

## Compare the classification of the quadratic function and the logistic function

### مقارنة تصنيف الدالة التربيعية والدالة اللوجستية

مريم مهدي عناد الخزعلي

أ.م. د. شروق عبد الرضا سعيد السباح

جامعة كربلاء/كلية الإدارة والاقتصاد/ قسم الإحصاء  
(بحث مستل من رسالة ماجستير)

#### الخلاصة:-

تم تطبيق دالة التمييز لتشخيص ثلاث أنواع من الأورام السرطانية وهي سرطان الثدي وسرطان العظم وسرطان الرئة لكثرة شيوعها في مجتمعنا حالياً، ولغرض دراسة هذا الموضوع تم تسجيل قيم المشاهدات لخمس متغيرات وهي (الجنس، العمر، مهنة المريض، حالة خروج المريض، فترة بقاء المريض في المستشفى) من عينة عشوائية بسيطة بحجم (270 مريض) درست في نموذجين وهي دالة التمييز التربيعية ودالة التمييز اللوجستية، واستخرجت النتائج باستعمال البرنامج الإحصائي Stata، وتم استعمال معيار خطأ التصنيف كمييار للمقارنة بين النماذج، بهدف الوصول إلى أفضل نموذج باقل خطأ تصنيفي ممكن، وتم الوصول إلى ان الدالة اللوجستية ذات تفوق على الدالة التربيعية من حيث قلة نسبة التصنيف الخاطئ.

#### Abstract:-

The function of discrimination has been applied to diagnose three types of tumors, namely breast cancer, bone cancer and lung cancer, for their prevalence in our society. For the purpose of this study, the values of observations were recorded for five variables (sex, age, occupation, prognosis and duration of hospitalization) the random sample of size (270 patients) studied in two models: Quadratic discrimination function and Logistic discrimination function. The results were extracted using the statistical program Stata. With a view to reaching Li best model the lowest rank possible error, was reached that Logistic discrimination functionis superior to Quadratic discrimination function in terms of the lack of proportion of the wrong classification.

#### Introduction and Problem

#### المقدمة ومشكلة البحث

يعد التمييز بين المشاهدات من الأساليب الشائعة الاستعمال وذلك لكثرة الظواهر التطبيقية التي يمكن أن يتم تحليلها من خلال أسلوب التمييز. إذ إن دالة التمييز الخطية المسندة على تركيب خطي للمتغيرات لكي تكون مثلى يجب أن تنتج اصغر احتمال لخطأ التصنيف علماً بان هناك افتراضات يجب بتوفرها حول البيانات المستخدمة في التحليل وهي أن يكون متجه المتغيرات التوضيحية ذا توزيع طبيعي متعدد المتغيرات ومصفوفات تباينات مشتركة متساوية ومتجهات متوسطات مختلفة في كل مجموعة من المجاميع؛ ولكن أحياناً نواجه اختلال في بعض الفرضيات وبالتالي تفقد دالة التمييز الخطية بعض خواصها المثلى. فعند عدم تحقق ثبات مصفوفة التباين والتباين المشترك يكون من الضروري استعمال دالة التمييز التربيعية (Quadratic Discriminate)، فضلاً عن ذلك ان دالة التمييز الخطية تفقد خواصها المثلى عندما تكون جميع او بعض المتغيرات الاضافية ذات طبيعة ثنائية او مصنفة فتكون دالة التمييز اللوجستية (Logistic Discrimination) البديل الامثل للتمييز الطبيعي. أما عملية التصنيف (Classification) هي العملية اللاحقة بعد تكوين الدالة المميزة حيث يتم الاعتماد على هذه الدالة بالتنبؤ وتصنيف المفردة الجديدة لإحدى المجموعات قيد الدراسة بأقل خطأ تصنيف ممكن.

**Search Hypothesis**

$H_0$ : عدم وجود مشاكل في النموذج الاحصائي

$H_1$ : وجود مشاكل في النموذج الاحصائي

**هدف البحث**

يهدف البحث إلى الوصول إلى أفضل تصنيف لبعض أنواع الأورام السرطانية على أساس معيار احتمال اقل خطأ تصنيفي ممكن باستعمال دالة التمييز التربيعية ودالة التمييز اللوجستية.

**Research Aim**

**الجانب النظري**

**1. دالة التمييز التربيعية**

**Quadratic Discriminant Function**

تعد هذه الطريقة من الطرق المعلمية إذ يستعمل هذا النوع من الدوال في حالة عدم تساوي مصفوفة التباين والتباين المشترك للمجموعات ويكون المجتمع قيد الدراسة ذا توزيع طبيعي متعدد المتغيرات وبمتجهات متوسطات مختلفة. وعليه فان مقياس التمايز  $V$  سيكون كالآتي [10]:

$$V = \frac{f_1(x_i)}{f_2(x_j)} \neq 1 \quad (1)$$

إذا كانت  $V > 1$  المشاهدة  $x$  تعود إلى المجموعة الأولى و إذا كانت  $V < 1$  فان المشاهدة  $x$  تعود إلى المجموعة الثانية وعشوائياً عدا ذلك.

