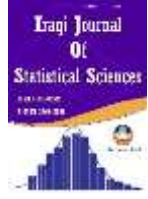




المجلة العراقية للعلوم الإحصائية

<http://stats.uomosul.edu.iq>



اختيار المتغيرات في نموذج الانحدار اللوجستي باستخدام خوارزمية اليراعات المضيئة المعدلة

هبة سليمان داؤد

قسم الاحصاء والمعلوماتية, كلية علوم الحاسوب والرياضيات , جامعة الموصل , الموصل , العراق

الخلاصة

يعتبر نموذج الانحدار اللوجستي هو الأكثر استخدامًا في العديد من التطبيقات ، وهو من النماذج الرئيسية في عائلة النماذج الخطية المعممة.. وكغيره من سائر نماذج الانحدار , قد يحتوي النموذج على متغيرات مستقلة كثيرة ما يؤثر سلباً على دقة النموذج وبساطته في تفسير النتائج. تهدف هذه الدراسة إلى استخدام خوارزمية اليراعات المضيئة ومقارنتها مع طرائق اخرى في اختيار المتغيرات في نموذج الانحدار الاسي باستخدام المحاكاة والبيانات الحقيقية . وأظهرت النتائج أنه بالمقارنة مع الطرائق الأخرى المستخدمة سابقاً، فإن الاسلوب المقترح يؤدي أداء أفضل ويساعد على خفض متوسط مربع الخطأ للنموذج.

معلومات النشر

تاريخ الاستلام : 25 كانون الاول 2023
تاريخ القبول : 28 نيسان 2024
تاريخ القبول: 5 أيار 2024
متاح على الانترنت 1 حزيران 2024

الكلمات المفتاحية:

اختيار المتغيرات، خوارزمية اليراعات المضيئة، المحاكاة، النموذج الانحدار الاسي.

المراسلة:

هبة سليمان داؤد

heba.sulaiman82@uomosul.edu.iq

DOI [10.33899/IQJOSS.2024.183255](https://doi.org/10.33899/IQJOSS.2024.183255) , ©Authors, 2024, College of Computer Science and Mathematics University of Mosul.
This is an open access article under the CC BY 4.0 license (<http://creativecommons.org/licenses/by/4.0/>).

1. مقدمة Introduction

تعالج التقنيات الجديدة النمو الهائل للبيانات، إذ تساعد هذه التقنيات الباحثين على نقل كميات هائلة من البيانات إلى المعلومات. قد تحتوي البيانات الضخمة على متغيرات غير ذات صلة أو زائدة عن الحاجة. ولذلك يفضل الباحثون اختيار المهم من هذه المتغيرات عن طريق اختيار مجموعة فرعية صغيرة من المتغيرات المهمة المتوفرة من مجموعات البيانات. تعتبر دراسة أي مشكلة أو ظاهرة من المجالات الاقتصادية، الاجتماعية، الطبية أو غيرها، من أهم أسس البحث العلمي. فالغاية الرئيسية من دراستها هي تحديد المعادلة الرئيسية التي تمثل تلك الظاهرة بدقة، وذلك عن طريق جمع البيانات المتعلقة بها من مختلف المصادر المتاحة. ومن ثم يتم تحليل تلك البيانات باستخدام تقنيات الإحصاء والتحليل الرياضي لتحديد العلاقات بين المتغيرات المختلفة وتصميم نماذج إحصائية تصف تلك العلاقات. وهذا يشكل المدخل الأساسي لفهمها بشكل أعمق وتحديد معالمها الرئيسية. ويشار إلى أن هذه العملية في علم الإحصاء بنمذجة الظواهر (Månsson, 2013). ومن بين جميع نماذج الانحدار الخطي المعممة، يمكن القول أن نموذج الانحدار اللوجستي هو أحد أشهر هذه النماذج ، حيث يتم استخدامه بشكل واسع في العديد من التطبيقات.

الانحدار اللوجستي هو احد نماذج المعتمد على التصنيف الثنائي من خلال التنبؤ باحتمالية حدوث نتيجة أو حدث. يقدم النموذج نتيجة ثنائية أو ثنائية التفرع تقتصر على نتيجتين محتملتين: نعم/لا، 1/0، أو صحيح/خطأ. يقوم الانحدار اللوجستي بتحليل العلاقة بين واحد أو أكثر من المتغيرات المستقلة ويصنف البيانات إلى فئات منفصلة. يتم استخدامه على نطاق واسع في النمذجة التنبؤية، حيث يقوم النموذج بتقدير الاحتمال الرياضي لما إذا كان المثل ينتمي إلى فئة معينة أم لا (Alharthi, Lee, & Algama, 2021).

