

مقارنة مقدرات M الحصينة مع مقدرات شرائح التمهيد التكعيبية لأنموذج المعاملات المتغيرة زمنياً للبيانات الطولية المترنة

الباحث علي سيف الدين عبد الحافظ

أ. د. ظافر حسين رشيد
كلية الادارة والاقتصاد - جامعة بغداد
قسم الاحصاء

المستخلاص

في هذا البحث تم مقارنة مقدرات (M) الحصينة لتقنية شرائح التمهيد التكعيبية لتلافي مشكلة الشواذ في البيانات أو تلوث الخطأ مع طريقة التقدير التقليدية لتقنية شرائح التمهيد التكعيبية ، باستعمال معيارين للمفاضلة بينهما هما (MADE) و لمختلف حجوم العينة ومستويات التباين، وذلك لتقدير دوال المعاملات المتغيرة زمنياً للبيانات الطولية المترنة، والتي تتصف بكون المشاهدات يتم الحصول عليها من n من القطاعات المستقلة كل واحد منها يقاس تكرارياً خلال مجموعة نقاط زمن محددة m ، أذ تكون القياسات المكررة داخل القطاعات مرتبطة على الأغلب ومستقلة بين القطاعات المختلفة.

المصطلحات الرئيسية للبحث/ المعاملات المتغيرة زمنياً، تقدير المرحلتين، تقديرات M الحصينة، تمهيد الشرائح التكعيبية، البيانات الطولية المترنة.



مجلة العلوم

الاقتصادية والإدارية

المجلد 19

العدد 73

الصفحات 413-398

بحث مستقل من أطروحة دكتوراه

المقدمة : 1

في الدراسات الطويلة المشاهدات غالباً يتم الحصول عليها من n من القطاعات المستقلة كل واحداً منها يقاس تكرارياً خلال مجموعة نقاط زمن محددة ، غالباً ما يتركز اهتمام هذه الدراسات على تقييم آثار الزمن (t) وكذلك مجموعة المتغيرات المستقلة $\chi_r(t)$ ، $r=1,2,\dots,d$ على نتيجة المتغير المعتمد $y(t)$ ، نفرض أن (t_{ij}) تمثل الزمن لقياسات χ_r للقطاع i^{th} ، وأن y_{ij} و χ_{ij} تمثل مشاهدات القطاع i^{th} للمتغير المعتمد والمستقل عند الزمن t_{ij} على التوالي، فأن مجموعة المشاهدات الطويلة تعطى كالتالي :-

حيث t_i هي عدد القياسات المتكررة للفيصل i^{th} ، على الرغم من أن القياسات هي مستقلة بين القطاعات المختلفة إلا أنها على الأغلب تكون مرتبطة داخل كل قطاع.

التحليل الأحصائي مع هكذا نوع من البيانات مهم بنمذجة منحنى المتوسط $(t)y$ والتأثيرات للمتغيرات المستقلة على $(t)y$ ، وتطوير التقدير وأجراءات الاستدلال ، وتحت إطار النماذج المعلميه مثل النماذج الخطية وغير الخطية ونماذج ذات التأثيرات المختلفة ، درست نظريات وطرائق التقدير للمعلم والاستدلالات بصورة موسعة، وتحت إطار النماذج اللامعلميه مع نقاط زمن تصميم ثابتة (Hart(1991) اعتمد طرائق (kernel) لتقدير التوقع $((t)y)$ بدون وجود المتغيرات المستقلة ، ولأخذ المتغيرات المستقلة بالحساب (Zeger and Diggle (1994) درسا الانموذج شبه المعلمي التالي:

$$Y_{ij} = \mu(t_{ij}) + X'_{ij} B + \varepsilon(t_{ij}) \dots \dots \dots (2)$$

حيث $B = (B_1, B_2, \dots, B_d)$ هو متجة ثوابت غير معلوم في R^d ، $\mu(t)$ هي دالة ممهدة محددة الى (t) ، $\varepsilon(t)$ عملية عشوائية بمتوسط (0) ، وحققوا اجراءات تكرارية حيث افترضوا اجراء (t) الذي يقدر في البداية $\mu(t)$ بطريقة (kernel) وثم تكرار التقديرات لـ β و $\mu(t)$ ، وأن الأمودج في (2) هو أكثر مرونة من النماذج الخطية التقليدية ، ويطلب لتأثيرات المتغيرات أن تكون ثابتة خلال الزمن ، هذا التقييد يمكن أن لا يكون واقعياً للكثير من الحالات العملية ، بعبارة أخرى حجوم العينة الفعلية فيأغلب الدراسات الطولية يمكن أن لا يكون حجمة كافية لدعم كامل للأمودج الامامي العام عندما تكون المتغيرات المستقلة ذات بعد عالي وهو ما يسمى مشكلة البعدية (Curse of Dimensionality) ، لذلك ولأجل تعليم عملي أكثر للأمودج (2) ، (Hoover et al 1998) أعتمدوا أنواع المعاملات المتغيرة التالية:

