

# التحليل المميز والانحدار اللوجستي بأسعمال المربعات الصغرى

## الجزئية (دراسة تجريبية (محاكاة))

أ.م.د. رباب عبد الرضا صالح البكري /جامعة بغداد / كلية الادارة والاقتصاد  
الباحث / محمد شاكر محمود العزي / جامعة بغداد / كلية الادارة والاقتصاد

تاريخ التقديم: 2018/2/2

تاريخ القبول: 2018/7/24

### المستخلص

بعد اسلوبی الانحدار اللوجستي الثنائی regression logistic Binary والدالة المميزة الخطية Linear discriminant function من اهم الاساليب الاحصائية المستخدمة في التصنيف والتبؤ، عندما تكون البيانات من النوع الثنائي (0,1) فإنه لا يمكن استخدام الانحدار الاعتيادي فلذلك نل JACK الى الانحدار اللوجستي الثنائی والدالة المميزة الخطية في حالة وجود مجموعتين او اكثر، وفي حالة وجود مشكلة التعدد الخطی Multicollinearity بين البيانات (ان البيانات يوجد فيها ارتباطات عالية بين المتغيرات) أصبح عدم الامکان في استخدام الانحدار اللوجستي والدالة المميزة الخطية، ولحل هذه المشكلة توجد عدة طرائق منها طريقة انحدار المربعات الصغرى الجزئية Partial least square regression لحل مشكلة التعدد الخطی.

وقد جرى في هذه البحث المقارنة بين الانحدار اللوجستي الثنائی regression logistic binary والدالة المميزة الخطية linear discriminant function عن طريق خطأ التصنيف. حيث تم توليد بيانات بمتغير استجابة (Y) نوع ثانی (0,1) تحتوي على مشكلة التعدد الخطی ويحوم عينات (50-100-150-200-250) ومتغيرات (20-40-5-10). حيث تمت معالجة مشكلة التعدد الخطی بأسعمال طريقة المربعات الصغرى الجزئية .Partial least square

وتوصل البحث الى ان الدالة المميزة الخطية linear discriminant function هي افضل في تصنيف البيانات من الانحدار اللوجستي الثنائی binary logistic regression ، اذ صنفت الدالة المميزة البيانات بشكل صحيح وأكثر دقة من الانحدار اللوجستي الثنائی.

**المصطلحات الرئيسية للبحث:** الدالة المميزة الخطية – الانحدار اللوجستي الثنائی – المربعات الصغرى الجزئية – مشكلة التعدد الخطی – نسبة التصنيف – محاكاة.



مجلة العلوم  
الاقتصادية والإدارية  
العدد 106 المجلد 24  
الصفحات 419-407

\*البحث مستقل من رسالة ماجستير



## التحليل المميز والانحدار اللوجستي باستعمال المربعات الصغرى الجزئية [دراسة تجريبية [محاكاة]]

### 1- المقدمة Introduction

ان اسلوب الانحدار اللوجستي من الاساليب المهمة التي تدخل في تحليل البيانات التي تكون فيها بيانات المتغير المعتمد ( $Y$ ) يكون فيه البيانات ثنائية (0،1)، ويكون الهدف الاساس من هذه الطريقة هو ايجاد أفضل نموذج يصف الحالة بين المتغير المعتمد والمتغير التوضيحي (المتغيرات التوضيحية). وان طريقة التحليل المميز هي من التحليلات الاحصائية التي تهتم بمتعدد المتغيرات حيث يتم بهذا التحليل استعمال مجموعة من المتغيرات للتمييز بين مجموعتين او أكثر بواسطة عدة دوال تميزية. سيتم اعتماد هذين الاسلوبين للعمل في هذا البحث، وللتخلص من مشكلة التعدد الخطى سيتم استعمال طريقة المربعات الصغرى الجزئية.

### 2- مشكلة البحث problem of research

لغرض استعمال بيانات من النوع ثانى الاستجابة والتي تعانى من وجود مشكلة التعدد الخطى من اجل تكوين انماذج احتمالي والذي يتم على اساسه تصنيف البيانات، حيث ثم استعمال اسلوبين من اساليب التصنيف لغرض ايجاد افضل انماذج احتمالي بأقل خطأ صنف ممكن.

### 3- هدف البحث objective of research

يتضمن هذا البحث المقارنة بين اسلوبين من اساليب التصنيف وهما الانحدار اللوجستي الثنائى والدالة المميزة الخطية بوجود مشكلة التعدد الخطى ولمعرفة مدى قابليتهم لتصنيف البيانات بأقل احتمال خطأ للتصنيف. وذلك باستعمال تجارب محاكاة عندما يتوزع الخطأ توزيع طبيعى لحجوم عينات وابعاد مختلفة.

