



المجلة العراقية للعلوم الإحصائية

<http://stats.uomosul.edu.iq>



طريقة المربعات الصغرى الموزونة الحصينة باستخدام صيغ مختلفة لمقدرات M (RWLSM) دراسة مقارنة

فاطمة محمد احمد¹ ، بشار عبدالعزيز الطالب²

^{1,2}قسم الاحصاء والمعلوماتية، كلية علوم الحاسوب والرياضيات، جامعة الموصل، الموصل، العراق

الخلاصة

تم في هذا البحث تقليل أو إستبعاد تأثير عدم تحقق فرض التوزيع الطبيعي للبيانات، بسبب وجود أنواع من القيم الشاذة فيها عند الرغبة في إختيار أفضل معادلة إنحدار بالطرق الحصينة، وتمّ تحقيق ذلك من خلال إدخال أوزان من طرق حصينة في التقدير واختبار حسانتها وملاءمتها للنموذج مسبقاً، ومن ثم إختيار الأوزان الناتجة من أعلى الطرق الحصينة كفاءة وإدخال هذه الأوزان في مراحل طرق إختيار أفضل معادلة إنحدار، فينتج عن ذلك نموذج يحقق صفتين في آن واحد وهما الحصانة وتقليل الأبعاد مقابل زيادة الكفاءة. وقد تمّ استخدام أسلوب المحاكاة على نماذج بأبعاد مختلفة وأحجام عينات مختلفة ونسب تلويث مختلفة في المتغير المعتمد مرة، وفي المتغيرات المستقلة مرة أخرى وفي الاثنین معاً مرة ثالثة، مع التركيز على دراسة احتمال تأثير وجود القيم الشاذة على المتغيرات التي ستبقى في الأنموذج والمتغيرات التي سيتم حذفها. ولتحقيق فكرة البحث تمّت مقارنة عدد من طرق التقدير الحصينة ومقارنة النتائج مع طريقة المربعات الصغرى الإعتيادية (OLS) وطريقة LASSO الحصينة المكيفة على بيانات تجريبية بأستخدام المحاكاة وكذلك على بيانات لعينة من مرضى التلاسيميا في محافظة نينوى.

معلومات النشر

تاريخ المقالة:
تاريخ الاستلام: 27 آذار 2023
تاريخ القبول: 25 ايار 2023
تم القبول في 28 ايار 2023
متاح على الانترنت 1 حزيران 2024
الكلمات الدالة:
المربعات الصغرى الموزونة، مقدرات M، اختيار المتغيرات، الكشف عن الشواذ، الانحدار الحصين.

المراسلة:

فاطمة محمد احمد
fatemamohammad738@gmail.com
[orcid](https://orcid.org/0000-0001-9148-4444)

DOI [10.33899/IQJOSS.2024.0183230](https://doi.org/10.33899/IQJOSS.2024.0183230) , ©Authors, 2024, 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/>).

1- المقدمة Introduction

إنّ تقليل عدد المتغيرات التوضيحية التي تُستخدَم في نماذج الإنحدار يعمل على التقليل من الجهد والوقت والتمن المستغرق في تقدير وتحليل وتفسير الأنموذج، ويركز على المتغيرات ذات التأثير المعنوي الحقيقي على متغير الاستجابة، وكذلك يضمن لنا سهولة التحليل والفهم لطبيعة العلاقة بين المتغير المعتمد وأهم المتغيرات المستقلة تأثيراً على المتغير التابع إذ إنّ أية زيادة في عدد المتغيرات التوضيحية غير المهمة في الأنموذج ستؤدي إلى ضياع الجهد دون جدوى، وعليه يجب أن يكون هناك توازن بين عملية تقليل عدد المتغيرات التوضيحية وبين زيادة عددها للحصول على نتائج تنبؤية دقيقة، ويكون من الأفضل اختيار المعادلة المُفسرة بأقل عدد من المتغيرات التوضيحية بحيث تكون هذه المتغيرات مُهمة ولها تأثير معنوي فعلي على المتغير التابع (Al-Subaihi, 2004).

وبرزت فكرة البحث التي تأخذ بنظر الاعتبار حصانة وكفاءة المقدرات ولهذا ستتصب آليات التحليل على تحسين طرق اختيار أفضل معادلة إنحدار بالطرق التقليدية ضد القيم الشاذة وذلك من خلال تطبيق مجموعة من دوال الوزن وطرق التقدير ومقدرات التباين والمقدرات الابتدائية البديلة ومن ثم أنتخاب أعلى المقدرات الموزونة كفاءة وتوظيف الاوزان المستحصلة منها في مراحل اختيار أفضل معادلة انحدار .

2- التقدير الحصين Robust Estimation

هي طرائق التقدير التي تعمل بشكل جيد ليس فقط تحت ظروف مثالية، ولكن أيضاً في ظلّ ظروف تمثّل خروجاً عن التوزيع أو الأ نموذج المفترض.

إنّ الهدف الأساسي من الإحصاء الحصين هو تطوير الإجراءات التي تبقى موثوقة وفعالة بشكل معقول في ظلّ الإنحرافات الصغيرة عن الأ نموذج. أي عندما يقع التوزيع الأساسي بالقرب من الأ نموذج المفترض.

الطرائق الإحصائية الحصينة هي امتداد للطرائق المعلمية ، مع الأخذ بنظر الإعتبار أنّ النماذج المعلمية هي أفضل تقريب للواقع ولكنها تفقد كفاءتها عند عدم تحقق الفروض التي تقوم عليها (Hurn and Mirosevich, 2008).

3- بعض طرائق التقدير الحصين Robust Estimation Methods

إنّ إتباع الطرائق التقليدية لتقدير معالم الأ نموذج تكون غير دقيقة في تحليل البيانات عند وجود القيم الشاذة أو وجود خلل في إحدى فرضيات الإنحدار أو أنّ توزيع الخطأ يكون غير طبيعي، حيث إنّ وجود قيمة شاذة واحدة في البيانات سوف يؤدي إلى خلل في خصائص مقدرات المربعات الصغرى وإنّ المقدر الحصين هو الذي يحافظ على الخصائص المرغوب بها للمقدرات عند خرق بعض فروض الإنحدار وسنلجأ إلى بعض طرائق التقدير الحصينة التي تمّ تطبيقها في هذه الرسالة:

1- مقدرات M-estimators

تعدّ طريقة التقدير M واحدة من أهم الطرائق الحصينة شائعة الإستخدام، إذ أشارت أغلب الدراسات إلى إنّ هذه الطريقة تعدّ من أكثر الطرائق الحصينة سواء في كفاءتها المكافئة لطريقة المربعات الصغرى عندما تتوزّع الأخطاء توزيعاً طبيعياً بوسط (صفر) وتباين σ^2 ، وتكون كفاءتها أعلى من كفاءة المربعات الصغرى عندما لا تتوزع الأخطاء توزيعاً طبيعياً أو عند وجود قيم شاذة في البيانات. وقد وسع (Huber, 1973) نتائجه للتقدير الحصين من معلمة الموقع إلى حالة الإنحدار الخطي. وقد اكتسبت هذه التقديرات شهرة أكثر من بقية المقدرات الحصينة الأخرى لأنّها أكثر مرونة، وكذلك توفر إمكانية تعميمها مباشرة إلى الإنحدار المتعدّد.

