

دمج البوتسtrap الاحتمالي الموزون مع طريقة لارس الحصينة لاختيار المتغيرات في نموذج الانحدار الخطّي بوجود مشكلة الأبعاد العالية والقيم الشاذة.

<https://doi.org/10.29124/kjeas.1547.19>

أ. د. باسم شلبيه مسلم⁽²⁾

زينه حكمت عبد المنعم الأمير⁽¹⁾

الجامعة العراقية / كلية الإدارة والاقتصاد

جامعة واسط / كلية الإدارة والاقتصاد

المستخلص

تم في هذا البحث اقتراح خوارزمية جديدة لاختيار المتغيرات المهمة في نموذج الانحدار بوجود مشكلة الأبعاد العالية والقيم الشاذة، من خلال توظيف ودمج أسلوب البوتسtrap الاحتمالي الموزون في طريقة Weighted Bootstrap probability - Robust Least Angle Regression - Selecting WBP-LARS اختصاراً (WBP-LARS)، ومقارنتها مع طريقة اختيار أخرى هي طريقة لارس الحصينة المعتمدة على أسلوب البوتسtrap العادي والمعروفة اختصاراً (B-LARS) تجريبياً بالمحاكاة وتطبيقياً بالاعتماد على بيانات حقيقة تتعلق بالقيمة السوقية لبعض المصارف الأهلية في سوق الأوراق المالية للمدة الزمنية 2010-2017. وقد تضمنت المقارنة في المحاكاة حالتين لعدد المتغيرات التوضيحية المطلوب اختيارها $K=5$ ، فضلاً عن حالتين عندما ($n < P$) ($n > P$) وأحجام عينات (50، 70، 20، 26)، وبقيمة ارتباط 0.95 وبنسب تلوث مختلفة (0.05، 0.10، 0.15) = α . وقد حُلِّمَ البحث إلى استنتاجات أهمها تحديد عدد المتغيرات المهمة بين الأعداد 7-10 أفضليّة طريقة (WBP-LARS) على طريقة (B-LARS). في حين ظهرت أفضليّة بسيطة لطريقة (B-LARS) على الطريقة المقترنة عندما ($P > n$). وحجم العينة بعيد عن عدد المتغيرات الكلية في النموذج، ولكن تقارب الكفاءة كلما اقترب حجم العينة بشكل كبير من عدد المتغيرات، ويمكن أن يعتمد عليها في عملية الاختيار للمتغيرات في هذه الحالة.

الكلمات المفتاحية : البوتسtrap الاحتمالي الموزون، البوتسtrap العادي ، الأبعاد العالية ، القيم الشاذة ، طريقة LARS الحصينة.

Abstract

In this research, a new algorithm was proposed to select the important variables in the regression model with the presence of two problems of high dimensions and outliers by employing and integrating the weighted bootstrap probability - Robust Least Angle Regression Selecting (WBP-LARS) and comparing it with another selection method. It is a method of impregnable Lars based on the regular bootstrap method, known as (B-LARS), empirically simulated and applied, based on real data related to the market value of some private banks in the stock market for the period 2010-2017. The comparison in the simulation included two cases for the required number of explanatory variables Choosing ($K = 5, K = 7$) as well as two cases when ($n > P$) ($n < P$) and sample sizes (50, 70, 20, 26) with a correlation value of 0.95 and with different contamination rates $\alpha = (0.05, 0.10, 0.15)$. The research concluded with conclusions, the most important of which is determining the number of important variables between the numbers 7-10, the preference of the (WBP-LARS) method over the (B-LARS) method when ($n > P$), while a slight preference for the (B-LARS) method appeared over the proposed method when ($n < P$) and the sample size is far from the total number of variables in the model, but the efficiency converges whenever the sample size is very close to the number of variables and can be relied upon in the selection process for the variables in this case.

Keywords: Weighted Bootstrap Probability, Classical Bootstrap, high dimensions, outliers, Robust LARS Method.

1- المقدمة:

تمهيد 1

يواجه الكثير من الباحثين المعوقات والمشكلات التي من الممكن أن تُصعب من عملية التحليل الإحصائي المتضمنة التقدير Estimation واختبار الفرضيات Testing Hypothesis لبيانات الظاهرة تحت دراسة. وتنتزع هذه المشكلات بحسب ظهورها، فقد تظهر في بيانات الظاهرة نفسها Data (أي طبيعة البيانات) كما هو الحال بوجود القيم الشاذة Outliers التي تظهر في قيم بعض المتغيرات أو كلها - موضوع الدراسة- والتي غالباً ما تكون حاضرة في الواقع التطبيقي ، وقد تظهر تلك المشكلات في هيكلية (طبيعة) الأنماذج الرياضي Model (أي في مكوناته)، والذي يمثل تلك البيانات، كما هو الحال في أنماذج الانحدار نتيجة لخرق أحد شروطه أو فرضياته الأساسية، ولا يقتصر حدوثها عند استعمال أسلوب تحليل

الانحدار فقط ، بل قد تحدث عند استعمال أغلب الأساليب الإحصائية في عملية تحليل البيانات إحصائياً، ولا سيما وجود القيم الشاذة.

ومن الشروط المهمة الواجب توفرها عند التحليل الإحصائي لبيانات ظاهرة بأسلوب تحليل الانحدار الخطّيتحقق شرط أن حجم العينة sample يكون أكبر من عدد المتغيرات التوضيحية explanatory variables بما فيها المتغير التوضيحي المرتبط بالحد الثابت ($P > n$) ، أو بعبارة أخرى أن يكون حجم العينة أكبر من عدد المعلمات في الأنماذج؛ لكي يتم تعريف معكوس المصفوفة X^{-1} ، ثم تطبيق طريقة المرربعات الصغرى الاعتيادية OLS بوصفها الطريقة الشائعة والمفضّلة لدى أغلب الباحثين مع تحقق بقية الفرضيات الأساسية الأخرى، ومنها: ثبات تجانس التباين للأخطاء العشوائية Error terms في الأنماذج، ما يعني خلوّ الأنماذج من مشكلة عدم تجانس التباين Heteroscedasticity، وكذلك عدم ترابط الأخطاء العشوائية لأنماذج فيما بينها مما يعني خلوه من مشكلة الارتباط الذاتي Autocorrelation ، فضلاً عن عدم وجود مشاكل في البيانات نفسها، ومنها وجود القيم الشاذة Outliers ، ما يعني الحصول على مقدرات كفؤة لها صفة أفضل مقدر خطّي غير متحيز Best Linear Unbiased Estimator (BLUE)، والتي يمكن الاعتماد عليها في عملية التنبؤ المستقبلي للظاهرة تحت الدراسة، ومن ثم تسهل عملية اختيار المتغيرات التوضيحية المهمة المؤثرة فعلاً على متغير الاستجابة (المعتمد) بالاستعانة بأحد أساليب الاختيار الشائعة: (أسلوب اختيار الانحدار المترجر المتسلسل) Stepwise Regression ، أسلوب الاختيار الامامي Forward ، وأسلوب الاختيار الخلفي Backward .

2-1 مشكلة البحث

في الواقع التطبيقي لا تتحقق شروط الأنماذج كلها GLR General Linear Regression ، فقد يكون ($P < n$) ، مما يسبّب مشكلة الأبعاد العالية high dimensions ، فضلاً عن ذاك قد تكون بيانات الظاهرة المدروسة تعاني من وجود القيم الشاذة outliers سواءً في متغير الاستجابة أم النقاط الرافعة leverage points في المتغيرات التوضيحية أو وجودهما معاً في البيانات ووجود هاتين المشكلتين في آنٍ واحدٍ.

يؤدي إلى فقدانطرائق الاعتيادية الشائعة كفاءتها لاختيار المتغيرات التوضيحية المهمة في الأنماذج التي تؤثر على متغير الاستجابة Variables Selection والتقدير Estimation لمعلماته، وبدوره يفقد الثقة بتلك المقدرات؛ نتيجة عدم واقعيتها وابتعادها عن الحقيقة. ومن هنا تظهر مشكلة البحث فالباحث في موضوع الانحدار يهتم بعملية التنبؤ المستقبلي للظاهرة تحت الدراسة، مما يفرض عليه البحث عن طرائق اختيار وتقدير مناسبة، تتصف بالدقة والكفاءة بوجود مشكلتي الأبعاد العالية والقيم الشاذة.

وقد نوقشت مشكلة البحث في كثير من البحوث المنشورة في خوارزميات الاختيار وطرائق التقدير لمعلمات أنماذج الانحدار المتعدد بوجود مشكلتي الأبعاد العالية والقيم الشاذة، ذكر منها على سبيل المثال في عام 2005 قدم الباحث khan وآخرون [14] بحثاً تضمن أسلوبين مختلفين للحصول على مقدرات LARS الحصين Robust؛ التعامل مع مشكلتي الأبعاد والقيم الشاذة في أنماذج الانحدار الخطّي GLR للحصول على أنماذج التنبؤ الخطّي؛ لأنّ طريقة LARS تقليدية. وفي عام 2008 قدم الباحثان Barrera & Aelst

خمسة بحوث اقترحوا فيها أسلوباً تمهدياً حصيناً لاختيار الحصين للمتغيرات المؤثرة في الأنماذج بواسطة استعمال مقدرات نقطية حصينة robust points estimators بوجود مشكلة الأبعاد العالية والقيم الشاذة . وفي عام 2009 قدمت الباحثة Midi وآخرون [17] بحثاً تضمنَت خوارزمية بوتسناب (تمهدية) حصينة Least robust bootstrap algorithm ، والتي تعتمد على طريقة المربعات الصغرى المشدبة Dynamic Robust Trimmed squares LTS ، والتي لا تتأثر بالقيم الشاذة، أطلق عليها تسمية DRBLTS (Bootstrap-LTS). وفي عام 2011 قدم الباحث Alfons وآخرون [4] بحثاً لوصف أسلوب تحديد المتغيرات التوضيحية المؤثرة على المتغير المعتمد، وتم تسميته Bootstrap Robust Least Angle Regression Selected BRLARS ، والتي اعتمدت على أسلوب بوتسناب العادي classical bootstrap Hettigoda اطروحة دكتوراه ناقش فيها انحدار Least Angle Regression (LAR) بالتفصيل، وبين أنه يشابه في أساس عمله لانحدار المرحلة الأمامي forward stagewise regression ولكنه يعاني من صعوبات في العمليات الحسابية ، وفي عام 2017 قدم الباحث السرائي [1] اطروحة دكتوراه ركز فيها على استعمال المقدرات الحصينة الجزئية في عملية التقدير لمعلمات أنماذج الانحدار الخطّي المتعدد بوجود المشكلتين دون التركيز على طرائق عملية اختيار المتغيرات التوضيحية المهمة. وفي عام 2021 قدم الباحث Lindskou وآخرون خمسة عشر بحثاً تضمنَ تطويراً لطريقة جديدة للكشف عن القيم الشاذة outliers في البيانات ذات الأبعاد العالية high dimensional ، وقد استعنوا بالتقنيات الحاسوبية المتطرورة لاستعمال نماذج الرسم . ويلاحظ أنَّ أغلب البحوث لم تتطرق إلى عملية اختيار المتغيرات، بل ركزت على عملية التقدير بطرائق جزئية حصينة فضلاً عن أنها لم تتطرق إلى أسلوب البوتستراب الاحتمالي الموزون واعتمدت البوتستراب العادي .

