

# استخدام خوارزمية (K-Means) للعنقدة في تنقيب البيانات (Data Mining) مع واقع تطبيقي

أ.م.د. قتيبة نبيل نايف / كلية الادارة والاقتصاد / جامعة بغداد  
الباحث / محي الدين خلف أيوب

## المستخلص

ان التقدم العلمي الكبير أدى الى الانتشار الواسع للمعلوماتية بحيث اصبحت المعلومات تتراكم بشكل هائل في قواعد بيانات كبيرة ، وهنا تكمن اهمية البحث في محاولة تنقيح وتبويب هذا الكم الهائل من البيانات وتصنيفها بحيث تؤدي الغرض المطلوب في استخراج المعلومات المخفية او في تصنيف البيانات بموجب علاقاتها ببعضها بغية الافادة منها لأغراض تقنية .

وان العمل بمصطلح التنقيب في البيانات (Data Mining) يعد ملائم في هذا المجال لأهمية البحث في استخدام خوارزمية (K-Means) لتصنيف البيانات بأسلوب تقني في واقع تطبيقي مع ما يمكن من ملاحظة التأثير في المتغيرات (v) من خلال تغيير حجم العينة (n) وكذلك عدد العناقيد (K) واثرها في عملية العنقدة في مراحل الخوارزمية ، من خلال تكوين عناقيد مثالية بحيث تحقق مجموعة بيانات جديدة ومفيدة تجيب عن كل الاستفسارات وبحسب صفات البيانات (Object) العاندة لخوارزمية البحث وبحسب متغيرات البحث (V) المطبقة في برنامج الجانب التطبيقي.

**المصطلحات الرئيسية للبحث /** العناصر ، تنقيب البيانات ، العنقدة ، التعليم الالي ، الخوارزمية.



مجلة العلوم  
الاقتصادية والإدارية  
العدد 91 المجلد ٢٢  
الصفحات ٣٨٩-٤٠٦

البحث مستل من رسالة ماجستير

### ١- المقدمة Introduction:

إن التطور الكبير في مجالات العلوم المختلفة وما صاحبه من تطور في تكنولوجيا المعلومات في عصر الاتصالات الحديثة والانترنت، أدى إلى زيادة كمية البيانات الرقمية إذ لم تعد وسائل التحليل الإحصائي التقليدية قادرة على التعامل معها بغية استخلاص معلومات تفي باتخاذ قرارات صائبة ودقيقة. كذلك فإن وجود كميات كبيرة من البيانات في قواعد البيانات (Data Base) وفي مخازن البيانات (Data warehouse) زادت الحاجة إلى تطوير أدوات تحليل البيانات لغرض استخلاص واكتشاف المعرفة (Knowledge Discloser) والافادة منها. ومن هنا ظهر اسم تنقيب البيانات (Data Mining) كأسلوب تقني يهدف إلى استخراج المعلومات المخفية (Hidden information) بأدوات تقنية وإحصائية حديثة. ومما تجدر الإشارة إليه أن أساليب تنقيب البيانات (Data Mining) تركز أيضاً على بناء التنبؤات المستقبلية في سلوك واتجاهات (Trends) المتغيرات مما يساعد على اتخاذ قرارات أكثر دقة وفي الوقت المناسب.

ومن هنا تكمن أهمية البحث في استخدام إحدى طرائق العقدة وهي خوارزمية (K-Means) في تنقيب البيانات DM في واقع تطبيقي من أجل استخلاص الفائدة والمعرفة من البيانات المطبقة، وبحسب صفات (Object) متغيرات البحث (V) للتجارب المبينة في الجانب التطبيقي.

### ٢- هدف البحث:

من الانظمة المعلوماتية الجديدة انظمة التنقيب في البيانات ضمن قواعد البيانات (Data Base) ومخازن البيانات (Data Warehouse)، ومن ادوات وسائل التنقيب في البيانات طريقة (K-Means) إذ يهدف البحث الى استخدام خوارزمية (K-Means) للعقدة في واقع تطبيقي للوصول الى المعلومة من خلال تصنيف صفات (Object) لمتغيرات (V) البحث التطبيقي فضلاً عن معرفة طبيعة العلائق الموجودة بين صفات البيانات والحصول على النتائج المرجوة لأفضل تصنيف للبيانات بالسرعة المطلوبة.

### البحث الأول / الجانب النظري:

#### ١- مفهوم تنقيب البيانات DM - (Data Mining):

ظهر مصطلح تنقيب البيانات في منتصف الستينات في الولايات المتحدة الأمريكية<sup>[2]</sup> وهو مصطلح يجمع بين الإحصاء وتكنولوجيا المعلومات وقواعد البيانات (Data Bases)، الذكاء الاصطناعي (Neural Intelligent)، التعليم الآلي (Machine Learning)، وظهرت عدة تعريفات لمفهوم تنقيب البيانات DM منها ((الاستكشاف الآلي أو الموتمت لأنماط شائعة وغير جلية مخفية في قاعدة البيانات))<sup>[3]</sup>. وعُرف ((عبارة عن تحليلات لكمية كبيرة من البيانات لغرض ايجاد قواعد وأمثلة ونماذج يمكن أن تستخدم لتفود وتدل صاحب القرار، وتتنبأ بالسلوك المستقبلي))<sup>[4]</sup>، وهناك تعريف آخر لمفهوم تنقيب البيانات ((تحليل لمجموعات كبيرة الحجم من البيانات المشاهدة للبحث عن علاقات محتملة، وتلخيص البيانات في أشكال جديدة لتكون مفهومة ومقيدة لمستخدميها))<sup>[26]</sup>.

وتسمى أحياناً (اكتشاف المعرفة) وهي عملية تحليل البيانات من منظورات مختلفة واستخلاص العلاقات فيما بينها وتلخيصها إلى معلومات مفيدة من خلال استكشاف قواعد جديدة من قواعد البيانات الكبيرة وكذلك اكتشاف نماذج (Model) جديدة وصولاً إلى معلومات ذات قيمة وذلك باستعمال مجموعة من الأدوات والتقنيات الحديثة أو الخوارزميات وغيرها من أدوات الإحصاء الاعتيادية والرسوم البيانية. ومن خلال التعريفات السابقة يمكن القول أن تنقيب البيانات (Data Mining) هي عملية اكتشاف المعرفة من قواعد البيانات (Data Base)، وتسمى أحياناً اكتشاف المعرفة (Knowledge discovering from Data- KDD).

## ٢- أهمية تنقيب البيانات:

أن تنقيب البيانات (DM) أسلوب تمكن من خلاله الوصول إلى المعلومات المخزونة في مستودع البيانات (Data Warehouse) ويتضمن استخدام التحليل الإحصائي (Statistical Analysis) لاكتشاف العلاقات الخفية في البيانات (Romney and Baul, 2009)<sup>(32)</sup>، كما يعتبر تنقيب البيانات أحد تكنولوجيا الذكاء الاصطناعي (Artificial Intelligent) فضلاً عن الأنظمة الخبيرة والشبكات العصبية، ويهدف إلى تمكين النشاط أو المنظمة من الاستغلال الأمثل لبياناتها فهي تحاول إيجاد المعلومات في مجاميع البيانات الكبيرة التي لا تعلم المنظمة أو النشاط بوجودها، كذلك إيجاد العلاقات وعمل التنبؤات .