وتبدأ عملية التقدير وفق هذه الطريقة (M-estimation) إبتداءً باستخدام طريقة المربعات الصغرى الموزونة تكرارياً (Iteratively Re-weighted Least Squares IRWLS) مع استخدام إحدى دوال وزن مقدرات M، وهي كثيرة، ومن أشهرها هي دوال (Huber, Hampel, Tukey's Bisquare)، علماً أنّ دالة Huber تضمن الحصول على مقدرات وحيدة (حل وحيد Unique Solution)، أما الطريقتان الأخرتان (Tukey's Bisquare و Hampel) فينتج عنهما عدة نهايات صغرى (تعدد الحلول المثلى)، وعندها يصبح من الضروري استخدام قيم إبتدائية جيّدة بالشكل الذي نضمن فيه حصول تقارب سريع في المقدرات. ومن الناحية الرياضية تقوم فكرة طريقة التقدير (M) على تحويل الأسلوب المتبع في طريقة المربعات الصغرى التي تهدف إلى تصغير المقدار الآتي :

$$\text{Min} \sum_{i=1}^n e_i^2 \quad (1)$$

أي أنّ

$$\text{Min} \sum_{i=1}^n \left(y_i - \sum_{j=1}^p X_{ji} \hat{\beta}_j \right)^2$$

أما طريقة مقدرات (M) فتهدف إلى تصغير المقدار

$$\text{Min} \sum_{i=1}^n \rho(e_i)$$

$$\text{Min } \sum_{i=1}^n \rho \left(\underline{y}_i - \sum_{j=1}^p \underline{x}_{ji} \hat{\beta}_j \right) \quad (2)$$

إذ تمثّل ρ دالة بدلالة الأخطاء، ولتصغير المعادلة (1) نشقها جزئياً بالنسبة لمتجه المقدرات $\hat{\beta}$ ومساواتها بالصفر، وكما يلي:

$$\sum_{i=1}^n x_i \Psi \left(\underline{y} - \underline{X} \hat{\beta} \right) \quad (3)$$

إذ تمثّل Ψ المشتقة الجزئية للدالة ρ بالنسبة للمعاملات في المعادلة (2)، وتمثّل منظومة مكونة من P من المعادلات، وتحلّ باستخدام إحدى الطرق العددية المعروفة أو طريقة المربعات الصغرى الموزونة (Weighted Least Squares Method)، ولايجاد مقدرات M يتم استخدام الصيغة الآتية:

$$\hat{\beta}_M = \left(\underline{X}' \underline{W} \underline{X} \right)^{-1} \underline{X}' \underline{W} \underline{Y} \quad (4)$$

إذ تمثّل W مصفوفة الأوزان، وهي مصفوفة قطرية ($n \times n$) عناصر قطرها الرئيسية معطاة بالصيغة الآتية:

$$W_i = \frac{\Psi(e_i)}{(e_i)} \quad (5)$$

$$W_i = \frac{\left[\Psi \left(\underline{y}_i - \sum_{j=1}^p \underline{x}_{ji} \hat{\beta}_j \right) \right]}{\left(\underline{y}_i - \sum_{j=1}^p \underline{x}_{ji} \hat{\beta}_j \right)} \quad (6)$$

إذ تمثّل β^0 القيم الابتدائية لمتجه معاملات النموذج، ويتم استخدامها لتحديد الأوزان، ويمكن استخدام مقدرات المربعات الصغرى كقيم ابتدائية وقد تم في هذه البحث استخدام مقدرات المربعات الصغرى الاعتيادية OLS ومقدرات المربعات الصغرى المبتورة أو المشذبة LTS، ومن التكرار الأول نجد قيمة $\hat{\beta}^1$ ، أما في التكرار الثاني فنستخدم $\hat{\beta}^1$ في إيجاد الأوزان التي ستستخدم لإيجاد $\hat{\beta}^2$ وهكذا تستمر عملية التكرار حتى نحصل على مقياس التقارب (Convergence) المعروف بالصيغة الآتية :

$$\text{Max} \left[\left| \hat{\beta}_j^{(r)} - \hat{\beta}_j^{(r-1)} \right| \right] < \delta \quad (7)$$

إذ تمثّل δ قيمة صغيرة جداً، r تمثّل عدد مرات التكرار، أي أنّ الحل يتوقف عندما يصبح الفرق المطلق بين المعلمات المقدرّة في المرحلة الحالية والمعلمات المقدرّة في المرحلة السابقة أصغر من القيمة المختارة δ أو يساويها، ولجعل مقدرات (M) تمتلك خاصية ثبات التباين Invariant Scale، فإنّ الدالة المطلوب تصغيرها هي:

$$\text{Min } \sum_{i=1}^n \rho \left(\underline{y}_i - \sum_{j=1}^p \underline{x}_{ji} \hat{\beta}_j \right) / \hat{\sigma} \quad (8)$$

حيث نقوم بإشتقاق الدالة (8) بالنسبة للمتجه $\hat{\beta}_j$ ومساواتها بالصفر فنحصل على دالة الوزن، وكما يلي:

$$\sum_{i=1}^n X_{ij} \Psi \left(y_i - x_i \beta_j \right) / \hat{\sigma} = 0 \quad (9)$$

ويمكن حلّ المعادلة أعلاه باستخدام الصيغة (2) إذ إنّ الأوزان يتم إيجادها وفق الصيغة الآتية:-

$$W_i = \frac{[\Psi(e_i/\hat{\sigma})]}{(e_i/\hat{\sigma})} = \frac{\left[\Psi \left(\frac{\underline{y}_i - \sum_{j=1}^p \underline{x}_{ji} \hat{\beta}_j}{\hat{\sigma}} \right) \right]}{\left(\frac{\underline{y}_i - \sum_{j=1}^p \underline{x}_{ji} \hat{\beta}_j}{\hat{\sigma}} \right)} \quad (10)$$

