

0-3: تمهيد

تعتبر الطرق الإستكشافية الإبتدائية ذات فائده في فهم طبيعه المعقده لعلاقة المتغيرات المتعدده. ويعتبر اسلوب فحص البيانات من اجل التوصل الي هيكل للتجمعات "الطبيعيه" من الأساليب الإستكشافية الهامه فهي تساعدنا في تحديد الأبعاد والتعرف على القيم الشاذه بالإضافة الي اقتراح روابط او فروض تصف العلاقة بين المتغيرات.

لا يضع تحليل التجميع أي فروض تخص عدد المجموعات أو هياكلها مما يجعله اسلوب بدائي لوصف التجمعات. وتتم عملية التجميع بناء على أوجه التماثل أو الإختلاف وتعد مقاييس التماثل أو البيانات في حساب التماثل مدخلات لإسلوب تحليل التجميع.

كثيراً ما نهتم بتجميع عدد من الوحدات في مجموعات متجانسة داخلياً في ضوء مشاهدات ذات بعد P مأخوذة من هذه الوحدات في المتغيرات x_1, x_2, \dots, x_p مثلاً تجميع أشخاص حسب قراءات أخذت منهم إلى مجموعات قد تمثل أجناس أو أعراق مختلفة. كذلك قد يهمننا تجميع عدد من المتغيرات في مجموعات قد ترتبط المتغيرات في كل منها بعامل معين.

1-3 المبحث الأول: مقاييس البيانات الرقمية Measures for numerical data

الهدف الأساسي من تحليل التجميع هو اكتشاف التجمعات الطبيعيه للمفردات أو المتغيرات لذا يجب علينا ان نتوصل أولاً الى مقياس كمي لقياس التوافق "التماثل" بين المفردات.

في أغلب الأحيان يمكننا تجميع المفردات بمجرد النظر في اشكال انتشارها "ذات بعدين او ثلاثه" وذلك على الرغم من عدم وجود تعريف دقيق لمفهوم التجميع الطبيعي .

وللإستفاده من المقدره العقليه على تجميع الأشياء المتشابهه فقد تم التوصل حديثاً الى اساليب بيانيه عديدة لرسم المشاهدات ذات الأبعاد المتعدده في بعدين فقط.

1-1-3: المسافة Distance

عند محاولة تجميع وحدات لابد أن يكون لدينا أولاً مقياس للمسافة بين أي وحدتين مثلاً A و B . وهناك عدة مقاييس للمسافة منها ، المسافة الإقليدية Euclidian Distance والإنحراف Deviance وهو مربع المسافة الإقليدية ومقياس مهالونوبس للمسافة Mahalanobis Distance

والذي يسمي أيضاً المسافة الإحصائية ،مسافة مانهاتن Manhattan Distance ، المسافة القصوى Maximum Distance ، مسافة منكوفسكي Minkowski Distance ، المسافة المتوسطة Average Distance. ففكر المسافة يشير إلى وجود عدد كبير من الأزواج غير المتماثلة.¹

وفيما يلي نتناول هذه المقاييس بإيجاز :

3-1-2: المسافة الإقليدية Euclidean Distance

المسافة الإقليدية معروفة لنا جميعاً فإذا نظرنا إلى النقطة $P = (x_1, x_2)$ في المستوى فإنه وفقاً لنظرية فيثاغورث يمكن إيجاد المسافة الخطية بينها وبين نقطة الأصل $0 = (0,0)$ ويرمز لهذه المسافة بالرمز $d = (0, p)$ كما يلي:

$$d = (0, p) = \sqrt{x_1^2 + x_2^2} \dots \dots \dots (1-3)$$

وعموماً إذا كان لدينا نقطة في فراغ ذي p من الأبعاد ، أي إذا كان لدينا النقطة $p = (x_1, x_2, \dots, x_p)$ فإن المسافة بينها وبين نقطة الأصل $0 = (0,0, \dots, 0)$ هي :

$$d(0, p) = (\sqrt{x_1^2 + x_2^2 + \dots + x_p^2})$$

ويمكن تعميم هذا المفهوم لقياس المسافة بين وحدتين A و B فإذا كانت \underline{X}_A و \underline{X}_B مشاهدات ذات بعد P (في P متغير) على الوحدتين A و B فإن المسافة الإقليدية بينهما تعرف:²

$$\bar{d}_\Delta = \sqrt{(\underline{X}_A - \underline{X}_B)'(\underline{X}_A - \underline{X}_B)} \dots \dots \dots (2-3)$$

3-1-3: الانحراف Deviance

يعرف الانحراف للمسافة بين وحدتين A و B كما يلي:

1. ريشلرلد جونسون ، ديرن وشرن ، التحليل الإحصائي للمتغيرات المتعددة من الوجهة التطبيقية ، تعريب . عبدالمرضي حامد عزام - دار المريخ للنشر، المملكة العربية السعودية - 1418هـ - 1998م. ص (855).

2- عزة أحمد - تجميع مستشفيات ولاية الخرطوم - بحث تكميلي لنيل درجة الماجستير في الإحصاء - جامعة النيلين - يناير 2009م ص(11)

$$(\underline{X}_A - \underline{X}_B)'(\underline{X}_A - \underline{X}_B) = \text{الانحراف}^3 \dots \dots \dots (3-3)$$

أي أنه مربع المسافة الإقليدية.