ويتضح من ذلك أن أهمية تنقيب البيانات ناجمة عن عملية استكشاف وتحليل كميات كبيرة من البيانات لغرض الحصول على علاقات ونماذج خفية تساعد على استخلاص المعلومات المفيدة والمساعدة لاتخاذ قرارات عمل استراتيجية كفيلة بزيادة أداء المنظمة أو النشاط.

١- صياغة المتغيرات وتحولها (Variables Construction and Translation): حيث يجب أن تصاغ المتغيرات الجديدة لبناء النماذج الفعالة التي تفي بالجانب التطبيقي .

٢- تكامل البيانات (Data Integration): إذ إن مجاميع البيانات في دراسة التنقيب عن البيانات من الممكن خزنها في قواعد البيانات متعددة الأغراض والتي تكون بحاجة إلى توحيدها في قاعدة بياناتية واحدة.

٣- تصميم وتنسيق البيانات (Data Formating and consisting): في هذه الخطوة إعادة ترتيب حقول البيانات وفقاً لنموذج التنقيب في البيانات.

## 3- أساليب تنقيب البيانات [Techniques of Data Mining] [11,23,26]:

تستخدم عملية تنقيب البيانات تقنيات وأساليب عديدة تتمكن من خلالها اكتشاف الاتجاهات والنماذج الخفية في مقادير كبيرة من البيانات ويمكن استخدام واحدة أو أكثر من هذه الأساليب، وهي كالآتي:

### أولاً: التصنيف (Classification):

يتم في التصنيف تحليل مجموعة البيانات لتكوين مجموعة من القواعد المتجمعة التي يمكن أن تستخدم لتصنيف بيانات التشغيل إلى مجموعات بحسب صفات معينة، أي: إيجاد المعلومات التي تتعلق بالخصائص المشتركة، وللتصنيف أدوات عديدة مثل:

أولاً: شجرة القرار (Decision Tree).

ثانياً: المجاور الأقرب (Nearest Neighbor).

ثالثاً: الانحدار (Regression).

### ثانياً: الاقتران (Association):

وهي القاعدة التي تتضمن علاقات اقتران ثابتة بين مجموعة من الأشياء في قاعدة البيانات، أي الاقتران بين حدوث حدث ما، وحدث حدث آخر وهي غالباً ما تسمى سلسلة السوق (Market Basket Analysis).

### ثالثاً: التحليل المتسلسل (Sequential Analysis):

وهو يشبه الاقتران ويوضع تحت مسمى تحليل الربط (Link Analysis) لكونه مرتبطاً بالزمن، فيبحث عن نماذج تحدث بالتتابع، أي: يتعامل مع البيانات التي تحدث في حالات منفصلة.

### رابعاً: التجميع أو العنقدة (Clustering):

وهي تقنية تجمع الكيانات المتشابهة سوياً وتفصلها عن البيانات غير المتشابهة وفي مجموعات مختلفة، وتعتمد بصورة أساسية على قياس المسافة وتعد تقنية المجاور الأقرب (Nearest Neighbor) شكلاً آخر للتجميع، إذ من الممكن أن تكون هناك مفاتيح مختلفة لاثنين من أدوات التنقيب على البيانات نفسها. يختلف التجميع (Clustering) عن التصنيف (Classification) بأنه في الأول لا تعرف ما ستكون عليه التجمعات عند البدء أو بأية صفة ستجمع البيانات.

وتستخدم في التجميع أدوات مثل:

أ- متوسط (K means).

ب- الشبكات العصبية (Neural Networks).

٤- أدوات أساليب تنقيب البيانات (Tools of Data Mining Methods) [15,25,27]:

هناك العديد من الأدوات (Tools) التي تستخدم في تقنيات تنقيب البيانات وهذه الأدوات متنوعة ولكل منها دور يخدم غرض معين ومن هذه الأدوات (Tow Crows, 1999: 10-15)<sup>[25]</sup>:

أولاً: أشجار القرار (Decision Trees):

وهي مشتقة من الإحصاء والذكاء الاصطناعي حيث تستخدم الارتباط في البيانات وتستخدم الاستدلال الإحصائي على قواعد العمل وتعد أساس بناء النموذج التنبؤي كما يمكن أن تستخدم في الشبكات العصبية كمدخلات.

ثانياً: الشبكات العصبية (Neural Networks):

وهي قريبة من أشجار القرار لكنها أصعب فهماً وتقدم نماذج ذات قوة تنبؤية أفضل وتتكون من طبقات المدخلات (Layers) وترتبط بعقد الطبقات المخفية (Hidden Layers) التي ترتبط بدورها بعقد طبقة مخفية أخرى وحتى طبقة المخرجات (Out Put Layer) حيث تضم واحد أو أكثر من المتغيرات التابعة.

ثالثاً: الانحدار (Regression):

يستخدم الانحدار في التنبؤ بالقيم الجديدة بالاعتماد على القيم الموجودة ويستخدم الانحدار الخطي للحالات البسيطة، أما الحالات المعقدة التي يصعب التنبؤ بها يستخدم الانحدار النسبي لأنها تعتمد على تفاعلات معقدة لمتغيرات متعددة.

رابعاً: التنبؤ (Predictor):

وتنبأ بالقيم المستقبلية غير المعروفة بالاعتماد على سلاسل تغير الزمن للمتنبات، فيؤخذ بالحسبان الخواص المميزة للزمن كتردد المدد والموسمية.

خامساً: استنتاج القاعدة (Rule Induction):

فيها يتم اشتقاق مجموعة من القواعد المستقلة وعلى خلاف أشجار القرار فهي لا تأتي من شجرة، وقد لا تغطي القواعد الممكنة كل الحالات الممكنة كما أنها قد تتعارض في تنبؤاتها.

سادساً: الجاور الأقرب (K- Nearest Neighbor):

أن فكرة الجاور الأقرب تؤسس على أن حل المشاكل الجديدة يكون عن طريق ملاحظة ومعرفة حلول المشاكل المشابهة والتي تم حلها مسبقاً.

سابعاً: التحليل التمييزي (Discriminate Analysis):

وهي أداة تصنيف تجد السطوح المتعددة التي تفصل الفئات ويكون النموذج الناتج سهل التغيير لأن كل ما على المستخدم هو تحديد جانب الخط الذي تقع عليه النقطة.

ثامناً: الإسناد (Boosting)<sup>[33]</sup>:

وتعتمد على أخذ عينات عشوائية متعددة من البيانات وبناء نموذج التصنيف لكل منها، ويتغير وضع البناء بالاعتماد على نتيجة النماذج السابقة ويكون التصنيف الأخير هو الفئة الأكثر تخصيصاً من قبل النماذج.<sup>[33]</sup> (Dasid H.A Firet,1996).

