

اثر المشاهدات الشاردة وذات القوة الرافعة في بناء فترات الثقة البيزية و بوتستراب

د. مزاحم محمد يحيى الهاشمي

muzahim63@gmail.com

المستخلص:

يهدف البحث إلى مقارنة كل من حدود الثقة لبوتستراب (bootstrap confidence intervals) مع حدود الثقة البيزية (Bayesian confidence intervals) لشرائح ممهدة (smoothing splines)، فضلاً عن حدود الثقة التقليدية، وذلك لتحديد أي من هذه الحدود هي الأفضل في ظل وجود نقاط شاردة وذات القوة الرافعة في البيانات. تم إجراء تجارب المحاكاة على أنموذجين. الأول: خطي في ظل وجود بيانات ملوثة بمشاهدات شاردة، وأخرى بمشاهدات ذات القوة الرافعة، و الانموذج الثاني اللاخطي في ظل وجود بيانات ملوثة بمشاهدات شاردة. وتم تنفيذ تجارب المحاكاة على حجوم لعينات مختلفة. واستخدام طريقة المربعات الصغرى الجزائية (Penalized Least Squares method) لتوفيق الانحدار اللامعلمي. وتم استخدام دالة الفاعلية المتقاطعة المعممة (Generalized Cross Validation function) (GCV) لاختيار قيمة التمهيد (The amount of smoothing).

الكلمات المفتاحية: المشاهدات الشاردة، المشاهدات ذات القوة الرافعة، طريقة المربعات الصغرى الجزائية، فترات الثقة البيزية، فترات الثقة لبوتستراب، فترات الثقة التقليدية.

This is an open access article under the CC BY 4.0 license
<http://creativecommons.org/licenses/by/4.0/>

The Effect of the Outliers and Leverage Points in the Construction of the Bayesian and Bootstrap Confidence Intervals

Abstract

The aim of this research is to compare the bootstrap confidence intervals with the Bayesian confidence intervals for smoothing splines as well as the traditional confidence intervals to determine which of these limits are best in the presence of Outliers and Leverage points in data. The simulation experiments were conducted on two models: the first was linear in the presence of data that was contaminated with outliers and the other with the Leverage points: The second model was nonlinear in the presence of data contaminated with outlying observations.

* مدرس / قسم الاحصاء والمعلوماتية / كلية علوم الحاسوب والرياضيات / جامعة الموصل .

Simulation experiments were carried out on different samples. The Penalized Least Squares method was used to fit the Nonparametric regression. The Generalized Cross Validation function (GCV) was used to select the amount of smoothing.

Key words: Outliers, Leverage points, Penalized Least Squares, Bayesian confidence intervals, Bootstrap confidence intervals, traditional confidence intervals.

المقدمة:

لغرض تقدير معالم المجتمع الإحصائي التي غالباً ما تكون مجهولة، يُلجأ إلى العينات للحصول على تقديرات لمعاملات المجتمع. وللتقدير أسلوبان، الأول: يسمى تقدير المعلمة بنقطة، والثاني يسمى تقدير المعلمة بفترة. ففي حالة تقدير المعلمة بنقطة نحصل على قيمة واحدة من العينة، وتستخدم هذه القيمة الواحدة كتقدير لمعلمة المجتمع المجهولة. أما في حالة تقدير المعلمة بفترة، فنحصل على فترة من العينة تتحدد بحدين هما: الحد الأدنى، والحد الأعلى، وإن هذه الفترة تحتوي على عدد غير محدد من القيم.

يتميز تقدير المعلمة بفترة بأنه يمكن معرفة مدى دقة التقدير؛ وذلك من خلال حساب احتمال أن يكون التقدير صحيحاً، لذا فإن فترات التقدير تسمى أيضاً " فترات الثقة " Confidence intervals؛ لأن هذه الفترات تعتمد إحصائياً على درجات أو مستويات ثقة معينة Confidence Levels. ويتم تحديد المستوى المطلوب من الثقة من الباحث مثل 95% أو 99 % وغيرها، بمعنى أن احتمال أن تكون فترة التقدير صحيحة هو 0.95 أو 0.99.

ومن العوامل التي تؤثر في تباعد حدود الثقة هو حجم العينة، ومستوى الثقة. وعادة ما يكون لحجم العينة الأثر الكبير في الوصول إلى تقدير أفضل لمعلمة المجتمع، إذ كلما كبر حجم العينة فإن ذلك سيقارب من حدود الثقة؛ لأنه يقلل من الانحراف المعياري، وهو بذلك دالة على كفاءة المقدر.

تعد دراسة القيم الشاردة (Outlying values)، فضلاً عن القيم ذات القوة الرافعة من المسائل المهمة في الإحصاء النظري والتطبيقي، لما لها من أثر في الوصول إلى دقة في النتائج. لقد أثير الكثير من النقاش حول قبول أو استبعاد هذه المشاهدات، حتى استقر الحال بقبولها في التحليل، خاصة إذا كانت هذه المشاهدات هي من صميم البيانات التي تحت الدراسة. وقد ظهرت الكثير من طرائق التحليل الحصينة، التي تختلف بشكل كبير عن الطرائق التقليدية، إذ أخذت بنظر الاعتبار وجود هذه المشاهدات في التحليل.

إنَّ البحث عن معلمة المجتمع غير المعلومة في ظل وجود المشاهدات الشاردة وذات القوة الرافعة تعد من المسائل المهمة، لما لهذه المشاهدات من أثر في تضليل النتائج، وخصوصاً حينما يتم التعامل مع حدود الثقة (الذي يفترض أن تشمل القيمة الحقيقية)، إذ تؤثر هذه المشاهدات في سعة فترات الثقة (كلما تقلصت فترة الثقة زادت دقة التقديرات)؛ لتغطي مشاهدات أكبر نتيجة لتطرف هذه المشاهدات عن تجمع بقية المشاهدات (Bulk of data)، مما يؤدي الى تباعد حدود الثقة الأمر الذي يؤدي إلى تشتت أكبر في النتائج.

دراسات سابقة:

استنتج كل من (Wang & Wahba, 1995) (في حالة عدم وجود المشاهدات الشاردة وذات القوة الرافعة). أنَّ فترات الثقة البيزية مشابهة لفترات الثقة لبوتستراب فيما يتعلق بمتوسط احتمال نطاق التغطية. كما توصل الباحثان إلى أنَّ فترات الثقة للبوتستراب الأفضل عندما تكون حجوم البيانات الصغيرة. فضلاً عن ذلك فإنَّ فترات الثقة لبوتستراب المئني (Bootstrap percentile-interval) أفضل من الأنواع الأخرى لحدود الثقة لبوتستراب.

هدف البحث:

يهدف البحث إلى مقارنة كل من حدود الثقة لبوتستراب (bootstrap confidence intervals) مع حدود الثقة البيزية (Bayesian confidence intervals) للشرائح الممهدة (smoothing splines)؛ فضلاً عن حدود الثقة التقليدية، لتحديد أيٍّ من هذه الحدود حصينة تجاه القيم الشاردة وذات القوة الرافعة في البيانات.

أهمية البحث:

إنَّ الوصول إلى الدقة في التقديرات، فضلاً عن البحث عن معلمة المجتمع غير المعلومة في ظل وجود المشاهدات الشاردة وذات القوة الرافعة تعد من المسائل المهمة في الإحصاء، لما لهذه المشاهدات من أثر في تضليل النتائج، وخصوصاً حينما يتم التعامل مع حدود الثقة (الذي يفترض أن تشمل القيمة الحقيقية، فضلاً عن انها دالة على كفاءة المقدر)، إذ من المتوقع أن تؤثر هذه المشاهدات في عرض فترات الثقة (كلما تقلصت فترة الثقة زادت دقة التقديرات)؛ لكونها تزيد من حالة التشتت، ثم في الوصول الى دقة في النتائج.

1. الجانب النظري

1.1: تقدير المربعات الصغرى الجزائية (Penalized Least Squares Estimation)

يعد الانحدار اللامعلمي (Nonparametric Regression) شكلاً من أشكال تحليل الانحدار، إذ لا توجد أية صيغ مسبقة للمتغيرات التوضيحية مع متغير الاستجابة، كما لا يفترض الانحدار اللامعلمي شكل العلاقة بين المتغيرات، ولكنه قادر على استنتاج العلاقات المعقدة من البيانات (Simonoff, 2012). وبذلك يختلف الانحدار اللامعلمي عن الانحدار الخطي والانحدار اللاخطي، فضلاً عن النماذج الخطية المعممة (Generalized Linear Model) بأن النماذج الثلاثة الأخيرة تتطلب توصيف أنموذج الانحدار، فضلاً عن أن أنموذج الانحدار هو ذو عدد محدد من المعلمات، في حين يستخدم الانحدار اللامعلمي حينما لا تتوفر معلومات مسبقة عن أنموذج الانحدار، أو عندما لا يمكن تمثيل البيانات بأنموذج ذي عدد محدد من المعلمات.

توفر طرائق التمهيد (Smoothing Methods) جسراً بين عدم وجود افتراضات لنماذج الانحدار اللامعلمية والافتراضات القوية لنماذج الانحدار المعلمية، وذلك من خلال إيجاد افتراضات ضعيفة نسبياً (Simonoff, 2012).

لقد أصبحت إجراءات الانحدار الجزئية نهجاً شائعاً جداً لتقدير الدوال المعقدة. فعلى سبيل المثال الصفائح الممهدة هي الحل لمسألة التقليل (minimization)، أو الوصول إلى الحد الأدنى في الفضاء الدالي. ففي أية مسألة تقليل، هناك العديد من الأسئلة التي يتم طرحها، مثل: هل الحل موجود؟ فإذا كانت الإجابة نعم، هل الحل فريد (unique)؟ وكيف يمكن الوصول إليه؟. فإذا كانت المشكلة تطرح في إطار إعادة إنتاج فضاء هيلبرت (الذي سيتم تناوله لاحقاً)، فذلك يعني ضمان وجود الحل، وأنه فريد، فضلاً عن أنه يأخذ صيغة بسيطة للغاية (Nosedal, et al., 2012).

تستخدم طريقة المربعات الصغرى الجزئية لتوفيق أنموذج الانحدار اللامعلمي؛ وذلك بتقليل مجموع مربعات الباقي، فضلاً عن حد الجزاء. إذ يتم حساب الشرائح الممهدة للصفائح الرقيقة (thin-plate smoothing spline) (تشير الصفائح الرقيقة إلى مقارنة للخصائص الفيزيائية للرقائق المعدنية التي تتضمن قدرة هذه الرقائق على التقوس والانحناء (Ruan, Zhixing, et al., 2013)) وذلك لتقريب الدوال المتعددة الممهدة للبيانات المشاهدة. وبهذه الطريقة يتم توفيق البيانات وفق أنموذج مرن بحيث إن عدد المعلمات المعنوية يمكن أن يكون أكبر من نقاط التصميم الوحيدة (unique design points). لذا فإنه كلما ازداد حجم العينة فإن فضاء الأنموذج يزداد أيضاً مما يؤدي إلى اختيار أنموذج من مجموعة كبيرة من النماذج المحتملة بحيث يحقق الاتزان بين تعقيد الأنموذج وأفضل توفيق (Wood, 2003).

يمكن بوساطة المربعات الصغرى الجزائية توفيق البيانات بعناية وبصورة متزنة بعيداً عن التوفيق غير الممهد (roughness of the fit) الناتجة عن وجود الكثير من العقد (Knots) (وهي النقاطات (junctions) متعددة الحدود من الدرجة n ، التي تؤدي إلى ظهور منحني ممهد خلال مجموعة من النقاط) في الأنموذج (رشيد و خمو، 2005) (Ruppert, et al., 2003). يعرف تقدير المربعات الصغرى الجزائية بأنه حل فريد في إيجاد السطح الذي يقلل (minimizes) المربعات الصغرى الجزائية لفئة كل السطوح التي تحقق الشروط المنتظمة الكافية (sufficient regularity conditions).

ليكن (\quad) متجه المتغيرات التوضيحية ذا البعد (d) من المصفوفة (\mathbf{X}) ذات البعد $(n \times d)$ ، وليكن i متجه المتغيرات التوضيحية ذا البعد p ، وان y_i هي القيم المشاهدة لمتغير الاستجابة المقترنة بكل من المتجهين $(\mathbf{x}_i, \mathbf{z}_i)$. وعلى افتراض أن العلاقة بين \mathbf{z}_i و y_i هي علاقة خطية، وأن العلاقة بين \mathbf{x}_i و y_i هي علاقة غير معلومة، فإنه يمكن توفيق البيانات باستخدام الأنموذج شبه المعلمي (Semiparametric Model) (هي نماذج إحصائية، إذ إن بعض فضاء المعلمة (parameter space) لها مكون واحد أو أكثر من الأبعاد غير النهائية (infinite-dimensional component)، وكالاتي (Nosedal, et al., 2012):

$$\varepsilon_i = y_i - f(\mathbf{x}_i) - \mathbf{z}_i \boldsymbol{\beta} \quad i = 1, \dots, n \quad (1)$$

إذ إن y_i متغير الاستجابة، وإن \mathbf{x}_i هي متغيرات التمهيد، و f هي الدالة الممهدة (smooth function)، وإن \mathbf{z}_i هي متغيرات الإنحدار، و $\boldsymbol{\beta}$ هو المتجه المعلمي (parametric vector) غير المعلوم ذو البعد p . وان $\varepsilon_i \sim N(0, \sigma^2)$ بشرط أنها مستقلة.

