

# مقارنة بعض الطرائق الحصينة لتقدير معلمات انحدار المربعات الصغرى الجزئية\*

أ.م.د. سجي محمد حسين  
م. رباب عبد الرضا صالح  
كلية الادارة والاقتصاد/ جامعة بغداد- قسم الاحصاء

## المستخلص

تعد تقنية تخفيض الابعاد واختيار المتغيرات من المواضيع المهمة في التحليل الاحصائي لنماذج متعدد المتغيرات، فعندما يرتبط إثنان او اكثر من المتغيرات التوضيحية في الإنحدار بعلاقة او عدة علاقات تامة او غير تامة ، تحدث مشكلة التعدد الخطي والتي فيها خرق لأحد الفروض الأساسية لطريقة المربعات الصغرى الاعتيادية مما يؤدي الى تقديرات غير دقيقة .

هناك طرائق عدة اقترحت لمعالجة هذه المشكلة نذكر منها طريقة المربعات الصغرى الجزئية (PLS) والتي تستعمل لتخفيض الأبعاد في تحليل الإنحدار، باستعمال تحويلات خطية تقوم بتحويل مجموعة من المتغيرات المرتبطة ارتباطاً عالياً، الى مجموعة من المتغيرات المستقلة الجديدة وغير المرتبطة تعرف بالمكونات ، وتكون هذه المكونات خطية متعامدة ومستقلة بعضها عن البعض الآخر. ان طريقة PLS تفشل في التعامل مع البيانات التي تتضمن وجود القيم الشاذة ،وعليه فان نجاح هذه الطريقة يتوقف على عدم وجود هذه القيم الشاذة التي لها تأثير غير مرغوب على النتائج، وللمحد من تواجد هذه القيم نلجأ الى استعمال الطرائق الحصينة .

في هذا البحث أستعمل خوارزمية PLSKURSD، والتي تطبق خوارزمية SIMPLS على مصفوفة التباين والتباين المشترك الحصين، فضلاً عن الطريقة المقترحة MPLSKURSD وهي تعديل الى طريقة PLSKURSD .

جرى مقارنة بين معلمات نموذج الانحدار الخطي بطريقة المربعات الصغرى الجزئية مع الطرائق الحصينة للمربعات الصغرى الجزئية من خلال تجارب محاكاة، اعتمدت على وجود أنواع عدة من القيم الشاذة من البيانات وبنسب مختلفة من التلوث ولحجوم عينات وابعاد متغيرات مختلفة.

**المصطلحات الرئيسية للبحث:** المربعات الصغرى الجزئية- مصفوفة التباين المشترك الحصينة- التفرطح- الاسقاطات.



مجلة العلوم

الاقتصادية والإدارية

المجلد 20

العدد ٧5

لسنة ٢٠١٤

الصفحات ٤١٣ - ٤٣١

\*مستل من اطروحة دكتوراه

## ١-١ المقدمة

أن أغلب التقنيات الاحصائية التقليدية تصمم بصورة خاصة لأبعاد قليلة من البيانات عندما يكون عدد المشاهدات أكبر من عدد المتغيرات لكن في بعض المسائل الخاصة بموعد البقاء أو تصنيف الورم لتنبؤ المريض كما في تحليل بيانات التعبير الجيني التي تتضمن اعداد كبيرة من المتغيرات (الجينات) التي ليس لها صلة بالتنبؤ إذ أن ادراج هذه المتغيرات في الانحدار والتصنيف يؤدي الى خسارة في الاداء،<sup>[9]</sup> حتى لو كان اسهام هذه المتغيرات قليلة وغيرها من التطبيقات التي تعتمد على الأبعاد العليا من البيانات فالتطبيقات الاحصائية في مثل هذه الحالات قليلة ومن هذه التطبيقات طريقة المربعات الصغرى الجزئية التي تحلل مثل هكذا البيانات .

تعد طريقة أنحدار المربعات الصغرى الجزئية (pls) partial least square إحدى تطبيقات أنحدار الخطي لحل مشكلة تخفيض الأبعاد ومن ثم التخلص من مشكلة التعدد الخطي، وأول من أستعمل هذه التطبيقات الباحث بعلم الاقتصاد Herman Wold عام (١٩٦٦) ثم طورت من قبل العالم نفس عام(1975) وقد تم تطبيق خوارزمية (NIPALS (Non-linear Iterative partial least squares (PLS1, PLS2) لمعالجة سلسلة من المصفوفات والتطبيقات في الاقتصاد القياسي ذات الأبعاد العليا لوحد أو أكثر من متغيرات الاستجابة فهي طريقة تتصل بالمصفوفة  $x$  والمتجه  $y$  او المصفوفة  $Y$  . وقد طورت هذه الطريقة واصبحت مشهورة ليس في مجال البحوث القياسية وإنما في بقية الحقول وخصوصا في مجال علوم الكيمياء (chemometrics) ، وبعدها استخدمت في مجالات الطب والفيزياء والعلوم الاجتماعية ، علوم صناعة الاغذية وغيرها من العلوم الأخرى<sup>[13]</sup> .

وقد قدم De Jong عام ١٩٩٣ خوارزمية بديلة عن هذه الخوارزمية وهي statistically inspired modification of the partial least squares (SIMPLS) وهي خوارزمية سريعة وكفوءة<sup>[8,3]</sup> بالمقارنة مع خوارزمية NIPALS (pls1, pls2) ، إذ انها تعطي نتائج pls1 نفسها في حالة متغير الاستجابة واحد ونتائجها تختلف عن pls2 في حالة وجود عدة متغيرات استجابة. وتوالت الدراسات ضمن طريقة المربعات الصغرى الجزئية فقد قدم الباحث Andersson بحثاً عن استخدام خوارزميات عدة الى المربعات الصغرى الجزئية فضلا عن الخوارزميات المذكورة انفاً، نذكر منها خوارزمية KERNAL المنسوبة الى Dayal ، وخوارزمية PLS-F المنسوبة الى Manne<sup>[1]</sup> وغيرها من الخوارزميات. ان طريقة المربعات الصغرى الجزئية PLS حساسة لوجود القيم الشاذة في مجموعة البيانات والتي تؤدي الى عدم دقة النتائج ولمعالجة القيم الشاذة ، تم استخدام الطرائق الحصينة للمربعات الصغرى الجزئية .