3-1-4: مسافة مهالانوبس Mahalanobis Distance

يتم إيجاد هذه المسافة عن طريق العلاقة :

$$d_{mah} = \sqrt{(X - Y)\Sigma^{-1}(X - Y)'} \dots \dots \dots (4-3)$$

Σ هي عبارة عن مصفوفة التباينات والتغايرات.⁴

3-1-5: المسافة القصوى Maximum Distance

ويطلق على هذه المسافة أيضاً المسافة الجزئية (sup distance) إذ نوجد أقصى مسافة بين المتغيرات فمثلاً إذا كان لدينا المتغيرين X,Y في d من الأبعاد المسافة القصوى بينهما توجد كما يلي:⁵

$$d_{\max}(x, y) = \max_{1 \leq k \leq d} |x_j - y_j| \dots \dots \dots (5-3)$$

3-1-6: مسافة منكوفسكي Minkowski Distance

المسافة الإقليدية ومسافة مانهاتن والمسافة القصوى عبارة عن ثلاث حالات خاصة من مسافة منكوفسكي وهذه المسافة تعرف بـ . :

$$d_{\min}(x, y) = \left(\sum_{j=1}^d |x_j - y_j|^r \right)^{\frac{1}{r}}, \dots, r \geq 1 \dots \dots \dots (6-3)$$

عندما:

3-عزة أحمد – مرجع سابق ص(12)

4 - ملرتن تي ويلز- وآخرون – تجميع البيانات – سلسلة إصدارات الجمعية الأمريكية الأمريكية (فرجينيا) – جمعية الصناعات والتطبيقات الرياضية(بنسلفينيا) – الولايات المتحدة الأمريكية – سبتمبر-2007م. ص(72)

5 - ملرتن تي ويلز- وآخرون – مرجع سابق ص (72)

$$r = 1, 2 \text{ and } \infty$$

تنتج المسافة الإقليدية مسافة مانهاتن والمسافة القصوى على التوالي.

Average Distance: 7-1-3: المسافة المتوسطة

هذه المسافة يمكن إيجادها من المسافة الإقليدية كآتي:

أفرض أن لدينا المتغيرين x, y ذات بعد d المسافة المتوسطة تصبح⁶:

$$d_{ave}(x, y) = \left(\frac{1}{d_j} \sum (x_j - y_j) \right)^{\frac{1}{2}} \dots \dots \dots (7-3)$$

⁶ - مارتن تي ويلز - وآخرون - مرجع سابق ص (74)

2-3 طرق التجميع والدمج Clustering Method and Amalgamation

وفي هذا المبحث سنتناول بإيجاز المفاهيم الأساسية المرتبطة بتجميع وحدات في ضوء مشاهدات

$$\underline{X}' = [x_1, x_2, \dots, x_p]$$
 مأخوذة منها.

توجد ثلاث طرق مختلفة يمكن استخدامها لعنقدة أو تجميع البيانات وهي:

1- التحليل العنقودي بواسطة الـ k متوسط : يستخدم عند معرفة عدد العناقيد المطلوبة حيث تدخل قبل اجراء التحليل.

2- التحليل العنقودي الهرمي: يستخدم عندما يكون عدد البيانات صغير.

3- التحليل العنقودي ذو الخطوتين : يستخدم عندما يكون عدد البيانات كبير جداً أو عندما تكون البيانات مختلطة ما بين المتغيرات المتصلة والمتغيرات الوصفية.

وسوف يقتصر الباحث على التحليل العنقودي الهرمي والتحليل العنقودي بواسطة الـ K متوسط.

1-2-3 Amalgamation : الدمج

مثلاً تتطلب عملية التجميع تحديد مقياس للمسافة نحتاج أيضاً لتحديد طريقة لدمج العناقيد. وهناك ثلاثة طرق للدمج هي طريقة الربط المفرد وطريقة الربط الكامل وطريقة الربط المتوسط والتي يتم تناولها بإيجاز كما يلي:

2-2-3 طريقة الربط المفرد : Single Linkage Method

تعتبر واحدة من أسهل طرق التحليل العنقودي - عرفت بواسطة فلوريك Florek (1951م) - والربط المفرد يعرف بعدة تعريفات منها الجار الأقرب ، طريقة الحد الأدنى توظف أقرب الجارين لدراسة او قياس الاختلاف بين مجموعتين، ويظهر كما مبين:

الشكل (1-3) الربط المفرد



المصدر : اعداد الدارس باستخدام برنامج MS-Word

في هذه الحالة ومن الشكل (1-3) يمكن ان نفترض ان a هي اقرب وحدة في المجموعة الاولى الي المجموعة الثانية و b هي اقرب وحدة في المجموعة الثانية الي المجموعة الاولى عليه اعتماد المسافة بين a و b لقياس الاختلاف بين المجموعتين يسمي الربط المفرد والصيغة الرياضية لهذه الطريقة تكون كما يلي:

أفرض أن لدينا C_i, C_j و C_k عبارة عن ثلاث مجموعات فالمسافة D بين C_k و $C_i \cup C_j$ يمكن الحصول عليها من صيغة لانس - ويليامز Lance - Williams كالآتي⁷:

$$\begin{aligned}
 & D(C_k, C_i \cup C_j) \\
 &= \frac{1}{2}D(C_k, C_i) + \frac{1}{2}D(C_k, C_j) - \frac{1}{2}|D(C_k, C_i) - D(C_k, C_j)| \\
 &= \min \{D(C_k, C_i), D(C_k, C_j)\} \dots \dots \dots (8-3)
 \end{aligned}$$