$$Y(t) = X'(t) B(t) + \varepsilon(t) \dots \dots \dots (3)$$



وأقترحوا صنف متعدد الحدود الموضعي (kernel) لتقدير $B(t)$ ، حيث $B(t) = (B_1(t), B_2(t), \dots, B_d(t))'$ هو متجة لدوال الممهدة خلال (t) ، $\varepsilon(t)$ كما معرفة في (2) ولجميع قيم (t) فأن $X(t)$ و $Y(t)$ مستقلان، وبشكل عام نلاحظ أن الأنماذج (3) هو أنماذج خطى بين $X(t)$ و $Y(t)$ عند كل زمن ثابت (t) ، أجراءات التقدير في (Hoover et al 1998) طورت الحالة الخاصة لمقدرات (kernel) والتي اعتمدت على عرض حزمة (pandwidth) واحد ومقدرات (spline) والتي اعتمدت على أكثر من معلمة تمهد، وأعتمدوا على خوارزمية (backfitting) لحل مشكلة تعدد الأبعاد وخاصة لتقنية الشرائح التمهيدية ، والتي عانت من بعض المشكلات وهي كثافة الحسابات والجهد البرمجي، وكعلاج للمشكلات السابقة (Fan and Zhang 2000) أقترحنا أجراء المرحلتين (Two Step) كبديل (Raw Estimate) إلى دوال المعالم $B(t)$ في (3) ، أو لا أحسبوا التقديرات الخام (Smooth Estimate) للحصول على المقدرات التمهيدية لدوال المعاملات بواسطة استعمال أحدى تقنيات التمهيد المعروفة، واستعملوا شرائح التمهيد كأحدى تلك التقنيات.

أن تقديرات دوال المعاملات (t) بأسلوب المرحلتين يتم الحصول عليه باستعمال طريقة المربعات الصغرى، وكما هو معلوم أن مقدرات الربعات الصغرى تمتلك بعض الخصائص الجيدة، وخاصة عندما الخطأ العشوائي يتبع التوزيع الطبيعي، ولكن المقدرات المعتمدة على المربعات الصغرى حساسة جداً إلى الشوائب في البيانات أو عند تلوث (contamination) الخطأ؛ لذلك فإن طرائق التقدير الحصينة مطلوبة أكثر، في هذا البحث سيتم الاعتماد على البيانات الطولية المتزنة (عندما تكون عدد القياسات للقطاعات متساوية وهي m) والمصاغة وفق الأنماذج (3) ، لأيجاد تقديرات شرائح التمهيد التكعيبية الحصين لدوال المعاملات بطريقة المرحلتين، وبالأعتماد على أسلوب M الحصين، ومقارنته مع طرائق التقدير التقليدية عن طريق تجربة محاكاة بنسبي تلوث مختلفة وحجمع عينة مختلفة ومستويات تباين مختلفة ولأغراض الملاعنة والتعريم، تم عرض كل الصيغ والمعادلات بدلالة d من المتغيرات التوضيحية، على الرغم من إستعمال دراسة المحاكاة لحالة ثنائية للمتغيرات (متغيرين فقط).

2. طريقة تدبير المرحلتين (Two-Step Estimation Method) $(1), (2), (10)$

لنفترض $t_r, r=1, 2, \dots, m$ هي نقاط زمن محددة، حيث تم جمع البيانات ، إذ ان m تمثل عدد القياسات المكررة لكل قطاع، لأن هناك عدد من المشاهدات التي جمعت في الزمن t_r ، فمن الممكن لهذا الثابت t_r أستعمال البيانات المجمعة هناك لمطابقة أنماذج (3) والحصول على المقدرات الخام (Raw estimates)

$$b(t_j) = (b_1(t_j), \dots, b_d(t_j))'$$

هذه هي المرحلة الأولى، عادةً المقدرات الخام هي غير ممهدة تحتاج إلى تمهيدها للحصول على المقدرات الممهدة إلى دوال المعاملات لذلك، في المرحلة الثانية لكل مركبة معطاة $r=1, 2, \dots, d$ نطبق تقنية تمهد إلى البيانات $\{b(t_j)\}_{j=1}^m$ ، وإن مرحلة تقديرات التمهيد (Smoothing estimates) هذه حاسمة لأنها تعطي مقدرات تمهدية لدوال معاملات التمهيد الأساسية، وإضافة إلى ذلك فإن مرحلة التمهيد ذو بعد واحد (one-dimensional) .



1.2 مرحلة التقديرات الخامسة (Raw Estimates Step)

ولتوضيح هذه المرحلة نفترض $t_j = 1, 2, \dots, m$ ، y_{ij} هي نقاط زمن محددة لمجموعة البيانات الطولية لكل نقطة زمن j ، لنفترض N_j هو مجموعة الفهارس للقطاع الى جميع مشاهدات y_{ij} عند j ، نجمع كل χ_{ij} و y_{ij} التي تقابل الفهرس للقطاعات في N_j ونشكل مصفوفة التصميم \tilde{X}_j ومتجه الاستجابة \tilde{Y}_j بالتتابع، الجدير بالاشارة بان البحث يعتمد على حالة البيانات المتزنة وبذلك فان N_j مجموعة الفهارس للقطاع الى جميع مشاهدات y_{ij} عند j ستكون متساوية ولجميع القطاعات. عندها فان صيغة أنموذج (3) عندما البيانات تجمع عند الزمن يتبع الأنموذج الخطى الآتى :

$$\left. \begin{aligned} \tilde{Y}_i(t_j) = & \tilde{X}_i(t_j) B(t_j) + \tilde{e}_i(t_j) \\ & j=1, 2, \dots, m \end{aligned} \right\} \dots \quad (4)$$

