

## استخدام تقنيات تنقيب البيانات لكشف التطفل في شبكات الحاسوب

أ.د. السمانى عبد المطلب أحمد - باحث ثانى

أ.إياد محمد مهيبوب غالب البريهي - باحث أول

جامعة النيلين، كلية علوم الحاسوب وتقانة

المعلومات، الخرطوم، جمهورية السودان

### الملخص

# 2

أمان وسرية المعلومات في الشبكات تعتبر القضية الرئيسية المؤثرة لكثير من الشركات والمؤسسات التي تستخدم كمية كبيرة من البيانات، وبالمقابل هناك العديد من الطرق المستخدمة لحماية الشبكات في الوقت الراهن مثل التشفير والـ VPN، وجدران الحماية، ولكن كل هذه الطرق تعتبر طرق استاتيكية جداً للحماية الفعالة ضد المهددات وعدد المهددات، وتستخدم تقنيات تنقيب البيانات لهذا الغرض حيث تطبق لكشف التطفل intrusion detection.

تهدف هذه الورقة لإستخدام تقنيات تنقيب البيانات (data mining) لكشف حالات الشذوذ في الشبكات تطبيقاً على عينتين عشوائيتين من مجموعة بيانات NLS-KDD Data Set، استخدمت الورقة تقنية التصنيف مصنف شجرة القرار (decision tree) التي تنفذ خوارزمية C4.5، نتيجة التجربة تعرض أن مصنف C4.5 أظهر نتائج فعالة لكشف التطفل في الشبكات الحاسوبية وأظهرت النتائج أنه كلما كانت كمية البيانات كبيرة تكون نسبة الخطأ أقل ودقة التنبؤ عالية.

الكلمات المفتاحية : تنقيب البيانات، التطفل، التصنيف، شجرة القرار، اكتساب المعلومات، Gain Ratio، gini index.

## مقدمة

أمن أنظمة الحواسيب الخاصة بنا والبيانات في خطر مستمر، ولقد أدى النمو الواسع النطاق للإنترنت وازدياد توفر أدوات وحيل التطفل ومهاجمة الشبكات للتوجه لكشف ومنع التطفل، مما أصبح عنصراً أساسياً وفعالاً في أنظمة الشبكات.

يمكن تعريف التطفل ( Intrusion ) بأنه مجموعة من الأحداث والمهددات التي تهدد السرية والسلامة أو توافر موارد الشبكة كحساب المستخدم، أنظمة الملفات، نواة النظام ... الخ[١].

يتضمن كشف التطفل تحديد مجموعة من الأحداث الخبيثة التي تهدد سلامة وسرية، وتوافر موارد المعلومات. تستند الطرق التقليدية لكشف التسلسل على المعرفة الواسعة لتوقعات الهجمات المعروفة.

تقنيات كشف التسلسل المعتمدة على تنقيب البيانات تستند عموماً على واحدة من فئتين، الكشف عن سوء الاستخدام misuse detection أو الكشف عن الشذوذ anomaly detection. في misuse detection يسمى كل instance في مجموعة البيانات إما 'normal' أو 'intrusion' ويتم تدريب خوارزمية التعلم على بيانات التدريب، وهذه التقنيات قادرة على إعادة تدريب نماذج كشف التسلسل تلقائياً على بيانات الإدخال المختلفة التي تشمل أنواع جديدة من الهجمات ما دام قد وصفت بشكل مناسب.

الكشف عن الشذوذ anomaly detection ، من ناحية أخرى ، يبني نماذج للسلوك الطبيعي normal behavior ، وبشكل تلقائي يتم الكشف عن أي انحراف عن ذلك السلوك.

أنظمة كشف التطفل أو أنظمة منع التطفل كليهما تعمل على مراقبة حركة مرور الشبكة ومنع الأنشطة الخبيثة، وغالبية أنظمة كشف ومنع التطفل تستخدم إما الكشف المعتمد على التوقيع signature-based detection أو الكشف المعتمد على الشذوذ anomaly-based detection .

أ. الكشف المعتمد على التوقيع signature-based detection : هذه الطريقة تستخدم للكشف عن التوقعات التي هي أنماط هجوم مكونة مسبقاً ومحددة سلفاً من قبل خبراء المجال، ونظام منع الاختراق المستندة إلى التوقيع تراقب حركة مرور الشبكة لمطابقتها مع هذه التوقعات، وإذا تم العثور مرة واحدة على تطابق فإن نظام كشف التسلسل سوف يبلغ عن وجود الشذوذ وسيتخذ نظام منع الاختراق إجراءات إضافية مناسبة.

ب. الكشف المعتمد على الشذوذ Anomaly-based detection : هذه الطريقة تبني نماذج لسلوك الشبكات الطبيعية تسمى ( تشكيلات profiles ) تستخدم بعد ذلك للكشف عن الأنماط التي تحيد كثيراً عن هذه التشكيلات، ومثل هذه الانحرافات قد تمثل الإختراقات الفعلية أو ببساطة السلوكيات الجديدة التي تحتاج أن تضاف إلى هذه التشكيلات. الميزة الرئيسية للكشف المعتمد على الشذوذ قدرته على

كشفت الاختراقات غير المألوفة التي لم يتم ملاحظتها، وعادة يجب على المحلل فرز الانحرافات للتأكد أيا منها تمثل الاختراقات الحقيقية.

ستتناول هذه الورقة استخدام تقنية التصنيف كتقنية من تقنيات تنقيب البيانات لتصنيف بيانات الشبكة لعينتين عشوائيتين من مجموعة البيانات NLS- KDD99 data set العينة الأولى تحتوي على ( 5,290,866 byte ) كبيانات تدريب و(946,848 byte) كبيانات اختبار وتحتوي العينة العشوائية الثانية على كمية بيانات ( 1,058064 byte ) كبيانات تدريب و ( 497,700 byte ) كبيانات اختبار يتم تدريب الخوارزمية C4.5 بواسطة أداة الويكا Weka Tools المصممة لتدريب واختبار خوارزميات تنقيب البيانات، يتم تقييم مصنف C4.5 بواسطة حساب دقة المصنف ونسبة الخطأ الناتجة من مصفوفة التضارب Confusion Matrix .

٢ - استخدام تنقيب البيانات في كشف التطفل :

لتعزيز أمن الشبكة هناك العديد من الخوارزميات المستخدمة في كشف التسلسل مثل التشفير، و VPN، وجدار الحماية،... الخ، وهذه الخوارزميات توفر الحماية بشكل ثابت عبر الشبكة، ولكن كشف التسلسل عملية ديناميكية ولذلك يجب توفير حماية ديناميكية عبر الشبكة [٤]. يمكن أن تساعد تقنيات تنقيب البيانات في تعزيز أداء أنظمة حماية ومنع التطفل بعدة طرق أهمها :

٢ - ١ خوارزميات تنقيب بيانات جديدة لكشف التطفل :

يمكن استخدام خوارزميات تنقيب البيانات للكشف المعتمد عن التوقيع والكشف المعتمد عن الشذوذ معاً.

