

تقنية ROBUST – LOESS في تحليل الانحدار

زينب حسن راضي

جامعة القادسية

علوم الحاسوب وتكنولوجيا المعلومات

الخلاصة :

تناولنا في هذا البحث دراسة استخدام تقنية robust-loess في تقدير دالة الانحدار اللامعلمي ، حيث أن هذه الطريقة تضيف الحصانة الى طريقة loess كونها طريقة غير حصينة بسبب اعتمادها استخدام المربعات الصغرى في التمهيد والتي تتأثر بوجود الشواذ ويتم تنفيذ robust-loess باستخدام تقنية Last Absolute Residuals وتقنية bisquare التي تضيف المتانة الى المربعات الصغرى الموزونة في loess .

The Technique of Robust – Loess in Regression Analysis

ZAINB HASSAN RATIY

Al Qadisiya University

Faculty of Computer Science and Information Technology

Abstract:

In this paper the robust-loess method is used to estimate the nonparametric regression function. The Loess is non-robust method and used in case of outliers where it bases on the less squares in regression which affects by presence of outliers. In this paper, the robust-loess has been implemented through applying the last absolute residuals and bi-square techniques which enhanced robustness to the weighted least squares in loess.

1 المقدمة

تحليل الانحدار هو اداة احصائية تقوم ببناء نموذج احصائي لتخمين العلاقة بين المتغيرات وقد استخدمنا في هذا الاسلوب تمهيد مخطط التشتت لانه يساعدنا في رؤية تلك العلاقة فهو يلخص متوسط البيانات باستخدام دالة تمهيد للنقاط اضافة الى انه يوفر لنا تقدير او تنبؤ لكل قيمة معطاة من x يطلق على هذه الطريقة اسم Loess وهي احدى ادوات تحليل الانحدار اللامعلمي وقد اضعنا لها تقنية robust التي استخدمت من قبل باحثين منهم [1] (2004) خلود يوسف خمو و [2] (2015) زينب حسن راضي و [3] William(1979) و [4] William and Devlin(1988) ، صيغة نموذج الانحدار اللامعلمي [1][2] :

$$y_i = g(x_i) + \varepsilon_i \quad i = 1, 2, \dots, n$$

حيث ε_i يمثل الخطأ العشوائي و $g(x_i)$ هي دالة الانحدار المجهولة والتي نريد تقديرها أو تمهيدها و y_i المتغير المعتمد

2 هدف البحث :

يهدف البحث الى اضافة استخدام تقنية Robust الى طريقة Loess لتمهيد الدالة اللامعلمية لاضافة المتانة والدقة في التقدير للوصول الى التقارب من المنحني الحقيقي.

3 البواقي Residuals :

من طبيعة تحليل الانحدار (سواء كان معلمي او لامعلمي) هو وجود البواقي التي تزيد من معرفتنا بمدى تشتت الخطأ العشوائي حول خط الانحدار (والتي هي الفروق بين القيمة التي نحسبها من نموذج الانحدار والقيمة الحقيقية) وتعرف بالصيغة التالية:

$$\varepsilon_i^\wedge = y_i - y_i^\wedge \quad i = 1, 2, \dots, n$$

حيث ان y_i قيم المتغير المعتمد نسبة الى المتغير المستقل x_i و y_i^\wedge قيمة متغير التنبؤ. ويتم العمل على عرض هذه

الخطوات المبينة ادناه توضح وصف خطوات Loess في عملية التمهيد:

(a) نحسب الانحدار الموزون لكل نقطة من نقاط البيانات ضمن المجال ،حيث الوزن يعطى باستخدام دالة الوزن الثلاثية :

$$w_i = \left(1 - \left| \frac{x - x_i}{d(x)} \right|^3\right)$$

حيث x تمثل قيمة التنبؤ المرتبطة بالتمهيد و x_i قيم نقاط الجوار الى x المعرفة في الفترة و $d(x)$ طول المسافة على المحور السيني لمعظم قيم المؤشر في الفترة . ويمتاز الوزن بالخصائص التالية :

- نقاط البيانات الممهدة تمتلك وزن اكبر ولها التأثير الاكبر على التمهيد
- نقاط البيانات الشاذة في المجال لها وزن صفر وليس لها تأثير على التمهيد.

(b) نستخدم انحدار المربعات الصغرى الموزونة ، اذا كان انحدار lowess تستخدم متعددة حدود من الدرجة الاولى ،اذا كان انحدار loess تستخدم متعددة حدود تربيعية

الفروق او الاختلافات باستخدام احدى الطرق مثل المدرج التكراري او box plot او استخدام مخطط التشتت البواقي.

