



التحليل التمييزي

+

(دراسة تطبيقية)

الدكتور / رحيم

المبحث الأول

التحليل التميزي The Discriminate ate analysis

مقدمة:

يلعب التحليل الإحصائي دوراً هاماً في تحليل وتقدير الطواهر الاجتماعية والطبيعية في المجتمع ويُعد التحليل الإحصائي أحد طرق البحث العلمي الذي يستخدم عند دراسة المشاكل الاجتماعية والصحية والاقتصادية وتم تقسيم التحليل الإحصائي إلى التحليل احادي وثنائي المتغيرات وبينى على حزمة من المتغيرات أو العوامل من الأساليب الإحصائية للتحليل متعدد المتغيرات هو أسلوب التحليل التميزي والذي يشاع استخدامه في المجالات الطبية حيث يهتم التحليل التميزي بكيفية التمييز بين مجموعتين أو أكثر من الأفراد أو الأشياء وتصنيف المفردات الجديدة على المجموعات التي سبق تعریفها ويعتمد أسلوب تحليل التمايز على الوصول إلى دالة تسمى دالة التمايز تعمل على زيادة الفروق بين متوسط المجموعات حيث كلما كان هناك تباعد بين متوسط المجموعات كلما كان التمييز كفاء وبالناتي يقل خطأ التصنيف ويعتبر التحليل التميزي بين مجموعتين أو أكثر من الأفراد أو الأشياء وتصنيف المفردات الجديدة على المجموعات التي سبق تعریفها ويعتبر التحليل التميزي استكشافياً بطبعته حيث يكتشف أسباب الاختلاف المشاهدة عندما لا تستطيع فهم العلاقات المسببة بدرجة كافية الدقة. (الجاعوني، غانم: ٢٠٠٧)

أهمية التحليل التميزي:

ترجع أهمية التحليل التميزي كأحد أساليب التحليل متعدد المتغيرات إلى مقدرته في التمييز بين مجموعتين أو أكثر من خلال مجموعة من المتغيرات ويتم ذلك بإنشاء دوال تمايز "Discriminate Function" تعمل على تعظيم الاختلاف أو الفروق بين المجموعات بأقل خطأ للتصنيف.

أنواع التحليل التميزي:

هناك ثلاث أنواع من التحليل التميزي تتمثل في: (النويري: ٢٠١٣)

- التحليل التميزي المباشر Direct discriminate analysis: حيث تدخل المتغيرات إلى التحليل دفعه واحدة دون إعطاء أي أهمية لأي متغير.
- التحليل التميزي الهرمي Hierarchical discriminant analysis: يتم فيها إدخال المتغيرات حسب رؤية الباحث.
- التحليل التميزي المتدرج Stepwise discriminant analysis: يتم إدخال المتغيرات للتحليل حسب معيار إحصائي يحدد أولوية إدخال المتغيرات إلى النموذج حيث يتم إضافة المتغيرات إلى الدوال التميزية واحد تلو الآخر حتى نجد أن إضافة متغيرات لا يعطي تمييزاً أفضل.

أهداف التحليل التميزي: -

هناك عدة أهداف للتحليل التميزي أهمها: -

- إنشاء دوال تميزية للفصل أو التمييز بين فئات المتغير التابع.

- تعمل هذه الدول على تعظيم الفروق بين المجموعات (فئات المتغير التابع).
- ترتيب المتغيرات التي تسهم بقدر كبير في التمييز أو توضيح الاختلافات بين المجموعات (فئات المتغير التابع).
- تصنيف المشاهدات الجديدة وتوزيعها على المجموعات (فئات المتغير التابع).
- الوصول إلى أقل نسبة خطأ للتوصيف - تقييم دقة التصنيف كنسبة مئوية.

شروط التحليل التميزي:

- ١- عدم تساوي متوسطات المجموعات (فئات المتغير التابع).
- ٢- تساوي مصفوفة التباين والتغاير بين المجموعتين.
- ٣- ان تكون المجموعات منفصلة وقابلة للتحديد.
- ٤- ان توزع المتغيرات التابعة والكمية توزيعاً طبيعياً.
- ٥- العينة تختار عشوائياً.
- ٦- استقلال المشاهدات؛ أي عدم وجود ارتباط بين المتغيرات المستخدمة في الدراسة أو ما يعرف بمشكلة حيث *Multicollinearity* حيث كلما كان هناك ارتباط بين المتغيرات كلما كان هناك صعوبة في تفسير نتائج تحليل التمايز وذلك صعوبة في تحديد المساهمة النسبية لكل متغير على حدة.
- ٧- عدم وجود قيمة متطرفة حيث أن تحليل التمايز أكثر حساسية وتتأثر بالقيم الشاذة ووجودها يبعد توزيع البيانات عن التوزيع الطبيعي.

الدالة التمييزية :Discriminate Function

تقوم الدالة التمييزية على فكرة أساسية وهي تقسيم الأشخاص إلى مجموعتين هما (مصاب أو غير مصاب) وذلك بالاعتماد على مجموعة من المتغيرات أو العوامل وتعمل الدالة على زيادة درجة التجانس بين مفردات المجموعة الواحدة وتقليل درجة التجانس بين المجموعتين وبالتالي تسهيل إمكانية تصنيف أي مشاهدة جديدة إلى إحدى المجموعتين بأقل خطأ للتصنيف كما تعمل الدالة على استبعاد المتغيرات التي ليس لها تأثير معنوي في التمييز والفصل بين المجموعتين.

ويتم حساب الدالة التمييزية كالتالي: -

في حالة تعدد المجموعات تتعدد الدوال التمييزية ولكننا سنقتصر على الدالة التمييزية بين مجموعتين فقط.

