

# المجلة العراقية للعلوم الإحصائية



http://stats.uomosul.edu.iq

# الكشف عن القيم المتطرفة في نموذج الانحدار الخطى مع التطبيق على بيانات تلوث مياه الآبار أطراف مدينة الموصل

الافضل في الكشف من بين الطرق التي تم استخدامها.

سجى مروان إسماعيل 🌕 صفوان ناظم راشد

قسم الاحصاء والمعلوماتية، كلية علوم الحاسوب والرباضيات، جامعة الموصل، الموصل، العراق

#### الخلاصة

# معلومات النشر

تاريخ المقالة: تواستلامه في 2

تم استلامه في 22 شباط 2022 تم القبول في 30 نيسان 2022

تم القبول في 9 ايار 2022 متاح على الإنترنت في 1 حزيران 2022

الكلمات الدالة:

القيم المتطرفة، الكشف التشخيص، المعالجة، طريقة الحذف، مقدر M الحصين ،مقدر MM الحصين ومقدر M الحصين الموزونة

#### المراسلة:

سجى مروان إسماعيل

 $\underline{saja.20csp120@student.uomosul}\\ \underline{.edu.iq}$ 

DOI: 10.33899/IQJOSS.2022.174334, @Authors, 2022, College of Computer Science and Mathematics, University of Mosul. This is an open access article under the CC BY 4.0 license (http://creativecommons.org/licenses/by/4.0/).

تهتم فكرة البحث في التعرف عن تأثير القيم المتطرفة على معلمات نموذج تحليل الانحدار الخطى المتعدد، حيث يتم

الكشف عن القيم المتطرفة وتشخيصها والموجودة في البيانات إن كانت في المتغيرات المستقلة أو المتغير المعتمد مما

يتسبب في التأثير على تقدير معلمات النموذج المدروس . وقد تم التعرف على أنواع البيانات المتطرفة وطرق معالجتها

للحصول على نموذج أفضل ذات كفاءة عالية أو التقليل من تأثير هذه القيم على النموذج ، وتم وضع معيار متوسط

مربعات الخطأ MSE لغرض المقارنة بين طرق المعالجة وتم تطبيقها على بيانات حقيقية مأخوذة من مركز بحوث

السدود والموارد المائية جامعة الموصل ، وظهرت النتائج توضح أفضلية الكشف عن القيم الشاذة بطريقة الرسم الصندوقي Box-Plot كذلك أفضلية طرق المعالجة بطريقة مقدر M الحصين الموزونة المقترحة من قبل (Shaker ، 2009)

## مقدمة

من المعلوم أن التقدير في الطرائق الاحصائية يعتمد على مجموعة من الفروض المهمة للحصول على نموذج انحدار دقيق ، وتعد معلومية التوزيع الاحتمالي للبيانات أحدى أهم الأحيان تأخذ البيانات الموزعة نمطاً مختلفاً وقد لا تتمثل بنمط معين من التوزيعات والسبب يعود أحياناً الى وجود القيم الشاذة (Outlier) وهو الأمر الذي يؤدي الى عدم التحقق في فروض المربعات الصغرى وعندها ستفقد خصائصها الجيدة وعليه يتم البحث عن طرق بديلة حصينة لمعالجة هذه المشكلة وتعطينا مقدرات كفوءة ، وقد تم البحث عن طرق للكشف عن القيم الشاذة التي ظهرت في البيانات منها طريقة الرسم الصندوقي ، وطريقة مقدر M الحصين الذي وطريقة حذف ستيودنت ، وقد وضعت أساليب عديدة لمعالجة هذه القيم عند تقدير معلمات نموذج الانحدار أهمها أسلوب الحذف ، طريقة مقدر M الحصين الذي يعالج الخلل في المتغيرات التوضيحية ، كذلك استخدام طريقة مقدر MM الحصين الذي يعالج الخلل في المتغيرات المستقلة والمتغير التابع ، كذلك طريقة مقدر M الحصين الموزونة المقترحة من قبل (شاكر ، 2009) .

تضمن هذا البحث جانبين ، الأول مفهوم القيم الشاذة وطرق الكشف عنها وتشخيصها ، وطرق كشف التأثير عن القيم الشاذة ، وأساليب معالجة القيم الشاذة .أما الجانب الثاني تمثل بتطبيق الطرائق الثلاثة للكشف على بيانات حقيقية المتمثلة بتلوث مياه الابار أطراف مدينة الموصل وتم الحصول على البيانات من مركز بحوث السدود والموارد المائية -جامعة الموصل وتم وضع معيار متوسط مربعات الخطأ MSE للمقارنة .