غالبية البيانات في الواقع التطبيقي الحقيقي تحتوي على مشاكل مثل مشكلة العدد الكبير من المتغيرات المستقلة المدروسة، وهي من المشاكل المعروفة لدى الباحثين الإحصائيين، وتؤثر سلباً على عملية التقدير. في بعض الحالات، يمكن أن تؤدي هذه المشكلة إلى تجاهل بعض المتغيرات التوضيحية المهمة. حيث أصبحت الاساليب التقليدية لاختيار المجموعات الجزئية غير جيدة في أداء وظيفتها حيث أصبحت أكثر تكلفة في حسابها، إضافة إلى ذلك فإن معايير المعلومات لاختيار المتغيرات مثل معيار أكاكي للمعلومات (Akaike information (AIC)) ومعيار بيز للمعلومات (Bayesian information criterion (BIC)) أصبحت غير عملية في اختيار المتغيرات التوضيحية وذلك بسبب تعقيدها الحسابي الذي ينمو بشكل طردي مع ازدياد عدد المتغيرات التوضيحية (Özkale & Arican, 2018).

يهدف هذا البحث إلى توظيف خوارزمية البراعات المضيفة المعدلة ومقارنتها مع طرائق إختيار المتغيرات التوضيحية في أنموذج الانحدار اللوجستي الأخرى باستخدام المحاكاة والبيانات الحقيقية، من خلال تسليط الضوء على عدد من العوامل التي قد تؤثر على جودة هذه الطرائق ووجوب استخدامها ضمن شروط معينة دون غيرها من الطرائق.

يعتبر العالم Yang أول من استخدم خوارزمية البراعات المضيفة عام 2007 وطورها في عام 2009 لاجراء التصنيفات للبيانات الهندسية وبعدها تم استخدامها في الكثير من المجالات الاحصائية مثل تنقيب البيانات و في جال تعلم الآلة وغيرها من المجالات.

2. نموذج الانحدار اللوجستي (LRM) Logistic Regression Model

يعرف الانحدار اللوجستي على انه نوع من انواع نماذج الانحدار اللاخطية الذي تكون فيه العلاقة بين المتغير التابع (الاستجابة) (y_i) ومجموعة من المتغيرات التوضيحية (x_1, x_2, \dots, x_k) علاقة غير خطية (Özkale & Arican, 2016)، إذ يكون فيها المتغير التابع (الاستجابة) متغير نوعي (Varathan & Wijekoon, 2018). قد يأخذ المتغير التابع في نموذج الانحدار اللوجستي صفتين فقط ويرمز لهاتين الصفتين بـ (0) او (1)، وهو ما يطلق عليه بالانحدار اللوجستي الثنائي (Binary logistic regression). اما فيما يخص المتغيرات التوضيحية، فيمكن ان تكون هذه المتغيرات مستمرة او متقطعة سواءا وصفية كانت او عددية.

يبنى نموذج الانحدار اللوجستي على فرض أساسي هو أن المتغير التابع الذي نهتم بدراسته هو متغير ثنائي الصفة ويتبع توزيع برنولي وفق الدالة الاحتمالية المعرفة بالصيغة الآتية (Özkale & Arican, 2018)

$$p(Y = y_i) = \pi_i^{y_i} (1 - \pi_i)^{1-y_i} \quad (1)$$

اذ ان: π_i : تمثل احتمال حدوث الاستجابة عندما $y_i = 1$ و $1 - \pi_i$: تمثل احتمال عدم حدوث الاستجابة عندما $y_i = 0$.

يمكن تعريف الاحتمال (π_i) رياضياً بدلالة المتغيرات التوضيحية والدالة اللوجستية وكما في الصيغة الآتية:

$$\pi_i = \frac{e^{X_i \beta}}{1 + e^{X_i \beta}} \quad (2)$$

اذ ان: β : متجه من المعلمات أبعاده $(p \times 1)$ و $X_i = \{x_{i1}, x_{i2}, \dots, x_{ip}\}$: متجه صفي من المتغيرات التوضيحية أبعاده $(1 \times p)$.

