



طريقة مقترحة لتعديل معلمة الموقع

أ. م. د. لقاء علي العلوي
كلية الادارة والاقتصاد- جامعة بغداد
قسم الاحصاء

المخلص

تقدير الموقع والتشتت في متعدد المتغيرات مع خاصيتي تساوي التغاير والانهيال الموجب دائما يكون صعب والمقدر الذي يحقق هذه الخواص هو مقدر القطع البيضوي الأصغر Minimum Volume Ellipsoid (MVE) Estimator . حساب (MVE) المضبوط غالبا يكون غير مقنع لذلك يتم إعادة ترتيب خوارزمية تقريبية لحسابه. في الانحدار هناك خوارزميات عديدة لحساب مقدرات الانهيال الموجب مثل اصغر وسيط للمربعات الذي يعيد حساب التقاطع في كل خطوة لتحقيق النتيجة. هذا الأسلوب يدعى بـ (تعديل الموقع). في هذا البحث اشرنا إلى تقنيات مشابهة تدعى تعديل الموقع، ممكن وضعها في حساب (MVE). لهذا الغرض تم استخدام اصغر قطع كروي Minimum Volume Ball (MVB) لغرض تقليل دالة الهدف لـ (MVE). تم اعتبار خوارزمية مضبوطة لحساب (MVE). وكبديل لتعديل الموقع (MVB) تم استخدام تعديل الموقع (L_1) في مرحلة ثانية، والذي ليس من الضروري ان يقلل دالة الهدف لـ (MVE) ولكنه يفرز مقدرات أكثر كفاءة لجزء الموقع. والمحاكاة تمت لمقارنة هذه الأنواع في تعديل الموقع.

Abstract

Estimating multivariate location and scatter with both affine equivariance and positive break down has always been difficult. A well-known estimator which satisfies both properties is the Minimum volume Ellipsoid Estimator (MVE) Computing the exact (MVE) is often not feasible, so one usually resorts to an approximate Algorithm. In the regression setup, algorithm for positive-break down estimators like Least Median of squares typically recomputed the intercept at each step, to improve the result. This approach is called intercept adjustment. In this paper we show that a similar technique, called location adjustment, Can be applied to the (MVE). For this purpose we use the Minimum Volume Ball (MVB). In order to lower the (MVE) objective function. An exact algorithm for calculating the (MVB) is presented. As an alternative to (MVB) location adjustment we propose (L_1) location adjustment, which does not necessarily lower the (MVE) objective function but yields more efficient estimates for the location part. Simulations Compare the two type of location adjustment.

1. المقدمة Introduction

ان وجود قيم شاذة outliers ضمن مجموعة البيانات تؤثر بشكل كبير في نتائج التحليل الاحصائي لتلك البيانات ويظهر هذا في حقيقة ان اغلب الطرائق الإحصائية تؤكد على افتراض ان البيانات المستخدمة في التحليل هي بيانات متجانسة homogeneous data أي ان كل نقاط البيانات الموجودة في مجموعة معينة تتبع نفس التوزيع المفترض لذلك طرائق بديلة طورت للتعامل مع الأخطاء في النموذج والتلويث في البيانات.

في التحليل الحصين مقدر اصغر قطع بيضوي Minimum Volume Ellipsoid Estimator (MVE) غالبا ما يستخدم في تقدير الموقع والتشتت، حيث يعرف (MVE) نسبة الى (Rousseeuw) (1985) كاصغر قطع بيضوي منظم يحوي على الاقل h من عناصر مجموعة البيانات $X = \{X_1, X_2, \dots, X_n\} \subset R^p$ ، حيث ان مقدر الموقع Locayion Estimator للـ (MVE) هو مركز ذلك القطع، بينما مقدر التشتت Scatter Estimator في الـ (MVE) يعود الى شكل المصفوفة، وبذلك يمكن تقليل مقدرات (MVE) بالصورة ادناه :

$$(\hat{\mu}, \hat{\Sigma}) = \underset{(\mu, \Sigma) \in R^p \times SPD(p)}{\operatorname{arg\,min}} d_h^2(\mu, S) \quad \dots \dots \dots (1)$$

|Σ|=1

حيث SPD(P) مجموعة كل المصفوفات المتماثلة الموجبة قطعاً Symmetric Positive definit matrix $S \in R^{p \times p}$. يطلق على \hat{S} مصفوفة الشكل Shape matrix لان \hat{S} تحدد شكل القطع البيضوي وليس الكمية ، طالما ضروريا $|\hat{S}| = 1$.

بواسطة $d_h^2(\mu, S)$ تعرف المسافة التربيعية h^{th} بين X_i و μ بواسطة المصفوفة S كما يلي

$$d_h^2(\mu, S) = \left\{ (X_i - \mu)' S^{-1} (X_i - \mu); \leq i \leq n \right\}_{(h)} = \left\{ \|X_i - \mu\|_S^2; 1 \leq i \leq n \right\}_{(h)}$$

$$= \operatorname{median}_i \|X_i - \mu\|_S^2 \quad \dots \dots \dots (2)$$

حيث الوسيط للاحصاءات المرتبة (Order Statistic h^{th}) . ولهذا مقدر (MVE) يكون بالصيغة

$$(\hat{\mu}, \hat{\Sigma}) = (\hat{\mu}, c(n, p, h) d_h^2(\hat{\mu}, \hat{S}) \hat{S}) \quad \dots \dots \dots (3)$$

حيث $c(n, p, h)$ هو عامل التصحيح Correction Factor لتكون $\hat{\Sigma}$ متسقة Consistent لـ Σ في النموذج الطبيعي وقيمته مساوية الى $c = 1 + \frac{15}{n-p}$.