أولاً: - حساب متوسطات المتغيرات في كل مجموعة وإيجاد الفرق بين متوسط: -

$$\bar{x}_{l(1)} = \begin{bmatrix} \bar{x}_1(1) \\ \bar{x}_2(1) \\ \vdots \\ \bar{x}_k(1) \end{bmatrix}$$

متوسطات المتغيرات في المجموعة الثانية: -

$$\bar{x}_{i(2)} = \begin{bmatrix} \bar{x}_{1(2)} \\ \bar{x}_{2(2)} \\ \vdots \\ \bar{x}_{k(2)} \end{bmatrix}$$

k عدد المتغيرات المستقلة

الفرق بين متوسط المتغير في المجموعتين:

$$(المسافة) d_i = \bar{x}_{i(1)} - \bar{x}_{i(2)} = \begin{bmatrix} \bar{x}_{11} - \bar{x}_{12} \\ \bar{x}_{21} - \bar{x}_{22} \\ \vdots & \vdots \\ \bar{x}_{k(1)} - \bar{x}_{k(2)} \end{bmatrix} = \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_k \end{bmatrix}$$

ثانياً: إيجاد التباين والتغير المشترك بين المجموعتين: -

$$S_{ii} = \sum x_i^2 - \frac{(\sum x_i)^2}{n}$$

$$S_{ij} = \sum x_i x_j - \frac{\sum x_i \sum x_j}{n}$$

\therefore التباين المشترك

$$V_{ii} = \frac{S_{ii} + S_{ii(2)}}{n_1 + n_2 - 2}$$

\therefore التغير المشترك

$$V_{ij} = \frac{S_{ij(1)} + S_{ij(2)}}{n_1 + n_2 - 2}$$

مصفوفة التباين والتغير المشترك بين المجموعتين.

$$V = \begin{bmatrix} V_{11} V_{12} V_{13} \dots \dots \dots V_{1k} \\ V_{21} V_{22} V_{23} \dots \dots \dots V_{2k} \\ \vdots & \vdots & \vdots \\ V_{k1} V_{k2} V_{k3} \dots \dots \dots V_{kk} \end{bmatrix}$$

وهي عبارة عن مصفوفة مربعة ومتماثلة والقطر الرئيسي لها يمثل التباين المشترك وبقي العناصر التغير المشترك.

بناء الدالة التمييزية:

تأخذ الدالة التمييزية بمعاملات معيارية الشكل التالي:

$$\hat{L} = \hat{\alpha}_1 x_1 + \hat{\alpha}_2 x_2 + \dots + \hat{\alpha}_k x_k$$

حيث

$$\hat{\alpha} = V^{-1} d$$

- 1

$$\begin{bmatrix} \hat{\alpha}_1 \\ \hat{\alpha}_2 \\ \vdots \\ \hat{\alpha}_k \end{bmatrix} = \begin{bmatrix} V_{11} V_{12} \dots \dots \dots V_{1k} \\ V_{21} V_{22} \dots \dots \dots V_{2k} \\ \vdots & \vdots & \vdots \\ V_{k1} V_{k2} \dots \dots \dots V_{kk} \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_k \end{bmatrix}$$

حيث

$\hat{\alpha}$ معاملات الدالة التمييزية المعيارية.

v^{-1} : معكوس مصفوفة التباين والتغير المشترك.

d_i : مصفوفة المسافة بين متوسط المتغيرات في كلا المجموعتين.

- الأهمية النسبية للعامل المؤثرة (المتغيرات المستقلة):

بعد قيام التحليل التميizi بإنشاء وتكونن الدوال التمييزية تظهر له ميزة إضافية وهي تحديد

الأهمية النسبية للمتغيرات المستقلة والممؤثرة في عملية التمييز والفصل بين المجموعات وترتيبها ويتم ذلك من

خلال استبعاد إشارات المعاملات المعيارية لدالة التمييز وصاحب أعلى قيمة هو الأكثر أهمية أما عن نسبة

المُساهمة في عملية التمييز تحدد من خلال مُعامل الارتباط القانوني "Canonical correlation"

اختبارات الدالة التمييزية:

لاختبار قدرة الدالة على التمييز والفصل بين المجموعات تستخدم الاختبارات الآتية: -

١ - اختبار F (F test)

ونذلك لاختبار قدرة الدالة على التمييز وعن طريق الفرضية التي تنص على ان الدالة ليس لديها

القدرة على التمييز (H_0) ضد الدالة لديها القدرة على التمييز (H_1) ويعتمد هذا الاختبار على قياس

الاختلافات بين المجموعات وداخل المجموعات بين المفردات ويتم ذلك من خلال تكون جدول تحليل التباين

التالي: -

Source	SS	Df	Ms	F
بين المجموعات Between x's	SSB	k-1	M _{SB}	M _{SB}
الخطأ Within x's	SSE	n-k	M _{SE}	M _{SE}
الكلي Total	SST	n-1		

حيث ان: -

١ - مجموع مربعات الأخطاء يحسب كالتالي: -

$$SSE = D^2 = \hat{\alpha}_1 d_1 + \hat{\alpha}_2 d_2 + \dots + \hat{\alpha}_k d_k$$

٢ - مجموع مربعات بين المتغيرات: -

$$SSB = \frac{n_1 n_2}{(n_1 + n_2)(n_1 + n_2 - 2)} \times (D^2)^2$$

٣ - مجموع مربعات الكلي: -

$$SST = SSB + SSE$$

و يتم الاختيار كالتالي:

١ - صياغة الفروض:

الدالة ليس لها قدرة على التمييز: H_0

الدالة لها القدرة على التمييز: H_1

٢ - القيمة المحسوبة:

$$F = \frac{M_{SB}}{M_{SE}}$$

٣ - القيمة الجدولية:

$$F(k-1, n-k)$$

٤ - القرار:

إذا كانت F المحسوبة أكبر من F الجدولية ترفض الفرض العدمي ونقبل بالفرض البديل ويكون للدالة قدرة عالية على التمييز والعكس صحيح.

