



Using Artificial Intelligence Techniques to Detect Extremism in Social Media Content

Noor Thamer Mahmood¹, Suhad Hatim Jihad¹, Sumar Mohamed Khaleel¹
Shahad Jassim Hassan¹

¹ Computer Center, University of Babylon,
nour.thamer95@uobabylon.edu.iq

توظيف تقنيات الذكاء الاصطناعي للكشف عن التطرف في محتوى منصات التواصل الاجتماعي

نور ثامر محمود سهاد حاتم جهاد سומר محمد خليل شهد جاسم حسن

¹ مركز الحاسبة الالكترونية-جامعة بابل ، nour.thamer95@uobabylon.edu.iq ، بابل ، العراق

Received: 18/8/2025

Published: 28/8/2025

ABSTRACT

Environments, posing significant threats to both social and psychological security. This extremism appears in various forms, including violence, terrorism, and extremist crimes, with detrimental impacts on individuals, especially younger generations. Consequently, there is a pressing need to adopt evidence-based solutions involving preventive, educational, and social interventions to address this growing phenomenon. As extremism escalates, the role of technology and artificial intelligence (AI) in combating it has become more crucial, positioning it as a key focus in behavioral health and social science research. In this context, the study presents a hybrid model for detecting harmful and extremist content by analyzing URLs on social media platforms, utilizing AI and machine learning techniques. The research evaluated several algorithms, including Logistic Regression, Support Vector Machines (SVM), Decision Trees, Random Forest, and K-Nearest Neighbors (KNN). The proposed hybrid model demonstrated outstanding performance, achieving an accuracy rate of 100% in detecting harmful content—outperforming individual models, particularly Logistic Regression, which achieved an accuracy of 81.7%. These findings highlight the hybrid model's effectiveness in reducing classification errors and enhancing the detection of digital threats.

Keywords: Digital extremism, artificial intelligence, machine learning, social media, counter-extremism.



الخلاصة

تتزايد مظاهر الكراهية والتطرف العنيف في المجتمعات، خاصة في البيئات الإلكترونية، مما يشكل تهديدًا كبيرًا للأمن الاجتماعي والنفسي. يتجلى هذا التطرف في أشكال متعددة، مثل العنف والإرهاب والجرائم المتطرفة، وله تأثير سلبي واسع على الأفراد، خصوصًا الأجيال الناشئة. لذا، أصبحت الحاجة ملحة لاعتماد حلول فعالة مبنية على الأدلة العلمية، تشمل التدخلات الوقائية والتعليمية والاجتماعية لمواجهة هذه الظاهرة.

ومع تصاعد التطرف، برزت أهمية الاستفادة من التكنولوجيا والذكاء الاصطناعي في مكافحته، مما جعله مجالًا رئيسيًا في أبحاث الصحة السلوكية والعلوم الاجتماعية. في هذا السياق، يقدم البحث نموذجًا هجينًا للكشف عن المحتوى الضار والمتطرف عبر تحليل عناوين URL في وسائل التواصل الاجتماعي باستخدام تقنيات الذكاء الاصطناعي والتعلم الآلي.

شمل البحث تجريب عدد من الخوارزميات، مثل الانحدار اللوجستي، وآلات المتجهات الداعمة، وأشجار القرار، وغابة القرار العشوائية، وأقرب الجيران (KNN). وقد أظهر النموذج الهجين أداءً متميزًا، محققًا دقة بلغت 100% في الكشف عن المحتوى الضار، متفوقًا على النماذج الفردية، خصوصًا الانحدار اللوجستي الذي سجل دقة 81.7%. هذه النتائج تبرز فاعلية النهج الهجين في تقليل أخطاء التصنيف وتعزيز القدرة على رصد التهديدات الرقمية.

الكلمات المفتاحية: التطرف الرقمي، الذكاء الاصطناعي، التعلم الآلي، وسائل التواصل الاجتماعي، مكافحة التطرف.