١.١.٢ مقدرات المربعات الصغرى العامة المقبولة مع اخطاء AR(1) (Feasible GLS Estimation With AR(1) Errors)

لتقدير معلم الأنماذج (4) تحت افتراض هيكل الارتباطات للأخطاء يتبع (1) AR كالتالي :

يمكننا تطبيق المربعات الصغرى العامة حيث مقدراتها ستكون كالتالي :

$$b_{GLS}(t_j) = \left(X_i'(t_j) \Omega^{-1} X_i(t_j) \right)^{-1} X_i'(t_j) \Omega^{-1} \tilde{Y}_i(t_j) \quad \dots \quad (6)$$

ان المشكلة في مقدرات GLS بأنها تفترض مصفوفة التباين المشترك Ω معلومة، بعبارة أخرى إن ρ معلوم وهذا نادراً ما هو معلوم من الناحية العملية.

ولتجنب افتراض GLS علينا ايجاد تقدير متسق لـ $\hat{\Omega}$ أي $\hat{\theta}$ واستعماله لايجاد تقدير لمعالم الانموزج (4)، إن هذا المقدر يدعى مقدرات المربيعات الصغرى العامة المقبولة (FGLS) والذي يمكن ايجاده بحسب الخطوات الآتية :

a. نجد اولاً مقدرات المربعات الصغرى الاعتيادية الى أنموزج (4) والذي سيكون كالتالي :

$$b_{OLS}(t_j) = \left(\tilde{X}'_i(t_j) \tilde{X}_i(t_j) \right)^{-1} \tilde{X}'_i(t_j) \tilde{Y}_i(t_j) \quad \dots \dots \quad (7)$$

ثم نجد الاخطاء باستعمال مقدرات OLS .



إذ إن :

$$\tilde{e}_i(t_j) = \tilde{Y}_i(t_j) - \tilde{X}_i(t_j) b_{OLS}(t_j) \quad \dots \dots .8)$$

b. وتحت افتراض الاخطاء هي عمليات عشوائية مشتركة فان مقدر ρ المشترك يمكن تقديره كالتالي :

$$\hat{\rho} = \frac{\sum_{i=1}^n \sum_{j=2}^m \tilde{e}_i(t_j) \tilde{e}_i(t_{j-1})}{\sum_{i=1}^n \sum_{j=1}^m \tilde{e}_i^2(t_j)} \quad \dots \dots .9)$$

c. إجراء تحويل للبيانات باستعمال **(Prais – Winsten) transformation**.

$${}^* Y_i(t_j) = \begin{bmatrix} \sqrt{1 - \hat{\rho}^2} \tilde{y}_i(t_1) \\ \tilde{y}_i(t_2) - \hat{\rho} y_i(t_1) \\ \tilde{y}_i(t_3) - \hat{\rho} y_i(t_2) \\ \vdots \\ \tilde{y}_i(t_m) - \hat{\rho} \tilde{y}_i(t_{m-1}) \end{bmatrix}, \quad i = 1, 2, \dots, m \quad \dots \dots .10)$$

$${}^* X_i(t_j) = \begin{bmatrix} \sqrt{1 - \hat{\rho}^2} \tilde{x}_i(t_1) \\ \tilde{x}_i(t_2) - \hat{\rho} \tilde{x}_i(t_1) \\ \tilde{x}_i(t_3) - \hat{\rho} \tilde{x}_i(t_2) \\ \vdots \\ \tilde{x}_i(t_m) - \hat{\rho} \tilde{x}_i(t_{m-1}) \end{bmatrix}, \quad i = 1, 2, \dots, m \quad \dots \dots .11)$$

d. وبتطبيق المربعات الصغرى الاعتيادية على البيانات المحولة فاننا نحصل على مقدرات FGLS كالتالي :

$$b_{FGLS}(t_j) = \left({}^* X'_{\cdot i}(t_j) {}^* X_{\cdot i}(t_j) \right)^{-1} {}^* X'_{\cdot i}(t_j) {}^* Y_{\cdot i}(t_j) \quad \dots \dots .12)$$

وان مقدرات FGLS هي تقاريباً أكثر كفاءة من مقدرات OLS عندما الاخطاء تتبع هيكل ارتباط AR(1).

Robust Raw Estimate 2.1.2 مقترح التقديرات الخام الحصينة

اقتصر أسلوب M الحصين او لاً من قبل (Huber)⁽⁶⁾، وتستند الفكرة ببساطة الى تقليل بعض الدوال للأخطاء بدلاً عن مجموع المربعات لها، والمقدر الحصين يحدد عن طريق الاختيار لدالة وزن، وان اسلوب مقدرات M بحاجة الى بعض التوسيع لتطبيقها على البيانات الطولية المتزنة الموصوفة في أنموذج (4) والتي تحتوي على n من القطاعات و m من القياسات المكررة لكل قطاع.



لأنموذج المعاملات المتغيرة زمنياً للبيانات الطولية المتزنة
إذ إن أهم ما يميز هذه البيانات هو ارتباطها ضمن القطاع، أي بمعنى أن الأخطاء مرتبطة، وهذا سينافي الافتراض لأسلوب مقدرات M وهو أن تكون الأخطاء غير مرتبطة، ولتلafi هذه المشكلة سيتم الاعتماد على البيانات المحولة في طريقة (Feasible GLS).

