



دراسة مقارنة بين طريقة GM6.IDRGP(RMVN) وطريقة GM6 لتحليل انحدار الخطى المتعدد

في ظل وجود القيم الشاذة واستعمال أسلوب المحاكاة

A comparative study between the two methods GM6.IDRGP(RMVN) and the GM6 method for analyzing the multiple linear regression model in the presence of outliers using the simulation method.

<https://doi.org/10.29124/kjeas.1549.16>

أ. م. د طارق عزيز صالح⁽¹⁾ فرحان يعقوب عذاب⁽²⁾

كلية الإدارة والاقتصاد/ جامعة واسط

المستخلص

يُعد الانحدار الخطى المتعدد أحد أكثر طرائق التحليل الإحصائى استعمالاً في العديد من المجالات العلمية. إذ يتم تقدير معلماته بالاعتماد على طريقة المربيات الصغرى الاعتيادية، التي تعطي أفضل تقدير خطى غير متخيّر في حالة تحقق فرضياتها. من أهم هذه الفرضيات هو أن يكون لها توزيع طبيعي للخطأ بمتوسط صفر وتبين ثابت. إذا كانت البيانات لا تفي ببعض الافتراضات، فقد تكون تقديرات العينة والناتج مُضللة. إذ يكون نموذج الانحدار الخطى حساساً لظهور القيم الشاذة ونقطات الرافعه. لذلك تم تطوير تقنيات إحصائية قادرة على التعامل مع القيم الشاذة أو اكتشافها. ما أدى إلى ظهور العديد من طرائق البديلة عن طريقة OLS مثل M-Huber و S و MM و LTS و IDRGP.RMVN و DRGP. وذلك اقتربنا في هذه الدراسة توظيف طريقة IDRGP.RMVN في خوارزمية طريقة GM6. ومقارنتها مع بعض طرائق الانحدار الحصين عن طريق دراسة محاكاة في تحديد أفضل طرائق.

المصطلحات / الانحدار الخطّي المتعدد، نقاط الرفع العالية، مقدّر DRGP و GM6 و IDRGP.RMVN

.GM6.IDRGP(RMVN)

Abstract

Multiple linear regression is one of the most widely used statistical analysis methods in many scientific fields. Its parameters are estimated based on the ordinary least squares method. Which gives the best unbiased linear estimate if its assumptions are met. The most important of these assumptions is that it has a normal distribution of error with a mean of zero and a constant variance. If the data does not meet certain assumptions, the sample estimates and results may be misleading. The linear regression model is sensitive to the appearance of outliers and leverage points. Therefore, statistical techniques have been developed capable of dealing with or detecting outliers. This led to the emergence of many alternative methods to the OLS method, such as M-Huber, s, LTS, and MM, which have high efficiency and breakdown points, but are affected by HLPs, which result in the problem of Masking and Swamping. The GM6 method is one of the methods that was developed in order to treat such problems through the use of a weight function, but the weight function depends on the individual diagnosis, which gives inaccurate results. In order to overcome this problem, comprehensive diagnostic methods have been proposed, such as the DRGP and IDRGP.RMVN methods. Therefore, in this study, we proposed to employ the IDRGP.RMVN method in the GM6 method algorithm. And comparing it with some hippocampal regression methods through a simulation study in determining the best methods.

Key wards: Multiple Linear Regression, Leverage Point, DRGP Estimator and IDRGP.RMVN, methods GM6 and GM6.IDRGP(RMVN).

Introduction

- المقدمة

اهتم العديد من المؤلفين والباحثين في مجال تقدیر المعلمات بإيجاد أفضل مقدّر لهذه المعلمات. لذلك، تم تطوير العديد من الطرائق لهذا الغرض. يُعد نموذج الانحدار الخطّي المتعدد أحد النماذج الرياضية الأكثر استعمالاً لتحليل بيانات الحقيقي. إن طريقة المرّبعات الصغرى الاعتيادية (OLS) هي طريقة شائعة الاستعمال في الانحدار الخطّي المتعدد بسبب خصائصها المثلثيّة. ومع ذلك، إذا كانت البيانات تحتوي على قيم شاذة، فإن تقدیرات OLS تصبح غير فعالة. يؤدّي هذا إلى استعمال طرائق تقدیر بديلة أكثر دقة وكفاءة من التقدیر (OLS). يطلق على هذه الأساليب بالطرائق الحصين (Robust estimation). الهدف من الإحصائيات الحصينة هو طرائق إحصائية جديدة لا تتأثر بالقيم الشاذة وتكون أكثر عملية للبيانات

