

# مقارنة مقدرات M الحصينة مع مقدرات شرائح التمهيد التكميبيية لأنموذج المعاملات المتغيرة زمنياً للبيانات الطولية المتزنة

الباحث علي سيف الدين عبد الحافظ

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

## المستخلص

في هذا البحث تم مقارنة مقدرات (M) الحصينة لتقنية شرائح التمهيد التكميبيية لتلافي مشكلة الشواذ في البيانات أو تلوث الخطأ مع طريقة التقدير التقليدية لتقنية شرائح التمهيد التكميبيية ، بأستعمال معيارين للمفاضلة بينهما هما (MADE)، (WASE) ولمختلف حجوم العينة ومستويات التباين، وذلك لتقدير دوال المعاملات المتغيرة زمنياً للبيانات الطولية المتزنة، 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 الاغلب ومستقلة بين القطاعات المختلفة.

**المصطلحات الرئيسية للبحث/** المعاملات المتغيرة زمنياً، تقدير المرحلتين، تقديرات M الحصينة، تمهيد الشرائح التكميبيية، البيانات الطولية المتزنة.



مجلة العلوم  
الاقتصادية والإدارية  
المجلد ١٩  
العدد 73  
الصفحات ٣٩٨-٤١٣

بحث مستل من أطروحة دكتوراه

## 1. المقدمة :

في الدراسات الطولية المشاهدات غالباً يتم الحصول عليها من  $n$  من القطاعات المستقلة كل واحداً منها يقاس تكرارياً خلال مجموعة نقاط زمن محددة ، وغالباً 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 للقياسات  $j^{th}$  للقطاع  $i^{th}$  ، وأن  $y_{ij}$  و  $\chi_{ij}$  تمثل مشاهدات القطاع  $i^{th}$  للمتغير المعتمد والمستقل عند الزمن  $t_{ij}$  على التوالي، فإن مجموعة المشاهدات الطولية تعطى كالآتي :-

$$\{ (t_{ij}, y_{ij}, \chi_{ij}) ; i=1,2,\dots,n; j=1,2,\dots,t_i \} \dots\dots(1)$$

حيث  $t_i$  هي عدد القياسات المتكررة للقطاع  $i^{th}$  ، 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 مهتم بنمذجة منحنى المتوسط لـ  $y(t)$  والتأثيرات للمتغيرات المستقلة على  $y(t)$  ، وتطوير التقدير وأجراءات الاستدلال ، وتحت أطار النماذج المعلمية مثل النماذج الخطية وغير الخطية ونماذج ذات التأثيرات المختلطة ، درست نظريات وطرائق التقدير للمعالم والاستدلالات بصورة موسعة، وتحت أطار النماذج اللامعلمية مع نقاط زمن تصميم ثابتة (Hart(1991) إعتد طرائق (kernel) لتقدير التوقع  $E(y(t))$  بدون وجود المتغيرات المستقلة ، ولأخذ المتغيرات المستقلة بالحسبان (Zeger and Diggle (1994) درساً لأنموذج شبة 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 ثابت غير معلوم في  $R^d$  ،  $\mu(t)$  هي دالة ممهدة محددة الى  $(t)$  ،  $\varepsilon(t)$  عملية عشوائية بمتوسط (0) ، وحققوا إجراءات تكرارية حيث أقترحوا إجراء (backfitting) الذي يقدر في البداية  $\mu(t)$  بطريقة (kernel) و ثم تكرار التقديرات لـ  $B$  و  $\mu(t)$  ، وأن الأنموذج في (2) هو أكثر مرونة من النماذج الخطية التقليدية ، ويتطلب لتأثيرات المتغيرات الحالات العملية ، بعبارة أخرى حجوم العينة الفعلي في أغلب الدراسات الطولية يمكن أن لا يكون واقعيّاً للكثير من كامل لأنموذج اللامعلمي العام عندما تكون المتغيرات المستقلة ذات بعد عالي وهو ما يسمى مشكلة البعدية (Curse of Dimensionality) ، لذلك ولأجل تعميم عملي أكثر للأنموذج (2) ، (Hoover et al (1998) أعتدوا أنموذج المعاملات المتغيرة التالي:

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



واقترحوا صنف متعدد الحدود الموضعي (kernel) لتقدير  $B(t)$ ، حيث

$$B(t) = (B_1(t), B_2(t), \dots, B_d(t))'$$

هو متجة للدوال الممهدة خلال  $(t)$ ،  $\varepsilon(t)$  كما معرفة في (2) ولجميع قيم  $(t)$  فإن  $\varepsilon(t)$  و  $X(t)$  مستقلان، وبشكل عام نلاحظ أن الأنموذج (3) هو أنموذج خطي بين  $X(t)$  و  $Y(t)$  عند كل زمن ثابت  $(t)$ ، إجراءات التقدير في (Hoover et al (1998)) طورت الحالة الخاصة لمقدرات (kernel) والتي اعتمدت على عرض حزمة (bandwidth) واحد ومقدرات (spline) والتي اعتمدت على أكثر من معلمة تمهيد، واعتمدوا على خوارزمية (backfitting) لحل مشكلة تعدد الابعاد وخاصة لتقنية الشرائح التمهيديّة، والتي عانت من بعض المشكلات وهي كثافة الحسابات والجهد البرمجي، وكعلاج للمشكلات السابقة (Fan and Zhang (2000)) اقترحوا إجراء المرحلتين (Two Step) كبديل لأجل تقدير  $B(t)$  في (3)، أولاً احتسبوا التقديرات الخام (Raw Estimate) الى دوال المعالم  $B(t)$  عن طريق مطابقة أنموذج خطي قياسي، ثانياً مهدوا المقدرات الخام (Smooth Estimate) للحصول على المقدرات التمهيديّة لدوال المعاملات بواسطة استعمال إحدى تقنيات التمهيد المعروفة، واستعملوا شرائح التمهيد كأحدى تلك التقنيات.

أن تقديرات دوال المعاملات  $B(t)$  بأسلوب المرحلتين يتم الحصول عليه باستعمال طريقة المربعات الصغرى، وكما هو معلوم أن مقدرات المربعات الصغرى تمتلك بعض الخصائص الجيدة، وخاصة عندما الخطأ العشوائي يتبع التوزيع الطبيعي، ولكن المقدرات المعتمدة على المربعات الصغرى حساسة جداً الى الشواذ 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 الأنموذج (3)، لأيجاد تقديرات شرائح التمهيد التكعيبية الحصين لدوال المعاملات بطريقة المرحلتين، وبلا اعتماد على أسلوب M الحصين، ومقارنته مع طرائق التقدير التقليدية عن طريق تجارب محاكاة بنسب تلويث مختلفة وحجوم عينة مختلفة ومستويات تباين مختلفة ولأغراض الملاءمة والتعميم، تم عرض كل الصيغ والمعادلات بدلالة  $d$  من المتغيرات التوضيحية، على الرغم من استعمال دراسة المحاكاة لحالة ثنائي المتغيرات (متغيرين فقط).

## 2. طريقة تقدير المرحلتين (Two-Step Estimation Method) (10)، (2)، (1)

لنفترض  $t_j, j=1, 2, \dots, m$ ، هي نقاط زمن محددة، حيث تم جمع البيانات، أذ ان  $m$  تمثل عدد القياسات المكررة لكل قطاع، لأن هناك عدد من المشاهدات التي جمعت في الزمن  $t_j$ ، فمن الممكن لهذا الثابت  $t_j$  استعمال البيانات المجمعة هناك لمطابقة أنموذج (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 الى تمهيدها للحصول على المقدرات الممهدة الى دوال المعاملات لذلك، في المرحلة الثانية لكل مركبة معطاة  $r=1, 2, \dots, d$  نطبق تقنيّة تمهيد الى البيانات  $\{(b_r(t_j), t_j), j=1, 2, \dots, m\}$ ، وإن مرحلة تقديرات التمهيد (Smoothing estimates) هذه حاسمة لأنها تعطي مقدرات تمهيدية لدوال معاملات التمهيد الاساسية، وإضافة الى ذلك فإن مرحلة التمهيد ذو بعد واحد (one-dimensional).