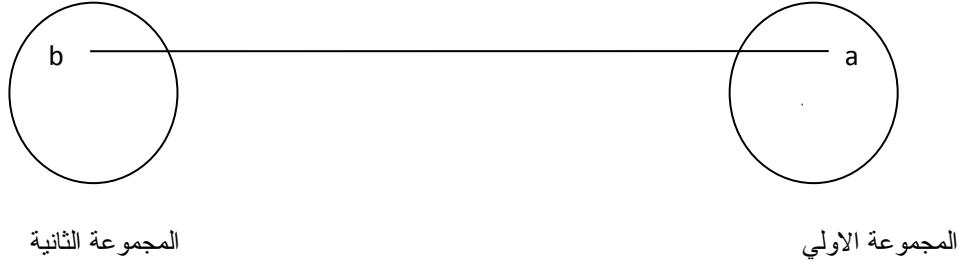
3-2-3: طريقة الربط الكامل Complete Link Method

هذه الطريقة أكثر حذراً إذ تعتبر المسافة بين أي مجموعتين هي المسافة بين أبعد وحدتين فيهما كما يلي:

الشكل (2-3) الربط الكامل⁸

7 - مارتن تي ويلز وآخرون- تجميع البيانات - مرجع سابق ص (118)

8 - د. ريتشارد جونسون ، ديرن وشرن، مرجع سابق ص (866)



المصدر : اعداد الدارس باستخدام برنامج MS-Word

فلنفرض ان a هي ابعاد وحدة في المجموعة الاولى من المجموعة الثانية و b هي ابعاد وحدة في المجموعة الثانية من المجموعة الاولى عليه يتم استخدام هاتين الوحدتين لقياس الاختلاف بين المجموعتين.

والصيغة الرياضية لهذه الطريقة تكون كما يلي:

أفرض أن لدينا C_i, C_j & C_k عبارة عن ثلاث مجموعات فالمسافة D بين $C_i \cup C_j$ و C_k يمكن الحصول عليها من صيغة لانس - ويليامز Lance - Williams كالآتي:⁹

$$\begin{aligned}
 & D(C_k, C_i \cup C_j) \\
 &= \frac{1}{2} D(C_k, C_i) + \frac{1}{2} D(C_k, C_j) + \frac{1}{2} |D(C_k, C_i) - D(C_k, C_j)| \\
 &= \max \{D(C_k, C_i), D(C_k, C_j)\} \dots \dots \dots (9-3)
 \end{aligned}$$

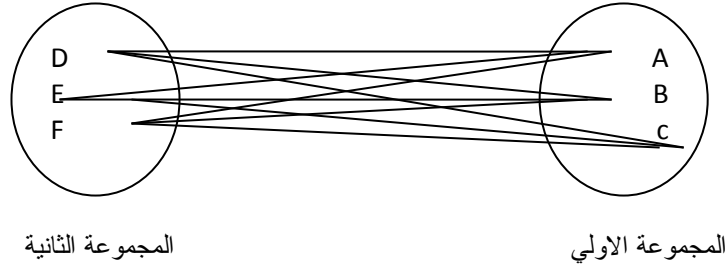
4-2-3 طريقة الربط المتوسط: Average Linkage Method

طريقة الربط المتوسط تنظر للمسافة بين مجموعتين على أنها متوسط المسافة بين جميع الأزواج التي ينتمي أحد عناصرها إلى إحدى المجموعتين بينما ينتمي العنصر الآخر إلى المجموعة الأخرى ، كما يلي:

الشكل (3-3) الربط المتوسط¹⁰

⁹ - مارتن تي ويلز وآخرون- تجميع البيانات - مرجع سابق ص (121).

¹⁰ - د. ريتشارد جونسون ، ديرن وشرن، مرجع سابق ص (866) .



المصدر : اعداد الدارس باستخدام برنامج MS-Word

في هذه الحالة اذا افترضنا (A,B,C) تنتمي للمجموعة الاولى و (D,E,F) تنتمي الي المجموعة الثانية حيث يتم حساب المسافات بين الوحدات في المجموعة الأولى مع المجموعة الثانية لإيجاد متوسط المسافة لاستخدامها في قياس الاختلاف بين المجموعتين

والصيغة الرياضية لهذه الطريقة تكون كما يلي:

أفرض أن لدينا $C_i, C_j & C_k$ عبارة عن ثلاث مجموعات فالمسافة D بين C_k و $C_i \cup C_j$ يمكن الحصول عليها من صيغة لانس – ويليامز Lance – Williams كالآتي:¹¹

$$D(C_k, C_i \cup C_j) = \frac{|C_i|}{|C_i| + |C_j|} D(C_k, C_i) + \frac{|C_j|}{|C_i| + |C_j|} D(C_k, C_j) \dots \dots \dots (10-3)$$

أفرض أن C, C' عبارة عن مجموعتين غير خاليتين عليه يمكن إيجاد المسافة بطريقة الربط المتوسط الآتي:

$$D(C, C') = \frac{1}{|C||C'|} \sum_{x \in C, y \in C'} d(x, y) \dots \dots \dots (.11-3)$$

أفرض أن C_1, C_2, C_3 عبارة عن ثلاث مجموعات غير خالية عليه أفرض أن :

$$D(C_i, C_j) = \frac{1}{n_i n_j} \sum (C_i, C_j), 1 \leq i \leq j \leq 3 \dots \dots \dots (12-3)$$

حيث:

¹¹ -مارتن تي ويلز وآخرون- تجميع البيانات – مرجع سابق ص (123)

$C_i, C_j \dots \sum(C_i, C_j)$ والمجموع الكلي لمسافات المجموعات $n_i = |C_i|$ و $n_j = |C_j|$