لإيجاد $(\hat{\sigma})$ في المعادلة أعلاه، والتي تمثل قيمة الانحراف المعياري، وأن هذه القيمة تقدر مرة واحدة فقط باستخدام القيم الأولية قبل البدء بالتكرار وهناك عدة صيغ لتقديرها منها:

$$\hat{\sigma} = 1.5 \text{ Med}|e_i| = \text{Med}|e_i|/0.6745 \quad (11)$$

$$\hat{\sigma} = 1.2 \text{ Med}|e_i| \quad (12)$$

$$\hat{\sigma} = 1.4825[\text{Med}|e_i - \text{Med}e_i|] = \text{Med}|e_i - \text{Med}(e_i)|/0.6745 \quad (13)$$

إذ تمثل e_i البواقي و Med يشير الى الوسيط، ولقد اقترح الباحثون عدداً من الدوال $\rho(\cdot)$ أو مشتقاتها $\Psi(\cdot)$ ، بحيث تجعل نتائج التقدير حسينية لا تتأثر بوجود الشواذ، وفيما يلي بعض الدوال المهمة لهذا النوع من المقدرات والمعرفة بدلالة الدالة $\Psi(e_i)$ ، وبافتراض أن وسيط الأخطاء المطلقة (MAD: Median Absolute Deviation) الوارد في الصيغة (13) أعلاه كمقدر للانحراف المعياري، وكما يلي:

$$\text{MAD} = \text{Med}|e_i - \text{Med}(e_i)|/0.6745$$

اذ أن 0.6745 تمثل وسيط التوزيع الطبيعي القياسي.

وقد أوجد الباحثون (Montgomery et. al., 2001) الصيغة القياسية للباقي Standardized Residuals باستخدام $e_{is} = \frac{e_i}{\text{MAD}}$ ، أن ثابت القطع (C: Tuning Constant) الذي يجعل التباين المقدر MAD غير متحيز تقريباً ل σ عندما يكون حجم العينة كبير والخطأ يتوزع طبيعياً (Hasan and Ridha, 2011).

بعض دوال الأوزان لمقدر M

A. دالة Tukey Bisquare (Beatone and Tukey,1974):

$$\Psi_{\text{Bisquare}}(e_{is}, c) = \begin{cases} \left[1 - \left(\frac{e_{is}}{c}\right)^2\right]^2 & \text{if } |e_{is}| \leq c \\ 0 & \text{otherwise} \end{cases} \quad (14)$$

اذ إن c تأخذ القيمة الافتراضية $c = 4.685$ ، دالة Huber (Huber,1964):

$$\Psi_{\text{Huber}}(e_{is}, c) = \begin{cases} 1 & \text{if } |e_{is}| \leq c \\ \frac{c}{|e_i|} & \text{otherwise} \end{cases} \quad (15)$$

اذ إن c تأخذ القيمة الافتراضية $c = 1.345$ ، دالة Hampel (Grubbs,1969):

$$\Psi_{\text{Hampel}}(e_{is}, c) = \begin{cases} 1 & \text{if } |e_{is}| \leq c \\ \frac{a}{|e_i|} & \text{if } a < |e_{is}| \leq b \\ \frac{a(c-|e_{is}|)}{|e_{is}|(c-b)} & b < |e_{is}| \leq c \end{cases} \quad (16)$$

اذ إن القيم الافتراضية لثوابت القطع (Tuning Constants) هي $a = 2$ $b = 4$ and $c = 8$

علماً أن ثابت التوليف (Tuning Constant) لكل دالة يستخدم لتعديل كفاءة المقدرات الناتجة لتوزيعات محدّدة ويحقق كفاءة تقريبية مقدارها (95 %) عندما تتبع الاخطاء التوزيع الطبيعي، وأن الإختيار الجيد لقيمة هذا الثابت يؤدي إلى زيادة حصانة المقدرات، لأن لهذا الثابت تأثيراً كبيراً على حصانة المقدرات، وإن قيمته تتراوح بين انحراف معياري واحد إلى انحرافين معياريين لقيم المشاهدات أو الاخطاء.

2- طريقة S الحصينة:

اقترح مقدر للانحراف المعياري للأخطاء بالاعتماد على فكرة (Rousseeuw and Yohai, 1984) و (Rousseeuw and Leroy, 1987) الذين اقترحوا مقدرات S باعتبارها الحل الذي يوجد أقل قيمة تشتت ممكن للأخطاء، أي أنه يقوم بإيجاد

$$\text{Min}_{\hat{\beta}} S(e_1(\hat{\beta}), e_2(\hat{\beta}), \dots, e_n(\hat{\beta})) \quad (17)$$

المربعات الصغرى التي تقوم بتصغير تباين الأخطاء، ويجب أن يكون واضحاً بأن مقدرات المربعات الصغرى الإعتيادية OLS يمكن النظر إليها كحالة خاصة أقل حسانة من مقدرات S، فعندما نبدأ بتصغير تباين الأخطاء فإن مقدرات S الحصينة ستقوم بتصغير مقدر التباين الحصين للأخطاء أي أنه يقوم بتصغير قيمة المعادلة التالية:

$$\frac{1}{n} \sum_{i=1}^n \rho\left(\frac{e_i}{\hat{\sigma}_e}\right) = b \quad (18)$$

إذ إن b قيمة ثابت تعرف على أنها $b = E_{\emptyset}[\rho(e)]$ ، حيث \emptyset تمثل التوزيع الطبيعي القياسي، وبعد أخذ الإشتقاق لهذه المعادلة وحلها نحصل على

$$\frac{1}{n} \sum_{i=1}^n \Psi\left(\frac{e_i}{\hat{\sigma}_e}\right) = b \quad (19)$$

حيث إن Ψ يتم استبدالها بدالة وزن مناسبة. كما هو الحال بالنسبة لأكثر مقدرات M، ومن أمثلة دوال الوزن والتي Huber أو Hampel أو biweight وغيرهم. وعلى الرغم من أن مقدرات S لها نقطة انهيار تبلغ 0.5 فإن كلفة ذلك هو أن هذه المقدرات تمتلك كفاءة منخفضة جداً حيث تبلغ كفاءتها مقارنة بطريقة المربعات الصغرى 30% فقط حسب ما بين كل من (Croux et al., 1994).