تستخدم طريقة المربعات الصغرى الاعتيادية لتقدير $f(\mathbf{x}_i)$ و $\boldsymbol{\beta}$ بوساطة تقليل مجموع مربعات البواقي الجزائية الآتي:

$$\frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i) - \mathbf{z}_i \boldsymbol{\beta})^2 \quad \dots (2)$$

لكون أن الفضاء الدالي (functional space) $f(\mathbf{x}_i)$ كبير جداً، فإنه يمكن استخدام طريقة المربعات الصغرى الجزائية التي يمكن من خلالها الحصول على التقدير الذي يوفق البيانات بشكل جيد، فضلاً عن أن له درجة معينة من التمهيد، وتعرف بالآتي:

$$S_\lambda(f) = \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i) - \mathbf{z}_i \boldsymbol{\beta})^2 + \lambda J_2(f) \quad \dots (3)$$

إنَّ الحدَّ الأول من الصيغة في أعلاه يمثل دقة التوفيق، وإنَّ الحدَّ الثاني المقترن بالوزن λ يمثل الجزء غير الممهّد لـ f ، الذي يعرف في معظم الحالات بأنه التكامل لمربع المشتقة الثانية لـ f ($\int \{f''(x)\}^2 dx$) (Nosedal, Storlie, & Christensen 2011).

تمثل المربعات الصغرى الجزائية حلاً وسطاً بين حسن المطابقة (goodness-of-fit) والتمهيد، ويتم التحكم في التوازن بين الاثنين بواسطة معلمة التمهيد λ ، التي تأخذ القيم من 0 إلى ∞ (Wang, 2011). والتي تتحكم بالمفاضلة بين التمهيد وحسن المطابقة. عندما تكون λ كبيرة فإنه من الصعوبة بمكان إجراء الجزاءات للتوفيق غير الممهدة. وعلى العكس من ذلك حينما تكون λ صغيرة، فإنَّ التركيز سيكون بشكل أكبر على حسن المطابقة.

لذا، فإنه يمكن تمثيل الصيغة (2) كتركيبية خطية (linear combination) من سلسلة من الدوال الأساسية. وبالتالي فإنَّ التقديرات النهائية للصفات الممهدة f تكون بالصيغة الآتية:

$$\hat{f}_\lambda(\mathbf{x}_i) = \theta_0 + \sum_{j=1}^d \theta_j \mathbf{x}_j + \sum_{j=1}^p \delta_j \beta_j(\mathbf{x}_j) \quad \dots (4)$$

إذ إنَّ j هي الدالة الأساسية، التي تعتمد على مكان وجود البيانات \mathbf{x}_i ، وإنَّ $\theta = \{\theta_0, \dots, \theta_d\}$ و $\delta = \{\delta_0, \dots, \delta_d\}$ هي المعلمات المطلوب تقديرها.

إذا كانت λ ثابتاً (λ) تمثل معلمة انحدار الحرف في انحدار الحرف (Bates, et al., 1987)، فإنه يمكن تقدير المعلمات (θ, δ, β) بواسطة إيجاد حل للنظام $n \times n$. كما يمكن اختيار معلمة التمهيد بواسطة تقليل دالة الفاعلية المتقاطعة المعممة (Generalized Cross Validation) (GCV) ويرمز لها بـ $V(\lambda)$. إذ تستخدم دالة الفاعلية المتقاطعة المعممة في البحث عن القيمة التقديرية $\hat{\lambda}$ المثالية لـ λ من البيانات. إذ إنَّ $\hat{\lambda}$ تقلل المقدار الآتي (Golub, Heath & Wahba, 1978).

$$V(\lambda) = \frac{n^{-1} \|(I - A(\lambda))y\|^2}{[n^{-1} \text{tr}(I - A(\lambda))]^2} \quad \dots (5)$$

إذ إنَّ $A(\lambda)$ هي مصفوفة القبعة (hat matrix) ذات بعد $n \times n$ لمصفوفة انحدار الحرف، وتسمى مصفوفة التمهيد (smoothing matrix) كما تسمى أيضاً (influence matrix)، وتعرف بالآتي (Bates, et al., 1987):

$$A(\lambda) = X(X^T X + n\lambda I)^{-1} X^T \quad \dots (6)$$

على فرض أن \mathcal{R}_m هو فضاء للدوال التي مشتقاتها الجزئية من الرتبة الكلية m ، وعلى فرض أن أنموذج البيانات هو الآتي:

$$y_i = f(\mathbf{x}_i) + \varepsilon_i, \quad i = 1, \dots, n \quad \dots (7)$$

وان $f \in \mathcal{R}_m$.

وإذا كانت λ ثابتاً. فإنه لغرض تقدير f يتطلب تقليل دالة المربعات الصغرى الجزئية

$$\frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i) - \mathbf{z}_i \boldsymbol{\beta})^2 + \lambda J_m(f) \quad \dots (8)$$

إن $m(f)$ هو حد الجزاء لفرض التمهيد على f إذ إنه يتناقص كلما زاد التمهيد.

هناك عدة طرائق لتعريف $m(f)$. بالنسبة لشرائح التمهيد ذات الصفائح الرقيقة، عندما تكون

\mathbf{X} ذات بعد d ، فإنه يمكن تعريف $J_m(f)$ على أنها (Wahba, 1990)

$$J_m(f) = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} \frac{m!}{\alpha_1! \dots \alpha_d!} \left(\frac{\partial^m f}{\partial x_1 \alpha_1 \dots \partial x_d \alpha_d} \right)^2 d_{x_1} \dots d_{x_d} \quad \dots (9)$$

إذ إن $i = \sum_{i=1}^n m$. وفقاً للتعريف في أعلاه، فإن بعض دوال $m(\cdot)$ سيكون لها جزاء

مقداره الصفر. ويسمى الفضاء الذي يمتد بوساطة مجموعة متعددي الحدود (polynomials)

ذات المساهمات الصفرية الجزاء الفضاء متعدد الحدود (polynomial space). إن البعد

للفضاء متعدد الحدود M هو دالة من البعد d والرتبة m للجزء التمهيدي.

عندما $d = 2$ و $m = 2$ ، فإن $m(\cdot)$ تعرف بالآتي:

$$J_2(f) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (f_{x_1 x_1}^2 + 2f_{x_1 x_2}^2 + f_{x_2 x_2}^2) d_{x_1} d_{x_2} \quad \dots (10)$$

إذ إن

$$f_{x_1 x_1}^2 = \left(\frac{\partial^2 f}{\partial x_1^2} \right)^2; \quad f_{x_1 x_2}^2 = \left(\frac{\partial^2 f}{\partial x_1 \partial x_2} \right)^2; \quad f_{x_2 x_2}^2 = \left(\frac{\partial^2 f}{\partial x_2^2} \right)^2$$

وإن M هي الحجم ل $\{1, x_1, x_2\} = 3$.

بشكل عام فإن m و d يجب أن تستوفيا الشرط بأن $2m - d > 0$. لإجل التبسيط، سيتم

التعامل مع $m = 2$ للصيغ والمعادلات في أدناه.

وقد برهن Duchon (1977) على أن f_λ لها الصيغة الآتية (لمزيد من التفاصيل، انظر

(Duchon (1977) [3].

$$f_\lambda(\mathbf{x}_i) = \theta_0 + \sum_{j=1}^d \theta_j \mathbf{x}_{ij} + \sum_{j=1}^d \delta_j \mathbf{E}_2(\mathbf{x}_i - \mathbf{x}_j) \quad \dots (11)$$

إذا إنَّ $\| \cdot \|_2(\mathbf{x}_i - \mathbf{x}_j) = \frac{1}{2^{3\pi}} \|\mathbf{x}_i - \mathbf{x}_j\|^2 \log (\|\mathbf{x}_i - \mathbf{x}_j\|)$ هي المعيار الاقليدي (Euclidean norm) عندما $d = 2$.

إذا كانت \mathbf{K} تحتوي على العناصر $\mathbf{x}_i - \mathbf{x}_j$ ، وإن \mathbf{T} تحتوي على العناصر $\mathbf{T}_{ij} = (\mathbf{X}_{ij})$ ، فإن الهدف هو في إيجاد المعلمات β ، θ و δ التي تقلل المقدار الآتي

$$S_\lambda(\beta, \theta, \delta) = \frac{1}{n} \|\mathbf{y} - \mathbf{T}\theta - \mathbf{K}\delta - \mathbf{Z}\beta\|^2 + \lambda \delta^T \mathbf{K}\delta \quad \dots (12)$$

يمكن الحصول على حل وحيد إذا كانت المصفوفة \mathbf{T} ذات رتبة كاملة (full rank) ، وإن δ يجب ان تستوفي الشرط الآتي $\mathbf{T}^T \delta = 0$.

إذا كانت $\alpha = \begin{pmatrix} \beta \\ \theta \end{pmatrix}$ ، وإن $\mathbf{X} = (\mathbf{T} \mathbf{Z})$ ، فإن الصيغة لـ S_λ ستكون على الآتي:

$$\frac{1}{n} \|\mathbf{y} - \mathbf{X}\alpha - \mathbf{K}\delta\|^2 + \lambda \delta^T \mathbf{K}\delta \quad \dots (13)$$

يمكن الحصول على المعلمات α و δ وذلك بحل المعادلات الآتية:

$$(\mathbf{K} + n\lambda \mathbf{I}_n)\delta + \mathbf{X}\alpha = \mathbf{y} \quad \dots (14)$$

$$\mathbf{X}^T \delta = \mathbf{0} \quad \dots (15)$$

يمكن حساب α و δ وذلك بتحلل QR¹ (QR decomposition) للمصفوفة \mathbf{X} ، وعلى الآتي (Wang, Yuedong)

$$\mathbf{X} = (\mathbf{Q}_1 \quad \mathbf{Q}_2) \begin{pmatrix} \mathbf{R} \\ \mathbf{0} \end{pmatrix} \quad \dots (16)$$

إذا إن \mathbf{Q}_1 مصفوفة ذات بعد $n \times p$ و \mathbf{Q}_2 مصفوفة ذات بعد $n \times (n - p)$ و \mathbf{R} هي مصفوفة مثلثية عليا ذات بعد $p \times p$ ، وإن $(\mathbf{Q}_1 \quad \mathbf{Q}_2)$ مصفوفة متعامدة (orthogonal matrix) ، وإن \mathbf{R} هي مصفوفة مثلثية عليا ، وإن $\mathbf{X}^T \mathbf{Q}_2 = \mathbf{R}^T \mathbf{Q}_1^T \mathbf{Q}_2 = \mathbf{0}$.

بما أن $\mathbf{Q}_1^T \delta = \mathbf{0}$ كما في المعادلة (15) ، فإن δ يجب ان تكون ضمن فضاء العمود لـ \mathbf{Q}_2 . لذا فإن δ يمكن صياغتها على أنها $\delta = \mathbf{Q}_2 \gamma$ للمتجه γ . وبالتعويض عن $\delta = \mathbf{Q}_2 \gamma$ في المعادلة السابقة وبالضرب بـ \mathbf{Q}_2^T ، نحصل على

1 يمكن تحلل أية مصفوفة مربعة A ، وذلك بأن $A = QR$ ، اذا ان Q هي مصفوفة متعامدة (orthogonal matrix) وان R هي المصفوفة المثلثية العليا (upper triangular matrix).

$$Q_2^T(K + n\lambda I)Q_2\gamma = Q_2^T y \quad \dots (17)$$

أو

$$Q_2\gamma = [Q_2^T(K + n\lambda I)]^{-1} Q_2^T y \quad \dots (18)$$

$$\delta = Q_2\gamma = Q_2[Q_2^T(K + n\lambda I)Q_2]^{-1} Q_2^T y \quad \dots (19)$$

يمكن الحصول على α بحل المعادلة الآتية:

$$R\alpha = Q_1^T[y - (K + n\lambda I)\delta] \quad \dots (20)$$

وإذا تم افتراض أن الأثر (trace) لـ $A(\lambda)$ هي درجات الحرية للنموذج، وأن الأثر لـ $(I - A(\lambda))$ هي درجات الحرية للخطأ، فإن σ^2 يمكن أن تأخذ الصيغة الآتية:

$$\hat{\sigma}^2 = \frac{RSS(\lambda)}{tr(I - A(\lambda))} \quad \dots (21)$$

إذ إن $RSS(\lambda)$ هو مجموع مربعات البواقي (Residuals Sum of Squares).