التي لا يتم توزيعها بشكل طبيعي. يمكن العثور على العديد من طرائق التقدير الحصينة مثل M و MM و LMS و LTS ومع ذلك، فإن هذه الأساليب تعاني من تأثير القيم الشاذة ونقاط الرافعه العالية (HLPs). التي تسبب بظهور ظاهرة الاخفاء والغمر Masking and Swamping التي تجعل من تقديرات هذه الطريقة غير دقيقة، مما يؤدي إلى عدم قدرتها على اكتشاف أنواع القيم الشاذة جميع و (HLPs) بدقة في مجموعة البيانات. الغمر (swamping) هو تشخيص قيمة نظيفة على انها قيم شاذة لكن بالحقيقة هي ليست بالقيمة الشاذة. اما ظاهرة الاخفاء (masking) هي تشخيص أحد القيمة الشاذة وتطغى على بقية القيم الشاذة التي لا يمكن تشخيصها بسبب تأثير القيمة المشخصة (Maroona and Yohai, 2006). إذ قدم Imon (2002) مقياساً جديداً للقضاء على الاخفاء والغمر. السمية بمقياس (GP). ومع ذلك، لم يتمكن مقياس GP من تحديد العدد الدقيق لنقطات الرافع العالية ولا يزال يعاني من تأثيرات الاخفاء والغمر. لذلك تم اقتراح خوارزمية جديدة من قبل Habshah و آخرون عام (2009) تسمى Diagnostic Generalized Robust Scale(DRGP). الغرض من هذا المقياس هو التشخيص الدقيق وتقليل تأثير هذه الظاهرة. لم تقمي الطريقة الأخيرة تماماً على تأثيرات ظاهرة الاخفاء والغمر، ولكنها قلل من هذا التأثير قدر الإمكان، لاسيما مع أحجام العينات الصغيرة وزيادة معدلات تلوث مجموعات البيانات. لمعالجة هذا الوضع في عام 2015، محمد وآخرون اقترحوا تحسين الطريقة السابقة من خلال طريقة جديدة تسمى Improvised DRGP (IDRGP) التي رفعت معدل التشخيص والتقليل من معدل تأثير الاخفاء والغمر إلى أدنى حد ممكن لكنها لم تتمكن من التخلص منها بشكل نهائي. قام كل من (Uraibi and Alhussieny, 2022) و بالاعتماد على طريقة RMD (RMVN) المقترحة من (Olive and Hawkins, 2010) بتطوير طريقة DRGP من خلال دمج RMVN مع RMD من MVE وكذلك استعمال فكرة (Mohammad et.al, 2015) بتحسين الطريقة (DRGP) من خلال طريقة جديدة أطلق عليها (IDRGP.RMVN) سمية Improvised DRGP. لتكون أكثر كفاءة في كشف ظاهر الاخفاء والغمر وتقليلها والحد من تثيرها. مما تقدم يمكن تلخيص مشكلة البحث بوجود قيم شاذة ونقاط رافعة عالية توثر على نتائج التقدير أما هدف هذا البحث هو إيجاد طرائق بديلة تمتاز بكافأة و نقطة انهيار عالية ولا التأثر بهذه القيم. في القسم 2 يتم وصف القيم الشاذة واهم الطرائق المستخدمة في تشخيصها اما في القسم 3 يتم شرح للطرائق الحصينة وفي القسم 4 تكون دراسة المحاكمات ونتائج المستحصلة عليها من المحاكاة.

2- القيم الشاذة في النموذج الانحدار الخطى المتعدد وطرائق تشخيصها

وقد صنفَ كلَّ من (Rousseeuw & Zomeren 1990)^[12] القيم الشاذة على ثلاث لـما يخصّ موقعها وتأثيرها مشاهدات الشاذة لمتغير (Y) وهي القيم الشاذة العمودية (Vertical Outliers) او (VO) وهي المشاهدات الشاذة أو البعيدة في تنسيق متغير المعتمد او متغير الاستجابة (y). مشاهدات الشاذة لمتغير (X) نقاط التأثير المرتفعة (HLPs)^[3] وهي المشاهدات الشاذة التي تكون في قيم المتغير (X) او البعيدة عن مجموعة المتغيرات التوضيحية وهي قيم خارجية توثر على المتغير (X) والتي تؤدي إلى انعطاف ميل الانحدار نحوها وتكون مرتبطة بالمتغيرات التوضيحية. وتكون على نوعين: الأولى تسمى (Bad-Leverage) او (BLPs) نقاط تأثير سيئة وهي المشاهدات التي تكون بعيدة عن خط الانحدار أو مركز بيانات المتغير التوضيحي وتكون ذات تأثير على القيم المحسوبة تقدير متغير المعتمد. أما الثانية فتسمى (Good-Leverage) او (GLPs) نقاط التأثير الجيدة وهي المشاهدات التي تكون بعيدة عن بيانات المتغير التوضيحي. وتبقى على مقربة من خط الانحدار وتسهم في كفاءة التقدير وليس لها تأثير على تقديرات المتغير المعتمد. وهناك نوع ثالث يكون في

المتغير المعتمد (Y) والمتغيرات التوضيحية (X) [11]. ان لتشخيص القيم الشاذة دور مهم في معرفة اين بقع تأثير تلك القيم بالتحديد ليتم التعامل معها بحرفية اذ تم وضع طرائق للكشف عنها سواء أكانت بالمحور العمودي او الافقى. لذلك اهتم الباحثون كثيرا في قضية الكشف عن HLP، مثلًا قدم (Hadi) عام (1992) طريقته المعروفة Hadi's Potential للكشف عن هذه المشاهدات المؤثرة على التقديرات.

Hat Matrix -1-2

هي مصفوفات الوزن وبذلك تسمى بـ (Hat Matrix) [5] ويتم ايجاد هذه المصفوفة في تحديد صروف المشاهدات لـ (X) التي تحتوي القيم الشاذة ونُعرّف تحليل الانحدار كمقياس للكشف عن وجود القيم الشاذة وتكون الصيغة لها كالتالي :

$$H = X(X'X)^{-1}X' \quad (1)$$

إذ أن عناصر القطر الرئيس لمصفوفة الوزن (H) تكون صيغتها كالتالي :

$$h_{ii} = x_i(X'X)^{-1}X'_i \quad (2)$$

لهذه العناصر خصائص مفيدة ولاسيما عندما تكون قيمها تتراوح ما بين (1 و 0) وتكون مجموع (h_{ii}) يساوي (p). ومن الممكن إضافة الإثبات (h_{ii}) المؤشر على احتوائها للقيم الشاذة لـ (LP) لـ (i^{th}) إذ تكون مقياس للمسافة بين قيم لـ (x) للحالات جميعها (n)، لذلك تشير القيمة الكبيرة لـ (h_{ii}) إلى أن الحالة (i^{th}) بعيد عن مركز مشاهدات المتغير (x) [11]

2-الباقي القياسي المحوسبة Deleted Studentized Residuals

طريقة أخرى للتشخيص العمودي هي بقایا ستیوپنت المتبقي القياسي أي إن المتبقي القياسي هو مجرد بقایا محوسبة مقسومة على الانحراف المعياري لأخطاء $\hat{\sigma}$ مضروب في h_{ii} . $\sqrt{1 - h_{ii}}$ مصفوفة القطرية لمصفوفة القاعدة [14]

$$t_i = \frac{e_i}{\hat{\sigma} \sqrt{1 - h_{ii}}} \quad i = 1, 2, \dots, n \quad (3)$$

إذ إن $|t_i|$ إذا كانت أكبر من 2.5 فتعد قيمة شاذة.