### 4- الجانب النظري

من الاساليب الإحصائية المهمة في متعدد المتغيرات والتي تستعمل في تحليل وتقويم العلاقات بين مجموعة من المتغيرات يكون فيه المتغير التابع (مقطع)، ليس دائماً يكون فيها المتغير التابع مستمراً وذلك لغرض الحصول على انماذج رياضي يوضح العلاقة بين مثل هكذا بيانات تستعمل اسلوب التحليل المميز والانحدار اللوجستي. وان هناك بعض الافتراضات الخاصة للأسلوبين أعلاه منها وان تكون المتغيرات التوضيحية مستقلة ولا يوجد أي ارتباط بينها ولمعالجة مثل هكذا حالات تستعمل عدة طرائق منها طريقة المركبات الرئيسية وطريقة المربعات الصغرى الجزئية وسيتم في هذا البحث التركيز على طريقة المربعات الصغرى الجزئية.

#### 1-4 المربعات الصغرى الجزئية<sup>(1)(5)(6)</sup> Partial least square

تعد طريقة المربعات الصغرى الجزئية أكثر الطرائق أهمية في الانحدار فهي تستعمل لتقليل عدد المتغيرات التوضيحية المرتبطة في الانماذج الى مركبات غير مرتبطة (خطية، متعامدة)، او عندما يكون عدد المتغيرات التوضيحية أكثر من عدد المشاهدات في التجربة. وطريقة المربعات الصغرى الجزئية مشابهة لطريقة المركبات الرئيسية وأسلوب انحدار الحرف لمعالجة مشكلة التعدد الخطى ولكنها تختلف في الحسابات فخوارزمية (Partial least square) تأخذ بنظر الاعتبار التباين المشترك ما بين م بين متغير (متغيرات) الاستجابة والمتغيرات التوضيحية، اما طريقة المركبات الرئيسية (PCA) تأخذ بنظر الاعتبار التباين بين المتغيرات التوضيحية فقط. يقوم بتحويل المتغيرات التوضيحية المرتبطة الى مركبات رئيسية تختلف في الحسابات فخوارزمية الحل بطريقة PLS طريقة تكرارية عند استعمالها تنتج سلسلة من النماذج ويتوقف الحل التكراري عندما نصل الى العدد الكلى من المركبات في الانماذج او عندما تكون الباقي مساوية للصفر.

ففي حال تساوي عدد المركبات مع عدد المتغيرات التوضيحية فإن النتائج ستكون متطابقة مع طريقة المربعات الصغرى. ولأجل تحديد عدد المركبات التي تصغر من خطأ التنبؤ تستعمل طريقة العبور الشرعي Cross-validation وبعد تحديد عدد المركبات تقدر معلم انماذج الانحدار لكل متغير.



## التحليل المميز والانحدار اللوجستي باستعمال المربعات الصغرى الجزئية [دراسة تجريبية [محاكاة]]

وان طريقة PLS اول من طبقها الباحث Wold عندما يكون هناك ارتباط عالي ما بين المتغيرات التوضيحية او عندما يكون عدد المتغيرات التوضيحية تفوق عدد المشاهدات. وتوجد عدة خوارزميات وأكثرها تداولاً خوارزمية NIPALS عام 1973 للباحث Wold وخوارزمية SIMPLS 1993 المنسوبة للعالم De Jong وخوارزمية KERNEL المنسوبة للعالم DAYAL وخوارزمية TRYGG latent structures (PLS1,PLS2) حيث ان خوارزمية PLS1 تستعمل عندما يكون متغير الاستجابة متوجه اما خوارزمية PLS2 تستعمل عندما يكون متغير الاستجابة مصفوفة، وخوارزمية PLS1 تعطي نفس نتائج خوارزمية SIMPLS. وسنركز في هذا البحث على خوارزمية SIMPLS.

### 2-4 خوارزمية SIMPLS Algorithm

#### 1- ايجاد مصفوفة التباين والتباين المشترك للبيانات

$$S_{xy} = \hat{X}\hat{Y} \quad (\text{X and Y are enteral}) \dots\dots\dots(1)$$

$$\tilde{X}_i = X_i - \bar{X} \dots\dots\dots(2)$$

$$\tilde{Y}_i = Y_i - \bar{Y} \dots\dots\dots(3)$$

2- إعادة تكرار الخطوات من (1-2) الى (6-2) لكل من  $h=1,2,3,\dots\dots\dots,k$

1-2 ايجاد اول متوجه مميز للمصفوفة  $S_{xy}^*$  ول  $\mathbf{h}^{th}$  من متوجهات الاوزان ( $r_h$ ) لطريقة المربعات الصغرى الجزئية (X-weight) مصفوفة الاوزان.

2-2 ايجاد  $\mathbf{h}^{th}$  من المركبات الخطية (X-score) مصفوفة الدرجات

$$t_h = \hat{X}r_h \quad \text{and normalized} \dots\dots\dots(4)$$

$$t_h = t_h / \|t_h\| \dots\dots\dots(5)$$

### 3-2 حساب $\mathbf{h}^{th}$ تحميلات X – loading على $t_h$

$$\hat{P}_h = \hat{X}t_h \dots\dots\dots(6)$$

### 4-2 خزن المتوجهات $t_h r_h, p_h$ في المصفوفات

$$\mathbf{R}_h = \{r_1, \dots, r_h\} \dots\dots\dots(7)$$

$$\mathbf{T}_h = \{t_1, \dots, t_h\} \dots\dots\dots(8)$$

$$\mathbf{P}_h = \{P_1, \dots, P_h\} \dots\dots\dots(9)$$

5-2 اذا كان  $h=h+1$  فيتم حساب المعادلة:

$$S_{xy}^h = (I_p - V_{h-1} * \hat{V}_{h-1}) \hat{S}_{xy}^{h-1} \dots\dots\dots(10)$$