3-1 هدف البحث

يهدف البحث إلى اقتراح طريقة اختيار للمتغيرات المهمة المؤثرة على متغير الاستجابة في أنماذج الانحدار الخطّي المتعدد بوجود مشكلتي الأبعاد العالية والقيم الشاذة وباحتين ($P < n$) ، من خلال الدمج بين أسلوب البوتستراب الاحتمالي الموزون Weighted Bootstrap Probability وطريقة لارس RIARS ، والتي يرمز لها بـ WBP-RLARS ومقارنتها مع طريقة لارس الحصينة المدمجة مع أسلوب البوتستراب العادي B-RLARS بأسلوب المحاكاة فضلاً عن إجراء تطبيق عملي على بيانات حقيقة تتعلق بالقيمة السوقية لبعض المصارف الأهلية في سوق الأوراق المالية للمدة الزمنية 2010-2017

4-1 هيكليّة البحث

ومن أجل بلوغ البحث غايته فقد تم تقسيمه على مباحث الأول منها تضمن المقدمة وأهمية البحث، وكذلك مشكلة البحث وهدفه وبعض الاستعراض المرجعي لموضوع البحث. أمّا المبحث الثاني فتضمن الجانب النظري للبحث متضمناً بعض المفاهيم الأساسية لمصطلحات موضوع البحث فضلاً عن استعراض طرائق الاختيار لمتغيرات أنماذج الانحدار الخطّي العام (المتعدد) GLR المعتمدة على أسلوب البوتستراب Bootstrap وهي :

- 1- Bootstrap- Least Angle Regression Selection (B-LARS).
- 2- Bootstrap- Robust Least Angle Regression Selecting(B-RLARS).
- 3- Weighted Bootstrap probability - Robust Least Angle Regression Selecting (WB-RLARS).

وتضمن المبحث أيضاً طريقة التقدير LARS العادية، وطريقة MM إحدى أنواع طريقة M . في حين تضمن المبحث الثالث الجانب التجاري والعملي. وأخيراً المبحث الرابع تضمن الاستنتاجات والتوصيات.

2- الجانب النظري

1-2 مفاهيم أساس:

1-1-2 مشكلة الأبعاد العالية **High dimensions**

عادة تحدث مشكلة الأبعاد العالية عند زيادة عدد المتغيرات التوضيحية في نموذج الانحدار، اي إن $n \geq p$) إذ تظهر ارتباطات وهمية بين المتغيرات التوضيحية ومتغير الاستجابة، مما يؤدي إلى الاستدلال الإحصائي الخاطئ. ويظهر تأثير المشكلة على معكوس المصفوفة (XX^{-1}) إذ يكون غير معرف، وعندها لا يوجد حل للنظام [1] [10].

2-1-2 القيم الشاذة: **Outliers**

يطلق مصطلح القيم الشاذة outliers على المشاهدات أو القيم التي تختلف بسلوكها عن أغلبية البيانات، أي إن المشاهدات التي لا تشبه بسلوكها السياق العام لبيانات العينة تحت الدراسة. بعبارة أخرى أنهما تكون مشاهدات غير متنسقة (inconsistent) مع المشاهدات الأخرى ، وقد تتوارد هذه القيم في بيانات متغير الاستجابة أو في بيانات المتغيرات التوضيحية أو تتوارد في كلها [16] [7] [9].

3-1-2 أسلوب البوتراب : **Bootstrap approach**

يعتبر أسلوب البوتراب من أساليب الإحصاء المهمة الذي يستند على فكرة إعادة المعينة لعدد محدد من المرات بحسب قناعة الباحث وهو على نوعين [12] [16] :

أولاً: البوتراب التقليدي (الباقي) **Classical Bootstrap**

يعتمد البوتراب التقليدي على تقديرات OLS للحصول على بقية البيانات الأصلية. ثم يتم الحصول على عينات bootstrap من خلال إعادة أخذ العينات المتباعدة من الانحدار الأصلي. وخوارزميته كالتالي :

1- تطبق طريقة OLS بالاعتماد على بيانات العينة الأصلية للمشاهدات للحصول على مقدرات المعلمات

$$\hat{Y}_i = f(X_i, \hat{\beta}_{ols})$$

2- استخرج الباقي $e_i = Y_i - \hat{Y}_i$ وإعطاء احتمال $\frac{1}{n}$ لكل قيمة .

3- اسحب عينة عشوائية من \hat{Y}_i مع الاستبدال ، أي e_i^* مأخوذ من e_i ومرفق بها \hat{Y}_i للحصول على

$$Y_i^{*b} = f(X_i, \hat{\beta}_{ols}) + e_i^{*b} \quad \text{قيمة الباقي } Y_i^{*b} \text{ أي إن:}$$

4- تطبيق طريقة OLS بالاعتماد على بيانات Y_i^{*b} على X الثابت للحصول عليها $\hat{\beta}_{ols}^{*b}$

5- كرر الخطوتين 3 و 4 للحصول على مرات B $\hat{\beta}_{ols}^{*1}, \dots, \hat{\beta}_{ols}^{*B}$ هي مكررات bootstrap .

ثانياً : بوتستراب الأزواج (المشاهدات) Pairs Bootstrap

تنص هذه الطريقة على أخذ ازواج من العينة الأصلية ذات الحجم n وهي :

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$$

ويحسب عدد B من العينات كأزواج من حجم العينة الأصلية نفسه بإرجاع ونحصل على عينات البوتستراب:

$$X^{*b} = (x_1^*, y_1^*), (x_2^*, y_2^*), \dots, (x_n^*, y_n^*) \quad , b = 1, 2, \dots, B$$

وبقدر عدد العينات البوتسترابية B يتم إيجاد مقدار المربعات الصغرى البوتسترابية لنجعل على :

$$b_b^* = b_1^*, b_2^*, \dots, b_B^*$$

إذ يمثل المتجه b_b^* بصفة عامة مقدرات بوتستراب أزواج المشاهدات لمعلمات الانحدار الخطّي ، وتكون تلك المقدرات ما يعرف بالتوزيع البوتسترابي والذي توقعه يمثل تقدير لموجة معلمات المجتمع . β

2-2 نموذج الانحدار الخطّي بوجود الأبعاد العالية ومشكلة القيم الشاذة.

يختلف نموذج الانحدار الخطّي ذو الأبعاد العالية عن نموذج الانحدار الخطّي العادي في عدد المتغيرات التوضيحية، إذ يكون فيه عدد المتغيرات التوضيحية أكبر من حجم العينة، أي إن $(p > n)$ عندما تظهر مشكلة الأبعاد العالية HD (High Dimensional) [10][7] ويعُرف كالتالي :

$$Y = X\beta + \varepsilon \quad , i = 1, 2, \dots, n \quad , j = 1, 2, \dots, p \quad , p > n$$

وبغض النظر عن هذين النوعين من النماذج فإن اهتمام الباحث ينصب على مجموعة من الأهداف، وهي: التقدير، والتتبؤ، و اختيار المتغيرات selection variable [14]. ويعد اختيار المتغيرات التوضيحية المهمة في مقدمتها؛ لأنّه يحدد المتغيرات المهمة في الأنماذج، وبقيها فيه وبهمل الأخرى، مما يؤدي إلى خفض عدد المتغيرات التوضيحية في الأنماذج بأكمله، ما يعني اعتماد المقدرات التي تدخل في المعادلة التقديرية فقط لغرض إجراء التتبؤ؛ ولذا نجد أنّ أغلب البحوث قد ركزت على عملية التقدير دون طرائق اختيار المتغيرات مع إن طرائق الاختيار لكل منها خصائصها التي تتميز بها عن الأخرى، والتي قد تفقدها عند وجود مشكلات أخرى، ومنها مشكلة وجود القيم الشاذة متغير الاستجابة والنقط الرافعة في المتغيرات

التوسيعية أو في كليهما؛ لذا اقترحت أيضاً طائق اختيار متغيرات حصينة للحفاظ على قوة الطريقة والاختيار، ومنها طريقة لارس الحصينة Robust LARS Method والتي تُعد من الطائق الجيدة في عملية الاختيار لتقارن مع الطيق المقتراح من قبل الباحثة بدمج أسلوب البوتسناب وهي كالتالي :

1-2-2 دمج أسلوب البوتسناب مع طيق لارس الحصينة Method

للحصول على طيق البوتسناب لارس الحصينة B-RLARS لا بد من عرض طيفي LARS العادي و LARS الحصينة فقد اقترح Efron وآخرون [8] في عام (2004) طيق انحدار الزاوية الصغرى باجراءات الاختيار المرحلي الأمامي وطيق لاسو LASSO (Tibshirani 1996) ، إذ يوفر SFS ، ترتيب أو ترتيبا تدخل فيه المتغيرات التوضعية أو يثبتها في نموذج الانحدار الخطى، وهذا الترتيب أو التسلسل للمتغيرات المهمة غالبا ما يكون نفسه في LASSO أو SFS، ولكن يتم الحصول عليه بطريقة فعالة حسابيا، وهذه ميزة تحسب لطيق LARS. وهناك ميزة أخرى تحسب لطيق هي أنه يمكن إيجاد الترتيب أو التسلسل للمتغيرات التوضعية المهمة في نموذج الانحدار من مصفوفة الارتباط البيانات دون الاعتماد على البيانات نفسها. للتوضيح نفرض توفر المتغيرات Y, X_1, \dots, X_K وبشكل معيارية standardized بالاعتماد على الوسط الحسابي والانحراف المعياري لـ كل متغير ، ولنفرض أن r_{jY} يمثل معامل الارتباط بين متغير الاستجابة Y والمتغير التوضعى X_j ، وخوارزمية طيق LARS كالتالي [8] :

: [11]

1- لتكن B مجموعة تتضمن المتغيرات المهمة أو الفعالة وفي أول الأمر $\emptyset = B$ وأيضاً ليكن $\emptyset = S_B$ يمثل موجّه الإشارات للارتباطات بين المتغيرات التوضعية .

