



بسم الله الرحمن الرحيم

جامعة السودان للعلوم والتكنولوجيا

كلية علوم الحاسوب وتقانة المعلومات

مقارنة خوارزميات التعلم الآلي علي بيانات حوادث المرور

Comparing The Machine Learning Algorithms By Using Traffic Accident Dataset

بحث تكميلي مقدم كمطلوب لنيل درجة بكالوريوس الشرف في علوم الحاسوب

أكتوبر 2017 م

بسم الله الرحمن الرحيم

جامعة السودان للعلوم والتكنولوجيا

كلية علوم الحاسوب وتقانة المعلومات

مقارنة خوارزميات التعلم الآلي علي بيانات حوادث

المرور

Comparing the Machine Learning
Algorithms By Using Traffic Accident
Dataset

أكتوبر 2017 م

بحث تكميلي مقدم كمطلوب لنيل درجة بكالوريوس الشرف في علوم
الحاسوب

إشراف: أ.علي الأمين

إعداد الطلاب

عبد الله عمر إبراهيم

محمد معاوية محمد

آية

قال تعالى : (فاذكروني أذكركم واشكروا لي ولا تكفرون) البقرة (152)

صدق الله العظيم

الحمد لله

الحمد لله الذي أنزل على عبده الكتاب . أظهر الحق بالحق وأخزى الأحزاب . وأتمّ نوره، وجعل كيد الكافرين في تباب . أرسل الرياح بشرى بين يدي رحمته وأجرى بفضل السحاب . وأنزل من السماء ماء، فمناه شجر، ومنه شراب . جعل الليل والنهار خلفه فتذكر أولو الألباب . نحمده تبارك وتعالى على المسببات والأسباب . وصلي اللهم وسلم وبارك علي جناب الحبيب المحبوب حبيبك المصطفى وعلي اله وصحبه وسلم سلاماً تاماً اللهم علمنا ما جهلنا وانفعنا بما علمتنا وزدنا علماً .

الإهداء

إلي الوالد الحبيب

- إلي قدوتي الأولي ونبراسي الذي ينير دربي.
- إلي من أعطاني ولم يزل يعطيني بالحدود.
- إلي من رفعت رأسي عليا افتخارا به.
- ولاكني لا املك إلي أن ادعوا الله عز وجل أن يبقيه ذخرا لنا.
- أطال الله عمرك فيما يحب ويرضي.

إلي الوالدة العزيزة

- فيا من علمتني أبجدية الحروف.
- ويا من علمتني الصمود مهما تبدلت الظروف.
- اخط لكي كلمات مدها حبر دمي.
- كلمات ملؤها شكر و عرفان
- كلمات تتردد علي كل لسان
- نعم إنها أمي الغالية..

الشكر والعرفان

لابد لنا ونحن نخطو خطواتنا الاخيريه في الحياة الجامعية من وقفه ونعود إلي أعوام قضيناها في رحاب الجامعة مع أساتذتنا الكرام الذين قدموا لنا الكثير باذلين بذلك جهودا كبيره في بناء الغد لتبعث الأمة من جديد . وقيل إن نمضي نقدم اسمي آيات الشكر والامتنان والتقدير والمحبة إلي الذين حملوا أقدس رسالة في الحياة.

"كن عالما .. فإن لم تستطع فكن متعلما، فإن لم تستطع فأحب العلماء ، فإن لم تستطع فلا تبغضهم"

إلي الذين مهدوا لنا طريق العلم والمعرفة...

إلي جميع أساتذتنا الأفاضل...

وأخص بالتقدير والشكر:

الأستاذ :علي الأمين.

وكذلك نشكر كل من ساعد علي إتمام هذا البحث وقدم لنا العون ومد لنا يد المساعدة وزودنا بالمعلومات اللازمة لإتمام هذا البحث.

المستخلص

إن استخدام تقنية التنقيب في البيانات يوفر للمؤسسات القدرة علي فهم جميع البيانات و التركيز علي أهم المعلومات الموجودة في قواعد البيانات, وترتكز تقنيات التنقيب علي البيانات علي بناء التنبؤات المستقبلية واكتشاف المعرفة و السلوك والاتجاهات ، مما يسمح باتخاذ القرارات الصحيحة في الوقت المناسب.

نموذج تحليل بيانات حوادث المرور هو نموذج يقوم علي تحليل بيانات الحوادث معتمدا علي بعض العوامل .

تعتمد فكرة المشروع علي تنقيب البيانات حيث يهدف المشروع إلي تحليل مجموعة من بيانات حوادث المرور لولاية الخرطوم وتدريب بعض خوارزميات تعلم الالة ومقارنة نتائج تلك الخوارزميات مع بعضها ، حيث استخدم في هذا البحث خوارزميات تعلم الآلة في تحليل البيانات.

Abstract

The use of data mining technology provides access to information and data in databases. Data mining techniques are based on future predictions, knowledge discovery and behavior, allowing for the right decisions to be made in a timely fashion.

Traffic Accident Data Analysis Model is a model based on the analysis of accident data based on certain factors.

The idea of the project is based on data mining. The project aims at analyzing a set of traffic accident data for the state of Khartoum and the training of some algorithms of learning algorithms and comparing the results of these algorithms with each other.

فهرس المصطلحات

المصطلح	شرح المصطلح
FP	False positive
TN	True negative
FN	False negative
TP	True positive

فهرس الجداول

الموضوع	رقم الجدول
جدول مقارنات الدراسات السابقة	الجدول رقم (2-2)
جدول الاقسام	جدول رقم (3-1-1)
جدول الايام	جدول رقم (3-1-2)
جدول الشهور	جدول رقم (3-1-3)
جدول الاعمار	جدول رقم (3-1-4)
جدول الرخص	جدول رقم (3-1-5)
جدول اعمار المصابين	جدول رقم (3-1-6)
مقارنة بين الخوارزميات	الجدول رقم (3-2-3)
تقييم الخوارزميات	جدول رقم (3-3-1)

فهرس الاشكال

الموضوع	رقم الشكل
أنواع ومهام تقنية تنقيب البيانات	شكل(2-1)
تقنية التصنيف	الشكل(3-2-1)
مصفوفة التضارب	الشكل رقم(2-2-3)

فهرس الصور

الموضوع	رقم الصورة
بيانات الحوادث بعد التعديل	صورة رقم (3-2-8)
بيانات حوادث المرور	صورة رقم (3-1-7)
نتيجة تدريب خوارزمية support vector machine	الصورة رقم (3,3,1)
نتيجة تدريب خوارزمية naïve bayes	الصورة رقم (3,3,2)
نتيجة تدريب خوارزمية 48 j	الصورة رقم (3,3,3)

فهرس المحتويات

iii	آية	
iv	الحمد لله	
v	الإهداء	
vi	الشكر والعرفان	
vii	المستخلص	
viii	Abstract	
ix	فهرس المصطلحات	
x	فهرس الجداول	
xi	فهرس الأشكال	
xii	فهرس الصور	
3	الباب الاول	
3	المقدمة	
1	المقدمة	1.1
1	مشكلة البحث	1.2
2	فرضيات البحث	1.3
2	أهمية البحث	1.4
2	أهداف البحث	1.5
2	منهجية البحث	1.6
2	حدود البحث	1.7
3	هيكلية البحث	1.8
4	الباب الثاني	2
5	الفصل الاول	
5	الإطار النظري	2.1
5	المقدمة	2.1.1
5	مفهوم التنقيب في البيانات	2.1.2
6	أهمية أسلوب تنقيب البيانات	2.1.3
6	أهداف أسلوب تنقيب البيانات	2.1.4
7	نماذج التنقيب في البيانات	2.1.5
7	أدوات التنقيب في البيانات	2.1.6
8	مراحل عملية التنقيب في البيانات	2.1.7
10	تعلم الآلة	2.1.8