وبأخذ اللوغاريتم للطرفين نحصل على الآتي:

$$G = \ln V = \ln f_1(x_i) - \ln f_2(x_j) = 0 \quad (2)$$

وفي حالة وجود  $P$  من المتغيرات  $(X_1, X_2, \dots, X_P)$  لكل مجموعة من المجموعتين فان مقياس التمييز يكون:

$$G = \ln f_1(x_1, \dots, x_P) - \ln f_2(x_1, \dots, x_P) \quad (3)$$

وفي حالة إن المعالم الدالة الاحتمالية غير معلومة فسوف تكون دالة التمييز التربيعية التقديرية  $g$  والتي قدرت معالمها بطريقة الإمكان الأعظم كالآتي [13]:

$$g = \hat{G} = \frac{1}{2} \log \frac{|S_1|}{|S_2|} - \frac{1}{2} (\bar{X}'_1 S_1^{-1} \bar{X}_1 - \bar{X}'_2 S_2^{-1} \bar{X}_2) + \bar{X}' (S_1^{-1} \bar{X}_1 - S_2^{-1} \bar{X}_2) - \frac{1}{2} \bar{X}' (S_1^{-1} - S_2^{-1}) \bar{X} \quad (4)$$

وسيتم تصنيف المشاهدات الجديدة بأنها تعود إلى المجموعة الأولى إذا كانت  $g > 0$  وإلى المجموعة الثانية إذا كانت  $g < 0$  وعشوائياً عدا ذلك.

**2. اشتقاق دالة التمييز التربيعية [16], [15]**

تكون الفرضية في حالة وجود مجموعتين من مجتمعين فرضية خاصة تمثل انتماء كل مشاهدة من المشاهدات إلى إحدى المجموعتين وهي:

$$H_0: \underline{X} \in f_2(x)$$

$$H_1: \underline{X} \in f_1(x)$$

وباستعمال أسلوب بيز لتقليل الخطورة التي تمثل نسبة الإمكان الأعظم والاحتمالات الأولية  $P_1$  و  $P_2$  وكلف التصنيف الخاطئ  $c_{21}$  و  $c_{12}$  تكون:

$$\frac{f_1(x)}{f_2(x)} < \frac{c_{12}P_2}{c_{21}P_1} \quad (5)$$

وعند امتلاك كل من  $f_1$  و  $f_2$  توزيعاً طبيعياً و اخذ اللوغاريتم للطرفين ينتج:

$$\ln \frac{|\Sigma_1|}{|\Sigma_2|} + (\underline{x}_{11} - \underline{\mu}_1)' \Sigma_1^{-1} (\underline{x}_{1i} - \underline{\mu}_1) - (\underline{x}_{2i} - \underline{\mu}_2)' \Sigma_2^{-1} (\underline{x}_{2i} - \underline{\mu}_2) < -2 \ln \frac{c_{12}P_2}{c_{21}P_1} \quad (6)$$

وبعد التبسيط تصبح دالة التمييز كما يأتي :

$$h(x) = \underline{x}'A\underline{x}b'x + c \langle \rangle T \quad (7)$$

إذ إن:

$$A = \Sigma_1^{-1} - \Sigma_2^{-1}$$

$$b = 2 \left( \Sigma_2^{-1} \underline{\mu}_2 - \Sigma_1^{-1} \underline{\mu}_1 \right)$$

$$c = \left( \underline{\mu}'_1 \Sigma_1^{-1} \underline{\mu}_1 - \underline{\mu}'_2 \Sigma_2^{-1} \underline{\mu}_2 + \ln \frac{|\Sigma_1|}{|\Sigma_2|} \right)$$

$$T = -2 \ln \frac{c_{12} P_2}{c_{21} P_1}$$

إذ تصنف المشاهدة  $x$  على أنها تنتمي للمجتمع  $f_1$  إذا كانت  $h(\underline{x}) < T$  وتنتمي إلى المجتمع  $f_2$  إذا كانت  $h(\underline{x}) > T$  وتعرف بدالة التمييز التربيعية.

أما عند تساوي مصفوفة التباين والتباين المشترك  $\Sigma_i$  تصبح الدالة بالشكل الآتي:

$$h^*(x) = b'^* x + c^* \langle \rangle T(8)$$

إذ أن:

$$b^* = 2 \Sigma (\underline{\mu}_2 - \underline{\mu}_1)$$

$$c^* = \underline{\mu}'_1 \Sigma^{-1} \underline{\mu}_1 - \underline{\mu}'_2 \Sigma^{-1} \underline{\mu}_2$$

وهي دالة التمييز الخطية.

الرمز  $\langle \rangle$  يعني تصنيف المشاهدة  $x$  على أنها تنتمي إلى المجتمع المميز  $f_1$  إذا كانت قيمتها اقل من  $T$  وتنتمي إلى المجتمع  $f_2$  إذا كانت قيمتها اكبر من  $T$ .

وعليه تكون كلف التصنيف الخاطئ تساوي دائماً واحد ، وان خطورة بيز تتحول إلى احتمال خطأ التصنيف لذلك فان الحل الذي يقلل من خطورة بيز يقلل أيضاً من احتمال خطأ التصنيف ولهذا فإن :

$$P_1 f_1(\underline{x}) < P_2 f_2(\underline{x}) \quad (9)$$

وعند التعميم واخذ مجموعات متعددة تكون الفرضية:

$$H_0: \underline{x} \in w_i \quad i = 1, \dots, g$$

وفي حالة الكلف تساوي واحد يكون الحل باختبار المجموعة التي تحقق ما يأتي كمجموعة تنتمي إليها المفردات  $\underline{x}$ :

$$P_k f_k > P_i f_i \quad (10)$$

$$\text{for } i \neq k, i = 1, \dots, g$$

وبعبارة أخرى فان لدينا  $g$  من الدوال التمييزية  $h_i(\underline{x})$  والحل هو اختبار الدالة التي تعطي أكبر قيمة ، وكما يأتي:

$$h_k(\underline{x}) = \max h_i(\underline{x}) \quad (11)$$

إذ أن:

$$h_i(\underline{x}) = P_i f_i(\underline{x})$$

وتوجد نظرية بهذا الخصوص تنص على إن دالة متزايدة رتيبة  $h_i(\underline{x})$  يمكن أن تحل بدلا من  $h_i(\underline{x})$  وفي حالة توزيع البيانات طبيعياً يكون [15]:

$$h_i(\underline{x}) = P_i N(\underline{x}_i \underline{\mu}_1, \Sigma_i) \quad (12)$$