هذا يعني أن:

$$\sum(C_i, C_j) = \sum_{x \in C_i, y \in C_j} d(x, y)$$

ومن المعادلات (11-3) و (12-3) نستنتج:

$$\begin{aligned} D(C_k, C_i \cup C_g) &= \frac{n_2}{n_2 + n_3} D(C_1, C_2) + \frac{n_3}{n_2 + n_3} D(C_1, C_3) \\ &= \frac{n_2}{n_2 + n_3} \cdot \frac{1}{n_1 n_2} \sum(C_1, C_2) + \frac{n_3}{n_2 + n_3} \cdot \frac{1}{n_1 n_2} \sum(C_1, C_3) \\ &= \frac{1}{n_1(n_2 + n_3)} \sum(C_1, C_2 \cup C_3) \end{aligned}$$

بعد ذلك:

$$\sum(C_1, C_2) + \sum(C_1, C_3) = \sum(C_1, C_2 \cup C_3)$$

وهذا يتحقق أيضاً للمعادلة (11-3)

هناك عدة طرق للتجميع توفرها النظرية الإحصائية و نتناول فيما يلي بشكل موجز إثنين من أهم هذه الطرق وهي طريقة التجميع المتدرجة وطريقة الـ (K) متوسط واللتين سيتم تطبيقهما في التجميع وسنفترض أن المطلوب هو تجميع وحدات وليس متغيرات وهي الحالة التي تهتمنا

3-2-5 طرق التجميع المتدرجة: Clustering Methods

هناك نوعان من طرق التجميع المتدرجة :

1- طريقة الفصل المتدرجة Divisive hierarchical methods ، حيث نبدأ بمجموعة تتضمن جميع العناصر ثم تقسم هذه المجموعات إلى مجموعتين فرعيتين بحيث تكون العناصر الموجودة في مجموعة منها "بعيدة" عن العناصر الموجودة في المجموعة الأخرى ، يتم بعد ذلك تقسيم كل

من هاتين المجموعتين إلى مجموعات فرعية غير مماثلة ، نستمر في ذلك حتى يكون لدينا عدداً من المجموعات الفرعية مساوياً لعدد العناصر ، أي حتى يكون كل عنصر مجموعة بنفسه.¹²

2- طريقة التجميع المتدرجة Agglomerative Hierarchical Methods وهي إحدى طرق الإدماج "linkage methods" المتتالية التي تلائم تجميع المفردات وكما تلائم تجميع المتغيرات وهذا لا يتحقق بالنسبة لجميع طرق الإدماج المتدرجة الأخرى ، وفيما يلي خطوات تنفيذ الطرق التجميع المتدرجة لإدماج المجموعات عند وجود N من العناصر (المفردات أو المتغيرات) وتسمى هذه الخطوات بالطريقة العامة:¹³

1- نبدأ بعدد N من المجموعات ، كل مجموعة بها عنصر واحد ونوجد مصفوفة المسافة (أو مصفوفة قيم معامل التماثل المستخدم) $D = d_{ik}$ وهي مصفوفة متماثلة أبعادها $(N \times N)$

2- نبحث في مصفوفة المسافة عن أقرب زوج من المجموعات (الزوج الأكثر تماثلاً). نفترض أن d_{uv} تشير الي المسافة بين الزوج الأكثر تماثلاً U, V .

3- ندمج المجموعة U مع المجموعة V ، ونستخدم الرمز $(U V)$ للإشارة الى المجموعة الجديدة ، ونعدل من عناصر مصفوفة المسافة على النحو التالي:

(a) نحذف الصفوف والأعمدة المناظرة للمجموعتين U, V .

(b) نضيف صفاً وعموداً جديدين يعطيان المسافة بين المجموعة الجديدة $(U V)$ والمجموعات الأخرى

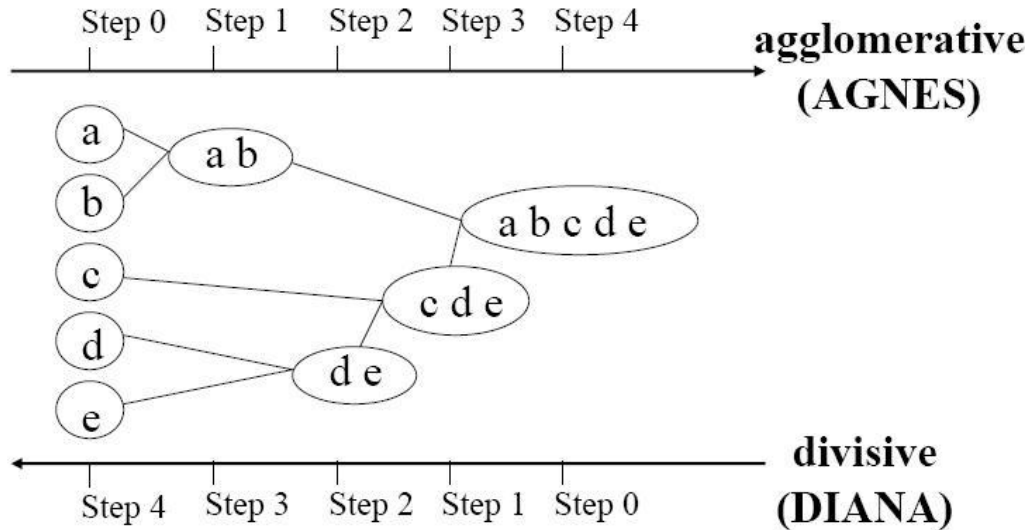
4- نكرر الخطوتين (a) ، (b) عدد $N-1$ من المرات (في النهاية تتجمع جميع العناصر في مجموعة واحدة) نقوم بتسجيل هوية المجموعات التي أدمجت والمسافات (أو قيم معامل التماثل المستخدم) التي تم عندها الإدماج.