إذا اردنا تقليل نقطة الانهيار للمقدر، القيمة لـ h في المعادلة (1) و (3) تكون $h = \frac{n+p+1}{2} \approx \frac{n}{2}$

لكن اذا (هذا ما يحدث غالبا) عرفنا كسر الشواذ على الاكثر α حيث $0 < \alpha < 1/2$. نستطيع العمل مع مقدر MVE (α) حيث $h = [n(1 - \alpha)]$. اختيار $\alpha = 1/4$ يكون جيد وبهذا نستطيع وضع خوارزمية ايجاد مقدر (MVE) كما يلي :

1. بعد تحديد حجم العينة (n) يتم تحديد العينات الجزئية بحجم ($p+1$) في المشاهدات ، أي نختار C_{p+1}^n من العينات .

2. لكل عينة جزئية يتم حساب متجه الاوساط ومصفوفة التشتت ليجري بعدها استخراج المسافات التربيعية .

3. ترتب قيم المسافات تصاعديا ويتم اختيار d_h^2 لكل عينة جزئية ، ليتم بعدها اختيار افضل عينة جزئية وهي العينة التي تمتلك اصغر دالة هدف (1).

4. بعد ايجاد مصفوفة التشتت المصححة يتم استخراج المسافات التربيعية وفقا لكافة المشاهدات ليتم بعدها الاسلوب التعاقبي في الحساب والذي يتم بزيادة مشاهدة واحدة الى حجم المجموعة الجزئية ليصبح ($p+2$) ثم لتصل بعدها الى (h) من المشاهدات .

اذا وجدنا نقاط h نستطيع تغطيتها بشكل بيضوي باصغر قطع عشوائي . في المعادلة (1) اصغر $d_h^2(\mu, S)$ تصبح صفر ولكنها لا تمتد لكافة المشاهدات الشرط هذا كافي لايجاد مقدرات (MVE) طالما

اصغر قطع هو الاصغر للعدد المنتهي $\binom{n}{h}$ للقطوع الموجبة بالضبط .

2. حساب المقدرات الحصينة للانحدار

Computation of the Robust Estimators in Regression

في نموذج الانحدار $\gamma_i = \beta' X_i + \alpha + \epsilon_i (i = 1, \dots, n)$ معلمة الميل Slope Parameter

هي $\beta \in R^{p-1}$ ومعلمة التقاطع intercept parameter هي $\alpha \in R$ لذلك فان مقدر اصغر وسيط للمربعات Least Median of Square (LMS) Estimator – (Rousseeuw 1985) بالنسبة لـ (α, β) يعرف بواسطة :

$$(\hat{\alpha}, \hat{\beta}) = \arg \min_{(\alpha, \beta) \in R \times R^{p-1}} \text{median}(y_i - \beta' X_i - \alpha)2 \quad \dots \dots \dots (4)$$

عادة، تقدير التقاطع $\hat{\alpha}$ يحسب مشروطا بقيمة تقدير الميل $\hat{\beta}$ لغرض تخفيض قيمة دالة الهدف. هذا الاسلوب في تعديل التقاطع intercept adjustment يتضمن تشريح مشكلة التصغير (4) الى جزئين:

$$\hat{\beta} = \arg \min_{\beta \in R^{p-1}} \text{median}(y_i - \beta' X_i - \hat{\alpha}(\beta))2 \quad \dots \dots \dots (5)$$

$$\hat{\alpha}(\hat{\beta}) = \arg \min_{\alpha \in R} \text{median}(y_i - \beta' X_i - \alpha)2 \quad \dots \dots \dots (6)$$

وهذا يعود لقولنا ان $\hat{\alpha}(\beta)$ هي تقدير الموقع لـ (LMS) احادي التغير لـ (n) من العناصر $\gamma_i - \beta'X_i$, $(i = 1,2,\dots,n)$. وهناك خوارزمية مضبوطة لتقدير موقع (LMS) احادي المتغير، طالما انها نقطة الوسط لاصغر فترة تحتوي h من المشاهدات . بصورة ملائمة لما ذكر اعلاه، نستطيع كتابة (1) بالشكل :

$$\hat{S} = \arg \min_{\substack{S \in SPD(P) \\ |S|=1}} d_h^2(\hat{\mu}(s), S) \quad \dots\dots\dots (7)$$

حيث لـ (s) المعطاة مع $|s|=1$ نضع :

$$\hat{\mu}(s) = \arg \min_{\mu \in R^p} d_h^2(\mu, S) \quad \dots\dots\dots (8)$$

$$\hat{\mu}(s) = \arg \min_{\mu \in R^p} \text{median} \|x_i - \mu\| S$$

$$\therefore \hat{\mu}(s) = \arg \min_{\mu \in R^p} \text{median} \|S^{-1/2} x_i - S^{-1/2} \mu\| \quad \dots\dots\dots (9)$$