2- نحدد قيمة أكبر قيمة مطلقة للارتباطات بين المتغيرات التوضعية وإشارتها وتحديد قيمة الارتباط r وكالاتي

$$X_m = \operatorname{argmax} |r_j| , s_m = \operatorname{sign}(r_m) , r = s_m r_m$$

3- ضم المتغير X_m إلى المجموعة B في الخطوة الأولى، وكذلك ضم إشارة الارتباط إلى الموجّه S_B أي إن :

$$B \leftarrow B \cup \{X_m\} , S_B \leftarrow S_B \cup \{s_m\}.$$

4- حساب المقدار الآتي :

إذ إن :

بموجّه عناصره كلها مسؤولة إلى الواحد الصحيح ومرتبته (عدد المتغيرات المهمة أو الفعالة مضروبا بـ 1_B)

$H_B = \text{diag}(S_B))$ أي إن H_B

: مصفوفة قطرية عناصرها عناصر الموجه S_B .

ويُثمن في هذه الخطوة احتساب الآتي :

$$\theta_B = b(H_B C_B H_B)^{-1} 1_B , \quad b_j = (H_B r_{jB})' \theta_B , \quad j \in B^c ,$$

: موجّه الارتباطات بين المتغير X_j والمتغيرات التوضيحية المهمة أو الفعالة.

: المجموعة المكملة للمجموعة B وتضمّ المتغيرات التوضيحية التي لم يتم اختبار أهميتها بعد.

وتتجدر الإشارة هنا أنه في حالة وجود متغير مهم أو فعال واحد فقط تم اختياره فإنه يكون :

$$b = 1 , \quad \theta = 1 , \quad b_j = r_{jm}$$

-5- لـكلّ متغير بتسلسل j إذ إن $(X_j \in B^c)$ نقوم باحتساب المسافتين الآتتين:

$$\delta_j^+ = \frac{r - r_{jY}}{b - b_j} , \quad \delta_j^- = \frac{r + r_{jY}}{b + b_j}$$

$$\delta_j = \min(\delta_j^+, \delta_j^-) \quad \text{ويُثمن اختيار القيمة الأقل، أي إن :}$$

$\delta = \min(\delta_j, X_j \in B^c)$ وبالاعتماد على قيمة δ يتم تحديد المسافة δ كالتالي:

إذا قيمة δ الأقل مساوية إلى δ^+ نضع $s_{\min} = +1$ ، وأي شيء آخر نضع -1 .

وفي هذه الخطوة يتم تعديل قيم الارتباط وكالآتي : $r \leftarrow r - \delta b$ ، $r_j = r_j - \delta b_j \quad \forall X_j \in B^c$.

-6- نعيد إجراء الخطوات 3، 4، 5 .

ويمكن كتابة مصفوفة الارتباطات للمتغيرات التوضيحية كالتالي :

$$C_X = [1 \ r_{12} \ \dots \ r_{1K} : \vdots \vdots \ r_{K1} \ r_{K2} \ \dots \ 1] \quad \dots(1)$$

وبافتراض أن r_{mY} يملك القيمة المطلقة الأعلى من بين قيم الارتباطات بين المتغيرات التوضيحية ومتغير الاستجابة، وايضاً نفرض أن $s_m = \text{sign}(r_{mY})$ يمثل إشارة الارتباط، أي إن

ووفقاً لما نقدم فإن المتغير X_m يصبح المتغير المهم أو الفعال الأول الواجب بقاؤه في الأنماذج ومن ثم يكون التبؤ الحالي لمتغير الاستجابة μ وبعبارة أخرى $\mu \rightarrow 0$ ، إذ يتطلب تعديله عن طريق التحرك على طول اتجاه $S_m X_m$ حتى مسافة معينة δ يمكن التعبير عنها من حيث الارتباطات بين المتغيرات ، وبحديد قيمة المسافة δ تقوم طريقة LARS بتحديد المتغير المهم الفعال الثاني الواجب وجوده في الأنماذج، وبمجرد أن يكون لدينا أكثر من متغير مهم أو فعال أو نشط ، يقوم LARS بتعديل التبؤ الحالي على طول الاتجاه المتساوي ، أي الاتجاه الذي له زاوية متساوية (ارتباط) مع المتغيرات المشتركة الفعالة جميعها أو المهمة الموجودة في الأنماذج. ومن هذه الآلية (اعتمادها على الزاوية المتساوية للمتغيرات المشتركة جميعها) اكتسبت الطريقة تسميتها [11] [8].

إن التحرك على طول هذا الاتجاه وفق الزاوية المتساوية يضمن أن الارتباط الحالي لكل متغير مهم أو فعال ينشط مع بقية المتغيرات الأخرى سينخفض بالتساوي، وقد وضح الباحث Khan وأخرون [14] أسلوب الاشتغال لتحديد المسافة δ لمتغير توضيحي مهم أو فعال في ملحق البحث سلسل A (يمكن الاطلاع عليه) وقيمتها كالتالي:

$$\delta(+X_j) = \frac{r - r_{jY}}{1 - r_{jm}} \quad \dots(2)$$

$$\delta(-X_j) = \frac{r + r_{jY}}{1 + r_{jm}} \quad \dots(3)$$

وعلى ضوء المعادلين (2) (3) يتم اختيار القيمة الأقلّ بينهما، وكذلك الحد الأدنى على المتغيرات غير المهمة أو الفعالة جميعها، ثم وفقاً لقيمة المسافة δ ، يتم تحديد المتغير المشترك الذي يدخل إلى الأنماذج مرفقاً للمتغير السابق في الاختيار .

واعتماداً على ما نقدم أشار الباحث Khan وأخرون أن طريقة LARS بشكلها الحالي تتأثر بتلوّث البيانات. وقد وضح ذلك بمثال يمكن الاطلاع عليه ، ولذا اقترحوا أسلوبين لزيادة الحصانة لطريقة LARS وسمّي الأسلوبين: بالإضافة، والتنظيف plug-in and cleaning مكونات طريقة LARS غير الحصينة non robust وهي: (المتوسط، والتباين، والارتباط) بمكونات حصينة robust، وأكثر مقاييس النزعة المركزية حصانة وأقل تأثيراً بالقيم الشاذة هو الوسيط median (med) في حين يُعدّ مقياس التشتت (mad) الأكثر حصانة. ولكن الحصول على الارتباط الحصين correlation robust يتوجّب توفر على بيانات بقدر d من الأبعاد، وهذا يعني تطلب وقتاً طويلاً بعكس المقياسين الآخرين. وأشار الباحث Khan وأخرون أن الباحث Huber ذكر في كتابه عام 1981 أن أساليب الأزواج الحصينة pairwise robust ليست متكافئة، وهذا يعني أنها حساسة للقيم المتطرفة ثنائية الأبعاد، لذا يُعدّ أحد الحلول للحصول على الارتباط الحصين باستعمال ارتباطات قوية مشتقة من مقدار التغاير المتساوي الأزواج. وهناك طريقتان لإيجاد الارتباط الحصين إحداهما طريقة مقدر M ثنائي المتغير، الذي عرّفه الباحث Maronna ، والثانية حساب الارتباط ثنائي المتغير من بيانات Windsorized ثنائية المتغير، وبقدر تعلق البحث سيُتم الاعتماد على الطريقة الثانية.

تُعدّ مهمة الحصول على مقدّر ارتباط حصين وحساب سريع من أهمّ المعوقات التي تواجه الباحثين بوجود بيانات الأبعاد العالية dimensions. في 1981 قدّم الباحث Huber فكرة البيانات أحادية البعـد ونذرريشن one-dimensional Windsorization وبالاعتماد على البيانات المحولـة إلى الصيغة القياسيـة يُتمّ حساب معاملات الارتباط التقليديـة، ثمّ جرت إعادة تدقـيقـة لأسلوب الباحث Huber لتقدير عناصر مصفوفة الارتباطـات ذات الأبعـاد العـالـية أو الكـبـيرـة، إذـ إنـ البيانات المحـولـة إلى الصـيـغـة الـقـيـاسـيـة لـبعد وـاحـدـ (متـغـيرـ واحدـ بـحـجمـ عـيـنةـ nـ) ليـتمـ الحصولـ عـلـيـهاـ كـالـآـتـيـ :

$$u_i = \frac{f_c((X_i - \text{med}(X_i)))}{\text{med}(X_i)} , i=1, 2, \dots, n \quad (4)$$

$$f_c = \min[\max(-c, X_i), c] \quad (5)$$

إذـ إنـ قيمةـ cـ يـتمـ اختيارـهاـ منـ قـبـلـ البـاحـثـ .

ونستنتجـ منـ الصـيـغـةـ (4ـ)ـ أنـ أـسـلـوبـ (Huberـ one-dimensional Windsorizationـ)ـ يـتـمـيزـ بـسـرـعـةـ الحـسـابـاتـ الـرـياـضـيـةـ لـتـحـوـيلـ الـبـيـانـاتـ بـشـكـلـ الصـيـغـةـ الـقـيـاسـيـةـ،ـ وـلـكـنـ فـيـ الـبـيـانـاتـ ذـاتـ الـأـزـواـجــ وـلـكـنـ تـفـقـدـ هـذـهـ المـيـزةـ فـيـ الـبـيـانـاتـ الثـانـيـةـ bivariate dataـ عـلـىـ الرـغـمـ مـنـ وـجـودـ هـذـهـ المـيـزةـ لـكـنـ لاـ يـتـغـلـبـ عـلـىـ تـأـيـرـ الـقـيـمـ الشـاذـةـ بـشـكـلـ نـهـائـيـ،ـ وـنـجـدـ أـنـ بـعـضـ قـيـمـ الـبـيـانـاتـ يـتـرـكـهاـ دـوـنـ تـغـيـيرـ.