## 1. نموذج الانحدار الخطى المتعدد:

تعد طريقة المربعات الصغرى الاعتيادية أحد الطرق في تقدير معلمات النموذج التي تأخذ العلاقة:

$$Y = X\beta + e \tag{1}$$

تشترط توفر عدد من الشروط من أجل الحصول على مقدرات دقيقة لمعلمات نموذج والذي يتمثل بالمعادلة الأتية:

$$\hat{\beta}_{\text{ols}} = (X'X)^{-1}X'Y \tag{2}$$

ومن هذه الشروط هي التوزيع الطبيعي لمتجه الاخطاء وتجانس تباين الاخطاء (Heteroscedastic ) إلا أن هذه الشروط لاتتحقق في كثير من الحالات والسبب يعود إلى وجود القيم الشاذة في متجه الأخطاء أو في قيم المتغيرات التوضيحية ، لذا فأن كشف هذه القيم ومعالجتها قبل تحليل البيانات أمر في غاية الأهمية من ألجل الحصول على نتائج واقعية دقيقة (Maronna and Martia ,2006) وقد وضعت أساليب للكشف عن هذه القيم منها:

# 2. بعض طرق الكشف عن المشاهدات الشاذة:

- a) الرسم الصندوقي : هو أحد الطرائق الاستكشافية الحديثة لتعيين القيم الشاذة وهي طريقة العرض بالرسم الصندوقي Box\_Plot بأستخدام خمس ملخصات المقترحة من قبل العالم (Tukey , 1977) ، وهو أحد الأساليب الرسومية التي توضح القيم ذات التطرف القوي (Extrem Outlier) والقيم ذات التطرف المعتدل وتحتويها البرامج والحزم الرسومية (الصائغ ,2013).
- b) فحص عناصر قطر المصفوفة (Hat- Matrix): هي أحدى الطرق الأحصائية التي من خلالها يتم الكشف عن القيم الشاذة في المتغيرات المستقلة والتي تسمى أحياناً بقيم قوة الرفع (Leverage values) (يوسف ،2015)إذا إن:

$$\mathbf{H}_{(n\times n)} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' = \begin{bmatrix} \mathbf{h}_{11} & \mathbf{h}_{12} & \dots & \mathbf{h}_{1n} \\ \mathbf{h}_{21} & \mathbf{h}_{22} & \dots & \mathbf{h}_{2n} \\ \vdots & \vdots & & \vdots \\ \mathbf{h}_{n1} & \mathbf{h}_{n2} & \dots & \mathbf{h}_{nn} \end{bmatrix}$$
(3)

. (Belsly and Walsch, 1980) فهذا يدل على ان المشاهدة (i فهذا يدل على ان المشاهدة  $h_{ii}>0.2$  or or  $h_{ii}>\frac{3p}{n}$  منظرفة فإذا كانت n

# ظرق تشخيص القيم الشاذة في المتغيرات المستقلة:

هناك ثلاث طرق لتشخيص المتغيرات المستقلة:

○ طريقة بليسلي وآخرون:أوضح (Belsly et al,1980 ) في هذه الطريقة عن المشاهدات الشاذة وبين فيما اذا كانت مؤثر ام لا على قوة الرافعة في المعادلة (3)
 كالأتي :

$$h_{ii} > \frac{2p}{n} \tag{4}$$

(Neter et al.,1990) إن الرافعة التي تزيد قيمتها عن 0.5 كبيرة

مربقة نيتر وآخرون: تشخيص القيمة الشاذة للمشاهدة i حيث اعتبر

كالاتى:

$$h_{ii} \ge 0.5 \tag{5}$$

oطريقة جون فوكس: اقترح (Fox,1997) دراسة كل الحالات التي تزيد قيم رافعاتها عن ثلاثة أضعاف متوسط الرافعة.