لنفرض أن $(Y_{ij}^* = \beta^* X_{ij}^* + \epsilon_{ij}^*)$ و $(\beta^* = \beta(t_j))$ ، ولتوسيع هذا الأسلوب، ان المربعات الصغرى الاحتيادية للبيانات المحولة تخفض مجموع مربعات الخطأ إلى أقل ما يمكن كالتالي :

$$\min \sum_{i=1}^n \sum_{j=1}^m e_{ij}^{*2} = \min \sum_{i=1}^n \sum_{j=1}^m \left(Y_{ij}^* - \beta^* X_{ij}^* \right)^2 \quad \dots \dots \dots .(13)$$

إذ إن :

y_{ij}^* : المشاهدة j للقطاع i للمتغير المعتمد.

X_{ij}^* : الصنف ij لمصفوفة التصميم X .

وإن :

$$X_{ijk}^* = \begin{bmatrix} * \\ \chi_{ij1}, \dots, \chi_{ijd} \end{bmatrix}, k = 1, 2, \dots, d$$

. β : متوجه معلم ذو بعد $dm * 1$

ان تقديرات M المطورة من قبل (Huber) والتي لها خاصية (Scale Invariant) تعتمد على فكرة ابدال مجموع مربعات الأخطاء e_{ij}^{*2} بدالة اخرى للأخطاء الهدف منها تقليل المقدر الآتي:

$$\sum_{i=1}^n \sum_{j=1}^m P \left(\frac{\left(y_{ij}^* - X_{ij}^* \beta^* \right)^2}{\sigma_e^2} \right) \quad \dots \dots \dots .(14)$$

وإن P هي دالة محدبة متتماثلة (symmetric convex function) وتقليل المقدار اعلاه تؤخذ المشتقة بالنسبة الى متوجه المعلم وجعلها مساوية الى الصفر كالتالي :

$$\sum_{i=1}^n \sum_{j=1}^m X_{ij}^* \Psi \left(\frac{\left(Y_{ij}^* - X_{ij}^* \beta^* \right)}{\sigma_e^2} \right) = 0 \quad \dots \dots \dots .(15)$$



إذ إن :

Ψ : المشتقة الجزئية إلى متوجه المعالم β للدالة P و $\Psi = P'$ وبهذا يكون هناك ($d m$) من المعادلات غير الخطية والتي يمكن حلها بعدة طرائق منها طريقة المربعات الصغرى الموزونة تكرارياً، (IWLS) ولإيجاد مقدرات M بالاعتماد على طريقة IWLS يتطلب حساب دالة الوزن وفيها يتم اعادة كتابة الصيغة (15) كما يلي :

$$\sum_{i=1}^n \sum_{j=1}^m W_{ij} \mathbf{X}_{ij}^* \left(\frac{\left(\mathbf{Y}_{ij}^* - \mathbf{X}_{ij}^* \boldsymbol{\beta} \right)}{\sigma_e^2} \right) = \mathbf{O} \quad \dots \dots \dots (16)$$

وبحل المعادلة اعلاه نحصل على تقديرات أسلوب M الحصين b_μ باستعمال IWLS .

إذ إن :

$$b_M = \left(\mathbf{X}' \mathbf{W} \mathbf{X} \right)^{-1} \mathbf{X}' \mathbf{W} \mathbf{Y} \quad \dots \dots \dots (17)$$

إذ إن :

b_M : متوجه ذو بعد ($dm * 1$)
 W : مصفوفة اوزان قطرية ببعد ($nm * nm$) تحسب عناصرها كما يأتي :

$$W_{ij} = \frac{\Psi \left(\frac{\left(\mathbf{y}_{ij}^* - \mathbf{X}_{ij}^* \boldsymbol{\beta} \right)}{\hat{\sigma}_e} \right)}{\Psi \left(\frac{\left(\mathbf{y}_{ij}^* - \mathbf{X}_{ij}^* \boldsymbol{\beta} \right)}{\hat{\sigma}_e} \right)} \quad \dots \dots \dots (18)$$

إذ إن :

$$\hat{\sigma}_e = 1.483 \left[Median \left| e_{ij}^* - Median(e_{ij}^*) \right| \right]$$

وتم استعمال دالة Ψ ومشتقها التالية :

دالة (Andrews)

$$\Psi(e_{ij}^*) = \begin{cases} A \sin\left(\frac{e_{ij}^*}{A}\right) & \left| e_{ij}^* \right| \leq A\pi \\ 0 & \left| e_{ij}^* \right| > A\pi \end{cases}$$

$A = 1.339$



2.2 مرحلة تحسين او تمهيد المقدرات الخام⁽¹⁾ (Refining Or Smoothing The Raw Estimates)

و قبل البدء بهذه المرحلة فان مركبات دوال المعاملات المقدرة في المرحلة الاولى نحصل عليها كالتالي :
 ولـ $r=1, 2, \dots, d$ نفرض أن :

$b_r(t_j)$ يحتوي على r^{th} من المركبات لتقديرات المرحلة الاولى $b_{(t_j)}$ عندئذ

$$b_r(t) = L'_r \begin{pmatrix} * & * \\ X' & X \end{pmatrix}^{-1} \begin{pmatrix} * & * \\ X' & Y \end{pmatrix} \quad \dots \dots \dots \text{19)$$

اڑ ان :

L_r : يعرف كمتوجه وحدة ذو بعد $(dm)^{th}$ فيه r من المدخلات (1) والباقي (0).
ومن المعلوم إن

$$E[b_r(t)] = \beta_r(t) \quad \dots \dots \text{QO}$$

وبذلك سيصبح التمهيد عند كل مركبة لـ (٢) ذو بعد واحد، حيث سنستعمل تقنية تمهيد شرائح التمهيد التكعيبية CSS، ولكن للبيانات التالية :

اڈ ان :

t_j : هي نقاط زمن التصميم.