٢-١-١ في خوارزميات الكشف المعتمد على التوقيع تصنف بيانات التدريب إما "طبيعية normal" او "متطفلة intrusion"، ويمكن بعد ذلك اشتقاق المصنف الذي يكتشف الاختراقات المعروفة في هذا المجال.

٢-١-٢ في خوارزميات الكشف المعتمد على الشذوذ يتم بناء نماذج السلوك الطبيعي normal (behavior) للبيانات وبشكل تلقائي يتم الكشف عن الانحرافات الهامة منها.

٢-٢ تحليل الربط Association، والارتباط correlation ونمط التمييز

تساعد هذه الخوارزميات في اختيار وبناء المصنفات التمييزية وتطبق لإيجاد العلاقة بين خصائص النظام system attributes التي تصف بيانات الشبكة، ومثل هذه المعلومات يمكن أن توفر الرؤيا الواضحة فيما يتعلق باختيار سمات مفيدة لكشف التسلسل.

## ٢-٣ تحليل البيانات المجدولة :Analysis of stream data

نظراً لطبيعة تنقل وحركة التطفلات والهجمات الخبيثة لا بد من إجراء الكشف عن التطفلات في بيئة البيانات المجدولة، وبالرغم من أن الحدث قد يكون طبيعياً من تلقاء نفسه إلا إنه اعتبر ضاراً إذا نظر إليه كجزء من تسلسل الأحداث، وبالتالي فمن الضروري دراسة ماهي تسلسلات الاحداث المتصادفة معاً بشكل متكرر ولا بد من إيجاد الأنماط المتسلسلة وتحديد القيم المتطرفة. طرق تعدين البيانات الأخرى ضرورية أيضاً لإيجاد التجمعات المتطورة وبناء نماذج تصنيف ديناميكية في جداول البيانات في الوقت الحقيقي لكشف التسلسل.

## ٢-٤ طرق تنقيب البيانات الموزعة :

التطفلات يمكن إطلاقها من عدة مواقع مختلفة لتوجه إلى عدة وجهات مختلفة، وفي هذه الحالة يمكن استخدام طرق تنقيب البيانات الموزعة لتحليل بيانات الشبكة من عدة مواقع للشبكة للكشف عن التهديدات الموزعة. [٣]

## ٣ - أنواع مهددات التطفل Type Intrusion Attacks

ستتناول هذه التجربة في هذه الورقة الكشف عن أربعة مهددات تهاجم الشبكات الحاسوبية، حيث سيتم الكشف عنها بواسطة تقنية التصنيف كتقنية من تقنيات تنقيب البيانات بواسطة تنفيذ خوارزمية C4.5 على قاعدة بيانات NLS-KDD99 Data Set.

## ٣-١ مهدد إنكار الخدمة (DoS) Denial of Service Attack

هو المهدد الذي يجعل مورد الحوسبة او الذاكرة مشغولة جداً او ممتلئة للتعامل مع أي طلبات مشروعة أو حرمان المستخدمين الشرعيين من الوصول إلى الخدمة .

## ٣-٢ مهدد المستخدم الرئيسي (U2R) User to Root Attack

هو المهدد الذي يتظاهر بأنه المستخدم الشرعي للنظام دون إذن، ثم يستغل نقاط ضعف النظام للحصول على صلاحية الوصول إلى مصدر التحكم بالنظام على سبيل المثال امتلاك صلاحيات تنزيل البرامج المخفية التي قد تتسبب في فقد النظام أو تنفيذ برامج هجومية كما لو أنها جزء من برامج النظام الأصلية.

## ٣-٣ مهدد محلي عن بعد (R2L) Remote To Local Attack

هم المهددين غير المصرح لهم خلال الشبكات ويحصلون على وصول محلي للشبكة كمستخدمي الأجهزة المحلية، ويمكنهم شن الهجمات من أي مكان على الانترنت، وفي حالة ان مستخدماً ما أصبح لديه حق الوصول إلى نظام المعلومات فإنه يمكنه استغلال نقاط الضعف في الجهاز المحلي لسرقة البيانات الهامة أو تدمير نظام المعلومات. أي خدمة تحتاج لكلمة سر للوصول إليها تعتبر هدفاً لهذا

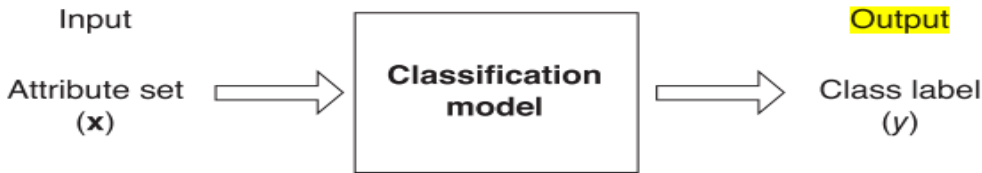
المهدد. هجمات R2L من أصعب الهجمات كشافاً لأنها تنطوي على مستوى الشبكة ومميزات مستوى المضيف.

### ٣-٤ مهدد التدقيق : PROBE ATTACK

يتم الهجوم من قبل المهاجمين باستخدام برامج فحص تلقائية لعدد كبير من عناوين ال IP Address الخاص بالشبكات من أجل العثور على نقاط ضعف يمكن استغلالها، فإذا عثر على نقاط ضعف لمرة واحدة فإنه يمكن للمهاجمين استغلالها لكسب الوصول إلى الشبكة والبدء بجمع المعلومات بدون إذن.[٤]

### ٤- تقنية التصنيف classification

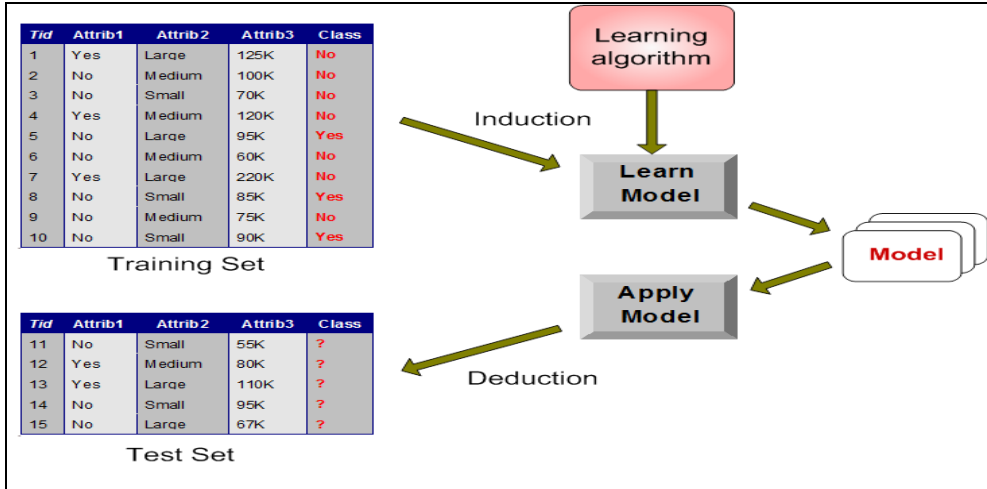
التصنيف هو مهمة تعليم دالة معينة  $f$  لربط مجموعة الخصائص  $(X)$  بفئة  $(class)$  معرفة مسبقاً تسمى  $Y$ ، تعرف الدالة  $f$  أيضاً بنموذج التصنيف  $classification\ model$ .