2- طريقة MM الحصينة (Rousseeuw and Leroy, 1987):

تقوم هذه الطريقة بتقدير معاملات الانحدار باستعمال تقدير (S الحصين)، وذلك بتصغير minimize مقياس الانحراف المعياري للبقايا من طريقة M. إن طريقة (MM) تهدف للحصول على مقدرات ذات قيم عالية الدقة أو أكثر كفاءة. قبل استعمال المشاهدات في الأنموذج وهي:

$$\sum_{i=1}^n \rho(e_{is})x_{ij} = 0 \quad \text{or} \quad (20)$$

$$\sum_{i=1}^n \rho\left(\left(\frac{y_i - \sum_{j=1}^k x_{ij}\hat{\beta}_j}{S_{MM}}\right)\right)x_{ij} = 0$$

إذ إن S_{MM} هو الانحراف المعياري الذي يتم الحصول عليه من بواقي تقدير S و ρ هي دالة Tukeys الموزونة

$$\rho(e_{is}) = \begin{cases} \frac{e_{is}^2}{2} - \frac{e_{is}^4}{2c^2} + \frac{e_{is}^6}{6c^2} & -c \leq e_{is} \leq c \\ \frac{c^2}{6} & e_{is} < -c \text{ or } e_{is} > c \end{cases} \quad (21)$$

3- طريقة المربعات الصغرى المشدبة (LTS): Least Trimmed Squares ethod:

إن العديد من الباحثين لم يدرك أن أداء المربعات الصغرى الإعتيادية OLS (Ordinary Least Squares) يمكن أن يكون ضعيفاً جداً عندما يكون شكل التوزيع الطبيعي للبيانات ذات ذيول (أطراف) ثقيلة (Heavy tails)، والتي تنشأ من القيم الشاذة حتى إذا كانت هناك قيمة شاذة واحدة فقط فانه سوف يكون لها تأثير كبير على تقديرات OLS. وللتغلب على هذه المشكلة سوف يتم استخدام مقدر حصين له نقطة انهيار عالية كبديل عن طريقة OLS. ومن البدائل المتوفرة مقدر LTS الذي اقترح من قبل (Rousseeuw and Yohai, 1984) الذي يكون لديه نقطة انهيار (Breakdown point) عالية تساوي $BP = \lceil \{n - k/2 + 1\} / n \rceil$ (Uraibi et al., 2009).

نحصل عليه بالشكل الآتي:

$$\text{Min} \sum_{i=1}^h e_i^2$$

إذ إنَّ e_i^2 تمثِّل مربعات حدِّ الخطأ. وهذه الطريقة تقلل من نسبة تأثير القيم الشاذة (α) في البيانات. لتثذيب (لقطع) نسبة α اقترح (Rousseeuw and Leroy, 1987) اختيار h حسب الصيغة الآتية:-

$$h = n/2 + (p + 1)/2 \quad (22)$$

إذ إنَّ p تمثِّل عدد المعلمات. ومن مميزات استخدام طريقة LTS السيطرة على مستوى التثذيب الذي يعتمد على تشابه النسبة المنوية المبتورة لنسبة القيم الشاذة. إذا كان هناك شك في إنَّ البيانات تحتوي على ما يقرب 10% من القيم الشاذة، فإنَّ LTS سوف تقطع 10% من الطرفين أي أنَّ عملية التقدير تتمُّ بما نسبته 80% من المشاهدات الأصلية (Arif, 2012).

4-طريقة LASSO المكيفة (Wang et al. 2013) (Zou. 2006):

$$\arg \min_{\hat{\beta}} \sum_{i=1}^n (1 - \exp \{-(y_i - X_i^T \beta)^2 / \gamma_n\}) + n \sum_{j=1}^d \lambda_{nj} \beta_j \quad (23)$$

ذ أن : γ_n : معلمة الضبط (التوليف) معلمة التنظيم λ_{nj} :

التطبيق العملي

1. جانب المحاكاة (الخوارزمية المقترحة للجانب التجريبي) Simulation Side

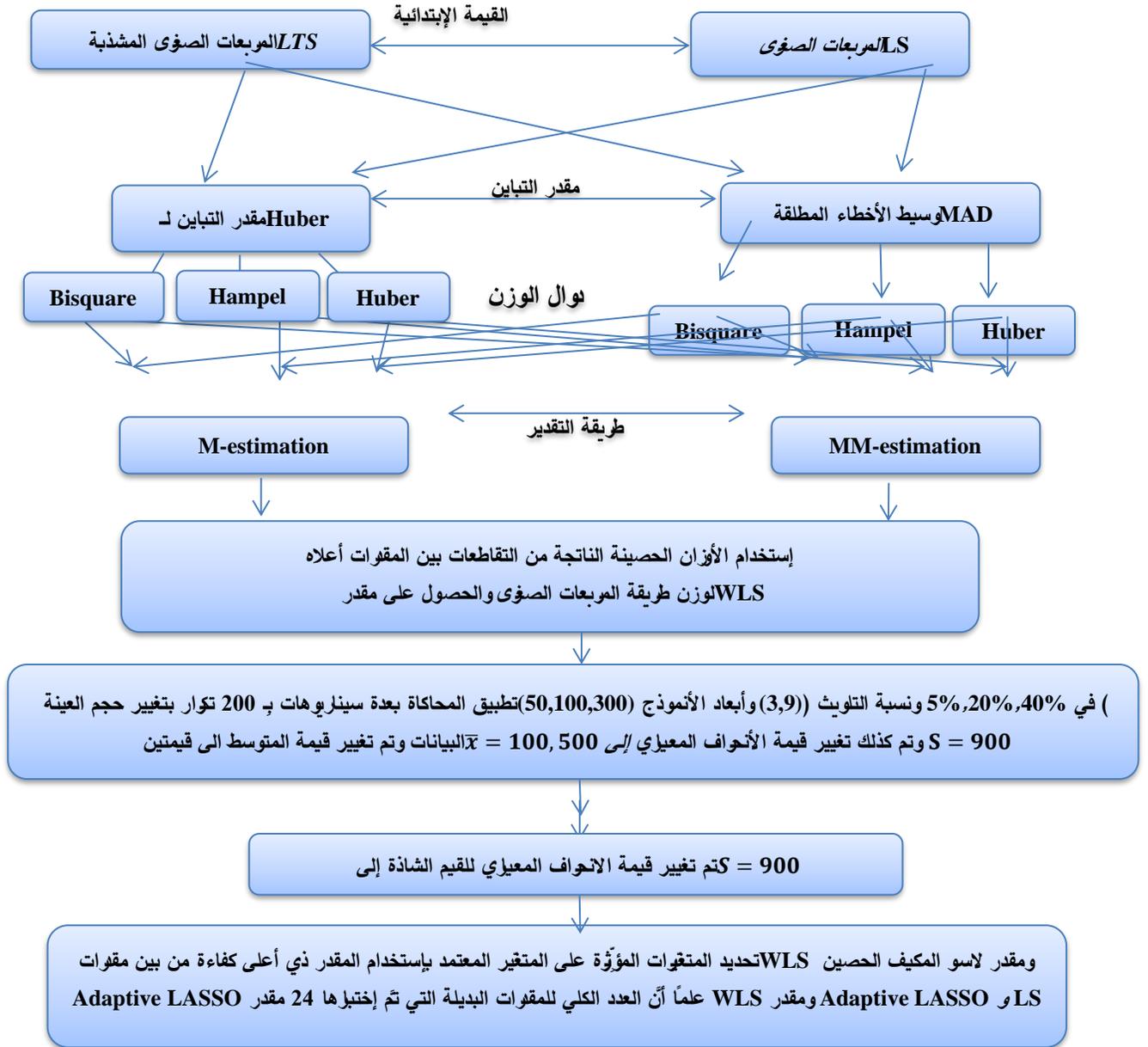
تقوم فكرة الخوارزمية التي تسعى البحث إلى تقديمها إلى تقدير المعلمات بأسلوب M وبدوال أوزان مختلفة من خلال إدخال كلِّ المتغيرات في الأنموذج وتطبيق عدة حالات تتمثِّل بتغيير المقدرات الابتدائية ومقدرات التباين والأوزان المستخدمة. وقد تمَّ اللجوء إلى إثنين من أهم المقدرات الحصينة، وهما مقدر M-estimation الحصين ضد القيم الشاذة في المتغير المعتمد (y-outliers)، وكذلك مقدر MM-estimation الذي يمتلك حصانة ضد القيم الشاذة في كلِّ من المتغيرين المعتمد، وكذلك المتغير المستقل (تدعى القيم الشاذة في فضاء المتغيرات المستقلة بالقيم الجاذبة أو الفعالة أو المخلة (X-leverage points)، حيث أنَّ مقدر MM-estimation هو عبارة عن مقدر M-estimation بدالة وزن Tukey's biweight معتمدة على أحد مقدرات S (S-estimator) وفي أدناه مخطط إنسيابي للخوارزمية المقترحة.