$b_r(t_j)$: هي الاستجابة عند نقاط زمن التصميم، مع التأكيد بأنها اي تقدير مستخرج من المرحلة الاولى.
ان نموذج الانحدار الامثلسي البسيط للبيانات السابقة سيكون كالآتي :

ونحن نريد تقدير الدالة الممهدة (

اڈ ان :

$f(t_j)$: هي اخطاء القياسات التي لايمكن شرحها بواسطة دالة الانحدار
 $b_r(t_j) \setminus t_j = t$ هي التوقع الشرطي لـ $f(t)$ رياضياً أي

$$f(t) = E(b_r(t_j) \mid t_j = t)$$



1.2.2 شرائح التمهيد التكعيبة

(Cubic Smoothing Splines) :

لإيجاد مطابقة شرائط التمهيد للألموزج (22) وبدون فقدان التعيم لو أفترض مدى الفترة f في الألموزج (22) هي فترة منتهية $[a, b]$ ولبعض الأعداد المنتهية a, b .
 فإن الجزء غير الممهد (Roughness Penalty) لـ f هو عادة يعرف كتكامل لمربع مشتقه (V) من المرات كالتالي :

$$\int_a^b \{f_h^{(V)}\}^2 dt$$

ولبعض $V \geq 1$ ، عندما مهم شرائح التمهيد لـ f في نموذج (22) يعرف كتقدير $\hat{f}_\lambda(t)$ لصيغة المربعات الصغرى الجزئية PLS التالية :

$$\sum_{j=1}^m \left[b_r(t_j) - f(t_j) \right]^2 + \lambda \int_a^b \left\{ f_{(t)}^{(V)} \right\}^2 dt \quad \dots \dots \dots \text{Eq3}$$

وخلال V^{th} من فضاء (Sobolev) المرتب $W_2^V[a, b]$ إن :

وإن اختيارنا لشريحة التمهيد التكعيبية لتقليل الصيغة (23) هو صعوبة الحصول على التكامل الذي يعرف الجزء غير الممهد فضلاً على وجود طريقة لحسابه في شريحة التمهيد التكعيبية.

وباستخدام نقاط الزمن (t_j) كعقد حيث نفترض أن $j = 1, 2, \dots, m$ ، T_j هي نقاط الزمن المحددة

$$a = T_1 < T_2 < \dots < T_m = b$$

فإن جميع العقد لشراحت التمهيد التكعيبية التي تقلل (23) عندما $V = 2$ تكون كالتالي :

$$h_L = T_{L+1} - T_L \quad , \quad L=1,2,\dots,m-1$$

نعرف ($A = (a_{LS})$) كمصفوفة ذات بعد $(m*(m-2))$ مع جميع مدخلاتها تكون (0) ما عدا عندما $L=1, 2, \dots, m-2$ فان

$$a_{L,L} = h_L^{-1}$$

$$a_{L^{\pm 1}, L} = -\left(h_{L^-}^{-1} + h_{L^{\pm 1}}^{-1}\right)$$

$$a_{I+2-I} = -h_{I+1}^{-1}$$



ونعرف C كمصفوفة ذات بعد $(m-2) \times (m-2)$ مع جميع مدخلاتها (0) ماعدا

$$c_{11} = (h_1 + h_2)/3$$

$$c_{21} = h_2/6$$

ولأن $L=1, 2, \dots, m-4$ فإن

$$c_{L,L+1} = h_{L+1}/6$$

$$c_{L+1,L+1} = (h_{L+1} + h_{L+2})/3$$

$$c_{L+2,L+1} = h_{L+2}/6$$

$$c_{m-3,m-2} = h_{m-2}/6$$

$$c_{m-2,m-2} = (h_{m-2} + h_{m-1})/3$$

وأخيراً نعرف G كمصفوفة وهي مصفوفة ذات بعد $(m \times m)$ غير الممهد التكعيبية كالتالي :

$$G = A C^{-1} A'$$

نفرض أن $f' = (f_1, \dots, f_m)'$ إذ إن

$$f_j = f(T_j) \quad , \quad j = 1, 2, \dots, m$$

وبذلك الجزء غير الممهد يمكن ان نعبر عنه كالتالي :

$$\int_a^b [f''(t)]^2 dt = f' G f \quad \dots\dots\dots \text{24}$$

لاجل ذلك : نحن فقط نشير الى G كمصفوفة غير الممهد (Roughness matrix)، وهذا يعني أن صيغة

في (23) يمكن كتابتها كالتالي :

$$\|b_r - Wf\|^2 + \lambda f' G f$$



إذا إن : $W = \begin{pmatrix} W_{jj} \\ \vdots \\ W_{1j} \end{pmatrix}$ هي مصفوفة حد ذات بعد $(m*m)$ مع $t_j = T_j$ إذا و 0 في الحالات الأخرى.