حيث $S^{1/2}$ تعرف بالجذر المتماثل لـ S أي ان $(S = S^{1/2} S^{1/2})$ مع $S^{1/2}$ المتماثلة . باستخدام مجموعة بيانات التحويلات $\{y_i = S^{-1/2} X_i, i = 1, \dots, n\}$ نحصل على:

$$\hat{\mu}(s) = S^{1/2} \arg \min_{\theta} \text{median} \|y_i - \theta\| \quad \dots\dots\dots (10)$$

نعرف ان $\hat{\mu}(\hat{S})$ مع \hat{S} المعرفة بـ (7) تساوي $\hat{\mu}$ ولذلك تبقى متساوية التغير affine equivariant . قيمة θ التي تصغر وسيط $\|y_i - \theta\|$ هي مركز الكرة مع مقدر اصغر قطع (غير صفري) الذي يحوي على الاقل h من مجموعة البيانات Y. وهذا يقود لمقدر اصغر قطع كروي Minimum Volume Ball Estimator (MVB) معرف من قبل (Rousseeuw 1984, PP877) علما ان هذا المقدر يتمتع بخصائص الحصانة (متساوي التغير ومتعامد وذا نقطة انهيار مساوية الى 50%) .

كذلك نحن نفترض مقدرات الانحدار الحصينة والمعرف مقياس الحصانة لها بتقليل انتشار البواقي. كمثال عندما نستخدم الانحراف القياسي المشذب Trimmed standard deviation فنحصل على مقدر المربعات المشذبة الصغرى Least Trimmed squares estimator (LTS) — (Rousseeuw & Leroy 1987) للعلم فان الغرض من اجراء تعديل الموقع في (MVE) هو لتخفيض دالة الهدف (1). هذا مشابه لتقنية تعديل التقاطع في الانحدار الحصين والذي يكون ذو فائدة.

3. حساب اصغر قطع كروي Computation of Minimum Volume Ball

لناخذ $Y = \{\gamma_1, \gamma_2, \dots, \gamma_n\} \subset R^p$ ، مقدر موقع (MVB) يعرف بـ

$$MVB(Y) = \arg \min_{\mu \in R^p} \text{median} \|\gamma_i - \mu\| \quad \dots \quad (11)$$

وثانية الوسيط يمثل الاحصاءة المرتبة h^{th}

الخطوات ادناه تعطي الخوارزمية المضبوطة لمقدر (MVB) .

1. نضع R_{best} الابتدائية مساوية الى $+\infty$

2. لاي عدد $2 \leq k \leq P+1$ ولاي مجموعة جزئية k وان

$$J = \{i_1, i_2, \dots, i_k\} \subset \{1, 2, \dots, n\} \text{ يكون:}$$

a. نضع $A_j = \text{Affinespan}\{\gamma_i : j \in J\}$ واذا كان $\dim(A_j) < k-1$ نرجع للخطوة 2 (أي

نحذف هذه المجموعة الجزئية J)

b. لذلك يكون $\dim(A_j) = k-1$. نحدد النقطة الوحيدة μ_j في A_j التي تقع في نفس البعد

الاقليدي Euclidean distance لكل y_j لكل $j \in J$ وذلك بحل $k \times k$ من نظام المعادلات

الخطية .

c. نحسب $R_j = \text{median}_i \|\gamma_i - \mu_j\|$ واذا كان $R_j \geq R_{best}$ نرجع للخطوة 2

d. نسجل $R_{best} = R_j$ وكذلك $\mu_{best} = \mu_j$

3. نسجل $MVB(Y) = \mu_{best}$

المشكلة هي في ايجاد k والمجموعة الجزئية $\{i_1, i_2, \dots, i_k\}$ حيث نبحت خلال كل المجموعات

الجزئية J الممكنة لكل الاحجام الممكنة ونحسب R_j .

طالما ان خوارزمية (MVB) المضبوطة تاخذ وقت طويل (وقت تكرار العملية يعتمد على n) لذلك وجدت خوارزمية تقريبية بديلة ، حيث الخطوة 2 تتعامل فقط مع المجموعة الجزئية N_{samp} ذات الحجم

($p+1$) (بواسطة المحاكاة نجد ان المجموعة الجزئية J التي تنتج R_{best} و μ_{best} ذات حجم ($p+1$)

باحتمال عالي).

بديل اخر يكون قريب جدا لخوارزمية المجموعة الجزئية ($p+1$) التقريبية يكون :

1. لكل $j=1, 2, \dots, n$ نحسب $R_j = \text{median}_i \|\gamma_i - \gamma_j\|$.

2. نضع $MVB(y)$ مساوية الى المشاهدة y_j لنحصل على اوطا R_j .

هذه الخوارزمية اقل كلفة ولكنها تبقى تحسب مقدرات متساوية التباين متعامدة وذات 50% نقطة انهيار.