وقد كان اول اقتراح للطرق الحصينة للمربعات الصغرى الجزئية من قبل العالم Wakelling and Macfie عام ١٩٩٢ حيث تضمن تطبيقات الانحدار الحصين الى المربعات الصغرى الجزئية باستعمال خوارزمية biweight وهي طريقة تكرارية تعتمد على إعادة الأوزان لكل مشاهدة بالاعتماد على حجم بواقي الانحدار لطريقة المربعات الصغرى الجزئية<sup>[16]</sup> وتوالت الدراسات حول حصانة المربعات الصغرى الجزئية فقد قدم الباحث ( Gonzales Javier. & et al عام ٢٠٠٩ ) أسلوباً جديداً لإزالة تأثير القيم الشاذة من نقاط البيانات لطريقة المربعات الصغرى الجزئية وتقدير حصين لمصفوفة التباين المشترك يتم حسابه بواسطة البحث عن الشواذ لاسقاطات المتغير الاحادي للبيانات يجمع بين الاتجاه العشوائي الى (Stahel \_ Donoho) والاتجاه المحدد الحاصل من تعظيم وتصغير معامل التفرطح او (التفلطح) لـ (Pena and Prieto) وقد اوضحت هذه الطريقة كفاءتها واعطائها افضل النتائج من الطرائق الأخرى وقد تم استعمال بيانات حقيقية اضافة الى اسلوب المحاكاة عند ظهور عدة انواع من الشواذ في البيانات وكذلك تمت المقارنة بين التقديرات تحت النموذج الذي يتوزع توزيع طبيعي وتوزيعات مختلفه من الأخطاء .



## الجزئية \*

## ٢-١ هدف البحث

ان الهدف من البحث هو الحصول على افضل تقدير لمعاملات الانحدار في حالة وجود التعدد الخطي بين المتغيرات التوضيحية فضلا عن وجود القيم الشاذة في البيانات باستعمال طريقة المربعات الصغرى الجزئية PLS بخوارزمية SIMPLS ومقارنتها مع الطرائق الحصينة باستعمال خوارزميتي PLSKURSD , MPLSKURSD بتجارب محاكاة لعدة انواع من الشواذ وبنسب وحجوم عينات وابعاد مختلفة .

## ١-٢ الجانب النظري

ان طريقة المربعات الصغرى الجزئية PLS تعد اداة لاسلوب الانحدار الخطي وضعت لربط العديد من الانحدارات لواحد او اكثر من متغيرات الاستجابة بواسطة متغيرات كامنة Latent variable عندما يكون للمتغيرات التوضيحية ارتباطاً عالياً وكذلك عندما يكون عدد المتغيرات التوضيحية تفوق عدد المشاهدات. وسيتم الاعتماد في البحث على خوارزمية SIMPLS حيث ان العوامل لـ PLS باستعمال خوارزمية SIMPLS تحدد باستعمال معيار تعظيم مصفوفة التباين المشترك.

## ٢-٢ انحدار المربعات الصغرى الجزئية Partial Least Squares Regression

في هذا البحث يتم الفرض بان المصفوفة  $X_{n,p}$  لها  $n$  من المشاهدات و  $p$  من المتغيرات ويرمز لها

$$X_{n,p} = (X_1, X_2, \dots, X_n)' \text{ حيث } X_i = (X_{i1}, X_{i2}, \dots, X_{ip}) \text{ الى } i^{th} \text{ من المشاهدات.}$$

والمصفوفة  $Y$  تتكون من  $n$  من المشاهدات و  $q$  من متغيرات الاستجابة ونرمز لها بالرمز

$$Y_{n,q} = (Y_1, Y_2, \dots, Y_n)' \text{ اما مجموعة البيانات المدمجة } (X_{n,p}, Y_{n,q}) \text{ تشير اليها بالرمز } Z_{n,m}$$

$$[8] \text{ حيث } m = p + q$$

فنموذج الانحدار الخطي :

$$y_i = \beta_0 + \beta'_{p,q} x_i + e_i \quad \dots (2-1)$$

$e_i$  : حد الخطأ الذي يشترط فيه:

$$\text{cov}(e_i) = \Sigma_e, E(e_i) = 0 \text{ لحجم } q$$

$$\beta_0 = (\beta_{01}, \dots, \beta_{0q}) \text{ , حد ثابت غير معلوم ببعد } q$$

$$\beta_{pq} \text{ تمثل المعالم غير المعلومة وهي مصفوفة الميل ببعد } pxq$$

## ٢-٢ خوارزمية SIMPLS The SIMPLS algorithm [14,8,3,2]

ان خوارزمية SIMPLS تفرض المتغيرات  $y_i$  و  $x_i$  حيث  $i=1,2,\dots,n$  لها علاقة من خلال الانموذج التالي الاتي:

$$x_i = \bar{x} + P_{p,k} \tilde{t}_i + g_i \quad \dots (2-2)$$

$$y_i = \bar{y} + A'_{q,k} \tilde{t}_i + f_i \quad \dots (2-3)$$

## الجزئية \*

$\bar{x}, \bar{y}$ : الوسط الحسابي للمتغيرات  $x_i$  و  $y_i$

$\tilde{t}_i$ : القياسات او المركبات او النقاط (scores or component) ببعد  $k$  حيث  $k \ll p$

$p_{p,k}$ : مصفوفة التحميل (x-loading)

$A'_{k,p}$ : مصفوفة الميل في انحدار  $y_i$  على  $\tilde{t}_i$

$g_i, f_i$ : البواقي.

تفترض الخوارزمية اولا تشكيل او بناء المركبات حيث يتم الحصول على  $k$  من المركبات التي تكون كتوافيق خطية من  $x$  من المتغيرات والتي لها اكبّر تباين مشترك مع التوافق الخطية من المتغيرات  $y$  وبصورة ادق نفرض  $\tilde{Y}_{n,q}, \tilde{X}_{n,p}$  تشير الى مصفوفة البيانات الممركزة (mean-centered data matrices) ان:-

$$\tilde{x}_i = x_i - \bar{x}$$

$$\tilde{y}_i = y_i - \bar{y} \quad \dots (2-4)$$

$q_a, r_a$ : متجهات اوزان pls الطبيعية (normalized) حيث  $(\|r_a\| = \|q_a\| = 1)$  تعرف على انها المتجهات التي تعظم لكل  $a$  حيث  $a = 1 \dots k$ .

$$\text{cov}(\tilde{Y}_{n,q} q_a, \tilde{X}_{n,p} r_a) = q'_a \frac{\tilde{Y}'_{q,n} \tilde{X}_{n,p}}{n-1} r_a = q'_a S_{y,x} r_a \quad \dots (2-5)$$

حيث ان  $S'_{yx} = S_{xy}$  هي مصفوفة التباين المشترك التجريبية empirical cross-covariance matrix بين المتغيرات  $X, Y$  وعناصر الـ  $\tilde{t}_i$  تعرف على انها التوافق الخطية للبيانات الممركزة (mean-centered data) اوتكافئ

$$\tilde{t}_{ia} = \tilde{X}'_i r_a$$

$$\tilde{T}_{n,k} = \tilde{X}_{n,p} R_{p,k} \quad \dots (2-6)$$

$$R_{p,k} = (r_1 \dots r_k)$$

حيث

ان تعظيم المعادلة (2-5) لها حل واحد ان اول زوج لمتجه الاوزان لـ  $(r_a, q_a)$  يتم الحصول عليه من يسار ويمين المتجه المفرد المميز الى  $S_{yx}$ . وهذا يعني بانته  $r_1$  المتجه المميز الى  $S_{yx} * S_{yx}$  حيث ابعادها  $(p \times p)$  و  $q_1$  هو المتجه المميز الى  $S_{yx} * S_{xy}$  حيث ابعادها  $(q \times q)$ . وللحصول على