إذ إن :

$N(\underline{x}_i \underline{\mu}_1, \Sigma_i)$  : هي دالة التوزيع الطبيعي متعدد المتغيرات

وفي حالة تحقق التوزيع الطبيعي فان الصيغة  $h_i(\underline{x})$  تكون :

$$h_i(\underline{x}) = - \left( \underline{x} - \underline{\mu}_i \right)' \Sigma_i^{-1} \left( \underline{x} - \underline{\mu}_i \right) - \ln |\Sigma_i| + 2 \ln P_i \quad (13)$$

$$i = 1, \dots, g$$

وهي دالة التمييز التربيعية المتعددة (Multiple Quadratic Discriminant Function) وعند تساوي مصفوفة التباين والتباين المشترك تصبح الدالة بالشكل الآتي:

$$h_i^*(\underline{x}) = 2\mu_i'\Sigma^{-1}\underline{x} - \mu_i'\Sigma^{-1}\mu_i + 2\ln P_i, i = 1, \dots, g \quad (14)$$

وهي دالة التمييز الخطية .

### 3. تقدير معالم النموذج [10],[12]

يتم تقدير  $\Sigma_i$  و  $\mu_i$  بطريقة الإمكان الأعظم لوجود خاصية الثبات إذ إن :

$$\hat{\mu}_i = \bar{x}_i = \frac{1}{n_i} \sum_{j=i}^{n_i} x_{ij} \quad (15)$$

$$\hat{\Sigma}_i = S_i^2 = \frac{1}{n_i} \sum_{j=i}^{n_i} (x_{ij} - \bar{x}_j)(x_{ij} - \bar{x}_j)^2 \quad (16)$$

ويمكن استبدال تقدير  $\Sigma_i$  المتحيز بتقدير  $\Sigma_1^{-1}$  الغير متحيز المشتق من توزيع (Wishert) وهو :

$$\hat{\Sigma}_1^{-1} = \left( \frac{ni - P - 3}{ni} \right) S_i^{-1} \quad (17)$$

إذ إن  $S_i$  هي تقدير الإمكان الأعظم لـ  $\Sigma_1^{-1}$  وكذلك بدلا من  $\Sigma_i$  يمكن استعمال التقدير غير متحيز لمسافة مهالانوس .

$$\hat{D}_{ij}^2 = \frac{n_i + n_j - p - 3}{n_i + n_j - p - 2} (\bar{x}_j - \bar{x}_i)' S^{-1} (\bar{x}_j - \bar{x}_i) - \frac{p(n_i - n_j)}{n_i n_j} \quad (18)$$

وإذا كانت  $\hat{D}_{ij}^2$  سالبة فتؤخذ صفراً .

أما بالنسبة لتقدير الاحتمالات الأولية فتقديرها يكون:

$$\hat{P}_i = p_i = \frac{n_i}{n} \quad (19)$$

إذ أن:

$n$  : حجم العينة الكلي

$n_i$  : حجم العينة للمجموعة  $i$ .

### Logistic discriminant Function

### 4. دالة التمييز اللوجستي

يعد التمييز اللوجستي من بين الكثير من الأساليب المقترحة في التمييز الإحصائي للطرق المعلمية المستندة على التوزيع الطبيعي متعدد المتغيرات ، والطرق اللامعلمية حرة التوزيع كطريقة الجار الأقرب (Nearest Neighbor) لهذا السبب غالباً ما تدعى هذه الطريقة بالمعلمية الجزئية أو شبه معلمية . وفقاً لذلك فإن النموذج اللوجستي احد الأدوات الأكثر جاذبية واستعمالاً في حل مسائل التمايز والانحدار اللذان يسيران نحو إحدى أو كلا الاتجاهين الآتيين :

1. تلخيص ووصف الفروقات بين المجتمعات.

2. تحديد موقع مفردات جديدة إلى مجاميعها أي التنبؤ بقيم متغير الاستجابة لمفردات جديدة.

ان هذه الطريقة من الطرائق المهمة وذلك لأنها تعالج مشاكل التمييز في حالة إن تكون المتغيرات (المقاييس) منقطعة أو مستمرة أو مختلطة فضلاً عن سهولتها في الاستخدام إذ أن عدد حدودها يساوي عدد حدود الدالة الخطية. وهناك نوعان من النموذج اللوجستي وهي (1) النموذج اللوجستي ثنائي الاستجابة (2) النموذج اللوجستي متعدد الاستجابة وسيكون تركيزنا على النوع الثاني .

### مميزات النموذج اللوجستي<sup>[5]</sup>

1. يتطلب القليل من الافتراضات حول التوزيع، فهو لا يفترض وجود علاقة خطية بين المتغير التابع والمتغيرات التوضيحية.
2. سهولة استخدامه، فعندما يتم تقدير معالم النموذج فإن تحديد موقع مفردة جديدة يتطلب حساب دالة خطية واحدة فقط.
3. تكون المعالم المقدره من النموذج ذات تفسير مباشر كدالة خطية للوغاريتم نسبة الارجحية .
4. المتغير التابع يجب ان يكون ثنائي التفرع بحيث يحتوي على فئتين (مصابين او غير مصابين على سبيل المثال)
5. يجب ان تكون الفئات محددة وشاملة بحيث ان كل مفردة تنتمي إلى فئة واحدة فقط .
6. النموذج اللوجستي لا يشترط ان تكون المتغيرات التوضيحية من النوع المستمر ولا تتبع التوزيع الطبيعي ولا ان تكون العلاقة بين المتغير التابع والمتغيرات التوضيحية خطية ولا يفترض تساوي التباين ضمن كل فئة . وهذا يجعل النموذج اللوجستي اكثر مرونة من بقية نماذج التنبؤ والتصنيف .
7. يجب ان يكون حجم العينة المستخدم في النموذج اللوجستي اكبر من حجم العينة المستخدم في النموذج الخطي، لأن معاملات نموذج اللوجستي يتم تقديرها باستخدام طريقة الامكان الاعظم وهي طريقة تحتاج إلى عينة كبيرة الحجم نسبياً .