1.2 مرحلة التقديرات الخام (Raw Estimates Step) <sup>(8)</sup>

ولتوضيح هذه المرحلة نفترض  $t_j, j=1, 2, \dots, m$  هي نقاط زمن محددة لمجموعة البيانات الطولية لكل نقطة زمن  $t_j$ ، لنفترض  $N_j$  هو مجموعة الفهارس للقطاع الى جميع مشاهدات  $y_{ij}$  عند  $t_j$ ، نجمع كل  $y_{ij}$  و  $X_{ij}$  التي تقابل الفهرس للقطاعات في  $N_j$  ونشكل مصفوفة التصميم  $\tilde{X}_j$  و متجه 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  $t_j$  ستكون متساوية ولجميع القطاعات. عندئذ فان صيغة أنموذج (3) عندما البيانات تجمع عند الزمن يتبع الأنموذج الخطي الاتي :

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

 1.1.2 مقدرات المربعات الصغرى العامة المقبولة مع اخطاء AR(1) <sup>(9)</sup>

(Feasible GLS Estimation With AR(1) Errors)

لتقدير معالم الأنموذج (4) وتحت افتراض هيكل الارتباطات للأخطاء يتبع AR(1) كالاتي :

$$\tilde{e}_i(t_j) = \rho \tilde{e}_i(t_{j-1}) + u_{ij} \quad \dots\dots\dots(5)$$

يمكننا تطبيق المربعات الصغرى العامة حيث مقدراتها ستكون كالاتي :

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

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

ولتجنب افتراض GLS علينا ايجاد تقدير متسق لـ  $\hat{\Omega}$  أي  $\hat{\rho}$  واستعمالة لاجاد تقدير لمعالم الأنموذج (4)، إن هذا المقدر يدعى مقدرات المربعات الصغرى العامة المقبولة (FGLS) والذي يمكن ايجاده بحسب الخطوات الاتية :

a. نجد اولاً مقدرات المربعات الصغرى الاعتيادية الى أنموذج (4) والذي سيكون كالاتي :

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

ثم نجد الاخطاء باستعمال مقدرات OLS .

إذ إن :

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

 b. وتحت افتراض الاخطاء هي عمليات عشوائية مشتركة فان مقدر  $\rho$  المشترك يمكن تقديره كالآتي :

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

 c. إجراء تحويل للبيانات باستعمال (Prais - Winsten) transformation <sup>(3)</sup>.

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

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

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

$$b_{FGLS}(t_j) = \left( X_i^*(t_j) X_i^*(t_j) \right)^{-1} X_i^*(t_j) Y_i^*(t_j) \quad \dots\dots\dots(12)$$

وان مقدرات FGLS هي تقاربياً اكثر كفاءة من مقدرات OLS عندما يتبع هيكل ارتباط AR(1).

### 2.1.2 مقترح التقديرات الخام الحصينة Robust Raw Estimate

اقترح أسلوب M الحصين اولاً من قبل (Huber) <sup>(6)</sup>، وتستند الفكرة ببساطة الى تقليل بعض الدوال للأخطاء بدلاً عن مجموع المربعات لها، والمقدر الحصين يحدد عن طريق الاختيار لدالة وزن، وان أسلوب 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 لاسلوب مقدرات M وهو أن تكون الأخطاء غير مرتبطة، ولتلافي هذه المشكلة سيتم الاعتماد على البيانات المحولة في طريقة (Feasible GLS).