وأن $\|a\|^2 = \sum_{j=1}^m a_j^2$ وتعرف عادة $L_2 - norm$ لذلك التعبير لشراح التمهيد التكعيبية

يتحقق عند العقد \hat{f}_λ ، تكون T_j ، $j = 1, 2, \dots, m$ كالتالي :

$$\hat{f}_\lambda = (\mathbf{W}'\mathbf{W} + \lambda G)^{-1} \mathbf{W}' b_r$$

وأن متجه التقديرات عند نقاط زمن التصميم هو

اذ ان :

$$A_\lambda = W(W'W + \lambda G)^{-1} W'$$

وعندما جميع نقاط زمن التصميم تستعمل كعقد فإن التقدير سيصبح كالآتي :

اُذْ اِنْ :

$$W = I_n$$

2.2.2 اختيار معلمة التمهيد

(Smoothing Parameter Selection)

إن أحد أفضل طائق اختيار معلمة التمهيد وأكثرها انتشاراً هو (Generalized Cross – GCV) ، ولا اختيار معلمة التمهيد (λ) للمهد الخطى (t_{λ}) ولأتموج الاملمي في الصيغة **Validation** ، عن طريق تصغير المعيار الآتى : (22)

$$GVC(\lambda) = \frac{m^{-1} \sum_{j=1}^m \left[b_{r,j} - \hat{b}_{r,j} \right]^2}{\left\{ 1 - \frac{t_r(A_\lambda)}{m} \right\}^2} = \frac{m^{-1} SSE_h}{\left(1 - \frac{df}{m} \right)^2} \quad \dots \dots \text{Eq7}$$

. **GCV** أقل تقابل التي **(λ)** التمهيد معلمة اختيار ويتم .



3.2.2 مقترن التقديرات المهددة الحصينة⁽⁷⁾ (Robust Smoothing Estimates)

ان طريقة تقدير الأنماذج الامعملي في (22) وهي طريقة شرائط التمهيد التكعيبية حساسة إتجاه وجود قيم شاذة وجعلها أكثر صرامة بوجود الشوائب يمكن اجراء الاساليب الحصينة في القسم (2.1.2)، إذ إن الاخطاء لأنموذج الامعملي في (22) ستكون كالآتي :

$$e_j = b_r(t_j) - \hat{b}_{r,j} \quad , \quad j = 1, 2, \dots, m$$

لذلك سيكون من الواجب تحصين معيار (λ) GCV كالتالي :

$$GCV_{Rob}(\lambda) = \frac{m^{-1} \sum W_j \left(b_{r,j} - \hat{b}_{r,j} \right)^2}{\left\{ 1 - t_r(A_\lambda) / m \right\}^2} \quad 28)$$

إذ إن :

W_j : هي دالة وزن تحتوي على (m) من العناصر ويمكن حسابها دون الحاجة الى التوسيع لأسلوب M في القسم (2.1.2).

3. المحاكاة (Simulation)

تم تنفيذ تجارب المحاكاة باستخدام ($n=10$) ويمثل عدد القطاعات مع ($m=5, m=10, m=15$) وتمثل القياسات المتكررة لكل قطاع. وبذلك سيكون لدينا ثلاثة جموم للعينات ($nm=50$) و($nm=100$) وأخيراً ($nm=150$) ، ولأنموذج التالي:

$$Y_{i,j} = X_{1,i}(t_j) B_1(t_j) + X_{2,i}(t_j) B_2(t_j) + e_i(t_j) \quad , \quad i=1, 2, \dots, n ; j=1, 2, \dots, m$$

إذ إن :

$\beta_r(t_j)$ ، $r = 1, 2$: هي دوال معاملات مهددة.

المتغيران التوضيحيان (t_j) و ($X_{2,i}(t_j)$) يتبعان التوزيع الطبيعي بمتوسط μ وتباین σ^2 ويتم توليدهما باستعمال طريقة (Box – Muller) وبصورة مستقلة لكل واحداً منها، أما الاخطاء العشوائية فيتم توليدها كالتالي:

1. متوجه الأخطاء ($e_i(t_j)$) يتبع التوزيع الطبيعي بمتوسط (0) وتباین σ^2 يتم عن طريق استعمال

طريقة (Box – Muller) ، وقد تم تناول ثلاثة مستويات للتباین:

* تباین عالي (High Noise)

$$\sigma = \left(\frac{1}{2} \right) * Function Range$$

* تباین متوسط (Medium Noise)

$$\sigma = \left(\frac{1}{4} \right) * Function Range$$

* تباین واطئ (Low Noise)

$$\sigma = \left(\frac{1}{8} \right) * Function Range$$



إذ إن :

 σ : هو الانحراف المعياري للخطأ .

2. أما التوزيع الآخر للخطأ العشوائي $e_i(t_j)$ فهو التوزيع الملوث ويستعمل في حالة تلوث البيانات بقيم شاذة وبنسبة 10% و 20% إذ تم توليد بيانات تتبع توزيع طبيعي بمتوسط (0) وتباين $(= 36\sigma^2)$.