حل أكثر من واحد، المركبات  $\tilde{T}_{n,k}$  تتطلب أن تكون غير مرتبطة orthogonal:

$$T'_a T_j = 0 \quad \forall j \neq a$$



## الجزئية \*

القيود المذكور انفا يفرض توليد سلسلة من الحلول المختلفة للمعادلة (٥-٢) هذا فضلا عن تجنب التعدد الخطي بين الانحدارات في الخطوة الثانية من الخوارزمية ، وبوجود هذا القيد اوزان المتجهات لخوارزمية *Simpls* الباقية وهي  $r_a, q_a$  ،  $(2 \leq a \leq k)$  يتم الحصول عليها بواسطة المتجهات المميزة  $S_{xy}^a S_{yx}^a$  و  $S_{xy}^a S_{yx}^a$  وبشكل متتالي حيث ان مصفوفة التباين -التباين المشترك الجديدة (deflated) يتم الحصول عليها من المعادلة الاتية :

$$S_{x,y}^a = (I_p - V_{a-1} V_{a-1}') S_{x,y}^{a-1} \quad \dots (2-7)$$

تبدأ هذه الخوارزمية مع  $S_{xy}^1 = S_{xy}$  وتكرر العملية حتى نحصل على  $k$  من المركبات والتي يتم بعدها تحديد عدد المركبات  $k$  بحسب طريقة *cross validation*

وتمثل  $(V_1, \dots, V_{a-1})$  بالمتعامد الطبيعي *orthonormal base* الى مصفوفة التحميل  $X, \dots, P_{a-1}$  حيث ان مصفوفة التحميل  $P_j$  تصف العلاقة الخطية بين المتغيرات  $X$  و  $j^{th}$  من المركبات وكالاتي :

$$P_j = \tilde{X}' T_j / T_j' T_j \quad \dots (2-8)$$

وبالتعويض عن  $T_j$  نحصل على الاتي :-

$$P_j = \tilde{X}' \tilde{X} r_j / (r_j' \tilde{X}' \tilde{X} r_j) = S_x r_j / (r_j' S_x r_j) \quad \dots (2-9)$$

والتي تمثل معامل انحدار المربعات الصغرى في انحدار  $\tilde{X}$  على المركبة  $T_j$

حيث:  $S_x$  هي مصفوفة التباين التجريبية للانحدارات  
الخطوة الثانية في الخوارزمية هي انجاز الانحدار المتعدد الخطي (MLR) لانحدار المركبات المستخلصة  
على متغيرات  $y$  الاصلية حيث صيغة انموذج الانحدار كالاتي:

$$y_i = \alpha_0 + A'_{q,k} \tilde{t}_i + f_i \quad \dots (2-10)$$

حيث  $E(f_i) = 0$  and  $cov(f_i) = \sum_f$

تقدير الانحدار الخطي المتعدد MLR

$$\hat{\alpha}_0 = \bar{y} - \hat{A}'_{q,k} \bar{\tilde{t}} \quad \dots (2-12)$$

$$\hat{A}_{k,q} = (S_t)^{-1} S_{ty} = (R'_{k,p} S_x R_{p,k})^{-1} R'_{k,p} S_{xy} \quad \dots (2-11)$$

$$S_f = S_y - \hat{A}'_{q,k} S_t \hat{A}_{k,q} \quad \dots (2-13)$$

حيث  $S_y$  &  $S_t$ : مصفوفة التباين الاولية للمتغيرات  $y$  &  $t$

**الجزئية \***

نلاحظ ان الانحدار الخطي المتعدد (multiple linear regression MLR) يشير الى انحدار المربعات الصغرى الكلاسيكية لمتغيرات  $x$  المتعددة، وعندما عدد المتغيرات المعتمدة أكبر من الواحد لمتغيرات  $y$  المتعدد يعرف بالانحدار الخطي لمتعدد المتغيرات (multivariate multiple linear regression

وبسبب  $\bar{t} = 0$  الحد الثابت  $\alpha_0$  يقدر بواسطة  $\bar{y}$

وعند تعويض  $\tilde{t}_i = R'_{k,p}(x_i - \bar{x})$  في المعادلة (2-3) نحصل على تقديرات المعامل لنموذج الانحدار الاصيلي في المعادلة (2-1) وكالاتي :

$$\hat{\beta} = R_{p,k} (R'_{K,p} S_x R_{p,K})^{-1} R'_{K,p} S_{xy} \quad \dots(2-14)$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}'_{q,p} \bar{x} \quad \dots(2-15)$$

واخيرا التقدير الى  $\Sigma_e$  هو

$$S_e = S_y - \hat{\beta} S_x \hat{\beta}' \quad \dots(2-16)$$

في حالة الاستجابة واحدة ( $q=1$ ) التقدير للمعالم يكون  $\hat{\beta}_{p,1}$  وتكتب على شكل متجه، بينما تباين الخطأ

$$S_e \text{ يتحقق كالاتي } \hat{\sigma}_e^2 = S_e^2$$

**٤-2 انحدار المربعات الصغرى الجزئية الحصينة Robust Partial Least Squares Regression**

على الرغم من ان المربعات الصغرى الجزئية (PLS) تعالج مشكلة تقليل الابعاد في متعدد المتغيرات الا انها تفشل في التعامل مع بيانات متعدد المتغيرات والتي تتضمن وجود القيم الشاذة التي تكون بعيدة عن تمركز البيانات او تكون غير منسجمة او مختلفة عن بقية البيانات وتحدث نتيجة اخطاء في القياس او التسجيل او تنتمي بعض المشاهدات الى مجتمع اخر ومن ثم ينتج تقدير غير دقيق ومن المهم جدا اكتشاف القيم الشاذة لان وجودها لها تأثير غير مرغوب به على التقديرات ولازالة هذه التأثيرات الغير مرغوب بها فقد افترض الباحثون استبدال الطرائق الاعتيادية بالطرائق الحصينة وهناك عدة طرائق لحصانة المربعات الصغرى الجزئية منها الطرائق التي تعتمد على حصانة مصفوفة التباين والتباين المشترك والطرائق التي تعتمد على الانحدار الحصين [25].

[12,11,10,5]

**5--2 خوارزمية PLSKURSD The PLSKURSD algorithm**

تعتمد خوارزمية PLSKURSD على تطبيق خوارزمية SIMPLS على مصفوفة التباين والتباين المشترك الحصين، ان تقدير مصفوفة التباين يتم بواسطة اسقاط البيانات في بعض الاتجاهات وايجاد القيم الشاذة على هذه الاتجاهات وحذفها من العينة واستعمال البيانات النظيفة لحساب مصفوفة التباين والتباين المشترك.

