

المحاضرة الأولى: أدوات تقسيم السوق - التحليل العنقودي -**1. تمهيد:**

يؤدي التحليل العنقودي دورا مهما في الكثير من الميادين مثل: علم النفس، إحصاء، العلوم الاجتماعية، بيولوجيا، ... الخ، وفي التسويق تقوم الشركات بجمع بيانات كثيرة عن العملاء الحاليين والمحتملين، وهنا يساعدها التحليل العنقودي في تقسيم العملاء إلى عدد صغير من المجموعات، بهدف استخدامها في تحليلات إضافية وأنشطة التسويق، فمثلا تصنيف الأشخاص، المناطق والمنتجات إلى مجاميع عنقودية ذات صفات متشابهة، يساعد في تحديد قطاعات السوق والشركات المنافسة.

إذن التحليل العنقودي هي أداة إحصائية تستخدم في مجال الأبحاث التسويقية لوضع المستهلكين الذين يتشابهون في نفس الصفات ضمن مجموعات محددة، مما يسهل على المؤسسة تصميم المنتجات بشكل يتناسب مع كل مجموعة، كما يسهل عليها وضع الخطط التسويقية المناسبة لكل مجموعة.

كما يستخدم التحليل العنقودي لتجميع أي مجموعة من المتغيرات ضمن مجموعات متجانسة ومنفصلة، ويطبق هذا الاتجاه في مراجعة وتنقيح الاستبيانات بناءً على إجابات المستجيبين على استمارات الاستبيان، حيث أن تجميع الأسئلة حسب المتوسطات بواسطة التحليل العنقودي يساعدنا في التعرف على الأسئلة الضعيفة وإعادة النظر فيها، وهذا يزيد حظوظ معدل الإجابات الجيدة بالنسبة لإجمالي أسئلة الاستبيان.

2. تعريف التحليل العنقودي:

التحليل العنقودي هو: "أحد الأساليب الإحصائية الرياضية لتقسيم عناصر المجتمع المدروس إلى عدة مجموعات متعلقة ومتجانسة داخليا (متشابهة) ومتباينة خارجياً عن بعضها البعض، أي أنه يهدف إلى جعل تباين العناصر داخل كل مجموعة أصغر ما يمكن، وجعل التباين بين المجموعات (بين مراكزها) أكبر ما يمكن".

التحليل العنقودي عبارة عن: "إجراءات تهدف إلى تصنيف مجموعة حالات Cases (أو متغيرات Variables) بطرق معينة وترتيبها داخل عناقيد Clusters، بحيث تكون الحالات المصنفة داخل عنقود معين متجانسة فيما يتعلق بخصائص محددة، وتختلف عن حالات أخرى موجودة في عنقود آخر".

3. أنواع التحليل العنقودي:

بصورة عامة يتفرع التحليل العنقودي إلى نوعين أساسيين، كل منهما ينقسم إلى طرق فرعية عديدة، هذين النوعين هما:

أ. التحليل العنقودي الهرمي Hierarchical

هذه الطريقة لا تتطلب المعرفة المسبقة لعدد العناقيد المسبقة التي سيتم تجميع العناصر على أساسها، وهي تناسب العينات الصغيرة نسبياً.

ب. التحليل العنقودي غير الهرمي Non- Hierarchical

في هذه الطريقة يكون عدد العناقيد محدد مسبقاً أو معين كجزء من إجراءات التحليل العنقودي، ومن بين أهم أشكالها: "طريقة ك-متوسطات (K-means)"، والتي تعتمد في التصنيف على تجميع العناصر المدروسة في k عنقود أولي، ثم حساب مراكز (النقطة المتوسطة) هذه العناقيد، والمسافات بين العناصر والمراكز، ثم نقل العناصر إلى العناقيد أقرب إليها.

ويعتبر أسلوب التحليل العنقودي الهرمي من الأساليب المفضلة في التحليل العنقودي، لأنه يعتمد على أسس بسيطة، ويعمل على عنقود مفردات العينة (n مفردة)، وبشكل متتالي، ضمن m عنقوداً، بواسطة دمج المفردات المتقاربة ضمن مجموعات متعلقة تسمى عناقيد، وبحيث يكون العنقود الأول C_1 أصغر وأبسط العناقيد، ويكون العنقود الأخير C_m أعدها وأشملها (لأنه يضم جميع المفردات والعنقيد الجزئية)، على أن يتألف كل عنقود من عدة

مجموعات متقاربة ومرتبطة مع بعضها بواسطة علاقات تحقق شروط التقارب المفضلة (حسب المتغير أو الصفة المدروسة).