4. تعديل الموقع بواسطة MVB (MVB) Location Adjustment by MVB

ايجاد الحل المضبوط لـ (MVB) الذي يقلل المشكلة (1) يكون غالبا غير مقتنع . لذلك تتم اعادة ترتيب خوارزمية المجموعة الجزئية (p+1) التقريبية. بسهولة يمكن تكييف هذه الخوارزمية بعدم اشتراك تعديل الموقع باستخدام مقدر اصغر قطع كروي. وهذا يقود الى الخوارزمية الاتية:

1. نحدد القيمة الابتدائية لـ R_{best} بـ $+\infty$
2. لكل مجموعة جزئية (p+1) حيث $J \subset \{1,2,\dots,n\}$ نعمل مايلي:
 - a. نحسب $\mu_j = \frac{1}{p+1} \sum_{i \in J} X_i$ وكذلك $C_j = \frac{1}{p} \sum_{i \in J} (\chi_i - \mu_j)(\chi_i - \mu_j)'$ واذا كان $|C_j| = 0$ نرجع للخطوة 2.
 - b. نحسب $S_j = |C_j|^{-1/p} C_j$ لتكون $|S_j| = 1$.
 - c. تحول مصفوفة البيانات X الى Y حيث $Y = \{S_j^{-1/2} X_i ; i = 1,2,\dots,n\}$
 - d. نحسب تقدير $\theta_j = MVB(Y)$ ونضع $R_j = \text{median} \|\gamma_i - \theta_j\|$
 - e. اذا كان $R_j = R_{best}$ نرجع للخطوة 2
 - f. نسجل التقدير الاخير $(\hat{\mu}, \hat{\Sigma})$ حيث $\hat{\mu} = S_{best}^{1/2} \theta_{best}$ وكذلك $\hat{\Sigma} = C(n, p, h) R_{best}^2 S_{best}$ نعرف ان اوطا قيمة لدالة الهدف (1) مساوية الى R_{best}^2 .

بزيادة الابعاد P، يصبح من غير الممكن اعتماد كل المجموعات الجزئية $\binom{n}{p+1}$ لذلك نستطيع ان نبقي نبحث باسلوب الاختيار العشوائي للمجموعات الجزئية N_{samp} ذات الحجم (p+1). نستطيع كذلك وضع تعديل (MVB) مرة واحدة كتحسين اخير لخوارزمية المجموعة الجزئية (P+1) العادية لمقدر (MVE). من الضروري الاشارة الى ان كل هذه النسخ لـ (MVE) المقارنة مع تعديل الموقع (MVB) هي نظريات متساوية التباين. لاننا نضع (MVB) المتساوية التباين والمتعامد للبيانات في مصفوفة (MVE).

5. المحاكاة Simulation:

خوارزمية المجموعة الجزئية (p+1) لـ (MVB) وضعت من قبل (Rousseeuw & Leroy 1987 pp 259-260) ووضعت ضمن برامج SAS, S-plus. هذا المبحث ضم دراستين للمحاكاة، الاولى ضمن متعدد المتغيرات والثانية في تحليل الانحدار. دراسة المحاكاة الاولى قارنت بين المقدرات الاتية من حيث الموقع والتشتت: مقدر المجموعة الجزئية (p+1) الاصيلي $(\hat{\mu}, \hat{\Sigma})$ ، مقدر المجموعة الجزئية (p+1) $(\tilde{\mu}, \tilde{\Sigma})$ مع تعديل (MVB) في كل خطوة، وفي النهاية مقدر المجموعة الجزئية (p+1) $(\tilde{\mu}, \tilde{\Sigma})$ مع تعديل (MVB) مرة واحدة في النهاية. تم توليد نوعين من البيانات الاول هو حالة البيانات الطبيعية حيث $X_i \approx N(0, I_p)$ لكل $i = 1, 2, \dots, n$. وفي الحالة الثانية هناك 20% من المشاهدات الملوثة (بإبدالها بـ $100e_1$ حيث e_1 هو المتجه الاحادي الاول first unit vector). هذا ينتج عنقود من الشواذ المتطرفة. تم اولا اعتماد $p=2, n=30$ مع $N_{samp} = 400$ وفي الحالة الثانية $p=3, n=400$ و $N_{samp} = 500$. القيم الملخصة لـ $m=500$ من التكرارات تم تسجيلها اضافة الى التحيز Bias ومتوسط مربعات الخطا mean squared error لمقدرات الموقع. علما ان:

$$Bias(\hat{\mu}) = \|\bar{\mu} - \mu\| = \left\| \left(\frac{1}{m} \sum_{k=1}^m \hat{\mu}^k \right) - \mu \right\| \quad \dots \dots \dots (12)$$

$$MSE = \frac{1}{m} \sum_{k=1}^m \|\hat{\mu}^k - \mu\|^2 \quad \dots \dots \dots (13)$$

حيث $\hat{\mu}^k$ هو تقدير الموقع لعينة المحاكاة k^{th} والمعلمة الحقيقية هي $\mu = 0$. لقياس انحراف التسطح deviation from sphericity لمصفوفة التشتت المقدرة $\hat{\Sigma}^{th}$ للعينة k^{th} ، يتم حساب

$$\phi_k = \frac{trac(\hat{\Sigma}^k / p)^p}{\det(\hat{\Sigma}^k)} \quad \dots \dots \dots (14)$$

والجدول (1) ادناه يسجل قيم $median In \phi_k$ نعرف ان المصفوفة $\hat{\Sigma}, \tilde{\Sigma}$ تختلف بعامل واحد فقط، لذلك

لها نفس انحراف التسطح.