اما الية عمل هذه الطريقة فانها تكون على ثلاث خطوات، حيث تفترض وبدون فقدان العمومية ان البيانات الاصلية لها وسط حسابي صفر ومصفوفة التباين  $S_z$  حيث تحول البيانات باستعمال المعادلة الآتية لمصفوفة البيانات  $Z=[X,Y]$  :-

$$\tilde{Z}_i = S_z^{-1/2} Z_i, \quad i = 1, 2, \dots, n \quad \dots(2-17)$$

وعندما تكون مصفوفة التباين مفردة (singular) نستخدم المعكوس العام (generalized inverse)

وبمكن تطبيق هذه الخوارزمية في حالة عدد المتغيرات أكبر من المشاهدات باستعمال خوارزمية svd [7] اما خطوات هذه الطريقة فهي كالاتي:

**الخطوة الاولى STEP ONE**



## الجزئية \*

يتم حساب الاتجاه الذي يعظم معامل التفلطح (Kurtosis) للاسقاطات وكذلك يتم حساب المسافات (  $r_i^{(j)}$  ) حيث  $j=1,2$  لكلا الاتجاهين (تعظيم وتصغير) للمتغير الاحادي الطبيعية للبيانات اول اتجاه نحصل عليه من حل الصيغة التالية

$$D1 = \arg \max 1/n \sum_{i=1}^n (d' \tilde{Z}_i)^4 \quad \dots (2-18)$$

$$s.t \quad d' d = 1$$

ونفس العملية تكرر لحساب الاتجاه الذي يصغر معامل التفرطح

$$D2 = \arg \min 1/n \sum_{i=1}^n (d' \tilde{Z}_i)^4 \quad \dots (2-19)$$

$$s.t \quad d' d = 1$$

نفرض ان

$$P_i^{(j)} = d'_j \tilde{z}_i \quad i=1,2,\dots,n \quad j=1,2 \quad \dots (2-20)$$

الذي يمثل اسقاطات القيم على هذين الاتجاهين و تم تحسب قياس الشواذ للمتغير الاحادي او المسافات  $r_i^{(j)}$  للمتغير الاحادي الطبيعية كالآتي

$$r_i^j = 1/\beta_p \frac{|P_i^{(j)} - \text{median}_i(P_i^{(j)})|}{MAD_i(P_i^{(j)})} \quad j=1,2 \quad \dots (2-21)$$

حيث ان

$$MAD_i(p_i^{(j)}) = \text{median}_i |p_i^{(j)} - \text{median}_i(p_i^{(j)})| \quad \dots (2-22)$$

حيث ان MAD : وسيط الانحراف المطلق Median Absolute Deviation و ان قيمة القطع  $\beta_p$  تعتمد على عدد المتغيرات (p) ويتم الحصول عليها من تجارب محاكاة لضمان في حالة غياب الشواذ بان نسبة المشاهدات الصحيحة التي صنفت بانها شواذ تقريبا هي 0.05 .

## الخطوة الثانية STEP TWO

حساب الاتجاهات العشوائية أيجاد الاتجاه العشوائي وذلك من خلال اسلوب المعاينة التطبيقية ومن ثم البحث عن القيم الشاذة في هذه الاتجاهات العشوائية كل اتجاه يولد بخطوتين وهي

1- يتم اختيار مشاهدتين عشوانيا من العينة ثم تحسب الاتجاه المعرف لهاتين المشاهدين ( اما ان يكون على اساس المسافة او الفرق بين المشاهدين ) ومن ثم اسقاط المشاهدات على هذا الاتجاه وتكرر الى  $h$  من المرات حيث  $p * 10, 2, \dots, h$  حيث  $p$  هو ( عدد المتغيرات ) .

٢- في هذه الخطوة يتم بناء مجموعة من  $k$  من العينات التطبيقية حيث يتم ترتيب الاسقاطات وتقسيم الى  $k$  من الفترات حيث  $k$  تؤخذ اما 3 او 5 تحدد مسبقا وكل فترة سيكون حجمها  $n/k$  لكل  $k$  من الفترات،  $1 < k < K$ ، سيتم اختيار عينة جزئية تتكون من  $p$  من المشاهدات تنتخب وبدون الارجاع ويحسب الاتجاه  $d_j$  الذي يستخدم لحساب الشواذ كما في الخطوة الاولى حيث ان الاسقاطات  $\tilde{p}_i^{(j)}$  هي:-

$$\tilde{P}_i^{(j)} = \tilde{d}_j \tilde{Z}$$

ومن خلالها سوف تحسب المسافات  $\tilde{r}_i$  للمتغير الاحادي الطبيعية :

$$\tilde{r}_i^j = 1 / \beta_p \frac{|\tilde{P}_i^{(j)} - \text{median}_i(\tilde{P}_i^{(j)})|}{MAD_i(\tilde{P}_i^{(j)})} \quad j = 1, 2, \dots, h * K \quad \dots(2-23)$$

## الخطوة الثالثة (التحقيق) STEP THREE (checked)

لكل مشاهدة (i) فان قياس الشواذ الطبيعي لها يتم الحصول عليه كما يلي

$$r_i = \max \{ r_i^1, r_i^2, \tilde{r}_1, \dots, \tilde{r}_i^j \} \quad \dots(2-24)$$

$$j = 1, 2, \dots, h * K$$

والتي تم الحصول عليها بموجب المعادلات (2-21) ، (2-23) ، هذه المشاهدات اذا كانت  $r_i > 1$  فانها تصنف على انها قيمة شاذة وتزال من العينة اذا كان عددهم اصغر من  $[n - [n + p + 1] / 2]$  عدا ذلك فان المشاهدات التي عددها  $[n - [n + p + 1] / 2]$  والتي لها اكبر قيم الى  $r_i$  فهي تعتبر شواذ وتزال من البيانات.

واخيرا نجد المجموعة  $U$  والتي تمثل المشاهدات التي لا تتضمن القيم الشاذة بعدها يتم حساب مسافة مهالنوبس بالاعتماد على المشاهدات الجيدة وكما يلي:-

$$\tilde{m} = \frac{1}{|U|} \sum_{i \in U} Z_i \quad \dots(2-25)$$

$$\tilde{S}_z = \frac{1}{|U| - 1} \sum_{i \in U} (Z_i - \tilde{m})(Z_i - \tilde{m})' \quad \dots(2-26)$$

$$V_i = (Z_i - \tilde{m})' \tilde{S}_z^{-1} (Z_i - \tilde{m}) \quad \dots(2-27)$$

حيث ان  $|U|$  تمثل عدد المشاهدات في المجموعة  $U$



**الجزئية \***

فإذا كانت  $V_i < \chi_{p,1-99}^2$  فإنها لا تصنف على أنها قيم شاذة وتوضع في مجموعة U وفي حالة

$\tilde{S}_z^{-1}$  غير موجودة سوف نستعمل المعكوس العام لها.