¹² www.norusis.com/pdf/spc-v13.pdf.

¹³ - د. ريتشارد جونسون ، ديرن وشرن ، مرجع سابق ص (867).

ويمكن عرض نتائج هاتين الطريقتين بيانياً في فراغ ذي بعدين في شكل بياني يعرف باسم "الندوجرام" Dendrogram.

الشكل (4-3) طريقة التجميع المتدرجة وطريقة الفصل المتدرج¹⁴



1. المصدر : فرح عبدالله محمد , (2010) ،"تصنيف الولايات السودانيه ذات الخصائص الديمغرافية المتشابهة باستخدام التحليل العنقودي للعام 2002 م" , جامعة السودان للعلوم والتكنولوجيا.

نلاحظ في الشكل اعلاه في حالة التجميع المتدرج في الخطوة Step0 اعلي الشكل ان كل مفردة تعتبر عنقود لوحدها وعند اجراء عملية التجميع ننقل الي الخطوة Step1 لنجد ان العنقودين a,b تم دمجها في عنقود واحد نسبة لوجود تقارب بينهما اما البقية (c,d,e) تظل في عناقيد لوحدها ، وفي الخطوة Step2 نلاحظ ان العنقودين (d,e) تم جمعها معا في عنقود واحد نسبة لدرجة التقارب بينهما والبقية (a,b) في عنقود و (c) في عنقود لوحدها ، في الخطوة Step3 تم دمج العنقود (c) مع العنقود (d,e) ليصير لدينا عنقودين فقط هما (c,d,e) و (a,d) اما في الخطوة Step4 تم تجميع اخر عنقودين من الخطوة Step3 في عنقود واحد كبير يجمع جميع العناقيد في الخطوة Step0 في عنقود واحد وبذلك تكون طريقة التجميع المتدرج وصد وصلت الي آخر خطوة.

كي -بينق زاهنق - وسائل التحليل العنقودي في البيانات البحثية -رسالة دكتوراة -جامعة ماكوري -سديني - أستراليا 2007م.14.

اما في خطوة الفصل المتدرج اسفل الشكل ففي الخطوة Step0 نجد ان هذه الطريقة تعتبر كل العناصر عبارة عن عنقود واحد كبيرة وعند الانتقال للخطوة الاولى في الطريقة Step1 نجد انه تم فصل العنصرين (a,b) عن العنقود الكبير لتصبح عنقود لوحدها والبقية تصبح (c,d,e) وفي عند الانتقال الي الخطوة Step2 نجد ان العنصر (c) انفصل عن العنقود الكبير واصبح عدد العناقيد 3 وهي (a,b) و (c) و (d,e) وفي الخطوة Step3 نجد ان العنقود (d,e) قد انفصل الي عنقودين هما (d) و (e) فاصبح عدد العناقيد اربعة وهي (a,b) و (c) و (d) و (e) وفي الخطوة الختامية لطريقة الفصل المتدرج تكون كل العناصر قد اصبحت في عناقيد كل علي حده.

ونلاحظ ان كلا الطريقتين تعملان عكس بعضهما الاولى تجمع بالتدرج والاخري تفصل

بالتدرج

6-2-3 مصفوفة القرابة¹⁵: Proximity Matrix

هي مصفوفة متماثلة مدخلاتها عبارة عن مربع المسافة الإقليدية لكل زوج من الحالات الجدول رقم (3-1) يمثل نموذج لمصفوفة القرابة ، ونلاحظ أن قيم القطر الرئيسي (0) مما يعني أنه لا توجد فروق بين الحالة ونفسها وكذلك القيمة (2.6) تعني أن المسافة بين (E) و (D) هي عبارة عن أصغر قيمة في المصفوفة مما يعني أنهما متماثلين بدرجة عالية ، وأيضاً نجد القيمة (12.8) تعني أن المسافة بين (F) و (G) هي أكبر المسافات الإقليدية مما يعني أنهما على درجة عالية من عدم التماثل.

1-د. محفوظ جودة ، التحليل الإحصائي المتقدم باستخدام SPSS دار وائل للطباعة والنشر - عمان الأردن الطبعة الأولى 2008م ص(100)

جدول (1-3) مصفوفة القرابة

samples	A	B	C	D	E	F	G
A	0	7.5	10.4	6.4	4.7	6.6	8.8
B	7.5	0	12.4	9.3	5.6	9.2	8.8
C	10.4	12.4	0	5.4	3.7	9.1	10.6
D	6.4	9.3	5.4	0	2.6	5.4	10.7
E	4.7	5.4	3.7	2.6	0	4.7	5.4
F	6.6	9.2	9.8	5.7	4.8	0	12.8
G	8.8	8.8	10.6	10.7	5.4	12.8	0

المصدر: المرجع نفسه، ص[19]

7-2-3: شكل الأنواع الجليدية: The Icicle Plot

هو عبارة عن شكل يبين الطريقة التي تتم بها دمج الحالات في كل خطوة من خطوات التحليل العنقودي ، ويسمي شكل الالواح الجليدية لأن لها نفس شكل الألواح الجليدية المتشكلة من مساقط الماء.¹⁶ ولكي نتمكن من فهمه لابد أن نقرأه من أسفل إلى أعلى.