شكل (١) التصنيف

تتم عملية بناء المصنف بخطوتين رئيسيتين هما التعليم (خطوة بناء نموذج التصنيف) والتصنيف (استخدام النموذج ليتنبأ بفئات البيانات غير المعروفة).

- أ. في الخطوة الأولى: يتم تدريب خوارزمية التصنيف (التعلم)  $classifying\ algorithm$  على بيانات التدريب  $training\ data$  المحتوية على سجلات معروفة لبناء المصنف الذي سيستخدم لفحص بيانات الفحص  $data\ test$  التي تحتوي على سجلات غير معروفة.
- ب. الخطوة الثانية: يقيم أداء المصنف بحساب عدد السجلات المتوقعة المصنفة بشكل صحيح والسجلات المصنفة بشكل خطأ فيما يسمى بمصفوفة التعارض  $confusion\ matrix$ ، يتم تقييم أداء نماذج المصنفات بالحصول على خوارزميات تصنيف تسعى للحصول على أعلى دقة و أقل نسبة خطأ عند تطبيقها على بيانات الاختبار  $data\ test$ . [٥]



شكل (٢) خطوات بناء مصنف البيانات

٤- ١ حساب مصفوفة التعارض confusion matrix

يتم جدولة عدد السجلات المصنفة بشكل صحيح وعدد السجلات المصنفة بشكل خطأ على شكل مصفوفة تعرف بمصفوفة التعارض.

		Predicted class	
		Class = 1	Class = 0
Actual class	Class = 1	TP	FN
	Class = 0	FP	TN

شكل (٣) مصفوفة التعارض الثنائية [١] binary confusion matrix

كل مدخل  $f_{ij}$  في مصفوفة التعارض يشير إلى عدد السجلات في  $i$  class المتوقع تكون في  $j$  class. مثلا  $f_{01}$  تشير إلى عدد السجلات في  $i=0$  class التي تم توقعها بشكل خطأ في  $j=1$  class.

اعتماداً على المصنوفة فإن :

مجموع السجلات المتوقعة بشكل صحيح هي  $(f_{11} + f_{00})$  .

مجموع السجلات المتوقعة بشكل خطأ هي  $(f_{01}+f_{10})$  .

تحسب الدقة Accuracy لتقييم أداء المصنفات اعتماداً على مصنوفة التعارض من المعادلة

التالية

$$Accuracy = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}} = \frac{f_{11} + f_{00}}{f_{11} + f_{10} + f_{01} + f_{00}}$$

وبشكل مكافئ يحسب معدل الخطأ من المعادلة التالية :

$$Arrore\ rate = \frac{\text{Number of wrong predictions}}{\text{Total number of predictions}} = \frac{f_{01} + f_{10}}{f_{11} + f_{10} + f_{01} + f_{00}}$$

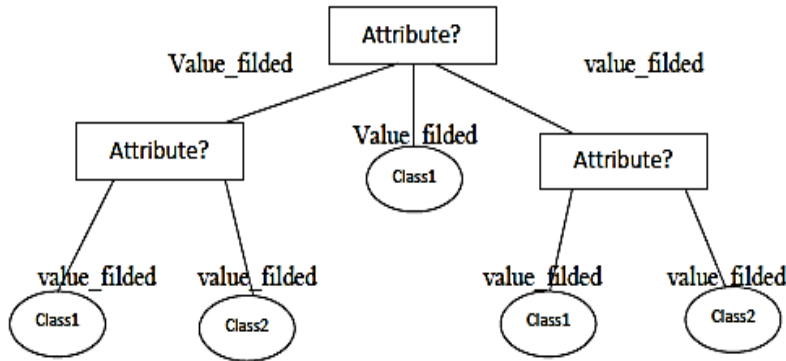
## تقنيات التصنيف Classification Techniques

٤- ٢

يتضمن التصنيف عدة تقنيات رئيسية ستناقش هذه الورقة أهم هذه التقنيات وهي تقنية شجرة القرار Decision Tree .

### ٤- ٢- ١- مصنف شجرة القرار Decision Tree classifier

أشجار القرار تصنف الحالات عن طريق فرزها على أساس قيم الصفة attribute. كل عقدة داخلية في شجرة القرار تمثل صفة اختبار test attribute، كل فرع يمثل قيمة العقدة، وكل عقدة طرفية (ورقة) تمثل مخرجات الاختبار (قيمة value)، أعلى عقدة في الشجرة هي جذر الشجرة، وتصنف الحالات انطلاقاً من عقدة الجذر في شجرة القرار، يتم تعيين كل عقدة ورقية كفئة class.



الشكل (١- ١٧) شجرة قرار التصنيف [4]

الخوارزميات ID3، C4.5، CART تستخدم أسلوب فرق تسد من أعلى إلى أسفل (top-down recursive divide-and-conquer manner) لتوليد أشجار القرار التي تبدأ بالتدريب على مجموعة بيانات التدريب والفئات (class labels) المرتبطة بها.

#### ٤- ٢- ١- بناء شجرة القرار How to Build a Decision Tree

من حيث المبدأ، هناك عدد من أشجار القرار يمكن بناؤها من مجموعة معينة من الصفات المتوفرة، بينما كثير من أشجار القرار تكون أكثر دقة من الأخرى، الحصول على الشجرة الأمثل مناسب حسابياً بسبب الكم الهائل من حجم عمليات البحث التي تقوم بها الخوارزمية، ومع ذلك فقد تم تطوير خوارزميات فعالة للبحث على الدقة بشكل معقول. واحدة من هذه الخوارزميات خوارزمية هانت Hunt's algorithm، التي هي أساس العديد من خوارزميات حث شجرة القرارات القائمة، بما في ذلك ID3، C4.5، CART. يعرض هذا القسم مناقشة ريفية المستوى لخوارزمية هانت [5].

#### ٤- ٢- ١- خوارزمية هانت Hunt's Algorithm

في خوارزمية هانت يتم توليد شجرة القرار باتباع أسلوب التكرار الذاتي عن طريق تقسيم سجلات التدريب إلى مجموعات فرعية بشكل متتابع.