وهذه الخطوات تكرر حتى لا تبقى قيم شاذة (او تصبح U مجموعة كل البيانات) .

ويمكن ان تستعمل هذه الخوارزمية عندما  $p > n$  حيث يتم تطبيق تجزئة القيم المفردة (SVD) على

البيانات لاجل تقليص مجال البيانات الى مجال جزئي مقسم بواسطة  $n$  [28,19] ولا يوجد معيار قيمة القطع للبيانات التي لا تتوزع توزيعاً طبيعياً وإنما يتم استعمال معيار بديل لقيم القطع

(الحصينة) تكون بالاعتماد على الجذر التربيعي لمسافة مهالنوبس  $\sqrt{\tilde{V}_i}$  الذي يكون توزيعاً مماثل أكثر من

$\tilde{V}$  وان المشاهدة تعتبر شاذة اذا تحقق مايلي:-

$$\sqrt{\tilde{V}_i} \geq \text{median}_i(\sqrt{\tilde{V}_i}) + 4.5 \text{MAD}_i(\sqrt{\tilde{V}_i}) \quad \dots(2-28)$$

وان  $\text{MAD}_i$  عرف في المعادلة (2-22)

**٦-٢ الطريقة المقترحة THE PROPOSED ALGORITHM MPLSKURSD**

يتم في هذه الخوارزمية ايجاد مصفوفة التباين والتباين المشترك الحصينة والتي توظف في خوارزمية SIMPLS لإيجاد مقدرات المعلمات في حالة وجود القيم الشاذة فضلاً عن وجود مشكلة التعدد الخطي، أذ تعتبر هذه الطريقة كتعديل لطريقة المربعات الصغرى الجزئية الحصينة لخوارزمية PLSKURSD حيث ان هذه الطريقة تعتمد على اربع خطوات وكالاتي :

١- في الخطوة الاولى يتم ايجاد قيم  $\tilde{X}$  الحصينه وذلك أولاً بحساب تقدير الوسيط متعدد الابعاد Column-wise median أو L1-median (وهو تقدير حصين للموقع في متعدد المتغيرات وله خصائص جيدة ويسمى بالوسيط المكاني (spatial median) وقد تم تعريفه من قبل العالم weber

1909) بأنه أي نقطه لها أقل مجموع لمسافة أفليديس من كل النقاط في مجموع البيانات [4] وتوجد عدة

خوارزميات لإيجاد l1-median وقد تم الاعتماد على خوارزمية Hossjer and Croux (1995) [6]

فالوسيط L1- (spatial median) الى مجموعة من البيانات  $X = \{x_1, \dots, x_n\}$  لكل  $x_i \in R^p$

هو متجه يرمز له بـ  $\hat{\mu}_{sm}$  والذي يحقق أقل مجموع وكالاتي :

$$\hat{\mu}_{sm} = \arg \min_{\eta \in R^p} \sum_{i=1}^n \|x_i - \eta\| \quad \|\cdot\| \text{ يشير الى الاقليديس الطبيعي.}$$

ويتم إيجاد قيم  $\tilde{y}$  الحصينة بأستعمال الوسيط حيث يعتبر من المقاييس الحصينة أيضاً.

ثانياً المتغيرات التوضيحية (X) يطرح منها الوسيط الخاص L 1\_median والمتغير (y) يطرح منه

الوسيط، المتغيرات الحصينة  $Z = [\tilde{X}, \tilde{y}]$  يتم إيجاد

$$\tilde{Z}_i = S_z^{-1/2} Z_i, \quad i = 1, 2, \dots, n$$

٢- في هذه الخطوة يتم ايجاد الاتجاه الذي يعظم ويصغر معامل التفرطح للبيانات الحصينة وحسب المعادلات

(٢-١٨)، (٢-١٩)، ثم يتم إيجاد أسقاطات القيم الحصينة على هذين الاتجاهين وبعدها يتم ايجاد قياس الشواذ

للمتغير الاحادي  $r_i(t)$  وحسب المعادلة (٢-٢١)

## الجزئية \*

٣- حساب قياس الشواذ للمتغير الاحادي  $\tilde{r}_i$  من خلال أسقاطات القيم الحصينة على الاتجاه العشوائي من خلال اسلوب المعاينة التطبيقية ومن ثم البحث عن القيم الشاذة في هذه الاتجاهات العشوائية وحسب المعادلة (٢-٢٣)

٤- في هذه الخطوة يتم التدقيق لكل مشاهدة (i) بعد ايجاد قياس الشواذ للمتغير الاحادي  $r_i$  فإذا كانت  $r_i > 1$  فان المشاهدة i تصنف على انها قيمة شاذة وتزال من العينة اذا كان عددهم اصغر من  $[n - [n + p + 1] / 2]$  ، عدا ذلك فان المشاهدات التي عددها  $[n - [n + p + 1] / 2]$  والتي لها اكبر قيم الى  $r_i$  فهي تعتبر شواذ وتزال من البيانات، وبعدها يتم حساب مسافة مهالويس الحصينة ( $V_i$ ) بالاعتماد على البيانات الخالية من القيم الشاذة وحسب المعادلة (٢-٢٧) فإذا كانت مسافة مهالويس أكبر من قيمة القطع  $\chi_{p,1-99}^2$  تصنف على انها قيم شاذة . وفي حالة  $\tilde{S}_z^{-1}$  غير موجودة سوف نستعمل المعكوس العام لها. وتكرر هذه الخطوات حتى لا تبقى قيم شاذة ونستطيع ان نستخدم هذه الطريقة في حالة عدد المتغيرات أكبر من المشاهدات أيضاً باستخدام خوارزمية SVD.

## ١-٣ الجانب التجريبي [15,5,8]

من اجل الوصول الى الهدف تم استعمال اسلوب المحاكاة للمقارنة بين تقدير معاملات الانحدار لطريقة المربعات الصغرى الجزئية بأستعمال خوارزمية SIMPLS وطريقة المربعات الصغرى الجزئية الحصينة لخوارزميتي MPLSKURSD, PLSKURSD لعدة انواع من الشواذ وبنسب التلوث المختلفة ولابعاد متغيرات وحجوم عينات مختلفة.

## ٢-٣ وصف تجربة المحاكاة

في هذه التجربة تم توليد العينات بحجوم (200,150,100,50,30) وعدد المتغيرات التوضيحية (p=5) في حالة  $n > p$  اما في حالة  $p > n$  تم توليد العينات بالحجم  $n=20$  وعدد المتغيرات التوضيحية (p=40,60,80,100,150) وكان عدد المركبات (2) في الحالتين اما مصفوفة التباين الى t (مصفوفة القياس) فهي قطرية diag [6,4] وتم استعمال نسب تلوث  $\epsilon$  مساوية الى (0.3, 0.2, 0.1) وقد كررت التجربة 1000 مرة للحصول على الدقة العالية. في هذه التجربة سنقارن بين نوعين من الشواذ بالاضافة الى المشاهدات الاعتيادية وقد تمت المقارنة بين متجه معاملات الانموذج للطرائق المختلفة ومتجه المعلمات المقدر  $\hat{\beta}_a^{(l)}$  حيث تم تكرار التجربة لـ 1000 بمعيار  $Norm(\beta)$  فقد تم حساب المقياس كالاتي :

$$MSE_a(\hat{\beta}) = 1/m \sum_{l=1}^m \|\hat{\beta}_a^{(l)} - \beta\|^2 = Norm(\beta) \dots(3-1)$$

حيث  $\|\|\|$  يشير الى اقليدس الطبيعي و متجه مقدرات الانحدار  $\hat{\beta}_a^{(l)}$  في التكرار l و m يشير الى عدد التكرارات.



