



# طرق التحليل الإحصائي

## متعدد المتغيرات

---

إعداد

الأستاذ الدكتور زياد رشاد الراوي

نشر بدعم من المعهد العربي للتدريب والبحوث الإحصائية

2017

الطبعة الأولى  
**2017م**  
حقوق الطبع محفوظة

المملكة الأردنية الهاشمية  
رقم الإيداع لدى دائرة المكتبة الوطنية  
(2017 /10 /5553)

❖ يتحمل المؤلف كامل المسؤولية القانونية عن محتوى مصنفه ولا يعبر هذا المصنف عن رأي دائرة المكتبة الوطنية أو أي جهة حكومية أخرى.  
❖ تم إعداد بيانات الفهرسة والتصنيف الأولية من قبل دائرة المكتبة الوطنية

لا يسمح بإعادة إصدار هذا الكتاب أو أي جزء منه أو تخزينه في نطاق استعادة المعلومات أو نقله بأي شكل من الأشكال، دون إذن خطي مسبق.

## فهرس المحتويات

5	تقديم
7	مدخل عام General Introduction
8	تصنيف طرق التحليل متعدد المتغيرات
8	تكوين متغيرات جديدة
9	عندما نبدأ التحليل
10	عمليات المصفوفات Matrix Operations
10	عملية الجمع والطرح في المصفوفات
12	عملية الضرب في المصفوفات
12	حالة ضرب المصفوفة بالقيمة الثابتة c
13	محددة المصفوفة  A  Determinant
15	طريقة لابلاس Laplace
17	بعض النظريات المفيدة في موضوع محددة المصفوفة
17	حالة المصفوفة المثلثية Triangular Matrix
18	معكوس المصفوفة Matrix Inverse
20	حل نظام معادلات خطية Solving system of linear equations
21	طريقة غاوس - جوردن Gouss - Jordan للمعادلات الخطية
24	طريقة الحذف Elimination Method لحل نظام المعادلات الخطية
25	قاعدة كرامير Cramer's rule لحل نظام المعادلات الخطية
26	القيم المميزة والمتجهات المميزة Eigen values & Eigen vectors
31	مستويات قياس المتغيرات Variables Levels of Measurment
31	المقياس الاسمي (Nominal Scale)
32	المقياس الرتبوي (الترتبي) (Ordinal Scale)
32	المقياس الفئوي (الفترى) (Interval Scale)
33	المقياس النسبي (Ratio Scale)
35	تحليل الانحدار والإرتباط المتعدد Multiple Regression & Correlation Analysis
35	أولاً: تحليل الانحدار الخطي البسيط
35	الفرضيات الخاصة بنموذج الانحدار البسيط
36	تقدير دالة الانحدار
37	معامل التحديد Coefficient of Determination
38	ثانياً: نموذج الانحدار الخطي المتعدد
40	تقدير معالم نموذج الانحدار الخطي المتعدد
41	سمات مقدرات طريقة المربعات الصغرى
42	تحليل التباين لنموذج الانحدار الخطي المتعدد
42	إختبار معنوية الانحدار
43	معامل التحديد $R^2$
43	معامل التحديد المصحح ( $adjusted R^2$ )
44	الإرتباط The Correlation
44	معامل الإرتباط الجزئي
47	إختيار أفضل معادلة إنحدار
55	تحليل المركبات الرئيسية (PCA) Principal Components Analysis
56	طبيعة المركبات الرئيسية
57	خطوات الحسابات
59	خواص المركبات الرئيسية
60	بعض إستخدامات المركبات الرئيسية

60	عيوب المركبات الرئيسية
60	تحليل الإنحدار بالمكونات الرئيسية
61	تحقيق التعامدية للمركبات الرئيسية
71	نموذج الإنحدار اللوجستي (LRM) Logistic Regression Model
71	مقدمة
75	ملاحظة مهمة (للمودج التجميعي) Additive Model
76	النموذج الضربي للإحتمالات A Multiplicative Model
76	العلاقة بين الأرجحية Odds والإحتمال Prob.
78	نموذج الإنحدار اللوجستي المتعدد Multiple Logistic Regression Model
80	تحليل النتائج
83	تحليل التباين متعدد المتغيرات (MANOVA) Multivariate Analysis of Variance
95	التحليل المميز Discriminant Analysis (DA)
96	أنواع الدوال التمييزية
99	الإختبارات المستخدمة في التحليل المميز
100	إحتمال خطأ التصنيف: the probability of misclassification
101	طريقة التعويض: Resubstitution Method
104	بعض الطرق اللامعلمية
104	طريقة الرتب
105	الجانب التطبيقي
108	قاعدة التصنيف
115	التحليل العنقودي (CA) Cluster Analysis
116	مقياس المسافة Distance Measures
117	الطريقة الهرمية Hierarchical clustering
118	الطريقة المركزية
118	طريقة الربط الفردي
126	تحليل الارتباط القويم (ارتباط المجموعات) Canonical Correlation Analysis
128	طرق تنفيذ تحليل الارتباط القويم
130	إختبار معنوية الارتباط القويم
131	تحليل نتائج المتغيرات القوية
135	مقاييس أخرى للتحليل
138	التحليل العاملي (FA) Factor Analysis
139	أهداف التحليل العاملي
139	النموذج العاملي Factor Model
140	الفروض الأساسية للتحليل العاملي
142	الإشتراكيات (Communalities) وطرق تقديرها
142	طرق تقدير الإشتراكيات
145	حساب مصفوفة الارتباط

## تقديم

إن التزايد المتسارع والملاحظ في الحاجة لإستخدام الطرق الإحصائية بمستوياتها وجوانبها المختلفة لتشمل الظواهر الحياتية منها والعلمية يصاحبه تطور وتنوع في أبعاد هذه الطرق الإحصائية لضمان شمول جميع هذه الظواهر وبما يناسب نوعية وحجم البيانات المتاحة من جهة، ومستويات وأبعاد التحليل الإحصائي المراد إنجازها من جهة أخرى.

وإنطلاقاً من رسالة المعهد العربي للتدريب والبحوث الإحصائية في العمل المتواصل على تطوير قدرات العاملين في مجال التحليل الإحصائي لأنواع مختلفة من البيانات، يقوم المعهد بتنظيم دورات تدريبية ويشجع على إنجاز دراسات وأبحاث في هذا الإتجاه وفي اتجاهات إحصائية أخرى. وفي هذا الإطار يأتي دعمه لإصدار هذا الكتاب من تأليف الأستاذ الدكتور زياد رشاد الراوي تحت عنوان " طرق التحليل الإحصائي متعدد المتغيرات ".

إن هذه الطرق من التحليل الإحصائي تعتبر القمة في تحليل البيانات لأنها تأخذ في الإعتبار التنوع في البيانات والتعدد في المتغيرات التي تمثلها. ولكونها يغلب عليها التعامل مع المصفوفات فقد ارتأى المؤلف تخصيص فصل موسع حول هذا الموضوع في بداية الكتاب. من ناحية أخرى، ونظراً لتعدد المتغيرات التي تتضمنها طرق التحليل هذه فقد يواجه الباحث حالة تعدد أنماط القياس لهذه المتغيرات، ولتدارك الأمر في هذا الجانب يتضمن الكتاب أيضاً موضوع القياس بمختلف أنواعه (الإسمي، الترتيبي، الفئوي، النسبي) مع توضيح ذلك بالنسبة لجميع المتغيرات التي تتضمنها كل طريقة من طرق التحليل المتعدد والتي تناولها هذا الكتاب .

لا يسعني إلا أن أقدم بالشكر الجزيل للأستاذ الدكتور زياد رشاد الراوي على هذه المساهمة العلمية المتميزة، وأمل أن يكون المعهد قد قدم بذلك مرجعاً هاماً في التحليل الإحصائي متعدد المتغيرات لمختلف الإحصائيين والمهتمين بذات الموضوع.

والله ولي التوفيق

عبد العزيز معلمي  
المدير العام



## مدخل عام General Introduction

إن المقصود بمتعدد المتغيرات هو التعامل مع حالة وجود أكثر من متغير واحد سواءً فيما يتعلق بالمتغيرات (العوامل) التوضيحية (وتسمى بالمستقلة أحياناً Independent) من جهة، أو متغيرات الإستجابة (وتسمى بالمعتمدة Dependent) من جهة أخرى.

والبيانات متعددة المتغيرات تظهر في جميع تفرعات العلوم تقريباً. وفي الغالب نجد أن جميع البيانات التي يتم جمعها من خلال الوحدة التجريبية Experimental Unit وتحليلها من قبل الباحثين يمكن تصنيفها بكونها بيانات متعددة المتغيرات. والمقصود بالوحدة التجريبية هنا هي أي حالة أو عنصر يمكن قياسه أو تقييمه بطريقة ما. والبيانات متعددة المتغيرات تظهر هنا متى ما قام الباحث بقياس أو تقييم أكثر من خاصية أو سمة واحدة لكل وحدة تجريبية. وهذه الخواص أو السمات تسمى عادة بالمتغيرات من قبل الإحصائيين.

وطرق التحليل متعدد المتغيرات في غاية الأهمية لكونها تساعد الباحثين في تكوين إستنتاج فيما يخص مجاميع كبيرة ومتداخلة ومعقدة أحياناً من البيانات تتضمن عدداً كبيراً من المتغيرات مأخوذة من عدد كبير من الوحدات التجريبية. إن ضرورة وفائدة إستخدام طرق التحليل متعدد المتغيرات تزداد بشكلٍ طردي مع زيادة عدد الوحدات التجريبية أو عدد المتغيرات المأخوذة عنها للتحليل.

وغالباً ما يكون الهدف من إستخدام التحليل متعدد المتغيرات هو تلخيص الكمية الكبيرة من البيانات من خلال عدد صغير نسبياً من المعلمات Parameters. وبالتالي فإن الوظيفة الرئيسية لغالبية الأساليب متعددة المتغيرات هي التبسيط.

من جانب آخر، فإن التحليل متعدد المتغيرات غالباً ما يرتبط مع إيجاد علاقات ما بين:

- 1) متغيرات الإستجابة Response Variables
- 2) الوحدات التجريبية Experimental Units
- 3) كل من متغيرات الإستجابة والوحدات التجريبية

إن غالبية أساليب التحليل متعدد المتغيرات تميل إلى كونها ذات طبيعة إستكشافية لحالة ما بدلاً من كونها تأكيدية لتلك الحالة. وهذا يعني ميلها إلى خلق الفرضيات الإحصائية بدلاً من اختبارها.

ولتوضيح هذه النقطة، إفترض حالة كون باحثٍ ما لديه (50) خمسون متغيراً مقاسة على أكثر من (2000) ألفي وحدة تجريبية. إن الطرق الإحصائية الإعتيادية تتطلب من الباحث البدء بإدراج بعض الفرضيات أولاً ومن ثم قيامه بجمع البيانات، ومن بعد ذلك إستخدام هذه البيانات لتثبيت الميل الإحتمالي لتبني (قبول) هذه الفرضيات أو نفي ثبوت صحتها (رفضها).

والحالة البديلة التي غالباً ما تظهر هنا هي حالة إمتلاك الباحث لكمية كبيرة من البيانات ويساوره الشعور فيما إذا كانت هنالك ثمة معلومات ذات أهمية ضمن هذه البيانات. إن أساليب التحليل متعدد المتغيرات غالباً ما تكون مفيدة في عملية إستكشاف ضمن البيانات لمحاولة الوصول إلى نوع من القناعة بأن ثمة معلومات ذات قيمة ومفيدة تنطوي عليها مجموعة البيانات هذه.

## تصنيف طرق التحليل متعدد المتغيرات

أحد الفروقات الأساسية ما بين الطرق متعددة المتغيرات يكمن في تصنيفها إلى صنفين رئيسيين من حيث أساليب التحكم في التحليل وهما:

- 1) أساليب تحكّم المتغيرات Variable – directed techniques وتشمل تلك التي تتعامل بشكلٍ رئيسي مع العلاقات التي من الممكن ظهورها ضمن متغيرات الإستجابة Response Variables التي يتم قياسها. ومثال ذلك:
  - التحليلات المعتمدة على مصفوفات معامل الارتباط Correlation Matrices
  - تحليل المركبات الرئيسية Principle Components Analysis
  - التحليل العاملي Factor Analysis
  - تحليل الارتباط القويم (إرتباط المجاميع) Canonical Correlation Analysis

- 2) أساليب التحكم الشخصي Individual – directed techniques وتشمل تلك التي تتعامل بشكلٍ رئيسي مع العلاقات التي من الممكن ظهورها ضمن الوحدات التجريبية و/أو الأشخاص الخاضعين للقياس. ومثال ذلك:
  - التحليل المميز Discriminant Analysis
  - التحليل العنقودي Cluster Analysis
  - تحليل التباين المتعدد Multivariate Analysis of Variance (MANOVA)

## تكوين متغيرات جديدة

قد نجد في كثير من الأحيان أنه من المفيد تكوين متغيرات جديدة لكل وحدة تجريبية ليكون بالإمكان عمل مقارنة فيما بينها بطريقة أكثر سهولة. هذه المتغيرات الجديدة عبارة عن دوال تتضمن جميع المتغيرات الأصلية المعتمدة في التجربة. إن العديد من الطرق متعددة المتغيرات تساعد الباحث في تكوين متغيرات جديدة تتسم بمزايا مرغوبة. مثل هذه الطرق هي المركبات الرئيسية، التحليل العاملي، تحليل الارتباط القويم، التحليل المميز القويم.



## عندما نبدأ التحليل

ما أن يبدأ الباحث التفكير في إجراء تحليل لمجموعة جديدة من البيانات، فإن أسئلة عديدة حول هذه البيانات يجب أن تؤخذ بالإعتبار. ومن هذه الأسئلة المهمة:

- 1) هل هنالك أية جوانب في البيانات يمكن إعتبارها غريبة أو غير إعتيادية؟
- 2) هل يمكن الإفتراض بأن البيانات تتوزع طبيعياً Normality؟
- 3) هل هنالك أي جانب منها خارج الطبيعي Abnormality؟
- 4) هل هنالك أية قيم شاردة (شاذة) Outliers في البيانات؟ نقصد هنا بالقيمة الشاذة لوحدة التجربة حيث المتغيرات المقاسة منها تبدو وكأنها غير متناسقة مع القياسات المأخوذة لوحدات أخرى.

وفي هذا الكتاب سنتناول المواضيع التالية:

- تحليل الإنحدار المتعدد
- المركبات الرئيسية
- الإنحدار اللوجستي
- تحليل التباين المتعدد
- التحليل المميز
- التحليل العنقودي
- الارتباط القويم
- التحليل العاملي

ولضرورة فهم التعامل مع المصفوفات والعمليات المرتبطة بها، نرى من الضروري البدء في تغطية هذا الموضوع قبل البدء في المواضيع الرئيسية أعلاه.

وبما أن المتغيرات التي سوف نتعامل معها تتبع مستويات مختلفة من القياس، لذلك سيتم تناول هذا الموضوع أيضاً بعد استعراض فصل المصفوفات.

## عمليات المصفوفات Matrix Operations

يتم تعريف المصفوفة بأبعادها التي تشير إلى عدد الصفوف وعدد الأعمدة ضمن هذه المصفوفة. فالمصفوفة  $A_{m,n}$  ترمز للمصفوفة  $A$  ذات أبعاد  $m$  من الصفوف و  $n$  من الأعمدة. أي أنها تتضمن  $m \times n$  من العناصر  $a_{ij}$  وتكتب بالصيغة التالية:

$$A_{m,n} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix}$$

ويمكن تدوير هذه المصفوفة بإستبدال الصفوف مع الأعمدة وتسمى بالمصفوفة المدورة Transpose of matrix ونرمز لها بالرمز  $A'_{m,n}$  وتكون بالصيغة التالية:

$$A'_{m,n} = \begin{bmatrix} a_{11} & a_{21} & \dots & a_{m1} \\ a_{12} & a_{22} & \dots & a_{m2} \\ \dots & \dots & \dots & \dots \\ a_{1n} & a_{2n} & \dots & a_{mn} \end{bmatrix}$$

**مثال:** إذا كانت لدينا المصفوفة:

$$A_{3,2} = \begin{bmatrix} 2 & 3 \\ 0 & -1 \\ 1 & 0 \end{bmatrix} \Rightarrow A'_{2,3} = \begin{bmatrix} 2 & 0 & 1 \\ 3 & -1 & 0 \end{bmatrix}$$

### عملية الجمع والطرح في المصفوفات

لنفترض المصفوفتين  $A_{m,n}$  و  $B_{m,n}$  حيث أن:

$$A_{m,n} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix}, \quad B_{m,n} = \begin{bmatrix} b_{11} & b_{12} & \dots & b_{1n} \\ b_{21} & b_{22} & \dots & b_{2n} \\ \dots & \dots & \dots & \dots \\ b_{m1} & b_{m2} & \dots & b_{mn} \end{bmatrix}$$

فإن ناتج جمع المصفوفتين أو طرحهما يتم بجمع أو طرح ما بين كل عنصرين متقابلين ( $a_{ij}$  و  $b_{ij}$ ) من المصفوفتين لتكوين مصفوفة جديدة  $C_{m,n}$  بعناصر  $c_{ij}$  وعلى النحو التالي:

$$A_{m,n} + B_{m,n} = \begin{bmatrix} (a_{11} + b_{11}) & (a_{12} + b_{12}) & \dots & (a_{1n} + b_{1n}) \\ (a_{21} + b_{21}) & (a_{22} + b_{22}) & \dots & (a_{2n} + b_{2n}) \\ \dots & \dots & \dots & \dots \\ (a_{m1} + b_{m1}) & (a_{m2} + b_{m2}) & \dots & (a_{mn} + b_{mn}) \end{bmatrix}$$

أما عملية الطرح فنتم باستبدال إشارة الجمع بإشارة الطرح. أي أن:

$$A_{m,n} - B_{m,n} = \begin{bmatrix} (a_{11} - b_{11}) & (a_{12} - b_{12}) & \dots & (a_{1n} - b_{1n}) \\ (a_{21} - b_{21}) & (a_{22} - b_{22}) & \dots & (a_{2n} - b_{2n}) \\ \dots & \dots & \dots & \dots \\ (a_{m1} - b_{m1}) & (a_{m2} - b_{m2}) & \dots & (a_{mn} - b_{mn}) \end{bmatrix}$$

**مثال:** لنفترض المصفوفتين:

$$A_{3,3} = \begin{bmatrix} 1 & -1 & 3 \\ 0 & 4 & -2 \\ 1 & 0 & 5 \end{bmatrix}, \quad B_{3,3} = \begin{bmatrix} 0 & 1 & 2 \\ 2 & 0 & 3 \\ 1 & 1 & 4 \end{bmatrix}$$

$$A_{3,3} + B_{3,3} = \begin{bmatrix} 1 & 0 & 5 \\ 2 & 4 & 1 \\ 2 & 1 & 9 \end{bmatrix}$$

$$A_{3,3} - B_{3,3} = \begin{bmatrix} 1 & -2 & 1 \\ -2 & 4 & -5 \\ 0 & -1 & 1 \end{bmatrix}$$

**ملاحظة:**

عند عملية الجمع أو الطرح، يشترط أن تكون المصفوفتان بنفس الأبعاد من حيث عدد الصفوف وعدد الأعمدة.

## عملية الضرب في المصفوفات

هنا يشترط أن يكون عدد صفوف إحدى المصفوفتين مساوياً لعدد أعمدة الأخرى. وتتم العملية بضرب عناصر العمود (j) من المصفوفة الثانية بعناصر الصف (i) من المصفوفة الأولى وبشكل متقابل ترتيبياً وجمع حاصل الضرب ليكون العنصر (ij) من المصفوفة الناتجة وبالشكل التالي:

$$A_{m,n} \times B_{n,p} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix} \times \begin{bmatrix} b_{11} & b_{12} & \dots & b_{1p} \\ b_{21} & b_{22} & \dots & b_{2p} \\ \dots & \dots & \dots & \dots \\ b_{n1} & b_{n2} & \dots & b_{np} \end{bmatrix}$$

$$= \begin{bmatrix} \sum a_{1i}b_{i1} & \sum a_{1i}b_{i2} & \dots & \sum a_{1i}b_{ip} \\ \sum a_{2i}b_{i1} & \sum a_{2i}b_{i2} & \dots & \sum a_{2i}b_{ip} \\ \dots & \dots & \dots & \dots \\ \sum a_{mi}b_{i1} & \sum a_{mi}b_{i2} & \dots & \sum a_{mi}b_{ip} \end{bmatrix}$$

مثال: لنفترض المصفوفتين:

$$A = \begin{bmatrix} 2 & 1 & 3 \\ 0 & -1 & 2 \\ 3 & 0 & 4 \\ 5 & 1 & -1 \end{bmatrix}, \quad B = \begin{bmatrix} -2 & 1 \\ 0 & 5 \\ 3 & 2 \end{bmatrix}$$

$$A_{4,3} \times B_{3,2} = \begin{bmatrix} 2 & 1 & 3 \\ 0 & -1 & 2 \\ 3 & 0 & 4 \\ 5 & 1 & -1 \end{bmatrix} \times \begin{bmatrix} -2 & 1 \\ 0 & 5 \\ 3 & 2 \end{bmatrix}$$

$$= \begin{bmatrix} (-4+0+9) & (2+5+6) \\ (0+0+6) & (0-5+4) \\ (-6+0+12) & (3+0+8) \\ (-10+0-3) & (5+5-2) \end{bmatrix} = \begin{bmatrix} 5 & 13 \\ 6 & -1 \\ 6 & 11 \\ -13 & 8 \end{bmatrix}$$

## حالة ضرب المصفوفة بالقيمة الثابتة c

في هذه الحالة ينتج لدينا مصفوفة بعناصر جديدة عبارة عن العناصر الأصلية مضروبة بـ القيمة الثابتة C.

$$C \times A_{m,n} = C \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix} = \begin{bmatrix} ca_{11} & ca_{12} & \dots & ca_{1n} \\ ca_{21} & ca_{22} & \dots & ca_{2n} \\ \dots & \dots & \dots & \dots \\ ca_{m1} & ca_{m2} & \dots & ca_{mn} \end{bmatrix}$$

مثال: لنفترض المصفوفة:

$$A = \begin{bmatrix} 3 & 1 & 2 \\ 0 & 4 & 1 \\ -1 & 2 & -3 \end{bmatrix}$$

$$5A = \begin{bmatrix} 15 & 5 & 10 \\ 0 & 20 & 5 \\ -5 & 10 & -15 \end{bmatrix}$$

### محددة المصفوفة |A| Determinant

لنفترض المصفوفة المربعة  $A_{n,n}$  حيث  $n = 3$  على سبيل المثال ولغرض التبسيط:

$$A_{3,3} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}$$

حيث أن المحددة |A| تكون عند إستخدامنا الصف الأول كعامل ضرب:

$$|A| = a_{11}M_{11} - a_{12}M_{12} + a_{13}M_{13}$$

$$= a_{11}C_{11} + a_{12}C_{12} + a_{13}C_{13}$$

$$C_{ij} = (-1)^{i+j} M_{ij}$$

حيث

وأن:

$$M_{11} = \begin{vmatrix} a_{22} & a_{23} \\ a_{32} & a_{33} \end{vmatrix}$$

بعد حذف الصف الأول والعمود الأول، وبشكل عام، فإن:

$$M_{ij} = |A^{-}|$$

أي أنها تساوي محددة ما يتبقى من المصفوفة A بعد حذف الصف (i) والعمود (j) (أي محددة ما سميناه  $A^{-}$ ).

### ملاحظة :

في أعلاه، إستخدمنا الصف الأول كعامل ضرب وإستطاعتنا إستخدام أي صف أو عمود لنفس الغرض ونختار الأسهل في الحسابات. وبصورة عامة، يمكن كتابة قيمة محددة المصفوفة  $A_{n,n}$  على الشكل التالي:

1- بإستخدام الصف (i) كعامل ضرب تكون:

$$\det(A) = |A| = \sum_{j=1}^n a_{ij} C_{ij}$$

2- بإستخدام العمود (j) كعامل ضرب تكون:

$$\det(A) = |A| = \sum_{i=1}^n a_{ij} C_{ij}$$

مثال: إفترض المصفوفة:

$$A_{3,3} = \begin{bmatrix} 3 & 4 & -1 \\ 2 & -4 & 5 \\ 0 & 1 & -6 \end{bmatrix}$$

وإستخدام العمود الأول تكون:

$$\begin{aligned} |A| &= 3 \begin{vmatrix} -4 & 5 \\ 1 & -6 \end{vmatrix} - 2 \begin{vmatrix} 4 & -1 \\ 1 & -6 \end{vmatrix} + 0 \begin{vmatrix} 4 & -1 \\ -4 & 5 \end{vmatrix} \\ &= 3[(-4)(-6) - (5)(1)] - 2[(4)(-6) - (-1)(1)] + 0.0 \\ &= 3(19) - 2(-23) \\ &= 103 \end{aligned}$$

ولو إستخدمنا الصف الأخير، يكون لدينا:

$$\begin{aligned} |A| &= 0 \begin{vmatrix} 4 & -1 \\ -4 & 5 \end{vmatrix} - 1 \begin{vmatrix} 3 & -1 \\ 2 & 5 \end{vmatrix} - 6 \begin{vmatrix} 3 & 4 \\ 2 & -4 \end{vmatrix} \\ &= 0.0 - 1[(3)(5) - (-1)(2)] - 6[(3)(-4) - (4)(2)] \\ &= 0.0 - 17 + 120 \\ &= 103 \end{aligned}$$

### ملاحظة:

لأجل التبسيط في عملية تحديد قيمة المحددة لأي مصفوفة مربعة، من المستحسن إختيار الصف أو العمود الذي يتضمن أصفاراً أكثر. وعادة ما يطلق على هذه الطريقة بالإختيار الذكي للصف أو العمود.

**مثال:** لو كانت لدينا المصفوفة التالية:

هنا من المستحسن إختيار العمود الثاني لكثرة الأصفار فيه. أي أنه سيكون لدينا:

$$|A| = 0 + 1 \begin{vmatrix} 1 & 0 & -1 \\ 1 & -2 & 1 \\ 2 & 0 & 1 \end{vmatrix} + 0 + 0$$

وهنا يتبقى لدينا فقط الحد الثاني وسنختار العمود الثاني أيضاً لنفس السبب ليكون لدينا:

$$|A| = 0 + 1(-2) \begin{vmatrix} 1 & -1 \\ 2 & 1 \end{vmatrix} + 0 \\ = -2 [1 + 2] = -6$$

### ملاحظة:

في حالة عدم وجود أصفار كثيرة في أحد الصفوف أو الأعمدة لتسهيل العملية الحسابية في إستخراج محددة المصفوفة، لا يزال لدينا خياراً آخر لتسهيل حصول ذلك وهو تطبيق طريقة "لابلاس Laplace" وكما هو مبين في أدناه.

### طريقة لابلاس Laplace

تستخدم هذه الطريقة لإحداث أصفار في صفٍ ما أو عمودٍ ما للمصفوفة في حالة عدم وجودها وذلك قبل حساب محددة المصفوفة، وعن طريق ضرب ذلك الصف بقيمة ثابتة وجمعها مع صف آخر أو إجراء نفس الشيء مع ذلك العمود وتكون هذه العملية مجزية خاصة في حالة القيم الكبيرة نسبياً لعناصر المصفوفة. وقد تتوالى هذه العملية لأكثر من مرة واحدة وحسب الضرورة والمثال التالي يوضح ذلك.

**مثال:** لو كانت لدينا المصفوفة A

$$A = \begin{bmatrix} 4 & 3 & 0 & -1 \\ 9 & 7 & 2 & -3 \\ 4 & 0 & 2 & -1 \\ 3 & -1 & 4 & 5 \end{bmatrix}$$

$$|A| = \begin{vmatrix} 4 & 3 & 0 & -1 \\ 9 & 7 & 2 & -3 \\ 4 & 0 & 2 & -1 \\ 3 & -1 & 4 & 5 \end{vmatrix}$$

وبضرب العمود الأخير بالعدد (4) وجمعه إلى العمود الأول يكون هذا مساوياً إلى:

$$= \begin{vmatrix} 0 & 3 & 0 & -1 \\ -3 & 7 & 2 & -3 \\ 0 & 0 & 2 & -1 \\ 23 & -1 & 4 & 5 \end{vmatrix}$$

والآن نقوم بضرب العمود الأخير بالعدد (3) وجمعه إلى العمود الثاني فيكون هذا مساوياً

إلى:

$$= \begin{vmatrix} 0 & 0 & 0 & -1 \\ -3 & -2 & 2 & -3 \\ 0 & -3 & 2 & -1 \\ 23 & 14 & 4 & 5 \end{vmatrix}$$

$$= 1 \begin{vmatrix} -3 & -2 & 2 \\ 0 & -3 & 2 \\ 23 & 14 & 4 \end{vmatrix}$$

لأننا إستخدمنا الصف الأول

ويمكننا إستخراج العدد (2) مضروباً بالعمود الأخير (وهذا يصح في حالة المحددة) ليكون

لدينا:

$$= 2 \begin{vmatrix} -3 & -2 & 1 \\ 0 & -3 & 1 \\ 23 & 14 & 2 \end{vmatrix}$$

وبضرب العمود الأخير بالعدد (3) وجمعه للعمود الثاني يصبح لدينا:

$$= 2 \begin{vmatrix} -3 & 1 & 1 \\ 0 & 0 & 1 \\ 23 & 20 & 2 \end{vmatrix}$$

$$= 2 (-1) \begin{vmatrix} -3 & 1 \\ 23 & 20 \end{vmatrix} = -2 [(-3)(20) - (1)(23)] = 166$$



### ملاحظة:

ليس من الضروري الإبقاء على التعامل مع الأعمدة وإنما يمكن التغيير نحو الصفوف أيضاً بشكلٍ مماثل وفي أي خطوة.

### بعض النظريات المفيدة في موضوع محددة المصفوفة

- (1) إذا كان كل عنصر ضمن صف أو عمود مضروباً في عدد ثابت  $k$ ، فإن محددة المصفوفة تكون قيمتها مضروبة في  $k$ .
- (2) قيمة محددة المصفوفة تساوي صفراً في حالة:
  - أ- كل عناصر أي صف أو أي عمود أصفاراً.
  - ب- كان صفان أو عمودان متشابهين.
  - ت- وجود علاقة نسبية بين أي صفين أو عمودين.
- (3) في حالة تغيير موقع عمودين أو صفين مع بعضهما، فإن قيمة المحددة تتغير إشارتها فقط.
- (4) لا تتغير قيمة المحددة إذا:
  - أ- تم كتابة الأعمدة صفوفاً أو الصفوف أعمدة.
  - ب- أضفنا لكل عنصر في صفٍ ما العنصر المقابل له في صفٍ آخر مضروباً في العدد  $k$ . (نفس الشيء ينطبق بالنسبة للأعمدة).

### ملاحظة:

إن وجود محددة للمصفوفة من عدمها تعتمد على وجود معكوس للمصفوفة من عدمها وكما سنرى ذلك لاحقاً.

### حالة المصفوفة المثلثية Triangular Matrix

تعرف المصفوفة المثلثية بأنها تلك التي تكون جميع عناصرها ضمن المثلث فوق المحور Diagonal أو الذي تحته أصفاراً. أي أنها مثل:

$$A = \begin{bmatrix} a_{11} & 0 & \dots & \dots & 0 \\ a_{21} & a_{22} & 0 & \dots & 0 \\ & & & & 0 \\ a_{n1} & a_{n2} & \dots & \dots & a_{nn} \end{bmatrix}$$

وفي مثل هذه الحالة، فإن قيمة محددة المصفوفة  $A$  عبارة عن حاصل ضرب جميع العناصر المحورية ( $a_{ii}$ ) فقط. أي أن:

$$|A| = \prod_{i=1}^n a_{ii} = a_{11}a_{22}\dots a_{nn}$$

مثال: لنفترض المصفوفة التالية:

$$A = \begin{bmatrix} 3 & 0 & 0 \\ -1 & 1 & 0 \\ 2 & -2 & 4 \end{bmatrix}$$

$$|A| = (3)(1)(4) = 12$$

### معكوس المصفوفة Matrix Inverse

لنفترض أن للمصفوفة  $A$  معكوس هو  $A^{-1}$  حيث أن:

$$A^{-1} = \text{adj}(A)/|A|$$

ولذلك فإن معكوس المصفوفة يتحدد بوجود محددة لها. وحيث أن:

$$\text{adj}(A) = C'$$

علماً أن:

$$C = \begin{bmatrix} c_{11} & c_{12} & \dots & c_{1n} \\ c_{21} & c_{22} & \dots & c_{2n} \\ \dots & \dots & \dots & \dots \\ c_{n1} & c_{n2} & \dots & c_{nn} \end{bmatrix}$$

$$C_{ij} = (-1)^{i+j} M_{ij} \quad \text{وأن:}$$

وأن  $M_{ij}$  هي محددة المصفوفة  $A$  بعد حذف الصف  $(i)$  والعمود  $(j)$  منها ومثلما سبق وذكرنا ذلك في بداية الفصل.

مثال: لنفترض المصفوفة  $A$

$$A = \begin{bmatrix} 3 & 2 & -1 \\ 1 & 6 & 3 \\ 2 & -4 & 0 \end{bmatrix}$$

$$M_{11} = \begin{vmatrix} 6 & 3 \\ -4 & 0 \end{vmatrix} = 12 \quad C_{11} = (-1)^{1+1} M_{11} = 12$$

$$M_{12} = \begin{vmatrix} 1 & 3 \\ 2 & 0 \end{vmatrix} = -6 \quad C_{12} = 6$$

$$M_{13} = \begin{vmatrix} 1 & 6 \\ 2 & -4 \end{vmatrix} = -16 \quad C_{13} = -16$$

$$M_{21} = \begin{vmatrix} 2 & -1 \\ -4 & 0 \end{vmatrix} = -4 \quad C_{21} = 4$$

$$M_{22} = \begin{vmatrix} 3 & 3-1 \\ 2 & 0 \end{vmatrix} = 2 \quad C_{22} = 2$$

$$M_{23} = \begin{vmatrix} 3 & 2 \\ 2 & -4 \end{vmatrix} = -16 \quad C_{23} = 16$$

$$M_{31} = \begin{vmatrix} 2 & -1 \\ 6 & 3 \end{vmatrix} = 12 \quad C_{31} = 12$$

$$M_{32} = \begin{vmatrix} 3 & -1 \\ 1 & 3 \end{vmatrix} = 10 \quad C_{32} = -10$$

$$M_{33} = \begin{vmatrix} 3 & 2 \\ 1 & 6 \end{vmatrix} = 16 \quad C_{33} = 16$$

وبذلك فإن:

$$C = \begin{bmatrix} 12 & 6 & -16 \\ 4 & 2 & 16 \\ 12 & -10 & 16 \end{bmatrix}$$

$$\text{adj}(A) = C' = \begin{bmatrix} 12 & 4 & 12 \\ 6 & 2 & -10 \\ -16 & 16 & 16 \end{bmatrix}$$

وحيث أن:

$$\begin{aligned} |A| &= 3 \begin{vmatrix} 6 & 3 \\ -4 & 0 \end{vmatrix} - 2 \begin{vmatrix} 1 & 3 \\ 2 & 0 \end{vmatrix} + (-1) \begin{vmatrix} 1 & 6 \\ 2 & -4 \end{vmatrix} \\ &= 3(12) - 2(-6) - 1(-16) \\ &= 36 + 12 + 16 = 64 \end{aligned}$$

وبذلك تكون معكوسة المصفوفة A كما يلي:

$$A^{-1} = \text{adj}(A) / |A| = (1/64) \begin{bmatrix} 12 & 4 & 12 \\ 6 & 2 & -10 \\ -16 & 16 & 16 \end{bmatrix}$$

$$= \begin{bmatrix} \frac{12}{64} & \frac{4}{64} & \frac{12}{64} \\ \frac{6}{64} & \frac{2}{64} & \frac{-10}{64} \\ \frac{-16}{64} & \frac{16}{64} & \frac{16}{64} \end{bmatrix}$$

**ملاحظة:**

إن جدوى إستخراج معكوس المصفوفة تظهر في غالبية عمليات التحليل متعدد المتغيرات والتي يمكن تمثيل كل أو جزء من بياناتها بمصفوفة. ونحن نعلم بأن جميع طرق التحليل متعدد المتغيرات تبرز فيها مثل هذه الحالة.

كما أن هنالك طرقاً مثل تحليل الإنحدار المتعدد وتحليل الإنحدار اللوجستي المتعدد تبرز فيهما حالة نظام المعادلات الخطية والتي يمكن إستخدام  $A^{-1}$  فيها إضافة إلى إمكانية اتباع طرق أخرى لحل مثل هذه المعادلات وإستخراج قيم المجاهيل فيها. وسوف نتناول ذلك من خلال ما يلي:

### حل نظام معادلات خطية Solving system of linear equations

لنفترض أنه لدينا المعادلتين الخطيتين:

$$2X - Y = 5$$

$$X + 2Y = -5$$

وهذه يمكن تمثيلها بالمصفوفات وكما يلي:

$$A \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 2 & -1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 5 \\ -5 \end{bmatrix}$$

وبالتالي فإننا نريد معرفة قيم X و Y اللتان تحققان صحة المعادلات هذه. وهذا يعني أن:

$$\begin{bmatrix} x \\ y \end{bmatrix} = A^{-1} \begin{bmatrix} 5 \\ -5 \end{bmatrix} = \begin{bmatrix} 2 & -1 \\ 1 & 2 \end{bmatrix}^{-1} \begin{bmatrix} 5 \\ -5 \end{bmatrix}$$

$$= \frac{1}{5} \begin{bmatrix} 2 & 1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} 5 \\ -5 \end{bmatrix} = \frac{1}{5} \begin{bmatrix} 5 \\ -15 \end{bmatrix} = \begin{bmatrix} 1 \\ -3 \end{bmatrix}$$

## طريقة غاوس – جوردن Gouss - Jordan للمعادلات الخطية

ويمكن تطبيق هذه الطريقة على المثال أعلاه بالشكل التالي:

1- نضع الصيغة التالية:

$$\left[ \begin{array}{cc|c} 2 & -1 & 5 \\ 1 & 2 & -5 \end{array} \right]$$

2- نبدأ بإجراء التحويلات على الجزء الأيسر من الصيغة هذه بحيث نصل بها للشكل التالي:

$$\left[ \begin{array}{cc|c} 1 & 0 & c_1 \\ 0 & 1 & c_2 \end{array} \right] = \left[ \begin{array}{cc|c} 1 & 0 & 1 \\ 0 & 1 & -3 \end{array} \right]$$

أي أننا نجري التحويلات المناسبة على المصفوفة A والتي تمثل الطرف الأيسر حتى نحولها إلى مصفوفة الوحدة (مصفوفة بمحور جميع عناصره "1" مع بقية العناصر خارجه أصفاراً). هذه التحويلات بالطبع تأخذ مداها على الجانب الأيمن تلقائياً. والنتيجة هي أن:

$$\begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} = \begin{bmatrix} 1 \\ -3 \end{bmatrix}$$

وفيما يلي الأسلوب المتبع لإجراء هذه العملية:

نبدأ بنفس الصيغة:

$$\left[ \begin{array}{cc|c} 2 & -1 & 5 \\ 1 & 2 & -5 \end{array} \right]$$

ولأننا بحاجة إلى الرقم "1" للعنصر  $a_{11}$ ، فإنه من المناسب إستبدال الصفين أحدهما مكان الآخر. بهذه العملية لم يتم تغيير أي شئ سوى أننا إعتبرنا المعادلة الثانية هي الأولى. أي أنه لدينا:

$$\left[ \begin{array}{cc|c} 1 & 2 & -5 \\ 2 & -1 & 5 \end{array} \right]$$

وبضرب الصف الأول بالمقدار (-2) وجمعه للصف الثاني نحصل على:

$$\left[ \begin{array}{cc|c} 1 & 2 & -5 \\ 0 & -5 & 15 \end{array} \right]$$

وبضرب الصف الثاني بالمقدار (-1/5) يكون لدينا:

$$\left[ \begin{array}{cc|c} 1 & 2 & -5 \\ 0 & 1 & -3 \end{array} \right]$$

وبضرب الصف الثاني بالمقدار (-2) وجمعه للصف الأول يكون لدينا:

$$\left[ \begin{array}{cc|c} 1 & 0 & 1 \\ 0 & 1 & -3 \end{array} \right]$$

وبما أننا توصلنا لما نريده فإن الحل هو:

$$\begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 1 \\ -3 \end{bmatrix}$$

**مثال آخر:** لنفترض أن مجموعة المعادلات الخطية التالية:

$$X + Y = 5$$

$$-2X - Y + 2Z = -10$$

$$3X + 6Y + 7Z = 14$$

هذه المعادلات تعطينا الصيغة التالية:

$$\left[ \begin{array}{cccc|c} 1 & 1 & 0 & 5 \\ -2 & -1 & 2 & -10 \\ 3 & 6 & 7 & 14 \end{array} \right]$$

وبضرب الصف الأول بالمقدار (2) وجمعه للصف الثاني، وكذلك بضرب الصف الأول

بالمقدار (-3) وجمعه للصف الثالث نحصل على:

$$\left[ \begin{array}{cccc|c} 1 & 1 & 0 & 5 \\ 0 & 1 & 2 & 0 \\ 0 & 3 & 7 & -1 \end{array} \right]$$

وبضرب الصف الثاني بالمقدار (-1) وجمعه للصف الأول، وكذلك بضرب الصف الثاني

بالمقدار (-2) وجمعه للصف الثالث نحصل على:

$$\left[ \begin{array}{cccc|c} 1 & 0 & -2 & 5 \\ 0 & 1 & 2 & 0 \\ 0 & 0 & 1 & -1 \end{array} \right]$$

وبضرب الصف الثالث بالمقدار (2) وجمعه للصف الأول، وكذلك بضرب الصف الثالث بالمقدار (-2) وجمعه للصف الثاني نحصل على:

$$\left[ \begin{array}{ccc|c} 1 & 0 & 0 & 3 \\ 0 & 1 & 0 & 2 \\ 0 & 0 & 1 & -1 \end{array} \right]$$

وبما أننا توصلنا لما نريده فإن الحل هو:

$$\begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} 3 \\ 2 \\ -1 \end{bmatrix}$$

### ملاحظة:

عندما يكون لدينا حلاً لنظام المعادلات، فإن ذلك يعني وجود معكوسة للمصفوفة A وهذا بدوره يعني حتمية وجود محددة لهذه المصفوفة. وفي حالات أخرى غير ذلك، فقد لا نجد حلاً لنظام المعادلات هذا. ومثال ذلك في الحالتين التاليتين:

(1) عند وجود حالة الترابط الخطي ما بين أية معادلتين وهذه الحالة تدعى dependent system ومثال ذلك لو كانت لدينا الصيغة:

$$\left[ \begin{array}{ccc|c} 2 & -4 & 4 & \\ 1 & -2 & 2 & 2 \end{array} \right]$$

فإنه بضرب الصف الأول بالمقدار (-1/2) وجمعه للصف الثاني، سيصبح لدينا:

$$\left[ \begin{array}{ccc|c} 2 & -4 & 4 & \\ 0 & 0 & 0 & 0 \end{array} \right]$$

وهذه عبارة عن معادلة واحدة بمتغيرين مجهولين وهنا يصعب إيجاد حلاً لها بالنسبة لقيم المتغيرين.