ليكن  $D_t$  مجموعة سجلات التدريب المرتبطة بالعقدة  $N$  و  $y = \{y_1, y_2, \dots, y_n\}$  عبارة عن أصناف التصنيف class labels:

١ - إذا كانت كل السجلات في  $D_t$  تنتمي إلى نفس الـ class  $y_t$  فإن  $t$  عقدة ورقية مرتبطة كـ class  $y_t$ .

٢ - إذا Dt تحتوي سجلات تنتمي إلى أكثر من class ، يتم اختيار شرط اختبار الصفة (attribute test condition) لتقسيم السجلات إلى مجموعات فرعية صغيرة. يتم توليد العقدة الابن (child node) لكل مخرج من مخرجات شرط الاختبار ويتم توزيع سجلات Dt لعقد الأبناء اعتماداً على المخرجات. يتم تطبيق الخوارزمية بشكل تكرار ذاتي لكل عقد الأبناء child node [٥]

#### ٤- ٢- ١- ٣- معايير اختيار الصفة Attribute Selection Measures

بناء شجرة القرارات، بجانب اختيار خوارزمية البناء يجب أن تؤخذ بعين الاعتبار معايير اختيار الصفة المناسب لاختيار معيار التقسيم الأفضل الذي يقسم بيانات معينة D، من المجموعات التدريبية المرتبطة بصف (class) إلى مجموعات فردية. ستناقش هذه الورقة ثلاثة معايير عامة لاختيار الصفة المناسبة للتقسيم الأفضل Information Gain ، Gain Ratio ، Gini Index .

#### ٤- ٢- ١- ٣- اكتساب المعلومة Information Gain

يعرف اكتساب المعلومة بالانخفاض المتوقع للمعلومات اللازمة لتصنيف الصفة A. الصفة التي تكون قيمة اكتساب المعلومة لها أعلى من قيمة الصفة الأخرى يتم إختيارها كصفة تقسيم للعقدة N في مجموعة البيانات D [١].

ولحساب اكتساب المعلومة لصفة منفصلة A في مجموعة البيانات D نحسب الآتي :

١- المعلومات المتوقعة اللازمة لتصنيف سجل في D أو ما يسمى Entropy(D) من المعادلة التالية :

$$\text{Info}(D) = \sum_{i=1}^m p_i \log_2(p_i) \quad \dots(1)$$

$$\text{أو Entropy}(D) = \sum_{i=1}^m p_i \log_2(p_i)$$

حيث  $p_i$  احتمالية انتماء السجل في D إلى class  $C_i$  وتحسب من العلاقة الآتية :

$$p_i = \frac{|C_{i,D}|}{|D|}$$

حيث  $C_{i,D}$  عدد السجلات في D المنتمية إلى class C ، D : العدد الكلي للسجلات .

تستخدم دالة اللوغاريتم بالاساس ٢ لان المعلومات ترمز بال bits .

- Info(D) هي متوسط كمية المعلومات اللازمة لتحديد فئة السجل i في D .

ب - نحسب المعلومات المتوقعة لل attribute A في D حيث A تحتوي على عدد v من القيم

المنفصلة حيث  $v=\{a_1, a_2, \dots, a_v\}$  من المعادلة التالية :

$$info_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times info(D_j) \dots \dots \dots (2)$$

$$Entropy_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times Entropy(D_j)$$

حيث:  $\frac{|D_j|}{|D|}$  وزن الجزء  $j$  بالنسبة لقيم الصفة  $A$ ،  $D_j$  عدد ال tuples التي تنتمي لل class التابعة لـ attribute  $A$  في  $D$ .

ج- لحساب معادلة Information Gain(A) نحسب الفرق بين المعادلة (1) والمعادلة (2)

$$Gain(A) = info(D) - info_A(D) \dots \dots (3)$$

No.	Duration	dst_host_same_src_port_rate	dst_host_srv_diff_host_rate	Flag	src_bytes	dst_bytes	Class
1	0.0	0.17	0.0	SF	491.0	0.0	Normal
2	0.0	0.88	0.0	SF	146.0	0.0	Normal
3	0.0	0.0	0.0	S0	0.0	0.0	Anomaly
4	0.0	0.03	0.04	SF	232.0	8153.0	Normal
5	0.0	0.0	0.0	SF	199.0	420.0	Normal
6	0.0	0.0	0.0	REJ	0.0	0.0	Anomaly
7	0.0	0.0	0.0	S0	0.0	0.0	Anomaly
8	0.0	0.0	0.0	S0	0.0	0.0	Anomaly
9	0.0	0.0	0.0	S0	0.0	0.0	Anomaly
10	0.0	0.0	0.0	S0	0.0	0.0	Anomaly
11	0.0	0.0	0.0	REJ	0.0	0.0	Anomaly
12	0.0	0.0	0.0	S0	0.0	0.0	Anomaly
13	0.0	0.12	0.03	SF	287.0	2251.0	Normal
14	0.0	1.0	0.2	SF	334.0	0.0	Anomaly
15	0.0	0.0	0.0	S0	0.0	0.0	Anomaly
16	0.0	0.0	0.0	S0	0.0	0.0	Anomaly
17	0.0	0.01	0.02	SF	300.0	13788.0	Normal
18	0.0	1.0	1.0	SF	18.0	0.0	Anomaly
19	0.0	0.02	0.03	SF	233.0	616.0	Normal

$$\text{Gain}(A) = \text{Entropy}(D) - \text{Entropy}_A(D).$$

يعرض الجدول عينة عشوائية من سجلات البيانات D اخذت من مجموعة البيانات الشبكية -nls  
.KDD99

جدول (١) عينة عشوائية لمجموعة سجلات nls\_KDD99 data set

لحساب Gain(flag) اعتماداً على الخاصية class يتم حساب المعادلات الآتية :  
١- نحسب Entropy (D) كما يلي :

$$\text{Entropy}(D) = -((7/19) \log_2(7/19) + (12/19) \log_2(12/19)) = 0.949 \text{ bits}$$

٢- نحسب EntropyA(D) لكل خاصية في D فمثلاً A= flag

S0=8	SF=9	REJ=2	Class
0	7	0	Normal
8	2	2	Anomaly

$$\text{Entropy}_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times \text{Entropy}(D_j)$$

$$\begin{aligned} &= - (8/19)((0/8) \log_2(0/8) - (8/8) \log_2(8/8)) + (-9/19)((7/9) \log_2(7/9) - ((2/9) \log_2(2/9))) + (-2/19) \\ &((0/2) \log_2(0/2) - (2/2) \log_2(2/2))) \\ &= - (8/19)(\log_2 1) + (-9/19)((7/9) \log_2(7/9) - (2/9) \log_2(2/9)) + (-2/19)(-\log_2 1) \\ &= 0 + 0.473(0.77 \log_2 0.77) - (0.2 \log_2 0.2) + 0 \\ &= 0.178 + 0.464 \\ &= 0.642 \text{ bits} \end{aligned}$$

٣- نحسب الـ information gain للخاصية flag :

$$\text{Gain}(flag) = \text{Entropy}(D) - \text{Entropy}_{flag}(D) = 0.949 - 0.642 = 0.307 \text{ bits}$$

التقسيم بواسطة معيار information gain يعتمد على اختبار عدد المخرجات أي أن هذا المعيار يفضل اختيار الـ attribute ذات عدد أكبر من القيم. فمثلاً الصفات ذات القيم الوحيدة unique attribute (مثلاً product\_id) سوف تكون محل اهتمام لهذا المعيار وبالتالي فإن التقسيمات الناتجة عنه سوف تكون كبيرة جداً ووحيدة pure (بعدد قيم الـ unique attribute) ولذلك فإن  $\text{infounique\_attribute}(D) = 0$  ، وبالتالي فإن هذا المعيار غير مناسب لعملية التصنيف لأنه ينتج عدداً كبيراً جداً من التقسيمات.