2.1: حدود الثقة:

تعد حدود الثقة من المقاييس الإحصائية التي لا يقتصر الاستفادة منها على تقدير معلمة المجتمع، بل تتعدى ذلك إلى إعطاء معلومات عن دقة التقدير. حينما تكون حدود الثقة متباعدة فهي دالة على الكثير من عدم اليقين بقيمة معلمة المجتمع، وعندما تكون حدود الثقة متقاربة فتدل على أن التقدير محدد بدقة كبيرة (Reed, Pileggi & Winkelman, 1973) (Terry & Kelley, 2012).

2.1.1: بناء حدود الثقة للبوستراب (Wang & Wahba, 1995)

لنكن $\{x_i\}_{i=1,n}$ نقاط تصميم ثابتة (fixed design points)، ولنكن $\hat{f}_{\hat{\lambda}}$ و $\hat{\sigma}^2$ تقديرات لكل من f و σ^2 من البيانات. وعلى افتراض أن $\hat{f}_{\hat{\lambda}}$ هي دالة حقيقية لـ f ، فإنه يمكن اتباع الخطوات الآتية لغرض بناء حدود الثقة للبوستراب:

1. توليد عينة بوستراب بحسب الصيغة الآتية:

$$y_i^* = \hat{f}_{\hat{\lambda}}(x_i) + \epsilon_i^*$$

إذ إن $\epsilon_i^* = (\epsilon_1^*, \dots, \epsilon_n^*)^T \sim N(0, \hat{\sigma}^2 I_{n \times n})$.

2. إيجاد تقدير للصفائح الممهدة \hat{f}_λ^* بالاستناد على عينة البوتستراب.
 3. افرض أنّ $f^*(x_i)$ تمثل متغيراً عشوائياً لقيم بوتستراب التقديرية عند x_i .
 4. إعادة الخطوات السابقة B من المرات.
 5. إيجاد D_i^* والقيم $\chi_{1-\alpha/2}$ و $\chi_{\alpha/2}$ ، إذ إنّ $D_i^* = (f_\lambda^*(x_i) - \hat{f}_\lambda(x_i))$
- يلاحظ من الخطوات في أعلاه أنه في كل x_i ، لدينا B من تقديرات البوتستراب لـ $\hat{f}_\lambda(x_i)$ وهي بذلك تحقق B من $(x_i)^*$.

2.1.2: بناء فترات الثقة البيزية:

اقترح (Wahba, 1983) أنموذج الشريحة (Spline model) كأنموذج بيزي (Bayesian model)، واقترح الصيغة الآتية لفترات الثقة البيزية لتمهيد تقديرات الشرائح:

$$\hat{f}_\lambda(x_i) = \pm z_{\alpha/2} \sqrt{\hat{\sigma}^2 a_{ii}(\lambda)}$$

إذ إنّ $a_{ii}(\lambda)$ هو العنصر i للمصفوفة $A(\lambda)$ ، وإنّ $z_{\alpha/2}$ هي الربيعي $1 - \alpha/2$ للتوزيع الطبيعي القياسي. وتفسر فترات الثقة على أنها فترات "عبر الدوال" (across the function) بدلا من الفترات النقطية (pointwise intervals).

2 الجانب العملي:

ينظر الى المحاكاة على أنها الملاذ الأخير الذي يتم توظيفه في حالة صعوبة الحصول على حلول ممكنة باستخدام منهجيات الحل التحليلي. إنّ الأساليب الحديثة المتقدمة في المحاكاة، نتيجة لتوفر البرامجيات، فضلاً عن تطور التقنيات. جعل من المحاكاة واحدة من أكثر الطرائق انتشاراً، والادوات الأكثر قبولا في الوصول إلى معلومات قيمة عن النماذج المستخدمة وعن المتغيرات الأكثر أهمية، فضلاً عن كيفية تفاعل هذه المتغيرات مع بعضها البعض.

تم اجراء تجارب المحاكاة على أنموذجين، الاول: خطي في ظل وجود بيانات ملوثة بمشاهدات شاردة، الآخر بمشاهدات ذات القوة الرافعة، والانموذج الثاني اللاخطي في ظل وجود بيانات ملوثة بمشاهدات شاردة. تم تنفيذ تجارب المحاكاة على حجوم لعينات مختلفة، وهي على الآتي: $n = 100, 90, 80, 70, 60, 50, 40, 30, 20, 15, 10, 5, 4$ وفيما يتعلق بتوصيف النماذج فقد تم إجراء تجارب المحاكاة بموجب الآتي:

1. أنموذج الانحدار الخطي:

a. توليد بيانات ملوثة بمشاهدات شاردة.

- المتغير التوضيحي: تم توليد مشاهدات المتغير التوضيحي وفق التوزيع المنتظم $U(0, 1)$.
- الخطأ العشوائي: تم توليد 90% من الأخطاء وفق التوزيع الطبيعي القياسي $N(0, 1)$ ، و 10% وفق التوزيع الطبيعي $N(0, 10)$ التي تعد على أنها مشاهدات شاردة، ويصطلح عليها حالة نقل المتوسط (Mean shift).
- b. توليد بيانات ملوثة بمشاهدات ذات القوة الرافعة (القيم الشاردة في فضاء المتغير التوضيحي)، فقد تم توليد البيانات وفقاً للآتي:
- المتغير التوضيحي: تم توليد 15% من مشاهدات المتغير التوضيحي وفق التوزيع الطبيعي $N(0, 10)$.
- الخطأ العشوائي: تم توليد 75% من الأخطاء وفق التوزيع الطبيعي القياسي $N(0, 1)$ ، و 25% وفق التوزيع الطبيعي $N(0, 10)$.

2. أنموذج الانحدار اللاخطي:

- المتغير التوضيحي: تم توليد مشاهدات المتغير التوضيحي وفق التوزيع المنتظم $U(0, 1)$.
- تم توليد 90% من الأخطاء وفق التوزيع الطبيعي القياسي $N(0, 1)$ ، و 10% وفق التوزيع الطبيعي $N(0, 10)$ ، التي تعد على أنها مشاهدات شاردة.
- الانموذج اللاخطي وصيغته هي $y = 5 \sin(3x) + e$.

ابتداءً يتم توفيق أنموذج التقدير الخطي واللاخطي للبيانات التي يتم توليدها والمتضمنة للمشاهدات الشاردة، فضلاً عن المشاهدات ذات القوة الرافعة باستخدام طريقة المربعات الصغرى الجزئية.

و بالاعتماد على التقديرات التي يتم الحصول عليها من توفيق أنموذجي التقدير الخطية و اللاخطي، تم توليد 1200 عينة بوتستراب (لمزيد من التفاصيل عن حجوم العينات في البوتستراب، انظر Efron & Tibshirani, 1994). بما أن بعضاً من عينات البوتستراب يمكن أن تتسبب في مشاكل معينة، لذا فإن هذه العينات يتم فرزها واستبعادها من قبل دالة الفاعلية المتقاطعة المعممة. لذا فإن الـ 200 عينة من البوتستراب هي لأجل تعويض النقص في العينات التي يتم استبعادها. ويتم توفيق جميع المتغيرات $(y_1, y_2, \dots, y_{1200})$ كمتغيرات استجابة، ثم يتم توفيق جميع عينات البوتستراب آنياً. وبالتالي إيجاد i^* والقيم $\chi_{\alpha/2}$ و $\chi_{1-\alpha/2}$.

يتضمن الانموذج قيد الدراسة متغير استجابة مستمراً، ومتغيراً توضيحياً الذي عُدَّ على أنه متغير تمهيدي (smoothing). إن رتبة الاشتقاق في الجزء (Order of Derivative in the Penalty Smoothing) هي (2) ولجميع التجارب، وهي ناتج حساب الصيغة الآتية $\max(2, \text{int}(\frac{d}{2}) + 1)$ ، إذ إن d هي عدد المتغيرات التمهيديّة (Smoothing Variables). وإن بُعِدَ الفضاء المتعدد الحدود (Dimension of Polynomial Space) هو (2) وهو يُمثل كلاً من الحد الثابت والمتغير التوضيحي. وفيما يتعلق بالانحراف المعياري (Standard Deviation) فيتم ايجاده من خلال الصيغة الآتية:

$$SD = \sqrt{\left(\frac{RSS(\lambda)}{\text{tr}(I - A)}\right)}$$

1.2 تحليل تجارب المحاكاة:

تم تحليل نتائج تجارب المحاكاة للانموذج الخطي بحجوم عينات مختلفة في ظل وجود المشاهدات الشاردة، فضلاً عن المشاهدات ذات القوة الرافعة، وللانموذج اللاخطي بحجوم عينات مختلفة في ظل وجود المشاهدات الشاردة. وتم تفسير النتائج وفقاً للرسوم البيانية، الجزء التمهيدي (Smoothing Penalty)، الانحراف المعياري، فضلاً عن دالة الفاعلية المتقاطعة المعممة.

الأشكال 1، 3، و 5 هي الرسوم البيانية لما ورد في الجداول (1) الى (3) على التوالي، وذلك لغرض زيادة الإيضاح، ولتسهيل تفسير النتائج.

الأشكال (2)، (4) و (6) تمثل المنحنيات التقديرية وحدود الثقة البيزية والبوتستراب والتقليدية باستخدام طريقة المربعات الصغرى الجزئية، ويلاحظ من هذه الرسوم دقة التقدير ولجميع تجارب المحاكاة.

1.3 مناقشة النتائج:

a. الانموذج الخطي في ظل وجود المشاهدات الشاردة:

من الشكل (2) يلاحظ أنّ حدود الثقة بوتستراب هي الأفضل من حدود الثقة البيزية والتقليدية في ظل وجود المشاهدات الشاردة ولحد حجم العينة 10، ثم تصبح حدود الثقة البيزية هي الأفضل. كما يلاحظ أنّ حدود الثقة البيزية أفضل من التقليدية ولجميع حجومات العينات.

يلاحظ من الشكل (1) والجدول (1) بأن الـ GCV تراوحت بين 0 حينما حجم العينة 4، والـ 61.8 حينما حجم العينة 60، كما يلاحظ تناغم تذبذب الـ SD مع الـ GCV. فيما يتعلق بالجزء التمهيدي، فيلاحظ أن قيمته قد ارتفعت بشكل كبير عند حجم العينة 20، وبشكل أكبر عند حجم العينة 60. أما ما يتعلق بـ (Loglambda) فهو ناتج اللوغاريتم الطبيعي لحاصل ضرب معلمة التمهيد بحجم العينة $\log_{10}(n\lambda)$.

b. الانموذج الخطي في ظل وجود المشاهدات ذات القوة الرافعة:

من الشكل (4) يلاحظ أن حدود الثقة بوتستراب هي الأفضل من حدود الثقة البيزية والتقليدية في ظل وجود المشاهدات الشاردة ولجميع حجوم العينات. كما يلاحظ أن حدود الثقة البيزية أفضل من التقليدية ولجميع حجوم العينات.

يلاحظ من الشكل (2) والجدول (2) بأن الـ GCV تراوحت بين 57 حينما حجم العينة 90، والـ 1719 حينما حجم العينة 40. فيما يتعلق بالجزء التمهيدي، فيلاحظ أنه مستقر عند الصفر ولجميع حجوم العينات.

c. الانموذج اللاخطي في ظل وجود المشاهدات الشاردة:

من الشكل (6) يلاحظ أن حدود ثقة البوتستراب أفضل من حدود الثقة البيزية والتقليدية ولجميع التجارب. كما يلاحظ أن هناك تشابهاً كبيراً بين حدود الثقة البيزية والتقليدية لجميع التجارب ما عدا حينما يكون عدد المشاهدات 4، إذ تتفوق حدود الثقة البيزية على كل من حدود الثقة بوتستراب والتقليدية.

يلاحظ من الشكل (3) والجدول (3) بأن الـ GCV تراوحت بين 0.7-3.07 للاحجام العينات من 100 الى حجم العينة 5، إلا أنها قفزت الى 54.3 عند حجم العينة 4، وبالرجوع الى الشكل (6) يلاحظ أن تشتت النقاط كان كبيراً، مما أدى الى عدم تمثيل النقاط بشكل جيد، وهذا انعكس بدوره على الـ SD، وكما هو موضح ذلك في الشكل (5). أما فيما يتعلق بالجزء التمهيدي فيلاحظ أن قيمته قد ارتفعت بشكل كبير عند حجم العينة 60 وبشكل أكبر عند حجم العينة 20، وعند الرجوع الى الشكل (5)، ويلاحظ أن المنحني يعاني من تذبذب عند حجم العينة 20، وتذبذب أكبر عند حجم العينة 60.