### 5. نموذج اللوجستي ثنائي الاستجابة [8]، [14]

يُطبق هذا النموذج عندما يأخذ المتغير العشوائي  $Y$  قيمتين فقط قد تدل على نمط الاستجابة لمؤثر ما (كالأدوية والمنشطات والسموم). في هذه الحالة تأخذ  $Y$  القيمتين 0 و 1 تشير إلى عدم حدوث وحدث الاستجابة على التوالي. كذلك يستعمل في مسائل التشخيص الطبي والتميز بين أنواع الأمراض. فالأورام السرطانية مثلا قد تكون خبيثة أو حميدة. تعد هذه الطريقة من الطرائق المهمة وذلك لأنها تعالج مشاكل التمييز في حالة ان تكون المتغيرات مستمرة أو متقطعة أو مختلطة فضلا عن سهولتها في الاستعمال إذ إن عدد حدودها يساوي عدد حدود الدالة الخطية، يتكون نموذج اللوجستي من متغير الاستجابة  $Y_i$  يتأثر بمجموعة من المتغيرات الموضحة ( $X_i$ 's) (*Explanatory variables*) على وفق علاقة تحتوي على مجموعة من المعالم ( $b_i$ 's) إذ أن متغير الاستجابة  $Y_i$  في التجارب الحياتية ثنائية الاستجابة الذي يعد نموذج اللوجستي احدها يتوزع حسب توزيع برنولي أي ان له مستويان هما الصفر والواحد.

$$Y_i \sim b(1, P_i)$$

إذ إن:

$P_i$ : هو احتمال الاستجابة عندما ( $Y = 1$ )

$1 - P_i$ : احتمال عدم الاستجابة عندما ( $Y = 0$ )

يعتمد  $P_i$  على مجموعة من المتغيرات التوضيحية ( $X_i$ 's) على وفق الصيغة الآتية التي يطلق عليه انموذج دالة الاستجابة اللوجستية:

$$P_i = \frac{e^{-x_i B}}{1 + e^{-x_i B}} \quad (20)$$

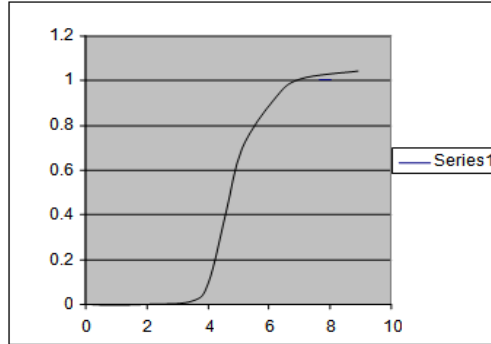
$$1 - P_i = \frac{1}{1 + e^{-x_i B}} \quad (21)$$

إذ ان:

$B$ : معالم مجهولة يتم تقديرها.

$x_i$ : تمثل المتغيرات التوضيحية.

وان الصيغة (20) تسمى دالة الاستجابة اللوجستية ولها مخطط في الشكل (1):



الشكل (1) يوضح العلاقة بين احتمال الاستجابة ( $P_i$ ) والمتغير الموضح ( $x_i$ )

من الشكل (1) نلاحظ عدم خطية العلاقة بين المتغير الموضح  $x_i$  واحتمال الاستجابة  $P_i$ ، ولغرض اجراء التحويل الخطي للدالة اللوجستية تمكن الباحث (*Berkson*) عام (1944) من تحويل العلاقة بين المتغيرات ( $x_i$ ) ومتغير الاستجابة ( $P_i$ ) إلى علاقة خطية وذلك برسم  $\log P$  بدلاً من  $\text{Probit } P$  مقابل ( $X$ ) كالاتي [14].

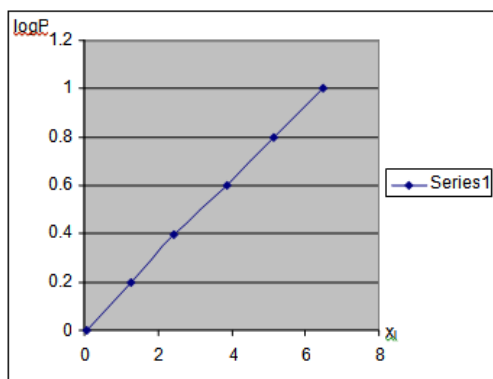
$$L_i = \ln\left(\frac{P_i}{1 - P_i}\right) = \log P_i = x_i b \quad (22)$$

إذ إن:

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

$L_i$ : يمثل تحويل ( $\log$ ) للاحتمال ( $P_i$ ).

والصيغة (22) تبين التحويل الخطي ( $\log$ ) كما في الشكل (2):



الشكل (2) يوضح العلاقة الخطية بين  $\log P_i$  و  $x_i$

$$L_i = \ln\left(\frac{P_i}{1-P_i}\right) = b_0 + b_1x_1 + \dots + b_px_p \quad (23)$$

اذ تكون  $\underline{b}$  موجه عمودي من رتبة  $[(P+1) \times 1]$  للمعالم المطلوب تقديرها

$$\underline{b} = \begin{pmatrix} b_0 \\ b_1 \\ \vdots \\ b_p \end{pmatrix}$$

$\underline{X}$  مصفوفة من رتبة  $[g \times (p+1)]$  للمتغيرات التوضيحية تكون :

$$\underline{X} = \begin{pmatrix} 1 & X_{11} & \dots & X_{1p} \\ 1 & X_{21} & \dots & X_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & X_{g1} & \dots & X_{gp} \end{pmatrix}$$

وبما ان التوزيع هنا هو برنولي ولحالتين فقط عندئذ يكون <sup>[9]</sup>:

$$\underline{L}_i = \underline{P}_i = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$$

$$\underline{L}_i^* = \ln\left(\frac{P_i + \frac{1}{2n}}{(1-P_i) + \frac{1}{2n}}\right) \quad (24)$$

اذ استعمل هنا معامل تصحيح الاستمرارية  $\frac{1}{2n}$  حيث اعطى نتائج اكثر دقة ويخطأ اقل <sup>[2]</sup>.