ولـعـلاـجـ هـذـاـ القـصـورـ فـيـ التـحـوـيلـ Windsorizationـ الـأـحـادـيـ يـمـكـنـ تـحـوـيلـ الـبـيـانـاتـ وـنـذـرـرـيـشنـ الثـانـيـ اـعـتـمـادـاـ عـلـىـ القـطـعـ النـاقـصـ الـأـوـلـيـ للـتـسـامـحـ الـأـوـلـيـ لـغـالـيـةـ الـبـيـانـاتـ أوـ التـفاـوتـ الـمـسـمـوحـ initial tolerance ellipseـ عـلـىـ هـذـهـ هـذـهـ القـطـعـ النـاقـصـ باـسـتـعـالـ التـحـوـيلـ ثـانـيـ المـتـغـيرـ الآـتـيـ:

$$u \left(\min \sqrt{\frac{c}{d(X)}}, 1 \right) X , \quad X = (X_1, X_2) \quad (6)$$

ـ d(X)ـ تمـثـلـ مـسـافـةـ مـهـانـوـبـسـ Mahalanobisـ distanceـ الـثـانـيـ الـأـوـلـيـ فـيـ مـصـفـوفـةـ الـارـتـبـاطـاتـ وـالـذـيـ يـمـثـلـ مـعـالـمـ الـارـتـبـاطـ التقـلـيدـيـ كـيـ يـتـمـ حـسـابـ بـيـانـاتـ المـعـدـلـةـ adjusted Windsorization dataـ وـالـذـيـ تـحـسـبـ بـالـصـيـغـةـ الـآـتـيـةـ [9ـ]:ـ

$$d(X) = [(X_1 - \bar{X}_2)^T R_0^{-1} (X_1 - \bar{X}_2)]^{0.5} \quad (7)$$

ـ ماـ تـقـدـمـ يـمـثـلـ طـرـيـقـةـ لـأـرـسـ الـحـصـيـنـةـ LARSـ وـأـخـتـصـارـاـ Robustـ الـبـوـتـسـتـرـابـ Bootstrap Approachـ فـيـ عـلـيـةـ اـخـتـيـارـ الـمـتـغـيرـاتـ التـوضـيـحـيـةـ الـمـهـمـةـ أوـ الـفـعـالـةـ فـيـ أـنـمـوذـجـ الـانـحدـارـ الـخـطـيـ بـجـوـدـ الـبـيـانـاتـ ذـاتـ الـأـبعـادـ الـعـالـيـةـ high dimensionsـ لـلـحـصـولـ عـلـىـ نـتـائـجـ أـكـثـرـ وـاقـعـيـةـ وـمـوـثـقـيـةـ إـذـ يـتـمـ تـولـيـدـ Mـ مـنـ عـيـنـاتـ الـبـوـتـسـتـرـابـ (ـ التـهـيـدـ)ـ مـنـ الـبـيـانـاتـ الـأـصـلـيـةـ وـفـقـ النـوعـ الـثـانـيـ مـنـ أـنـوـاعـ الـبـوـتـسـتـرـابـ ،ـ وـنـسـتـعـالـ بـعـدـ ذـلـكـ طـرـيـقـةـ RLARSـ التـسـلـسـلـ الـمـقـابـلـ لـكـلـ مـنـ الـمـتـغـيرـاتـ التـوضـيـحـيـةـ لـكـلـ عـيـنةـ مـنـ عـيـنـاتـ الـبـوـتـسـتـرـابـ وـتـصـنـيـفـ الـمـتـغـيرـاتـ بـشـكـلـ

رتب من التسلسل 1 إلى التسلسل V من حيث الأهمية، ثم لكل متغير مشترك يمكن أن تأخذ المتوسط على هذه الرتب M ، والمتغيرات المشتركة m التي تمتلك أصغر متوسط الرتب تكون المجموعة المخفضة أو المجموعة المختزلة set reduced. وعليه يمكن وضع خوارزمية لتطبيق B-RLARS لاختيار المتغيرات في نموذج الانحدار بوجود البيانات ذات الأبعاد العالية وكالآتي [14][11] :

1- توليد عينات بعدد M من عينات البوتستراب Bootstrap samples باختيار زوج البيانات مع الإرجاع من البيانات الأصلية، أي إن :

نفرض أن هناك بيانات للمتغيرات التوضيحية ومتغير الاستجابة ($Y, X_1, X_2, X_3, \dots, X_K$) فتكون أزواج المشاهدات كالآتي [12][17] :

$$(8) \quad (Y_i, X_{i1}, X_{i2}, \dots, X_{ik}) = (Y_i, X_i), \quad (X_i = X_{i1}, X_{i2}, \dots, X_{ik}), i = 1, 2, \dots, n$$

وبالاعتماد على تعريف البيانات في (8) ستكون البيانات الأصلية (أزواج البيانات)، التي سيتم سحب عينات البوتستراب منه كالآتي:

$$(9) \quad (Y_1, X_1), (Y_2, X_2), \dots, (Y_n, X_n)$$

وعلى ضوء بيانات الأزواج المعرفة في (9) يتم سحب عينات البوتستراب مع الإعادة وكالآتي:

$$(10) \quad (\dot{Y}_{1m}, \dot{X}_{1m}), (\dot{Y}_{2m}, \dot{X}_{2m}), \dots, (\dot{Y}_{nm}, \dot{X}_{nm}), m = 1, 2, \dots, M$$

2- تطبيق طريقة RLARS على كل عينة واحدة من عينات البوتستراب للحصول على تسلسل أو رتب المتغيرات التوضيحية بحسب الأهمية.

3- ترتيب المتغيرات أو وضعها بتسلسل بحسب أهميتها بالاعتماد على نتائج الخطوة 2.

4- إعادة إجراء الخطوات 1,2,3 لـ M من عينات البوتستراب.

5- يتم حساب الوسط الحسابي لرتب كل متغير توضيحي، والتي حصل عليها من الخطوات 1,2,3,4 وهي بقدر M من عينات البوتستراب.

6- ترتيب المتغيرات التوضيحية بحسب أهميتها بحسب أقيم الأوساط الحسابية للرتب التي حصلت عليها من الخطوة 5، ويكون التسلسل الأفضل للمتغير الذي يملك أصغر وسط حسابي من بينها.

7- يتم تحديد المجموعة المخفضة أو المجموعة المختزلة من المتغيرات التوضيحية المهمة التي تم ترتيبها في الخطوة 6 وبعدد 5 أو 7 من المتغيرات ولا يتعدى 10 متغيرات.

1-2-2 الطريقة المقترنة لاختيار المتغيرات Suggest method for selection

تستند فكرة الطريقة المقترنة لاختيار المتغيرات المهمة في نموذج الانحدار، والحصول على المجموعة المختزلة بتوظيف أسلوب البوتستراب الموزون الاحتمالي، ودمجه مع طريقة لارس الحصينة Weighted Bootstrap with probability RLARS و اختصاراً WBP-RLARS. ولذا فعملية كتابة الخوارزمية للطريقة المقترنة تستوجب عرض أسلوب البوتستراب (التمهيد) الموزون الاحتمالي RLARS ، و اختصاراً WBP ، ودمجهما مع خوارزمية Weighted Bootstrap with probability وكالآتي [18] :

يعد هذا البوتستراب تعديلا على طريقة التشخيص قبل البوتستراب Diagnostic -Before Bootstrap المقدم من قبل الباحث Imon وآخرون [12] لوقاية إجراء البوتستراب من القيم الشاذة Outliers ، وفيها يتم تحديد شيئين: الأول (تحديد أوزان للمشاهدات استعملت دالة هامبل المرجحة Hamble's weighting Psi في تعريف الأوزان weights لنقط البيانات كلها، والثاني (تحديد الاحتمالات لـلكل زوج من البيانات بالاعتماد على الأوزان المستخرجة)؛ ولذا فالأوزان تستخرج بحسب الصيغة الآتية:

$$\phi(u)_{Hample} = \begin{cases} u & 0 \leq \text{abs}(u) \leq a \\ a \text{ sign}(u) & a \leq \text{abs}(u) \leq b \\ a(c - \text{abs}(u)) / (c - b) \text{ sign}(u) & b \leq \text{abs}(u) \leq c \\ 0 & c \leq \text{abs}(u) \end{cases}$$

$$\{ a = 1.31, b = 2.039, c = 4 \} \quad (11)$$

وبهذا يتوقع أن تأخذ القيم الشاذة outliers values أوزاناً تبعاً لخطورتها .

ويتم حساب هذه الأوزان من البوابي القياسي بطريقة مربعات الوسيط الصغرى Least median squares وفق الصيغة (LMS) [19][4] وكالآتي :

$$u_i = \frac{r_i}{\text{MAD}(r_i)}, \quad i = 1, 2, \dots, n$$

$$(12)$$

واعتماداً على الصيغة (12) تستخرج الأوزان وكالآتي :

$$w_i(u_i) = \frac{\phi(u_i)}{u_i}, \quad i = 1, 2, \dots, n$$

$$(13)$$

وبناءً على الصيغة (13) يتم استخراج الاحتمالات المقابلة لـلكل زوج من المشاهدات وفق الصيغة الآتية :

$$p_i = \frac{w_i}{\sum_{i=1}^n w_i} \quad i = 1, 2, \dots, n \quad (14)$$

واعتماداً على الاحتمالات المعرفة في المعادلة 14 يتم تخصيص كل احتمال لزوج من البيانات، معنى أن الاحتمالات p_1, p_2, \dots, p_n مُخصصة لكل زوج من البيانات $(Y_1, X_1), (Y_2, X_2), \dots, (Y_n, X_n)$ ، وبطبيعة الاحتمالات التي قد تكون صفرًا عندما تكون هناك قيمة شاذة يتم تقسيم البيانات على جزأين: البيانات الباقية R ، والبيانات المحذوفة D ، إذ يجب أن يكون للبيانات الباقية قيمة احتمالية أكبر من الصفر ووفقاً لذلك يمكن تعريف البيانات كالتالي :

$$X = [X^R \ X^D], Y = [Y^R \ Y^D] \quad (15)$$

وبافتراض أنَّ وجَهَ مقدَّرات المَعْلَمَات لِأَنْمُوذِجِ الانحدار هو $\hat{\beta}^D$ بعد حذف D من البيانات الشاذة، والذي يستخرج بالاعتماد على البيانات الباقية لمُتَغَيِّرِ الاستجابة Y^R والمُتَغَيِّرِات التوضيحية X^R ، ويمكن توضيح خوارزمية طريقة **WBP** وكالآتي :

1- تقدير مَعْلَمَاتِ الأَنْمُوذِجِ اعتماداً على البيانات الكلية من دون الحذف بطريقة LMS للحصول على أوزان هامبل Hample ثم تحديد الأوزان الأكبر من الصفر، والتي تقابل زوج البيانات الباقية؛ لكي نحصل على وجَهَ المقدَّرات $\hat{\beta}^D$ ، ثم استخراج القيم التنبئية لمشاهدات مُتَغَيِّرِ الاستجابة وكالآتي:

$$\hat{Y}_i^D = g(x_i, \hat{\beta}^D) \quad (16)$$

2- نستخرج الباقي residuals من خلال العلاقة الآتية :

$$\hat{\varepsilon}_i^D = Y_i - g(x_i, \hat{\beta}^D) \quad (17)$$

3- سحب مجموعة من الباقي الجديد \hat{Y}_i^D من الباقي المستخرجة من الخطوة 2 مع الإرجاع باستعمال البوتستراسب، واستخراج الاحتمالات الجديدة وفقاً للصيغة 14 ، وباعتماد الباقي الجديد، ومنها استخراج القيم التنبئية لمُتَغَيِّرِ الاستجابة الجديدة \hat{Y}_i^D إذ إنَّ :

$$\hat{Y}_i^D = g(x_i, \hat{\beta}^D) + \hat{\varepsilon}_i^D \quad (18)$$

4- استخراج وجَهَ المقدَّرات الجديد $\hat{\beta}^D$ بالاعتماد على بيانات البوتستراسب المُتَغَيِّرِ الاستجابة والمُتَغَيِّرِات التوضيحية المقابلة لها .