10.....	أنواع التعلم.....	2.1.9
11.....	الدراسات السابقة.....	2.2
11.....	المقدمة.....	2.2.1
13.....	الباب الثالث.....	3
13.....	منهجية البحث.....	
14.....	الفصل الأول.....	
14.....	منهجية البحث.....	3.1
14.....	المقدمة.....	3.1.1
19.....	القسم الثاني.....	
19.....	التصنيف.....	3.2
19.....	المقدمة.....	3.2.1
19.....	التصنيف.....	3.2.2
19.....	تقنية التصنيف.....	3.2.3
20.....	تقنيات التصنيف.....	3.2.4
22.....	الفصل الثالث.....	
22.....	التدريب والتقييم.....	3.3
22.....	مقدمة.....	3.3.1
22.....	تدريب الخوارزميات علي البيانات.....	3.3.2
25.....	تقييم خوارزميات التصنيف.....	3.3.3
26.....	مقارنة الخوارزميات في التقييم.....	3.3.4
27.....	الباب الرابع.....	4
28.....	النتائج.....	4.1
28.....	التوصيات.....	4.2
28.....	الخاتمة.....	4.3
29.....	المصادر والمراجع.....	

الباب الأول

المقدمة

1.1 المقدمة

تمثل مشكلة الحوادث المرورية من المشكلات الاجتماعية والاقتصادية والصحة العامة المرتبطة بالتنمية . حيث تسبب الحوادث في مقتل أكثر من مليون شخص ، كما يصاب أكثر من 15 مليون شخص بجروح من جراء تلك الحوادث علي الطرق كل عام ، تتحمل الدول النامية والدول ذات الاقتصاديات المحدودة العبء الأكبر حيث تمثل الحوادث احد قضايا التنمية التي تؤثر تأثيرا غير مناسباً علي الفقراء في الدول المنخفضة الدخل والمتوسطة الدخل . وتستنزف الحوادث المرورية عادة من 1 إلى 3 في المائة من إجمالي الناتج المحلي لأي دول.

وينجم عن الحوادث المرورية وفاة (1200000) شخص سنويا كما يصاب 50 مليون شخصا بالإعاقة بسبب حوادث المرور ، تعد حوادث المرور السبب الرئيسي الثاني للوفيات للفئة العمرية من 5 إلى 29 سنة ، كما تعد السبب الرئيسي الثالث في الوفيات في الفئة العمرية ما بين 30 و44 سنة ، وتقدر منظمة الصحة العالمية إن عدد الوفيات سيزداد بنسبة 80% في الدول النامية وذات الدخل المتدنية بحلول 2020 م إذا لم تتخذ إجراءات فورية للتصدي لهذه الحوادث وأسبابها سيفقد العالم يوميا أكثر من (3000) شخص من جراء حوادث المرور أما علي صعيد إقليم الشرق الأوسط في المنظمة الصحة العالمية والذي يشمل معظم الدول العربية ، فإنه يتوفي أكثر من (130) ألف شخص سنويا. [1]

أصبحت الحوادث المرورية تمثل وبشكل كبير هاجسا وقلقا لكافة أفراد المجتمع ، وأصبحت واحدة من أهم المشكلات التي تستنزف الموارد المادية والطاقات البشرية وتستهدف المجتمعات في أهم مقومات الحياة والذي هو العنصر البشري باضا فه إلي ما تكبده من مشاكل اجتماعية ونفسية وخسائر مادية ، مما أصبح لزاما إيجاد حلول ومقترحات ووضعها موضع التنفيذ للحد من هذه الحوادث أو علي اقل تقدير معالجة أسبابها والتخفيف من أثارها السلبية[2].

1.2 مشكلة البحث

بالرغم من تطور التقنيات في تحليل البيانات إلا أننا نجد انه لا توجد أي دراسة علمية منهجية تستخدم للحد من هذه الحوادث ووضع استراتيجيات مناسبة للسلامة المرورية ،ونجد أن معظم الطرق المستخدمة في تحليل بيانات حوادث المرور إحصائية.

1.3 فرضيات البحث

كلما زادت كمية البيانات المستخدمة في التحليل كلما كانت نسبة الخطأ اقل ودقة تنبؤ المصنف عالية.

1.4 أهمية البحث

تعد هذه الدراسة من الدراسات القليلة التي تنبه لموضوع أهمية تحليل المعطيات المرورية وتمثل أهمية البحث في معرفة أسباب الحوادث وتقليل الحوادث. لذلك لا بد من الاستفادة القصوى من بيانات حوادث المرور المأخوذة من الإدارة العامة للمرور.

1.5 أهداف البحث

يهدف هذا البحث إلي:-

- تحليل بيانات حوادث المرور بطريقة منهجية علمية.
- جعل الالة قادرة علي التعلم.
- تحليل عميق للبيانات لكي تساعد في حل المشكلة.

1.6 منهجية البحث

- يتم في هذا البحث استخدام المنهج الوصفي والاعتماد علي بيانات حوادث المرور المقدمة من الإدارة العامة للمرور.
- كما تتم الإشارة إلي الأسلوب الاستدلالي في تحليل معطيات المرور من اجل تحليل بيانات الحوادث وما يرتبط بها من عوامل مثل رخصة السائق وحالة السائق والفئة العمرية وزمان الحادث ومكان الحادث ونوعية المركبة.
- استخدام تقنية machine learning في التحليل.

1.7 حدود البحث

بيانات حوادث المرور لولاية الخرطوم من يناير إلي ابريل 2017 م.

1.8 هيكلية البحث

يتضمن هذا البحث بالاضافه إلى هذا الفصل الفصول.

الباب الثاني : يتضمن نبذه عامه عن الإطار النظري ل مفهوم تنقيب البيانات وأساليبها ,
والدراسات السابقة.

الباب الثالث : يتضمن المنهجية المستخدمة في البحث.

الباب الثاني

الإطار النظري والدراسات

السابقة

الفصل الاول

1.9 الإطار النظري

1.9.1 المقدمة

يتميز عصرنا الراهن عصر الانترنت والاقتصاد الرقمي بوجود كميات كبيرة للبيانات حتى أصبح من المستحيل علي المحللين استخلاص معلومات ذات معني باللجوء فقط إلى المداخل التقليدية للتحليل التمهيدي للبيانات.

مع وجود كميات كبيرة من البيانات المخزنة في قواعد البيانات ومخازن البيانات إزدادت الحاجة إلي تطوير أدوات تمتاز بالقوة لتحليل البيانات واستخراج المعلومات والمعارف منها ، من هنا ظهر مايسمى بالتنقيب في البيانات كتقنية تهدف إلي استخراج المعرفة من كميات هائلة من البيانات.

وتعتبر تقنية التنقيب في البيانات من التقنيات الحديثة التي فرضت نفسها في بقوة في عصر المعلوماتية، واستخدامها يوفر للشركات والمنظمات والمؤسسات الحكومية في جميع المجالات القدرة علي اكتشاف المعلومات الموجودة في قواعد البيانات، كما تركز تقنية التنقيب في البيانات علي بناء التنبؤات المستقبلية مما يسمح باستخدام القرارات المناسبة في الوقت المناسب.