1. المقدمة:

التطرف هو ظاهرة اجتماعية وفكرية تهدد استقرار المجتمعات، حيث يتمثل في تبني أفكار وآراء متشددة ترفض التعايش مع الرؤى المختلفة وتؤدي إلى الكراهية والعنف [1]. يأخذ التطرف أشكالًا متعددة، منها التطرف الديني، السياسي، والاجتماعي، وهو مرتبط برفض القيم الإنسانية الأساسية مثل التسامح والاحترام المتبادل [2]. غالبًا ما يؤدي التطرف إلى تهديد الأمن والسلام الاجتماعي، إذ يخلق انقسامات عميقة داخل المجتمعات، مما يهيئ الظروف لنشوء العنف المنظم والإرهاب [3].

شهدت السنوات الأخيرة تصاعدًا في حوادث العنف المرتبطة بالتطرف، لا سيما في البيئات التي تشهد استقطابًا اجتماعيًا واسعًا. على سبيل المثال، في عام 2021، أثارت حادثة قتل عائلة مسلمة في لندن، أونتايريو، صدمة واسعة، حيث ارتكبت الجريمة بدوافع كراهية دينية [4]. كما وقعت حادثة أخرى في صالون لتصفيف الشعر في تورونتو عام 2020، حيث تبين أن الجاني كان متأثرًا بحركة "العزوبية غير الطوعية" [5] (Incel) تشير هذه الحوادث إلى تزايد الأيديولوجيات المتطرفة وتأثيرها على الأفراد، خاصة في المجتمعات التي تشهد استقطابًا سياسيًا وثقافيًا متزايدًا [6].

لعبت التكنولوجيا ووسائل التواصل الاجتماعي دورًا رئيسيًا في انتشار التطرف، حيث أصبحت منصات مثل "ستورمفرونت" والمنديات الخاصة بحركات مثل "اليمين البديل" بؤرًا لتبادل الأفكار المتطرفة وتجنيد الأفراد [7]. توفر هذه المنصات بيئة خصبة لإنشاء مجتمعات افتراضية تعزز الشعور بالانتماء إلى جماعات متطرفة، مما يسهل نشر الفكر المتشدد عبر الإنترنت [8]. في هذا السياق، أصبحت منصات مثل فيسبوك وتويتر ويوتيوب تلعب دورًا مزدوجًا، حيث تساهم في نشر المحتوى المتطرف من جهة، وتحاول مكافحته عبر أدوات الذكاء الاصطناعي من جهة أخرى. ومع ذلك، فإن الخوارزميات الموصية في هذه المنصات

قد تسهم في تضخيم المحتوى المتطرف عبر استهداف المستخدمين بمحتوى يعزز توجهاتهم الفكرية المتشددة، مما يزيد من احتمالية الانزلاق نحو التطرف.[9]

في ظل هذه التحديات، ظهرت جهود بحثية وتقنية للحد من انتشار التطرف الرقمي، تشمل تطوير أنظمة ذكاء اصطناعي قادرة على اكتشاف وتصنيف المحتوى المتطرف، وتعزيز التعاون بين الحكومات وشركات التكنولوجيا لمكافحة الخطاب العنيف. لا يزال هذا المجال يشهد تطورات مستمرة، مع التركيز على إيجاد حلول توازن بين حرية التعبير وحماية الأمن الاجتماعي.

