



اسم المقال: استخدام المصنف C4.5 في تمييز سمة الكائن - دراسة مقارنة

اسم الكاتب: م. غيداء عبدالعزيز الطالب، م. رائد عبدالقادر الدباغ، م.م. نعمة عبدالله الفحري

<https://political-encyclopedia.org/library/3070>

تاريخ الاسترداد: 2025/05/10 02:07 +03

الموسوعة السياسية هي مبادرة أكاديمية غير هادفة للربح، تساعد الباحثين والطلاب على الوصول واستخدام وبناء مجموعات أوسع من المحتوى العلمي العربي في مجال علم السياسة واستخدامها في الأرشيف الرقمي الموثوق به لإغناء المحتوى العربي على الإنترنت.

لمزيد من المعلومات حول الموسوعة السياسية – Encyclopedia Political – يرجى التواصل على

info@political-encyclopedia.org

استخدامكم لأرشيف مكتبة الموسوعة السياسية – Encyclopedia Political يعني موافقتك على شروط وأحكام الاستخدام

<https://political-encyclopedia.org/terms-of-use>

تم الحصول على هذا المقال من موقع مجلة تنمية الراذدين كلية الإدارة والاقتصاد / جامعة الموصل ورفده في مكتبة الموسوعة السياسية مستوفياً شروط حقوق الملكية الفكرية ومتطلبات رخصة المشاع الإبداعي التي يتضمن المقال تحتها.



استخدام المصنف C4.5 في تمييز سمة الكائن دراسة مقارنة

نعمة عبدالله الفخرى	رائد عبدالقادر الدباغ	غيداء عبدالعزيز الطالب
مدرس مساعد-قسم نظم المعلومات	مدرس-قسم نظم المعلومات	مدرس-قسم علوم الحاسوبات
كلية الادارة والاقتصاد-جامعة الموصل	كلية الادارة والاقتصاد-جامعة الموصل	كلية علوم الحاسوبات والرياضيات-جامعة الموصل

المستخلص

ان تعدين البيانات فعالية الحصول على المعرفة لتحقيق هدف أساس وهو اكتشاف الحقائق الخفية (Hidden Facts) التي تتضمنها قواعد البيانات وذلك من خلال استخدام تقنيات متعددة تشمل على الذكاء الاصطناعي، التحليلات الاحصائية، تقنيات ونمذجة البيانات ... الخ. فإن عملية تعدين البيانات تولد نماذج وعلاقات واضحة في البيانات والتي تساعده على توقيع النتائج في المستقبل. وقد ظهرت العديد من الخوارزميات التي في هذا المجال، وترتبط عليها مقارنة بين هذه الخوارزميات لاختبار الخوارزمية المناسبة في الحصول على نتائج أفضل.

وقد هدف البحث الى استخدام المصنف C4.5 وربطها مع الشبكة العصبية نوع - Back Propagation (BP) وذلك لتكوين نموذج تصنيف يحمل خواص الطرفين، فضلا عن مقارنة النتائج المستحصلة مع نتائج التصنيف باستخدام الحزمة البرمجية الجاهزة Minitab. وتوصل البحث الى ان المعادلات الخاصة بالمصنف C4.5 كانت أفعى في الاداء وخاصة بعد ربطها بالشبكة العصبية BP لازالة التناقض والتشویش الموجود في البيانات، كما عززت النتائج من افضلية استخدام لغات البرمجة مقارنة بنتائج التطبيق الجاهزة.

١. المقدمة

ان التطور الحاصل في موضوع تعدين البيانات (Data Mining) في المجالات والصناعات المختلفة أدى الى ظهور العديد من الخوارزميات، الأمر الذي جعل من الأهمية اختيار خوارزمية تعدين مناسبة للحصول على نتائج أفضل وذلك بسبب تنوع البيانات واختلافها، مما يعمل جيداً على بيانات معينة قد لا يعمل بنفس الجودة على بيانات أخرى. وتأخذ عملية تعدين البيانات الاعتبارات الآلية (الفخرى، ٢٠٠٣، ١-٢):

أولاً: تمثيل المعرفة باستخدام خوارزميات خاصة، وهي خوارزميات واسعة ومتنوعة، فمنها ما يعمل بأسلوب شجرة القرار (Decision Tree)، ومنها ما يستخدم قاعدة اذا ... فان (if – then – rule) ... الخ .

ثانياً: كيف تستطيع الخوارزمية الوصول الى أعلى مقياس اعتماداً على فضاء البحث المتوفر لديها.

وبناءً على ذلك ظهرت عمليات مقارنة بين الخوارزميات و على مديات مختلفة من البيانات كمحاولة لوضع خصائص لهذه البيانات، ثم مطابقة تلك الخوارزميات مع خصائص تلك البيانات، وهذه المعلومات تساعد في تعدين البيانات لصنع قرارات ذكية في اختيار الخوارزمية الملائمة لملفات البيانات .