(2) عند غياب حالة التناسق ما بين المعادلات مما يعني عدم صحة معادلة أو أكثر عند نفس القيم للمجهول. وهذه الحالة تدعى inconsistent system. ومثال ذلك لو كانت لدينا الصيغة التالية:

$$\left[ \begin{array}{ccc|c} 2 & -5 & 10 & \\ 2/5 & -1 & 7 & \end{array} \right]$$

فإنه بضرب الصف الأول بالمقدار  $(-1/5)$  وجمعه للصف الثاني، سيصبح لدينا:

$$\left[ \begin{array}{ccc|c} 2 & -5 & 10 & \\ 0 & 0 & 5 & \end{array} \right]$$

وهذا يعني أن  $0.0 = 5$  وهو غير منطقي وبالتالي لا يوجد حلاً لهذا النظام.

### طريقة الحذف Elimination Method لحل نظام المعادلات الخطية

وهذه الطريقة تعتمد أسلوب حل كل معادلتين أنيتين بطريقة الحذف التوافقي للمتغيرات والإنتهاء بمتغير واحد نجد قيمته ومن ثم الرجوع عكسياً لتعويض هذه القيمة في معادلة ذات متغيرين تتضمن متغيراً واحداً إلى جانب هذا المتغير ليتم تحديد قيمة المتغير الثاني. وتتكرر هذه العملية على معادلة أخرى تتضمن متغيراً ثالثاً إلى جانب هذين المتغيرين لغرض إيجاد قيمته. وتستمر هذه العملية تباعاً حتى الإنتهاء من تحديد قيم جميع المتغيرات في مجموعة المعادلات هذه.

وفيما يلي مثالاً لتوضيح مجريات هذه الطريقة:

مثال: لنفترض نظام المعادلات الخطية التالية:

A  $2X + y - 3Z = -7$

B  $3X - 2Y + Z = 11$

C  $-2X - 3Y - 2Z = 3$

ملاحظة:

لقد تم تسمية المعادلات هذه بالحروف A و B و C لغرض تسهيل الرجوع إليها وسنستمر بتسمية المعادلات الجديدة بنفس الأسلوب.

ومن أجل إيجاد الحل، سنتبع الخطوات التالية بالنسبة لهذه المجموعة:

1) حذف المتغير Y من المعادلتين A و B بعد ضرب المعادلة A بالمقدار (2) والذي سيتم التنويه بذلك إزاء المعادلة نفسها. أي أننا سنحذف المتغير Y من المعادلتين

2A و B بطريقة الجمع أو الطرح لتصبح لدينا المعادلة الجديدة D:

(2A)  $4X + 2y - 6Z = -14$

B  $3X - 2Y + Z = 11$



بالجمع

$$D \quad 7X \quad - 5Z = -3$$

(2) حذف المتغير Y من المعادلتين A و C بنفس الأسلوب:

$$(3A) \quad 6X + 3y - 9Z = -21$$

$$C \quad -2X - 3Y - 2Z = 3$$

بالجمع

$$E \quad 4A \quad - 11Z = -18$$

(3) حذف المتغير X من المعادلتين D و E بنفس الأسلوب:

$$(-4D) \quad -28X + 20Z = 12$$

$$(7E) \quad 28X - 77Z = -126$$

بالجمع

$$-57Z = -114$$

ومنها نجد أن:  $Z = 2$

(4) وبالتعويض عن هذه القيمة في المعادلة D نحصل على  $X = 1$ .

(5) وبالتعويض هاتين القيمتين بالمعادلة A نحصل على  $Y = -3$ .

(6) وبذلك تكون نتيجة الحل النهائي لهذا النظام هي:

$$\begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = \begin{bmatrix} 1 \\ -3 \\ 2 \end{bmatrix}$$

### قاعدة كرامير Cramer's rule لحل نظام المعادلات الخطية

لنفترض المعادلات الخطية التالية بصيغة المصفوفات:

$$\begin{bmatrix} 1 & 0 & 2 \\ -3 & 4 & 6 \\ -1 & -2 & 3 \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix} = \begin{bmatrix} 6 \\ 30 \\ 8 \end{bmatrix}$$

وهذا يعني أن المصفوفة A هي:

$$A = \begin{bmatrix} 1 & 0 & 2 \\ -3 & 4 & 6 \\ -1 & -2 & 3 \end{bmatrix}$$

وبعد ذلك يتم إحلال المتجه مكان كل عمود من أعمدة المصفوفة A ليكون لدينا المصفوفات الجديدة التالية:

$$C = \begin{bmatrix} 6 \\ 30 \\ 8 \end{bmatrix}$$

فعند إحلال C مكان العمود الأول أو الثاني أو الثالث، تكون لدينا المصفوفات A<sub>1</sub> و A<sub>2</sub> و A<sub>3</sub> على التوالي:

$$A_1 = \begin{bmatrix} 6 & 0 & 2 \\ 30 & 4 & 6 \\ 8 & -2 & 3 \end{bmatrix}$$

$$A_2 = \begin{bmatrix} 1 & 6 & 2 \\ -3 & 30 & 6 \\ -1 & 8 & 3 \end{bmatrix}$$

$$A_3 = \begin{bmatrix} 1 & 0 & 6 \\ -3 & 4 & 30 \\ -1 & -2 & 8 \end{bmatrix}$$

وبالتالي فإن إستخراج قيم المتغيرات X<sub>1</sub> , X<sub>2</sub> , X<sub>3</sub> ستكون بالشكل التالي:

$$X_1 = |A_1|/|A| = (-40)/(44) = -10/11$$

$$X_2 = |A_2|/|A| = (72)/(44) = 18/11$$

$$X_3 = |A_3|/|A| = (152)/(44) = 38/11$$

### القيم المميزة والمتجهات المميزة Eigen values & Eigen vectors

في تحليل متعدد المتغيرات، غالباً ما تكون لدينا مصفوفة مربعة A<sub>n,n</sub> ونحتاج إلى تحديد القيم والمتجهات المميزة لهذه المصفوفة.

وهذا واضح في طرق المكونات الرئيسية Principle components وإستخداماتها في طرق تحليل أخرى وفي مقدمتها التحليل العاملي Factor analysis. ولأننا نتعامل في هذه الطرق مع مصفوفة التباين - التباين المشترك أو مصفوفة معاملات الارتباط، فإن المصفوفة تكون في العادة مربعة ومتماثلة بنفس الوقت.

### ملاحظة:

لاحظ أن التماثل لا يتحقق إلا مع المصفوفة المربعة. ولذلك يمكن الإكتفاء بالقول مصفوفة متماثلة.

ولتوضيح كيفية تحديد القيم والمتجهات المميزة لمصفوفة ما، سنعرض ذلك من خلال المثال التالي:

**مثال :** لنفترض أنه لدينا مصفوفة التباين - التباين المشترك A التالية:

$$A = \begin{bmatrix} 5 & -2 \\ -2 & 2 \end{bmatrix}$$

وعلينا حل المعادلة:

$$\text{Det}(A - \lambda I) = 0.0$$

بالنسبة إلى  $\lambda$  والتي تمثل لنا متجه القيم المميزة للمصفوفة A.

وهذا يعني لنا العمل على:

$$\begin{vmatrix} 5-\lambda & -2 \\ -2 & 2-\lambda \end{vmatrix} = 0.0$$

وبالنتيجة يكون لدينا:

$$(5 - \lambda)(2 - \lambda) - (-2)(-2) = 0.0$$

$$\lambda^2 - 7\lambda + 6 = 0.0$$

أو:

وبحل المعادلة هذه نحصل على:

$$(\lambda - 1)(\lambda - 6) = 0.0$$

$$\lambda = 6 \text{ أو } \lambda = 1$$

أي أنه لدينا:

وهاتان القيمتان هما التعبير عن القيمة المميزة الأولى (الأكبر)  $\lambda_1 = 6$  والقيمة المميزة

الثانية  $\lambda_2 = 1$ .

### ملاحظة:

يجب أن يبقى في بالنا أن:

$$\lambda_1 \geq \lambda_2$$

$$\lambda_1 + \lambda_2 = \text{trace}(A) = 7 \quad \text{و}$$

ولغرض تحديد المتجه المميز مقابل كل قيمة مميزة، فإننا نستخدم الآتي:

$$\begin{pmatrix} 5-\lambda_1 & -2 \\ -2 & 2-\lambda_1 \end{pmatrix} \begin{pmatrix} a_{11} \\ a_{12} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$\begin{pmatrix} 5-6 & -2 \\ -2 & 2-6 \end{pmatrix} \begin{pmatrix} a_{11} \\ a_{12} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$\begin{pmatrix} -1 & -2 \\ -2 & -4 \end{pmatrix} \begin{pmatrix} a_{11} \\ a_{12} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

حيث أن  $\begin{pmatrix} a_{11} \\ a_{12} \end{pmatrix}$  هو المتجه المميز المقابل للقيمة المميزة  $\lambda_1$

وبإستخدام أيّ من المعادلتين (المتجانستين):

$$-a_{11} - 2a_{12} = 0.0$$

$$-2a_{11} - 4a_{12} = 0.0$$

$$a_{11} = -2a_{12}$$

سنحصل على:

$$(a_{11})^2 + (a_{12})^2 = 1$$

وبالتنسيق مع شرط كون

سنحصل على:

$$\begin{pmatrix} a_{11} \\ a_{12} \end{pmatrix} = \begin{pmatrix} -2/\sqrt{5} \\ 1/\sqrt{5} \end{pmatrix}$$

وبتطبيق نفس الإسلوب مع  $\lambda_2 = 1$  نحصل على المتجه المميز الثاني المقابل لها وهو:

$$\begin{pmatrix} a_{21} \\ a_{22} \end{pmatrix} = \begin{pmatrix} 1/\sqrt{5} \\ 2/\sqrt{5} \end{pmatrix}$$

**ملاحظة 1:**

من المهم هنا ولغرض تطبيق طرق تحليل متعدد المتغيرات والتي تعتمد على هذه القيم، أن نتحقق لدينا حالة التعامدية Orthogonality ما بين هذه الإتجاهات (المتجهات المميزة) وهذا يعني:

$$\sum_j a_{ij}^2 = 1 \quad , \quad i=1, 2 \quad -1$$

$$\sum_j a_{1j}a_{2j} = 0.0 \quad -2$$

وهذان الشرطان متحققان في مثالنا أعلاه.

## ملاحظة 2:

في حالة إستخدامنا مصفوفة التباين-التباين المشترك، فإن:

$$\sum_i \lambda_i = \sum_i \text{var}(X_i)$$

أما في حالة إستخدامنا مصفوفة الارتباط، فإن:

$$\sum_i \lambda_i = p$$

وإن  $p =$  عدد المتغيرات ضمن المصفوفة

والمثال أعلاه يمكن أن يتطابق مع الحالة الأولى.

والمثال التالي ينطبق على إستخدام الحالة الثانية:

**مثال:** لنفترض المصفوفة:

$$A = R = \begin{bmatrix} 1 & 0.7 \\ 0.7 & 1 \end{bmatrix}$$

أي أنها تمثل مصفوفة الارتباط بين متغيرين  $X_1$  و  $X_2$  وأن  $r_{12} = 0.7$

وهذا يعني لنا العمل على:

$$\begin{vmatrix} 1-\lambda & 0.7 \\ 0.7 & 1-\lambda \end{vmatrix} = 0.0$$

وبالنتيجة يكون لدينا:

$$(1 - \lambda)(1 - \lambda) - 0.49 = 0.0$$

$$\lambda^2 - 2\lambda + 1 = 0.49 \quad \text{أو:}$$

وبحل المعادلة هذه نحصل على القيمتان المميزتان:

$$\lambda_2 = 0.3 \quad \text{و} \quad \lambda_1 = 1.7$$

ولغرض تحديد المتجه المميز مقابل كل قيمة مميزة، فإننا نستخدم الآتي:

$$\begin{pmatrix} 1-\lambda_1 & 0.7 \\ 0.7 & 1-\lambda_1 \end{pmatrix} \begin{pmatrix} a_{11} \\ a_{12} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$\begin{pmatrix} 1-1.7 & 0.7 \\ 0.7 & 1-0.7 \end{pmatrix} \begin{pmatrix} a_{11} \\ a_{12} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$\begin{pmatrix} -0.7 & 0.7 \\ 0.7 & -0.7 \end{pmatrix} \begin{pmatrix} a_{11} \\ a_{12} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

حيث أن  $\begin{pmatrix} a_{11} \\ a_{12} \end{pmatrix}$  هو المتجه المميز المقابل للقيمة المميزة  $\lambda_1$

$$\begin{pmatrix} a_{11} \\ a_{12} \end{pmatrix} = \begin{pmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \end{pmatrix} \quad \text{وهنا نحصل على المتجه المميز الأول :}$$

وفي حالة استخدامنا القيمة المميزة  $\lambda_2 = 0.3$  نحصل على:

$$\begin{pmatrix} 0.7 & 0.7 \\ 0.7 & 0.7 \end{pmatrix} \begin{pmatrix} a_{21} \\ a_{22} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

ومنها سنحصل على المتجه المميز الثاني:

$$\begin{pmatrix} a_{21} \\ a_{22} \end{pmatrix} = \begin{pmatrix} 1/\sqrt{2} \\ -1/\sqrt{2} \end{pmatrix}$$

وهما متعامدان.

## مستويات قياس المتغيرات Variables Levels of Measurement

بطبيعة الحال، عندما نتعامل مع أي نوع من الطرق الإحصائية، فإننا نتعامل مع متغيرات تختلف في تصنيفها من حيث طبيعة القيم التي تشير إليها. ونقصد بذلك، طبيعة مستوى القياس الذي تم بموجبه تحديد قيمة معينة لأي وحدة إحصائية ضمن المتغير الواحد. أي أن كل متغير يخضع لمستوى محدد من القياس تتعامل معه جميع وحداته عند تحديد قيمها. وعندما نستعرض العديد من الطرق الإحصائية سنجد طريقة ما تتطلب متغيرات بمستوى قياس محدد واحد لجميع المتغيرات وطرق أخرى تستوعب متغيرات بأكثر من مستوى قياس واحد.

إن ضرورة توضيح هذا الموضوع هنا هو كون طرق التحليل الإحصائي متعدد المتغيرات التي سنتناولها ضمن هذا الكتاب تتباين أحياناً حسب طبيعة مستويات قياس المتغيرات التي تتعامل معها. وبالتالي فإننا سنتناول بشئ من التفصيل طبيعة مستويات القياس المعروفة هذه مع أمثلة مناسبة لغرض التوضيح. وبذلك سوف نشير إلى طبيعة مستوى القياس لكل متغير يتم استخدامه ضمن أي طريقة تحليل في هذا الكتاب.

وبشكل عام، فإنه من المعروف بأن هنالك أربعة أنواع من مستويات القياس والتي سنذكرها تباعاً حسب القوة (درجة التطور) التي تتميز بها وهي:

### المقياس الإسمي (Nominal Scale)

وهو أدنى مستويات القياس ويناسب المتغيرات الكيفية أو النوعية التي تتطلب تصنيف الأفراد إلى مجموعات منفصلة للتمييز بينهم في سمة معينة، ويكون الهدف من عملية القياس في هذه الحالة هو التصنيف الذي يراعي الفروق النوعية بين الأفراد. والأعداد المستخدمة في هذا المستوى من القياس تعد بمثابة رموز بسيطة تستخدم كأسماء لفئات أو مجموعات منفصلة وتمييزة. التصنيفات في هذه الحالة مختلفة وغير متكررة والأرقام (الأعداد) لم توضع إلا لسهولة التعامل مع المجموعات وليس لها أي دلالة رقمية. فالبيانات (الأعداد) في هذه الحالة فقط تصنف البيانات ولا تعطي لها أي ترتيب.

ولا نستطيع أن نجري أي عمليات حسابية على هذه الأعداد ولكن بالإمكان عد الأفراد في كل فئة. ومن أمثلة متغيرات هذا المستوى: النوع (ذكور أو إناث) وقد نرسم العدد (1) للذكور والعدد (2) للإناث فنعرف مثلاً أنه لدينا 14 ذكور، 16 إناث بالنسبة لمتغير النوع وهكذا للأمثلة أخرى مثل الجنسية، والديانة، والحالة الاجتماعية، أو حسب مناطق السكن (جنوب، شرق، شمال، غرب)، أو ألوان السيارات (أخضر، أسود، أبيض) فيتم تصنيف السيارات بحسب ألوانها ولا يمكن ترتيبها.

## 1) المقياس الرتبوي (الترتبي) (Ordinal Scale)

هذا المقياس يصنّف البيانات كما هو حال المقياس السابق لكن يضيف إليها خاصية الترتيب، بحيث أنه يمكن وضع التصنيفات في ترتيب واضح متسلسل. وبذلك يعتبر أكثر تطوراً من المقياس الإسمي. ومن الأمثلة الواضحة على هذا النوع من المقاييس هي المقاييس الخاصة بالتقييم (Rating Scale) أو مقاييس لكرت (Likert Scale) المستخدمة في تصنيف أجوبة أسئلة الإستبانة. فهو يفيد الترتيب بين الأفراد أو وحدات المتغير ولكن ليس من الضروري أن تكون الفروق في مقدار أو درجة الخاصية بين كل رتبتين متجاورتين منتظمة. عندما تكون المشاهدات مختلفة فقط من فئة إلى أخرى، بل فيما يمكن ترتيبها بالنسبة إلى معيار معين فقد يقال عنها بأنها مقياس رتبوي. المرضى في دور النقاهة قد يتم وصفهم ( غير متحسنين/ متحسنين/ متحسنين جداً). الأشخاص قد يتم تصنيفهم طبقاً لحالتهم الإقتصادية والإجتماعية مثل (واطيئ/متوسط/مرتفع)، ودرجة الذكاء عند الأطفال قد تكون (فوق المعدل/ في المعدل/ تحت المعدل). في كل واحد من هذه الأمثلة نرى أن أعضاء الفئة الواحدة يعتبرون متساويين لكن أعضاء إحدى الفئات يعتبرون أوطأ، أسوأ، أو أصغر من أعضاء الفئة الأخرى وهي التي تحمل تبعاً نفس العلاقة مع فئة ثالثة. فمثلاً المريض المتحسن جداً يكون أكثر صحة من مريض آخر تم تصنيفه كمتحسن بينما المريض المتحسن يكون أحسن حالاً من الآخر غير المتحسن. وبالطبع من المستحيل الإستنتاج أن الفرق بين أعضاء إحدى الفئات وأخرى مجاورة في الأسفل يساوي الفرق بين أعضاء تلك الفئة وأعضاء الفئة المجاورة لها في الأعلى. أي أن درجة التحسن بين غير المتحسن والمتحسن قد لا تكون نفسها بين المتحسن والمتحسن جداً. إن عمل الأرقام المعطاة لبيانات رتبوية هو لترتيب (أو تدرج to rank) المشاهدات من الأوطأ إلى الأعلى وهذا بدوره رتبوي. لهذا لا معنى للعمليات الحسابية في هذا المستوى من القياس، على الرغم من القدرة على إجرائها، وذلك لأن نتائج العمليات الحسابية لا تعكس مقدار الكم للصفة المراد قياسها.

## 2) المقياس الفئوي (الفترّي) (Interval Scale)

المقياس الفئوي (الفترّي) يعتبر أكثر تطوراً من المقياس الإسمي أو الرتبوي لأنه مع هذا المقياس ليس بالإمكان ترتيب القياسات فقط، لكن المسافة بين أي قياسين تكون معروفة. فنحن نعرف مثلاً، أن المسافة بين القياس 20 والقياس 30 يكون مساوياً للمسافة بين القياس 30 والقياس 40. والمقدرة في عمل هذا، يقود إلى إستخدام وحدة المسافة ونقطة الصفر وكلاهما عفوي. فنقطة الصفر المختارة ليست صفراً حقيقياً لكونها لا تشير إلى الغياب الكلي للمقدار الذي يكون مقاساً. فالفروق أو المسافات المتساوية على هذا المقياس متساوية تدل على مقادير متساوية من الخاصية التي نقيسها، ولذا يمكن جمع هذه المسافات أو طرحها أو ضربها مع مراعاة أنه لا يوجد لمقياس المسافة صفر حقيقي أو مطلق (يدل على إنعدام الشيء أو عدم



وجود الخاصية). وربما أحسن مثال لمقياس الفترة هو الطريقة التي تقاس بها درجة الحرارة. فوحدة القياس هي الدرجة ونقطة المقارنة هي " درجة الصفر " المختارة عفويًا. فدرجة الصفر مئوي لا تعني إنعدام الحرارة من الوجود، كما أن الفرق ما بين درجتَي الحرارة 25 و 28 هو نفسه الفرق ما بين درجتَي الحرارة 81 و 84. ومن جانب آخر نجد أن درجة الحرارة الصفر بمقياس الفهرنهايت هي 32 ، كما أن الصفر الجامعي في المساق هو 35. إن مقياس الفترة يختلف عن المقياس الإسمي والمقياس الرتبوي بكونه مقياس مقادير حقيقية ويستخدم للبيانات الكمية.

### 3) المقياس النسبي (Ratio Scale)

يحتفظ هذا النوع من المقاييس بمزايا الثلاثة أنواع السابقة، فهو يصنّف، يرتّب ويوضح المسافات بشكل متساوي وموزون. وبالإضافة لذلك، فإن الخاصية التي ينفرد بها مقياس النسب هي نقطة الصفر الحقيقية أي صفر مطلق يناظر بالفعل إنعدام الخاصية والسمة المقاسة. ويمكن إجراء جميع العمليات الحسابية الأساسية على هذه المقاييس ومن أمثلتها مقياس الوزن، والحجم، والطول والمسافة وغيرها. وسمي بمستوى القياس النسبي، لأن نسبة الأرقام إلى بعضها ذات معنى ودلالة. العمليات الرياضية (الحسابية) والمقارنات عند هذا المقياس لها معنى ويمكن للباحثين إجراء عمليات القسمة والضرب دون تغيير في الخصائص. وتمثل المسطرة مثالاً بسيطاً للمقياس النسبي، فالفرق فيها بين نقاط القياس متساو في العرض، وهناك نقطة صفر حقيقية على المسطرة التي تجعل أي قياس تحت الصفر ليس بذي معنى. ولذلك يمكن تصنيفه بأنه أعلى مستوى للقياس. وهذا المقياس ينطبق دائماً على قياس المتغير الكمي والذي سيتم تناوله تالياً. وبذلك قد نطلق عليه "المقياس الكمي".

وإنطلاقاً من هذا الإستعراض لمستويات القياس للمتغيرات، نلاحظ أن المتغيرات يمكن أن تكون، بحسب مستوى قياسها، مترية (كمية)، مثل (درجات الامتحان، السرعة، الفئات العمرية، مستويات الدخل،...) وهذه تشمل المقياسين (النسبي Ratio) و (الفئوي Interval)، أو لامتريّة (نوعية) مثل (اللون، الجنس، مستوى التعليم، الرتبة العسكرية....) وهذه تشمل المقياسين (الإسمي Nominal) و (الرتبوي Ordinal).

كما ينقسم المتغير الكمي إلى متغير متصل (مستمر Continuous) والذي يمكن أن تأخذ مقاديره أي قيمة، ومتغير منفصل (متقطع Discrete) تنحصر مقاديره في قيم العدد (الأعداد الصحيحة). وتصنيف المتغير الكمي إلى متصل أو منفصل يعتمد على أداة القياس المستخدمة في قياسه وليس على القيم التي يظهر بها ونتعامل بموجبها. ومثالاً على ذلك قيم درجات الحرارة حيث أننا نتعامل في العادة مع أعداد صحيحة منفصلة (1 ، 2 ، 3 ،.....) ويظهر ذلك بوضوح عند مراقبة نشرات الأنواء الجوية. ولكن مقياس درجة الحرارة لا ينتقل

من الدرجة (1) إلى الدرجة (2) مباشرةً وإنما يقرأ أي جزء ما بين الدرجتين، وبالتالي فإن درجة الحرارة هي متغير متصل ولو تعاملنا مع قيمه المنفصلة.

#### **ملاحظة:**

من خلال إستعراضنا لمستويات القياس هذه، والعلم بأن غالبية الطرق الإحصائية في التحليل المتعدد التي سيتم تناولها في هذا الكتاب تعتمد مصفوفة التباين أو مصفوفة معامل الارتباط في العمليات الحسابية الخاصة بها والتي تتعامل حصراً مع المتغيرات الكمية، يصبح من الواضح جداً بأن المقياس النسبي (Ratio scale) هو المقياس المناسب لجميع المتغيرات التي تدخل في أي مصفوفة وقد يكون من الممكن التعامل مع المقياس الفئوي (Interval scale) وسوف نراعي ذلك عند تناول هذه الطرق.

## تحليل الإنحدار والإرتباط المتعدد Multiple Regression & Correlation Analysis

### أولاً: تحليل الإنحدار الخطي البسيط

عند دراسة العلاقة بين المتغيرات فإن الخطوة الأولى تبدأ في تحديد المتغيرات التي تدخل في هذه العلاقة. فإذا كانت هذه العلاقة بسيطة أي بين إثنين من المتغيرات، ففي هذه الحالة عادة ما نفكر بأحد المتغيرين بأنه المتغير السببي ويوصف بأنه المتغير المستقل ( $x$ ) والمتغير الآخر بأنه المتغير التابع أو متغير الاستجابة ( $y$ ) أي أن  $y$  هو دالة  $x$  لكن القيمة المشاهدة إلى  $y$  لا يمكن أن ترتبط بعلاقة خطية مضبوطة مع القيمة المشاهدة إلى  $x$  في كل محاولة من المحاولات. لذلك نلجأ إلى إضافة حد جديد يسمى حد الخطأ أو الخطأ العشوائي الذي يمثل الفشل بالنسبة للقيمة المشاهدة إلى  $y$  في أن تكون مساوية للعلاقة الخطية ( $\beta_0 + \beta_1 x$ ) فتصبح الصيغة الملائمة بالشكل التالي:

$$y_i = \beta_0 + \beta_1 x_i + e_i \quad i = 1, 2, \dots, n$$

حيث أن:

$y_i$  : تمثل متغير الاستجابة (المعتمد) في المشاهدة ( $i$ ) وهو متغير كمي (نسبي أو فئوي).

$x_i$  : تمثل المتغير التوضيحي (المستقل) في المشاهدة ( $i$ ) وهو متغير كمي (نسبي أو فئوي).

$e_i$  : يمثل حد الخطأ في المشاهدة ( $i$ ).

### الفرضيات الخاصة بنموذج الإنحدار البسيط

وأهمها الفرضيات الخاصة بحد الخطأ  $e_i$  والتي تشمل الآتي:

- $e_i$  متغير عشوائي
- يتوزع توزيع طبيعي
- وسطه الصفر  $E(e_i)=0$
- تباينه ثابت  $\sigma^2$  لكل قيم  $x_i$
- أي أن  $e_i \sim N(0, \sigma^2)$
- المتغير العشوائي  $e_i$  مستقل عن قيم  $x_i$  أي أن  $cov(x_i, e_i) = 0.0$
- الخطأ العشوائي في أي محاولة يكون مستقلاً عن الأخطاء العشوائية في محاولة أخرى أي أن:

$$cov(e_i, e_j) = 0.0 \quad i \neq j$$

أما الفرضيات الخاصة بالمتغير المستقل  $x_i$  فتنحصر عموماً بكون قيمه ثابتة.

إستنتاجات خاصة بالمتغير المعتمد (  $y$  )

- بما أن  $E(e_i) = 0$  فإن  $E(y_i) = \beta_0 + \beta_1 x_i$
- بما أن  $v(e_i) = \sigma^2$  لكل قيم  $i$  فإن  $v(y) = \sigma^2$  لكل قيم  $i$
- $y_i$  متغير عشوائي يتبع التوزيع الطبيعي بمتوسط مقداره  $\beta_0 + \beta_1 x_i$  وتباين ثابت مقداره  $\sigma^2$ . أي أن:  
 $y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$
- بما أن قيم الاخطاء العشوائية غير مرتبطة ببعضها البعض، لذلك فإن  $y_i$  و  $y_j$  غير مرتبطة أيضاً، أي أن:  
 $Cov(y_i, y_j) = 0$

### تقدير دالة الانحدار

لإيجاد معادلة الانحدار التقديرية لابد من إيجاد تقديرات معاملات الانحدار  $(\beta_0, \beta_1)$  وذلك باستخدام إحدى طرق التقدير المعروفة والتي تعطي نتائج تقديرية للمعالم تحمل الكثير من الصفات المرغوبة في التقديرات الاحصائية، ومن هذه الطرق:

- طريقة المربعات الصغرى Least Squares Method
- طريقة الإمكان الأعظم Maximum Likelihood

وسوف نقتصر هنا على استخدام الطريقة الأولى وهي طريقة المربعات الصغرى العادية ordinary least squares.

إن أساس طريقة المربعات الصغرى يعتمد على تقدير قيم المعالم المجهولة لنموذج الانحدار  $\beta_0$  و  $\beta_1$  والتي تجعل مجموع مربعات الأخطاء العشوائية في نهايتها الصغرى.

وبإتباع إجراءات النهايات الصغرى يمكن حساب التقديرات  $b_0$  و  $b_1$  وذلك باستخدام أسلوب التفاضل الجزئي بالنسبة إلى كل من  $B_0$  و  $B_1$  ثم نجعل هذه المشتقات الجزئية مساوية للصفر. وبحل المعادلات الآتية بالنسبة إلى  $\beta_0$  و  $\beta_1$  نحصل على ما يلي:

$$b_0 = \bar{y} - b_1 \bar{x}$$

$$b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

ولغرض تحليل التباين ومعرفة معنوية نموذج الانحدار الخطي البسيط يمكننا استخدام العلاقات الآتية:

مجموع المربعات الكلي: مجموع مربعات الانحدار + مجموع مربعات الخطأ البسيط.

أي أن:

$$SST = SSR + SSE$$

حيث أن:

$$SST = \sum (y_i - \bar{y})^2 = \sum y_i^2 - n\bar{y}^2$$

$$SSR = \sum (\hat{y}_i - \bar{y})^2 = \sum \hat{y}_i^2 - n\bar{y}^2 = b_1(\sum x_i y_i - n\bar{x}\bar{y})$$

$$SSE = \sum (y_i - \hat{y}_i)^2 = SST - SSR$$

وبذلك فإن جدول تحليل التباين سوف يكون كالآتي:

S.V.	df	SS	MS	F
Regression	1	$SSR = b_1(\sum x_i y_i - n\bar{x}\bar{y})$	$MSR = SSR/1$	$MSR/MSE$
Error	n-2	$SSE = \text{by Sub.}$	$MSE = SSE/(n-2)$	
Total	n-1	$SST = \sum y_i^2 - n\bar{y}^2$		

ومن جدول تحليل التباين أعلاه يمكن أن يتضح لنا مدى معنوية نموذج الإنحدار أو هل أن المتغير المستقل له تأثير معنوي (جوهري) على المتغير المعتمد وذلك بإختبار الفرضية التالية:

$$H_0 : \beta_1 = 0.0$$

$$H_1 : \beta_1 \neq 0.0 \quad \text{مقابل الفرضية البديلة:}$$

ويتم ذلك بإستخدام اختبار F والذي تتضح قيمته في الجدول أعلاه.

أو إستخدام اختبار t حسب الصيغة التالية:

$$T = b_1 / \sqrt{\text{Var}(b_1)}$$

### معامل التحديد Coefficient of Determination

يمكننا الحصول على معامل التحديد الذي يمثل نسبة الإنحرافات الكلية الموضحة أو المشروحة بواسطة معادلة الإنحدار التقديرية أو نسبة مساهمة معادلة الإنحدار التقديرية في تفسير أو شرح الإنحرافات الكلية عن قيم  $y$  حول الوسط الحسابي  $\bar{y}$  وهو يأخذ الصيغة التالية:

$$R^2 = SSR/SST = \sum (\hat{y}_i - \bar{y})^2 / \sum (y_i - \bar{y})^2 \quad , \quad 0 \leq R^2 \leq 1$$

## ثانياً: نموذج الإنحدار الخطي المتعدد

إن نموذج الإنحدار الخطي البسيط يقتصر فقط على تحليل العلاقة بين متغيرين أحدهما مستقل والآخر معتمد ولكن في الواقع أن أغلب الظواهر الإقتصادية والإجتماعية وغيرها بحاجة إلى وضع المتغير المعتمد كدالة لأكثر من متغير مستقل. ففي مثل هذه العلاقات يجب تمثيلها بنماذج خطية متعددة وأن المعادلة التي تمثل هذه العلاقة تنتمي لنموذج الإنحدار الخطي المتعدد وفيه نفترض وجود علاقة خطية بين أحد المتغيرات  $y$  (المتغير المعتمد) وعدد  $k$  من المتغيرات المستقلة  $(X_1, X_2, \dots, X_k)$  والتي تدعى بالمتغيرات التوضيحية بسبب إستخدامها في توضيح التباعد في المتغير  $y$ . وكل هذه المتغيرات تظهر بمقياس كمي نسبي في الغالب ولكن يمكن التعامل أيضاً مع المقياس الفئوي. بالإضافة إلى ذلك، فإن بعض المتغيرات المستقلة يمكن أن تظهر بمقياس رتبوي أو حتى إسمي ويطلق عليها عندئذٍ بالمتغيرات الوهمية *Dummy Variables* وهذا سيضيف بعض الصعوبات في تفسير النتائج. وسنتناول مثالاً بهذه الحالة نهاية هذا الفصل لتوضيح طريقة تحليل النتائج المرتبطة بها. ويمكن صياغة نموذج الإنحدار الخطي المتعدد بالشكل التالي:

$$y = f[(x_1, x_2, \dots, x_k), e]$$

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + e_i, \quad i = 1, 2, \dots, n$$

حيث نلاحظ وجود عدد  $k$  من المتغيرات المستقلة و  $(k+1)$  من المعامل و  $n$  من المشاهدات.

ولغرض تبسيط عرض النموذج الخطي المتعدد سوف نستخدم أسلوب المصفوفات في عرض مشاهدات المتغيرات المستقلة والمتغير المعتمد إضافة إلى إستخدام هذا الأسلوب عند تقدير معالم النموذج الخطي المتعدد. لذلك يمكن التعبير عن النموذج أعلاه وكما يلي:

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1k} \\ 1 & x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nk} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_n \end{bmatrix} + \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix}$$

حيث أن:

$Y$  قيمة مشاهدات المتغير المعتمد وهو متجه ذو رتبة  $(n,1)$ .

$X$  مصفوفة مشاهدات المتغيرات المستقلة ذو رتبة  $(n,k+1)$ .

$\beta$  قيمة لمعالم النموذج الخطي المتعدد المطلوب تقديرها وهو متجه ذو رتبة  $(k+1,1)$ .

$e$  قيمة للأخطاء وهو متجه ذو رتبة  $(n,1)$ .

ولغرض التفكير بأسلوب ملائم لتقدير معالم النموذج يجب علينا النظر في الفروض الأساسية الخاصة بمتجه الأخطاء العشوائية والتي يفترض أن تنطبق على جميع المشاهدات أي أن:

$$e_i \sim N(0, \sigma^2)$$

وهي أن قيم الأخطاء تتوزع توزيعاً طبيعياً بتوقع (متوسط) قدره:

$$E(\mathbf{e}) = E \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix} = \begin{bmatrix} E(e_1) \\ E(e_2) \\ \vdots \\ E(e_n) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

وتباين ثابت  $\sigma^2 I_n$  ويمكن وضعه بالصيغة التالية:

$$\begin{aligned} \text{Var}(\mathbf{e}) &= E(\mathbf{e} \cdot \mathbf{e}') = E \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix} [e_1 \ e_2 \ \dots \ e_n] \\ &= E \begin{bmatrix} e_1^2 & e_1 e_2 & \dots & e_1 e_n \\ e_2 e_1 & e_2^2 & \dots & e_2 e_n \\ \vdots & \vdots & \ddots & \vdots \\ e_n e_1 & e_n e_2 & \dots & e_n^2 \end{bmatrix} \\ &= \begin{bmatrix} E(e_1^2) & E(e_1 e_2) & \dots & E(e_1 e_n) \\ E(e_2 e_1) & E(e_2^2) & \dots & E(e_2 e_n) \\ \vdots & \vdots & \ddots & \vdots \\ E(e_n e_1) & E(e_n e_2) & \dots & E(e_n^2) \end{bmatrix} \\ &= \begin{bmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_n^2 \end{bmatrix} = \sigma^2 I_n \end{aligned}$$

وذلك لأنه، وحسب الافتراض بأن التباين ثابت، أي:

$$\sigma_1^2 = \sigma_2^2 = \dots = \sigma_n^2 = \sigma^2$$

بالإضافة إلى ذلك هناك فروض أخرى يجب توفرها في نموذج الإنحدار الخطي المقصود هي:

- (1) عدم وجود علاقة خطية محددة أو تامة بين المتغيرات المستقلة.
- (2) يجب أن يكون عدد المشاهدات أكبر من عدد المتغيرات المستقلة.

## تقدير معالم نموذج الإنحدار الخطي المتعدد

سنقوم باستخدام طريقة المربعات الصغرى لتقدير معالم النموذج لأنها تهدف إلى جعل مجموع مربعات الأخطاء أقل ما يمكن. (الرموز في أدناه هي مصفوفات مثل  $X$  والبقية متجهات):

$$y = X\beta + e$$

$$e = y - \hat{y} = y - Xb$$

$$\begin{aligned} e'e &= (y - Xb)'(y - Xb) \\ &= y'y - y'Xb - b'X'y + b'X'Xb \\ &= y'y - 2b'X'y + b'X'Xb \end{aligned}$$

$$\frac{\partial e'e}{\partial b} = -2X'y + 2X'Xb = 0.0$$

وبذلك نحصل على:

$$X'y = X'Xb$$

$$b = (X'X)^{-1} X'y$$

ويمكننا إيجاد مكونات العلاقة السابقة عن طريق إجراء العمليات الإعتيادية على المصفوفات وكما يلي:

$$X'X = \begin{bmatrix} n & \sum x_1 & \sum x_2 & \cdots & \sum x_k \\ \sum x_1 & \sum x_1^2 & \sum x_1 x_2 & \cdots & \sum x_1 x_k \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \sum x_k & \sum x_k x_1 & \sum x_k x_2 & \cdots & \sum x_k^2 \end{bmatrix}$$

$$X'y = \begin{bmatrix} \sum y \\ \sum x_1 y \\ \sum x_2 y \\ \vdots \\ \sum x_k y \end{bmatrix}$$



## سمات مقدرات طريقة المربعات الصغرى

وتمتاز مقدرات المربعات الصغرى بصفات كثيرة إلا أن أهمها:

1- عدم التحيز unbiasedness أي أن:

$$E(\hat{\beta}) = E(b) = \beta$$

وهذا واضح من خلال حقيقة كون:

$$\hat{\beta} = b = (X'X)^{-1} X'y$$

$$E(y) = X\beta$$

وبالتالي يكون لدينا:

$$\begin{aligned} E(\hat{\beta}) &= E((X'X)^{-1} X'y) \\ &= ((X'X)^{-1} X'E(y)) \\ &= (X'X)^{-1} X'X\beta \\ &= \beta \end{aligned}$$

2- أقل تباين minimum variance

لكي نثبت أن التقديرات بطريقة المربعات الصغرى لها أقل تباين، لابد أن نقوم أولاً بتقدير التباين للمتجه  $\hat{\beta} = b$  وكما يلي:

$$\begin{aligned} \text{var}(\hat{\beta}) &= \text{var}(b) = E(\hat{\beta} - \beta)^2 \\ &= E(\hat{\beta} - \beta)(\hat{\beta} - \beta)' \\ &= E[(X'X)^{-1} X'e][e'X(X'X)^{-1}] \\ &= (X'X)^{-1} X'E(e e')X(X'X)^{-1} \\ &= \sigma^2 I_n (X'X)^{-1} X X (X'X)^{-1} \\ &= \sigma^2 I_n (X'X)^{-1} \end{aligned}$$

وبذلك يكون توزيع المتجه  $\hat{\beta} = b$  هو طبيعي وكما يلي:

$$\hat{\beta} \sim N(\beta, \sigma^2 I_n (X'X)^{-1})$$

ولكي نتأكد من أن هذا التقدير لتباين  $\hat{\beta}$  هو الأقل، سنلاحظ ذلك من خلال الآتي:

لنفترض أن  $b^*$  تمثل أي تقدير آخر غير متحيز للمتجه  $\beta$  حصلنا عليه بطريقة أخرى غير طريقة المربعات الصغرى:

$$b^* = [(X'X)^{-1} X' + D]y$$

حيث أن D مصفوفة ثوابت وبذلك من الواضح أن يكون:

$$\begin{aligned}\text{Var}(b^*) &= [(X'X)^{-1} X' + D] \text{var}(y) [(X'X)^{-1} X' + D]' \\ &= \sigma^2 I_n [(X'X)^{-1} + DX (X'X)^{-1} + (X'X)^{-1} X' D' + DD']\end{aligned}$$

وحتى لو افترضنا أن  $DX = 0.0$  فإن ذلك يبقى مساوياً إلى:

$$= \sigma^2 I_n [(X'X)^{-1} + DD']$$

وهذا يثبت أن تقديرات المربعات الصغرى لها أقل تباين.