5- إعادة إجراء الخطوات 3,4 ، لعدد من المرات M من عينات البوتستراسب للحصول على M من وجَهَات المقدَّرات المَعْلَمَات . $\hat{\beta}_m^D, m = 1, 2, \dots, M$

ووفقًا لخوارزمية WBP يمكن توظيفها مع طريقة RLARS للحصول على طريقة اختيار مقتراحه حصينة أخرى تعتمد على أسلوب البوتستراب الاحتمالي with Probability Weighted Bootstrap LARS ونسمى WBP-RLARS وحيث أنها هي خليط من خوارزمية طريقي WBP و LARS وكالآتي :

1- توليد عينات بعدد M من عينات البوتستراب samples للحصول على زوج البيانات من المتغير الاستجابة، والمتغيرات التوضيحية باتباع خطوات خوارزمية WBP الموضحة سلفاً، أي الحصول على الآتي :

$$(\dot{Y}_{1m}, \dot{X}_{1m}), (\dot{Y}_{2m}, \dot{X}_{2m}), \dots, (\dot{Y}_{nm}, \dot{X}_{nm}) , \quad m = 1, 2, \dots, M \quad (19)$$

2- تطبيق طريقة RLARS على كل عينة واحدة من عينات البوتستراب للحصول على تسلسل أو رتب المتغيرات التوضيحية بحسب الأهمية.

3- ترتيب المتغيرات أو وضعها بتسلسل بحسب أهميتها بالاعتماد على نتائج الخطوة 2.

4- إعادة إجراء الخطوات 1,2,3 لـ M من عينات البوتستراب.

5- يتم حساب الوسط الحسابي لرتب كل متغير توضيحي والتي حصل عليها من الخطوات 1,2,3,4 وهي بقدر M من عينات البوتستراب.

6- ترتيب المتغيرات التوضيحية بحسب أهميتها بحسب أقيم الأوساط الحسابية للرتب التي حصلت عليها من الخطوة (5) ، ويكون التسلسل الأفضل للمتغير الذي يملك أصغر وسطاً حسابياً من بينها.

7- يتم تحديد المجموعة المختزلة من المتغيرات التوضيحية المهمة، التي تم ترتيبها في الخطوة 6 وبعد 5 أو 7 من المتغيرات ولا يتعدى 10 متغيرات.

3-2 طريقة التقدير MM

اقترح الباحث Yohai في عام 1987 تقديرًا ، هو مزيج من تقدير قيمة الانهيار العالية high breakdown value والتقدير الفعال efficient estimation estimation ، وهي طريقة مطورة عن طريقة تقدير M ، ويكون مقدار MM هو الحل للمعادلة الآتية [19][4]:

$$\sum_{i=1}^n T(u_i) \left(\frac{Y_i - \sum_{j=0}^K X_{ij} \hat{\beta}_j}{s_{MM}} \right) X_{ij} = 0 \quad (20)$$

: S : يمثل الانحراف المعياري المستحصل عليه من بوافي تقدير s_{MM}

: $T(u_i)$: تتمثل دالة توكي ثنائية الوزن Tukey's biweight function وتساوي :

$$T(u_i) = \begin{cases} \frac{u_i^2}{2} - \frac{u_i^4}{2(1.547)^2} + \frac{u_i^6}{6(1.547)^2} & , -1.547 \leq u_i \leq 1.547 \\ \frac{(1.547)^2}{6} & u_i < -1.547 \text{ or } u_i > 1.547 \end{cases} \quad (21)$$

$$u_i = \frac{\varepsilon_i}{\sigma_s}, \quad \hat{\sigma}_s = \frac{\text{median} |\varepsilon_i - \text{median}(\varepsilon_i)|}{0.6745} \quad (22)$$

ووفقاً لما تقدم يمكن وضع خوارزمية طريقة تقدير MM وكالآتي [21] :

1- إيجاد تقدير معلمات أنموذج الانحدار بطريقة OLS .

2- اختبار الفرضيات الأساسية لأنموذج الانحدار الخطي العادي . Classical regression

3- إجراء اختبار للكشف عن وجود القيم الشاذة في البيانات .

4- حساب قيم الباقي residual بالاعتماد على مقدرات S وفق المعادلة .

5- حساب قيمة . $\hat{\sigma}_s$

6- حساب قيمة $i = 1, 2, \dots, n$. $u_i = \frac{\varepsilon_i}{\sigma_s}$

7- حساب قيمة الوزن w_i وفق الصيغة الآتية :

$$w_i = \begin{cases} \left[1 - \left(\frac{u_i}{4.685} \right)^2 \right]^2 & , |u_i| \leq 4.685 \\ 0 & |u_i| > 4.685 \end{cases}, i = 1, 2, \dots, n \quad (23)$$

8- حساب مقدر معلمات أنموذج الانحدار $\hat{\beta}_{MM}$ بطريقة المربعات الصغرى الموزونة WLS بالاعتماد على الأوزان المحسوبة بالخطوة 7 .

9- إعادة إجراء الخطوات من 5- 8 للحصول على تقارب بين قيم المقدرات .

10- إجراء اختبار معنوية للمعلمات؛ لتحديد أيّ من المتغيرات التوضيحية لها تأثير على متغير الاستجابة .

3- المحاكاة والتطبيق العملي

1-3 المحاكاة :

1-1-3 خطّة المحاكاة :

تم تتنفيذ برنامج المحاكاة باستعمال لغة R. إذ تضمن البرنامج توليد البيانات بالاعتماد على نموذج الانحدار الخطّي المتعدد الآتي :

$$Y = X\beta + \epsilon$$

إذ أن X هي مصفوفة المُتغيّرات التوضيحية ذات بعد $p \times n$ سيُتم توليدها في كلّ مرّة خلال توزيع طبيعي متعدد المُتغيّرات بمتوازن قيمته صفرًا، وتبين $\rho^{ij} = \sigma^2$ ، أي إنّ :

$$X \sim N(0, \rho^{ij}) \quad , \quad \rho = 0.95 \quad (n = 50,70) , \quad \rho = 0.50 \quad n = 20,26$$

أمّا الأخطاء العشوائية ϵ فقد تم توليدها وفقاً للتوزيع الطبيعي بمتوازن صفر وتبين σ^2 ، والذي تم افتراضه . 1

وانّ القيم الأولية للمعلمات كالتالي:

$$\beta = [3,0,3,0,0,3,0,0,0,3,0,0,0,0,3^{15}, 0,0, \dots, 0^{13}]_{1 \times 28}$$

إذ تم اختيار خمسة مُتغيّرات بقيم لا صفرية ، والمتبقي من المُتغيّرات التوضيحية بقيم صفرية كما موضح في فيما تقدّم.

وقد تم تحديد عدد المُتغيّرات التوضيحية $p = 28$ مُتغيّراً (مساوٍ لعدد المُتغيّرات التوضيحية في التطبيق العملي). وأمّا أحجام العينات n فكانت حالات عدّة منها ما هو أكبر من عدد المُتغيّرات التوضيحية $n > p$ كما في أحجام (50,70) ، وبعضها الآخر كانت أصغر من عدد المُتغيّرات التوضيحية $n < p$ كما في أحجام (20,26) ، وأيضاً لضمان وجود القيم الشاذة والنقط الرافعة تم تثبيت نسب التأوث $\alpha = 0.05, 0.10, 0.15$.

وقد تمت عملية التلويث كالتالي:

1. تلويث المُتغيّرات اللاصفرية ($X_1, X_3, X_6, X_{10}, X_{15}$) في مصفوفة المُتغيّرات التوضيحية بـ α من النقاط الرافة من خلال جمع القيم الأصلية بالقيمة 20.
2. تلويث المُتغيّر المعتمد بـ α من القيم الشاذة تتبع توزيع مربع كاي بدرجة حرية 50.
3. تلويث بيانات المُتغير الاستجابة والمُتغيّرات التوضيحية في آنٍ واحدٍ، ووفقاً للخطوتين المذكورين آنفًا .

تكرار توليد البيانات (100) مرّة و تخصيص 50 دورة لعينات البوتاستراب مع مراعاة أحجام العينات المفترضة.

لاختبار كفاءة الطريقة المقترحة (WBP-RLARS) سيُتم تشغيل (100) دورة محاكاة لغرض مقارنتها مع نتائج الطريقة B-RLARS ، ولأنّ الطريقتين تقومان بوضع تسلسل لأهميّة المُتغيّرات التوضيحية بالنسبة

لمتغير الاستجابة؛ لذا تم اختيار أول (5) متغيرات توضيحية من كل طريقة على التوالي، وذلك بحساب معدل عدد المتغيرات الصحيحة (Right Selection)، التي تم اختيارها و معدن عدد المتغيرات التي تم اختيارها خطأ (False Selection) ، هذا الأخير قد يتضمن الـ Over fit والـ Under fit . إن السبب الرئيس لاختيار أول المتغيرات الخمسة أو السبعة يعود إلى إننا وضعنا آلية لاحتساب عدد متغيرات الـ (Under fit) و الـ (Over fit) من عدد المتغيرات الـ (28) وفقاً لنسب التلوث (Under fit) = α ، وقد اعتمدت الباحثة مقياس الاختيار الخاطئ False Selection (0.05, 0.10, 0.15) ، والاختبار الصائب Right Selection لطريقة الاختيار كمقياس للمقارنة في تحديد الأفضل منها في عملية اختيار المتغيرات المهمة، والمؤثرة على متغير الاستجابة، فالطريقة الأفضل هي تلك الطريقة التي ظهر أقل اختيارات خاطئاً وأكبر اختياراً صائباً لتحديد المتغيرات المؤثرة في الأنماذج.