٥ - اختبار ويلكس لامدا (Wilks Lambda) (A):

تأخذ الفروض الشكل الآتي:

الدالة ليس لها مقدرة على التمييز

$$H_0: \mu_1 = \mu_2$$

الدالة لها القدرة على التمييز

القيمة المحسوبة

$$\Lambda = \prod_{i=1}^k \frac{1}{1+\lambda_i}$$

λ الجذر الكامن (eigenvalues) لكل المتغيرات.

k عدد المتغيرات

القرار: تتحصر قيمة

$$0 \leq \Lambda \leq 1$$

إذا كان

. معناها تساوي متواسطات المجموعتين وبالتالي عدم مقدرة الدالة على التمييز والفصل.

. معناها عدم تساوي متواسطات المجموعتين والدالة لها القدرة عالية على التمييز.

إذا اقتربت قيمة Λ من الواحد دليل على عدم مقدرة الدالة على التمييز وإذا اقتربت من الصفر دليل على قدرة الدالة على التمييز.

وستستخدم إحصائية "ويلكس لمدا" لاختبار معنوية المتغيرات الداخلية في النموذج حيث يتم الإبقاء على المتغيرات لها أدنى قيمة لإحصائية Wilk's Lambda وأعلى قيمة لـ F .

٣- اختبار هوتلنج - Lawley test (T^2)

إحصاء هوتلنج تأخذ الشكل الآتي:

$$T^2 = \sum_{i=1}^s \lambda_i$$

حيث أن

λ : eigenvalues الجذور المميزة للمتغيرات

s : عدد المتغيرات

وتعادل إحصائية هوتلنج قيمة F من جدول تحليل التباين ويمكن تحويله إلى قيمة لها توزيع F تقريبي صيغته كالتالي:

$$F = \frac{n_1 + n_2 - k - 1}{(n_1 + n_2 - 2)k} * T^2$$

والقيمة الجدولية:

$$F_{\alpha}(k - 1, n_1 + n_2 - k - 1)$$

إذا كانت F المحسوبة أكبر من F الجدولية رفض الفرض العدلي وقبول البديل بين للدالة قدرة عالية على التمييز.

نقطة الفصل (القطع) :Cut Of Point

بعد تكوين الدالة التمييزية واختبار قدرتها على التمييز والفصل بين المجموعتين يبدأ الاستخدام الثاني لها وهو كيفية تصنيف المشاهدة الجديدة إلى أي المجموعتين تتبع ويتم ذلك من خلال الخطوات الآتية:

١- تحديد نقطة الفصل وهي تمثل متوسط المتوسطين:

$$\bar{\bar{L}} = \frac{\bar{L}_{(1)} + \bar{L}_{(2)}}{2}$$

حيث ان

$\bar{\bar{L}}$: نقطة الفصل.

$\bar{L}_{(1)}$ متوسط القيم التمييزية للمجموعة الأولى.

$\bar{L}_{(2)}$ متوسط القيم التمييزية للمجموعة الثانية.

قاعدة التصنيف :Classification Rule

من خلال هذه القاعدة يمكن تصنيف أو التبع بانتماء مفردة جديدة لإحدى المجموعتين بأقل خطأ
تصنف على النحو التالي:

١) إذا كان $\bar{L}_{(1)} > \bar{L}_{(2)}$

وإذا كانت القيمة التمييزية للمفردة الجديدة أكبر من نقطة الفصل تصنف ضمن المجموعة الأولى
وإذا كانت القيمة التمييزية للمفردة الجديدة أقل من نقطة الفصل تصنف ضمن المجموعة الثانية وإذا ساوت
نقطة الفصل تصنف عشوائياً ضمن أي مجموعة من المجموعتين.

٢) إذا كان $\bar{L}_{(1)} < \bar{L}_{(2)}$

وإذا كانت القيمة التمييزية للمفردة الجديدة أعلى من نقطة الفصل تصنف ضمن المجموعتين
وإذا كانت أقل تصنف ضمن المجموعة الأولى وإذا ساوت معها تصنف عشوائياً ضمن أي مجموعة في
المجموعتين.

أخطاء التصنيف:

يقصد بأخطاء التصنيف وضع المفردة في مجموعة غير مناسبة لها أي وضع مفردة في مجموعة ما ولكن
هي تتبع لمجموعة أخرى ويغير خطأ التصنيف عامل مهم عند الحكم على كفاءة الدالة التمييزية.
هناك نوعان من أخطاء التصنيف هما:

١- خطأ التصنيف الظاهري.

ويحسب من جدول التصنيف التالي.

المجموعة	تابع المجموعة الأولى (١)	تابع المجموعة الثانية (٢)	مجموع
الأولى (١)	n_{11}	n_{12}	n_1
الثانية (٢)	n_{21}	n_{22}	n_2

n_{11} : عدد المفردات من المجموعة الأولى والتي تم تصنيفها في نفس المجموعة وبالتالي هي صنفت بطريقة
صحيحة.

n_{12} : عدد المفردات من المجموعة الأولى والتي تم تصنيفها خطأ في المجموعة الثانية.

n_{21} : عدد المفردات التي تتبع للأصل إلى المجموعة الثانية وتم تصنيفها خطأ في المجموعة الأولى.

n_{22} : عدد المفردات في المجموعة الثانية التي تم تصنيفها في نفس المجموعة وبالتالي هي صنفت بطريقة
صحيحة.