اخيرا، القيمة المتوسطة لدالة الهدف (1) لـ (m) من التكرارات تم ادراجها ايضا ضمن الجدول (1) والذي نلاحظ من خلاله ان تعديل (MVB) يخفض دالة الهدف لـ (MVE) مقارنة بخوارزمية (MVE) ذات المجموعة الجزئية (p+1) الاصلية وخاصة عندما يتم التعديل في كل خطوة. في الحقيقة، بواسطة

الخوارزمية نعرف ان $(\tilde{\mu}, \tilde{\Sigma})$ دائما تنتج قيم اوطا لدالة الهدف عن $(\tilde{\mu}, \tilde{\Sigma})$ والتي بتكرارها تنتج قيم

اوطا عن $(\tilde{\mu}, \tilde{\Sigma})$. نلاحظ من جدول (1) ايضا ان التحيز ومتوسط مربعات الخطا اضافة الى

$median_k In \phi_k$ يبقى نفسه. حتى عندما يتم التعديل في كل خطوة فاننا نحصل على الدقة بصعوبة. بينما

وقت الحسابات يزداد فجأة.

جدول رقم (1)

التحيز ومتوسط مربعات الخطا لمقدرات الموقع مع وسيط $In\phi_k$ لمقدرات التشتت اضافة الى متوسط

قيمة دالة الهدف (MVE) المحسوبة باستخدام تعديل الموقع (MVB) .

		البيانات الطبيعية			20% البيانات الملوثة		
		$(\hat{\mu}, \hat{\Sigma})$	$(\tilde{\mu}, \tilde{\Sigma})$	$(\tilde{\tilde{\mu}}, \tilde{\tilde{\Sigma}})$	$(\hat{\mu}, \hat{\Sigma})$	$(\tilde{\mu}, \tilde{\Sigma})$	$(\tilde{\tilde{\mu}}, \tilde{\tilde{\Sigma}})$
p=2 n=30	$Bias(\hat{\mu})$	0.0079	0.0079	0.0077	0.028	0.015	0.018
	$MSE(\hat{\mu})$	0.234	0.229	0.231	0.259	0.228	0.252
	$med_k In\phi_k$	0.594	0.586	0.594	0.408	0.364	0.408
	$Ave_k obj_k$	1.016	0.878	0.985	1.508	1.293	1.453
p=3 n=40	$Bias(\hat{\mu})$	0.018	0.023	0.023	0.037	0.035	0.028
	$MSE(\hat{\mu})$	0.295	0.266	0.319	0.310	0.301	0.353
	$med_k In\phi_k$	0.921	0.786	0.921	0.702	0.719	0.702
	$Ave_k obj_k$	2.150	1.897	2.108	2.991	2.616	2.885

دراسة المحاكاة الثانية تمت لمشابهة النتائج اعلاه لشكل او اطار الانحدار ، حيث تم مقارنة مقدر المربعات المشدبة الصغرى (LTS) مقارنة مع او بدون تعديل التقاطع. في حالة وحيد المتغير (LTS) بحسب بالضبط من خلال خوارزمية

(Rousseeuw & Leroy 1987) . تم توليد ثلاث حالات مختلفة ، في الاولى النموذج معطى بالشكل:

$$\gamma_i = \beta_1 \chi_{i1} + \beta_2 \chi_{i2} + \beta_3 \chi_{i3} + \alpha + e_i$$

حيث $i = 1, 2, \dots, 40$ مع $\beta_1 = \beta_2 = \beta_3 = \alpha = 1$ حيث $e_i \approx N(0, 1)$ والمتغيرات التوضيحية

مولدة باستقلال (غير معتمدة) بالشكل $X_{ij} \approx N(0, 10)$ لكل $j = 1, 2, 3$. في الحالة الثانية تبديل اول (8)

نقاط بالشواذ في الاتجاه y (20% من التلوث) اي ان $e_i \approx N(10, 1)$ لكل $i = 1, 2, \dots, 8$. وفي

الحالة الثالثة تبديل بالشواذ في اتجاه X اي ان $X_{ij} \approx N(100, 10)$ مع $e_i \approx N(0, 1)$ لكل i

$i = 1, 2, \dots, 8$. التحيز ومتوسط مربعات الخطا المعرفة في (12) و (13) سجلت لمعالم الانحدار

ومعالم التقاطع باستخدام خوارزمية المجموعة الجزئية (p) بدون تعديل التقاطع ، ومع تعديل التقاطع في

كل خطوة . واخيرا باستخدام تعديل التقاطع فقط في الخطوة النهائية.

من جدول (2) نرى ان تعديل التقاطع في الحقيقة يخفض دالة الهدف لـ (LTS) خصوصا اذا اجري

التعديل في كل خطوة ، حيث التحيز و (MSE) للمعاملات لا تتغير معنويا. هذا يؤيد نتائج المحاكاة للموقع

والتشتت في التجربة الاولى الخاصة بمتعدد المتغيرات (جدول (1)).

جدول رقم (2)
التحيز ومتوسط مربعات الخطا لمقدرات الميل والتقاطع مع متوسط قيمة دالة الهدف لـ (LTS) المحسوبة
بتعديل التقاطع لانحدار (LTS).