تأسيساً على ما نقدم، فقد هدف البحث الى استخدام المصنف C4.5 بوصفه أحد

خوارزميات نوع شجرة القرارات وربطها مع الشبكة العصبية – Back Propagation (BP) وذلك لتكوين نموذج تصنيف يحمل خواص الطرفين، فضلاً عن مقارنة النتائج المستحصلة مع نتائج التصنيف للحزمة البرمجية الجاهزة Minitab^(*) سعياً لتحقيق فرضية البحث ومفادها "ان استخدام البرمجيات الجاهزة في التصنيف Classification يعتمد على اسلوب محدد باستخدام احدى خوارزميات التصنيف دون مقارنة النتائج مع بقية الخوارزميات، الامر الذي يقلل من أهمية استخدام لغات البرمجة في كتابة البرامج الأكفاء والأفضل في التصنيف".

لقد اعتمد البحث في أسلوبه جانبيين، الأول يمثل الجانب النظري و الذي يتم من خلاله وصف المصنف C4.5، فضلاً عن استخدام الشبكات العصبية الاصطناعية لغرض توليد البيانات بعد تغذية البرنامج بخصائص الحالات قيد الدراسة. في حين تناول الجانب العملي وصفاً لنتائج تطبيق البرنامج على الحالات المولدة ومقارنتها مع نتائج التصنيف باستخدام الحزمة البرمجية الجاهزة Minitab .

وتم استخدام لغة البرمجة VB الاصدار 6.0 في كتابة البرامج كافة المستخدمة لتنفيذ خطوات الخوارزمية وربطها مع الشبكة العصبية BP .

٢. الجانب النظري

١-٢ تعدين البيانات Data Minig

يتميز عصرنا الحالي باستخدام تكنولوجيا البيانات المتطرفة لحفظ و استرجاع البيانات وبكميات هائلة و التي توصف بمستودعات البيانات Data Warehousing. ان توفر هذه البيانات فتح الباب امام مجموعة من الموضوعات المتخصصة في ادارة تلك البيانات كان من ابرزها موضوع تعدين البيانات Data Minig والذي يعد من الاساليب المهمة للحصول على معلومات مفيدة من البيانات (Ibrahim, 1999,9). وقد عرف العلماء مصطلح تعدين البيانات على انه "جزء من عملية اكتشاف المعرفة في قواعد البيانات و التي تتم باستخدام طرائق متعددة هدفها تكوين نماذج من البيانات" .

^(*) الحزمة الجاهزة Minitab Release 13.0 والتي تعمل في بيئة التوافذ، و هي إحدى برامجيات لوائح العمل في مجال الإحصاء و الرياضيات .

وبصورة اخرى، فقد تم تعريف مصطلح تعداد البيانات بأنه "عملية اكتشاف المعرفة وطريقة تحليلها من زوايا مختلفة و تلخيصها وتحويلها الى ما يسمى بـ (معلومات - معلومات) لتوضيح امام صانعي القرار للعمل على اساسها في مجالات عده مثل المحاسبة، الاتصالات، استخدام اوسع في المجالات الطبية، ... الخ و بشكل يعمل على زيادة الدخل او تقليل الكلف او كليهما معاً".

وفي ضوء ما سبق من تعریفات، فإنه يمكن القول بأن تعداد البيانات هو فعالية الحصول على المعرفة لتحقيق هدف اساس وهو اكتشاف الحقائق المخفية Hidden Facts التي تتضمنها قواعد البيانات وذلك من خلال استخدام تقنيات متعددة تشمل على الذكاء الاصطناعي، التحليلات الاحصائية، تقنيات ونمذجة البيانات. ان عملية تعداد البيانات تولد نماذج وعلاقات واضحة في البيانات والتي تساعدها على توقع النتائج في المستقبل.

اذن، لتهجين البيانات ادوات، وبرامجيات تعداد البيانات هي احدى هذه الادوات والتي تستخدم لتحليل البيانات من قبل المستخدم و تلخيص العلاقات التي تعرفها، فضلا عن دورها في ايجاد روابط بين تقنية المعلومات و الانظمة التحليلية المسئولة لتمثيل العلاقات بين نماذج البيانات المخزونة

(www.Anderson.Ucla.edu)

وتقييماً، فإن تعداد البيانات هو عبارة عن إجراء أو معالجة لإيجاد الروابط بين مجاميع (الحقول/السجلات) في قواعد البيانات الكبيرة والتي تشتهر فيها برامجيات أخرى فضلاً عن برامجيات تعداد البيانات كالبرام吉ات الاحصائية، والشبكات العصبية، الخ. وبشكل عام فإنه يمكن القول بأن أي من العلاقات الآتية مهمة في مجال تعداد البيانات:

١. الاصناف Classes: وتستخدم عادة لوضع البيانات المخزونة في مجاميع تم تحديدها مسبقاً لبناء نموذج بالاعتماد على بعض المتغيرات المستقلة .
٢. العناقيد Clusters : و تستخدم لوضع البيانات في مجاميع اعتماداً على العلاقات المنطقية. بعبارة اخرى، فان الخوارزميات المستخدمة للتصنيف في هذه الطريقة تسعى لتقسيم البيانات الى مجاميع (عناقيد) بحيث ان السجلات المشابهة تقع في المجموعة نفسها وهذه المجاميع يجب ان تكون مختلفة عن بعضها قدر الامكان.
٣. الروابط Associations : و هي تعرف العلاقات الخاصة بتهجين البيانات، اذ إن الخوارزميات المستخدمة فيها تنشئ قواعد لربط الحوادث التي تظهر سوية في البيانات .
٤. النماذج المتسلسلة: يتم تعداد البيانات لتوقع سلوك واتجاهات النماذج المستحصلة .(www.Anderson.Ucla.edu)

٢-٢ خوارزميات التصنيف Classification Algorithms

ان طرق تعداد Data Mining Methods هي عبارة عن مجموعة الاجراءات و الخوارزميات المصممة لتحليل البيانات المخزونة Data Baise، و هي تتعامل مع

عدة عوامل أهمها، أولاً: الدقة بين النماذج المترکونة والبيانات المتوفرة، وثانياً: تمثيل المعرفة باستخدام خوارزميات خاصة (www.towercrows.com). وبما ان ايجاد شكل موحد يمثل المعرفة لجميع البيانات أمر صعب المنال، فقد أوجد الباحثون أشكالاً للتصنيف غير معتمدة على تمثيل المعرفة أي انها مصنفات للاغراض العامة وهي تقع في نوعين، الأول يطلق عليها المصنفات التي تعمل بمحرر و الثاني يعمل بدون محرر (Ibrahim, 1999,11-23)

٢-٤-١ C4.5 المصنف C4.5 Classification Algorithm

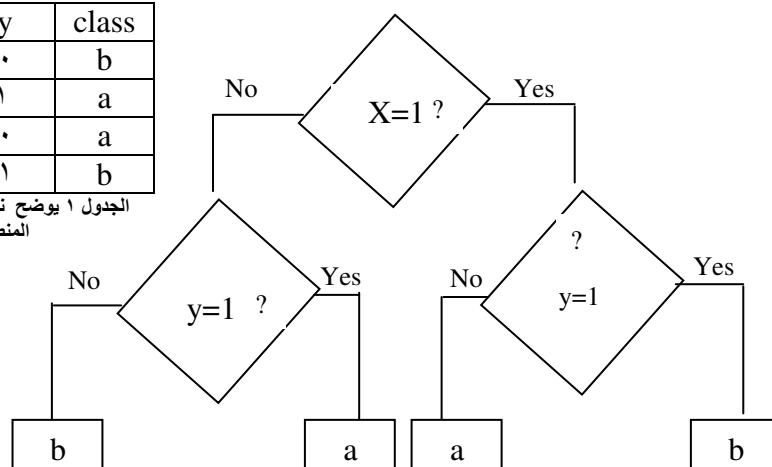
ان خوارزمية C4.5 هي احدي انواع الخوارزميات التي تعمل باسلوب شجرة القرار Decision Tree، حيث تمثل فيها تمثيل مجاميع القرارات و هذه القرارات تولد قوانين التصنيف الخاصة بمجاميع البيانات، وبشكل آخر فانها تعمل على تصنیف الحالات Instances الى فئات مختلفة، وهي من طرائق التصنيف شائعة الاستخدام.

تقع هذه الخوارزمية تحت مجموعة المصنفات التي تعمل بمحرر و التي تصنف حالات معينة الى عدد من الفئات (التصانيف) باسلوب فرق – تسد، اذ انها تجزي المسألة المعقدة الى مسائل ابسط ثم يتم الاستدعاء الذاتي للدالة نفسها و لكل اجزاء المسألة، ويجمع حلول المسائل المجزأة يتم الحصول على حل المسألة المعقدة . (Ali and Abraham , 2002, 2-3)

ويمكن تشبيه عمل شجرة القرار في هذه الخوارزمية بعمل بوابة XOR المنطقية المتمثل بجدول الحقيقة الآتي : (keller,2000,8)

x	y	class
.	.	b
.	١	a
١	.	a
١	١	b

الجدول ١ يوضح نتائج بوابة XOR المنطقية



الشكل ١ عمل شجرة القرارات للمصنف C4.5

اذا ان : العقد Nodes تمثل قرارات الخصائص Decisions On Attributes . الاوراق Leaves تمثل الاصناف Class.

٢-٢-٢ مراحل عمل الخوارزمية C4.5

يمكن توضيح عمل الخوارزمية C4.5 بالمراحل الآتية:

١. تحديد التصنيف لكائن معين من خلال تنفيذ العلاقة الآتية:

$$I(p,n) = \frac{p}{p+n} \log \left[\frac{p}{p+n} \right] - \frac{n}{p+n} \log \left[\frac{n}{p+n} \right] \quad \dots \dots \dots (1)$$

اذا ان : P, n تمثلان اصناف مختلفتين.
 $P \neq n$

٢. تعين خاصية معينة و اختيارها و L_k من القيم وذلك باستخدام العلاقة الآتية:

$$E(A,p,n) = \sum_{i=1}^k \frac{p_i + n_i}{p+n} \cdot I(p_i, n_i) \quad \dots \dots \dots (2)$$

اذا ان : A هي الخاصية التي تم اختيارها
 p_i, n_i عدد الحالات لكل صنف من الشجيرات الناتجة عن شجرة القرار و تكون مرتبطة مع الجزء I .