تاسعاً: الخوارزميات الجينية (Genetic Algorithms):

سميت بذلك لأنها تتبع نموذج نشو الإحياء الذي يتنافس فيه أعضاء النشو الواحد من النماذج لتتقدم خصائصها في النشو اللاحق من النماذج إلى أن يتم إيجاد النموذج الأفضل. أن المعلومات موجودة في الكروموسومات التي تضم خواص بناء النموذج، فالخوارزميات الجينية تعمل أساساً كطريقة لإنجاز البحث الموجه عن النماذج الجيدة.



اذ ان  $S_f$  الانحراف المعياري  
وأن قياس المسافة الاقليدية لكل زوج من الصفات يكون :

$$d(i, j) = \sqrt{(X_{i1} - X_{j1})^2 + (X_{i2} - X_{j2})^2 + \dots + (X_{in} - X_{jn})^2} \quad (4)$$

حيث  $i = (X_{i1}, X_{i2}, \dots, X_{in})$

$J = (X_{j1}, X_{j2}, \dots, X_{jn})$

وهناك مقياس (Manhattan) لقياس المسافة ويكون بالشكل الآتي:

$$d(i, j) = |X_{i1} - X_{j1}| + |X_{i2} - X_{j2}| + \dots + |X_{in} - X_{jn}| \quad (5)$$

وكلا المقياسين للمسافة Euclidean و Manhattan يحققان الآتي:

١- ليست سالبة  $d(i, j) \geq 0$

٢- المسافة لنفس الصفة تساوي صفر  $d(i, j) = 0$

٣- المسافة دالة متماثلة  $d(i, j) = d(j, i)$

٤-  $d(i, j) \leq d(i, h) + d(h, j)$  حيث  $h$  صفة أخرى.

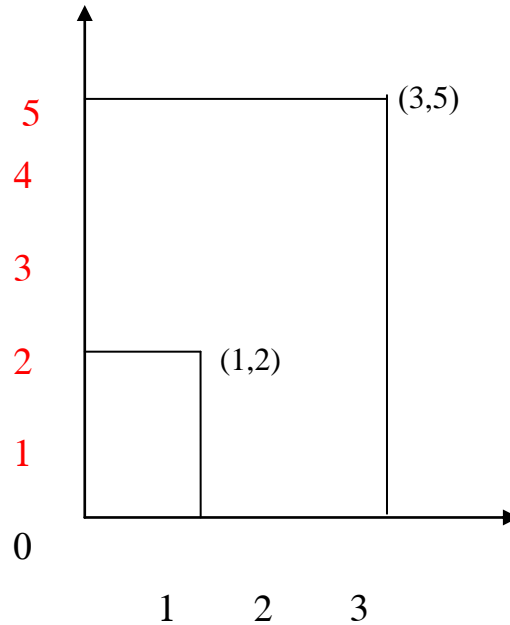
ثانياً: مقياس المتغيرات الثنائية Binary variables:

ويستعمل هذا المقياس لحساب عدم التشابه بين صفات البيانات من خلال قياس المتغيرات الثنائية المتماثلة أو غير المتماثلة وكالاتي:

$$\text{Euclidean distance} = (2^2 + 3^2)^{1/2} = 3.61$$

$$\text{Manhattan distance} = (2+3) = 5$$

وإن الشكل (١) الآتي يبين قياس المسافة الاقليدية ومانهاتن (Manhattan) لـ ٢ من الصفات :



الشكل (١-١) يبين قياس المسافة لصفتين في خوارزمية العنقدة

يكون المتغير الثنائي يملك صفتين فقط هما 0, 1، حيث 0 يمثل أن المتغير غير موجود، و 1 يمثل أن المتغير حاضر.

ويوجد نوعان من المتحولات الثنائية:

$$d(i,j) = \frac{r+s}{q+r+s+t} \quad (\text{Symmetric})$$

أذا أن  $r, s$  عدد مرات الاختلاف.

و  $q$  عدد مرات التشابه (الغرضان يجتمعان بنفس الصفة).

و  $t$  عدد مرات التشابه (الغرضان لا يجتمعان بنفس الصفة).

$$d(i,j) = \frac{r+s}{q+r+s} \quad (\text{Asymmetric})$$

وبعد تمثيل البيانات نقوم بتشكيل مصفوفة عدم التشابه.

ثالثاً: قياس المتغيرات الرتبوية (Ordinal variables):

وهي تشبه المتحولات المطلقة (Absolute variables) ولكن في هذه الحالة يؤخذ الترتيب بعين العناية مثلاً درجات التقدير (دكتوراه، ماجستير، بكالوريوس، دبلوم، إعدادية) فتعطي الرتب (1, 2, 3, 4, 5) فيصبح لدينا مجال التصنيف  $\{1, \dots, m\}$ ، حيث  $m=5$ . وفيما يأتي خطوات تشكيل المصفوفة:

١- نحول المجال  $\{1, m\}$  إلى المجال  $\{0, 1\}$  من خلال العلاقة الآتية:

$$Z_{if} = \frac{r_{if} - 1}{M_f - 1}$$

٢- نقوم بمعالجة المتحولات الناتجة وكأنها متحولات المجال.

٣- نقوم بحساب المسافات.

٤- ثم نشكل مصفوفة عدم التشابه.

أذا أن  $f$  تملك  $M_f$  حالة مرتبة من  $1, \dots, M_f$

و  $f$  تملك  $i$  th من الصفات  $x_{if}$

$$r_{if} \in \{1, \dots, m\}$$

٦- توزيع  $(F)^{[9]}$ :

أذا كان هناك متغيرات  $(Q, W)$  كل منهما يتوزع بصورة مستقلة توزيع  $(X^2)$  مربع كاي، وإذا قسم كل منهما على درجة الحرية المناظرة له فإن النسبة الناتجة من قسمة أحدهما على الأخر تسمى نسبة التباين، وأن التوزيع لها هو توزيع  $(F)$  وفي حالة اختبار تساوي الأوساط فإننا قد أوجدنا نسبة مجموع المربعات بين ومجموع المربعات داخل أي ان صيغة:

$$F = \frac{\frac{1}{t-1} \sum_j n(\bar{y}_{i.} - \bar{y}_{..})^2}{\frac{1}{n-t} \sum_i \sum_j (y_{ij} - \bar{y}_{i.})^2}$$

اذ أن البسط والمقام هما تقديران غير متحيزين لتباين المجتمع  $\sigma^2$  وأن هذه النسبة أيضاً لها توزيع F ويمكن تبسيط هذه النسبة كما يأتي:

$$F = \frac{\text{تقدير التباين بين المجموع}}{\text{تقدير التباين داخل المجموع}}$$