4. أساليب التحليل العنقودي الهرمي:

يستخدم التحليل العنقودي الهرمي لعنقدة مفردات العينة أسلوبيين عمليين ما:

أ. أسلوب التجميع Agglomerative Technique

يفترض هذا الأسلوب من البداية أن كل مفردة من مفردات العينة تشكل عنقوداً خاصاً بها، ثم يتم دمج أي مفردتين متقاربتين في عنقود خاص (أول)، ثم نضيف إليهما أي مفردة ثالثة متقاربة مع ذلك العنقود فيشكل لدينا عنقود ثانٍ، وهكذا نتابع إضافة المفردات واحدة بعد الأخرى إلى بعضها أو إلى العناقيد السابقة، مع تحديد العلاقات بينها ضمن العناقيد، حتى نحصل على العنقود الأخير، الذي يضم جميع مفردات العينة (n مفردة) مع العلاقات التي ترتبط بينها، ويعتمد هذا الأسلوب على مصفوفة التقارب بين مفردات العينة حسب المسافات المحسوبة.

ب. أسلوب التقسيم Divisive Technique

يفترض هذا الأسلوب من البداية أن جميع مفردات العينة (n مفردة) تشكل عنقوداً واحداً شاملاً، ثم تتم تجزئته إلى عناقيد جزئية متباينة تتضمن عدداً أقل من المفردات، وبعد فرز هذه العناقيد وتحديد العلاقات بينها، تتم تجزئتها إلى عناقيد أصغر فأصغر، وتتابع هذه العملية حتى يتكون عنقود خاص لكل مفردة من مفردات العينة أو نتوقف عن التقسيم عند حد معين.

5. متطلبات التحليل العنقودي التجميعي:

إن أسلوب التجميع والتقسيم يعتمدان على عدد عناصر العينة المدروسة (n مفردة أو مشاهدة)، وعلى طبيعة المتغيرات المستقلة X_1, X_2, \dots, X_p المستخدمة في عملية العنقدة، على أن تنظم بيانات مفردات العينة حسب المتغيرات النظامية (معيارية أو ثنائية)، وتوضع في جدول مناسب كما يلي:

نموذج جدول البيانات اللازمة للتحليل العنقودي

المتغيرات/المفردات	X_1	X_2	$X_3 \dots$	$X_i \dots$	X_p
1	X_{11}	X_{12}	X_{13}	X_{1i}	X_{1p}
2	X_{21}	X_{22}	X_{23}	X_{2i}	X_{2p}
.					
.					
J	X_{j1}	X_{j2}	X_{j3}	X_{ji}	X_{jp}
.					
.					
n	X_{n1}	X_{n2}	X_{n3}	X_{ni}	X_{np}

إذا كانت المتغيرات X_1, X_2, \dots, X_p كمية، فإننا نقوم بحساب عناصر مصفوفة التباعد (Dissimilarity(D) ، وهي عبارة عن المسافات التي تفصل بين مفردات العينة. أما إذا كانت المتغيرات نوعية أو مختلطة، فإننا نقوم بحساب عناصر مصفوفة أخرى مصفوفة التشابه أو التقارب (Similarity(S)، وهي أوزان التكرارات التي تقابل المفردتين المتشابهتين (j, k).

5. حساب مصفوفة التباعد (حالة متغيرات كمية):

تتألف مصفوفة التباعد من قيم المسافات بين أزواج مفردات العينة، وتحسب من قيم المتغيرات المستخدمة في عملية العنقدة، لذلك نفترض أنه لدينا n مفردة هي $1, 2, 3, \dots, j, \dots, n$ ، نريد تصنيفها ضمن عناقيد حسب قيم المتغيرات المؤثرة عليها، والتي سنرمز لها بـ X_1, X_2, \dots, X_p ، وسنستخدم قيم هذه المتغيرات لحساب عناصر مصفوفة التباعد أو مصفوفة المسافات، والتي سنرمز لها بـ D ونكتبها كمايلي:

- فإذا كان المتغير X ثنائياً، أي يتألف من حالتين فقط (نجاح وفشل)، فإننا نفترض أنه يأخذ القيمة (1) إذا تحققت حالة النجاح، ويأخذ القيمة (0) إذا تحققت حالة الفشل.