#### ٤- ٢- ١- ٣- ٢- نسبة الكسب Gain Ratio

نسبة الكسب امتداداً للكسب المعلوماتي يتخلص من نقاط ضعفه ويطبق نوعاً من التطبيع على الكسب المعلوماتي باستخدام معلومات التقسيم "split information" لقيمة معرفة بشكل قياسي كما في المعادلة الآتية [١] :

$$SplitInfo_A(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2 \left( \frac{|D_j|}{|D|} \right).$$

هذا المعيار يأخذ في الاعتبار حساب عدد المخرجات الناتجة لاختبار الصفة على سبيل المثال، في خوارزمية شجرة القرارات C4.5 يتم استخدام معيار تقسيم Gain Ratio لتحديد الانقسام الأفضل ويعرف هذا المعيار بالمعادلة الآتية: [4] :

$$GainRatio(A) = \frac{Gain(A)}{SplitInfo(A)}.$$

الصفة التي يكون عندها نسبة الكسب أعلى قيمة يتم اختيارها كأفضل صفة تقسيم للمجموعة D [1].  
ولتوضيح ذلك نسبة الكسب للصفة flag من الجدول (1)

$$\begin{aligned} Splitinfo_{flag}(D) &= -(8/19) * \log_2(8/19) - (2/19) * \log_2(2/19) - (9/19) * \log_2(9/19) \\ &= 0.525 + 0.342 + 0.510 \\ splitinfo(flag) &= 1.377 \text{ bits} \\ Gainratio(flag) &= gain(flag) / splitinfo(flag) \\ &= 0.307 / 1.377 \\ &= 0.2229 \text{ bits} . \end{aligned}$$

#### ٤- ٢- ١- ٣- مؤشر جيني Gini Index

معيار مؤشر جيني يحاول تجنب مشكلة الكسب المعلوماتي Information Gain بتقييد شروط الاختبار إلى تقسيمات ثنائية فقط بمعنى يتم حساب مجموع الشوائب الموزونة لكل قسم ، ويستخدم هذا المعيار مع خوارزمية CART .

Gini index يقيس شوائب D و مجموعة سجلات التدريب أو أقسام البيانات حسب المعادلة التالية :

$$Gini(D) = 1 - \sum_{i=1}^m P_i^2 .$$

حيث  $P_i = |C_{i,D}| / |D|$

إذا الصفة A. تقسم قاعدة البيانات D إلى  $D_1$  و  $D_2$  فإن الـ gini index يقسم D إلى التقسيم التالي :

$$Gini_A(D) = \frac{|D_1|}{|D|} Gini(D_1) + \frac{|D_2|}{|D|} Gini(D_2).$$

D1 تمثل مجموعة من السجلات في D تحقق الشرط  $A \leq \text{split\_point}$

D2 تمثل مجموعة من السجلات في D تحقق الشرط  $A > \text{split\_point}$

تخفيض الشوائب المتوقع من الانقسام الثنائي على صفة A منفصلة أو مستمرة القيمة تحسب وفق المعادلة التالية :

$$\Delta Gini(A) = Gini(D) - Gini_A(D).$$

ولتوضيح ذلك لنفترض مجموعة التدریب D كما في الجدول (1) فإن عدد السجلات عندما ال class = normal تساوي 7 وعندما ال class = anomaly تساوي 12 ، لإنشاء العقدة N التي تمثل جذر الشجرة لسجلات بيانات التدریب D نحسب الآتي :

$$\begin{aligned} \text{Gain}(D) &= 1 - (7/19)^2 - (12/19)^2 \\ &= 1 - 0.1357 - 0.39889 = 0.46541 \text{ bits} . \end{aligned}$$

لإيجاد التقسيم المناسب لـ D باستخدام الخاصية flag يتم تقسيم مجموعة التدریب D إلى مجموعتين D1 و D2 كما يلي:

D1 = {S0, SF} و D2 = {REJ} عدد السجلات في D1 = 17 و عدد السجلات في D2 = 2

$$\begin{aligned} \text{Gini}_{\text{tag}}(D) &= (17/19)\text{Gini}(D1) + (2/19)\text{Gini}(D2) \\ &= (17/19)(1 - (7/17)^2 - (9/17)^2) + (2/19)(1 - (2/2)^2 - (0/2)^2) \\ &= 0.894(1 - 0.169 - 0.280) + 0.105(1 - 1 - 0) \\ &= 0.527 \text{ bits} . \end{aligned}$$

الصفة التي لها قيمة Gini index أقل من قيمة الصفات الأخرى يتم إختيارها كصفة تقسيم لمجموعة بيانات التدریب D [1].

٤- استخدام تقنية التصنيف لكشف التطفل في قاعدة بيانات NLS-KDD99 data set

استخدم العديد من الباحثين بيانات KDD Cup 1999 لبناء نظم كشف التطفل وأظهرت الدراسات السابقة وجود بعض المشاكل الكامنة في هذه البيانات، إن التحديد المهم لهذه البيانات هو العدد الهائل للسجلات الزائدة بمعنى أن 78% من سجلات التدریب و 75% من سجلات الاختبار متكررة، هذه البيانات تعاني من بعض المشاكل وقد لا تكون مثلى للشبكات الفعلية الموجودة، ويمكن أن يكون محاكاة الهجوم ضمن واحد من الأصناف الأربعة (Dos , R2L , U2R , Probe).

مجاميع البيانات المتولدة KDD Train ، KDD Test شملت (125973 , 22544) سجلاً على التوالي. اقترحت بيانات NSL-KDD من قبل (Tavallae et al.) حل مشاكل بيانات KDD المذكورة سابقاً تعتبر NSL-KDD التي تحوي في كل سجل اتصال TCP على ٤١ ميزة مع عنوان يوضح

هل هذا الاتصال هو اتصال اعتيادي أو نوع من أنواع التطفل ، وهناك ٣٨ ميزة رقمية و ٣ ميزة رمزية ، يأتي فوائد NSL-KDD مقارنة بمجموعة بيانات KDD الأصلية .