2-1-3 نتائج المحاكاة :

لإجراء المحاكاة تم تكوين حالتين: الأولى عندما يكون حجم العينة أكبر من عدد المتغيرات، والثانية عندما يكون حجم العينة أصغر من حجم العينة، وهذه الحالة تمثل مشكلة الأبعاد العالية high dimensions مع قيمة ارتباط 0.95 ونسبة تلوث 0.15 للقيم الشاذة والنقط الرافعة معاً ولخمس متغيرات وسبع متغيرات عند الاختيار وتفاصيل الحالتين كالتالي :

الحالة الأولى : (n (50, 70) > p = 28

الجدول(1) يعرض نتائج الاختيار الصحيح و الاختيار الخاطئ Right Selection و الاختيار الخاطئ False Selection لطريقي الاختيار ولجمي (50، 70) بحسب نسب التلوث لأول خمس متغيرات فقط.

$n > p = 28, \rho = 0.5, \alpha = 0.05$									
مكان تلوث البيانات		القيم الشاذة في المتغير Y		النقط الرافعة في المتغيرات X ^s		قيم شاذة ونقاط رافعة في ان واحد Outliers & leverage points			
طريقة الاختيار للمتغيرات Selection method	n	اختيار صحيح للمتغيرات Right Selection	اختيار خاطئ للمتغيرات False Selection	اختيار صحيح للمتغيرات Right Selection	اختيار خاطئ للمتغيرات False Selection	اختيار صحيح للمتغيرات Right Selection	اختيار خاطئ للمتغيرات False Selection	اختيار صحيح للمتغيرات Right Selection	اختيار خاطئ للمتغيرات False Selection

B- RLARS	50	4.2	0.8	2.1	2.9	2.1	2.9
WB- RLARS		5	0	5	0	5	0
B- RLARS	70	5	0	2.1	2.9	2.1	2.9
WBP- RLARS		4.1	0.9	4.1	0.9	4.1	0.9
<i>n > p = 28 , ρ = 0.5 , , α = 0.10</i>							
B- RLARS	50	4.1	0.9	1.1	3.9	2.1	2.9
WB- RLARS		5	0	4.1	0.9	5	0
B- RLARS	70	4.1	0.9	1.1	3.9	2.1	2.9
WBP- RLARS		4.1	0.9	4.1	0.9	4.1	0.9
<i>n > p = 28 , ρ = 0.5, , α = 0.15</i>							
B- RLARS	50	3.1	1.9	1	4	2.1	2.9
WB- RLARS		5	0	5	0	5	0
B- RLARS	70	4.1	0.9	1	4	2.9	2.1
WBP- RLARS		4.1	0.9	4.1	0.9	4.1	0.9

يلاحظ أن الجزء الأول من الجدول (1) عندما يكون $n > p = 28$, $\rho = 0.5$, $\alpha = 0.05$ على طريقة WBP-RLARS في الاختيار الصحيح للمتغيرات الخمسة المهمة في مختلف أمكنة التلوث سواء في القيم الشاذة أم النقاط الرافعه أو كليهما مجتمعة عند حجم عينة 50 ، وينطبق الحال نفسه عند حجم عينة 70 باستثناء تفوق بسيط لطريقة B-RLARS على طريقة WB-RLARS في مكان التلوث في القيم الشاذة لـ Y . أما الجزء الثاني من الجدول (1) عندما يكون $n > p = 28$, $\rho = 0.5$, $\alpha = 0.10$ تفوق الطريقة المقترنة WBP-RLARS على طريقة B-RLARS في الاختيار الصحيح للمتغيرات الخمسة في مختلف أمكنة التلوث سواء في القيم الشاذة أم النقاط الرافعه أو كليهما مجتمعة عند حجم عينة 50 ، وينطبق الحال نفسه عند حجم عينة 70 علمًا أن بارتفاع نسبة التلوث بالقيم الشاذة للمتغير الاستجابة Y حصرًا، وازيد حجم العينة تتساوى كفاءة الطرفيتين مع بعض. وبعبارة أخرى كلما زاد حجم العينة مع زيادة نسبة التلوث في القيم الشاذة تقارب أو تساوى كفاءة الطرفيتين في هذه الحالة فقط. وقد تم التأكيد من ذلك مع أحجام عينات أكبر لم تذكر هنا في هذا البحث. أما الجزء الثالث من الجدول (1) عندما يكون $n > p = 28$, $\rho = 0.5$, $\alpha = 0.15$ ينطبق ماثم التوصل إليه في تحليل نتائج الجزء الثاني من الجدول (1) .

الحالة الثانية : (n(20, 26) < p = 28)

الجدول(2) يعرض نتائج الاختيار الصحيح و الاختيار الخاطئ Selection Right Selection و لطريقي الاختيار ولجمي عينة (20، 26) بحسب نسبة تلوث 0.15 لأول خمسة وسبعة متغيرات .

Outliers & leverage points						
		قيم شاذة ونقاط رافعة في ان واحد			اختيار خمسة متغيرات	
		K=5			K=7	
طريقة الاختيار للمتغيرات	n	اختيار صحيح للمتغيرات	اختيار خاطئ للمتغيرات	اختيار صحيح للمتغيرات	اختيار خاطئ للمتغيرات	اختيار خاطئ للمتغيرات
Selection method		Right Selection	False Selection	Right Selection	False Selection	False Selection
B-RLARS	20	2	3	2		5
WBP-RLARS		0.2	4.8	0.3		6.7

B-RLARS	26	3.9	1.1	3.9	3.1
WBP- RLARS		2.2	2.8	3.1	3.9

يلاحظ من الجدول (2) سوء أداء الطرقين B-RLARS و WBP-RLARS في اختيار أول خمسة متغيرات مهمة K=5 من بين 28 بعد ترتيب أهميتها عند حجم عينة 20 مع أفضلية لطريقة B-RLARS ، ويرتفع أداؤهما عند حجم عينة 26 معبقاء الأفضلية للطريقة نفسها، وبزيادة عدد أول المتغيرات المراد اختيارها من 5 إلى 7 يتحسن أداء الطرقين ويقارباً مع بعض عند حجم العينة 26 ، ولكن أداؤهما يكون سيئاً عند حجم العينة 20 معبقاء الأفضلية لطريقة B-RLARS . وتتجدر الإشارة أن التجربة كانت أكثر من حالة تلوث وكان تحليل النتائج متواافق لما تقدم.

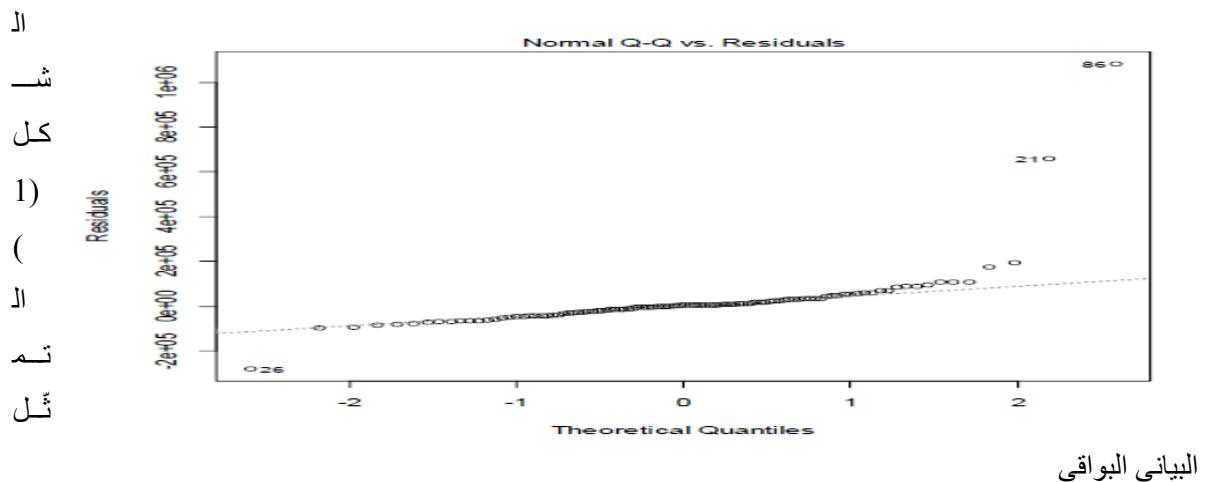
3-1-3 نتائج التطبيق العملي :

تم في هذا المبحث تحليل بيانات عن القيمة السوقية كمتغير استجابة Y لبعض المصارف الأهلية، وعددتها 13 مصرف وهي: (الأهلي، آشور، بابل، بغداد، إيلاف، الاستثمار، العراقي، الخليج، المنصور، سومر، التجاري والمصرف الوطني)، وهي عينة من عدد المصارف الكلي البالغ 48 ، وقد جمعت البيانات من النشرات المالية لسوق العراق للأوراق المالية المنشورة على الموقع الرسمي لسوق الأوراق المالية في العراق ، إذ تم الحصول على 28 نسبة مالية وهي: (معدل دوران السهم (X₁) نسبة العائد على السهم (X₂) ، نسبة الملكية(X₃) ، مكرر الأرباح(X₄) ، نسبة التداول(X₅) ، القيمة الدفترية(X₆) ، سعر الإغلاق السنوي(X₇) ، معدل السعر السنوي(X₈) ، أعلى سعر نقد(X₉) ، أدنى سعر نقد(X₁₀) ، نسبة كفاية راس المال(X₁₁) ، نسبة السيولة(X₁₂) ، نسبة التشغيل(X₁₃) ، نسبة صافي الأرباح إلى رأس المال(X₁₄) ، نسبة صافي الدخل إلى حقوق المساهمين(X₁₅) ، معدل دوران إجمالي الأصول(X₁₆) ، معدل دوران الأصول المتداولة(X₁₇) ، نسبة الإيرادات إلى إجمالي الودائع(X₁₈) ، نسبة المصروفات إلى الإيرادات(X₁₉)،نسبة المصروفات إلى إجمالي الودائع(X₂₀) ، نسبة الودائع إلى (رأس المال + الاحتياطي) (X₂₁) ، نسبة (رأس المال + الاحتياطي) إلى الموجودات(X₂₂) ، مؤشر السيولة (X₂₃) = (صافي رأس المال العامل / الموجودات) * 100 ، مؤشر السيولة(X₂₄) = (الأصول المتداولة / إجمالي الموجودات) * 100 ، مؤشر الرفع (X₂₅) = (حقوق المساهمين / الموجودات) * 100 ، مؤشر الربح قبل الضريبة/ إجمالي الموجودات) * 100 ، مؤشر الرفع (X₂₇) = (إجمالي الموجودات/ إجمالي المطلوبات) * 100 ، مؤشر الرفع (X₂₈) = (حقوق الملكية/ الموجودات الثابتة) * 100) ، وفقاً لذلك يكون أنموذج الانحدار لعينة الدراسة كالتالي :

