

مقارنة بين بعض طرائق تقدير حجم العينة لتقدير معلمات أنموذج التصنيف في حالة وجود بيانات عالية الابعاد باستعمال المحاكاة^(*)

أ.د. دجلة إبراهيم مهدي
كلية الادارة والاقتصاد
جامعة بغداد

dr.dejela.mahdi@gmail.com

الباحث: باسم سعدون جاسم
مديرية التربية ديالى
وزارة التربية

blasm5517@gmail.com

المستخلص:

تم في هذا البحث استعمال عدة طرائق لتحديد حجم العينة الامثل لتقدير معلم البيانات ذات الابعاد العالية (High dimensional: HDD) التي يكون فيها عدد المتغيرات التوضيحية اكبر من حجم العينة ($n < P$). وهذه الطرائق هي طريقة متباينة بونفروني وهي حالة خاصة من التقرير الطبيعي وطريقة متباينة بيرشتاين. يتم تقدير انموذج الانحدار اللوجستي الثنائي اللاخطي بحجم عينة كل طريقة في حالة بيانات عالية الابعاد باستعمال طريقة الذكاء الاصطناعي وهي طريقة الشبكة العصبية الاصطناعية (ANN) كونها تعطي تقدير عالي الدقة بما يتناسب مع نوع البيانات ونوع الدراسة الطبية. يتم توظيف القيم الاحتمالية التي تم الحصول عليها من الشبكة العصبية الاصطناعية في حساب مؤشر اعادة التصنيف الصافي (NRI)، تم كتابة برنامج لهذا الغرض باستعماله لغة البرمجة الاحصائية (R) حيث تم الاعتماد على معيار متوسط اكبر خطأ مطلق (MME) لمؤشر شبكة اعادة التصنيف الصافي (NRI) للمقارنة بين طرائق تحديد حجم العينة وبوجود عدد المعلمات الافتراضية مختلفة في ظل قيمة هامش خطأ معين (ϵ) . للتحقق من اداء الطرائق باستعمال معايير المقارنة اعلاه حيث كانت اهم الاستنتاجات هي ان طريقة متباينة بيرشتاين هي الافضل في تحديد حجم العينة الامثل باختلاف عدد المعلمات الافتراضية وقيمة هامش الخطأ.

الكلمات المفتاحية: متباينة بونفروني، الشبكة العصبية الاصطناعية، متباينة بيرشتاين، متوسط اكبر خطأ مطلق.

A comparison between some methods of estimating the sample size to estimate the parameters of the classification model in the case of high-dimensional data using simulations

Researcher: Balasim Saadoun Jassim
Diyala education directorate
Ministry of Education

Prof. Dr. Dejela Ibrahim Mahdi
College of Administration and Economics
University of Baghdad

Abstract:

In this research, several methods were used to determine the optimum sample size to estimate the parameters of high-dimensional data (HDD). Where the number of explanatory variables is greater than the sample size ($P > n$). These methods are the Bonferroni inequality method a special case of normal approximation and the

(*) بحث مستقل من رسالة ماجستير.

Bernstein's inequality method. The non-linear logistic regression model is estimated in the sample size of each method in the case of high-dimensional data using the artificial intelligence method, which is the artificial neural network method (ANN), as it gives a high-precision estimate commensurate with the data type and type of medical study. The probabilistic values obtained from the artificial neural network are used in the calculation of the net reclassification index. A program was written for this purpose using the statistical programming language (R) where the mean maximum absolute error criterion (MME) of the net reclassification network index (NRI) was used to compare the methods of specifying the sample size and the presence of a number of different default parameters under the given margin of error value (ϵ). To verify the performance of the methods using the comparison criteria above where the most important conclusions were that the Bernstein's inequality method is the better in determining the optimal sample size according to the number of default parameters and the error margin value.

Keywords: Bonferroni inequality, Artificial neural network, Bernstein's inequality, Mean maximum absolute error.

١. المقدمة:

زاد اهتمام الباحثين في السنوات الاخيرة بتحديد حجم العينة الامثل للحصول على دقة وتقدير كافية للحصول على معلم عالية الدقة وذلك لتقييم عدد كبير من الاختبارات في مجال التشخيص في ان واحد.

من اجل دراسة اداء دقة التشخيص للحصول على مجموعة بيانات مناسبة بحجم عينة معقول. قد لا يمكننا الحصول على حجم العينة ذات كفاءة احصائية لتحقيق نتائج معنوية ذات دقة عالية، من ناحية اخرى لا يمكن اجراء دراسة بحجم عينة كبيرة جدا يؤدي الى فقدان الدقة في النتائج المرجوة وكذلك الجهد والوقت والكلفة المادية. فهناك عدة طرق لحساب حجم العينة الامثل في مجالات العلوم الاجتماعية والاقتصادية والطبية وخصوصا في مجال الطب التشخيصي.

