

مقارنة طريقتي تقدير الدالة اللامعلمية لبيانات عنقودية عن كريات الدم البيضاء لمرضى اللوكيميا

أ.م.د. سجي محمد حسين / كلية الادارة والاقتصاد / جامعة بغداد
م.م. حلا كاظم عبيد / كلية الادارة والاقتصاد / الجامعة العراقية

تاريخ التقديم : 2016/8/16
تاريخ القبول : 2016/11/9

المستخلص:

البيانات العنقودية تظهر في الكثير من العلوم الاجتماعية والصحية والسلوكية. ويتميز هذا النوع من البيانات بوجود الارتباط بين مشاهداتها. ويمكن التعبير عن العنقدة من حيث العلاقة بين القياسات على الوحدات ضمن المجموعة نفسها. تم في هذا البحث تقدير الدالة اللامعلمية للبيانات العنقودية باستعمال طريقة المقدرات اللبية غير المرتبطة ظاهريا *The Seemingly Unrelated Kernel Estimators* , وطريقة المربعات الصغرى المعممة لمقدرات الشريحة التمهيدية *The Generalized Least Squares Smoothing Spline Estimators* وتم تطبيق الطريقتين المذكورتين على بيانات مرضى اللوكيميا وتمت المقارنة بين الطريقتين عن طريق معيار المقارنة *MAE* و *MSE* ووضحت النتائج التطبيقية كفاءة مقدر المربعات الصغرى المعممة لمقدرات الشريحة التمهيدية.

المصطلحات الرئيسية للبحث/ البيانات العنقودية ، الطريقة المقدرات اللبية غير المرتبطة ظاهريا،
وطريقة المربعات الصغرى المعممة لمقدرات الشريحة التمهيدية، *MAE*، *MSE*



مجلة العلوم
الاقتصادية والإدارية
العدد 97 المجلد 23
الصفحات 394-419

*البحث مستل من اطروحة دكتوراه



مقارنة طريقتي تقدير الدالة اللامعلمية لبيانات عنقودية عن كريات الدم البيضاء لمرضى اللوكيميا

المقدمة:

يطلق على مجموعة من القياسات المشتركة من المعلومات او الاشياء ضمن عنقود واحد (cluster) التي تمت هيكلتها في عناقيد بالبيانات العنقودية. حيث ان مجموعة من الاشياء المشتركة تسمى عنقود والبيانات ضمن العنقود الواحد تكون متماثلة (مرتبطة) اكثر مع بعضها البعض من تلك الموجودة في العناقيد الاخرى^[3]. استعملت البيانات العنقودية في الكثير من مجالات الطب الحيوي وعلم الاوبئة والظواهر الزراعية والاقتصادية وغيرها ، ومما تجدر الاشارة اليه ان البيانات العنقودية تشبه البيانات الطولية من حيث اعتمادية البيانات داخل العنقود الواحد (القطاع بالنسبة للبيانات الطولية) على بعضها بعضا^[5]. ومن هنا تاتي اهمية الارتباط الموجود داخل المشاهدات لتقدير الدوال حيث ان استخدام الاسلوب المعلمي لا يكون دائما مرغوب فيه لتقدير بعض الدوال بسبب ان شكل الدالة الملائم للبيانات غير معروف مسبقا او نتيجة لوجود بعض المعوقات لذلك يتم استخدام الاسلوب اللامعلمي لتقدير (تمهيد) الدالة اللامعلمية. لقد تطورت البحوث في الاونة الاخيرة حول استخدام الانحدار اللامعلمي عندما تكون الافتراضات المعلمية غير متحققة، والانحدار اللامعلمي يسمح بمرونة اكبر للدوال المعتمدة على المتغيرات الناتجة من البيانات. لقد تطرقت البحوث السابقة الى حالة تقدير البيانات العنقودية بالطرائق اللامعلمية والطرائق الشبه معلمية واعتمدت حالة اهمال الارتباط داخل بيانات العنقود الواحد الخاصية التي تميز البيانات العنقودية بشكل خاص . حيث ان المقدرات اللبية الموضوعية حققت كفاءة اكثر باهمال الارتباط داخل العناقيد (حتى وان كان الارتباط في مصلحة الدراسة). حيث ان مصفوفة الارتباط داخل العنقود تلعب دورا مهما في كفاءة التقدير للدوال غير المعلومة $\theta(x)$ في نموذج الانحدار اللامعلمي

$$y = \theta(x) + \varepsilon \quad (1)$$

وتتشا البيانات العنقودية في كثير من المجالات ومنها المجالات الطبية ويعد سرطان الدم احدى المجالات التي يمكن ان تعد البيانات الماخوذة من المرضى المصابين به احدى انواع البيانات العنقودية. حيث اعتبر كل شخص (مريض) كعنقود وتؤخذ له قياسات دورية لنسبة كرات الدم البيضاء ونسبة الهيموغلوبين وتعد القياسات المتكررة في حالات (اوقات) مختلفة وحدات المستوى الاول وتتمثل وحدات المستوى الثاني بالاشخاص.

يتميز سرطان الدم بوجود عدد كبير من الخلايا البيضاء الناضجة نسبيا ولكنها تمتلك شكل غير طبيعي فضلا عن ان خلايا الدم البيضاء في العادة تحتاج الى شهور لتنضج ولكن في هذه الحالة خلايا الدم البيضاء تنتج بسرعة عالية مما يؤدي الى وجود عدد من الخلايا البيضاء غير الطبيعية غير المسيطر عليها التي من غير الممكن التكهن بسلوكها. لهذا النوع من البيانات يمكن تقديرها باستعمال النموذج (1).

هدف البحث:

يهدف البحث الى تقدير الدالة اللامعلمية للبيانات العنقودية باستخدام طريقة المقدرات اللبية غير المرتبطة ظاهريا (seemingly unrelated kernel regression) وطريقة المربعات الصغرى المعممة لمقدر الشريحة التمهيدية (generalized least squares smoothing spline estimator) مع الاخذ بنظر العناية الارتباط الموجود داخل العنقود الواحد حيث ان مصفوفة الارتباط داخل العنقود الواحد تلعب دورا مهما في تحليل البيانات العنقودية مع اخذ حالة الارتباط الموجود داخل العناقيد ومقارنتها مع حالة الاستقلالية لبيانات كريات الدم البيضاء لمرضى اللوكيميا من خلال معيار المقارنة MAE و MSE .



مقارنة طريقتي تقدير الدالة اللامعلمية لبيانات عنقودية عن كريات الدم البيضاء لمرضى اللوكيميا

الاستعراض المرجعي:

في عام 2000 درس كل من LIN, X. & CARROLL, R. J. طريقة الانحدار اللبي لمتعدد الحدود الموضوعي مع متغير واحد لبيانات عنقودية باستخدام المعادلات التقديرية على افتراض ان كل عنقود يشمل عدة مشاهدات ولتكن $m < \infty$ حيث يتم قياس المؤشر في حالة غياب خطأ القياس للانحدار العشوائي بتجاهل هيكلية الارتباط داخل العنقود واثبت ان المقدر اللبي الموضوعي لمتعدد الحدود يكون اكثر كفاءة في حالة تجاهل الارتباط الموجود داخل العنقود الواحد واستخدموا الطريقة المقترحة في تحليل بيانات نقص المناعة البشرية لخلايا ال CD₄ . واكد ان استعمال الطريقة اللبية المحددة لا تسفر بالضرورة عن مقدر جيد عند حساب الارتباط داخل العناقيد.^[12]

وفي عام 2001 استخدم كل من الباحثين LIN, X. and CARROLL R. J. النماذج شبه المعلمية الخطية العامة للبيانات العنقودية باستخدام المعادلات التقديرية المعممة GEE وطبقت النتائج في حالة كون عدد المشاهدات لكل عنقود محدودة وعدد العناقيد كبير ومتوسط النتيجة μ يمثل بالشكل

حيث $g(\mu) = X^T \beta + \theta(T)$ دالة ربط و X و T متغيرات و β متجه المعالم غير المعلومة و $\theta(T)$ دالة ممهدة غير معلومة^[11]

وفي عام 2002 قام كل من الباحثين Welsh, A. H., Lin, X. and Carroll, R. J. بدراسة مقدر الشريحة الذي يمتلك وزنا اكثر وتباين اقل من الطريقة اللبية والتباين المحاذي يكون اقل عند حساب الارتباط بشكل صحيح. حيث ان الطريقة اللبية تسلك سلوك مختلف عن طريقة الشرائح التمهيدية في حالة اخذ هيكلية الارتباط الموجود داخل العناقيد في حين ان كلا الطريقتين الشرائحية واللبية تكونان متكافئتين في حالة اهمال الارتباط داخل العناقيد.^[17]

وفي عام 2005 قام الباحث مناف يوسف حمود باستعراض بعض الطرائق اللامعلمية وشبه المعلمية مع بعض الطرائق المقترحة لتقدير دالة الكثافة الاحتمالية كذلك استعراض الطرائق الالهة والخاصة بتقدير المعلمة التمهيدية ومقارنة تلك الطرائق من خلال استخدام اسلوب المحاكاة باستعمال توزيعات وحجوم عينات ومستويات تباين مختلفة للمتغير X .^[1]