- لا تشمل سجلات زائدة في مجموعة التدريب ولن تميل المصنفات باتجاه سجلات أكثر حدوثاً.
  - عدد السجلات المختارة من كل مجموعة : مستوى الصعوبة يتناسب عكسياً ونسبة السجلات في مجموعة KDD الأصلية، ونتيجة لذلك نسب تصنيف طرائق تعليم الآلة المتميزة تختلف بمدى واسع، مما يجعل زيادة الفعالية امتلاك تقييم دقيق لتقنيات تعليم مختلفة.
  - عدد السجلات في التدريب ومجاميع الاختبار معقول مما يجعل من المحتمل إجراء تجارب على المجموعة الكاملة دون الحاجة للاختبار العشوائي لنسبة ضئيلة نتيجة لذلك سيكون تقييم نتائج البحوث المختلفة ثابتة ومشابهة.
  - ٤- ٦ تدريب خوارزمية C4.5 على NLS-KDD
  - استخدمت أداة Weka 3.7.9 لتدريب خوارزمية C4.5 على بيانات التدريب NLS-KDD Train لبناء مصنف البيانات وتم اختباره على بيانات الاختبار NLS-KDD Test ثم تم تقييم أداء المصنف من خلال حساب دقة التصنيف ونسبة الخطأ في التصنيف يمكن قياس دقة المصنفات بحساب الآتي:
  - نسبة الضبطية (Precision) : نسبة سجلات الـ class A المصنفة بشكل صحيح إلى السجلات المصنفة في الـ class A .
  - أ - نسبة الاستدعاء Recall : نسبة سجلات الـ class A المصنفة بشكل صحيح إلى عدد السجلات في الـ class A .
  - ب - الإيجابية الكاذبة False positive (FP) أو الإنذار الكاذب false alarm : أي يتم تصنيف السجلات على أنها attacks وهي في الحقيقة normal
  - ت - السلبية الكاذبة False Negative (FN) : تصنيف السجلات على أنها طبيعية وهي في الحقيقة مهددات . وهذه الهجمات هي الهدف من عملية الكشف.
  - ث - الإيجابية الصادقة True Positive (TP) : تصنيف السجلات على أنها normal وهي في الحقيقة normal .
  - ج - السلبية الصادقة False Negative (FN) : تصنيف السجلات على أنها attacks وهي في الحقيقة attacks .
- يتم حساب دقة المصنف اعتماداً على معدل الكشف Detected Rate و معدل الإنذار الكاذب False Alarm Rate .

استخدمت هذه الورقة مجموعة عشوائية من بيانات NLS-KDD وكانت نتائج التدريب كما يأتي :

Subset: CDDTrain-20Percent

34: dst_host_srv_rate	35: dst_host_diff_srv_rate	36: dst_host_same_srv_port_rate	37: dst_host_diff_host_rate	38: dst_host_serror_rate	39: dst_host_srv_serror_rate	40: dst_host_percent_rate	41: dst_host_srv_error_rate	42: class
Numeric	Numeric	Numeric	Numeric	Numeric	Numeric	Numeric	Numeric	Nominal
25.0	0.17	0.03	0.17	0.0	0.0	0.0	0.05	0.0(normal)
1.0	0.0	0.6	0.88	0.0	0.0	0.0	0.0	0.0(normal)
25.0	0.1	0.05	0.0	0.0	1.0	1.0	0.0	0.0(anomaly)
255.0	1.0	0.0	0.0	0.04	0.02	0.02	0.0	0.0(normal)
255.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0(normal)
48.0	0.07	0.07	0.0	0.0	0.0	0.0	1.0	1.0(anomaly)
9.0	0.04	0.05	0.0	0.0	1.0	1.0	0.0	0.0(anomaly)
15.0	0.06	0.07	0.0	0.0	1.0	1.0	0.0	0.0(anomaly)
23.0	0.09	0.05	0.0	0.0	1.0	1.0	0.0	0.0(anomaly)
13.0	0.05	0.06	0.0	0.0	1.0	1.0	0.0	0.0(anomaly)
12.0	0.05	0.07	0.0	0.0	0.0	0.0	1.0	1.0(anomaly)
13.0	0.05	0.07	0.0	0.0	1.0	1.0	0.0	0.0(anomaly)
218.0	1.0	0.0	0.12	0.03	0.0	0.0	0.0	0.0(normal)
28.0	1.0	0.0	1.0	0.2	0.0	0.0	0.0	0.0(anomaly)
1.0	0.0	0.07	0.0	0.0	1.0	1.0	0.0	0.0(anomaly)
2.0	0.01	0.06	0.0	0.0	1.0	1.0	0.0	0.0(anomaly)
255.0	1.0	0.0	0.01	0.02	0.0	0.0	0.0	0.0(normal)
48.0	1.0	0.0	1.0	1.0	0.0	0.0	0.0	0.0(anomaly)
255.0	1.0	0.0	0.02	0.03	0.0	0.0	0.02	0.0(normal)
255.0	1.0	0.0	0.01	0.04	0.0	0.0	0.0	0.0(normal)
23.0	0.09	0.05	0.0	0.0	1.0	1.0	0.0	0.0(anomaly)
17.0	0.07	0.06	0.0	0.0	0.99	1.0	0.0	0.0(anomaly)
255.0	1.0	0.0	0.01	0.02	0.0	0.0	0.0	0.0(normal)
1.0	0.0	0.06	1.0	0.0	0.0	0.0	0.0	0.0(normal)
2.0	0.01	0.06	0.0	0.0	1.0	1.0	0.0	0.0(anomaly)
25.0	0.1	0.05	0.0	0.0	0.53	0.0	0.02	0.35(normal)
13.0	0.05	0.07	0.0	0.0	1.0	1.0	0.0	0.0(anomaly)
255.0	1.0	0.0	0.02	0.04	0.0	0.0	0.56	0.5(normal)
255.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0(normal)
255.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0(normal)
53.0	1.0	0.0	1.0	0.51	0.0	0.0	0.0	0.0(anomaly)
58.0	0.23	0.04	0.01	0.0	1.0	1.0	0.0	0.0(anomaly)
255.0	1.0	0.0	0.11	0.01	0.0	0.0	0.0	0.0(normal)
1.0	0.0	0.12	0.38	0.0	0.0	0.0	0.29	1.0(anomaly)
25.0	0.08	0.01	0.0	0.0	0.0	0.0	0.0	0.0(normal)
5.0	0.12	0.05	0.05	0.0	0.0	0.0	0.0	0.0(normal)
255.0	1.0	0.0	0.01	0.0	0.0	0.0	0.0	0.0(normal)

٤ - ٥ - ١- نتائج تدريب الخوارزمية :