ويحسب الخطأ الظاهري كما يلي:

$$P_{12} = \frac{n_{12}}{n_1}$$

P_{12} نسبة المفردات التي تتبع للمجموعة الأولى وصنفت خطأ للثانية.

$$P_{21} = \frac{n_{21}}{n_2}$$

P_{21} نسبة المفردات التي تتبع للمجموعة الثانية وصنفت خطأ في الأولى. ويمكن حساب معدل الخطأ
الظاهري باستخدام المعادلة

$$\frac{n_{12} + n_{21}}{n_1 + n_2}$$

٢- الخطأ الحقيقي: يمثل نسبة التصنيف الخاطئ في المجتمع:

$$P_{12} = P_{21} = F \left[\frac{-\sqrt{D^2}}{2} \right]$$

حيث F دالة التوزيع الطبيعي المعياري، D إحصائية Mahalanobis. تحسب القيمة بين القوسين ويحسب الاحتمال المقابل لها من جدول التوزيع الطبيعي المعياري وكلما اقترب الاحتمال من الصفر دل على صنف وانخفاض خط التوصيف وبالتالي قدرة الدالة على التمييز والتصنيف أما إذا كان الاحتمال قريب من الواحد يدل على ارتفاع خط التوصيف وانخفاض قدره الدالة على التمييز والتصنيف.
الدالة التمييزية بمعاملات غير معيارية: -

تأخذ الشكل التالي:

$$y = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_k x_k$$

y : الدالة التمييزية غير المعيارية.

b_0 : ثابت التمييز.

b_n 's: معاملات التمييز غير المعيارية.

x 's: المتغيرات غير المعيارية.

وللحكم على جودة النموذج التمييزي من خلال مُعامل الارتباط القانوني Canonical correlation حيث ان القيم المرتفعة لمعامل الارتباط القانوني تكون مؤشر لجودة التوفيق العالي للنموذج التمييزي وبتربيع قيمة مُعامل الارتباط القانوني تحصل على قيمة معامل التحديد " R^2 " الذي يحدد نسبة مُساهمة المتغيرات المستقلة في التمييز والتصنيف.

المبحث الثاني الجانب التطبيقي

تمهيد:

بالرغم من الدور الذي يلعبه التحليل احادي المتغير أو ثانوي المتغير في تفسير وتحليل كثير من الظواهر الاقتصادية والاجتماعية والطبية إلا أنه عندما يتعلق الأمر بعدد كبير من المتغيرات فلا بد من اللجوء إلى التحليل متعدد المتغيرات ومن أهم أساليب التحليل متعدد المتغيرات والشائع استخدامه في المجالات الطبية هو أسلوب التحليل التمييزي والذي يقوم بدوره بالتمييز وفصل الأشخاص إلى مجموعتين رئيسيتين هما أما (مصاب أو غير مصاب) بمرض السكري وذلك على عينة من ٣٥٠ شخص منهم ١٦٨ مصاب و ١٨٢ غير مصاب بهدف الوصول إلى دالة تمييزية من خلالها يتم تصنيف الأشخاص أو المشاهدات الجديدة على احدى المجموعتين بناءً على فرضيات معينة.

١- متغيرات البحث:

تتمثل متغيرات البحث في متغير تابع نوعي ثانوي القيمة (ع) غير مصاب (٠) ومصاب (١) ومجموعة من العوامل المؤثرة (المتغيرات المستقلة) وهي.

١. الوراثة (x_1) تأخذ (٠ لا يوجد، ١ يوجد).

٢. الوزن. (x_2)

٣. ضغط الدم (x_3) تأخذ (٠ طبيعي، ١ مرتفع).

٤. الغمر (x_4) .

٥. النوع (x_5) تأخذ (٠ أنثى، ١ ذكر).

٦. التدخين (x_6) تأخذ (٠ لا يدخن، ١ مدخن).

٧. ممارسة الرياضة (x_7) تأخذ (٠ لا يمارس، يمارس رياضة).

٨. مرض الثغرس (x_8) تأخذ (٠ لا يوجد، ١ يوجد).

٩. الكوليسترونول (x_9) تأخذ (٠ لا يوجد، ١ يوجد).

١٠. الحالة الاجتماعية (x_{10}) تأخذ (٠ أعزب، ١ متزوج).

١١. أمراض القلب والكلري (x_{11}) تأخذ (٠ لا يوجد، ١ يوجد).

ولاستخدام تحليل التمييزي مجموعة من الافتراضات لابد من توافرها وهي:

١- اختبار التوزيع الطبيعي للبيانات: -

نظرًا أن حجم العينة يزيد عن ٣٠ مفردة طبقاً لنظرية النهاية المركزية فإن البيانات تتبع التوزيع الطبيعي ولا داعي لإجراء اختبار الطبيعي.

٢- اختبار تساوي متوسطي المجموعتين.

بالنظر إلى الجدول رقم (١) التالي:

(١)

Sig	المجموعة		المتغير
	الثاني (غير المصابين) (١)	الأولى (المصابين) (١)	
	المتوسط	المتوسط	
0.00	0.3352	0.7381	x_1
0.00	73.7692	100.3274	x_2
0.00	0.3187	0.7202	x_3
0.00	38.2363	48.333	x_4
0.014	0.5385	0.6667	x_5
0.00	0.6154	0.3512	x_6
0.00	0.6044	0.3512	x_7
0.891	0.5549	0.5476	x_8
0.00	0.4121	0.7143	x_9
0.011	0.6099	0.7381	x_{10}
0.00	0.2802	0.5179	x_{11}