وعندما تكون اوساط ( المجموعات ) غير متساوية فإن تقدير التباين من مجموع مربعات (بين) مساوياً الى  $\sigma^2 + c$  حيث  $C > \text{Zero}$  وبالنتيجة فإن نسبة التباين الى قيمة F ستكبر ومن ثم سنكون في شك في تساوي الاوساط المجموعات او العينات .  
ويمكن صياغة مجموع المربعات الكلية ومجموع المربعات بين المجموعات ومجموع المربعات داخل المجموعات بالشكل الاتي:

$$\sum_i \sum_j (y_{ij} - \bar{y}_{..})^2 = \sum_i \sum_j y_{ij}^2 - \frac{(\sum \sum y_{ij})^2}{N} \quad (1-1)$$

$$\sum_i \sum_j (y_{ij} - \bar{y}_{..})^2 = \sum_i \sum_j y_{ij}^2 - cf$$

$$\sum_i \sum_j (\bar{y}_{i.} - \bar{y}_{..})^2 = \frac{\sum y_{i.}^2}{ni} - \frac{y_{..}^2}{N}$$

$$\sum_i \sum_j (y_{ij} - \bar{y}_{i.})^2 = \sum (\sum y_{ij}^2 - \frac{y_{ij}^2}{ni}) \quad (2-1)$$

ويبين الجدول (1-1) تقدير التباين بين المجموعات وداخل المجموعات وتقدير F لمصدر التباين (S.O.V):

مصادر التباين (S.o.v)	درجة الحرية DF	مجموع المربعات SS	متوسط المربعات MS	F
تباين المجموعات Between groups	t-1	$SSG = \sum_{i=1}^t \frac{y_{i.}^2}{ni} - \frac{y_{..}^2}{N}$	$MSG = \frac{SSG}{t-1}$	$\frac{MSG}{MSW}$
داخل المجموعات Within groups	N-t	$SSW = \sum \left  \sum y_{ij}^2 - \frac{y_{i.}^2}{ni} \right $	$MSW = \frac{SSW}{N-t}$	
الكل total	N-1	$\sum_i \sum_j y_{ij}^2 - \frac{(\sum \sum y_{ij})^2}{N}$		

جدول (1-1) : تحليل التباين لمعيار واحد



## 6- خوارزمية K-Means<sup>[31]</sup>:

التعريف البسيط لعنقدة K-Means هو تصنيف البيانات إلى مجموعات من الأشياء تؤسس على صفات (Attribute) في K عدد من المجموعات، حيث K عدد صحيح موجب. أن خوارزمية K-Mean للعنقدة تؤسس على المركز (Center Based) وهو خوارزمية تقنية تقوم بحل مشاكل العنقدة بمراحل (Steps) متعددة.

وتعرف K-Mean بأنها معدل (Centroid) النقاط في المجموعات وتطبق بين الأبعاد (Dimension) ذات المجال (Space) المستمر.

وإن خوارزمية K-mean يظهر لنا أجزاء (مجموعات) ذات كفاية معقولة من التباين لبياناتها ضمن الصف الواحد ومن الممكن تأمين ذلك بالتحليل الرياضي والتجارب العملية، أضف إلى ذلك أن K-Means ينتج برامج سهلة وهي اقتصادية من الناحية الحسابية ولذلك فأنها سهلة في العمليات والبرامج التي تستخدم العينات الكبيرة.

وخوارزمية K-Mean هي إحدى الخوارزميات التي يكون فيها تحليل البيانات يحقق مجموعة بيانات جديدة ومفيدة بسبب بساطة خوارزمتها، كونها متينة نسبياً، كفاية جيدة، تجيب على كل الاستفسارات لمجموعات البيانات المختلفة.

أن تحليل خوارزمية K-Mean في العنقدة تحدد من خلال:

- 1- تكوين عناقيد مثالية (Clustering membership optimal): أي كل نقطة هي عضو (member) في العنقود وتكون الأقرب إلى المركز.
- 2- بلوغ الأمثلية للمحتوى: أي أن كل نقاط العناقيد تحقق الأمثلية إلى المركز (Centroid) بالتشابه والتقارب بينها.

## أساسيات خوارزميات K-Means<sup>[31,7]</sup>:

أن تقنية العنقدة لـ K-Means هي إحدى الخوارزميات البسيطة والتي تتبع عدة مراحل (Steps) وكما في النقاط أدناه:

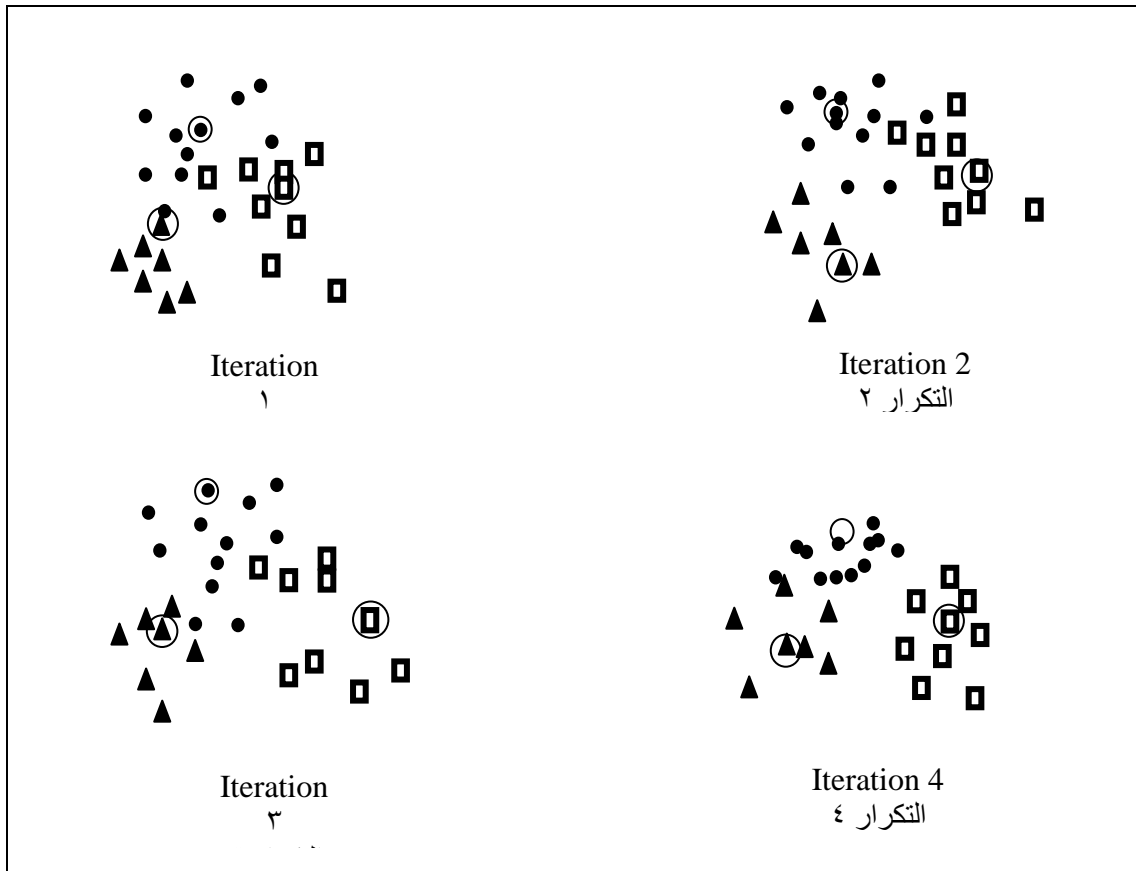
Step 1: اختيار k من النقاط والتي تمثل المراكز (Centroid) الأولية حيث K تمثل أيضاً عدد العناقيد المطلوبة.