جدول (١) البيانات قبل تصنيفها بواسطة المصنف

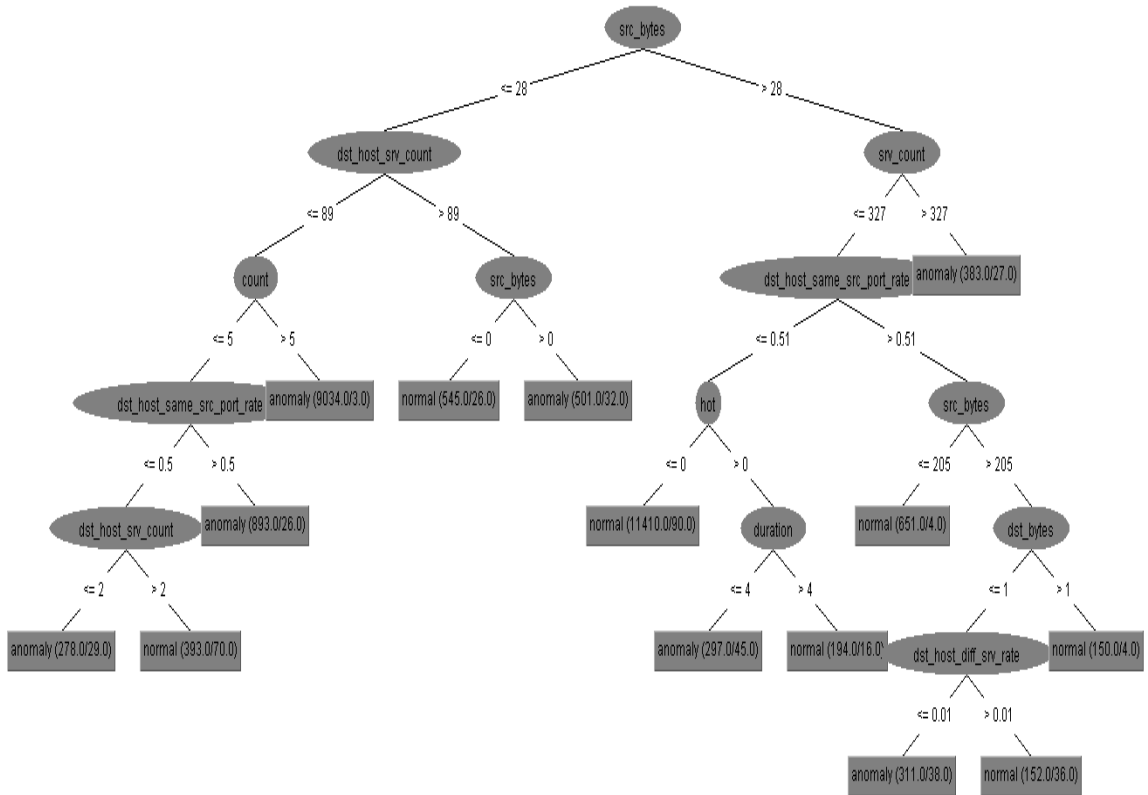
Subset: KDDTrain-20Percent\_predicted

35: dst_host_diff_srv_rate	36: dst_host_same_srv_port_rate	37: dst_host_diff_host_rate	38: dst_host_serror_rate	39: dst_host_srv_serror_rate	40: dst_host_percent_rate	41: dst_host_srv_error_rate	42: prediction margin	43: predicted class	44: class
Numeric	Numeric	Numeric	Numeric	Numeric	Numeric	Numeric	Numeric	Nominal	Nominal
0.06	0.0	0.0	0.0	0.0	0.0	1.0	1.0	1.0(anomaly)	anomaly
0.06	0.0	0.0	0.0	0.0	0.0	1.0	1.0	1.0(anomaly)	anomaly
0.04	0.61	0.02	0.0	0.0	0.0	0.0	0.0	1.0(normal)	normal
0.0	1.0	0.28	0.0	0.0	0.0	0.0	0.0	0.97214(anomaly)	anomaly
0.17	0.03	0.02	0.0	0.0	0.0	0.83	0.71	-0.95209(normal)	anomaly
0.0	0.01	0.03	0.01	0.0	0.0	0.0	0.0	0.99862(normal)	normal
0.72	0.0	0.0	0.0	0.0	0.0	0.72	0.04	0.99862(normal)	normal
0.0	0.0	0.0	0.01	0.01	0.0	0.0	0.02	-0.99862(normal)	anomaly
0.0	0.01	0.02	0.0	0.0	0.0	0.0	0.0	0.99862(normal)	normal
0.08	0.02	0.0	0.0	0.0	0.0	0.0	0.0	-0.99496(normal)	anomaly
0.01	0.0	0.0	0.0	0.0	0.0	0.66	0.32	0.99496(normal)	anomaly
0.03	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.99862(normal)	normal
0.07	0.0	0.0	0.0	0.0	0.95	6.44	0.0	1.0(anomaly)	anomaly
0.05	0.03	0.04	0.0	0.0	0.77	0.0	0.0	-0.98496(normal)	anomaly
0.0	0.01	0.0	0.0	0.0	0.0	0.0	0.0	0.99862(normal)	normal
0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.99862(normal)	normal
0.0	0.03	0.05	0.0	0.0	0.0	0.0	0.0	0.99862(normal)	normal
0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.987673(normal)	normal
0.07	0.0	0.0	0.0	0.0	0.0	1.0	1.0	1.0(anomaly)	anomaly
0.05	0.0	0.0	0.0	0.0	0.0	1.0	1.0	1.0(anomaly)	anomaly
0.01	0.0	0.0	1.0	1.0	0.0	0.0	0.0	1.0(anomaly)	anomaly
0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.99862(normal)	normal
0.0	0.03	0.02	0.0	0.0	0.0	0.0	0.0	0.99862(normal)	normal
0.06	0.0	0.0	0.0	0.0	0.0	1.0	1.0	1.0(anomaly)	anomaly
0.06	0.0	0.0	0.0	0.0	0.0	0.0	1.0	1.0(anomaly)	anomaly
0.0	0.01	0.01	0.0	0.0	0.0	0.0	0.0	0.99862(normal)	normal
0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.99862(normal)	normal
0.0	1.0	0.14	0.0	0.0	0.0	0.0	0.0	1.0(anomaly)	anomaly
0.0	0.83	0.0	0.0	0.0	0.0	0.0	0.0	0.987673(normal)	normal
0.01	0.0	0.0	0.0	0.0	0.0	0.07	0.07	0.0(normal)	anomaly
0.06	0.01	0.0	0.0	0.0	0.0	0.0	0.0	0.99862(normal)	normal
0.0	0.05	0.02	0.05	0.0	0.0	0.0	0.0	0.99862(normal)	normal
0.58	0.99	0.0	0.0	0.0	0.0	0.01	0.0	-1.0(anomaly)	normal

جدول رقم (٢) البيانات بعد تصنيفها بواسطة المصنف

يبين الجدول (2) نتيجة اختبار المصنف C4.5 على بيانات الاختبار حيث يوضح الجدول الصف الحقيقي (class nominal) والصف المتوقع بواسطة المصنف (predicted class nominal) الذي يستخدم كصف معتمد لتصنيف هذه السجلات ، كما يبين الشكل أن السجلات التي صنفت بالفعل على أنها طبيعية normal هي في الأصل طبيعية والسجلات التي صنفت على أنها مهددات هي في الأصل مهددات بالفعل، ولكن توجد نسبة خطأ في تصنيف بعض السجلات التي صنفت على أنها مهددات وهي بالحقيقية سجلات طبيعية وكذلك سجلات صنفت على أنها طبيعية وهي بالحقيقة مهددات .