لنفرض أن  $Y_{ij} = y_i(t_j)$  و  $X_{ij} = \chi_i(t_j)$  ومتجه المعالم  $\beta = \beta(t_j)$ ، ولتوضيح هذا الاسلوب، أن المربعات الصغرى الاعتيادية للبيانات المحولة تخفض مجموع مربعات الخطأ إلى أقل ما يمكن كالآتي:

$$\min \sum_{i=1}^n \sum_{j=1}^m e_{ij}^2 = \min \sum_{i=1}^n \sum_{j=1}^m \left( y_{ij} - \chi_{ij} \beta \right)^2 \quad \dots\dots\dots(3)$$

إذ إن:

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

$\chi_{ij}^*$ : الصف  $jz$  لمصفوفة التصميم  $\chi$ .

وإن:

$$X_{ijk}^* = \left[ \chi_{ij1}^*, \dots, \chi_{ijd}^* \right], \quad k = 1, 2, \dots, d$$

$\beta$ : متجه معالم ذو بعد  $1 \times dm$ .

أن تقديرات M المطورة من قبل (Huber) والتي لها خاصية (Scale Invariant) تعتمد على فكرة ابدال مجموع مربعات الأخطاء  $e_{ij}^2$  بدالة أخرى 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 متماثلة (symmetric convex function) ولتقليل المقدر اعلاه تؤخذ المشتقة بالنسبة إلى متجه المعالم وجعلها مساوية إلى الصفر كالآتي:

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



إذ إن :

$\Psi$  : المشتقة الجزئية الى متجه المعالم  $\beta$  للدالة  $P$  و  $\Psi = P'$  وبهذا يكون هناك  $(dm)$  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 إعادة كتابة الصيغة (١٥) كما يلي :

$$\sum_{i=1}^n \sum_{j=1}^m W_{ij} X_{ij}^* \left( \frac{(Y_{ij}^* - X_{ij}^* \beta)}{\sigma_e^2} \right) = 0 \quad \dots\dots\dots(16)$$

وبحل المعادلة اعلاه نحصل على تقديرات أسلوب M الحصين  $b_{\mu}$  باستعمال IWLS .  
إذ إن :

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

إذ إن :

$b_M$  : متجه ذو بعد  $(dm*1)$

$W$  : مصفوفة اوزان قطرية ببعد  $nm*nm$  () تحسب عناصرها كما يأتي :

$$W_{ij} = \frac{\Psi \left( \frac{(y_{ij}^* - X_{ij}^* \beta)}{\hat{\sigma}_e} \right)}{\left( \frac{(y_{ij}^* - X_{ij}^* \beta)}{\hat{\sigma}_e} \right)} \quad \dots\dots\dots(18)$$

إذ إن :

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

وتم استعمال دالة P ومشتقتها  $\Psi$  التالية :

دالة (Andrews)

$$\Psi(e_{ij}^*) = \begin{cases} A \sin(e_{ij}^*/A) & |e_{ij}^*| \leq A\Pi \\ 0 & |e_{ij}^*| > A\Pi \end{cases}$$

$$A = 1.339$$



## 2.2 مرحلة تحسين أو تمهيد المقدرات الخام<sup>(1)</sup>

### (Refining Or Smoothing The Raw Estimates)

وقبل البدء بهذه المرحلة فإن مركبات دوال المعاملات المقدره في المرحلة الاولى نحصل عليها كالآتي :  
ولـ  $r=1,2, \dots, d$  نأفرض أن :

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

$$b_r(t) = L'_r \left( X' X \right)^{-1} X' Y \quad \dots\dots\dots(9)$$

إذ إن :

$L_r$  : يعرف كمتجه وحدة ذو بعد  $(dm*1)$  فيه  $r^{th}$  من المدخلات (1) والباقي (0).

ومن المعلوم إن

$$E[b_r(t)] = \beta_r(t) \quad \dots\dots\dots(10)$$

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

$$(b_r(t_j), t), j = 1, 2, \dots, m \quad \dots\dots\dots(11)$$

إذ إن :

$t_j$  : هي نقاط زمن التصميم.

$b_r(t_j)$  : هي الاستجابة عند نقاط زمن التصميم، مع التأكيد بانها اي تقدير مستخرج من المرحلة الاولى.