ونستعمل بعد ذلك طريقة المربعات الصغرى العامة (GLS) *Generalized Least Squares* التقديرية توهي:

$$\underline{b}^* = (X'X)^{-1}X'\underline{L}_i^* \quad (25)$$

$$\hat{L}_i = \hat{b}_0 + \hat{b}_1X_1 + \dots + \hat{b}_pX_p \quad (26)$$

وسوف يتم تصنيف المشاهدات الجديدة بأنها تعود إلى المجموعة الأولى إذا كانت  $\hat{L}_i$  اكبر من الصفر وإلى المجموعة الثانية إذا كانت  $\hat{L}_i$  اصغر من الصفر وعشوائياً عدا ذلك.

### 6. نموذج اللوجستي متعدد الاستجابة

يعد نموذج اللوجستي متعدد الاستجابة امتداداً طبيعياً لنموذج ثنائي الاستجابة . وهي الحالة التي يكون فيها المفردة أو المشاهدة محددة بوحدة من  $g$  من القيم الممكنة تدعى بفئات أو مجاميع الاستجابة . فقد تكون ذات طبيعة ترتيبية مثل درجة خطورة أو شدة مرض معين . كما يمكن أن تكون مجاميع الاستجابة ذات طبيعة اسمية كأصناف الدم مثلاً أو نوع واسطة النقل المفضلة في السفر (سيارة ، طائرة ، قطار ، .....).

نفرض وجود مجتمع كبير يحتوي على  $g$  من المجاميع المختلفة  $(\pi_1, \dots, \pi_g)$  وسحبت عينة عشوائية منه . فأن عدد المفردات الداخلة ضمن العينة وتنتمي إلى المجموعة  $\pi_g$  يتبع توزيع متعدد الحدود *Multinomial Distribution* وان احتمال وقوع  $Y$  من المشاهدات ضمن العينة يكون كالآتي<sup>[11]</sup>:

$$P(Y_1 = y_1, \dots, Y_g = y_g) = \binom{n}{y} P_1^{y_1}, \dots, P_g^{y_g} \quad (27)$$

إذ إن:

$(P_1, \dots, P_g)$ : تمثل نسبة المجاميع في المجتمع

$$\binom{n}{y} = \frac{n!}{y_1! y_2! \dots y_p!} \quad (28)$$

نستنتج من ذلك إن التوزيع متعدد الحدود يتولد كنتيجة لأسلوب المعاينة. فإذا رمزنا لاحتمال حدوث الاستجابة  $j$  بـ  $P_j(j = 1, \dots, g)$  فعند قيم المشاهدات  $i$  ذات متجه الخصائص  $x_i = (x_{i1}, \dots, x_{ip})$  يكون النموذج بالشكل الآتي:

$$P_j(x_i) = \frac{\exp(b'_j x_i)}{\sum_{j=1}^g \exp(b'_j x_i)}, j = 1, \dots, g - 1 \quad (29)$$

ويكون  $P_g = 0$  إذ يمثل  $g$  مجموعة الأساس  $b_j$  يمثل متجه المعالم الخاصة بالمجموعة  $j$  إذ يوجد  $g-1$  من المتجهات المختلفة ضمن النموذج . وبذلك فإن النموذج يحتوي على  $(p+1)(g-1)$  من المعالم المجهولة ويمكن كتابة النموذج بصيغة لوغاريتم نسبة الإمكان الأعظم كما يلي :

$$\log \left( \frac{P_j(x_i)}{P_g(x_i)} \right) = b'_j x_i, j = 1, \dots, g - 1 \quad (30)$$

### 7. تقدير معلمات النموذج اللوجستي<sup>[3],[11]</sup>

لاحظ *Anderson (1972)* إن تطبيق النموذج اللوجستي يخضع للمعادلة الآتية :

$$\log \left( \frac{P_j(x_i)}{P_g(x_i)} \right) = b_0 + b'x \quad (31)$$

وهو تركيب خطي للمتغير  $x$  يناسب المتغيرات الثنائية والثلاثية عندما  $x$  يأخذ القيم  $(0, 1, 2)$  ، بافتراض إن نسبة احتمال اللوغاريتمي هو تركيب خطي للمتغير  $x$  .

من أجل تطبيق النموذج اللوجستي يجب تحويل المتغير  $x$  إلى متغيرين ثنائيين  $x_1$  و  $x_2$  ، إذ إن  $x = 0$  إذا فقط إذا كان  $x_1 = 0$  و  $x_2 = 0$ ؛  $x = 1$  إذا فقط إذا كان  $x_1 = 1$  و  $x_2 = 0$ ؛  $x = 2$  إذا فقط إذا  $x_1 = 0$  و  $x_2 = 1$  .  
 على نحو مماثل، يمكن استبدال كل متغير مع أكثر من  $2 \geq r$  مستويات من المتغيرات الثنائية  $r - 1$  .  
 وان النموذج اللوجستي في تحليل التمايز يكون :

$$T_1(x) = \frac{\exp(b_0 + b'x)}{1 + \exp(b_0 + b'x)} \quad (32)$$

إذ إن  $T_1(x)$  تمثل الاحتمالات اللاحقة لعضوية المجموعة من خلال نمذجة العلاقة بين متغير الاستجابة الفئوية الذي نرسم له بالرمز  $G$  وبين متجه المتغيرات التوضيحية  $X$  .