وأن  $V_{h-1}$  هو متعامد طبيعي . orthonormal

$$S_{xy}^h = S_{xy}^{h-1} - P_{h-1}(\hat{P}_{h-1}, P_{h-1})^{-1} \hat{P}_{h-1} S_{xy}^{h-1} \dots\dots\dots(11)$$

وأن  $(V_1, \dots, V_{h-1})$  متعامد (orthogonal) الى مصفوفة التحميل X ( $P_1, \dots, P_{h-1}$ )

3- يتم الحصول على مقدرات المعلومات لإنموذج الانحدار الخطى

$$\hat{B} = R_{ph}(\hat{R}_{h*p} * S_x * R_{p*1}) \hat{R}_{h,p} S_{xy} \dots\dots\dots(12)$$

$$\widehat{B}_0 = \bar{Y} - \widehat{B}_{1,p} \bar{X} \dots\dots\dots(13)$$



حيث ان:  
 $R$ : مصفوفة اوزان –  $X$  –  $p^*h$  (X-weights) بأبعاد  $p \times h$  حيث ان  $p$  عدد المتغيرات التوضيحية وان  $h$  عدد المركبات.

$n^*h$  هي مصفوفة الدرجات  $P_h$  هي مصفوفة التحميل  $S_{xy}$ : مصفوفة التباين والتباين المشترك الى  $X$  و  $S_x$ : مصفوفة التباين الى  $X$ .

## 5- الانحدار اللوجستي logistic regression

بعد انموذج الانحدار اللوجستي من النماذج الإحصائية المهمة في تحليل البيانات اذ ان الهدف الأساسي من معظم الدراسات هو تحليل وتقويم العلاقات بين مجموعة من المتغيرات للحصول على صيغة نستطيع من خلالها ان نصف النموذج ويستعمل انموذج الانحدار اللوجستي لوصف العلاقة بين بين متغير الاستجابة من النوع المتقطع والمتغيرات التوضيحية ويكون على نوعين انموذج الانحدار اللوجستي ثانوي الاستجابة وانموذج الانحدار اللوجستي متعدد الاستجابة وسنركز على النوع الأول من الانحدار اللوجستي.

### 1-5 انموذج الانحدار اللوجستي ثانوي الاستجابة

ان من خصائص الانحدار اللوجستي ثانوي الاستجابة ان المتغير التابع ( $y$ ) متغير الاستجابة يتبع توزيع بيرنولي ويأخذ القيم (0) و (1) اي بمعنى ان (1) باحتمال قدره ( $p$ ) احتمال ناجح، وباحتمال فشل ( $1-p$ ) قدره (0) وبمعنى اخر احتمال حدوث الاستجابة (1) واحتمال عدم حدوث الاستجابة (0). ان نموذج الانحدار الخطي المتعدد MLR يكون كالتالي:

$$Y = XB + \epsilon \quad (14)$$

حيث ان:

$Y$  : موجه المتغير المعتمد بأبعاد ( $n \times 1$ )  
 $X$  : مصفوفة المتغيرات التوضيحية ( $(n \times k+1) \times (k+1) \times 1$ )  
 $B$  : موجه معلمات دالة الانحدار ( $(k+1) \times 1$ )  
 $\epsilon$  : موجه الأخطاء العشوائية والذي يشترط به تحقق الشروط الآتية:

-1  $\epsilon_i$  تتوزع توزيع طبيعي

-2  $E(\epsilon_i) = 0$

-3  $Cov(\epsilon_i, \epsilon_j) = 0$

-4  $Var(\epsilon_i) = \sigma^2$

$n$ : حجم العينة.

$k$ : عدد المعلمات

في حالة وجود متغير توضيحي واحد فان متوسط قيم المشاهدة  $Y$  عند متوجه معين للمتغير  $x$  فهو  $E(y|x)$  وبذلك يمكن كتابة النموذج على النحو التالي:

$$E(y|x) = \beta_0 + \beta_1 x \quad (15)$$

ان الطرف اليمين لهذا النموذج يأخذ قيم  $(-\infty, +\infty)$  ، لكن عندما يكون المتغير ( $y$ ) ثانيا فان الانموذج اعلاه لا يكون ملائما لأن:

$$E(y/x) = P_r(y=1) = p \quad (16)$$



وفي هذه الحالة يكون الطرف اليمين محصور بين  $(0, 1)$ ، وهذا يعني ان الانموذج يكون غير قابل للتطبيق احصائيا. للتخلص من هذه المشكلة سنقوم بادخال تحويلة رياضيا على المتغير التابع  $y$ . وبما ان  $(0 \leq p \leq 1)$  وان  $\frac{p}{1-p}$  هو مقدار موجب محصور بين  $(0, \infty)$  أي  $(0 \leq \frac{p}{1-p} \leq \infty)$  وبأخذ اللوغاريتم الطبيعي للاسas  $(e)$  للمقدار  $(\frac{p}{1-p})$  فان مجال القيمة تصبح محصور بين  $(-\infty, +\infty)$  وتكون كالاتي  $\log_e(\frac{p}{1-p}) \leq \infty$ ، وبالنهاية يمكن كتابة انموذج الانحدار في حالة وجود متغير واحد وكالاتي:

اما اذا كان لدينا اكثر من متغير توضيحي فتصبح صيغته كالتالي:

اذ ان:

$$i = 1, 2, 3, \dots, n \quad j = 1, 2, 3, \dots, k$$

$\beta_1, \beta_2, \dots, \beta_k$ : موجه للمعلم المطلوب تقديرها.

$x_{ij}$ : متغيرات توضيحية.

اما بالنسبة  $\frac{p}{1-p}$  نسبة افضلية النجاح ( odds of success ) او نسبة الافضلية للحدث المرغوب به وصيغته الرياضية هي كالتالي:

$$\frac{P(Y=1)}{1-P(Y=1)} = e^{\beta_0 + \sum_{j=1}^k \beta_j x_{ij}} \dots \dots \dots (19)$$

وصيغة احتمالات الاستجابة لأنموذج الانحدار اللوجستي تكتب كالتالي:

$$p = \frac{1}{1 + (e^{\beta_0 + \beta_1 x_i})^{-1}} \dots \dots \dots (20)$$

وان المقدار  $(\frac{p}{1-p})$  يسمى لوغاریتم نسبة افضلية النجاح (logs odds of success).

## ٥-افتراضات الخاصة بالانحدار الوجستي

ان الانحدار اللوجستي لا يتطلب افتراضات كثيرة فقط يتطلب عدم وجود ارتباط بين المتغيرات التوضيحية وان حجم المشاهدات كبيرة في كل مجموعة يفترض انها تكون اكبر من خمس مرات من عدد المعلومات المستعملة في الانموذج النهائي.

### 3- تقدیر احتمالات الاستجابة 3-5 (8)(4)(2)Estimation of response probabilities

ان صيغة احتمال الاستجابة لأنموذج الانحدار اللوجستي الثاني يمكن الحصول عليها وذلك بادخال ال (e) على المعادلة التالية:

وذلك صيغة احتمال الاستجابة تصبح كالتالي:



ويمكن ايضا ان نكتب المعادلة اعلاه بالصيغة الآتية:

وباستعمال الطرق الجبرية للاستدلال على صحة احتمال الاستجابة لاموزج الانحدار اللوجستي الثنائي يكون كالتالي:

وهذا يعني ان احتمال متغير الاستجابة ( $y$ ) تأخذ القيمة (1) ويكون كالتالي:

و عند القيمة (0) فإن احتمال متغير الاستجابة (y) يكون على النحو التالي:

وبما ان مجموع الاحتمالات يساوي (1) فأن:

## 6- التحليل التمييزي Discriminant analysis :

يعد التحليل التميزي من الأساليب الإحصائية المهمة في متعدد المتغيرات التي تهتم بتفريق (تميز) بين مجموعتين (مجموعتين) أو أكثر من خلال إيجاد تواافق خطية للمتغيرات التوضيحية تعرف بذالة التميز التي يمكن عن طريقها الفصل او تميز بين مجموعتين او أكثر. وهناك عدة دوال تميزية منها الدالة المميزة الخطية والتي قدمت من قبل العالم فيشر عام (1936) والتي تستعمل في حالة تساوي مصفوفة التباينات المشتركة للمجاميع كافة وان كل مجموعة تعبر عن مجتمع طبيعي لمتعدد المتغيرات اما في حالة عدم تساوي مصفوفة التباينات المشتركة نجأ الى الدالة المميزة التربيعية. ان من الافتراضات الخاصة بالتحليل المميز بأن المجاميع المدروسة ذات توزيع طبيعي متعدد متغيرات ، كذلك يتطلب ان تكون حجم العينة كبيرة حيث ان المجاميع المختلفة تتضمن على الأقل 20 مشاهدة لكل متغير توضيحي وأيضا يشترط عدم وجود ارتباط بين المتغيرات التوضيحية وعدم وجود قيم شاذة بينها.

1-6 **الدالة المتميزة الخطية لاثنتين** The linear discriminant function for two groups (9)(7)(3)

تسمى هذه الدالة بدالة فيشر (Fisher Function) توزع فيها المشاهدات توزيعاً طبيعياً. لنفترض ان لدينا مجتمعين ونريد المقارنة بينهما ولنفرض ان هذين المجتمعين لهما نفس مصفوفة التباين والتباين المشترك ( $\Sigma_1 = \Sigma_2$ ) ولهم موجه متوسطات ( $\mu_1, \mu_2$ ) بالتتابع، وتم اختيار عينتين عشوائيتين

لكل من المجتمعين.  $(\underline{x}_{21}, \underline{x}_{22}, \underline{x}_{23}, \dots, \underline{x}_{2n_2})$  و  $(\underline{x}_{11}, \underline{x}_{12}, \underline{x}_{13}, \dots, \underline{x}_{1n_2})$