1.4 الاستنتاجات:

1. تفوق حدود الثقة بوتستراب على حدود الثقة البيزية والتقليدية لتجارب المحاكاة ولحجوم العينات أكبر، او يساوي الـ 10 حينما تكون البيانات ملوثة بالنقاط الشاردة للانموذج الخطي.

2. أظهرت حدود الثقة بوتستراب تفوقاً على حدود الثقة البيزية والتقليدية لتجارب المحاكاة ولجميع العينات حينما تكون البيانات ملوثة بالنقاط ذات القوة الرافعة للانموذج الخطي.
3. تقدم حدود الثقة بوتستراب على حدود الثقة البيزية والتقليدية لتجارب المحاكاة وبحجوم العينات اكبر، او يساوي الـ 5 حينما تكون البيانات ملوثة بالنقاط الشاردة للانموذج اللاخطي.
4. لوحظ أنَّ حدود الثقة البيزية تتفوق على حدود الثقة التقليدية لجميع حجوم العينات ولجميع تجارب البحث.

المصادر

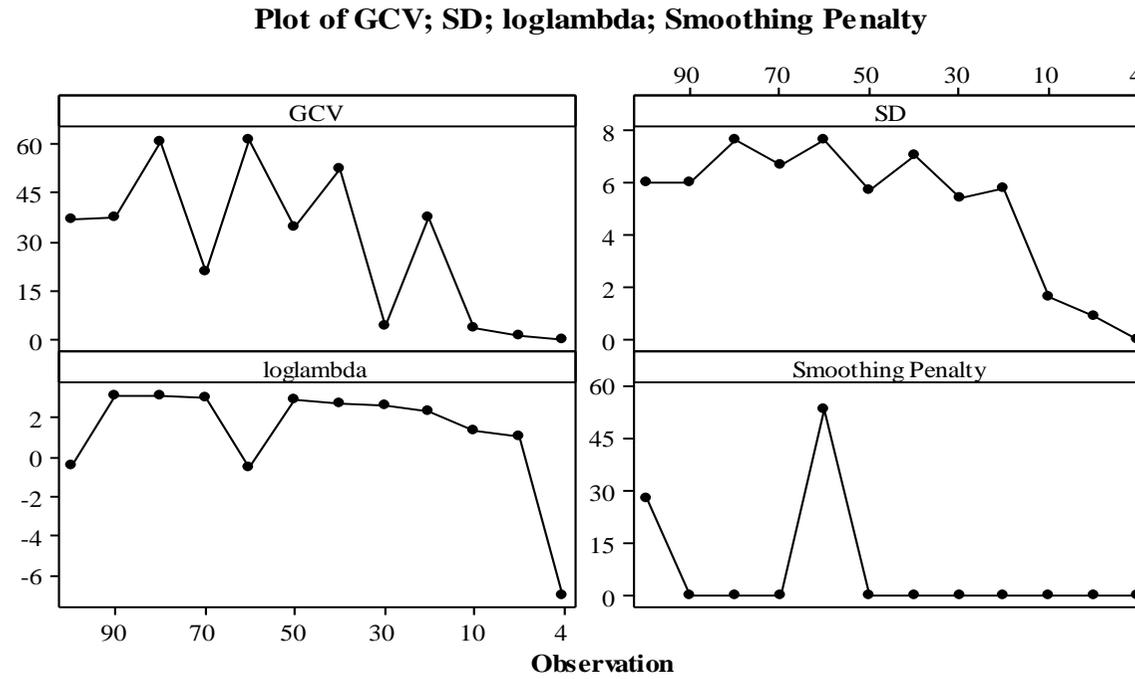
1. رشيد، ظافر حسين وخلود يوسف خمو، (2005)م " مقارنة الطرائق الشرائحية لتقدير منحى الانحدار اللامعلمي " ، المجلة العراقية للعلوم الاحصائية، العدد 8، ص-ص: 62-40.
2. Bates, D. M., Lindstrom, M. J., Wahba, G., & Yandell, B. S. (1987). Gcvpack-routines for generalized cross validation: Gcvpack-routines for generalized. *Communications in Statistics-simulation and Computation*, 16(1), 263-297.
3. Duchon, J. (1977). Splines minimizing rotation-invariant seminorms in Sobolev spaces. In *Constructive theory of functions of several variables* (pp. 85-100). Springer, Berlin, Heidelberg.
4. Golub, G. H., Heath, M., & Wahba, G. (1979). Generalized cross-validation as a method for choosing a good ridge parameter. *Technometrics*, 21(2), 215-223.
5. Efron, B., & Tibshirani, R. J. (1994). *An introduction to the bootstrap*. CRC press.
6. Nosedal-Sanchez, A., Storlie, C. B., Lee, T. C., & Christensen, R. (2012). Reproducing kernel Hilbert spaces for penalized regression: A tutorial. *The American Statistician*, 66(1), 50-60.
7. Nosedal-Sanchez, A., Storlie, C. B., Lee, T. C., & Christensen, R. (2012). Reproducing kernel Hilbert spaces for penalized regression: A tutorial. *The American Statistician*, 66(1), 50-60.
8. Reed, A. H., Cannon, D. C., Pileggi, V. J., & Winkelman, J. W. (1973). Use of confidence intervals to assess precision of normal range estimates. *Clinical biochemistry*, 6, 29-33.
9. Ruppert, D., Wand, M. P., & Carroll, R. J. (2003). *Semiparametric regression* (No. 12). Cambridge university press.

10. Simonoff, J. S. (2012). *Smoothing methods in statistics*. Springer Science & Business Media.
11. Ruan, Z., Guo, H., Liu, H., & Yan, S. (2013). Glacier surface velocity estimation in the West Kunlun Mountain range from L-band ALOS/PALSAR images using modified synthetic aperture radar offset-tracking procedure. *Journal of Applied Remote Sensing*, 8(1), 084595.
12. Wahba, G. (1983). "Bayesian" confidence intervals" for the cross-validated smoothing spline. *Journal of the Royal Statistical Society. Series B (Methodological)*, 133-150.
13. Wahba, G. (1990). *Spline models for observational data*. Society for industrial and applied mathematics.
14. Wang, Y., & Wahba, G. (1995). Bootstrap confidence intervals for smoothing splines and their comparison to Bayesian confidence intervals. *Journal of Statistical Computation and Simulation*, 51(2-4), 263-279.
15. Wang, Y. (2011). *Smoothing splines: methods and applications*. CRC Press.
16. Wang, Y. (2011). *Smoothing splines: methods and applications*. Chapman and Hall/CRC.
17. Wood, S. N. (2003). Thin plate regression splines. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 65(1), 95-114.
18. Terry, L., & Kelley, K. (2012). Sample size planning for composite reliability coefficients: Accuracy in parameter estimation via narrow confidence intervals. *British Journal of Mathematical and Statistical Psychology*, 65(3), 301.

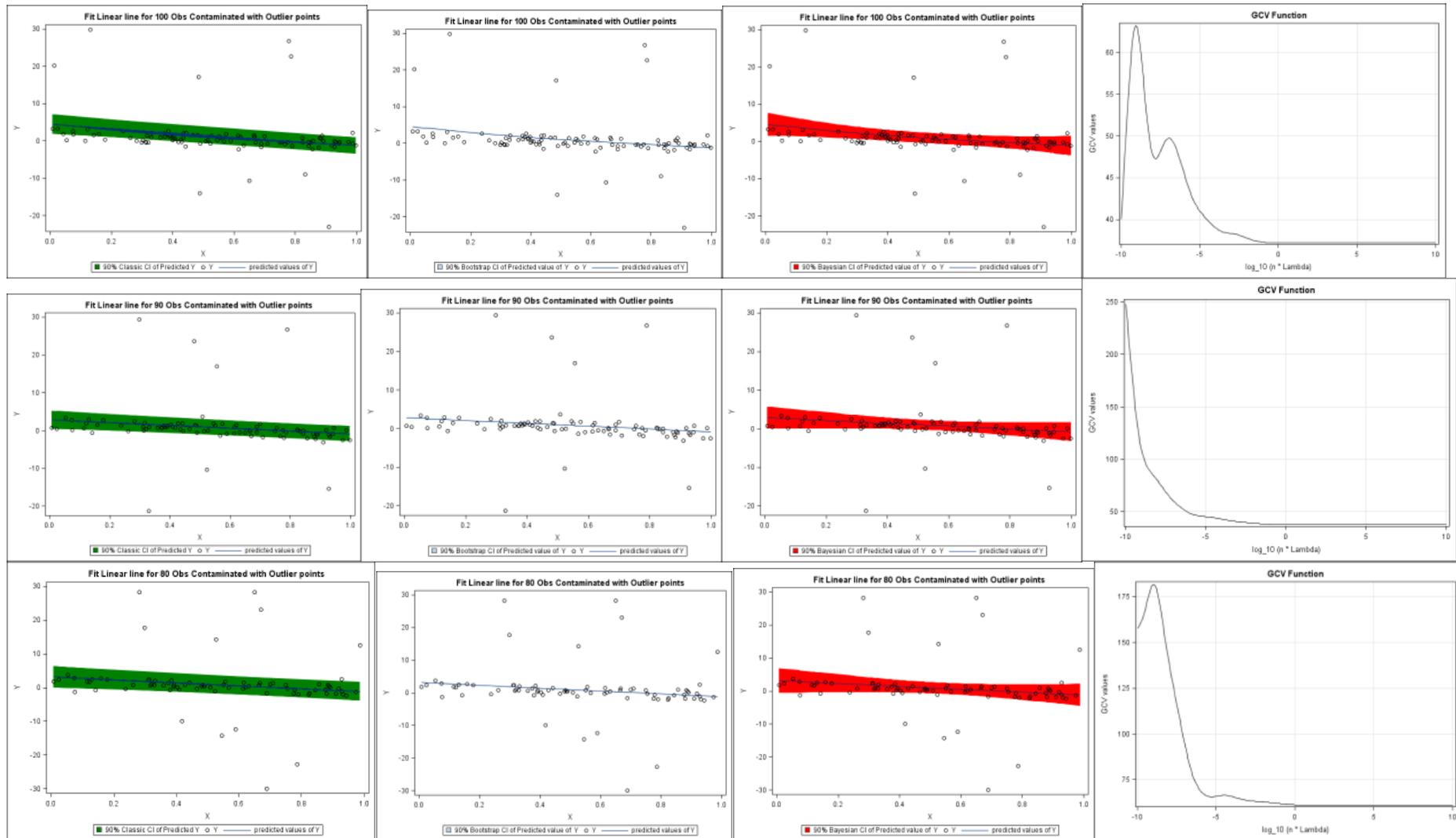
الملحق

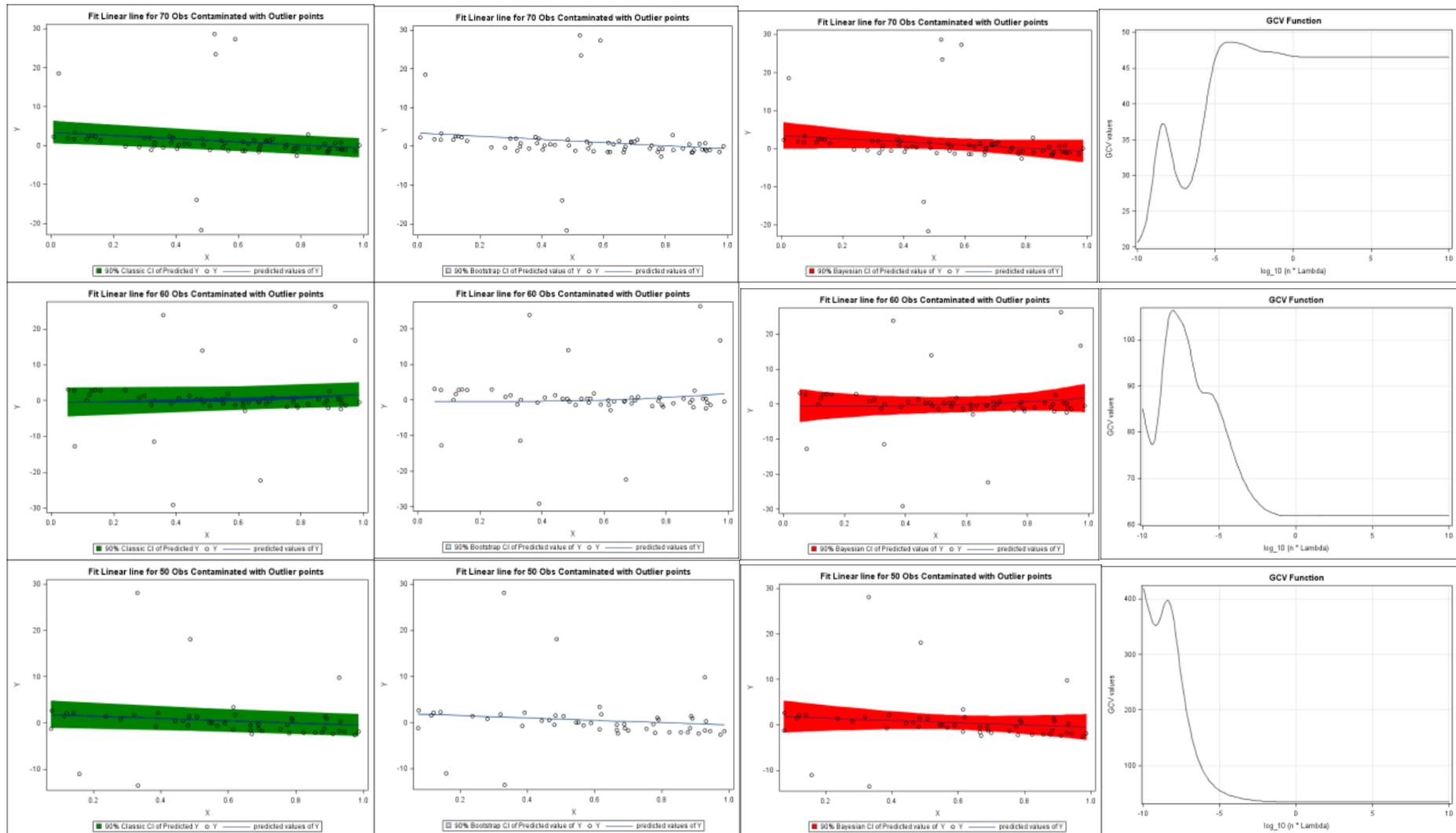
الجدول (1): يوضح الإحصاءات للتقديرات النهائية للانموذج الخطي في ظل وجود المشاهدات الشاردة، ولحجوم عينات مختلفة