علماً أن متغير الاستجابة  $G$  هو متجه من المجموعات أو الفئات المختلفة  $(G_1, \dots, G_g)$  .

وعليه يكون نموذج لوغاريتمي خطي للاحتمالات اللاحقة للمتغير مستقل مع المتغير  $x$  مقارنة مع متغير أساس  $x_0$  هو :

$$\log \left[ \frac{T_1(x)}{T_1(x_0)} \right] = b'(x - x_0) \quad (33)$$

في كثير من الأحيان المرض قيد الدراسة يكون نادراً إلى حد ما ، بحيث  $T_2(x)$  و  $T_2(x_0)$  قريبة من الواحد ، وبالتالي نسبة الفرق تقريبا تكون :

$$T_1(x)/T_1(x_0) \quad (34)$$

ومن خلال التحويل اللوجستي نحصل على :

$$\begin{aligned} \text{logit}[T_1(x)] &= \log[T_1(x)/T_2(x)] \\ &= b_0 + b'x \end{aligned} \quad (35)$$

ولربط  $x$  بالاحتمال اللاحق  $T_1(x)$  يكون التمييز الاحتمالي بالصيغة الآتية :

$$\text{probit}[T_1(x)] = \Phi^{-1}[T_1(x)] \quad (36)$$

إذ إن  $\Phi$  هي دالة التوزيع الطبيعي القياسي .

عملياً يمكن استعمال المعادلة الآتية للنماذج اللوجستية والاحتمالية وهي :

$$\text{logit}[T_1(x)] \cong c \Phi^{-1}[T_1(x)] \quad (37)$$

إذ إن  $c = \sqrt{8/\pi} = 1.6$  .

### Classification

### 8. عملية التصنيف

إن عملية التصنيف (Classification) هي العملية اللاحقة بعد تكوين دالة التمييز حيث يتم الاعتماد على هذه الدالة بالتنبؤ وتصنيف المفردة الجديدة لإحدى المجموعات قيد الدراسة بأقل خطأ تصنيف ممكن، ويشترط تساوي التباينات للمجموعات قيد البحث، وهناك تمييز خطي في حالة مجموعتين، وتمييز خطي في حالة أكثر من مجموعتين، أما التمييز غير الخطي فيستخدم في حالة عدم تساوي التباينات<sup>[1]</sup>.

### احتمال خطأ التصنيف

إن خطأ التصنيف هو عامل مهم لإثبات كفاءة الدالة التمييزية، أي أن الدالة التمييزية التي تعطي أقل خطأ تصنيف هي الدالة الأكثر كفاءة وتكون الأفضل من بين دوال التمييز. وهناك نوعين من احتمال خطأ التصنيف<sup>[7]</sup>:  
أولاً: احتمال خطأ التصنيف  $P_{ij}$  وهو احتمال تصنيف المشاهدة إلى المجموعة  $j$  وهي تعود إلى المجموعة  $i$  .

$$P_{ij} = \Phi\left(-\frac{\delta}{2}\right) \quad (38)$$

ثانياً: احتمال خطأ التصنيف  $P_{ji}$  وهو احتمال تصنيف المشاهدة إلى المجموعة  $i$  وهي تعود إلى المجموعة  $j$ .

$$P_{ji} = \Phi\left(-\frac{\delta}{2}\right) \quad (39)$$

إذ أن:

$\Phi$  : تمثل دالة التوزيع الطبيعي القياسي .

$\delta$  : تمثل جذر مقياس مسافة مهالانوبس  $D^2$  .

إذن :

$$\delta^2 = D^2 = (\bar{X}_1 - \bar{X}_2)' S^{-1} (\bar{X}_1 - \bar{X}_2)$$

وسيكون احتمال خطأ التصنيف كالآتي:

$$P_{ij} = P_{ji} = \Phi\left(-\frac{D}{2}\right) \quad (40)$$

إذ إن  $D$  هو جذر مقياس مهالانوبس .

### الجانب التطبيقي

تم جمع البيانات من خلال اخذ عينة عشوائية بسيطة حجمها 270 من طيلات المرضى الراقدين في مستشفى (مدينة مرجان الطبية) في محافظة بابل لإمراض السرطان للتمييز بين اشخاص مصابين بأورام الثدي والعمود الفقري (سرطان العظم) والجهاز التنفسي (سرطان الرئة) وذلك لكثرة المصابين بهذه الانواع خلال عام 2016، ولغرض تحليل البيانات تم استعمال البرنامج الاحصائي Stata وكذلك الحزمة الاحصائية SPSS، وتم دراسة المتغيرات التالية لكل مجموعة: (الجنس، العمر، مهنة المريض، تشخيص المرض، حالة خروج المريض، فترة بقاء المريض في المستشفى).

## Definition of variables

## تعريف المتغيرات

اعتمدنا في تكوين دالة التمييز على عدد من المتغيرات تم جمعها عن كل مشاهدة (مريض) من مشاهدات العينة ، وحيث إن جزء من المتغيرات هي متغيرات متقطعة وأخرى متصلة .

أ- المتغير المتعمد (Y) والذي يمثل نوع أو تشخيص المرض:

سرطان الثدي = 1 سرطان العظم = 2 سرطان الرئة = 3

ب- المتغيرات التوضيحية تمثل ما يلي:

1-  $X_1$  = الجنس: شملت العينة كلا الجنسين وكانت نسبتهم إلى مجموع العينة متفاوتة ويمثل متغير ثنائي (ذكر=1 ، أنثى=2)

2-  $X_2$  = العمر: تتراوح أعمار المرضى المشمولين ضمن العينة بين (4 – 95) سنة.