$$h_{ii} > \frac{3p}{p} \tag{6}$$

طريقة بواقي ستيودنت المحذوفة: \_ تمثل أحدى طرق الكشف عن القيم الشاذة في المتغير التابع ، والتي يتم الحصول عليها بايجاد القيمة المعيارية للبواقي المحذوفة (Deleted Residual)، حيث ان البواقي المحذوفة  $d_i$  للمشاهدة يساوي الفرق بين قيم ( $y_i$ ) الفعلية والقيم المقدرة لها ( $y_i$ ) باستخدام نموذج الانحدار الخطي الذي تقديره باستبعاد المشاهدة (i)، مما يجعل تحليل البواقي أكثر فاعلية في الكشف عن المشاهدات القاصية في المتغير التابع ( $y_i$ ) وتستند التوصل إلى أفضل معادلة تم الاعتماد عليها في الدراسة حيث تم حساب بواقي ستيودنت أستناداً إلى المعادلة أدناه والتي تتبع توزيع  $y_i$  بدرجة حرية ( $y_i$ ) وتستند في حسلبها على الخطأ  $y_i$ 0 ومجموع مربعات الخطأ  $y_i$ 2010 فضلاً عن قيم الرفع  $y_i$ 3 لمصفوفة الـ Hat-Matrix (المطيري، 2010):

$$d_{i}^{*} = e_{i} \left[ \frac{n - k - 1}{SS_{Res}(1 - h_{ii}) - e_{i}^{2}} \right] \sim t(n - k - 1)$$
(7)

#### طرق تشخيص القيم الشاذة في المتغير المعتمد:

يتم تشخيص القيم الشاذة لمتغير التابع بمقارنة القيمة المطلقة لباقي ستيودنت المحذوفة  $d_i^*$  بقيمة توزيع t عند درجة حرية t ومستوى معنوية t حيث تعتبر الحالة t حالة شاذة لابد من دراستها وتحديد مدى تأثير ها على مقدرات المربعات الصغرى(المطيري ,2010) .

$$\left|\mathbf{d}_{i}^{*}\right| = \mathbf{t}_{\alpha/2}, \mathbf{n} - \mathbf{k} - 1 \tag{8}$$

#### 3. المشاهدات الشاذة المؤثرة وطرق الكشف عنها:

هنالك مقاييس يتم من خلالها معرفة إذا كانت المشاهدات الشاذة مؤثرة أم لا من هذه لمقاييس:

# a. مقياس DFFITS لتأثير القيم الشاذة:

يستخدم مقياس DFFITS لقياس اثر المشاهدة i على القيمة المقدرة ، وتم الاعتماد على الصيغة أدناه لقياس أثر المشاهدة i على القيم المقدرة

$$(DFFITS)_{i} = d_{i}^{*} \left[ \frac{h_{ii}}{1 - h_{ii}} \right]^{\frac{1}{2}}$$
 (9)

# ♦ طرق كشف التأثير لقيمة DFFITS على النموذج:

يمكن تشخيص المشاهدة i باعتبارها مؤثرة على نتائج نموذج تحليل الانحدار الخطى المتعدد وهناك طرق كشف منها:-

# وطريقة بليسلى وآخرون (Belsley et al., 1980):

أختبر بليسلي الحالة مؤثرة على نتائج تحليل الانحدار الخطي المتعدد أستناداً إلى عدد معلمات النموذج p الموضح بالصيغة الأتية:

$$|DFFITS| > 2\sqrt{\frac{p}{p}}$$
 (10)

#### oطريقة شاترجي وهادي(Chatterjee and Hadi, 1988):

اقترح كل من شاترجي وهادي معيار لمقارنة القيمة المطلقة لـ DFFITS بقيمة اكبر قليلاً من القيمة التي اقترحها بليملي وإخرون تعرف كالآتي:

$$\left| \text{DFFITS} \right| > 2\sqrt{\frac{p}{n-p}} \tag{11}$$

#### أ. مقياس COVRATIO لتأثير القيم الشاذة:

نستخدم مقياس الاثر على الأخطاء المعيارية Influence on standard Error الذي طوره بليسلي اثر على حالة مصفوفة تباين تغاير معاملات الانحدار المقدرة. وبعد مراحل من التطور تم التوصل إلى معادلة حسابية تعتمد على قيم الرافعة  $h_{ii}$  وقيم بواقي ستيودنت المحذوفة  $d_i^*$  وعدد معلمات النموذج  $d_i^*$  وعدد المتغيرات  $d_i^*$  وقق الصيغة الأتية:

COVRATIO = 
$$\frac{1}{(1 - h_{ii}) \left(\frac{(n - k - 1) + d_{i}^{*2}}{n - k}\right)^{p}}$$
 (12)

نلاحظ في المعادلة تزيد قيمة COVRATIO بزيادة قيمة الرافعة وانخفاض بواقي ستيودنت المحذوفة ويكون ذلك مؤثر جيد لاكتشاف قيم COVRATIO المؤثرة على الأخطاء المعيارية.