3-طرائق التقدير الحصينة

عمل كثير من الباحثين في مجال الانحدار على طرائق التقدير البديلة للمرجعات الصغرى الاعتيادية عندما لا تتحقق طبيعة البيانات من الافتراضات الأساسية. إذ ناقش Alma (2011) [4] التقدير الحصين هو طريقة مهمة لتحليل البيانات الملوثة بالقيم الشاذة، الهدف الرئيس هو توفير مقدرات غير متخيزة ومرغوبة بوجود تلك القيم. من خلال استعمال خصائص الكفاءة ونقاط

الانهيار ونقطة الرافعة العالية لتحديد أداء التقييمات الحصينة. يتمثل أحد أهداف المقدرات الحصينة في الحصول على نقطة انهيار عالية محدودة للعينة والتي حددتها Donoho & Huber ، في عام 1983 [7]. تم تصميم طائق التقدير الحصينة بإذ لا تتأثر بشكل كبير بالقيم الشاذة. سوف نتطرق لبعض طائق التقدير الأكثر شيوعاً بشيء من التفصيل.

Break down Point

3-نقطة الانهيار

تُعدّ نقطة الانهيار من المعايير المهمة التي تقيس حصانة المقدرات. وتكون من الركائز الرئيسية في الحصانة إذ تُعدّ نقطه الفصل بين المقدّر ذي الفائدة من عديم الفائدة فكلاًما كانت نقطة انهيار المقدّر عالية كلما كانت المقاومة أكبر. وقد وصفها الباحث (Hampel) عام (1971) بأنّها (درجة حصانة المقدّر عندما تحتوي البيانات على قيم شاذة). وعرفها Donoho و Huber عام (1983)، Leroy و Rousseeuw عام (1987) [14] و Maronna و آخرون عام (2006) بأنّها هي أصغر كمية من التلوّث التي قد تتسبّب في أن يأخذ المقدّر قيمًا شاذةً كبيرةً بشكل تعسفي. كما عرفها Alma, (2011) [4] بأنّها النقطة أو النسبة المئوية المحدّدة للتلوّث في البيانات التي يَتَمّ عندها إغراق أي إحصائيات اختبار أولًا.

وعليه يمكن وصف نقطة الانهيار بأنّها الحدّ الذي يوضح مدى مقاومة المقدّر للقيم الشاذة ويمكن التعبير عنها بالشكل الرياضي الآتي التي تشير فيه كُلُّ من T إلى مقدرات الانحدار و Z إلى بيانات العينة و n إلى حجم العينة و m إلى النقطة الشاذة :

$$BP(T, Z) = \min \left\{ \frac{m}{n} : bais(m; T, Z) = \infty \right\} \quad \dots (4)$$

طريقة IDRGP(RMVN)

-1-3

هذه الطريقة المقترحة من قبل (Uraibi;Haraj;2022) [10] تدعو إلى توظيف مصفوفة الموقع والقياس Reweighted DRGP(MVE) في خوارزمية Multivariate Normal Estimators (RMVN) بدلًا من MVE. فضلًا عن الحال الخطوات التي أضافها Mohammad al et 2015 للتأكد من صحة التشخيص. إن مصفوفة RMVN هي من اقتراح Olive and Hawkins (2010) لإعادة ترجيح مقدّر متعدد المتغيرات الطبيعي من خلال خوارزمية سريعة ومتّسقة وتملك نقطة انهيار عالية. إذ قام الباحثان بالتركيز على خوارزمية ذات الخطوات الخمسة في المرحلتين الأولى و الثانية التي تتطلّب حساب مصفوفة DGK للبيان والتباين المشترك الذي اقترحها Devlin al et.(1981) و مصفوفة الموقع والقياس الـ Median Ball(MB) المفترضة من قبل Olive (2004). وبعدها يتمّ حساب مقدّر Fast FCH (consistent and high breakdown) الذي اقترحه Olive and Hawkins (2010) وفي الأخير يتمّ إعادة الترجيح مرة أخرى للحصول على المقدّر النهائي لهذه الطريقة. لقد وَظَفَ الباحثان (Uraibi;Haraj;2022) مقدرات هذه الطريقة لاقتراح خوارزمية جديدة بعنوان DRGP(RMVN) التي يمكن وصفها كالتالي:

1- حساب الوسط الحسابي ومصفوفة التباين والتباين المشترك كمقدرين تبدأ بهما الخوارزمية إذ يرمز لهما بالرموز

$T_{0,1}, C_{0,1}$ على التوالي.

2- استخراج مصفوفة التباين والتباين المشترك DGK: بعده يتم حساب مسافة مهالانوبيز (MD) بالاعتماد على المقدّرات في الخطوة (1) ثم نرتب القيم MD تصاعدياً للحصول على وسيط قيم MD كعبة للحصول على مصفوفة جيدة من المشاهدات تناظر تلك الصفوف التي قيم MD لها أقل من قيمة العتبة. ونقوم بتكرار الخطوتين السابقتين خمس مرات وفي كلّ مرة تحسب قيم MD ل الكامل مجموعة البيانات للحصول على ($T_{5,1}, C_{5,1}$) كمقدّرين للموضع والقياس حصينين.

3- إيجاد مقدّري الموضع والقياس لمصفوفة (MB) خوارزمية هذه المصفوفة تبدأ مع مقدّر الوسيط ومصفوفة الوحدة كمصفوفة للتباين والتباين المشترك كمراحل بداية لها. تختلف هذه المصفوفة عن سابقتها في أنها تستعمل في حساب MD الوسيط وليس الوسط الحسابي ونعيد الخطوات الخمس السابقة نفسها للحصول على ($T_{5,2}, C_{5,2}$)

4- مقدّرات FCH: الخطوة الأولى في خوارزمية هذه الطريقة هي حساب المسافة التقليدية بين مقدّري الموضع في الخطوتين السابقتين ($T_{5,1}, T_{5,2}$) على التوالي. فإذا كانت قيمة المسافة أقل من آخر قيمة للعتبة (وسيط مسافة مهالانوبيز في الخطوة الخامسة) في مصفوفة MB تبدأ FCH لاختيار أحد هذين المقدّرين كالتالي:

$$T_{FCH} \begin{cases} T_{5,2} & \text{if } \sqrt{|C_{5,1}|} < \sqrt{|C_{5,2}|} \\ T_{5,1} & \text{otherwise} \end{cases} \quad (5)$$