3-  $X_3$  = مهنة المريض: أعطينا رموز لكل مهنة وذلك لتسهيل تحليل البيانات

(طفل=1، ربة بيت أو كاسب=2 ، طالب =3 ، عاجز=4 ، متقاعد=5 ، موظف =6)

4-  $X_4$  = حالة خروج المريض: (وفاة=1 ، إحالة=2 ، خرج على مسؤوليته=3 ، تحسين=4)

5-  $X_5$  = فترة بقاء المريض في المستشفى.

## اختبار التوزيع الطبيعي لكل مجتمع

نختبر البيانات لمعرفة ما إذا كانت المتغيرات التوضيحية للمجموعات الثلاث لإمراض السرطان تتوزع طبيعياً أم لا باستخدام (*Kolmogorov-Smirnov*) وبموجب البرنامج الإحصائي SPSS وحسب الفرضية التالية:

$H_0$ : البيانات تتبع التوزيع الطبيعي

$H_1$ : البيانات لا تتبع التوزيع الطبيعي

الجدول (1) نتائج اختبار البيانات للتوزيع الطبيعي

Variables	Kolmogorov-Smirnov	
	statistic	Sig.
الجنس ( $X_1$ )	0.40	.000
العمر ( $X_2$ )	0.14	.000
المهنة ( $X_3$ )	0.34	.000
حالة خروج المريض ( $X_4$ )	0.48	.000
فترة بقاء المريض ( $X_5$ )	0.17	.000

في الجدول (1) أظهرت نتائج قيم اختبار (*Kolmogorov-Smirnov*) إن مستوى المعنوية لجميع المتغيرات أقل من 0.05 المستوى المعتمد في هذه الدراسة وعليه نرفض فرضية العدم القائلة إن البيانات تتبع التوزيع الطبيعي. ولكون حجم البيانات 270 مشاهدة يمكن إن نعد إن البيانات تقترب من التوزيع الطبيعي حسب نظرية (الغاية المركزية).

## حساب احتمال التصنيف الصحيح

### The probability of correct classification

إن عملية التصنيف قد تؤدي إلى الوقوع فيما يعرف بخطأ التصنيف (*misclassification*) وهو احتمال تصنيف مفردة معينة إلى المجموعة الأولى بينما هي في الحقيقة تعود للمجموعة الثانية أو الثالثة وبالعكس وهناك نوعان من الخطأ الأول يسمى بنسبة الخطأ الظاهري والثاني نسبة الخطأ الحقيقي. وفي هذه الدراسة تم التركيز على احتساب النوع الأول من الخطأ لجميع مشاهدات العينات الثلاث والجدول التالي يمثل التصنيف بحسب دالة التمييز التربيعية ودالة التمييز اللوجستية.

تصنيف دالة التمييز التربيعية

الجدول (2) يبين تصنيف الدالة التربيعية

التصنيف الصحيح (Y)	التصنيف			المجموع
	المجموعة الاولى	المجموعة الثانية	المجموعة الثالثة	
المجموعة الاولى	87	0	5	92
	94.57	0.00	5.43	%100
المجموعة الثانية	1	52	1	54
	1.85	96.30	1.85	%100
المجموعة الثالثة	47	2	75	124
	37.90	1.61	60.48	%100
المجموع	135	54	81	270
	50.00	20.00	30.00	

يبين الجدول (2) مدى دقة النتائج النهائية لتصنيف الدالة التربيعية إذ يتبين إن (87) حالة من المجموعة الأولى وبنسبة 94.57 % قد تم تصنيفها بشكل صحيح وإن (5) حالات وبنسبة 5.43% تم تصنيفها بشكل خاطئ إذ صنفت بأنها ضمن المجموعة الأولى وهي تعود إلى المجموعة الثالثة وفي المجموعة الثانية يتبين إن (52) حالة وبنسبة 96.30% قد تم تصنيفها بشكل صحيح وبناءً عليه فإن هناك حالة واحدة وبنسبة (1.85) تم تصنيفها بشكل خاطئ إذ صنفت بأنها ضمن المجموعة الأولى وهي تعود إلى المجموعة الثانية و كذلك هناك حالة واحدة وبنسبة (1.85) صنفت بشكل خاطئ إذ تم تصنيفها ضمن المجموعة الثانية وهي تعود إلى المجموعة الثالثة. وفي المجموعة الثالثة يتبين إن (75) حالة وبنسبة 60.81% تم تصنيفها بشكل صحيح وبناءً عليه فإن هناك (47) حالة وبنسبة 37.90% تم تصنيفها بشكل خاطئ إذ صنفت ضمن المجموعة الثالثة وهي تعود إلى المجموعة الأولى وكذلك هناك حالتان وبنسبة 1.61% صنفت بشكل خاطئ إذ صنفت بأنها ضمن المجموعة الثالثة وهي تعود إلى المجموعة الثانية، وكننتيجة عامة فقد دلت النتائج بأن الحالات المصنفة تصنيفاً صحيحاً كانت (214) حالة أي ما نسبته 79.26% من حالات العينة البالغة (270) .