لذلك فان مسألة تقدير حجم العينة في هذه الحالة أصبح أصعب واعتقد لذلك لابد من استعمال طرق حساب حجم عينة تتناسب مع البيانات عالية الابعاد وكذلك تتناسب مع معايير الدقة، اذ اقترح (Pencina et al.,, ٢٠٠٨) مؤشر اعادة التصنيف الصافي Net reclassification index: NRI (اذ يستعمل هذا المؤشر في تحسين دقة التشخيص عند اضافة علامات حيوية جديدة الى العلامات الحيوية الاساسية (٦)). والحصول على انموذج جديد يضمن المتغيرات الاساسية بإضافة الى العلامات الجديدة، ولتكن العلامات الحيوية الجديدة على سبيل المثال الجينات التي تخص بروتين معين الخاصة بمرضى سرطان الثدي حيث يتم اضافة هذه العلامات الى العلامات الحيوية الاساسية مثلما العمر والجنس وعدد ضربات القلب ودرجة الحرارة للجسم وضغط الدم ومعدل التنفس وغيرها لغرض تحسين دقة التشخيص للمصابين بهذا المرض وللحكم عن كل هذه العملية تستعمل معيار دقة وهو مؤشر اعادة التصنيف الصافي (NRI) لذلك لابد من تحديد حجم عينة مناسب قياس دقة التشخيص باستعمال هذا المؤشر.

٢. هدف البحث: يهدف البحث الى جانبين اساسيين هما:

أولاً. تحديد أفضل حجم عينة ليعطي أفضل النتائج من الممكن الاعتماد عليها في اتخاذ القرار الاحصائي والطبي.

ثانياً. ايجاد أفضل طريقة لتقدير معلمات النموذج والتي تمثل أفضل المعالم الافتراضية للبيانات
عالية الابعاد.

٣. الجانب النظري:

٣-١. طرائق تحديد حجم العينة (Sample Size Determination Methods): وفي هذا البحث سوف نستعمل عدة طرائق لتحديد حجم العينة الامثل في ظل وجود بيانات ذات الابعاد العالية لتقدير مؤشر اعادة التصنيف الصافي (Net reclassification index) (NRI) باستعمال أنماذج الانحدار اللوجستي اللاخطي الثنائي ويمكن سرد هذه الطرق على النحو الآتي:

٣-١-١. طريقة التقرير الطبيعي (Normal Approximation Method): تفترض هذه الطريقة أن المقدر يمتلك توزيعاً طبيعياً محاذياً (Asymptotical normal) :

$$\sqrt{n}(\hat{\beta}_j - \beta_j) \sim N(0, \sigma_j^2), \quad k = 1, 2, \dots, p \quad \dots (1)$$

وحيث ان متوجه الأخطاء يفترض ان يتوزع توزيعاً طبيعياً بمتوسط صفر وتبين ثابت (σ_j^2) يمكن تحقيق ذلك من خلال خواص تقدير المعلمات. (Jiang and Li, 2018: 2)، تستخدم مثل هذه النتائج المقاربة لحساب حجم العينة عندما يكون حجم العينة n كبير فاذا كان لدينا:

$$Pr\left(\frac{\sqrt{n}|(\hat{\beta}_j - \beta_j)|}{\sigma_j} > Z_{\frac{\alpha}{2}}\right) < \alpha \quad \dots \dots \dots (2)$$

$$n^* = \frac{Z_{\alpha/2}^2 \sigma_j^2}{\epsilon^2} \quad \dots \dots \dots (3)$$

اذ ان:

n^* : حجم العينة.

$Z_{\alpha/2}$: القيمة الجدولية للتوزيع الطبيعي القياسي التي تقابل مستوى المعنوية (قيمة الثقة) والتي تستخرج من جداول التوزيع الطبيعي القياسي.

σ_j^2 : تباينات المجتمع غير المعلومة.

ان استعمال الصيغة (3) لحساب حجم العينة تضمن ان خطأ التقدير لمعلمة المجتمع لـ (β_k) ضمن حدود حد الخطأ التقدير(ϵ) ، باحتمالية ($1 - \alpha$) هذه الصيغة هي مناسبة إذا كنا ندرس فقط معلمة واحدة $\epsilon = 1$ (يعني مؤشر حيوى واحد).

٣-١-١-١. متباعدة بونفروني (Bonferroni Inequality): نستعمل متباعدة بونفروني من اجل التوسيع في حساب حجم العينة للصيغة (3) في حالة اكثراً من معلمة واحدة ($P > 1$) حيث تم تعريفها من قبل (Janos Galambos) التي تمثل مجموع الحوادث الاحتمالية التي تتوزع حسب توزيع برنولي للعزوم. (Simes, 1986: 751) اذ ان:

$$Pr\left(\max_j |\hat{\beta}_j - \beta_j| > \epsilon\right) \leq \sum_{j=1}^p Pr\left(|\hat{\beta}_j - \beta_j| > \epsilon\right) \dots (4)$$

حيث ان مجموع الاحتمالات لقيم العلية لخطا التقدير اكبر من ϵ وهي قيمة صغيرة جدا حيث ان $0 > \epsilon$. وحدود احتمال الخطأ لكل معلمة تساوي $jth \alpha/p$. لـ α/p من المقدرات حسب الصيغة الاحتمالية التالية:

$$Pr\left(\frac{\sqrt{n}|\hat{\beta}_j - \beta_j|}{\sigma_j} > Z_{\alpha/(2p)}\right) < \alpha/p \quad \dots \dots \dots (5)$$