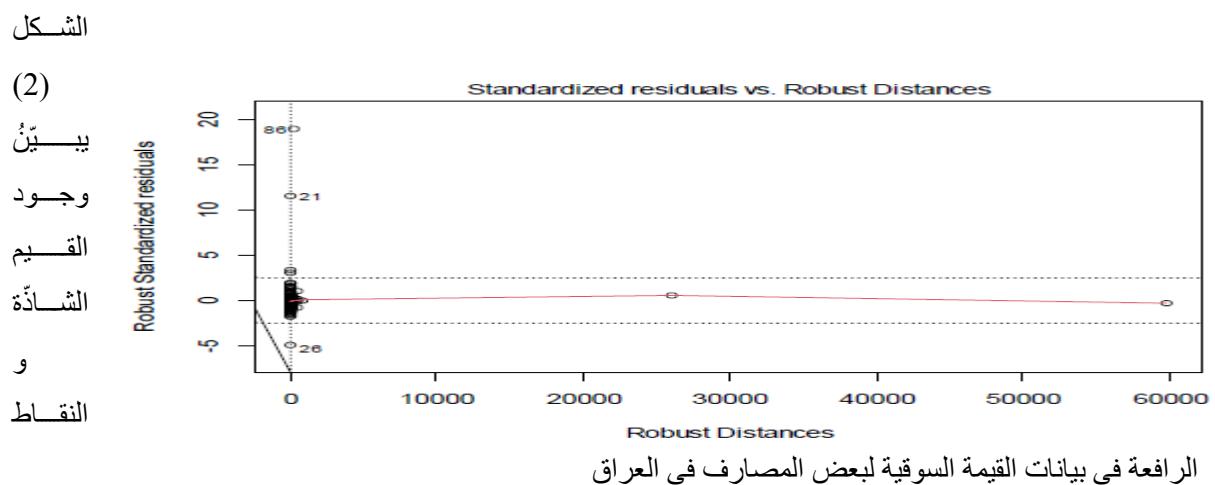
$$Y = \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_{28} X_{28} + \epsilon$$

إذ إن : (ε) يمثل حد الخطأ العشوائي الذي يتوزع توزيعاً طبيعياً.

وبرسم البوافي يتضح من الشكل (1) أن التوزيع الطبيعي لهذه البوافي تم خرقه ببعض القيم الشاذة، أبرزها القيم المفهرسة (21، 86)، الأمر الذي يدفعنا إلى استعمال الطرائق الحصينة بدلاً من الطرائق التقليدية كما موضح في الشكل الآتي:



واعتماداً على المحاكاة تبين أن تأثير النقاط الرافعه كان كبيراً على طريقة B-RLARS ، ولذلك لا بدّ من فحص وجود هذه النقاط في البيانات. والشكل (2) يُظهر العلاقة بين البوافي القياسية ومسافة مهالنوبيز الحصينة، وقد تبين أن هناك قيم شاذة outliers والظاهرة خارج الخطين الأفقيين اللذان يمثلان القيمتين (2.5، -2.5)، كذلك وجود نقطتين رافعتين ظهرتا بعيدة جداً عن نقطة القطع، التي تمثل الخط العمودي إلا أنهما ليستا الوحيدين إذ هناك ثمان نقاط رافعة أخرى تجاوزت خط القطع العمودي هذا. وقد تم استعمال Hat.matrix للكشف عن النقاط الرافعه والتي تبين أن نسبتها (15%) تقربياً وكما موضح بالشكل الآتي :



ووفقاً لحجم العينة الكلي البالغ 104 وبوجود 28 متغيراً توضيحياً، وأيضاً نتائج المحاكاة تم استعمال طريقة WBP-RLARS؛ لأنها الطريقة الأمثل في عملية اختيار المتغيرات عندما يكون حجم العينة أكبر من عدد المتغيرات التوضيحية وقد تم إجراء الآتي :

-1 تحويل المتغيرات إلى متغيرات قياسية حصينة

$$X_i = \frac{x_i - \text{med}(x_i)}{\text{mad}(x_i)} \quad Y_i = \frac{y_i - \text{med}(y_i)}{\text{mad}(y_i)}, \quad i = 1, 2, \dots, n.$$

-2 اختيار أول K من المتغيرات في تسلسل الأهمية الناتج من الطريقة WBP-RLARS إذ إن

. K=7 و K=5

-3 استخلاص مصفوفة فرعية $X_{(n \times k)}$ من مصفوفة المتغيرات التوضيحية $(X_{(n \times p)})$ ، إذ إن $k < p$.

-4 تقدير المعلمات الحصينة باستعمال طريقة MM.

-5 إيجاد قيمة p-value لـ كل معلمة مقدرة بهدف اختبار الفرضيات الآتية :

$$H_0: \beta_j = 0 \quad v.s \quad H_1: \beta_j \neq 0, \quad j = 1, 2, \dots, K$$

وقد كانت النتائج كالتالي :

الجدول (3) يوضح اختيار الطريقة WBP-RLARS للمتغيرات التوضيحية المؤثرة في القيمة السوقية بعض المصادر الأهلية في العراق عندما K=7 و K=5

n=104 , p=28 , K=5				
المتغير	$\hat{\beta}$	Se($\hat{\beta}$)	t-value	p-value
الحد الثابت	0.21	0.07	1.51	0.76
X_9	0.310	0.04	8.04	2.0E-12***
X_{28}	0.233	0.13	1.18	0.073
X_{24}	-0.231	0.08	-2.77	0.007**
X_7	0.004	0.03	0.20	0.85
X_{21}	-0.145	0.09	-1.60	0.11
n=104 , p=28 , K=7				
الحد الثابت	0.29	0.09	3.31	0.001**
X_9	0.310	0.04	8.17	2.0E-12***

X ₂₈	0.249	0.08	3.24	0.002 **
X ₂₄	-0.184	0.05	-3.52	0.0007***
X ₇	-0.135	0.07	-1.83	0.0700
X ₂₁	-0.231	0.09	-2.66	0.009 **
X ₂₇	0.037	0.04	0.78	0.435
X ₁₉	-0.014	0.01	-1.54	0.128

ويلاحظ من الجدول (3) أنَّ الطريقة اختارت أول خمس مُتغيرات هي { X₂₁, X₂₈, X₂₄, X₇, X₂₁ } من حيث الأهمية عندما K=5 وتمثل أعلى سعر نقد(X₉), مؤشر الرافعه المالية (X₂₈), مؤشر السيولة (X₂₄), سعر الإغلاق السنوي(X₇), نسبة الودائع إلى (رأس المال + الاحتياطي) (X₂₁) وفقاً للنموذج النظري الآتي:

$$Y = \beta_0 + \beta_9 X_9 + \beta_{28} X_{28} + \beta_{24} X_{24} + \beta_7 X_7 + \beta_{21} X_{21} + \epsilon$$

وقد تمَّ تقدير المعلمات الحصينة للنموذج المذكور آنفًا باستعمال طريقة MM وكما موضح في الجدول ويتبَّع من اختبار t وقيمة p-value رفض الفرضية الصفرية فيما يخص المعلمات β_9 و β_{24}

وعند توسيع عدد المُتغيرات المراد اختيارها، أي جعل K=7 يثُم إضافة مُتغيرين آخرين بحسب أهميَّتهما، ومن ثُمَّ تكون المُتغيرات المختارة هي: { X₉, X₂₈, X₂₄, X₇, X₂₁, X₂₇, X₁₉ }, أي إضافة نسبة المصرفوفات إلى الإيرادات(X₁₉), مؤشر الرافعه المالية (X₂₇) الذي يساوي (إجمالي الموجودات/ إجمالي المطلوبات)*100. ليكون النموذج النظري الثاني كالتالي:

$$Y = \beta_0 + \beta_9 X_9 + \beta_{28} X_{28} + \beta_{24} X_{24} + \beta_7 X_7 + \beta_{21} X_{21} + \beta_{27} X_{27} + \beta_{19} X_{19} + \epsilon$$

وقد أظهرت نتائج الاختبار وقيم p-value رفض الفرضية الصفرية للمعلمات β_9 و β_{28} و β_{24} و β_{21} و β_{27} وبلاحظ أنه على الرغم من دخول المُتغيرين X₁₉ و X₂₇ إلى الأنماذج لكنهما لم يؤثرا على مُتغير الاستجابة للقيمة السوقية للمصارف في عينة البحث، ولكن دخولهما في الأنماذج ساعد في الكشف عن مُتغيرات مهمة أخرى هي X₂₈ و X₂₁ لتكون المعادلة التقديرية النهائية التي توضح العلاقة بين مُتغير الاستجابة للقيمة السوقية وبعض المُتغيرات التوضيحية المهمة وفق الآتي :

$$\hat{Y} = 0.29 + 0.31 X_9 + 0.249 X_{28} - 0.184 X_{24} - 0.231 X_{21}$$

ولتأكيد نتائج تجربة المحاكاة عندما n = 28 < p تمَّ اختزال المدة الزمنية للمصارف - عينة البحث - إلى سنتين فقط، بحيث يصبح حجم العينة 26 مشاهدة، وبنطبيق الطريقتين لاختيار المُتغيرات عند اختيار كان اختيار طريقة WBP-RLARS للمُتغيرات { X₂₀, X₄, X₁₈, X₁₁, X₁₇, X₆, X₂₂ }, وهي مُتغيرات

تختلف تماماً عن التي اختارتها عندما كان عدد المشاهدات أكبر من عدد المُتغيرات في التجارب السابقة. أما طريقة B-RLARS فقد اختارت $\{X_{13}, X_1, X_{26}, X_{21}, X_6, X_{11}, X_7\}$ وقد اشتركت مع اختيارات WBP-RLARS في المُتغيرين (X_7, X_{21}) عندما $p = n > 28$ ، وهذا الاختيار يؤكّد نتائج المحاكاة وأفضلية طريقة WBP-RLARS على طريقة B-RLARS عندما تكون هناك مشكلة الأبعاد العالية high outliers، أي $n < p$ وجود القيم الشاذة dimensions.