وفي عام 2007 قام الباحث عمر عبد المحسن بدراسة امكانية تقدير النموذج التجميعي (Generalized Additive Model) لتحليل النماذج الشبه معلمية فضلا عن التحليل اللامعلمي واقترح تقدير GAM الحصينة عبر دالة وزن معينة للتحليل المعلمي وشبه المعلمي في حالة وجود تلويف في البيانات ووضح استعمال الشرائح التمهيدية بشكل جزئي (partially) (شبه معلمي) افضل من استعمالها بشكل تام (completely) (لامعلمي) عند تحليل GAM.^[2]

وفي عام 2008 قام كل من Wang, Y. G. & Zhao, Y. بتحليل البيانات العنقودية لنماذج الانحدار على اساس الرتبة واقترحوا طريقة رتبة wilxon الموزونة لحساب الارتباط داخل العناقيد ولأحجام مختلفة من العناقيد. واجراء دراسة المحاكاة لمقارنة المقدرات المختلفة لعدد من الطرائق حول هيكلية الارتباط بوجود او عدم وجود القيم المتطرفة وان الطرائق المقترحة تظهر أداء جيد في حالة حساب الارتباط الموزون داخل العناقيد فضلا عن الطرائق الحصينة المتضمنة عدم الحساسية للقيم الشاذة وهيكلية الارتباط.^[16]

وفي عام 2010 قام كل من Ibrahim, N. A. & Suliadi بتحليل البيانات العنقودية / الطولية باستخدام المعادلات التقديرية GEE مع مهادت الشريحة الطبيعية التكميلية باستخدام المحاكاة مع اخذ حالة الارتباط الموجود داخل العنقود بعين العناية حيث اظهر ان GEE-Smoothing spline لها خصائص افضل من GEE-Local Polynomial Kernel الذي تناولها Lin & Carroll 2000 الذي اهمل الارتباط الموجود داخل العناقيد للحصول على مقدرات كفوءة في حين برهن Ibrahim, N. A. & Suliadi كفاءة المقدرات الناتجة من حيث التحيز والكفاءة والاتساق في حالة الارتباط الحقيقي للبيانات داخل العناقيد.^[13]

وفي عام 2013 قام كل من Ma, S., Song, Q. & Wang, L. بتطوير منهجية عامة لاختيار المتغيرات انيا وتقدير العناصر المجهولة للنماذج التجميعية الخطية جزئيا (additive partially linear models) باختصار (APLMS) للبيانات العنقودية



مقارنة طريقتي تقدير الدالة اللامعلمية لبيانات عنقودية عن كريات الدم البيضاء لمرضى اللوكيميا

واقترحوا في الخطوة الاولى المربعات الصغرى للحصول على التقدير المعلمي اما المركبات اللامعلمية فيتم تقديرها على اساس متعدد الحدود للشرائح التمهيدية polynomial spline smoothing ويعد هذا الاسلوب مرن وسهل التطبيق على ارض الواقع وركزوا على اخذ الجزء الخطي في النموذج ولتحسين الكفاءة يجب على المرء ان يختار هيكليّة الارتباط المناسبة للبيانات المدروسة.^[7]

الجانب النظري:

تم في هذا البحث عرض بعض طرائق تقدير الدالة اللامعلمية للبيانات عنقودية وكما يأتي:

(1) طريقة مقدر النواة غير المتصل ظاهريا: The Seemingly Unrelated Kernel Estimator

نفرض n من العناقد ب m_i من المشاهدات داخل كل عنقود حيث $i = 1, 2, \dots, n$ وان j^{th} من المشاهدات $j = 1, 2, \dots, m_i$ في كل عنقود مؤلف من متغيرات الاستجابة Y_{ij} ومتغير مفرد X_{ij} التي تمثل القياسات المتكررة داخل كل عنقود عبر اوقات مختلفة . بالنسبة للبيانات الطولية حيث ان كل عنقود يمثل قطاع، اما بالنسبة للبيانات العائلية فان كل عائلة تمثل عنقود والمشاهدات داخل العنقود تمثل مختلف افراد الاسر من العائلة نفسها كما موضح في النموذج في المعادلة رقم 1 ، فان لأي مصفوفة تباين وتباين مشترك V فان معدل مقدر النواة غير المتصل ظاهريا ويرمز لها بالرمز $\hat{\theta}_k(x)$ ويساوي لأي نقطة وعندما تكون رتبة مقدر النواة غير المتصل ظاهريا $p = 0$ فان

$$\hat{\theta}_k(x) = k'_{wh}(x) \{I + (\tilde{V}^{-1} - \tilde{V}^d) k_w\}^{-1} \tilde{V}^{-1} Y \quad (2)$$

حيث ان: k_{wh} : متجه من الرتبة $nm \times 1$ ويساوي:

$$k_{wh}(x) = \left\{ \sum_{i=1}^n \sum_{j=1}^{m_i} k_h(X_{ij} - x) V^{jj} \right\}^{-1} \{k_h(X_{11} - x), \dots, k_h(X_{nm} - x)\}'$$

وان $Y = [\underline{y}_1, \underline{y}_2, \dots, \underline{y}_n]'$ مصفوفة من الدرجة $nm \times nm$

$$y_i = (y_{1i} \ y_{2i} \ \dots \ y_{im})' \quad \text{و } nm \times 1 \text{ متجه من الرتبة}$$

وان k_w عبارة عن مصفوفة $nm \times nm$ $k_w = \{k_{wh}(X_{11}), \dots, k_{wh}(X_{nm})\}'$

$$\tilde{V}^d = \text{diag}(V^d, V^d, \dots, V^d); V^d = \text{diag}(V^{-1}) = \text{diag}(V^{jj}) \quad \text{وان}$$

$$\tilde{V} = \text{diag}(V_1, \dots, V_n)$$

ولمتجه $\hat{\theta}(X)$ من الرتبة $nm \times 1$ يحتوي كل المقدرات $\hat{\theta}_k(x)$ لكل نقاط التصميم X

$$\hat{\theta}_k(X) = \{\hat{\theta}_k(X_{11}) \dots \hat{\theta}_k(X_{nm})\}' \quad \text{حيث ان:}$$

فان:

$$\hat{\theta}_k(X) = \{I + K_w(\tilde{V}^{-1} - \tilde{V}^d)\}^{-1} k_w \tilde{V}^{-1} Y \quad (3)$$



وحيث ان

$$k_{wh}(x)' = \delta_1' \{ \tilde{X}(x)' k_{dh}(x) \tilde{V}^d \tilde{X}(x) \}^{-1} \tilde{X}(x)' k_{dh}(x)$$

$$\delta_1 = (1, 0, \dots, 0)'$$

وان $\tilde{X}(x)$ مصفوفة من الرتبة $nm \times p + 1$ ب $\{(n-1)^{i+j}\}$ من الصفوف

$$\tilde{X}(x) = \begin{bmatrix} 1 & (X_{11} - x) & (X_{12} - x)^2 & \dots & (X_{1m} - x)^p \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & (X_{n1} - x) & (X_{n2} - x)^2 & \dots & (X_{nm} - x)^p \end{bmatrix}$$

$$k_{dh}(t) = \text{diag} \{ (X_{11} - x) / h, \dots, (X_{nm} - x) / h \}$$

$$\text{حيث: } k_h(s) = h^{-1} k_h(s/h)$$

$k_h(s)$: تمثل دالة الكيرنل وباستخدام دالة ال Gaussian فان:

$$k(.) = (2\pi)^{-0.5} \text{Exp} \left(-\frac{u^2}{2} \right) \quad (4)$$

و h يمثل عرض الحزمة (Bandwidth) التي سيتم حسابها لاحقا
ومصفوفة ال V من الدرجة $nm \times nm$ تحسب كالآتي:

$$V = A^{1/2} R(\alpha) A^{1/2} \quad (5)$$

A : تمثل مصفوفة الانحراف المعياري من الدرجة $nm \times nm$

$$A = \begin{bmatrix} a_{i1}^{1/2} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & a_{im}^{1/2} \end{bmatrix}$$

$$\text{حيث ان: } a_{i1} = \sqrt{\text{var } y_{i1}}$$

وان $R(\alpha)$: تمثل مصفوفة الارتباطات من الدرجة $nm \times nm$ [11],[7]

$$R(\alpha) = \begin{bmatrix} 1 & \rho_{i1} & \dots & \rho_{imi} \\ \vdots & \ddots & \ddots & \vdots \\ \rho_{imi} & \dots & \dots & 1 \end{bmatrix}$$



مقارنة طريقتي تقدير الدالة اللامعلمية لبيانات متقودية عن
كريات الدم البيضاء لعرضى اللوكيميا

$$R(\alpha) = \text{corr}(y_{ij}, y_{ij+1}) \quad i = 1, \dots, n-1 \quad \text{و}$$

$$\hat{\rho}_{ij} = \frac{y_{ij} - m(y_{ij})}{\sqrt{v(y_{ij})}}$$

(2) طريقة المربعات الصغرى المعممة لمقدر الشريحة التمهيدية: The Generalized Least

Squares Smoothing Spline Estimator

يراد من دراسة الشرائح التمهيدية مطابقة البيانات بشكل جيد من جهة ويكون الممهد على قدر من التمهيد من جهة اخرى ولتقدير المعلمة الممهدة في النموذج (1) باستخدام ممهدات الشريحة من الدرجة p ولأي مصفوفة تباين عاملة V والذي يقلل المعادلة

$$\frac{1}{n} \sum \{Y_i - \theta(X_i)\}' V^{-1} \{Y_i - \theta(X_i)\} + \lambda \int \{\theta^{(p)}(X)\}^2 dX \quad (6)$$