# ❖ طربقة كشف التأثير لقيمة (COVRATIO):

اقترح (Belsley et al. , 1980) مقارنة قيمة COVRATIO بالقيمة  $\frac{3p}{n}$  لتشخيص اثر الحالة رقم (i) على الأخطاء المعيارية لمعاملات الانحدار أو COVRATIO فارج هذه الفترة فان المشاهدة رقم (i) تعتبر مؤثرة على قيم الأخطاء المعيارية لمعاملات نموذج الانحدار أعتماداً على الفترة الانية:

$$1 - \frac{3p}{n} < \text{COVRATIO} < 1 + \frac{3p}{n} \tag{13}$$

# c. مقياس Cook's Distance نتأثير القيم الشاذة:\_

يستخدم مقياس (مسافة كوك) لقياس أثر المشاهدة (i) على كل معاملات نموذج الانحدار المقدرة، وقد تم اعتماد مقياس مسافة كوك على القيمة  $e_i$ ,  $h_{ii}$  عندما تكون احد هاتان القيمتان كبيرة أو كلاهما فان قيمة مسافة كوك ستصبح كبيرة أيضاً.

$$D_{i} = \frac{e_{i}^{2}}{p * MSE} \left( \frac{h_{ii}}{(1 - h_{ii})^{2}} \right)$$
 (14)

# ❖ طربقة كشف التأثير لقيمة مسافة كوك 1: \_\_\_\_

طريقة فوكس: حيث اقترح فوكس طريقة في عملية الكشف القيمة الشاذة مؤثرة على قيم معاملات الانحدار وبخلاف ذلك تكون القيمة غير مؤثرة من خلال الصيغة الأثبية:

$$D_{i} > \frac{4}{n - p} \tag{15}$$

#### d. مقياس (DFBETAS) لتأثير القيم الشاذة:

يستخدم مقياس DFBETAS لقياس الفرق بين معاملات الانحدار المقدرة باستخدام كل المشاهدات وقيم معاملات الانحدار المقدرة بعد حذف المشاهدة رقم (i) في كل مرة. وهناك معادلة يجب حسابها في كل مرة يتم فيها توفيق النموذج بعد حذف المشاهدة رقم i مستنداً إلى الخطأ المعياري  $S_i$  وعنصر القطر  $S_i$  من مصفوفة  $S_i$  وفي الصيغة الاتية :

DFBETAS<sub>k(i)</sub> = 
$$\frac{b_k - b_{k(i)}}{S_i \sqrt{(X'X)_{kk}^{-1}}}$$
 for k=0,1,2,...,p (16)

♣ طریقـــة کشــف التـــأثیر

#### لقيمة(DFBETAS):

لتشخيص الحالات المؤثرة على قيمة معامل الانحدار.

صربقة نيتر وأخرون ( Neter

:(et al,1990

اقترح نيتر معياراً لتحديد الحالات المؤثرة في حالة العينات الصغيرة والمتوسطة  $|DFBETAS_{k(i)}| > 1$  أما في حالة العينات الكبيرة

$$\left| \text{DFBETAS}_{k(i)} \right| > \frac{2}{\sqrt{n}}$$

4. طرق معالجة القيم الشاذة:

ه. طريقــة حــذف المثــاهدات

الشاذة : تستخدم هذه الطريقة أذا كان حجم العينة كبيرا (إسماعيل،2001) حيث يتم بناء نموذج جديد في كل مرة يتم فيها حذف مشاهدة شاذة (i) لحين الحصول على نموذج يبلغ فيه قيمة الخطأ (MSE) أقل ما يمكن فضلاً عن قيمة معامل التحديد  $\mathbb{R}^2$  والمختبر الاحصائى  $\mathbb{R}$ .

b. Robust M Estimation : طربقة مقدر M الحصين

هو أحد طرق المعالجة الذي يتم من خلاله معالجة القيم المتطرفة في المتغير المعتمد وذالك بأستعمال الاساليب الحصينة التي أقترحها (Huber,1973) .وتم الاعتماد على الصيغة أدناه يمكن الحصول على معلمات النموذج وكما يلى:

$$\hat{\boldsymbol{\beta}}_{M} = (\mathbf{X}'\mathbf{W}_{M}\mathbf{X})^{-1}\mathbf{X}'\mathbf{W}_{M}\mathbf{Y} \tag{17}$$

هدف أسلوب M الحصين هو أعطاء أوزان صغيرة للمشاهدات غير الاعتيادية (المتطرفة ) من خلال عناصر القطر للمصفوفة القطرية W ، ولتطبيق أسلوب W يتطلب الأمر مقدر ابتدائي وأسلوب تكراري للوصول في النهاية إلى تقارب في مقدرات W للمعلمة W ، ويدعى هذا الأسلوب بأسلوب المربعات الصغرى الموزونة (RousseeandLeroy,1987).

