

## اكتشاف عناوين URL الضارة باستخدام التعلم الآلي

سوسن خضر\*

(تاريخ الإيداع ٢٠٢٤ / ٧ / ٩ - تاريخ النشر ٢٠٢٤ / ١٠ / ٢٤)

### □ ملخص □

رابط URL الخبيث هو عنوان ويب يقود المستخدمين إلى محتوى ضار و عادةً ما يتم تصميم عناوين URL الخبيثة لخداع المستخدمين لفتح مواقع تحتوي على برامج ضارة وخبيثة تعرض بياناتهم وأجهزتهم للخطر. هناك العديد من الهجمات التي يمكن استغلالها باستخدام عناوين URL الخبيثة، مثل هجمات التصيد الاحتيالي وهجمات البرامج الضارة والعديد من الهجمات الأخرى التي تخدع المستخدمين عبر الإنترنت، ومن هنا تأتي الأهمية الكبرى لاكتشاف عناوين URL الخبيثة من أجل الحفاظ على خصوصية المستخدمين ومنع استغلالهم. مع التطور الكبير الذي يشهده العالم في أنظمة الذكاء الاصطناعي، أصبح استغلال القدرة الكبيرة لتقنيات التعلم الآلي على تحليل واكتشاف الأنماط داخل عناوين URL الخبيثة أمرًا بالغ الأهمية لتعزيز القدرة على اكتشاف التهديدات والاستجابة لها في الوقت المناسب وبطريقة فعالة. واستجابة لهذه المتطلبات يهدف البحث إلى استغلال القدرات التحليلية القوية للتعلم الآلي وتوظيفها في عملية تحليل عناوين URL وتحسين دقة وكفاءة أنظمة اكتشاف عناوين URL الخبيثة، مما يؤدي بدوره إلى تعزيز قدرة هذه الأنظمة ضد التهديدات السيبرانية. تم تصميم النموذج المقترح باستخدام شبكة عصبية اصطناعية متعددة الطبقات من نوع Multi-Layers Perceptron وهي عبارة عن خوارزمية تعلم آلي خاضعة للإشراف تسمح للنموذج بتعلم الأنماط والتمثيلات المعقدة في البيانات. تم تدريب النموذج على مجموعة بيانات تحتوي على 651191 عنوان URL مصنفة إلى (حميدة (Benign)، مشوهة (Defacement)، برامج ضارة (Malware)، احتيال (Phishing)). حقق النموذج المقترح دقة تصنيف بنسبة (93%) متفوقًا على خوارزميات شجرة القرار (Decision Tree) (82%)، والغابة العشوائية (Random Forest) (82%)، والانحدار اللوجستي (Logistic Regression) (80%) و(SVM) (80%). تكشف نتائج البحث أن تطبيق التعلم الآلي يعزز بشكل كبير من اكتشاف عناوين URL الضارة. من خلال تحليل متعمق للميزات المتنوعة، توضح الدراسة فعالية النموذج المقترح في تصنيف عناوين URL الضارة والحميدة. تكمن أصالة هذا البحث في التصميم المبتكر لشبكتنا العصبية المقترحة التي تتضمن أربع طبقات كثيفة (Dense Layers) وهندسة الميزات والتي تضمنت استخراج الميزات من عناوين URL المجردة ومعالجة عدم التوازن في مجموعة البيانات باستخدام طريقة SMOTE ومن بعدها إجراء عملية Standardization للميزات وهي تقنية تستخدم في معالجة البيانات مسبقًا حيث يتم إعادة قياس الميزات بحيث يكون لها خصائص التوزيع الطبيعي القياسي بمتوسط (0) وانحراف معياري (1).

الكلمات المفتاحية: التعلم الآلي، عناوين URL الضارة، الذكاء الاصطناعي

\* ماجستير في علوم الحاسب - الجامعة الافتراضية السورية

## Detect malicious URLs using machine learning

Sawsan Kheder\*

(Received 9/7/2024.Accepted24 /10/2024)

### □ABSTRACT □

A malicious URL is a web address that leads users to malicious content. Malicious URLs are usually designed to trick users into opening sites that contain malicious and malicious programs that put their data and devices at risk. There are many attacks that can be exploited using malicious URLs, such as phishing attacks, malware attacks, and many other attacks that deceive users online. Hence, the great importance of detecting malicious URLs in order to maintain users' privacy and prevent their exploitation. With the great development witnessed by the world in artificial intelligence systems, exploiting the great ability of machine learning techniques to analyze and discover patterns within malicious URLs has become crucial to enhance the ability to detect and respond to threats in a timely and effective manner. In response to these requirements, the research aims to exploit the powerful analytical capabilities of machine learning and employ them in the process of analyzing URLs and improving the accuracy and efficiency of malicious URL detection systems, which in turn enhances the ability of these systems against cyber threats. The proposed model was designed using a multi-layer artificial neural network of the Multi-Layers Perceptron type, which is a supervised deep learning algorithm that allows the model to learn complex patterns and representations in data. The model was trained on a dataset containing 651,191 URLs classified into (Benign, Defacement, Malware, Phishing). The proposed model achieved a classification accuracy of (93%), outperforming the Decision Tree (82%), Random Forest (82%), Logistic Regression (80%) and SVM (80%) algorithms. The research results reveal that the application of machine learning significantly enhances the detection of malicious URLs. Through an in-depth analysis of various features, the study demonstrates the effectiveness of the proposed model in classifying malicious and benign URLs. The originality of this research lies in the innovative design of our proposed neural network that includes four dense layers and feature engineering, which included extracting features from abstract URLs and addressing the imbalance in the dataset using the SMOTE method, and then performing the feature standardization process, which is a technique used in data preprocessing where the features are rescaled to have the characteristics of a standard normal distribution with a mean (0) and a standard deviation (1).

**Keywords:** Machine Learning, Malicious URLs, Artificial Intelligence

---

\*Master of computer science, Syrian Virtual Univ.

## 1- مقدمة:

في هذا العصر الذي يتميز بالانتشار الواسع للإنترنت في كل مكان، ومع زيادة عدد مستخدمي الإنترنت بشكل كبير جدًا، أدى ذلك إلى تصاعد تطور التهديدات السيبرانية التي أصبحت تشكل تحديًا أساسيًا للأنظمة الرقمية [1]. ومع الزيادة الكبيرة في الاعتماد على الإنترنت سواء من قبل الأفراد أو المنظمات أو الحكومات، وصل خطر وقوع هؤلاء المستخدمين ضحية للهجمات إلى مستويات غير مسبوقة [2]. إن الانتشار الواسع والمستمر لهذا النوع من التهديدات يتطلب التقدم والتطوير المستمر لمنهجيات الحماية، الأمر الذي يتطلب الانتقال من أساليب الحماية التقليدية إلى أساليب أكثر نكاءً تتكيف مع التغيرات والتطورات المستمرة في أنواع الهجمات وأساليبها. يركز هذا البحث بشكل خاص على اكتشاف عناوين URL الضارة لأنها تعتبر البوابة الأساسية للمستخدمين للوصول إلى محتويات الإنترنت [3]. يستغل المهاجمون هذه الثغرة من خلال إخفاء الروابط الضارة في عناوين URL أخرى تبدو غير ضارة أو عناوين شرعية ومعروفة بالنسبة للمستخدمين. أظهرت طرق الكشف التقليدية، مثل الكشف القائم على التوقيع والأنظمة القائمة على القواعد، قيودًا في قدرتها على مواكبة التطور الديناميكي والمستمر للهجمات [4]. تعمل عناوين URL الضارة كنقطة ضعف قوية تدفع المهاجمين إلى اختراق أنظمة ومعلومات المستخدمين الحساسة وتطوير أشكال مختلفة من أساليب الهجوم. غالبًا ما تبدو عناوين URL الضارة شرعية وغير ضارة وبالتالي تستغل ثقة المستخدمين في عناوين الويب الشرعية [5]. مع الانتشار الأخير للهجمات التصيد الاحتمالي وبرامج الفدية، تظهر الأهمية الكبرى لمعالجة هذا النوع من الثغرات الأمنية. واستجابة لهذه التحديات، يظهر التعلم الآلي، بقدراته التحليلية الهائلة وكفاءته في اكتشاف الأنماط المخفية داخل مجموعات البيانات الكبيرة، كنهج استراتيجي ومبتكر لتعزيز أنظمة الكشف عن الهجمات [6]. مع استمرار نمو حجم عناوين URL التي يتم الوصول إليها يوميًا عبر الإنترنت بشكل كبير، فإن هذا يجعل من غير العملي أن يكون المسح اليدوي آلية دفاع موثوقة. يتعمق هذا البحث في تطبيق تقنيات التعلم الآلي في الكشف عن أكثر من نوع من الهجمات بناءً على عناوين URL الضارة [7] [8]. من خلال الاستفادة من قدرات التعلم الآلي، يسعى هذا البحث إلى تعزيز دقة تصنيف عناوين URL الحميدة والخبيثة وبالتالي تعزيز كفاءة اكتشافها والتخفيف من خطر وقوع المستخدمين ضحية لهذه الهجمات.