حيث ان $\theta^{(p)}(X)$ قابلة للاشتقاق P من المرات لدالة الانحدار، وان الحد الثاني من المعادلة (6) يمثل مقدار عدم التمهيد فان مقدر الشريحة الممهد بالصيغة الاتية:

$$\hat{\theta}_s(X) = (\tilde{V}^{-1} + n\lambda\psi)^{-1} \tilde{V}^{-1} Y \quad (7)$$

λ : معلمة تمهيد التي تتغير بتغير الدوال وليس بتغير العقد وتكون $\lambda > 0$ [9] وان الـ

$$Y = [y_1, y_2, \dots, y_n]'$$

$$y_i = (y_{1i} \ y_{2i} \ \dots \ y_{im})' \quad \text{متجه من الرتبة } nm \times 1 \quad \text{و}$$

$$\text{ومصفوفة الـ } \tilde{X}(x) \text{ من الدرجة } nm + 1 \text{ [15],[8],[10]}$$

ولتقليل المعادلة (6) وباستخدام معلمة تمهيد ثابتة λ لجميع الدوال المستمرة والمختلفة تقود الى استخدام الشرائح التكميلية مع عقد عند نقاط البيانات (اي تستخدم البيانات كعقد عند نقاط التصميم) حيث ان قيمة الـ $p = 2$ في المعادلة (6) [4] ولإيجاد مصفوفة الـ ψ فان:

$$\psi = CQ^{-1}C' \quad (7)$$

حيث ان C مصفوفة من الدرجة $n \times n - 2$

$$i = 1, \dots, n; \quad j = 2, \dots, n-2$$

$$i = 1, \dots, n-1$$

$$C_{n \times n-2} = [C_{ij}]$$

$$L_i = x_{(i+1)} - x_i$$

$$C_{ij} = 0 \quad \text{if } |i - j| \geq 2$$

$$C_{j-1,j} = L_{j-1}^{-1}$$



مقارنة طريقتي تقدير الدالة الامعلمية لبيانات عنقودية عن
كريات الدم البيضاء لمرض اللوكيميا

$$C_{jj} = -L_{j-1}^{-1} - L_j^{-1}$$

ومصفوفة الـ Q من الدرجة $n - 2 \times n - 2$

$$Q = [q_{ij}] \quad i, j = 2, \dots, N - 1 \quad (8)$$

حيث ان :

$$q_{ij} = 0 \quad \text{if } |i - j| \geq 2$$
$$q_{i,i+1} = q_{i+1,i} = L_i / 6 \quad i = 2, \dots, N - 2$$
$$q_{ii} = (L_{i-1} + L_i) / 3 \quad i = 2, \dots, N - 1$$

$$\tilde{V} = \text{diag}(V, V, \dots, V) \quad ; V = D^{1/2} R(\alpha) D^{1/2} \quad (9)$$
$$D = \text{diag} [\sqrt{\text{var}(y)}]$$

ولنفرض ان نقاط التصميم x_1, x_2, \dots, x_m المرتبة تصاعديا تمثل جميع نقاط المشاهدات في العناقد $x_{ij}, i = 1, 2, \dots, n, j = 1, 2, \dots, m$ تستخدم كعقد حيث ان $a < x_1 < x_2 < \dots < x_m < b$ حيث ان a, b تحدد مدى فترة الزمن . تستعمل جميع نقاط التصميم ويجب ان تكون العقد اصغر من حجم العينة^[14]

معيار العبور الشرعي العام (General Cross Validation):

يتم استخدام معيار العبور الشرعي العام لاختيار معلمة التمهيد من خلال ابدال القيمة الموضحة في

المعادلة (6) والتي تساوي

$$\sum_{i=1}^n (I - A_{ii})^{-1} (I - A_{ii})^{-1}$$

بالقيمة

$\{1 - (nm)^{-1} \text{trace}(A)\}^2$ تكون صيغته كالآتي

$$GCV(\eta) = \frac{\sum_{i=1}^n \sum_{j=1}^m \{Y_{ij} - \hat{\theta}(X_{ij})\}^2}{\{1 - (nm)^{-1} \text{trace}(A)\}^2} \quad (11)$$

ويتم اختيار معلمة التمهيد الـ η التي تقابل اصغر قيمة لـ $GCV(\lambda)$ ^[17]

معيار المقارنة:

يستخدم معيار المقارنة للتعرف على افضل طريقة للتقدير ، من خلال المقارنة باقل متوسط مربعات الخطا او اقل متوسط القيمة المطلقة للأخطاء واللذان تكون صيغة كل منهما بالشكل الآتي:

$$MSE[\hat{\theta}(X)] = \text{mean}[\theta(X) - \hat{\theta}(X)]^2 \quad (12)$$

$$MAE[\hat{\theta}(X)] = \text{mean}|\theta(X) - \hat{\theta}(X)| \quad (13)$$



مقارنة طريقتي تقدير الدالة الالاعلمية لبيانات عنقودية عن كريات الدم البيضاء لمرض اللوكيميا

الجانب التطبيقي:

يعرف سرطان الدم (اللوكيميا) بأنه عبارة عن تكاثر غير محكوم لخلايا الدم وعادة ما يكون هذا التكاثر لخلايا الدم البيضاء، حيث ان عددها في الدم يزيد عن الحد الطبيعي وهو من 5 الاف الى 10 الاف خلية لكل ملم. ففي سرطان الدم يكون عدد الخلايا من (15-30) الف خلية لكل ملم . وقد يصل هذا التكاثر الى اكثر من 100 الف خلية لكل ملم او حتى يكون اقل من 5 الاف خلية لكل ملم (اي اقل من العدد الطبيعي). ويعود اختلاف هذا العدد الى مرحلة سرطان الدم ومدى تقدمه في الجسم ، فخلايا كرات الدم البيضاء غير طبيعية تنتج في نخاع العظم وهو مصنع ومصدر خلايا كرات الدم جميعها في الجسم فيؤدي تكاثرها في نخاع العظم الى خروجها من الدم . ففي تكاثرها في المراحل الاولى يؤدي ذلك الى تقليل عدد خلايا كرات الدم المصنعة وقبل ان تخرج خلايا كرات الدم البيضاء الى خارج الدم مما يؤدي الى وجود نقص مرحلي في كرات الدم البيضاء في الدم في بداية المرض، ولكن سرعان ما يصدر نخاع العظم هذه الخلايا السرطانية الى الدم ليزيد بذلك عدد كرات الدم البيضاء في الدم. وهذه الزيادة عادة ما تكون مضطربة. ومن اثار هذه الزيادة يصاب المريض بنقص الهيموغلوبين ومن ثم بفقر في الدم فيبدأ يعاني من تعب من اي مجهود يقوم به ودوران وكثرة حب النوم وكسل شديد. ومع استمرار المرض يعاني المريض من نقص في الصفائح الدموية فيصبح دمه يتأخر في التجلط وتكون عدة كدمات على جلده نتيجة عدم قدرة الجسم على وقف التجلطات^[6]

تم جمع بيانات عن مرض اللوكيميا لـ 100 شخص واخذت ثلاثة قياسات لكل شخص لثلاث اوقات مختلفة من مستشفى بغداد التعليمي /مدينة الطب، واخذت نسبة التغيرات لعدد كريات الدم البيض كمتغيرات توضيحية ومن ثم اخذت نسبة الهيموغلوبين في الدم كمتغيرات استجابة، ترتبط هذه القياسات فيما بينها مكونه مصفوفة الارتباط التي تلعب دورا مهما في تقدير دالة الانحدار الالاعلمية لنسبة كريات الدم البيضاء حيث تكون مصفوفة الارتباط في حالة الاستقلالية لثلاث اوقات لمتغير الاستجابة بالشكل الاتي:

$$R(\alpha) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

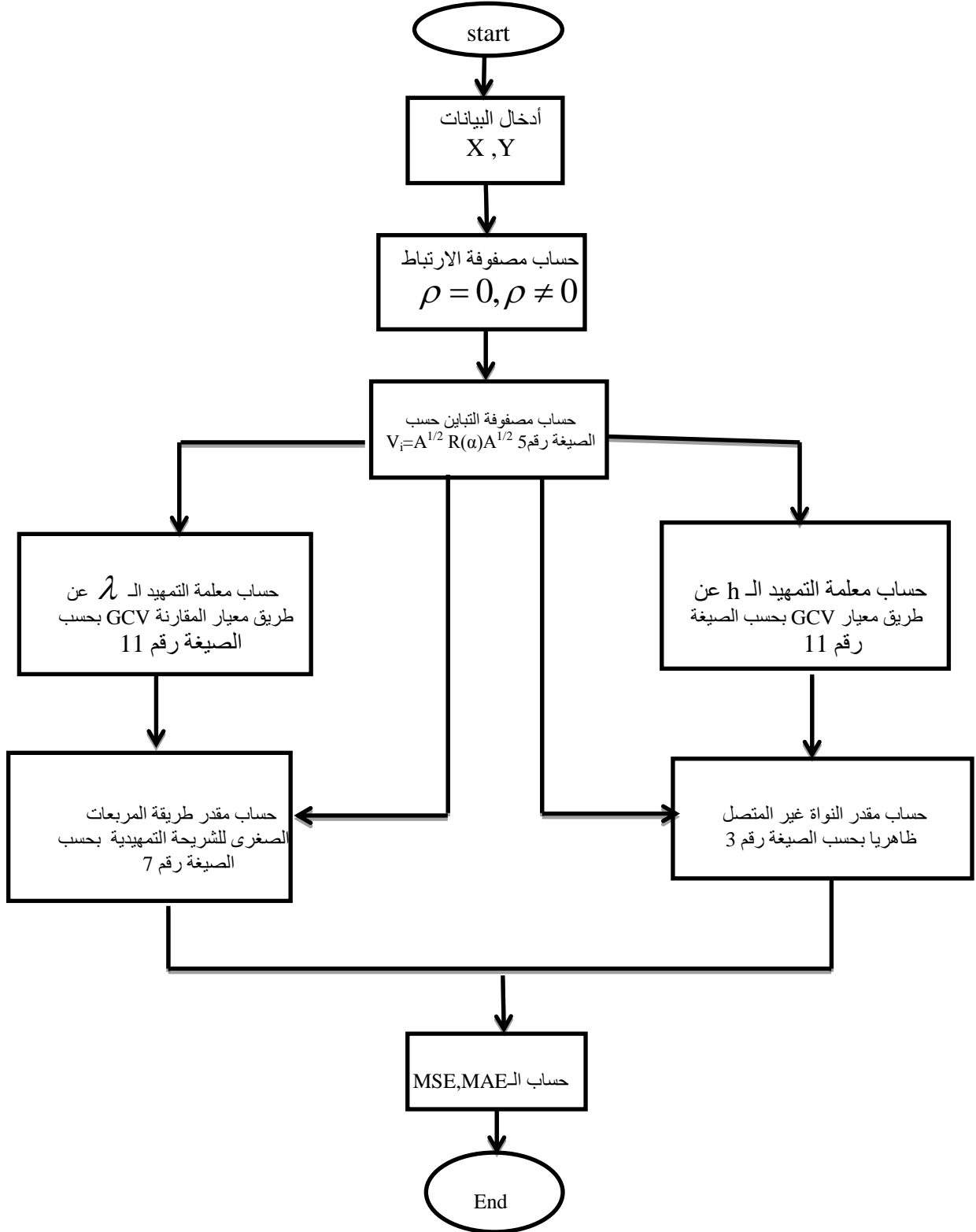
وفي حالة الارتباط الحقيقي بين العناقد فان:

$$R(\alpha) = \begin{bmatrix} 1 & 0.57 & 0.415 \\ 0.57 & 1 & 0.584 \\ 0.415 & 0.584 & 1 \end{bmatrix}$$



مقارنة طريقتي تقدير الدالة الامعلمية لبيانات عنقودية عن كريات الدم البيضاء لمرضى اللوكيميا

مخطط انسيابي يوضح العمل الرياضي





مقارنة طريقتي تقدير الدالة الالمعلمية لبيانات عنقودية عن كريات الدم البيضاء لمرضى اللوكيميا

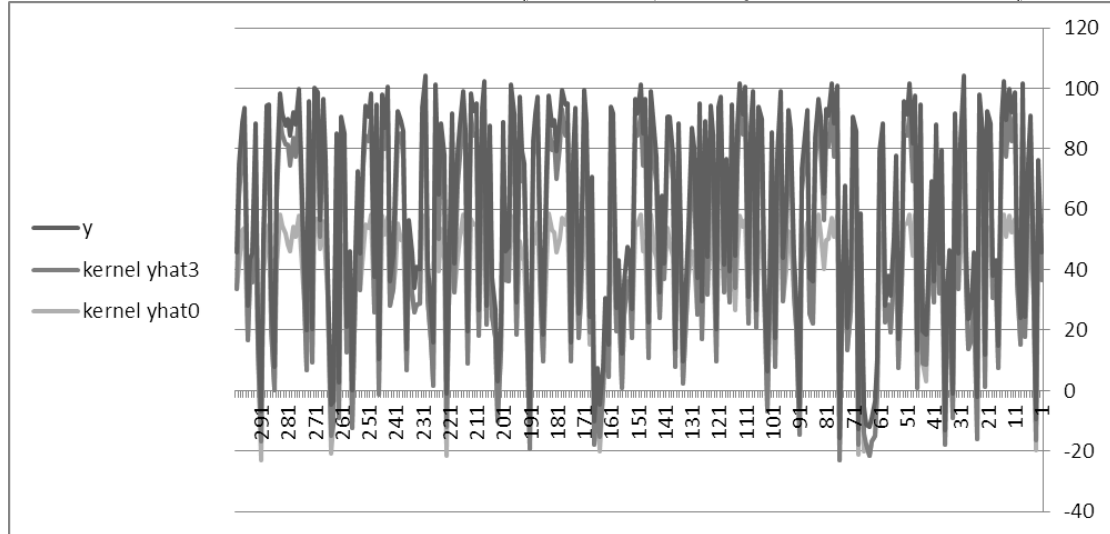
تحليل النتائج

تم تطبيق الطرائق المذكورة في الجانب النظري لبيان افضل طريقة لتقدير الدالة الالمعلمية التي تحقق اقل متوسط مربعات الخطأ و اقل قيمة لمتوسط الاخطاء المطلقة باستخدام برنامج بلغة الماتلاب (MATLAB) كما موضح في الملحق رقم (2) والجدول رقم (1) يمثل قيم MSE , MAE لطريقة المقدرات اللبية غير المرتبطة ظاهريا، وطريقة المربعات الصغرى المعممة لمقدرات الشريحة التمهيدية في حالتى الارتباط والاستقلالية ونلاحظ من الجدول ان اقل متوسط مربعات الخطأ و اقل قيمة لمتوسط الاخطاء المطلقة كانت لطريقة المربعات الصغرى للشريحة التمهيدية في حالة وجود الارتباط.

جدول رقم (1) يبين قيم MSE, MAE لحالتى الاستقلالية والارتباط لكل من طريقتى (المقدرات اللبية غير المرتبطة ظاهريا، وطريقة المربعات الصغرى المعممة لمقدرات الشريحة التمهيدية)

Methods	حالة الاستقلالية	حالة الارتباط
Kernel	Mse=(30.60145)	Mse=(15.21513)
	Mae=(1156.274)	Mae=(283.6377)
Spline	Mse=(1.685315)	Mse=(1.681781)
	Mae=(4.290639)	Mae=(4.284745)

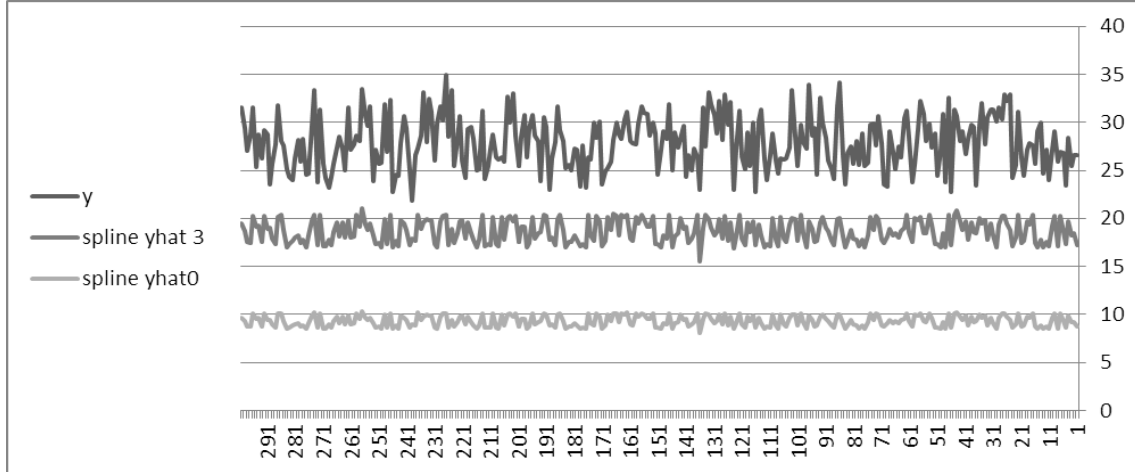
شكل رقم (1): يمثل قيم y الحقيقية و \hat{y}_0 يمثل تقدير النموذج باستخدام المقدرات اللبية غير المرتبطة ظاهريا في حالة الاستقلالية و \hat{y}_3 يمثل قيم المقدرات في حالة الارتباط .