ولتحقيق أكثر عمومية الى احتمال الخطأ والحصول على صيغة لحساب حجم العينة في حالة البيانات عالية الابعاد. نقوم بتعظيم صيغة حجم العينة رقم (٣) باستعمال هذه المتباينة حيث نحصل على الصيغة الآتية:

$$n^* = \frac{Z_{\alpha/(2p)}^2 v}{\epsilon^2} \quad \dots \dots \dots (6)$$

اذ ان:

$Z_{\alpha/(2p)}$: تمثل القيمة المعيارية بمستوى ثقة معين والتي تستخرج من جداول التوزيع الطبيعي المعياري.

$$v = \max_j \sigma_i^2$$

اذ ان:

v : تمثل القيم العظمى لبيانات المجتمع لكل j .

٣-٢. طريقة المتباينات (Inequalities Method): في بعض الأحيان يكون تطبيق التقرير الطبيعي باستخدام نظرية الغاية المركزية غير مرغوب فيه وقد لا يكون التوزيع احادي القيمة المتماثل مناسباً لشكل التوزيع وذلك لوصف المعلمات المقدرة في العينات غير المحدودة وخاصة عندما تكون المعلمات الحقيقة ضمن الحدود الطبيعية .

٣-٢-١. متباينة بير شتاين (Bernstein inequality): هي طريقة لحساب حجم العينة تم اقتراها من قبل الباحث (Van der Laan and Bryan) في عام (٢٠٠١)م باستعمال (Van der Laan and Bryan, 2001: 445) .(Bernstein inequality)

$$Pr(|\tilde{\beta}_j - \beta_j| > \epsilon) = pr\left(\left|\sum_{i=1}^n [\psi_{ij} - \beta_j]\right| > n\epsilon\right) \leq 2 \exp\left(\frac{-n\epsilon^2}{2v_j + 2\frac{M\epsilon}{3}}\right)$$

ولتحقيق حدود احتمال الخطأ تحت مستوى معنوية (α).

نستعمل أسلوب (Bonferroni Correction) يعتبر أحد الطرق المهمة للتتصدي لمشكلة الاختبارات المتعددة.

من خلال (Bonferroni Correction) يمكن ان نعمل على تجزئة لمستوى المعنوية الاختبار لكل معلمة في النموذج وتقليل عدد المعالم التي تكون تحت مستوى المعنوية المخصص لكل منها وجعلها مساوية للصفر تماما، من خلال ذلك نحدد النموذج الملائم للبيانات تحت الدراسة.

وذلك بأخذ الطرف الايمن للمتباعدة الذي يساوي تقدير احتمال الخطأ (α/P) .(Van Der vaart and Wallner, 1996: 102)

$$2\exp\left(\frac{-n\epsilon^2}{2\nu_j + 2\frac{M\epsilon}{3}}\right) = \frac{\alpha}{p} \quad \dots\dots\dots (7)$$

من المعادلة اعلاه نحصل على صيغة لحساب حجم العينة التالية:

$$n^* = \frac{2\nu_j + \frac{2M\epsilon}{3}}{\epsilon^2} \left(\log p + \log \frac{2}{\alpha} \right) \quad \dots\dots\dots (8)$$

اذ ان:

p : تمثل عدد المعلمات. $M > 0$: قيمة ثابتة اكبر من الصفر.

ν_j : وتمثل مقدار الخطأ المسموح به في تقدير المتوسط أو النسبة (ويساوي مضروب القيمة المعيارية عند مستوى ثقة معين في الخطأ المعياري للمتوسط أو مضروباً في الخطأ المعياري للنسبة في حالة تقدير النسبة).

٣-٣. **نموذج الانحدار اللوجستي الثاني:** الانحدار اللوجستي هو أحد عناصر مجموعة من النماذج تسمى بمجموعة النماذج الخطية العامة ويستعمل نموذج الانحدار اللوجستي لوصف العلاقة بين متغير الاستجابة الذي يكون من النوع المتقطع والمتغيرات التوضيحية تكون من النوع المتقطع او المستمر او الخلط بينهما كما يعتبر من الاساليب المهمة في تصنیف البيانات وتخصیصها لفئات معینة وتكون فيه العلاقة بين متغير الاستجابة (y) والمتغيرات التوضيحية (x_i) ايًّا كان نوعها كمياً او نوعياً غير خطية حيث يكون متغير الاستجابة (y) ثانٍ الاستجابة يتوزع برنولي (Bernoulli distribution) مفترضاً احدى القيمتين (١,٠)، وان احتمال حدوث الاستجابة (π_i) واحتمال عدم حدوث الاستجابة ($\pi_i - 1$). كما يعتبر حالة خاصة من نماذج الانحدار العامة (GLRM) فأن دالة الكثافة الاحتمالية يمكن كتابتها على النحو الاتي: (Scott, 2005: 5)

$$p(Y=y_i) = \pi_i^{y_i} (1-\pi_i)^{1-y_i} \quad \dots\dots\dots (9)$$

اذ ان:

$$y_i = 0, 1$$

y_i : متغير الاستجابة الثنائي.