#### c. طربقة مقدر MM الحصين:

هو أحد الأساليب الحصينة ذات الخصائص الجيدة والأكثر استخداما ، وهو أحد طرق المعالجة يستخدم لمعالجة القيم المتطرفة في المتغيرات المستقلة والمتغير التابع وهو مقدر يجمع بين الكفاءة النسبية المحاذية العالية لمقدر M مع نقطة انهيار عالية لنوع معين من مقدرات S ، حيث أن مقدر MM يقوم بحساب تباين الأخطاء من مقدر ابتدائي ذو نقطة انهيار عالية الا وهو مقدر S-Estimation ، بينما أسلوب M الذي كان يستخدم المربعات الصغرى الموزونة يأخذ الاخطاء من مقدر ابتدائي وهو مقدر المربعات الصغرى الاعتيادية الذي تكون نقطة انهياره 0% (شاكر ,2017).

وقد وصف (Yohai, 1987) ثلاث مراحل للحصول على مقدر MM:

المعادلة مقدر ابتدائي ذو نقطة انهيار عالية مثل مقدر S ، إذ يرمز له بالرمز  $\widetilde{\beta}$  ومن ثم استخدام أخطاء هذا المقدر من خلال المعادلة

$$e_i(\widetilde{\beta}) = y_i - x_i'\widetilde{\beta}$$

يتم أستخدام أخطاء المقدر  $\widetilde{\beta}$  والتي هي  $S(e_1(\widetilde{\beta})...\cdot e_n(\widetilde{\beta}))$  ويرمز لها بالرمز  $S(e_1(\widetilde{\beta})...\cdot e_n(\widetilde{\beta}))$  ويرمز لها  $S(e_1(\widetilde{\beta})...\cdot e_n(\widetilde{\beta}))$  ويرمز لها  $S(e_1(\widetilde{\beta})...\cdot e_n(\widetilde{\beta}))$  ويرمز لها بالرمز  $S(e_1(\widetilde{\beta})...\cdot e_n(\widetilde{\beta}))$  والتي هي هذه المرحلة يرمز لها بالرمز  $S(e_1(\widetilde{\beta})...\cdot e_n(\widetilde{\beta}))$  ويرمز لها بالرمز  $S(e_1(\widetilde{\beta})...\cdot e_n(\widetilde{\beta}))$ 

3. يتم في هذه المرحلة ايجاد مقدر MM من خلال

$$\begin{split} j=1,\dots,k\;,\quad \sum_{i=1}^n x_{ij} \psi_1(\frac{y_1-x_i'\beta}{S_n}) &= 0\\ \psi_1(e) &= \frac{\partial \rho_1(e)}{\partial e} \end{split}$$
   
 \(\psi\_i)

. (Yohai , 1987) في المرحلة الثالثة يجب أن لا تكون نفس دالة  $ho_0$  ، ولكن يجب أن تحقق الشروط الثلاثة (Yohai , 1987) .

والصيغة أدناه تمثل مقدر معلمات النموذج MM الحصين التي تم الاعتماد عليها في العمل.

$$\hat{\beta}_{MM} = (X'W_{MM}X)^{-1}X'W_{MM}X \tag{18}$$

# d. طريقة مقدر M الحصين الموزونة (R.M.W):

تعتبر هذه الطريقة أحد طرق معالجة القيم الشاذة الموجودة في متجه الأخطاء الأ أنه في حالة وجود مثل هذه القيم في مصفوفة المتغيرات التوضيحية يصبح هذا الأسلوب غير قادر على معالجة هذه القيم والتخلص من أثرها ، لذا أقترح (شاكر،2009) أولاً تعديل القيم المنظرفة الموجودة في مصفوفة المتغيرات التوضيحية باستعمال مصفوفة الأوزان لطريقة المربعات الصغرى الموزونة (weighted least squares (W.L.S) ومن ثم معالجة القيم المتطرفة الموجودة في متجه متغير الاستجابة من خلال استخدام (متجه أخطاء المربعات الصغرى الموزونة ) باستعمال أسلوب M الحصين ثانياً، وأخيراً إيجاد المقدرات الجديدة بعد التعديل الأخير وهذه المقدرات سيطلق عليها اسم M الحصينة الموزونة (Robust M-Weighted Estimator(R.M.W).