يتم حساب مقدّر القياس أيضاً بالاعتماد على الشرط نفسه ولكن المقدّر هنا يضرب بثابت أو معامل تصحيح كالتالي:

$$C_{FCH} \begin{cases} \frac{MED(MD_i^2(T_{5,1}, C_{5,1}))}{\chi_{(p,0.5)}^2} \times C_{5,1}, & \text{if } \sqrt{|C_{5,1}|} < \sqrt{|C_{5,2}|} \\ \frac{MED(MD_i^2(T_{5,2}, C_{5,2}))}{\chi_{(p,0.5)}^2} \times C_{5,2}, & \text{otherwise} \end{cases} \quad (6)$$

5- مقدّر RFCH: بعد الحصول على المقدّرين ($T_{1,RFCH}, C_{1,RFCH}$) من الخطوة السابقة يتم حساب مسافة مهالانوبيز مرة أخرى ل الكامل المشاهدة ثمّ نكرر الخطوة السابقة لإيجاد ($T_{2,RFCH}$) مع مراعاة معامل التصحيح و كالتالي:

$$C_{2,RFCH} = \frac{MED(MD_i^2(T_{1,RFCH}, C_{1,RFCH}))}{\chi_{(p,0.975)}^2} \times C_{1,RFCH}, \quad (7)$$

6- مقدّري الـ RMVN: خوارزمية هذه الطريقة تسعى أولاً إلى إيجاد مصفوفة مشاهدات جديدة للمتغيرات بالاعتماد على مقدّر الموضع والقياس السابقة و كالتالي:

افرض ان $S^0 = \sum_{j=1}^{n_1} X_k$ إذ ان

$$X_k = \left\{ X_k : MD_i(T_{2,RFCH}, C_{2,RFCH}) \leq MED(MD_i(T_{2,RFCH}, C_{2,RFCH})) \right\}$$

Where $k = 1, \dots, n_1$ and $i = 1, \dots, n$

لـ افرض ان $\{Q^{(1)} = \min\{0.5 \times 0.975 \times n/s^0, 0.995\}$ للحصول على مقدار القياس الأول وكالاتي RMVN

$$C_{RMVN}^{(1)} = \frac{MED(D_i^2(T_{RFCH}, C_{RFCH}))}{\chi^2_{(p, Q^{(1)})}} \quad (8)$$

و لتحسين أداء هذا المقدار اقترح Olive and Hawkins (2010) إعادة ترجيح مقدار القياس السابق

(8) مع الجزء الثاني من البيانات n_2 للحصول على المقدار النهائي $C_{RMVN}^{(2)}$.
 مما نصل إلى أن مقدار معلمة الموقع $T_{RMVN} = T_{RFCH}$ وأن مقدار معلمة القياس هو

- حساب مسافة مهالانوبير الحصينة:

$$RMD_i(RMVN) = \sqrt{(X - T_{RMVN}(X))' C_{RMVN}^{(2)}^{-1} (X - T_{RMVN}(X))}, \quad (9)$$

8- تحديد المشاهدات المشكوك بها بالمقارنة مع العتبة، $RMD_i(RMVN) > \sqrt{\chi^2_{p, 0.95}}$

و وضعها في المصفوفة D الموصوفة في طريقة GP ، أما بقية المشاهدات فتووضع في المصفوفة R.

9- نحسب p_{ii} كما في المعادلة (10) لمعرفة المشاهدات في المصفوفة D هل تحتوي على LP أم لا ، فإذا تم تشخيصها تبقى في المصفوفة D وبخلافه تنقل إلى المصفوفة R.

$$p_{ii} = \begin{cases} w_{ii}^{(-D)} & \forall i \in D, \\ \frac{w_{ii}^{(-D)}}{1 - w_{ii}^{(-D)}} & \forall i \in R, \end{cases} \quad (10)$$

أما الخطوات التي اقترحها Mohammad (2015) قد استخدمها Uraibi;Sawsan;2022 لتحسين أداء الطريقة المقترحة التي اطلق عليها IDRGP(RMVN) فلاحظ Mohammad (2015) أن تشخيص HLP في الخطوة الثانية من خوارزمية DRGP(RMVN) التي وضعها في مصفوفة جزئية أطلق عليها D لم يتم التحقق منها بشكل صحيح. لذلك اقترح إضافة خطوة جديدة إلى الخوارزمية من خلال التشخيص باستعمال matrix hat. تم مقاطعة التشخيصين للتحقق من التشخيص الأول، إذ يقارن فيها بين ما تم تشخيصه ك HLP بوصفها نتيجة نهاية للخوارزمية الأولى و التي افترضها

بالمصفوفة الجزئية D_2 و قارنها بما تم تشخيصه في الخطوة الثانية التي انتجت المصفوفة الجزئية D أي المقارنة بين D و D_2 كالتالي:

- الحالـة الأولى: إذا كانت المشاهـدات المشـخصـة كـ HLP هي ذاتـها في D و D_2 فالخطـوة الأخيرة لـخوارـزمـية الإعلـان عن هذا التشـخيص وـتـوقفـ.
- الحالـة الثانية: إذا كان عدد المشـاهـدـات المشـخصـة كـ HLP في D_2 هي أكـبرـ من عدد تلكـ المشـخصـة في D ، عندـئـذـ تـعملـ الخـوارـزمـية على إـعادـةـ المشـاهـدـاتـ التي لم يـتمـ تشـخيصـهاـ فيـ D_2 إلىـ المـصـفـوفـةـ R وـاحـدةـ ثـلـوـ الآـخـرـىـ علىـ انـيـتمـ حـاسـبـ p_{ii} لـكـلـ وـاحـدةـ مـنـهـاـ فيـ كـلـ مـرـةـ وـ تـقـارـنـ بـقـيمـةـ نـقـطـةـ الـقطـعـ أوـ العـتـبةـ فيـ المعـادـلـةـ (11)ـ فـإـذـاـ تـجاـوزـتـ قـيمـةـ العـتـبةـ تـضـافـ إـلـىـ المشـاهـدـةـ المشـخصـةـ وـ بـخـالـفـهـ تـبـقـىـ معـ المشـاهـدـاتـ النـظـيفـةـ.