٣. الحصول على النسبة النهائية للمعلومات من المعادلة الآتية:

$$Gain(A,p,n) = I(p_i, n_i) - E(A,p,n) \quad \dots \dots \dots (3)$$

(Ali and Abraham, 2002, 3; Hidalgo and others, 2002, 327)

١-٢ الشبكات العصبية Neural Network

تميزت الشبكات العصبية ومنذ انتشارها في منتصف الثمانينيات بامكانياتها العالية في اجراء المعالجة المتوازنة وعدم حاجتها الى علاقات معقدة في عملها، بل تحتاج الى بعض الامثلة لتعلمها فقط، ومن ثم سهولة الاجابة . وتعطي الاجابة الصحيحة للدخلات التي ليست ضمن فقرة التدريب مع استبعاد ما يميل منها الى التناقض والتشویش.

واعتماداً على نوعية التطبيق فقد تم اختيار شبكة الانتشار الخلفي – Back Propagation (BP) في بناء النموذج، والتي تعد احدى الشبكات الواسعة الاستخدام والكافحة في طرق التعلم والتي تستخدم لتدريب الشبكات متعددة الطبقات. وتعتمد هذه الشبكة في خوارزمياتها على القاعدة المعروفة باسم الانحدار التدريجي Stepwise Regression لمربع معدل الاوزان. ان تدرج الخطأ واوزان الشبكة يعطي الاتجاه الذي يتزايد فيه الخطأ بأسرع ما يمكن.

اما المعادلات المستخدمة في الإخراج و تصحيح الاوزان فتعطى بالاتي

$$\left. \begin{array}{l} a = f(w_1 * h) \\ h = f(w_2 * e) \end{array} \right\} \dots\dots\dots(4)$$

اذ أن : a تمثل متجه الإفراج، e تمثل متجه الادخال، h تمثل الطبقة المخفية و (W₁ , W₂) مصفوفتي الاوزان .

ومعادلة التفعيل Sigmoid function هي :

$$f(x) = 1/(1 + \exp(-c.x)) \dots\dots\dots(5)$$

$c > 0$

ويتم احتساب الاوزان بطريقة ما بحيث يكون الخطأ أقل ما يمكن و ان دالة الخطأ تتمثل بالعلاقة الآتية :

$$E = 1/2 * \sum_i (z_i - a_{i1})^2 \dots\dots\dots(6)$$

اذ إن : z_i تمثل القيمة النهائية للدالة قيد التدريب.

a_i تمثل القيمة الحقيقة (اخراج الشبكة)

ومن العلاقة ٦ تمثل قيمة L أحسن قيمة لدالة الخطأ للشبكة و اذا كانت $E=0$ فان الشبكة تعمل بصورة أدق .

اما الأوزان فتتغير وفق المعادلتين الآتتين:

$$\left. \begin{array}{l} \Delta_{ij}^1 = \alpha * \varepsilon_i * a_i (1-a_i) * h_j \\ \Delta_{ij}^2 = \alpha * \sum_m \varepsilon_m * a_m * (1-a_m) w_{ij}^1 * h_j * (1-h_i) * e_i \end{array} \right\} \dots\dots\dots(7)$$

وتعطى قيمة الخطأ بالعلاقة الآتية:

$$\varepsilon_i = z_i - a_i = \text{Final Value} - \text{Network Value} = \text{error} \dots\dots\dots(8)$$

(kinnerbrock ,1995 ,40-41)

وعليه يمكن وصف خوارزمية التعليم للشبكة BP بالخطوات الآتية :

١. احتساب قيم الأوزان كافة و لأعداد عشوائية .

٢. اختيار نموذج ادخال – اخراج عشوائي للدالة قيد التعلم، وحساب قيم h_i للطبقة المخفية.

٣. لقيم الادخال e_i والقيم النهائية للشبكة z_i قيم تصحيح الاوزان وفقاً للمعادلة ٧.

٤. العودة الى الخطوة ٢ .

الجانب العملي

٢-١ عينة الدراسة

قبل البدء في عرض النتائج المستحصلة من تطبيق المصنفات، لا بد من وصف البيانات التي استخدمت في البحث، فقد تم اختيار بذور النباتات (الحنطة، الشعير، الرز،...) لتنفيذ النموذج المقترن وذلك لما تحمله هذه البذور من خصائص وصفات مثل (نسبة الكربوهيدرات، الدهون، البروتين، الألياف،... الخ). وقد تم اعتماد نسبة وجود البروتين في البذور بوصفها صفة لدعم التصنيف فضلاً عن الصفات الأخرى، وتمييز أصناف بذور النباتات عن بعضها البعض والحصول على التصانيف الصحيحة.

والجدول الآتي يوضح أنواع الحبوب مصنفة وفقاً لنسب البروتين الموجود في كل منها.