و تكون صيغة الدالة المميزة الخطية لفيشر تكتب بالشكل الآتي:

## حیث ان:



وأن مقياس (Mahalanobis) يعتمد بالأساس على القياسات ذات المسافة القليلة بين قيم المتغيرات للمشاهدات الجديدة وقيم المتوسطات للمتغيرات لكل مجموعة يمكن كتابتها بالشكل الآتي:

$$D_i^2 = \left( \bar{x}_1 - \bar{x}_i \right)' S_p^{-1} \left( \bar{x}_1 - \bar{x}_i \right) \dots \dots \dots (30)$$

إذ أن :

**x<sub>i</sub>**: موجة متوسطات المتغيرات لكل قيمة من المجموعة (i)

**$S_p^{-1}$**  : هو معكوس مصفوفة التباین والتباين المشترك المقدرة داخل العینتین.

## ٧- **تصنيف البيانات** classification of data <sup>(١١)</sup>

لنفرض ان:

## لدينا عينة بحجم (n)

ان عدد المشاهدات من النوع (0) هي  $n_1$

وأن عدد المشاهدات من النوع (1) هي  $n_2$

وان لدينا بيانات من نوع ثانوي (binary) وكانت لدينا تصنیف البيانات كما مبين في الجدول رقم (1) الآتي:

### **جدول رقم (1) يبين تصنیف البيانات**

|                         |                       |                       |
|-------------------------|-----------------------|-----------------------|
| <b>البيانات الثانية</b> | <b>(0)</b>            | <b>(1)</b>            |
| <b>(0)</b>              | <b>A<sub>11</sub></b> | <b>A<sub>12</sub></b> |
| <b>(1)</b>              | <b>A<sub>21</sub></b> | <b>A<sub>22</sub></b> |

وتكون معايير التصنيف بالشكل الآتي:

### **١- نسبة التصنيف الصحيحة تحسب بالشكل الآتي:**

ومن ثم يمكن حساب معيار خطأ التصنيف الصحيح بالشكل الآتي:

نسبة التصنيف الصحيح = معيار خطأ التصنيف .....(32)

2- نسبة التصنيف الجزئية على مستوى (0) و (1) وتحسب بالشكل الآتي:

١- نسبة التصنيف الكلية على مستوى (٠) و (١) وتحسب بالشكل الآتي:



8- الجانب التجريبى

## المحاكاة 1-8

تعرف المحاكاة ب أنها تصميم أنماذج افتراضي مشابه للأنماذج الحقيقية لغرض معرفة سلوك الانموزج للوصول الى الهدف المطلوب وبمعنى اخر هو إعادة صياغة الواقع الفعلي بواقع افتراضي باستعمال صيغ ونماذج معينة وفق أسلوب تكراري لعدة مرات من التكرارات لغرض الوصول الى نتائج تجريبية تهدف الى المقارنة بينها وبين نتائج الجانب التطبيقي.

ان الاعتماد على تجارب المحاكاة في المجال الاحصائي بصورة عامة وفي مجالات التقدير والتباين والتصنيف بصورة خاصة كونه يقارن بين حجوم العينات وعدد المتغيرات في الدراسات ومن ثم التوصل الى افضل الطرائق تحت الدراسة.

-8 المحاكاة لتجربة وصف

ان موضوع الرسالة سيتم توليد بيانات بمتغيرات عدد (P=15) بحجم عينة (25، 50، 100، 150، 250، 400) باستعمال برنامج Matlab (R2015b) وبالاعتماد على المعادلة الآتية:

$$x_{p-1} = N(0,0.1) + x_1 \dots \dots \dots \quad (3-120)$$

$$Y = \text{randi}([0,1], n, 1) \dots \dots \dots \quad (3-121)$$

ان هذه البيانات تعانى من مشكلة تعدد خطى.

**1-2-8 في حالة حجم العينة (100 – 150 – 250 – 400) وبعد مرکبات (2):**

#### **الجدول (4) بين نسب التصنيف**

**في المحاكاة للدالة المميزة بحجم عينة (100 – 150 – 250 – 400) ومركبات (2)**

| $h=2$ n=100        | (0)   | (1)   | نسبة التصنيف الصحيحة |
|--------------------|-------|-------|----------------------|
| (0)                | 38    | 9     | 80.9%                |
| (1)                | 11    | 42    | 79.2%                |
| نسبة التصنيف الكلي | 49%   | 51%   | %80                  |
| $h=2$ n=150        | (0)   | (1)   | نسبة التصنيف الصحيحة |
| (0)                | 59    | 18    | 76.6%                |
| (1)                | 18    | 55    | 75.3%                |
| نسبة التصنيف الكلي | 51.4% | 48.6% | %76                  |
| $h=2$ n=250        | (0)   | (1)   | نسبة التصنيف الصحيحة |
| (0)                | 90    | 37    | 70.9%                |
| (1)                | 45    | 78    | 63.4%                |
| نسبة التصنيف الكلي | 54%   | 46%   | 67.2%                |
| $h=2$ n=400        | (0)   | (1)   | نسبة التصنيف الصحيحة |
| (0)                | 130   | 72    | 64.6%                |
| (1)                | 82    | 116   | 58.6%                |
| نسبة التصنيف الكلي | 53%   | 47%   | %61.5                |