### تحليل التباين لنموذج الانحدار الخطي المتعدد

سبق وأن أوضحنا في الانحدار الخطي البسيط أن مجموع مربعات الانحرافات الكلية يمكن أن يجرأ إلى قسمين هما مجموع مربعات الانحدار ومجموع مربعات الأخطاء حيث:

$$SST = SSR + SSE$$

$$SST = y'y - n\bar{y}^2$$

$$SSR = \hat{y}'\hat{y} - n\bar{y}^2 = (xb)'xb - n\bar{y}^2$$

$$= b'x'x(x'x)^{-1}x'y - n\bar{y}^2$$

$$= b'x'y - n\bar{y}^2$$

$$SSE = SST - SSR$$

$$= y'y - b'x'y$$

وبذلك يكون جدول تحليل التباين كما يلي:

S.V.	df	SS	MS	F
Regression	k-1	$SSR = b'x'y - n\bar{y}^2$	$MSR = SSR/(k-1)$	$MSR/MSE$
Error	n-k	$SSE = y'y - b'x'y$	$MSE = SSE/(n-k)$	
Total	n-1	$SST = y'y - n\bar{y}^2$		

### إختبار معنوية الانحدار

يمكننا أن نميز نوعين من الإختبارات:

1. إختبار معنوية الانحدار في حالة النموذج الكامل (جميع المعلمات بضمنها  $\beta_0$ ):

$$H_0: \beta = 0$$

$$H_1: \beta \neq 0$$

وتقارن قيمة  $F = \frac{\hat{\beta}'x'x\hat{\beta} / k}{SSE / (n - k)}$  مع  $f(k, n - k)$  الجدولية.

2. إختبار معنوية الإنحدار في حالة النموذج المختزل (بدون  $\beta_0$ ):

$$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$$

$$H_1: \beta_i \neq 0, \quad i=1,2,\dots,k$$

وُتقارن قيمة  $F = \frac{\hat{\beta}'x'x\hat{\beta}/(k-1)}{SSE/(n-k)}$  مع  $f(k-1, n-k)$  الجدولية وهو في العادة ما

يعتمد لأن  $\beta_0$  لا تشكل شيئاً مهماً في مساهمة المتغيرات المستقلة في النموذج.

## معامل التحديد $R^2$

لاحظنا سابقاً أن نموذج الإنحدار الخطي المتعدد يمثل العلاقة بين المتغير المعتمد وعدد من المتغيرات المستقلة. كما أن مجموع المربعات الكلي لهذا النموذج يتكون من مجموع مربعات الإنحدار للمتغيرات المستقلة مضافاً إليها مجموع مربعات الأخطاء (غير المُفسرة). وإذا إفترضنا أن معادلة الإنحدار تمثل هذه العلاقة تمثيلاً جيداً فإنه من الضروري أن تكون نسبة مجموع مربعات الإنحدار إلى مجموع المربعات الكلي كبيرة وهذا ما يسمى بمعامل التحديد  $R^2$  والذي يمثل نسبة مساهمة المتغيرات المستقلة في تفسير التباين في المتغير المعتمد  $Y$  وحسب الصيغة التالية:

$$SST = SSR + SSE$$

$$1 = SSR/SST + SSE/SST$$

$$R^2 = 1 - SSE/SST$$

وبالتالي فإن:

$$R^2 = SSR/SST = \frac{b'x'y - n\bar{y}^2}{y'y - n\bar{y}^2}$$

## معامل التحديد المصحح ( $\text{adjusted } R^2$ )

يمتاز معامل التحديد  $R^2$  بأنه لو أضيف متغير مستقل إلى النموذج فإن قيمته سترتفع حتى وإن لم يكن للمتغير المضاف من الأهمية ما يستحق معها إدخاله في النموذج. ولذا ولغرض الحصول على معيار أفضل لقياس مدى قابلية مجاميع مختلفة من المتغيرات لتحليل العلاقة قيد الدراسة وبنفس الوقت الأخذ بنظر الإعتبار عدد المتغيرات المشمولة، فإنه يتم حساب ما يسمى بمعامل التحديد المصحح بموجب الصيغة التالية:

$$\bar{R}^2 = 1 - \frac{\sum e_i^2 / (n-k)}{\sum (y_i - \bar{y})^2 / (n-1)}$$

وبذلك فإن العلاقة بين معامل التحديد المصحح ومعامل التحديد غير المصحح هي:

$$\bar{R}^2 = 1 - \frac{n-1}{n-k}(1-R^2)$$

ويلاحظ أن قيمة  $\bar{R}^2$  سوف تنخفض عند إضافة متغير مستقل إذا لم تؤدي هذه الإضافة إلى تقليص قيمة  $(1-R^2)$  بما يعوض عن الزيادة التي تحصل في  $\frac{n-1}{n-k}$  نتيجة لإرتفاع قيمة  $K$ .  
وبعبارة أخرى، من الأفضل عدم إضافة متغير إلى النموذج إذا تسببت إضافته إلى تخفيض قيمة  $\bar{R}^2$ .

### الإرتباط The Correlation

كما سبق وذكرنا في فرضيات نموذج الإنحدار الخطي المتعدد أنه يجب أن تكون العلاقة بين المتغيرات التوضيحية معدومة حتى لا تنتج لدينا مشكلة تعدد الإرتباط الخطي Multicollinearity في الإنحدار. وبالمقابل، يجب أن تكون هناك علاقة قوية بين المتغيرات التوضيحية من جهة والمتغير المعتمد من جهة أخرى لكي تكون معادلة الإنحدار الخطي المتعدد كفوءة في تفسير الظاهرة المدروسة.

ولدراسة هذه الإرتباطات يمكن تحديد نوعين أساسيين هما:

1. إرتباط المتغير المستقل  $X$  مع المتغير المعتمد  $y$  وهو ما يسمى بالإرتباط البسيط.
2. إرتباط المتغيرات المستقلة بمجموعها مع المتغير المعتمد وهو ما يسمى بالإرتباط الخطي المتعدد.

ويمكن التعبير عن الإرتباط الخطي البسيط بالصيغ التالية:

$$\begin{aligned} r_{xy} &= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}} \\ &= \frac{s_{xy}}{\sqrt{s_{xx}s_{yy}}} \\ &= \frac{Cov(x, y)}{\sqrt{V(x)V(y)}} \end{aligned}$$

### معامل الإرتباط الجزئي

يمثل معامل الإرتباط الجزئي صافي الإرتباط بين المتغير المعتمد والمتغير المستقل بعد حذف التأثير المشترك لباقي المتغيرات المستقلة على كل من المتغير المعتمد والمتغير المستقل (أي بعد تثبيت المتغيرات الأخرى). مثلاً :

$r_{y1.2}$  أو  $r_{yx_1, x_2}$  يعني الارتباط الجزئي بين  $Y$  و  $X_1$  بعد حذف تأثير  $X_2$  على كل من  $X_1$  و  $Y$  ، ويستخدم لتحديد الأهمية النسبية للمتغيرات المستقلة المختلفة في الإنحدار المتعدد علماً أن قيمته، مثلما هي لأي معامل ارتباط، تنحصر ما بين (+1 و -1) ويأخذ إشارة المعلمة المناظرة ويتم احتسابه بموجب الصيغة التالية:

$$r_{y1.2} = \frac{r_{y1} - r_{12}r_{y2}}{\sqrt{(1 - r_{12}^2)(1 - r_{y2}^2)}}$$

أو يمكن حسابه من جدول تحليل التباين بالصيغة:

$$|r_{y1.2}| = \sqrt{\frac{SSR(X_1 | X_2)}{SST - SSR(X_2)}}$$

أما الارتباطات الجزئية مع إضافة متغيرات مستقلة أخرى فتحتسب معاملاتها على نفس النمط وكما يلي:

في حالة ثلاثة متغيرات مستقلة تكون:

$$r_{y1.23} = \frac{r_{y1.3} - r_{12.3}r_{y2.3}}{\sqrt{(1 - r_{12.3}^2)(1 - r_{y2.3}^2)}} = \frac{r_{y1.2} - r_{13.2}r_{y3.2}}{\sqrt{(1 - r_{13.2}^2)(1 - r_{y3.2}^2)}}$$

أو

$$|r_{y1.23}| = \sqrt{\frac{SSR(X_1 | X_2, X_3)}{SST - SSR(X_2, X_3)}}$$

وفي حالة أربعة متغيرات مستقلة تكون:

$$r_{y1.234} = \frac{r_{y1.34} - r_{12.34}r_{y2.34}}{\sqrt{(1 - r_{12.34}^2)(1 - r_{y2.34}^2)}}$$

أو

$$|r_{y1.234}| = \sqrt{\frac{SSR(X_1 | X_2, X_3, X_4)}{SST - SSR(X_2, X_3, X_4)}}$$

وبصورة عامة فإن معامل الارتباط الجزئي بين المتغيرين (i و j) بعد جعل جميع تأثيرات المتغيرات الأخرى ثابتة هو:

$$r_{ij.(all\ other\ variables)} = \frac{-C_{ij}}{\sqrt{C_{ii}C_{jj}}}$$

حيث أن  $C_{ii}$  و  $C_{jj}$  و  $C_{ij}$  هي عناصر معكوس مصفوفة معامل الارتباط.

ولأجل توضيح ذلك، لنفترض ثلاثة متغيرات مستقلة  $X_1, X_2, X_3$  إلى جانب المتغير المعتمد  $y$ . في هذه الحالة، تكون لدينا مصفوفة الارتباط:

$$R = \begin{vmatrix} 1 & r_{12} & r_{13} & r_{1y} \\ r_{21} & 1 & r_{23} & r_{2y} \\ r_{31} & r_{32} & 1 & r_{3y} \\ r_{y1} & r_{y2} & r_{y3} & 1 \end{vmatrix}$$

ولذلك فإن معكوس هذه المصفوفة:

$$R^{-1} = \begin{vmatrix} C_{11} & C_{12} & C_{13} & C_{1y} \\ C_{21} & C_{22} & C_{23} & C_{2y} \\ C_{31} & C_{32} & C_{33} & C_{3y} \\ C_{y1} & C_{y2} & C_{y3} & C_{yy} \end{vmatrix}$$

وبذلك، وعلى سبيل المثال، فإن:

$$r_{23.1y} = \frac{-C_{23}}{\sqrt{C_{22}C_{33}}}$$

أما بالنسبة لمعامل الارتباط المتعدد والذي يرمز له برمز  $R_{y.12\dots k}$  فإنه يقيس مدى قرب سطح الإنحدار من النقاط المشاهدة، وهذا يعني أن معامل الارتباط المتعدد يقيس التأثير المشترك لجميع المتغيرات المستقلة على المتغير التابع. ويُعرّف معامل الارتباط المتعدد على النحو التالي:

$$R_{y.12\dots k} = \sqrt{\frac{SSR}{SST}} = \sqrt{\frac{b'x'y - n\bar{y}^2}{y'y - n\bar{y}^2}}$$

أيضاً يمكن كتابته (في حالة متغيرين مستقلين وثلاثة متغيرات مستقلة على سبيل المثال) كما يلي:

$$R_{y.12} = \left[ 1 - (1 - r_{y1}^2)(1 - r_{y2.1}^2) \right]^{1/2}$$

ويمكن عمل الاختبارات التالية:

1. بالنسبة للإرتباط المتعدد:

$$H_0: r_{y.123 \dots p} = 0$$

$$H_1: r_{y.123 \dots p} > 0$$

ونرفض  $H_0$  إذا كان:

$$\left( \frac{r_{y.123 \dots p}^2}{1 - r_{y.123 \dots p}^2} \right) \left( \frac{n-p}{p-1} \right) > f_{\alpha, p-1, n-p}$$

2. بالنسبة للإرتباط الجزئي (حالة ثلاثة متغيرات مستقلة):

$$H_0: r_{y2.13} = 0$$

$$H_1: r_{y2.13} \neq 0$$

باستخدام إحصاءة الاختبار:

$$T = \frac{r_{y2.13}}{\sqrt{\frac{1 - r_{y2.13}^2}{n-3}}} \sim t_{n-3}$$

### إختيار أفضل معادلة إنحدار

إن أحد أصعب مسائل تحليل الإنحدار هي إختيار مجموعة المتغيرات المستقلة المتضمنة في النموذج. وهناك عدد من الطرق التي تساعد في إيجاد أفضل مجموعة من المتغيرات المستقلة. وهذه الطرق بصورة عامة يفضل استخدامها مع الحاسب الآلي لأنها تحتاج إلى عمليات حسابية مطولة جداً وخاصة في حالة وجود عدد كبير من المتغيرات المستقلة. ومن هذه الطرق ما يلي:

1. طريقة الخطوات المتسلسلة Stepwise
2. طريقة كل الإنحدارات الممكنة All possible regression.
3. طريقة الإختيار الأمامي أو المباشر Forward method.
4. طريقة الحذف المعاكس Backward method.

وسنشرح هنا طريقة الإختيار الأمامي أو المباشر Forward method وتتلخص بما يلي:

نبدأ المعادلة بدون أي متغير مستقل ثم نختار المتغيرات المستقلة التي تدخل للمعادلة واحداً بعد الآخر ونتوقف عن الإختيار عندما تقل قيمة F الجزئية عن قيمة f الجدولية المقابلة.

وأول المتغيرات الذي يدخل المعادلة هو المتغير الذي له أعلى قيمة F محسوبة وتزيد عن قيمة f الجدولية. المتغير الثاني الذي يدخل المعادلة أعلاه هو المتغير الذي له أعلى قيمة F جزئية بوجود المتغير الأول المنتخب بالخطوة الأولى وتزيد عن قيمة f الجدولية المعينة لتلك الخطوة. وهكذا نستمر بإضافة المتغير الذي له أعلى قيمة F جزئية وتزيد عن f الجدولية إلى أن نصل إلى أعلى قيمة F جزئية تقل عن f الجدولية فعند ذلك نتوقف عن الإضافة. وستطبق هذه الطريقة على المثال التالي ببيانات حقيقية لتجربة مختبرية: (1)

المشاهدات	Y	X <sub>1</sub>	X <sub>2</sub>	X <sub>3</sub>	X <sub>4</sub>
1	78.5	7	26	6	60
2	74.5	1	29	15	52
3	104.3	11	56	8	20
4	87.6	11	31	8	47
5	95.9	7	52	6	33
6	109.2	11	55	9	22
7	102.7	3	71	17	6
8	72.5	1	31	22	44
9	93.1	2	59	18	22
10	115.9	21	47	4	26
11	83.8	1	40	23	34
12	113.5	11	66	9	12
13	109.4	10	68	8	12

والذي يمثل تجربة مختبرية لدراسة تأثير بعض العناصر الفلزية على درجة التوصيل الحراري لقضيب معدني.

حيث أن:

Y تمثل الحرارة المنبعثة

X<sub>1</sub> كميات الومينات الكالسيوم الثلاثي

X<sub>2</sub> كميات سليكات الكالسيوم الثلاثي

X<sub>3</sub> الومينات الكالسيوم الرباعي الحديدي

X<sub>4</sub> كمية سليكات الكالسيوم الثنائي

وجميع هذه المتغيرات بقياس كمي نسبي.



الخطوة الأولى:

أ- نحسب قيمة F لانحدار Y على كل من  $X_1$  و  $X_2$  و  $X_3$  و  $X_4$  وبشكل منفرد كما في الجداول التالية:

تحليل التباين لانحدار Y على  $X_1$

S.V	DF	SS	MS	F
R( $X_1$ )	1	1950.0769	1950.0769	
ERROR( $X_1$ )	11	1265.6667	115.0629	<b>12.60</b>
TOTAL	12	2715.7631		

تحليل التباين لانحدار Y على  $X_2$

S.V	DF	SS	MS	F
R( $X_2$ )	1	1809.9268	1809.9268	
ERROR( $X_2$ )	11	906.3363	82.3942	<b>21.96</b>
TOTAL	12	2715.7631		

تحليل التباين لانحدار Y على  $X_3$

S.V	DF	SS	MS	F
R( $X_3$ )	1	776.3626	776.3626	
ERROR( $X_3$ )	11	1939.9005	176.3091	<b>4.40</b>
TOTAL	12	2715.7631		

تحليل التباين لانحدار Y على  $X_4$

S.V	DF	SS	MS	F
R( $X_4$ )	1	1831.8962	1831.8962	
ERROR( $X_4$ )	11	883.8669	80.3515	<b>22.80</b>
TOTAL	12	2715.7631		

ب- إن أول متغير يُنتخب ليدخل المعادلة هو  $X_4$  (كمية سليكات الكالسيوم الثنائي) لأنه أعلى قيمة F وهي تزيد عن قيمة:

$$F=22.80 > f_{0.01, (1,11)} = 3.23$$

أي أن كمية سليكات الكالسيوم الثنائي له التأثير الأوضح على درجة التوصيل الحراري.

ملاحظة:

لو كانت قيمة أعلى F أقل من قيمة f الجدولية هذه، فعندئذ نتوقف ونقول بأنه لا يوجد أي متغير مستقل له أي تأثير معنوي على معدل Y.

### الخطوة الثانية:

أ- ولإنتخاب المتغير الثاني لإدخاله في المعادلة التي تضم المتغير المنتخب الأول  $X_4$ ، نحسب قيمة F الجزئية لكل متغير آخر بوجود المتغير  $X_4$ . أي نحسب قيم F الجزئية التالية:

$$F(X_1 | X_4) = \frac{MSR(X_1 | X_4)}{MSE(X_1, X_4)}$$

$$F(X_2 | X_4) = \frac{MSR(X_2 | X_4)}{MSE(X_2, X_4)}$$

$$F(X_3 | X_4) = \frac{MSR(X_3 | X_4)}{MSE(X_3, X_4)}$$

قيمة F الجزئية للمتغير  $X_1$  بوجود  $X_4$

S.V	DF	SS	MS	F
R(X <sub>1</sub> , X <sub>4</sub> )	2	2641.0010		
R(X <sub>4</sub> )	1	1831.8962		
R(X <sub>1</sub>   X <sub>4</sub> )	1	809.1048	809.1048	<b>108.22</b>
ERROR(X <sub>1</sub> , X <sub>4</sub> )	10	74.7621	7.4762	
TOTAL	12	2715.7631		

قيمة F الجزئية للمتغير  $X_2$  بوجود  $X_4$

S.V	DF	SS	MS	F
R(X <sub>2</sub> , X <sub>4</sub> )	2	1846.8830		
R(X <sub>4</sub> )	1	1831.8962		
R(X <sub>2</sub>   X <sub>4</sub> )	1	14.9888	19.9888	<b>0.1725</b>
ERROR(X <sub>2</sub> , X <sub>4</sub> )	10	868.8801	86.8880	
TOTAL	12	2715.7631		

قيمة F الجزئية للمتغير  $X_3$  بوجود  $X_4$

S.V	DF	SS	MS	F
R(X <sub>3</sub> , X <sub>4</sub> )	2	2590.0251		
R(X <sub>4</sub> )	1	1831.8962		
R(X <sub>3</sub>   X <sub>4</sub> )	1	708.1289	708.1289	<b>40.29</b>
ERROR(X <sub>3</sub> , X <sub>4</sub> )	10	175.7380	17.5738	
TOTAL	12	2715.7631		

ب- نختار المتغير المستقل الذي له أعلى F جزئية والتي هي (108.22) والخاصة بالمتغير  $X_1$  (كميات الومينات الكالسيوم الثلاثي) التي تزيد عن القيمة الجدولية المقابلة لها حيث أن:

$$F=108.22 > f_{0.01, (1,10)} = 3.23$$

لذا فإن  $X_1$  نختارها وندخلها في المعادلة التي تحتوي على  $X_4$  ثم نحسب هذه المعادلة التي كانت:

$$\hat{Y} = 103.097 + 1.940X_1 - 6.614X_4$$

### الخطوة الثالثة:

نحسب قيمة F الجزئية لكل من المتغيرات المستقلة الباقية بوجود  $X_1, X_4$  أي نحسب F  $(X_2|X_1, X_4)$  ،  $F(X_3|X_1, X_4)$  كما في الجدولين التاليين:

قيمة F الجزئية للمتغير  $X_2$  بوجود  $X_1, X_4$

S.V	DF	SS	MS	F
R( $X_1, X_2, X_4$ )	3	2667.7904		
R( $X_1, X_4$ )	2	2641.0110		
R( $X_2   X_1, X_4$ )	1	26.77	26.7744	<b>5.02</b>
ERROR( $X_1, X_2, X_4$ )	9	47.9727	5.3303	
TOTAL	12	2715.7631		

قيمة F الجزئية للمتغير  $X_3$  بوجود  $X_1, X_4$

S.V	DF	SS	MS	F
R( $X_1, X_3, X_4$ )	3	2664.9270		
R( $X_1, X_4$ )	2	2641.0110		
R( $X_3   X_1, X_4$ )	1	23.9160	23.9160	<b>4.23</b>
ERROR( $X_1, X_3, X_4$ )	9	50.8361	5.6485	
TOTAL	12	2715.7631		

بما أن المتغير  $X_2$  له أعلى قيمة F جزئية (5.02) حيث:

$$F = 5.02 > f_{0.01, (1, 9)} = 3.23$$

لذا نختار  $X_2$  ليضاف للمعادلة السابقة لتصبح لدينا المعادلة الجديدة:

$$\hat{Y} = 71.6480 + 1.4519X_1 + 0.4161X_2 - 0.2365X_4$$

#### الخطوة الرابعة:

نحسب قيمة F الجزئية للمتغير الباقي  $X_3$  مع وجود  $X_1, X_2, X_4$  في المعادلة أي أننا نحسب

قيمة  $F(X_3 | X_1, X_2, X_4)$  وحسب الآتي:

قيمة F الجزئية إلى  $X_3$  بوجود  $X_1, X_2, X_4$

S.V	DF	SS	MS	F
$R(X_1, X_2, X_3, X_4)$	4	2667.8995		
$R(X_1, X_2, X_4)$	3	2667.7904		
$R(X_3   X_1, X_2, X_4)$	1	0.1091	0.1091	<b>0.02</b>
$ERROR(X_1, X_2, X_3, X_4)$	8	47.8637	5.9830	
TOTAL	12	2715.7631		

وبما أن قيمة F الجزئية هنا إلى  $X_3$  صغيرة (0.02) وهي أقل من الجدولية المقابلة لها، أي أن:

$$F = 0.02 < f_{0.01, (1, 8)} = 3.46$$

وهذا يعني أن تأثير المتغير  $X_3$  (الومينات الكالسيوم الرباعي الحديدي) غير معنوي لذا فإنه لا يدخل إلى معادلة الإنحدار.

وبناءً على ذلك، فإن أفضل معادلة إنحدار في هذه الحالة هي:

$$\hat{Y} = 71.6480 + 1.4519X_1 + 0.4161X_2 - 0.2365X_4$$

وباعتماد هذا النموذج، فإنه بالإمكان إحتساب قيمة معامل التحديد  $R^2$  وفقاً لذلك لتكون

$$R^2 = SSR/SST = 2667.7904/2715.7631 = 0.9823$$

أي أن حوالي 98.23 % من التباينات الموجودة في المتغير المعتمد  $Y$  (الحرارة المنبعثة) تعود أسبابها إلى تأثير المتغيرات الثلاثة  $X_1$  (كمية الومينات الكالسيوم الثلاثي) و  $X_2$  (كمية سليكات الكالسيوم الثلاثي) و  $X_4$  (كمية سليكات الكالسيوم الثنائي).

وقد نقوم بتحديد قيمة معامل التحديد المصحح  $\bar{R}^2$  وفقاً للصيغة:

$$\begin{aligned}\bar{R}^2 &= 1 - \frac{n-1}{n-k} (1 - R^2) \\ &= 1 - \frac{12-1}{12-3} (1 - 0.9823) \\ &= 0.9784 \quad (97.84\%) \end{aligned}$$

وقد لا نلاحظ تغيير كبير في القيمتين لكون عدد المتغيرات المستقلة المعتمدة ليس كبير جداً وبعيداً عن الواحد. أي أنه لا فروقات كبيرة ما بين درجة الحرية الكلية  $(n - 1)$  ودرجة الحرية للخطأ التجريبي  $(n-k)$ .

والمثال التالي يتضمن بيانات فعلية لتجربة بمتغيرات مستقلة بعضها ذات قياس كمي والبعض الآخر ذات قياس غير كمي لغرض إستعراض كيفية التعامل مع مثل هذه المتغيرات. والنتائج مأخوذة عن دراسة في جامعة بوسطن الأمريكية<sup>(2)</sup> عام 2013.

**مثال:**

في دراسة للعلاقة ما بين ضغط الدم  $Y$  كمتغير معتمد ودرجة السمنة Body Mass Index (BMI) كمتغير مستقل  $X_1$  أظهرت وجود ارتباط موجب بمعنوية عالية من خلال نموذج الإندار البسيط لبيانات مأخوذة عن عينة بحجم  $(n = 3,539)$  حيث كانت النتيجة:

$$\hat{Y} = 108.28 - 0.67 (\text{BMI})$$

وفي ضوء ذلك تم تطوير الدراسة لنموذج إندار متعدد بإعتماد المتغيرات التالية كمتغيرات مستقلة:

$$\text{BMI} = X_1$$

$$\text{العمر} = X_2$$

$$X_3 = \text{الجنس} \quad (1 = \text{ذكر} / 0 = \text{أنثى})$$

$$X_4 = \text{تناول أدوية الضغط} \quad (1 = \text{نعم} / 0 = \text{كلا})$$

وبذلك يتضح لنا أن المتغيرين  $X_1$  و  $X_2$  بمقياس كمي نسبي فيما كون المتغيران  $X_3$  و  $X_4$  بمقياس إسمي.

والنتيجة كانت:

$$\hat{Y} = 68.15 + 0.58 X_1 + 0.65 X_2 + 0.94 (X_3=1) + 6.44 (X_4=1).$$

أو يمكننا كتابتها بالشكل التالي:

$$\hat{Y} = 68.15 + 0.58(\text{BMI}) + 0.65(\text{Age}) + 0.94(\text{Male}) + 6.44(\text{Trt})$$

والتي تشير (من خلال إختبار T) إلى معنوية عالية (p-value = 0.0001) لمعاملات الإندار بالنسبة إلى جميع المتغيرات عدا متغير الجنس الذكوري  $X_3$  حيث نجد أن (p-value = 0.1133).

وهنا يظهر لنا أن العمر Age هو الأعلى معنوية من بين المتغيرات المستقلة، يتبعه في ذلك BMI ومن ثم المعالجة وأخيراً الجنس الذكوري. كما نلاحظ أن قيمة معامل الإندار بالنسبة إلى المتغير BMI قد أصبحت 0.58 مقابل 0.67 أي بإنخفاض قدره 0.13 نتيجة تأثير المتغيرات الثلاثة الأخيرة على قيمة ضغط الدم.

وفيما يخص تفسير النتائج وفقاً لما جاء في أعلاه، فإنها تشير إلى أن زيادة BMI وحدة واحدة يرتبط بزيادة في ضغط الدم بمقدار 0.58 من الوحدات بتثبيت قيم العمر والجنس الذكوري وحالة المعالجة. كذلك يمكننا القول أن الذكور لديهم ضغط دم أعلى نسبياً من الإناث بنحو 0.94 من الوحدات بتثبيت قيمة BMI والعمر وحالة المعالجة.

وباعتماد معادلة الإندار أعلاه نستطيع، بطبيعة الحال، من تقدير قيمة ضغط الدم كدالة لمجموعة المتغيرات المستقلة الأربعة وبالشكل الذي نريده. وعلى سبيل المثال، نجد أن القيمة التقديرية لضغط الدم لدى شخص ذكر عمره 50 سنة ويتسم بمقياس بدانة BMI = 25 ولا يخضع لعلاج هي:

$$\hat{Y} = 68.15 + 0.58 (25) + 0.65 (50) + 0.94 (1) + 6.44 (0) = 116.09$$

بينما هذه القيمة لدى أنثى تخضع للعلاج وبنفس العمر ومقياس البدانة ستكون:

$$\hat{Y} = 68.15 + 0.58 (25) + 0.65 (50) + 0.94 (0) + 6.44 (1) = 121.59.$$

## تحليل المركبات الرئيسية

### Principal Components Analysis (PCA)

إن أول من طرح موضوع طريقة تحليل المركبات الرئيسية هو كارل بيرسون Karl Pearson وذلك عام 1901 لأهميتها آنذاك للمختصين في مجال علم الأحياء القياسي Biometrics. أعقبه هوتلنك Hottling عام 1931 بوصف طرق عملية في هذا الجانب.

إن من أهم أسباب تطبيق تحليل المركبات الرئيسية هو لإستخدامها أداة لفحص البيانات متعددة المتغيرات. ويمكن من خلالها تكوين متغيرات جديدة تسمى علامات المركبات الرئيسية Principal Components Scores حيث يمكن حساب قيمها. إن فحص وضعية الشكل الناتج غالباً ما يعطينا الإنطباع بتحقق حالة اللاتبيعي في البيانات التي نحن بصدد تحليلها.

وفي تحليل المركبات الرئيسية يتم إستخدام أسلوب رياضي يقوم على أساس تحويل مجموعة من المتغيرات التوضيحية المترابطة فيما بينها إلى مجموعة جديدة من المتغيرات غير المترابطة (أو المتعامدة Orthogonal) تدعى المركبات الرئيسية. وكل واحدٍ من هذه المتغيرات الجديدة عبارة عن توليفة رياضية خطية تضم جميع المتغيرات التوضيحية الأصلية ويمكن استخدامها كمدخلات لبرامج الرسومات والأشكال البيانية وطرق تحليل أخرى لمتعدد المتغيرات.

يمكن تنفيذ التحليل هذا باستخدام مصفوفة (التباين-التباين المشترك) -Variance-Covariance Matrix أو مصفوفة الارتباط Correlation Matrix للمتغيرات التوضيحية. وإن نوع المصفوفة المفضل استخدامها يعتمد في الغالب على طبيعة المتغيرات قيد التحليل. فإذا كانت هذه المتغيرات بوحدات متشابهة، يمكن إستخدام مصفوفة (التباين-التباين المشترك). أما إذا كانت الحالة عكس ذلك، فمن الأجدر إعتداد مصفوفة الارتباط.

بالنسبة للمتغيرات الجديدة (المركبات الرئيسية) يمكننا، أحياناً وليس دائماً، إعطاء تفسير لها. ولذلك لا يمكننا دائماً التوقع بأن نكون قادرين على تفسيرها. بل أنه عندما نعطي تفسيراً لها فإن ذلك يعتبر شيئاً إيجابياً إضافياً وغير متوقع لوظيفتها الرئيسية وهي استخدامها أداة في متابعة التحليل بطرق أخرى سواءً أمكن تفسيرها أم لا.

ولا يخفى كون تحليل المركبات الرئيسية في العادة مفيداً للباحثين الذين يريدون تقسيم مجموعة الوحدات التجريبية الرئيسية إلى مجاميع فرعية حيث أن الوحدات المتمثلة إلى حدٍ ما تكون ضمن المجموعة الفرعية الواحدة. وفي هذه الحالة، فإن علامات المركبات الرئيسية Principal Components Scores يمكن استخدامها بمثابة مدخل إلى برامج العنقدة Clustering programs مختصرين الوقت والجهد بهذا الإتجاه. والأكثر من ذلك، فإن

علامات المركبات الرئيسية يمكن، بل ويفضل، توظيفها في المساعدة للتحقق من نتائج برامج العنقدة.

### طبيعة المركبات الرئيسية :

المركبات الرئيسية عبارة عن توليفات خطية من جميع المتغيرات التوضيحية  $X_1, X_2, \dots, X_p$  تحدها المتجهات المميزة ( $a_i$ ) Eigen Vectors والتي ترتبط بالقيم المميزة Eigen Values ( $\lambda_i$ ) (أو الجذور المميزة Characteristic Roots) الناتجة من مصفوفة التباين – التباين المشترك أو مصفوفة الارتباط. ولا بد من أن نذكر هنا بأن عدد هذه المركبات الرئيسية في الأصل هو بعدد المتغيرات المستقلة. وصيغتها الرياضية بشكل عام هي:

$$PC_i = a_{i1}X_1 + a_{i2}X_2 + \dots + a_{ip}X_p, \quad i = 1, 2, \dots, p$$

ويمكن التعبير عن جميع هذه المركبات الرئيسية بصيغة المصفوفات وهي

$$\begin{bmatrix} PC_1 \\ PC_2 \\ \vdots \\ PC_p \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1p} \\ a_{21} & a_{22} & \dots & a_{2p} \\ \vdots & \vdots & & \vdots \\ a_{p1} & a_{p2} & \dots & a_{pp} \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{bmatrix} = A'X$$

حيث أن كل عمود من المصفوفة A يمثل المتجه المميز الذي يقابل القيمة المميزة المرتبط بها والناتج عنها كما سنرى ذلك لاحقاً.

ومن الجدير بالذكر أن المركبة الرئيسية الأولى ستكون عبارة عن توليفة خطية من المتغيرات التوضيحية بمعاملات المتجه المميز المقابل للقيمة المميزة الأولى  $\lambda_1$  وهي أكبر قيمة مميزة حيث أن هذه القيم تكون عند مقارنتها مع بعضها بالشكل التالي:

$$\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots \geq \lambda_p$$

ومع أن عدد المركبات الرئيسية، وكما ذكرنا، هي بعدد المتغيرات التوضيحية، إلا أننا لا نستخدم سوى عدد قليل منها سواءً في دالة الإنحدار المتعدد الخطية أو في طرق أخرى في التحليل متعدد المتغيرات. وهذه الحالة تدخل ضمن عملية تخفيض اتجاهات التحليل الإحصائي Reduction of Dimensionality للبيانات ومعطياتها كما يبدو ذلك واضحاً في بعض طرق التحليل المتعدد وبالأخص في التحليل العاملي. وهذا التخفيض عادةً ما يتم وفقاً لقرار إستبعاد المركبات الرئيسية الضعيفة من حيث تباينها المتمثل بالقيمة المميزة  $\lambda$  وبالأخص من تكون قيمتها أقل من الواحد.



وفيما يتعلق بدالة الإنحدار المتعدد الخطية، فإن الدافع لإستخدامها هو معالجة ظهور حالة التعدد الخطي (الإرتباط المتعدد) Multicollinearity فيما بين المتغيرات التوضيحية وذلك وفقاً للنموذج الآتي:

$$Y = \alpha_0 + \alpha_1 PC_1 + \alpha_2 PC_2 + \dots + \alpha_p PC_p$$

وهنا نلاحظ إحلال المركبات الرئيسية محل المتغيرات التوضيحية حيث أنها تتسم بخصوصية التعامدية والتي تعكس عدم وجود أي إرتباط فيما بينها. وهذه السمة ضرورية لأنها تحقق تقديرات صحيحة لمعاملات النموذج.

### خطوات الحسابات:

(1) نقوم باحتساب مصفوفة التباين – التباين المشترك أو مصفوفة الإرتباط للمتغيرات التوضيحية وذلك طبقاً لكل من الحالتين التاليتين  
أ- إذا كانت وحدات قياس المتغيرات التوضيحية متشابهة، فإننا نتجه إلى حساب مصفوفة التباين – التباين المشترك وهي:

$$S = \begin{bmatrix} V_{11} & V_{12} & \dots & V_{1p} \\ V_{21} & V_{22} & \dots & V_{2p} \\ \vdots & & & \\ V_{p1} & V_{p2} & \dots & V_{pp} \end{bmatrix}$$

حيث أن:

$$V_{ii} = \frac{S_{ii}}{n-1} \quad , \quad V_{ij} = \frac{S_{ij}}{n-1}$$

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

$$S_{ij} = \sum X_i X_j - \frac{(\sum X_i)(\sum X_j)}{n}$$

ب- وإذا كانت وحدات القياس هذه مختلفة، فيستحسن في هذه الحالة تحويل المتغيرات التوضيحية إلى الحالة القياسية بوسط حسابي (0) وتباين (1) وهذا يتطابق تماماً مع إستخدامنا لمصفوفة الإرتباط بدلاً من مصفوفة التباين وهي:

$$R = \begin{bmatrix} 1 & r_{12} & \dots & r_{1p} \\ r_{21} & 1 & \dots & r_{2p} \\ \vdots & & & \\ r_{p1} & r_{p2} & \dots & 1 \end{bmatrix}$$

(2) إيجاد القيم (الجذور) المميزة  $\lambda$  من خلال حل المعادلة المميزة للمصفوفة  $R$  وهي:

$$|R - \lambda I| = \begin{vmatrix} 1 - \lambda & r_{12} & \cdots & r_{1p} \\ r_{21} & 1 - \lambda & \cdots & r_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ r_{p1} & r_{p2} & \cdots & 1 - \lambda \end{vmatrix} = 0$$

والتي ستكون:

$$\lambda_1 > \lambda_2 > \lambda_3 > \cdots > \lambda_p$$

(3) إيجاد المتجه المميز الأول  $\underline{a}_1$  المقابل للقيمة المميزة الأولى  $\lambda_1$  من المعادلة

$$(R - \lambda_1)\underline{a}_1 = 0$$

ويتم إختيار قيم عناصر هذا المتجه المميز بشرط أن يتحقق لدينا:

$$\underline{a}'_1 \underline{a}_1 = 1$$

وفي هذا السياق، فقد نعلم

المتجه المميز القياسي بمثابة المتجه المميز الأول وهو:

$$\underline{a}'_1 = \left[ \frac{a_1}{\sqrt{\sum a_i^2}} \quad \frac{a_2}{\sqrt{\sum a_i^2}} \quad \cdots \quad \frac{a_p}{\sqrt{\sum a_i^2}} \right]$$

$$= [a_{11} \quad a_{12} \quad \cdots \quad a_{1p}]$$

وبذلك تكون المركبة الرئيسية الأولى بالصيغة:

$$PC_1 = a_{11}X_1 + a_{12}X_2 + \cdots + a_{1p}X_p$$

(4) إيجاد المتجه المميز الثاني  $\underline{a}_2$  المقابل للقيمة المميزة الأولى  $\lambda_2$  من المعادلة:

$$(R - \lambda_2)\underline{a}_2 = 0$$

وبذلك تكون المركبة الرئيسية الأولى بالصيغة:

$$PC_2 = a_{21}X_1 + a_{22}X_2 + \cdots + a_{2p}X_p$$

ويتم إختيار قيم عناصر هذا المتجه المميز بشرط أن يتحقق لدينا:

$$\underline{a}'_2 \underline{a}_2 = 1 \quad , \quad \underline{a}'_1 \underline{a}_2 = 0$$

والذي هو شرط إعتبار  $PC_2$  و  $PC_1$  متعامدتين.

(5) نعيد الإجراءات في الخطوة (4) أعلاه لتحديد المتجه المميز الثالث  $\underline{a}_3$  بشرط أن

يتحقق لدينا

$$\underline{a}'_3 \underline{a}_3 = 1 \quad , \quad \underline{a}'_1 \underline{a}_3 = 0 \quad , \quad \underline{a}'_2 \underline{a}_3 = 0$$

والذي هو شرط إعتبار  $PC_3$  متعامدة مع كل من  $PC_2$  و  $PC_1$

ونستمر هكذا حتى نكمل جميع المركبات الرئيسية واحدةً بعد أخرى وبشرط تحقق حالة تعامد كل واحدة جديدة مع كل ما سبقها.

## خواص المركبات الرئيسية:

من المفيد هنا معرفة الجوانب التالية حول بعض خصائص القيم (الجزور) المميزة وهي  
 (1) في حالة إستخدامنا لمصفوفة التباين-التباين المشترك فإن:

$$\sum \lambda_i = \text{Trace}(S) = \sum V_{ii}$$

(2) في حالة إستخدامنا لمصفوفة الارتباط فإن:

$$\sum \lambda_i = \text{Trace}(R) = p$$

(3) نفس الشيء بالنسبة لحاصل ضرب القيم المميزة ليكون مساوياً لمحددة المصفوفة S  
 أي أن:

$$|S| = \prod_{i=1}^p \lambda_i$$

وهي خاصية مهمة جداً في تفسيرات المركبات الرئيسية سيما إذا ما عرفنا أن  
 القيمة المميزة هي بمثابة التباين للمركبة الرئيسية المقابلة لها.

(4) إن الأهمية النسبية للمركبة الرئيسية في وصف النموذج تقاس بما يلي

$$\frac{\text{Var}(PC_i)}{\sum \text{Var}(PC_i)} = \frac{\lambda_i}{\sum \lambda_i}$$

ولذلك لو كانت هناك قيمتين مميزتين أو أكثر متساوية في القيمة، فإنها تكون  
 متساوية في الأهمية النسبية.

(5) وحيث أن المركبات الرئيسية متعامدة وبدون أي إرتباط فيما بينها، فإن مصفوفة  
 التباين لها تكون بالصيغة التالية:

$$V(PC) = \begin{bmatrix} V(PC_1) & 0 & \dots & 0 \\ & V(PC_2) & \dots & 0 \\ & & & \vdots \\ & & & V(PC_p) \end{bmatrix}$$

$$= \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ & \lambda_2 & \dots & 0 \\ & & & \vdots \\ & & & \lambda_p \end{bmatrix}$$

وبالتالي من الواضح أن يكون لدينا

$$|V(PC)| = \prod_{i=1}^p \lambda_i$$

## بعض استخدامات المركبات الرئيسية:

1. إن استخدام المركبات الرئيسية في التحليل المتعدد يساهم في تحديد العوامل الأولية من حيث أهميتها في التحليل إضافة إلى إعطاء فكرة عن مجاميع المتغيرات التي تشكل معاً بعداً خاصاً بها وكما يظهر ذلك في التحليل العاملي حيث أن كل عامل يمكن أن يحمل مسمى يعكس سمة هذه المجاميع.
2. حيث أن المركبات الرئيسية متعامدة فيما بينها جميعاً، فإن استخدامها في تحليل الانحدار المتعدد سوف يستبعد حالة التعدد الخطي (الإرتباط المتعدد) وأثره على دقة التقدير مما يؤدي إلى دقة التنبؤ.

## عيوب المركبات الرئيسية:

1. ليس للمكونات الرئيسية أي تفسير منطقي واضح في أغلب الأحيان.
2. تعتمد هذه الطريقة على افتراض التعدد الخطي ما بين المتغيرات التوضيحية. وبذلك فإن استخدامها في حالة عدم وجود هذه الحالة، سيعطي نتائج تبتعد نسبياً عن الصواب.
3. تعتمد هذه الطريقة على حدس شخصي في بقاء المكونة أو خروجها من التحليل وهو متعلق بقدر القيم (الجذور) المميزة ولا توجد أية وسيلة يمكن إعتماها بشكل ثابت في هذا الموضوع.
4. إن إستبعاد بعض المكونات الرئيسية من التحليل يعني إهمال جزء من المعلومات التي قد تكون مفيدة في التحليل النهائي.
5. إن تطبيق هذه الطريقة على القيم المعيارية يعطي نتائج مختلفة عن ما لو تم تطبيقها على المتغيرات الأصلية.
6. الإختبارات المختلفة للمتجهات المميزة (المعدلة) يعطي نوعاً ما أيضاً نتائج مختلفة.

## تحليل الإنحدار بالمكونات الرئيسية:

وكما ذكرنا سابقاً، فإن الغرض من استخدام المكونات الرئيسية في تحليل الإنحدار الخطي المتعدد هو إستبعاد الأثر السلبي للتعدد الخطي في حالة وجوده مع المتغيرات التوضيحية (المستقلة).