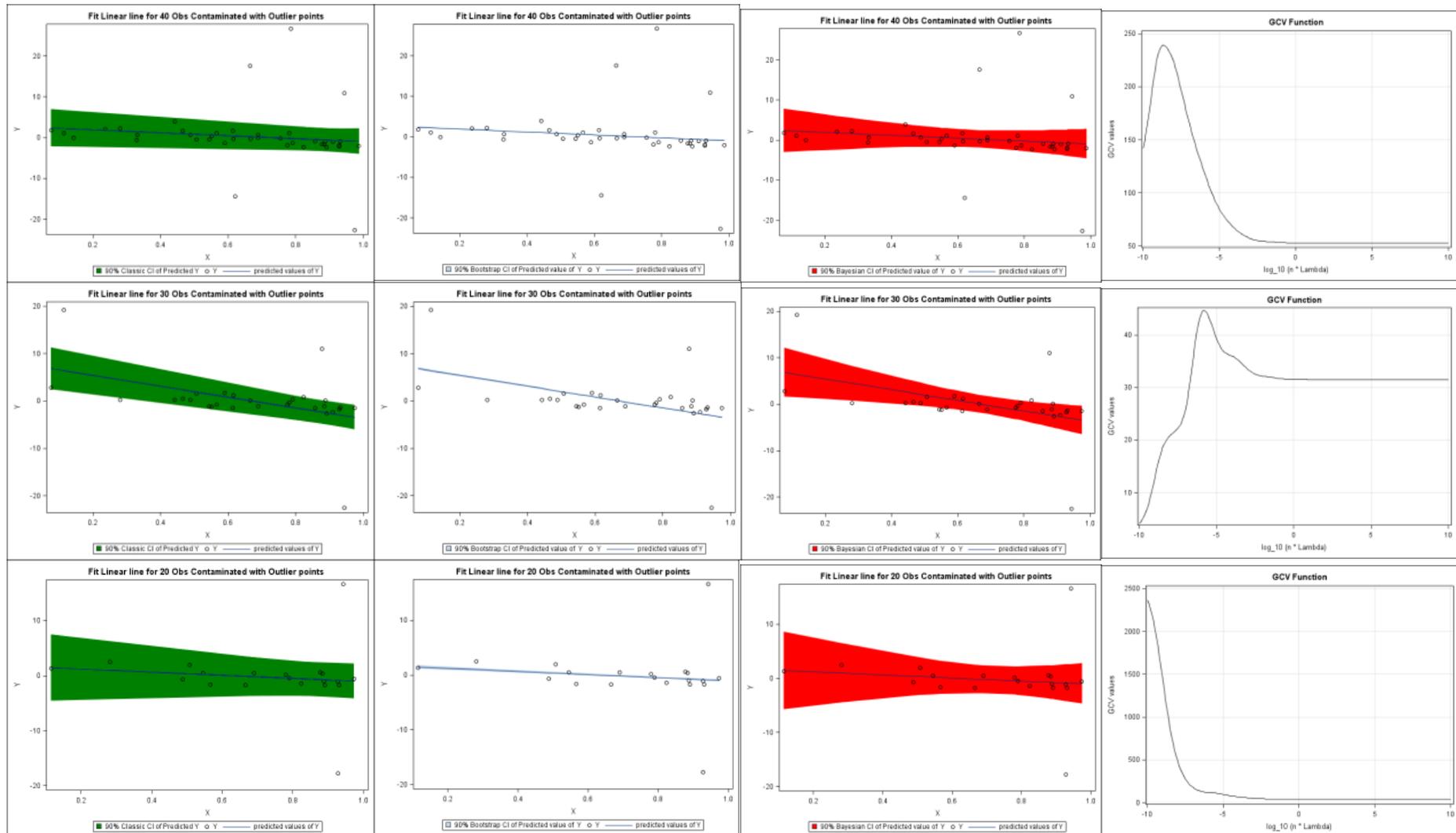
NO of observations	loglambda	SMOOTHING PENALTY	Residual	DF	SD	GCV
100	-0.3187	27.7405	3547.9786	2.3713	6.0284	37.224314
90	3.2303	0.0000	3232.5208	2.0001	6.0608	37.568138
80	3.1777	0.0000	4645.6375	2.0001	7.7175	61.086662
70	3.1484	0.0000	3073.5613	2.0001	6.7231	20.636316
60	-0.4702	53.3494	3430.1515	2.3026	7.7104	61.823631
50	2.9526	0.0000	1596.5690	2.0001	5.7673	34.647872
40	2.8188	0.0000	1905.6808	2.0001	7.0816	52.789193
30	2.7059	0.0000	824.0716	2.0001	5.4251	4.087004
20	2.4289	0.0000	612.3235	2.0001	5.8325	37.797773
10	1.4288	0.0000	21.9219	2.0001	1.6554	3.425383
5	1.1420	0.0000	2.1475	2.0001	0.8461	1.193091
4	-6.9773	0.0000	0.0000	3.5107	0.0000	0.0000

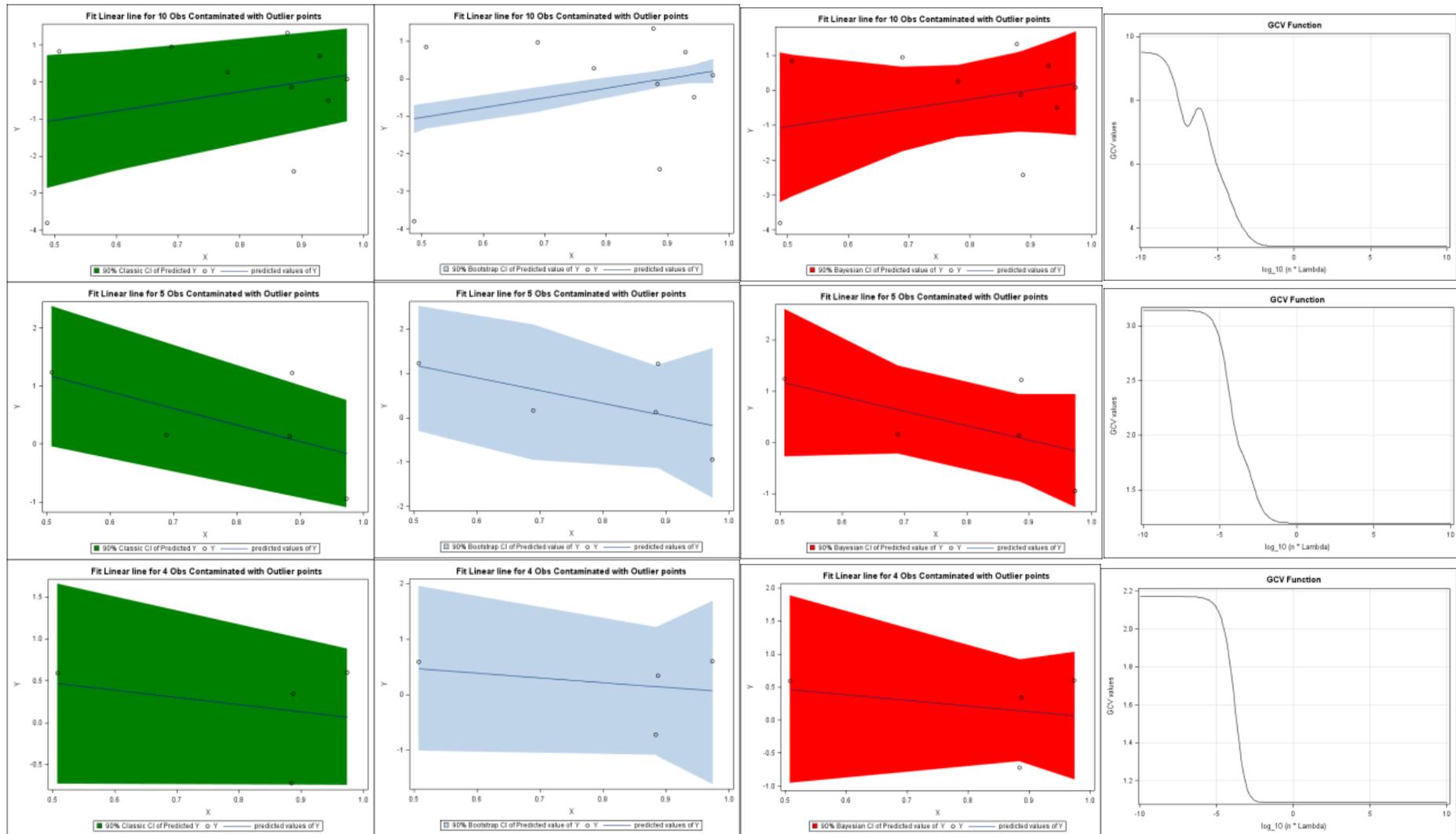


الشكل (1): يوضح الرسوم البيانية لإحصاءات التقديرات النهائية الواردة في الجدول (1) للنموذج الخطي في ظل وجود المشاهدات الشاردة، ولحجوم عينات مختلفة







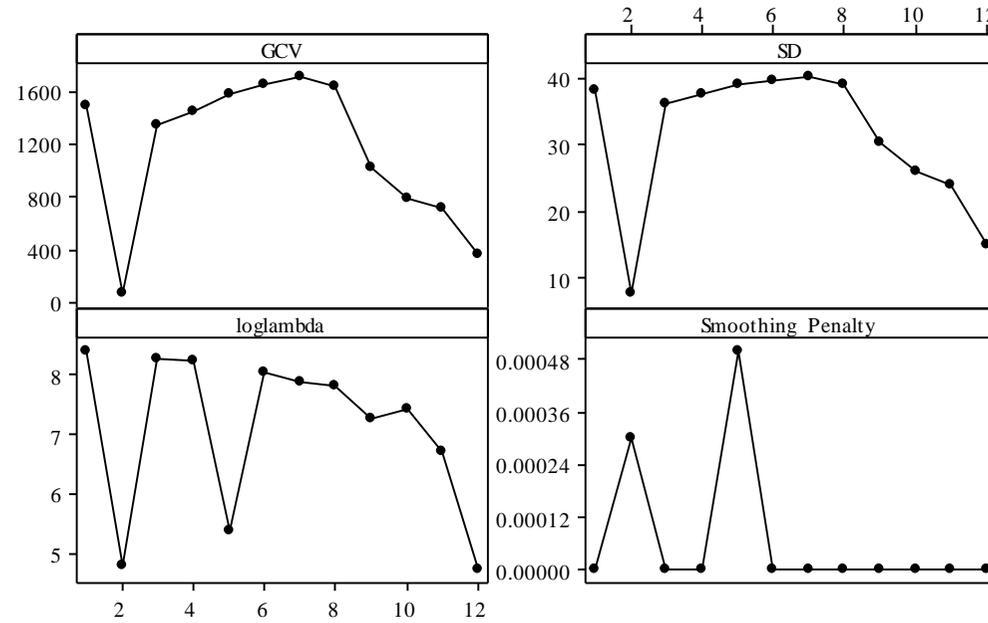


الشكل (1): يوضح الرسوم البيانية لكل من حدود الثقة بوتستراب والبيزية والكلاسيكية، فضلاً عن الرسم البياني لدالة الصلاحية التقاطع المعممة للانموذج الخطي في ظل وجود المشاهدات الشاردة ولحجوم عينات مختلفة.