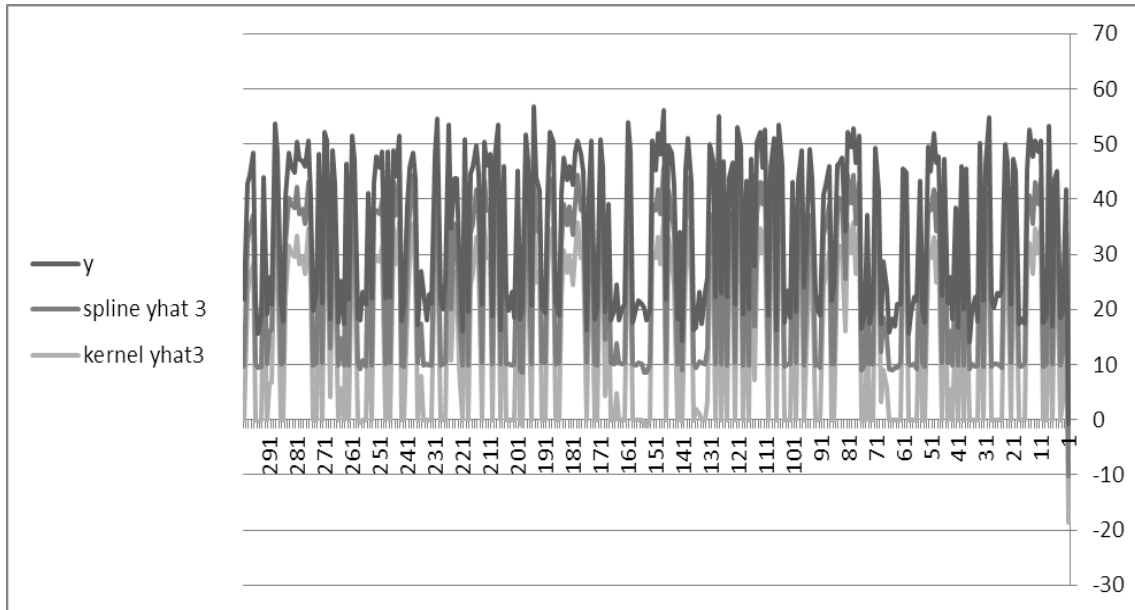


مقارنة طريقتي تقدير الدالة الالمعلمية لبيانات عنقودية عن كريات الدم البيضاء لعرضى اللوكيميا

شكل رقم (2): يمثل قيم y الحقيقية يمثل قيم y الحقيقية و \hat{y}_{00} يمثل تقدير النموذج باستخدام مقدرات المربعات الصغرى المعممة لممهد الشريحة في حالة الاستقلالية و \hat{y}_{03} يمثل قيم المقدرات في حالة الارتباط



شكل رقم (3): y يمثل القيم الحقيقية للمتغير المعتمد و \hat{y}_{03} يمثل تقدير النموذج باستخدام مقدرات المربعات الصغرى المعممة لممهد الشريحة في حالة الارتباط و \hat{y}_{3} يمثل تقدير النموذج باستخدام قيم المقدرات اللبية غير المرتبطة ظاهريا في حالة الارتباط .



من ملاحظة الاشكال الثلاثة فنرى من الشكل (1) مقدرات طريقة مقدر النواة غير المتصل ظاهريا من المقدرات Seemingly Unrelated Kernel مقارنة مع قيم y الحقيقية حيث تقترب قيم \hat{y}_{3} في حالة الارتباط من القيم الحقيقية اكثر من قيم المقدر الـ \hat{y}_{0} في حالة الاستقلالية .



مقارنة طريقتي تقدير الدالة اللامعلمية لبيانات عنقودية من كريات الدم البيضاء لمرضى اللوكيميا

اما في الشكل (2) فنلاحظ بان مقدرات طريقة المربعات الصغرى المعممة لمقدر الشريحة التمهيدية (Generalized least squares smoothing spline) مقارنة مع قيم y الحقيقية حيث تقترب قيم \hat{y}_{03} في حالة الارتباط من القيم الحقيقية اكثر من قيم المقدر الـ \hat{y}_{00} في حالة الاستقلالية. اما في الشكل (3) فيوضح مقدرات مقدر النواة غير المتصل ظاهريا Seemingly Unrelated Kernel و مقدرات طريقة المربعات الصغرى المعممة لمقدر الشريحة التمهيدية generalized least squares smoothing spline في حالة ارتباط البيانات مقارنة مع قيم y الاصلية حيث ان مقدرات مقدرات طريقة المربعات الصغرى المعممة لمقدر الشريحة التمهيدية Generalized least squares smoothing spline تقترب اكثر من القيم الحقيقية مقارنة مع قيم مقدرات النواة غير المتصل ظاهريا Seemingly Unrelated Kernel في حالة الارتباط لكلا المقدرين .

الاستنتاجات :

- 1- ان كل من الطريقتين (طريقة النواة غير المتصل ظاهريا) و(طريقة المربعات الصغرى لممهد الشريحة التمهيدية) حققت اقل MSE واقل MAE عند حساب مصفوفة الارتباط داخل العناقيد بعين العناية بالمقارنة مع حالة استقلال البيانات .
- 2- حققت طريقة المربعات الصغرى لممهد الشريحة التمهيدية اقل MSE واقل MAE مقارنة مع طريقة ممهد النواة غير المتصل ظاهريا في الحالتين (حالة الارتباط الحقيقي للبيانات) و(حالة الاستقلالية) .

التوصيات :

- 1- يمكن استخدام الطرائق المذكورة في البحث في حالة الامراض غير المسيطر عليها (النادرة) التي يكون تطورها غير معلوم .
- 2- يمكن اخذ حالة البيانات عنقودية غير المتزنة وتطبيق الطرائق المذكورة انفا .

المصادر:

المصادر العربية:

- 1- حمود د مناف يوسف (2005) ((مقارنة المقدرات اللامعلمية لتقدير دالة الكثافة الاحتمالية)) اطروحة دكتوراه في الاحصاء كلية الادارة والاقتصاد/جامعة بغداد.
- 2- علي د.عمر عبد المحسن (2007) ((مقارنة مقدرات النماذج التجميعية المعممة باستخدام الشرائح التمهيدية عند تحليل الانحدار اللامعلمي وشبه المعلمي)) اطروحة دكتوراه في الاحصاء كلية الادارة والاقتصاد/ جامعة بغداد.

المصادر الأجنبية:

- 3-Galbraith, S. Daniel, J. A .& Vissel ,B.(2010)((A Study of Clustered Data and Approaches to Its Analysis)),The Journal of Neuroscience ,30(32) : 10601-10608.
- 4-Garrett_Fitzmaurice,_Marie_Davidian,_Geert_Verbeke & Geert _Molenberghs .(2009) ((Longitudinal Data Analysis)) Chapman& Hall/ CRC.
- 5-[Http://www.jormedic.com/component/mailto/?tmpl=com](http://www.jormedic.com/component/mailto/?tmpl=com)
- 6-[Https://en.wikipedia.org/wiki/Cluster_analysis](https://en.wikipedia.org/wiki/Cluster_analysis)
- 7- Ibrahim, N, A.& Suliadi,2010a ((Analyzing Longitudinal Data Using Gee-Smoothing Spline)) . WSEAS Transactions on Systems and Control
- 8-Liang ,K.-Y. and Zeger, S.L. (1986) ((Longitudinal data analysis using generalized linear models)).Biometrika,73,13-22.



مقارنة طريقتي تقدير الدالة الالعملية لبيانات عنقودية من كريات الدم البيضاء لعرضى اللوكيميا

- 9- Lin, D. & Ying, Z. (2001) ((semiparametric and nonparametric regression analysis of longitudinal data)). J.Amer. Statist. Assoc. ,96,103-126.
- 10- Lin, X. Carroll, R. J.(2001) ((Semiparametric Regression For Cluster data), Biometrika,88,No.4;1179-1185.
- 11-Lin,X.&Carroll,R. J.(2000) (Nonparametric function estimation for clustered data when the predictor is measured without/with error) J.Amer. Statist. Assoc.,95,520-534al Data)).
- 12- Linton, O.B., Mammen, E., Lin, X. and Carroll, R.J.(2003) ((Accounting for correlation in marginal longitudinal nonparametric regression)) 2nd Seattle Symmp. Biostatistics.
- 13-Ma, S., Song., Q. and Wang,L.,2013,((Simultaneous variable selection and estimation in semi parametric modeling of longitudinal/clustered data))Bernoulli 19,252-274
- 14-Ruckstuhl,A.,Welsh,A.H.&Carroll,R.J.(2000).((Nonparametric function estimation of the relationship between two repeatedly measured variables)).Statist. Sinica, 10,51-71.
- 15- Wang, N.(2003)((Marginal nonparametric Kernel regression accounting for within-subject correlation)),Biometrika,90,43-52.
- 16- Wang, Y. G.& Zhao, Y.(2008)((Weighted rank regression for clustered data analysis)). [Biometrics](#). 64,:39-45.
- 17- Welsh, A. H., Lin, X.& Carroll, R. J.(2002).((Marginal longitudinal non parametric regression Locality and efficiency of spline and kernel methods)).