π_i : احتمال حدوث الاستجابة عندما $y_i = 1$.
فإن توقع متغير الاستجابة يمثل احتمال حدوث الاستجابة.

$$E(y_i) = p(Y=1) = \pi_i \quad \dots\dots\dots (10)$$

فيكون تباين متغير الاستجابة حسب توزيع برنولي (Bernoulli distribution) هو.

$$V(y_i) = \pi_i(1-\pi_i) \quad \dots\dots\dots (11)$$

ولتكن x_1, x_2, \dots, x_k مجموعة من المتغيرات التوضيحية وان n تمثل عدد المشاهدات لهذه المتغيرات التي تكون المصفوفة X .

$$X = (x_{ij})_{n \times k}$$

اذ ان:

$i = 1, 2, \dots, n$ تمثل حجم العينة.
 $p = k + 1$ عدد المتغيرات التوضيحية ، $j = 1, 2, \dots, k+1$ يمثل عدد المعلمات.

فإذا كان $y_i = [y_1, y_2, \dots, y_n]$ عينه عشوائية من المتغير الثنائي الاستجابة وأن $\{0, 1\}$. وبذلك تكون صيغة أنموذج الانحدار اللوجستي كالتالي:

$$y_i = \pi_i + e_i \quad \dots \dots \dots (12)$$

اذ ان:

π_i تمثل دالة الانحدار اللوجستي (احتمال الاستجابة).

$$\pi_i = \frac{e^{\beta_0 + \sum_{j=1}^p X_{ij}\beta_j}}{1 + e^{\beta_0 + \sum_{j=1}^p X_{ij}\beta_j}} \quad \dots \dots \dots (13)$$

اذ ان:

β_j : متوجه المعلمات من درجة $(1 \times p)$ اذا كان:

$$i = 1, 2, \dots, n, j = 1, 2, \dots, p$$

$(n \times p)$: مصفوفة من المتغيرات التوضيحية من درجة $X_{ij} = \{x_{i1}, x_{i2}, \dots, x_{ij}\}$

n : يمثل عدد المشاهدات.

p : يمثل عدد المتغيرات التوضيحية.

حيث ان حد الخطأ العشوائي (e_i) يتوزع برنولي بمتوسط صفر كما في الصيغة:

$$e_i = y_i - \pi_i \quad \dots \dots \dots (14)$$

وبأخذ التوقع لطيفي المعادلة (14) نحصل على:

$$E(e_i) = \pi_i - \pi_i = 0 \quad \dots \dots \dots (15)$$

أما تباين حد الخطأ العشوائي فانه يساوي تباين متغير الاستجابة الثنائي.

$$V(e_i) = \pi_i(1 - \pi_i) \quad \dots \dots \dots (16)$$

ولذا فإن حد الخطأ العشوائي يكون له متوسط صفر وتباین $(1 - \pi_i)\pi_i$ ومن الملاحظ

ان تباين حد الخطأ يعتمد على قيم احتمال الاستجابة (π_i). (Shen and Gao, 2008: 4).

٤-٣. التحويل الخطى الى دالة الانحدار اللوجستي الثنائى: بسبب الانحناءات السلبية في معلمات الانحدار اللوجستي والتي تؤثر هذه الانحناءات على خصائص مقدرات المعلمات وقيم حدوث الاستجابة التي يمكن التنبؤ بها لذلك يلجأ الكثير من الاحصائيين الى التحويل الخطى من خلال استعمال تحويل دالة (Logit function) من اجل ازالة انحناء معلماتها حيث تكون نتائج الاختبار مضللة بسبب ان المعلمات لا تتوزع طبيعيا وتكون متحيزه وتبایناتها لا تكون اقل ما يمكن. وقد قام الباحث Berkson عام (١٩٤٤) م بايجاد علاقة لوغاريتمية من اجل تحويل العلاقة بين المتغيرات التوضيحية و احتمال حدوث الاستجابة (π_i) إلى علاقة خطية عن طريق

تحويل الاحتمال $\text{pr}(\mathbf{Y} = \mathbf{1}|\mathbf{X})$ إلى دالة يكون مدها من $(-\infty, +\infty)$ وهذه الدالة تسمى نسبة الارجحية والتي تعطى بالصيغة الآتية:

$$\frac{\text{pr}(\mathbf{Y} = \mathbf{1}|\mathbf{X})}{\text{pr}(\mathbf{Y} = \mathbf{0}|\mathbf{X})} = \frac{\text{pr}(\mathbf{Y} = \mathbf{1}|\mathbf{X})}{1 - \text{pr}(\mathbf{Y} = \mathbf{1}|\mathbf{X})} \dots \dots \dots \quad (17)$$