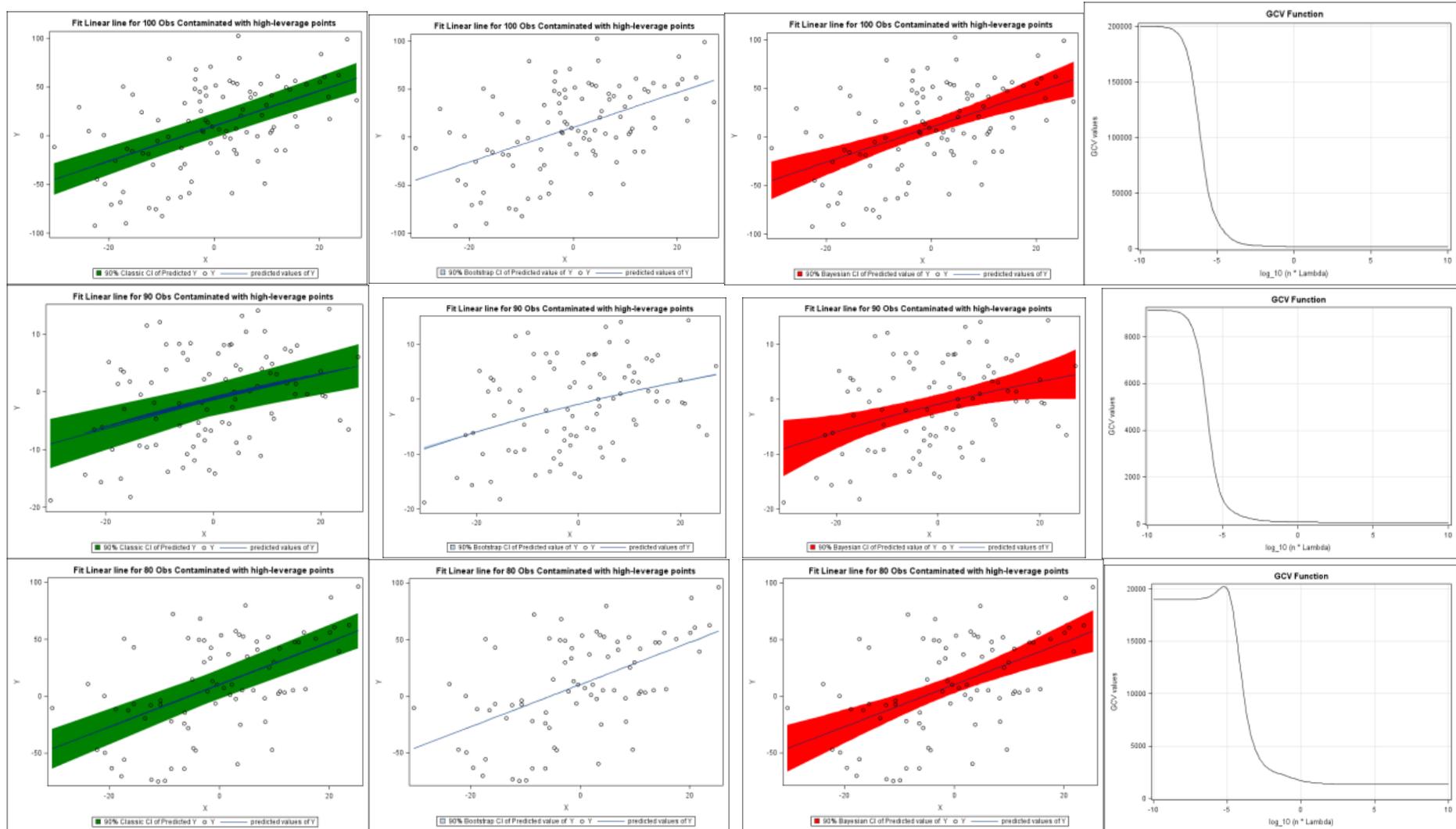
الجدول (2): يوضح الاحصاءات للتقديرات النهائية للانموذج الخطي في ظل وجود المشاهدات ذات القوة الرافعة ولحجوم عينات مختلفة

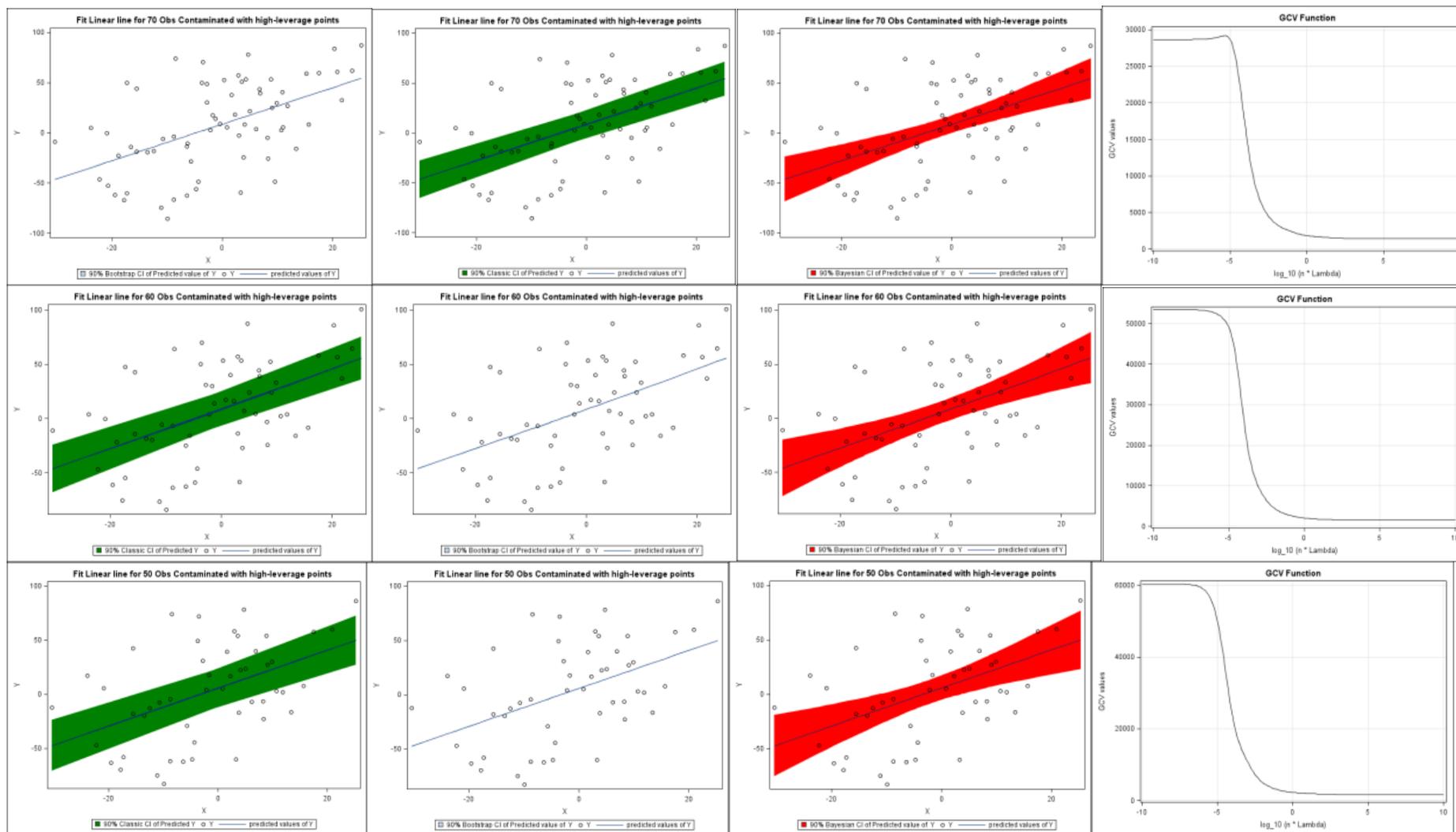
NO of observations	loglambda	SMOOTHING PENALTY	Residual	DF	SD	GCV
100	8.4041	0.0000	143428.561	2.0001	38.2565	1493.42637
90	4.8019	0.0003	4929.3933	2.3314	7.4985	57.566274
80	8.2682	0.0000	102752.408	2.0001	36.2952	1349.71252
70	8.2258	0.0000	95948.9183	2.0001	37.5635	1452.51794
60	5.3884	0.0005	88187.2184	2.0741	39.0181	1576.92778
50	8.0470	0.0000	76176.7415	2.0001	39.8374	1653.14797
40	7.8807	0.0000	62057.5531	2.0001	40.4116	1719.05042
30	7.8290	0.0000	43030.7973	2.0001	39.2023	1646.59402
20	7.2660	0.0000	15605.2598	2.0001	30.2979	1025.96014
15	7.4396	0.0000	8871.0250	2.0001	26.1227	787.372401
10	6.7249	0.0000	4579.9019	2.0001	23.9269	715.623698
5	4.7226	0.0000	652.0314	2.0001	14.7429	362.248500

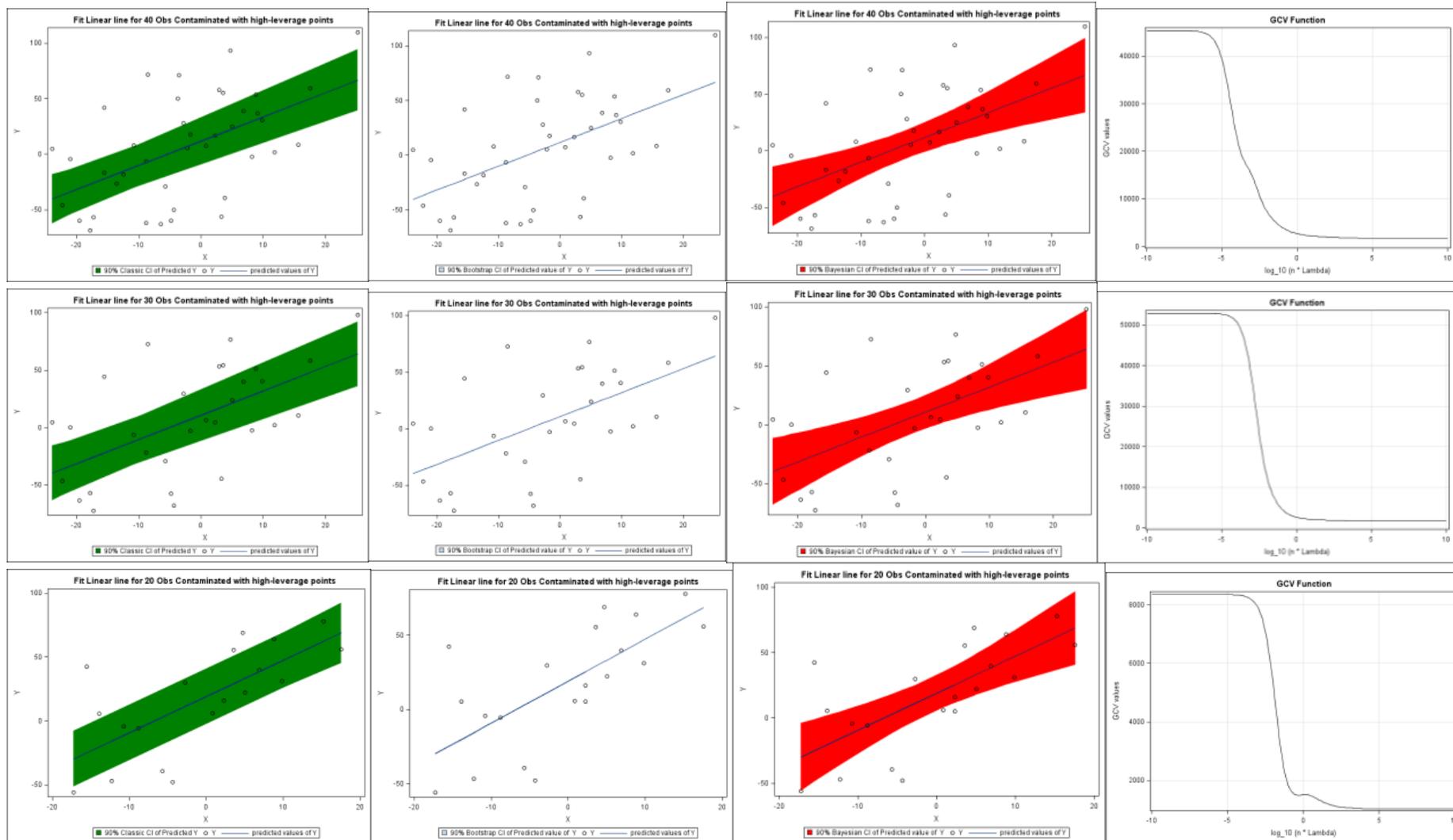
Plot of GCV; SD; loglambda; Smoothing Penalty