ان أنموذج الانحدار اللامعلمي البسيط للبيانات السابقة سيكون كالآتي :

$$b_r(t_j) = f(t_j) + \varepsilon_j, j = 1, 2, \dots, m \quad \dots\dots\dots(12)$$

ونحن نريد تقدير الدالة الممهدة  $f(t_j)$

إذ إن :

$\varepsilon_j$  : هي اخطاء القياسات التي لايمكن شرحها بواسطة دالة الانحدار  $f(t_j)$ .

رياضياً  $f(t)$  هي التوقع الشرطي لـ  $b_r(t_j) \setminus t_j = t$  أي

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





## 1.2.2 شرائح التمهيد التكعيبية

(Cubic Smoothing Splines) :

لايجاد مطابقة شرائح التمهيد لأنموذج (22) وبدون فقدان التعميم لو افترض مدى الفترة لـ  $f$  في الأنموذج (22) هي فترة منتهية  $[a, b]$  ولبعض الاعداد المنتهية  $b, a$ . فان الجزاء غير الممهد (Roughness Penalty) لـ  $f$  هو عادة يعرف كتكامل لمربع مشتقة  $(V)$  من المرات كالاتي :

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

ولبعض  $V \geq 1$  ، عندما ممهد شرائح التمهيد لـ  $f$  في أنموذج (22) يعرف كتقليل  $\hat{f}_\lambda(t)$  لصيغة المربعات الصغرى الجزائية PLS التالية :

$$\sum_{j=1}^m [b_r(t_j) - f(t_j)]^2 + \lambda \int_a^b \{f_h^{(V)}\}^2 dt \quad \dots\dots\dots 23)$$

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

إذ إن :

$\lambda > 0$  هي معلمة التمهيد .

وان اختيارنا لشرائح التمهيد التكعيبية لتقليل الصيغة (23) هو صعوبة الحصول على التكامل الذي يعرف الجزاء غير الممهد فضلاً على وجود طريقة لحسابه في شرائح التمهيد التكعيبية.

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

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

فان جميع العقد لشرائح التمهيد التكعيبية التي تقلل (23) عندما  $V = 2$  تكون كالاتي :

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

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

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

$$a_{L+1,L} = -(h_L^{-1} + h_{L+1}^{-1})$$

$$a_{L+2,L} = -h_{L+1}^{-1}$$

مقارنة مقدرات M الحصينة مع مقدرات شرائح التمهيد التكميلية  
 لأنموذج المعاملات المتغيرة زمنياً للبيانات الطولية المتزنة

ونعرف  $C = (c_{LS})$  كمصفوفة ذات بعد  $((m-2)*(m-2))$  مع جميع مدخلاتها (0) ماعدا  
 $c_{11} = (h_1 + h_2)/3$   
 $c_{21} = h_2/6$

ولـ  $L=1,2,\dots,m-4$  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  $(m*m)$  غير الممهدة التكميلية كالآتي :

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

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

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

وبذلك الجزاء غير الممهدة يمكن ان نعبر عنه كالآتي :

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

لاجل ذلك : نحن فقط نشير الى  $G$  كمصفوفة غير الممهدة (Roughness matrix)، وهذا يعني أن صيغة PLS في (23) يمكن كتابتها كالآتي :

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



مقارنة مقدرات M الحسنة مع مقدرات شرائح التمهيد التكميلية  
لأنموذج المعاملات المتغيرة زمنياً للبيانات الطولية المتزنة

إذ إن :  $W = (W_{jj})$  هي مصفوفة حدث ذات بعد  $(m * m)$  مع  $W_{jj} = 1$  إذا  $t_j = T_j$  و  $0$  في الحالات الأخرى.