	البيانات القياسية			20% شواذ عمودية			20% شواذ أفقية		
	$(\hat{\beta}, \hat{\alpha})$	$(\tilde{\beta}, \tilde{\alpha})$	$(\tilde{\tilde{\beta}}, \tilde{\tilde{\alpha}})$	$(\hat{\beta}, \hat{\alpha})$	$(\tilde{\beta}, \tilde{\alpha})$	$(\tilde{\tilde{\beta}}, \tilde{\tilde{\alpha}})$	$(\hat{\beta}, \hat{\alpha})$	$(\tilde{\beta}, \tilde{\alpha})$	$(\tilde{\tilde{\beta}}, \tilde{\tilde{\alpha}})$
$10^2 \times Bias(\hat{\beta})$	0.142	0.079		0.331	0.262	0.331	0.125	0.172	0.125
$10^2 MSE(\hat{\beta})$	0.463	0.460	0.463	0.422	0.430	0.422	0.414	0.385	0.414
$10^2 Bias(\hat{\alpha})$	1.677	1.749	1.70	0.227	0.275	0.786	2.010	1.428	2.720
$MSE(\hat{\alpha})$	0.129	0.124	0.126	0.134	0.123	0.122	0.020	0.014	0.027
$Ave_k obj_k$	0.810	0.781	0.799	1.195	1.148	1.174	1.194	1.14	1.174
								7	

6. تعديل L_1 Adjustment

مما تقدم نرى ان تعديل (MVB) يخفض قيمة دالة الهدف لـ (MVE) لكن هذا التعديل ليس من الضروري ان يزيد من كفاءة المقدر للعينة النهائية لكنه يبقى يستلزم وقت اطول في الحساب ، لهذا يتم ابدال تعديل (MVB) بتعديل اخر يستلزم وقت اقصر في الحساب . مقدر الموقع L_1 يكون ملائم طالما له نفس خصائص مقدر (MVB) بالنسبة لخاصية تساوي التباين ونقطة الانهيار العالية اضافة الى كفاءته الاحصائية المحسوبة بوقت اقصر. لمجموعة البيانات المعطاة ذات p من الابعاد $Y = \{y_1, y_2, \dots, y_n\}$ فان مقدر L_1 أي $\mu_L(Y)$ يكون حل لمشكلة التقليل:

$$\mu_L = \arg \min_{\mu \in R^p} \sum_{i=1}^n \|y_i - \mu\| \quad \dots \dots \dots (16)$$

علما ان الخوارزمية الاسرع لمقدر L_1 مقدمة من قبل (Hossjer & Croux 1995) . الان يمكن عمل تعديل الموقع من خلال مقدر L_1 ولاجل ذلك نطبق الخوارزمية المقدمة في المبحث الثالث وتبديل الخطوة d بالاتي :

d. نحسب تقدير L_1 ، $\theta_j = \mu_L(Y)$ ونحسب كذلك $R_j = \text{median}\|y_i - \theta_j\|$. أي شيء غير ذلك يبقى نفسه. بالتاكيد ممكن حساب المقدرات سوية في حالة البحث في كل المجموعات الجزئية ذات الحجم $(p+1)$ من خلال N_{samp} وتسحب عشوائيا المجموعات الجزئية $(p+1)$. كذلك ممكن وضع تعديل الموقع

L_1 مرة واحدة في النهاية . هذا يقود الى مقدر الخطوتين two – stage estimator المعروف بـ :

نحسب مقدرات (MVE) أي $(\hat{\mu}, \hat{\Sigma})$ وتبدل بعدها $\hat{\mu}$:-

$$\tilde{\mu} = \arg \min_{\mu \in R^p} \sum_{i=1}^n \|x_i - \mu\| \quad \dots\dots\dots (17)$$

نقطة الانهيار والخصائص التقريبية لهذا المقدر مقدمة من قبل (Hossjer & Croux 1995) علما ان كل هذه النسخ لـ (MVE) مع تعديل الموقع L_1 تبقى متساوية التباين . ولغرض تطبيق المقدر تم اعادة دراسة المحاكاة الاولى المقدمة في المبحث الخامس لمقارنة خوارزمية المجموعة الجزئية (p+1) العشوائية لـ (MVE) محسوبة مع او بدون تعديل الموقع L_1 . الجدول (3) ادناه يوضح النتائج التي تم التوصل اليها، وكما هو متوقع فان هذا النوع من تعديل الموقع لا يقلل دالة الهدف لـ (MVE) كذلك فان تأثير تعديل L_1 على التحيز $Bias(\hat{\mu})$ ومتوسط مربعات الخطأ $MSE(\hat{\mu})$ وكذلك $median_k In\phi_k$ تبقى نفسها عندما يجري التعديل في كل خطوة او فقط في النهاية. النسخة الاخيرة تعرف بـ $(\tilde{\mu}, \tilde{\Sigma})$ وهي الاسرع في الحساب.

جدول رقم (3)

التحيز ومتوسط مربعات الخطأ لمقدرات الموقع مع وسيط $In\phi_k$ لمقدرات التشتت ومتوسط قيمة دالة الهدف لـ (MVE) المحسوبة باستخدام تعديل L_1

		البيانات الطبيعية			20% البيانات الملوثة		
		$(\hat{\mu}, \hat{\Sigma})$	$(\tilde{\mu}, \tilde{\Sigma})$	$(\tilde{\mu}, \tilde{\Sigma})$	$(\hat{\mu}, \hat{\Sigma})$	$(\tilde{\mu}, \tilde{\Sigma})$	$(\tilde{\mu}, \tilde{\Sigma})$
p=2 n=30	$Bias(\hat{\mu})$	0.0079	0.0062	0.0034	0.0211	0.4082	0.4015
	$MSE(\hat{\mu})$	0.234	0.093	0.094	0.259	0.277	0.272
	$med_k In\phi_k$	0.594	0.531	0.594	0.408	0.404	0.406
	$Ave_k obj_k$	1.016	1.055	1.332	1.508	1.677	2.178
p=3 n=40	$Bias(\hat{\mu})$	0.0178	0.0036	0.0057	0.0369	0.4605	0.4601
	$MSE(\hat{\mu})$	0.295	0.092	0.093	0.310	0.328	0.328
	$med_k In\phi_k$	0.921	0.845	0.921	0.702	0.693	0.702
	$Ave_k obj_k$	2.150	2.093	2.461	2.991	3.109	3.695