$$\text{Median}(P_{ii}) \pm 3\text{MAD}(P_{ii}) \quad (11)$$

- الحالـةـ الثالثـةـ: إذاـ كانـ عـدـدـ المشـاهـدـاتـ المشـخصـةـ كـ HLP ـ فيـ D_2 ـ هيـ أـقـلـ منـ عـدـدـ تلكـ المشـخصـةـ فيـ D ـ،ـ أيـ تـمـ تشـخيصـ مشـاهـدـاتـ جـديـدةـ لـمـ تـشـخصـ سـابـقاـ،ـ عـندـئـذـ تـعـملـ الخـوارـزمـيةـ عـلـىـ الدـمـجـ بـيـنـ D ـ وـ D_2 ـ ثـمـ يـتمـ حـاسـبـ p_{ii} ـ لـكـلـ مشـاهـدـةـ فيـ المـصـفـوفـةـ الجـديـدةـ وـ تـقـارـنـ بـقـيمـةـ العـتـبةـ فيـ المعـادـلـةـ (23-2)ـ فـإـذـاـ تـجاـوزـتـ قـيمـةـ العـتـبةـ تـضـافـ إـلـىـ المشـاهـدـةـ المشـخصـةـ وـ بـخـالـفـهـ تـبـقـىـ معـ المشـاهـدـاتـ النـظـيفـةـ.

طـرـيقـةـ مـقـدرـ GM6 2-3

ان طـرـيقـةـ GM6ـ منـ الطـرـائقـ الحـصـينـةـ شـاسـعـةـ الـاستـعـمالـ لـأـنـهـاـ تـحـتـويـ عـلـىـ نـقـطـةـ انـهـيـارـ عـالـيـةـ تـسـاوـيـ تـقـرـيبـاـ 50%ـ وـكـفاءـةـ عـالـيـةـ وـتـأـثـيرـ مـحـدـودـ فيـ النـمـوذـجـ العـادـيـ وـالـتـيـ تـمـ اـقـتـراـحـهاـ منـ قـبـلـ العـالـمـانـ (Coakleyـ وـ Hettmanspergerـ)ـ فيـ عـامـ (1993)ـ وـ تـسـتـعـمـلـ المـرـبـعـاتـ الـأـقـلـ تـشـذـبـاـ كـمـقـدـرـ اـولـيـ فيـ خـوارـزمـيةـ GM6ـ.ـ يـعـرـفـ مـقـدـرـ GMـ بـأـنـهـ حلـ المـعـدـلـاتـ العـادـيـةـ التـيـ يـتـمـ تـقـديـمـهاـ مـنـ خـالـلـ:

$$\sum_{i=1}^n \pi_i \Psi \left\{ \frac{y_i - x_i^t \hat{\beta}}{\hat{\sigma} \pi_i} \right\} x_i = 0 \quad \dots (12)$$

إـذـاـ $\psi = p$ ـ هوـ مشـتقـ منـ دـالـةـ إـعادـةـ تـنـازـلـيـ (ـدـالـةـ الـوزـنـ).

$i = 1, 2, \dots, n$ ـ هوـ عـنـصـرـ الـوزـنـ الـأـولـيـ لـلـمـصـفـوفـةـ الـقـطـرـيـةـ W ـ

$\hat{\sigma}$ هو تقدير المقياس.

$\hat{\beta}$ هو متجه تقديرات المعلمات.

إن الهدف الرئيس لمقدرات GM^[8] يعتمد على تقليل نقاط الرفع العالية التي تحتوي على مخلفات كبيرة. في هذا الصدد يجب اكتشاف القيم الشاذة ونقط الرفع العالية في بداية الامر باستعمال طرائق تشخيص HLPs وب مجرد اكتشافها. بحسب تقليل اثارها لزيادة كفاءة مقدرات GM. يستعمل GM6 مسافات مهالانوبيز الحصينة (RMD) استنادا إلى الحد الأدنى من حجم او الحد الأدنى من محدد التغير (MCD) للكشف عن نقاط التأثير العالية. Ellipsoid(MVE)

قيمة النقطة الفاصلة كما اقترحها (Rousseeuw and Leroy (1987) هي $\sqrt{\chi^2_{p,0.5}}$. سيئ تحديد المشاهدة (i) التي

تنوافق مع RMD_i والتي تتجاوز نقطة القطع كنقطة رفع عالية.

خوارزمية مقرر GM6

الخطوة 1: احسب القيم المتبقية (r_i) بناءً على مقدر المربعات الأقل تشدبيا (LTS).

الخطوة 2: حساب المقياس المقدر (σ) للبواقي،

$$s = (1.4826)\left(1 + \frac{5}{(n-p-1)}\right)(\text{median}(|r_i|))$$

الخطوة 3: حساب القيم المتبقية المعيارية (e_i). إذ

الخطوة 4: نحسب الوزن الأولى ، والمشار إليه بالرمز π_i . إذ

$$\pi_i = \min \left\{ 1, \frac{\chi^2_{0.95,k}}{RMD(MVE)} \right\}$$

الخطوة 5: حساب دالة التأثير المحدود،

الخطوة 6: احسب خطوة واحدة نيوتن رافسون للحصول على تقديرات GM6

طريقة المقترحة GM6.IDRGP(RMVN) -3-3

من خلال اقتراح (Uraibi; Haraj) [10] في عام 2022 لطريقة IDRGP(RMVN) والتي استطاعت أن تقلل تأثير ظاهرة الاخفاء والغمر في البيانات. نقوم بتوظيف هذه الطريقة في اطار خوارزمي واحد مع GM6 مع تعديل خوارزمية GM6 لتسوّب الاوزان المستخدمة كافة في الطريقة كالتالي :

1- حساب بوافي (r_i) الانحدار بناء على مقدّر حصين (MM).

2- حساب المقدّر (S) للواقي.

$$s = (1.4826)\left(1 + \frac{5}{(n - p - 1)}\right)(\text{median}(|r_i|))$$

3- حساب معلمات المقدّر الحصين $\hat{\beta}_{Rob}$.