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

$\hat{f}_\lambda$  ، يتحقق عند العقد  $T_j$  ،  $j = 1, 2, \dots, m$  وتكون كالآتي :

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

وأن متجه التقديرات عند نقاط زمن التصميم هو

$$b_{r,\lambda}^* = A_\lambda b_r \quad \dots\dots\dots (25)$$

إذ إن :

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

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

$$b_{r,\lambda}^* = \hat{f}_\lambda = (I_n + \lambda G)^{-1} b_r \quad \dots\dots\dots (26)$$

إذ إن :

$$W = I_n$$

## 2.2.2 اختيار معلمة التمهيد

### (Smoothing Parameter Selection)

إن احد أفضل طرائق اختيار معلمة التمهيد وأكثرها انتشاراً هو (GCV - Generalized Cross

Validation) ، ولاختيار معلمة التمهيد  $(\lambda)$  للممهد الخطي  $\hat{f}_\lambda(t)$  ولأنموذج اللامعلمي في الصيغة (٢٢)

، عن طريق تصغير المعيار الآتي :

$$GVC(\lambda) = \frac{m^{-1} \sum_{j=1}^m [b_{r,j} - b_{r,j}^*]^2}{\left\{1 - \frac{t_r(A_\lambda)}{m}\right\}^2} = \frac{m^{-1} SSE_h}{\left(1 - \frac{df}{m}\right)^2} \quad \dots\dots\dots (27)$$

ويتم اختيار معلمة التمهيد  $(\lambda)$  التي تقابل أقل GCV .



### 3.2.2 مقترح التقديرات الممهدة الحصينة <sup>(7)</sup> (Robust Smoothing Estimates)

ان طريقة تقدير الأنموذج اللامعلمي في (22) وهي طريقة شرائح التمهيد التكميلية حساسة إتجاه وجود قيم شاذة وجعلها أكثر صرامة بوجود الشواذ يمكن اجراء الاساليب الحصينة في القسم (2.1.2)، إذ إن الاخطاء للأنموذج اللامعلمي في (22) ستكون كالآتي :

$$e_j = b_r(t_j) - b_{r,j}^* \quad , \quad j = 1, 2, \dots, m$$

لذلك سيكون من الواجب تحصيل معيار  $(\lambda)$  GCV كالآتي :

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

إذ إن :

$W_j$  : هي دالة وزن تحتوي على  $(m)$  من العناصر ويمكن حسابها دون الحاجة الى التوسيع لاسلوب 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 ثلاث حجوم للعينات  $(nm=50)$  و  $(nm=100)$  وأخيراً  $(nm=150)$  ، وللأنموذج 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 ممهدة.

المتغيران التوضيحيان  $X_{1,i}(t_j)$  و  $X_{2,i}(t_j)$  يتبعان التوزيع الطبيعي بمتوسط  $\mu$  وتباين  $\sigma^2$  ويتم توليدهما باستعمال طريقة (Box – Muller) وبصورة مستقلة لكل واحداً منهما، أما الاخطاء العشوائية فيتم توليدها كالآتي:

١. متجه الأخطاء  $e_i(t_j)$  يتبع التوزيع الطبيعي بمتوسط (٠) وتباين  $\sigma^2$  يتم عن طريق استعمال

طريقة (Box – Muller) ، وقد تم تناول ثلاث مستويات للتباين:

\* تباين عالي (High Noise)

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

\* تباين متوسط (Medium Noise)

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

\* تباين واطى (Low Noise)

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



مقارنة مقدرات M الحصينة مع مقدرات شرائح التمهيد التكميلية  
لأنموذج المعاملات المتغيرة زمنياً للبيانات الطولية المتزنة

إذ إن :

$\sigma$  : هو الانحراف المعياري للخطأ  $e$ .

٢. أما التوزيع الآخر للخطأ العشوائي  $e_i(t_j)$  فهو التوزيع الملوث ويستعمل في حالة تلوث البيانات بقيم شاذة وبنسب 10% و 20% إذ تم توليد بيانات تتبع توزيع طبيعي بمتوسط (٠) 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 توليده مباشرة من خلال استخدام الأنموذج في دراسة المحاكاة. ولتقييم أداء طرائق التقدير لمرحلة التقدير الخام ومرحلة التقدير الخام الحصينة وكذلك مرحلة التمهيد ومرحلة التمهيد الحصينة تم استعمال المعايير التالية:

١. متوسط الانحرافات المطلقة للأخطاء (Mean Absolute Deviation Error):

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

٢. متوسط مربعات الخطأ الموزون (Weighted Average Squared Error):

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

إذ إن :

$(range(\beta_r))$  هو المدى الى دالة  $\beta_r(t_j)$ .