شكل (3-5) الألواح الجليدية¹⁷

		Case																	
		5	8	6	7	3	2	4	1										
عدد العناقيد	1	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	x	
	2	x	x	x	x	x	x	x	x	x	x	x	x		x	x		x	
	3	x	x	x	x		x	x	x	x	x	x	x		x	x		x	
	4	x	x	x	x		x		x	x	x	x	x		x	x		x	
	5	x	x	x	x		x		x	x	x	x		x		x	x		x
	6	x	x	x	x		x			x	x	x		x		x	x		x
	7	x		x	x		x			x	x	x		x		x	x		x
	8	x		x	x		x			x	x	x		x		x			x

المصدر: المرجع نفسه ، ص [20]

3-2-8: خطوات التجميع¹⁸: The Agglomeration Schedule:

شكل الألواح الجليدية لا يمكن من معرفة أصغر مسافة إقليدية بين الحالتين أو الحالات التي تم دمجها في عنقود لذلك علينا اللجوء لما يسمى بخطوات التجميع جدول (3-3) فنجد في هذه الخطوة عمود يسمى المعامل Coefficient به قيم المسافات الإقليدية قيم (التمائل) التي أستخدمت لتكوين العناقيد. أما عمود (Stage Cluster First appears) يبين لنا الخطوة التي يظهر فيها العنقود لأول مرة .

¹⁷ www.nourisis.com/pdf/spc/-v13pd-

¹⁸ - محفوظ جودة ، التحليل الإحصائي المتقدم باستخدام SPSS مرجع سابق ص (100)

جدول (2-3) خطوات التجميع

المرحلة	الانضمام للعناقيد		المعاملات	الظهور اول مرة		المرحلة التالية
	العنقود 1	العنقود 2		العنقود 1	العنقود 2	
1	4	5	2.566	0	0	4
2	1	11	2.694	0	0	5
3	3	10	3.903	0	0	4
4	3	4	4.023	3	1	6
5	1	8	4.279	2	0	11
6	3	15	4.336	4	0	13
7	2	9	5.205	0	0	14

المصدر : المرجع نفسه، ص [21]

9-2-3: الدينودوجرام The Dendrogram

هو عبارة عن شكل بياني يوضح المسافات التي تم فيها تكوين العناقيد ويظهر ذلك في الشكل رقم (7-3) ويقرأ من الشمال إلى اليمين الخطوط العمودية تشير إلى العناقيد المدمجة فموقع الخط يشير إلى المسافة التي تمت فيها عملية الدمج وتكوين العناقيد ، فأول عمود رأسي يقابل أصغر مسافة إقليدية ولكي نفهم شكل الدينودوجرام لابد لنا ان نعرف مايسمى بالـ n-tree :

أفرض ان $D=(x_1,x_2,\dots,x_n)$ وأن J مجموعة جزئية من D تحقق الشروط التالية:

$$J \in D \quad .1$$

$$\text{Empty set } \Phi \in J \quad .2$$

For all J .3

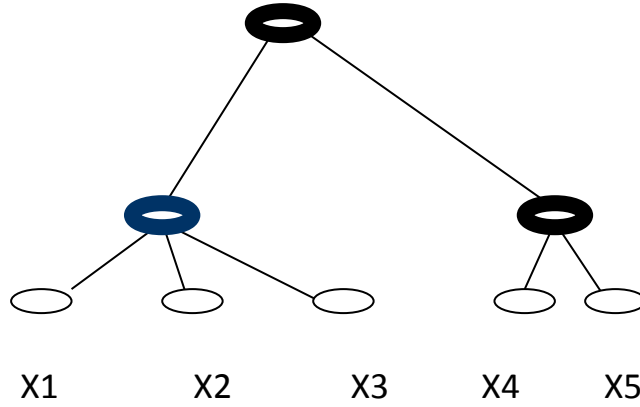
$$i=1,2,3,\dots,n \{x_i\} \in J$$

if $A, B \in B \cap J$, .4

$$\text{then } A \in \{\Phi, A, B\}.$$

الـ n-tree في الشكل (6-3)العقد الطرفية المرسومة بدائرة فقط من غير تظليل تبين البيانات

بصورة مفردة ،أما العقد المرسومة بدوائر مظلة المجموعات أو العناقيد.



X1 الي X5 = الوحدات التي تتكون منها العناقيد

المصدر : المرجع نفسه ، ص [22]

من ذلك نستطيع أن نصف شكل الديندوجرام بأنه شجرة تعرف با لـ (5-tree) يحقق الشرط التالي:

$$h(A) \leq h(B) \Leftrightarrow A \subseteq B \dots \dots \dots (13-3)$$

لكل مجموعة جزئية من البيانات A و B بشرط $A \cap B \neq \Phi$

عندما $h(A) \& h(B)$ تعرف الارتفاعات لـ (A,B) على التوالي .

من الشكل أعلاه أفرض أن (h_{ij}) عبارة عن إرتفاع زوج البيانات (X_i, X_j) فالقيمة الأصغر لـ (h_{ij}) تمثل أكبر معامل تماثل أو تشابه بين (X_i, X_j) وكذلك أكبر قيمة لـ (h_{ij}) تمثل أصغر قيمة لمعامل التماثل أو تشابه بين $((X_i, X_j))$.