4- إيجاد معلمتى الموقف والقياس باستعمال RMVN.

5- حساب مسافة مهالانوبيز الحصينة RMD المعادلة (10) و تحديد المشاهدات المشكوك بها و وضعها في المصفوفة D الموصوفة في طريقة GP ، أمّا بقية المشاهدات فتوضع في المصفوفة R.

6- نحسب p_{ii} كما في المعادلة (9) لمعرفة المشاهدات في المصفوفة D هل تحتوي على LP أم لا ، فإذا تم تشخيصها تبقى في المصفوفة D و بخلافه تنقّل إلى المصفوفة R.

7- أية قيمة من قيم p_{ii} تتجاوز قيمة القطع في معادلة (10) تعد نقطة رافعة.

8- تضرب المشاهدات كُلّها بوزن أولي ل π_i حيث:

$$\pi_i = \min\left\{1, \frac{p_{ii}}{RMD(RMVN)}\right\},$$

9- نجد قيم المعلمات المقدّرة للبيانات المرجحة باستعمال طريقة المربيّات الصغرى $\hat{\beta}_{OLS}$ تتوقف الخوارزمية وتكون $|\hat{\beta}_{Rob} - \hat{\beta}_{OLS}| \leq 0.0001$ هي معلمة المقدّر الحصينة لطريقة GM6 إذا كانت لطريقة.

10- إذا لم يتحقق الشرط السابق نحسب بوافي الانحدار r_i^* و العودة إلى الخطوة 2.

دراسة أسلوب المحاكاة -4

ليكن لدينا نموذج الانحدار الخطّي المتعدد الآتي :-

$$y = X\beta + e \quad (15)$$

حيث أن X هي مصفوفة المتغيرات التوضيحية ذات بعد $(1 + n) \times p$ المولدة من توزيع طبيعي متعدد المتغيرات

بمتوسط قيمته صفرًا وتبين $\rho^{[i-j]} = \sigma$ أي أن

$$X \sim MVTN(0, \rho^{[i-j]})$$

إذ $n = \{30, 50, 90, 130\}$ بمعنى (4) متغيرات في دراسة المحاكاة. و $p = \{4\}$

$$\rho = 0.5$$

هي متتجه الوحدة لمعلمات هذا النموذج مع الحد الثابت كالتالي:

$$\beta = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}_{p \times 1}$$

اما α فهو متتجه الأخطاء العشوائية للنموذج بمتوسط قيمته صفر وانحراف معياري قيمته 2.

قام الباحث بتلويث بيانات كل دراسة محاكاة بنسب مختلفة $(0.15, 0.10, 0.05)$ من القيم الشاذة والنقط الرافعة. إذ تم

توليد عينة بحجم n مشاهدة وتلويث حد الخطأ العشوائي بـ α من القيم الشاذة من خلال توليد مشاهدات تتبع توزيع χ^2 بدرجة حرية (25). أما مصفوفة المتغيرات التوضيحية فقد تم تلويتها بـ α من BLP بوجود HLP واحدة. وذلك من خلال ضرب الصفوف الأولى من المتغير الثاني إلى المتغير الرابع بالرقم 10 و بما يتواافق مع نسب التلويث، اخيراً تضرب أعظم قيمة للمتغير الأول وما يناظرها في المتغير بـ بالرقم 50.

استعمل الباحث عدداً من المعايير التي ممكن توظيفها للمقارنة. المعيار الأول هو القيمة المطلقة لتحيز المعلمات المقدرة عن المعلمات المفترضة في دراستنا هذه. المعيار الثاني هو معيار متوسط مربعات الخطأ MSE والثالث هو متوسط انحرافات القيم المطلقة لبواقي الانحدار عن متوسطها MAE وأخيراً الوسيط المعياري لأنحرافات القيم المطلقة لبواقي الانحدار عن متوسطها MAD . عليه تكون الطريقة التي تجمع فيها أقل قيمة للمعايير كلها هي الطريقة الأفضل.

الجدول رقم (1) متوسط تقدير معلمات طائق التقدير $GM6, GM6.IDRGP.RMVN$ لهذه المعلمات لبيانات بحجم (30,50,90,130) مشاهدة ملوثة بـ $(0.15, 0.10, 0.05)$ من القيم الشاذة والنقط الرافعة بوجود نقطة رافعة عالية جداً.

	α	n	$\hat{\beta}_0$	$\hat{\beta}_1$	$\hat{\beta}_2$	$\hat{\beta}_3$	$\hat{\beta}_4$

GM6	0.05	30	1.817	3.755	- 0.232	0.647	0.715
GM6.IDRGP. RMVN			0.98	0.907	0.976	0.977	0.974
GM6	50	1.614	2.798	0.291	0.683	0.717	
GM6.IDRGP. RMVN		0.846	1.158	1.058	0.868	0.961	
GM6	90	1.359	2.467	0.283	0.593	0.591	
GM6.IDRGP. RMVN		0.76	0.904	1.11	0.88	0.998	
GM6	130	1.139	1.902	0.299	0.623	0.621	
GM6.IDRGP. RMVN		0.696	1.027	0.939	0.992	0.932	
GM6	0.10	30	2.293	3.786	- 0.207	0.695	0.288
GM6.IDRGP. RMVN			1.021	0.895	0.968	1.007	0.964
GM6	50	2.102	3.25	- 0.044	0.576	0.339	
GM6.IDRGP. RMVN		0.959	1.061	1.044	0.793	1.028	
GM6	90	1.605	2.97	0.029	0.441	0.316	
GM6.IDRGP. RMVN		0.851	1.122	0.931	0.986	0.931	

GM6	130	1.5	2.437	0.073	0.352	0.318
GM6.IDRGP. RMVN		0.772	0.959	0.914	1.01	0.968
GM6	30	2.28	2.937	0.463	0.435	0.563
GM6.IDRGP. RMVN		1.028	1.042	0.866	0.997	0.908
GM6	50	1.881	1.948	0.288	0.501	0.307
GM6.IDRGP. RMVN		0.937	1.081	0.816	0.919	0.851
GM6	90	1.553	1.749	0.331	0.307	0.4
GM6.IDRGP. RMVN		0.869	1.027	0.982	0.921	1.004
GM6	130	1.416	2.014	0.322	0.381	0.337
GM6.IDRGP. RMVN		0.754	1.01	0.945	0.986	0.909

الجدول (2) متوسط التحيّز معلمات طرائق التقدير GM6,GM6.IDRGP.RMVN لبيانات بحجم (30,50,90,130) مشاهدة ملوثة بـ (0.15,0.10,0.05) من القيم الشاذة و النقاط الرافعه بوجود نقطة رافعة عاليه جدا.