**التحليل المميز والانحدار اللوجستي باستعمال المربعات الصغرى  
الجزئية [دراسة تجريبية [محاكاة]]**

**الجدول (5) يبين نسب التصنيف في المحاكاة  
للانحدار اللوجستي الثنائي بحجم عينة 100 – 150 – 250 – 400 ومرکبات (2)**

| $h=2$ n=100        | (0)   | (1)   | نسبة التصنيف الصحيحة |
|--------------------|-------|-------|----------------------|
| (0)                | 21    | 26    | 49%                  |
| (1)                | 23    | 30    | 44.7%                |
| نسبة التصنيف الكلي | 44%   | 56%   | %51                  |
| $h=2$ n=150        | (0)   | (1)   | نسبة التصنيف الصحيحة |
| (0)                | 40    | 37    | 51.9%                |
| (1)                | 38    | 35    | 47.9%                |
| نسبة التصنيف الكلي | 52%   | 48%   | %50                  |
| $h=2$ n=250        | (0)   | (1)   | نسبة التصنيف الصحيحة |
| (0)                | 64    | 63    | 50.4%                |
| (1)                | 62    | 61    | 49.6%                |
| نسبة التصنيف الكلي | 50.4% | 49.6% | %50                  |
| $h=2$ n=400        | (0)   | (1)   | نسبة التصنيف الصحيحة |
| (0)                | 102   | 100   | 50.5%                |
| (1)                | 100   | 98    | 49.5%                |
| نسبة التصنيف الكلي | 50.5% | 49.5% | %50                  |

نلاحظ من الجدولين (4) (5) أعلاه عندما يكون عدد المرکبات (2) وبجوم عينات (100 – 150 – 250 – 400) ان نسب تصنیف الدالة المميزة الخطیة كانت اعلى من نسب التصنيف لانحدار اللوجستی الثنائی وبالتالي فأن خطأ التصنیف للدالة المميزة الخطیة اقل من خطأ التصنیف لانحدار اللوجستی وهذا يعني ان الدالة المميزة الخطیة افضل من الانحدار اللوجستی.

2-2-2 في حالة حجم العينة (100 – 150 – 250 – 400) وبعد مرکبات (4):  
**جدول (6) يبين نسب التصنيف في المحاكاة**

**للدالة المميزة بحجم عينة 100 – 150 – 250 – 400 ومرکبات (4)**

| $h=4$ n=100        | (0)   | (1)   | نسبة التصنيف الصحيحة |
|--------------------|-------|-------|----------------------|
| (0)                | 42    | 11    | 79.2%                |
| (1)                | 10    | 37    | 84%                  |
| نسبة التصنيف الكلي | 58%   | 42%   | %79                  |
| $h=4$ n=150        | (0)   | (1)   | نسبة التصنيف الصحيحة |
| (0)                | 52    | 23    | 69.3%                |
| (1)                | 21    | 54    | 72%                  |
| نسبة التصنيف الكلي | 48.6% | 51.4% | %70.7                |
| $h=4$ n=250        | (0)   | (1)   | نسبة التصنيف الصحيحة |
| (0)                | 84    | 47    | 46.1%                |
| (1)                | 50    | 69    | 57.9%                |
| نسبة التصنيف الكلي | 53.6% | 46.4% | %61.2                |
| $h=4$ n=400        | (0)   | (1)   | نسبة التصنيف الصحيحة |
| (0)                | 122   | 80    | 60.4%                |
| (1)                | 81    | 117   | 59.1%                |
| نسبة التصنيف الكلي | 50.8% | 49.2% | %59.7                |



**التحليل المميز والانحدار اللوجستي باستعمال المربعات الصغرى  
الجزئية [دراسة تجريبية [محاكاة]]**

**جدول (7) يبين نسب التصنيف في المحاكاة  
للانحدار اللوجستي الثنائي بحجم عينة (100 – 150 – 250 – 400) ومرکبات (4)**

| <b>نسبة التصنيف الصحيحة</b> | <b>(1)</b> | <b>(0)</b> | <b>نسبة التصنيف الكلي</b> |
|-----------------------------|------------|------------|---------------------------|
| 54.7%                       | 24         | 29         | (0)                       |
| 44.7%                       | 21         | 26         | (1)                       |
| %50                         | 45%        | 55%        | نسبة التصنيف الكلي        |
| <b>نسبة التصنيف الصحيحة</b> | <b>(1)</b> | <b>(0)</b> | <b>نسبة التصنيف الكلي</b> |
| 50.7%                       | 37         | 38         | (0)                       |
| 49.3%                       | 37         | 38         | (1)                       |
| %50                         | 49.4%      | 50.6%      | نسبة التصنيف الكلي        |
| <b>نسبة التصنيف الصحيحة</b> | <b>(1)</b> | <b>(0)</b> | <b>نسبة التصنيف الكلي</b> |
| 52..7%                      | 62         | 69         | (0)                       |
| 47.6%                       | 56         | 63         | (1)                       |
| %50                         | 47.2%      | 52.8%      | نسبة التصنيف الكلي        |
| <b>نسبة التصنيف الصحيحة</b> | <b>(1)</b> | <b>(0)</b> | <b>نسبة التصنيف الكلي</b> |
| 50.5%                       | 100        | 102        | (0)                       |
| 49.5%                       | 98         | 100        | (1)                       |
| 50%                         | 49.5%      | 50.5%      | نسبة التصنيف الكلي        |