2. الاعمال البحثية ذات العلاقة

يعرض هذا القسم من الفصل أهم الأبحاث العلمية التي تناولت تحليل محتوى منصات التواصل الاجتماعي . حيث تناقش أهم التقنيات التي تم تبنيها في الدراسات الحديثة دور الذكاء الاصطناعي، وبالأخص تقنيات التعلم الآلي والتعلم العميق، في مكافحة التطرف الإلكتروني من خلال تحليل المحتوى الرقمي، وتصنيف النصوص المتطرفة، والتنبؤ بالسلوك العنيف للأفراد عبر الإنترنت. ركزت دراسة [10] على استرجاع وتحليل البيانات ذات الطابع الإرهابي والتطرفي، مع تقديم خوارزميات مستقلة عن اللغة لتحديد المخاطر المحتملة بناءً على بنية الشبكات الاجتماعية والتفاعلات بين المستخدمين، مما يتيح التعرف على الأفراد الخطرين حتى مع محدودية الوصول إلى المحتوى الذي ينشرونه. أما دراسة [11] فقدت نموذج Hist Gradient Boosting الذي حقق دقة 89.06% في تصنيف الأفراد العنيفين باستخدام بيانات من قاعدة PIRUS، كما تم استخدام SHAP و Permutation Feature Importance لتوضيح كيفية تأثير الميزات المختلفة على قرارات النموذج، مما يعزز فهم ديناميكيات التطرف العنيف.

من جهة أخرى، تناولت دراسة [12] مراجعة شاملة لآليات الكشف عن التطرف في وسائل التواصل الاجتماعي، حيث تم تحليل 64 دراسة باستخدام منهجية PRISMA، وأكدت الدراسة على الحاجة إلى تطوير أداة تلقائية متاحة للجمهور لجمع البيانات وتصنيفها وفقاً لمنهجيات تحقق دقيقة، بهدف تحسين دقة تحليل النصوص المتطرفة. في حين ركزت دراستهم الأخرى [13] على تطوير مجموعة بيانات MIWS (Merged ISIS/Jihadist-White Supremacist)، وهي مجموعة بيانات متوازنة متعددة الفئات تستهدف تصنيف النصوص المتطرفة إلى فئات مثل الدعاية، والتجنيد، والتطرف. تم اختبار بيانات MIWS باستخدام نماذج التعلم العميق BERT و RoBERTa و DistilBERT، حيث حقق BERT أعلى دقة (F1-score = 0.72)، مما يثبت فعالية هذه النماذج في تصنيف المحتوى المتطرف عبر أيديولوجيات مختلفة. على صعيد مراجعات الأدبيات، أشارت [14] إلى تصنيف الأبحاث في مجال كشف التطرف الإلكتروني ضمن ثلاث فئات رئيسية: التحليل، والكشف، والتنبؤ، مع إبراز تحديات مثل نقص مجموعات البيانات الموثوقة، وضعف التعاون البحثي، وتطور الخطاب المتطرف. أما دراسة [15] فقد قدمت تصنيفاً شاملاً لأبحاث خطاب الكراهية، مشيرةً إلى نقص البيانات العامة المستخدمة في هذه الأبحاث، بينما ركزت دراسة [16] AI- على كشف خطاب الكراهية في السياق العربي، مع مناقشة التحديات المرتبطة بتحليل النصوص باللغة العربية. وأخيراً، قدمت دراسة [17] نموذجاً للكشف عن الهجمات الاجتماعية المستندة إلى عناوين URL، والتي تُستخدم في الهندسة الاجتماعية لنشر التطرف. تم اختبار نماذج تعلم عميق تشمل LSTM و CNN و CharacterBERT، حيث حقق الأخير أعلى دقة وصلت إلى 99.65%، مما يبرز إمكانياته في كشف الهجمات الإلكترونية المرتبطة بالمحتوى المتطرف.