ومعادلة الإنحدار تكون وفقاً لذلك:

$$Y = \alpha_0 + \alpha_1 PC_1 + \alpha_2 PC_2 + \dots + \alpha_p PC_p$$
$$= \alpha_0 + \alpha_1 (a_{11}X_1 + a_{12}X_2 + \dots + a_{1p}X_p) + \alpha_2 (a_{21}X_1 + a_{22}X_2 + \dots + a_{2p}X_p) + \dots + \alpha_p (a_{p1}X_1 + a_{p2}X_2 + \dots + a_{pp}X_p)$$

$$\begin{aligned}
& \alpha_0 + (\alpha_1 a_{11} + \alpha_2 a_{21} + \dots + \alpha_p a_{p1}) X_1 \\
& + (\alpha_1 a_{12} + \alpha_2 a_{22} + \dots + \alpha_p a_{p2}) X_2 \\
= & \vdots \\
& + (\alpha_1 a_{1p} + \alpha_2 a_{2p} + \dots + \alpha_p a_{pp}) X_p
\end{aligned}$$

بينما معادلة الانحدار بصيغتها الأصلية تكون:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p$$

لذلك ستكون معاملات الانحدار الأصلية عند العودة إليها بالتحليل حسب الآتي:

$$\begin{aligned}
\beta_0 &= \alpha_0 \\
\beta_1 &= \alpha_1 a_{11} + \alpha_2 a_{21} + \dots + \alpha_p a_{p1} \\
\beta_2 &= \alpha_1 a_{12} + \alpha_2 a_{22} + \dots + \alpha_p a_{p2} \\
&\vdots \\
\beta_p &= \alpha_1 a_{1p} + \alpha_2 a_{2p} + \dots + \alpha_p a_{pp}
\end{aligned}$$

### تحقيق التعامدية للمركبات الرئيسية:

إنطلاقاً من وضع المركبات الرئيسية بصيغة المصفوفة:

$$\begin{bmatrix} PC_1 \\ PC_2 \\ \vdots \\ PC_p \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1p} \\ a_{21} & a_{22} & \dots & a_{2p} \\ \vdots & \vdots & & \vdots \\ a_{p1} & a_{p2} & \dots & a_{pp} \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{bmatrix} = A' X$$

ونبدأ بالمركبة الرئيسية الأولى:

$$PC_1 = a_{11} X_1 + a_{12} X_2 + \dots + a_{1p} X_p = \underline{a_1}' \underline{X}$$

وتباين هذه المركبة يكون:

$$V(PC_1) = V(\underline{a_1}' \underline{X}) = a_1' S a_1$$

وهي ما نريد تعظيمه من خلال إدخال ما يسمى بمضاعف لاكرانج Lagrange Multiplier إلى عملية أخذ المشتقة لهذا التباين مع الأخذ بنظر الإعتبار القيود المصاحبة لهذه المركبة. وفي أدناه توضيحاً لذلك قبل الخوض في العملية المطلوبة للمركبة الرئيسية الأولى.

لو إفترضنا أن المقصود بتعظيمه هي الدالة  $f(x)$  ولدينا القيد  $g(x) = C$  . وبعد إدخال مضاعف لاكرنج سيتم إستحداث دالة جديدة تحتوي جميع هذه الأمور وهي:

$$h(x, \lambda) = f(x) - \lambda[g(x) - C]$$

لاحظ أننا عند إستخدام مضاعف لاكرنج لم نزيد ولم ننقص شيئاً من الدالة  $f(x)$  لأننا طرحنا منها ما قيمته صفرأ بسبب القيد  $g(x) = C$  .

ثم نأخذ المشتقة لطرفي المعادلة أعلاه بالنسبة إلى  $x$  فيكون لدينا:

$$\frac{\partial h(x, \lambda)}{\partial x} = \frac{\partial f(x)}{\partial x} - \lambda \frac{\partial g(x)}{\partial x} = 0$$

**ملاحظة:**

إذا كان المطلوب تعظيم الدالة  $f(x)$  علينا أن نحرص بأن تكون  $\lambda$  هي الأعظم والعكس بالعكس.

وبالنسبة للمركبات الرئيسية، سنقوم بعرض توضيحي لكل مركبة تحديداً:

**المركبة الرئيسية الأولى:**

تكون لدينا الدالة:

$$\phi_1 = \underline{a_1}' S \underline{a_1} - \lambda(\underline{a_1}' \underline{a_1} - 1)$$

ويكون الحل المطلوب لهذه المعادلة بأخذ مشتقتها بالنسبة إلى  $a_1$  ومساواتها للصفر هو الآتي:

$$\frac{\partial \phi_1}{\partial a_1} = 2S a_1 - 2\lambda a_1 = 0$$

$$(S - \lambda I) a_1 = 0$$

ويتم حل هذه المعادلة على أساس أن  $a_1 \neq 0$  وهذا يعطينا أن:

$$|S - \lambda I| = 0$$

وهذا يعكس لنا بأن  $\lambda$  تمثل القيمة المميزة لمصفوفة التباين  $S$  وأن  $a_1$  تمثل المتجه المميز المناظر للقيمة المميزة  $\lambda$  والذي علينا تحديده من خلال العودة إلى النتيجة أعلاه

$$(S - \lambda I)a_1 = 0$$

$$Sa_1 - \lambda Ia_1 = 0$$

$$Sa_1 = \lambda Ia_1$$

وبضرب الطرفين من اليسار بالقيمة  $a'_1$  يكون لدينا:

$$a'_1 Sa_1 = a'_1 \lambda Ia_1 = \lambda a'_1 a_1 = \lambda a'_1 a_1 = \lambda = \lambda_1$$

ويتم هنا حل هذه المعادلة بالنسبة إلى  $a_1$  والتي تناظر القيمة المميزة العظمى  $\lambda_1$ .

### ملاحظة:

لقد إستخلصنا من أعلاه بأن تباين المركبة الرئيسية الأولى هو  $a'_1 Sa_1 = \lambda_1$  وهو الأعظم ومتوازناً مع الشرط  $a'_1 a_1 = 1$  وبعدها ننتقل للمركبة الرئيسية الثانية والتي هي:

$$PC_2 = a_{21} X_1 + a_{22} X_2 + \dots + a_{2p} X_p = a'_2 X$$

وتباين هذه المركبة يكون:

$$V(PC_2) = V(a'_2 X) = a'_2 S a_2$$

وهي ما نريد تعظيمه من خلال إدخال مضاعف لاكرانج Lagrange Multiplier إلى عملية أخذ المشتقة لهذا التباين مع الأخذ بنظر الإعتبار القيود المصاحبة لهذه المركبة. وهي

$$a'_2 a_2 = 1$$

مع تثبيت شرط التعامدية مع المركبة الرئيسية الأولى وهو  $a'_2 a_1 = 0$ .

### المركبة الرئيسية الثانية:

تكون لدينا الدالة:

$$\phi_2 = a'_2 Sa_2 - \lambda(a'_2 a_2 - 1) - 2\mu(a'_2 a_1 - 0)$$

ويكون الحل المطلوب لهذه المعادلة بأخذ مشتقتها بالنسبة إلى  $a_2$  ومساواتها للصفر هو الآتي:

$$\frac{\partial \phi_2}{\partial a_2} = 2Sa_2 - 2\lambda a_2 - 2\mu a_1 = 0$$

$$(S - \lambda I)a_2 = \mu a_1$$

ويتم حل هذه المعادلة على أساس  $a_2 \neq 0$ :

$$Sa_2 - \lambda a_2 = \mu a_1$$

وبضرب الطرفين من اليسار بالقيمة  $a'_1$  يكون لدينا:

$$Sa'_1 a_2 - \lambda a'_1 a_2 = \mu a'_1 a_1$$

$$0 - 0 = \mu$$

$$\mu = 0$$

وبذلك يكون لدينا:

$$(S - \lambda I)a_2 = 0$$

وبما أن  $a_2 \neq 0$  سيكون لدينا:

$$|S - \lambda I| = 0$$

وبضرب طرفي المعادلة  $(S - \lambda I)a_2 = 0$  بالقيمة  $a'_1$  يكون لدينا:

$$a'_2 Sa_2 - a'_2 \lambda a_2 = 0$$

$$a'_2 Sa_2 - \lambda a'_2 a_2 = 0$$

$$a'_2 Sa_2 = \lambda = \lambda_2$$

ويتم هنا حل هذه المعادلة بالنسبة إلى  $a_2$  والتي تناظر القيمة المميزة العظمى الثانية  $\lambda_2$ .

وبعدنا ننتقل للمركبة الرئيسية الثالثة والتي هي:

$$PC_3 = a_{31} X_1 + a_{32} X_2 + \dots + a_{3p} X_p = a'_3 X$$

وتباين هذه المركبة يكون:

$$V(PC_3) = V(a'_3 X) = a'_3 S a_3$$

وهي ما نريد تعظيمه من خلال إدخال مضاعف لاكرانج Lagrange Multiplier إلى عملية أخذ المشتقة لهذا التباين مع الأخذ بنظر الاعتبار القيود المصاحبة لهذه المركبة. وهي:

$$a'_3 a_3 = 1 \text{ مع تثبيت شرط التعامدية مع المركبة الرئيسية الأولى وهو } a'_3 a_1 = 0 \text{ ومع}$$

$$\text{المركبة الرئيسية الثانية } a'_3 a_2 = 0.$$



### المركبة الرئيسية الثالثة:

تكون لدينا الدالة:

$$\phi_3 = a'_3 S a_3 - \lambda(a'_3 a_3 - 1) - 2\mu_1(a'_3 a_1 - 0) - 2\mu_2(a'_3 a_2 - 0)$$

ويكون الحل المطلوب لهذه المعادلة بأخذ مشتقتها بالنسبة إلى  $a_3$  ومساواتها للصفر هو الآتي:

$$\frac{\partial \phi_3}{\partial a_3} = 2S a_3 - 2\lambda a_3 - 2\mu_1 a_1 - 2\mu_2 a_2 = 0$$

$$(S - \lambda I) a_3 - \mu_1 a_1 - \mu_2 a_2 = 0$$

ويتم حل هذه المعادلة على أساس  $a_3 \neq 0$ .

وبضرب الطرفين من اليسار بالقيمة  $a'_1$  يكون لدينا:

$$a'_1(S - \lambda I) a_3 - a'_1 \mu_1 a_1 - a'_1 \mu_2 a_2 = 0$$

$$(S - \lambda I) a'_1 a_3 - \mu_1 a'_1 a_1 - \mu_2 a'_1 a_2 = 0$$

وهذا يعطينا  $\mu_1 = 0$

ولو ضربنا نفس الطرفين من اليسار بالقيمة  $a'_2$  يكون لدينا:

$$a'_2(S - \lambda I) a_3 - a'_2 \mu_1 a_1 - a'_2 \mu_2 a_2 = 0$$

$$(S - \lambda I) a'_2 a_3 - \mu_1 a'_2 a_1 - \mu_2 a'_2 a_2 = 0$$

وهذا يعطينا  $\mu_2 = 0$  وبذلك يكون لدينا:

$$(S - \lambda I) a_3 = 0$$

وبما أن  $a_3 \neq 0$  سيكون لدينا:

$$|S - \lambda I| = 0$$

وبضرب طرفي المعادلة  $(S - \lambda I) a_3 = 0$  من اليسار بالقيمة  $a'_3$  يكون لدينا:

$$a'_3 S a_3 - a'_3 \lambda a_3 = 0$$

$$a'_3 S a_3 - \lambda a'_3 a_3 = 0$$

$$a'_3 S a_3 = \lambda = \lambda_3$$

وهكذا لبقية المركبات الرئيسية حيث نستمر بنفس النهج التصاعدي لمضاعف لاكرنج.

## مثال 1 (حالة متغيرين)

البيانات التالية تمثل الدخل الشهري  $X_1$  بعملة معينة و  $X_2$  تمثل الخدمة بالسنوات و  $Y$  تمثل الزيادة التشجيعية لعينة من خمسة عاملين. ومقياس جميع المتغيرات هنا هو كمي (نسبي أو فنوي).

العامل	Y	$X_1$	$X_2$
1	8	102	4
2	9	104	5
3	7	101	7
4	9	93	1
5	10	100	3

سنستخدم هنا مصفوفة الارتباط وهي:

$$R = \begin{bmatrix} 1 & 0.7483 \\ 0.7483 & 1 \end{bmatrix}$$

$$|R - \lambda I| = \begin{vmatrix} 1 - \lambda & 0.7483 \\ 0.7483 & 1 - \lambda \end{vmatrix} = 0$$

$$(1 - \lambda)^2 - 0.56 = 0$$

$$\lambda_1 = 1.7483$$

$$\lambda_2 = 0.2517$$

وعلينا تحديد قيم عناصر المتجه  $\underline{a_1}$  من خلال الآتي:

$$(R - \lambda I) \underline{a_1} = 0$$

$$(R - 1.7483 I) \begin{pmatrix} a_1 \\ a_2 \end{pmatrix} = 0$$

$$\begin{pmatrix} 1 - 1.7483 & 0.7483 \\ 0.7483 & 1 - 1.7483 \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \end{pmatrix} = 0$$

$$\begin{pmatrix} -0.7483 & 0.7483 \\ 0.7483 & -0.7483 \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \end{pmatrix} = 0$$

ومن الواضح هنا أن  $a_1 = a_2$  ولذلك لو إفترضنا القيمة (1) لأحدهما ستكون الأخرى مساوية لها وبنفس القيمة. أي أنه يصبح لدينا:

$$\begin{bmatrix} a_1 \\ a_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

وهذه القيم الأولية يمكن إستخدامها لتحديد قيم عناصر المتجه المميز الذي يتسم بشرط مجموع مربع قيم العناصر مساوياً إلى 1. ولتحقيق ذلك، نتبع التحويل التالي:

$$\underline{a}'_1 = \begin{bmatrix} a_{11} \\ a_{12} \end{bmatrix} = \begin{bmatrix} \frac{a_1}{\sqrt{\sum a_i^2}} \\ \frac{a_2}{\sqrt{\sum a_i^2}} \end{bmatrix} = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix}$$

وهي قيم عناصر المتجه القياسي المعتمد الأول.

والآن علينا تحديد قيم عناصر المتجه القياسي الثاني  $\underline{a}_2$  من خلال إتباع نفس الإسلوب ولكن بإستخدام القيمة المميزة الثانية وحسب الآتي:

$$(R - \lambda I) \underline{a}_2 = 0$$

$$(R - 1.7483I) \begin{pmatrix} a_1 \\ a_2 \end{pmatrix} = 0$$

$$\begin{pmatrix} 1-0.2517 & 0.7483 \\ 0.7483 & 1-0.2517 \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \end{pmatrix} = 0$$

$$\begin{pmatrix} 0.7483 & 0.7483 \\ 0.7483 & 0.7483 \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \end{pmatrix} = 0$$

ومن الواضح هنا أن  $a_1 = -a_2$  ولذلك لو افترضنا القيمة (1) لأحدهما ستكون الأخرى

هي:

(1) أي أنه يصبح لدينا:

$$\begin{bmatrix} a_1 \\ a_2 \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

وهذه القيم الأولية يمكن إستخدامها لتحديد قيم عناصر المتجه المميز الذي يتسم بشرطي مجموع مربع قيم العناصر مساوياً إلى 1 وحاصل ضرب المتجهين يساوي 0.0.

ولتحقيق ذلك، نتبع التحويل التالي:

$$\underline{a}'_1 = \begin{bmatrix} a_{11} \\ a_{12} \end{bmatrix} = \begin{bmatrix} \frac{a_1}{\sqrt{\sum a_i^2}} \\ \frac{a_2}{\sqrt{\sum a_i^2}} \end{bmatrix} = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \end{bmatrix}$$

وهي قيم عناصر المتجه القياسي المعتمد الثاني.

وحيث أن  $\frac{1}{\sqrt{2}} = 0.707$  ، فإن المركبتين الرئيسيتين ستكونا بالشكل التالي:

$$PC_1 = 0.707X_1 + 0.707X_2$$

$$PC_2 = 0.707X_1 - 0.707X_2$$

ومن المناسب هنا تجميع جميع القيم المرتبطة بالمركبات الرئيسية التي تم تحديدها وهو ما يتعلق بالمتجهات والقيم المميزة مع الأهمية النسبية لكل مركبة. كل ذلك سيكون ضمن الجدول التالي:

	PC <sub>1</sub>	PC <sub>2</sub>
X <sub>1</sub>	0.707	0.707
X <sub>2</sub>	0.707	- 0.707
$\lambda$	1.7483	0.2517
الأهمية النسبية	87.4%	12.6%

وهنا لو قررنا الإبقاء على المركبات الرئيسية المهمة فقط، فإننا سنكتفي بالمركبة الأولى لكون القيمة المميزة تزيد عن الواحد. وهذا الإجراء هو ما نسميه بتقليص الإتجاهات Reduction of Dimensionality في التحليل.

ولغرض إعطاء توضيح أكثر لهذا الموضوع، سنأخذ المثال التالي لبيانات حقيقية:

## مثال 2:

في دراسة G. R. Bryce في جامعة Brigham Young حول احتمالية وجود علاقة ما بين تصميم خوذة الرأس للاعب كرة القدم الأمريكية وجروح الرقبة، تم تسجيل 6 قياسات (متغيرات مستقلة) لثلاثة مجاميع من اللاعبين وبواقع 30 لاعب لكل مجموعة وهي كما يلي<sup>(3)</sup>:

المجموعة (1) : تمثل لاعبي المرحلة الثانوية لكرة القدم

المجموعة (2) : تمثل لاعبي المرحلة الجامعية لكرة القدم

المجموعة (3) : تمثل من هم بالعمر الجامعي وليسوا من لاعبي كرة القدم

أما القياسات (المتغيرات المستقلة) الست لاعبين فهي كما يلي:

$$X_1 = \text{عرض الرأس عند أكبر عمق}$$

$$X_2 = \text{محيط الرأس}$$

$$X_3 = \text{المسافة ما بين مقدمة الرأس ومؤخرته عند مستوى العين}$$

$$X_4 = \text{المسافة ما بين العين وقمة الرأس}$$

$$X_5 = \text{المسافة ما بين الأذن وقمة الرأس}$$

$$X_6 = \text{عرض الفك}$$

وجميع هذه المتغيرات بمقياس كمي نسبي وبنفس الوحدات.

**تنويه:** في هذا المثال تم إعتداد المجموعتين الثانية والثالثة فقط لوجود التجانس العمري بينهما.

ونبدأ العمل بطبيعة الحال بمصفوفة التباين والتباين المشترك - Variance - Covariance Matrix للمتغيرات المستقلة الست وهي:

$$S = \begin{bmatrix} 0.320 & 0.602 & 0.149 & 0.044 & 0.107 & 0.209 \\ & 2.629 & 0.801 & 0.666 & 0.103 & 0.377 \\ & & 0.458 & 0.111 & -0.013 & 0.120 \\ & & & 1.474 & 0.252 & -0.054 \\ & & & & 0.488 & -0.036 \\ & & & & & 0.324 \end{bmatrix}$$

ومجموع العناصر القطرية لهذه المجموعة يمثل مجموع التباين والذي يمثل بدوره مجموع القيم المميزة لهذه المصفوفة. أي أن:

$$\sum_{j=1}^6 S_{jj} = \sum_{i=1}^6 \lambda_{ii} = 0.320 + 2.629 + 0.458 + 1.474 + 0.488 + 0.324 = 5.743$$

ومن خلال التحليل تم تحديد القيم المميزة Eigen values الست مع المتجهات المميزة Eigen vectors للقيم المعتمدة وهي الأولى والثانية فقط كونهما توضحان نسبة 81.8% من التباين وكلاهما أكبر من الواحد، وكانت كما يلي:

Eigen value القيمة المميزة	Proportion of Variance نسبة التباين	Cumulative proportion النسبة المتجمعة	Eigen Vectors المتجهات المميزة		
			المتغير	a <sub>1</sub>	a <sub>2</sub>
3.323	0.579	0.579	X <sub>1</sub>	0.207	-0.142
1.374	0.239	0.818	X <sub>2</sub>	0.873	-0.219
0.476	0.083	0.901	X <sub>3</sub>	0.261	-0.231
0.325	0.057	0.957	X <sub>4</sub>	0.326	0.891
0.157	0.027	0.985	X <sub>5</sub>	0.066	0.222
0.088	0.015	1.000	X <sub>6</sub>	0.128	-0.187

ووفقاً لذلك فإننا سنكتفي بالمركبة الرئيسية الأولى PC<sub>1</sub> والثانية PC<sub>2</sub> وتكون بالصيغة التالية:

$$PC_1 = 0.207X_1 + 0.873X_2 + 0.261X_3 + 0.362X_4 + 0.066X_5 + 0.128X_6$$

$$PC_2 = -0.142X_1 - 0.219X_2 - 0.231X_3 + 0.891X_4 + 0.222X_5 - 0.187X_6$$

ومن الملاحظ هنا أن X<sub>2</sub> (محيط الرأس) له تأثير واضح ضمن المركبة الرئيسية الأولى بمعامل (0.873) ومثل ذلك بالنسبة إلى X<sub>4</sub> (المسافة ما بين العين وقمة الرأس) ضمن المركبة الرئيسية الثانية بمعامل (0.891). وهذا متوقع حدوثه لكون هذين المتغيرين لهما أكبر تباين (2.629 بالنسبة إلى X<sub>2</sub> و 1.474 بالنسبة إلى X<sub>4</sub>) بالمقارنة مع تباينات بقية المتغيرات.

وجدير بالذكر هنا أنه لو كانت هذه المتغيرات بتباينات متقاربة، لكننا قد لاحظنا تقارب كبير في معاملاتهما ضمن أي من المركبتين. ومن جهة أخرى، لو كانت تباينات هذين المتغيرين X<sub>2</sub> و X<sub>4</sub> كبيرة جداً نسبياً، لكننا نلاحظ أن المركبة الرئيسية الأولى PC<sub>1</sub> تساوي إلى حدٍ ما المتغير X<sub>2</sub>، كما أن المركبة الرئيسية الأولى PC<sub>2</sub> تساوي إلى حدٍ ما المتغير X<sub>4</sub>.

# نموذج الانحدار اللوجستي

## Logistic Regression Model (LRM)

### مقدمة

يستخدم الانحدار اللوجستي في الغالب لنمذجة احتمال عائدة وحدة تجرية لمجموعة معينة استناداً إلى معلومة مأخوذة من تلك الوحدة. مثل هذه النماذج يمكن استخدامها لأغراض التمييز. وفي حالة البطاقة الائتمانية، يمكننا نمذجة احتمال كون شخص ما ذو خصائص ديموغرافية معينة يقع ضمن مجموعة الجيدين من ناحية خطورة الائتمان. وبعد تطوير هذا النموذج، بالإمكان استخدامه للتنبؤ بالمجموعة التي ينتمي إليها شخصاً جديداً وفقاً لسماته الديموغرافية المحددة. والشخص الذي يعطي عنه النموذج احتمالاً يزيد عن 0.5 يحدد بأنه ينتمي إلى مجموعة الجيدين بالنسبة لخطورة الائتمان.

وطالما أنه نموذج انحدار، فهذا يعني أن هناك متغيرات توضيحية تقابل متغير معتمد لكن هذا المتغير المعتمد قد يكون بقياس رقمي عادي Numerical أو مقياس فئوي إسمي Nominal. ولأننا نتحدث عن احتمال، فإن هذا هو احتمال انتساب المتغير المعتمد لفئة معينة. ولذلك، فإنه في حالة كون القياس رقمي فإننا نلجأ إلى تقسيم القيم بحدود مناسبة لغرض تحديد الفئات بموجبها.

إن نماذج الانحدار اللوجستي تعتبر أسلوباً جيداً للوقوف على كيفية تأثير عددٍ من العوامل على ظهور مشاهدة ثنائية القياس Binary كمتغير معتمد (إستجابة Response). ونقصد بالمتغير الثنائي Binary هو أن هذا المتغير يأخذ قيمتين محتملتين. والأمثلة على ذلك الوفيات (حي/متوفي) والحالة المرضية (سليم/ مريض) والإجابة (نعم/ كلا) والجنس (ذكر/ أنثى) وهكذا.

وفي بعض الأحيان يمكن أن يحصل لدينا متغير تابع (معتمد) مستمر القياس وليس ثنائياً مثل مدة المكوث في المستشفى جراء الولادة. وفي هذه الحالة فقد يتم تقسيم القيم إلى أصغر أو يساوي 48 ساعة مقابل أكثر من 48 ساعة ليصبح القياس عندئذٍ ثنائياً.

وبالنسبة للملاحظات الثنائية، فقد نجد أسلوب الترميز الرقمي مناسباً باستخدام القيمتين (صفر/1) لكل مشاهدة ويفضل استخدام العدد (1) ليعكس وجود الظرف (الحالة) والعدد (صفر) لغيابها. وكمثال على ذلك، ينسب الرقم (1) للمريض والرقم (صفر) لغير المريض (السليم).

ولتوضيح أهمية استخدام نموذج الانحدار اللوجستي، افترض أنه لدينا مشاهدة ثنائية  $Y$  (المرض : نعم / كلا ) بصفة الإستجابة ومتغير مؤثر آخر  $X$  (التعرض : نعم / كلا). وهنا يمكن تسمية  $X$  بأنه عامل الخطورة وتسمية  $Y$  بأنه عامل الإستجابة تجاه هذه الخطورة. وفي مثل هذه الحالة، يمكننا استخدام جدول تقاطعي ( $2 \times 2$ ) لتقدير الخطورة النسبية Relative Risk وعلى النحو التالي:

		$Y$	
		Yes (1)	No (0)
$X$	Yes(1)	a	b
	No(0)	c	d

وحيث أن  $d, c, b, a$  تمثل تكرارات لهذه التقاطعات وأن حجم العينة  $n$  تكون:

$$n = a + b + c + d$$

$$RR = \frac{a/(a+b)}{c/(c+d)}$$

في حالة الدراسات الفوجية Cohort أو ما نسميه الدراسات التتبعية Prospective.

كذلك يمكن استخدامها لتقدير نسبة الأرجحية Odds Ratio (OR) في حالتها التتبعية أو التقاطعية Case – Control. وتحسب نسبة الأرجحية هذه بأنها

$$OR = \frac{(a)(d)}{(b)(c)}$$

والآن لنفترض أن  $Y$  متغير ثنائي بينما  $X$  متغير مستمر. وفي هذه الحالة، لا يمكن استخدام الجدول التقاطعي ( $2 \times 2$ ).

وفي مثل هذه المسألة، نستخدم الانحدار اللوجستي البسيط في حالة متغير مؤثر (مستقل) واحد  $X$  أو الانحدار اللوجستي المتعدد في حالة وجود أكثر من متغير مؤثر (مستقل) واحد مثل  $X_1, X_2, \dots, X_p$  كمتغيرات مستقلة.

والآن دعنا نتناول نموذج الانحدار اللوجستي البسيط والذي يتضمن متغير مستقل واحد  $X$ .



ونموذج الإنحدار اللوجستي يكون مناسباً لدالة الإستجابة  $Y$  مقابل  $X$  (أو متعدد  $X_1, X_2, \dots$ ) مع منحنى يعرف بشكل  $S$  أي  $S$  - Shaped Curve ويعرف بالدالة اللوجستية Logistic Function ، ويعبر عنها بالآتي:

$$E(Y|X) = \exp(\beta_0 + \beta_1 X) / [1 + \exp(\beta_0 + \beta_1 X)]$$

حيث أن  $E(Y|X)$  هي القيمة المتوقعة إلى  $Y$  عند وجود  $X$ . وأن "exp" تعبر عن الأساس إلى اللوغاريتم الطبيعي ( $\log_e = \ln$ ) وأن ( $e = \exp = 2.71828$ ). كما أن  $\beta_0$  و  $\beta_1$  عبارة عن معامل الإنحدار وهما معلمتان يتم العمل على تقديرهما في النموذج.

ولأن معدل دالة الإستجابة  $E(Y|X) = p$  عندما يكون  $Y$  متغيراً ثنائياً Binary حينما نستخدم القيمة (صفر) للتعبير عن الفشل والقيمة (واحد) للتعبير عن النجاح والذي تكون  $p$  عبارة عن احتمال ظهوره. ولذلك فإننا نستخدم في العادة صيغة النموذج التالية :

$$P = \exp(\beta_0 + \beta_1 X) / (1 + \exp(\beta_0 + \beta_1 X))$$

ولكي نحصل على نموذج خطي بالنسبة إلى المعلمات  $\beta_0$  و  $\beta_1$  فإننا نحتاج لإستخدام التحويل اللوجستي هنا ويسمى Logistic transformation حيث أن:

$$\text{Logit} = \log_e \left( \frac{P}{1-P} \right) = \ln \left( \frac{P}{1-P} \right)$$

وعلى سبيل المثال، لو كان احتمال ظهور النجاح لمشاهدة ما هو  $p = 0.1$  فإن:

$$\text{Logit} = \ln \left( \frac{0.1}{0.9} \right) = -2.1972$$

وفيما يلي توضيح لكيفية الحصول على نموذج خطي في  $\beta_0$  و  $\beta_1$ . حيث لدينا:

$$P = \exp(\beta_0 + \beta_1 X) / (1 + \exp(\beta_0 + \beta_1 X))$$

$$\begin{aligned} 1 - p &= 1 - \exp(\beta_0 + \beta_1 X) / (1 + \exp(\beta_0 + \beta_1 X)) \\ &= 1 / (1 + \exp(\beta_0 + \beta_1 X)) \end{aligned}$$

وبالتالي فإن:

$$\frac{P}{1-P} = \exp(\beta_0 + \beta_1 X)$$

$$\ln \left( \frac{P}{1-P} \right) = \ln[\exp(\beta_0 + \beta_1 X)] = \beta_0 + \beta_1 X \quad \dots\dots\dots(1)$$

وجدير بالذكر هنا أننا نطلق تعبير "نسبة الأرجحية" Odd's Ratio على المقدار  $\frac{p}{1-p}$

$$OR = \frac{p}{1-p} \quad \text{أي أن:}$$

وجدير بالملاحظة أن  $\text{logit} = \log(OR)$  هو بمثابة المتغير المعتمد في النموذج الخطي كما في المعادلة (1) أعلاه.

#### ملاحظة:

أيما يذكر اللوغاريتم  $\log$  فإننا نعني اللوغاريتم الطبيعي  $\ln$ .

ولو افترضنا أن المتغير  $X$  ثنائي القياس Binary بحيث يكون  $X=0$  عند تواجده و  $X=1$  عند عدم تواجده لتوصلنا إلى حقيقة كون معدل التغير في نسبة الأرجحية  $OR$  تستند على تقدير المعلمة  $\beta_1$  أي  $b_1$ .

وفيما يلي توضيح لكيفية حدوث ذلك لغرض التثبيت لهذه الحقيقة. حيث أن:

$$\begin{aligned} \text{Log}(OR) &= \log[(\text{odds} | X=1)/(\text{odds} | X=0)] \\ &= \log(\text{odds} | X=1) - \log(\text{odds} | X=0) \\ &= \text{logit}(X=1) - \text{logit}(X=0) \\ &= (\beta_0 + \beta_1 X) - (\beta_0) \quad , \quad X=1 \\ &= \beta_1 \end{aligned}$$

وبذلك فإن المماس  $\beta_1$  هو عبارة عن  $\log(OR)$  ولذلك فإن:

$$OR = \exp(\beta_1)$$

في مثل هذه الحالات (أي حالة كون  $X$  متغير ثنائي القياس Binary).

والآن يتم تركيزنا على أهمية نموذج الإنحدار اللوجستي من أنه يفحص العلاقة ما بين متغير مستقل واحد أو أكثر و  $\log(OR)$  المتغير المعتمد (الإستجابة) ثنائي القياس.

وإذا ما نظرنا إلى  $\log(OR)$  يترأى لنا أننا نتبع طريقة معقدة لتفسير البيانات، ولكن هذا يؤدي إلى أسهل أسلوب لتفسير البيانات والتي تتماشى مع قواعد الاحتمالات.

ولتوضيح هذا الجانب، دعنا نفترض البيانات التالية:

GA	28	29	30	31	32	33	34
Prob(BF)	0.60	0.62	0.64	0.66	0.68	0.70	0.72

**ملاحظة:**

GA يمثل فترة الحمل بالأسابيع وهو المتغير المستقل وذو قياس نسبي كونه كمي متصل

BF يمثل رضاعة الطفل من صدر أمه بعد مغادرته المستشفى وهو متغير معتمد.

وثنائي القياس ( yes =1 , No = 0 ) ومقياسه إسمي.

وعند استخدامنا لنموذج خطي للربط ما بين المتغيرين، فقد نقول أن:

$$\text{Prob}(\text{BF}) = \beta_0 + \beta_1 (\text{GA})$$

أي أنه من الواضح كون Prob(BF) (وهي إحتمال رضاعة الطفل من صدر أمه) تزداد وفق دالة خطية بالنسبة إلى GA (مدة الحمل بالأسابيع). وعند تحليل هذه البيانات وفقاً لذلك، نجد أن:

$$\text{Prob}(\text{BF}) = 0.04 + 0.02 (\text{GA})$$

وعلى هذا الأساس، فإن رضيعاً مع فترة حملة 30 إسبوعاً سيرضع من صدر أمه بعد مغادرتها المستشفى بإحتمال 0.64 وفقاً لهذا النموذج التجميعي.

**ملاحظة مهمة (لنموذج التجميعي) Additive Model**

نلاحظ هنا أننا إستخدمنا نموذج تجميعي في الإحتمالات. وهذا قد يقودنا إلى نتائج غير منطقية في بعض الأحيان وغير محسوبة خاصة إذا ما إقترب الإحتمال إلى % 0.0 أو إلى 100% لأننا نتوقع أن نجتاز هذين الحدين فيكون لدينا إحتمال أكثر من 100% أو أقل من الصفر (إحتمال سالب) وكلاهما غير منطقي بالنسبة للإحتمالات. دعنا الآن نجري بعض التغيير في نموذجنا الخطي للآتي:

$$\text{Prob}(\text{BF}) = 0.04 + 0.03 (\text{GA})$$

وهذا سوف يعطينا النتائج التقديرية التالية:

GA	28	29	30	31	32	33	34
Prob(BF)	0.88	0.91	0.94	0.97	1.0	1.03	1.06

وبطبيعة الحال، سنجد صعوبة في إعطاء تفسير عما يعنيه الإحتمال 1.06 وهذا هو السبب الرئيسي الذي يدفعنا لتجنب استخدام النموذج التجميعي.

وبشكل عام، فإنه ينصح بمحاولة تجنب النموذج التجميعي هذا ما لم يكن هنالك سبب قوي (قناعة) لتوقع كون جميع قيم الإحتمالات المتوقعة ستكون ضمن المدى ( 80% - 20%).

ونعني بتجنب النموذج التجميعي هو تبني النموذج الضربي للإحتمالات وهو التالي.

### النموذج الضربي للإحتمالات A Multiplicative Model

قد يكون من الأجدر اعتماد النموذج الضربي في حالة الإحتمالات مع كونه قد يعاني أيضاً من نفس المشاكل كما هي الحال في النموذج التجميعي ولكن بشكل أقل. وهنا تعتمد عملية الضرب بدلاً من عملية الجمع عند تغيير قيمة الإحتمال. وفيما يلي مثالاً على ذلك.

لنفترض النتيجة التقديرية التالية:

GA	28	29	30	31	32	33	34
Prob(BF)	0.01%	0.03%	0.09%	0.27%	0.81%	2.43%	7.29%

وفي هذا المثال، فإن كل إسبوع إضافي في مدة الحمل يؤدي إلى مضاعفة إحتمال الرضاعة الطبيعية Prob(BF) ثلاث مرات. كما نلاحظ هنا أن النموذج الضربي لا يؤدي إلى إحتمال أقل من الصفر ولكنه قد يؤدي إلى إحتمال أكثر من 1.0. ولغرض تطبيقه علينا أن تكون لدينا القناعة القوية بأن توقع جميع الإحتمالات ستكون بقيم صغيرة، ولتكن أقل من 0.20.

### العلاقة بين الأرجحية Odds والإحتمال Prob.

لنفترض حالة الربح والخسارة لفريق كرة قدم. فإذا كانت الأرجحية 3 إلى 1 بجانب فريقك فهذا يعني أنك تتوقع الربح (الفوز) أن يحصل 3 مرات بقدر عدد الخسارة. سنتعامل مع الحالة هذه إضافة إلى الأرجحية 4 إلى 1 ونرى كيف يمكننا تحويل الأرجحية إلى إحتمال وبالعكس. ففي الحالة الأولى يكون الربح ثلاثة مرات مع كل أربعة أشواط أي أن إحتمال الفوز هو 0.75. أما إذا كان إحتمال الفوز هو 0.20 فهذا يعني أنك تتوقع مرة واحدة للفوز مقابل أربع مرات خسارة مع كل خمسة أشواط لعب.

وهنا نستطيع وضع الصيغة التالية لهذه العلاقة ما بين نسبة الأرجحية OR والإحتمال Prob

$$OR = Prob/(1 - Prob) = P/(1 - P)$$

$$Prob = OR / (1 + OR)$$

والآن دعنا نرى تطبيق نموذج الإنحدار اللوجستي البسيط بمتغير مستقل واحد ومن خلال البيانات التالية:

البيانات الأولية		القيم المتوقعة		
GA = X	Prob (BF) = Y	Log (Odds)	Odds(BF)	Prob(BF)
28	2/6= 0.333	- 0.57	0.57	0.362
29	2/5= 0.400	0.01	1.01	0.503
30	7/9 =0.778	0.59	1.80	0.643
31	7/9 =0.778	1.16	3.20	0.762
32	16/20=0.80	1.74	5.70	0.851
33	14/15=0.933	2.32	10.15	0.910

وفي حالة تطبيق نموذج الإنحدار اللوجستي البسيط يكون لدينا:

$$P = Y = \exp(\beta_0 + \beta_1 X) / [1 + \exp(\beta_0 + \beta_1 X)]$$

$$\text{Logit}(P) = \log [P/(1 - P)] = \beta_0 + \beta_1 X$$

وبالتالي سنحصل على النموذج المتوقع:

$$\text{Logit}(P) = \log [P/(1 - P)] = - 16.72 + 0.577 X$$

وفي حالة كون  $X = GA = 30$  ، فإن هذا النموذج يعطينا:

$$\text{Logit}(P) = \log [P/(1 - P)] = \log(\text{OR})$$

$$= - 16.72 + 0.577(30) = 0.59$$

ولتحويل هذه القيمة إلى Odds يكون لدينا:

$$\text{OR} = \exp(0.59) = 1.80$$

وأخيراً تحويل ذلك إلى الإحتمال (رجوعاً) يكون لدينا:

$$\text{Prob.} = 1.80 / (1 + 1.80) = 0.643$$

وهذه القيمة المتوقعة للإحتمال تعتبر قريبة بشكلٍ معقول من الإحتمال الحقيقي وهو (0.778).

ومن المفيد أيضاًلقاء نظرة على نسبة الأرجحية المتوقعة (BF) OR حيث نلاحظ من الجدول بأن نسبة أي OR لصفين متتاليين هي 1.78.

$$3.20/1.80 = 1.78 \quad \text{ومثالاً على ذلك نجد أن}$$

$$5.70/3.20 = 1.78 \quad \text{وكذلك}$$

وهذه النتيجة ليست من قبيل المصادفة وإنما هي واقع قيمة المقدار

$$\exp(\beta_1) = \exp(0.577) = 1.78$$

وهذه هي الصفة العامة لنموذج الإنحدار اللوجستي. فقيمة الميل في نموذج الإنحدار اللوجستي يمثل  $\log(OR)$  أو لوغارتيم نسبة الأرجحية والتي تمثل (الزيادة/النقصان) في الخطورة عندما تزداد قيمة المتغير المستقل وحدة واحدة.

### نموذج الإنحدار اللوجستي المتعدد Multiple Logistic Regression Model

في هذا النموذج يكون لدينا متغيرين مستقلين أو أكثر  $X_1, X_2, \dots, X_k$ . وبالتالي، يكون النموذج وفقاً للمعادلة:

$$Y = P = \exp(\beta_0 + \sum \beta_i X_i) / [1 + \exp(\beta_0 + \sum \beta_i X_i)]$$

$$\text{Logit}(P) = \log [P/(1 - P)] = \beta_0 + \sum \beta_i X_i$$

ومن المزايا الحيدة لهذا النموذج أنه يسمح بوجود المتغيرات المستقلة (التوضيحية) بقياسات مختلفة بالنسبة إلى  $X_1, X_2, \dots, X_k$ . فقد يكون لدينا متغير أو أكثر فنوي القياس (سواءً إسمي أو رتبي أو فترتي). وسنرى ذلك من خلال المثال التالي والذي يتضمن دراسة لتحديد مدى تأثير عوامل مختلفة عديدة ومتنوعة القياس تجاه حالة البدانة لدى المرأة الأردنية.

#### المثال :

في دراسة لتحديد مدى تأثير عدد من العوامل تجاه حالة البدانة لدى المرأة الأردنية<sup>(4)</sup>، تم إعتبار مؤشر البدانة والمتمثل بالمقياس BMI بمثابة المتغير المعتمد Y حيث:

BMI يرمز إلى Body Mass Index ويتم حسابه وفقاً للآتي

$$BMI = \text{Body weight (Kg)} / [\text{Body High (meter)}]^2 = \text{Kg/m}^2$$

وتعتبر المرأة بدينة في حالة كون المقياس  $BMI \geq 30 \text{ Kg/m}^2$  مع مجموعة العوامل المؤثرة X والتي هي بمثابة المتغيرات المستقلة وتشمل:

AGE : ويمثل عمر المرأة بالسنوات ومقياسه كمي نسبي.

ED : ويمثل مستوى التعليم (أساسي أو أقل / ثانوي ED1 / عالي ED2) ومقياسه كمي فنوي ثلاثي.

PAR : ويمثل حجم الولادات ( حد أعلى 1 / PAR1 (2-3) / PAR2 (≥ 4) ) ومقياسه كمي فنوي ثلاثي.

REG : ويمثل المنطقة الجغرافية (المركز / شمال REG1 / جنوب REG2) ومقياسه إسمي ثلاثي.

RES : ويمثل تصنيف منطقة السكن (حضري/ ريفي RES1) ومقياسه إسمي ثنائي.

BAD : ويمثل سكن البادية (خارج البادية / البادية BAD1) ومقياسه إسمي ثنائي.

WEL : ويمثل الوضع المادي ( ليست فقيرة / فقيرة WEL1) ومقياسه إسمي ثنائي.