ملحق رقم (1)

القيم التالية تمثل قيم المتغيرات التوضيحية x المتمثلة بعدد كريات الدم البيضاء (wbc) لكل مريض خلال 3 زيارات متتابة ، اما فيما يخص بيانات متغير الاستجابة y المتمثل بهيموغلوبين الدم (hb)

التسلسل	قيم المتغير التوضيحي x			قيم متغير الاستجابة y	
	الزيارة الاولى	الزيارة الثانية	الزيارة الثالثة	الزيارة الاولى	الزيارة الثانية
العنقود الاول	63.7	1.4	1.15	9.38	7.80
العنقود الثاني	2.9	6.9	10.1	8.1	8.6
العنقود الثالث	24.3	16.8	18.2	7.27	13.3
العنقود الرابع	5.97	4.75	8.9	8.69	7.98
العنقود الخامس	0.8	1.16	7.15	6.1	8.71
العنقود السادس	3.16	1.82	1.57	8.15	8.25
العنقود السابع	8.84	8.08	7.8	6.76	6.1
العنقود الثامن	0.320	0.4	0.277	8.8	7.5
العنقود التاسع	9.50	1.56	0.754	8.82	8.26
العنقود العاشر	4.6	7.7	8.68	7.9	8.8
العنقود الحادي عشر	0.393	0.430	0.644	6.86	9.42
العنقود الثاني عشر	1.24	0.859	0.7	9.6	6.65
العنقود الثالث عشر	0.2	0.20	0.41	7.6	9.1
العنقود الرابع عشر	1.5	2.4	11.3	12.2	13.1
العنقود الخامس عشر	0.114	56.9	1.56	12.1	10.8



مقارنة طريقتي تقدير الدالة الالعملية لبيانات عنقودية من كريات الدم البيضاء لعرضى اللوكيميا

8	5.6	8.3	0.15	0.5	0.94	العنقود السادس عشر
10.6	10.2	7.2	1.5	6.0	11.7	العنقود السابع عشر
11.0	6.9	8.5	2.9	3.1	4.8	العنقود الثامن عشر
9.8	9.4	7.5	5.3	5.3	5.7	العنقود التاسع عشر
6.4	7.9	6.8	1.7	0.5	1.3	العنقود العشرون
6.2	8.6	8.9	6.2	1.7	0.9	العنقود الواحد والعشرون
11.0	10.8	10.7	35.7	11.1	13.7	العنقود الثاني والعشرون
9.8	9.10	7.5	2.93	2.15	2.1	العنقود الثالث والعشرون
8.00	6.1	7.1	0.981	0.027	0.44	العنقود الرابع والعشرون
14.4	12.5	18.1	3.74	5.47	21.5	العنقود الخامس والعشرون
11.3	12.0	12.7	0.478	1.43	5.15	العنقود السادس والعشرون
14.5	12.7	12.7	13.6	9.04	9.04	العنقود السابع والعشرون
9.9	10.2	10.0	9.3	1.85	9.9	العنقود الثامن والعشرون
12.0	12.2	12.2	5.82	7.01	4.70	العنقود التاسع والعشرون
12.9	10.3	13.1	0.147	3.05	0.148	العنقود الثلاثون
8.69	12.5	13.1	0.818	2.40	2.5	العنقود الواحد والثلاثون
11.3	12.9	11.9	6.08	4.00	5.02	العنقود الاثنان والثلاثون
12.7	13.2	12.8	6.21	7.01	1.56	العنقود الثلاث والثلاثون
8.0	7.11	7.9	6.80	13.3	17.69	العنقود الاربع والثلاثون
13.6	12.5	12.6	5.38	4.07	4.96	العنقود الخمس والثلاثون
9.63	7.4	9.6	3.85	68.6	7.38	العنقود الستة والثلاثون
6.85	6.0	4.9	11.9	11.6	22.83	العنقود السبعة والثلاثون
8.9	8.89	11.0	1.46	2.77	2.96	العنقود الثمانية والثلاثون
4.0	7.1	10.2	1.7	1.9	5.39	العنقود التسعة والثلاثون
7.7	9.2	10.6	0.653	0.96	1.50	العنقود الاربعون
10.6	5.19	6.89	3.54	4.19	6.01	العنقود الواحد والاربعون
11.2	10.8	10.2	5.00	3.65	3.55	العنقود الاثنان والاربعون
7.96	8.25	8.3	6.24	7.21	5.81	العنقود الثلاث والاربعون
7.4	8.9	9.8	0.26	2.85	49.7	العنقود الاربع والاربعون
6.90	10.8	10.9	1.34	2.29	11.3	العنقود الخمس والاربعون
5.7	8.0	5.6	0.2	0.15	0.5	العنقود الستة والاربعون
12	11.5	18.2	14.65	10.7	13.8	العنقود السبعة واربعون
9.6	10.4	6.75	0.81	1.80	0.129	العنقود الثمانية والاربعون
12.0	10.8	12.5	6.83	2.40	2.74	العنقود التسعة والاربعون
8.8	10.2	10.2	0.09	0.150	0.151	العنقود الخمسون
8.3	7.25	7.17	0.93	0.708	0.693	العنقود الواحد والخمسون
9.8	11.5	11.6	0.71	0.8	0.7	العنقود الاثنان والخمسون
5.68	9.7	7.98	2.40	8.9	4.75	العنقود الثالث والخمسون
12.2	9.5	9.5	5.1	4.32	11.7	العنقود الرابع والخمسون
10.8	11.7	9.6	3.49	4.27	2.9	العنقود الخامس والخمسون
11.3	11.0	12.4	4.71	7.0	59.8	العنقود السادس والخمسون
12.5	11.3	12.1	42.9	13.82	8.4	العنقود السابع والخمسون
9	10.4	9.5	4.18	5.2	6.2	العنقود الثامن والخمسون
8.3	7.6	5.5	10.2	8.3	7.6	العنقود التاسع والخمسون
9.5	10.2	6.2	2.3	1.3	1.2	العنقود الستون
9.1	10.1	9.9	2.0	1.9	3.0	العنقود الواحد والستون
11.8	10.7	10.8	17.8	12.6	14.1	العنقود الاثنان والستون
7	10.2	11.4	2.1	8.1	20.8	العنقود الثلاث والستون
7.6	7.9	7.7	18.7	14.3	21.9	العنقود الاربع والستون



مقارنة طريقتي تقدير الدالة الالعملية لبيانات عنقودية من كريات الدم البيضاء لعرضى اللوكيميا

10.4	10.5	9.4	25.2	24.4	26.1	العنقود الخمس والستون
7.7	9.7	6.6	5.6	15	22.8	العنقود الستة والستون
7.1	7.6	8.7	3.0	39.4	24.5	العنقود السبعة والستون
7.2	7.2	10.1	0.5	3.4	3.8	العنقود الثمانية والستون
5.4	5.1	5.2	61.9	53.6	25.9	العنقود التسعة والستون
7.1	7.3	6.1	0.4	1.3	1.0	العنقود السبعون
8.5	6.5	10	0.4	0.3	1.6	العنقود الواحد والسبعون
11	10.4	10.7	11.2	5.8	6.8	العنقود الاثنان والسبعون
6.6	8.2	7.3	0.5	8.9	10.0	العنقود الثلاث والسبعون
13.0	12.3	11.0	11.9	1.3	3.59	العنقود الاربع والسبعون
9.9	7.9	9.6	6.2	2.5	8.4	العنقود الخمس والسبعون
6.4	6	7.5	2.5	11	29.3	العنقود الستة والسبعون
7.5	6.2	8.4	0.2	0.06	0.2	العنقود السبعة وسبعون
10.5	10.3	11.1	1.5	0.7	1.5	العنقود الثمانية والسبعون
8.5	6.2	8.5	0.9	0.3	0.3	العنقود التسعة والسبعون
9.9	9.4	10.3	2.4	1.6	1.6	العنقود الثمانون
8.2	9	7.8	2	2.5	1.7	العنقود الواحد والثمانون
6.4	7.4	8.7	1.3	1.2	3.6	العنقود الاثنان والثمانون
7	8.1	9.2	0.8	1.1	1.3	العنقود الثلاث والثمانون
8.2	8.3	6.6	0.1	0.07	0.1	العنقود الاربع والثمانون
8.8	9.3	8.3	3.41	3.4	2.3	العنقود الخمسة والثمانون
7.7	8.8	14.1	14.1	11	7.7	العنقود الستة والثمانون
11.6	11.7	11.7	8.2	6.9	6.9	العنقود السبعة والثمانون
10.4	10.7	6.9	0.7	0.58	0.6	العنقود الثمانية والثمانون
7.81	8.3	7.4	1.77	1.91	1.4	العنقود التسعة والثمانون
4.6	5.2	7.3	3.9	1.5	3.3	العنقود التسعون
9.8	9.2	9.34	32.7	38.3	20.6	العنقود الواحد والتسعون
8.9	10.2	9.8	9.6	10.7	8.1	العنقود الاثنان والتسعون
8.7	5.3	13	1.1	3.1	5.5	العنقود الثلاث والتسعون
9.5	9.7	6.9	4.52	2.75	1.35	العنقود الاربعة والتسعون
6.2	10.8	11.9	4.16	1.75	1.15	العنقود الخمس والتسعون
18.1	10.6	9.7	9.1	8.6	3.7	العنقود الستة والتسعون
11.3	11.8	18.1	0.94	1.0	5.92	العنقود السبعة والتسعون
9.5	9.4	10.3	1.08	0.17	0.05	العنقود الثمانية والتسعون
10.7	11.6	9.5	3.38	4.26	2.8	العنقود التسعة والتسعون
12.1	9.4	9.4	5.0	4.21	11.6	العنقود المئة

ملحق رقم (2)

لحساب طريقه ال kernel للبيانات الحقيقية

```
clc
clear
data=xlsread('Rdata.xlsx','data');
n=size(data,1);
subj=data(:,1);x=data(:,2);y=data(:,3);
corry0=[1 0 0;0 1 0;0 0 1];
corry1=[1 0.5700 0.415;0.5700 1 0.584;0.415 0.584 1];
yy1=[y(1:n/3)];
```



