصائص عالية الكفاءة وأكثر مقاومة لنقاط الانعطاف العالية (HLPs)، يتم تلخيص الطريقة من خلال عدة خطوات على النحو الاتي :

1- يتم حساب البواقي الاولى (r_i) من مقدر S، ومن ثم أيجاد مقياس للبواقي ($\hat{\tau}$) وكالاتي :

$$r_i = y_i - \hat{y}_i, \quad i=1, 2, \dots, n \quad \dots (8)$$

$$\hat{\tau} = 1.4826(1 + 5 / (n - p)) \text{ Median } |r_i| \quad \dots (9)$$

2- يتم أيجاد مصفوفة الاوزان القطرية (W)، وعناصر القطر هي الاوزان (w_i) من خلال المعادلة الاتية :

$$w_i = \min \left[1, \left\{ \frac{x_{(0.95, p+1)}^2}{RMD^2} \right\} \right], \quad i=1, 2, \dots, n \quad \dots (10)$$

RMD : تمثل مسافة (مهن لوبس) (Mahalanobis Distance) الحصينة بالاعتماد على المقدار الأدنى للمقياس الإهليجي (MVE) .

3- حساب دالة التأثير Ψ^* للبواقي المعيارية، وذلك بحل المعادلة الاتية :

$$A = \text{diag } \Psi^* \left(\frac{r_i}{\hat{\tau} * w_i} \right) \quad \dots (11)$$

Ψ^* : مشتقة من دالة التأثير لـ (Huber's) .

4- عليه فإن مقدر (GM2) والذي نرسم له بالرمز ($\hat{\beta}_{GM2}$) يمكننا الحصول عليه عن طريق اشتقاق خطوة واحدة لطريقة (Newton Raphson) ويكون بالصيغة الاتية :

$$\tilde{\alpha}_{GM2} = \hat{\beta}_0 + (X' A X)^{-1} X' W \Psi \left(\frac{r_i}{w_i \hat{\tau}} \right) \hat{\tau} \quad \dots (12)$$

[1] (Alguraibawi, et.al, 2015:P311)



المقارنة بين بعض الطرائق الحصينة في ظل وجود مشكلتي تعدد العلاقة الخطية وتقاط الانعطاف العالية

7- طرائق التقدير:

لتقدير معلمات أنموذج الانحدار الخطي المتعدد سيتم الاعتماد في هذا البحث على بعض الطرائق الحصينة:

7-1- طريقة (RJRR):

A-Robust Jackknife Ridge Regression based on MM- estimator(RJRMM)

B -Robust Jackknife Ridge Regression based on GM2-estimator(RJRGM2)

نظرياً اقترحت طريقة الـ جاكنايف (Jackknife method) من الباحث (Quenouille) عام (1949 م) لتقليل التحيز والتباين في التقدير، وان المبدأ الاساسي لهذه الطريقة هو استبعاد جزء من البيانات عند إجراء كل عملية تقدير بحيث أن الجزء المستبعد اما يهمل او يؤخذ بنظر العناية عند التقدير من خلال احتسابه، في عام (1986 م) اقترح (Singh) وآخرون أسلوباً للتحايل على التحيز في انحدار الحرف بالاعتماد على تقنية الـ جاكنايف من خلال الصيغة الاتية:

$$y_{-i} = X_{-i}\beta + \varepsilon^* \quad \dots (13)$$

من خلال أنموذج الانحدار الخطي أعلاه بصيغة الـ جاكنايف يتم استبعاد قيمة من المتجه y_{-i} يقابلها استبعاد صف بالكامل من المصفوفة X_{-i} حيث ليس من الضروري ان تكون المصفوفة X_{-i} كاملة الرتبة للأعمدة، أما بالنسبة لـ β تمثل متجه المعلمات لانموذج الانحدار. ε^* يمثل حد الخطأ العشوائي حيث ان:

$$E(\varepsilon^*) = 0$$

$$\text{Cov}(\varepsilon^*) = \sigma^2 I_{n-1}$$

هنا يكون الأنموذج المختزل لانحدار الحرف الاعتيادي (RR) لـ (β) بالصيغة الاتية:

$$\hat{\beta}_{RR(-i)} = (X'_{-i}X_{-i} + kI_p)^{-1}X'_{-i}y_{-i} \quad \dots (14)$$

X_{-i} : الجزء او الصف i المستبعد من المصفوفة X .

y_{-i} : الجزء او القيمة i المستبعدة من المتجه y .