أما دوال المعاملات فهي كالتالي :

$$\beta_1(t) = \sin(4\pi t)$$

$$\beta_2(t) = \cos(0.5\pi t)$$

أما المتغير المعتمد فيتم توليدته مباشرةً من خلال استخدام الأنماذج في دراسة المحاكاة. ولتقييم إداء طرائق التقدير لمرحلة التقدير الخام ومرحلة التقدير الخام الحصينة وكذلك مرحلة التمهيد الحصينة تم استعمال المعايير التالية:

1. متوسط الانحرافات المطلقة للأخطاء

$$MADE = (dm)^{-1} \sum_{j=1}^m \sum_{r=1}^2 \frac{|\beta_r(t_j) - \hat{\beta}_r(t_j)|}{range(\beta_r)}$$

:(Weighted Average Squared Error)

2. متوسط مربعات الخطأ الموزون

$$WASE = (dm)^{-1} \sum_{j=1}^m \sum_{r=1}^2 \frac{\{\beta_r(t_j) - \hat{\beta}_r(t_j)\}^2}{range^2(\beta_r)}$$

إذ إن :

 $\cdot \beta_r(t_j)$ هو المدى إلى دالة $(range(\beta_r))$

وتم تكرار جميع تجارب المحاكاة (Replicates = 200) مرة لكل تجربة وتم وضع جميع النتائج في الجداول من رقم (1) إلى (4).

جدول (1) معايير تقدير Two step لحالة عدم استعمال أسلوب الحصانة في المرحلة الأولى والثانية، ولجميع حجم العينة ولجميع مستويات التباين

ϵ	method	n	m	WASE			MADE		
				$\sigma = \frac{1}{2}$	$\sigma = \frac{1}{4}$	$\sigma = \frac{1}{8}$	$\sigma = \frac{1}{2}$	$\sigma = \frac{1}{4}$	$\sigma = \frac{1}{8}$
10 %	CSS	10	5	214.22	71.06	26.03	13.07	6.52	5.32
		10	10	47.12	11.12	13.52	4.29	2.31	3.63
		10	15	7.19	5.32	5.60	2.37	1.47	1.60
20 %	CSS	10	5	221.43	97.07	138.10	11.24	7.96	9.22
		10	10	10.74	5.67	19.11	2.89	1.92	3.37
		10	15	2.18	1.95	2.03	1.36	1.01	1.14



جدول (2) معايير تقدير Two Step لحالة استعمال أسلوب الحصانة M في المرحلة الأولى فقط ، ولجميع حجوم العينة ولجميع مستويات التباين

ϵ	method	n	m	Rob	WASE			MADE		
					$\sigma = 1/2$	$\sigma = 1/4$	$\sigma = 1/8$	$\sigma = 1/2$	$\sigma = 1/4$	$\sigma = 1/8$
10 %	CSS	10	5	M	163.12	99.05	26.27	11.24	6.68	15.64
		10	10	M	36.27	11.78	12.88	3.87	2.44	3.38
		10	15	M	2.26	7.66	2.57	1.09	2.08	1.23
20 %	CSS	10	5	M	335.80	187.45	213.42	15.59	11.55	11.89
		10	10	M	52.69	19.68	27.06	5.52	3.57	4.08
		10	15	M	2.16	1.90	2.01	1.33	0.99	1.11

4. تحليل النتائج:

لحالة عدم استعمال أسلوب الحصانة في المرحلة الأولى والثانية ، ومن نتائج جدول (1) ، قيم معياري (MADE) سجلاً انخفاضاً ترافق مع زيادة حجم العينة (الزمن) ولكن حالي التلوث، هذا وأن المعياريين سجلاً زيادة عند

مستوى التباين العالي ($1/2$) = σ والواطي ($1/8$) = σ مقارنة مع مستوى التباين المتوسط ($1/4$) = σ على الأغلب .

وللحالة استعمال أسلوب الحصانة M في المرحلة الأولى فقط، ومن نتائج جدول (2) ، قيم معياري (MADE) سجلاً انخفاضاً ترافق مع زيادة حجم العينة (الزمن) ولكن حالي التلوث ، وبمقارنة النتائج مع نتائج جدول (1)

نلاحظ ارتفاع قيم المعياريين عند تلوث 10% ومستوى تباين ($1/4$) = σ وبنسبة أكبر عند تلوث 20% إذ سجلاً ارتفاعاً عند حجم عينة ($m=5, n=10$) و ($m=10, n=10$) ، هذا وأن المعياريين سجلاً زيادةً عند مستوى التباين العالي

($1/2$) = σ والواطي ($1/8$) = σ مقارنة مع مستوى التباين المتوسط ($1/4$) = σ على الأغلب .

وللحالة استعمال أسلوب الحصانة M في المرحلة الثانية فقط، ومن نتائج جدول (3) ، قيم معياري (MADE) سجلاً انخفاضاً ترافق مع زيادة حجم العينة ولكن حالي التلوث، وبمقارنة النتائج مع نتائج جدول (1) نلاحظ

انخفاض قيم المعياريين عند تلوث 10% و 20% ، هذا وأن المعياريين سجلاً زيادةً عند مستوى التباين العالي ($1/2$) = σ

والواطي ($1/8$) = σ مقارنة مع مستوى التباين المتوسط ($1/4$) = σ على الأغلب .