وتعتبر مرحلة استكشاف المعرفة من البيانات من أكثر المراحل تعقيدا. والمرتبطة إلي حد بعيد بعملية تطوير أخرى مهمة جدا هي مستودعات البيانات، حيث إن الكثير من الشركات والمنظمات تستخدم عملية التنقيب في قواعد البيانات بشكل منهجي ومنظم بوصفها تشكل جوهر العمل الذي يعتمد عليه في تفعيل النشاط وتحقيق الميزة التنافسية.

1.9.2 مفهوم التنقيب في البيانات

ظهر مصطلح التنقيب في البيانات في منتصف التسعينات في الولايات المتحدة الأمريكية، كنوع يجمع ما بين الإحصاء وتكنولوجيات (قواعد البيانات، الذكاء الاصطناعي، التعلم الآلي) .

توجد عدة تعريفات لهذا المفهوم حيث يمكن تعريفها بأنها " الاستكشاف الآلي والمؤتمن لأنماط شائعة وغير جلية مخفية في قاعدة بيانات معينة وأيضا تعرف بأنها،[3]" إجراءات تحليل دقيقة وذكية، تفاعلية و تسلسلية، تسمح لمسيري النشاطات عند استخدام هذه الإجراءات باتخاذ قرارات والقيام بأعمال ملائمة في صالح النشاط المسؤولين عنة والمؤسسة التي يعملون فيها [4]" وإنما أيضا " عبارة عن تحليل لمجموعات كبيرة الحجم من البيانات المشاهدة للبحث عن علاقات محتملة وتلخيص للبيانات في أشكال جديدة لتكون مفيدة ومفهومه لمستخدمها. [5]" من خلال التعريفات السابقة يمكن القول أن

التنقيب في قواعد البيانات يهدف إلى استخلاص المعلومات المخبأة فيها، واستخدامه يوفر للمؤسسات في جميع المجالات علي استكشاف والتركيز علي أهم المعلومات في قواعد البيانات بالإضافة إلى كثرة البيانات الموجودة والمخزنة في ما يسمى بقواعد البيانات أصبحت تقنية التنقيب في البيانات محل تساؤل من عديد من الباحثين للاستفادة منها ، ومع زيادة انتشار مستودعات التخزين الضخمة ما يسمى (*data warehouses*)، أصبح من الضروري إيجاد تقنيات وطرق ووسائل لاستخلاص المعلومات والمعرفة من هذه البيانات المكثفة واستغلالها في حل المشاكل واتخاذ القرارات ، باستخدام تطبيقات الحاسوب الحديثة والتي تعتبر تكنولوجيا حديثة ذكية قائمة علي جعل الحاسوب " يفكر كما يفكر الإنسان ويفعل كما يفعل الإنسان " جاءت فكرة الكشف والتنقيب علي هذه البيانات بطرق ذكية للمساعدة في حل المشاكل واتخاذ القرارات وتعتبر خطوة من خطوات استكشاف المعرفة من قواعد البيانات.

1.9.3 أهمية أسلوب تنقيب البيانات

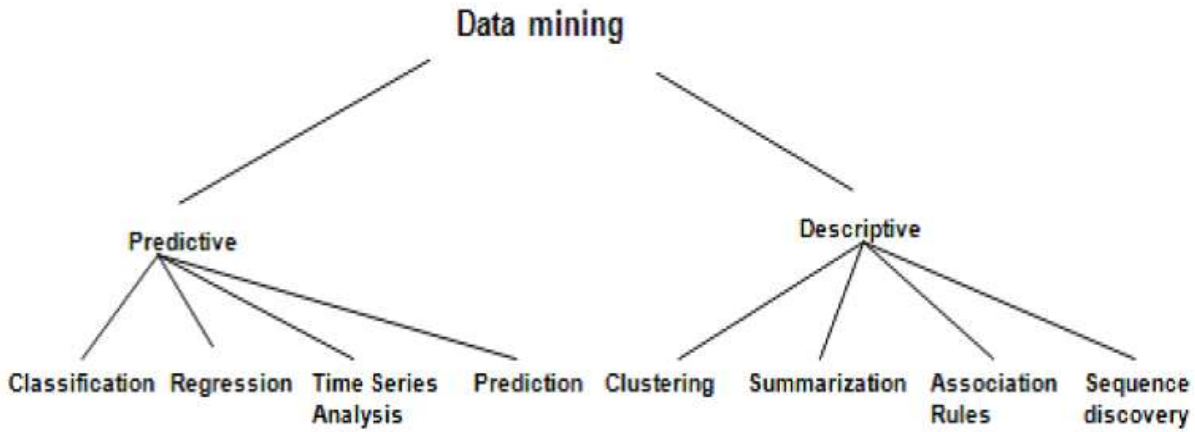
- عملية تحليلية للاستكشاف والبحث في بيانات ضخمة وهائلة لاستخراج أنماط مفيدة وإيجاد العلاقات ومدى الارتباط بين عناصرها.
- من أجل تحليل البيانات للحصول علي علاقات جديدة غير متوقعة.
- نقل عالم الأنظمة الواقعي إلى عالم افتراضي يمارس فيه متخذي القرار التحليل واختبار الفرضيات علي شاشات الحاسوب ذات القدرة الرسمية العالية الدقة إلى أن يصلوا إلى ما يطمحون إليه من فهم وقناعات قبل اتخاذ القرارات بشأن ما يدرسون من أنظمة.
- التنبؤ ومن ثم استنتاج إجابات مقدرة تقديرا إحصائيا لكميات البيانات.

1.9.4 أهداف أسلوب تنقيب البيانات

- إن التنقيب في قواعد البيانات يهدف إلى انتزاع واستخلاص أنماط مفيدة، وهي تكنولوجيا حديثة، أصبحت مهمة في ظل التطور السريع وانتشار استخدام قواعد البيانات.
- استخدامها يوفر للمؤسسات وأجهزة الأمن في جميع المجالات القدرة علي استكشاف، والتركيز علي أهم المعلومات في قواعد البيانات.
- تركز تقنيات التنقيب علي بناء التنبؤات المستقبلية واستكشاف السلوك والاتجاهات، مما يسمح بتقدير القرارات الصحيحة واتخاذها في الوقت المناسب.
- تجيب تقنيات التنقيب علي العديد من الأسئلة ، وفي وقت قياسي ، وخاصة تلك النوعية من الأسئلة التي يصعب الإجابة عليها، إن لم يكن مستحيلا ، باستخدام تقنيات الإحصاء الكلاسيكية ، والتي كانت وان وجدت فإنها تستخدم وقتا طويلا والعديد من الإجراءات.

1.9.5 نماذج التنقيب في البيانات

- النماذج التنبؤية (predictive data mining): يحاول إيجاد أفضل التنبؤات اعتمادا علي المعطيات ويستخدم هذا التنقيب علي المعلومات القديمة لتوقع ما سيحدث في المستقبل.
- النماذج الوصفية: (Descriptive data mining) تعتمد علي إعادة تنظيم البيانات والتنقيب في أعماقها لاستخراج المؤشرات الموجودة فيها.