4) تمهيد Loess

هي طريقة لامعلمية (بمعنى ليس لها تحديدات اولية لوصف شكل العلاقة بين المتغيرات) وهي اختصار لمصطلح الانحدار الموضوعي الموزون (Cleveland(1979)^[3] والمصطلحان Lowess و Loess مشتقان من مصطلح locally weight scatter plot smoothing وهما الطريقتان مختلفتان بالموديل المستخدم للانحدار، حيث تستخدم دالة Lowess دالة أنحدار متعددة حدود خطية اما Losee فتستخدم دالة متعددة حدود تربيعية. ان عملية التمهيد هنا تعتبر موضعية لان كل قيمة ممهدة تحدد بواسطة جوار الاقرب للبيانات الواقعة ضمن الفترة و موزونة لان دالة انحدار الوزن تعرف لكل نقاط البيانات الواقعة ضمن الفترة. وبأستطاعتنا ان نستخدم دالة الوزن الحصينة لجعلها مقاومة للقيم الشاذة او المتطرفة كما في طريقة Robust-Loess

• يتم حساب البواقي من الصيغة التالية :

$$\varepsilon_i = y_i - \hat{y}_i$$

• حساب الاوزان الحصينة لكل نقطة من

نقاط الفترة ، والاوزان تعطى بواسطة

دالة bisquar :

$$w_i = \begin{cases} 1 - \left(\frac{\varepsilon_i}{MAD}\right)^2 & |\varepsilon_i| \leq MAD \\ 0 & |\varepsilon_i| > MAD \end{cases}$$

حيث ان ε_i تمثل البواقي أو قيمة الخطأ العشوائي لكل نقاط الفترة، MAD هو المتوسط المطلق لانحرافات البواقي الذي هو مقياس لكمية انتشار البواقي

$$MAD = median(|\varepsilon|)$$

ان قيم الاوزان الحصينة تتغير مع تغير قيمة الخطأ العشوائي مقارنة مع قيمة المتوسط المطلق ، فاذا كانت قيمة الخطأ العشوائي (ε_i) صغيرة مقارنة مع MAD فإن قيم الاوزان الحصينة تكون قريبة الى الواحد ، اما اذا كانت قيمة الخطأ العشوائي (ε_i) كبيرة فإن الاوزان الحصينة MAD تساوي صفر ونقاط

(c) قيمة التمهيد تعطى بواسطة انحدار الموزون لقيمة التنبؤ ، اما اذا كان حساب التمهيد يتضمن نفس العدد للجوار الاقرب لنقاط البيانات على جانبي نقاط البيانات الممهدة فإن دالة الوزن متماثلة و اذا كان عدد نقاط بيانات جوارالتغير غير متماثل حول نقاط البيانات الممهدة فإن دالة الوزن هي دالة غير متماثلة.

1-4 تمهيد Robust Loess:

أن تقنية Losee هي تقنية غير حصينة لانها تعتمد الية المربعات الصغرى التي تتاثر بوجود القيم المتطرفة والتي تؤثر على قيمة التمهيد وتجعلها لاتعكس سلوك الجزء الاكبر من قيم البيانات ، لذا نلجأ الى استخدام الاجراء الحصين robust لمعالجة هذا التطرف في البيانات الذي يضيف متانة الى اوزان المربعات الصغرى في خطوات Loess اضافة الى استخدامها تقنية bisquare والتي تعيد اوزان النقاط الاساسية الى البواقي ، فاذا كانت البواقي كبيرة (انحراف كبير في النموذج) فان اوزان هذه النقاط تكون منخفض وبالعكس. ويمكن وصف هذه الطريقة بعدة خطوات :

البيانات المرتبطة تستبعد من حساب التمهيد .

• نعيد تمهيد البيانات بأستخدام الاوزان الحصينة حيث ان قيمة التمهيد النهائية تستخدم كلا من الانحدار الموضوعي الموزون والاوزان الحصينة

• نكرر الخطوات السابقة للحصول على افضل النتائج .