	α	N	$bias(\hat{\beta}_0)$	$bias(\hat{\beta}_1)$	$bias(\hat{\beta}_2)$	$bias(\hat{\beta}_3)$	$bias(\hat{\beta}_4)$
GM6	0.05	30	2.514	29.019	4.738	1.622	1.862
GM6.IDRGP. RMVN			0.093	0.394	0.535	0.327	0.294
GM6		50	0.718	9.824	1.42	0.458	0.349

GM6.IDRGP. RMVN			0.189	0.478	0.302	0.3	0.198
GM6	90	0.231	4.87	0.914	0.343	0.316	
GM6.IDRGP. RMVN		0.087	0.104	0.123	0.145	0.1	
GM6	130	0.052	1.454	0.598	0.274	0.234	
GM6.IDRGP. RMVN		0.11	0.091	0.084	0.104	0.069	
GM6	30	2.185	22.34	3.634	1.65	1.278	
GM6.IDRGP. RMVN		0.095	0.403	0.567	0.619	0.356	
GM6	50	1.786	16.148	1.975	0.833	0.778	
GM6.IDRGP. RMVN		0.058	0.283	0.29	0.359	0.186	
GM6	90	0.477	9.813	1.308	0.481	0.542	
GM6.IDRGP. RMVN		0.052	0.131	0.146	0.17	0.115	
GM6	130	0.304	3.967	0.977	0.459	0.515	
GM6.IDRGP. RMVN		0.124	0.773	0.588	0.238	0.276	
GM6	0.1 5	4.133	22.238	0.714	0.609	0.597	
GM6.IDRGP. RMVN		0.129	0.332	0.265	0.326	0.207	

GM6	50	1.006	7.863	0.607	0.367	0.63
GM6.IDRGP. RMVN		0.063	0.263	0.191	0.139	0.229
GM6	90	0.354	2.506	0.507	0.543	0.429
GM6.IDRGP. RMVN		0.041	0.11	0.074	0.088	0.066
GM6	130	0.222	2.563	0.494	0.419	0.476
GM6.IDRGP. RMVN		0.081	0.098	0.067	0.066	0.071

الجدول (3) المتوسط العام لمتوسط مربعات الخطأ معلمات طرائق التقدير GM6,GM6.IDRGP.RMVN لهذه المعلمات لبيانات بحجم (30,50,90,130) مشاهدة ملوثة بـ (0.15,0.10,0.05) من القيم الشاذة و النقاط الرافعه بوجود نقطة رافعة عالية جدا.

	α	N	$Mse(\hat{\beta}_0)$	$Mse(\hat{\beta}_1)$	$Mse(\hat{\beta}_2)$	$Mse(\hat{\beta}_3)$	$Mse(\hat{\beta}_4)$
GM6	30		3.526	29.02 3	5.391	2.39	2.66
GM6.IDRGP. RMVN			0.174	0.59	0.743	0.52	0.455
GM6	50		1.097	9.827	1.633	0.652	0.531
GM6.IDRGP. RMVN			0.093	0.245	0.313	0.292	0.221
GM6	90		0.362	4.871	0.961	0.396	0.365
GM6.IDRGP. RMVN			0.098	0.134	0.156	0.179	0.126
GM6	13		0.107	1.454	0.615	0.294	0.252

GM6.IDRGP. RMVN		0	0.115	0.108	0.104	0.123	0.086
GM6	0.1	30	3.399	22.34 6	4.214	2.299	1.634
GM6.IDRGP. RMVN			0.185	0.597	0.805	0.841	0.532
GM6		50	2.34	16.15	2.17	1.025	0.949
GM6.IDRGP. RMVN			0.094	0.366	0.384	0.451	0.255
GM6		90	0.649	9.814	1.344	0.518	0.57
GM6.IDRGP. RMVN			0.065	0.166	0.183	0.205	0.144
GM6		13	0.387	3.968	0.991	0.475	0.527
GM6.IDRGP. RMVN		0	0.072	0.114	0.144	0.101	0.102
GM6	0.1	30	5.237	22.24 3	1.029	0.925	0.931
GM6.IDRGP. RMVN			0.225	0.512	0.408	0.476	0.335
GM6		5	1.375	7.864	0.714	0.445	0.708
GM6.IDRGP. RMVN			0.099	0.334	0.249	0.19	0.282
GM6		90	0.464	2.507	0.523	0.557	0.444

GM6.IDRGP. RMVN			0.055	0.141	0.1	0.112	0.09
GM6	13 0	0.292	2.563	0.503	0.427	0.484	
GM6.IDRGP. RMVN		0.088	0.113	0.081	0.08	0.086	

الجدول(4) قيم MSE,MAE,MAD للطرائق GM6, GM6.IDRGP.RMVN لبيانات بحجم (30,50,90,130) مشاهدة ملوثة بـ 0.15,0.10,0.050 من القيم الشاذة و النقط الرافعة بوجود نقطة رافعة عالية جدا.