شكل(2-1)

وتنقسم إلي نوعان :

- النماذج التنبؤية
- النماذج الوصفية

1.9.6 أدوات التنقيب في البيانات

- السلاسل الزمنية
- التصنيف
- التنبؤ
- التلخيص
- التجزئة
- تحليل الارتباط
- الكشف عن التغيرات أو الانحرافات

1.9.7 مراحل عملية التنقيب في البيانات

يمكن تلخيص مراحل و خطوات عملية التنقيب في البيانات كما يلي:

- فهم طبيعة الأعمال
- يعتبر المطلب الأول لاكتشاف المعرفة هو فهم المشاكل و المسائل التي تواجهها الأعمال .
و بمعنى آخر ، كيف يمكن تحقيق المنفعة الأعظم من التنقيب في البيانات، مما يتطلب وجود صيغة واضحة ومحددة لأهداف الأعمال.
- فهم البيانات
- تعتبر مسألة معرفة ماهية وطبيعة البيانات عامل مهم في نجاح عملية التنقيب في البيانات و اكتشاف المعرفة حيث أن معرفة البيانات بصورة جيدة تعني مساعدة المصممين على استخدام الخوارزميات أو الأدوات المستخدمة للمسائل المحددة بدقة عالية.
وهذا يقود إلى تعظيم فرص النجاح بالإضافة إلى رفع الفعالية و الكفاءة لنظام اكتشاف المعرفة .ولاتحتاج عملية التنقيب في البيانات إلى تجميع البيانات في مستودع البيانات، أما إذا كاف مستودع البيانات موجود في الأنظمة، فمن الأفضل عدم احتكار المستودع بشكل مباشر لغرض التنقيب في البيانات.
- و يمكن تلخيص الخطوات الضرورية لعملية فهم البيانات كالآتي:
- تجميع البيانات
- و هي الخطوة الموجهة نحو تحديد مصدر البيانات في الدراسة بما في ذلك استخدام البيانات العامة الخارجية مثل) المرور و الضرائب وغيرها.(
- توصيف البيانات
- وهي الخطوة التي تركز على توصيف محتويات الملف الواحد من الملفات أو الجداول.
- جودة البيانات و تحقيقها
- هذه الخطوة تحدد ما إذا كان تقليل أو إهمال بعض البيانات غي الضرورية أو كونها رديئة الجودة وقد لا تصلح في الدراسة .لان النموذج الجيد يحتاج إلي بيانات جيدة مما يتوجب أن تكون البيانات صحيحة وذات مضمون دقيق.
- التحليل الاسترشادي للبيانات

تستخدم الأساليب مثل الإظهار المرئي أو التصور أو عملية التحليل المباشر التي تؤدي إلى إجراء التحليل الأولي للبيانات تعتبر هذه الخطوة مهمة وضرورية لأنها تركز على تطوير الفرضيات المتعلقة بالمشكلة قيد الدراسة.

• تهيئة البيانات

وتشمل الخطوات التالية:

• الاختيار وتعني اختيار المتغيرات المتوقعة و حجم العينة.

• صياغة المتغيرات و تحويلها

حيث يجب دائما أن تصاغ المتغيرات الجديدة لبناء النماذج الفعالة.

• تكامل البيانات حيث أن مجاميع البيانات في دراسة التنقيب عن البيانات من الممكن خزنها في قواعد بيانات متعددة لأغراض التي تكون بحاجة توحيدها في قاعدة بياناتية واحدة.

• تصميم وتنسيق البيانات حيث تتعلق هذه الخطوة في إعادة ترتيب حقول البيانات كما يتطلب في نموذج التنقيب في البيانات.

• صياغة نماذج الحل و ثبوتها

إن بناء و صياغة نموذج الحل السليم و الدقيق يتم من خلال عملية الخطأ و الصواب، حيث كثيرا ما تحتاج مثل هذه العملية إلى مساعدة المتخصصين في التنقيب عن البيانات بهدف اختبار و فحص مختلف البدائل للحصول على أفضل نموذج لحل المشكلة قيد الدراسة.

• التقييم و تحليل نتائج النموذج

حالما يتم صياغة النموذج و التحقق من ثباته ك صدقة ، تجري مباشرة عملية التحقق من ثبات حزمة البيانات التي يتم تغذيتها بواسطة النموذج وبما أن نتائج هذه البيانات معروفة ، لذا فان النتائج المتوقعة تقارن مع النتائج الفعلية في ثبات حزمة البيانات قيد التشغيل وتؤدي هذه المقارنة أو المفاضلة إلى التحقق من دقة النموذج

• نشر و توزيع النموذج

حيث تشتمل هذه الخطوة على نشر و توزيع النموذج داخل المنظمة لمساعدة عملية صنع القرار . وان النموذج الصالح يجب أيضا أن يحقق الرضا لدي المستخدمين طالما أن اختيار النموذج لا بد أن يتم من خلا الدراسة الاسترشادية أو نموذج مصغر من الدراسة الشاملة.

1.9.8 تعلم الآلة

هو أحد فروع الذكاء الاصطناعي التي توفر القدرة على التعلم لأجهزة الحاسوب. تهتم بتصميم وتطوير الخوارزميات والتقنيات التي تسمح للحاسوب بامتلاك خاصية "التعلم". هناك مستويين من التعلم: الاستقرائي والاستنتاجي. يقوم التعلم الاستقرائي باستنتاج قواعد وأحكام عامة من البيانات الضخمة.

يتداخل "علم الآلة" مع علم إحصاء الحوسبة. Computational Statistics ويهتم بصنع التنبؤات من خلال استخدام الحاسب، يرتبط علم التحسين الرياضي Mathematical Optimization ، الذي يركز على اختيار البديل الأفضل من بين العديد من البدائل المتاحة، كما يوفر الكثير من الوسائل والنظريات والتطبيقات لتعلم الآلة [6].

1.9.9 أنواع التعلم

- التعلم المراقب.
- التعلم الغير المراقب.

سينتظر هذا البحث علي نوع من أنواع تعلم الآلة وهو التعلم المراقب الذي يهدف إلي ربط المدخلات والمخرجات وذلك عن طريق إعطاء أمثلة للمدخلات (input) والمخرجات (output) المرغوبة للآلة من قبل المعلم.

من أهم الأدوات المستخدمة في عملية التعلم المراقب هي أداة التصنيف وهو النوع الأكثر استخداما في تعلم الآلة. في هذا النوع يكون الدخل مصنفا إلى نوعين أو أكثر. وهدف عملية التعلم إنتاج نموذج يستطيع تصنيف أي دخل جديد إلى نوع أو أكثر من الأنواع المعرفة سابقاً. مثال على هذا النوع، عملية تصنيف البريد الإلكتروني وعملية التعرف على الوجوه.

الفصل الثاني

1.10 الدراسات السابقة

1.10.1 المقدمة

سنتناول في هذا الباب الدراسات والبحوث التي قدمت في مجال التنقيب علي البيانات في تحليل بيانات المرور، وتتبع كيفية تتحاياها للبيانات والإستفاده منها في كثير من الاستخدامات وفي حل الكثير من المشاكل.

• S. Krishnaveni، (2011)، العمل مع بعض نماذج التصنيف للتنبؤ بالإصابات التي حدثت في حادث سير في نيجيريا ومقارنتها[7]

يستخدم هذا البحث على نهج (neural networks) القائمة في حين أن تحليل (decision trees) البيانات يمكن استخدامها للعمل على الحد من المجزرة على الطرق السريعة. تم تصنيف البيانات في بيانات مستمرة وفئوية حيث تم تحليل البيانات المستمرة باستخدام تقنية (neural networks) الاصطناعية والبيانات الفئوية باستخدام تقنية (decision trees). وأظهرت النتائج أن نهج (decision trees) تفوق على (neural networks) بمعدل خطأ أقل ومعدل دقة أعلى. يعتمد هذا البحث على ثلاثة أسباب مهمة للحوادث بسبب انفجار الإطارات، وفقدان السيطرة والإكثار من السرعة.