من الملاحظ من خلال جدول (١) ان قيمة Sig أقل من ٠.٠٥ وبالتالي معنوية الفرق بين متوسطي كل متغير في المجموعتين أو خلال فئات المتغير التابع فيما عدا المتغير الثامن.
كما ظلّاحظ أن المتosteatas الأعلى في المجموعة الأولى للمصابين في كلٍّ من x_2, x_4, x_6, x_9 ، وهما (الوزن والعمر والوراثة والكوليسترون وضغط الدم وذلك أمر بديهي فإن المصابين بمرض السكري نسبة الضغط العالي والوزن الزائد والعمر المتقدم أما المتosteatas الأعلى في المجموعة الثانية وهي مجموعة غير المصابين تتبع في x_2, x_4, x_6 .
ومن خلال إحصائية ويلكس لاما Wilks's lambda distribution التالي:

Wilks' Lambda

Test of Function(s)	Wilks' Lambda	Chi-square	Df	Sig.
1	.365	346.661	8	.000

حيث ان الفروض تصاغ:

$$\begin{aligned} H_0 \mu_1 &= \mu_2 \\ H_1 \mu_1 &\neq \mu_2 \end{aligned}$$

حيث Sig (000) أقل من ٠.٠٥ رفض الفرض العدلي وقبول الفرض البديل وبالتالي هناك اختلاف بين متوسطي المجموعتين كما ان إحصاء ويلكس لاما تساوي ٣٦٥ وهي تقترب من الصفر دليل على وجود اختلاف بين متوسطي المجموعتين وهذا يعني أن الدالة التمييزية لديها القدرة على التمييز وتصنيف المشاهدات إلى مجتمعها الحقيقي.

- اختبار فرضية تجانس التباين بين المجموعتين:

حيث تصاغ الفروض الإحصائية كالتالي:

$$\begin{aligned} H_0 \Sigma_1 &= \Sigma_2 \\ H_1 \Sigma_1 &\neq \Sigma_2 \end{aligned}$$

ويستخدم اختبار Box's M كانت النتائج كالتالي:

Log Determinants (٣)

Y	Rank	Log Determinant
N	8	1.230
Y	8	.976
Pooled within-groups	8	1.244

Test Results (٤)

Box's M	47.341
F	1.283
df1	36
df2	402161.071
Sig.	0.119

ومن الملاحظ أن قيمة Sig (0.119) أكبر من 0.05 وبالتالي قبول الفرض العدلي بتساوي مصفوفة التباين والتغير للمجموعتين وبالتالي تتحقق افتراض تجانس التباين بين المجموعتين.

٤- اختبار معنوية العوامل المؤثرة (المتغيرات المستقلة) في النموذج التميزي:

تم اختبار معنوية جميع العوامل المؤثرة في النموذج التميزي لمعرفة أهمية كل متغير ومدى إسهامه في عملية التمييز والتصنيف وكانت كالتالي (٥)

Tests of Equality Group of Means (٥)

	Wilks' Lambda	F	df1	df2	Sig.
x1	.857	58.046	1	348	.000
x2	.477	381.840	1	348	.000
x3	.839	66.816	1	348	.000
x4	.873	50.479	1	348	.000
x5	.983	6.052	1	348	.014
x6	.930	26.082	1	348	.000
x7	.936	23.829	1	348	.000
x8	1.000	.019	1	348	.891
x9	.908	35.422	1	348	.000
x10	.981	6.594	1	348	.011
x11	.941	21.827	1	348	.000

ومن الملاحظ أن جميع المتغيرات تتمتع بمعنى عالي حيث ان Sig (000) أقل من 0.05 لجميع المتغيرات ماعدا المتغير الثامن وذلك يدل على المتغيرات لها تأثير معنوي كبير في عملية التمييز بين المجموعتين ومن ثم توصلنا إلى النموذج التحليل التميزي مناسب لبيانات مرضى السكري.

١) تكوين الدالة التمييزية:

لإنشاء الدالة التمييزية تحدد أولاً المتغيرات الداخلة في تكوين الدالة التمييزية حيث ان اختزال عدد المتغيرات في نموذج التمييز يفيد في قياس المتغيرات ذات العلاقة المعنوية وذات الصلة الأكبر بالموضوع محل الدراسة وللتعرف على المتغيرات ذات القوة التمييزية المعنوية والتي تغطي أقل خطأ تصنيف (James, 1985). وهناك عدة معايير إحصائية للإدخال والحذف وهي الإبقاء على التمييز صاحب القيمة الأكبر لـ F وأقل قيمة إحصائية ويلكس لاما Wilks's lambda كما في الجدول الآتي:

Variables Entered/Removed^{a,b,c,d} (٥)

Step	Entered	Wilks' Lambda				Exact F				Sig.
		Statistic	df1	df2	df3	Statistic	df1	df2		
1	x2	.477	1	1	348.000	381.840	1	348.000	.000	
2	x4	.441	2	1	348.000	219.946	2	347.000	.000	
3	x3	.417	3	1	348.000	160.993	3	346.000	.000	
4	x1	.400	4	1	348.000	129.221	4	345.000	.000	
5	x7	.390	5	1	348.000	107.676	5	344.000	.000	
6	x5	.380	6	1	348.000	93.114	6	343.000	.000	
7	x6	.370	7	1	348.000	83.038	7	342.000	.000	
8	x9	.365	8	1	348.000	74.142	8	341.000	.000	

المصدر: SPSS V23

وبلغ من الجدول أنه تم استبعاد ثلاثة متغيرات وتم الإبقاء على ثمان متغيرات التي لها قدرة أعلى في التمييز والفصل بين المجموعتين المصابين وغير المصابين والتي لها أعلى قيمة F وأقل قيمة Wilks's lambda وتم الاختيار بناء على الاختيار التدريجي على ٨ خطوات وبالتالي المتغيرات الداخلة للنموذج هي $x_2, x_4, x_3, x_1, x_7, x_5, x_6, x_9$.