Method	α	n	MSE	MAE	MAD
GM6	0.05	30	23.224	2.807	3.024
GM6.IDRGP.RMVN			2.128	1.123	1.588
GM6	50	50	16.919	2.507	2.893
GM6.IDRGP.RMVN			1.607	1.045	1.557
GM6	90	90	8.936	1.8	1.961
GM6.IDRGP.RMVN			0.996	0.859	1.241
GM6	130	130	6.506	1.544	1.605
GM6.IDRGP.RMVN			0.71	0.732	1.037
GM6	0.10	30	33.138	3.453	4.068
GM6.IDRGP.RMVN			2.335	1.17	1.608
GM6	50	50	22.163	2.896	3.282
GM6.IDRGP.RMVN			1.7	1.076	1.549

GM6	90	11.044	2.147	2.215
GM6.IDRGP.RMVN		1.187	0.935	1.39
GM6	130	8.412	1.904	1.852
GM6.IDRGP.RMVN		0.904	0.82	1.242
GM6	30	59.576	4.712	5.274
GM6.IDRGP.RMVN		5.252	1.705	2.257
GM6	50	35.835	3.869	4.215
GM6.IDRGP.RMVN		2.177	1.201	1.68
GM6	90	15.088	2.697	2.631
GM6.IDRGP.RMVN		1.473	1.036	1.513
GM6	130	9.866	2.26	2.045
GM6.IDRGP.RMVN		1.191	0.938	1.381

يعرض الجدول (3,2,1) نتائج المحاكاة لتقدير معلمات الطرائق GM6, GM6.IDRGP.RMVN لـ (30,50,90,130) مشاهدة ملوثة بنسبي تلوث مختلف (0.15,0.10,0.05) من النقاط الرافعة والقيم الشاذة بوجود نقطة رافعة عالية. الملاحظ أن متواسط مقدرات طريقة المقترحة GM6.IDRGP.RMVN هي الأقرب لقيم المعلمات الأولية المستعملة في دراسة المحاكاة الا وهو الواحد الصحيح لكل معلمة، كما نلاحظ ان طريقة المقترحة حققت أقل قيم من قيم \hat{bias} أقل من الواحد الصحيح. حصول الطريقة المقترحة GM6.IDRGP.RMVN على قيم $Mse(\hat{\beta})$ أقل بكثير من طريقة GM6. مما يشير إلى دقة تقدير معلمات الطريقة المقترحة هذا من جانب دقة تقدير المعلمات.

أمّا من جانب المعايير التي تخص المقارنة بين النماذج كما أسلفت فلنجاً الباحث إلى حساب متواسط القيمة المطلقة للتحيز لكل طريقة وعرضها في الجداول (4) الذي بين أن MSE و MAE و MAD الطريقة المقترحة تتفاوت عند ثبات نسبة التلوث وزيادة حجم العينة. وأنّ الطريقة المقترحة حققت أقل قيم من قيم المعايير الثلاثة مما يؤكد تفوق طريقة المقترحة GM6 على طريقة GM6.IDRGP.RMVN

الاستنتاجات: من خلال نتائج البحث تم التوصل إلى ما يأتي:

1- نلاحظ أنَّ الطريقة المقترنة (GM6.IDRGP.RMVN) كانت لها الأفضلية عند حجوم العينات الصغيرة(30) والمتوسطة(90,50) والكبيرة(130) بالاعتماد على معايير المقارنة (MAD وMAE وMSE).

النوصيات: من خلال نتائج البحث نوصي بالأتي:-

1- يُتم الاعتماد على الطريقة المقترنة (GM6.IDRGP.RMVN) لأنَّ لها الأفضلية بالتقدير ومقاومة النقاط الرفع العالية ولحجوم العينات جميعاً HLPs

المصادر

- 1- A A. S. Hadi,(1992) “A new measure of overall potential influence in linear regression,” Computational Statistics & Data Analysis, vol. 14, no. 1, pp. 1–27.
- 2- Ali, D. A and Habshah. M,(2020)” On the Robust Parameter Estimation Method for Linear Model with Autocorrelated Errors in the Presence of High Leverage Points and Outliers in the Y-Direction” Malaysian Journal of Mathematical Sciences 14(3): 505-517.
- 3- Bagheri, A., and Midi , H. (2009) “Robust Estimations as a Remedy for Multicollinearity Caused by Multiple High Leverage Points” Journal of Mathematics and Statistics,NO. 5(4), PP 311-321.
- 4- Alma, Ö. G. (2011) “Comparison of Robust Regression Methods in Linear Regression” Int. J. Contemp. Math. Sciences,NO. 6(9), PP 409-421.
- 5- Chave, A. D., and Thomson , D. J. (2003) “A Bounded Influence Regression Estimator Based on the Statistics of the Hat Matrix ” Journal of the Royal Statistical Society: Series C (Applied Statistics),NO.52 (3),PP 307-322.
- 6- Coakley.Clint W and Thomas P. Hettmansperger(1993)” A Bounded Influence, High Breakdown, Efficient Regression Estimator” Journal of the American Statistical Association, Vol. 88, No. 423 , pp. 872- 880.
- 7- D.L. Donoho and P.J. Huber(1983), The notion of breakdown point. In: Bickel PJ, Doksum KA, Hodges JL Jr (Editors), A Festschrift for Erich L. LehmannWadsworth, Belmont,(pp 157-184).
- 8- Habshah. M ,Shelan.S, Jayanthl.T, & Mohammed.A(2021):” Simple and Fast Generalized - M (GM) Estimator and Its Application to Real Data Set” Sains Malaysiana 50(3) 859-867.

- 9- Hassan S. Uraibi (2009) "DYNAMIC ROBUST BOOTSTRAP ALGORITHM FOR LINEAR MODEL SELECTION USING LEAST TRIMMED SQUARES"**Universiti Putra Malaysia, in Fulfilment of the Requirements for the Degree of Master Science.
- 10-Hassan Uraibi and Sawsan Haraj (2022)" GROUP DIAGNOSTIC MEASURES OF DIFFERENT TYPES OF OUTLIERS IN MULTIPLE LINEAR REGRESSION MODEL"** Malaysian Journal of Science 41 (Special Issue 1): 23-33.
- 11-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.
- 12-Rousseeuw, P. J. and Leroy, A. M. (2005) Robust Regression and Outlier Detection,** John wiley & sons(p218).
- 13- Rousseeuw, P.J., Van Zomeren, B.C., 1990. Unmasking multivariate outliers and leverage points.** Journal of the American Statistical Association 85, 633–639.
- 14-Yohai , V.J. (1987) "High breakdown point and high efficiency robust estimates for regression"** The Annals of Statistics, NO. 15 (20), PP 642-656.