وتم استخدام ثلاثة أحجام عينات $n = 50, 100, 300$ ، وتم اختبار بُعدين لأنموذج الانحدار كما هو واضح في الخوارزمية أعلاه وهما $m = 3, 9$ وتم استخدام ثلاثة نسب للتلوين $\alpha = 5\%, 20\%, 40\%$ وتمت عملية التلوين باستخدام طريقة تلوين صف كامل وهو الأسلوب الذي إتبعه (Toka et.al., 2021) والذي يطلق عليه casewise وتم تغيير الوسط الحسابي للقيم الملوثة مرتين، حيث أن بيانات المحاكاة قد تم توليدها بوسط حسابي = صفر وعند تلوين البيانات تم تبديل الوسط الحسابي بقيمتين 100, 500 وتم تغيير قيمة الانحراف المعياري من $s = 100$ إلى $s = 900$ وبذلك أصبحت القيم الملوثة مختلفة من حيث قيمة المتوسط والتباين وكل الملوّثات قد تم توليدها من التوزيع الطبيعي وتم تكرار توليد $B = 200$ عينة لأختبار كل طريقة من طرق التقدير المقترحة وكانت سيناريوهات التلوين تتضمن تلوين y مره والمتغيرات المستقلة X'S مرة،

وكليهما معاً مرة ثالثة وبلغ إجمالي نماذج المحاكاة التي تم توليدها في دراسة المحاكاة هو 561,600، حيث تم توليدها وفقاً للتالي : $561,600 = 26 \times 3 \times 200 \times 2 \times 3 \times 2 \times 3$ أنموذج رياضي مقدر (علماً أن كل مقدر حصين يحتاج إلى ما أقصاه 20 تكراراً ليصل الى التقارب، أي أن عدد النماذج المقدره فعلياً قد يصل إلى $561,600 \times 20 = 11,232,000$ أنموذجاً مقدرًا).

المخطط الانسيابي الوارد في البحث عبر عن الخوارزمية المقترحة حيث اعتمدت الباحثة على مفهوم الخوارزمية الذي يمثل مجموعة من الخطوات الرياضية والمنطقية والمتسلسلة اللازمة لحل مشكلة البحث لتحسين عملية اختيار المتغيرات في معادلة الانحدار وتم رسم المخطط الانسيابي Flow Chart للتعبير عن خطوات الخوارزمية.

المخطط الانسيابي للطريقة المقترحة لإختيار أفضل مقدر لاختيار أفضل معادلة إنحدار



لغرض اختيار النموذج الأفضل تم تنفيذ على بيانات مولدة باستخدام برنامج R حيث تم اختيار الطريقة التي أثبتت أعلى كفاءة (أقل RMSE) في تجارب المحاكاة وتمت مقارنتها مع طريقة Adaptive LASSO الحصينة مقارنة بالطريقة الأكثر شيوعاً، وقد ثبت بأنها قد تغلبت على طريقة LASSO بعد اختيار أعلى دوال المربعات الصغرى الموزونة كفاءةً.

جدول (1) : الطرق التي حققت أفضل أداء (أقل RMSE) في حالة وجود قيم شاذة في المتغير المعتمد

لنموذج بثلاثة متغيرات						
أكفاً طريقة	قيمة المتوسط	حجم العينة	محور الشواذ	نسبة التلويث		
WLS.LS.Bisquare.Mad.M	عند قيمة وسط حسابي = 100	50	Y-outlier	6%		
WLS.LTS.HAMPEL.Huber.M				20%		
				40%		
WLS.LS.Bisquare.Mad.M		100		5%		
		20%				
		40%				
WLS.LS.Hampel.Huber.M		300		5%		
WLS.LS.Hampel.Mad.M				20%		
WLS.LS.Bisquare.Mad.M				40%		
لم تتفوق أي من الطرق على نظيراتها		عند قيمة وسط حسابي = 500		50	Y-outlier	6%
WLS.LS.Hampel.Mad.M						20%
WLS.LTS.HAMPEL.Huber.M						40%
WLS.LS.Hampel.Huber.M	100		5%			
WLS.LS.Bisquar.Mad.M			20%			
WLS.LS.Bisquar.Mad.M			40%			
WLS.LS.Hampel.Huber.M	300		5%			
WLS.LS.Bisquar.Mad.M			20%			
WLS.LS.Hampel.Huber.M			40%			
WLS.LS.Hampel.Mad.M	عند قيمة وسط حسابي = 100		50	Y-outlier		6%
WLS.LTS.HAMPEL.Huber.M						20%
						40%
WLS.LS.HAMPEL.Huber.M		100	5%			
		20%				
		40%				
WLS.LS.HAMPEL.Huber.M		300	5%			
WLS.LS.Bisquar.Mad.M			20%			
			40%			
WLS.LS.HAMPEL.Huber.M		عند قيمة وسط حسابي = 500	50		Y-outlier	6%
WLS.LS.Bisquar.Mad.M						20%
						40%
WLS.LS.HAMPEL.Huber.M	100		5%			
WLS.LTS.HAMPEL.Huber.M			20%			
WLS.LS.Bisquar.Mad.M			40%			
WLS.LS.HAMPEL.Huber.M	300		5%			
WLS.LS.Bisquar.Mad.M			20%			
			40%			