## 2- الدراسات السابقة:

لقد أصبح تطبيق تقنيات الذكاء الاصطناعي في مجال الأمن السيبراني موضوعًا ذا أهمية كبيرة، وخاصة في الكشف عن عناوين URL الضارة. وقد استعرضت الدراسات السابقة العديد من الأساليب التي تهدف إلى تطبيق نماذج الذكاء الاصطناعي للكشف عن عناوين URL الضارة والتخفيف من المخاطر المرتبطة بها. وقد تناولت هذه الدراسات العديد من الجوانب مثل جمع البيانات وهندسة الميزات وبناء النماذج، والتي تهدف إلى تعزيز دقة وكفاءة تحديد عناوين URL الضارة. يستعرض هذا القسم مجموعة من الدراسات السابقة التي تناولت موضوع الكشف عن عناوين URL الضارة.

يناقش البحث [9] اكتشاف عناوين URL الضارة باستخدام تقنيات التعلم الآلي، ويركز على تطبيق التعلم الآلي الكمي (Quantum Machine Learning) في مجال الأمن السيبراني. استخدم البحث خوارزميات تعلم آلي مختلفة (Logistic Regression و Decision Tree و SVM والشبكات العصبية) لاكتشاف عناوين URL الضارة. تبدأ المنهجية بالتحليل الاستكشافي للمفاهيم الأساسية في التعلم الآلي، ووصف مجموعات البيانات، ومعالجة البيانات مسبقًا، ونتائج التجارب باستخدام خوارزميات مختلفة. قارنت المنهجية نماذج التعلم الآلي الكلاسيكية بناءً على بعض

مقاييس الأداء (Precision, Recall, Accuracy). يستكشف البحث تطبيق التعلم الآلي الكمي، وخاصة المصنف الكمومي المتغير (VQC) Variable Quantum Classifier. وناقش تكييف مجموعات البيانات للخوارزميات الكمومية، واستخدام البوابات والدوائر الكمومية، وتقييم النماذج الكمومية بالمقارنة مع نماذج التعلم الآلي الكلاسيكية.

يناقش البحث [10] اكتشاف عناوين URL الضارة باستخدام نموذج يعتمد على الشبكة العصبية المتوازية (Parallel Neural Network). ويقدم نموذجًا للشبكة العصبية المتوازية لانتقاط معلومات عناوين URL مثل المعلومات الدلالية والبصرية وهذا يلغي هندسة الميزات اليدوية. يتضمن النهج خطوات مختلفة مثل استخراج الميزات والتصنيف والجمع بين الميزات المرئية والدلالية باستخدام تقنية الانتباه (Attention Technique). تعتمد هذه الدراسة على تقنيات التعلم العميق للتعلم التلقائي واستخراج وتمثيل ميزات عناوين URL دون الاعتماد على هندسة الميزات اليدوية. تبدأ المنهجية المقترحة بتضمين الأحرف لاستخراج المعلومات الدلالية وتحويل عناوين URL إلى صور بدرجات الرمادي. يعتمد النموذج المقترح على استخدام شبكة عصبية متكررة مستقلة (IndRNN) لمعالجة مشكلة التدرج المتلاشي في شبكات RNN التقليدية. استخدمت المنهجية أيضًا تقنية الانتباه لاستخراج ميزات مفيدة من مجموعة البيانات لتصنيف البيانات عن طريق ترجيح بيانات الإدخال.

في البحث [11]، تم اقتراح نظام للكشف عن عناوين URL الضارة. تحتوي مجموعة البيانات على 26054 عينة بيانات من عناوين URL تم جمعها من مصادر مختلفة. تتضمن المنهجية استخراج 41 ميزة من عناوين URL. ثم استخدمت الدراسة اختبار ANOVA لاختيار أهم الميزات وحصلت على 17 ميزة من 41 ميزة. تم استخدام مجموعة البيانات لتدريب مصنف خوارزمية XGBoost. يتضمن تحليل الأداء تقييم دقة واستقرار النموذج المقترح باستخدام التحقق المتبادل k-fold والذي يحتوي على 10-Folds.

ركز البحث [12] على تطبيق التعلم الآلي في الكشف عن عناوين URL الضارة. وقد استخدم خوارزمية مستوحاة من البيولوجيا والتعلم الآلي لتحسين الميزات. استخدمت الدراسة مصنفات التعلم الآلي المختلفة للكشف عن عناوين URL الضارة. تتضمن المنهجية المقترحة ثلاث مراحل: جمع البيانات وخوارزميات التعلم الآلي ومجموعة البيانات. تبدأ المنهجية باختيار الميزة لجمع الميزات التي سيتم استخدامها لتدريب مصنفات التعلم الآلي. بعد ذلك يتم تحسين الميزات باستخدام تحسين سرب الجسيمات Particle Swarm Optimization (PSO). وقد استخدمت خوارزمية

Naive Bayes و SVM للتدريب والاختبار. استخدمت الدراسة بعض مقاييس الأداء مثل الدقة ومعدل الإيجابية الحقيقية (TPR) ومعدل الإيجابية الكاذبة (FPR) والدقة والاستدعاء (Recall) لتحديد الأداء لكل مصنف.

يقدم البحث [13] نموذجًا للكشف عن عناوين URL الضارة باستخدام التعلم الجماعي (Ensemble Learning). يستخدم النموذج خطوات مختلفة مثل استخراج الميزات ومعالجة البيانات المسبقة وتقنيات TF-IDF. يحتوي هذا النموذج على ثلاث مكونات رئيسية:

(1) جمع الميزات: يتم في هذه المرحلة استخراج ميزات عنوان URL.

(2) معالجة البيانات: استخدمت هذه المرحلة TF-IDF (Term Frequency – Inverse Document Frequency) لتمثيل الميزات.

(3) مرحلة التصنيف: تُستخدم خوارزمية الغابة العشوائية لتحديد ما إذا كان عنوان URL ضارًا أم حميدًا.

يناقش البحث [14] أهمية أمن الإنترنت في الآونة الأخيرة، وخاصة مع ظهور التطبيقات والمستخدمين. ويسلط الضوء على التحول إلى المنصات الرقمية بسبب جائحة كوفيد-19 والتهديدات الأمنية المتزايدة مثل تسريبات البيانات. يتم استكشاف تقنيات التعلم الآلي لتصنيف مواقع الويب لتعزيز أمان المستخدم. يتم شرح خوارزميات مختلفة مثل أشجار القرار (DT)، والانحدار اللوجستي (LR)، و SVM، والانحدار التدريجي العشوائي في سياق الأمن السيبراني. تضمنت المنهجية الاستفادة من مجموعة بيانات "المواقع الضارة" من Kaggle، والتي كان بها في البداية اختلال في التوازن متحيز نحو المواقع الضارة. لمعالجة هذا الخلل، تم تطبيق ثلاث تقنيات لموازنة البيانات: نقص العينة، والإفراط في أخذ العينات، و SMOTE. بعد موازنة مجموعة البيانات، تم تنفيذ التحقق المتبادل K Fold لتقييم أداء النموذج.