Step2: تخصيص النقاط وحسب الصفات (Object) إلى أقرب مركز للعنقود.

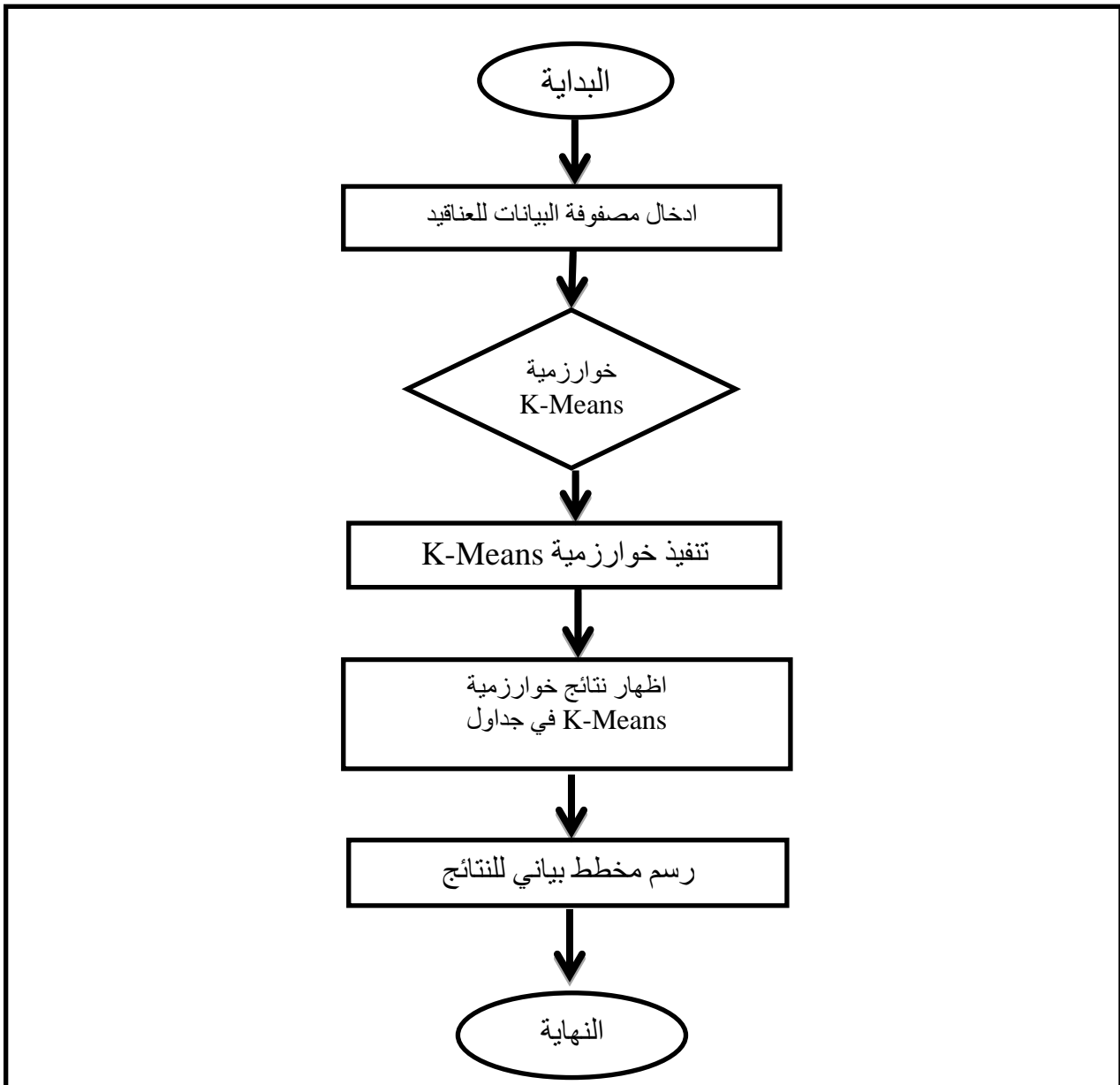
Step3: يتم احتساب مراكز العناقيد (K-Centroid) من جديد.

Step4: نعيد الخطوة (2) حتى تكون كل المراكز متماثلة أو متقاربة مع بعضها.

في الشكل (1-2) يبين عمل (K-mean) والذي يوضح عملية العنقدة مع ثلاثة مراكز (Centroid) حيث أن العنقود النهائي يكون رابع تكرار .



الشكل (1-2) يبين استعمال خوارزمية K-Means  
ولغرض توضيح المخطط الانسيابي لخوارزمية (K-Means) يبين لنا الشكل (1-3) الخوارزمية وبحسب  
المخطط الانسيابي لها .



الشكل (1-3) يوضح المخطط الانسيابي لخوارزمية K-Means

## المبحث الثاني / الجانب التطبيقي:

### ١- المقدمة Introduction:

في الجانب التطبيقي تناولت الدراسة البيانات الخاصة بمتغيرات الدالة والتي تمثل الملاك التدريسي المتغيرات (V) او المشاهدات لتربية محافظة الانبار والمبينة في الملاحق .  
تم سحب عينة مكونة من (١٠٠) مشاهدة اي ان حجم العينة (n=100) حيث  $X=\{X_1, X_2, \dots, X_n\}$  وان عدد المتغيرات التي وقع الاختيار عليها من عينة البحث هي  $m=5$  علماً ان كل متغير من هذه المتغيرات V يمثل عدد من الصفات الخاصة بكل من متغيرات البحث للبيانات في الملحق.  
في البرنامج التطبيقي تم تطبيق خوارزمية (K-Means) على عينة الدراسة واطهرت النتائج كما في الجدول (2-1) .