الجدول ٢

نسب البروتين الموجودة في بذور بعض النباتات

النسبة المئوية للبروتين

نوعية الحبوب
الرز الخام (الشلب)
الحنطة
حبة الشوفان الكاملة
الشعير
ذرة صفراء حلوة
ذرة بيضاء
الشيلم

١٠ - ٩
١١.٥ - ١٠.٥
١١.٨ - ١١.٦
١٢.٠ - ١١.٨
١٢.٣٠ - ١٢.١٠
١٣.٠٠ - ١٢.٤٠
١٣.٥ - فما فوق

المصدر: الفخري، واحمد صالح خلف، ١٩٨٣، ٢٩.

٢ - ٢ الخوارزميات المستخدمة في البحث

١. المصنف C4.5

وهي إحدى أنواع المصنفات التي تعمل بمشرف، وهي من نوع شجرة القرار، اذ تستخدم هذه الخوارزمية مبدأ تقسيم من الأعلى إلى الأسفل (Top – Down) وذلك وصولاً إلى الحل الأمثل، و من مميزاتها انها تتعامل مع الارقام، الخصائص، الفئيم المفقودة والبيانات المشوشة. فضلاً عن وصفها بأنها من أفضل خوارزميات التصنيف واكثرها دقة و سرعة في الوصول الى الحل النهائي. (Llora and et al, 2001, 4).

٢. الشبكات العصبية BP

تعتمد فكرة الشبكات العصبية الاصطناعية على ايجاد منظومة حسابية لها القدرة على التكيف و التعديل عن طريق التعلم و ذلك سعياً لايجاد دوال الربط بين المدخلات و المخرجات، او استنتاج قرار مبني على آلاف الاحتمالات و العلاقات التي تشكل بدورها ملفاً تاريخياً تبني من خلاله العلاقة او الدالة (العيدي ، ٢٠٠٠ ، ٧٥).

وتعد شبكة الانتشار الخلفي Back Propagation من أكثر الشبكات العصبية شيوعاً واستخداماً، إذ يجري تعديل اوزان الشبكة وتحسين ادائها من خلال دالة التعليم للوصول الى افضل نتيجة او نتيجة مقاربة. عليه اعتمد البحث هذا النوع من الشبكات BP في الحصول على النتائج وقد تم شرح المعادلات المستخدمة في الجانب النظري من البحث.

٢-٣ مراحل تنفيذ البرنامج المصمم

بعد وصف الطرائق المستخدمة في التصنيف فقد أعتمد البحث لغة البرمجة V.B في كتابة برامجه المتمثلة بربط نتائج الشبكة العصبية BP مع المصنف C4.5 وذلك للحصول على النسبة النهائية للمعلومات التي يتم التصنيف من خلالها الى فئات وذلك باعتماد البيانات المشار اليها في الجدول ٢.

الجدول ٣

جدول يوضح نتائج تنفيذ معادلات المصنف C4.5 من برنامج الـ V. Basic

التصنيف	نسبة المعلومات النهائية Gain	نسبة المعلومات المستحصلة من المعادلة E(APn)	عدد العناقيد	ت
حبة الشوفان الكاملة	١١.١	٢٣.٩٦٢٤	١٠	١
حنطة	١٠.٤	٢٢.٥٥٣٠	١٠	٢
حبة الشوفان الكاملة	١١.١	٢٤.٠٢٨٤	١٠	٣

التصنيف	نسبة المعلومات النهائية Gain	نسبة المعلومات المستحصلة من المعادلة E(APn)	عدد العناقيد	ت
حنطة	١١.٠١	٢٤.١٣٧٩	١٠	٤
رز خام	١٠.١	٢٢.١٩٧٤	١٠	٥
رز خام	٩.٧	٢٢.٩٣٧٨	٨	٦
حنطة	١٠.٧	٢٣.٣٠١٨	٨	٧
رز خام	١٠.٣	٢٥.٩٨٨٩	٨	٨
حنطة	١١.١	٢٤.٢٠٧٠	٨	٩
حنطة	١١.١	٢٣.٩٦٢٩	٦	١٠
حنطة	١٠.٧	٢٣.٢٤٥٤	٦	١١
رز خام	٩	٢٥.١٤٢٨	٦	١٢
رز خام	١٠	٢٣.٦٢١٩	٦	١٣
حبة الشوفان	١١	٢٤.٠١٤٩	٦	١٤
رز خام	١٠.١	٢٣.٣٠٥٣	٦	١٥
رز خام	١٠.٠١	٢٢.٠٧٨١	٤	١٦
رز خام	١٠.١	٢٣.٣٠١٨	٤	١٧
رز خام	٩	٢١.٦٩٧٨	٤	١٨
حبة الشوفان	١١	٢٤.٠١٦٥	٤	١٩
رز خام	٩	٢١.٥٦٨٠	٣	٢٠
رز خام	١٠	٢٢.٠٧٨١	٣	٢١
رز خام	١٠	٢٣.٥٨٦٨	٣	٢٢
رز خام	٩.٧٥	٢٢.٩٣٧٨	٨	٢٣
رز خام	٩.٥	٢٢.٤٢٤٥	٨	٢٤
رز خام	٩.٤٥	٢٢.٢٤٦٧	٨	٢٥
رز خام	٩.٧٧	٢٢.٨٥٥٧	٨	٢٦
رز خام	٩.٣٥	٢٢.٩٣٨٤	٨	٢٧