مقارنة طريقتي تقدير الدالة الالعملية لبيانات منقودية من
كريات الدم البيضاء لعرضى اللوكيميا

```
yy3=[y(((n/3)+1):2*(n/3))];  
yy5=[y((2*(n/3)+1):n)];  
usubj=unique(subj);  
nsubj=length(usubj);  
xmin=min(x);  
xmax=max(x);  
xfit=linspace(xmin*1.05,xmax*.95,101)';  
[ux,flag1,flag2]=unique(x);  
nux=length(ux);  
dpoly=2;  
dy1=std(yy1);  
dy2=std(yy3);  
dy3=std(yy5);  
dy=[dy1 dy2 dy3];  
sdy=diag(dy);  
dsd01= repmat(dy1,n/3);  
dsd1=diag((dsd01));  
dsd02= repmat(dy2,n/3);  
dsd2=diag((dsd02));  
dsd03= repmat(dy3,n/3);  
dsd3=diag((dsd03));  
dsdy=[dsd1' dsd2' dsd3'];  
dsdy=diag(dsdy);  
eyey=eye(n/3);  
output=['h' 'gcv' 'df'];  
output1=['vh' 'vgcv' 'vdf'];  
cdsdy0=kron(eyey,corry0);  
V0=dsdy*cdsdy0*dsdy;  
Vd10=inv(diag(diag(V0)));  
Vd20=diag(diag(inv(V0)));  
[yhat00,hgcv00,vhgcv00,B00,gcf01,gcf02]=glpfit(data,V0,Vd10,Vd20,[0,dpoly,1]);  
mse00=mean((y-yhat00).^2);  
mae00=mean(abs(y-yhat00));  
m00=[mse00 mae00];  
xlswrite('D:\KHala.xlsx',m00,'m00','A1:B1');  
xlswrite('D:\KHala.xlsx',yhat00,'yhat0');  
xlswrite('D:\KHala.xlsx',output,'hgcv0','A1:C1');  
xlswrite('D:\KHala.xlsx',hgcv00,'hgcv0','A2:C2');  
xlswrite('D:\KHala.xlsx',output1,'vhgcv0','A1:C1');  
xlswrite('D:\KHala.xlsx',vhgcv00,'vhgcv0','A2:C10');  
saveas(gcf01,'figure1','png')  
saveas(gcf02,'figure2','png')
```




```
cdsdy1=kron(eyey,corry1);
V1=dsdy*cdsdy1*dsdy;
Vd11=pinv(diag(diag(V1)));
Vd21=diag(diag(pinv(V1)));
[yhat01,hgcv01,vhgcv01,B01,gcf03,gcf04]=glpfit(data,V1,Vd11,Vd21,[0,dpoly,1]);
mse01=mean((y-yhat01).^2);
mae01=mean(abs(y-yhat01));
m01=[mse01 mae01];
xlswrite('D:\KHala.xlsx',m01,'m01','A1:B1');
xlswrite('D:\KHala.xlsx',yhat01,'yhat1');
xlswrite('D:\KHala.xlsx',output,'hgcv1','A1:C1');
xlswrite('D:\KHala.xlsx',hgcv01,'hgcv1','A2:C2');
xlswrite('D:\KHala.xlsx',output1,'vhgcv1','A1:C1');
xlswrite('D:\KHala.xlsx',vhgcv01,'vhgcv1','A2:C10');
saveas(gcf03,'figure3','png')
saveas(gcf04,'figure4','png')
```

لحساب دالة الkernel للبيانات الحقيقية وفق البرنامج الاتي

```
function [yhat,hgcv,vhgcv,B,gcf1,gcf2]=glpfit(data,V,Vd1,Vd2,params,xfit,kstr)
n=size(data,1);
subj=data(:,1);
usubj=unique(subj);
nsubj=length(usubj);
x=data(:,2);
y=data(:,3);
[ux,~,flag2]=unique(x);
nux=length(ux);
xrange=max(ux)-min(ux);
if nargin<2||isempty(V),
    V=dsdy*cdsdy*dsdy;
end
V=V(flag2,flag2);
if nargin<3||isempty(Vd1),
    V=sdy*cdsdy*sdy;
    Vd1=inv(diag(diag(V)));
end
V=V(flag2,flag2);
Vd1=Vd1(flag2,flag2);
if nargin<4||isempty(Vd2),
    V=sdy*cdsdy*sdy;
    Vd2=diag(diag(inv(V)));
end
```



```
V=V(flag2,flag2);
Vd2=Vd2(flag2,flag2);
if nargin<5||isempty(params),
    h=0;
    dpoly=1;
    indfig=0;
elseif length(params)==1,
    h=params;
    dpoly=1;
    indfig=0;
elseif length(params)==2,
    h=params(1);
    dpoly=params(2);
    indfig=0;
elseif length(params)==3,
    h=params(1);
    dpoly=params(2);
    indfig=params(3);
end
if nargin<7||isempty(kstr),
    kstr='(1/sqrt(2*pi))*exp(-.5*t.^2)';
end
if h~=0,
    nh=1;
else
    nh=10;
    hmin=.5*(dpoly+1)*xrange/nux;
    hmax=xrange/8;
    vh = logspace(log10(hmin),log10(hmax),nh);
end
if nh>1,
    for ii=1:nh,
        h=vh(ii);
        for jj=1:nux,
            temp=x-ux(jj);
            t=(temp/h)*h^(-1);
            W=diag(eval(kstr));
            X=ones(n,1);
            for r=1:dpoly,
                X=[X,temp.^r];
            end
            AA=0;
            BB=[];
```



```
for kk=1:nsubj,
    flagi=(subj==usubj(kk));
    Wi=W(flagi,flagi);
    Xi=X(flagi,:);
    yi=y(flagi);
    ni=length(yi);
    Vi=V(flagi,flagi);
    Vd1i=Vd1(flagi,flagi);
    temp=Wi.*diag(diag(pinv(Vi)));
    AA=AA+Xi'*temp*Xi;
    BB=[BB,Xi'*Wi];
end
temp=pinv(AA)*BB;
iii=eye(size(temp));
temp1=pinv(iii+temp*(Vd1-Vd2))*temp*Vd1;
A(jj,:)=temp1(1,:);
end
B=A(flag2,:);
yhat=B*y;
df(ii)=trace(B);
gcv(ii)=mean((y-yhat).^2)/(1-df(ii)/n)^2;
end
vhgcv=[vh(:),gcv(:),df(:)];
[gcv,temp]=min(gcv);
h=vh(temp);
hgcv=vhgcv(temp,:);
end
for jj=1:nux,
    temp=x-ux(jj);
    t=(temp/h)*h^(-1);
    W=diag(eval(kstr));
    X=ones(n,1);
    for r=1:dpoly,
        X=[X,temp.^r];
    end
    AA=0;
    BB=[];
    for kk=1:nsubj,
        flagi=(subj==usubj(kk));
        Wi=W(flagi,flagi);
        Xi=X(flagi,:);
        yi=y(flagi);
        ni=length(yi);
```



مقارنة طريقتي تقدير الدالة الامعلمية لبيانات منقودية من
كريات الدم البيضاء لمرضى اللوكيميا

```
Vi=V(flagi,flagi);
  Vd1i=Vd1(flagi,flagi);
  temp=Wi.*pinv(Vi);
  AA=AA+Xi'*temp*Xi;
  BB=[BB,Xi'*Wi];
end
temp=pinv(AA)*BB;
iii=eye(size(temp));
temp1=pinv(iii+temp*(Vd1-Vd2))*temp*Vd1;
A(jj,:)=temp1(1,:);
end
efit=A*y;
B=A(flag2,:);
yhat=B*y;
df=trace(B);
dfa=xrange/h/sqrt(2*pi);
resid=y-yhat;
gcv=mean(resid.^2)/(1-df/n)^2;
hsig2=sum(resid.^2)/(n-df);
hgcv=[h,gcv,df];
uyfit=efit;
uysig=sqrt(hsig2*diag(A*A'));
if nh==1,
  vhgcv=hgcv(1:3);
end
if nargin<6||isempty(xfit),
  xfit=ux;
  yfit=uyfit;
  ysig=uysig;
else
  yfit=spline(ux,uyfit,xfit);
  ysig=spline(ux,uysig,xfit);
end
fits=[xfit,yfit,ysig];
if indfig==1,
  figure(1);
  clf;
  subplot(2,2,1);
  plot(x,y,'o',fits(:,1),fits(:,2),'r-','...',
    fits(:,1),[fits(:,2)+1.96*fits(:,3),fits(:,2)-1.96*fits(:,3)],'g-');
  xlabel('x');
  ylabel('y');
```



```
title(['(a) ',num2str(dpoly),'-th order GLP fit']);
subplot(2,2,2);
plot(x,resid/sqrt(hsig2),'o');
xlabel('x');
ylabel('standardized residual');
title('(b) Standardized residuals vs x');
subplot(2,2,3);
plot(yhat,resid/sqrt(hsig2),'o');
xlabel('yhat');
ylabel('standardized residual');
title('(c) Standardized residuals vs fits');
subplot(2,2,4);
plot(y,resid/sqrt(hsig2),'o');
xlabel('y');
ylabel('standardized residual');
title('(d) Standardized residuals vs responses');
gcf1=gcf;
if nh>1,
    figure(2);
    clf;
    subplot(2,2,1);
    plot(log10(vhgcv(:,1)),vhgcv(:,2),'-o');
    xlabel('log_{10}(h)');
    ylabel('GCV');
    title('(a) GCV vs log_{10}(h)');
    subplot(2,2,2);
    plot(log10(vhgcv(:,1)),vhgcv(:,3),'-o');
    xlabel('log_{10}(h)');
    ylabel('df');
    title('(b) df vs log_{10}(h)');
end
gcf2=gcf;
end
clearvars -except yhat hgcv vhgcv B gcf1 gcf2
end

لحساب طريقه spline
clc
clear
data=xlsread('Rdata.xlsx','data');
n=size(data,1);
subj=data(:,1);x=data(:,2);y=data(:,3);
corry0=[1 0 0;0 1 0;0 0 1];
corry3=[1 0.5700 0.415;0.5700 1 0.584;0.415 0.584 1];
```