جدول (2): الطرق التي حققت أفضل أداء (أقل RMSE) في حالة وجود قيم إنعطاف (شواذ) في المتغيرات المستقلة

لنموذج بثلاثة متغيرات							
نسبة التلويث	محور الشواذ	حجم العينة	قيمة المتوسط	أكفاً طريقة			
X-leverage		50	عند قيمة وسط حسابي = 100	WLS.LTS.HAMPEL.Huber.M	6%		
					20%		
					40%		
					100	5%	
						20%	
						40%	
		300	5%				
			20%				
			40%				
		WLS.LS.Hampel.Mad.M LS					6%
							20%
							40%
X-leverage		50	عند قيمة وسط حسابي = 500	WLS.LTS.HAMPEL.Huber.M	6%		
					20%		
					40%		
					100	5%	
						20%	
						40%	
		300	5%				
			20%				
			40%				
		LS					6%
							20%
							40%
لنموذج بتسعة متغيرات							
X-leverage		50	عند قيمة وسط حسابي = 100	WLS.LTS.HAMPEL.Huber.M	6%		
					20%		
					40%		
					100	5%	
						20%	
						40%	
		300	5%				
			20%				
			40%				
		X-leverage		50	عند قيمة وسط حسابي = 500	WLS.LTS.HAMPEL.Huber.M	6%
							20%
							40%
100	5%						
	20%						
	40%						
300	5%						
	20%						
	40%						

جدول (3) : الطرق التي حققت أفضل أداء (أقل RMSE) في حالة وجود قيم شاذة في المتغيرين X و Y

لنموذج بثلاثة متغيرات						
WLS.LTS.HAMPEL.Huber.M	عند قيمة وسط حساسي = 100	50	X-leverage Y-outlier	6%		
				20%		
40%						
WLS.LS.Bisquar.Mad.M		100		50	5%	
					20%	
40%						
WLS.LTS.HAMPEL.Huber.M		300		50	5%	
					20%	
WLS.LS.Bisquar.Mad.M				40%		
WLS.LS.Hampel.Huber.M		عند قيمة وسط حساسي = 500		50	X-leverage Y-outlier	6%
						20%
40%						
WLS.LTS.HAMPEL.Huber.M	100		50	5%		
				20%		
40%						
WLS.LS.Bisquar.Mad.M	300		50	5%		
				20%		
40%						
WLS.LTS.HAMPEL.Huber.M	عند قيمة وسط حساسي = 500		50	X-leverage Y-outlier		6%
						20%
40%						
WLS.LS.Bisquar.Mad.M		100	50		5%	
					20%	
40%						
WLS.LTS.HAMPEL.Huber.M		300	50		5%	
					20%	
40%						
LS						40%
لنموذج بتسعة متغيرات						
WLS.LTS.HAMPEL.Huber.M		عند قيمة وسط حساسي = 100	50		X-leverage Y-outlier	6%
	20%					
40%						
WLS.LS.Hampel.Mad.M	100		50	5%		
				20%		
40%						
LS	300		50	5%		
				20%		
40%						
WLS.LTS.HAMPEL.Huber.M	عند قيمة وسط حساسي = 500		50	X-leverage Y-outlier		6%
						20%
40%						
WLS.LS.Bisquar.Mad.M		100	50		5%	
					20%	
40%						
WLS.LTS.HAMPEL.Huber.M		300	50		5%	
					20%	
40%						
WLS.LS.Bisquar.Mad.M		عند قيمة وسط حساسي = 500	50		X-leverage Y-outlier	6%
						20%
40%						
WLS.LTS.HAMPEL.Huber.M	100		50	5%		
				20%		
40%						
WLS.LS.Bisquar.Mad.M	300		50	5%		
				20%		
40%						
LS						40%

وعند النظر إلى الجداول الثلاثة (1) و (2) و (3) أعلاه الذين يمثلون ملخص نهائي لكافة تجارب المحاكاة لكل احجام النماذج والعينات ونسب التلويث ودوال الوزن وقيمتي المتوسط، حيث تظهر الطرق الأعلى كفاءة لكل حالة وتعطي الباحثين خارطة طريق حول طريقة

التقدير المناسبة لكل حالة. وإذا ودنا إستخلاص طريقة متفوقة إجمالاً لكل حالة من حالات توليخ البيانات بالقيم الشاذة وكما هو مبين فيما يلي:

A. كانت طريقة **LS.Bisquare.Mad.M** أكثر الطرق التي حققت أعلى كفاءة (اقلهم بقيمة RMSE) من بين كل الطرق التي

تمت مقارنتها في دراسة المحاكاة عند وجود شواذ في المتغير المعتمد y-outliers في اغلب الحالات التي تمت مقارنتها

B. كانت طريقة **LTS.HAMPEL.Huber.M** أكثر الطرق التي حققت أعلى كفاءة (اقلهم بقيمة RMSE) من بين كل الطرق

التي تمت مقارنتها في دراسة المحاكاة عند وجود شواذ في المتغيرات المستقلة X-leverage points في اغلب الحالات التي

تمت مقارنتها

C. كانت طريقة **LTS.HAMPEL.Huber.M** أكثر الطرق التي حققت أعلى كفاءة (اقلهم بقيمة RMSE) من بين كل الطرق

التي تمت مقارنتها في دراسة المحاكاة عند وجود شواذ في المتغير المعتمد y-outliers وفي المتغيرات المستقلة X-leverage

points في اغلب الحالات التي تمت مقارنتها

والنقاط الواردة في أعلاه تعطي خارطة طريق موجزة يمكن للباحثين إستخدامها لاختيار طريقة التقدير الملائمة في حالة عدم توفر معلومات واضحة حول نسب الشواذ في البيانات ونوعها.

2. الجانب التطبيقي Applied Side

جمع البيانات (Al-Nuaimi, 2005) :

طبقت هذه الدراسة على بيانات جمعت من مستشفى ابن الأثير التعليمي للولادة والأطفال في الموصل ، من مرضى مصابين بفقر دم البحر الأبيض المتوسط من نوع بيتا او ما يسمى الثلاثيميا ، وكان عدد المشاهدات (150) حيث أن المتغير المعتمد أو المستجيب في الدراسة مثل عمر العظم مقياساً بالشهر (Month) Age of the Bone كما واختير عدد من المتغيرات التي يعتقد بأنها تؤثر فيه بعد مراجعة أطباء اختصاصيين وممارسين في مرض الثلاثيميا.