### 3- أهمية البحث وأهدافه:

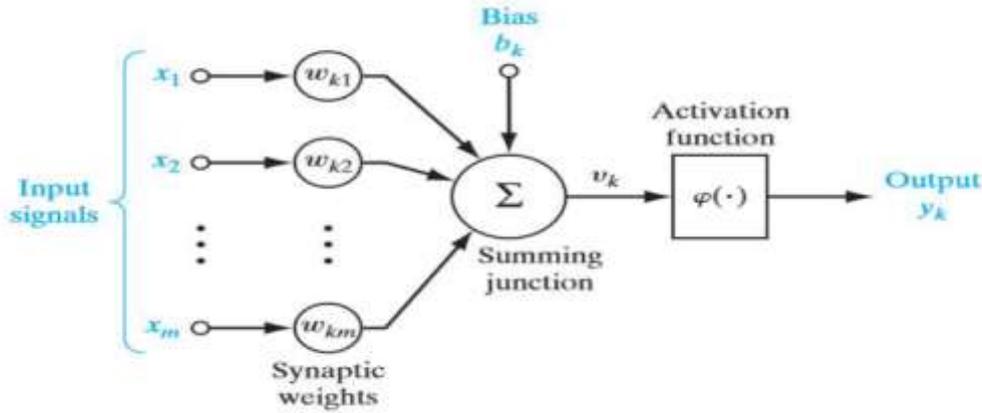
الهدف الأساسي من البحث هو تطوير نظام قوي ودقيق وفعال للكشف عن عناوين URL الضارة باستخدام شبكات عصبية متعددة الطبقات من نوع Multi-Layers Perceptron (MLP) وهي عبارة عن خوارزمية تعلم آلي خاضعة للإشراف تسمح للنموذج بتعلم الأنماط والتمثيلات المعقدة في البيانات. تسعى هذه الدراسة إلى معالجة قيود طرق الكشف التقليدية عن عناوين URL من خلال الاستفادة من قدرات MLP لتعزيز دقة الكشف. تشمل الأهداف المحددة تصميم وتنفيذ نموذج قائم على MLP، وتدريب النموذج على مجموعات بيانات شاملة تضم عناوين URL الحميدة والضارة، وتقييم أداء النموذج المقترح لتحقيق الأداء الأمثل. في النهاية، يطمح البحث إلى المساهمة في مجال الأمن السيبراني من خلال توفير حل قابل للتطوير وفعال للكشف عن تهديدات عناوين URL، وبالتالي تعزيز السلامة العامة وموثوقية استخدام الإنترنت.

### 4- طرق البحث ومواده:

#### 1.4 - الشبكات العصبية الاصطناعية:

الشبكات العصبية الاصطناعية هي طريقة من طرق الذكاء الاصطناعي التي يتعلم من خلالها الحاسب الآلي معالجة البيانات بطريقة مستوحاة من الدماغ البشري. تستخدم الشبكات العصبية الاصطناعية عقدًا أو عصبونات مترابطة في بنية متعددة الطبقات تشبه الدماغ البشري. تتكون الشبكات العصبية الاصطناعية من عقد أو عصبونات أو وحدات معالجة متصلة ببعضها البعض لتكوين شبكة من العقد تشبه عمل الخلايا العصبية البيولوجية أو الهياكل الإلكترونية [15]. يتم تحديد السلوك العام للشبكة من خلال هذا الاتصال، ويتم استخدام النموذج الرياضي الخاص بالشبكة لمعالجة المعلومات بناءً على طريقة الاتصال فيما بينها. تنتظم العقد العصبية الاصطناعية في طبقات متوازية طبقة إدخال وطبقة إخراج، وتعتبر هاتان الطبقتان الرئيسيتان، ويوجد بين هاتين الطبقتين أيضًا مجموعة من الطبقات المخفية وترتبط كل هذه الطبقات ببعضها البعض لتكوين بنية الشبكة العصبية الاصطناعية [16]. ولأن الشبكات العصبية الاصطناعية تشبه الدماغ البشري، فهذا يعني أنها تتغير وتتطور باستمرار، وهذا يعني أن الشبكات العصبية الاصطناعية قادرة على التعلم وطريقة تحليلها للبيانات تتغير، وبالتالي ستتجنب الأخطاء السابقة. تستقبل الشبكة

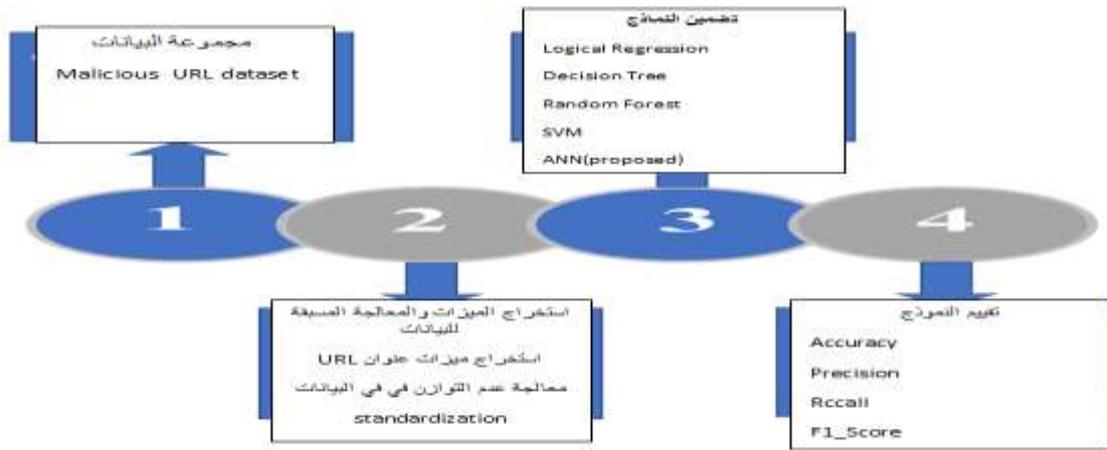
العصبية الاصطناعية مجموعة من المدخلات وتعالجها من خلال مجموعة من الأوزان المرتبطة بمدخلات كل خلية عصبية، وبعد ذلك يتم استخدام دالة رياضية تسمى دالة التنشيط (Activation Function)، والتي تنتج على مخرجاتها قيمة تمثل مخرجات الشبكة العصبية. يوضح الشكل (1) مبدأ عمل الشبكات العصبية الاصطناعية.



الشكل (1) مبدأ عمل الشبكات العصبية الاصطناعية

#### 2.4- المنهجية المقترحة:

تعتمد المنهجية المقترحة على بناء شبكة عصبية اصطناعية لتصنيف عناوين المواقع الضارة، ويمثل الشكل (2) الهيكل العام للمنهجية المقترحة.



الشكل (2) الهيكل العام للمنهجية المقترحة

#### 1.2.4 - مجموعة البيانات:

لقد استخدمنا مجموعة بيانات تتضمن 651191 عنوان URL تم الحصول عليها من موقع Kaggle. يتم تصنيف عناوين URL ضمن مجموعة البيانات إلى أربع فئات ( Benign URL, Defacement URL, Malware URL, Phishing URL). يوضح الشكل (3) جزء من مجموعة البيانات المستخدمة

url	type
br-icloud.com.br	phishing
mp3raid.com/music/krizz_kaliko.html	benign
bopsecrets.org/rexroth/cr/1.htm	benign
http://www.garage-pirene.be/index.php?option=com_content&view=article&id=70&vsig70_0=15	defacemer
http://buzzfil.net/m/show-art/ils-etaient-loin-de-s-imaginer-que-le-hibou-allait-faire-cesti-quand-ils-filmaient-2.html	benign
espn.go.com/nba/player/_/id/3457/brandon-rush	benign
yourbittorrent.com/?q=anthony-hamilton-soulife	benign
http://www.pashminaonline.com/pure-pashminas	defacemer
allmusic.com/album/crazy-from-the-heat-r16990	benign
corporationwiki.com/Ohio/Columbus/frank-s-benson-P3333917.aspx	benign
http://www.ikenmijnkunst.nl/index.php/exposities/exposities-2006	defacemer
myspace.com/video/vid/30602581	benign

الشكل (3) جزء من مجموعة البيانات المستخدمة

#### 2.2.4 - استخراج الميزات:

عناوين URL هي سلاسل من الأحرف التي توفر عنوان أو موقع الموارد على الإنترنت، مثل مواقع الويب أو الملفات أو الخدمات. يتضمن استخراج الميزات من عناوين URL تقسيم عنوان URL إلى مكوناته وتحديد العناصر الرئيسية التي يمكن استخدامها لتطبيقات مختلفة [17]. في منهجيتنا المقترحة، يتم تطبيق مجموعة من العمليات على مجموعة البيانات المستخدمة، واستخراج 19 ميزة لكل عنوان URL ضمن مجموعة البيانات وإنشاء مجموعة البيانات الجديدة التي تحتوي على ميزات كل عنوان URL وتصنيفه. يوضح الجدول (1) الميزات المستخرجة ووصف لكل ميزة.