### ٢- المتغيرات والصفات:

اولاً: الجنس : ويتضمن المتغير صفتين هما ( ذكر ، وانثى) وقد اعطي لهما الرمز (٢,١) على الترتيب.  
ثانياً: الحالة العلمية الشهادة: ويشمل متغير الحالة العلمية الصفات ( دكتوراه ، ماجستير ، بكالوريوس ، دبلوم ، اخرى ) وقد تم تمثيلها بالرموز (١,٢,٣,٤,٥) على الترتيب.  
ثالثاً: التخصص: هذا المتغير يشمل الاختصاصات التالية ( اسلامية، تربية رياضية، تاريخ، جغرافية ، انكليزية، عربية، فيزياء، رياضيات، كيمياء، علوم حياة، فنية) وقد مثلت بالرموز (١,٢,٣,٤,٥,٦,٧,٨,٩,١٠,١١) على الترتيب.  
رابعاً: الدرجة الوظيفية : ويشمل متغير الدرجة الوظيفية الدرجات ( الاولى ، الثانية ، الثالثة ، الرابعة ، الخامسة، السادسة ، السابعة، الثامنة ) وقد تم تمثيلها بالرموز (١,٢,٣,٤,٥,٦,٧,٨) على التوالي.  
خامساً: العنوان الوظيفي : ويشمل { مدرس اقدم اول (معلم جامعي اقدم اول)، مدرس اقدم (معلم جامعي اول)، مدرس اول ( معلم جامعي ) ، مدرس ثاني (معلم اقدم اول)، مدرس ثالث (معلم اقدم ثاني)، مدرس رابع (معلم ثاني)، معلم ثالث ، معلم رابع ، معلم خامس ، مدرس ، معلم } وتتمثل بالرموز من (١,٢,٣,٤,٥,٦,٧,٨,٩,١٠) على الترتيب.

### ٣- المقاييس المستخرجة:

تم استخراج بعض المقاييس الاحصائية من جداول البرنامج ومن اهمها:  
اولاً: مقياس اختبار F-Test وهذا المقياس كما معلوم يقيس اختبار نسبة تباين مجموعتين او عينتين بدرجتي حرية n-1 لكل منهما اي لدرجتي حرية البسط والمقام، وان مقياس F-Test يبين تشتت او تباين البيانات في العناقد المختلفة لبيانات العينة او متغيرات البحث.

ويشترط في هذا لمقياس ان تكون الاخطاء مستقلة من مفردة الى اخرى وتتوزع طبيعياً  $(e_i \sim N(0, \sigma^2))$  اما بالنسبة لصيغة F-test :

$$OR \quad F = \frac{MSC}{MSE} F = \frac{MSR}{MSE}$$

$$MSC=SSC/(n-1) \quad MSR=SSR/(n-1) \quad \text{اذ ان}$$

ثانياً: مقياس MSE (Mean square error)

ان مقياس MSE له اكثر من صيغة وان الصيغة المطبقة في خوارزميات تنقيب البيانات موضوع الدراسة

$$MSE = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{(n-p)} \quad \text{هي:}$$

اذ ان n-p هي درجة حرية الخطأ  
P عد الصفوف

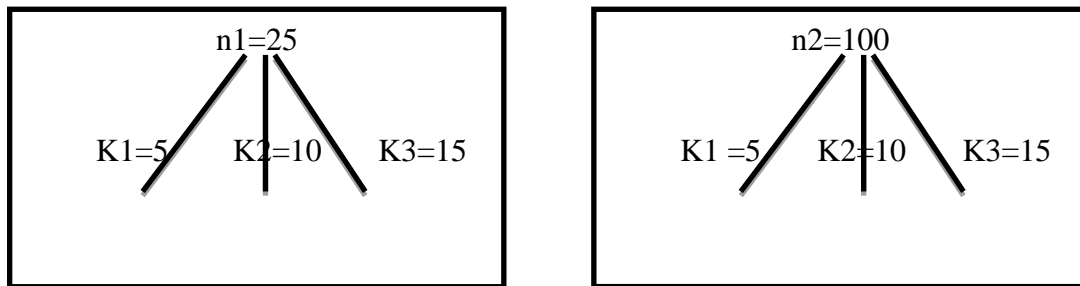
ثالثاً: عدد التكرارات (Iteration History) ويقصد بها عدد مرات اجراء العمليات للخوارزمية لحين الحصول على المتوسط الامثل وكلما كان عدد التكرارات اقل يكون التباين بين البيانات اصغر، ومن ثم فان الوصول الى الامثلية يكون بطريق اقصر.

#### ٤- أسلوب العمل:

تم اختيار حجوم عينات تكون ملائمة لعمليات تنقيب البيانات بحيث تتراوح من العينات الصغيرة الى العينات الكبيرة وهي ( $n_1=25$ ،  $n_2=100$ ) وقد اجريت عملية العنقدة وذلك باختيار عدد عناقيد البيانات بشكل موازي لحجم العينات اي تكون العناقيد ( صغيرة ، متوسطة ، كبيرة ) ايضاً وكالاتي ( $K_3=15, K_2=10, K_1=5$ ) على الترتيب ولكل عينة، وذلك لملاحظة مدى تأثير عملية عنقدة البيانات بكل من حجم البيانات وعدد العناقيد.

#### ٥- تحليل النتائج (Analyze the results):

تم تحليل النتائج لخوارزمية (K-Means) وكما فيما يأتي:  
تم تطبيق خوارزميات العنقدة (K-Means) وبحسب العناقيد K لكل عينة n وكما فيما يأتي:



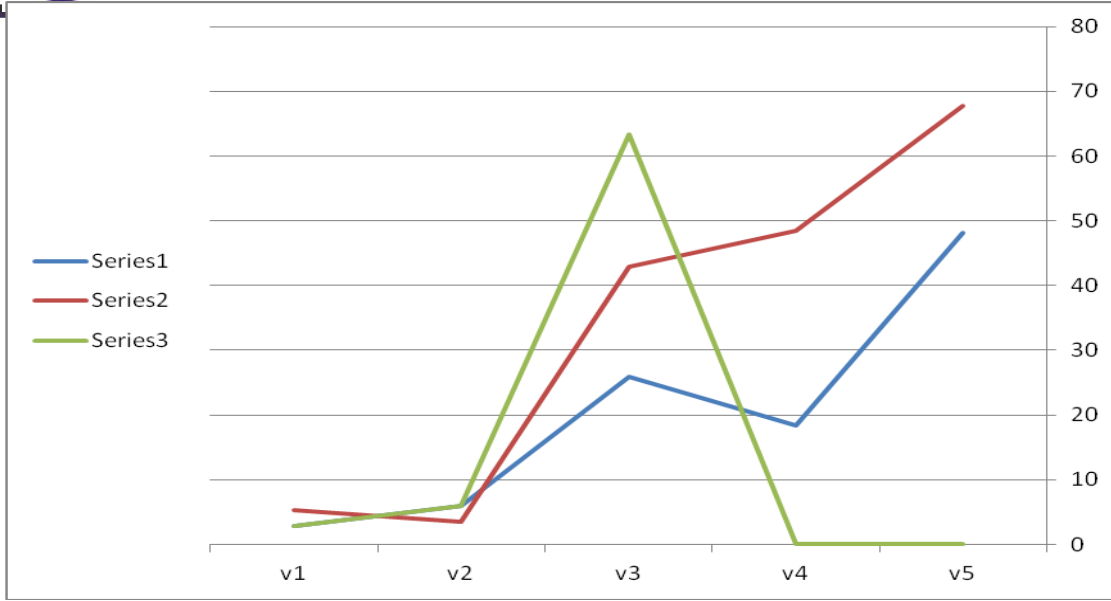
شكل (2-1) يبين توزيع العناقيد (K) على التجريبتين