بإمكاننا إعادة كتابة المعادلة (14) للأنموذج المختزل لانحدار الحرف الاعتيادي (RR) بالشكل القانوني (Canonical form) لـ (α) بالشكل الاتي:

$$\hat{\alpha}_{RR(-i)} = (Z'_{-i}Z_{-i} + kI_p)^{-1}Z'_{-i}Y_{-i} \quad \dots (15)$$

نأخذ كل من z_i ، y_i لمتجه الاعمدة Z و Y بالتنسيق، تكون الصيغة بالشكل الاتي:

$$\hat{\alpha}_{RR(-i)} = (Z'Z - z'_i z_i + kI)^{-1}(Z'y - z_i y_i) \quad \dots (16)$$

[1] (Alguraibawi, et.al, 2015:PP308-309)



المقارنة بين بعض الطرائق الحصينة في ظل وجود مشكلتي تعدد العلاقة الخطية ونقاط الانعطاف العالية

ومن خلال تطبيق نظرية (Binomial Inverse Theorem) نحصل على :

$$\begin{aligned} &= \left[B^{-1} + \frac{B^{-1}z_i'z_iB^{-1}}{1 - z_i'B^{-1}z_i} \right] \times (Z'y - z_iy_i) \\ &= (B^{-1}Z'y - B^{-1}z_iy_i + \frac{B^{-1}z_i'z_iB^{-1}}{1 - z_i'B^{-1}z_i}Z'y - \frac{B^{-1}z_i'z_iB^{-1}}{1 - z_i'B^{-1}z_i}z_iy_i) \\ &= \hat{\alpha}_{RR} - B^{-1}z_iy_i \left(1 + \frac{z_i'B^{-1}z_i}{1 - z_i'B^{-1}z_i} \right) + \frac{B^{-1}z_i'z_i\hat{\alpha}_{RR}}{1 - z_i'B^{-1}z_i} \\ &= \hat{\alpha}_{RR} - \frac{B^{-1}z_i(Y_i - z_i'\hat{\alpha})}{1 - z_i'B^{-1}z_i} \\ \hat{\alpha}_{RR(-i)} &= \hat{\alpha}_{RR} - \frac{B^{-1}z_i e_i}{1 - h_i} \quad \dots(17) \end{aligned}$$

[2] (Duran & Akdeniz, 2010:P268)

أذ أن

$$h_i = z_i'(Z'Z)^{-1}z_i \quad , \quad e_i = (Y_i - z_i'\hat{\alpha})$$

e_i : موجة البواقي .

عندها مقدر ال جاكنايف يكون وفق الشكل الآتي :

$$p_i = n\hat{\alpha}_{RR} - (n - 1)\hat{\alpha}_{RR(-i)} \quad \dots (18)$$

المعادلة أعلاه تسمى بالقيم غير الأكيدة (الوهمية) (Pseudo - values) لتقدير ال جاكنايف ، من المعادلة (17) و (18) يكون الوسط الحسابي لهذه التقديرات بالصيغة الآتية :

$$\bar{P} = \frac{1}{n} \sum_{i=1}^n P_i = \hat{\alpha}_{RR} + \frac{(n-1)}{n} B^{-1} \sum_{i=1}^n \frac{z_i e_i}{(1 - h_i)} \quad \dots (19)$$

ان القيم غير الأكيدة في المعادلة (18) تعرف انها متماثلة فيما يتعلق بالملاحظات ، في حين ان الأنموذج غير متوازن بشكل عام وينعكس نقص التوازن في المسافات (h_i) ، بالإضافة الى ان التباين ($\hat{\alpha}_{RR(-i)} - \hat{\alpha}_{RR}$) هو دالة متزايدة لـ h_i ، يمكننا تعريف القيم غير الأكيدة الموزونة كما في الصيغة الآتية :



المقارنة بين بعض الطرائق الحصينة في ظل وجود مشكلتي تعدد العلاقة الخطية ونقاط الانعطاف العالية

$$Q_i = \hat{\alpha}_{RR} + n(1 - h_i)(\hat{\alpha}_{RR} - \hat{\alpha}_{RR(-i)}) \quad \dots (20)$$

هنا تكون تقديرات (JRR) (Jackknife Ridge Regression) متماثلة، أي مقدر $\hat{\alpha}_{JRR}(k)$ يعطى وفق بالصيغة الآتية :

$$\hat{\alpha}_{JRR}(k) = \bar{Q} = \frac{1}{n} \sum Q_i = \hat{\alpha}_{RR} + B^{-1} \sum z_i e_i \quad \dots (21)$$

[7] (Singh, et.al, 1986:344)

يمكن تبسيط النموذج (JRR) بالصيغة الآتية :

$$\begin{aligned} \hat{\alpha}_{JRR}(k) &= (I + kB^{-1})\hat{\alpha}_{JRR}(k) \quad \dots (22) \\ &= (I - k^2B^{-2})\hat{\alpha} \end{aligned}$$