الجدول (1) الميزات المستخرجة من عناوين URL

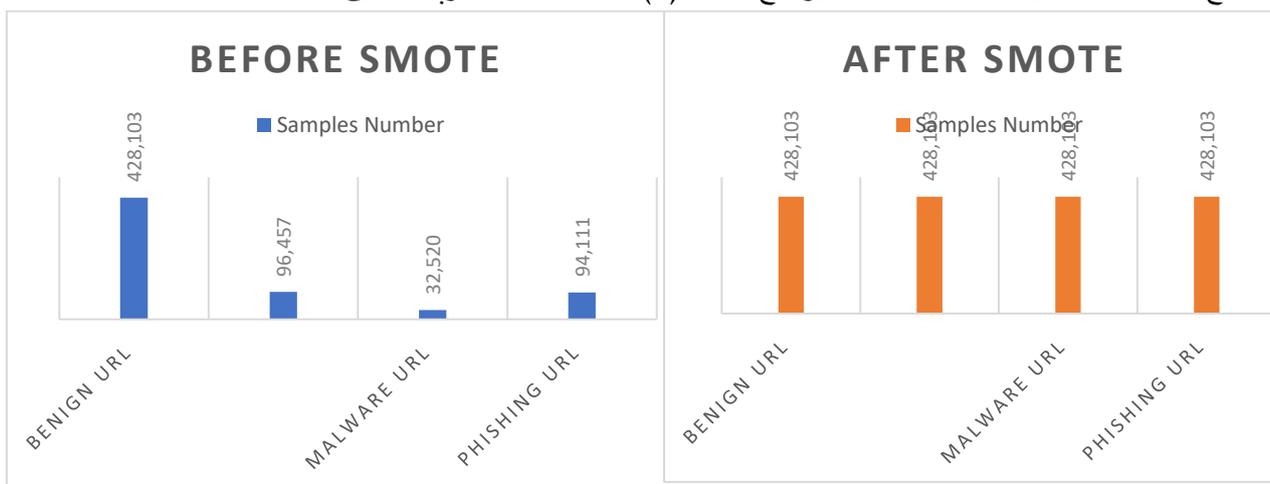
الوصف	
يشير إلى ما إذا كان عنوان URL يستخدم عنوان IP كجزء من المجال بدلاً من اسم المضيف	Use_of_ip
يقوم بحساب عدد النقاط (dots) في عنوان URL مما يوفر معلومات حول بنية المجال	Count.
يقوم بحساب عدد مرات ظهور "www" في عنوان URL	Count-www
عدد مرات ظهور الرمز "@" في عنوان URL	Count@
عدد المسارات أو المجلدات المحددة في مسار عنوان URL	Count_dir
عدد مرات ظهور المجالات المضمنة داخل عنوان URL، والتي قد تكون مؤشراً على محاولات التصيد الاحتيالي	Count_embed_domain
يشير إلى ما إذا كان عنوان URL قد خضع لاختصار URL، وهو أمر شائع في عناوين URL الضارة	Short_URL
عدد مرات ظهور "https" في عنوان URL	Count-https



بشكل جيد على فئة الأغلبية ولكنها ضعيفة على فئة الأقلية، حيث قد يتعلم النموذج التنبؤ دائماً بفئة الأغلبية بسبب عدد العينات الكبير لهذه الفئة ضمن مجموعة البيانات [19]. يعد معالجة هذا الخلل في التوازن أمراً بالغ الأهمية لضمان اكتشاف النموذج بدقة لفئة الأقلية. تُستخدم تقنيات - مثل إعادة أخذ العينات (أخذ عينات زائدة من فئة الأقلية (Oversampling) أو أخذ عينات أقل من فئة الأغلبية (Undersampling)) - بشكل شائع للتخفيف من التحديات التي تفرضها مجموعات البيانات غير المتوازنة [20].

تحتوي مجموعة البيانات المستخدمة في هذا البحث على 428103 عينة لفئة عناوين URL الحميدة (Benign URL)، و96457 عينة لفئة عناوين URL المشوهة (Defacement URL)، و32520 عينة لفئة عناوين URL الضارة (Malware URL)، و94111 عينة لفئة عناوين URL الاحتيالية (Phishing URL). وبالتالي، فإن مجموعة البيانات لدينا غير متوازنة. وبالتالي، في هذه المرحلة، تتضمن الخطوة الحاسمة تصحيح الخلل الموجود في مجموعة البيانات.

عادةً ما يتم تفضيل الإفراط في أخذ العينات (Oversampling) على أساليب إنقاص العينات (Undersampling) لأن إنقاص العينات يؤدي إلى تجاهل الحالات التي قد تحتوي على معلومات مهمة [21]. لقد استخدمنا طريقة SMOTE (Synthetic Minority Oversampling Technique) وهي طريقة أخذ عينات زائدة تنتج عينات اصطناعية لفئة الأقلية [22]. يوضح الشكل (5) عدد العينات قبل وبعد تطبيق تقنية SMOTE.



الشكل (5) عدد العينات قبل وبعد تطبيق تقنية SMOTE

يعد ضمان الاتساق (Consistency) في بيانات الإدخال الرقمية أمراً بالغ الأهمية لتحسين أداء خوارزميات التعلم الآلي. لتحقيق هذا الاتساق، من الضروري ضبط البيانات إلى نطاق موحد [23]. لقد استخدمنا StandardScaler وهي تقنية معالجة مسبقة للبيانات تُستخدم في الذكاء الاصطناعي والإحصاء. والهدف منها هو تحويل قيم ميزات مجموعة البيانات بحيث يكون متوسطها 0 وانحرافها المعياري 1. تهدف هذه العملية إلى جعل جميع السمات على مقياس مماثل [24]. يقوم StandardScaler بقياس كل ميزة ( $x_i$ ) في مجموعة البيانات عن طريق طرح المتوسط ( $\mu$ ) من تلك الميزة ثم القسمة على الانحراف المعياري ( $\sigma$ ):

$$\text{Standardized feature} = \frac{x_i - \mu}{\sigma}$$

الهدف الأساسي لاستخدام StandardScaler هو:

1- يساعد StandardScaler في ضمان أن جميع الميزات في مجموعة البيانات لها مقاييس مماثلة. وهذا مهم لأن نماذج الكفاءة الاصطناعي حساسة لحجم ميزات الإدخال. إذا كانت الميزات بمقاييس مختلفة، فقد يتسبب هذا في إعطاء النموذج أهمية أكبر للميزات ذات المقاييس الأكبر.

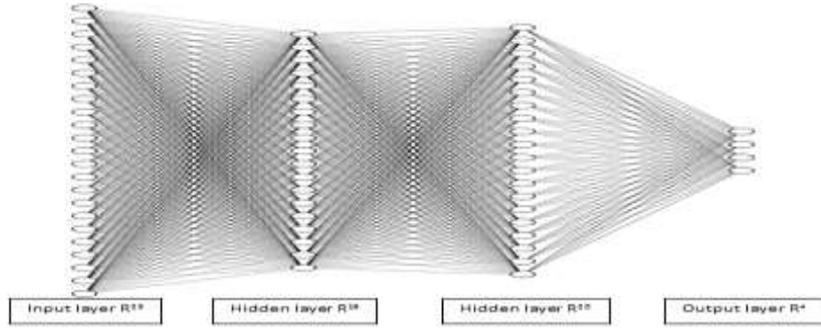
2- بالنسبة لنماذج التعلم الآلي والتعلم العميق فإن وجود ميزات على نطاق مماثل يؤدي إلى تقارب أسرع في الوصول الى الاداء الأمثل.

3- يمكن أن يساعد StandardScaler في التخفيف من المشكلات المتعلقة بالتذبذبات الرقمية والصعوبات في الحوسبة عندما تكون للميزات مقاييس مختلفة.

يعد StandardScaler مفيدًا بشكل خاص في الحالات التي تحتوي فيها الميزات على قيم مختلفة من حيث الحجم، مثل عند التعامل مع مجموعات بيانات تحتوي على ميزات مثل العمر والدخل ودرجة الحرارة. يضمن توحيد الميزات عدم هيمنة الميزات ذات المقاييس الأكبر على أداء النموذج، مما يؤدي بدوره إلى نموذج أكثر استقرارًا ويحسن الأداء العام للنموذج.

#### 4.2.4 - تصميم النموذج:

في هذا البحث تم تطوير نموذج لتصنيف عناوين المواقع باستخدام تقنيات التعلم الآلي. يوضح الشكل (6) بنية النموذج المقترح.