حتى نحصل على دالة مدها من $(-\infty, +\infty)$ نأخذ اللوغاريتم الطبيعي لنسب الارجحية (Odde) ويكون بالصيغة الآتية:

$$\text{logit } (\pi_i) = \ln \left(\frac{\pi_i}{1 - \pi_i} \right) \dots \dots \dots \quad (18)$$

$$\ln \left(\frac{\pi_i}{1 - \pi_i} \right) = \ln \left(e^{\beta_0 + \sum_{j=1}^p X_{ij} \beta_j} \right) \dots \dots \dots \quad (19)$$

$$\ln \left(\frac{\pi_i}{1 - \pi_i} \right) = \beta_0 + \sum_{j=1}^p X_{ij} \beta_j \dots \dots \dots \quad (20)$$

ويمكن اعادة كتابة الأنماذج كما في الشكل الآتي:

$$\ln \left(\frac{\pi_i}{1 - \pi_i} \right) = \beta_0 + \beta_u \mathbf{U} + \beta_1 \mathbf{x}_{i1} + \beta_2 \mathbf{x}_{i2} + \dots + \beta_j \mathbf{x}_{ij}, \dots \dots \dots \quad (21)$$

اذ ان:

$$i = 1, 2, \dots, n, \quad j = 1, 2, \dots, p$$

U : متغير الاساسي

$X_{i1}, X_{i2}, \dots, X_{ij}$: تمثل المشاهدات الخاصة بالمتغيرات التوضيحية.

$\beta_0, \beta_1, \beta_2, \dots, \beta_j$: معلمات نموذج الانحدار اللوجستي.

وهذا يعني ان نموذج الانحدار اللوجستي الثاني أصبح خطياً باستعمال تحويل دالة (Logit) (Kleinbaum and Klein, 2002: 17-20)

٥-٣. الشبكة العصبية الاصطناعية ((ANN)) (Artificial Neural Net Work) الشبكات العصبية الاصطناعية (ANN) عبارة عن شبكات من الخلايا العصبية الاصطناعية متصلة ببعضها البعض لتعلم كعنصرو معالجة بسيطة من أجل القيام بمهمة محددة. ان المعرفة التي اكتسبتها الشبكة من بيئتها تشبه معالجة الخلايا العصبية في الدماغ البشري، حيث يتم الحصول على المعلومات من عملية التعلم وتخزينها من نقاط القوة في الاتصال بين الخلايا العصبية، والتي تسمى بالأوزان ايضاً.

هناك انواع مختلفة من الشبكات فمنها شبكات متعددة الطبقات (Multilayer Net Work) والتي استعملت بشكل واسع لحل مشاكل الانحدار اذ يحتوي نموذج الشبكة العصبية متعددة الطبقات على طبقة إدخال (عقد المصدر)، او طبقة مخفية واحدة او أكثر (عقد حسابية) وطبقة الإخراج (عقد حساب). حيث ان طبقة الإدخال لها وظيفة استقبال إشارات الداخلة وأن الطبقات المخفية مسؤولة عن معالجة ونشر الإشارات المستقبلة (الداخلة) وابراجها. (Haykin, 2008: 27) فقد تكون لدى بعض الشبكات العصبية متعددة الطبقات اتصالات خاصة تسمى (Skip Layer)،

حيث يكون بعض إشارات الإدخال اتصال مباشر بطبقة الإخراج أي تنفيذ الطبقة المخفية.
(Ripple, 1996: 144)

$$y_i = \phi_i \left(b_i + \sum_{l=1}^m W_{kl} X_l + \sum_{i=1}^n W_{il} \cdot \phi_k \left(\sum_{j=1}^m W_{kj} X_j + b_j \right) \right) \dots \dots \dots (22)$$

تبين المعادلة (٩) الشبكة العصبية الاصطناعية مع تخطي طبقة الاتصالات حيث تتم معالجة نتيجة عقدة الإدخال (k) بواسطة عقدة الإخراج (l) وسوف يتم استعمال دالة التنشيط اللوجستية والتي تأخذ الصيغة الآتية:

$$\phi(z) = \frac{1}{1 + e^{-z}} \dots \dots \dots (23)$$

واعتماداً على دالة التنشيط (Activation function) سوف تكون مخرجات الشبكة العصبية قيم ثنائية محصورة بين (٠، ١) وكما ان هذه الشبكة تمثل الاتصال المباشر بين عقدة الإدخال وعقدة الإخراج وتسمى هذه الشبكة بالشبكة العصبية أحادية الطبقة المخفية مع تخطي طبقة الاتصالات. (Habibnia and Maasoumi, 2019: 6).