علما ان معاملات الانحدار الاصلية لـ (JRR) هي :

$$\hat{\beta}_{JRR} = \gamma \hat{\alpha}_{JRR}(k)$$

وان مقدر (JRR) متحيز بالنسبة للمعلمة α ومقدار التحيز هو :

$$Bias(\hat{\alpha}_{JRR}(k)) = -k^2B^{-2}\hat{\alpha} \quad \dots (23)$$

ومصفوفة التباين لمقدر (JRR) كالآتي :

$$Var(\hat{\alpha}_{JRR}(k)) = \sigma^2(I - k^2B^{-2})\Lambda^{-1}(I - k^2B^{-2})' \quad \dots (24)$$

ومصفوفة متوسط مربعات الخطأ لمقدر (JRR) بالشكل الآتي :

$$\begin{aligned} MSE(\hat{\alpha}_{JRR}(k)) &= Var(\hat{\alpha}_{JRR}(k)) + [Bias(\hat{\alpha}_{JRR}(k))][Bias(\hat{\alpha}_{JRR}(k))]' \\ &= \sigma^2(I - k^2B^{-2})\Lambda^{-1}(I - k^2B^{-2})' + k^4B^{-2}\hat{\alpha}\hat{\alpha}'B^{-2} \quad \dots (25) \end{aligned}$$

ان الصيغة العامة لمقدر (RJRR) (Robust Jackknife Ridge Regression) لـ α يكون كالآتي :

$$\begin{aligned} \hat{\alpha}_{RJRR}(k) &= [I + kB]\tilde{\alpha}_{RRR} \quad \dots (26) \\ &= [I + kB^{-1}][I - kB^{-1}]\tilde{\alpha} \\ &= (I - k^2B^{-2})\tilde{\alpha} \end{aligned}$$

وان :

$$k = \frac{p\tilde{\sigma}^2}{\tilde{B}'\tilde{B}} \quad \dots (27)$$

وان :

$$\tilde{\sigma}^2 = \frac{(Y - X\tilde{B})'(Y - X\tilde{B})}{\tilde{B}'\tilde{B}} \quad \dots (28)$$



المقارنة بين بعض الطرائق الحصينة في ظل وجود مشكلتي تعدد العلاقة الخطية ونقاط الانعطاف العالية

علما إن معاملات الانحدار الاصلية لمقدر (RJRR) — α تعطى بالشكل الآتي :

$$\hat{\beta}_{RJRR} = \gamma \hat{\alpha}_{RJRR}(k)$$

أن مقدر (RJRR) متحيز للمعلمة الأصلية، ومقدار التحيز هو :

$$\begin{aligned} Bias(\hat{\alpha}_{RJRR}(k)) &= E[\hat{\alpha}_{RJRR}(k)] - \alpha \\ &= E[(1 - k^2 B^{-2})\tilde{\alpha}] - \alpha \\ &= (1 - k^2 B^{-2})E[\tilde{\alpha}] - \alpha \\ &= (1 - k^2 B^{-2})\alpha - \alpha \\ &= -k^2 B^{-2} \alpha \end{aligned} \quad \dots (29)$$

والتباين لمقدر (RJRR) يعطى بالشكل الآتي :

$$\begin{aligned} Var(\hat{\alpha}_{RJRR}(k)) &= E[\hat{\alpha}_{RJRR}(k) - E(\hat{\alpha}_{RJRR}(k))][\hat{\alpha}_{RJRR}(k) - E(\hat{\alpha}_{RJRR}(k))] \\ &= E[\hat{\alpha}_{RJRR}(k) - \alpha][(\hat{\alpha}_{RJRR}(k) - \alpha)]' \quad \dots (30) \\ &= (1 - k^2 B^{-2}) \Omega (1 - k^2 B^{-2})' \end{aligned}$$

ومتوسط مربعات الخطأ (MSE) لمقدر (RJRR) يكون بالشكل الآتي :

$$\begin{aligned} MSE(\hat{\alpha}_{RJRR}(k)) &= Var(\hat{\alpha}_{RJRR}(k)) + [Bias(\hat{\alpha}_{RJRR}(k))][Bias(\hat{\alpha}_{RJRR}(k))] \\ &= (1 - k^2 B^{-2}) \Omega (1 - k^2 B^{-2})' + k^4 B^{-2} \alpha \alpha' B^{-2} \end{aligned} \quad \dots (31)$$