إيجاد الدالة التمييزية:-

أولاً: من جدول رقم (١) تم إيجاد متوسطات المتغيرات في كل المجموعتين.

ثانياً: إيجاد الفرق بين متوسط كل متغير في كل المجموعتين.

$$d = \bar{x}_{i(1)} - \bar{x}_{i(2)} = \begin{bmatrix} \bar{x}_{11} - \bar{x}_{12} \\ \bar{x}_{21} - \bar{x}_{22} \\ \vdots \\ \bar{x}_{9(1)} - \bar{x}_{9(2)} \end{bmatrix} = \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_9 \end{bmatrix}$$

ثالثاً: مصفوفة التباين والتغيير بين المجموعتين.

$$V = \begin{bmatrix} x_1 & x_2 & x_3 & x_4 & x_5 & x_6 & x_7 & x_9 \\ x_1 & 0.284 & 2.969 & 0.50 & 1.446 & 0.049 & 0 & -0.41 & -0.41 & 0.57 \\ x_2 & 2.969 & 337.465 & 3.136 & 62.014 & 0.293 & 0 & -1.640 & -1.528 & 2.410 \\ x_3 & 0.050 & 3.136 & 0.251 & 0.960 & -0.007 & -0.047 & -0.024 & 0.035 \\ x_4 & 1.446 & 62.014 & 201.452 & 201.45 & -0.348 & -1.571 & -0.814 & 0.916 \\ x_5 & 0.049 & 0.293 & -0.007 & -0.348 & 0.241 & 0.053 & -0.013 & -0.014 \\ x_6 & -0.041 & -1.640 & -0.047 & -1.571 & 0.053 & 0.251 & 0.016 & -0.041 \\ x_7 & -0.041 & -1.528 & -0.024 & -0.814 & -0.013 & 0.016 & 0.250 & -0.035 \\ x_9 & 0.057 & 2.410 & 0.035 & 0.914 & -0.014 & -0.041 & -0.035 & 0.247 \end{bmatrix}$$

وبذلك الدالة التمييزية بمعاملات معيارية: -

$$L = \hat{\alpha}_1 x_1 + \hat{\alpha}_2 x_2 + \hat{\alpha}_3 x_3 + \dots + \hat{\alpha}_9 x_9$$

حيث ان

$$\hat{\alpha} = v^T d$$

$$\hat{\alpha} = \begin{bmatrix} 0.160 \\ 0.797 \\ 0.267 \\ 0.295 \\ 0.268 \\ -0.203 \\ -0.187 \\ 0.154 \end{bmatrix}$$

$$\hat{L} = 0.160 x_1 + 0.797 x_2 + 0.267 x_3 + 0.295 x_4 + 0.268 x_5 - 0.203 x_6 - 0.187 x_7 +$$

$$0.154 x_9$$

ولتحديد الأهمية النسبية للعوامل المؤثرة ونسبة المساهمة في التمييز والتنبؤ في التمودج التميزي:

لتحديد أكثر العوامل أثر على مستوى الإصابة ومساهمة العامل في التمييز والتصنيف مكان كالتالي:

العامل	المتغير	الأهمية النسبية (معامل الارتباط القانوني التميزي)
0.797	x_2	0.794
0.295	x_4	0.289
0.268	x_5	0.10
0.267	x_3	0.332
-0.203	x_6	-0.208
-0.187	x_7	-0.198
0.16	x_1	0.310
0.154	x_9	0.242

ولمعرفة أهم العوامل المؤثرة تتظر لعمود المعاملات المعيارية ($\hat{\alpha}_i$) حيث انه القيمة

المطلقة الكبيرة يقابلها العامل الأكثر أهمية في التأثير على الإصابة وتكون هذه الأهمية موجبة أو

سالبة.

للحظ أن أكثر المتغيرات أهمية (x_2) الوزن ثم (x_4) العمر ثم (x_3) النوع ثم (x_6) ضغط الدم ثم (x_9) التدخين ثم (x_1) ممارسة رياضة ثم (x_7) الوراثة ثم (x_5) الكوليسترول.

صيغة المعاشرة:

صيغة المعاشرة:

أما الأهمية النسبية فإن (2x) الوزن يساهم بنسبة أكبر في عملية تمييز المجموعتين %٦٧٩,٤ يليه ضغط الدم %٣٣,٢ والوراثة %٣١ العمر %٢٨,٩ وكوليسترون %٢٤,٢ والتدخين %٢٠,٨ وممارسة الرياضية %١٩,٨ وأخيراً النوع %.١٠.

ولاختبار قدرة الدالة على التمييز:

١) باستخدام جدول تحليل التباين وختبار F.

• ومن خلال الفروض الآتية:

الدالة ليس لها قدرة على التمييز H_0

الدالة لها القدرة على التمييز H_1

• إيجاد قيمة مجموع مربعات الأخطاء (داخل المجموعات).

$$SSE = D^2 = \alpha_1 d_1 + \alpha_2 d_2 + \alpha_3 d_3 + \alpha_4 d_4 + \alpha_5 d_5 + \alpha_6 d_6 + \alpha_7 d_7 + \alpha_9 d_9$$

$$= [0.160 \ 0.797 \ 0.267 \ 0.295 \ 0.268 - 0.203 - 0.187 \ 0.154]$$

$$\begin{bmatrix} 0.4029 \\ 26.55 \\ 0.4015 \\ 10.09 \\ 0.1282 \\ -0.2642 \\ -0.2532 \\ 0.3022 \end{bmatrix}$$

$$SSE = 24.50$$

مجموع مربعات بين المتغيرات

$$SSB = \frac{n_1 n_2}{(n_1+n_2)(n_1+n_2)} \times (D^2)^2 = \frac{168 \times 182}{(168+182)(168+182-2)} * 24.5^2 = 150.68$$