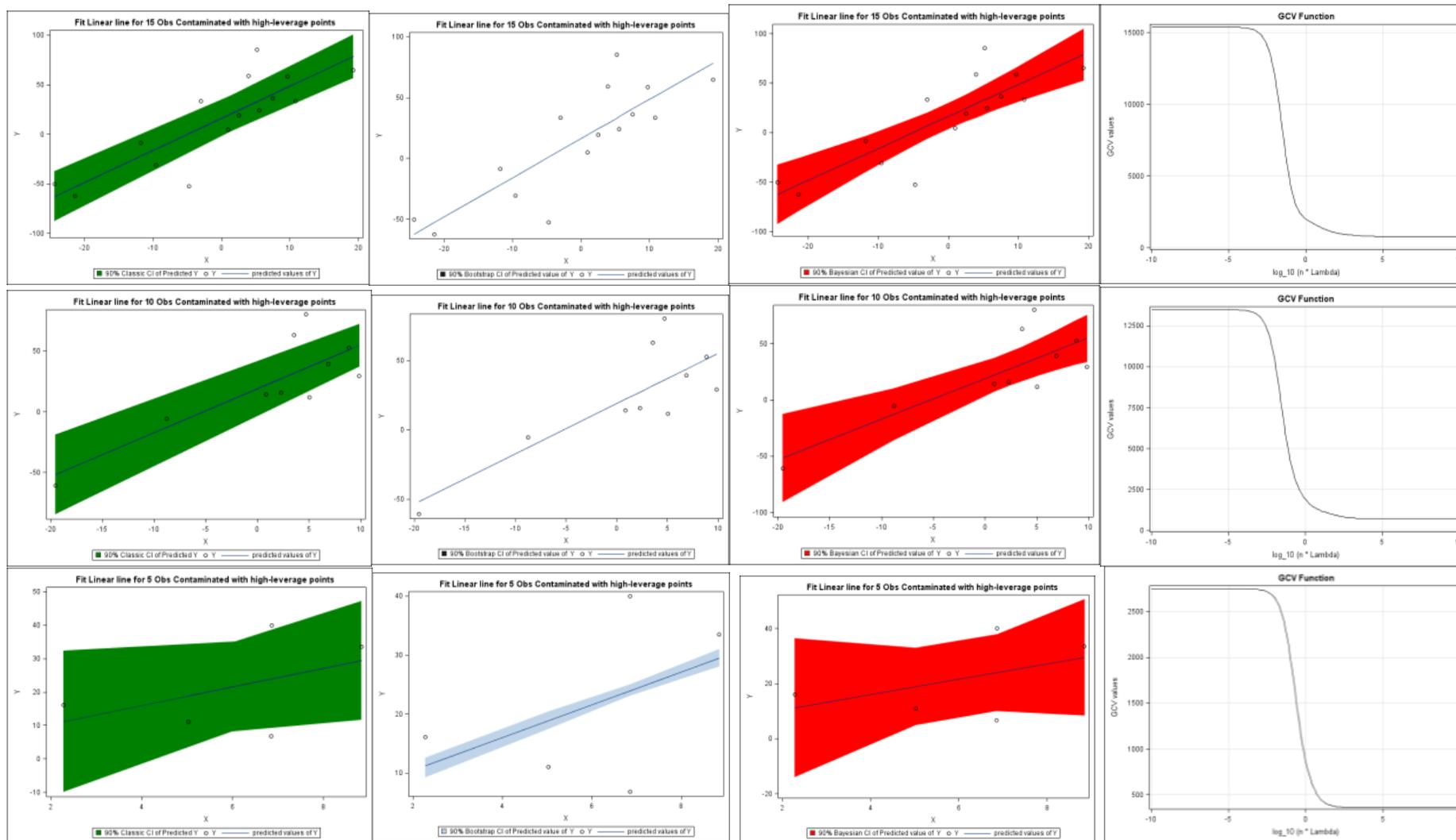


الشكل (3): يوضح الرسوم البيانية لاحصاءات التقديرات النهائية الواردة في الجدول (2) للانموذج الخطي في ظل وجود المشاهدات ذات القوة الرافعة ولحجوم عينات مختلفة







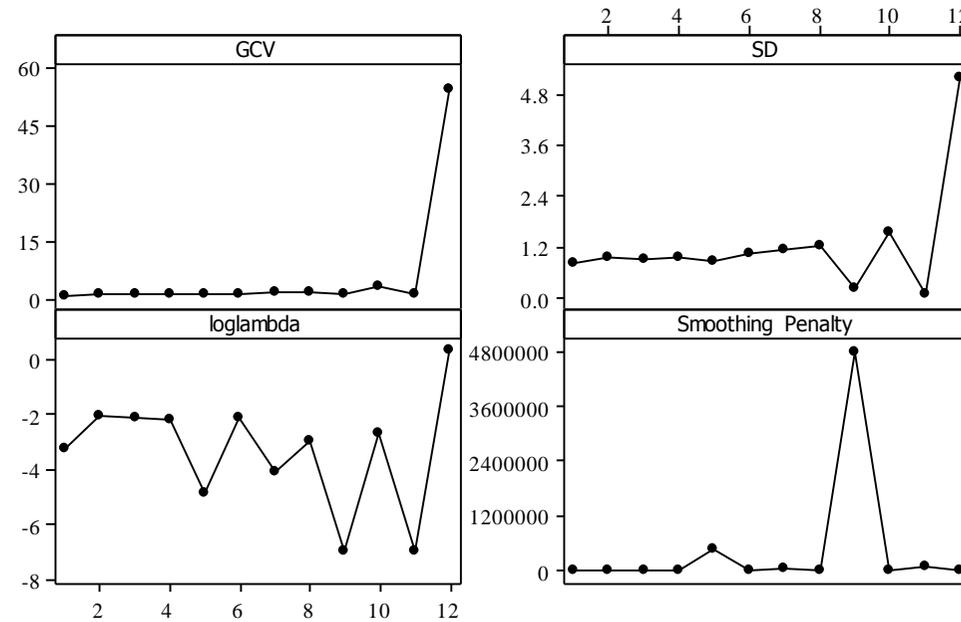


الشكل (4): يوضح الرسوم البيانية لكل من حدود الثقة بوتستراب والبيزية والكلاسيكية، فضلاً عن الرسم البياني لدالة الصلاحية التقاطع المعممة للانموذج الخطي في ظل وجود المشاهدات ذات القوة الرافعة ولحجوم عينات مختلفة.

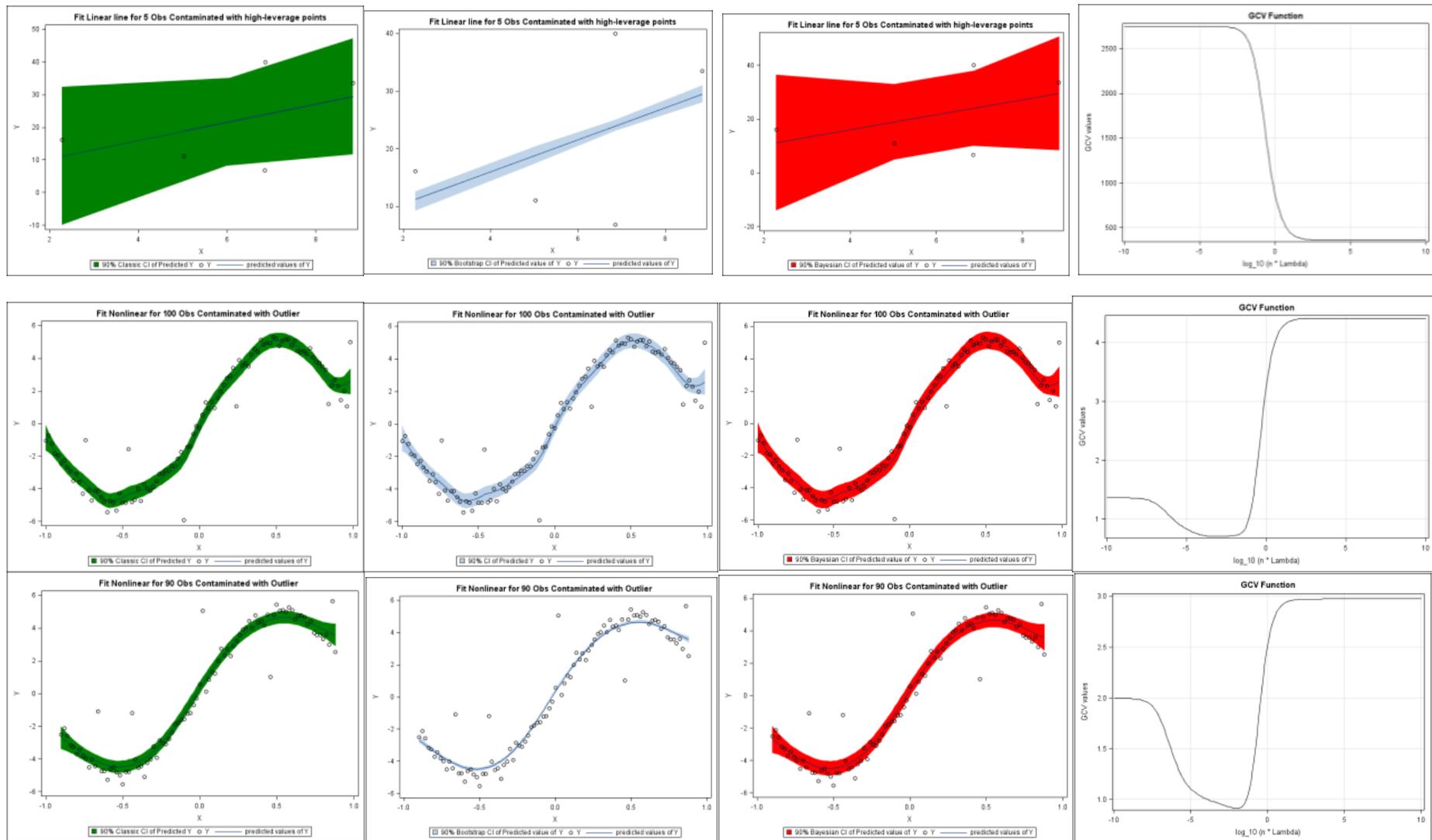
الجدول (3): يوضح الاحصاءات للتقديرات النهائية للانموذج اللاخطي في ظل وجود المشاهدات الشاردة

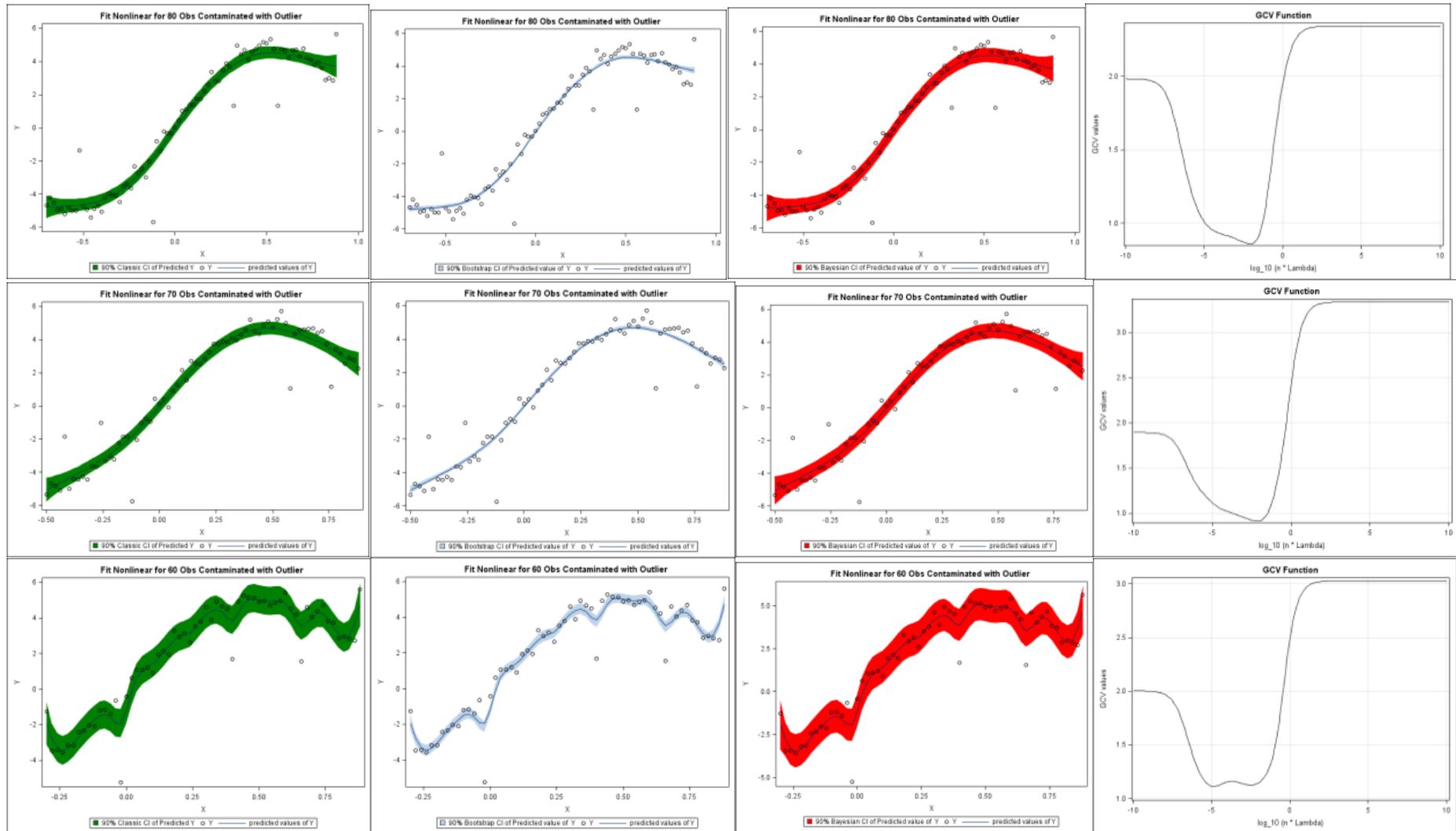
NO of observations	loglambda	SMOOTHING PENALTY	Residual	DF	SD	GCV
100	-3.2311	7055.6430	53.5584	13.0682	0.7849	0.708726
90	-2.0069	1215.1608	70.9038	6.3704	0.9208	0.912419
80	-2.0935	1006.5407	58.8419	6.0172	0.8918	0.860036
70	-2.1718	1103.3531	54.1068	5.5919	0.9165	0.913078
60	-4.8731	486287.3070	30.3667	19.5572	0.8665	1.114034
50	-2.1157	647.0535	50.4123	4.1750	1.0489	1.200351
40	-4.0394	50432.7360	38.1933	8.6672	1.1041	1.556276
30	-2.9831	2655.8819	38.9258	4.1340	1.2267	1.745452
20	-6.9773	4796272.1017	0.0352	19.1236	0.2005	0.905041
10	-2.6935	580.6901	17.3717	2.4862	1.5205	3.077003
5	-6.9773	92441.7743	0.0001	4.9823	0.0568	0.907823
4	0.3291	0.0013	54.2786	2.0001	5.2097	54.284200