CONT : ويمثل إستخدام طرق موانع الحمل (كلا/ طرق تقليدية CONT1 / طرق حديثة CONT2) ومقياسه إسمي ثلاثي.

SMOK : ويمثل التدخين (كلا / نعم SMOK1) ومقياسه إسمي ثنائي.

WORK : ويمثل صفة العمل (لا تعمل / تعمل WORK1) ومقياسه إسمي ثنائي.

وتم تطبيق نموذج الإنحدار اللوجستي المتعدد التالي:

$$\text{Logit } Y = \log\left(\frac{P}{1-P}\right)$$

$$= \beta_0 + \beta_1 AGE + \beta_2 ED1 + \beta_3 ED2 + \beta_4 PAR1 + \beta_5 PAR2 + \beta_6 REG1 + \beta_7 REG2 + \beta_8 RES1 + \beta_9 BAD1 + \beta_{10} WEL1 + \beta_{11} CONT1 + \beta_{12} CONT2 + \beta_{13} SMOK1 + \beta_{14} WORK1$$

$$P = \text{Prob}(y=1|X) = \text{Prob}(BMI \geq 30 \text{ Kg/m}^2)$$

و X هنا تمثل مجموعة المتغيرات التوضيحية (المستقلة).

**ملاحظة:**

جدير بالتنويه أن أي متغير ظهر في النموذج أعلاه (عدا العمر AGE لكونه كمي) يأخذ القيمة (1) عند تحقق وجوده والقيمة (0) عدا ذلك في حال كونه ثنائي. أما في حالة كونه ثلاثي (ولنأخذ متغير مستوى التعليم ED على سبيل المثال) فإن القيم التي تحتسب لمجاميعه تكون بالشكل التالي:

ED1 يأخذ القيمة (1) في حالة المستوى ثانوي والقيمة (0) عدا ذلك

ED2 يأخذ القيمة (1) في حالة المستوى جامعي والقيمة (0) عدا ذلك

وفي حالة كون كليهما بالقيمة (0) فيعني أن مستوى التعليم بدون أو ابتدائي

وهكذا بالنسبة لبقية المتغيرات الثلاثية.

وجدير بالذكر أن الدراسة موضوع البحث اعتمدت بيانات من مسوحات ديموغرافية وصحية لثلاثة سنوات 2002 و 2009 و 2012 للنساء بعمر (15 – 49) سنة ولحالتين منهن (جميع النساء المشمولات بالمسح / النساء المتزوجات فقط). كما تم تطبيق الدراسة لكل مسح على حدة بالإضافة إلى دمج المسوحات الثلاثة.

ولغرض التوضيح فيما نتناوله من تحليل لنتائج النموذج اللوجستي، سنكتفي بعرض نتائج حالة مسح عام 2009 وللنساء المتزوجات بعمر (15 – 49) سنة نظراً لشموليته لجوانب عديدة من التحليل لنتائج هذا النموذج. والنتائج كانت حسبما مبين في الجدول التالي:

المتغير	مسح عام 2009 للمتزوجات بعمر (15 – 49)	
	$\beta$	OR
AGE	0.077 ***	1.08
ED1	-0.041	0.96
ED2	-0.562 ***	0.57
PAR1	0.166	1.18
PAR2	0.56 **	1.75
REG1	0.336 **	1.40
REG2	0.531 ***	1.70
RES1	0.068	1.07
BAD1	0.077	1.08
WEL1	0.00	1.00
CONT1	-0.329 *	0.72
CONT2	-0.117	0.98
SMOK1	-0.528 **	0.59
WORK1	0.02	1.02

\*P-value < 0.05 \*\* P-value < 0.01 \*\*\* P-value < 0.001

## تحليل النتائج

سنتناول أولاً المتغيرات النوعية:

### 1) بالنسبة للمتغيرات غير المعنوية

في حالة المتغير الثلاثي ED "مستوى التعليم" (بدون أو ابتدائي/ ثانوية ED1 /جامعية ED2) نجد أن ED1 والذي يمثل المستوى الثانوي غير معنوي، وهذه النتيجة تشير إلى أنه عند مقارنة



امرأتين من حيث مقياس البدانة إحداهما بمستوى تعليم ابتدائي أو دون ذلك والأخرى بمستوى تعليم ثانوي، مع تثبيت بقية المتغيرات لهما فإننا نتوقع مقياس بدانة لهما بنمط واحد أي كلاهما فوق 30 أو كلاهما دون ذلك. بمعنى آخر، فإن المرأة بمستوى تعليم ثانوي لا تختلف عن أخرى بمستوى تعليم ابتدائي تجاه مقياس البدانة عندما تكونا متساويتان بالنسبة للمتغيرات الأخرى.

وبنفس الوقت، لو نظرنا لحالة المتغير الثنائي RES "تصنيف منطقة السكن" (حضري/ريفي) RES1 نجد أن RES1 والذي يمثل الريف غير معنوي، وهذه النتيجة تشير أيضاً إلى أنه عند مقارنة امرأتين من حيث مقياس البدانة إحداهما من سكنة منطقة حضرية والأخرى من سكنة منطقة ريفية، مع تثبيت بقية المتغيرات لهما فإننا نتوقع مقياس بدانة لهما بنمط واحد أي كلاهما فوق 30 أو كلاهما دون ذلك. بمعنى آخر، فإن المرأة من سكنة المنطقة الحضرية لا تختلف عن أخرى من سكنة المنطقة الريفية تجاه مقياس البدانة عندما تكونا متساويتان بالنسبة للمتغيرات الأخرى. أي أنهما يتعرضان للبدانة بإحتمال متقارب جداً.

وهكذا لبقية المتغيرات غير المعنوية.

## 2) بالنسبة للمتغيرات المعنوية بإشارة سالبة إلى تقدير $\beta$

ولنأخذ متغير المستوى التربوي ED على سبيل المثال ومقياسه رتبوي ثلاثي (أساسي أو دون ذلك/ ثانوي ED1 / عالي ED2). وأن ED2 الذي يمثل مستوى التعليم العالي وذو معنوية عالية جداً حيث أن (P-value < 0.001). وهذا يشير إلى أن المرأة التي لديها مستوى تعليم عالي يتوقع أن تكون أقل عرضة للبدانة مقارنة مع التي لديها مستوى تعليم أساسي أو دون ذلك عند تماثل جميع المتغيرات الأخرى للمرأتين. وبشكل أدق، فإن قيمة نسبة الأرجحية (OR = 0.57) تعني أن المرأة التي لديها مستوى تعليم عالي ينخفض لديها احتمال البدانة إلى 0.57 من نفس الإحتمال مع تلك التي لديها مستوى تعليم أساسي أو دون ذلك. أي أن هذا الإحتمال سينخفض إلى النصف تقريباً مع زيادة مستوى التعليم على نحو الشكل المبين هنا. وفي ضوء التفسير الآخر لهذه القيمة فإن مقابل كل امرأة بدينة ضمن مجموعة المستوى التعليمي العالي سننتوقع أن نجد امرأتان تقريباً بهذا الوضع ضمن مجموعة المستوى التعليمي الأساسي أو أقل.

## 3) بالنسبة للمتغيرات المعنوية بإشارة موجبة إلى تقدير $\beta$

وهنا علينا أن نتوقع إتجاهاً عكسياً لما توصلنا إليه في الفقرة (2) أعلاه.

لنأخذ المتغير PAR ويمثل حجم الولادات ومقياسه فنوي ثلاثي (حد أعلى 1 / (2-3) PAR1 / (PAR2 ≥ 4)). وأن PAR2 والذي يمثل عدد الأطفال (4 ≥) ونو معنوية عالية حيث أن (P-value < 0.01). وهذا يشير إلى أن المرأة التي لديها أربعة أطفال فأكثر يتوقع أن تكون أكثر عرضة للبدانة مقابل التي لديها طفل واحد في الأكثر عند تماثل جميع المتغيرات

الأخرى للمرأتين. وبشكل أدق، فإن قيمة نسبة الأرجحية (OR = 1.75) تعني أن المرأة التي لديها طفل واحد أو بدون وقد تكون بدينة بإحتمال معين، سيزداد إحتمال أن تصبح بدينة بمقدار 75% عندما يكون لديها أربعة أطفال. أي أن هذا الإحتمال سيتضاعف تقريباً مع زيادة عدد الأطفال. وتفسيراً آخر لهذه القيمة يفيد بأنه مقابل كل امرأة بدينة ضمن مجموعة ذوي طفل واحد أو بدون سنتوقع إمرأتان تقريباً تتسمان بالبدانة ضمن مجموعة ذوي أربعة أطفال فأكثر.

لكننا، من جانب آخر، نلاحظ عدم معنوية PAR1 والذي يمثل حالة الأمومة لطفلين أو ثلاثة مع نسبة أرجحية (OR = 1.18) والتي تعني أن المرأة التي لديها طفل واحد أو بدون وقد تكون بدينة بإحتمال معين، سيزداد إحتمال أن تصبح بدينة بمقدار 18% فقط عندما يكون لديها طفلان أو ثلاثة. وهذا يعتبر صغيراً لا تأثير له بتفسير المعنوية.

وفي مثل حالة هذا المتغير حيث أحد المستويات التعليمية (العالى ED2) له تأثير معنوي واضح والآخر (الثانوي ED1) بتأثير لا معنوي، فإننا لا نحتاج إلى المقارنة بينهما من خلال مقياس الخطورة النسبية (RR) والذي يستساغ أحياناً استخدامه للمقارنة ما بين خطورة مستويين وبشكل معنوي تجاه المتغير المعتمد نسبة إلى مستوى المقارنة. ولتوضيح هذا الجانب، دعنا نأخذ المتغير REG والذي يمثل المنطقة الجغرافية (المركز / شمال REG1 / جنوب REG2) ومقياسه إسمي ثلاثي. والمركز هنا، بطبيعة الحال، يمثل مستوى المقارنة. ومن الجدول أعلاه، نجد أن نسبة الأرجحية OR لكل من الشمال (REG1) والجنوب (REG2) هي 1.40 بمستوى معنوية (P-value < 0.01) و 1.70 بمستوى معنوية (P-value < 0.001) على التوالي. وبذلك فإن الجنوب يعتبر أكثر خطورة من الشمال تجاه حدوث البدانة مقارنة بالمركز. ولكي نقارن بين خطورة هاتين المنطقتين، فإننا نحسب الخطورة النسبية RR لهما وحسب ما يلي:

$$RR = 1.70/1.40 = 1.22$$

ومعنى ذلك أن منطقة الجنوب تعتبر، تقديرياً، 1.22 مرة أكثر خطورة لحدوث البدانة من الشمال مقارنة بالمركز.

## تحليل التباين متعدد المتغيرات

### Multivariate Analysis of Variance (MANOVA)

إن تحليل التباين متعدد المتغيرات MANOVA عبارة عن متعدد متغيرات تعميمي لحالة تحليل التباين الأحادي ANOVA والذي هو عبارة عن أسلوب يستخدم لمقارنة أوساط عدد من المجتمعات حول متغير مقياس واحد. وعندما تكون هنالك عدد من المتغيرات المقاسة على كل وحدة تجريبية، فإنه بإمكاننا تطبيق أسلوب ANOVA بشكل منفرد على كل واحد من هذه المتغيرات. على سبيل المثال، لو كان هنالك 15 متغير فإن بإمكان الباحث تنفيذ 15 تحليلاً منفصلاً بواقع تحليل ANOVA واحد منفرد لكل متغير. وعلى أية حال، هذا ليس معقولاً ولكن هذا ما نلاحظه في الغالب.

بالنسبة للإحصائيين، فإنه لديهم اعتراضين رئيسيين تجاه التحليل الإفرادي هذا وهما:

- (1) إن هذه المجتمعات قد تكون مختلفة بالنسبة لبعض المتغيرات ولكنها ليست كذلك بالنسبة لمتغيرات أخرى. والباحث هنا يواجه حيرة حول قرار منطقي في أي من هذه المجتمعات يمكن إختبارها مختلفة وأياً منها متجانس. إن تحليل التباين متعدد المتغيرات يمكن أن يساعد الباحث حينها في إجراء المقارنة بين هذه المجتمعات آخذاً في الإعتبار جميع المتغيرات المشمولة بالتحليل في نفس الوقت وبعملية واحدة وبالتالي اتخاذ القرار المنطقي المناسب بشأن حقيقة إختلاف هذه لمجتمعات من عدمها.
- (2) ليست هنالك حماية (صيغة) كافية تجاه عمل الأخطاء من النوع الأول Type I Error عند عمل التحليل واحداً بعد آخر وبشكل منفرد (تذكر من مبادئ الإحصاء كون الخطأ من النوع الأول يحدث عند رفض فرضية العدم  $H_0$  وهي صحيحة). وهذا ناتج عن حقيقة أنه كلما زاد عدد المتغيرات لدى الباحث للتحليل، كلما زادت امكانية ظهور واحد في الأقل من هذه المتغيرات ليعطي زيادة في مستوى الثقة الإحصائية Statistical Significance. بمعنى آخر، كلما زاد عدد المتغيرات الخاضعة للتحليل كلما زاد احتمال ظهور أحد التحليلات بفرق معنوي  $\{Pr(p\text{-value} < 0.05) \rightarrow 1\}$

بطبيعة الحال، على الباحث أن يكون حذراً لتفادي حالة زيادة احتمال الخطأ من النوع الأول. وعليه أن يكون أكثر ثقة عندما يدعي (يستنتج) بأن مجتمعين أو أكثر لديها أوساط حسابية مختلفة بالنسبة لأحد المتغيرات وأن لا يكون أحداً غيره قد يستنتج عكس ذلك من خلال تطبيق نفس التحليل وعلى نفس البيانات.

إنه من الضروري تطبيق تحليل التباين متعدد المتغيرات متى ما كانت هناك مقارنة ما بين مجتمعين أو أكثر على أساس عدد كبير من المتغيرات. فإذا ما أظهر تحليل التباين متعدد المتغيرات فروقات معنوية فإن الباحث يكون واثقاً من أن هذه الفروقات حقيقية.

أما في حالة كون تحليل التباين متعدد المتغيرات لم يظهر أية فروقات معنوية، فإن على الباحث أن يكون شديد الحذر في إعطاء أي استنتاج عند عمل التحليل المنفرد لكل متغير لأن هذه التحليلات قد لا تحدد أي شيء أكثر من " إيجابية كاذبة " false positive.

## T مقابل T<sup>2</sup>

لكي نرى طبيعة تحليل MANOVA، دعنا نبدأ من نقطة الصفر في هذا الإتجاه والذي موضوعه الرئيسي إيجاد إحصاءة إختبار للفرضيات حول الأوساط الحسابية للمجتمعات Population Means وبإختلاف صيغها من ثنائية إلى أكثر من ذلك.

## إختبار t

ويستخدم إختبار Studentized t لإختبار الفرضية :

$$H_0 : \mu_1 = \mu_2$$

$$H_1 : \mu_1 \neq \mu_2 \text{ or } \mu_1 > \mu_2 \text{ or } \mu_1 < \mu_2$$

ولغرض تسلسل الأفكار وتسهيل عملية الربط، فإننا نستخدم  $H_1 : \mu_1 \neq \mu_2$  حيث أن  $\mu_i$  هي الوسط الحسابي الحقيقي للمجتمع (i). ولإجراء هذا الإختبار نستخدم إحصاءة الإختبار:

$$T = \frac{(\bar{X}_1 - \bar{X}_2)}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t_{n_1+n_2-2}$$

$$S_p^2 = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}$$

وتسمى  $S_p^2$  بالتباين المدمج Pooled Variance وأن  $\bar{X}_1$  هي الوسط الحسابي لعينة بحجم  $n_1$  من المجتمع الأول و  $\bar{X}_2$  هي الوسط الحسابي لعينة بحجم  $n_2$  من المجتمع الثاني وليس من الضرورة أن تكون  $n_1 = n_2$ . كما أن  $S_1^2$  و  $S_2^2$  هما التباين من العينتين وبمثابة التقدير للتباين  $\sigma_1^2$  و  $\sigma_2^2$  للمجتمعين.

وعلينا أن نتذكر بأن كلاً من  $\mu_1$  و  $\mu_2$  هما بقيمة واحدة (أي 1x1).

والآن إذا ما كان أيّاً من المجتمعين يتضمن مجموعتين أو أكثر، فهذا يعني أن وسطي المجتمعين سيكونا  $\mu_1$  و  $\mu_2$  وكلٍ منهما عبارة عن متجه وكان يكون بحجم (px1).

أي أن:

$$\underline{\mu}_1 = \begin{bmatrix} \mu_{11} \\ \mu_{12} \\ \vdots \\ \mu_{1p} \end{bmatrix}, \quad \underline{\mu}_2 = \begin{bmatrix} \mu_{21} \\ \mu_{22} \\ \vdots \\ \mu_{2p} \end{bmatrix}$$

وبالتالي فإننا بصدد إختبار الفرضية:

$$H_0 : \underline{\mu}_1 = \underline{\mu}_2$$

$$H_1 : \underline{\mu}_1 \neq \underline{\mu}_2$$

وهنا يصبح لدينا موضوع تعدد المتغيرات. وبالتالي، فإن إختبار t لم يعد ممكناً إستخدامه هنا وإنما نستخدم إختبار هوتلنك  $T^2$  والذي هو ما يقابل إختبار t عندما تكون الأوساط بصيغة متجهات (px1).

فبينما تكون الحالة الأولى بمتغير منفرد لكلا المجتمعين مفترضين

$$X_1 \sim N(\mu_1, \sigma_1^2)$$

$$X_2 \sim N(\mu_2, \sigma_2^2) \quad , \quad \sigma_1^2 = \sigma_2^2 = \sigma^2$$

نجد في حالة متعدد المتغيرات أن:

$$\underline{X}_1 = \begin{bmatrix} X_{11} \\ X_{12} \\ \vdots \\ X_{1p} \end{bmatrix}, \quad \underline{X}_2 = \begin{bmatrix} X_{21} \\ X_{22} \\ \vdots \\ X_{2p} \end{bmatrix}$$

ولذلك فإننا نفترض:

$$\underline{X}_1 \sim N_p(\underline{\mu}_1, \Sigma_1)$$

$$\underline{X}_2 \sim N_p(\underline{\mu}_2, \Sigma_2) \quad , \quad \Sigma_1 = \Sigma_2$$

$$\Sigma_1 = (\sigma_{ij}^{(1)})_{p \times p} \quad , \quad \Sigma_2 = (\sigma_{ij}^{(2)})_{p \times p}$$

$$(\sigma_{ij}^{(1)})_{p \times p} = (\sigma_{ij}^{(2)})_{p \times p} = \Sigma_{p \times p}$$

وبعد أن نحسب مصفوفة التباين المدمج  $\hat{\Sigma}$  وحسب الصيغة:

$$\hat{\Sigma} = \frac{(n_1 - 1)\hat{\Sigma}_1 + (n_2 - 1)\hat{\Sigma}_2}{n_1 + n_2 - 2}$$

فإنه بالإمكان إجراء إختبار للفرضية أعلاه باستخدام إحصاءة الإختبار  $T^2$  بالصيغة التالية:

$$T^2 = \frac{n_1 n_2}{n_1 + n_2} (\bar{X}_1 - \bar{X}_2)' \hat{\Sigma}^{-1} (\bar{X}_1 - \bar{X}_2)$$

حيث أن  $\bar{X}_1$  و  $\bar{X}_2$  هما تقديران لوسطي المجتمعين  $\mu_1$  و  $\mu_2$  على التوالي وأن:

$$\bar{X}_1 = \begin{bmatrix} \bar{X}_{11} \\ \bar{X}_{12} \\ \vdots \\ \bar{X}_{1p} \end{bmatrix}, \quad \bar{X}_2 = \begin{bmatrix} \bar{X}_{21} \\ \bar{X}_{22} \\ \vdots \\ \bar{X}_{2p} \end{bmatrix}$$

ولغرض التبسيط في إستكمال قرار الرفض أو القبول للفرضية:  $H_0: \mu_1 = \mu_2$  ، فإنه بالإمكان استخدام توزيع F لهذا الغرض وعلى النحو التالي:

رفض  $H_0$  إذا كان:

$$\frac{(n_1 + n_2 - p - 1)}{p(n_1 + n_2 - 2)} T^2 > f_{p, (n_1 + n_2 - p - 1)}$$

## ANOVA مقابل MANOVA

من الممكن استخدام ANOVA بدلاً من إختبار t في حالة وجود مجتمعين أو أكثر قيد فرضية الإختبار وأن حجم  $\mu_i$  هو (1x1) لجميع المجتمعات. أي أننا نختبر الفرضية:

$$H_0: \mu_1 = \mu_2 = \dots = \mu_k, \quad k \geq 2$$

$$H_1: \text{at least } \mu_i \neq \mu_j, \quad i \neq j$$

$$\sigma_i^2 = \sigma^2, \quad i = 1, 2, \dots, k \quad \text{مفترضين أن:}$$

وعملية الإختبار هذه تتم من خلال جدول تحليل التباين الآتي:

S.V.	df	SS	MS	F
Between	k-1	SSB	MSB=SSB/(k-1)	MSB/MSE
Within(Error)	N -k	SSE	MSE=SSE/(N-k)	
Total	N -1	TSS		

حيث أن  $N = \sum_{i=1}^k n_i$  وأن  $n_i$  هو حجم العينة من المجتمع (i). ومكونات الجدول أعلاه تحتسب بالشكل التالي:

$$TSS = \sum_{i=1}^k \sum_{j=1}^{n_i} X_{ij}^2 - \frac{X_{..}^2}{N}$$

$$X_{..} = \sum_i \sum_j X_{ij}$$

$$SSB = \frac{\sum_i X_{i.}^2}{n_i} - \frac{X_{..}^2}{N}$$

$$SSE = TSS - SSB$$

والبيانات التي نتعامل معها ستكون بالشكل التالي:

Populations (groups)			
1	2	.....	k
$X_{11}$	$X_{21}$	.....	$X_{k1}$
$X_{12}$	$X_{22}$	.....	$X_{k2}$
.	.		.
.	.		.
$X_{1n1}$	$X_{2n2}$	.....	$X_{knk}$
$X_{1.}$	$X_{2.}$	.....	$X_{k.}$

**ملاحظة:**

في حالة المجموعتين ( $k=2$ ) فإنه يمكن استخدام إحصاءة الاختبار F أو T وفي هذه الحالة فإن القيم الحسابية  $T^2 = F$  وهذه تقودنا إلى مقارنة القيم الجدولية

$$t_{\alpha/2, (n1+n2-2)}^2 = f_{\alpha, (1, n1+n2-2)}$$

وباختصار، فإذا كانت ANOVA تتعامل مع اختبار الفروقات ما بين أوساط حسابية للمجتمعات (إثنان أو أكثر) فإن MANOVA تتعامل مع اختبار الفروقات ما بين متجهات الأوساط الحسابية (إثنان أو أكثر).

بمعنى آخر، وبلغة المصفوفات، فإن ANOVA تتعامل مع متجه أوساط حسابية بأبعاد (1x1) لأي مجموعة (مجتمع)، بينما MANOVA تتعامل مع متجه أوساط حسابية بأبعاد (px1) لأي مجموعة (مجتمع) وأن p يمثل عدد المتغيرات المعتمدة الناتجة من خلال التجربة.

## ملاحظة:

يجدر بنا أن نعلم بأن أيًا من هاتين الطريقتين لا تعطينا إستنتاجاً مباشراً عن أي من هذه الأوساط سيكون مختلفاً عن الآخر في حالة تثبيت وجود فروقات بشكلٍ معنوي ما بين هذه الأوساط. ولتغلب على هذه الإشكالية في الطريقتين، يمكننا استخدام أساليب أخرى لهذا الغرض.

ومن أجل التوضيح، فإننا نستخدم ANOVA لإختبار الفرضية:

$$H_0 : \mu_1 = \mu_2 = \dots = \mu_k$$

حيث كل  $\mu_i$  عبارة عن قيمة مفردة (1x1).

بينما نستخدم MANOVA لإختبار الفرضية:

$$H_0 : \underline{\mu}_1 = \underline{\mu}_2 = \dots = \underline{\mu}_k$$

حيث كل  $\underline{\mu}_i$  عبارة عن قيمة مفردة (px1). أي أن:

$$\underline{\mu}_i = \begin{bmatrix} \mu_{i1} \\ \mu_{i2} \\ \vdots \\ \mu_{ip} \end{bmatrix}$$

وأن: k تمثل عدد المجاميع (المجموعات).

P هي حجم المتجه والذي يعكس عدد المتغيرات المعتمدة التي يتم قياسها خلال التجربة.

## كيفية اختبار الفرضية $H_0$

كما نتذكر الإختبار لهذه الفرضية في حالة ANOVA حيث أن TSS للمتغير المعتمد يتم تقسيمه إلى SSB و SSE. ومع التعدد في المتغير المعتمد وإعتماد أسلوب MANOVA، فإنه لا زال بالإمكان احتساب SSB و SSE والتي ستكون عبارة عن مصفوفة (pxp). بالإضافة إلى ذلك، فإنه بالإمكان احتساب مجموع الضرب المتبادل (المتقاطع) Sum of Cross Products وتجزئته إلى SSB و SSE المتجمع Pooled. وهنا في حالة MANOVA، لعله من المناسب توصيف الآتي من المصفوفات (مفترضين متغيرين معتمدين  $X_1$  و  $X_2$  لغرض التوضيح الأسهل لهذه العملية):

$$W = \text{Pooled within groups (SSCP)} = \begin{bmatrix} SSW_1 & SCP_w \\ SCP_w & SSW_2 \end{bmatrix}$$



والرموز هذه تعني:

SSCP = (Sum of Squares for Cross Product) مجموع مربعات الضرب المتبادل

مجموع المربعات المدمج ضمن:

$SSW_1$  = Pooled SS within groups for  $X_1$

$SSW_2$  = Pooled SS within groups for  $X_2$

$SCP_w$  = Pooled within group Sum of Product for  $X_1$  &  $X_2$

$$B = \begin{bmatrix} SSb_1 & SCP_b \\ SCP_b & SSb_2 \end{bmatrix}$$

$SSb_i$  = Between – groups SS for  $X_i$

$SCP_b$  = Between – groups Sum of Cross Products of  $X_1$  &  $X_2$

$$T = \begin{bmatrix} SS_1 & SCP_{12} \\ SCP_{12} & SS_2 \end{bmatrix}$$

$SS_i$  = The Total SS for  $X_i$

$SCP_{12}$  = The Total Sum of Cross Products of  $X_1$  &  $X_2$

ومن المعلوم أن  $T = B + W$  (Total = Between + Within).

ولأجل توضيح كيفية احتساب هذه العناصر للمصفوفات  $T$ ,  $B$ ,  $W$  دعنا نفترض

مجموعتين  $G_1$ ,  $G_2$  ومتغيرين معتمدين لكلٍ منهما  $X_1$  و  $X_2$  بالبيانات التالية:

G1		G2	
$X_1$	$X_2$	$X_1$	$X_2$
8	3	4	2
7	4	3	1
5	5	3	2
3	4	2	2
3	2	2	5
$\sum X = 26$	18	14	12
$\sum X^2 = 156$	70	42	38
$\sum X_1 X_2 = 95$		31	

$$\begin{aligned}\sum X_{i1} &= 26 + 14 = 40 \\ \sum X_{i2} &= 18 + 12 = 30 \\ \sum X_{i1}^2 &= 156 + 42 = 198 \\ \sum X_{i2}^2 &= 70 + 38 = 108 \\ CP_i &= 95 + 31 = 126\end{aligned}$$

ومن البيانات هذه يمكننا حساب العناصر التالية:

$$\begin{aligned}SSW_1 &= \left[ 156 - \frac{(26)^2}{5} \right] + \left[ 42 - \frac{(14)^2}{5} \right] = 23.6 \\ SSW_2 &= \left[ 70 - \frac{(18)^2}{5} \right] + \left[ 38 - \frac{(12)^2}{5} \right] = 14.4 \\ SCP_w &= \left[ 95 - \frac{(26)(18)}{5} \right] + \left[ 31 - \frac{(14)(12)}{5} \right] = -1.2 \\ W &= \begin{bmatrix} 23.6 & -1.2 \\ -1.2 & 14.4 \end{bmatrix}\end{aligned}$$

كذلك فإن:

$$\begin{aligned}SSb_1 &= \left[ \frac{(26)^2}{5} + \frac{(14)^2}{5} \right] - \frac{(40)^2}{10} = 14.4 \\ SSb_2 &= \left[ \frac{(18)^2}{5} + \frac{(12)^2}{5} \right] - \frac{(30)^2}{10} = 3.6 \\ SCP_b &= \left[ \frac{(26)(18)}{5} + \frac{(14)(12)}{5} \right] - \frac{(40)(30)}{10} = 7.2 \\ B &= \begin{bmatrix} 14.4 & 7.2 \\ 7.2 & 3.6 \end{bmatrix}\end{aligned}$$

وبالتالي وعن طريق الجمع نجد أن:

$$T = B + W = \begin{bmatrix} 14.4 & 7.2 \\ 7.2 & 3.6 \end{bmatrix} + \begin{bmatrix} 23.6 & -1.2 \\ -1.2 & 14.4 \end{bmatrix} = \begin{bmatrix} 38 & 6 \\ 6 & 18 \end{bmatrix}$$

**ملاحظة:**

جدير بالذكر أن عناصر المصفوفة T يمكن احتسابها بشكل مستقل وعلى النحو التالي:

$$SST_1 = 198 - \frac{(40)^2}{10} = 38$$

$$SST_2 = 108 - \frac{(30)^2}{10} = 18$$

$$SCP_t = (95 + 31) - \frac{(40)(30)}{10} = 6$$

ولأننا نعمل على مجموعتين، فإن الفرضية التي نحن بصدد عمل الإختبار لها هي:

$$H_0 : \underline{\mu}_1 = \underline{\mu}_2$$

وأن إحصاءة الإختبار العامة في حالة المجموعتين هي:

$$F = \frac{(1 - \Lambda)/t}{\Lambda/(N - t - 1)} \sim f_{t, (N-t-1)}$$

حيث  $t = 2$  وهي عدد المتغيرات المعتمدة  $X_1$  و  $X_2$

$N = 10$  المجموع الكلي لعدد المشاهدات لكل مجموعة

والرمز  $\Lambda$  يشير إلى ما يعرف (ولكس لامبدا Wilk's Lambda)

$$\Lambda = \frac{|W|}{|T|} = \frac{338.4}{648} = 0.52222$$

وحيث أن مثالنا هذا فيه  $N=10$  و  $t=2$  فإن:

$$F = \frac{(1 - 0.52222)/2}{0.52222/(10 - 2 - 1)} = 3.202$$

وهذه القيمة المحسوبة لإحصاءة الإختبار هي أقل من القيمة الجدولية المقابلة لها وهي:

$$f_{t, (N-t-1)} = f_{2, 7} (0.05) = 4.74$$

فإننا لا نرفض الفرضية  $H_0$  ونستنتج أن لا فروقات معنوية بين المجموعتين بالنسبة إلى  
أوساط مجموعة المتغيرات المعتمدة  $X_1$  و  $X_2$ .

### تنويه:

إن ما أوردناه ضمن هذا المثال لا يمثل الحالة العامة لإحتساب الإحصاءة  $F$  والتي سيتم  
تناولها تالياً.

## الحالة العامة لإحتساب الإحصاءة $F$

لو إفترضنا الأبعاد التالية في التجربة وهي:

$$K = \text{عدد المجاميع}$$

$$P = \text{عدد المتغيرات المعتمدة في كل مجموعة ومقياسها كمي (نسبي أو فنوي)}$$

$$N = \text{العدد الكلي للملاحظات}$$

وفي ضوء ذلك، ستكون قيمة الإحصاءة  $F$  حسب الصيغة التالية

$$F = \frac{\left(1 - \Lambda^{1/b}\right) / df_1}{\Lambda^{1/b} / df_2}$$

حيث أن:

$$df_1 = P(K - 1)$$

$$df_2 = ab - c$$

$$a = N - K - \frac{P - K + 2}{2}$$

$$b = \sqrt{\frac{P^2(K - 1)^2 - 4}{P^2 + (K - 1)^2 - 5}}$$

$$c = \frac{P(K - 1) - 2}{2}$$

مع ملاحظة أن قيمة  $b = 1$  في حالة عدم تحقق كون  $(P^2 + (K - 1)^2 - 5 > 0.0)$ .

وإذا ما طبقنا هذه الصيغة العامة على المثال أعلاه حيث  $(N=10, P=2, K=2)$

$$\text{نجد أن } (df_1 = 2, df_2 = 7, c = 0, b = 1, a = 7)$$

وبالتالي فإن الحسابات التي اعتمدت تكون متوافقة كلياً مع الحالة العامة.

وسوف نرى أبعاد تطبيق الحالة العامة من خلال المثال التالي لبيانات حقيقية:

**مثال:**

في تجربة زراعية لمعرفة ما إذا توجد فروقات معنوية ما بين 4 نوعيات من التربة (سطحية، رملية، ملحية، طينية) بالنسبة لزراعة نوع جديد من بذور الذرة في ضوء 3 من

المتغيرات المعتمدة (الناتج، كمية مياه الري، كمية المبيد الحشري) واستخدام 8 نباتات لكل نوع من التربة<sup>(5)</sup>.

$$K = \text{عدد المجاميع} = 4$$

$$P = \text{عدد المتغيرات المعتمدة} = 3$$

$$N = \text{العدد الكلي للملاحظات} = 32$$

وبالتالي سيكون لدينا، وفقاً للصيغ أعلاه، القيم التالية:

$$a = N - K - \frac{P - K + 2}{2} = 32 - 4 - \frac{3 - 4 + 2}{2} = 27.5$$

$$b = \sqrt{\frac{P^2(K-1)^2 - 4}{P^2 + (K-1)^2 - 5}} = \sqrt{\frac{3^2(4-1)^2 - 4}{3^2 + (4-1)^2 - 5}} = 2.434$$

$$c = \frac{P(K-1) - 2}{2} = \frac{3(4-1) - 2}{2} = 3.5$$

$$df_1 = P(K-1) = 3(4-1) = 9$$

$$df_2 = ab - c = (27.5)(2.434) - 3.5 = 63.43$$

ومن خلال حسابات مصفوفات التباين والتباين المشترك الثلاثة W و B و T وجدنا الآتي:

$$W = \begin{bmatrix} 4058 & 714 & -273 \\ 714 & 2834 & 123 \\ -273 & 123 & 113 \end{bmatrix}$$

$$B = \begin{bmatrix} 911 & 63 & 163 \\ 63 & 122 & 24 \\ 163 & 24 & 32 \end{bmatrix}$$

$$T = \begin{bmatrix} 4969 & 777 & -110 \\ 777 & 2956 & 147 \\ -110 & 147 & 145 \end{bmatrix}$$

وبالتالي نجد أن:

$$\Lambda = \frac{|W|}{|T|} = 0.489$$

$$F = \frac{1 - (0.489)^{1/2.434}}{(0.489)^{1/2.434}} \left( \frac{63.43}{9} \right) = 2.405 > f_{9, 63.43} (0.05) = 2.032$$

وبالتالي يتم رفض الفرضية:

$$H_0 : \underline{\mu}_1 = \underline{\mu}_2 = \underline{\mu}_3 = \underline{\mu}_4$$

مستنتجين بوجود فروق معنوية ما بين اتجاهات الأوساط الحسابية الأربعة والتي يتضمن كلاً منها ثلاثة أوساط حسابية للمتغيرات المعتمدة. أي أن اتجاهات أوساط (الناتج، مياه الري، مبيد الحشرات) تختلف بشكلٍ معنوي ما بين نوعيات التربة (سطحية، رملية، ملحية، طينية).

## التحليل المميز

### Discriminant Analysis (DA)

إن التحليل المميز يستخدم بشكل رئيسي لتصنيف الأفراد أو الوحدات التجريبية إلى إثنين أو أكثر من المجتمعات المحددة بشكل منفرد لا تداخل فيما بينها. ولأجل تطوير قاعدة تمييز لغرض تصنيف Classifying الوحدات التجريبية لواحدة من عددٍ من الفئات المحتملة، فعلى الباحث أن يحصل على عينة عشوائية من الوحدات التجريبية من كل فئة محتملة للتصنيف. وبعد ذلك، فإن التحليل المميز يهيئ طرقاً تمكن الباحث من بناء قواعد يمكن استخدامها لتصنيف وحدات تجريبية أخرى ضمن واحدة من الفئات التصنيفية.

ومثالاً على ذلك، لو أن شركة معينة متخصصة في منح بطاقات إئتمانية Credit Cards فإنها حتماً ترغب في أن تكون قادرة على تصنيف طلبات الحصول على البطاقة إلى مجموعتين من الأفراد (1- أفراداً جيدين من ناحية خطورة الإئتمان) و (2- أفراداً غير جيدين من ناحية خطورة الإئتمان). والشركة تمنح المجموعة الأولى بطاقات الإئتمان فيما لا تمنحها للمجموعة الثانية. ولأجل مساعدة الشركة في هذا الجانب وتحديد المجموعتين، فإنها قد تأخذ في الاعتبار عدد من الخصائص الديموغرافية التي يمكن قياسها لدى كل فرد. فالشركة قد تأخذ، على سبيل المثال، المستوى التعليمي، الراتب الشهري، المديونيات، سجل الماضي في الإئتمان كمؤشرات تنبؤية بشأن استحقاق البطاقة. والشركة بعد ذلك تحاول استخدام هذه المعلومات لفردٍ محدد للمساعدة في اتخاذ قرار منحه البطاقة من عدمها. الشركة بحاجة لتوظيف طريقة من طرق التحليل متعدد المتغيرات لمساعدتها في تصنيف الأفراد لأي من المجموعتين، وهذه الطريقة هي طريقة التحليل المميز.

بالنسبة لهذا المثال، على الشركة أن تجمع هذه البيانات من أفراد معروفين لديها بأنهم ينتمون للمجموعة الأولى وإعادة أخذ نفس البيانات من أفراد معروفين لديها أيضاً بأنهم ينتمون للمجموعة الثانية. ومن ثم بإستطاعة الشركة تصنيف طالبي البطاقة الجدد لأي من المجموعتين بإستخدام القاعدة الناتجة من تطبيق التحليل المميز على الأفراد المعروفين لديها.

من المهم هنا أن يكون في بالنا احتمال ظهور خطأ في التصنيف وباحتمال معين يتم تشبيته من خلال إعادة تصنيف الأفراد المعروفين لديها في ضوء تطبيق القاعدة الناتجة عليهم. ويمكن أن نحكم على جودة قاعدة التصنيف الناتجة بتناسب عكسي مع مقدار احتمال خطأ إعادة التصنيف لكامل المجموعة.

إن التحليل المميز هو أحد المواضيع المهمة في التحليل المتعدد المتغيرات multivariate analysis والذي يهتم في كيفية التمييز بين مجموعتين أو أكثر. إن الفكرة الأساسية من التمييز discriminate هو التفريق ما بين المجتمعات المتداخلة أو المتشابهة ولها نفس الخصائص أو

الصفات. بمعنى آخر، لنفرض انه لدينا مجتمعين أو اكثر ولدينا عينة تحتوي على مجموعة من المشاهدات من كل مجتمع. إن وظيفة التحليل المميز هي إيجاد دالة يمكن بواسطتها تصنيف أو تمييز المشاهدات الجديدة بالنسبة لمجمعاتها الأصلية.

إن التحليل المميز يختلف عن تحليل الإنحدار في أن المتغير المعتمد في التحليل التمييزي هو متغير ذو مقياس إسمي nominal variable وهو من المتغيرات النوعية (يأخذ قيمتين , 0 , 1) في حالة التمييز ما بين مجموعتين

$Y = 1$  إذا كانت المشاهدة تعود للمجتمع الأول.  
 $Y = 0$  إذا كانت المشاهدة تعود للمجتمع الثان.

بينما المتغير المعتمد في تحليل الإنحدار هو على الأكثر متغير مستمر (وهو من المتغيرات الكمية) ويتشابه التحليلين بأن كلاهما يهدف لإيجاد علاقة بين المتغير المعتمد والمتغيرات المستقلة.