من الصيغة (26) نعوض المعلمة $(\tilde{\alpha})$ بالمقدر (MM) الحصين حيث تكون الصيغة النهائية لمقدر (RJMM) كالآتي :

$$\begin{aligned} \hat{\alpha}_{RJMM}(k) &= [I + kB]\tilde{\alpha}_{MM} \\ &= [I + kB^{-1}][I - kB^{-1}]\tilde{\alpha}_{MM} \\ &= (1 - k^2 B^{-2})\tilde{\alpha}_{MM} \end{aligned} \quad \dots (32)$$

وايضا من للمعادلة (26) نعوض المعلمة $(\tilde{\alpha})$ بالمقدر الحصين (GM2) و تكون الصيغة النهائية للمقدر (RJGM2) كالآتي :

$$\begin{aligned} \hat{\alpha}_{RJGM2}(k) &= [I + kB]\tilde{\alpha}_{GM2} \\ &= [I + kB^{-1}][I - kB^{-1}]\tilde{\alpha}_{GM2} \\ &= (1 - k^2 B^{-2})\tilde{\alpha}_{GM2} \end{aligned} \quad \dots (33)$$

[1] (Alguraibawi, et.al, 2015:PP309-313)



المقارنة بين بعض الطرائق الحصينة في ظل وجود مشكلتي تعدد العلاقة الخطية ونقاط الانعطاف العالية

8- الجانب التجريبي

1-8- مفهوم المحاكاة The Concept of Simulation

يعد أسلوب المحاكاة من الأساليب العلمية الرصينة باعتباره أسلوباً للاختبار يقوم على أساس إعطاء صورة طبق الأصل لظاهرة حقيقية قبل تطبيق التجربة على بيانات واقعية، لغرض الاستفادة من هذه الصورة في دراسة خواص تلك الظاهرة ومميزاتها، حيث يعتمد أسلوب المحاكاة على إثبات البرهان الرياضي نظرياً. [1] (الجشعبي، 2007: 34)

إن أسلوب المحاكاة هو عملية تمثيل سلوك الظاهرة الحقيقية قيد الدراسة بشكل تمثل النموذج المدروس للواقع الحقيقي، ويمكن أن يعتمد أسلوب المحاكاة في إثبات صحة طريقة معينة يكون من الصعب إثبات صحتها نظرياً، وكلما كانت النتائج دقيقة دل ذلك على أن أسلوب المحاكاة يكون أكثر قرباً للواقع المدروس. [2] (شهاب، 2017: 43)

2-8- مراحل تطبيق تجارب المحاكاة

Stages of the application of simulation experiments

لقد تضمنت تجارب المحاكاة لهذا البحث كتابة عدد من البرامج بلغة (R) في توليد البيانات، واقترحت ثلاث حجوم للينة (n=20, n=50, n=100) ونسب تلوث مختلفة (5%, 15%, τ)، وقد تم استعمال الصيغة الآتية في توليد المتغيرات التوضيحية :-

$$X_{ij} = p v_{ik} + (1 - p^2)^{1/2} v_{ij} , \quad i = 1, 2, \dots, n , \quad j = 1, 2, \dots, k$$

v_{ij} : الإعداد العشوائية المولدة والتي تتبع التوزيع الطبيعي القياسي .

v_{ik} : يمثل قيم العمود الأخير من اعمدة المتغيرات المولدة .

k : يمثل عدد المتغيرات المرتبطة .

i : يمثل عدد المشاهدات .

ρ : يمثل قيمة الارتباط بين المتغيرات التوضيحية في النموذج المدروس، الذي يأخذ القيم (0.99, 0.95, 0.90) . [1] (Alguraibawi, et.al, 2015:P313)

ويكون النموذج كالاتي :

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_k X_{ik} + \varepsilon_i \quad ..(34)$$

$$i = 1, 2, \dots, n$$

يتم توليد الأخطاء العشوائية وفقاً للتوزيع الطبيعي بمتوسط (صفر) وتباين (σ^2) أي أن :

$$\varepsilon_i \sim (0, \sigma^2) , \quad i = 1, 2, \dots, n$$

وبالنسبة لقيم الانحراف المعياري فهي (0.5, 1, 1.5) .

يتم تحديد القيم الافتراضية للمعلمات من خلال دراسات سابقة وبافتراض أن عدد المعالم هي $k = 6$ وفقاً للقيم التالية :-

$$\beta_1 = 5 , \beta_2 = 3 , \beta_3 = \sqrt{6} , \beta_4 = 0 , \beta_5 = 0 , \beta_0 = 0 \quad [4] \text{ (Kan, et.al, 2013:P648)}$$

ثم يتم تقدير معلمات النموذج الانحدار الخطي المتعدد (MLR) وفق طرائق التقدير التي تم عرضها في الجانب النظري من البحث وهي كما يلي :-

1- طريقة انحدار الحرف (RR) لمقدر ال جاكنايف بالاعتماد على مقدر (MM) (RJMM) .