ومن ملاحظة النتائج في جدول التحليل (2-1) الاتي : يتبين لدينا مدى العلاقة بين حجم العنقود (K) لكل عينة (n) من عينتي البحث مع قيمة الاختبار (F-test) ولكل متغير (v) من متغيرات الدراسة :

F-test حجم العينة	V1	V2	V3	V4	V5
25-5	2.836	5.918	25.928	18.338	48.095
25-10	5.333	3.489	42.924	48.523	67.807
25-15	2.806	5.914	63.351	0.000	0.000
100-5	9.341	45.747	77.670	117.451	161.996
100-10	5.663	936.660	64.152	227.274	280.926
100-15	4.600	24.874	102.294	250.906	334.852

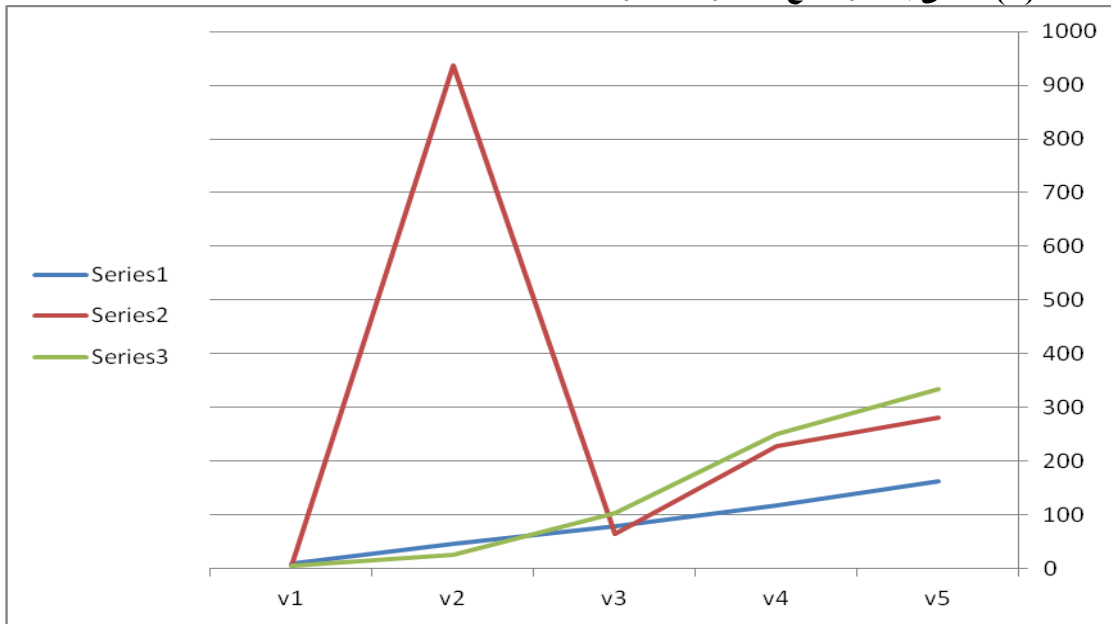
جدول (2-1) يبين قيمة F-Test للمتغيرات لكل عنقود ولحجوم العينات المختلفة

وعند ملاحظة جدول (2-1) يتبين معرفة تأثير F-Test بكل من (n) حجم العينة و (K) حجم العنقود وكذلك ملاحظة مدى تأثير كل من (F,n,K) لكل من المتغيرات الخمسة لأغراض العنقدة من خلال الاشكال المبينة بيانياً فيما يأتي:



شكل (٢-٢) يبين تجارب (n=25) للعناقد K المختلفة

من ملاحظة الشكل السابق تبين لنا ان قيمة (F) تتأثر بكل متغير من المتغيرات الأخوذة في البحث وللتجارب الثلاثة الاولى عند (n=25) ونلاحظ ان قيمة (F) الاعلى تعود للمتغير الخامس (V5) في حين ان قيم (F) الادنى تعود للمتغير الاول (V1) للتجربتين الاولى والثانية اما التجربة الثالثة فقد تميز المتغير الثالث بتقديره (F) الاعلى بالمقارنة مع المتغيرات الأخرى.



شكل (٢-٣) يبين تجارب (n=100) للعناقد K المختلفة

من ملاحظة الشكل (٢-٣) السابق تبين ان قيمة (F) اخذت قيمةً متقاربةً للتجارب الثلاثة عند المتغير V3 في حين انها اخذت قيمةً مرتفعةً جداً عند المتغير V2 ثم انخفضت قيمة (F) لتتقارب مع بعضها عند المتغير الثالث V3 والمتغير الرابع V4 والمتغير الخامس V5.

نستنتج من الاشكال السابقة أن قيمة (F) في المتغير الثالث V3 لحجوم العينات للتجارب (n=25) أخذت المتغير الثالث V3 كقيمة اعلى من بقية المتغيرات في حين ان في التجربة (n=100) تكون قيمة (F) هي الاعلى في المتغير الثاني V2 وهذا يفسر ان افضل عنقدة بالنسبة للتجربة (5-25) هي للمتغير V3 في حين ان المتغير V2 كان هو الافضل بالنسبة للتجربة (10-100). وبالعودة الى النتائج في الجدول (2-1) والجدول السابق ان خوارزمية العنقدة قد اختلفت بحسب ظروف التجربة المتمثلة بحجم العينة n وعدد العناقيد K وعدد التكرارات بغية الوصول الى العنقدة الامثل فضلاً عن عادية كل مفردة او عنصر الى كل عنقود وهذه تختلف من مرحلة الى اخرى ومن متغير الى اخر.

### الاستنتاجات :

- 1- تجب ملاحظة مدى تأثر العنقدة الاقل بظروف التجربة (عدد المفردات) فضلاً عن عدد العناقيد المطلوبة K.
- 2- ان قيم (F) قد تغيرت فيما بينها من الاعلى الى الاقل فنلاحظ ان اعلى قيمة لـ (F) كانت للتجربة (100-10).  
10 (المتغير الثاني (V2) والتي تمثل اقل قيمة لـ MSE. واما بالنسبة للتجربة (25-5) فان قيمة (F) اخذت اقل قيمة عند المتغير الاول (V1) والتي تكون فيها قيمة MSE اعلى من التجربة السابقة اي ان عنقدة البيانات بالنسبة للمتغير (V2) هي الافضل .
- 3- يمكن للوصول الى العنقدة الامثل بالمرور بأعداد مختلفة من التكرارات بغية الوصول الى الامثلية المطلقة وبحسب صفات (Object) متغيرات البحث (V).
- 4- ان خوارزمية (K-Means) اظهرت الاشكال انه عند تغير ظروف التجربة تكون نتائج الخوارزمية متغيرة من حيث النوع (المفردات التي يمتلكها كل عنقود) والعدد (عدد المفردات داخل العنقود).
- 5- في جدول خوارزمية (K-Means) المثلى يكون التغيرات (F) بين العناقيد المختلفة في اقصى ما يمكن ويكون عندها MSE اقل ما يمكن ، بينما يكون التغيرات (F) داخل كل عنقود في اقل ما يمكن.