## الجزئية \*

بالاعتماد على النموذج الابتدائي (٢-٢) و (٢-٣) وعلى فرض ان الوسط الحسابي يساوي صفر حيث النموذج يتبع التوزيع الطبيعي وغير المرتبط تم توليد البيانات التالية:

١- المشاهدات الاعتيادية (Regular observation) وهي المشاهدات التي تكمن من ضمن الغالبية من المشاهدات التي تقع على خط الانحدار كالاتي:

$$t \propto N_a(0_a, \sum_t)$$

$$x = I_{p,a}t + N_p(0_p, 0.1 I_p), \quad p > a \quad \dots(3-2)$$

$$y = q^T t + N(0,1)$$

حيث  $I_{p,a} = 1$  عندما  $i=j$  و  $I_{p,a} = 0$  ماعدا ذلك .

$q$  متجه من الواحدات ذو بعد  $ax1$

$I_p$  مصفوفة الوحدة ذو بعد  $p \times p$

اما بالنسبة للبيانات الملوثة فقد تم توليد مجموعة تتكون من  $[100 * (1 - \epsilon)]$  من المشاهدات وفقا للنموذج اعلاه وتم اضافة  $[100 * \epsilon]$  من المشاهدات الملوثة التي تولد من خلال نماذج التلوث ادناه لتكون انواع التلوث التالية

٢- الشواذ العمودية  $vertical\ outliers$  وهي النقاط التي تكون بعيدة عن خط الانحدار بينما اسقاطاتها تكون ضمن المشاهدات الجيدة .

$$y_\epsilon = q^T t + N(30, 0.3) \quad \dots(3-3)$$

٣- نقاط الانعطاف جيدة ( $good\ leverage\ point$ ) وهي النقاط التي تقع على خط الانحدار لكنها بعيدة عن تجمع الغالبية من المشاهدات التي تحدد خط الانحدار .

$$x_\epsilon = t_\epsilon I_{a,p} + N_p((0_a, 3_{p-a}), 0.01 I_p)$$

$$t_\epsilon \propto N_a(3_a, \sum_t)$$

$$y_\epsilon = q^T t + N(30, 0.3) \quad \dots(3-4)$$

## نتائج المحاكاة

تم تقدير المعلمات للطريقتين المذكورة في الجانب النظري وتم حساب  $(\hat{\beta})$   $MSE_\alpha$  عند جميع حجوم العينات ولمختلف اعداد المتغيرات التوضيحية وبنسب تلوث مختلفة والتي تم ذكرها وكانت النتائج كما في الجداول ادناه:



## الجزئية \*

جدول رقم (١) يوضح قيم MSE لمتجه المعلمات المقدرة عند (n=30) و (p=5) عند نسب التلوث (0.3,0.2,0.1)

الطرائق الشواذ	PLS	PLSKUR	MPLSKUR
No contamination لاتوجد شواذ	3.93198	6.705675	6.29413
10%Vertical outliers الشواذ العمودية	17.85582	6.683227	7.938724
10%tGood leverage point نقاط الرفع الجيدة	26.62007	6.15471	6.179438
20%Vertical outliers الشواذ العمودية	30.9934	11.32486	9.957953
20%Good leverage point نقاط الرفع الجيدة	28.49926	6.120336	5.882657
30%Vertical outliers الشواذ العمودية	52.44824	28.51931	26.12804
30%Good leverage point نقاط الرفع الجيدة	29.05623	6.703405	7.145314

جدول رقم (٢) يوضح قيم MSE لمتجه المعلمات المقدرة عند (n=50) و (p=5) عند نسب التلوث = (0.3,0.2,0.1)

الطرائق الشواذ	PLS	PLSKUR	MPLSKUR
No contamination لاتوجد شواذ	3.143267	4.455928	4.145776
10%Vertical outliers الشواذ العمودية	12.74727	4.238654	4.396087
10% Good leverage point نقاط الرفع الجيدة	27.27451	4.212243	4.217679
20% Vertical outliers الشواذ العمودية	22.78322	4.551482	4.438952
20% Good leverage point نقاط الرفع الجيدة	29.07668	4.707521	4.964694
30% Vertical outliers الشواذ العمودية	30.0315	10.47879	11.1865
30% Good leverage point نقاط الرفع الجيدة	29.26588	4.551511	4.448104



## الجزئية \*

جدول رقم (٣) يوضح قيم MSE لمتجه المعلمات المقدرة عند (n=100) و (p=5) عند نسب التلوث = (0.3,0.2,0.1)

الطرائق الشواذ	PLS	PLSKUR	MPLSKUR
No contamination لاتوجد شواذ	2.784376	2.944539	2.955385
10% Vertical outliers الشواذ العمودية	7.436278	3.101319	3.100473
10% Good leverage point نقاط الرفع الجيدة	27.10441	2.959698	2.945401
20% Vertical outliers الشواذ العمودية	12.25034	3.202687	3.128556
20% Good leverage point نقاط الرفع الجيدة	29.06709	3.078286	3.070211
30% Vertical outliers الشواذ العمودية	18.85743	4.1549	3.713399
30% Good leverage point نقاط الرفع الجيدة	29.88383	3.466043	3.560048

جدول رقم (٤) يوضح قيم MSE لمتجه المعلمات المقدرة عند (n=150) و (p=5) عند نسب التلوث = (0.3,0.2,0.1)

الطرائق الشواذ	PLS	PLSKUR	MPLSKUR
No contamination لاتوجد شواذ	2.564526	2.619275	2.616338
10% Vertical outliers الشواذ العمودية	5.421788	2.669105	2.666211
10% Good leverage point نقاط الرفع الجيدة	27.49836	2.646156	2.651284
20% Vertical outliers الشواذ العمودية	9.084631	2.828726	2.82872
20% Good leverage point نقاط الرفع الجيدة	28.98738	2.711855	2.704045
30% Vertical outliers الشواذ العمودية	11.94638	3.17712	2.992124
30% Good leverage point نقاط الرفع الجيدة	29.57321	3.069245	3.128813



## الجزئية \*

جدول رقم (٥) يوضح قيم MSE لمتجه المعلمات المقدرة عند (n=200) و (p=5) عند نسب التلوث = (0.3,0.2,0.1)