وبنفس طريقة المربعات الصغرى الموزونة الاعتيادية يمكن إيجاد M الحصين الموزون المقترح (شاكر،2009)، نحصل على المعادلة الأتية:

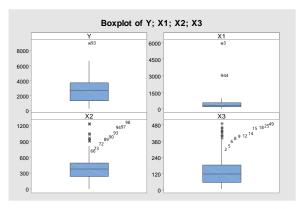
$$\therefore \hat{\beta}_{\text{RMW}} = (X'W_{\text{RMW}}X)^{-1}X'W_{\text{RMW}}Y \tag{19}$$

المعادلة (19) تمثل صيغة أسلوب M الحصين الموزون المقترح (R.W.M ) والذي يتم بواسطته معالجة التطرف الموجود في المتغيرات التوضيحية أو متغير الاستجابة أو كليهما معاً.

#### 5. الجانب التطبيقي:

في هذا الجانب من البحث تم تطبيق ما ورد في الجانب النظري على بيانات الدراسة تتمثل عن تلوث مياه الابار والعناصر الداخلة في التركيب والتي تتمثل بمتغيرات X وعناصر الحرى تتمثل بالمتغير Y والتي تحمل اعلى نسبة من الشواذ (2020) ومن هذه العناصر التي تم اعتمادها في الدراسة كمتغيرات مستقلة (الكالسيوم ( $(Ca^{2+})$ ) والبيكربونات ( $(Ca^{2+})$ ) والمغنيسيوم ( $(Mg^{+})$ ) ، والاملاح الذائبة ((D.S)) كمتغير معتمد. ومن خلال دراستنا استنتجنا بأن العلاقة بين كل متغير مستقل مع متغير معتمد علاقة طردية موجبة.

1. في الخطوة الاولى تم اختبار وجود القيم الشاذة في كلا المتغيرين ( المتغيرات المستقلة والمتغير التابع ) بأستعمال الرسم الصندوقي وهو أفضل طرق الكشف كما موضح في الشكل الاتي:



الشكل (1) يوضح القيم الشاذة في المتغرات المستقلة X والمتغير التابع Y

من خلال الشكل (1) تبين إن المتغير  $X_2$  ظهرت فيه قيمة شاذة واحدة أما في المتغير  $X_1$  ظهرت فيه قيمتين شاذتين كذلك المتغير  $X_2$  ظهرت فيه تسعة قيم شاذة والمتغير  $X_3$  يحتوى احدى عشر قيمة شاذة .

- 2. الخطوة الثانية معالجة القيم الشاذة : يتم معالجة القيم الشاذة بعدة طرق منها
- طريقة الحذف حيث تبين أن افضل نموذج تم الحصول عليه عند حذف المشاهدتين (3)،(93).
   نتائج النموذج المقدر بعد حذف المشاهدة رقم (3)(93):

$$\hat{\mathbf{Y}} = 0.0185 + 0.4988\mathbf{x}_1 + 0.3342\mathbf{x}_2 + 0.6651\mathbf{x}_3$$

$$R^2 = 79.88\%$$
 MSE = 0.1851 F = 134.98  $VIF_{x_1} = 1,16; VIF_{x_2} = 1.12; VIF_{x_3} = 1.06$  D.W = 1.47435

حيث تم الحصول على أقل MSE عند حذف المشاهدتين (3)،(3) مقارنتة بباقى النماذج.

• أما المعالجة بطريقة مقدر M الحصين : حيث تم اختيار نقطة انهيار مختلفة ولكل نقطة لها ثابت توليف معين (شاكر،2017) ، وتم التوصل الى افضل نموذج مقدر عند نقطة انهيار 50% كما موضح في الجدول الاتي:

تحديد باستخدام مقدر Mالحصين	ى لتباين الاخطاء ومعامل اا	، النموذج والجذر التربيع	جدول(1):مقدر معلمات
-----------------------------	----------------------------	--------------------------	---------------------

Obs	$\hat{eta}_{ m o}$	$\hat{oldsymbol{eta}}_{\scriptscriptstyle 1}$	$\hat{eta}_2$	$\hat{eta}_{\scriptscriptstyle 3}$	$\mathbb{R}^{2}$	MSE	7
1	-0.0076	0.8951	0.2052	0.6272	0.9910	0.8730	2971.80

حيث ظهرت قيمة ال MSE=0.8730 وهي قيمة صغيرة مقارنة بباقي النماذج كذلك قيمة  $^2$  = 0.9910 وهي قيمة تفسيرية تبين أن متغيرات الانموذج (الكالسيوم ،والبيكربونات ، والمغنيسيوم) قد فسرت بنسبة (0.991) من أجمالي التغيرات الحاصلة بالمتغير التابع (1.0.5) الأملاح الذائبة الكلية) .