مجموع المربعات الكلية:

$$SST = SSB + SSE = 150.68 + 24.5 = 175.18$$

جدول تحليل التباين

Source	Ss	Df	Ms	F
بين المجموعات	150.68	k-1 7	٢١,٥٢٥	
داخل المجموعات	24.5	n-k 342	٠,٠٧٢	٣٠٠,٥
الكلي	175.18	n-1 349		

القيمة المحسوبة:

$$F = 300.5$$

القيمة الجدولية:

$$F_{0.05}(7, 342) = 1.40$$

F القيمة المحسوبة أكبر F الجدولية . رفض الفرض العدلي وقبول الفرض البديل فإن للدالة قدرة عالية على التمييز والفعل بين المجموعتين.

٢) اختبار ويلكس لاماda :Wilks' lambda

تصاغ الفروض كالتالي:

$$H_0 \mu_1 = \mu_2 \quad \text{الدالة ليس لها قدرة على التمييز}$$

$$H_1 \mu_1 \neq \mu_2 \quad \text{الدالة لها القدرة على التمييز}$$

Wilks' Lambda

Test of Function(s)	Wilks' Lambda	Chi-square	Df	Sig.
1	.365	346.661	8	.000

أولاً: قيمة إحصائية ويلكس لاماda . وهي أقرب للصفر وذلك دليل على القدرة العالية للدالة على التمييز. وكما نلاحظ أن Sig (00) أقل من .٠٠٥ وبالتالي رفض الفرض العدلي وقبول الفرض البديل فإن للدالة قدرة على التمييز والفصل بين المجموعتين. كما أن في جدول:

Function	Eigenvalue	% of Variance	Cumulative %	Canonical Correlation
1	1.739 ^a	100.0	100.0	.797

ومن الملاحظ قيمة الجذر الكامن $\lambda = 1,739$ وتشير إلى أن نسبة التباين المفسر بين مجموعتي المصابين وغير المصابين والتي تعود للفرق بينها في النموذج التميزي الوحيدة وجمع قيمة معامل الارتباط القانوني .، .٧٩٧ معامل الارتباط بين مجموعة العوامل المؤثرة ونموذج التمييز الوحيد وبترتيب هذه القيمة تحصل على ٦٣,٥ % وهذا يعني نسبة مُساهمة العوامل المؤثرة في التباين والاختلاف في التمييز بين المجموعتين.

للاستخدام الثاني للنموذج التميزي وهو التصنيف فكانت النتائج كالتالي:

	غير مصاب (٠)	مصاب (١)	مجموع
غير مصاب (٠)	١٦٩	١٣	١٨٢
مصاب (١)	٢٠	١٤٨	١٦٨

٦٠,٩٦ % النسبة الإجمالية.

حيث ان النموذج التمييزي الذي يتكون من ثمان متغيرات هما $x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_9$ ، قام بالتصنيف الصحيح ١٦٩ مفردة من غير المصابين بمرض السكري وتصنيف غير صحيح ١٣ أي نسبة تصنيف صحيح ٦٩٢٪ في الحالات بشكل صحيح. أما بالنسبة للمصابين قام تصنيف الصحيح ١٤٨ مفردة وغير الصحيح ٢٠ مفردة بنسبة إجمالية صحيحة ٨٨,١٪ وإن الدقة الإجمالية للتصنيف هي ٩٠,٦٪ بخطأ ٤٪.

نسبة الخطأ الظاهري:

للمجموعة الأولى:

$$p_{12} = \frac{n_{12}}{n_1} = \frac{13}{182} = 0.0714$$

للمجموعة الثانية:

$$p_{21} = \frac{n_{21}}{n_1} = \frac{20}{168} = 0.119$$

تقدير النموذج التمييزي بمعاملات غير معيارية:

يتم تقدير النموذج التمييزي بمعاملات غير معيارية كما يلي:

المتغير	b
x1	0.324
x2	0.063
x3	0.581
x4	0.022
x5	0.551
x6	-0.420
x7	-0.386
x9	0.325
(Constant)	-6.974

$$\hat{y} = -6.974 + 0.324x_1 + 0.063x_2 + 0.581x_3 + 0.022x_4 + 0.551x_5 - 0.420x_6 - 0.386x_7 + 0.325x_9$$

وهذا النموذج فعال وقدر على التصنيف الصحيح للمرضى بنسبة ٩٠,٦٪ وأقل خطأ تصنيف ٤٪.
تصنيف المشاهدات الجديدة:

١- إيجاد نقطة الفصل:

حساب متوسط المجموعة الأولى (المصابين)

حساب متوسط المجموعة الثانية (غير المصابين)

$$\bar{y}_1 = 1.369$$

نقطة الفصل

$$\bar{y}_2 = 1.263$$

$$\bar{y} = \frac{1.369 - 1.263}{2} = 0.052$$

$$\bar{y}_1 > \bar{y}_2 \therefore$$

إذا كانت المُفردة الجديدة أكبر في النقطة الفاصلة تصنف المُفردة إلى المجموعة الأولى وإذا كانت أقل في النقطة الفاصلة تصنف إلى المجموعة الثانية.