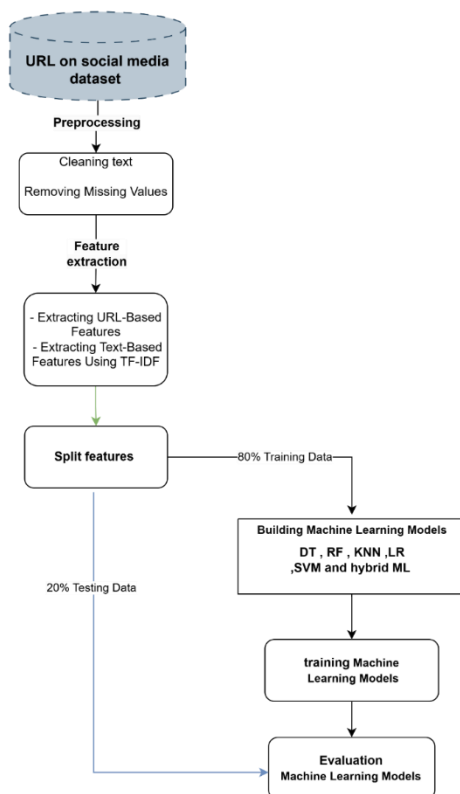


توضح هذه المراجعة أن الدراسات السابقة تفتقر إلى تحليل مفصل للمجموعات البيانية، والمناهج التقنية، والمقاييس التقييمية، وأساليب التحقق من البيانات، مما يشير إلى فجوة بحثية تتعلق بتطوير مناهج تصنيف متقدمة وأدوات موثوقة لرصد التطرف عبر الإنترنت. لذلك، تفتح هذه الدراسة المجال لأبحاث مستقبلية تستهدف تحليل شامل للبيانات، وتطوير تقنيات كشف وتصنيف أكثر دقة، وتقديم أساليب تحقق محسنة، وتصميم أدوات متاحة للعامة لرصد التطرف الإلكتروني بشكل فعال.

3. المنهجية المقترحة

يستعرض هذا الفصل المراحل الأساسية للجوانب العملية والمنهجية المعتمدة في بناء النظام المقترح، والذي يمر بعدة مراحل أساسية لتمكينه من تنفيذ مهمة تحليل النصوص وتصنيف محتوى منصات التواصل الاجتماعي بكفاءة عالية. تعتمد منهجية هذا البحث على خطوات جوهرية لضمان بناء نظام فعال ودقيق، وتتمثل في:

1. جمع مجموعة البيانات اللازمة لتدريب واختبار نموذج التعلم الآلي.
 2. تنفيذ عمليات المعالجة الأولية للبيانات لتجهيزها للمرحلة التالية.
 3. تحويل البيانات النصية إلى تمثيل عددي (Form Vector) باستخدام تقنية (TF-IDF)
 4. تقسيم البيانات إلى مجموعتي تدريب واختبار لضمان دقة النموذج وقياس أدائه.
- بعد تنفيذ هذه الخطوات، يتم تدريب نماذج التعلم الآلي باستخدام بيانات التدريب، ومن ثم اختبار قدرتها على التنبؤ بالمحتوى الضار في بيانات الاختبار. يوضح الشكل (3.1) المخطط الانسيابي المقترح للهيكل العام للنظام، والذي يبرز تدفق العمليات من مرحلة إدخال البيانات إلى تصنيف المحتوى بكفاءة وذكاء.



الشكل 3.1: النظام المقترح للتعليم الآلي

3.1 مجموعة البيانات

تم إنشاء مجموعة بيانات الروابط الضارة مقابل الروابط الآمنة لمعالجة مشكلة اكتشاف الروابط الخبيثة، حيث تم جمع البيانات من مصادر متعددة على الإنترنت ومنصات التواصل الاجتماعي، مثل PhisTank. تأتي البيانات على شكل ملف بصيغة CSV وتحتوي على سجلات لعناوين URL مصنفة إما كآمنة (0) أو ضارة (1). تتألف المجموعة من 450,176 رابطاً، حيث تمثل الروابط الآمنة 77%، بينما تشكل الروابط الضارة 23%. تتضمن البيانات أربعة أعمدة رئيسية:

- **URL:** يحتوي على الرابط الكامل.
- **Label:** يحدد ما إذا كان الرابط ضاراً أم لا.
- **Score:** يُعطى القيمة 1 للروابط الضارة و 0 للروابط الآمنة.