يكون الهدف الرئيس من الانحدار اللوجستي الثنائي هو تفسير التغير في قيم المتغير التابع من خلال تفسير حدوث الاستجابة باحتمال (π_i) او عدم حدوث الاستجابة باحتمال $(1 - \pi_i)$. بناءً على ذلك، وكما هو معروف عند بناء نموذج الانحدار، فإن من خلال المعادلة الأخيرة (2) يتضح ان العلاقة بين المتغير التابع والمتغيرات التوضيحية علاقة غير خطية وغالباً ما تأخذ الدالة اللوجستية شكلاً منحنياً. ويلجأ الاحصائيون غالباً إلى التحويل الخطي لهذه النماذج لإزالة انحناءات معلماتها وذلك لتأثير هذه الانحناءات السلبية في حالة وجودها على خصائص المقدرات اذ بالإمكان افتراض علاقة معينة تربط بين المتغير التابع والمتغيرات التوضيحية الأخرى، لذلك تم اقتراح تحويل دالة اللوجت (Logit Function) التي تقوم بتحويل علاقة الانحدار اللاخطية بين المتغيرات التوضيحية ودالة احتمال الاستجابة (π_i) في نموذج الانحدار اللوجستي الى علاقة انحدار خطي، وذلك من خلال اخذ اللوغارتم الطبيعي للمقدار $(\frac{\pi_i}{1-\pi_i})$ (Steyerberg, Borsboom, van Houwelingen, Eijkemans, & Habbema, 2004; Varathan & Wijekoon, 2018) كما مبين في المعادلات الآتية:

$$\text{logit}(\pi_i) = \ln \frac{\pi_i}{1-\pi_i} \quad (3)$$

$$\text{logit}(\pi_i) = \ln \frac{e^{X_i\beta}}{1+e^{X_i\beta}} \quad (4)$$

$$\text{logit}(\pi_i) = \ln(e^{X_i\beta}) = X_i\beta = (B_0 + \sum_{j=1}^k B_j x_{ij}) \quad (5)$$

اذ ان: $\beta_0, \beta_1, \dots, \beta_p$: معالم مجهولة يتم تقديرها.

دالة الإمكان الأعظم لنموذج الانحدار اللوجستي الذي يتبع توزيع برنولي تكون بالصيغة الآتية

$$y_i \sim \text{Bernoulli}(\pi_i) \quad (6)$$

إذا:

$$L(\beta) = \prod_{i=1}^n \pi_i^{y_i} (1 - \pi_i)^{1-y_i} \quad (7)$$

$$L(\beta) = \prod_{i=1}^n (1 - \pi_i) \exp \left[\sum_{i=1}^n y_i \ln \left(\frac{\pi_i}{1-\pi_i} \right) \right] \quad (8)$$

وحسب خاصية التحويل اللوجستي (دالة اللوجستيك) فان

$$\ln \frac{\pi_i}{1-\pi_i} = X_i\beta \quad (9)$$

فان $L(\beta)$ تساوي:

$$L(\beta) = \prod_{i=1}^n (1 - \pi_i) \exp \left[\sum_{i=1}^n y_i (X_i\beta) \right] \quad (10)$$

وبأخذ اللوغاريتم إلى دالة الإمكان.

$$\ln L(\beta) = \sum_{i=1}^n \ln(1 - \pi_i) + \left[\sum_{i=1}^n y_i (X_i\beta) \right] \quad (11)$$

إذا:

$$\ln L(\beta) = \left[\sum_{i=1}^n y_i (X_i\beta) \right] - \sum_{i=1}^n \ln(1 + e^{X_i\beta}) \quad (12)$$

وبأخذ المشتقة الأولى الى لوغاريتم دالة الإمكان الأعظم ثم مساواة المشتقة بالصفر.

$$\frac{\partial \log L(\beta, X)}{\partial \beta} = 0 \quad (13)$$

فان المعادلات الناتجة من المشتقة الأولى هي معادلات غير خطية والتي ليس لها حل واضح, لذلك يتم حل هذه المعادلات عن طريق الطرق العددية والتي منها الطريقة العددية الأكثر شيوعاً هي خوارزمية نيوتن رافسون .

3. خوارزمية اليراعات المضيئة (FFA) Firefly Algorithm

في السنوات الأخيرة أصبح الاهتمام متزايد بتصميم خوارزميات التحسين المستوحاة من الطبيعة وتطويرها، حاول الباحثون إيجاد الإلهام من مصادر مختلفة في الطبيعة مثل النحل والنمل واليراعات والاسماك والطيور والنباتات وأنظمة الامواج والأنهار. يعد ذكاء السرب أداة مهمة لحل العديد من المشكلات المعقدة في البحث العلمي، إذ تمت دراسة خوارزميات ذكاء السرب على نطاق واسع حيث تم تطبيقها بنجاح على مجموعة متنوعة من مشكلات التحسين المعقدة نظراً لتمتعها بالبساطة والمرونة والكفاءة العالي (Yang, 2010).

تعتمد معظم خوارزميات التحسين المستوحاة من الطبيعة على ذكاء السرب، وتشكل الخوارزميات القائمة على ذكاء السرب جزءاً كبيراً من الخوارزميات المعاصرة، وأصبحت هذه الخوارزميات مستخدمة على نطاق واسع في التحسين وتحليل البيانات وكذلك في التعلم الآلي والذكاء الاصطناعي. وتعد خوارزمية اليراعات المضيئة (Firefly Algorithm FA) واحدة من أحدث أساليب ذكاء السرب الجديدة واغوى خوارزميات التحسين التي تم تطويرها لأول مرة من قبل الباحث Yang في بداية عام 2008.