• Naina et. Al. (2016) استخدم نموذج التصنيف، مع 3 Dichotomiser التكرارية وخوارزمية (decision trees). كما تقارن الخوارزمية الحالية مع تعزيز خوارزمية C 4.5 مع استخدام أداة Weka أنها أظهرت أن من المفيد عندما نستخدم مجموعات البيانات الكبيرة والنتائج هي مؤثرة جدا[8].

- K. Geetha 2015 هذه الدراسة تعمل على بيانات حوادث المرور من Tamilnadu city. والهدف الرئيسي من هذه الدراسة هو تقليل عدد حوادث الطرق .تتم إدارة بيانات حوادث المرور في شكل نص أو تنسيقات رقمية بطريقة غير مصنفة[9].
- قدم معاوية وأسماء هذه الدراسة في السودان 2015 بعمل تحليل بيانات الحوادث والهدف من الدراسة حوادث الطرق وعلاقتها مع حوادث الشباب واستخدم الشبكة العصبية في البحث[10].
- Miao Chong هذه الدراسة تعمل علي بيانات حوادث المرور في الولايات المتحدة مأخوذة من (NASS) والهدف من البحث جعل لغة الآلة ذكية في كشف نوعية الهدف اعتمادا علي بعض العوام

جدول يوضح المقارنة بين الدراسات السابقة

الدراسة	السنة	الأدوات والتقنيات	نقاط ضعف الدراسة
1	2011	neural networks	لا يوجد
2	2016	Weka وخوارزمية (decision trees)	قلة البيانات
3	2015		لا يوجد
4	2015	الشبكة العصبية	لا يوجد
5	2015	neural networks-support vector machine-	لا يوجد
البحث المقترح	2017	Weka(support vector machine-شجرة القرار-naïve Bayes)	فقدان بعض البيانات

الجدول رقم(2-2)

الباب الثالث

منهجية البحث

الفصل الأول

1.11 منهجية البحث

1.11.1 المقدمة

هذا القسم يتحدث عن المنهجية المتبعة في البحث من فهم المشاكل التي تواجهها المؤسسة وكذلك في فهم البيانات بصورة جيدة ومعرفة ماهية البيانات , ويوضح كيفية تجميع البيانات، ويوضح معني كل متغير من تلك البيانات كما تم استخدام منهجية ال CRISP-DM في التحليل.

• طبيعة الأعمال

من خلال دراستنا إلي البيانات المأخوذة من الإدارة العامة للمرور نجد أن عملية التحليل المستخدمة في البيانات عبارة عن تحليل تقليدي للبيانات أي أنها تستخدم الطريقة الإحصائية في تحليل بياناتها وهنا تكمن المشكلة لأنها لا تستفيد من البيانات لذلك لابد من إيجاد بديل في عملية التحليل للاستفادة من هذه البيانات بصورة كبيرة , لذلك تمكنت هذه الدراسة من إيجاد الحل والاستفادة من البيانات بصورة كبيرة حيث استخدمت تعلم الآلة في التحليل .

• فهم البيانات

من الخطوات الضرورية لعملية فهم البيانات يمكن تلخيصها في الآتي:-

■ تجميع البيانات

لقد تم تجميع هذه البيانات من الإدارة العامة للمرور من سجلات الحوادث لولاية الخرطوم وذلك من الفترة من يناير إلي ابريل 2017 م.

■ توصيف البيانات

تتكون البيانات من 12 متغير سنقوم بوصف كل متغير علي حدة .

• الأقسام : ونقصد بها أقسام محليات ولاية الخرطوم وهي بيانات نصية.

اسم المحلية	#
الخرطوم	1
الجبيل	2
امبدة	3
شرق النيل	4
كرري	5
بحري	6
امدرمان	7

جدول رقم (3-1-1)

اليوم : ونقصد بها أيام الأسبوع.

#	اليوم
1	الجمعة
2	السبت
3	الأحد
4	الاثنين
5	الثلاثاء
6	الأربعاء
7	الخميس

جدول رقم(2-1-3)

الشهور : ونقصد بها الشهور التي سجل فيها الحدث.

#	الشهر
1	يناير
2	فبراير
3	مارس
4	أبريل

جدول رقم(3-1-3)

نوع الحادث : يقصد به نوع الحدث الذي تم تسجيله وينقسم إلي(1 حالة وفاة, 2 تعني أذي جسيم).

عمر المتهم

#	عمر المتهم
2	الأعمار من 11 سنة إلي 20 سنة
3	الأعمار من 21 سنة إلي 30 سنة
4	الأعمار من 31 سنة إلي 40 سنة
5	الأعمار من 41 سنة إلي 50 سنة
6	الأعمار من 51 فما فوق

جدول رقم(4-1-3)

زمن الحادث: ويقصد به زمن وقوع الحادث الذي تم تسجيله وينقسم إلي(ص : صباحا , م : مساء).

الجنس : ويقصد به نوع الشخص ذكر أو انثي.

رخصة السائق : ويقصد بها الرخصة التي يستخدمها السائق.

#	نوع الرخصة
أجنبية	يحمل رخصة أجنبية
م خ	يحمل رخصة ملاكي الخرطوم
ع خ	يحمل رخصة عمومي الخرطوم
ع و	يحمل رخصه عمومي ولائي
م و	يحمل رخصه ملاكي ولائي
غير محدد	لا يمكن تحديد إذا كان يحمل رخصة
لا يملك	لا يملك رخصة قيادية

جدول رقم (3-1-5)

نوع المركبة - : ويقصد بها نوع العربة التي تسببت في الحادث ومن أمثلتها (موتر , صالون , حافلة , ركشة , وجميع الموديلات).

الشارع : ويقصد به الطريق الذي وقع فيه الحادث مثل (شارع مدني , شارع عبيد ختم , عمر الصول , وغيرها).

الموقع : يقصد به موقع الحادث.

عمر المصاب : هي نفس أعمار المتهمين بالاضافه إلي ("1" تعني الأعمار من 10 فما تحت).

#	عمر المتهم
2	الأعمار من 11 سنة إلي 20 سنة
3	الأعمار من 21 سنة إلي 30 سنة
4	الأعمار من 31 سنة إلي 40 سنة
5	الأعمار من 41 سنة إلي 50 سنة
6	الأعمار من 51 فما فوق
1	الأعمار من 10 فما تحت

جدول رقم (3-1-6)

The image shows a screenshot of a Microsoft Excel spreadsheet. The spreadsheet contains a list of data points, likely related to a study or survey. The columns are labeled with letters from V to S, and the rows are numbered from 1 to 278. The data points include names of individuals or groups, their locations, and various numerical values. The spreadsheet is displayed in a window titled 'بيانات خزانة' (Data Warehouse) and 'Microsoft Excel'.