2- طريقة انحدار الحرف لمقدر ال جاكنايف بالاعتماد على مقدر (GM2) (RJGM2) .

وقد تم الاعتماد على المقياس الاحصائي متوسط مربعات الخطأ (MSE) للنموذج لمعرفة أي النماذج أفضل في تمثيل البيانات وكان عدد مرات تكرار التجربة (R=1000) مرة .



المقارنة بين بعض الطرائق الحصينة في ظل وجود مشكلتي تعدد العلاقة الخطية ونقاط الانعطاف العالية

3—8 تحليل نتائج تجربة المحاكاة

Analyzing the results of the simulation experiment

جدول (1) تقدير المعلمات و (MSE) للطريقتين ، ولكافة حجوم العينات عندما تكون نسبة التلوث ($\tau=5\%$) وقيمة الانحراف المعياري ($\sigma = 1$) .

Coffe	N	Parameters Methods	$\bar{\beta}_1$	$\bar{\beta}_2$	$\bar{\beta}_3$	$\bar{\beta}_4$	$\bar{\beta}_5$	MSE
$\rho = 0.90$	n=20	RJMM	0.213876	0.20999	0.20863	0.10212	0.12416	0.014047
		RJGM2	0.31949	0.27230	0.24319	0.06064	0.06622	0.01249
	n=50	RJMM	0.25065	0.17109	0.05731	0.10252	0.09237	0.00957
		RJGM2	0.48167	0.29385	-0.0299	0.00786	0.02436	0.00753
	n=100	RJMM	0.12598	0.11411	0.11846	0.10658	0.10446	0.00582
		RJGM2	0.21053	0.15107	0.16763	0.09859	0.09065	0.00551
$\rho = 0.95$	n=20	RJMM	0.19441	0.21181	0.21690	0.11222	0.13002	0.01186
		RJGM2	0.29125	0.26430	0.24253	0.07102	0.07328	0.01119
	n=50	RJMM	0.24999	0.17245	0.06736	0.11688	0.10058	0.00883
		RJGM2	0.50334	0.30061	-0.0621	0.01252	0.02111	0.00728
	n=100	RJMM	0.12718	0.11955	0.12334	0.11595	0.11441	0.00541
		RJGM2	0.19522	0.14513	0.16643	0.10926	0.10178	0.00522
$\rho = 0.99$	n=20	RJMM	0.16647	0.23724	0.25986	0.10714	0.12521	0.01026
		RJGM2	0.289963	0.26008	0.26526	0.03962	0.07908	0.01018
	n=50	RJMM	0.28012	0.18168	0.05966	0.12924	0.09955	0.00810
		RJGM2	0.66635	0.37027	-0.2073	-	-	0.00711
	n=100	RJMM	0.138178	0.13343	0.13719	0.13260	0.13108	0.00499
		RJGM2	0.18417	0.13747	0.17326	0.11980	0.10954	0.00496



المقارنة بين بعض الطرائق الحصينة في ظل وجود مشكلتي تعدد العلاقة الخطية ونقاط الانعطاف العالية

جدول (2) تقدير المعلمات و (MSE) للطريقتين، ولكافة حجوم العينات عندما تكون نسبة التلوث ($\tau=5\%$) وقيمة الانحراف المعياري ($\sigma = 1.5$).

Coffe	N	Parameters	$\hat{\beta}_1$	$\hat{\beta}_2$	$\hat{\beta}_3$	$\hat{\beta}_4$	$\hat{\beta}_5$	MSE
		Methods						
$\rho = 0.90$	n=20	RJMM	0.18473	0.19654	0.20110	0.07738	0.10546	0.01520
		RJGM2	0.28340	0.25953	0.23804	0.04137	0.05015	0.01266
	n=50	RJMM	0.12111	0.10701	0.10109	0.09288	0.09013	0.01596
		RJGM2	0.22294	0.15039	0.12809	0.09257	0.08784	0.01593
	n=100	RJMM	0.08740	0.07964	0.08395	0.07675	0.07484	0.00753
		RJGM2	0.15491	0.10867	0.13112	0.07943	0.07174	0.00711
$\rho = 0.95$	n=20	RJMM	0.17208	0.20816	0.22058	0.08088	0.10686	0.01611
		RJGM2	0.26637	0.26073	0.24855	0.04025	0.04742	0.01383
	n=50	RJMM	0.12702	0.11308	0.10811	0.10324	0.10132	0.01221
		RJGM2	0.22407	0.14417	0.12259	0.09836	0.09658	0.01156
	n=100	RJMM	0.09509	0.08913	0.09298	0.08756	0.08596	0.00673
		RJGM2	0.15100	0.10812	0.13389	0.08668	0.07878	0.00642
$\rho = 0.99$	n=20	RJMM	0.14250	0.25549	0.29369	0.05336	0.08750	0.01938
		RJGM2	0.28771	0.27179	0.30409	- 0.03643	0.03812	0.01960
	n=50	RJMM	0.30866	0.16922	-0.0168	0.08828	0.04639	0.01259
		RJGM2	0.79406	0.42444	-0.3773	- 0.11089	- 0.11054	0.01111
	n=100	RJMM	0.10998	0.10489	0.10967	0.10491	0.10274	0.00695
		RJGM2	0.15097	0.10279	0.14727	0.09297	0.07936	0.00693