أثبتت الخوارزمية أنها فعالة وذات أداء جيد في حل مشكلات التحسين المختلفة. تم إيجاد خوارزمية اليراعات من محاكاة السلوك الاجتماعي لليراعات المضيئة على أساس جاذبية الفلاش (الأضواء الساطعة) من خلال تمثيل ميزة بعض الخصائص الواضحة لليراعات وكيفية التفاعل معها، إذ أن وميض اليراعة هو نظام إشارة يستخدم لجذب يراعة أخرى (Long Zhang, Shan, & Wang, 2016; Li Zhang, Srisukkhom, Neoh, Lim, & Pandit, 2018).

حيث يمكن حساب المسافة بين اثنين من اليراعات في المواقع والمسافة الديكارتية والتي يمكن حسابها باستخدام المعادلة الآتية :

$$r_{ij} = \|x_i - x_j\| = \sqrt{\sum_{d=1}^D (x_{id} - x_{jd})^2} \quad (14)$$

إذ أن x_i : هو موقع اليراعة i إذ $x_i = \{x_{1i}, x_{2i}, \dots, x_{di}\}$ ، وأن x_j : موقع اليراعة j إذ $x_j = \{x_{1j}, x_{2j}, \dots, x_{dj}\}$ ، عدد الأبعاد D ، وأن $d \in D$

يمكننا تلخيص آلية عمل خوارزمية اليراعات (FA) بالخطوات الآتية (Yang, 2010):

1- جميع اليراعات للجنسين، إذ يمكن أن تتجذب كل يراعة إلى كل اليراعات الأخرى. إذ أن اليراعات الأقل جاذبية (إشراقاً) تتجذب إليها اليراعات الأكثر جاذبية (إشراقاً).

2- تتناسب جاذبية اليراعة مع شدة الضوء الذي يتناقص كلما زادت المسافة عن اليراعات الأخرى.

3- يتم تحديد جاذبية اليراعة من خلال موقعها داخل مساحة البحث.

4- تؤدي القيمة الأفضل لوظيفة اللياقة في موقع معين إلى زيادة جاذبية اليراعة.

لكل فراشة شدة ضوء أو سطوع يتم استخدام قيمته لتقييم جودتها. إن سطوع اليراعة i في موقع معين x نستطيع أن نشير إليه بالآتي:

$$I(x_i) = f(x_i) \quad (15)$$

حيث أن شدة ضوء اليراعة تتناسب طردياً مع سطوعها وترتبط بالقيم الموضوعية. عند المقارنة بين اليراعات، تتجذب اليراعة التي لها شدة ضوء منخفضة نحو اليراعة الأخرى ذات الضوء الأعلى، شدة ضوء اليراعة تعتمد على I_o من الضوء المنبعث من اليراعة والمسافة r_{ij} بين زوج من اليراعات. يمكن وصف شدة الضوء $I(r)$ من خلال دالة متناقصة بشكل رتيب لـ r_{ij} والتي يمكن صياغتها كالآتي:

$$I(r) = I_o e^{-(\gamma r_{ij})^2} \quad (16)$$

γ : هو عامل امتصاص تأثير الضوء.

ونظراً لأن الجاذبية لكل فراشة تتناسب مع شدة الضوء التي تراها اليراعات المجاورة، لذلك يجب السماح للجاذبية بالتنوع باختلاف درجة الإمتصاص، حيث يمكن تحديد الشكل الرئيسي لتباين الجاذبية Z بالمعادلة التالية: (Xu, Yu, Chen, & Zuo, 2018)

$$Z(r) = Z_o e^{-(\gamma r_{ij})^2} \quad (17)$$

إذ أن $Z(r)$ تمثل دالة جاذبية اليراعة عند المسافة r و Z_o هي الجاذبية الأولية لليراعة عند مسافة $(r = 0)$

ويمكن أن تكون ثابتة. عند التنفيذ Z_o تساوي الواحد ولمعظم المشاكل. حيث يتم تحديث الحركة للفراشات حسب المعادلة الآتية:

$$x_i^{(t+1)} = x_i^{(t)} + Z_o e^{-\gamma r_{ij}^2} (x_j^{(t)} - x_i^{(t)}) + \alpha_i \mathcal{E}_i^t \quad (18)$$

إذ أن α_i : هو معامل التوزيع العشوائي. \mathcal{E}_i^t : متجه لأرقام عشوائية مأخوذة من توزيع Uniform.