نلاحظ من الجدولين رقم (6) (7) أعلاه عندما يكون عدد المرکبات (2) وبحجم عينات (100 – 150 – 250 – 400) ان نسب تصنيف الدالة المميزة الخطية كانت اعلى من نسب التصنيف للانحدار اللوجستي الثنائي ومن ثم فأن خطأ التصنيف للدالة المميزة الخطية اقل من خطأ التصنيف للانحدار اللوجستي وهذا يعني ان الدالة المميزة الخطية افضل من الانحدار اللوجستي.

**3-2-8 في حالة حجم العينة (100 – 150 – 250 – 400) وبعد مرکبات (7):**  
**جدول (8) يبين نسب التصنيف في المحاكاة**

**للدالة المميزة بحجم عينة (100 – 150 – 250 – 400) ومرکبات (7)**

| <b>نسبة التصنيف الصحيحة</b> | <b>(1)</b> | <b>(0)</b> | <b>نسبة التصنيف الكلي</b> |
|-----------------------------|------------|------------|---------------------------|
| 68.6%                       | 11         | 24         | (0)                       |
| 81.5                        | 53         | 12         | (1)                       |
| 77%                         | 64%        | 36%        | نسبة التصنيف الكلي        |
| <b>نسبة التصنيف الصحيحة</b> | <b>(1)</b> | <b>(0)</b> | <b>نسبة التصنيف الكلي</b> |
| 65.4%                       | 27         | 51         | (0)                       |
| 66.7%                       | 48         | 24         | (1)                       |
| %66                         | 50%        | 50%        | نسبة التصنيف الكلي        |
| <b>نسبة التصنيف الصحيحة</b> | <b>(1)</b> | <b>(0)</b> | <b>نسبة التصنيف الكلي</b> |
| 66.4%                       | 40         | 79         | (0)                       |
| 64.1%                       | 84         | 47         | (1)                       |
| %65.2                       | 49.6%      | 50.4%      | نسبة التصنيف الكلي        |
| <b>نسبة التصنيف الصحيحة</b> | <b>(1)</b> | <b>(0)</b> | <b>نسبة التصنيف الكلي</b> |
| 62.6%                       | 73         | 122        | (0)                       |
| 60.5%                       | 124        | 81         | (1)                       |
| 61.5%                       | 49.2%      | 50.8%      | نسبة التصنيف الكلي        |



## التحليل المميز والانحدار اللوجستي باستعمال المربعات الصغرى الجزئية [دراسة تجريبية [محاكاة]]

**جدول (9) يبين نسب التصنيف في المحاكاة  
للانحدار اللوجستي الثنائي بحجم عينة (100 – 150 – 250 – 400) ومرکبات (7)**

| نسبة التصنيف الصحيحة | (1)   | (0)   | نسبة التصنيف الكلي |
|----------------------|-------|-------|--------------------|
| 28.6%                | 25    | 10    | (0)                |
| 70.8%                | 46    | 19    | (1)                |
| 56%                  | 71%   | 29%   | نسبة التصنيف الكلي |
| نسبة التصنيف الصحيحة | (1)   | (0)   | نسبة التصنيف الكلي |
| 52.6%                | 37    | 41    | (0)                |
| 47.2%                | 34    | 38    | (1)                |
| %50                  | 47.4% | 52.6% | نسبة التصنيف الكلي |
| نسبة التصنيف الصحيحة | (1)   | (0)   | نسبة التصنيف الكلي |
| 48.3%                | 63    | 56    | (0)                |
| 52.7%                | 69    | 62    | (1)                |
| %50                  | 52.8% | 47.2% | نسبة التصنيف الكلي |
| نسبة التصنيف الصحيحة | (1)   | (0)   | نسبة التصنيف الكلي |
| 48.7%                | 100   | 95    | (0)                |
| 51.2%                | 105   | 100   | (1)                |
| 50%                  | 51.2% | 48.8% | نسبة التصنيف الكلي |

نلاحظ من الجدولين رقم (8) (9) أعلاه عندما يكون عدد المرکبات (2) وبحجم عينات (100 – 150 – 250 – 400) ان نسب تصنیف الدالة المميزة الخطية كانت اعلى من نسب التصنیف للانحدار اللوجستي الثنائي وبالتالي فان خطأ التصنیف للدالة المميزة الخطية اقل من خطأ التصنیف للانحدار اللوجستي وهذا يعني ان الدالة المميزة الخطية افضل من الانحدار اللوجستي.