المقارنة بين بعض الطرائق الحصينة في ظل وجود مشكلتي تعدد العلاقة الخطية ونقاط الانعطاف العالية

جدول (3) تقدير المعلمات و(MSE) للطريقتين ، ولكافة حجومات العينات عندما تكون نسبة التلوث $(\tau=15\%)$ وقيمة الانحراف المعياري $(\sigma = 1)$.

Coffe	N	Parameters Methods	$\bar{\beta}_1$	$\bar{\beta}_2$	$\bar{\beta}_3$	$\bar{\beta}_4$	$\bar{\beta}_5$	MSE
$\rho = 0.90$	n=20	RJMM	0.21387	0.20999	0.20863	0.10212	0.12416	0.01404
		RJGM2	0.31949	0.27230	0.24319	0.06064	0.06622	0.01249
	n=50	RJMM	0.25065	0.17109	0.05731	0.10252	0.09237	0.00957
		RJGM2	0.48167	0.29385	-0.0299	0.00786	0.02436	0.00753
	n=100	RJMM	0.12598	0.11411	0.11846	0.10658	0.10446	0.00582
		RJGM2	0.21053	0.15107	0.16763	0.09859	0.09065	0.00551
$\rho = 0.95$	n=20	RJMM	0.19441	0.21181	0.21690	0.11222	0.13002	0.01186
		RJGM2	0.29125	0.26430	0.24253	0.07102	0.07328	0.01119
	n=50	RJMM	0.24999	0.17245	0.06736	0.11688	0.10058	0.00883
		RJGM2	0.50334	0.30061	-0.0621	0.01252	0.02111	0.00728
	n=100	RJMM	0.12718	0.11955	0.12334	0.11595	0.11441	0.00541
		RJGM2	0.19522	0.14513	0.16643	0.10926	0.10178	0.00522
$\rho = 0.99$	n=20	RJMM	0.16647	0.23724	0.25987	0.10714	0.12521	0.01026
		RJGM2	0.28996	0.26008	0.26526	0.03962	0.07908	0.01018
	n=50	RJMM	0.28012	0.18168	0.05966	0.12924	0.09955	0.00810
		RJGM2	0.66635	0.37027	-0.2073	-0.02308	-0.02941	0.00711
	n=100	RJMM	0.13817	0.13343	0.13719	0.13260	0.13108	0.00499
		RJGM2	0.18417	0.13747	0.17326	0.11980	0.10954	0.00496



المقارنة بين بعض الطرائق الحصينة في ظل وجود مشكلتي تعدد العلاقة الخطية ونقاط الانعطاف العالية

جدول (4) تقدير المعلمات و(MSE) للطريقتين ، ولكافة حجومات العينات عندما تكون نسبة التلوث $(\tau=15\%)$ وقيمة الانحراف المعياري $(\sigma = 1.5)$.

Coffe	N	Parameters Methods	$\hat{\beta}_1$	$\hat{\beta}_2$	$\hat{\beta}_3$	$\hat{\beta}_4$	$\hat{\beta}_5$	MSE
$\rho = 0.90$	n=20	RJMM	0.18504	0.19528	0.20032	0.07817	0.10580	0.02321
		RJGM2	0.28380	0.25765	0.23790	0.04353	0.04924	0.02215
	n=50	RJMM	0.21261	0.13888	0.02280	0.07539	0.06197	0.01385
		RJGM2	0.45060	0.26984	-0.0915	-0.01845	-0.00559	0.01145
	n=100	RJMM	0.08733	0.07976	0.08373	0.07653	0.07480	0.00776
		RJGM2	0.15466	0.10977	0.12987	0.12987	0.07952	0.00747
$\rho = 0.95$	n=20	RJMM	0.17286	0.20711	0.21896	0.08197	0.10722	0.02049
		RJGM2	0.26706	0.25900	0.24766	0.04342	0.04619	0.02019
	n=50	RJMM	0.23127	0.14732	0.02410	0.08570	0.06555	0.013187
		RJGM2	0.50844	0.29377	-0.1412	-0.02800	-0.01964	0.01123
	n=100	RJMM	0.09504	0.08914	0.09292	0.08743	0.08588	0.00735
		RJGM2	0.15087	0.10857	0.13355	0.08652	0.07867	0.00719
$\rho = 0.99$	n=20	RJMM	0.14458	0.25377	0.28908	0.05508	0.08939	0.01822
		RJGM2	0.29000	0.27218	0.29967	-0.03164	0.03520	0.01854
	n=50	RJMM	0.30642	0.16823	-0.0143	0.08929	0.04773	0.01259
		RJGM2	0.78643	0.41721	-0.3690	-0.10590	-0.10728	0.01113
	n=100	RJMM	0.10998	0.10489	0.10967	0.10491	0.10274	0.00695
		RJGM2	0.15097	0.10279	0.14727	0.09297	0.07936	0.00693