يعتمد تأثير هذه الحركة العشوائية في المعلمة α_i فيما إذا تم إختياره ليكون كبيراً فإن الحل x_i سيتحرك بشكل عشوائي مبتعداً عن الموقع، بخلاف إذا كان α_i صغيرة جداً، فستتحرك في الموقع وقد تصبح ضئيلة مقارنة بالحركة نحو اليراعات الأكثر إشراقاً.

في BFFA، تُستخدم وظيفة النقل لتعيين مساحة بحث مستمرة إلى مساحة ثنائية، وتم تصميم عملية التحديث لتبديل مواقع النجوم بين 0 و 1 في مساحات البحث الثنائية. من أجل بناء هذا المتجه الثنائي، وظيفة النقل في المعادلة (18) يمكن استخدامها، حيث يكون الحل الجديد مقيداً بالقيم الثنائية فقط

$$x_i^t = \begin{cases} 1 & \text{if } T(x) > \alpha \\ 0 & \text{otherwise} \end{cases} \quad (19)$$

إذ أن $\alpha \in [0, 1]$ هي عبارة عن رقم عشوائي وإن $T(x)$ هي دالة تحويل. إن دالة التحويل تعرف بالشكل الآتي:

$$T_{BCSA}(x_i^t) = \frac{1}{1 + e^{10(x_i^t - 0.5)}}, \quad (20)$$

في هذا البحث تم اقتراح استخدام دالة تحويل متغيرة خلال الزمن. أي ان دالة التحويل هذه سوف تتغير خلال تكرار الحل. تم هذا الاقتراح من خلال اضافة معلمة تحكم وهي θ ، اذا احتاج هذه المعلمة الى قيمة عليا وقيمة دنيا لها من خلال المعادلة الخاصة بها وهي:

$$\theta = \theta_{\min} + (\theta_{\max} - \theta_{\min}) e^{-t}, \quad (21)$$

وعليه سوف تصبح دالة التحويل المقترحة بالشكل التالي:

$$T_{TV}(x_i^t) = \frac{1}{1 + e^{10(x_i^t - 0.5)/\theta}}, \quad (22)$$

من أجل إتمام هدف البحث وتحقيقه، وبالاعتماد على هذه التقنية، فإن كل عنصر (براعة) في المجموعة سيكون لديه d من المواقع التي تمثل عدد المتغيرات التوضيحية في النموذج الانحدار اللوجستي. بناءً على ذلك، فإن توظيف خوارزمية البراعات المضيفة تكون وفق الخطوات التالية:

الخطوة الأولى: تحديد حجم المجموعة (عدد البراعات) وهو 30 فراشة، حيث إن كل فراشة سيكون له متجه من عدد المتغيرات المستقلة فضلاً عن ذلك تحديد عدد التكرارات داخل خوارزمية البراعات المضيفة حيث استقرت النتائج عند التكرار 500.

الخطوة الثانية: توليد القيم الأولية التي تحتاجها الخوارزمية، التي ستمثل القيم الأولية الافتراضية ، فإن توليدها سيكون من التوزيع المنتظم المستمر وفق الفترة $[0,1]$.

الخطوة الثالثة: لغرض اختيار القيم المثلى، تم الاعتماد على Fitness Function وفق الصيغة الآتية:

$$\text{Fitness Function} = \min \left[\frac{\sum_{i=1}^n (y_i - \hat{m}(\mathbf{X}))^2}{n} \right] \quad (23)$$

الخطوة الرابعة: بالاعتماد على أقل قيمة تحصل عليها أي فراشة وفق المعادلة (22) يتم تحديث مواقع باقي البراعات.

الخطوة الخامسة: نستمر بالحل لحين الوصول الى أعلى تكرار للخوارزمية، الذي تم تحديده بالخطوة الأولى والذي سيمثل الحل الأمثل.

| | | | | |
|-------|-------|-------|-----------|-------|
| x_1 | x_2 | | x_{p-1} | x_p |
| 1 | 0 | | 1 | 0 |

الشكل 1: آلية اختيار المتغيرات حسب خوارزمية البراعات المضيفة

3- معايير تقييم طرائق اختيار المتغيرات Criteria for Evaluating Methods for Selecting Variables

3-1 معايير تقييم دقة التنبؤ

اولاً: خطأ التنبؤ (PE) (Prediction Error)

ويعرف بأنه مربع الفرق بين القيمة الحقيقية لمتغير الاستجابة والقيمة التنبؤية المرافقة له، ويعرف رياضياً بالمعادلة التالية :

$$PE = (y - \hat{y})^T (y - \hat{y}) \quad (24)$$

وبالاعتماد على هذا المعيار يتم تحديد الطريقة الأفضل التي تعطي اقل قيمة مقارنة بالطرائق الأخرى.