تصنيف دالة التمييز اللوجستية

الجدول (3) يبين تصنيف الدالة اللوجستية

التصنيف الصحيح (Y)	التصنيف			المجموع
	المجموعة الاولى	المجموعة الثانية	المجموعة الثالثة	
المجموعة الاولى	83	1	8	92
	90.22	1.09	8.70	%100
المجموعة الثانية	1	52	1	54
	1.85	96.30	1.85	%100
المجموعة الثالثة	42	2	80	124
	33.87	1.61	64.52	%100
المجموع	126	55	89	270
	46.67	20.37	33.96	

يبين الجدول (3) إلى مدى دقة النتائج النهائية لتصنيف الدالة اللوجستية إذ يتبين إن (83) حالة من المجموعة الأولى وبنسبة 90.22% قد تم تصنيفها بشكل صحيح وبناءً عليه فإن هناك حالة واحدة وبنسبة 1.09% تم تصنيفها بشكل خاطئ إذ صنفتم ضمن المجموعة الأولى وهي تعود إلى المجموعة الثانية وكذلك يتبين إن هناك (8) حالات وبنسبة 8.70% تم تصنيفها بشكل خاطئ إذ صنفتم بأنها ضمن المجموعة الأولى وهي تعود إلى المجموعة الثالثة وفي المجموعة الثانية يتبين إن (52) حالة وبنسبة 96.30% قد تم تصنيفها بشكل صحيح وبناءً عليه فإن هناك حالة واحدة وبنسبة (1.85) تم تصنيفها بشكل خاطئ إذ صنفتم بأنها ضمن المجموعة الأولى وهي تعود إلى المجموعة الثانية وكذلك هناك حالة واحدة وبنسبة (1.85) صنفتم بشكل خاطئ إذ تم تصنيفها ضمن المجموعة الثانية وهي تعود إلى المجموعة الثالثة. وفي المجموعة الثالثة يتبين إن (80) حالة وبنسبة 64.52% تم تصنيفها بشكل صحيح وبناءً عليه فإن هناك (42) حالة وبنسبة 33.87% تم تصنيفها بشكل خاطئ إذ صنفتم ضمن المجموعة الثالثة وهي تعود إلى المجموعة الأولى وكذلك هناك حالتان وبنسبة 1.61% صنفتم بشكل خاطئ إذ صنفتم بأنها ضمن المجموعة الثالثة وهي تعود إلى المجموعة الثانية، وكنتيجة عامة فقد دلت النتائج بأن الحالات المصنفة تصنيفاً صحيحاً كانت (215) حالة أي ما نسبته 79.62% من حالات العينة البالغة (270).

### Conclusions

### الاستنتاجات

1. عند اختبار نسبة التصنيف الصحيح تبين ان الدالة اللوجستية افضل من الدالة التربيعية اذ بلغت نسبتها (79.62) في حين ان الدالة التربيعية فقد بلغت نسبتها (79.26).
2. عند المقارنة بين افضل نموذج للتصنيف تبين ان النموذج اللامعلمي افضل النماذج في تمثيل بيانات الدراسة كون البيانات كانت تعاني من مشاكل في التوزيع الطبيعي وتجانس التباينات.

### Recommendations

### التوصيات

1. نوصي بتطبيق التصنيف على وفق دالة التمييز التربيعية ودالة التمييز اللوجستية في مجالات أخرى غير الطبية.
2. نوصي باستخدام الدوال التمييزية الأخرى (دالة كيرنل، دالة الرتبة، الدالة اللوجستية التربيعية، نايف بيز للتصنيف).
3. توسيع عدد المتغيرات في الدراسة لتشمل (التدخين، الوراثة، تناول المشروبات الكحولية، التعرض المكثف للأشعة الضارة، الإصابة بعدوى فيروسية أو بكتيرية،...) من خلال استمارة شاملة كل المعلومات.
4. من الضروري تطوير مراكز الإحصاء في المستشفيات فضلاً عن معلومات جديدة للمرضى المصابين.
5. ضرورة تشخيص المرض وإجراء العمليات الجراحية اللازمة في وقت مبكر.

### المصادر

1. الجبوري، شلال و حمزة، صلاح، "تحليل متعدد المتغيرات"، دار الكتب لجامعة بغداد، (2000).
2. Ager, J.W, JR and Brent, S.B, "An index of agreement between ahypothesized partial order and an empirical Rank order", American statistical Association, volume 73, number 364, Theory and method section, (1978).
3. Anderson, J.A., "Separate Sample Logistic Discrimination", Biometrika, 59, 19-35, (1972).
4. Berkson J., "Application of the Logistic function to Bioassay", JASA, vol. pp.357-365, (1944).
5. Burns, Robert and Burns, Richard, "Business Research Methods and Statistics using SPSS", Five extra advanced chapters, chapter 24 Logistic Regression: pp 568 -575, (2008).
6. Constanza, M.C. and Afifi, A.A., "Comparison of Stopping Rules in Forward Stepwise Discriminant Analysis", JASA. Vol. 74, No.368, pp777 – 785, (1979).
7. Hardle, W. and Simar, L., "Applied Multivariate Statistical Analysis", Berlin and Louvain-la-Neuve, Germany,p332, (2003).
8. He, X. and Fung, W.K., "High breakdown estimation for multiple with applications to discriminant analysis", Journal of Multivariate analysis, 72, 151-162, (2000).
9. Kiang, M.Y., "A Comparative Assessment of Classification Methods". Decision Support Systems, 35, 441-454, (2003).
10. Kshlr Sagar M.A., "Multivariate Analysis", Marcel Dekker, Inc, New York, (1972).

11. McLachlan, G. J., "Discriminant Analysis and Statistical Pattern Recognition", New York: Wiley, (1992).
12. Nussbaumer, H. J., "Fast Fourier Transform and Convolution Algorithms" Springer – Verlag Berlin Heidelberg, (1982).
13. Rausch, J. R., "A comparison of Linear, quadratic, and mixture models for Discriminant Analysis under non Normal Continuous Predictors", <http://www.stat.notre dame.edu/www/research/reports>, (2005).
14. Schlotzhauer, D.C, "Some issues in using proc logistic for Binary logistic regression", the technical for SAS software users vol. 2. No. 4. (1993).
15. Therrien, C.W., "Decision Estimation and classification", John wiley & Sons New York, (1989).
16. Young, T.Y., "Classification, Estimation and Pattern Recognition" American Elsevier Publishing Company Inc, New York, (1974).