٣-٦. مقاييس دقة التشخيص والتصنيف

(Measures of Diagnostic and classification accuracy):

يعد التشخيص والتصنيف أساساً في الممارسة الطبية في حالة مرض معين فقد تكون الحاجة إلى التمييز بين الوضع الحالي للمرض وحالة غياب المرض، فقد يشمل التصنيف على أكثر من فئتين من المرضى يلزم معاجلتها بشكل منفصل، وغالباً ما تستعمل المؤشرات الحيوية (Biomarker) للتنبؤ بحالة مرض معين وذلك بالاعتماد على مجموعة من الأدوات الاحصائية لتقدير دقة التشخيص (Li and Fine, 2013: 383)، وسوف نستعمل أحد هذه المقاييس المهمة وهو مؤشر إعادة التصنيف الصافي (NRI) ومن ثم الحصول على أفضل حجم عينة يعطي أدق تشخيص للمؤشرات الحيوية الجديدة المضافة إلى المؤشرات الحيوية الأساسية بهدف التشخيص أو الفحص الطبي. (Pencina D'Agostino and Steyerberg, 2011: 12).

٣-٧. مؤشر إعادة التصنيف الصافي (Net reclassification index: NRI):
هو مقاييس تم اقتراحه من قبل (Pencina and others) في العام (٢٠٠٨) لتقدير التنبؤ بالعلامة الجديدة المضافة إلى الأنماذج الأساسية في حالة وجود فئات متعددة. (Pencina D'Agostino and Steyerberg, 2011: 12)
والصيغة التقديرية لمؤشر (NRI) هي على النحو الآتي:

$$\hat{S}_j = \sum_{m=1}^M \frac{\omega_m}{n_m} \sum_{i=1}^n I\{\hat{p}_{mi}(\mathcal{M}_2) = \max \hat{p}_i(\mathcal{M}_2), \hat{p}_{mi}(\mathcal{M}_1) \neq \max \hat{p}_i(\mathcal{M}_1), Y_i = m\} \dots \dots \dots (24)$$

اذ ان:

$$\hat{p}_i(\mathcal{M}_j) = (\hat{p}_{1i}(\mathcal{M}_j), \hat{p}_{2i}(\mathcal{M}_j), \dots, \hat{p}_{Mi}(\mathcal{M}_j))$$

ويعتبر أفضل مقاييس دقة تصنیف الذي يمتلك أكبر قيمة لمؤشر إعادة التصنيف الصافي (NRI).

الجانب التجريبي

٤. **مفهوم المحاكاة:** عبارة عن اسلوب يتضمن استخدام نموذج رياضي نظري و مشابهته بالأنموذج الحقيقي الذي يمثل المشكلة المدروسة أو كما عرفه (Naylor) انه اسلوب يستخدم للتحكم بالتجربة المراد دراستها ومعرفة صفاتها وخصائصها باستخدام القابلية الرياضية والمنطقية المتوفرة في الحاسبة الالكترونية أو كما عرفه آخرون اسلوب يتم من خلاله ايجاد نموذج بديل مماثل للنموذج الحقيقي من دون محاولة الحصول على النموذج الحقيقي نفسه. (بتال، ٢٠٠٨: ١٧٥)

٤-١. **مراحل بناء تجارب المحاكاة:** لتحليل البيانات باستعمال نماذج المحاكاة لا بد من تحديد مراحل بناء المحاكاة والتي تكون كما يلي:

٤-١-١. **توليد المتغيرات التوضيحية:** يتم توليد المتغيرات التوضيحية وفق التوزيع الطبيعي متعدد المتغيرات وفق الصيغة الآتية:

$$X_1, X_2, \dots, X_p \sim \text{MN}(\mathbf{0}, \Sigma)$$

اذ ان:

$$\Sigma = (\sigma_{ij})_{p \times p}, \quad \sigma_{ij} = 0.5^{|i-j|}, \quad 1 \leq i, j \leq p$$

اما بالنسبة لمتغير الاساسي يتم توليد من التوزيع الطبيعي القياسي وكما يأتي:

$$U \sim \text{N}(0, 1)$$

٤-٢. **تحديد المعلمات الأولية:** تم وصف القيم الافتراضية لتجارب المحاكاة وفقاً للجدول التالي:

الجدول (١): القيم الافتراضية للمعلمات الأولية

(Case one)	عدد المعلمات = p	(parameter)	$\beta = (\beta_0, \beta_u, \beta_1, \dots, \beta_p)$
1	$p = 75$	$\beta = (1.5, 1.5, 1.5, 0.5, 0.5, 0.5), \beta_j = 0, j = 7, \dots, 75$	
2	$p = 150$	$\beta = (1.5, 1.5, 1.5, 0.5, 0.5, 0.5), \beta_j = 0, j = 7, \dots, 150$	
3	$p = 350$	$\beta = (1.5, 1.5, 1.5, 0.5, 0.5, 0.5), \beta_j = 0, j = 7, \dots, 350$	
4	$p = 550$	$\beta = (1.5, 1.5, 1.5, 0.5, 0.5, 0.5), \beta_j = 0, j = 7, \dots, 550$	

٤-٣. **حساب متغير الاستجابة (Response Variable):** يتم حساب متغير الاستجابة الاحتمالي بالاعتماد على أنموذج الانحدار اللوجستي وعلى النحو الآتي:

$$y = \log \frac{p}{1-p} \dots \dots \dots (25)$$

$$\log \frac{p}{1-p} = \beta_0 + \beta_u U + \beta_1 X_1 + \dots + \beta_p X_p \dots \dots \dots (26)$$