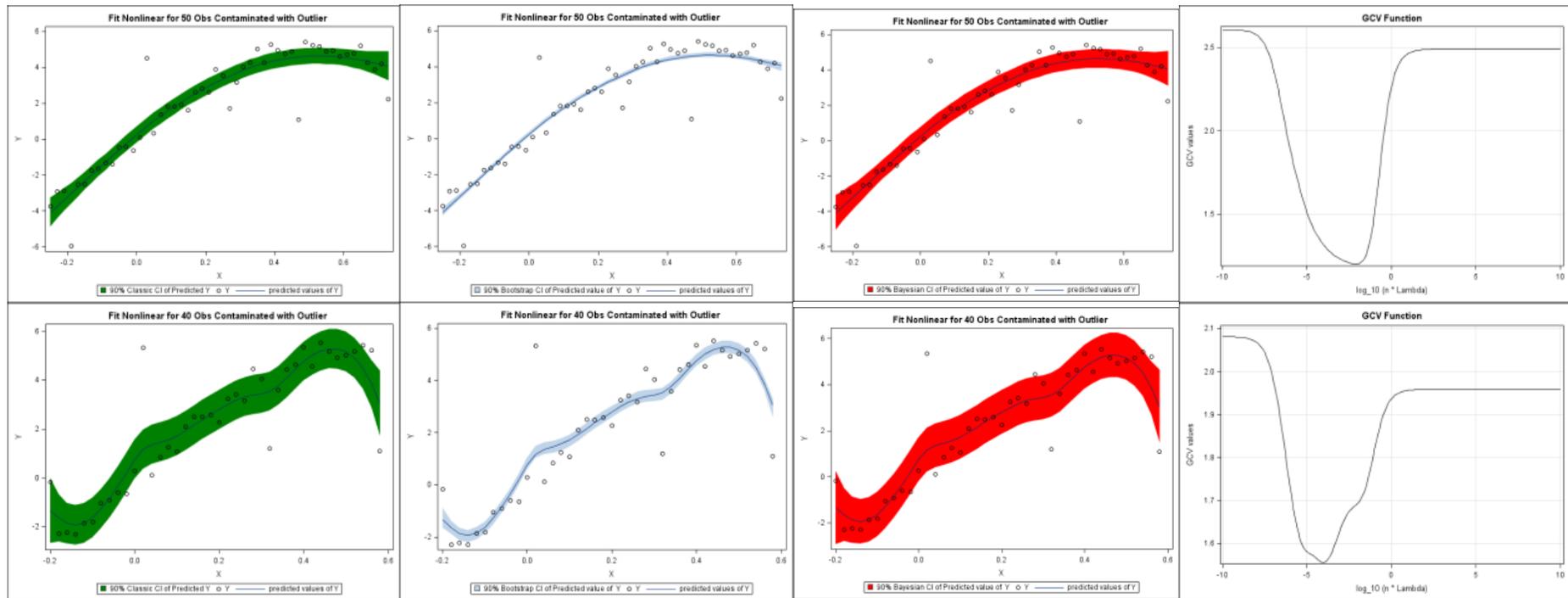
Plot of GCV; SD; loglambda; Smoothing Penalty

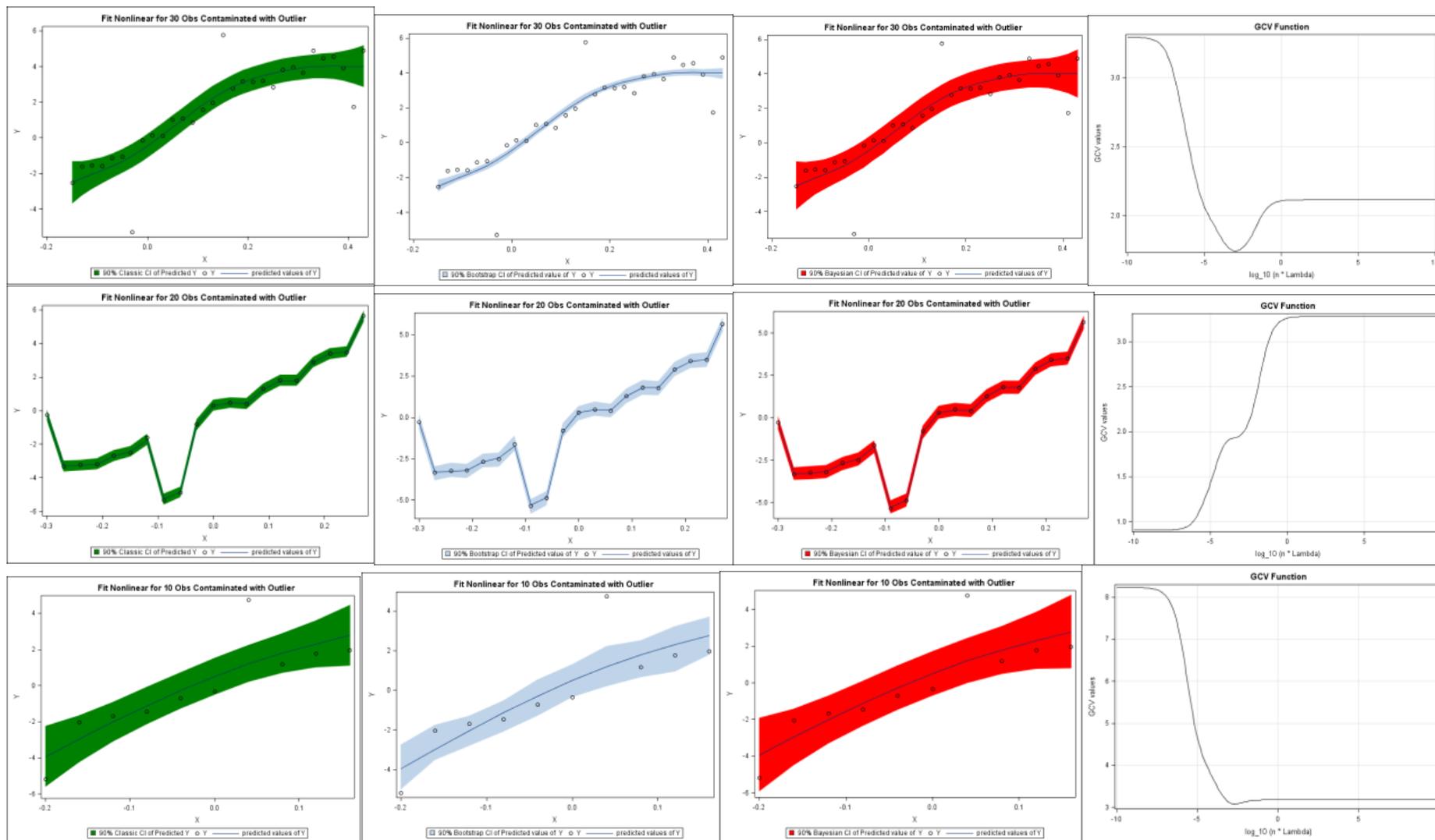


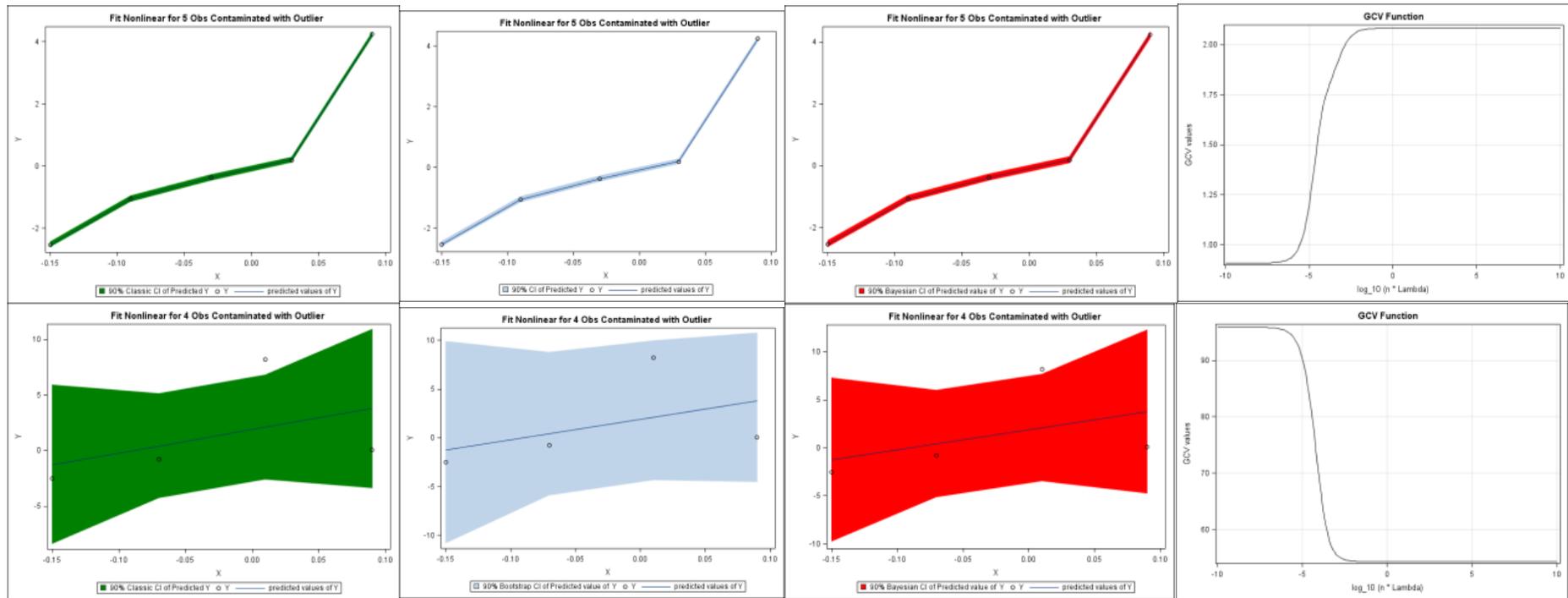
الشكل (5): يوضح الرسوم البيانية لإحصاءات التقديرات النهائية الواردة في الجدول (3) للنموذج اللاخطي في ظل وجود المشاهدات الشاردة، ولحجوم عينات مختلفة











الشكل (6): يوضح الرسوم البيانية لكل من حدود الثقة بوتستراب والبيزية والكلاسيكية، فضلاً عن الرسم البياني لدالة الصلاحية التقاطع المعممة للنموذج اللاخطي في ظل وجود المشاهدات الشاردة ولحجوم عينات مختلفة.