فمثلاً - إذا كانت هناك سيدة ($x_5 = 0$) لديها مرض السكري وراثة ($x_1 = 1$) والوزن 100 ($x_2 = x_4 = 70$) والعمر ($x_3 = 1$) ولديها ضغط الدم ($x_6 = 0$) ولا يدخن ($x_7 = 0$) ولا تمارس رياضة ($x_8 = 0$) وليس لديها كوليسترونول ($x_9 = 0$)

$$\hat{y} = -6.974 + 0.324(1) + 0.063(100) + 0.581(1) + 0.022(70) + 0.551(0) - 0.420(0) \\ - 0.386(0) + 0.325(0) = 1.771$$

: القيمة التمييزية للمُشاهدة الجديدة أكبر من نقطة الفصل وبالتالي تصنف ضمن المجموعة الأولى (المصابين) لمرض السكري.

وإذا كان هناك رجل ($x_5 = 1$) ليس لديه مرض وراثي ($x_1 = 0$) وزنه ($x_2 = 65$) وضغط الدم عادي ($x_3 = 0$) وعمره ($x_4 = 60$) يدخن ($x_6 = 1$) ولا يمارس رياضة ($x_7 = 0$) وكوليسترونول عالي ($x_9 = 1$) القيمة التمييزية .
لـ.

$$\hat{y} = -6.974 + 0.324(0) + 0.063(65) \\ + 0.581(0) + 0.022(60) + 0.551(1) \\ - 0.420(1) - 0.386(0) + 0.325(1) \\ = -1.103$$

القيمة التمييزية للمُشاهدة أقل من نقطة الفصل وبالتالي تصنف ضمن المجموعة الثانية (غير المصابين). وبالتالي فالنموذج التميزي المقدر بنسبة مُساهمة العوامل المؤثرة فيه ٦٣,٥ % وكفاءة النموذج في النصف ٩٠,٦ % أما الحساسين أي تصنف غير المصاب على أنه غير مصاب تمثل ٩٢,٩ % أما النوعية تصنف المصاب على أنه مصاب بنسبة ٨٨,١ % ونسبة خطأ تصنيف ٤,٩ %.

النتائج والتوصيات:

يهدف البحث إلى استخدام التحليل التميزي كأحد أساليب التحليل متعدد الحدود لتحديد أهم العوامل المؤثرة بالإصابة بمرض السكري وذلك من خلال متغير تابع نوعي (مصاب أو غير مصاب) مجموعة من المتغيرات (العوامل) المستقلة وهي عامل الوراثة والوزن وضغط الدم والعمر والنوع والتدخين وممارسة رياضة ومرض التقرس والكوليسترونول والحالة الاجتماعية وأمراض القلب والكلى.

وتم التوصل إلى النتائج التالية:

١ - بعد التأكيد من توافر افتراضات أسلوب تحليل التمييز وهي شرط طبيعة البيانات وشرط عدم تساوي متواسطات المجموعتين وتساوي مصفوفة التباين والتغاير بين المجموعتين ومعنى غالبية العوامل المؤثرة توصلت إلى

ملائمة أسلوب التحليل التمييزي لبيانات مرض السكري أي يمكن استخدامه في تمييز وتصنيف المفردات

الجديدة إلى مصابين أو غير مصابين وفقاً لمجموعة العوامل المستقلة.

٢ - باختبار معنوية العوامل المؤثرة تم استبعاد ثلاثة عوامل (متغيرات) هي أمراض النقرس، الحالة الاجتماعية وأمراض القلب أما باقي المتغيرات لها معنوية عالية في أسلوب تحليل التمييزي.

٣ - وبالتالي فإن الدالة التمييزية للفصل والتمييز بين المجموعتين بمعاملات معيارية هي.

$$\hat{L} = 0.160x_1 + 0.797x_2 + 0.267x_3 + 0.295x_4 + 0.268x_5 - 0.203x_6 - 0.187x_7 + 0.154x_9$$

٤ - أكثر العوامل المؤثرة وأهمها على الإصابة بمرض السكري هو الوزن ثم العمر ثم ضغط الدم ثم التدخين وممارسة الرياضة والوراثة والكوليسترول.

٥ - أكثر العوامل مُساهمة في التمييز بين المجموعتين هو الوزن بنسبة ٧٩,٤٪ يليه ضغط الدم بنسبة ٦٣,٢٪ يليه الوراثة ٦٣٪ ثم العمر ٢٨,٩٪ والكوليسترول بنسبة ٢٤,٢٪ ثم التدخين بنسبة ٢٠,٨٪ وممارسة الرياضة ١٩,٨٪ وأخيراً النوع بنسبة ١٠٪.

٦ - العوامل المؤثرة في الإصابة بمرض السكري تساهم بنسبة ٦٣,٥٪ من التمييز والتصنيف بين المجموعتين.

٧ - النموذج التمييزي ذو كفاءة عالية التصنيف بنسبة ٩٠,٦٪ وحساسية ٩٢,٩٪ ونوعية ٨٨,١٪.

٨ - نسبة خطأ التصنيف صغيرة ٤٪.

الوصيات:

يوصي الباحث بـ:

١ - التوسيع في استخدام التحليل التمييزي كأحد أساليب التحليل متعدد المتغيرات في المجالات الاقتصادية والاجتماعية.

- ٢- استخدام التحليل التمييزي لتحديد العوامل المؤثرة في الإصابة بمرض السكري مع إضافة مُتغيرات أخرى كالنظام الغذائي وتناول الكحوليات وغيرها.
- ٣- استخدام نموذج الدالة التمييزية في التشخيص المبكر.
- ٤- الوصول للوزن المثالي تقادياً من الإصابة بمرض السكري.
- ٥- الاهتمام بالتحليل الإحصائي وابراز الدور الهام له في الجانب الطبي.
- ٦- استخدام أساليب تحليل مُتعدد المعدات لوس كالانحدار اللوجستي والانحدار المُتعدد وتحليل التباين في تحديد العوامل المؤثرة في الإصابة بمرض السكري ومقارنة النتائج بالتحليل التمييزي.