9- مناقشة تجارب المحاكاة في حالة وجود تلوث ($\tau = 5\%$, 15%)

- من خلال النتائج المبينة بالجدول (1) (2) (3) (4) نلاحظ ما يلي :-
- 1- يتضح لنا ان قيمة متوسط مربعات الخطأ (MSE) للأنموذج تتناقص كلما زاد حجم العينة ولجميع طرائق التقدير المدروسة .
 - 2- في حالة في حالة وجود تلوث ($\tau = 5\%$) وتوزيع الاخطاء التوزيع الطبيعي بمتوسط (0) وانحراف معياري ($\sigma = 1$) نلاحظ من خلال مقياس المقارنة (MSE) للجدول (1) ان طريقة (RJGM2) افضل طريقة عند احجام العينات ($n=20, n=50, n=100$) وعندما تكون قيمة معامل الارتباط (0.99) ($p = 0.90, 0.95$) وذلك لانها حققت اقل قيمة ل (MSE) للأنموذج.
 - 3- وايضاً عند وجود تلوث ($\tau = 5\%$) وتوزع الاخطاء التوزيع الطبيعي بمتوسط (0) وانحراف معياري ($\sigma = 1.5$) للجدول (2) ومن خلال المقياس (MSE) للأنموذج نلاحظ ان طريقة (RJGM2) هي افضل طريقة عند احجام العينات ($n=20, n=50, n=100$) عندما تكون قيمة معامل الارتباط ($0.90, 0.95$) ($p = 0.90$)، ايضاً حققت طريقة (RJGM2) الأفضلية من خلال معيار المقارنة (MSE) للأنموذج وعند احجام العينات ($n=50, n=100$) عندما تكون قيمة معامل الارتباط ($p = 0.99$)، بينما حققت طريقة (RJMM) اقل قيمة ل (MSE) عند حجم العينة ($n=20$) .
 - 4- وفي حالة وجود التلوث ($\tau = 15\%$) وتوزيع الاخطاء بالتوزيع الطبيعي بمتوسط (0) وانحراف معياري ($\sigma = 1$) نلاحظ من خلال مقياس المقارنة (MSE) للجدول (3) ان طريقة (RJGM2) افضل طريقة عند احجام العينات ($n=20, n=50, n=100$) عندما تكون قيمة معامل الارتباط ($0.90, 0.95, 0.99$) ($p = 0.99$) وذلك لانها حققت اقل قيمة ل (MSE) للأنموذج، ايضاً حققت طريقة (RJGM2) للجدول (4) الأفضلية لا قل قيمة لمعيار المقارنة (MSE) عند احجام العينات ($n=20, n=50, n=100$) عند قيمة معامل الارتباط ($0.90, 0.95$) ($p = 0.90$)، وايضاً الأفضلية عندما تكون قيمة معامل الارتباط ($p = 0.99$) لا حجام العينات ($n=50, n=100$)، بينما حققت طريقة (RJMM) الأفضلية عند حجم العينة ($n=20$) كونها حققت اقل قيمة لمتوسط مربعات الخطأ (MSE) .

10-الاستنتاجات والتوصيات

10-1- الاستنتاجات

- 1- اظهرت نتائج المحاكاة لهذا البحث ان طريقة (RJGM2) هي الافضل، و تمتلك اقل قيمة لمتوسط مربعات الخطأ (MSE) مقارنة مع الطريقة الأخرى .
- 2- ايضاً اثبتت طريقة (RJGM2) اكثر كفاءة في حالة زيادة حجوم العينات ونسب التلوث العالية .
- 3- اثبتت طريقة (RJMM) كفاءتها عند حجوم العينات الصغيرة .