الطرائق الشواذ	PLS	PLSKUR	MPLSKUR
No contamination لاتوجد شواذ	2.468142	2.498064	2.504395
10% Vertical outliers الشواذ العمودية	4.565388	2.555479	2.553579
10% Good leverage point نقاط الرفع الجيدة	27.33878	2.570506	2.567473
20% Vertical outliers الشواذ العمودية	7.58146	2.602076	2.605444
20% Good leverage point نقاط الرفع الجيدة	29.00699	2.623537	2.621256
30% Vertical outliers الشواذ العمودية	9.250026	2.688	2.748775
30% Good leverage point نقاط الرفع الجيدة	29.69553	3.007342	3.29134

جدول رقم (٦) يوضح قيم MSE لمتجه المعلمات المقدرة عند (n=20) و (p=40) عند نسب التلوث = (0.3,0.2,0.1)

الطرائق الشواذ	PLS	PLSKUR	MPLSKUR
No contamination لاتوجد شواذ	90.76676	83.67835	83.62133
10 % Vertical outliers الشواذ العمودية	203.7097	130.7868	129.3845
10% Good leverage point نقاط الرفع الجيدة	76.85447	76.39718	76.36648
20% Vertical outliers الشواذ العمودية	255.2324	121.5665	119.8704
20% Good leverage point نقاط الرفع الجيدة	43.28528	42.79096	42.77652
30% Vertical outliers الشواذ العمودية	399.5773	163.3382	159.6765
30% Good leverage point نقاط الرفع الجيدة	51.60886	51.29083	51.2872



## \*الجزئية\*

جدول رقم (٧) يوضح قيم MSE لمتجه المعلمات المقدرة عند (n=20) و (p=60) عند نسب التلوث = (0.3, 0.2, 0.1)

الطرائق الشواذ	PLS	PLSKUR	MPLSKUR
No contamination لا توجد شواذ	66.45336	58.18977	58.18797
10% Vertical outliers الشواذ العمودية	174.4606	95.6495	95.23426
10% Good leverage point نقاط الرفع الجيدة	58.34411	57.55186	57.5446
20% Vertical outliers الشواذ العمودية	281.1992	123.3543	122.2643
20% Good leverage point نقاط الرفع الجيدة	52.75271	51.99841	51.98885
30% Vertical outliers الشواذ العمودية	421.7863	175.031	173.1078
30% Good leverage point نقاط الرفع الجيدة	74.58289	74.13526	74.10261

جدول رقم (٨) يوضح قيم MSE لمتجه المعلمات المقدرة عند (n=20) و (p=80) عند نسب التلوث = (0.3, 0.2, 0.1)

الطرائق الشواذ	PLS	PLSKUR	MPLSKUR
No contamination لا توجد شواذ	70.61393	62.59702	62.54708
10% Vertical outliers الشواذ العمودية	177.7266	98.43487	97.7031
10% Good leverage point نقاط الرفع الجيدة	57.87697	56.80069	56.78465
20% Vertical outliers الشواذ العمودية	320.3175	168.6173	167.5111
20% Good leverage point نقاط الرفع الجيدة	99.34089	98.55477	98.51187
30% Vertical outliers الشواذ العمودية	397.5902	166.802	165.8302
30% Good leverage point نقاط الرفع الجيدة	79.74979	79.03737	78.97919



## الجزئية \*

جدول رقم (٩) يوضح قيم MSE لمتجه المعلمات المقدرة عند (n=20) و (p=100) عند نسب التلوث = (0.3, 0.2, 0.1)

الطرائق الشواذ	PLS	PLSKUR	MPLSKUR
No contamination لاتوجد شواذ	74.42864	66.54818	66.44204
10% Vertical outliers الشواذ العمودية	173.8288	100.4384	100.4113
10% Good leverage point نقاط الرفع الجيدة	63.48287	62.27119	62.24431
20% Vertical outliers الشواذ العمودية	303.9499	155.7178	155.3580
20% Good leverage point نقاط الرفع الجيدة	97.46692	96.38357	96.34031
30% Vertical outliers الشواذ العمودية	417.1883	188.8668	188.0394
30% Good leverage point نقاط الرفع الجيدة	109.4832	108.5655	108.5351

جدول رقم (١٠) يوضح قيم MSE لمتجه المعلمات المقدرة عند (n=20) و (p=150) عند نسب التلوث = (0.3, 0.2, 0.1)

الطرائق الشواذ	PLS	PLSKUR	MPLSKUR
No contamination لاتوجد شواذ	147.6072	141.1587	141.1408
10% Vertical outliers الشواذ العمودية	230.5816	169.0253	168.9759
10% Good leverage point نقاط الرفع الجيدة	143.02	141.3591	141.3234
20% Vertical outliers الشواذ العمودية	339.5476	215.9492	215.4868
20% Good leverage point نقاط الرفع الجيدة	167.7448	166.3475	166.3294
30% Bad leverage point نقاط الرفع السيئة	158.4073	156.3259	156.265
30% Good leverage point نقاط الرفع الجيدة	155.5631	154.501	154.4384