توفر هذه البيانات مصدراً مهماً لدراسة خصائص الروابط الضارة، وتطوير خوارزميات التصنيف، وتحليل أنماط التهديدات الإلكترونية لتعزيز الأمن السيبراني. كما تساهم في عزل المحتوى الضار على منصات التواصل الاجتماعي، خاصة المحتوى المستخدم في عمليات التصيد الاحتيالي والاحتيال الإلكتروني.

مجموعة البيانات متاحة على موقع Kaggle عبر الرابط التالي:

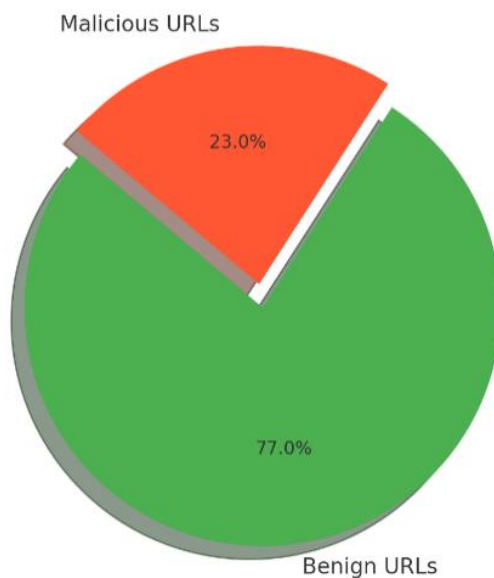
Malicious and Benign URLs Dataset

يوضح الجدول 1 والأشكال (2-3) محتوى هذه البيانات بشيء من التفصيل.

الجدول 1: وصف مجموعة البيانات

الوصف	القيمة
إجمالي الروابط	450,176
الروابط الآمنة	77% تقريبًا 346,635
الروابط الضارة	23% تقريبًا 103,541

Distribution of Benign and Malicious URLs



الشكل (2): توزيع عناوين URL في مجموعة البيانات.

url_df - DataFrame

Index	Unnamed: 0	url	label	result
0	0	https://www.google.com	benign	0
1	1	https://www.youtube.com	benign	0
2	2	https://www.facebook.com	benign	0
3	3	https://www.baidu.com	benign	0
4	4	https://www.wikipedia.org	benign	0
5	5	https://www.reddit.com	benign	0
6	6	https://www.yahoo.com	benign	0
7	7	https://www.google.co.in	benign	0
8	8	https://www.qq.com	benign	0
9	9	https://www.amazon.com	benign	0
10	10	https://www.taobao.com	benign	0

الشكل 3: عينة من مجموعة البيانات ووصف الميزات.

3.2 المعالجة المسبقة واستخراج الميزات

قبل تدريب نماذج التعلم الآلي، تمر مجموعة البيانات بسلسلة من العمليات الأساسية لضمان دقة التصنيف وكفاءة النموذج. تتضمن هذه العمليات المعالجة المسبقة، واستخراج الميزات، وإعداد البيانات للتدريب والاختبار، مما يعزز قدرة النموذج على التعرف على الروابط الضارة بفعالية.

١. المعالجة المسبقة للبيانات:

- تبدأ هذه المرحلة بتحميل البيانات وفحصها للتأكد من خلوها من المشكلات التي قد تؤثر على أداء النموذج. تشمل الخطوات الأساسية:
- التأكد من عدم وجود قيم مفقودة: يتم التحقق من أي بيانات غير مكتملة أو فارغة وإزالتها أو استبدالها عند الضرورة.
- تنقية البيانات: تتضمن إزالة أي رموز أو علامات غير ضرورية قد تؤثر على تحليل النصوص داخل الروابط.
- تحويل جميع الروابط إلى تنسيق موحد: مثل تحويل الأحرف الكبيرة إلى صغيرة لتجنب التكرارات غير الضرورية.