٤ - ٥ - ٢ مصنف شجرة القرار الناتج :



شكل (٥) مصنف شجرة القرار الناتج لقاعدة بيانات NLS-KDD data set

#### ٤- ٦- تقييم المصنف Evaluation Classifier

مصنوفة التضارب الناتجة لعملية اختبار المصنف على ٢٠% Of NLS-KDD Data Set Test

		Predicted class	
		Class = normal	Class = anomaly
Actual class	Class = normal	١٨٧٩	٢٧٣
	Class = anomaly	٣٩٩٦	٥٧٠٢

شكل رقم (٦) مصنوفة التضارب الناتجة لإختبار المصنف C4.5 على ٢٠% Of NLS-KDD Data Set Test

$$\begin{aligned} \text{Accuracy} &= (1879 + 5702) / (1879 + 273 + 3996 + 5702) \\ &= (7581) / (11850) \\ &= 0.639 = 63.9 \% \end{aligned}$$

$$\begin{aligned} \text{Error rate} &= (3996 + 273) / (1879 + 273 + 3996 + 5702) \\ &= (4269) / (11850) \\ &= 0.360 \\ &= 36\% \end{aligned}$$

$$\begin{aligned} \text{true positive rate (TPR) or sensitivity} &= TP / (TP + FN) \\ &= 1879 / (1879 + 273) = 87 \% \end{aligned}$$

$$\begin{aligned} \text{true negative rate (TNR) or specificity} &= TN / (TN + FP). \\ &= 5702 / (5702 + 3996) = 59\% \end{aligned}$$

$$\begin{aligned} \text{false positive rate (FPR)} &= FP / (TN + FP), \\ &= 3996 / (5702 + 3996) = 41\% \end{aligned}$$

$$\begin{aligned} \text{false negative rate (FNR)} &= FN / (TP + FN). \\ &= 273 / (1879 + 273) = 13\% \end{aligned}$$

$$\begin{aligned} \text{Precision} &= tp / (tp + fp) \\ &= 1879 / (1879 + 3996) = 32\% \end{aligned}$$

$$\begin{aligned} \text{Recall} &= tp / (tp + fn) \\ &= 1879 / (1879 + 273) = 87\% \end{aligned}$$

تلخص نتائج التجربة 1 في الجدول التالي :

TPR	FP	TP	FN	Precisoin	Recall	Errore rate	Accurcy
87 %	41%	58%	13%	32%	87%	36%	63.9 %

شكل رقم (٧) نتائج اختبار المصنف C4.5 على بيانات 20% Of NLS-KDD Data Set Test

مصنوفة التضارب الناتجة لعملية اختبار المصنف على Full NLS- KDD 99 Data Set Test

		Predicted class	
		Class = normal	Class = anomaly
Actual class	Class = normal	٩٤٤٨	٢٦٣
	Class = anomaly	٣٩٠٠	٨٩٣٣

شكل رقم (٦) مصنوفة التضارب الناتجة لإختبار المصنف C4.5 على Full NLS-KDD 99 Data

Set Test

نتائج التجربة 2 :

TP	FP	TN	FN	Precisoin	Recall	Errore rate	Accurcy
97 %	59%	69.6%	2.7%	70.8%	97%	36%	٨١,٥٣٣٩%

شكل رقم (٩) نتائج اختبار المصنف C4.5 على ب على Full NLS-KDD 99 Data Set Test

٥ - الخلاصة CONCLUSION :

إذا أردنا استخدام الشبكة فإن هناك حاجة ملحة لكشف هجوم أنظمة الشبكات، وفي هذه الورقة استخدمت واحدة من أهم تقنيات استخراج البيانات وهي C4.5 للكشف عن الشذوذ في الشبكة. نتيجة التجربة المعروضة سابقاً أظهرت خوارزمية C4.5 نتيجة فعالة في كل من كشف الشذوذ ومعدل الانذار الكاذب في مجموعة البيانات المتوفرة، وأظهرت النتائج أنه كلما كانت كمية البيانات كبيرة تكون نسبة الخطأ أقل ودقة تنبؤ المصنف عالية حيث كانت دقة المصنف % 63.9 عندما كانت كمية بيانات الاختبار (946,848 byte) وعندما زادت كمية بيانات إلى ( 497,700 byte ) ارتفعت دقة المصنف إلى %٨١,٥٣٣٩، بينما كانت كمية بيانات التدريب ثابتة، هذه النتائج تدل قدرة المصنف العالية للتعلم كلما زاد حجم كمية البيانات.

٦- مراجع:

1. Data Mining: Concepts and Techniques Third Edition - Jiawei Han - University of Illinois at Urbana – Champaign - Micheline Kamber Jian Pei Simon Fraser University- Morgan Kaufmann is an imprint of Elsevier.-
2. Data Mining for Network Intrusion Detection Paul Dokas, Levent Ertöz, Vipin Kumar, Aleksandar Lazarevic, Jaideep Srivastava, Pang-Nig Tan ,Computer Science Department, 200 Union Street SE, 4-192, EE/CSC Building ,University of Minnesota, Minneapolis, MN 55455, USA ,aleks@cs.umn.edu ,srivasta@cs.umn.edu ,[kumar@cs.umn.edu](mailto:kumar@cs.umn.edu) .
3. Zirkle, L., “What is host-based intrusion detection? “Virginia Tech CNS. SANS -Institute Resources, Intrusion Detection FAQ, Hyperlink ID FAQ, 2000.
4. ISSN (Print): 2279-0047 ISSN (Online): 2279-0055 - International Journal of Emerging -Technologies in Computational and Applied Sciences (IJETCAS) [www.iasir.net](http://www.iasir.net), Differentiating Network Attacks using C4.5 Algorithm with Multiboosting, V.Balaji1, Varalakshmi.K2 # Department of Computer Science and Engineering, SRM University, SRM Nagar, Kattankulathur – 603203, Kancheepuram District, Tamil Nadu, INDIA.
5. Introduction to Data Mining- by Tan, Steinbach, Kumar – 2004



# جامعة الناصر

## AL-NASSER UNIVERSITY