### أنواع الدوال التمييزية

1. الدالة المميزة الخطية
2. الدالة المميزة التربيعية
3. الدالة المميزة اللوجستية

### 1. الدالة المميزة الخطية Linear Discriminant Function

تستخدم هذه الدالة عندما تكون المجتمعات المدروسة ذات توزيع طبيعي متعدد المتغيرات بمتجهات متوسط مختلفة و مصفوفة تباين مشترك متساوية وهناك حالتان:

- 1- حالة مجموعتين (مجتمعين)
- 2- حالة عدة مجاميع (مجتمعات)

### حالة المجموعتين

نفرض لدينا عينة مسحوبة من مجتمعين يتوزعان توزيعاً طبيعياً بمتوسطين  $\mu_1$  ,  $\mu_2$  ومصفوفة تباين مشترك  $\Sigma$  لكل مجتمع. أي أن:

$$\underline{X}_1 \sim N_p(\underline{\mu}_1 , \Sigma )$$

$$\underline{X}_2 \sim N_p(\underline{\mu}_2 , \Sigma )$$

يمكن صياغة دالة بالإعتماد على مقاييس من هذه القيم وأن هذه الدالة تمكننا من إختيار أي مشاهدة و تحديد المجتمع الذي تعود إليه. إن المتغير العشوائي  $X$  له عادة دالة كثافة إحصائية إما  $f_1$



$(x, \theta_1)$  أو  $f_2(x, \theta_2)$  و  $\theta_i$  هنا ترمز لأي معلمة يختلف التوزيع بموجبها. وبذلك، فإن هذا يعني أن:

$$f_i(x, \mu_i) = \frac{1}{(2\pi)^{\frac{p}{2}} (\Sigma)^{\frac{1}{2}}} e^{-\frac{1}{2}(x-\mu_i)' \Sigma^{-1} (x-\mu_i)}, \quad i=1,2$$

وبذلك فإن نسبة الإمكان الأعظم هنا تكون:

$$\frac{f_1(x, \mu_1)}{f_2(x, \mu_2)} = \frac{\frac{1}{(2\pi)^{\frac{p}{2}} (\Sigma)^{\frac{1}{2}}} e^{-\frac{1}{2}(x-\mu_1)' \Sigma^{-1} (x-\mu_1)}}{\frac{1}{(2\pi)^{\frac{p}{2}} (\Sigma)^{\frac{1}{2}}} e^{-\frac{1}{2}(x-\mu_2)' \Sigma^{-1} (x-\mu_2)}} \geq \lambda$$

أو أن:

$$\exp\left[-\frac{1}{2}\left\{(x-\mu_1)' \Sigma^{-1} (x-\mu_1) - (x-\mu_2)' \Sigma^{-1} (x-\mu_2)\right\}\right] \geq \lambda$$

وعند تفكيك هذه المعادلة وإضافة وطرح المقدار  $\mu_2' \Sigma^{-1} \mu_1$  يصبح لدينا:

$$\exp\left[\underline{x}' \Sigma^{-1} (\underline{\mu}_1 - \underline{\mu}_2) - \frac{1}{2} \left\{(\underline{\mu}_1 + \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 - \underline{\mu}_2)\right\}\right] \geq \lambda$$

وبأخذ اللوغاريتم الطبيعي للطرفين:

$$\left[\underline{x}' \Sigma^{-1} (\underline{\mu}_1 - \underline{\mu}_2) - \frac{1}{2} \left\{(\underline{\mu}_1 + \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 - \underline{\mu}_2)\right\}\right] \geq \log \lambda$$

وعندما تكون  $\lambda = 1$  فإن:

$$\left[\underline{x}' \Sigma^{-1} (\underline{\mu}_1 - \underline{\mu}_2) - \frac{1}{2} \left\{(\underline{\mu}_1 + \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 - \underline{\mu}_2)\right\}\right] \geq 0$$

وبإستخدام مقدرات الإمكان الأعظم إلى كل من  $\Sigma, \mu_2, \mu_1$  وتعويضها هنا، يصبح لدينا:

$$W = \left[\underline{x}' S^{-1} (\bar{x}_1 - \bar{x}_2) - \frac{1}{2} (\bar{x}_1 + \bar{x}_2)' S^{-1} (\bar{x}_1 - \bar{x}_2)\right] \geq 0$$

$$\bar{x}_1 = \sum_{j=1}^{n_1} x_{1j} / n_1$$

حيث أن:

$$\bar{x}_2 = \sum_{j=1}^{n_2} x_{2j} / n_2$$

$$S = \left[ \sum_{j=1}^{n_1} (\underline{x}_{1j} - \bar{x}_1)(\underline{x}_{1j} - \bar{x}_1) + \sum_{j=1}^{n_2} (\underline{x}_{2j} - \bar{x}_2)(\underline{x}_{2j} - \bar{x}_2) \right] / (n_1 + n_2 - 2)$$

وأن جزئي المعادلة  $w$  أعلاه يتكون من دالة التمييز الخطية:

$$y = \underline{x}' S^{-1} (\bar{x}_1 - \bar{x}_2)$$

ونقطة الفصل أو التمييز  $z$ :

$$z = \frac{1}{2} D^2 = \frac{1}{2} (\bar{x}_1 + \bar{x}_2)' S^{-1} (\bar{x}_1 - \bar{x}_2)$$

حيث أن  $D^2$  هي إحصاءة مهلونوبيس Mahalanobis.

وأن دالة التمييز الخطية ممكن أن تكتب على شكل الدالة الخطية التالية:

$$y = \underline{x}' \underline{C} \quad \text{or} \quad y = \underline{C}' \underline{x}$$

حيث:

$$\underline{C} = S^{-1} (\bar{x}_1 - \bar{x}_2)$$

والآن يمكن أن نستخلص الآتي من دالة التمييز الخطية:

$$\bar{y}_1 = \bar{x}_1' S^{-1} (\bar{x}_1 - \bar{x}_2)$$

$$\bar{y}_2 = \bar{x}_2' S^{-1} (\bar{x}_1 - \bar{x}_2)$$

وبالتالي يمكننا التعبير عن نقطة الفصل بالآتي:

$$z = \frac{\bar{y}_1 + \bar{y}_2}{2}$$

ولو إفتراضنا أن  $\bar{y}_1 < \bar{y}_2$  فإن خطة التصنيف للمشاهدة الجديدة  $y$  هو أنها:

تعود للمجموعة الأولى في حالة كون  $y \leq z$

تعود للمجموعة الثانية في حالة كون  $y > z$

أو يمكن إستخدام  $w$  أعلاه بالكامل لغرض التصنيف للمشاهدة  $x$  بكونها:

تعود للمجموعة الأولى في حالة كون  $w \leq 0$

تعود للمجموعة الثانية في حالة كون  $w > 0$

وجدير بالذكر أن المصفوفة S تستخرج بالشكل التالي (مفترضين  $n = 3$ ):

$$S = \begin{bmatrix} V_{11} & V_{12} & V_{13} \\ & V_{22} & V_{23} \\ & & V_{33} \end{bmatrix}$$

حيث أن:

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

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

### الإختبارات المستخدمة في التحليل المميز

مع أن كفاءة الدالة المميزة الخطية تقاس وفقاً لنسبة التصنيف الصحيح للمشاهدات حسب مجاميعها الأصلية، إلا أنه من الممكن التنبؤ مسبقاً بشئ عن هذه الكفاءة لأنها تزداد وفقاً لزيادة الفرق ما بين أوساط المجاميع من جهة، وتقارب قيم مصفوفات التباين – التباين المشترك لهذه المجاميع من جهة أخرى. وفي لغة الإحصاء ومفهومه، فإن ذلك يعني تطبيق إختبارات إحصائية ونستعين بنتائجها لتقييم كفاءة الدالة المميزة الخطية. وفيما يلي الإختبارات المناسبة في هذا الجانب.

1. إختبار معنوية الفرق بين الأوساط عن طريق إختبار الفرضية:

$$H_0 : \underline{\mu}_1 = \underline{\mu}_2$$

$$H_1 : \underline{\mu}_1 \neq \underline{\mu}_2$$

ويتم الإختبار بإستخدام إختبار F الذي يعتمد على إحصاءة هوتلنك  $T^2$  Hotteling والتي تكون:

$$T^2 = \frac{n_1 n_2}{n_1 + n_2} D^2$$

حيث أن  $D^2$  هي إحصاءة مهلونوبيس Mahalanobis والتي صيغتها

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

وبالتالي، فإن إختبار F سيكون:

$$F = \frac{n_1 + n_2 - p - 1}{n_1 + n_2 - 2} T^2 \sim f_{p, n_1 + n_2 - p - 1}$$

والذي نريد نتيجته الرفض بمعنوية عالية لعكس زيادة نسبة التصنيف الصحيح. حيث:

$n_1$  : حجم العينة الأولى

$n_2$  : حجم العينة الثانية

$P$  : عدد المتغيرات

2. إختبار تساوي مصفوفتي التباين- التباين المشترك للمجموعات من خلال الفرضية:

$$H_0 : \Sigma_1 = \Sigma_2$$

$$H_1 : \Sigma_1 \neq \Sigma_2$$

وأن إحصاءة الإختبار لهذه الفرضية  $Q$  وهي:

$$Q = MC^* \sim \chi^2_{(k-1)(p-1)}$$

حيث أن:

$$M = \ln \frac{|S|^{n_1+n_2}}{|S_1|^{n_1}|S_2|^{n_2}}$$

$$\text{or } M = \sum_{i=1}^2 n_i \ln|S| - \sum_{i=1}^2 n_i \ln|S_i|$$

$$S = \frac{A_1 + A_2}{n_1 + n_2 - 2}$$

$$A_1 = (n_1 - 1)S_1$$

$$A_2 = (n_2 - 1)S_2$$

$$C^* = 1 - \frac{2p^2 + 3p - 1}{6(p+1)(k-1)} \left[ \sum \frac{1}{n_i} - \frac{1}{\sum n_i} \right]$$

وأن:

$$k = \text{عدد المجاميع}$$

$$p = \text{عدد المتغيرات}$$

إن الفرضية أعلاه والتي تنص على تساوي مصفوفتي التباين - التباين المشترك بين المجموعتين هي أهم فرضية هنا لأن عدم تحققها يعني أننا لا يمكن أن نستخدم الدالة التمييزية الخطية وإنما نستخدم الدالة التمييزية التربيعية.

### إحتمال خطأ التصنيف: the probability of misclassification

هو إحتمال تصنيف مشاهدة معينة إلى المجموعة الأولى بينما هي تعود في الحقيقة إلى المجموعة الثانية و بالعكس. نفترض لحساب خطأ التصنيف أن حجم العينة يكون كبير لذلك فإننا نضمن كون توزيع المشاهدات يقترب من التوزيع الطبيعي (حسب نظرية الحد المركزي). حيث

أن هذا الخطأ يعتمد على أن توزيع العينة هو التوزيع الطبيعي أو يقترب من التوزيع الطبيعي. هذا الاحتمال يكون:

$$P_{12} = P(\text{classifying } x \text{ to be from group(1)/ } x \text{ is from group(2)}) \\ = \phi(-D/2)$$

حيث  $D^2$  هي إحصاء مهالونوبيس.

ويتم إيجاد هذه القيمة من جداول التوزيع الطبيعي القياسي. إن خطأ التصنيف هو عامل مهم لإثبات كفاءة الدالة المميزة. والتي تعطي أقل خطأ تصنيف هي الدالة الأكثر كفاءة و تكون الأفضل من بين دوال التمييز.

وبالإمكان أيضاً استخدام طريقة إعادة التعويض في هذا الجانب وكما هو في أدناه.

### طريقة التعويض : Resubstitution Method

تستخدم هذه الطريقة لإيجاد احتمال خطأ التصنيف وأن أسلوب هذه الطريقة يعتمد على أنه لو كان  $n_j$  يمثل عدد المشاهدات التي تعود للمجموعة (j) وأن  $n_{ij}$  هو عدد المشاهدات في المجموعة (j) وصنفت وفق دالة التمييز على انها تعود للمجموعة (i)، فإن تقدير احتمال خطأ التصنيف في هذه الحالة وبشكل عام سيكون:

$$P_{ij} = \frac{n_{ij}}{n_j}$$

وفي حالة المجموعتين التي نحن فيها الآن، فإن احتمال خطأ التصنيف الكلي للدالة المميزة سيكون:

$$\hat{p} = \frac{n_{12} + n_{21}}{n_1 + n_2}$$

أما متوسط احتمال خطأ التصنيف سيكون:

$$\hat{p} = \frac{\hat{p}_{21} + \hat{p}_{12}}{2}$$

### حالة عدة مجاميع

في حالة كون مسألة التمييز بين أكثر من مجموعتين (k من المجاميع)، يتم التصنيف عن طريق المقارنة بين كل مجموعتين وتكون لذلك عدة دوال مميزة  $y_{ij}$  وعددها  $C_2^k$  وتكتب على النحو التالي:

$$y_{ij} = x'S^{-1}(\bar{x}_i - \bar{x}_j)$$

ولو كانت لدينا ثلاثة مجاميع ولكل مجموعة  $p$  من المتغيرات ( $p \geq 2$ ) فإنه من المناسب حساب:

$$\begin{aligned} \underline{C}_1 &= S^{-1}(\underline{\bar{x}}_1 - \underline{\bar{x}}_2) \\ \underline{C}_2 &= S^{-1}(\underline{\bar{x}}_1 - \underline{\bar{x}}_3) \\ \underline{C}_3 &= S^{-1}(\underline{\bar{x}}_2 - \underline{\bar{x}}_3) \end{aligned}$$

وبذلك يكون لدينا دوال التمييز التالية:

$$Y_{12} = \underline{x}' \underline{C}_1, \quad Y_{13} = \underline{x}' \underline{C}_2, \quad Y_{23} = \underline{x}' \underline{C}_3$$

ثم نجد ما يلي:

$$\begin{aligned} W_{12} &= \left[ \underline{x}' S^{-1}(\underline{\bar{x}}_1 - \underline{\bar{x}}_2) - \frac{1}{2}(\underline{\bar{x}}_1 + \underline{\bar{x}}_2)' S^{-1}(\underline{\bar{x}}_1 - \underline{\bar{x}}_2) \right] \\ W_{13} &= \left[ \underline{x}' S^{-1}(\underline{\bar{x}}_1 - \underline{\bar{x}}_3) - \frac{1}{2}(\underline{\bar{x}}_1 + \underline{\bar{x}}_3)' S^{-1}(\underline{\bar{x}}_1 - \underline{\bar{x}}_3) \right] \\ W_{23} &= \left[ \underline{x}' S^{-1}(\underline{\bar{x}}_2 - \underline{\bar{x}}_3) - \frac{1}{2}(\underline{\bar{x}}_2 + \underline{\bar{x}}_3)' S^{-1}(\underline{\bar{x}}_2 - \underline{\bar{x}}_3) \right] \end{aligned}$$

فتكون العلاقة بينها:

$$W_{23} = W_{13} - W_{12}$$

وتكون قاعدة التصنيف إذا كانت لدينا عدد المتغيرات ( $p \geq 2$ ) لكل مجموعة حسب الآتي:

تصنف المشاهدة  $x$  لواحدة من المجاميع التالية:

- ضمن المجموعة (1) إذا كانت  $W_{12} > 0$  و  $W_{13} > 0$
- ضمن المجموعة (2) إذا كانت  $W_{13} > 0$  و  $W_{13} > W_{12}$
- ضمن المجموعة (3) إذا كانت  $W_{13} > 0$  و  $W_{12} > W_{13}$

أما إذا كان لدينا متغير وحيد لكل مجموعة ( $p=1$ ) وأن أوساط المجاميع رتبت للسهولة بالشكل التالي  $\bar{X}_1 < \bar{X}_2 < \bar{X}_3$  فإن قواعد التصنيف للمشاهدة  $x$  تكون كالآتي:

- ضمن المجموعة (1) إذا كانت  $x < \frac{1}{2}(\underline{\bar{x}}_1 + \underline{\bar{x}}_2)$
- ضمن المجموعة (2) إذا كانت  $\frac{1}{2}(\underline{\bar{x}}_1 + \underline{\bar{x}}_2) \leq x \leq \frac{1}{2}(\underline{\bar{x}}_2 + \underline{\bar{x}}_3)$
- ضمن المجموعة (3) إذا كانت  $x > \frac{1}{2}(\underline{\bar{x}}_1 + \underline{\bar{x}}_2)$

ومصفوفة S لهذه الحالة تستخرج بالشكل التالي:

$$S = \begin{bmatrix} V_{11} & V_{12} & V_{13} \\ & V_{22} & V_{23} \\ & & V_{33} \end{bmatrix}$$

$$V_{ii} = \frac{S_{ii}(1) + S_{ii}(2) + S_{ii}(3)}{n_1 + n_2 + n_3 - 3}, \quad S_{ii} = \sum x_i^2 - \frac{(\sum x_i)^2}{n}$$

$$V_{ij} = \frac{S_{ij}(1) + S_{ij}(2) + S_{ij}(3)}{n_1 + n_2 + n_3 - 3}, \quad S_{ij} = \sum x_i x_j - \frac{\sum x_i \sum x_j}{n}$$

## 2. دالة التمييز التربيعية Quadratic Discriminant Function

وكما ذكرنا سابقاً، يستخدم هذا النوع من الدوال في حالة عدم تساوي مصفوفة التباين - التباين المشترك بين المجموعات. و تقدر معالم هذه الدالة بطريقة MLE بافتراض أن حجم العينة كبير بحيث يصبح من الممكن أن نفترض بأن المشاهدات تقترب من التوزيع الطبيعي (النظرية المركزية).

وفي حالة المجتمعين، فإن تقدير  $\mu_1$  هو  $\bar{X}_1$  و أن تقدير  $\mu_2$  هو  $\bar{X}_2$  وتقدير  $\Sigma_1$  هو  $S_1$  وتقدير  $\Sigma_2$  هو  $S_2$ . وأن مقياس التمييز هو V حيث

$$V = \frac{f_1(x_i)}{f_2(x_i)} > \text{or} < 1$$

$$G = \ln V = \ln f_1(x_i) - \ln f_2(x_i) > \text{or} < 0$$

ملاحظة:

في حالة وجود p من المتغيرات لكل مجموعة فإن G هذه بطبيعة الحال ستكون:

$$G = \ln f_1(x_1, x_2, \dots, x_p) - \ln f_2(x_1, x_2, \dots, x_p)$$

ووفقاً لما جاء في أعلاه، فإن دالة التمييز التربيعية التقديرية في حالة متغيرين ستكون:

$$\hat{G} = \frac{1}{2} \ln \frac{S_2}{S_1} - \frac{1}{2} (\bar{x}_1' S_1^{-1} \bar{x}_1 - \bar{x}_2' S_2^{-1} \bar{x}_2) + \bar{x}' (S_1^{-1} \bar{x}_1 - S_2^{-1} \bar{x}_2) - \frac{1}{2} \bar{x}' (S_1^{-1} - S_2^{-1}) \bar{x}$$

وبالتالي، فإن التصنيف للمشاهدة x سيكون

- تعود للمجموعة الأولى إذا كانت  $\hat{G} > 0$
- تعود للمجموعة الثانية إذا كانت  $\hat{G} < 0$

### 3. دالة الإنحدار اللوجستية للتمييز Logistic Regression Discriminant Function

تستخدم في حالة كون توزيع البيانات غير التوزيع الطبيعي. أي عندما يكون توزيع المتغيرات التوضيحية من عائلة التوزيع الأسّي. وتعتمد هذه الطريقة على الاحتمالات السابقة واللاحقة لمجتمع  $(Y_2, Y_1)$  وتكون دالة التمييز التالية:

$$Z = \ln \frac{p(x | Y_1)}{p(x | Y_2)}$$

أما قاعدة التصنيف للمشاهدة  $x$  سيكون:

- تعود للمجموعة الأولى  $Y_1$  إذا كانت  $Z > 0$
- تعود للمجموعة الثانية  $Y_2$  إذا كانت  $Z < 0$

#### بعض الطرق الالاعلمية

تستخدم عندما يكون التوزيع غير طبيعي أو تكون هناك قيم شاذة مثل:

1. طريقة الرتب
2. الطرق الالاعلمية باستخدام التقديرات الحصينة (طريقة HUBER)
3. طريقة الدمج ما بين طريقة الرتب و طريقة HUBER

#### طريقة الرتب

تستخدم بغض النظر عن توزيع البيانات سواء كان طبيعي أو غير طبيعي ويلجأ إليها الباحث عند عدم توفر الفرضيات الخاصة بالدالة المميزة فيتم استخدام تحويلات الرتب للبيانات الأصلية. أي إستبدال البيانات الأصلية برتبها وبعدها يتم تطبيق الطرق السابقة (الخطية أو التريبيعية).

#### 1. طريقة HUBER

يتم إستبدال البيانات الأصلية ببيانات مشذبة. ولغرض تشذيب البيانات يتم حساب المتوسطات ومصفوفة التباين- التباين المشترك حيث يتم إستبدال  $\bar{X}$  و  $S$  بأوزان جديدة تعتمد على إحصاءة مهالونوبيس ( $D_i$ ) وحسب الآتي:

$$\bar{X}^r = \frac{\sum w_i x_i}{\sum w_i}$$
$$S^r = \frac{\sum w_i^2 (x_i - \bar{X}^r)(x_i - \bar{X}^r)'}{\sum w_i^2}$$
$$w_i = \begin{cases} \frac{2}{D_i} & \text{if } D_i > 2 \\ 1 & \text{if } D_i \leq 2 \end{cases}$$



وبعد الحصول على البيانات المشدبة يتم تطبيق الطرق السابقة (الخطية، التربيعية).

## 2. طريقة الدمج ما بين طريقة الرتب وطريقة HUBER

نلجأ إلى إجراء تحويلين على البيانات الأصلية:

أولاً: تشذيب البيانات وجعلها فريية من التوزيع الطبيعي بإستخدام أسلوب HUBER بالطريقة السابقة.

ثانياً: أخذ الرتب بالطريقة السابقة للبيانات المشدبة وبعدها تؤخذ البيانات الجديدة وتعالج على الطرق السابقة (الخطية، التربيعية).

### الجانب التطبيقي

مثال 1:

سُحبت عينة عشوائية مؤلفة من 12 رياضي من مجتمع طبيعي وأجرى عليهم اختباري اللياقة والكفاءة وذلك لغرض تصنيفهم إلى مهرة أو غير مهرة وكانت النتائج كما في الجدول أدناه حيث أن:

$X_1$  تمثل إختبار اللياقة و  $X_2$  تمثل إختبار الكفاءة:

المهرة		غير المهرة	
$X_2$	$X_1$	$X_2$	$X_1$
33	60	35	57
36	61	36	59
35	64	38	59
38	63	39	61
40	65	41	63
		43	65
		41	59

الحل:

يجب أن نتحقق من شرطي الدالة التمييزية الخطية:

1. البيانات تتوزع توزيعاً طبيعياً
2. مصفوفة التباين - التباين المشترك متساوية للمجتمعين.
3. ومن ثم سوف نقوم بإجراء إختبار لإثبات تساوي مصفوفة التباين - التباين المشترك للمجتمعين من خلال الفرضية

$$H_0 : \Sigma_1 = \Sigma_2$$

$$H_1 : \Sigma_1 \neq \Sigma_2$$

وأن إحصاءة الإختبار لهذه الفرضية Q وهي:

$$Q = MC^* \sim \chi^2_{(k-1)(p-1)}$$

حيث أن:

$$M = \sum_{i=1}^2 n_i \ln|S| - \sum_{i=1}^2 n_i \ln|S_i|$$

$$C^* = 1 - \frac{2p^2 + 3p - 1}{6(p+1)(k-1)} \left[ \sum \frac{1}{n_i} - \frac{1}{\sum n_i} \right]$$

$$S = \begin{pmatrix} 7.92 & 5.68 \\ 5.68 & 6.29 \end{pmatrix}$$

$$S_1 = \frac{\begin{pmatrix} S_{11} & S_{12} \\ S_{12} & S_{22} \end{pmatrix}}{n_1 - 1} = \begin{pmatrix} 7.3 & 4.2 \\ 4.2 & 4.3 \end{pmatrix}$$

$$S_2 = \frac{\begin{pmatrix} S_{11} & S_{12} \\ S_{12} & S_{22} \end{pmatrix}}{n_2 - 1} = \begin{pmatrix} 8.33 & 6.66 \\ 6.66 & 7.61 \end{pmatrix}$$

$$|S| = 17.56 \quad , \quad \ln|S| = 2.865$$

$$|S_1| = 13.7 \quad , \quad \ln|S_1| = 2.62$$

$$|S_2| = 19.035 \quad , \quad \ln|S_2| = 2.946$$

$$M = 12(2.865) - [5(2.62) + 7(2.946)]$$

$$= 0.685$$

$$C^* = 1 - \frac{2(2)^2 + 3(2) - 1}{6(2+1)(2-1)} \left[ \left( \frac{1}{3} + \frac{1}{7} \right) - \frac{1}{12} \right]$$

$$= 0.8131$$

$$MC^* = (0.685)(0.8131) = 0.534$$

$$\chi^2_{1,0.05} = 3.841$$

ومن الواضح أن القيمة المحتسبة اقل بكثير من القيمة الجدولية:

$$H_0 : \Sigma_1 = \Sigma_2$$

لذا فالقرار هو عدم رفض (قبول) الفرضية، ومعنى ذلك عدم وجود فرق معنوي بين  $\Sigma_1$  و  $\Sigma_2$ .  
نلاحظ بأن شرطي الدالة التمييزية الخطية متحقق. ولإيجاد دالة التميز الخطية وهي:

$$y = \underline{x}'S^{-1}(\underline{\bar{x}}_1 - \underline{\bar{x}}_2)$$

علينا أن نحسب القيم  $\bar{x}_1$  و  $\bar{x}_2$  و  $S$  لكل مجموعة واحتساب النتائج بموجبها وكما يلي:

### المجموعة الأولى:

$$\sum X_1 = 182 \quad , \quad \sum X_1^2 = 6654$$

$$\sum X_2 = 313 \quad , \quad \sum X_2^2 = 19611$$

$$\sum X_1 X_2 = 11410$$

$$\bar{X}_1 = \begin{pmatrix} 36.4 \\ 62.6 \end{pmatrix}$$

$$S_{11} = \sum x_1^2 - \frac{(\sum x_1)^2}{n} = 6654 - \frac{(182)^2}{5} = 29.2$$

$$S_{22} = \sum x_2^2 - \frac{(\sum x_2)^2}{n} = 19611 - \frac{(313)^2}{5} = 17.2$$

$$S_{12} = \sum x_1 x_2 - \frac{(\sum x_1)(\sum x_2)}{n} = 11410 - \frac{(182)(313)}{5} = 16.8$$

### المجموعة الثانية:

$$\sum X_1 = 273 \quad , \quad \sum X_1^2 = 10697$$

$$\sum X_2 = 423 \quad , \quad \sum X_2^2 = 25607$$

$$\sum X_1 X_2 = 16537$$

$$\bar{X}_2 = \begin{pmatrix} 39 \\ 60.42 \end{pmatrix}$$

$$S_{11} = \sum x_1^2 - \frac{(\sum x_1)^2}{n} = 10697 - \frac{(273)^2}{7} = 50$$

$$S_{22} = \sum x_2^2 - \frac{(\sum x_2)^2}{n} = 25607 - \frac{(423)^2}{7} = 45.71$$

$$S_{12} = \sum x_1 x_2 - \frac{(\sum x_1)(\sum x_2)}{n} = 16537 - \frac{(273)(423)}{7} = 40$$

ومن هذه النتائج نستخرج القيم التالية والتي تمثل التباينات المدمجة:

$$V_{11} = \frac{S_{11}(1) + S_{11}(2)}{n_1 + n_2 - 2} = \frac{29.2 + 50}{5 + 7 - 2} = 7.92$$

$$V_{22} = \frac{S_{22}(1) + S_{22}(2)}{n_1 + n_2 - 2} = \frac{17.2 + 45.71}{5 + 7 - 2} = 6.291$$

$$V_{12} = \frac{S_{12}(1) + S_{12}(2)}{n_1 + n_2 - 2} = \frac{16.8 + 40}{5 + 7 - 2} = 5.68$$

وبالتالي فإن مصفوفة التباين الكلي تكون:

$$|S| = 17.5623$$

$$S^{-1} = \frac{adj(S)}{|S|} = \frac{\begin{pmatrix} 6.291 & -5.68 \\ -5.68 & 7.92 \end{pmatrix}}{17.5623} = \begin{pmatrix} 0.358 & -0.323 \\ -0.323 & 0.451 \end{pmatrix}$$

**الدالة المميزة**

$$y = \underline{x}' S^{-1} (\underline{\bar{x}}_1 - \underline{\bar{x}}_2) = \underline{x}' C^*$$

$$C^* = S^{-1} (\underline{\bar{x}}_1 - \underline{\bar{x}}_2)$$

$$= \begin{pmatrix} 0.358 & -0.323 \\ -0.323 & 0.451 \end{pmatrix} \begin{pmatrix} -2.6 \\ 2.18 \end{pmatrix} = \begin{pmatrix} -1.634 \\ 1.822 \end{pmatrix} = \begin{pmatrix} C_1 \\ C_2 \end{pmatrix}$$

$$y = (x_1 \quad x_2) \begin{pmatrix} -1.634 \\ 1.822 \end{pmatrix}$$

$$y = -1.634x_1 + 1.822x_2$$

**نقطة الفصل**

$$\bar{y}_1 = \underline{\bar{x}}_1' S^{-1} (\underline{\bar{x}}_1 - \underline{\bar{x}}_2) = (36.4 \quad 62.6) \begin{pmatrix} -1.634 \\ 1.822 \end{pmatrix} = 54.5796$$

$$\bar{y}_2 = \underline{\bar{x}}_2' S^{-1} (\underline{\bar{x}}_1 - \underline{\bar{x}}_2) = (39 \quad 60.42) \begin{pmatrix} -1.634 \\ 1.822 \end{pmatrix} = 46.359$$

$$z = \frac{\bar{y}_1 + \bar{y}_2}{2} = \frac{54.5796 + 46.359}{2} = 50.469$$

**قاعدة التصنيف**

المشاهدة  $x$  تعود للمجتمع الأول إذا كانت  $y - z > 0$

المشاهدة  $x$  تعود للمجتمع الثاني إذا كانت  $y - z \leq 0$

فلو افترضنا المشاهدة  $\underline{X} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$  لنرى لأي مجتمع تعود، وهنا يجب أن نجد قيمة الدالة

المميزة:

$$y = -1.634x_1 + 1.822x_2$$

$$= -1.634(1) + 1.822(1) = 0.188$$

$$y - z = 0.188 - 50.469 = -50.281 < 0$$

إذا هذه المشاهدة يتم تصنيفها بأنها تعود للمجتمع الثاني.

### أهمية كل متغير

بالإمكان إظهار أهمية كل متغير من خلال تطبيق المقياس الآتي:

$$C_i^* = C_i \sqrt{V_{ii}}$$

$$C_1^* = C_1 \sqrt{V_{11}} = -1.634 \sqrt{7.92} = -4.598$$

$$C_2^* = C_2 \sqrt{V_{22}} = 1.822 \sqrt{6.291} = 4.569$$

ونرى بأن المتغيرين لهما نفس الأهمية إذ الفرق قليل جداً.

### إيجاد خطأ التصنيف

$$\begin{aligned} D^2 &= (\bar{x}_1 - \bar{x}_2)' S^{-1} (\bar{x}_1 - \bar{x}_2) \\ &= (-2.6 \quad 2.18) \begin{pmatrix} 0.358 & -0.323 \\ -0.323 & 0.451 \end{pmatrix} \begin{pmatrix} -2.6 \\ 2.18 \end{pmatrix} = 8.22 \\ &2.867 = D \end{aligned}$$

وبذلك فإن احتمال خطأ التصنيف أعلاه هو صغير جداً ويعطي إنطباعاً عن كفاءة عالية لدالة التمييز.

وللوقوف على مستوى هذه الكفاءة من جانب إستنتاجي إحصائي، نقوم بإختبار الفرضية التالية:

$$H_0 : \underline{\mu}_1 = \underline{\mu}_2$$

$$H_1 : \underline{\mu}_1 \neq \underline{\mu}_2$$

ويتم الإختبا بإستخدام إختبار F الذي يعتمد على إحصاءة هوتلنك  $T^2$  Hotteling والتي تكون:

$$T^2 = \frac{n_1 n_2}{n_1 + n_2} D^2 = \frac{(5)(7)}{5+7} (8.22) = 23.975$$

وبالتالي، فإن إختبار F سيكون:

$$\begin{aligned} F &= \frac{n_1 + n_2 - p - 1}{n_1 + n_2 - 2} T^2 \\ &= \frac{5+7-2-1}{5+7-2} (23.975) = 21.577 \end{aligned}$$

وبالمقارنة مع القيمة الجدولية  $f_{2,9,0.05} = 4.26$  فإننا نرفض  $H_0$  مستنتجين فرق معنوي كبير بين الوسط الحسابي للمجتمعين مما يتيح القول بمستوى كفاءة عالي لدالة التمييز أعلاه.

## مثال 2:

سُحبت عينة عشوائية من 30 طالب لإجراء دراسة باستخدام التحليل المميز للتمييز بين الطلاب للصفوف الثلاث (المجاميع)  $C_1, C_2, C_3$  للمرحلة الرابعة لقسم الأحصاء و ذلك بالإستناد على أربعة إختبارات وهي إختبارات متعدد  $(X_1)$ ، الإستدلال  $(X_2)$ ، العمليات  $(X_3)$ ، التصميم  $(X_4)$ . وكانت النتائج كما في الجدول أدناه والتي تمثل الأوساط الحسابية:

متغير الإختبار	$C_1$	$C_2$	$C_3$
$X_1$	64.5	60.5	58
$X_2$	86.4	81.8	83.1
$X_3$	75.2	73.4	71.4
$X_4$	81.9	88	86.6
حجم العينة	8	12	10

وأن مصفوفة التباين - التباين المشترك والمقدرة من العينة كانت كالآتي:

$$S = \begin{bmatrix} 1.00 & 0.5849 & 0.1774 & 0.1974 \\ 0.5849 & 1.00 & 0.2094 & 0.2170 \\ 0.1774 & 0.2094 & 1.00 & 0.2910 \\ 0.1974 & 0.2170 & 0.2910 & 1.00 \end{bmatrix}$$

## الدالة المميزة

لدينا ثلاثة مجاميع لذلك توجد لدينا ثلاث دوال مميزة وهي  $y_{12}, y_{13}, y_{23}$  وكما يلي:

$$y_{12} = \underline{x}' S^{-1} (\underline{\bar{x}}_1 - \underline{\bar{x}}_2) = \underline{x}' C_1^*$$

$$C_1^* = S^{-1} (\underline{\bar{x}}_1 - \underline{\bar{x}}_2)$$

$$|S| = 0.55406$$

$$S^{-1} = \frac{adj(S)}{|S|}$$

$$= \frac{1}{0.55406} \begin{pmatrix} 0.8508 & -0.4786 & -0.03503 & -0.0539 \\ -0.4786 & 0.86527 & -0.07559 & -0.0713 \\ -0.03503 & -0.07559 & 0.62194 & -0.15768 \\ -0.0539 & -0.0713 & -0.15768 & 0.62603 \end{pmatrix}$$

$$= \begin{pmatrix} 1.53557 & -0.8638 & -0.06322 & -0.09728 \\ -0.8638 & 1.56169 & -0.13642 & -0.128686 \\ -0.06322 & -0.13642 & 1.12251 & -0.2845 \\ -0.09728 & -0.128686 & -0.2845 & 1.12989 \end{pmatrix}$$

$$\begin{aligned} C_1^* &= S^{-1}(\bar{x}_1 - \bar{x}_2) \\ &= \begin{pmatrix} 1.53557 & -0.8638 & -0.06322 & -0.09728 \\ -0.8638 & 1.56169 & -0.13642 & -0.128686 \\ -0.06322 & -0.13642 & 1.12251 & -0.2845 \\ -0.09728 & -0.128686 & -0.2845 & 1.12989 \end{pmatrix} \begin{pmatrix} 4 \\ 4.6 \\ 1.8 \\ -6.1 \end{pmatrix} \\ &= \begin{pmatrix} 2.6484 \\ 4.26755 \\ 2.8755 \\ -8.3855 \end{pmatrix} \end{aligned}$$

$$\begin{aligned} y_{12} &= (x_1 \quad x_2 \quad x_3 \quad x_4) \begin{pmatrix} 2.6484 \\ 4.26755 \\ 2.8755 \\ -8.3855 \end{pmatrix} \\ &= 2.6484 x_1 + 4.26755 x_2 + 2.8755 x_3 - 8.3855 x_4 \end{aligned}$$

$$y_{13} = \underline{x}' S^{-1}(\bar{x}_1 - \bar{x}_3) = \underline{x}' C_2^*$$

$$\begin{aligned} C_2^* &= S^{-1}(\bar{x}_1 - \bar{x}_3) \\ &= \begin{pmatrix} 1.53557 & -0.8638 & -0.06322 & -0.09728 \\ -0.8638 & 1.56169 & -0.13642 & -0.128686 \\ -0.06322 & -0.13642 & 1.12251 & -0.2845 \\ -0.09728 & -0.128686 & -0.2845 & 1.12989 \end{pmatrix} \begin{pmatrix} 6.5 \\ 3.3 \\ 3.8 \\ -4.7 \end{pmatrix} \\ &= \begin{pmatrix} 7.3476 \\ -0.3746 \\ 4.7415 \\ -7.4485 \end{pmatrix} \end{aligned}$$

$$y_{13} = (x_1 \quad x_2 \quad x_3 \quad x_4) \begin{pmatrix} 7.3476 \\ -0.3746 \\ 4.7415 \\ -7.4485 \end{pmatrix}$$

$$= 7.3476 x_1 - 0.3746 x_2 + 4.7415 x_3 - 7.4485 x_4$$

$$y_{23} = \underline{x}' S^{-1} (\bar{x}_2 - \bar{x}_3) = \underline{x}' C_3^*$$

$$C_3^* = S^{-1} (\bar{x}_2 - \bar{x}_3)$$

$$= \begin{pmatrix} 1.53557 & -0.8638 & -0.06322 & -0.09728 \\ -0.8638 & 1.56169 & -0.13642 & -0.128686 \\ -0.06322 & -0.13642 & 1.12251 & -0.2845 \\ -0.09728 & -0.128686 & -0.2845 & 1.12989 \end{pmatrix} \begin{pmatrix} 2.5 \\ -1.3 \\ 2 \\ 1.4 \end{pmatrix}$$

$$= \begin{pmatrix} 4.699 \\ -4.462 \\ 1.866 \\ 0.9369 \end{pmatrix}$$

$$y_{23} = (x_1 \quad x_2 \quad x_3 \quad x_4) \begin{pmatrix} 4.699 \\ -4.462 \\ 1.866 \\ 0.9369 \end{pmatrix}$$

$$= 4.699 x_1 - 4.462 x_2 + 1.866 x_3 + 0.9369 x_4$$

## نقطة الفصل

نقطة التمييز بين المجموعة الأولى والثانية هي:

$$Z_{12} = \frac{\bar{y}_1 + \bar{y}_2}{2}$$

$$\bar{y}_1 = \underline{\bar{x}}_1' S^{-1} (\bar{x}_1 - \bar{x}_2) = (64.5 \quad 86.4 \quad 75.2 \quad 81.9) \begin{pmatrix} 2.6484 \\ 4.2675 \\ 2.8755 \\ -8.3855 \end{pmatrix} = 68.99$$

$$\bar{y}_2 = \underline{\bar{x}}_2' S^{-1} (\bar{x}_1 - \bar{x}_2) = (60.5 \quad 81.8 \quad 73.4 \quad 88) \begin{pmatrix} 2.6484 \\ 4.2675 \\ 2.8755 \\ -8.3855 \end{pmatrix} = -17.55$$

$$Z_{12} = \frac{68.99 - 17.55}{2} = 25.72$$



نقطة التمييز بين المجموعة الأولى والمجموعة الثالثة:

$$Z_{13} = \frac{\bar{y}_1 + \bar{y}_3}{2}$$

$$\bar{y}_1 = \underline{\bar{x}}_1' S^{-1} (\underline{\bar{x}}_1 - \underline{\bar{x}}_3)$$

$$= (64.5 \quad 86.4 \quad 75.2 \quad 81.9) \begin{pmatrix} 7.3476 \\ -0.3746 \\ 4.7415 \\ -7.4485 \end{pmatrix} = 188.08$$

$$\bar{y}_3 = \underline{\bar{x}}_3' S^{-1} (\underline{\bar{x}}_1 - \underline{\bar{x}}_3) = (58 \quad 83.1 \quad 71.4 \quad 86.6) \begin{pmatrix} 7.3476 \\ -0.3746 \\ 4.7415 \\ -7.4485 \end{pmatrix} = 88.53$$

$$Z_{13} = \frac{188.08 + 88.53}{2} = 138.305$$

نقطة التمييز بين المجموعة الثانية والثالثة:

$$Z_{23} = \frac{\bar{y}_2 + \bar{y}_3}{2}$$

$$\bar{y}_2 = \underline{\bar{x}}_2' S^{-1} (\underline{\bar{x}}_2 - \underline{\bar{x}}_3)$$

$$= (60.5 \quad 81.8 \quad 73.4 \quad 88) \begin{pmatrix} 4.699 \\ -4.462 \\ 1.866 \\ 0.9369 \end{pmatrix} = (58 \quad 83.1 \quad 71.4 \quad 86.6)$$

$$\begin{pmatrix} 4.699 \\ -4.462 \\ 1.866 \\ 0.9369 \end{pmatrix} = 116.11$$

$$Z_{23} = \frac{138.70 + 116.11}{2} = 127.405$$

### قاعدة التصنيف

لكي نجد كل مشاهدة تعود إلى أي مجتمع نجد:

$$W_{12} = y_{12} - Z_{12}$$

$$W_{13} = y_{13} - Z_{13}$$

$$W_{23} = y_{23} - Z_{23}$$

تقع ضمن المجموعة الأولى إذا كان:

$$W_{13} > 0 \quad , \quad W_{12} > 0$$

تقع ضمن المجموعة الثانية إذا كان:

$$W_{13} > W_{12} \quad , \quad W_{13} > 0$$

تقع ضمن المجموعة الثالثة إذا كان:

$$W_{12} > W_{13} \quad , \quad W_{13} > 0$$

وفي ضوء ذلك، نستطيع عمل جدول تصنيفي للمشاهدات لنقف على نسبة التصنيفات الصحيحة (أي تصنيف المشاهدة لمجموعتها الأصلية) والتصنيفات الخاطئة (أي تصنيف المشاهدة ضمن أي من المجاميع الأخرى غير مجموعتها الأصلية).

## التحليل العنقودي

### Cluster Analysis (CA)

التحليل العنقودي مشابه للتحليل المميز من ناحية كونه يستخدم لتصنيف مجموعة من الأفراد أو الوحدات التجريبية إلى مجاميع فرعية معرفة بشكل محدد وبلا تقاطع. الفرق بينهما هو أن التحليل المميز يمكن استخدامه عندما يكون لدى الباحث عينات عشوائية وكل واحدة منها مسحوبة مسبقاً من إحدى المجاميع الفرعية المحددة. بينما التحليل العنقودي يتعامل مع مسألة التصنيف عندما تكون الوحدات التجريبية غير معروفة مسبقاً لأي مجموعة فرعية تنتمي في الأصل.

ولكي نفهم بداية فكرة التحليل العنقودي، لنفترض أن صاحب شركة تجارية لبيع البضائع والمستلزمات يمتلك بيانات ما تم جمعها من المستهلكين الذين يبتاعون من شركته. وهذه البيانات قد تتضمن متغيرات عديدة عن المستهلك مثل: العمر، المستوى التعليمي، مستوى الدخل، الحالة الزوجية، الحالة الوظيفية، عدد الأطفال دون سن الخامسة، وعدد هم ما بين - 13 6، وعدد من هم 14 سنة فأكثر.

وصاحب الشركة هذه ربما يرغب في تصنيف هؤلاء المستهلكين إلى مجاميع متباينة مستخدماً بياناته هذه وذلك لغرض إعلان سياسة العرض لبضائع معينة اعتماداً على طبيعة هذه المجاميع والتي تسمى (العناقيد Clusters).

هذه المجاميع بطبيعة الحال، توحى بتجانس ما بين المنتمين لكل مجموعة مع فروقات واضحة ما بين المنتمين لأي مجموعتين مختلفتين. وبشكل عام، فإن الأسلوب العنقودي من شأنه تصنيف وحدات العينة إلى مجاميع (عناقيد Clusters) غير معروفة مسبقاً.

ومن الجدير بالذكر هنا أنه يجب عدم الخلط ما بين التحليل العنقودي والتحليل التمييزي. فالعنقودي لا يعتمد على أية معطيات تصنيفية مسبقة، بينما التمييزي يعتمد على كون عدد المجاميع ونوعيتها معرفة مسبقاً ويتم جمع البيانات بموجبها ثم بعد ذلك تكون المهمة أن نقف على صحة عائدة المشاهدة لمجموعتها من عدمه. وفي حالة العدم، نقف على لأي مجموعة هي أقرب.