٤-٤. **توليد مؤشر اعادة التصنيف الصافي:** بعد ان تم توليد الانموذج الاحتمالي في الفقرة (٤-٣) يتم توظيف قيمه الاحتمالية لحساب مؤشر اعادة التصنيف الصافي الذي يمثل الخلط بين نموذجين لوجيستيين أحدهما يتضمن المتغير الاساسي ومتغير الاستجابة الاحتمالي الذي يرمز له $\hat{p}(\mathcal{M}_1)$ والذي يأخذ الصيغة الآتية:

$$\hat{p}_i(\mathcal{M}_1) = \log \frac{p}{1-p} \quad \dots \dots \quad (27)$$

$$\log \frac{p}{1-p} = \beta_0 + \beta_u U \quad \dots \dots \quad (28)$$

والآخر يتضمن المتغير الاساسي بالإضافة الى المتغيرات التوضيحية ومتغير الاستجابة الذي يرمز له $\hat{p}(\mathcal{M}_{2j})$ الذي يأخذ الصيغة الآتية:

$$\hat{p}_i(\mathcal{M}_{2j}) = \log \frac{p}{1-p} \quad \dots \dots \quad (29)$$

$$\log \frac{p}{1-p} = \beta_0 + \beta_u U + \beta_1 X_1 + \dots + \beta_p X_p \quad \dots \dots \quad (30)$$

٤-٥. مناقشة تجارب المحاكاة: تم استعمال برنامج لغة البرمجة الإحصائية (R) حيث تم الاعتماد على معيار متوسط اكبر خطأ مطلق (MME) لمؤشر إعادة التصنيف الصافي (NRI) وكانت النتائج كالتالي:

الجدول (٢): مقارنة طرائق تحديد حجم العينة

NRI	$(P = 75)$		
	Methods	Bonf. Inq.	Berns. Inq.
$(\varepsilon = 0.1)$	n^*	1158	1655
	$MAE \hat{\epsilon}$	0.006074626	0.005433601
	(MME)	0.02074634	0.02327051
$(\varepsilon = 0.05)$	n^*	4632	6512
	$MAE \hat{\epsilon}$	0.002983268	0.001549753
	(MME)	0.01916864	0.00746901

الجدول (٣): مقارنة طرائق تحديد حجم العينة

NRI	$(P = 150)$		
	Methods	Bonf. Inq.	Berns. Inq.
$(\varepsilon = 0.1)$	n^*	1287	1798
	$MAE \hat{\epsilon}$	0.00872644	0.001934587
	(MME)	0.03047043	0.03465371
$(\varepsilon = 0.05)$	n^*	5149	7076
	$MAE \hat{\epsilon}$	0.00209387	0.001082647
	(MME)	0.01151935	0.006489755

الجدول (٤): مقارنة طرائق تحديد حجم العينة

NRI	$(P = 350)$		
	$\beta = (1.5, 1.5, 1.5, 0.5, 0.5, 0.5)$, $\beta_o = 0.5, \beta_u = 0.5$	Bonf. Inq.	Berns. Inq.
Methods			
$(\varepsilon = 0.1)$	n^*	1446	1973
	$MAE \hat{\varepsilon}$	0.006740543	0.002581634
	(MME)	0.0293	0.02326463
$(\varepsilon = 0.05)$	n^*	5786	7765
	$MAE \hat{\varepsilon}$	0.001964876	0.001197568
	(MME)	0.01093737	0.01370418

الجدول (٥): مقارنة طرائق تحديد حجم العينة

NRI	$(P = 550)$		
	$\beta = (1.5, 1.5, 1.5, 0.5, 0.5, 0.5)$, $\beta_o = 0.5, \beta_u = 0.5$	Bonf. Inq.	Berns. Inq.
Methods			
$(\varepsilon = 0.1)$	n^*	1532	2065
	$MAE \hat{\varepsilon}$	0.004702194	0.003010647
	(MME)	0.0508137	0.02102493
$(\varepsilon = 0.05)$	n^*	6127	8131
	$MAE \hat{\varepsilon}$	0.002059838	0.00110434
	(MME)	0.01076622	0.005256992

نتائج المحاكاة: نلاحظ من نتائج الجدول (٢) و (٣) عندما تكون عدد المعلمات ($P = 75$)، ($P = 150$) وقيم المعلمات الافتراضية اعلاه في الجدول بناءً على أفضل طريقة لتحديد حجم العينة عندما تكون قيمة هامش الخطأ ($\varepsilon = 0.1$) هي طريقة (Berns. Inq.) كونها تملك أعلى متوسط أكبر خطأ مطلق (MME)، وتعتبر طريقة (Bonf. Inq.) أفضل طريقة لتحديد حجم العينة عندما تكون قيمة هامش الخطأ ($\varepsilon = 0.05$) كونها تملك أعلى متوسط أكبر خطأ مطلق (MME).