من الجدول اعلاه نرى ان $median_k In\phi_k$ لمقدر مصفوفة التشتت يبقى نفسه في جميع الحالات. اما بالنسبة لمقدرات الموقع، فهي تعتمد على فيما اذا كانت البيانات طبيعية او ملوثة. للبيانات الطبيعية (غير الملوثة) نرى ان تعديل L_1 يقدم تحيز قليل وتخفيض في متوسط مربعات الخطأ طالما ان L_1 له افضل كفاءة احصائية ، لكن للبيانات الملوثة فان (MSE) تبقى نفسه بينما التحيز يصبح عالي جدا.

7. مقارنة مع خوارزمية الحل المقنع

A comparison with the feasible Solution Algorithm

مما تقدم نرى ان الحسابات المضبوطة لـ (MVE) تكون كالآتي : اعتبر كل المجموعات الجزئية الممكنة ذات الحجم (h) وتدعى بالعينات النصفية half samples ويحسب القطع لاصغر شكل بيضوي يحوي كل النقاط المحسوبة لـ (MVE) هو الشكل البيضوي الخاص بتلك العينة النصفية المثلى: طالما العدد الكلي للعينات النصفية الممكنة كبير جدا فالحسابات المضبوطة تكون غير مقنعة في التطبيق العملي، ماعدا حجم العينة الصغير جدا لذلك عادة ما يتم اعادة ترتيب الخوارزميات التقريبية. في هذا المبحث سوف نركز على خوارزمية Feasible Solution Algorithm (FSA) والمعروفة والمستخدمة في حل مقدر الانحدار اصغر وسيط للمربعات Least Median of square regression (Hawkins 1993 a).

FSA تبدأ بالاختيار العشوائي للعينة النصفية ليتم بعد ذلك ابدال النقاط في العينة النصفية مع النقاط غير العائدة لها، وهذا يقلل قيمة دالة الهدف، وهذا معناه ان القطع لاصغر شكل بيضوي يحوي كل النقاط في العينة النصفية. اذا لم نلاحظ انخفاض أكثر فالعينة النصفية الحاصل عليها يطلق عليها الحل المقنع Feasible Solution . الخوارزمية تعتمد على العدد الكلي N_{fsa} للبدائيات العشوائية. هذا العدد يكون كبير بدرجة كافية خاصة عند تعامل FSA مع مجموعة بيانات قلقة Nasily data set (يبدأ البرنامج مع $N_{fsa} = 50$).

وقت الحسابات الذي يستلزمه FSA اكثر كثيرا من خوارزمية المجموعة الجزئية (p+1). كلفة المقارنة لوقت الحسابات صعب، طالما معالم التذبذب Tuning parameters تختار من قبل المستخدم (N_{samp}, N_{fsa}, h) ، واذا كان حجم العينة والابعاد تزداد فان هذا الاختلاف في وقت الحسابات يصبح كبير جدا (علما اننا لا نتعامل مع دالة التعديل في كل خطوة). من ملاحظة الجدول (4) نرى ان FSA نجحت في ايجاد اوطا قيمة لدالة الهدف خلال المقدرات المعتمدة . حيث تم هنا اعادة دراسة المحاكاة للمبحث الخامس ومنه يمكن القول ان FSA حققت هدفها (حققت قيمة واطنة لدالة الهدف (1)).

على كل حال الفوائد الاحصائية بقيت محددة، حيث لجزء التشتت لا يوجد اختلاف معنوي مع $(\hat{\mu}, \hat{\Sigma})$ او مع $(\tilde{\mu}, \tilde{\Sigma})$ اما لجزء الموقع فهناك تحسن البسيط في (MSE) عندما $p=3$ فان هناك تحيز عالي وايضا MSE وهذا ضمن حالة التعامل بالبيانات الملوثة، وذا بسبب ان الانهيار يحصل 8 مرات في 500 تكرار. زيادة قيمة N_{fsa} الى 100 يلغي هذا الانهيار، ولكن يضاعف وقت الحسابات. عند الحصول على حل FSA لمشكلة (MVE) يمكن ملاحظة التحسن في تقدير الموقع باضافة (MVB) او تعديل (L_1) في النهاية. الجدول (4) يسجل نتائج المحاكاة لـ FSA مع تعديل الموقع باستخدام (MVB) او L_1 .