وبالتالي، فإن أسلوب التحليل العنقودي هدفه تجميع (تصنيف) المشاهدات وفق مجاميع Clusters بحيث كل مجموعة تحتوي على مشاهدات متجانسة قدر الإمكان بالنسبة لمتغيرات التعنقد المستخدمة. ولكي نبدأ تطبيق هذا الأسلوب، علينا إتباع الخطوات التالية:

1) تحديد مقياس التشابه Measure of Similarity

2) تحديد نوع أسلوب العنقدة المستخدم

3) تقرير عدد العناقيد

4) تفسير النتائج

والمهم هنا في مقياس التشابه هو الحصول على نمط بناء مجاميع سهلة لبيانات واسعة بإعتماد مبدأ التقارب أو التشابه. ولتحديد هذا المقياس، نجد أن الإختيار يعتمد على نوعية مقياس المتغيرات (إسمي، رتبي، فئوي، نسبي) أو طبيعتها (متقطعة، مستمرة، ثنائية). ومن الممكن استخدام أحد القاييس التالية:

- 1- قياسات المسافة Distance Measures.
  - 2- معامل الترابط (العلاقة) Association Coefficient
- وفيما يلي توضيح ذلك.

### مقياس المسافة Distance Measures

وهذا هو ما نعبر عنه بمسافة مهالانوبس Mahalanobis Distance ما بين أية مشاهدين (نقطتين)  $X_r$  ,  $X_s$  والتي تكون:

$$d(X_r, X_s) = d_{rs} = \sqrt{(X_r - X_s)' \Sigma^{-1} (X_r - X_s)}$$

حيث أن  $\Sigma$  تمثل مصفوفة التباين والتباين المشترك ويمكن إحلال تقدير مناسب لها.

#### ملاحظة:

في الغالب نستخدم مقياس المسافة لهذا الغرض Mahalanobis Distance وقد نطلق عليه أيضاً المسافة الإقليدية Euclidean Distance وبالتالي فإن المقياس للمتغيرات يجب أن يكون كمياً (نسبي أو فئوي). وفي أدناه مثلاً على ذلك.

#### مثال:

لنفترض المثال البسيط التالي لغرض توضيح كيفية استخدام أسلوب مقياس المسافة من خلال عينة من 6 أشخاص والبيانات التي تم جمعها عنهم تتمثل بمتغيرين هما:

$X_1$  : مقدار الدخل السنوي بالدولار

$X_2$  : مستوى التعليم بالسنوات

جدول البيانات الأولية

الشخص (المشاهدة)	مقدار الدخل (ألف دولار)	مستوى التعليم (سنة)
1	5	5
2	6	6
3	15	14
4	16	15
5	25	20
6	30	19

وبإستخدام مقياس المسافة Euclidean Distance بالنسبة لهذه البيانات، نحصل على مصفوفة التشابه (التمائل) Similarity Matrix التالية :

المشاهدة	1	2	3	4	5	6
1	0.0	2	181	221	625	821
2	2	0.0	145	181	557	745
3	181	145	0.0	2	136	250
4	221	181	2	0.0	106	212
5	625	557	136	106	0.0	26
6	821	745	250	212	26	0.0

والسؤال هو كيف بالإمكان إستخدام هذه البيانات الأولية ومقاييس المسافة في تحديد العناقيد؟. هنالك طريقتين رئيسيتين لعمل التحليل العنقودي وهما:

- 1- الطريقة الهرمية Hierarchical clustering
- 2- الطريقة اللاهرمية Non-hierarchical clustering

وفيما يلي سنتطرق للطريقة الهرمية لكونها المفضلة عموماً.

### الطريقة الهرمية Hierarchical clustering

من خلال ملاحظة جدول البيانات الأولية نجد أن المشاهدين الأولى والثانية متقاربتين مثلما هما المشاهدين الثالثة والرابعة. في البداية بالطبع لدينا عناقيد بعدد المشاهدات ويمكن أن نبدأ بأي زوج منهما بداية الأمر ولنفتراض أننا نختار المشاهدين الأولى والثانية كخطوة أولى. وبذلك يتم الدمج لهاتين المشاهدين في عنقود واحد وبالتالي يصبح لدينا خمسة عناقيد وفقاً للبيانات الأولية والمسافات.

الخطوة التالية هو تكوين مصفوفة جديدة من مقياس المسافات وفقاً لهذه العناقيد الخمسة.

ولأن العنقود رقم (1) المتشكل من المشاهدين الأولى والثانية يتضمن مشاهدين، فإنه يتوجب علينا إستخدام بعض القواعد لتحديد المسافات ما بين العناقيد الجديدة في مثل هذه الحالة.

وهناك عدد من القواعد شائعة الإستخدام ومنها:

- 1- الطريقة المركزية Centroid method
- 2- طريقة الربط الفردي Single Linkage method
- 3- طريقة الربط المتكامل Complete Linkage method
- 4- طريقة الربط المتوسطي Average Linkage method
- 5- طريقة وارد Ward's method

وفيما يلي عرض بسيط لأول طريقتين.

### الطريقة المركزية

وفقاً لهذه الطريقة، فإن كل مجموعة (عقود) Cluster تم استبدالها بمعدل المشاهدات لتلك المجموعة عوضاً عن القيم الأصلية. وعلى سبيل المثال، عندما ندمج المشاهدين الأولى والثانية لتكوين العقود الأول والذي يكون مركزاً وسطياً بينهما. بمعنى إننا نستخدم معدل قيم هاتين المشاهدين. وهذا يعني أن العقود الأول لديه معدل مستوى التعليم ما يساوي  $(5 + 6)/2 = 5.5$  yr. ونستمر هكذا لحين الوصول إلى العقود الأخير والذي يمكن أن يضم جميع المشاهدات.

### طريقة الربط الفردي

في هذه الطريقة، فإن المسافة ما بين عقودين يتم تمثيلها بأقل مسافة ما بين جميع الأزواج المحتملة للمشاهدات في كلا العقودين. وأحد الأمثلة على هذه الطريقة هو أسلوب /طريقة أقرب الجوار Nearest neighbor method والتي تتضمن الخطوات التالية:

- أ- إبدأ مع N من العناقيد حيث كل عقود يتضمن مشاهدة واحدة.
  - ب- إدمج أقرب نقطتين وفقاً لإحدى طرق مقاييس المسافة المعتمد.
  - ت- إعتد التباعد dissimilarity ما بين هذا العقود الجديد وأي نقطة أخرى بمثابة أصغر مسافة بين هاتين النقطتين في العقود وهذه النقطة الأخرى.
  - ث- الإستمرارية بدمج العناقيد الأكثر قرباً لبعضهما، وبالتالي سينخفض عدد العناقيد واحداً مع كل خطوة. إن التباعد ما بين أي عقودين هو دائماً المسافة ما بين أقرب عقودين.
- ولذلك فإن طريقة "أقرب الجوار" تبدأ أيضاً مع N من العناقيد حيث كل عقود يتضمن مشاهدة واحدة، وتستمر بدمج النقاط والعناقيد حتى تنتهي العملية بعقود واحد يتضمن جميع المشاهدات.

وإزاء ذلك، فإنه من الواضح أن العدد المناسب من العناقيد التي ننتهي عندها يقع ما بين عددها عند البداية وعددها عند النهاية. وهناك بعض الطرق لتحديد مكان التوقف لعملية الدمج هذه بضمنها النظرة المنطقية في هذا الشأن.

إحدى هذه الطرق التي تساعد في وقف عملية الدمج هي من خلال بناء شكل الشجرة الهرمي. وهذا الشكل يتضمن فروعاً تربط نقاط البيانات وتعكس بالترتيب الذي تأخذه النقاط لضمها للعناقيد. وطول فروعها يجب أن تتناسب مع المسافات ما بين النقاط والعناقيد عندما يتم دمج النقاط بالعناقيد. وسيتم توضيح هذه الآلية من خلال المثال التالي.

مثال:

لنفترض مصفوفة المسافات الإقليدية (مصفوفة التباعد) ما بين 6 من المشاهدات. وأعطينا هذا المثال بهذه القيم لتسهيل فهم عملية الإجراءات المتخذة وفقاً لما تم تبيانه في أعلاه:

العنقود (المشاهدة)	1	2	3	4	5	6
1	0.0	0.31	0.23	0.32	0.26	0.25
2		0.0	0.34	0.21	0.36	0.28
3			0.0	0.31	(0.04)	0.07
4				0.0	0.31	0.28
5					0.0	0.09
6						0.0

وقبل البدء برسم الشكل الشجري الهرمي، دعنا نتعامل في عنقدة هذه البيانات وفقاً لطريقة الربط الفردي وحسب الخطوات التالية:

في البداية، وكما قلنا، نعتبر عدد العناقيد مساوياً لعدد المشاهدات (النقاط). فإذا افترضنا الرمز  $C$  لمجموعة العناقيد، ففي هذه الحالة نبدأ بالمجموعة  $C_0$  وتكون الآتي:

$$C_0 = \{ [1], [2], [3], [4], [5], [6] \}$$

أي أننا لدينا الآن 6 عناقيد. وفي ضوء ذلك تكون خطوات العمل التالية كما يلي:

1- بملاحظة مصفوفة المسافات أعلاه، نجد أن النقطتين الأقرب لبعضهما هما (3) و (5) حيث أن المسافة هي (0.04) وبذلك فإن الخطوة الأولى هي دمج هاتين النقطتين في عنقود واحد وبذلك تصبح لدينا مجموعة العناقيد  $C_1$  وهي

$$C_1 = \{ [1], [2], [3, 5], [4], [6] \}$$

وبذلك يترتب علينا صياغة مصفوفة مسافات جديدة وفقاً إلى  $C_1$  وهي الآتي:

العنقود	[1]	[2]	[3, 5]	[4]	[6]
[1]	0.0	0.31	0.23	0.32	0.25
[2]		0.0	0.34	0.21	0.28
[3, 5]			0.0	0.31	(0.07)
[4]				0.0	0.28
[6]					0.0

## ملاحظة:

إن المصفوفة الجديدة هي عبارة عن المصفوفة السابقة بعد حذف العمود والصف للعنقود المدمج مع سابقه. ولأننا قمنا بدمج العنقود [5] مع العنقود [3] فإن المصفوفة أعلاه عبارة عن المصفوفة الأصلية محذوفاً منها عمود وصف [5].

وهذه المسافات محسوبة على أساس مقارنة النقطة [1] مع [5, 3] ونختار الأصغر من بين 0.23 و 0.26 فوضعنا 0.23 الأصغر. وعلى نفس المنوال، تم تحديد المسافات الجديدة ما بين بقية النقاط (العناقيد) جميعاً لنرى الوضع أعلاه.

2- ومن ملاحظة ذلك، نجد أن أقرب مسافة هي ما بين [6] و [5, 3] وهي (0.07) وبذلك يتم دمج هذين العنقودين لتكون لدينا المجموعة  $C_2$  التالية:

$$C_2 = \{ [1], [2], [3, 5, 6], [4] \}$$

وفي ضوء ذلك يترتب علينا صياغة مصفوفة مسافات جديدة والتي يمكن تحديدها من خلال الخطوة السابقة وهذا هو الجانب الإيجابي في طريقة الربط المنفرد والتي تغنينا عن الرجوع للمصفوفة الأولية. وعلى هذا الأساس، وفي ضوء  $C_2$  ستكون المصفوفة الجديدة للمسافات كما يلي:

العنقود	[1]	[2]	[3, 5, 6]	[4]
[1]		0.31	0.23	0.32
[2]			0.28	(0.21)
[3, 5, 6]				0.28
[4]				

3- وهذه المصفوفة تنتج لنا تقارب [2] مع [4] حيث المسافة الأقل وهي (0.21) وبذلك تكون مجموعة العناقيد الجديدة  $C_3$  بالشكل التالي:

$$C_3 = \{ [1], [2, 4], [3, 5, 6] \}$$

وفي ضوء  $C_3$  هذه ستكون المصفوفة الجديدة للمسافات كما يلي:

العنقود	[1]	[2, 4]	[3, 5, 6]
[1]	0.0	0.31	(0.23)
[2, 4]		0.0	0.28
[3, 5, 6]			0.0



4- وهذه المصفوفة تشير إلى أقرب مسافة ما بين [1] و [3, 5, 6] وهي (0.23). ولذلك سيتم الدمج بينهما لتصبح مجموعة العناقيد الجديدة  $C_4$  بالشكل التالي:

$$C_4 = \{ [3, 5, 6, 1], [2, 4] \}$$

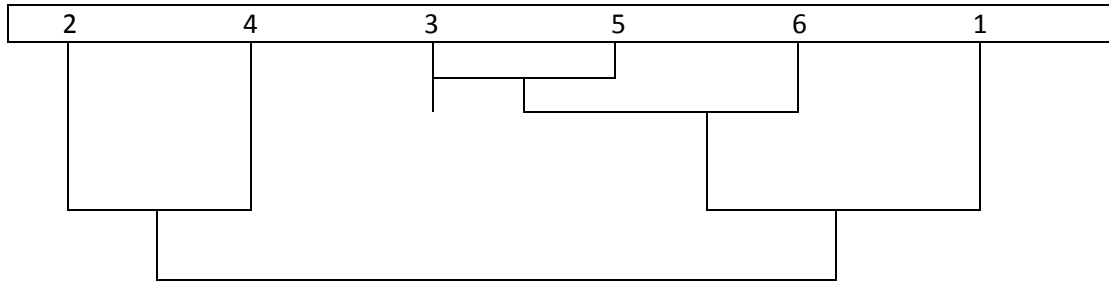
وأن مصفوفة المسافات في ضوء  $C_4$  ستكون:

العنقود	[3, 5, 6, 1]	[2, 4]
[3, 5, 6, 1]	0.0	0.28
[2, 4]		0.0

5- إن عملية الدمج الأخيرة هي بالتالي ما بين هاتين المجموعتين فتكون لدينا مجموعة (عنقود) واحد وهي:

$$C_5 = ( [1, 2, 3, 4, 5, 6] )$$

وفي ضوء هذه النتائج يمكننا تنفيذ الشكل البياني للشجرة الهرمية وعلى النحو المبين في أدناه:

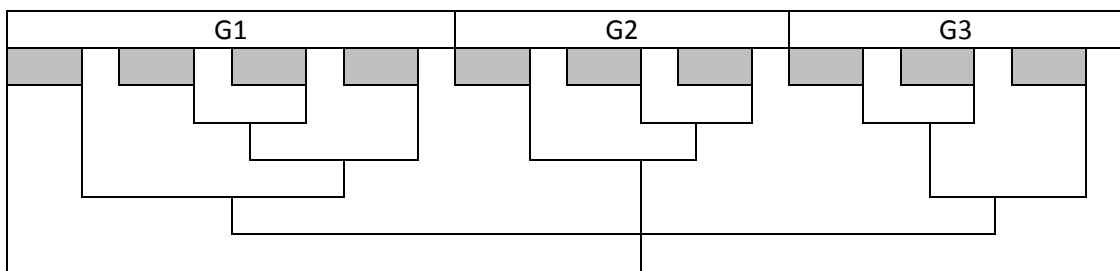


وهذا الشكل يساعد في إتخاذ قرار التوقف حيث قد نرى التوقف عند العنقودين

$$[2, 4] \text{ و } [3, 5, 6, 1]$$

$$\text{أو ما قبل ذلك وعند العناقيد الثلاثة } [1] \text{ و } [2, 4] \text{ و } [3, 5, 6]$$

ولكي يكون الموضوع أكثر وضوحاً حول إختيار توقف العنقدة، لنفترض أن حصيلتنا بالشكل التالي:



وهذا يمثل الشكل البياني المثالي الوضوح للشجرة الهرمية. وفي ضوء ذلك فقد نختار الإبقاء على ثلاثة عناقيد كإختيار معقول وهي  $G_1$  و  $G_2$  و  $G_3$

وبالعودة لمثالنا الأول حيث البداية بالمصفوفة:

المشاهدة	1	2	3	4	5	6
1	0.0	(2)	181	221	625	821
2		0.0	145	181	557	745
3			0.0	2	136	250
4				0.0	106	212
5					0.0	26
6						0.0

وعند تطبيق طريقة الربط الفردي، تكون لدينا النتائج التالية:

1- الخطوة الأولى دمج [1] مع [2] فتكون المجاميع  
 $C_1 = \{ [1, 2], [3], [4], [5], [6] \}$

2- وباعتماد المصفوفة الجديدة للمساقات وهي:

	[1, 2]	[3]	[4]	[5]	[6]
العنقود					
[1, 2]	0.0	181	221	625	821
[3]		0.0	(2)	136	250
[4]			0.0	106	212
[5]				0.0	26
[6]					0.0

ومنها نجد إدماج [3] و [4] لتصبح لدينا المجموعة:

$C_2 = \{ [1, 2], [3, 4], [5], [6] \}$

3- وباعتماد ذلك، تكون لدينا المصفوفة الجديدة:

العنقود	[1, 2]	[3, 4]	[5]	[6]
[1, 2]	0.0	181	625	821
[3, 4]		0.0	136	250
[5]			0.0	(26)
[6]				0.0

4- ومنها نجد إدماج [5] مع [6] لتصبح لدينا المجموعة:

$$C_3 = \{ [1, 2] , [3, 4] , [5,6] \}$$

وبذلك تنتج لدينا مصفوفة مسافات جديدة وهي:

العنقود	[1, 2]	[3, 4]	[5, 6]
[1, 2]	0.0	181	625
[3, 4]		0.0	(136)
[5, 6]			0.0

ومنها نجد أن [5, 6] أقرب إلى [3, 4] فيتم الدمج بينهما لتصبح لدينا المجموعة:

$$C_4 = \{ [1, 2] , [3, 4, 5,6] \}$$

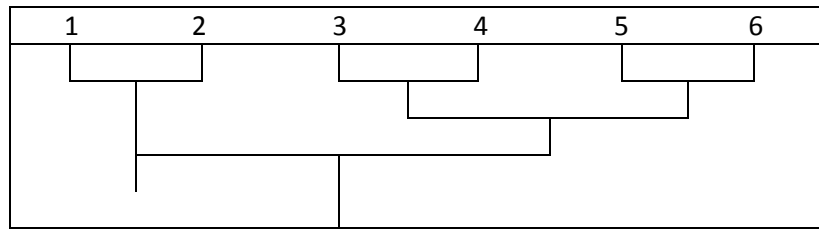
والنتيجة النهائية لمصفوفة المسافات تكون :

العنقود	[1, 2]	[3, 4, 5, 6]
[1, 2]	0.0	181
[3, 4, 5, 6]		0.0

5- وأخيراً ليس لدينا سوى دمج هاتين المجموعتين (العنقودين) لتصبح لدينا المجموعة الأخيرة وبواقع عنقود واحد وهو:

$$C_4 = \{ [1, 2, 3, 4, 5, 6] \}$$

والشكل الشجري الهرمي لهذه النتيجة يكون:



ونرى أن نتوقف عند العناقيد:

$$C_3 = \{ [1, 2] , [3, 4] , [5,6] \}$$

وفيما يلي مثالاً لبيانات حقيقية لواقع الجرائم المختلفة (القتل، الإغتصاب، السرقة،.....، سرقة السيارات) المسجلة في 6 مدن أمريكية (أطلنطا، بوسطن، شيكاغو، دالاس، دنفر، ديترويت) محسوبة لكل 100,000 من السكان حيث تم إحتساب مصفوفة المسافات ما بين هذه المدن وفقاً لبيانات الجرائم فكانت حسب الآتي<sup>(3)</sup>:

العنقود (المدن)	1	2	3	4	5	6
أطلنطا	0	536	516	590	693	716
بوسطن	536	0	447	833	915	881
شيكاغو	516	447	0	924	1073	971
دالاس	590	833	924	0	527	464
دنفر	693	915	1073	527	0	(358)
ديترويت	716	881	971	464	(358)	0

أصغر قيمة للمسافات محصورة بين قوسين ( )

ولكون أصغر مسافة في الجدول أعلاه هي (358) وهي المسافة ما بين المدينتين دنفر (رقم 5) و دترويت (رقم 6)، فإننا ندمج هاتين المدينتين بإعتبارهما يشكلان عنقوداً واحداً من حيث سجل الجرائم المرتكبة. وبعد عملية الدمج بين هاتين المدينتين وما تبع ذلك من الحذف والتعويض، أصبح لدينا الجدول التالي:

العنقود (المدن)	1	2	3	4	[ 5 , 6 ]
1	0.0	536	516	590	693
2	536	0.0	(447)	833	881
3	516	(447)	0.0	924	971
4	590	833	924	0.0	464
[ 5 , 6 ]	693	881	971	464	0.0

ولكون أصغر مسافة في الجدول أعلاه هي (447) وهي المسافة ما بين المدينتين بوسطن (رقم 2) و شيكاغو (رقم 3)، فإننا ندمج هاتين المدينتين بإعتبارهما يشكلان عنقوداً واحداً من حيث سجل الجرائم المرتكبة. وبعد عملية الدمج بين هاتين المدينتين وما تبع ذلك من الحذف والتعويض، أصبح لدينا الجدول التالي:

العنقود (المدن)	1	[ 2 , 3 ]	4	[ 5 , 6 ]
1	0.0	516	590	693
[ 2 , 3 ]	516	0.0	833	881
4	590	833	0.0	(464)
[ 5 , 6 ]	693	881	(464)	0.0

وبنفس الطريقة ندمج مدينة دالاس (رقم 4) مع المدينتين دنفر (رقم 5) وديترويت (رقم 6) ليصبح لدينا الجدول التالي:

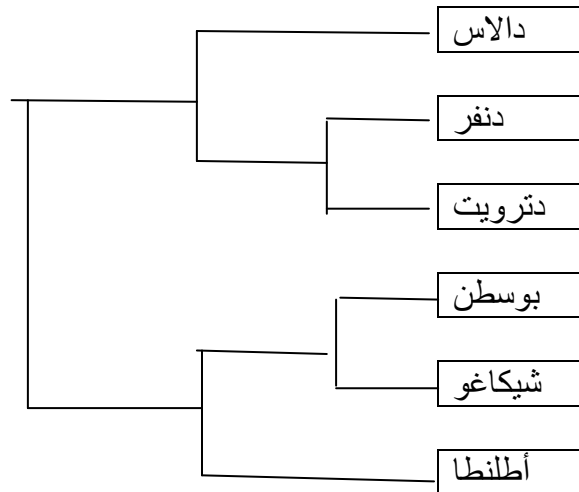
العنقود (المدن)	1	[2 , 3]	[4 , 5 , 6]
1	0.0	(516)	590
[ 2 , 3 ]	(516)	0.0	833
[4 , 5 , 6 ]	590	833	0.0

وهنا يتم دمج مدينة أطلنطا (رقم 1) مع المدينتين بوسطن (رقم 2) وشيكاغو (رقم 3) ليصبح لدينا الجدول التالي :

العنقود (المدن)	[1 , 2 , 3]	[4 , 5 , 6]
[ 1 , 2 , 3 ]	0.0	590
[4 , 5 , 6 ]	590	0.0

والذي أصبح لدينا بموجبه عنقودين (أطلنطا، بوسطن، شيكاغو) و (دالاس، دنفر، دترويت).

وخلاصة التحليل تشير إلى أن دنفر أكثر قرباً إلى دترويت ومن بعد ذلك كانت بوسطن أكثر قرباً إلى شيكاغو. وبعد الدمج وجدنا دالاس تقترب لمجموعة دنفر/ دترويت ثم بعد ذلك نجد أطلنطا تقترب لمجموعة بوسطن/ شيكاغو. والشكل التالي يوضح ذلك:



## تحليل الارتباط القويم (ارتباط المجموعات) Canonical Correlation Analysis

إن تحليل الارتباط القويم هو بمثابة تعميم للارتباط المتعدد Multiple Correlation في مسائل الانحدار المتعدد. ففي الارتباط المتعدد يكون لدينا متغير معتمد أحادي  $Y$  يرتبط بمتغيرين توضيحيين أو أكثر  $(X_1, X_2, \dots, X_p)$  ليتضح لدينا مدى علاقته بها. أما في الارتباط القويم، فيكون لدينا عدد  $q$  من المتغيرات المعتمدة  $(Y_1, Y_2, \dots, Y_q)$  بدلاً من الواحد ونريد معرفة مدى وشكل ارتباطاتها مع مجموعة المتغيرات  $(X_1, X_2, \dots, X_p)$ . أي أننا نعني بالارتباط القويم هنا بأنه الارتباط بين مجموعتين من المتغيرات إحداها مجموعة  $(X_1, X_2, \dots, X_p)$  والأخرى مجموعة  $(Y_1, Y_2, \dots, Y_q)$  عن طريق إيجاد أكبر ترابط خطي للمتغيرات في إحدى المجموعتين مع التراكيب الخطية للمتغيرات في المجموعة الثانية والمجموعتين ذات توزيع مشترك. ووفقاً لذلك، فإن تطبيق هذا التحليل يتطلب تقسيم المتغيرات الناتجة إلى المجموعتين. هذا التقسيم يكون على أساس طبيعة المتغير وليس على أساس التحري والفحص للبيانات. مثال ذلك يمكن دراسة الارتباط ما بين سمات معينة لدى الآباء وسمات أخرى لدى الأبناء للوقوف على مدى الترابط ما بين الآباء والأبناء فيما يخص مجال معين يتم جمع البيانات حولها من خلال متغيرات موصوفة مسبقاً للمجتمعين. كما أن القياس المناسب لكلا المجموعتين من المتغيرات هو النسبي أو الفئوي فقط.

في بعض الأحيان قد يجد الباحث صعوبة في جمع بيانات لمتغيرات محددة تغطي مجال معين ولكنه بحاجة لدراسة هذا المجال. هذه الصعوبة قد تكون نتيجة عدم توفر البيانات أو كلفتها العالية. ما الحل إذا ونحن نعرف أن علم الإحصاء لا يقف عاجزاً أمام مثل هذه الحالات وغيرها، على الباحث هنا اللجوء إلى ما يسمى بالمتغيرات المساعدة Auxiliary Variables لتكون بديلاً عن المتغيرات الأصلية. والتاريخ يزخر بأثلة من هذا الإجراء. والتاريخ الإسلامي يذكر أن معركة كانت على وشك أن تدور بين المسلمين وجيش الروم وأراد القائد معرفة تعداد جيش الروم بشكل غير مباشر فلجأ إلى معرفة ذلك تقديراً من خلال عدد أرغفة الخبز أو الذبائح التي يستهلكها ذلك الجيش والذي كان ممكناً معرفتها إلى حد ما.

إن تحليل الارتباط القويم قد يساهم في هذا الجانب من خلال توصيف مجموعة متغيرات أخرى كمتغيرات مساعدة مقابلة للمتغيرات الأصلية وجمع البيانات عنها ممكناً. إن درجة الثقة بإمكانية اعتماد هذه المتغيرات المساعدة بديلاً عن المتغيرات الأصلية في الدراسة تتناسب طردياً مع قيمة الارتباط القويم الذي نحصل عليه بين مجموعتي المتغيرات.

أحد التساؤلات الأساسية التي على تحليل الارتباط القويم الإجابة عنها هو ما إذا كان بالإمكان استخدام المتغيرات في إحدى المجموعتين للتنبؤ بالمتغيرات في المجموعة الأخرى.

فعندما يكون ذلك ممكناً، فإن ذلك يعني بأن تحليل الارتباط القويم يعمل على تلخيص العلاقات ما بين مجموعتي المتغيرات من خلال تكوين متغيرات جديدة من كلٍّ من مجموعتي المتغيرات الأصلية.

ومن الجدير بالذكر أن مفهوم الارتباط القويم ظهر في الفترة 1936/1935 من قبل الإحصائي هوتلنك Hotelling حيث ذكر أن تحليل الارتباط المتعدد ما هو إلا حالة خاصة من الارتباط القويم. وفي عام 1940 كان فيشر أول من استخدم الارتباط القويم لتحليل الجداول التوافقية ذات الإتجاهين ( rxc ) ذات فئات مرتبة.

وعند التطبيق هنا لو افترضنا وجود مجموعة متغيرات  $(X_1, X_2, \dots, X_p)$  ومجموعة متغيرات أخرى  $(Y_1, Y_2, \dots, Y_q)$  وأن  $(q < p)$  فإنه ستكون لدينا حصيلة إرتباطات ثنائية بعدد  $(p + q)$  ننطلق منها في عملية تحليل الارتباط القويم. هذه العملية تشمل تحديد الارتباط ما بين المتغيرات القويمة  $V_i$  و  $U_i$  حيث:

$$U_1 = a_{11} X_1 + a_{12} X_2 + \dots + a_{1p} X_p$$

$$U_2 = a_{21} X_1 + a_{22} X_2 + \dots + a_{2p} X_p$$

⋮

$$U_r = a_{r1} X_1 + a_{r2} X_2 + \dots + a_{rp} X_p$$

وكذلك:

$$V_1 = b_{11} Y_1 + b_{12} Y_2 + \dots + b_{1q} Y_q$$

$$V_2 = b_{21} Y_1 + b_{22} Y_2 + \dots + b_{2q} Y_q$$

⋮

$$V_r = b_{r1} Y_1 + b_{r2} Y_2 + \dots + b_{rq} Y_q$$

وأن  $(r)$  هنا يمثل أصغر الأعداد  $(p$  و  $q)$ . وهذه العلاقات الخطية يتم تحديدها بحيث نحصل على أعلى إرتباط ما بين  $U_1$  و  $V_1$ . كذلك يكون لدينا الارتباط الأعلى التالي ما بين  $U_2$  و  $V_2$  وهكذا. بمعنى آخر فإن كل زوج من المتغيرات القويمة  $(U_1, V_1)$  و  $(U_2, V_2)$  و..... و  $(U_r, V_r)$  يمثل إتجهاً مستقلاً في العلاقة ما بين مجموعتي المتغيرات  $(X_1, X_2, \dots, X_p)$  و  $(Y_1, Y_2, \dots, Y_q)$ . ولأن أول زوج  $(U_1, V_1)$  من المتغيرات القويمة يملك أعلى إرتباط، فإنه يعتبر الأكثر أهمية. وأن الزوج الثاني الذي يليه وهو  $(U_2, V_2)$  له ثاني أعلى إرتباط، وبالتالي فإنه يليه في الأهمية ويكون ثاني أهم زوج وهكذا بالنسبة

لجميع أزواج المتغيرات القوية إلى أن تأتي على الزوج الأخير وهو  $(U_r, V_r)$  والذي هو الأقل أهمية لكونه ذو الارتباط الأصغر.

### طرق تنفيذ تحليل الارتباط القويم

لعله من السهل برمجة الحسابات المتعلقة بتحليل الارتباط القويم لتنفيذها في الحاسبة شرط أن يكون البرنامج مناسب للتعامل مع المصفوفات الجبرية لأن الأساس في العملية الحسابية هو البدء بتحديد مصفوفات الارتباط ما بين متغيرات كل مجموعة من المتغيرات على انفراد بالإضافة إلى مصفوفة الارتباط ما بين المجموعتين. ونعني بذلك لو كانت لدينا مجموعتي المتغيرات  $X_1, X_2, \dots, X_p$  و  $Y_1, Y_2, \dots, Y_q$  فإنه سيكون لدينا مصفوفة الارتباطات التربيعية ما بين جميع هذه المتغيرات والتي ستكون بأبعاد  $(p+q) \times (p+q)$  وفقاً لما يلي:

$$\begin{array}{c}
 X_1 \quad X_2 \quad \dots \quad X_p \quad Y_1 \quad Y_2 \quad \dots \quad Y_q \\
 \left[ \begin{array}{cc}
 \begin{array}{l} p \times p \text{ matrix} \\ \text{for } X\text{'s Variables} \\ A \end{array} & \begin{array}{l} p \times q \text{ matrix} \\ \text{for } XY\text{'s Variables} \\ C \end{array} \\
 \begin{array}{l} q \times p \text{ matrix} \\ \text{for } YX\text{'s Variables} \\ C' \end{array} & \begin{array}{l} q \times q \text{ matrix} \\ \text{for } Y\text{'s Variables} \\ B \end{array}
 \end{array} \right]
 \end{array}$$

ومن هذه المصفوفة الرئيسية يمكننا تكوين المصفوفة  $B^{-1}C'A^{-1}C$  ذات الأبعاد  $q \times q$  لغرض استخدامها في تحديد القيم المميزة  $\lambda_1 > \lambda_2 > \dots > \lambda_r$  لكون  $(r=q < p)$  والتي تعتبر تربيعة لقيم الارتباطات القوية ما بين المتغيرات القوية المقابلة لها. وهذا بطبيعة الحال يتم من خلال حل مجموعة المعادلات الناتجة عن محددة المصفوفة

$$|B^{-1}C'A^{-1}C - \lambda I| = 0 \quad \dots \dots \dots (1)$$

ولغرض استكمال الصورة التي تكون عليها المتغيرات القوية، فإننا نحتاج إلى تحديد المتجهات المميزة  $b_1, b_2, \dots, b_r$  والتي هي عبارة عن الأوزان القوية لمجموعة المتغيرات  $Y$ 's ضمن المتغيرات القوية  $V$ 's، وهذا يتم من خلال معادلة المصفوفات:

$$(B^{-1}C'A^{-1}C - \lambda I)b = 0$$



والتي يمكن تحويلها للصيغة:

$$(C'A^{-1}C - \lambda B)\mathbf{b} = 0 \quad \dots\dots\dots (2)$$

أما المتجهات  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_r$  والتي هي عبارة عن الأوزان القوية لمجموعة المتغيرات  $X$ 's ضمن المتغيرات القوية  $U$ 's ، فيتم تحديدها من خلال معادلة المصفوفات:

$$\mathbf{a}_i = A^{-1}C \mathbf{b}_i \quad \dots\dots\dots(3)$$

ومن الجدير بالذكر أن هذه الحسابات في أعلاه تتم جميعها باستخدام القيم المعيارية للمتغيرات الأصلية والتي تكون أوساطها صفراً وانحرافها المعياري وحدة واحدة.

ووفقاً لهذه الأوزان القوية  $\mathbf{a}_i$  و  $\mathbf{b}_i$  تكون المتغيرات القوية للمجموعتين على الشكل الآتي:

$$U_i = \mathbf{a}'_i \mathbf{X} = (a_{i1}, a_{i2}, \dots, a_{ip}) \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{bmatrix} = a_{i1}x_1 + a_{i2}x_2 + \dots + a_{ip}x_p \quad , i=1,2,\dots,r$$

$$V_i = \mathbf{b}'_i \mathbf{y} = (b_{i1}, b_{i2}, \dots, b_{iq}) \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_q \end{bmatrix} = b_{i1}y_1 + b_{i2}y_2 + \dots + b_{iq}y_q \quad , i=1,2,\dots,r$$

بطبيعة الحال، عندما نتكلم عن الارتباط القويم ما بين زوج المتغيرات القوية، فإنه يكون معلوماً لدينا أن الارتباط بينهما ينعكس من عدد المشاهدات  $n$  وهو حجم العينة المستخدمة والتي تكون بالشكل التالي:

$$\begin{array}{cccccc} x_{11} & x_{21} & \dots & x_{p1} & y_{11} & y_{21} & \dots & y_{q1} \\ x_{12} & x_{22} & \dots & x_{p2} & y_{12} & y_{22} & \dots & y_{q2} \\ \vdots & \vdots & & \vdots & \vdots & \vdots & & \vdots \\ \vdots & \vdots & & \vdots & \vdots & \vdots & & \vdots \\ x_{1n} & x_{2n} & \dots & x_{pn} & y_{1n} & y_{2n} & \dots & y_{qn} \end{array}$$

وحيث أنه سيكون لدينا العدد الكلي للارتباطات القوية ( $r$ ) فإنه ليس من المنطقي أن نعتمدها جميعاً للتحليل وإنما سنعتمد عدداً محدوداً جداً وذلك وفقاً لمعنويتها في إختبار مناسب لها وكما هو في القسم التالي.

## إختبار معنوية الإرتباط القويم

من البديهي أنه طالما أن مربع أي إرتباط قويم يساوي القيمة المميزة المقابلة، فهذا يعني أن الإرتباط القويم الأول يعتبر الأعلى ويليه الثاني وهكذا فإن الأخير هو الأصغر. لذلك لو إفترضنا أننا وجدنا الإرتباط القويم الثاني غير معنوي، فلا حاجة لنا بإجراء الإختبارات للثالث وما بعده. وطريقة الإختبار هي الإختبار التقريبي الذي تم إقتراحه من قبل بارتليت (1947) Bartlett لمعرفة عدد الإرتباط القويم المعنوية وهي بسيطة وسهلة الإستخدام.

هذه الطريقة تعتمد بداية على إحصاءة الإختبار:

$$\Delta_0^2 = -\left\{n - \frac{1}{2}(p+q+1)\right\} \sum_{i=1}^r \ln(1 - \lambda_i)$$

والتي تتبع توزيع مربع كاي بدرجات حرية مساوية إلى (pq) أي  $(\chi_{pq}^2)$ . فإذا كانت معنوية (أي أن  $\Delta_0^2 > \chi_{pq}^2$ )، فإنه يكون لدينا واحداً في الأقل من الإرتباطات القويمية وهو الأول معنوياً وهو ما يكون في الغالب.

ووفقاً لمعنوية الإرتباط القويم الأول، فإننا نستمر لفحص الإرتباط الذي يليه وفقاً لنفس الصيغة ولكن بعد استبعاد ما يتعلق بأثر الإرتباط الأول من إحصاءة الإختبار ووفقاً لما يلي:

$$\Delta_1^2 = -\left\{n - \frac{1}{2}(p+q+1)\right\} \sum_{i=2}^r \ln(1 - \lambda_i)$$

والتي تتبع توزيع مربع كاي بدرجات حرية مساوية إلى (p-1)(q-1) أي  $(\chi_{(p-1)(q-1)}^2)$  (أي أن  $\Delta_1^2 > \chi_{(p-1)(q-1)}^2$ )، فإنه يكون لدينا الإرتباطات القويمية الأول والثاني معنويان. وهذا يدفعنا للتطبيق التالي لإحصاءة الإختبار بعد استبعاد ما يتعلق بأثر هذان الإرتباطان القويان منها وهكذا حتى نصل إلى آخر إرتباط معنوي. وبشكل عام فإننا لو وجدنا عدد (j) منها معنوياً فإننا نطبق الصيغة التالية:

$$\Delta_j^2 = -\left\{n - \frac{1}{2}(p+q+1)\right\} \sum_{i=j+1}^r \ln(1 - \lambda_i)$$

والتي تتبع توزيع مربع كاي بدرجات حرية مساوية إلى (p-j)(q-j) أي  $(\chi_{(p-j)(q-j)}^2)$ .

ومن الجدير بالذكر أنه لو حصل لدينا عدم معنوية في أي خطوة مما ذكرنا في أعلاه، فإنه يجدر بنا التوقف وإعتبار الإرتباطات القويمية التالية أيضاً ليست معنوية. كما نلاحظ أن حجم العينة (n) هنا له أثر واضح في قيمة دالة الإختبار وبالتالي في معنويته.

## تحليل نتائج المتغيرات القويمة

بعد الوقوف على نتائج الإختبارات أعلاه وحصولنا على نتائج معنوية، فإن الخطوة التالية هو الخروج بنتائج تحليلية لهذه النتائج ودلالة القياسات التي تعكسها. وبالتأكيد فإن ذلك يشمل جميع القيم المعنوية للإرتباطات القويمة ولو أن التركيز سيكون على الأول وهو الأكبر وقد نكتفي بذلك حتى في حالة كون قيمته غير معنوية.

فإذا إعتبرنا المتغيرات القويمة:

$$U_i = a_{i1} X_1 + a_{i2} X_2 + \dots + a_{ip} X_p$$

و

$$V_i = b_{i1} Y_1 + b_{i2} Y_2 + \dots + b_{iq} Y_q$$

فإنه من الممكن توضيح  $U_i$  بالنسبة إلى أي من مجموعة المتغيرات  $X_1, X_2, \dots, X_p$  التي لها أكبر أوزان قويمة  $a_{ij}$  وكذلك فيما يخص  $V_i$  بالنسبة إلى مجموعة المتغيرات  $Y_1, Y_2, \dots, Y_q$  التي لها أكبر أوزان قويمة  $b_{ij}$  بغض النظر عن إشارة هذه الأوزان موجبة كانت أم سالبة.

ومن المفيد هنا أن نذكر بأنه قد نجد الوزن القويم  $a_{i1}$  بإشارة موجبة بينما في الوقت نفسه يكون معامل الإرتباط (أو ما سنطلق عليه بالمعاملات التركيبية Structure Coefficients) ما بين  $U_i$  والمتغير  $X_1$  سالبة (أي عكس ذلك). وتفسيرنا لظهور هذه الحالة العكسية يعود إلى وجود إرتباط عالي ما بين المتغير  $X_1$  وأي عدد من المتغيرات الأخرى في مجموعة  $X$  مما يدفع إلى ظهور هذه الحالة. أي هنالك حالة التعدد الخطي (multicollinearity) ما بين هذه المجموعة من المتغيرات وتأثيرها هنا مشابه لتأثيرها على تقدير المعاملات في نموذج الإنحدار المتعدد. وفي مثل هذه الحالة، لا يمكن الوقوف على الحجم الحقيقي لإسهام هذا المتغير بشكل مستقل تجاه المتغير القويم بل هنالك تداخل من تأثيرات المتغيرات الأخرى التي له معها إرتباطات واضحة. ولتجاوز هذه الحالة، فإننا نرى أهمية النظر إلى معامل الإرتباط (المعاملات التركيبية) ما بين المتغير القويم وكل متغير أولي مقابل له بدلاً من الأوزان القويمة  $a_{ij}$  أو  $b_{ij}$ . والمعاملات التركيبية  $S_{xi}$  و  $S_{yi}$  يمكن حسابها حسب الصيغة التالية:

$$S_{xi} = R_{xx} (a_i) = A (a_i)$$

$$S_{yi} = R_{yy} (b_i) = B (b_i)$$

وفي ضوء هذه المعادلات وفي حالة كون مصفوفة الإرتباط لأي مجموعة عبارة عن مصفوفة الوحدة (Identity matrix)، فإنه من الواضح أن المعامل التركيبي لأي متغير

أصيل يكون مساوياً للوزن القويم لذلك المتغير. وبشكل عام، فإن مربع المعامل التركيبي لأي متغير يمثل نسبة مساهمته في تفسير التباين الحاصل في المتغير القويم.