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

جدول (١) معايير تقدير Two step لحالة عدم استعمال أسلوب الحصانة في المرحلة الاولى والثانية، ولجميع

حجوم العينة ولجميع مستويات التباين

€	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



## مقارنة مقدرات M الحصينة مع مقدرات شرائم التمهيد التكميلية

## لأنموذج المعاملات المتغيرة زمنياً للبيانات الطولية المتزنة

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

€	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) ، (WASE) سجلاً انخفاضاً ترافق مع زيادة حجم العينة (الزمن) ولكلا حالتي التلوث، هذا وأن المعيارين سجلاً زيادة عند

مستوى التباين العالي  $\sigma=(1/2)$  والواطي  $\sigma=(1/8)$  مقارنة مع مستوى التباين المتوسط  $\sigma=(1/4)$  على الاغلب .

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

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

$\sigma=(1/2)$  والواطي  $\sigma=(1/8)$  مقارنة مع مستوى التباين المتوسط  $\sigma=(1/4)$  على الاغلب.

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

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

والواطي  $\sigma=(1/8)$  مقارنة مع مستوى التباين المتوسط  $\sigma=(1/4)$  على الاغلب .

ولحالة استعمال أسلوب الحصانة M في المرحلة الأولى والثانية معاً ، ومن نتائج جدول (4) ، قيم معياري (MADE) ، (WASE) سجلاً انخفاضاً ترافق مع زيادة حجم العينة ولكلا حالتي التلوث ، وبمقارنة النتائج مع نتائج جدول (1) نلاحظ انخفاض قيم المعيارين 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 المعيارين سجلاً زيادةً عند

مستوى التباين العالي  $\sigma=(1/2)$  والواطي  $\sigma=(1/8)$  مقارنة مع مستوى التباين المتوسط  $\sigma=(1/4)$  على الاغلب .

ومن خلال متابعة الجداول من (1) الى (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 ولكلا المعيارين، وأظهرا المعيارين MADE و WASE تطابق في تفضيل ومقارنة طرائق التقدير.

## 5. الاستنتاجات والتوصيات :



### لأنموذج المعاملات المتغيرة زمنياً للبيانات الطولية المتزنة

- ١- أستعمال أسلوب M الحصين في المرحلة الأولى فقط أظهر تقدماً عند تلوث 10% على الطريقة التقليدية وأخفق عند تلوث 20%.
- ٢- أستعمال أسلوب M الحصين في المرحلة الثانية فقط أظهر تقدماً عند تلوث 10% و 20% على الطريقة التقليدية.
- ٣- أستعمال أسلوب M الحصين في المرحلة الأولى و الثانية أظهر تقدماً عند تلوث 10% و 20% على الطريقة التقليدية عدا حالة واحدة عند تلوث 20% وحجم عينة (m=5, n=10).
- ٤- بصورة مطلقة أفضل النتائج لطريقة تقدير CSS عند تلوث 10% هي أستعمال أسلوب الحصانة M في المرحلة الثانية فقط وعند تلوث 20% هي استعمال أسلوب الحصانة M في المرحلة الأولى والثانية معاً.
- ٥- أستعمال أحد المعيارين MADE و WASE يكون كافي في تفضيل ومقارنة طرائق التقدير.
- ٦- وجوب تقصي حجم العينة (عدد نقاط الزمن) ، مع مستوى التباين (Noise) المسموح به ، لتأثيرهما على قيم المعايير.
- ٧- وجوب تحري وأختيار معلمة التمهيد وفق معايير أختيار معالم 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.