- أما إذا كان للمتغير X أكثر من حالتين، فإننا نعتبر كل حالة مستقلة عن الحالات الأخرى، ونعرف عليها متغيرات ثنائية جديدة ($X_1', X_2' \dots X_p'$)، نفترض أن كل متغير جديد X_i' ، يأخذ القيمة (1) عندما تتحقق حالة معينة، ويأخذ القيمة (0) عندما لا تتحقق تلك الحالة (أي كل الحالات الأخرى). ولحساب نسبة التقارب S_{jk} بين أي مفردتين (ز و k) نقوم بتطبيق الصيغة التالية:

$$S_{jk} = \frac{\text{عدد الأزواج المتشابهة}}{\text{عدد المتغيرات}} = \frac{a+d}{P}$$

حيث a عدد المتغيرات X_i' التي تحصل عندها المفردتين ز و k على نفس القيمة 1، و b هو عدد المتغيرات X_i' التي تحصل عندها المفردتين ز و k على نفس القيمة 0 في نفس الوقت. وأخيراً نحصل على عناصر مصفوفة التقارب S ونكتبها كما يلي:

$$S = \begin{matrix} & \begin{matrix} 1 & 2 & 3 & k & n \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ \vdots \\ j \\ \vdots \\ n \end{matrix} & \begin{bmatrix} S_{11} & S_{12} & S_{13} & S_{1k} & S_{1n} \\ S_{21} & S_{22} & S_{23} & S_{2k} & S_{2n} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ S_{j1} & S_{j2} & S_{j3} & S_{jk} & S_{jn} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ S_{n1} & S_{n2} & S_{n3} & S_{nk} & d_{nn} \end{bmatrix} \end{matrix}$$

ومن خواص هذه المصفوفة إنها مصفوفة مربعة من المرتبة $n \times n$ ، ومتناظرة لأن المقاييس المستخدمة لحسابها متناظرة (أي أن: $S_{jk} = S_{kj}$)، كما أن عناصر القطر الرئيسي فيها ($S_{jj} = 1$)، لذلك يمكننا كتابة مصفوفة التقارب S على شكل مصفوفة مثلثية سفلى (لتمييزها عن مصفوفة التباعد D) كما يلي:

$$S = \begin{matrix} & \begin{matrix} 1 & 2 & 3 & k & n \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ \vdots \\ j \\ \vdots \\ n \end{matrix} & \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ S_{21} & 1 & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ S_{j1} & S_{j2} & 1 & 0 & 0 \\ \vdots & \vdots & \vdots & 1 & 0 \\ S_{n1} & S_{n2} & S_{n3} & \dots & 1 \end{bmatrix} \end{matrix}$$

ملاحظة (1): كل عناصر مصفوفة التقارب S تأخذ قيم كسرية في المجال $[0, 1]$. عكس عناصر مصفوفة التباعد D التي تأخذ قيمها في المجال $[0, +\infty[$.

ملاحظة (2): يمكن حساب عناصر مصفوفة التقارب S_{jk} انطلاقاً من عناصر مصفوفة التباعد d_{jk} من خلال العلاقة:

$$S_{jk} = \frac{1}{1 + d_{jk}} \quad 0 \leq S_{jk} \leq 1$$

ملاحظة (3):

يمكن معالجة المتغيرات النوعية أو المختلطة دون تحويلها إلى متغيرات ثنائية، حيث نستبدلها بأرقامها وندخلها إلى الجدول الأساسي لبيانات العينة، مثلاً: النوع X_1 (0 أنثى، 1 ذكر)؛ X_2 الدخل (1 منخفض، 2 متوسط، 3 مرتفع)؛ X_3 الحالة الاجتماعية (1 أعزب، 2 متزوج، 3 مطلق، 4 أرمل)؛ X_4 العمر بالسنوات (1 أقل من 30، 2 لـ $[30, 50[$ ، 3 لـ 50 فأكثر)؛ ...