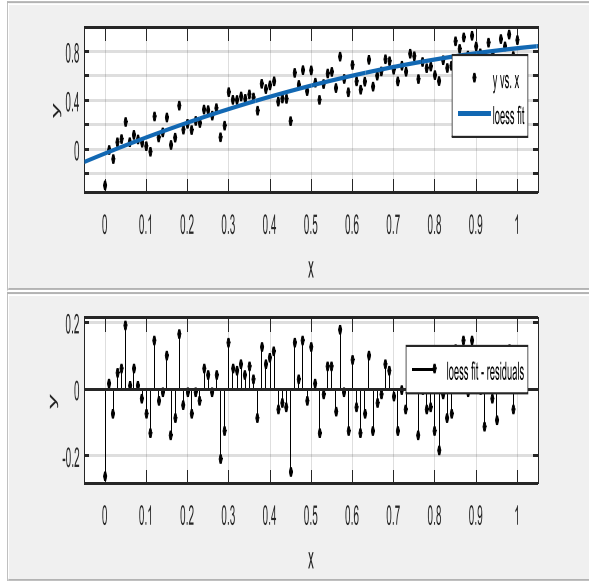
• تم تطبيق robust بأستخدام تقنيتين هما (bisquare) و Last

Absolut Residual ، وسوف نرملها بالرمز (LAR) حيث أن (LAR) هي طريقة لايجاد منحنى يقلل الفروق المطلقة للبواقي بدلا من تقليل مربعات الاختلافات ولذلك النقاط التي تمتلك قيم عالية يكون لها تأثير أقل على التمهيد أما bisquare فهي تقلل مجموع مربعات الاوزان حيث أن الوزن يعطى لكل نقطة من نقاط البيانات حسب بعدها أو قربها من المنحنى وبالتالي فالنقطة التي تمتلك وزن اكبر تأخذ القيمة صفر وبالعكس وبذلك فإن هذه التقنية تسعى في نفس الوقت الى أيجاد منحنى مناسب لمعظم نقاط البيانات إضافة الى تقليل تأثير القيم المتطرفة أن وجدت.

(5) التطبيق :

ولتطبيق ما تم ذكره اعلاه استخدمنا برنامج MATLAB 2015 ، حيث قمنا بتوليد المتغير العشوائي x_i و الاخطاء العشوائية ε_i بأستخدام الاوامر المتاحة في البرنامج ، إضافة الى توليد المتغير المعتمد من خلال جمع دوال المتغير التوضيحي مع متجه الاخطاء العشوائية ، واعتمدنا دالة اختبار متعددة الحدود من درجات مختلفة لتوضيح تمهيد البيانات .

في هذا البحث اعتمدنا معيار مجموع مربعات البواقي الحاصل بسبب التمهيد SSE (التي تعرف بأنها الاختلافات بين قيمة المشاهدة وقيمة التنبؤ) واحصائية معامل الارتباط R^2 في الطريقتين أعلاه ، ان استخدام SSE اداة تشخيص مفيدة من اجل تحديد ماذا كان المنحنى ممهد بشكل يشمل اغلبية البيانات حيث ان القيمة الصغيرة له تشير الى ان النموذج يمتلك خطأ عشوائي صغير وبالعكس، اما R-Square فهي مقياس لمدى نجاحنا في وصف او بيان العلاقة بين المتغيرات او هي مربع الارتباطات بين متغير الاستجابة ومتغير



شكل رقم (1) : معادلة الانحدار
لطريقة تمهيد loess

اما لطريقة robust-loess by
bisquar فقد حصلنا النتائج
التالية:

$$\beta_0 = 0.9887$$

$$\beta_1 = 1.138$$

$$\beta_2 = - 0.7756$$

Goodness of fit:

$$SSE = 0.9687$$

$$R\text{- Square} = 0.5908$$

$$\text{Adjusted R-Square} = 0.5824$$

$$\text{RMSE} = 0.09942$$

والشكل رقم (2) يوضح معادلة الانحدار
للنتائج اعلاه.

التنبؤ وهذه الاحصائية تاخذ قيمة بين 0 و 1

$$SSE = \sum (y_i - \hat{y}_i)^2$$

$$R^2 = \frac{\sum (y_i - \bar{y})^2}{\sum (y_i - \hat{y}_i)^2}$$

ووفقا لمعادلة الانحدار التالية والتي يتوزع
فيها المتغير x توزيعا طبيعيا بمتوسط
0.5 وتباين 0.293 :

$$g(x_i) = \beta_0 + \beta_1 x + \beta_2 x^2$$

فقد كانت النتائج لطريقة Loess كالتالي

$$\beta_0 = -0.03475$$

$$\beta_1 = 1.36$$

$$\beta_2 = - 0.4992$$

Goodness of fit:

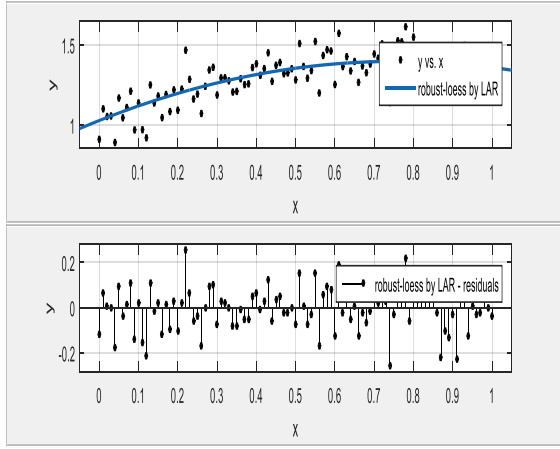
$$SSE = 0.9101$$

$$R\text{- Square} = 0.8772$$

$$\text{Adjusted R-Square} = 0.8747$$

$$\text{RMSE} = 0.09637$$

والشكل رقم (1) يوضح معادلة الانحدار
للنتائج اعلاه.



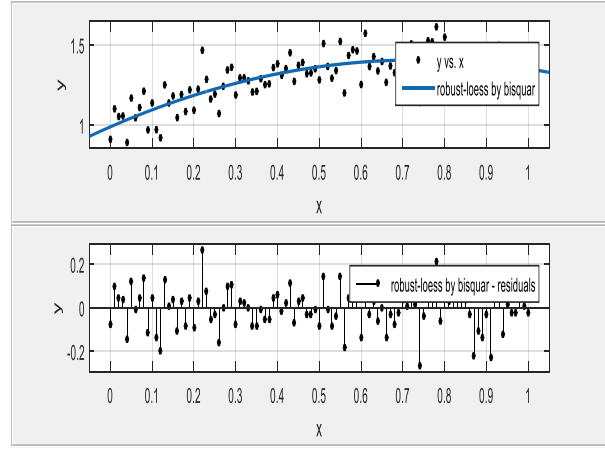
شكل رقم (3) : معادلة الانحدار لطريقة
robust-loess by (LAR)

(6) الاستنتاج :

1- أن استخدام loess يمثل نهجا مرنا للغاية للغاية لمشكلة البيانات أضافة الى أننا لانحتاج الى تحديد اولي لشكل العلاقة بين المتغيرات المستقلة وامعتمدة كذلك سهولة ووضوح تركيبها .

2- تعتمد طريقة loess استخدام المربعات الصغرى في التمهيد لذا نستطيع اعتبارها طريقة عامة لكل من الاسويين المعلمي واللامعلمي ولكن بنفس الوقت يواخذ عليها تأثيرها بالقيم المتطرفة لان المربعات الصغرى تتأثر بوجود الشواذ.

3- يمثل أضافة استخدام robust حلا لمشكلة القيم المتطرفة.



شكل رقم (2) : معادلة الانحدار لطريقة
robust-loess by bisquare

بينما اظهرت طريقة robust-

loess by (LAR) النتائج التالية:

$$\beta_0 = 1.028 , \beta_1 = 0.9745$$

$$\beta_2 = - 0.6433$$

Goodness of fit:

$$SSE = 0.9504$$

$$R\text{- Square} = 0.5985$$

$$\text{Adjusted R-Square} = 0.5903$$

$$RMSE=0.09848$$

والشكل رقم (3) يوضح معادلة الانحدار
للنتائج اعلاه.

regression and smoothing”
Journal of American
"statistical association
(1979) issue368 ، vol.74
-4 William S. Cleveland and
Locally "sudan J. Devlin
weighted regression:
"Approach analysis by local
Journal of American
vol.83 "statistical association
'pp596-610 ، NO.403 ،
(1988).
-5 "Learn Matlab" Matlab
. version 6 release 12

(7) التوصيات :
نوصي باستخدام تقنية robust-loess في
التمهيد لما تضيفه من متانة الى طريقة
المربعات الصغرى حيث انها تعيد اوزان
النقاط الى البواقي وبأستطاعتنا اضافة
robust الى الشرائح spline او عندما
يكون توزيع الخطأ العشوائي غير طبيعي
(ملوث).

المصادر :

1- خلود يوسف خمو "مقارنة اساليب بيز
مع طرائق اخرى لتقدير منحنى
الانحدار" أطروحة دكتوراه في الاحصاء
،جامعة بغداد ، كلية الادارة واقتصاد
(2004).

2- زينب حسن راضي "تمهيد دالة
الانحدار اللامعلمي بطرائق تمهيد
متنوعة" رسالة ماجستير في الاحصاء
الرياضي ،جامعة القادسية ،كلية علوم
الحاسوب والرياضيات (2015).

-3 William S.Cleveland
" Robust locally weighted "