الشكل (6) بنية النموذج المقترح

يوضح الشكل (7) ملخص للنموذج المقترح الممثل في الشكل (6).

Layer (type)	Output Shape	Param #
dense (Dense)	(None, 23)	460
dense_1 (Dense)	(None, 19)	456
dense_2 (Dense)	(None, 20)	400
dense_3 (Dense)	(None, 4)	84

Total params: 1400 (5.47 KB)  
 Trainable params: 1400 (5.47 KB)  
 Non-trainable params: 0 (0.00 Byte)

الشكل (7) ملخص للنموذج المقترح.

يمثل هذا الملخص الموضح في الشكل (7) النموذج المقترح المكون من أربع طبقات متصلة بالكامل (Fully Connected Layers).

تحتوي الطبقة الأولى على 23 خلية عصبية، وكل خلية عصبية متصلة بالكامل بطبقة الإدخال وهي تتضمن 460 بارامتر قابل للتدريب.

تحتوي الطبقة الثانية على 19 خلية عصبية، متصلة بالكامل بالطبقة السابقة التي تحتوي على 23 خلية عصبية وهي تتضمن 456 بارامتر قابل للتدريب. تحتوي الطبقة الثالثة على 20 خلية عصبية، متصلة بالكامل بالطبقة السابقة التي تحتوي على 19 خلية عصبية وهي تتضمن 400 بارامتر قابل للتدريب. تحتوي الطبقة الرابعة على 4 خلايا عصبية، متصلة بالكامل بالطبقة السابقة التي تحتوي على 20 خلية عصبية وهي تتضمن 84 بارامتر قابل للتدريب. تمثل هذه الطبقة طبقة الإخراج وتتضمن 4 خلايا عصبية لأن النموذج يحل مشكلة تصنيف مكونة من 4 فئات. بالتالي يكون إجمالي البارامترات القابلة للتدريب هو 1400 بارامتر.

يتم ادخال بيانات التدريب الى النموذج على شكل دفعات (Batches) بحيث يكون حجم الادخال هو  $(n\_samples \times n\_features)$  حيث تمثل  $n\_samples$  حجم الدفعة (Batch Size) والتي تم ضبطها إلى  $(Batch\ Size = 32)$  و  $n\_features$  تمثل عدد الميزات وهي عبارة عن 19 ميزة موضحة في الجدول (1)، بالتالي يكون حجم كل إدخال هو  $(32 \times 19)$ .

تم استخدام دالة تنشيط Rectified Linear Unit (RELU) للطبقات الثلاث الأولى ودالة تنشيط Softmax لطبقة الإخراج الأخيرة. تحتوي طبقة الإخراج على 4 وحدات لمهمة تصنيف بأربع فئات (أربعة أنواع من عناوين URL).

ReLU: دالة تنشيط غير خطية تستخدم بشكل أساسي في الطبقات المخفية لتقديم عدم الخطية ومنع مشكلة التدرج المتلاشي. تقوم بإخراج المدخلات مباشرة إذا كانت موجبة وإلا فإنها تخرج صفرًا. Softmax: دالة تنشيط تستخدم في طبقة الإخراج لمهام التصنيف متعددة الفئات. تقوم بتحويل درجات التنبؤ إلى توزيع احتمالي على الفئات بحيث يمثل الاحتمال الأعلى الفئة التي تنتمي إليها عينة الدخل.

#### 5.2,4 - تقييم النموذج:

إن مصفوفة الارتباك (CM) Confusion Matrix التي تشير إلى مصفوفة الخطأ [25] عبارة عن تخطيط جدول مصمم لتمثيل أداء النموذج، غالبًا ما يتم استخدامه في التعلم الخاضع للإشراف [26]. يتم توضيح بنية مصفوفة الارتباك في الشكل (8).

		Predicted Values	
		1	0
Actual Values	1	TP	FN
	0	FP	TN

الشكل (8) بنية مصفوفة الارتباك

حيث:

TP (True Positive): عندما يقوم النموذج بتنبؤات إيجابية صحيحة.

TN (True Negative): عندما يقوم النموذج بتنبؤ سلبي صحيح.

FP (False Positive): عندما يقوم النموذج بتنبؤ إيجابي غير صحيح.

FN (False Negative): عندما يقوم النموذج بتنبؤ سلبي غير صحيح.

يتم استخدام هذه المعلمات لحساب recall, precision, accuracy.

$$Recall = \frac{TP}{TP + FN}$$

توضح Recall أنه "بالنسبة لجميع الفئات الإيجابية، ما عدد الفئات التي تم التنبؤ بها بشكل صحيح". [27]

$$Precision = \frac{TP}{TP + FP}$$

توضح قيمة Precision أنه "من بين جميع الفئات الإيجابية المتوقعة، كم عدد الفئات الإيجابية فعلياً". [28]

يتم تعريف الدقة (Accuracy) على أنها النسبة المئوية للتنبؤات الصحيحة التي يحققها النموذج بالنسبة لجميع التنبؤات [22]. يتم إعطاء معادلة الدقة كما يلي:

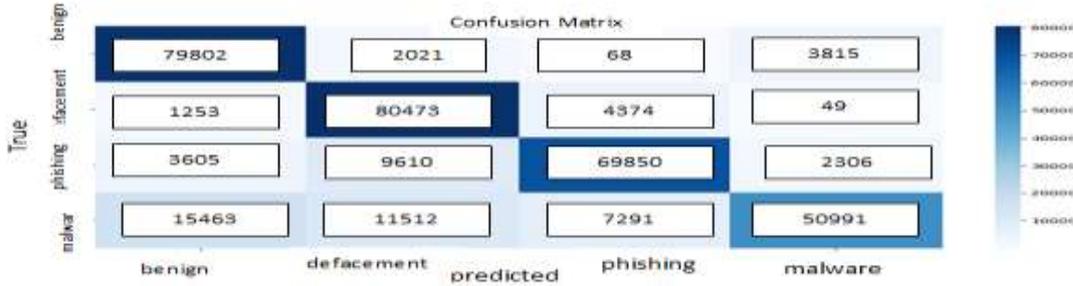
$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

F1\_Score هو مقياس يجمع بين Precision و Recall، وبالتالي يوفر تقييماً شاملاً لأداء النموذج. يتم تعريف F1\_Score على أنها المتوسط التوافقي ل Precision و Recall، وتؤدي إلى قيمة واحدة تمثل كل من الإيجابيات الخاطئة والسلبيات الخاطئة. تعطر معادلة F1\_Score كما يلي:

$$F1\_Score = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

## 5- النتائج والمناقشة:

سيتم في هذا القسم استعراض نتائج المحاكاة التي تم الوصول إليها ومقارنة نتائج النموذج المقترح مع كل من خوارزميات شجرة القرار، الغابة العشوائية، الانحدار اللوجستي و SVM. يمثل الشكل (9) مصفوفة الارتباك لخوارزمية شجرة القرار.



الشكل (9) مصفوفة الارتباك لخوارزمية شجرة القرار

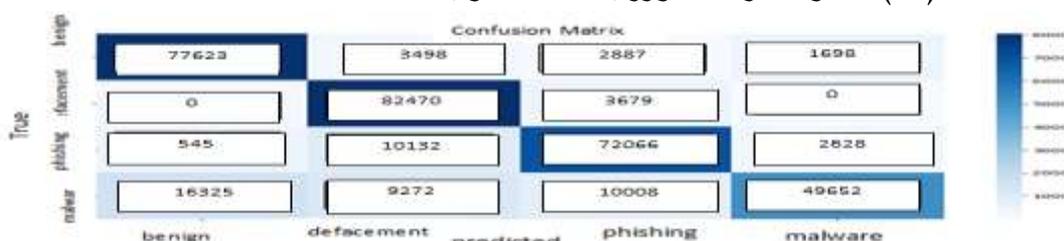
تمثل عناصر القطر الرئيسي التوقعات الصحيحة لكل فئة. على سبيل المثال، يشير العنصر الموجود في أعلى اليسار (79802) إلى 79802 توقعًا صحيحًا لفئة عنوان Benign. تمثل العناصر التي تكون خارج القطر الرئيسي التصنيفات الخاطئة. على سبيل المثال، يشير العنصر الموجود في الصف الثاني، العمود الأول (1253) إلى أن 1253 حالة من فئة عنوان Defacement تم تصنيفها خطأً على أنها فئة عنوان Benign. يوضح الجدول 2 قيم Precision و Recall و F1\_Score لكل فئة لخوارزمية شجرة القرار.