ثانياً: معايير تقييم دقة اختيار المتغيرات

بما ان الطرق المقترحة بصورة عامة تعمل على اختيار المتغيرات، لذلك من المهم تقييم وقياس قدرة هذه الطرائق وجودتها في كيفية اختيار المتغيرات المهمة. ولذلك، تم الاعتماد على معيارين في دراستنا لهذا الغرض وبالشكل التالي:

(1) معيار التقييم "C"

هو معيار التقييم الذي يرمز له بـ (C) والذي يعرف بأنه عدد المعاملات الحقيقية ذات القيم الصفرية والتي تم تقديرها بشكل صحيح على انها ذات قيم صفرية.

(2) معيار التقييم "I"

معيار التقييم الذي يرمز له بـ (I) وهو يعرف على انه عدد المعاملات الحقيقية ذات القيم غير الصفرية والذي تم تقديرها بشكل غير صحيح على انها ذات قيم صفرية. تعتمد جودة طرائق الجزاء من ناحية معايير تقييم دقة اختيار المتغيرات على من يعطي اعلی قيمة لـ (C) واقل قيمة لـ (I) .

3- نتائج المحاكاة Simulation Results

لقد تم تصميم تجربة ومحاكاتها باستخدام لغة البرمجة (R) حيث تم توليد المتغير (y_i) في انموذج انحدار كاوس المعكوس، حيث تم استخدام اسلوب مونت كارلو (Mont Carlo) في المحاكاة حيث تم تعيين قيم حجم العينات (n) حيث تم استخدام ثلاث احجام من العينات وهي (30,100,150,250) وذلك لأجل دراسة المقارنة وفق العينات باختلاف أنواعها. سوف تتم المقارنة مع كل من طريقة معيار بيز ومعيار اكاكي.

اولاً : تم توليد بيانات المتغير y التي تتبع انموذج الانحدار اللوجستي وكالاتي :

$$\hat{Y} = \frac{e^{x_i\beta}}{1 + e^{x_i\beta}}$$

ثانياً : تم توليد مصفوفة المتغيرات التوضيحية X ذات ابعاد $(n \times p)$ التي تتبع التوزيع الطبيعي المتعدد (Multivariate Normal Distribution) كالاتي :

$$X \sim MN(\mu, M)$$

حيث ان M هي مصفوفة التباين المشترك، حيث ان $M_{ij} = r^{|i-j|}$ عندما $(i, j = 1, 2, \dots, p)$ حيث ان المتغيرات التوضيحية تكون مرتبطة.

ثالثاً : تم تكرار التجربة (100) مرة وذلك لغرض تقليل التحيز في تجارب مونت كارلو (Mont Carlo).

رابعاً : تم توليد بيانات نموذج انحدار بواسون تبعاً لقيم متجه معاملات الانحدار β الذي ابعاده $(1 \times p)$ وكانت قيم متجه معاملات الانحدار β كالاتي $\beta = (1.8, 2.5, 1, -4, -7, 0, \dots, 0)^T$, حيث ان المعلمات غير الصفرية عددها $q = 5$, وان المعلمات الصفرية تساوي $p - q$. اذ تم اعتبار $p=10, 50, 100$

الجدول الآتية توضح النتائج العملية:

جدول (1) : معدل معايير تقييم طرائق الاختيار عندما n=30

| p | Method | PE | C | I |
|-----|--------|--------|---|---|
| 10 | AIC | 24.503 | 1 | 0 |
| | BIC | 22.955 | 2 | 0 |
| | FFA | 17.722 | 5 | 0 |
| 50 | AIC | 22.879 | 3 | 0 |
| | BIC | 21.331 | 3 | 0 |
| | FFA | 16.098 | 5 | 0 |
| 100 | AIC | 22.112 | 2 | 1 |
| | BIC | 20.564 | 3 | 0 |
| | FFA | 15.331 | 5 | 0 |

جدول (2) : معدل معايير تقييم طرائق الاختيار عندما n=100

| p | Method | PE | C | I |
|-----|--------|--------|---|---|
| 10 | AIC | 23.465 | 1 | 0 |
| | BIC | 21.917 | 3 | 0 |
| | FFA | 16.684 | 5 | 0 |
| 50 | AIC | 21.841 | 3 | 0 |
| | BIC | 20.293 | 4 | 0 |
| | FFA | 15.06 | 5 | 0 |
| 100 | AIC | 21.074 | 2 | 1 |
| | BIC | 19.526 | 3 | 0 |
| | FFA | 14.293 | 5 | 0 |