Variable Name	أسم المتغير	
Real Age (Month)	العمر الحقيقي (مقاساً بالشهر)	X1
Onset of Disease (Month)	عمر المريض عند ظهور المرض (مقاساً بالشهر)	2X
c.m)(Enlargement of liver	تضخم الكبد (مقاساً بالسنتيمتر)	3X
Hemoglobin	هيموكلوبين الدم	4X
Packed cell volume	مكداس الدم (خلايا الدم المضغوط)	5X
Reticulocyte	الخلايا الشبكية	6X
Normooblast	الأرومة الحمراء	7X
Fetal Hemoglobin	الهيموكلوبين الجنيني	8X
Number of Blood units	عدد وحدات الدم	9X
Onset of Blood Trans fusion To According	بداية نقل الدم حسب العمر (مقاساً بالشهر)	10X

جدول (4): كفاءة النماذج النهائية لكل الطرق المقارنة

الطريقة	X.Interc ept.	X ₁	X ₂	X ₃	X ₄	X ₅	X ₆	X ₇	X ₈	X ₉	X ₁₀	R ²	R ² _{adj}	AIC	p - value
LS(كل الانحدارات الممكنة بدون أوزان)	1	1	0	0	1	0	0	0	0	1	0	0.81 7	0.81 4	1274.4 78	2.2e- 16
LS(الحذف العكسي بدون أوزان) LS (الاختيار الامامي بدون أوزان) LS (الانحدار التدريجي بدون أوزان)	1	1	1	0	1	1	0	1	0	1	0	0.83 6	0.82 91	836.72	2.2e- 16
WLS.LS.Bisquare. Mad.M (كل الانحدارات الممكنة)	1	1	0	0	1	0	0	0	0	1	0	0.88 2	0.88	1205.1 05	2.2e- 16
WLS.LS.Bisquare. Mad.M (طريقة الحذف العكسي) WLS.LS.Bisquare. Mad.M (طريقة الإختيار الامامي) WLS.LS.Bisquare. Mad.M (الانحدار التدريجي)	1	1	1	1	1	1	0	1	0	1	0	0.90 36	0.89 89	729.29	2.2e- 16
Robust Variable Selection Exponential Squared Loss	1	1	0	0	1	0	0	0	0	0	0	0.74 15	0.72 29	- 181.90 97	Not Availa ble

وبالنظر إلى الجدول (4) أعلاه نلاحظ أن الطرق المختلفة قد أنتجت نماذجاً مختلفة، فمنها ما قام بإستبعاد أغلب المتغيرات مثل طريقتي كل الانحدارات الممكنة وطريقة LASSO اللتان أبقينا 3 و 2 متغيرات في النموذج على التوالي وهذا خلاف ماتمت مناقشته مع أطباء الاختصاص الذين أكدوا بأن المتغيرات العشرة الموجودة في النموذج مهمة وتؤثر على الإصابة بمرض التلاسيميا وعليه فإن طرق الاختيار الامامي والحذف العكسي والانحدار التدريجي التي تم وزنها بالأوزان الحصينة والتي حققت أعلى كفاءة في دراسات المحاكاة في غالبية الحالات وهي طريقة WLS.LS.Bisquare.Mad.M التي جمعت بين الحصانة والكفاءة فضلاً عن إتفاقها مع الرؤية الطبية لهذا المرض.

ومن ناحية أخرى فإن المقارنة بين مقاييس الكفاءة بين نماذج تختلف في عدد متغيراتها غير مجدية لأن تلك المقاييس ترتبط بعدد المتغيرات الموجودة في النموذج، فعلى سبيل المثال فإن قيمة معامل التحديد تزداد بزيادة عدد المتغيرات وذلك لأن كل متغير مضاف إلى النموذج يمتلك مجموع مربعات موجب وبالتالي فإن إضافته ستؤدي إلى إضافة قيمة موجبة إلى البسط، وبما أن مقام حساب معامل التحديد يحتوي على مجموع المربعات الكلي وبالتالي فهو قيمة ثابتة وزيادة البسط ستؤدي إلى زيادة القيمة الكلية للكسر وعليه تميل الباحثة إلى اعتماد نتائج المحاكاة التي تم فيها تنفيذ عملية التقدير مئات الالاف من المرات واعتماد النموذج النهائي الذي نتج من طريقة WLS.LS.Bisquare.Mad.M كنموذج نهائي يتم إعتماده لتشخيص أهم المتغيرات المؤثرة على الإصابة بالتلاسيميا. ونفس الامر ينطبق على معيار AIC والذي ترتبط عملية حسابه أيضاً بعدد المتغيرات في النموذج كون أن زيادة عدد المتغيرات سيطرر رقماً موجباً من البسط $(sse \times (n - p - 1))$ وبذلك فإن طرح رقم صغير أو كبير من n سيكون له تأثير كبير على قيمة لوغارتم المقدار فنتج لدينا أحياناً قيم سلبية تجعل قيمة AIC سالبة. أما لو كانت النماذج التي تتم مقارنتها من نفس البعد (نفس عدد المتغيرات) لأمكننا مقارنتها باستخدام تلك المعايير ولكن هذا لا يتحقق في طرق إختيار أفضل معادلة إندار التي تعمل على الموازنة بين تقليل عدد المتغيرات في النموذج وبين زيادة كفاءة النموذج المقدر.