الإرتفاع في الديندوجرام يحقق الشرط التالي:

$$h_{ij} \leq \max \{h_{ik}, h_{jk}\} \forall j, k \in \{1, 2, \dots, n\} \dots \dots \dots (14-3)$$

ويمكن عرض شكل الديندوجرام رياضياً بالمعادلة التالية:

19 - مارتن تي ويلز وآخرون- تجميع البيانات - مرجع سابق ص (112)

$$C: [(0, \infty) \rightarrow E(D)]$$

التي تحقق:

$$c(h) \subseteq c(h') \rightarrow \text{if } h \leq h'$$

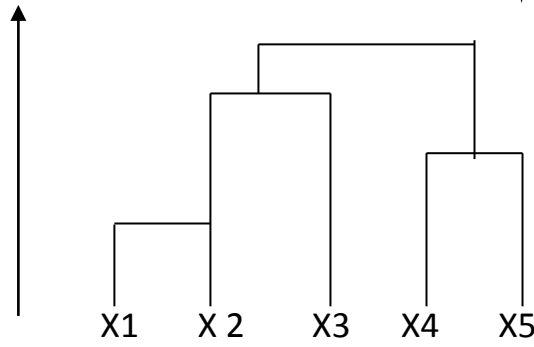
$c(h)$ is eventually in $D \times D$

$$c(h + \delta) = c(h) \text{ for some small } \delta > 0$$

عندما يكون لدينا المجموعة D و $E(D)$ عبارة

عن مجموعة من العلاقات المتكافئة في D ويمكن تمثيل ذلك بالشكل التالي:

شكل (7-3) ديندوجرام لخمس مجموعات²⁰



المصدر : المرجع نفسه ، ص [23]

X1 الي X5 = الوحدات التي تتكون منها العناقيد

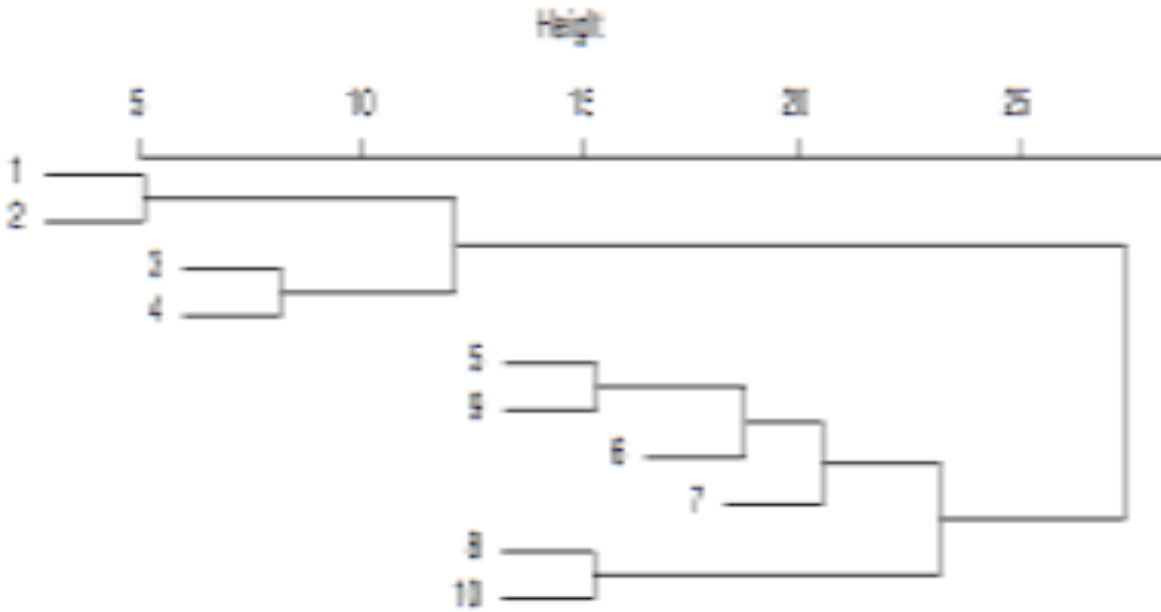
المعادلة (c) تحوي المعلومات في شكل الديندوجرام ويمكن بيان ذلك من الشكل (8-3) كالاتي:

$$c(h) = \begin{cases} \{(i, i) : i = 1, 2, 3, 4, 5\} \text{ if } 0 \leq h < 1, \\ \{(i, i) : i = 3, 4, 5\} \cup \{(i, j) : i, j = 1, 2\} \text{ if } 1 \leq h < 2 \\ \{(3, 3)\} \cup \{(i, j) : i, j = 1, 2\} \cup \{(i, j) : i, j = 4, 5\} \text{ if } 2 \leq h < 3 \dots \dots \dots (15-3) \\ \{(i, j) : i, j = 4, 5\} \cup \{(i, j) : i, j = 1, 2, 3\} \text{ if } 3 \leq h < 4 \\ \{(i, j) : i, j = 1, 2, 3, 4, 5\} \text{ if } 4 \leq h \end{cases}$$

شكل (8-3) الديندوجرام²¹

²⁰ - مارتن تي ويلز وآخرون- تجميع البيانات - مرجع سابق ص (112)

²¹ www.uga.edu/strata/software/pdf/cluster



المصدر : المرجع نفسه، ص [24]

3-2-10: طريقة K من المتوسطات

هي إحدى طرق التجميع غير المتدرجة Nonhierarchical Clustering Methods لتجميع المفردات وليس المتغيرات في K من المتوسطات، اقترح MacQueen عام 1976م اسم K من المتوسطات (K-means) لوصف طريقته وتعتبر من أبسط طرق العنقدة وأكثرها شيوعاً وذات كفاءة عالية ، وتستخدم ما يسمى النقاط المركزية (Centroids) التي تضع كل مفردة في المجموعة التي يكون وسطها الحسابي أقرب لها ، وتتكون هذه الطريقة في أبسط صورها من ثلاث خطوات²² هي:

1- تقسيم المفردات إلى K من المجموعات الأولية.