الجدول (2) قيم Precision و Recall و F1\_Score لكل فئة لخوارزمية شجرة القرار

	Precision	Recall	F1_Score
Benign URL	0.80	0.93	0.86
Defacement URL	0.78	0.93	0.85
Malware URL	0.86	0.82	0.84
Phishing URL	0.89	0.60	0.72

نلاحظ من الجدول (2) أن خوارزمية شجرة القرار تعمل بشكل جيد من حيث قيمة precision، وخاصة بالنسبة لعناوين URL الحميدة (Benign) والمشوهة (Defacement). وهذا يعني أنها تحدد بشكل صحيح معظم عناوين URL التي تنتمي إلى كل فئة. قيمة Recall جيدة أيضًا بالنسبة لمعظم الفئات، وخاصة عناوين URL الحميدة والمشوهة. وهذا يشير إلى أنها تكتشف نسبة عالية من الإيجابيات الحقيقية ضمن كل فئة. تواجه شجرة القرار صعوبة في التعامل مع عناوين URL للتصيد الاحتمالي (Phishing URLs)، حيث تظهر قيمة precision أعلى من Recall وهذا يعني أنها تحدد بدقة معظم عناوين URL التي لا تمثل تصيد احتيالي ولكنها تقوت جزءًا كبيرًا من عناوين URL للتصيد الاحتمالي الحقيقية. وقد يكون هذا بسبب تعقيد صياغة عناوين URL للتصيد الاحتمالي والميزات المحدودة التي تستخدمها شجرة القرار وهذه تعتبر مشكلة يمكن أن تؤدي إلى وقوع العديد من المستخدمين كضحية لعمليات التصيد الاحتمالي. بشكل عام، يمكن القول أن شجرة القرار تمثل نهجًا جيدًا للكشف عن عناوين URL الضارة، وخاصةً للفئات الحميدة والمشوهة. ومع ذلك، هناك حاجة إلى مزيد من التحليل لمعالجة ضعف الأداء مع عناوين URL للتصيد الاحتمالي.

يمثل الشكل (10) مصفوفة الارتباك لخوارزمية الغابة العشوائية.



الشكل (10) مصفوفة الارتباك لخوارزمية الغابة العشوائية

الجدول (3) يوضح قيم Precision, Recall and F1\_Score لكل فئة لخوارزمية الغابة العشوائية.

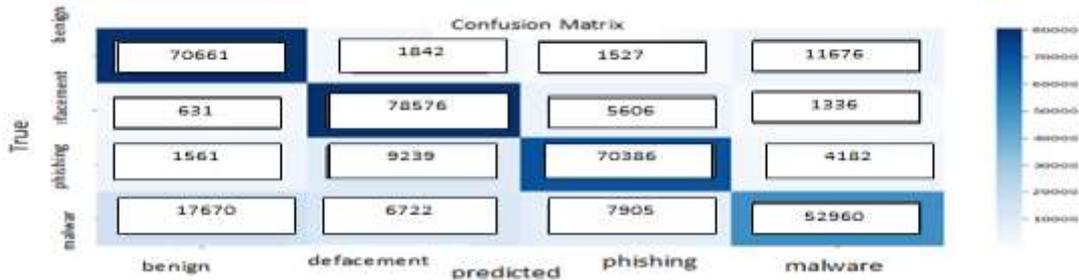
الجدول (3) قيم Precision و Recall و F1\_Score لكل فئة لخوارزمية الغابة العشوائية

	Precision	Recall	F1_Score
Benign URL	0.82	0.91	0.86
Defacement URL	0.78	0.96	0.86
Malware URL	0.81	0.84	0.83
Phishing URL	0.92	0.58	0.71

نلاحظ من الجدول (3) أن خوارزمية الغابة العشوائية تظهر قيمة precision مماثلة لشجرة القرار لمعظم الفئات، مما يشير إلى أنها تحدد بدقة معظم عناوين URL التي تصنفها. تظهر كلتا الخوارزمتين قيمة Recall مرتفعة لمعظم الفئات باستثناء عناوين URL الخاصة بالتصيد الاحتمالي حيث نلاحظ أنه على غرار شجرة القرار، تواجه الغابة العشوائية صعوبات في التعامل مع عناوين URL الخاصة بالتصيد الاحتمالي، حيث تظهر قيمة precision أعلى

ولكن قيمة Recall أقل. وهذا يشير إلى أن كلتا الخوارزميتين تواجهان تحديات مماثلة في الكشف بدقة عن محاولات التصيد الاحتيالي. من الجدير بالذكر أن خوارزمية الغابة العشوائية تظهر تحسناً في قيمة Recall لعناوين URL الخاصة بالثشويه (Defacement URLs) مقارنة بشجرة القرار (96% مقابل 93%). وهذا يشير إلى أنها تكتشف بنجاح نسبة أعلى من الإيجابيات الحقيقية في هذه الفئة والسبب في ذلك هو ما يسمى بتأثير المجموعة (Ensemble Effect) حيث تكون الغابة العشوائية، باعتبارها مجموعة من أشجار القرار، أقل عرضة للملاءمة المفرطة (Overfitting) وتوفر أداءً أكثر عمومية مقارنة بشجرة قرار واحدة. بشكل عام، تظهر الغابة العشوائية نتائج واعدة، وخاصة فيما يتعلق بكشف عناوين URL المشوهة. ومع ذلك، تشترك مع خوارزمية شجرة القرار في تحدي تحديد عناوين URL التصيد الاحتيالي بدقة.

يمثل الشكل (11) مصفوفة الارتباك لخوارزمية الانحدار اللوجستي.



الشكل (11) مصفوفة الارتباك لخوارزمية الانحدار اللوجستي

الجدول (4) يوضح قيم Precision, Recall and F1\_Score لكل فئة لخوارزمية الانحدار

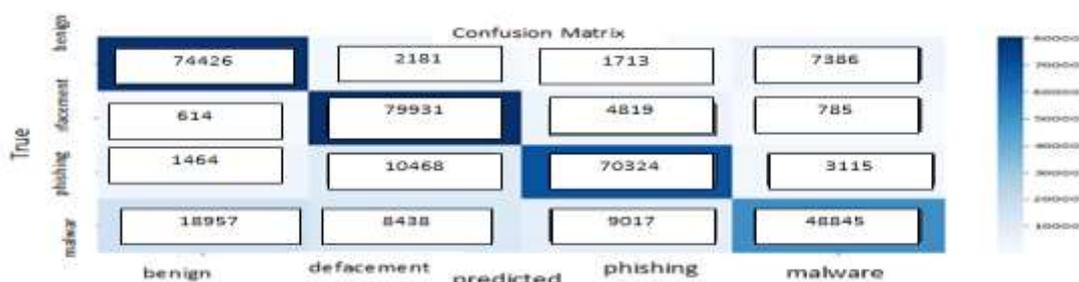
اللوجستي.

الجدول (4) قيم Precision, Recall and F1\_Score لكل فئة لخوارزمية الانحدار اللوجستي

	Precision	Recall	F1_Score
Benign URL	0.78	0.82	0.80
Defacement URL	0.82	0.91	0.86
Malware URL	0.82	0.82	0.82
Phishing URL	0.75	0.62	0.68

نلاحظ من الجدول (4) أن خوارزمية الانحدار اللوجستي تظهر قيم Precision و Recall جيدة لمعظم الفئات، ولكن على غرار الخوارزميات الأخرى، يواجه الانحدار اللوجستي صعوبة في اكتشاف عناوين URL الاحتيالية. وهذا يسلط الضوء على التحدي الكبير الذي تواجهه معظم الخوارزميات في تحديد محاولات التصيد الاحتيالي بدقة. نلاحظ كذلك ان الانحدار اللوجستي حقق قيم Precision و Recall لعناوين التصيد الاحتيالي اقل من خوارزمية الغابة العشوائية والسبب في ذلك أن حدود القرار الخطية للانحدار اللوجستي تحد من قدرته على التقاط العلاقات المعقدة بين الميزات مقارنة بنماذج المجموعة (Ensemble Models) مثل الغابة العشوائية.

يمثل الشكل (12) مصفوفة الارتباك لخوارزمية SVM.