صورة رقم (7-1-3)

• جودة البيانات و تحقيقها

يحتاج النموذج الجيد إلي بيانات جيدة مما يتوجب أن تكون البيانات صحيحة وذات مضمون دقيق لذلك توجب علينا حذف بعض البيانات التي تعتبر غير ضرورية أو مهمة وقمنا بدمج بعض البيانات وذلك للتقليل من البيانات في عملية التصنيف كما موضح في شكل رقم (7-1-3)

■ تهيئة البيانات

تم اختيار المتغيرات من البيانات وتنسيقها وإعادة ترتيب كل الحقول وجعل البيانات جاهزة لعملية التحليل كما موضح في الصورة رقم (8-2-3) .

id	departing date	date	type	time	age	driving lic	gender	bus	street	place	tabaco	al	injured	dead
1	الخميس 11/1	JAN	SUN	2 AM	4	ع	male	stroller	ساحل	لقطع الشرف	عابر	yes	no	
2	الخميس 11/1	JAN	MON	2 AM	4	ع	male	public	العرف	مخمة القعدة	عابر	yes	no	
3	الخميس 11/1	JAN	FRI	2 PM	5	ع	male	stroller	البنكية	كثري المشاة	عابر	yes	no	
4	الخميس 11/1	JAN	SUN	2 PM	3	ع	male	stroller	مدرسة كتيف	عند	تسليم	yes	no	
5	الخميس 11/1	JAN	SUN	2 PM	3	ع	male	stroller	عبد خلد	المركب الشرفي	عابر	yes	no	
6	الخميس 11/1	JAN	THU	1 AM	NA	ع	male	public	الروافض	للبيع ماله	عابر	yes	no	
7	الخميس 11/1	JAN	THU	1 AM	3	ع	male	stroller	كثري عابري	مطلة بنار	عابر	no	yes	
8	الخميس 11/1	JAN	FRI	1 PM	3	ع	male	stroller	عرب	الشمية	عابر	no	yes	
9	الخميس 11/1	JAN	SAT	2 PM	3	ع	male	stroller	الدارر القوي	القوية الأويبة	عابر	no	yes	
10	الخميس 11/1	JAN	SUN	2 PM	5	ع	male	public	عبد خلد	لقطع الشرف	عابر	yes	no	
11	الخميس 11/1	JAN	MON	2 PM	5	ع	male	stroller	عبد خلد	أويبة وسف	عابر	yes	no	
12	الخميس 11/1	JAN	TUE	1 AM	3	ع	male	NA	الروافض	أويبة وسف	عابر	no	yes	
13	الخميس 11/1	JAN	WED	2 PM	4	ع	male	stroller	أويبة وسف	أويبة وسف	عابر	yes	no	
14	الخميس 11/1	JAN	THU	2 PM	3	ع	male	stroller	أويبة وسف	أويبة وسف	عابر	yes	no	
15	الخميس 11/1	JAN	THU	2 PM	3	ع	male	stroller	أويبة وسف	أويبة وسف	عابر	yes	no	
16	الخميس 11/1	JAN	SAT	2 PM	5	ع	male	stroller	عبد خلد	كثري عابري	عابر	yes	no	
17	الخميس 11/1	JAN	SUN	1 PM	4	ع	male	stroller	عبد خلد	كثري عابري	عابر	yes	no	
18	الخميس 11/1	JAN	MON	2 AM	4	ع	male	stroller	عبد خلد	كثري عابري	عابر	yes	no	
19	الخميس 11/1	JAN	SUN	2 AM	4	ع	male	stroller	عبد خلد	كثري عابري	عابر	yes	no	
20	الخميس 11/1	JAN	TUE	2 AM	3	ع	male	stroller	عبد خلد	كثري عابري	عابر	yes	no	
21	الخميس 11/1	JAN	WED	2 PM	5	ع	male	stroller	عبد خلد	كثري عابري	عابر	yes	no	
22	الخميس 11/1	JAN	SUN	2 PM	3	ع	male	stroller	عبد خلد	كثري عابري	عابر	yes	no	
23	الخميس 11/1	JAN	SUN	2 PM	3	ع	male	stroller	عبد خلد	كثري عابري	عابر	yes	no	
24	الخميس 11/1	JAN	MON	1 PM	NA	NA	male	stroller	عبد خلد	كثري عابري	عابر	no	yes	
25	الخميس 11/1	JAN	MON	2 PM	5	ع	male	public	عبد خلد	كثري عابري	عابر	yes	no	
26	الخميس 11/1	JAN	TUE	2 PM	3	ع	male	stroller	عبد خلد	كثري عابري	عابر	yes	no	
27	الخميس 11/1	JAN	TUE	2 PM	4	ع	male	stroller	عبد خلد	كثري عابري	عابر	yes	no	
28	الخميس 11/1	JAN	TUE	2 PM	3	ع	male	stroller	عبد خلد	كثري عابري	عابر	yes	no	
29	الخميس 11/1	JAN	TUE	2 PM	3	ع	male	stroller	عبد خلد	كثري عابري	عابر	yes	no	
30	الخميس 11/1	JAN	TUE	1 AM	NA	NA	male	stroller	عبد خلد	كثري عابري	عابر	yes	no	

صورة رقم (8-3-2)

القسم الثاني

1.12 التصنيف

1.12.1 المقدمة

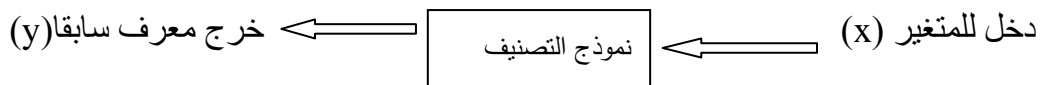
ينقسم تعلم الآلة إلى عدة أنواع منها التعلم الخاضع للإشراف ويسمى أيضاً التعلم التنبؤي في هذا النوع من التعلم يتم تدريب الآلة باستخدام دخل معروف الخرج مسبقاً ، مثلاً مجموعة من رسائل البريد الإلكتروني المصنفة مسبقاً إلى مهمة أو غير مهمة، والمطلوب تعلم كيفية ربط الدخل مع الخرج ليصبح بالإمكان مستقبلاً التنبؤ بالخرج من أجل أي دخل جديد. يندرج تحت هذا النوع أنواع فرعية من التعلم بحسب الخرج المطلوب من نظام تعلم الآلة ومن أمثلة هذا النوع هو التصنيف، ومن أهم التقنيات المستخدمة فيه وهي شجرة القرار، naïve Bayes، support vector machine.

1.12.2 التصنيف

وهو النوع الأكثر استخداماً في تعلم الآلة. في هذا النوع يكون الدخل مصنفاً إلى نوعين أو أكثر. وهدف عملية التعلم إنتاج نموذج يستطيع تصنيف أي دخل جديد إلى نوع أو أكثر من الأنواع المعروفة سابقاً.

1.12.3 تقنية التصنيف

التصنيف هو مهمة تعليم دالة معينة f لربط مجموعة الخصائص (x) بفئة (class) معرفة مسبقاً تسمى (y) تعرف الدالة (f) بنموذج التصنيف.



الشكل (1-2-3)

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

أ- يتم تدريب خوارزمية التصنيف (التعلم) علي بيانات التدريب المحتوية علي سجلات معروفة لبناء المصنف الذي يستخدم لفحص البيانات التي تحتوي علي البيانات .

ب- يتم اداء المصنف بحساب عدد السجلات المتوقعة المصنفة بشكل صحيح والسجلات المصنفة بشكل خاطئ فيما يسمى بمصفوفة التضارب يتم تقييم اداء المصنفات بالحصول علي اعلي دقة واقل نسبة خطأ عند تطبيقها علي بيانات الاختبار.