• المعالجة بطريقة مقدر MM الحصين: تبين من هذه الطريقة التي تعالج مشكلة القيم الشاذة في المتغيرات المستقلة والمتغير التابع لتقليل أثر القيم الشاذة كما تم دراستها في الجانب النظري وفق المعادلة (18) حيث تم أختيار أنسب معالجة للقيم بناءا على الدراسات السابقة (شاكر 2017)عند نقطة أنهيار 50%وكفاءة قدرها 3.42 بمعامل توليف قدره 1.547 وكانت النتائج كما موضح في الشكل الاتي:

جدول(2): يوضح معاملات الانحدار ومتوسط مربعات الخطأ و معامل التحديد

$\hat{oldsymbol{eta}}_{ m o}$	$\hat{eta}_{_1}$	$\hat{eta}_2$	$\hat{eta}_3$	$\mathbb{R}^2$	F	MSE
0.0906	1.2398	0.1891	0.5353	0.9626	1295.69	1.5638

حيث تم الحصول على نموذج تكون فيه قيمة الMSE=1.5638 وقيمة أقل ما يمكن وقيمة معامل التحديد  $\mathbf{R}^2 = \mathbf{0.9626}$  وهي قيمة تفسيرية تبين أن متغيرات الثلاثة (الكالسيوم، البيكربونـات ،المغنيسيوم) للنموذج قد فسرت ما نسبته (0.9626) من أجمالي التغيرات الحاصلة بالمتغير التابع (الامـلاح الذائدة. $\mathbf{C}$ T.D.S).

• المعالجة بطريقة مقدر M الحصين الموزونة: في هذه الطريقة تم اختبار معلمات نموذج الانحدار M الموزونة استناداً إلى المعادلى (19) وفق أوزان وختبرات مختلفة فقد تم التوصل الى أفضل نموذج تم الحصول عليه عند أضافة الوزن  $W_{RMW4}$  حيث كانت قيمة الـ MSE = 0.022 وهي قيمة صغيرة مقارنة ببياقي النماذج كما في الجدول الاتي:

 $\frac{4}{n-p}$  عند الموزونة عند M الحصين الموزونة عند n-p

Dependent variable	Coefficients	Std.Error	T	p-value
Constant	0.048	0.049	0.981	0.329
$X_1(Ca^{2+})$	0.587	0.121	4.869	0.000
$X_2(HCO_3^-)$	0.347	0.047	7.368	0.000
$X_3(Mg^+)$	0.659	0.047	13.953	0.000

# المقارنة بين طرق معالجة القيم الشاذة:

بعد عدة طرق معالجة اجرية على بيانات التلوث تم التوصول الى اربع نماذج أنحدار يتم المقارنة بين هذه النماذج للحصول على افضل مقدر ذو كفاءة عالية تكون فيه قيمة الMSE أقل مايمكن كما هو مبين في الجدول الاتي:

	طريقة الحذف/(3+39)	الحصين /عند نقطة Mمقدر	الحصين/عنده MMمقدر	الحصينة Mمقدر			
	طريقة الكدف/(د+دو)	أنهيار (50%)	نقطة الانهيار وكفاءة (3.42)	الموزونة/0.038			
R-Squer	0.7988	0.9910	0.9528	0.9772			
MSE	0.1851	0.8730	1.5638	0.022			
F	134.98	2971.805	1295.696	117.646			

جدول (4): نتائج المقارنة بين طرق معالجة القيم الشاذة

من ملاحظة النتائج في الجدول (4) يمكننا القول بأن R.M.W قد عالج نسبة كبيرة من القيم الشاذة مقارنة بكل من طريقة الحذف ومقدر M الحصين ومقدر MM الحصين الحصين .