والمثال التالي مأخوذ عن كتاب الطرق الإحصائية متعددة المتغيرات لمؤلفه B. F. J. Manly ويتضمن خمسة متغيرات جينية وتمثل مجموعة X وأربع متغيرات بيئية وتمثل مجموعة Y لنوع معين من الفراشات<sup>(6)</sup>. وقد تم جمع البيانات من 16 مستعمرة بولاييتين من الولايات الأمريكية ( كاليفورنيا وأوريغون). والجدول التالي يبين القيم المعيارية لمجموعي المتغيرات، ولغرض الترتيب المريح، فإنه يفضل أن تكون مجموعة المتغيرات الأكبر عدداً ممثلة بالمجموعة X والأخرى الأقل عدداً بالمجموعة Y.

المتغيرات الجينية					المتغيرات البيئية			
X <sub>1</sub>	X <sub>2</sub>	X <sub>3</sub>	X <sub>4</sub>	X <sub>5</sub>	Y <sub>1</sub>	Y <sub>2</sub>	Y <sub>3</sub>	Y <sub>4</sub>
-0.42	-0.56	0.32	0.29	-0.02	-0.56	1.07	0.12	-0.35
-0.42	1.18	0.13	-0.66	-0.38	-0.45	-0.58	-0.82	1.02
-0.42	-0.16	0.88	-0.26	-0.02	-0.53	0.00	0.12	0.47
-0.42	-0.43	0.04	-0.21	0.91	-0.54	0.00	0.12	0.47
-0.42	-0.83	-0.98	-0.06	1.65	-0.54	0.00	0.12	0.47
-0.42	-0.69	0.04	-0.36	1.37	-0.60	-0.93	0.27	0.65
-0.42	-0.96	-0.33	-0.06	0.91	-0.41	-0.50	0.27	0.65
1.99	1.84	1.99	-1.31	-1.22	-0.50	-1.29	0.59	0.56
2.96	2.51	1.25	-1.16	-1.59	-0.52	-1.29	0.59	0.56
-0.42	-0.83	-1.17	1.44	-0.48	-0.21	-0.65	0.27	0.19
-0.18	-0.43	1.43	-0.91	0.45	-0.12	-0.43	0.59	0.56
-0.42	-0.03	-0.42	0.74	-0.39	-0.04	2.14	0.43	-0.26
-0.42	0.24	-0.33	-0.21	0.35	0.14	0.42	0.74	-0.44
0.30	-0.03	-0.14	-0.96	0.91	-0.04	-0.50	1.21	-0.08
-0.42	-0.29	-1.07	1.64	-1.22	2.00	1.00	-2.08	-1.45
-0.42	-0.56	-1.63	2.03	-1.22	2.92	1.57	-2.55	-3.00

ومنها يتم تحديد مصفوفة معامل الارتباط التالية ما بين مجموعتي المتغيرات.

1.000	0.855	0.618	-0.532	-0.506	∴	-0.203	-0.530	0.295	0.221
0.855	1.000	0.615	-0.548	-0.597	∴	-0.190	-0.410	0.173	0.246
0.618	0.615	1.000	-0.824	-0.127	∴	-0.573	-0.550	-0.536	0.593
-0.532	-0.548	-0.824	1.000	-0.264	∴	0.727	0.699	-0.717	-0.759
-0.506	-0.597	-0.127	-0.264	1.000	∴	-0.458	-0.138	0.438	0.412
...	...	...	...	...	∴	...	...	...	...
-0.203	-0.190	-0.573	0.727	-0.458	∴	1.000	0.568	-0.828	-0.936
-0.530	-0.410	-0.550	0.699	-0.138	∴	0.568	1.000	-0.479	-0.705
0.295	0.173	-0.536	-0.717	0.438	∴	-0.828	-0.479	1.000	0.719
0.221	0.246	0.593	-0.759	0.412	∴	-0.936	-0.705	0.719	1.000

وهذا يمثل المصفوفة:

$$\begin{bmatrix} A & \vdots & C \\ \dots & \dots & \dots \\ C' & \vdots & B \end{bmatrix}$$

وبعد تطبيق المعادلة (1) في أعلاه نتوصل أولاً إلى القيم المميزة التالية:

$$0.7731, 0.5570, 0.1694, 0.0472$$

ومنها يتم تطبيق المعادلتين (2 ، 3) لـ **a** ومن ثم قيم **b** لـ **a** ليتم تحديد قيم الأوزان القوية **b** ومن ثم قيم **a** المقابلة لنخرج بالمعادلات الخطية التالية التي تمثل كامل أزواج المتغيرات القوية بقيم قياسية:

$$U_1 = -0.675X_1 + 0.909X_2 + 0.367X_3 + 1.442X_4 + 0.269X_5$$

$$V_1 = -0.114Y_1 + 0.619Y_2 - 0.693Y_3 + 0.048Y_4$$

$$U_2 = -1.087X_1 + 3.034X_2 + 2.216X_3 + 3.439X_4 + 2.928X_5$$

$$V_2 = -0.777Y_1 + 0.980Y_2 - 0.562Y_3 + 0.928Y_4$$

$$U_3 = 1.530X_1 + 2.049X_2 + 2.231X_3 + 4.916X_4 + 3.611X_5$$

$$V_3 = -3.654Y_1 - 0.601Y_2 - 0.565Y_3 - 0.0483.623Y_4$$

$$U_4 = 0.284X_1 - 2.331X_2 - 0.867X_3 - 1.907X_4 - 1.133X_5$$

$$V_4 = 1.594Y_1 + 0.860Y_2 + 1.599Y_3 + 0.742Y_4$$

ومن خلال القيم المميزة  $\lambda_i$  في أعلاه يتم تحديد الارتباطات القوية  $R_{ci}$  والتي تساوي

الجذر التربيعي للقيم المميزة حيث ستكون:

$$R_{c1} = 0.879, R_{c2} = 0.746, R_{c3} = 0.412, R_{c4} = 0.217$$

ومع أن قيم الارتباطين الأول والثاني تبدوان كبيرة نسبياً إلا أنه بتطبيق إختبار بارتليت لم تكن أيها ذات معنوية بمستوى (0.05). ومع أننا يمكن أن نكتفي بإختبار الأولى لأنها غير معنوية ونتوقف، إلا أننا سنذكر هنا جميع نتائج الإختبار هذا حيث أن

$$\Delta_0^2 = 27.85 \text{ ( بدرجات حرية 20 ) مقارنة مع } \chi_{20}^2 = 31.4104 \text{ الجدولية}$$

$$\Delta_1^2 = 11.53 \text{ ( بدرجات حرية 12 ) مقارنة مع } \chi_{12}^2 = 21.0261 \text{ الجدولية}$$

$$\Delta_2^2 = 2.57 \text{ ( بدرجات حرية 6 ) مقارنة مع } \chi_6^2 = 12.5916 \text{ الجدولية}$$

$$\Delta_3^2 = 0.52 \text{ ( بدرجات حرية 2 ) مقارنة مع } \chi_2^2 = 5.99147 \text{ الجدولية}$$

وقد يبدو من الغريب أن نحصل على نتيجة غير معنوية مع قيمة الارتباط القويم الأول الذي هو كبير بشكل واضح. والسبب، وكما ذكرنا سابقاً، هو نتيجة مباشرة لتأثير حجم العينة ( $n=16$ ) الصغير نسبياً. ومع ذلك لعله من المفيد أن نضع عدم المعنوية هذه جانباً لغرض المضي في استكمال تفسير النتائج بالنسبة للارتباط القويم الأول ( $RC_1$ ).

وقبل البدء في ذلك لعله من المفيد أن نذكر هنا قيم معامل الارتباط (المعامل التركيبية) ما بين الزوج الأول من المتغيرات القوية ( $U_1, V_1$ ) وكل متغير أولي مقابل لهما سوية مع الأوزان القوية المرافقة في الجدول التالي:

المعامل بالنسبة إلى $U_1$ مع مجموعة X			المعامل بالنسبة إلى $V_1$ مع مجموعة Y		
المتغير	الوزن القويم	الارتباط $S_{xi}$ (المعامل التركيبية) بين $U_1$ و X	المتغير	الوزن القويم	الارتباط $S_{yi}$ (المعامل التركيبية) بين $V_1$ و Y
$X_1$	-0.675	-0.57	$Y_1$	-0.114	0.77
$X_2$	0.909	-0.39	$Y_2$	0.619	0.85
$X_3$	0.367	-0.70	$Y_3$	-0.693	-0.86
$X_4$	1.442	0.92	$Y_4$	0.048	-0.78
$X_5$	0.269	-0.36			

ومن خلال هذا الجدول وبالنسبة إلى  $U_1$  نلاحظ أن الوزن القويم بالنسبة إلى  $X_1$  هو الوحيد بإشارة سالبة مما نستطيع القول بأن اتجاه قيم  $X_1$  معاكسة تماماً لإتجاه قيم المتغيرات الأخرى في المجموعة. ومن جهة أخرى، وبالنسبة إلى  $V_1$  فإننا نلاحظ أوزان قوية عالية موجبة مع  $Y_2$  (0.619) وعالية سالبة مع  $Y_3$  (-0.693) مما يمكننا قوله بأن القصور في قيم المتغير الجيني  $X_1$  في مستعمرة ما للفراشات يتزامن مع تزايد في قيم المتغير البيئي  $Y_2$  وتناقص قيم المتغير البيئي  $Y_3$ .

وبالنظر إلى معامل الارتباط ما بين المتغير القويم  $U_1$  ومجموعة المتغيرات X نلاحظ أن  $U_1$  يرتبط إيجابياً بشكل واضح مع المتغير  $X_4$  (0.92) وسلبياً مع بقية المتغيرات الأربعة الأخرى والتي أعلاها مع  $X_3$  (-0.70). وبناءً على ذلك يمكننا القول بأن  $U_1$  يؤشر ارتفاع قيمة المتغير  $X_4$ . وهذا بطبيعة الحال يعتبر مغايراً بعض الشيء عما يمكننا قوله من خلال النظر إلى الأوزان القوية لمجموعة المتغيرات X ضمن الدالة  $U_1$ . وبشكل عام، فإن التفسير في ضوء معامل الارتباط (المعامل التركيبية) يبدو أفضل هنا مما هو عليه مع الأوزان القوية.

وجدير بالذكر هنا أن هنالك مشكلة حقيقية في تفسير النتائج بالنسبة إلى المتغيرات القوية في حالة وجود ارتباط عالي نسبياً ما بين المتغيرات الأصلية وهذا ما وجدناه هنا بالفعل حيث

معامل الارتباط العالي نسبياً ما بين المتغيرات في المجموعة  $X$  وكما يتضح لنا من الجدول الأخير.

### مقاييس أخرى للتحليل

هنالك مقياسان آخران يمكن حسابهما من البيانات وقد يضيفان بعض الجوانب الأخرى من تفسير النتائج وهما:

#### 1- معامل كفاية الجودة ( $A_d$ ) ( Adequacy Coefficient )

وهذا المعامل يشير إلى نسبة التباين الكلي الحاصل لمجموعة المتغيرات الأصلية والمفسرة من قبل المتغير القويم لتلك المجموعة. ويتم حساب هذه المعامل والتي تساوي معدل مجموع مربعات المعامل التركيبية لمجموعة المتغيرات الأصلية عند متغير قويم معين ووفقاً للصيغة التالية:

$$A_d(U_{xi}) = \frac{\sum_{r=1}^p S_{xir}^2}{p} \quad (100)$$

$$A_d(V_{yi}) = \frac{\sum_{r=1}^q S_{yir}^2}{q} \quad (100)$$

ولأننا إعتدنا معامل الارتباط القويم الأول، فإن معامل كفاية الجودة لمجموعتي المتغيرات هما حسب الآتي:

$$A_d(U_{x1}) = \frac{\sum_{r=1}^5 S_{x1r}^2}{5} (100) = \frac{(-0.57)^2 + (-0.39)^2 + (-0.70)^2 + (0.92)^2 + (-0.36)^2}{5} = 0.40$$

$$A_d(V_{y1}) = \frac{\sum_{r=1}^4 S_{y1r}^2}{4} (100) = \frac{(0.77)^2 + (0.85)^2 + (-0.86)^2 + (-0.78)^2}{4} = 0.67$$

وهذا يعني أن المتغير القويم  $U$  يوضح ما نسبته 40% من التباين الكلي لمجموعة المتغيرات الجينية  $X$ . ووفقاً لقيم  $S_{x1r}^2$  فإنه بإستطاعتنا ترتيب المتغيرات الجينية  $X$  حسب قوة تأثيرها في المتغيرات البيئية  $Y$  فيكون المتغير  $X_4$  هو الأكثر تأثيراً (المرتبة الأولى) يليه المتغير  $X_3$  ثم  $X_1$  ثم  $X_2$  ثم الأقل تأثيراً  $X_5$ . كما أن المتغير القويم  $V$  يوضح ما نسبته 67% من التباين الكلي لمجموعة المتغيرات البيئية  $Y$  لهذه الفراشات.

## 2- معامل الإفاضة ( Redundancy Coefficient ) $Rd_{y|x}$

وهذا المعامل يمثل نسبة التباين الحاصلة في متغيرات مجموعة معينة والمفسرة من قبل متغيرات المجموعة الأخرى . أي أنه يفيد في معرفة مدى تأثير مجموعة متغيرات أصلية  $X$  في مجموعة المتغيرات الأخرى  $Y$  (والعكس بالعكس). وقيمة هذا المعامل تساوي معدل مجموع مربعات معاملات الارتباط المتعدد وقيمه ضمن المجال [ 0 , 1 ] ويأخذ القيمة 1.0 لأي زوج من المتغيرات القوية عندما تكون نسبة التباين المشترك لها مساوياً إلى ( 100% ) أي عندما يكون الارتباط القوي عند هذا الزوج مساوياً إلى (1.0). وقد تحسب هذه القيمة أيضاً كما في الصيغة التالية:

$$Rd_{y|x} = (Rc_1)^2 \cdot A_d(V_{y1}) = (0.879)^2 \cdot (0.67) = 0.52$$

وهذا يعني أن مجموعة المتغيرات الجينية  $X$  توضح ما نسبته 52% من تباين مجموعة المتغيرات البيئية  $Y$  للفرشات.

### ملاحظة:

نود التذكير هنا بأن الحسابات أعلاه مبنية على أساس إطلاق المجموعة  $X$  للمجموعة ذات العدد الأكبر من المتغيرات وعددها (p) بالنسبة للمجموعتين و (q < p).

وفي هذا السياق، نود أن نذكر هنا نتائج مثال تطبيقي آخر لإستخدام تحليل الارتباط القويم من خلال دراسة أدت نتائجها في حينه إلى قيام وزارة التعليم العالي العراقية لإعادة النظر في بعض أسس القبول في كليات المجموعة الطبية<sup>(7)</sup>. و خلاصة فكرة الدراسة، التي أجريت عام 1998، أن وزارة التعليم العالي قد حددت شروط خاصة للقبول في كليات المجموعة الطبية (كلية الطب، كلية طب الأسنان، كلية الصيدلة وكلية الطب البيطري) بتحديد مواد المفاضلة وهي (الفيزياء، الكيمياء، الأحياء) وأن لا تقل درجة الطالب في أي منها عن 70% في المرحلة الثانوية (التوجيهي)، والتي تم تطبيقها منذ العام الدراسي 1972/1971، منطلقين من الافتراض بأن هذه المواد الثلاثة لها تأثير على تحصيل الطالب في مواد السنة الأولى في هذه الكليات وبالتالي في مسيرته للسنوات الأخرى. ولذلك تم تصميم الدراسة موضوع البحث لغرض التأكد من منطقيته هذا الافتراض من عدمه ومن خلال نتائج الارتباط القويم بين درجات الطالب في مواد المرحلة الثانوية ( مجموعة المتغيرات  $X$  ) ودرجاته في مواد السنة الأولى من الكلية المنتسب إليها (مجموعة المتغيرات  $Y$ ) ولثلاث سنوات أكاديمية (1986/85، 1990/89، 1994/93). وسنكتفي بتناول بيانات كلية الطب فقط لغرض التوضيح فيما تشيراليه النتائج والتي تظهر في الجدول التالي:



السنة الدراسية						مواد المرحلة الثانوية	المتغير
1994/1993		1990/1989		1986/1985			
ترتيب العلاقة	قيمة % (Sxi) <sup>2</sup>	ترتيب العلاقة	قيمة % (Sxi) <sup>2</sup>	ترتيب العلاقة	قيمة % (Sxi) <sup>2</sup>		
6	34	2	50	6	14	لغة عربية	X <sub>1</sub>
3	49	5	24	5	24	لغة إنكليزية	X <sub>2</sub>
2	61	1	52	3	52	رياضيات	X <sub>3</sub>
5	38	3	46	2	56	أحياء	X <sub>4</sub>
1	62	4	37	4	50	كيمياء	X <sub>5</sub>
4	46	6	19	1	59	فيزياء	X <sub>6</sub>

وفي ضوء النتائج هذه، فإن مدى تحقق صحة إفتراض الوزارة يتحقق من خلال ترتيب علاقة متقدم لمواد المفاضلة (فيزياء، كيمياء، أحياء) في جميع هذه السنوات. أي أن تأخذ الترتيب (1، 2، 3). ولكن الذي حصل أن هذه المواد أخذت الترتيبات (1، 4، 2) و (6، 4، 3) و (4، 1، 5) للسنوات الثلاث على التوالي. أي أن مادة الرياضيات أخذت الترتيب المتقدم بدلاً من الكيمياء لعام 1986/85، ومادتي الرياضيات واللغة العربية بدلاً من الفيزياء والكيمياء لعام 1990/89، ومادتي الرياضيات واللغة الإنكليزية بدلاً من الفيزياء والأحياء لعام 1994/93. ووفقاً لذلك، يصبح من المؤكد عدم صحة الإفتراض الذي إعتدته الوزارة بشأن منطقية مواد المفاضلة.

## التحليل العاملي

### Factor Analysis (FA)

التحليل العاملي عبارة عن أسلوب غالباً ما يستخدم في تكوين متغيرات جديدة تلخص جميع المعلومات التي من الممكن توفرها في المتغيرات الأصلية. وعلى سبيل المثال، في حالة إعطاء إمتحان لطلبة مرحلة تعليمية محددة في عدد من المواد مثل القراءة والإملاء والرياضيات والعلوم في الوقت الذي يحصل فيه الطلبة في النهاية على تقدير عالي، متوسط، أو واطئ كنتيجة نهائية لجميع المواد. فإذا ما كان هذا الذي سيحصل، فإننا نستطيع القول بأن نتائج الإمتحان هذه يمكن توضيحها من خلال تحديد سمة أو عامل مشترك لجميع نتائج الإمتحانات الأربعة. في هذا المثال، قد يبدو منطقياً لإفترض مثل هذه السمة (العامل) بمثابة تعبير عن الذكاء أو الأداء العام.

والتحليل العاملي يستخدم أيضاً لدراسة العلاقات التي من الممكن وجودها ما بين المتغيرات التي تم قياسها ضمن مجموعة البيانات. ومثلما هي الحال في تحليل المركبات الرئيسية، فإن التحليل العاملي هو من أساليب تحكم المتغيرات.

أحد الأهداف الرئيسية للتحليل العاملي يكمن في تحديد ما إذا كانت المتغيرات الناتجة تعكس أنماطاً من العلاقات مع بعضها البعض بحيث يمكن تقسيم المتغيرات إلى مجاميع جزئية من المتغيرات تكون مترابطة بشكل واضح فيما بينها ضمن مجموعة بعينها في حين نجد هذه المتغيرات أقل ترابطاً ضمن المجاميع الأخرى. ولذلك فإن التحليل العاملي غالباً ما يستخدم لدراسة البناء الترابطي ما بين المتغيرات ضمن مجموعة البيانات. وهذه المتغيرات كمية وبالتالي فإنها ذات مقياس نسبي أو فئوي.

أحد أوجه الشبه ما بين تحليل المركبات الرئيسية والتحليل العاملي هو أن التحليل العاملي يمكن استخدامه أيضاً لتكوين متغيرات جديدة غير مترابطة مع بعضها البعض. مثل هذه المتغيرات تسمى الدرجات العاملية Factor Scores.

يتميز التحليل العاملي بإحدى الإيجابيات مقارنة بتحليل المركبات الرئيسية عند تكوين المتغيرات الجديدة وهي أن هذه المتغيرات الجديدة التي تتكون من خلال التحليل العاملي هي، بشكل عام، قابلة للتفسير بشكل أسهل بكثير مقارنة بتلك المكونة من خلال تحليل المركبات الرئيسية. وهذا يعني أنه إذا ما أراد الباحث تكوين مجموعة متغيرات جديدة أصغر عدداً وقابلة للتفسير والإستنتاج وتلخص معظم المعلومات في المتغيرات الأصلية، فإنه يجب أخذ التحليل العاملي في الإعتبار حيث عدد العوامل الناتجة تكون أقل عدداً من المتغيرات التي نبدأ التحليل بها.

## أهداف التحليل العاملي

- (1) التعرف على الأنماط البينية  
فإذا كانت لدينا مصفوفة إرتباطات بين مجموعة من المتغيرات تمثل خصائص معينة فإن أسلوب التحليل العاملي يكشف عن الأنماط المنفصلة للعلاقات البينية التي تتضمنها المتغيرات ويحدد علاقة كل متغير بتلك الأنماط ودرجة هذه العلاقة.
- (2) الإختصار في وصف البيانات  
إذا كان لدينا مجموعة كبيرة من المشاهدات التي تخص عدد كبير من المتغيرات، فإنه يمكن أن يتم تركيز هذه البيانات ضمن عدد قليل من العوامل تقوم مقام المتغيرات العديدة في إجراء الوصف والمقارنة.
- (3) بناء مقاييس التقدير  
كثيراً ما يتطلب الأمر تصميم مقياس لتقدير سلوك الأفراد في مجال معين، ويستلزم ذلك إعطاء أوزان معينة للخصائص التي يتضمنها هذا المقياس. والتحليل العاملي يحقق هذا الهدف بتصنيفه لهذه الخصائص والمتغيرات في صورة عوامل مستقلة.
- (4) تحويل البيانات  
يستعمل التحليل العاملي في تحويل البيانات إلى صورة أخرى تتوفر فيها بعض الشروط التي يمكن من خلالها تطبيق أساليب إحصائية أكثر جدوى على هذه البيانات. ومثال ذلك في تحليل الإنحدار المتعدد فإن التقدير الصحيح للمعاملات يتطلب عدم وجود حالة التعدد الخطي. ولكن في حالة وجود مثل هذه الحالة، فإنه يمكن استخدام التحليل العاملي وتحويلها إلى أقل عدداً من العوامل غير المترابطة يمكن إحلالها محل المتغيرات الأصلية في معادلة الإنحدار.

## النموذج العاملي Factor Model

إذا افترضنا نظام متعدد المتغيرات يتضمن عدد  $P$  من الإستجابات (المتغيرات العشوائية)

$[X_1, X_2, \dots, X_p]$  تتبع توزيعاً طبيعياً متعدداً، فإنه يمكن أن نكتب النموذج التالي

للإستجابات:

$$\begin{aligned} X_1 &= a_{11}F_1 + a_{12}F_2 + \dots + a_{1m}F_m + e_1 \\ &\vdots \\ X_p &= a_{p1}F_1 + a_{p2}F_2 + \dots + a_{pm}F_m + e_p \end{aligned}$$

أو بشكل عام:

$$X_i = a_{i1}F_1 + a_{i2}F_2 + \dots + a_{im}F_m + e_i \quad \dots \dots \dots (1)$$

وبصيغة المصفوفات تكون:

$$\underline{X} = \underline{A}\underline{F} + \underline{e} \quad \dots\dots\dots(2)$$

$$\underline{X}' = [X_1, X_2, \dots\dots\dots, X_p]$$

$$\underline{F}' = [F_1, F_2, \dots\dots\dots, F_p]$$

$$\underline{e}' = [e_1, e_2, \dots\dots\dots, e_p]$$

$$A = \begin{bmatrix} a_{11} & \dots\dots & a_{1m} \\ \vdots & \dots\dots & \\ a_{p1} & \dots\dots & a_{pm} \end{bmatrix}$$

حيث أن:

$\underline{X}$  : موجه المتغيرات العشوائية

$\underline{F}$  : موجه العوامل المشتركة Common Factors

$A$  : مصفوفة تحميلات العوامل Factor Loadings

$\underline{e}$  : الموجه العشوائي

### الفروض الأساسية لتحليل العاملي

#### 1- الفرضية الأولى

تعتمد هذه الفرضية على أساس وجود إرتباطات بين المتغيرات قيد الدراسة نتيجة وجود عوامل مشتركة فيما بينها، ويكون النموذج العاملي إلى (p) من المتغيرات المشاهدة لعينة حجمها (n) على أساس تكوين دالة خطية إلى (q) من العوامل (بالقيمة المعيارية) وكما يلي:

$$Z_{ij} = a_{1j}F_1 + a_{2j}F_2 + \dots\dots\dots + a_{qj}F_q + d_j u_{ij} \quad \dots\dots\dots (3)$$

حيث أن:

$Z_{ij}$  = القيمة المعيارية للمشاهدة (i) بالنسبة للمتغير (j)

$a_{qj}$  = تحميل العامل (بالنسبة للمتغير (j) وهي أوزان مرافقة لقيم العوامل المشتركة)

$F_q$  = القيمة المعيارية للعامل (q) المشترك المحدد

$u_{ij}$  = القيمة المعيارية للمفردة (i) للعامل المشترك المحدد

$d_j$  = معامل يمثل الوزن المرافق لقيمة العامل الفريد (العامل الخاص بمتغير واحد).

ويقسم التباين الكلي للمتغيرات إلى ثلاثة أنواع (ضمن هذه الفرضية) وهي:

1. التباين المشترك *Common Variance* : وهو الجزء الذي يرتبط مع بقية المتغيرات الأخرى من خلال العوامل المشتركة.
  2. التباين الخاص *Specific Variance* : وهو الجزء الذي لا يرتبط مع بقية المتغيرات بل مع المتغير نفسه.
  3. تباين الخطأ *Error Variance* : وهو الجزء الناتج من خلال حدوث أخطاء عند سحب العينة أو قياسها أو تغيرات أخرى ترتبط بالشخص الذي يسحب المشاهدة.
- ويمكن التعبير عن التباين الكلي  $\sigma_j^2$  للمتغير وأجزائه كما يلي:

$$\sigma_j^2 = (\sigma_{j1}^2 + \sigma_{j2}^2 + \dots + \sigma_{jq}^2) + (\sigma_{js}^2) + (\sigma_{je}^2) \quad \dots \dots \dots (4)$$

$$= (\text{Common Variance}) + (\text{Specific Variance}) + (\text{Error Variance})$$

ويساهم كلاً من التباين المشترك والتباين الخاص في تكوين التباين الثابت (المعتمد) كما يفترض بأن تباين الخطأ لا يرتبط بالتباين الثابت.

وبقسمة طرفي المعادلة (4) على  $\sigma_j^2$  يصبح لدينا:

$$\frac{\sigma_j^2}{\sigma_j^2} = \frac{(\sigma_{j1}^2 + \sigma_{j2}^2 + \dots + \sigma_{jq}^2)}{\sigma_j^2} + \frac{\sigma_{js}^2}{\sigma_j^2} + \frac{\sigma_{je}^2}{\sigma_j^2} \quad \dots \dots \dots (5)$$

$$1 = a_{j1}^2 + a_{j2}^2 + \dots + a_{jq}^2 + S_j^2 + e_j^2 \quad \dots \dots \dots (6)$$

حيث أن الجذر التربيعي للتباينات المشتركة  $a_{j1}^2, a_{j2}^2, \dots, a_{jq}^2$  هي احتمالات العوامل والتي تمثل مقدار الارتباط للمتغير (j) بكل عامل.

## 2- الفرضية الثانية

تقوم هذه الفرضية على أساس أن معامل الارتباط بين متغيرين (j , i) يعود إلى طبيعة تشبعهما بالعوامل المشتركة ومدى هذا التشبع. أي أن معامل الارتباط بين متغيرين يساوي مجموع حاصل ضرب احتمالات المتغيرات بالعوامل المشتركة بينهما. أي أن

$$r_{ij} = a_{i1}a_{j1} + a_{i2}a_{j2} + \dots + a_{iq}a_{jq} \quad \dots \dots \dots (7)$$

حيث أن هذه المعادلة إنما تعبر عن العوامل المتعامدة (orthogonal) ويمكن كتابتها بصيغة المصفوفات وبالشكل التالي:

$$R = AA'$$

حيث أن:

$R$  = مصفوفة الارتباط

$A$  = مصفوفة احتمالات العوامل

$$R = \begin{bmatrix} r_{11} & r_{12} & \cdots & r_{1p} \\ r_{21} & r_{22} & \cdots & r_{2p} \\ \vdots & \vdots & \cdots & \vdots \\ r_{p1} & r_{p2} & \cdots & r_{pp} \end{bmatrix}$$

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1q} \\ a_{21} & a_{22} & \cdots & a_{2q} \\ \vdots & \vdots & \cdots & \vdots \\ a_{p1} & a_{p2} & \cdots & a_{pq} \end{bmatrix}$$

### الإشتراكيات (Communalities) وطرق تقديرها

إن قيم الإشتراكيات (الشيوع) يرمز لها بالرمز  $h_j^2$  حيث  $0 \leq h_j^2 \leq 1$  والتي تكون قيمتها

$$h_j^2 = a_{j1}^2 + a_{j2}^2 + \dots + a_{jq}^2, \quad j = 1, 2, \dots, q$$

أي أنها تمثل مجموع مربعات تحميلات (تشبعات) كل متغير وهو بمثابة نسبة التباين الذي تفسره العوامل المشتركة الناتجة من تحليل مصفوفة الارتباط  $R$ .

وفي ضوء ذلك، يمكننا وضع الصيغة لتالية للمعادلة (6)

$$1 = h_j^2 + S_j^2 + e_j^2 = R_{jj} + e_j^2$$

وهنا سنعمل على تقدير قيم الإشتراكيات لتحل كعناصر تقديرية في مصفوفة الارتباط. فإذا كانت العناصر القطرية لمصفوفة الارتباط تساوي الواحد صحيح فإننا نطلق عليها تسمية "مصفوفة الارتباطات الكاملة" ( Complete correlation matrix ). أما إذا تم استبدال العناصر القطرية بقيم الإشتراكيات  $h_j^2$  فإنها تسمى "مصفوفة الارتباطات المخفضة" (Reduced correlation matrix )

### طرق تقدير الإشتراكيات

#### (1) الإرتباط الأكبر (maximum correlation)

وفي هذه الطريقة يتم اعتماد أكبر معاملات الارتباط بين متغيراً موضوع التقدير وبقيّة المتغيرات الأخرى لتمثل القيم التقريبية للإشتراكيات ويفضل استخدام هذه الطريقة في حالة مصفوفة ارتباطات مكونة من عدد كبير من المتغيرات.

## (2) ثلاثى الأبعاد (third dimension)

ووفقاً لهذه الطريقة يتم تقدير اشتراكية أي متغير (j) بالصيغة الثلاثية:

$$h_j^2 = \frac{r_{ji}r_{jk}}{r_{ik}}$$

حيث أن (i) و (k) هما المتغيران اللذان لهما أعلى ارتباط مع المتغير (j). ووفقاً لهذه الطريقة يتم تقليل تأثير الارتباطات العالية.

## (3) معدل الارتباطات (mean of correlations)

ويتم من خلال استخراج معدل معامل الارتباط لذلك المتغير (j) مع بقية المتغيرات وعلى النحو التالي:

$$h_j^2 = \sum_{j \neq i} \frac{r_{ji}}{n-1}$$

## (4) مربع الارتباط المتعدد (Square Multiple Correlation (SMC)

وتعتبر هذه الطريقة من أكثر الطرق استخداماً لتقدير قيم الارتباطات  $h_j^2$  وبالإعتماد على مصفوفة الارتباط  $R$  حيث يتم تحديد معكوس هذه المصفوفة  $R^{-1}$  واستخدامها لتقدير الارتباطات على النحو التالي:

$$h_j^2 = SMC = 1 - \frac{1}{r_{jj}}$$

حيث أن  $r_{jj}$  تمثل العنصر القطري لمعكوس مصفوفة الارتباطات لذلك المتغير (j).

ولغرض توضيح استخدام البعض من هذه الطرق، سنستخدم مثلاً بثلاثة متغيرات لغرض التبسيط.

**مثال:**

لنفترض أن لدينا مصفوفة الارتباطات التالية للمتغيرات  $X_3$  ,  $X_2$  ,  $X_1$  :

$$R = \begin{bmatrix} 1 & 0.61 & 0.40 \\ 0.61 & 1 & 0.48 \\ 0.40 & 0.48 & 1 \end{bmatrix}$$

فإن تقدير الارتباطات الثلاثة ستكون وفقاً للطريقة الأولى كما يلي:

$$h_1^2 = 0.61$$

$$h_2^2 = 0.61$$

$$h_3^2 = 0.48$$

ووفقاً للطريقة الثانية فإنها:

$$h_1^2 = \frac{(0.61)(0.40)}{0.48} = 0.508$$

$$h_2^2 = \frac{(0.61)(0.48)}{0.40} = 0.732$$

$$h_3^2 = \frac{(0.40)(0.48)}{(0.61)} = 0.314$$

وبالنسبة للطريقة الثالثة:

$$h_1^2 = \frac{(0.61) + (0.40)}{2} = 0.550$$

$$h_2^2 = \frac{(0.61) + (0.48)}{2} = 0.545$$

$$h_3^2 = \frac{(0.40) + (0.48)}{2} = 0.440$$

وبالنسبة للطريقة الرابعة نحتاج أولاً إلى

تحديد معكوس المصفوفة وحسب الآتي:

$$adjR = \begin{bmatrix} 0.77 & -0.42 & -0.11 \\ -0.42 & 0.84 & -0.24 \\ -0.11 & -0.24 & 0.63 \end{bmatrix}$$

$$|R| = 0.47$$

$$R^{-1} = \begin{bmatrix} 0.77 & -0.42 & -0.11 \\ -0.42 & 0.84 & -0.24 \\ -0.11 & -0.24 & 0.63 \end{bmatrix} / 0.47 = \begin{bmatrix} 1.64 & -0.89 & -0.23 \\ -0.89 & 1.78 & -0.51 \\ -0.23 & -0.51 & 1.34 \end{bmatrix}$$

وبالتالي يكون لدينا:

$$h_1^2 = 1 - \frac{1}{r_{11}} = 1 - \frac{1}{1.64} = 0.390$$

$$h_2^2 = 1 - \frac{1}{r_{22}} = 1 - \frac{1}{1.78} = 0.438$$

$$h_3^2 = 1 - \frac{1}{r_{33}} = 1 - \frac{1}{1.34} = 0.254$$



ونلاحظ أن هنالك تبايناً واضحاً في قيم الإشتراكيات المقدره بهذه الطرق الأربع وهذا بطبيعة الحال بسبب كون عدد المتغيرات صغير وتتوقع أن لا نجد مثل هذا الحجم من الفروقات عند التعامل مع عدد كبير نسبياً من المتغيرات.

### حساب مصفوفة الارتباط

بطبيعة الحال، نحن بحاجة إلى وجود مصفوفة الارتباط ويمكن تحديد هذه المصفوفة بالشكل الإعتيادي أو من خلال تحميلات العوامل والتي سنوضحها فيما يلي:

لنفترض أنه لدينا حالة أربعة متغيرات وهي  $X_1$  ,  $X_2$  ,  $X_3$  ,  $X_4$  وحصلنا على عاملين معنويين  $F_1$  ,  $F_2$  بالتحميلات التالية من قيم معيارية:

$$F = \begin{matrix} & F_1 & F_2 \\ \begin{matrix} 0.66 & 0.27 \\ 0.79 & 0.33 \\ 0.34 & 0.64 \\ 0.17 & 0.78 \end{matrix} \end{matrix}$$

$$F * F' = \begin{bmatrix} 0.50 & 0.61 & 0.40 & 0.32 \\ 0.61 & 0.73 & 0.48 & 0.39 \\ 0.40 & 0.48 & 0.52 & 0.55 \\ 0.32 & 0.39 & 0.55 & 0.63 \end{bmatrix}$$

وهذه المصفوفة عبارة عن مصفوفة الارتباط للمتغيرات الأربعة مع احلال الإشتراكيات لمواقع العناصر القطرية في المصفوفة. ولكي نثبت ذلك دعنا نحسب قيم الإشتراكيات الأربعة بالشكل الإعتيادي كونها تساوي مجموع مربعات تحميلات العوامل وهي:

$$h_1^2 = (0.66)^2 + (0.27)^2 = 0.50$$

$$h_2^2 = (0.79)^2 + (0.33)^2 = 0.73$$

$$h_3^2 = (0.34)^2 + (0.64)^2 = 0.52$$

$$h_4^2 = (0.17)^2 + (0.78)^2 = 0.63$$

حيث أن الإشتراكية لمتغير ما عبارة عن نسبة التباين لذلك المتغير (باعتباره متغيراً معتمداً) التي تم توضيحها من خلال العوامل المشتركة  $F_1$  و  $F_2$  (باعتبارها متغيرات مستقلة ومتعامدة). وبذلك يمكننا كتابة المعادلات الخطية التالية بمثابة معادلات إحدار خطية:

$$\begin{aligned}
X_1 &= a_{11}F_1 + a_{12}F_2 + d_1U_1 = 0.66F_1 + 0.27F_2 + 0.70U_1 \\
X_2 &= a_{21}F_1 + a_{22}F_2 + d_2U_2 = 0.79F_1 + 0.33F_2 + 0.52U_2 \\
X_3 &= a_{31}F_1 + a_{32}F_2 + d_3U_3 = 0.34F_1 + 0.64F_2 + 0.69U_3 \\
X_4 &= a_{41}F_1 + a_{42}F_2 + d_4U_4 = 0.17F_1 + 0.78F_2 + 0.60U_4
\end{aligned}$$

$$d_i = \sqrt{1-h_i^2} \text{ حيث أن}$$

ومن المعلوم أن مربع معامل الانحدار الجزئي بالوحدات المعيارية هو عبارة عن التباين الموضح للمتغير التابع من خلال المتغير المستقل. وبذلك فإنه من المعادلة الخطية الأولى في أعلاه نجد أن  $(0.66)^2 = 43.66\%$  وهذا يعني أن حوالي 44% من تباين المتغير  $X_1$  تم تفسيره من خلال العامل الأول  $F_1$  وأن  $(0.27)^2 = 7.29\%$  يعني أن حوالي 7% من تباين المتغير  $X_1$  تم تفسيره من خلال العامل الثاني  $F_2$ . وأن مجموعهما حوالي  $43.66 + 7.29 = 51\%$  وهي نسبة التباين الكلية التي تم تفسيرها من خلال العاملين  $F_1$  و  $F_2$ .

كما نلاحظ، ومن خلال تحميلات العوامل، أن العامل الأول هو المحدد المهم بالنسبة للمتغيرين  $X_1$  و  $X_2$  كما أن العامل الثاني هو المحدد المهم بالنسبة للمتغيرين  $X_3$  و  $X_4$ .

والمثال التالي<sup>(8)</sup> يوضح نتائج التحليل العاملي لبيانات دراسة تربوية تتضمن 9 متغيرات كمية تمثل إختبارات تربوية ونفسية واقتصر الإختيار على العوامل الثلاثة الأولى لمعنويتها حيث كانت النتائج لتحميلات هذه العوامل كما في الجدول التالي:

تحميلات العوامل			المتغيرات	
$F_1$	$F_2$	$F_3$	ما تمثله	رمزها
0.83	0.16	0.13	المتناسبات اللفظية	$X_1$
0.02	0.89	- 0.04	سلاسل الأرقام	$X_2$
0.01	0.43	0.71	ذاكرة الأشكال	$X_3$
0.73	0.21	0.01	فهم الأمثال	$X_4$
0.79	0.11	0.15	الإستدلال اللفظي	$X_5$
0.02	0.53	0.68	سلاسل الأشكال	$X_6$
0.02	0.84	0.12	ذاكرة الأرقام	$X_7$
0.07	0.38	0.84	إدراك الأشكال	$X_8$
0.03	0.91	- 0.01	عمليات حسابية	$X_9$

ومن خلال النتائج أعلاه نلاحظ التحميلات العالية لهذه العوامل وما تقابله من متغيرات ضمن كل عامل من هذه العوامل الثلاثة ليستقر الرأي على إعطاء سمة لكل عامل وفقاً لما يتناسب والمتغيرات التي تقابل أكبر التحميلات فيه لتكون كما يلي:

- العامل الأول (  $F_1$  ) والذي أطلق عليه الباحث عامل (الإدراك اللفظي) لتمييزه في:
  - المتناسبات اللفظية 0.83
  - الإستدلال اللفظي 0.79
  - فهم الأمثال 0.73
  
- العامل الثاني (  $F_2$  ) والذي أطلق عليه الباحث عامل (الإدراك الذهني) لتمييزه في:
  - عمليات حسابية 0.91
  - سلاسل الأرقام 0.89
  - ذاكرة الأرقام 0.84
  
- العامل الثالث (  $F_3$  ) والذي أطلق عليه الباحث عامل (الإدراك الشكلي) لتمييزه في:
  - إدراك الأشكال 0.84
  - ذاكرة الأشكال 0.71
  - سلاسل الأشكال 0.68

## المصادر

- 1) Johnson, Dallas E. (1998) "Applied Multivariate Methods for Data Analysts" ; Duxbury Press.
- 2) Boston University, School of Public Health (2013) "Multiple Linear Regression Analysis – seventh examination of the Framingham Offspring Study".
- 3) Rencher, Alvin C. (2002) "Methods of Multivariate Analysis " 2<sup>nd</sup> Ed. ; Wiley Interscience.
- 4) Al-Nsour, Mohannad & Arbaji, Ali (2014) " Obesity and Related Factors Among Jordanian Women of Reproductive Age (Based on Three DHS Surveys, 2002-2012) "; Middle East Health Observatory for Research and Studies (MEHORS).
- 5) Zaiontz, Charles (2017) "Real Statistics Using Excel"; [www.real-statistics.com](http://www.real-statistics.com).
- 6) Manly, Bryan F. J. (2004) " Multivariate Statistical Methods ; A Primer " ; 3<sup>rd</sup> Ed. ; Chapman & Hall/CRC.
- 7) الكبيسي، مائل كامل (1998) " إستخدام الارتباط القويم في دراسة العلاقة بين درجات مواد المفاضلة في القبول ودرجات المواد العلمية للسنة الأولى في كليات المجموعة الطبية" /رسالة ماجستير في الإحصاء/ إشراف أ.د. زياد الراوي/الجامعة المستنصرية.
- 8) العلاق، مهدي محسن اسماعيل (1982) "استخدام التحليل العاملي(طريقة الإمكان الأعظم) في تحليل وتفسير بعض نتائج المسح الجيولوجي في العراق/ رسالة ماجستير في الإحصاء/جامعة بغداد.