### 9- الاستنتاجات

نستنتج من هذا البحث الآتي:

- 1- عند تكرار المحاولة 1000 مرة في الدالة المميزة الخطية كما اقل حجم العينة كانت نتائج خطأ التصنیف اقل، اما في الانحدار اللوجستي الثنائي فان زيادة او قلة حجم العينة لا يؤثر على خطأ التصنیف.
- 2- عند المقارنة بين الطريقيتين باستخدام معيار خطأ التصنیف تم التوصل الى ان الدالة المميزة الخطية هي افضل في تصنیف البيانات من الانحدار اللوجستي الثنائي ولجميع العينات والمتغيرات. وهذا يعني ان نموذج الدالة المميزة افضل في التنبؤ من الانحدار اللوجستي الثنائي لأنها أعطت اقل خطأ تصنیف

### 10- التوصيات

بناءاً على ما تم التوصل اليه من الاستنتاجات نوصي بالآتي:

- 1- اجراء دراسات إحصائية في حالة وجود قيمة شاذة ومشكلة تعدد خطى واجراء دراسات مقارنة بين الانحدار اللوجستي والتحليل المميز بعد ان تتم المعالجة من وجود الشوافذ بأحد الطرق الحصينة والارتباطات بين المتغيرات (طريقة المربعات الصغرى الجزئية)
- 2- نوصي بتوسيع انحدار المربعات الصغرى الجزئية وذلك في حالة عدد المتغيرات التوضيحية اكبر من عدد المشاهدات ومن ثم استعمالها في المقارنة بين الانحدار اللوجستي والدالة المميزة الخطية سواء كانت خطية او تربيعية.



## 11-المصادر

- 1- البكري، رباب عبد الرضا (2015)، مقارنة بعض الطرائق الخطية لمعالجة مشكلة التعدد الخطي في النماذج مع تطبيق عملي، رسالة دكتوراه - كلية الادارة والاقتصاد - جامعة بغداد.
- 2- التميمي، رعد فاضل حسن، (2013)، "الانحدار والسلسل الزمنية أساليب احصائية تطبيقية متقدمة باستخدام برنامج Minitab، كتاب، كلية الادارة والاقتصاد، الجامعة المستنصرية.
- 3- الحمداني، بسمة رشيد، 2014، "تميز الكادر الطبي حسب معرفتهم للتصنيف الدولي (ICD-10)) باستعمال الدالة المميزة"، رسالة ماجستير في جامعة بغداد، كلية الادارة والاقتصاد.
- 4- عباس، علي خضير، 2012، "استخدام أنموذج الانحدار اللوجستي في التنبؤ بالدول ذات المتغيرات الاقتصادية التابعة النوعية"، مجلة كركوك للعلوم الادارية والاقتصادية، مجلد 2، العدد 2.
- 5- Abdi, hervi, 2010," Partial least squares regression and projection on latent structure regression (PLS Regression)", John Wiley & Sons.
- 6- Banh V. Nguyen, David M. Rocke, 2004, " On partial least squares dimension reduction for microarray-based classification: a simulation study, Division of Biostatistics, School of Medicine, University of California.
- 7- Boaz Nadler, Ronald R.Coifman, 2005, " Partial least squares, Beer's law and the net analyte signal: statistical modeling and analysis", Department of Mathematics, Yale University.
- 8- Erik Brorson, Asterios Geroukis, 2014, "A comparison between discriminant analysis and logistic regression using principal components", Department of Statistics, Uppsala University, Uppsala University.
- 9- Leo H. Chiang, Evan L. Russell, Richard D. Braatz, 2000, " Fault diagnosis in chemical processes using Fisher discriminant analysis, discriminant partial least squares, and principal component analysis", Department of Chemical Engineering, universty of Illinois.
- 10- Norliza Adnan, Ahmad Maizah Hura, Adnan Robiah, 2006, "A Comparative Study on Some Methods For Handling Multicollinearity Problems, Department of Mathematics, uivirsty of malysia.
- 11- Tormod Næs and Bjorn-Helge Mevik, "Understanding the collinearity problem in regression and discriminant analysis", Journal of Chemo metrics, P<sup>(413-426)</sup>, 2001.



## **discriminate analysis and logistic regression by use partial least square**

### **Abstract**

The method binary logistic regression and linear discriminant function of the most important statistical methods used in the classification and prediction when the data of the kind of binary (0,1) you can not use the normal regression therefore resort to binary logistic regression and linear discriminant function in the case of two group in the case of a Multicollinearity problem between the data (the data containing high correlation) It became not possible to use binary logistic regression and linear discriminant function, to solve this problem, we resort to Partial least square regression.

In this, search the comparison between binary logistic regression and linear discriminant function using error Category. Where the data has been generating a variable response (Y) binary data (0,1) containing Multicollinearity problem by the samples (50-100-150-250-400) and the variables (5-10-15). Multicollinearity problem has been processed using a method partial least square The research found that linear discriminant function It is the best in the classification of data from binary logistic regression classified as linear discriminant function the data correctly and more accurate than binary logistic regression.

**Keyword:** linear discriminant function- binary logistic regression- partial least square– multicollinearity problem – ratio of classification – simulation.