2- ضع كل مفردة من المفردات الموجودة في المجموعة التي يكون وسطها الحسابي أقرب ما يكون لها (عادة ماتستخدم المسافة الإقليدية لحساب المسافة سواء تم ذلك باستخدام المشاهدات الفعلية

²² - www.ranger.uta.edu/chqding/paper/k-means.com.

أو المشاهدات المعيارية) نعيد حساب الوسط الحسابي للمجموعة التي أضيفت إليها المفردة الجديدة وللمجموعة التي فقدت منها المفردة.

3- نكرر الخطوة الثانية إلى أن نتوقف عملية توزيع المفردات على المجموعات ، وبدلاً من أن نبدأ في الخطوة الأولى بتجزئة جميع المفردات إلى K من المتوسطات من المجموعات الأولية ، فإنه يمكننا البدء بتحديد k ، (نقاط الأساس) ثم نقوم بتنفيذ الخطوة الثانية.

ويعتمد التوزيع النهائي للمفردات على المجموعات إلى حد ما ، على التقسيم الأولي المستخدم أو على الاختيار الأولي لنقاط الأساس .

توجد بعض الحجج القوية التي تؤيد عدم تحديد عدد المجموعات K مسبقاً²³ ومن بينها:

1- إذا وقعت نقطتين أو أكثر من نقاط الأساس ، دون تعمد ، في مجموعة واحدة، فإنه من الصعب التمييز بين المجموعات الناجمة عن ذلك.

2- قد يؤدي وجود مشاهدة شاذة في الحصول على مجموعة واحدة على الأقل متناثرة المفردات.

3- وحتى إذا تكون المجتمع من K من المجموعات ، فإن طريقة المعاينة المستخدمة يمكن أن تؤدي إلى عدم احتواء العينة على بيانات من المجموعة الأصغر. في هذه الحالة، يؤدي توزيع البيانات على K من المجموعات إلى وجود مجموعات ليس لها معنى.

المثال التالي يوضح طريقة التجميع بواسطة الـ (K) متوسط:

أفرض أننا نستطيع قياس المتغيرين X_1 ، X_2 لكل مفردة من المفردات كما في الجدول التالي:

جدول (3-3) قيم المتغيرين X_1 ، X_2

²³د. ريتشارد جونسون ، ديرن وشرن، مرجع سابق ص (893) .

المفردة	X1	X2
A	5	3
B	-1	1
C	1	-2
D	-3	-2

المصدر : (فرح عبدالله ، [2010])

الهدف هو تقسيم هذه المفردات لمجموعتين أي $K=2$ بحيث يكون قرب المفردات الموجودة في أي مجموعة من بعضها البعض أكبر من قربها من مفردات المجموعات الأخرى.

- نقوم بتقسيم المفردات عشوائياً لمجموعتين (AB) والمجموعة (CD).

- نقوم بحساب إحداثيات المركز (\bar{X}_1, \bar{X}_2) لكل مجموعة كالاتي:

جدول (4-3) احداثيات مركز X_2, X_1

المجموعة	\bar{X}_1	\bar{X}_2
AB	2	2
CD	-1	-2

المصدر : المرجع نفسه ، ص [26]

نقوم بحساب المسافة الإقليدية بين كل مفردة وبين إحداثيات نقطتي المركز بين كل مفردة وبين إحداثيات نقطتي المركز ونعيد توزيع كل مفردة على المجموعة الأقرب لها ، وإذا انتقلت مفردة ما من مجموعة إلى أخرى يجب تعديل نقاط المركز المناظرة قبل مواصلة العمل كالاتي:

$$d^2(A,(AB)) = (5-2)^2 + (3-2)^2 = 10$$

$$d^2(A,(CD)) = (5+1)^2 + (3+2)^2 = 61$$

حيث أن المفردة A أقرب إلى المجموعة (AB) منها للمجموعة (CD) فتبقي في مجموعتها.

وبمواصلة حساب المسافات المربعة :

$$d^2(B,(AB)) = (-1-2)^2 + (1-2)^2 = 10$$

$$d^2(B,(CD)) = (-1+1)^2 + (1+2)^2 = 9$$

وبالتالي نقوم بوضع المفردة B في المجموعة (CD) لنحصل على المجموعة (BCD) كما نحصل على نقاط المركز الجديدة التالية:

جدول (5-3) احداثيات المركز الجديدة لـ X_2, X_1

المجموعة	\bar{X}_1	\bar{X}_2
A	5	3
CDB	-1	-1

المصدر: المرجع نفسه ، ص [27]

ومرة أخرى نقوم بفحص كل مفردة لمعرفة إمكانية إعادة توزيع المفردات وبحساب المسافات الإقليدية نحصل على مايلي:

جدول (6-3) توزيع المجموعات النهائي²⁴

المجموعة	A	B	C	D
A	0	40	41	89
CDB	52	4	5	5

المصدر : المرجع نفسه ، ص [27]

من ذلك يكون كل مفردة قد تم وضعها في المجموعة التي يكون مركزها أقرب ما يكون لها بالتالي تتوقف العملية أي أننا نحصل في النهاية على المجموعتين A و (BCD).

²⁴ ريتشارد جونسون ، ديرن وشرن، مرجع سابق ص (887)