مقارنة طريقتي تقدير الدالة الامعلمية لبيانات منقودية عن
كريات الدم البيضاء لمرضى اللوكيميا

```
yy1=[y(1:n/3)];
yy3=[y((n/3)+1:2*(n/3))];
yy5=[y(2*(n/3)+1:n)];
usubj=unique(subj);
nsubj=length(usubj);
xmin=min(x);
xmax=max(x);
xfit=linspace(xmin*1.05,xmax*.95,101)';
[ux,flag1,flag2]=unique(x);
nux=length(ux);
dy1=std(yy1);
dy2=std(yy3);
dy3=std(yy5);
dy=[dy1 dy2 dy3];
sdy=diag(dy);
dsd01= repmat(dy1,100);
dsd1=diag((dsd01));
dsd02= repmat(dy2,100);
dsd2=diag((dsd02));
dsd03= repmat(dy3,100);
dsd3=diag((dsd03));
dsdy=[dsd1' dsd2' dsd3'];
dsdy=diag(dsdy);
eyey=eye(n/3);
output=['spar' 'gecv'];
output1=['vspar' 'vgecv'];
cdsdy0=kron(eyey,corry0);
V0=dsdy*cdsdy0*dsdy;
Vd10=diag(diag(inv(V0)));
[yhat00,spargcv00,vspargcv00,A00,gcf01,gcf02]=ssfit(data,[],1,xfit,Vd10);
mse00=mean((y-yhat00).^2);
mae00=mean(abs(y-yhat00));
m00=[mse00 mae00];
xlswrite('D:\SHala.xlsx',m00,'m00','A1:B1');
xlswrite('D:\SHala.xlsx',yhat00,'yhat0');
xlswrite('D:\SHala.xlsx',output,'spargcv0','A1:B1');
xlswrite('D:\SHala.xlsx',spargcv00,'spargcv0','A2:B2');
xlswrite('D:\SHala.xlsx',output1,'vspargcv0','A1:B1');
xlswrite('D:\SHala.xlsx',vspargcv00,'vspargcv0','A2:B10');
saveas(gcf01,'figure1','png')
saveas(gcf02,'figure2','png')
cdsdy1=kron(eyey,corry1);
V1=dsdy*cdsdy1*dsdy;
```



مقارنة طريقتي تقدير الدالة اللامعلمية لبيانات منقودية عن كريات الدم البيضاء لمرضى اللوكيميا

```
Vd11=diag(diag(inv(V1)));  
[yhat01,spargcv01,vspargcv01,A01,gcf03,gcf04]=ssfit(data,[],1,x,Vd11);  
mse01=mean((y-yhat01).^2);  
mae01=mean(abs(y-yhat01));  
m01=[mse01 mae01];  
xlswrite('D:\SHala.xlsx',m01,'m01','A1:B1');  
xlswrite('D:\SHala.xlsx',yhat01,'yhat1');  
xlswrite('D:\SHala.xlsx',output,'spargcv1','A1:B1');  
xlswrite('D:\SHala.xlsx',spargcv01,'spargcv1','A2:B2');  
xlswrite('D:\SHala.xlsx',output1,'vspargcv1','A1:B1');  
xlswrite('D:\SHala.xlsx',vspargcv01,'vspargcv1','A2:B10');  
saveas(gcf03,'figure3','png')  
saveas(gcf04,'figure4','png')
```

لحساب دالة السبلاين للبيانات الحقيقية

```
function [yhat,spargcv,vspargcv,A,gcf1,gcf2]=ssfit (data, spar, indfig, xfit, V)  
subj=data(:,1);  
usubj=unique(subj);  
nsubj=length(usubj);  
x=data(:,2);  
y=data(:,3);  
n=length(x);  
[knots,~,flag2]=unique(x);  
nknot=length(knots);  
E=zeros(n,nknot);  
if nargin<3||isempty(indfig),  
    indfig=0;  
end  
if nargin<5||isempty(V),  
    V=eye(n);  
end  
for i=1:n,  
    E(i,flag2(i))=1;  
end  
S=E'*E;  
S=diag(diag(S).^(-1/2));  
Gmat=splmat(knots);  
[~,D]=eig(S*Gmat*S);  
d=diag(real(D));  
if nargin<2||isempty(spar),  
    nspar=15;
```



مقارنة طريقتي تقدير الدالة الالمعلمية لبيانات منقودية من
كريات الدم البيضاء لعرضى اللوكيميا

```
d0=sort(d);
dmin=d0(3);
dmax=max(d);
sparmax=((nknot-2)/(nknot-4)-1)/dmin;
sparmin=((nknot-2)/(4-2)-1)/dmax;
vspar = logspace(log10 (sparmin),log10 (sparmax), nspar)';
else
vspar=spar;
nspar=1;
end
for k=1:nspar,
A=E*pinv(E'*V*E+(n)*vspar(k)*Gmat)*E'*V;
yhat=A*y;
gcv(k)=mean((y-yhat).^2)/(1-trace(A)/n)^2;
end
vspargcv=[vspar(:),gcv(:)];
[gcv,temp]=min(gcv);
spar=vspar(temp(1));
spargcv=[spar,gcv];
A=E*pinv(E'*V*E+(n)*spar*Gmat)*E'*V;
yhat=A*y;
resid=y-yhat;
SSE=sum((y-yhat).^2);
hsig2=SSE/(n-trace(A));
ysig=sqrt(hsig2*diag(A*A'));
if nargin<4||isempty(xfit),
fits=[x,yhat,ysig];
else
[xx,xflag]=unique(x);
yyhat=yhat(xflag);
yysig=ysig(xflag);
yfit=spline(xx,yyhat,xfit);
ysig=spline(xx,yysig,xfit);
fits=[xfit,yfit,ysig];
end
if indfig==1,
figure(1);
clf;
[~,temp2]=sort(fits(:,1));
ffits=fits(temp2,:);
plot(x,y,'o',ffits(:,1),ffits(:,2),'r-',...
ffits(:,1),[ffits(:,2)-1.96*ffits(:,3),ffits(:,2)+1.96*ffits(:,3)],'g-');
gcf1=gcf;
if nspar>1,
```




مقارنة طريقتي تقدير الدالة الامعلمية لبيانات منقودية من كريات الدم البيضاء لعرضى اللوكيميا

```
figure(2);
clf;
yyhat=y-resid;
resid=resid/std(resid);
subplot(2,2,1)
plot(log10(vspargcv(:,1)),vspargcv(:,2),'-o')
xlabel('log_{10}(\lambda)')
ylabel('GCV score')
title('(b) GCV curve')
subplot(2,2,2)
plot(yyhat,resid,'o')
xlabel('SS fit')
ylabel('Standardized residual')
title('(a)')
subplot(2,2,3)
plot(y,resid,'o')
xlabel('Response')
ylabel('Standardized residual')
title('(c)')
subplot(2,2,4)
hist(resid)
xlabel('Standardized residual')
ylabel('Frequency')
title('(d)')
end
gcf2=gcf;
end
clearvars -except yhat spargcv vspargcv A gcf1 gcf2
end
برنامج لحساب مصفوفة الـ بساي في دالة السبلين
function [K,Q,R,h]=splmat(t)
[n,m]=size(t);
if n==1,n=m;
t=t';
end
h=t(2:n)-t(1:n-1);
hin=h.^(-1);
Q=diag([hin(1:n-2);0],-1)-diag([0;hin(1:n-2)+hin(2:n-1);0],0);
Q=Q+diag([0;hin(2:n-1)],1);
Q=Q(2:n-1,:);
R=diag(h(2:n-2),-1)/6+diag(h(1:n-2)+h(2:n-1))/3+diag(h(2:n-2),1)/6;
K=Q*inv(R)*Q';
clearvars -except K Q R h
end
```



Compared of estimating two methods for nonparametric function to cluster data for the white blood cells to leukemia patients

Abstract:

We can notice cluster data in social, health and behavioral sciences, so this type of data have a link between its observations and we can express these clusters through the relationship between measurements on units within the same group.

In this research, I estimate the reliability function of cluster function by using the seemingly unrelated Kernel Estimators method and the Generalized Least Squares Smoothing Spline Estimators method, and I applied these two methods on Leukemia patients and made a comparison between the two methods by using MSE and MAE comparison standard, the empirical results showed the efficiency of the Generalized Least Squares Smoothing Spline Estimators method.

Keywords: cluster data, the seemingly unrelated Kernel Estimators method, and the Generalized Least Squares Smoothing Spline Estimators method, MSE, MAE.