## الجزئية \*

## تحليل النتائج والاستنتاجات

- في حالة عدد المشاهدات اكبر من المتغيرات ( $n > p$ ) يتضح مايلي :-
- ١- يتضح من الجدول رقم (١) والذي يتضمن ( $n=30$   $p=5$ ) مايلي  
- عند عدم وجود الشواذ كانت الطريقة PLS هي الأفضل .  
- عند الشواذ العمودية كانت الطريقة plskursd هي الأفضل عند نسب التلوث 0.1 والطريقة mplskursd هي الأفضل عند نسب التلوث 0.2,0.3  
- عند نقاط الرفع الجيدة كانت الطريقة plskursd هي الأفضل عند نسب التلوث 0.3,0.1 و mplskursd هي الأفضل عند نسب التلوث 0.2
  - ٢- يتضح من الجدول رقم (٢) والذي يتضمن ( $n=50$   $p=5$ ) مايلي :  
- عند عدم وجود الشواذ كانت الطريقة PLS هي الأفضل .  
- عند الشواذ العمودية كانت الطريقة plskursd هي الأفضل عند نسب التلوث 0.3,0.1, والطريقة mplskursd هي الأفضل عند نسب التلوث 0.2  
- عند نقاط الرفع الجيدة كانت الطريقة plskursd هي الأفضل عند نسب التلوث 0.2,0.1 و mplskursd هي الأفضل عند نسب التلوث 0.3
  - ٣- يتضح من الجدول رقم (٣) والذي يتضمن ( $n=100$   $p=5$ ) مايلي :  
- عند عدم وجود الشواذ كانت الطريقة PLS هي الأفضل .  
- عند الشواذ العمودية كانت الطريقة mplskursd هي الأفضل في كل الحالات نسب التلوث  
- عند نقاط الرفع الجيدة كانت الطريقة mplskursd هي الأفضل عند نسب التلوث 0.1,0.2 وان plskursd هي الأفضل عند نسب التلوث 0.3
  - ٤- يتضح من الجدول رقم (٤) والذي يتضمن ( $n=150$   $p=5$ ) مايلي :  
- عند عدم وجود الشواذ كانت الطريقة PLS هي الأفضل .  
- عند الشواذ العمودية كانت الطريقة mplskursd هي الأفضل في كل الحالات نسب التلوث .  
- عند نقاط الرفع الجيدة كانت الطريقة plskursd هي الأفضل عند نسب التلوث 0.1,0.3 وان mplskursd هي الأفضل عند نسب التلوث 0.2 .
  - ٥- يتضح من الجدول رقم (٥) والذي يتضمن ( $n=200$   $p=5$ ) مايلي :  
- عند عدم وجود الشواذ كانت الطريقة PLS هي الأفضل .  
- عند الشواذ العمودية كانت الطريقة mplskursd هي الأفضل عند نسب التلوث 0.1 و plskursd هي الأفضل في الحالات الأخرى  
- عند نقاط الرفع الجيدة كانت الطريقة mplskursd هي الأفضل عند نسب التلوث و plskursd هي الأفضل عند نسب التلوث 0.1,0.2,0.3
  - ٦- يتضح من الجدول رقم (٦ و٧ و٨ و٩ و١٠) نجد مايلي:-  
- عند عدم وجود الشواذ والشواذ العمودية ونقاط الرفع الجيدة كانت الطريقة mplskursd هي الأفضل في جميع الجداول.
- ويمكن ان نلخص النتائج بما يلي:
- ١- ان طريقة PLS هي الأفضل في تقدير المعلمات في حالة عدم وجود الشواذ عندما عدد المشاهدات اكبر من المتغيرات وذلك لامتلاكها اقل متوسط مربعات الخطأ (MSE).
  - ٢- ان طريقة PLSKURSD هي الأفضل تقريبا بالمقارنة مع MPLSKURD في تقدير المعلمات عند احتواء البيانات على القيم الشاذة عندما عدد المشاهدات اكبر من المتغيرات
  - ٣- ان طريقة MPLSKURSD هي الأفضل في تقدير المعلمات في حالة عدد المتغيرات اكبر من المشاهدات وعند عدم وجود الشواذ وفي الشواذ العمودية ونقاط الانعطاف (الرفع) الجيدة .
  - ٤- اشارت النتائج الى ان قيم متوسط مربعات الخطأ تتناقص بزيادة حجم العينة ولكلا الطرائق في التقدير.



## \*الجزئية\*

## التوصيات

نوصي باستخدام الخوارزمية الحصينة المقترحة MPLSKURSD في حالة عدد المتغيرات اكبر من المشاهدات لكل نسب التلوث واستعمال PLSKURSD في حالة المشاهدات اكبر من المتغيرات لبعض نسب التلوث

## المصادر

- 1- Andersson, M., (2009) "Acomarison of nine PLS1 algorithms" J. chemometrics; 23:518-529
- 2- Collins,B., 2010 " partial least squares regression "october 13
- 3-De Jong, S., 1993. SIMPLS: an alternative approach to partial least squares regression. Chemometrics and Intelligent Laboratory Systems, 18: 251–26
- 4- Fritz,H., Filzmoser,P. , Croux,C., 2010" A comparison of algorithms for the multivariate l1\_median"  
<http://www.statistik.tuwien.ac.at>
- 5-Gonzales, J., & et al. (2009) "A robust partial least squares regression method with application" J . Chemometrics; 23:78-90
- 6- Hossjer and Croux (1995) "Generalizing Univariate Signed Rank Statistics for Testing and Estimating a Multivariate Location Parameter", Non-parametric Statistics, 4, 293
- 7-Hubert, M., Rousseeuw, P.J., Verboven,S (2002), "A fast robust method for principal components with applications to chemometrics" J . Chemometrics and Intelligent Laboratory Systems, 60,101-111.
- 8- Hubert, M., & Vanden , K., (2003) "Robust Methods for Partial Least Squares Regression "9<sup>th</sup> October J. chemometrics p 537-549.
- 9- - Nguyen,D.V. and Rocke ,D.M.,(2002) "Tumor classification by partial least squares using microarray gene expression data" Bioinformayics vol.18 no .1 pages 39-50.
- 10-Pena,D. and Prieto,F.J., (1997) "Robust Covariance Matrix Estimation and Multivariate Outlier Detection" working paper 97-08, statistics and Econometrics series 04
- 11- Pena,D. and Prieto,F.J.,(2001) "Multivariate Outlier Detection and Robust Covariance Matrix Estimation" Technometrics,vol.43,no.3
- 12- Pena,D., and Prieto,F.J., (2007 )"Combining Random and Specific Direction for Outlier Detection and Robust Estimation in High-Dimensional Multivariate Data" journal of computational and graphical statistics, ,vol.16, number 1,pages 228-254
- 13-Rosipal ,R., Kr`amer ,N.,(2006)" Overview and Recent Advances in Partial Least Squares" Springer-Verlag Berlin Heidelberg pp. 34–
- 14-Vanden Branden. K. & Hubert. M,(2003) "Robustness Properties of a robust PLS regression method" December 10. J.Analytica Chimica p229-241.
- 15-verardi .v and croux .V(2009) "Robust Regression in stata" K.u.leuven Facutly of Business and Economices.
- 16- Wakeling .J.N and Macfie. H.J.H "A Robust PLS Procedure" (1992) journal of chemometrics, vol.6,189-198.
- 17- Yuditskaya,s., (2010)" An Overview of Methods in Linear Least-Squares Regression"pattern recognition and Analysis ,November 4



## Comparison of some robust methods to estimate parameters of partial least squares regression (PLSR)

### Abstract

The technology of reducing dimensions and choosing variables are very important topics in statistical analysis to multivariate. When two or more of the predictor variables are linked in the complete or incomplete regression relationships, a problem of multicollinearity are occurred which consist of the breach of one basic assumptions of the ordinary least squares method with incorrect estimates results.

There are several methods proposed to address this problem, including the partial least squares (PLS), used to reduce dimensional regression analysis. By using linear transformations that convert a set of variables associated with a high link to a set of new independent variables and unrelated with each other, which are called, the components. These components are orthogonal and independent from each other.

The method of partial least squares PLS is failed in dealing with data that consist of the presence of Outliers values and hence the success of this method depends on the absence of such outliers values that have undesirable effect on the results. In order to reduce the presence of these values, we resorted to use the robust methods.

In this research a method of PLSKURSD that applied SIMPLS algorithms on variance-covariance robust matrix. Also the proposed method MPLSKURSD are used which is a modified method to the PLSKURSD method. parameters linear regression model by partial least squares(PLS) is compared with modalities robust partial least squares through the simulation experiments depends on the presence of several types of outlier values of data for different rates of pollution, volumes of samples, and variables dimensions

**Keyword:** Partial Least Squares; Robust Covariance Matrix; kurtosis; Projection