الشكل (12) مصفوفة الارتباك لخوارزمية SVM

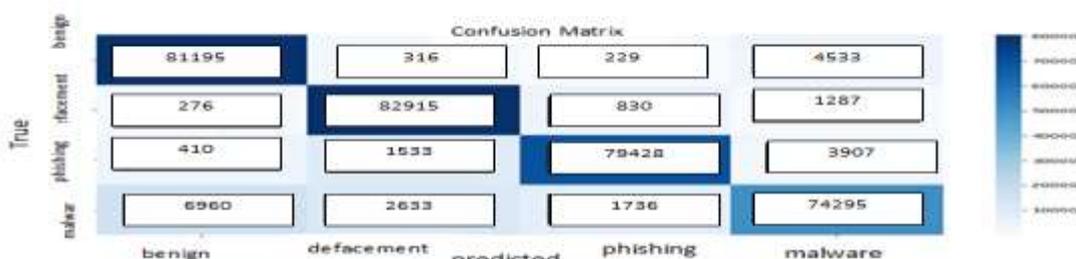
يوضح الجدول (5) قيم Precision, Recall and F1\_Score لكل فئة لخوارزمية SVM.

الجدول (5) قيم Precision, Recall and F1\_Score لكل فئة لخوارزمية SVM

	Precision	Recall	F1_Score
Benign URL	0.78	0.87	0.82
Defacement URL	0.79	0.93	0.85
Malware URL	0.82	0.82	0.82
Phishing URL	0.81	0.57	0.67

نلاحظ من الجدول (5) أن خوارزمية SVM تظهر أداء قريب من الخوارزميات الأخرى لمعظم الفئات من حيث قيم Precision و Recall. على غرار الخوارزميات الأخرى، تواجه SVM صعوبات في اكتشاف عناوين URL التصيد الاحتيالي.

يوضح الشكل (13) نتائج مصفوفة الارتباك للنموذج المقترح.



الشكل (13) نتائج مصفوفة الارتباك للنموذج المقترح

يوضح الجدول (6) قيم Precision, Recall and F1\_Score لكل فئة للنموذج المقترح.

الجدول (6) قيم Precision, Recall and F1\_Score لكل فئة للنموذج المقترح

	Precision	Recall	F1_Score
Benign URL	0.91	0.94	0.93
Defacement URL	0.95	0.97	0.96
Malware URL	0.97	0.93	0.95
Phishing URL	0.88	0.87	0.88

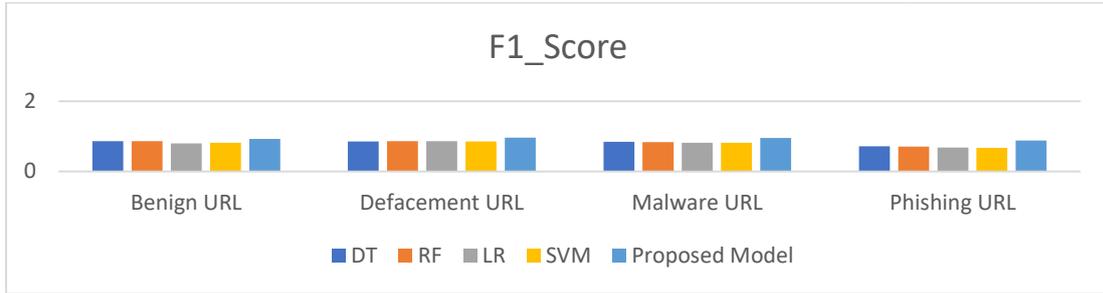
نلاحظ من الجدول (6) ما يلي:

- قيم Precision و Recall عالية: يحقق النموذج قيم Precision و Recall ممتازين لجميع الفئات، مما يشير إلى قدرة النموذج المقترح على تحديد الإيجابيات والسلبيات بدقة. وهذا أمر بالغ الأهمية عند التعامل مع مهام مثل اكتشاف عناوين URL الضارة حيث يمكن أن يكون للإيجابيات الخاطئة عواقب وخيمة.

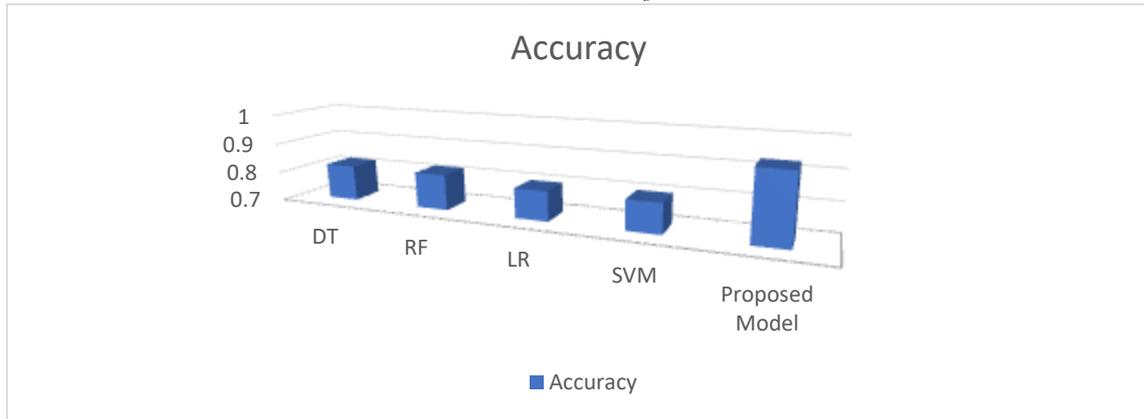
• أداء متوازن: نلاحظ أيضًا أن قيم  $F1\_Score$  عالية لجميع الفئات، مما يشير إلى أن النموذج لا يضحى بمقياس واحد من أجل الآخر. وهذا أمر بالغ الأهمية لضمان الأداء الموثوق لأنواع مختلفة من عناوين URL.

• نلاحظ أن النموذج المقترح يحقق دقة 0.88 وتذكرًا 0.87 لعناوين URL التصيد الاحتمالي، والتي واجهت الخوارزميات الأخرى صعوبات في اكتشافها. وهذا يشير إلى أن البنية المقترحة مناسبة لمعالجة مهام التصنيف المعقدة.

بمقارنة جميع النتائج السابقة وحساب قيمة الأداء (Accuracy) الإجمالية لكل خوارزمية، تم الحصول على النتائج الموضحة في الشكل (14) والشكل (15).



الشكل (14) قيم  $F1\_Score$  التي حققتها جميع الخوارزميات بالنسبة لجميع الفئات



الشكل (15) قيمة الدقة (Accuracy) الإجمالية لجميع الخوارزميات السابقة

نلاحظ من الشكل (14) والشكل (15) ما يلي:

• يتفوق النموذج المقترح على جميع الخوارزميات الأخرى من حيث قيم  $Precision$  و  $Recall$  وذلك بالنظر لقيم  $F1\_Score$ ، التي تمثل المتوسط التوافقي لكل من  $Precision$  و  $Recall$ ، مما يشير إلى فعاليته في تصنيف فئات URL المختلفة بدقة. هذه نتائج ممتازة بالنظر إلى الطبيعة الصعبة لاكتشاف عناوين URL الاحتمالية.

• شجرة القرار والغابة العشوائية: تُظهر هذه النماذج أداءً متشابهًا عبر معظم المقاييس، مع ميزة طفيفة للغابة العشوائية في اكتشاف عناوين URL الاحتمالية.

• الانحدار اللوجستي: يواجه الانحدار اللوجستي صعوبة في اكتشاف عناوين URL الاحتمالية والبرامج الضارة، مما يوضح حدوده في التعامل مع العلاقات المعقدة بين الميزات.

• SVM: يُظهر SVM أداءً جيدًا، ولكنه كذلك يواجه صعوبة في اكتشاف عناوين URL الاحتمالية.

بالنظر إلى كل فئة وبناءً على جميع النتائج السابقة يمكننا القول بأن:

- عناوين URL الحميدة والمشوهة: تحقق جميع النماذج أداءً عاليًا في هذه الفئات، مع تفوق النموذج المقترح على بقية الخوارزميات.
- عناوين URL للبرامج الضارة: يُظهر النموذج المقترح والغابة العشوائية وشجرة القرار أداءً جيدًا، في حين أن دقة SVM والانحدار اللوجستي أقل.
- عناوين URL للتصيد الاحتيالي: هذه الفئة الأكثر صعوبة في التصنيف، حيث نلاحظ أن النموذج المقترح يتفوق بشكل كبير على جميع الخوارزميات حيث حقق دقة عالية، يليه شجرة القرار والغابة العشوائية.