#### 6. الأستنتاجات:

كل ما تم الحصول عليه من نتائج وأستنتاجات وكفاءة مقدر وفق معيار المقارنة حسب طبيعة البيانات المدروسة حالياً والمأخوذة من مركز بحوث السدود والموارد المائية. أن أفضل طريقة لكشف المشاهدات الشاذة في بيانات نموذج الانحدار المتعدد باستخدام الرسم الصندوقي Box-plot. وأن معالجة القيم الشاذة أدى الى تحسين كبير جداً في أداء نموذج الانحدار الخطي المتعدد بالتنبؤ بقيم المتغير التابع بشكلها العام وتم التوصل الى أفضل نموذج عند معالجة القيم الشاذة بطريقة الحذف عند حذف المشاهدتين (3)،(92) حيث كانت قيمة MSE أقل مايمكن عند هذه الحالتين ، أما بطريقة مقدر M الحصين ظهر أفضل نموذج عند نقطة

انهيار قدرها 50% أما بطريقة مقدر 
$$M$$
 الحصين الموزونة حيث ظهر النموذج المفضل عند  $\frac{4}{n-p}$  وظهر تقوق لـ  $M$  الحصين بكفاءة  $3.42$  ونقطة أنهيار  $n-p$ 

قدرها %50 وهذا كان واضحا لحصوله على أقل MSE ، كما تبين أن أفضل طريقة لمعالجة القيم الشاذة والتي اظهرت كفاءتها العالية عن باقي الطرق هي

طريقة مقدر M الحصين الموزونة عند 
$$\frac{4}{n-p}$$
 عن باقي المختبرات.

#### Reference

- 1. Shaker ,Saleh Muayad,(2009)," Improving the hippocampal M method in estimating multiple linear regression model parameters", Iraqi Journal of Statistical Science, No.16,Pp.219-242.
- 2. Shaker ,Saleh Muayad,(2017)," Proposed robust methods for the median analysis of the linear regression model and their comparison with the ordinary least squares estimators using simulation",PHD. Thesis, University of Mosul, Mosul, Iraq.
- 3. AL-Mutery, Abed Al-Aziz Mnahe, (2010)," Methods for discovering anomalous and affecting observations on linear regression",King Saud University, AL-Reyad.
- 4. AL-Saeg, Mumen Amer Hsan,(2013)," The effect of outliers on the results of some statistical hypotheses",BSc. Thesis, University of Mosul, Mosul, Iraq.
- 5. Yusef, Isaam Al-deen Yusef Abd Alla,(2020)," The effect of outliers on the parameters of the multiple linear regression analysis model",PHD. Thesis, AL-Sudan University.
- Al-Youzbakey, K.T. and Sulaiman, A.M. (2020). "Hydrochemical Evaluation for Al-Sada Area Wells and their Suitability for Agricultural Usages", Journal of Umm Al-Qura University for Applied Science, Dams and Water Resources Researches Center, University of Mosul, Mosul, Iraq.
- 7. Belsley, D. et al. (1980). "Regression Diagnostics: Identifying Infuential Data and sources of Collinearity", Wiley, New York, p:105.
- 8. Chatterjee, S. and Hadi, A. S.(1988). "Sensitivity Analysis in Linear Regression", New York: john Wiley.
- 9. Fox, John,(1997). "Applied Regression Analysis, Linear Models, and Related Methods", Sage publications.
- 10. Neter, J. et al. (1990). "Applied Linear statistical Models: Regression, Analysis of Variance, and Experimental Designs". (3rd edition). Irwin, Homewood, IL 60430, Boston, MA 02116.
- 11. Rousseew P.J. and Leroy, A.M. (1987). Robust Regression and Outlier Detection. Wiley-Interscience, New York.
- 12. Tukey, J.W.(1977). "Exploratory Data Analysis", Addison-Wesley reading, MA.
- 13. Yohai, V.J. (1987). "High breakdown-point and high efficiency estimates for regression", The Annals of Statistics 15, 642-65.

# Detection Of Outliers In The Linear Regression Model With Application To Well Water Pollution Data On The Outskirts Of The City Of Mosul

Saja Marwan Esmaeel & Safwan Nathem Rashed Department of Informatics & Statistic, College of Computer Science, & Mathematics University of Mosul, Mosul, Iraq

**Abstract :**The research idea is concerned with identifying the effect of outliers on the parameters of the multiple linear regression analysis model. Where the outliers values that are present in the data are detected and diagnosed if they are in the independent variables or the dependent variable, which causes an impact on the estimation of the parameters of the studied model. The extreme data types and methods of processing them have been identified to obtain a better model with high efficiency or reduce the impact of These values on the model, The MSE standard was developed for the purpose of comparing treatment methods and was applied to real data taken from the Dams and Water Resources Research Center, University of Mosul. Suggested by (Shaker, 2009) is the best in detection among the methods that have been used.

**Keyword:** Outliers, Diagnostic Detection, Treatment, Elimination Method, Hippocampal M Estimator, Hippocampal MM Estimator and Weighted Hippocampal M Estimator