ثم نعرف على بيانات المفردتين (k و j) متغير جديد Z، يأخذ القيمة (1) إذا كانت إجابة المفردة j مساوية لإجابة المفردة k، ويأخذ القيمة (0) إذا كانت إجابة j مختلفة عن إجابة k، ثم نقوم بحساب عناصر مصفوفة التقارب s_{jk} من خلال حساب المتوسط الحسابي لقيم المتغير Z بالمنسبة لكل مفردتين k و i بتطبيق العلاقة التالية:

$$s_{jk} = \frac{\sum_i z_{ijk}}{p} \quad i=1 \dots P; \quad J, k=1..n$$

7. خطوات التحليل العنقودي الهرمي التجميعي:

إن عملية العنقدة الهرمية التجميعية لـ n مفردة (مشاهدات أو متغيرات)، تتألف من الخطوات التالية:

1- تبدأ العنقدة الهرمية التجميعية من اعتبار كل مفردة تشكل عنقوداً خاصاً، ومن مصفوفة متناظرة للمسافات (أو للتشابه)، نرسم لها بـ: $D = [d_{jk}]$

2- نبحث في مصفوفة المسافات D عن أصغر عنصر فيها (أصغر مسافة)، ومنه نحدد العنقودين الأكثر قرباً أو تشابهاً من بين العناقيد المدروسة، وذلك من أجل دمجها وتشكيل عنقود جديد منهما.

3- ندمج العنقودين u و v في عنقود واحد، ونرمز له بـ: (uv)، ثم نقوم بتحديث العناصر المتقاطعة في مصفوفة المسافات D كما يلي:

أ. نحذف العمودين والسطرين المقابلين للعنقودين u و v من المصفوفة.

ب. نضيف سطراً وعموداً جديدين، ونضعهما في مكان u أو v، أو في آخر أو أول المصفوفة D.

ج. نقوم بحساب شعاع عناصر العنقود الجديد (uv) من عناصر العنقودين u و v، ونرمز لذلك الشعاع بـ: $d_{(uv)}$ ، وهو عبارة عن المسافات المتقاطعة بين العنقود الجديد (uv) المتبقية في المصفوفة، ثم نضعها في العمود والسطر المخصصين لـ (uv) في المصفوفة، علماً بأن عناصر العنقود الجديد (uv) تحسب باستبدال كل عنصرين متقابلين في العنقودين u و v بأصغرهما أو بأكبرهما أو بمتوسطهما، وذلك حسب الربط المستخدم في عملية العنقدة (المنفرد أو التام أو المتوسط على الترتيب).

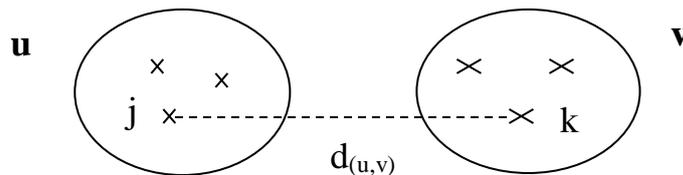
د. نكرر الخطوات (2) و (3) حتى (n-1) مرة حتى تصبح جميع المفردات في عنقود واحد وفي كل مرحلة نقوم بتسجيل مواصفات وخواص العناقيد المدمجة ومستوياتها الجديدة (عناصر المسافات أو التشابه) في الأماكن المخصصة للدمج.

هـ. نقوم برسم مخطط الأغصان لكل مرحلة من هذه المراحل، وأخيراً ندمجها في المخطط العنقودي العام.

أ. طريقة الربط المنفرد Single Linkage

تبدأ من اعتبار كل مفردة من المفردات تشكل عنقوداً خاصاً، وتعتمد على مصفوفة المسافات D (أو التشابه S) بين أزواج المفردات المدروسة، ويتم تشكيل العناقيد فيها من دمج العناقيد الأكثر تقارباً (الجوار الأقرب).

ولتحديد العنقودين الأكثر تقارباً نقوم بدراسة عناصر المصفوفة D، ونحدد أصغر عنصر فيها، ولنفترض أنه كان يقابل العنقودين u و v، لذلك نقوم بدمج هذين العنقودين في عنقود جديد ونرمز له بـ: (u, v)، ولحساب شعاع المسافات للعنقود الجديد (u, v)، نستبدل كل عنصرين متقابلين من (u) و (v) بأصغرهما (بأقربهما)، أي نطبق العلاقة: $d_{(u,v)} = \min[d_{jk}] \quad j \in u, k \in v$ حيث j ينتمي إلى u و k ينتمي إلى v.



ولحساب شعاع المسافات بين العنقود الجديد وأي عنقود آخر (أو مفردة w) نطبق العلاقة التالية:

$$d_{(u,v)w} = \min [d_{uw}, d_{vw}]$$

حيث أن: d_{uw} و d_{vw} هما المسافتان بين العنقود الأكثر تقارباً u مع العنقود w، والعنقود الأكثر تقارباً v مع العنقود w على الترتيب.