يتبع

ما قبله

رز خام	٩.٥٩	٢٢.٨٨٢٢	٨	٢٨
حنطة	١١.٢٠	٢٤.٤٨٦٥	٨	٢٩
حنطة	١٠.٩٧	٢٤.٠٣٧١	٨	٣٠
حنطة	١١.١٤	٢٤.٣٦٥٧	٨	٣١
حنطة	١١.١٦	٢٤.٤١٠٣	٨	٣٢
حنطة	١٠.٩٠	٢٣.٩٢٤١	٨	٣٣
رز خام	١٠.٠٢	٢٢.٣١٢٢	٨	٣٤
شعير	١١.١١	٢٣.٩٦٥٠	٨	٣٥
رز خام	١٠.٣١	٢٦.٠٦٦٤	٨	٣٦

حطة	١١.٠٠	٢٣٠٨٣٧	٨	٣٧
حطة	١١.٢	٢٤٤٨٦	٤	٣٨
حطة	١٠.٩	٢٤٠٣٧١	٤	٣٩
حطة	١١.٠٨	٢٤١٣٧٩	٤	٤٠
حطة	١١.١٨	٢٤٤٤٠٢	٣	٤١
حطة	١١.١٨	٢٤٤٢٦٦	٣	٤٢
حطة	١٠.٨٠	٢٣٣٥٦٧	٣	٤٣
رز خام	١٠.٢٤	٢٢٣٧٢٦	٣	٤٤
رز خام	١٠.٢٥	٢٢٣٧٢٦	٣	٤٥
رز خام	٩.٥٨	٢١٥١٢٤	٢	٤٦
رز خام	٩.٦١	٢١٥٥٨٢	٢	٤٧
ذرة بيضاء	١٢.٧٢	٢٩٠٢٠٢	٢	٤٨

الجدول ٤

النتائج من برنامج Minitab بعد تنفيذ المعادلات

التصنيف	نسبة المعلومات النهائية	عدد العناقيد	ت
رز خام	٩ ≈ ٨.٩	١٠	١
رز خام	٩.٢	١٠	٢
ذرة بيضاء	١٣.٦	١٠	٣
رز خام	٩.٣	١٠	٤
حطة	١٠.٨	١٠	٥
حطة	١٠.١	٨	٦
حطة	١٠.٥	٨	٧
شيلم	١٤.٥	٨	٨
رز خام	٩.٥	٨	٩