ونلاحظ من نتائج الجدول رقم (٤) عندما تكون عدد المعلمات ($P = 350$) وقيم المعلمات الافتراضية اعلاه في الجدول بناءً على أفضل طريقة لتحديد حجم العينة عندما تكون قيمة هامش الخطأ ($\varepsilon = 0.1$) هي طريقة (Bonf. Inq.) كونها تملك أعلى متوسط أكبر خطأ مطلق (MME)، وتعتبر طريقة (Berns. Inq.) أفضل طريقة لتحديد حجم العينة عندما تكون قيمة هامش الخطأ ($\varepsilon = 0.05$) كونها تملك أعلى متوسط أكبر خطأ مطلق (MME). ونلاحظ من نتائج جدول رقم (٥) عندما تكون عدد المعلمات ($P = 550$) وقيم المعلمات الافتراضية اعلاه في الجدول بناءً على أفضل طريقة لتحديد حجم العينة عندما تكون قيمة هامش الخطأ ($\varepsilon = 0.1, \varepsilon = 0.05$) هي طريقة (Bonf. Inq.) كونها تملك أعلى متوسط أكبر خطأ مطلق (MME).

٥. الاستنتاجات والتوصيات

- ١-٥. الاستنتاجات: من خلال تجارب المحاكاة وما تم عرضه من نتائج استنتج الباحث ما يأتي:
١. من خلال الجانب التجريبي اثبتت النتائج أن أفضل طريقة تحديد حجم العينة هي طريقة (Bonf. Inq.) وعند قيمة هامش الخطأ ($\alpha = 0.1, 0.05$).
 ٢. هناك علاقة عكسية بين حجم العينة وهامش الخطأ اذ نلاحظ انه كلما قل هامش الخطأ زادت حجم العينة والعكس صحيح.
 ٣. تم الحصول على أفضل تقدير للمعلمات الافتراضية للبيانات ذات الابعاد العالية من خلال طريقة (Bonf. Inq.) لكونها تمتلك أعلى متوسط أكبر خطأ مطلق (MME).
- ٤-٥. التوصيات: على ضوء الاستنتاجات التي توصلنا اليها من خلال البحث يمكن ادراج التوصيات الآتية:
١. نوصي باستعمال طريقة (Berns. Inq.) لتحديد حجم العينة في حالة البيانات ذات الابعاد العالية وللدراسات وخصوصاً في الجانب الطبي وفي الطب التشخيصي.
 ٢. نوصي باستعمال مؤشر اعادة التصنيف الصافي (NRI) للدراسات الطبية الخاصة بالطب التشخيصي للأمراض.
 ٣. استعمال اساليب مختلفة غير الشبكات العصبية الاصطناعية لتقدير انموذج الانحدار اللوجستي الثنائي وبالتالي تقدير مؤشر اعادة التصنيف الصافي (NRI).

المصادر

أولاً. المصادر العربية:

١. بتال، احمد حسين، احمد، عصام كامل، خضير، البراء عبد الوهاب (٢٠٠٨)، استخدام المحاكاة في تدريس الانحدار الخطي البسيط، مجلة جامعة الانبار للعلوم الصرفة، العدد الثالث، المجلد الثاني.

ثانياً. المصادر الأجنبية:

1. Haykin, S, (2008), Redes Neurais.2nd Edition ,Principiose pratica.
2. Habibnia, A., & Maasoumi, E., (2019), Forecasting in Big Data Environments: an Adaptable and Automated Shrinkage Estimation of Neural Networks (AAShNet) p.1-26
3. Jiang B., & Li J., (2018), Sample size determination for high dimensional parameter estimation with application to biomarker identification. Computational Statistics and Data Analysis, 118, 54–65.
4. Li, J., Jiang, B., and Fine, J., (2013), Multicategory reclassification statistics for assessing improvements in diagnostic accuracy. Biostatistics 14, 382–394.
5. Kleinbaum, D. G., & Klein, M., (2002), Logistic Regression A Self-Learning Text. 3d ,Springer.
6. Pencina, M., D'Agostino, R., and Vasan, R., (2008), evaluating the added predictive ability of a new marker: from area under the ROC curve to reclassification and beyond. Statistics in Medicine 27, pp- 157–172
7. Pencina, M., D'Agostino, R., and Steyerberg, E., (2011), Extensions of net reclassification improvement calculations to measure usefulness of new biomarkers. Statistics in Medicine 30, pp- 11–21.

8. R. J. Simes, (1986), An Improved Bonferroni procedure for multiple tests of significance. *Biometrika*, Vol. 73, No. 3, pp. 751-754.
9. Ripley, B. D., (1996), *Pattern Recognition and Neural Network*. Cambridge University.
10. Shen, J., & Gao, S., (2008), A Solution to Separation and Multicollinearity in Multiple Logistic Regression. *Journal of Data Science : JDS*, 6(4), PP-515-531.
11. Scott, M., (2002), *Applied Logistic Regression Analysis*. sage.
12. Van Der Vaart, A. and Wellner, J. A., (1996), *Weak Convergence and Empirical Processes*. New York: Springer.
13. Van der Laan, M., and Bryan, J., (2001), Gene expression analysis with the parametric bootstrap. *Biostatistics* 2, PP- 445-461.