10-2- التوصيات

- 1- استعمال طريقة (RJGM2) في تقدير معلمات أنموذج الانحدار الخطي المتعدد (MLR) وباختلاف احجام العينات لما تبديه من كفاءة ومرونة في التطبيق .
- 2- اعتماد طريقة (RJGM2) في تقدير معلمات أنموذج الانحدار الخطي المتعدد (MLR) في حالة احجام العينات الكبيرة ونسب تلوث العالية .
- 3- اعتماد طريقة (RJMM) عند حجوم العينات الصغيرة .



المقارنة بين بعض الطرائق الحصينة في ظل وجود مشكلتي تعدد العلاقة الخطية وتقاط الانعطاف العالية

11- المصادر العربية

1. الجشعمي, حسين علي عبد الله 2007 م، " مقارنة لبعض المقدرات الحصينة لمعالم النماذج اللاخطية " اطروحة دكتوراه في الاحصاء كلية الادارة والاقتصاد الجامعة المستنصرية .
2. شهاب ضمياء حامد (2017) م. " مقارنة بعض طرائق التقدير الحصينة مع أسلوب بيز في تقدير دالة الانحدار اللوجستي مع تطبيق عملي " اطروحة دكتوراه في الاحصاء كلية الادارة والاقتصاد الجامعة المستنصرية.
3. كاظم، أموري هادي و مسلم، باسم شليبية 2002م، " القياس الاقتصادي المتقدم النظرية والتطبيق "، مكتبة دنيا الامل ، بغداد.

12- المصادر الاجنبية

- 1.Alguraibawi, M., Midi, H., & Rana, L. S. (2015) "Robust Jackknife Ridge Regression to Combat Multicollinearity and High Leverage Points in Multiple Linear Regressions" Economic Computation and Economic Cybernetics Studies and Research, NO. 4, PP 305-322.
- 2.Duran, E. A ., and Akdeniz, F. (2010) "Efficiency of the Modified Jackknifed Liu-Type Estimator" Statistical Papers, NO. 53, PP 265-280.
- 3.Kamruzzaman, MD, and Imon, A.H.M.R. (2002) "High leverage point, another source of multicollinearity" Pakistanian Journal of Statistics, NO. 18, PP 435-448.
- 4.Kan, B., Alpu, Ö., & Yazıcı B. (2013) "Robust Ridge and Robust Liu Estimator for Regression Based on the LTS Estimator" Journal of Applied Statistics, NO. 40 (3), PP 644-655.
- 5.Lukman, A. F., Osowole, O. I. & Ayinde, K.(2015) "Two Stage Robust Ridge Method in a Linear Regression Model" Journal of Modern Applied Statistical Methods, NO. 14(2), PP 53-67.
- 6.Midi, H., and Mohammed .M.A (2015) "The Identification of Good and Bad High Leverage Points in Multiple Linear Regression Model" Mathematical Methods and System in Science and Engineering, PP 147-153.
- 7.Singh, B., Chaubey, Y. P., & Dwivedi, T. D. (1986) "An almost unbiased ridge estimator" The Indian Journal of Statistics, NO.48 (3), PP 342-346.



Comparison of some robust methods in the presence of problems of multicollinearity and high leverage points

Ghufran Ismail Kamal / Assistant Prof- College Of Administration & Economics- University Of Baghdad-Statistics section - ghufran62@gmail.com .

Researcher Saif alemam saadi khazaal / College Of Administration & Economics- University Of Baghdad-Statistics section- saif.alemams@gmail.com

Abstract

The multiple linear regression model of the important regression models used in the analysis for different fields of science Such as business, economics, medicine and social sciences high in data has undesirable effects on analysis results . The multicollinearity is a major problem in multiple linear regression. In its simplest state, it leads to the departure of the model parameter that is capable of its scientific properties, Also there is an important problem in regression analysis is the presence of high leverage points in the data have undesirable effects on the results of the analysis , In this research , we present some of the robust methods in the multiple linear regression model These methods include the (Jackknife Ridge regression) methods based on the (MM) estimator and the (GM2) estimator (Modified Generalized M-estimator) . Using the Monte Carlo simulation, the two methods were compared in accordance with the comparison criterion, the mean squares error (MSE) and sample sizes ($n = 20, n = 50, n = 100$) and different pollution ratios ($\tau = 5\%, 15\%$) , The comparison shows that (RJGM2) is the best method for estimating the parameters of the multiple linear regression model, which has the lowest value for mean squares error (MSE) compared with the rest of the other estimations.

Keywords : Multiple Linear Regression , Multicollinearity, high leverage point, Jackknife ridge regression, MM-estimator, GM2-estimator.