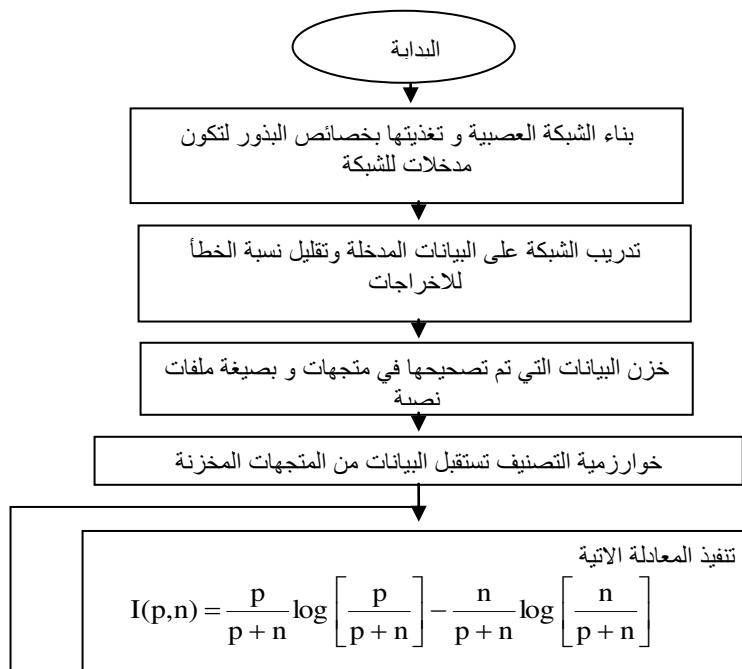
يتبع ←

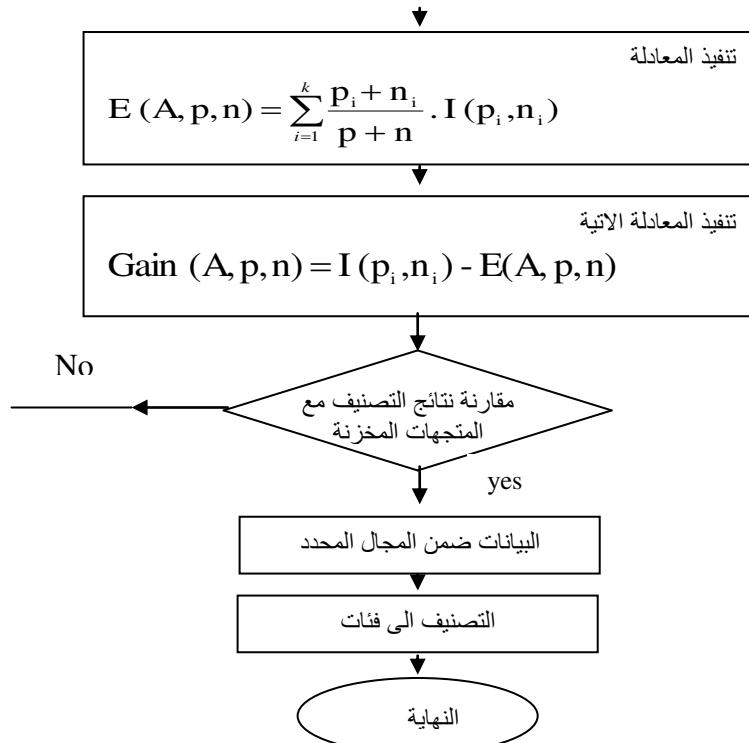
← ماقبله

رز خام	٩.٢	٦	١٠
حطة	١٠.٥	٦	١١
شيلم	١٣.٧	٦	١٢
رز خام	٩.٢	٦	١٣
رز خام	٩.٣	٦	١٤
شعير	١١.٦	٦	١٥
ادنى من المستوى	٨.٨	٤	١٦
رز خام	١٠.٢	٤	١٧
حطة	١٠.٥	٤	١٨
حطة	١٠.٦	٤	١٩
رز خام	٩	٣	٢٠

ادنى من المستوى	٧.١	٣	٢١
ادنى من المستوى	٦.٦	٣	٢٢
رز خام	٩.٧	٨	٢٣
ادنى من المستوى	٢.٧	٨	٢٤
ذرة بيضاء	١٢.٦	٨	٢٥
ادنى من المستوى	٧.٥	٨	٢٦
ادنى من المستوى	٧.٧-	٨	٢٧
اعلى من الحد الطبيعي	٣٩.٤	٨	٢٨
اعلى من الحد الطبيعي	٣٩.٤	٨	٢٩
اعلى من الحد الطبيعي	٤٢.٥	٨	٣٠
اعلى من الحد الطبيعي	٤٣.٢	٨	٣١
اعلى من الحد الطبيعي	٣٦.١	٨	٣٢
اعلى من الحد الطبيعي	٣٢.٤	٨	٣٣
اعلى من الحد الطبيعي	٤٢.٧	٨	٣٤
اعلى من الحد الطبيعي	٤٤.٦	٨	٣٥
اعلى من الحد الطبيعي	٣٨.٨	٨	٣٦
اعلى من الحد الطبيعي	٤٢.٧	٨	٣٧
اعلى من الحد الطبيعي	٤٤.٦	٨	٣٨
اعلى من الحد الطبيعي	٤٦.١	٨	٣٩
اعلى من الحد الطبيعي	٤١.٢	٤	٤٠
اعلى من الحد الطبيعي	٣٧.٤	٤	٤١

ولفهم عمل البرنامج المصمم يمكن تتبع المخطط الانسيابي في الشكل ٢ ، والذي يبيّن الخطوات كافة التي تمر بها الخوارزمية وصولاً إلى النتائج.





٢-٤

مخطط اسيابي يبين مراحل البرمجة المستخدمة في البحث

يمثل نتائج تنفيذ البرنامج المصمم والمكتوب بلغة V.B والمتضمن استخدام المصنف C4.5، في حين يظهر الجدول الثاني جدول ٤ نتائج تنفيذ الحزمة البرمجية الجاهزة Minitab Release 13.0.

من ملاحظة نتائج التصنيف في الجدول ٣ تبين بأن النتائج المستحصلة من تطبيق البرنامج المعد من قبل الباحثين تقع جميعها ضمن فئات التصنيف الصحيحة، بعبارة أخرى فإن نسبة المعلومات النهائية Gain والتي تمثل بالنتيجة نسبة البروتين الموجود في البذور قيد التصنيف، ولجميع الحالات التي تم تنفيذ البرنامج عليها تقع ضمن فئات الجدول ٢ ولا توجد أي حال شاذة في نتائج تنفيذ الخوارزمية C4.5. وهذه النتيجة تشير إلى توصل الخوارزمية C4.5 إلى حالات التصنيف الصحيحة أو المقاربة في جميع الحالات .

في حين تشير نتائج الجدول ٤ إلى وجود بعض من الحالات غير الصحيحة (حالات شاذة) اذ كانت النسبة النهائية للبروتين دون المستوى المطلوب او انها أعلى بكثير من المستوى. هذا يعني ان استخدام الحزمة الجاهزة Minitab قد لا يؤدي الى نتائج صحيحة في التصنيف، ناهيك عن استخدام بعض المعالجات (معالجة رياضية) على النتائج المستحصلة من تطبيق دالة التصنيف للحصول على النتائج النهائية وهذا

ما ينسجم مع فرضية البحث ويفرز من صحة النتائج المستحصلة من تطبيق البرنامج المصمم بلغة البرمجة V.B.