جدول (3) : معدل معايير تقييم طرائق الاختيار عندما $n=150$

| p | Method | PE | C | I |
|-----|--------|--------|---|---|
| 10 | AIC | 21.687 | 2 | 0 |
| | BIC | 20.139 | 3 | 0 |
| | FFA | 14.906 | 5 | 0 |
| 50 | AIC | 20.063 | 4 | 0 |
| | BIC | 18.515 | 4 | 0 |
| | FFA | 13.282 | 5 | 0 |
| 100 | AIC | 19.296 | 3 | 1 |
| | BIC | 17.748 | 3 | 0 |
| | FFA | 12.515 | 5 | 0 |

جدول (4) : معدل معايير تقييم طرائق الاختيار عندما $n=250$

| p | Method | PE | C | I |
|-----|--------|--------|---|---|
| 10 | AIC | 20.649 | 1 | 0 |
| | BIC | 19.101 | 2 | 0 |
| | FFA | 13.868 | 5 | 0 |
| 50 | AIC | 19.025 | 3 | 0 |
| | BIC | 17.477 | 4 | 0 |
| | FFA | 12.244 | 5 | 0 |
| 100 | AIC | 18.258 | 3 | 1 |
| | BIC | 16.71 | 3 | 0 |
| | FFA | 11.477 | 5 | 0 |

- سيتم تحليل وتفسير نتائج تجربة المحاكاة تبعاً لمعايير دقة التنبؤ ومعايير دقة اختيار المتغيرات. من خلال ملاحظة الجدول (1) و (2) و (3) و (4) الذي يوضح قيم معايير كل من (PE, C, I) للطرائق BIC و AIC والطريقة المقترحة FFA يمكن استخلاص ما يلي:
- 1- عندما تتغير قيمة معلمة التشتت ويغض النظر عن قيمة حجم العينة، يتبين أن طريقة (FFA) أعطت أقل قيم (PE) حيث بلغ مقدار التحسن بالتنبؤ بالاعتماد على المعيار (PE) بمقدار 35.14% و 30.86% عند ($n=50$) و $\tau = 1.5$ مقارنة بـ (AIC و BIC) على الترتيب.
 - 2- عندما يتغير حجم العينة ويغض النظر عن قيمة معلمة التشتت، أعطت طريقة (FFA) أفضل النتائج مقارنة بالطرائق الأخرى حيث تحسن التنبؤ بالاعتماد على المعيار (PE).
 - 3- بالاعتماد على معايير اختيار المتغيرات، فقد امتلكت طريقة (FFA) أعلى قيم (C) الذي هو عدد المعاملات الحقيقية ذات القيم الصفرية والتي تم تقديرها بشكل صحيح على أنها ذات قيم صفرية، وأعطت أقل قيم (I) الذي يعرف أنه عدد المعاملات الحقيقية ذات القيم غير الصفرية والذي تم تقديرها بشكل غير صحيح على أنها ذات قيم صفرية.
 - 4- ظهرت طريقة AIC كأسوأ طريقة في اختيار المتغيرات لأنها تعطي أعلى قيم لـ (PE) وكذلك كأسوأ طريقة في اختيار المتغيرات كونها تميل إلى اختيار متغيرات توضيحية غير مهمة.

5- الجانب التطبيقي Application Part

في هذا الجانب، يتم إجراء مقارنة بين أداء الطريقة المقترحة ومقدرات أخرى عن طريق استخدام البيانات الحقيقية. لغرض اتمام الفائدة المرجوة من البحث والطريقة المقترحة، تم التطبيق على بيانات تحتوي على تعدد خطي بين المتغيرات التوضيحية والتي أخذت من بيانات استخدمت من قبل (النعي، اسوان محمد طيب، 2005) حول مرض التلاسيميا الذي يصاب به الأطفال وبحجم 150 مريض. وقد تم اختيار عشرة متغيرات توضيحية وهي: العمر الحقيقي للطفل (بالشهر) (X1)، عمر المريض عند المرض مقاساً (بالشهر) (X2)، تضخم الكبد مقاساً (بالسنتمتر) (X3)، هيموكلوبين الدم (X4)، مكداس الدم (خلايا الدم المضغوطة) (X5)، الخلايا الشبكية (X6)، ارومة حمراء (X7)، الهيموكلوبين الجيني (X8)، عدد وحدات الدم (X9)، بداية نقل الدم حسب العمر مقاساً (بالشهر) (X10). في حين يمثل متغير الاستجابة وهو متغير ثنائي الصفة: العمر من العظم مقاساً بالشهر أكبر من أو يساوي 60 و العمر من العظم مقاساً بالشهر أقل من 60.

تم إجراء تقييم لنموذج الانحدار اللوجستي باستخدام طرائق اختيار المتغيرات المشار إليها من خلال حساب قيم متوسط مربعات الخطأ وكذلك عدد المتغيرات المستقلة التي تم اختيارها. توضح النتائج الملخصة في الجدول رقم 5 أن الأسلوب المقترح FFA تفوقت في الأداء على الطرائق الأخرى، حيث حققت أدنى قيمة لـ MSE وأقل عدد من المتغيرات المستقلة التي تم اختيارها.