## 6- الاستنتاجات والتوصيات:

في هذا البحث، تمت دراسة مجموعة من الخوارزميات المختلفة والتحقق من فعاليتها في تصنيف عناوين URL الضارة إلى أربع فئات (حميدة، ومشوهة، وبرامج ضارة، وتصيد احتيالي). قمنا بتنفيذ وتقييم خمسة نماذج مختلفة: شجرة القرار، والغابة العشوائية، والانحدار اللوجستي، و SVM، الشبكة العصبية الاصطناعية المقترحة. تم تقييم كل نموذج على باستخدام مجموعة من مقاييس الأداء وهي Precision، Recall، و F1\_Score، و Accuracy. أظهرت النتائج الرئيسية لهذه الدراسة تفوق بنية الشبكة العصبية المقترحة بشكل كبير على جميع النماذج الأخرى عبر جميع المقاييس وخاصة للفئات الصعبة مثل عناوين URL للتصيد الاحتيالي. وهذا يوضح إمكانات النموذج المقترح للكشف الدقيق عن عناوين URL الضارة. حققت شجرة القرار والغابة العشوائية أداءً جيدًا لمعظم الفئات. أظهر الانحدار اللوجستي و SVM صعوبات في التعامل مع العلاقات المعقدة داخل مجموعة البيانات، مما أثر على أدائهما في الكشف. يمكن أن يركز العمل المستقبلي على عدة مجالات لتطوير بنية الشبكة العصبية الاصطناعية المقترحة وذلك من خلال الخوض بشكل أعمق في تحسين هذه البنية حيث يمكن أن يشمل ذلك استكشاف بنى الشبكات العصبية المختلفة مثل الشبكات العصبية الالتفافية (CNNs) أو الشبكات العصبية المتكررة (RNNs) أو النماذج القائمة على المحولات مثل BERT. كذلك يمكن للباحثين استكشاف الميزات الخاصة بعناوين URL بشكل أكبر والتي يمكن أن تعزز أداء النموذج للكشف عن عناوين URL الضارة. قد يتضمن ذلك استكشاف الميزات المتعلقة ببنية عنوان URL أو المحتوى أو البيانات الوصفية.

## 7- المراجع:

- [1] Kebande, V. R., & Awad, A. I. (2024). Industrial Internet of Things Ecosystems Security and Digital Forensics: Achievements, Open Challenges, and Future Directions. *ACM Computing Surveys*, 56(5), 1-37.
- [2] Gangwar, S., & Narang, V. (2022). A Survey on Emerging Cyber Crimes and Their Impact Worldwide. In *Research Anthology on Combating Cyber-Aggression and Online Negativity* (pp. 1583-1595). IGI Global.
- [3] Kumar, M., Darshan, S. S., & Yarlagadda, V. (2023). Introduction to the cyber-security landscape. In *Malware Analysis and Intrusion Detection in Cyber-Physical Systems* (pp. 1-21). IGI Global.
- [4] Ibor, A. E., Oladeji, F. A., & Okunoye, O. B. (2018). A survey of cyber security approaches for attack detection prediction and prevention. *International Journal of Security and its Applications*, 12(4), 15-28.

- [5] Telo, J. (2022). Supervised Machine Learning for Detecting Malicious URLs: An Evaluation of Different Models. *Sage Science Review of Applied Machine Learning*, 5(2), 30-46.
- [6] Shah, V. (2021). Machine Learning Algorithms for Cyber security: Detecting and Preventing Threats. *Revista Espanola de Documentacion Cientifica*, 15(4), 42-66.
- [7] Allioui, H., & Mourdi, Y. (2023). Exploring the full potentials of IoT for better financial growth and stability: A comprehensive survey. *Sensors*, 23(19), 8015.
- [8] Idrissi, I., Boukabous, M., Azizi, M., Moussaoui, O., & El Fadili, H. (2021). Toward a deep learning-based intrusion detection system for IoT against botnet attacks. *IAES International Journal of Artificial Intelligence*, 10(1), 110.
- [9] Reyes-Dorta, N., Caballero-Gil, P., & Rosa-Remedios, C. (2024). Detection of malicious URLs using machine learning. *Wireless Networks*, 1-18.
- [10] Yuan, J., Chen, G., Tian, S., & Pei, X. (2021). Malicious URL detection based on a parallel neural joint model. *IEEE Access*, 9, 9464-9472.
- [11] Chen, Y. C., Ma, Y. W., & Chen, J. L. (2020, July). Intelligent malicious URL detection with feature analysis. In *2020 IEEE Symposium on Computers and Communications (ISCC)* (pp. 1-5). IEEE.
- [12] Heryanto, A., Ab Razak, M. F., Raffei, A. F. M., Phon, D. N. E., Kasim, S., & Sutikno, T. (2020). A malicious URLs detection system using optimization and machine learning classifiers. *Indonesian Journal of Electrical Engineering and Computer Science*, 17(3), 1210-1214.
- [13] Alsaedi, M., Ghaleb, F. A., Saeed, F., Ahmad, J., & Alasli, M. (2022). Cyber threat intelligence-based malicious URL detection model using ensemble learning. *Sensors*, 22(9), 3373.
- [14] Ul Hassan, I., Ali, R. H., Ul Abideen, Z., Khan, T. A., & Kouatly, R. (2022). Significance of machine learning for detection of malicious websites on an unbalanced dataset. *Digital*, 2(4), 501-519.
- [15] Dastres, R., & Soori, M. (2021). Artificial neural network systems. *International Journal of Imaging and Robotics (IJIR)*, 21(2), 13-25.
- [16] Islam, M., Chen, G., & Jin, S. (2019). An overview of neural network. *American Journal of Neural Networks and Applications*, 5(1), 7-11.
- [17] Rastogi, V., Shao, R., Chen, Y., Pan, X., Zou, S., & Riley, R. D. (2016, February). Are these Ads Safe: Detecting Hidden Attacks through the Mobile App-Web Interfaces. In *NDSS*.
- [18] Joshi, A., Kanwar, K., & Vaidya, P. (2022, July). Attribute Selection, Sampling, and Classifier Methods to Address Class Imbalance Issues on Data Set Having Ratio Less Than Five. In *Proceedings of Third International Conference on Computing, Communications, and Cyber-Security: IC4S 2021* (pp. 261-276). Singapore: Springer Nature Singapore.
- [19] Wang, X., Wu, Z., Lian, L., & Yu, S. X. (2022). Debaised learning from naturally imbalanced pseudo-labels. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (pp. 14647-14657).
- [20] Mohammed, R., Rawashdeh, J., & Abdullah, M. (2020, April). Machine learning with oversampling and undersampling techniques: overview study and experimental results. In *2020 11th international conference on information and communication systems (ICICS)* (pp. 243-248). IEEE.

- [21] Kim, M., & Hwang, K. B. (2022). An empirical evaluation of sampling methods for the classification of imbalanced data. *PLoS One*, 17(7), e0271260.
- [22] Mohammed, A. J., Hassan, M. M., & Kadir, D. H. (2020). Improving classification performance for a novel imbalanced medical dataset using SMOTE method. *International Journal of Advanced Trends in Computer Science and Engineering*, 9(3), 3161-3172.
- [23] L'heureux, A., Grolinger, K., Elyamany, H. F., & Capretz, M. A. (2017). Machine learning with big data: Challenges and approaches. *Ieee Access*, 5, 7776-7797.
- [24] Misra, P., & Yadav, A. S. (2019, March). Impact of preprocessing methods on healthcare predictions. In *Proceedings of 2nd International Conference on Advanced Computing and Software Engineering (ICACSE)*.
- [25] Markoulidakis, I., Kopsiaftis, G., Rallis, I., & Georgoulas, I. (2021, June). Multi-class confusion matrix reduction method and its application on net promoter score classification problem. In *The 14th pervasive technologies related to assistive environments conference* (pp. 412-419).
- [26] Roth, K., Pemula, L., Zepeda, J., Schölkopf, B., Brox, T., & Gehler, P. (2022). Towards total recall in industrial anomaly detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (pp. 14318-14328).
- [27] MacEachern, S. J., & Forkert, N. D. (2021). Machine learning for precision medicine. *Genome*, 64(4), 416-425.
- [28] Fu, G., Sun, P., Zhu, W., Yang, J., Cao, Y., Yang, M. Y., & Cao, Y. (2019). A deep-learning-based approach for fast and robust steel surface defects classification. *Optics and Lasers in Engineering*, 121, 397-405.