يتم جدولة عدد السجلات المصنفة بشكل صحيح وعدد السجلات المصنفة بشكل خاطئ علي شكل مصفوفة تسمى مصفوفة التعارض .

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

الشكل (2-2-3)

كل مدخل f_{ij} في مصفوفة التعارض يشير إلي عدد السجلات في class I المتوقعة أن تكون في ال class 0 فمثلا f_{01} تشير إلي عدد السجلات في ال class 0 التي تم توقعها بشكل خاطئ في class 1 .

اعتمادا علي المصفوفة فان:

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

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

1.12.4 تقنيات التصنيف

يتضمن التصنيف عدة تقنيات رئيسية سيناقدش هذا البحث ثلاثة تقنيات :

1. شجرة القرار

هنالك عدة خوارزميات في شجرة القرار في هذا البحث سوف يستخدم خوارزمية

.j48

2. Naïve Bayes

سوف يستخدم خوارزمية naïve Bayes

3. Super vector machine

سوف يستخدم خوارزمية support vector machine

	Decision Trees	Naïve Bayes	SVM
Accuracy in general	**	*	****
Speed of learning with respect to number of attributes and the number of instances	***	****	*
Speed of classification	****	****	****
Tolerance to missing values	***	****	**
Tolerance to irrelevant attributes	***	**	****
Tolerance to redundant attributes	**	*	***
Tolerance to highly Interdependent attributes (e.g. parity problems)	**	*	***
Dealing with discrete/binary/continuous attributes	****	***(not continuous)	** (not discrete)
Tolerance to noise	**	***	**
Dealing with danger of overfitting	**	***	**
Attempts for incremental learning	**	****	**
Explanation ability/transparency of knowledge/classifications	****	****	*
Model parameter handling	***	****	*

الجدول رقم (3-2-3)

الجدول رقم (3-2-3) يوضح بعض المقارنات للخوارزميات و علي هذا الاساس سوف يتم اختيار الخوارزمية الافضل في اداة عملية التصنيف.

الفصل الثالث

1.13 التدريب والتقييم

1.13.1 مقدمة

في هذا القسم يتم تدريب البيانات عن طريق إدخالها في weka وبعد ذلك يتم تقييم أداء المصنف.

1.13.2 تدريب الخوارزميات علي البيانات

تدريب خوارزمية supper vector machine علي بيانات حوادث المرور :

- استخدمت أداة weka 3.8.1 لتدريب خوارزمية support vector machine علي بيانات الحوادث لعمل (train) لبناء مصنف البيانات وتم اختباره علي بيانات الاختبار (test) ثم تم تقييم أداة المصنف من خلال حساب دقة التصنيف ونسبة الخطأ في التصنيف ويمكن حساب دقة المصنفات بحساب الآتي :-
 - نسبة الضبطية (precision) نسبة سجلات الclass المصنفة بشكل صحيح إلي السجلات المصنفة في ال class .
 - نسبة الاستدعاء (recall) نسبة سجلات الclass المصنفة بشكل صحيح إلي عدد إلي السجلات في ال class
 - الايجابية الكاذبة (FP) أي يتم تصنيف البيانات علي أنها كاذبة وهي في الحقيقة صحيحة
 - السلبية الكاذبة (fn) تصنيف علي أنها صحيحة وهي في الحقيقة مهددات .
 - الايجابية الصادقة (TP) تصنيف السجلات علي أنها صحيحة وهي في الحقيقة صحيحة.
 - السلبية الصادقة (fn) تصنيف السجلات علي أنها مهددات وهي في الحقيقة مهددات .
- يتم حساب دقة المصنف اعتمادا علي معدل الكشف ومعدل الإنذار الكاذب .
استخدم هذا البحث بيانات حوات المرور وكانت نتائج التدريب كما يلي :

Predicated class			
Actual class		Class dead	Class injured
	Class dead	851	0
	Class injured	0	367

بعد ادخال البيانات إلي weka واختيار التدريب المناسب وهنا تم اختيار متجة آلة الدعم تم الحصول علي الصورة رقم (3,3,1).

Time taken to build model: 0.48 seconds

=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances	1188	97.5369 %
Incorrectly Classified Instances	30	2.4631 %
Kappa statistic	0.9238	
Mean absolute error	0.0246	
Root mean squared error	0.1569	
Relative absolute error	7.433 %	
Root relative squared error	38.5741 %	
Total Number of Instances	1218	

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	0.993	0.090	0.977	0.993	0.985	0.925	0.951	0.975	injured
	0.910	0.007	0.971	0.910	0.939	0.925	0.951	0.902	dead
Weighted Avg.	0.975	0.073	0.975	0.975	0.975	0.925	0.951	0.960	

=== Confusion Matrix ===

```
a  b  <-- classified as
956  7  |  a = injured
 23 232 |  b = dead
```

صورة رقم (3,3,1)

تدريب خوارزمية naïve Bayes علي بيانات حوادث المرور

- استخدمت أداة weka 3.8.1 لتدريب خوارزمية naïve Bayes علي بيانات الحوادث لعمل (train) لبناء مصنف البيانات وتم اختباره علي بيانات الاختبار (test) ثم تم تقييم أداة المصنف من خلال حساب دقة التصنيف ونسبة الخطأ و معدل الإنذار الكاذب في التصنيف ويمكن حساب دقة المصنفات بحساب الآتي :-
 - نسبة الضبطية (precision) نسبة سجلات ال class المصنفة بشكل صحيح إلي السجلات المصنفة في ال class .
 - نسبة الاستدعاء (recall) نسبة سجلات ال class المصنفة بشكل صحيح إلي عدد إلي السجلات في ال class
 - الايجابية الكاذبة (FP) أي يتم تصنيف البيانات علي أنها كاذبة وهي في الحقيقة صحيحة
 - السلبية الكاذبة (fn) تصنيف علي أنها صحيحة وهي في الحقيقة مهددات .
 - الايجابية الصادقة (TP) تصنيف السجلات علي أنها صحيحة وهي في الحقيقة صحيحة.
 - السلبية الصادقة (fn) تصنيف السجلات علي أنها مهددات وهي في الحقيقة مهددات .
- يتم حساب دقة المصنف اعتمادا علي معدل الإنذار الكاذب .
- استخدم هذا البحث بيانات حوات المرور وكانت نتائج التدريب كما يلي :

Predicated class				
Actual class		Class dead	Class injured	
		Class dead	851	1
		Class injured	61	367

بعد ادخال البيانات إلي weka واختيار التدريب المناسب وهنا تم اختيار خوارزمية naïve bayes تم الحصول علي الصورة رقم (3,3,2).

```
=== Stratified cross-validation ===
=== Summary ===
```

```
Correctly Classified Instances      1186          97.3727 %
Incorrectly Classified Instances     32           2.6273 %
Kappa statistic                     0.919
Mean absolute error                 0.0306
Root mean squared error             0.1544
Relative absolute error             9.2373 %
Root relative squared error        37.9496 %
Total Number of Instances          1218
```

```
=== Detailed Accuracy By Class ===
```

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	0.991	0.090	0.976	0.991	0.984	0.920	0.953	0.975	injured
	0.910	0.009	0.963	0.910	0.935	0.920	0.953	0.943	dead
Weighted Avg.	0.974	0.073	0.974	0.974	0.973	0.920	0.953	0.968	

```
=== Confusion Matrix ===
```

```
 a  b  <-- classified as
954  9  |  a = injured
 23 232 |  b = dead
```