جدول 5: نتائج الجانب التطبيقي

| Method | MSE | Variables |
|--------|--------|-----------|
| AIC | 39.561 | 7 |
| BIC | 37.248 | 6 |
| CSA | 27.931 | 4 |

6- الاستنتاجات Conclusion

تشير النتائج التي تم الحصول عليها من خلال المحاكاة والبيانات الحقيقية في انموذج الانحدار اللوجستي إلى أن استخدام أسلوب FFA يؤدي إلى نتائج ممتازة عند استخدام معيار MSE و PE، مما يجعله موثوقاً للمستخدمين في التنبؤ بالنتائج وتقييم النماذج الإحصائية. وبالإضافة إلى ذلك، يبدو أن حجم العينة n له تأثير كبير على قيم PE، حيث تتخفف قيمه PE عند زيادة حجم العينة، مما يعني زيادة الدقة. وعلى الجانب الآخر، عند زيادة قيمة عدد المتغيرات المستقلة نلاحظ أيضاً انخفاضاً في قيمة PE. ويجدر بالذكر أن استخدام معيار MSE مع يؤدي إلى نتائج أفضل في التنبؤ بالنتائج وتقييم النماذج الإحصائية. علاوة على ذلك ان الأسلوب المقترح ابدى قوته باختيار اقل عدد من المتغيرات المستقلة.

Reference

1. Al-Naimi, Aswan Muhammad Tayyab Rashid, 2005, "Testing Variables in Letter Regression," unpublished master's thesis, College of Computer Science and Mathematics, University of Mosul, Iraq.
2. Alharthi, A. M., Lee, M. H., & Algamal, Z. Y. (2021). Gene selection and classification of microarray gene expression data based on a new adaptive L1-norm elastic net penalty. *Informatics in Medicine Unlocked*, 24. doi:10.1016/j.imu.2021.100622
3. Månsson, K. (2013). Developing a Liu estimator for the negative binomial regression model: method and application. *Journal of Statistical Computation and Simulation*, 83(9), 1773-1780.
4. Özkale, M. R., & Arıcan, E. (2016). A new biased estimator in logistic regression model. *Statistics*, 1-21. doi:10.1080/02331888.2015.1123711
5. Özkale, M. R., & Arıcan, E. (2018). A first-order approximated jackknifed ridge estimator in binary logistic regression. *Computational Statistics*, 34(2), 683-712. doi:10.1007/s00180-018-0851-6
6. Steyerberg, E. W., Borsboom, G. J., van Houwelingen, H. C., Eijkemans, M. J., & Habbema, J. D. (2004). Validation and updating of predictive logistic regression models: a study on sample size and shrinkage. *Stat Med*, 23(16), 2567-2586. doi:10.1002/sim.1844
7. Varathan, N., & Wijekoon, P. (2018). Liu-Type logistic estimator under Stochastic Linear Restrictions. *Ceylon Journal of Science*, 47(1). doi:10.4038/cjs.v47i1.7483
8. Xu, H., Yu, S., Chen, J., & Zuo, X. (2018). An Improved Firefly Algorithm for Feature Selection in Classification. *Wireless Personal Communications*, 102(4), 2823-2834. doi:10.1007/s11277-018-5309-1
9. Yang, X.-S. (2010). *Nature-inspired metaheuristic algorithms*: Luniver press.
10. Zhang, L., Shan, L., & Wang, J. (2016). Optimal feature selection using distance-based discrete firefly algorithm with mutual information criterion. *Neural Computing and Applications*, 28(9), 2795-2808. doi:10.1007/s00521-016-2204-0
11. Zhang, L., Srisukham, W., Neoh, S. C., Lim, C. P., & Pandit, D. (2018). Classifier ensemble reduction using a modified firefly algorithm: An empirical evaluation. *Expert Systems with Applications*, 93, 395-422. doi:10.1016/j.eswa.2017.10.001

Variable Selection In Logistic Regression Model Using Modified Firefly Algorithms

Heba Suleiman Dawood

Department of Statistics and Informatics, College of Computer Science and Mathematics, University of Mosul, Mosul, Iraq

Abstract: The logistic regression model is considered the most widely used in many applications, and it is one of the main models in the family of generalized linear models. Like other regression models, the model may contain many independent variables, which negatively affects the accuracy of the model and its simplicity in interpreting the results. This study aims to use the modified firefly algorithm and compare it with other methods for selecting variables in an exponential regression model using simulation and real data. The results showed that compared to other previously used methods, the proposed method performs better and helps reduce the mean square error of the model.

Keyword: Selection of Variables, Firefly Algorithm , Simulation , Exponential Regression Model.