4- الاستنتاجات والتوصيات :

1-4 الاستنتاجات :

بالاعتماد على نتائج تجربة المحاكاة والتطبيق العملي يمكن استنتاج الآتي:

- نستنتج أن لحجم العينة وعدد المُتغيرات التوضيحية تأثيره على بيان أفضلية طريقة اختيار المُتغيرات التوضيحية المهمة والمؤثرة فعلاً على مُتغير الاستجابة، فضلاً عن عدد المُتغيرات المراد اختيارها بحسب أهميتها إذ أبرزت نتائج الطريقتين عند $K=7$ ضماناً لاختيار المُتغيرات المهمة أكثر $K=5$ سواء عندما يكون حجم العينة أكبر أم أصغر من عدد المُتغيرات التوضيحية في الأنماذج فنجد أن طريقة- WBP- RLARS أثبتت كفاءتها على طريقة B-RLARS باختيار المُتغيرات المهمة في العينات من حجم 50.
- تقارب WBP-RLARS مع B-RLARS في كفاءتها مع حجم العينات 70 وتبقى الأفضلية للطريقة WBP-RLARS مع ازدياد أحجام العينات، وهذا الحال هو نفسه عند $K=7$ ، ولكن هذا التقارب يظهر في وجود القيم الشاذة فقط، أي في مُتغير الاستجابة حسراً ولكن الأفضلية تكون لطريقة WBP-RLARS بوجود النقاط الرافعية فقط أو وجود القيم الشاذة والنقاط الرافعية معاً.
- 3- أفضلية طريقة B-RLARS عندما $p < n$ عند حالة اختيار $K=5$ وبحجم عينة 20 على طريقة- WBP- RLARS التي يكون أداؤها في الاختيار سيئاً جداً في أحجام العينات الصغيرة ولكن أداؤها يتحسن كثيراً بتزايد حجم العينة كما حجم عينة 26 ، ويقترب كثيراً من أداء طريقة B-RLARS عند $K=7$.
- 4- اعتماداً على نتائج التطبيق العملي نستنتج أفضلية WBP-RLARS كانت مع حجم العينة $n=104$ الأكبر من عدد المُتغيرات التوضيحية $p=28$ وعدد مُتغيرات مختاراة بحسب الأهمية $K=7$ ولكن هذه الأفضلية تظهر لطريقة B-RLARS عند حجم عينة $n=26$ وعدد مُتغيرات $p=28$ على حساب طريقة- WBP- RLARS التي تعطي اختيارات مختلفة مع اختيارات الطريقة الأخرى التي تشترك معها ببعض الاختيارات، وهذا ما يتواافق مع نتائج المحاكاة وأن المُتغير X_9 والذي يمثل أعلى سعر نقد للمصارف، هو المؤثر الأكثر تأثيراً وب يأتي بالمرتبة الثانية المُتغير X_{28} ، الذي يمثل مؤشر الرفع للمصارف (حقوق الملكية / الموجودات الثابتة) *100 وباتجاه طردي. في حين جاء بالمرتبة الثالثة المُتغير X_{24} مؤشر السيولة: (الأصول المتداولة / إجمالي الموجودات) *100 وبالمرتبة الرابعة X_{21} نسبة (رأس المال + الموجودات) إلى الموجودات وباتجاه عكسي.

2-4 التوصيات :

اعتماداً على الاستنتاجات توصي الباحثة بالآتي:

- 1- اعتماد عدد المتغيرات المهمة من 7-10 حسراً في الأحوال كأها؛ كونه يعطي ضماناً لاختيار المتغيرات المهمة فعلاً.
- 2- استعمال خوارزمية WBP-RLARS في عملية الاختيار عندما $p > n$ ، والقيم الشاذة في متغير الاستجابة بوجود النقاط الرافعه في المتغيرات التوضيحية أو القيم الشاذة والنقط الرافعة معاً في بيانات أنموذج الانحدار الخطّي المتعدد، وبأي حجم عينة بما يتحقق الشرط المذكور آنفًا، وكما ويمكن استعمال طريقة-B-RLARS في حالة وجود القيم الشاذة في متغير الاستجابة فقط.
- 3- استعمال خوارزمية B-RLARS في حالة $p < n$ عند عملية اختيار المتغيرات المهمة ولا سيما عند أحجام العينات المبتدأة من عدد المتغيرات التوضيحية حسراً، ويمكن استعمال طريقة WBP-RLARS أيضاً لاختيار المتغيرات المهمة في الأنماذج عندما يقترب حجم العينة من عدد المتغيرات التوضيحية في الأنماذج، إذ تقترب كفاءة الطريقتين من بعضهما عند هذه الحالة.
- 4- اهتمام المصارف الأهلية في العراق على المحافظة على أعلى سعر نقد، ومؤشر الرفع لحقوق الملكية مقسموماً على الموجودات مرتفعاً؛ لأنّ له التأثير الكبير في رفع القيمة السوقية للمصارف في حين يجب عليها أن تحافظ على مستوى مناسب لمؤشر السيولة فيها، وكذلك المتغير X_{21} الذي يمثل نسبة رأس المال + الاحتياطي) إلى الموجودات، كونهما يعملان على خفض القيمة السوقية للمصارف.

المصادر:-

أولاً : المصادر العربية

- 1- السrai ،علي حميد يوسف (2017) " استعمال المقدرات الحصينة الجزئية لنموذج الانحدار الخطّي في ظلّ وجود مشكلتي الأبعاد والقيم الشاذة مع تطبيق عملي " ، اطروحة دكتوراه ، كلية الإدراة والاقتصاد - جامعة بغداد.

ثانياً : المصادر الأجنبية

- 1- Alamgir, Salahuddin, Ali Amjad (2013) " Split Sample Bootstrap Method" World Applied Sciences Journal 21 (7): 983-993.
- 2- Alfons, A., Ates ,N. Y. ,Groenen P J. F.(2022) " A Robust Bootstrap Test for Mediation Analysis" Organizational Research Methods, journals. saspub. Com/home/otm, Vol. 25(3) 591–617.
- 3- Alfons, A., Baaske, W.E., Filzmoser, P., Mader W., Wieser, R. (2011)" Robust variable selection with application to quality of life research" 20(1), 65–82.

- 4-** Barrera, M. S., Aelst S. V. (2008)" Robust model selection using fast and robust bootstrap" Computational Statistics and Data Analysis 52 , 5121–5135.
- 5-** Blaine, Bruce E. (2018). "Winsorizing." The SAGE Encyclopedia of Educational Research, Measurement, and Evaluation, 1817-1818.
<http://libguides.sjfc.edu/citations>.
- 6-** Buhlmann , Peter. Geer, Sara (2011) " Statistics for High Dimensional Data Methods, Theory and Applications". Springer.
- 7-** Efron, B., Hastie, T., Johnstone I., Tibshirani, R. (2004)" LEAST ANGLE REGRESSION" The Annals of Statistics, vol. 32, No.2, 407-541.
- 8-** Farnè, Matteo, Vouldis Angelos T. (2018)" A methodology for automatised outlier detection in high-dimensional datasets: an application to euro area banks' supervisory data" Working Paper Series, European Central Bank, No 2171.
- 9-** Filzmoser, Peter, Nordhausen, Klaus, (2021) " Robust linear regression for high-dimensional data: An overview " WIREs Comput Stat.;13:e1524.
- 10-** Hettigoda, Sandamala, (2016) "Computation of Least Angle Regression coefficient profiles and LASSO estimates.". Electronic Theses and Dissertations. Paper 2404. <https://doi.org/10.18297/etd/240>.
- 11-** Imon, A. H. M. R, Ali M. M.(2005)" Bootstrapping Regression Residuals" Journal of Korean Data & Information Science Society, Vol. 16, No. 3, pp. 665-682.
- 12-** Januaviani, T.M.A.,GusrianiN., Joebaedi, K., Supian, S., Subiyanto,(2019) " The Best Model of LASSO With The LARS (Least Angle Regression and Shrinkage) Algorithm Using Mallow's Cp" World Scientific News 116 , 245-252 .
- 13-** Khan J. A., Aelst, S. V., Zamar, R. H.(2007) " Robust Linear Model Selection Based on Least Angle Regression " Journal of the American Statistical Association ,Vol. 102, No. 480, 1289-1299.

- 14-** Lindskou, M., Tvedebrink, T., Eriksen, P.S., Morling, N. (2021)" Detecting Outliers in High-Dimensional Data with Mixed Variable Types Using Conditional Gaussian Regression Models"arXiv:2103.02366v3 [math.ST].
- 15-** Mami, A. M., Jaber, A. M., Almabrouk, O. S. (2020)" Applying Bootstrap Robust Regression Method on Data with Outliers" International Journal of Sciences: Basic and Applied Research (IJSBAR) Volume 49, No 1, pp 143-160.
- 16-** Midi. H., Uraibi, H.S., Talib, B.A. (2009)" Dynamic Robust Bootstrap Method Based on LTS Estimators " European Journal of Scientific Research, Vol.32 No.3, pp.277-287.
- 17-** Mohamad, M, Ramli, N. M., Ghani, N. A. M.(2016) " Weighted Split Sample Bootstrap for Regression Models with High Dimensional Data" Indian Journal of Science and Technology, Vol 9(28) .
- 18-** Ramli, Norazan M., Midi Habshah, Imon, A. H. M. R. (2009)" Estimating Regression Coefficients using Weighted Bootstrap with Probability" Wseas Transactions on Mathematics, Issue 7, Volume 8, 362-371.
- 19-** Rana Sohel, Midi, Habshah, Imon, A. H. M. R. (2012) " Robust Wild Bootstrap for Stabilizing the Variance of Parameter Estimates in Heteroscedastic Regression Models in the Presence of Outliers" Mathematical Problems in Engineering Volume 2012, Article ID 730328, 14 pages.
- 20-** Susanti, Yuliana, Pratiwi , Hasih, Sulistijowati, Sri, Liana ,Twenty (2014) " M ESTIMATION, S ESTIMATION, AND MM ESTIMATION IN ROBUST REGRESSION" International Journal of Pure and Applied Mathematics , Volume 91 No. 3, 349-360.