٤- الاستنتاجات

تم التوصل من خلال نتائج تنفيذ البرنامج المصمم والنتائج المستحصلة من تطبيق الحزمة البرمجية الجاهزة Minitab الى جملة من الاستنتاجات ندرجها بالاتي:

١. ان المعادلات الخاصة بالمصنف C4.5 كانت اكفاء في الاداء وهذا ما اظهرته نتائج تنفيذ البرنامج بلغة V.B .
٢. ان ربط نتائج الشبكة العصبية BP بخوارزمية التصنيف ادى الى ازالة التناقض والتشویش الموجود في البيانات، اذ تعمل الشبكة العصبية ومن خلال مصفوفات الاوزان الى ازالة هذه الحالات من البيانات .
٣. من ملاحظة المعادلات المستخدمة في المصنف C4.5، فان الحدود التي تحتوي على دالة اللوغاريتم قد تؤدي الى عدم الوصول الى حل نهائي وذلك في الحالات التي يكون فيها دليل الدالة Log سالبة، وقد تمت معالجة هذه الحالات الجاهزة نتيجة لمثل هذه الحالات .
٤. لاتعد النتائج المستحصلة من تطبيق الحزمة الجاهزة Minitab نتائجاً نهائية، اذ تتطلب عملية التصنيف القيام بالعديد من العمليات الحسابية للحصول على النتيجة النهائية .
٥. ان كانت الحالات التي تم تنفيذ البرنامج المصمم عليها حصلت على تصنيف، وذلك من خلال الحصول على نسبة المعلومات النهائية ضمن حدود الفئات الواردة في الجدول ٢ في حين ظهرت العديد من الحالات غير المصنفة لدى استخدام الحزمة البرمجية الجاهزة Minitab.

المراجع

اولاً- المراجع باللغة العربية

١. محمود خليل ابراهيم العبيدي، "الشبكات العصبية الاصطناعية" مجلة ابحاث الحاسوب، المجلد الرابع، العدد الاول، ٢٠٠٠ .
٢. عبدالله قاسم الفوري، و احمد صالح خلف، "بذور المحاصيل انتاجها ونوعيتها"، دار الكتب للطباعة والنشر، ١٩٨٣ .
٣. نعمة عبدالله الفوري، "استخلاص نموذج بياني من قاعدة بيانات باستخدام خوارزميتي K-means "، رسالة ماجستير، كلية علوم الحاسوب والرياضيات / جامعة الموصل ، ٢٠٠٣ .
٤. محاضرات مأخوذة من شبكة المعلومات الدولية (الانترنت) :

1- Data Mining Glossary Courses .

www.towcrows.com/glossary.html/2004.

2- Data Mining : What is data mining .

www.Anderson.Ucla.edu/faculty/jason.f.Fraud/teacher/technologies/2004.

ثانياً- المراجع باللغة الاجنبية

1. Ali, S., and Abraham A. (2002), "An Empirical comparison of kernel selection for support vector machine", Gippsland school of computing and information technology, Monash university, Victoria-Australia.
2. Gómez Hidalgo J. M., Mana López M., and Puertas Sanz E. (2002), "Evaluating cost-sensitive un solicited Bulk Email categorization", Journees internationales d'Analyse statistique des Données Textuelles (JADT) - <http://citeseer.nj.nec.com>
3. Ibrahim R. S. (1999), "Data Mining of machine learning performance data" (M.Sc. thesis), RMIT university-Australia.
<http://goanna.cs.rmit.edu.au/~vc/papers/ibrahim-mbc>.
4. Kinnebrock W.,(1995),"Neural Networks: Fundamentals, Applications, Examples", Galgotia publications Pvt. LTD, NewDelhi .
5. Keller F. (2000), "Introduction to machine learning", Journal of Machine Learning Research (JMLR).
http://www.aai.org/AI_Topics/html/mahine.html
6. Llora X., and Garrell J. M. (2001), "Knowledge-Independent data mining with fine-grained parallel evolutionary algorithms", Enginyeria I Arquitectura La Salle, Uinveristat Ramon Llull.
<http://gal4.ge.uivc.edu/~xllora/curriculum/cu>

ABSTRACT

The Use of The Classification Algorithms C4.5 in Distinguishing

Data mining is an action to obtain the knowledge to achieve the main purpose that is detecting hidden facts which contain database by using various techniques that include artificial intelligent, statistic analysis, techniques and modeling.

Data mining method brings models and obvious relations in data which help to expect future results. Algorithm appears in this field which indicates comparisons between algorithm to choose the suitable one to have the best results.

This paper aims at using classification algorithms. C4.5 and connect it with neural nets (Back propagation BP) to form a classification model hold a two methods characteristics. In addition to comparing the total results with the classification results by using Minitab.

This paper results in classified equation C4.5 which is sufficient in performance especially after connecting it with the neural net BP to erase the contradictions in data.

The results also confirm the importance to use programming language in achieving the best in comparison with the results of ready made applications.