جدول رقم (4)

استخدام تعديل الموقع لـ FSA لكل الحالات محسوبة ومتوسط مربعات الخطا لمقدر الموقع ووسيط $In\phi_k$ لمقدرات التشتت اضافة الى متوسط قيمة دالة الهدف لـ (MVE)

		البيانات الطبيعية			20% البيانات الملوثة		
		FSA	+MVB	+L ₁	FSA	+MVB	+L ₁
p=2 n=30	$Bias(\hat{\mu})$	0.028	0.028	0.022	0.027	0.026	0.388
	$MSE(\hat{\mu})$	0.242	0.242	0.092	0.228	0.227	0.261
	$med_k In\phi_k$	0.606	0.606	0.606	0.385	0.385	0.385
	$Ave_k obj_k$	0.823	0.823	1.293	1.257	1.257	2.119
p=3 n=40	$Bias(\hat{\mu})$	0.006	0.006	0.011	0.411	0.411	0.784
	$MSE(\hat{\mu})$	0.248	0.248	0.098	10.09	10.14	7.580
	$med_k In\phi_k$	0.973	0.973	0.973	0.682	0.682	0.682
	$Ave_k obj_k$	1.519	1.519	2.241	2.365	2.365	3.752

المدهش انه لا يوجد اختلاف نشاهده بين حلول FSA و FSA+MVB . في الحقيقة كلا الاسلوبين غالبا (ولكن ليس دائما) نفس النتيجة . اذا FSA وجدت MVE المضبوطة، عندها تعديل (MVB) لا يغير التقدير. اختبار فيما اذا تعديل (MVB) يغير مقدر الموقع يكون لذلك غير ضروري لكن غير كافي وهو شرط ايجاد (MVE) المضبوط. استخدام تعديل (L_1) يفرز قيمة مقاربة الى $(\tilde{\mu}, \tilde{\Sigma})$ في جدول (3): زيادة قيمة دالة الهدف، لكن الكفاءة الاحصائية الافضل مفاصة بواسطة MSE.

8. النتائج

في هذا البحث لاحظنا ان مقدر اصغر كروي (MVB) ممكن استخدامه لتعديل الموقع لمقدر اصغر بيضوي (MVE) . هذا التعديل غالبا ما يقلل دالة الهدف في (MVE) ولكن له تاثير قليل على التحيز Bias ومتوسط مربعات الخطا MSE للمتوسط $\hat{\mu}$. من جانب اخر فان تعديل الموقع بالاعتماد على (L_1) ليس ضروريا ان يقلل دالة الهدف لـ (MVE) ولكنه يبرهن على كفاءة $\hat{\mu}$ للبيانات الطبيعية Normal data.

لغرض تقليل تاثير الشواذ المتطرفة على تعديل (L_1) ممكن استخدام مقدر (L_1) معاد الوزن Reweighted (L_1) estimator بدلا عنه وكما في حالة مقدر (L_1) غير معاد الوزن، فان هذا النوع من التعديل لا يخفض دالة الهدف لـ (MVE) بينما التحيز للمقدر الناتج يكون مشابه لما هو في مقدر (MVE)، بينما MSE له يكون افضل حتى عند التعامل بالبيانات الملوثة. في المبحث الخامس لاحظنا ان تعديل الموقع بـ (LTS) يخفض قيمة دالة الهدف ولكن ليس له تاثير كبير على متوسط مربعات الخطا للمعاملات لكنه يفرز مقدر تقاطع بتحيز واطى وكفاءة اكثر. الغاية اذا من استخدام تعديل (L_1) هي في انه لا يكلف وقت حسابات عالي، يملك منحني تحيز واطى ويعطي مقدرات اكثر كفاءة لمعلمة الموقع في حالة البيانات الطبيعية متعددة المتغيرات.

استخدام تعديل (MVB) او ما يسمى بـ Feasible Solution Algorithm الموضحة في المبحث الثامن يعطي قيم اوطا لدالة الهدف المتعلقة بمقدرات اصغر قطع بيضوي (MVE) ولكن له حسابات عالية، علما ان اوطا القيم لدالة الهدف لا تولد بالضرورة دقة احصائية عالية. اخيرا من الضروري الاشارة الى ان هناك خوارزميات عديدة وجديدة لحسابات مقدرات عالية الانهيار، كالخوارزمية المعتمدة لحساب مقدر اصغر محدد تباين مشترك Minimum Covariance determinant (MCP) ولكن حساب مقدر (MVE) يبقى هو الاصعب.

المصادر

1. العلوي، لقاء علي، (مقارنة مقدرات التباين المشترك الحصينة في تحليل المركبات الرئيسية) رسالة دكتوراه- كلية الادارة والاقتصاد- جامعة بغداد.
2. Croux,C.,& Haesbroeck, (2002), "Finite – sample efficiencies of Estimators for the Minimum volume Ellipsoid", Journal of statistical computing and simulation, 72,585-596.
3. Croux, C., Van Aelst, S., & Dehon. C., (2003), " Bounded Influence Regression using High Breakdown Scatter Matrix", Annals of the Institute of statistical Mathematics, 55,256-285.
4. Hawkins, D.M.(1993a), (the Feasible set Algorithm for Least Median of squares Regression", Computational statistics and data Analysis, 16,81-101.
5. Hossjer, O., & Corux, C.(1995), "Generalizing univariate signed Rank statistics for Testing and estimating a multivariate Location parameter", Nan parametric statistics, 4,293-308.
6. Rousseeuw, P.J.(1985), "Multivariate Estimation with High Breakdown point", in Mathematical statistics and Application, VolB, eds.W.Grossmann, G. Pflug, I. Vincze and w.werts, (Dordrecht:Reidel), 283-297.
7. Rousseeuw, P.J. & Leroy, A.M.(1987), "Robust Regression and outlier Detection", New york: John wiley.