الاستنتاجات

- بناءً على ماتم التوصل إليه في الجانبين التجريبي والتطبيقي وصل البحث إلى مجموعة من الاستنتاجات، وكما يلي:
1. حافظت طريقة المربعات الصغرى الموزونة على خصائصها وإمكانية تطبيقها في كل الحالات وللنماذج التي تعاني من كل أنواع الشواذ سواء أكانت في المتغير المعتمد أو المتغيرات المستقلة أو كليهما معاً.
 2. لم تكن طريقة LASSO الحصينة موقفة في كل الحالات، فضلاً عن أنها تستبعد الكثير من المتغيرات التي قد يكون بعضها مهماً، بالإضافة إلى أنها تعاني من بعض المشاكل الحسابية المتعلقة بالتقارب حيث توقفت الدالة عن التنفيذ في بعض الحالات.
 3. تمكن البحث ومن بين 561,600 تجربة محاكاة من تثبيت بعض المقدرات الأعلى كفاءة والتي ينصح الباحثين باستخدامها حسب موقع الشواذ وفي حالتها النماذج بالابعد القليلة أو العالية لحصول مقدراتها على أقل قيمة لجذر متوسط مربعات الخطأ، وكما يلي:
- A. في حالة وجود الشواذ في المتغيرات المستقلة (القيم الفعالة Leverage Points) يفضل عند إيجاد مقدر M الحصين استخدام مقدر المربعات الصغرى المشدبة (LTS) كقيمة ابتدائية واستخدام دالة Hampel كدالة وزن، ومقدر Huber كمقدر للتشتت، وطريقة M في التقدير كما تم تطبيقه في طريقة LTS.HAMPEL.Huber.M.
- B. في حالة وجود الشواذ في المتغير المعتمد (Y-outliers) يفضل عند إيجاد مقدر M الحصين استخدام مقدر المربعات الصغرى (LS) كقيمة ابتدائية واستخدام دالة Bisquare كدالة وزن، ومقدر وسط الأخطاء المطلقة كمقدر للتشتت، وطريقة M في التقدير كما هو الحال في مقدر LS.Bisquare.Mad.M الذي تم استخدامه في البحث.
- C. في حالة وجود الشواذ في المتغير المعتمد والمتغيرات المستقلة (القيم الفعالة Leverage Points والقيم الشاذة y-outliers) يفضل عند إيجاد مقدر M الحصين استخدام مقدر المربعات الصغرى المشدبة (LTS) كقيمة ابتدائية واستخدام دالة Hampel كدالة وزن، ومقدر Huber كمقدر للتشتت، وطريقة M في عملية التقدير كما تم بيانه في المقدر LTS.HAMPEL.Huber.M.
- بناءً على ماورد في أعلاه نستنتج أن الكفاءة في حساب الأوزان يمكن طريقة المربعات الصغرى (الموزونة WLS) من إنتاج مقدرات كفوءة كما تم إثباته وبالتالي يجعل استخدام الطرق التقليدية في إختيار أفضل معادلة إنحدار ممكناً في تحديد أهم المتغيرات المؤثرة على المتغير المعتمد وكما تم إثباته في الجانبين التجريبي والتطبيقي حيث تم إختزال النماذج بشكل جمع بين الكفاءة والحصانة.

References

1. Al-Subaihi, Ali, A. (2004). " Variable Selection in Multivariable's Regression Using SAS/IML ", Institute of Public Admin station, Saudi Arabia.
2. Hurn , D. and Mirosevich , V. , (2008) . "Modern Robust Statistical Methods" , The American Psychologist Association , Vol. 63, No. 7, PP. 591–601.
3. Huber,P.J.(1973). "Robust estimation of location parameter" , Annals of Mathematical Statistics, 35: 73-101.
4. Montgomery D. C., Peck E. A., and Vining G. C. (2001). " Introduction to Linear Regression Analyses", 3rd ed., New york: John Wiley&Sons Inc, Library of congress, U.S.A.
5. Hasan, T. A. and Ridha, M. S. (2011). "Using Robust Regression to Find the Most Appropriate Model to Represent Meteorological Data in the City of Erbil during (1998-2010)", Journal of Administration and Economics, thirty four years, 39 year, Issue: 39.
6. Beaton, A.E., and Tukey, J.W., (1974). "The Fitting of Power series Meaning Polynomials, Illustrated on band – Spectroscopic Data", Techno metrics, 16, PP 147 – 185.
7. Huber,P.J. (1964). "Robust Estimation of location parameter".Ann .math .stat.35, PP 73-101.

8. Grubbs, F.E. (1969). "Procedures for Detecting Outlying observation", *Technometrics*.
9. Croux, C. Rousseeuw, P. J. and Hossjer, O. (1994). "Generalized S-Estimators" *Journal of the American Statistical Association*, 89(428), pp.1271-1281.
10. Rousseeuw, P. J. and Leroy, A. M. (1987). *Robust Regression and Outlier Detection*, Wiley-Interscience, New York (series in Applied probability and statistics), 329 Pages. ISBN-471-852333.
11. Rousseeuw P.J. and Yohai V.J. (1984). "Robust regression by means of S-estimators. in *Robust and Nonlinear Time Series Analysis*", eds. J. Franke. W. Härdle and R.D. Martin. *Lecture Notes in Statistics* 26, New York: Springer Verlag. <ftp://ftp.win.ua.ac.be/pub/preprints/84/Robreg84.pdf>
12. Uraibi, H. S., Talib, B. A., and Yousif, J. H., (2009). "Linear Regression Model Selection Based on Robust Bootstrapping Technique", *American Journal of Applied Sciences*, 6 (6), PP.1191-1198.
13. Arif, A. L. (2012). "Using Types of Robust Regression to Deal with Outlying Values in Simple Linear Regression", *Dhi Qar Journal of Agricultural Research, College of Science - Dhi Qar University, Iraq*, 1 (2), PP 93-103.
14. Wang, X., Yunlu J., Mian H., and Heping Z. (2013). "Robust Variable Selection with Exponential Squared Loss." *Journal of the American Statistical Association* 108(502):632–3. <https://doi.org/10.1080/01621459.2013.766613>.
15. Zou, H. (2006). "The Adaptive Lasso and Its Oracle Properties." *Journal of the American Statistical Association* 101 (476): 1418–29. <https://doi.org/10.1198/016214506000000735>.
16. Toka, O., Çetin, M., and Arslan, O. (2021). "Robust regression estimation and variable selection when cellwise and casewise outliers are present", *Hacettepe Journal*, Volume 50 (1), 289 – 303.
17. Al-Nuaimi, A. M. (2005). "Variable Selection in Ridge Regression", Un-Published MSc. Thesis, Faculty of Computer Science and Mathematics, University of Mosul – Iraq.

Robust Weighted Least Squares Method using different schemes of M-estimators (RWLSM), A comparative Study

Fatima M. Selevany , Bashar A. Al-Talib

Department of Statistics and Informatics, College of Computer Science and Mathematics, University of Mosul, Mosul, Iraq

Abstract: In this research, it was reducing or excluding the effect of not satisfying the assumption of normal distribution of the data, due to the presence of types of outlying values in it when we wish to choose the best regression equation by robust methods, and this was achieved by introducing weights from the robust methods in the estimate and testing their robustness and suitability for the model in advance, And then selecting the weights resulting from the highest efficient robust methods and introducing these weights in the stages of selecting best regression equation, which results in a model that achieves two characteristics at the same time, which are robustness and reducing dimensions in return for increasing efficiency. The simulation approach was used on models with different dimensions, different sample sizes, and different contamination rates in the dependent variable once, in the independent variables again, and in both together, with a focus on studying the possible impact of the presence of outliers on the variables that will remain in the model and the variables that will be deleted. To achieve the idea of the paper, a number of robust estimation methods were studied, and the results were compared with the ordinary least squares method (OLS) and the robust adaptive LASSO method on experimental data using simulation, as well as on data for a sample of thalassemia patients in Nineveh province.

Keyword: Variables Selection ,(M) M-estimators ,Weighted Least Squares ,Robust Regression ,Outliers Detection.