وللحالة استعمال أسلوب الحصانة M في المرحلة الأولى و الثانية معاً ، ومن نتائج جدول (4) ، قيم معياري (MADE) سجلاً انخفاضاً ترافق مع زيادة حجم العينة ولكن حالي التلوث ، وبمقارنة النتائج مع نتائج جدول (1)

نلاحظ انخفاضاً في قيم المعياريين عند استعمال أسلوب M الحصين في المرحلة الأولى والثانية أظهر تقدماً عند تلوث 10% ، وعند 20% سجلاً زيادةً عند حجم عينة ($m=5, n=10$) مقارنة مع الطريقة التقليدية، هذا وأن المعياريين سجلاً

زيادةً عند مستوى التباين العالي ($1/2$) = σ والواطي ($1/8$) = σ مقارنة مع مستوى التباين المتوسط ($1/4$) = σ على الأغلب .

ومن خلال متابعة الجداول من (1) إلى (4) فإن أفضل النتائج لطريقة تقدير CSS عند تلوث 10 هي استعمال أسلوب الحصانة M في المرحلة الثانية فقط وعند تلوث 20 هي استعمال أسلوب الحصانة M في المرحلة الأولى والثانية معاً ، إذ أفرزت أفضل النتائج ولكن المعياريين، وأظهرها المعياريين MADE و WASE تطابق في تفضيل ومقارنة طرائق التقدير.



5. الاستنتاجات والتوصيات :

- 1- استعمال اسلوب M الحصين في المرحلة الاولى فقط أظهر تقدماً عند تلوث 10% على الطريقة التقليدية وافق عند تلوث 20%.
- 2- استعمال اسلوب M الحصين في المرحلة الثانية فقط أظهر تقدماً عند تلوث 10% و 20% على الطريقة التقليدية.
- 3- استعمال اسلوب M الحصين في المرحلة الأولى و الثانية أظهر تقدماً عند تلوث 10% و 20% على الطريقة التقليدية عدا حالة واحدة عند تلوث 20% و حجم عينة ($m=5, n=10$) .
- 4- بصورة مطلقة أفضل النتائج لطريقة تقدير CSS عند تلوث 10% هي استعمال اسلوب الحصانة M في المرحلة الثانية فقط و عند تلوث 20% هي استعمال اسلوب الحصانة M في المرحلة الاولى والثانية معاً .
- 5- استعمال أحد المعيارين MADE و WASE يكون كافي في تفضيل و مقارنة طرائق التقدير.
- 6- وجوب تقصي حجم العينة (عدد نقاط الزمن) مع مستوى التباين (Noise) المسموح به ، لتاثيرهما على قيم المعايير.
- 7- وجوب تحري وأختيار معلمة التمهيد وفق معايير اختيار معلم التمهيد مثل معيار GCV لما لها من أثر على قيم المعايير للمفاضلة بين الطرائق.

المصادر

- 1- Fan, J. and Zhang, J. (2000), "Two – Step Estimation of Functional Linear Models with Applications to Longitudinal Data", Journal of Royal statistical society, vol. 62, no. 2, pp. 303-322.
- 2- Fan, J. and Zhang, W. (2008), "Statistical Methods with Varying Coefficient Models" Statistics and Interface, vol. 1, pp. 179-195.
- 3- Greene, W. H. (2003), "Econometric Analysis" Fifth Edition, Prentice Hall, New Jersey.
- 4-Hart, T.D. (1991), " kernel regression estimation with time series errors ", J. Roy. Stat. Soc. Ser. B, 53, 173-187.
- 5- Hoover, D. R., Rice, J. A., Wu, C. O. and Yang, L. (1998), " Non Parametric Smoothing Estimates of time - Varying Coefficient Models with Longitudinal Data", Biometrika, vol. 85, no. 4, pp. 809-822.
- 6- Huber, P. J. (1981), "Robust Statistics", John Wiley & Sons, New York.
- 7- Kovac, A. (2002), "Robust Nonparametric Regression and Modality", <http://Maths.bris.ac.uk/~Maxak/>
- 8- Senturk, D. and Muller, H. G. (2008), "Generalized varying coefficient models for longitudinal data", Biometrika, vol. 95, Iss.3, pp.653-666.
- 9- Wooldridge, J. M. (2002), "Introductory econometrics: A modern approach", Cambridge, MA: MIT press.
- 10- Wu, C. O., Tian, X. and Yu, J. (2010), "Non parametric estimation for time-varying transformation models with longitudinal data", Journal of non parametric statistics, vol. 22, no. 2, pp. 133-147.
- 11-Zeger, S.L. & Diggle, P.J. (1994), " semiparametric models for longitudinal data with application to CD4 cell numbers in HIV seroconverters" Biometrics, 50, 689-699.



Comparison Robust M Estimate With Cubic Smoothing Splines For Time-Varying Coefficient Model For Balance Longitudinal Data

Abstract

In this research, a comparison has been made between the robust estimators of (M) for the Cubic Smoothing Splines technique, to avoid the problem of abnormality in data or contamination of error, and the traditional estimation method of Cubic Smoothing Splines technique by using two criteria of differentiation which are (MADE, WASE) for different sample sizes and disparity levels to estimate the chronologically different coefficients functions for the balanced longitudinal data which are characterized by observations obtained through (n) from the independent subjects, each one of them is measured repeatedly by group of specific time points (m), since the frequent measurements within the subjects are almost connected and independent among the different subjects.

Keywords/ Time varying coefficient- Two step estimation- Robust M estimation- Cubic Splines smoothing- Balance longitudinal data.