الصورة رقم (3,3,2)

تدريب خوارزمية j48 علي بيانات حوادث المرور :

استخدمت أداة weka 3.8.1 لتدريب خوارزمية j48 علي بيانات الحوادث لعمل (train) لبناء مصنف البيانات وتم اختباره علي بيانات الاختبار (test) ثم تم تقييم أداة المصنف من خلال حساب دقة التصنيف ونسبة الخطأ في التصنيف ويمكن حساب دقة المصنفات بحساب الآتي :-

- نسبة الضبطية (precision) نسبة سجلات الclass المصنفة بشكل صحيح إلي السجلات المصنفة في ال class .
- نسبة الاستدعاء (recall) نسبة سجلات الclass المصنفة بشكل صحيح إلي عدد إلي السجلات في ال class
- الايجابية الكاذبة (FP) أي يتم تصنيف البيانات علي أنها كاذبة وهي في الحقيقة صحيحة
- السلبية الكاذبة (fn) تصنيف علي أنها صحيحة وهي في الحقيقة مهددات .
- الايجابية الصادقة (TP) تصنيف السجلات علي أنها صحيحة وهي في الحقيقة صحيحة.

- السلبية الصادقة (fn) تصنيف السجلات علي أنها مهددات وهي في الحقيقة مهددات .
يتم حساب دقة المصنف اعتمادا علي معدل الإنذار الكاذب .
استخدم هذا البحث بيانات حوات المرور وكانت نتائج التدريب كما يلي :

Actual class	Predicated class		
		Class dead	Class injured
	Class dead	956	7
Class injured	23	232	

بعد ادخال البيانات إلي weka واختيار التدريب المناسب وهنا تم اختيار لخوارزمية j48 وقد تم الحصول علي الصورة رقم (3,3,3).

```
Time taken to build model: 0.06 seconds
--- Stratified cross-validation ---
--- Summary ---
Correctly Classified Instances      1118          97.5368 %
Incorrectly Classified Instances      30           2.4632 %
Kappa statistic                      0.9238
Mean Absolute error                  0.0481
Root Mean Squared error              0.1882
Relative Absolute error               14.9018 %
Root Relative Squared error          20.1468 %
Total Number of Instances           1148
--- Detailed Accuracy By Class ---

```

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	ROC Area	Class
Weighted Avg.	0.915	0.007	0.971	0.950	0.938	0.925	0.927	0.950	dead
	0.975	0.071	0.975	0.975	0.975	0.925	0.927	0.951	injured

```
--- Confusion Matrix ---
 * 0 1 - classified as
 000 1 0 = injured
 23 211 0 = dead
```

صورة رقم (3,3,3)

1.13.3 تقييم خوارزميات التصنيف

1. التقييم بالنسبة لخوارزمية j48

Actual class	Predicated class		
		Class dead	Class injured
	Class dead	956	7
Class injured	23	232	

accuracy	recall	FP	Fn	TP	TN	Precision
97.5%	99.2%	9%	0.7%	99.2 %	81.8%	97.6%

2. التقييم بالنسبة لخوارزمية supper vector machine

Predicated class			
Actual class		Class dead	Class injured
	Class dead	851	0
	Class injured	0	367

Accuracy	recall	FP	Fn	TP	TN	Precision
100%	70%	0%	0%	100%	100%	100%

3. التقييم بالنسبة لخوارزمية naïve bayes

Predicated class			
Actual class		Class dead	Class injured
	Class dead	850	1
	Class injured	61	306

accuracy	recall	FP	Fn	TP	TN	Precision
94.9%	99.9%	16.6%	0.11%	99.8%	26.2%	93.3%

1.13.4 مقارنة الخوارزميات في التقييم

	accuracy	recall	FP	FN	TP	TN	Precision
Naïve Bayes	94.9%	99.9%	16.6%	0.11%	99.8%	26.2%	93.3%
J48	97.5%	99.2%	9%	0.7%	99.2 %	81.8%	97.6%
Supper vector machine	100%	70%	0%	0%	100%	100%	100%
The pest algorithm	Supper vector machine	Naïve Bayes	Supper vector machine	Supper vector machine	Supper vector machine	Supper vector machine	Supper vector machine

جدول رقم (1-3-3)

الباب الرابع

النتائج
والتوصيات

النتائج والتوصيات

2.1 النتائج

- من خلال مقارنة الخوارزميات التي تم التوصل اليها في الجدول رقم () نجد ان
- تم تحليل البيانات بصورة منهجية وأصبحت جاهزة لاستنتاج علاقات أخرى.
 - اثبتت خوارزمية supper vector machine كفاءتها في تحليل البيانات بمعدل دقة 100% ومعدل خطأ 0%

2.2 التوصيات

- نسبةً للإمكانيات المتاحة والزمن لم يتمكن الدارسون من عمل تحليل عميق للبيانات نسبة لمحدودية المتغيرات. لذا نوصي الدارسون ب:
- المحاولة بقدر الإمكان علي الحصول علي البيانات المفقودة التي يمكن أن تؤثر في عملية بناء النموذج مثل (حالة السائق والإضاءة ونوع الموديل) .
 - زيادة عدد البيانات.
 - تجربة عدة خوارزميات في عملية التحليل.

2.3 الخاتمة

ومالنا في الختام إلا أن نقول الحمد لله الذي هدانا و ما كنا لنهتدي لولا أن هدانا الله ،أملين أن يستمر البحث حتى يصبح للجهة المعنية نظاما خاصا للتقريب في البيانات ، وان يستفيد من بعدنا من بحثنا في ما أصبنا وان يستفيدوا مما أخطانا فيه.

المصادر والمراجع

لجنة الامم المتحدة: تحسين السلامة المرورية علي الصعيد العالمي) وضع الاهداف الاقليمية والوطنية للخد من هذه الحوادث المرورية علي [] الطرق. 2012

بشير عباس, العلاق , الادارة الرقمية: المجالات والتطبيقات, مركز الامارات للدراسات الاستراتيجية, ابوظبي 2005 [2]

[3] Bazsalica M., Naim P., Data mining pour le Web, éd. Eyrolles, Paris, 2001

عبد الستار العلي, عامر ابراهيم قنديلجي, غساف العمرم, المدخل إلى إدارة المعرفة، دار الدسنة للنشر و التوزيع و الطباعة، الطبعة الاولى [4]

عمان. 2006,

[5] Hand d., Mannila H., Smyth R., Principles of Data Mining, MIT Press, London.

[6] Krishnaveni ans Dr. M. Hemalatha, "A perspective analysis of traffic Accident Using Data Mining Techniques", International Journal of computer Application

[7]Naina Mahajan and Bikram Pal Kaur, "Analysis of Factors of Road Traffic Accidents using Enhanced Decision Tree Algorithm "International Journal of Computer Applications (0975 – 8887) Volume 135 – No.6, February 2016

[8]manner K. Geetha and C. Vaishnavi, "Analysis on Traffic Accident Injury Level Using Classification", International Journal of Advanced Research in Computer Science and Software Engineering, Volume 5, Issue 2, February 2015, ISSN: 2277 128X.

[9]دراسة حوادث الطرق والسلامة المرورية, السودان 2015

[10] Traffic Accident Analysis Using Machine Learning Paradigms USA 2004.

[11]<http://www.alyaum.com/article/4054373>