### التوصيات :

- 1- يمكن اعتماد نتائج خوارزميات العنقدة في التنقيب عن البيانات ولاسيما عندما تكون اعداد البيانات كبيرة جداً حيث تتمكن النتائج من الوصول الى المفردة المطلوبة بعدد المرور بالطرق الاقصر لهذه النتائج.
- 2- بالإمكان زيادة عدد المتغيرات (V) للوصول الى عنقدة أمثل تحقق نتائج بحث وتصنيف في قواعد البيانات بأقصر وقت ولكل متغير من متغيرات البحث.
- 3- يمكن استعمال خوارزمية العنقدة مع خوارزمية الذكاء الصناعي (Artificial Intelligence) وهذا ما يسهم في تقليل الجهد والوقت في الوصول الى المعلومة .
- 4- يمكن استعمال خوارزميات عنقدة اخرى مثل (Ward) او (Seeded) لأغراض المقارنة والبحث عن الامثلية (Optimization) .
- 5- ان تقنيات العنقدة يمكن ان تتعامل مع بيانات الحكومة الالكترونية والذكية وهذا ما يتطلب اعتماد نتائج العنقدة في الوصول الى قواعد البيانات الضخمة والتي من الممكن الدخول من خلالها في دائرة البحث اللامتناهية والتي تحتاج الى ادلة بحث وهذا ما تقدمه نتائج خوارزميات العنقدة .



## استخدام خوارزمية (K-Means) للتعقيد البيانات [Data Mining] مع واقع تطبيقي

### المصادر العربية :

- ١- مصطفى؛ فؤاد عبيد؛ "تقنيات التنقيب في قواعد البيانات واستكشاف في المعلومات المخبأة فيها" 1995.
- ٢- حسن؛ أ. طيار؛ شلاب؛ عمار؛ "التنقيب في البيانات واتخاذ القرارات"، كلية العلوم الاقتصادية، جامعة لوت ، 20؛(1995).
- ٣- العلق؛ بشير عباس؛ "الادارة الرقمية والتطبيقات"، مركز الامارات للدراسات والبحوث الذاتية 83، 2005.
- ٤- العلي؛ عبد القادر، قنديل؛ عامر ابراهيم، الطمري؛ غسان "المدخل الى دراسة المعرفة وادارة السيرة" عمان، 175 2006.
- ٥- السامرائي؛ لمياء عبد الصمد؛ " قواعد المعلومات"؛ العراق جامعة بغداد 20 (1988).
- ٦- العمران؛ حمد بن ابراهيم والعبيدي؛ هديل شوكت : "الوعي المعلوماتي والحكمة"، الرياض مكتبة الرشيد ص 88-89 (2007).
- ٧- رحيمية؛ وليد عبد الله : "استخدام التحليل العنقودي وتحليل الانحدار في تشخيص امراض القلب"، الجامعة المستنصرية؛ 1995.
- ٨- صديق؛ رضوان وعبد العزيز؛ غيداء : " تقييم صحة العنقدة " مجلة الراقدين لعلوم الحاسبات والرياضيات المجلد 5 العدد2 (2008)
- ٩- المشهداني، كمال علوان : "تصميم وتحليل التجارب-استخدام الحاسوب-" ، مكتب الجزيرة للطباعة والنشر، بغداد (٢٠١٠).

### المصادر الاجنبية:

- ١٠- Houston Andrea L.& others, "Medical Data Mining" on the internet Research on acancer Information system Artificial intelligence Review 2000.
- ١١- Edelstein, Herb, "Mining for gold" information week. April, 1997.
- ١٢- Remach an Ran M.puspa, "Mining for Gold wipro technologies" 2001
- ١٣- Luan Jing, "Data Mining And knowledge Management" in Higher Education Air Canada 2002.
- ١٤- Ahola Jussi & Rinta-Runsala Esa; "Data Mining case studies" in customer profiling 2001.
- ١٥- spitiopoulou, Myra & phole carsten: "Data Mining to measure and Improve the success of web sites" 2002.
- ١٦- Friedman I.H "Data Mining and Statistics" what in the connection the 29<sup>th</sup> sum passion on the information couputing science and Statistics, key note speech Houston (1997).
- ١٧- Little R.Rubin D, "statistical analysis with missing data widely": New York (1987).
- ١٨- Hoes J.kauper M. "Data mining concept and technique" morgan kaufmann, san Francisco .A (2006).
- ١٩- Efrom B, Boot strap method Another Look at the Jackknife Amals of statistics 7, 1-26 (1979).
- ٢٠- Efrom B,cous G: Alies Lrely Look at the Boost strap the Jackife and cress-ualidation American statistic 37,36-48 (1988).
- ٢١- Fland, D. uanni laH. Smuth R. "principle of data Mining" Mit press London 2001.





## استخدام خوارزمية (K-Means) للعثور على تنقيب البيانات [Data Mining] مع واقع تطبيقي

- ٢٢- source, Ranach and ran M. psuhpa, "Mining for gold wipro Technologies" Decemper 2001.
- ٢٣- Ramch M. push: "Data Mining for gold" wipro technologies 2001.
- ٢٤- Soni, Sanjay and Tang zhohui and young Jim performance study of Microsoft "Data Mining Algorithms Microsoft" 2001.
- ٢٥- Two crows corporation, "Introduction to Data Mining and knowledge discovery" third Edition 1999.
- 26- Atre shaka, Detining todays "Data mining" Excutive update Business intelligence Advisory service 2001.
- 27- Brand Estelle & Gerristen, Rob "Data Mining" solution Magazine 1998.
- 28- "Discovering knowledge in Data Mining" By Daniel T.Larose 2005.
- 29- "Data Mining" k- Dustering problem By Eham karousi university of Agder 2002.
- 30- "Data Mining concept and techniques" second Edition By Jiawi Han university of lionois at urbanchampaign Michline kamer.
- 31- Halkid, y.batistakis M.vozor giannis "clustering Algorithm and usability measures 2001.
- 32- Romney\_ Marshall B. and Baul John Standart "Accounting information system 9<sup>th</sup> ed" .upper saddle river prentice Hall (2009).
- 33- Dasid H.A Firt occurrence of common Jerm in mathematical statistic the American, 121-133 199 ٦.



## User (K-Means) for clustering in Data Mining with application

The great scientific progress has led to widespread Information as information accumulates in large databases is important in trying to revise and compile this vast amount of data and, where its purpose to extract hidden information or classified data under their relations with each other in order to take advantage of them for technical purposes.

And work with data mining (DM) is appropriate in this area because of the importance of research in the (K-Means) algorithm for clustering data in fact applied with effect can be observed in variables by changing the sample size (n) and the number of clusters (K) and their impact on the process of clustering in the algorithm.

**Key words :** object ,data mining, clustering ,machine learning ,algorithm