٢. استخراج الميزات: (Feature Extraction)

لتحليل الروابط وتصنيفها بدقة، يتم تحويل كل رابط إلى تمثيلات رقمية تستند إلى خصائص معينة، مما يساعد النموذج في التمييز بين الروابط الضارة والأمنة. تشمل الميزات المستخرجة:

التحليل الهيكلي لعناوين URL:

- التحقق من وجود عنوان IP (having_ip_address): يتم فحص ما إذا كان الرابط يحتوي على عنوان IP بدلاً من اسم نطاق، حيث إن الروابط التي تستخدم عناوين IP بدلاً من أسماء النطاقات غالباً ما تكون ضارة.
- عدد النقاط: (count_dot) يتم حساب عدد النقاط في الرابط، حيث إن الروابط الضارة تميل إلى استخدام عدد أكبر من النقاط لتضليل المستخدمين.
- التحقق من وجود (count_www) "www" و (count_com) ".com". يمكن أن تشير بعض الأنماط غير العادية إلى روابط خادعة.
- عدد الأرقام: (count_digits) الروابط التي تحتوي على أعداد كبيرة من الأرقام قد تكون مشبوهة، حيث تُستخدم هذه التقنية في التصيد الاحتيالي.
- طول الرابط: (url_length) يتم قياس طول الرابط، حيث إن الروابط الضارة تميل إلى أن تكون طويلة جداً لإخفاء نواياها الحقيقية.
- استخدام خدمات اختصار الروابط: (shortening_service) يتم التحقق مما إذا كان الرابط يستخدم خدمة اختصار، حيث تُستخدم هذه الخدمات غالباً لإخفاء الروابط الضارة.
- تحليل الرموز الخاصة داخل الرابط:
يتم احتساب عدد مرات ظهور الأحرف الخاصة مثل \$, # , + , ~ , ! , @ , & , = , ? , / , _ , - , % ، حيث إن الروابط الخبيثة غالباً ما تحتوي على أنماط معينة من هذه الرموز في محاولة لإخفاء طبيعتها الحقيقية أو لجذب المستخدمين للنقر عليها.

تحليل النصوص باستخدام TF-IDF:

يتم استخدام تقنية تحليل تردد الكلمات المعكوس (TF-IDF) لاستخراج الأنماط النصية المهمة من عناوين الروابط، حيث يساعد TfidfVectorizer في تحويل الروابط إلى بيانات رقمية بناءً على الكلمات الأكثر شيوعاً، مما يمكن النموذج من التعرف على المصطلحات المستخدمة في الروابط الخبيثة مقارنةً بالروابط الآمنة.

٣. إعداد البيانات للتدريب والاختبار:

- بعد استخراج جميع الميزات، يتم تقسيم البيانات إلى مجموعتين:
 - مجموعة التدريب: (80%) تُستخدم لتدريب نموذج التعلم الآلي على التعرف على الروابط الضارة.
 - مجموعة الاختبار: (20%) تُستخدم لاختبار دقة النموذج في تصنيف الروابط غير المرئية سابقاً.
- يضمن هذا الإعداد أن يكون النموذج قادراً على التعميم، مما يقلل من احتمالية حدوث فرط التكيف (Overfitting) ويزيد من قدرته على اكتشاف الروابط الضارة بدقة في البيانات الحقيقية.



3.3 التعلم الآلي

يعتمد النظام المقترح على أفضل خمسة خوارزميات تعلم آلي من حيث الدقة وسرعة التنفيذ، مما يضمن تحقيق تصنيف دقيق للروابط المنشورة على منصات التواصل الاجتماعي. في هذه الدراسة، تم إجراء تجارب باستخدام خوارزميات متقدمة تشمل:

- أقرب الجيران (K-Nearest Neighbors – KNN)
- أشجار القرار (Decision Trees)
- الغابات العشوائية (Random Forests)
- الانحدار اللوجستي (Logistic Regression)
- مصنف متجهات الدعم (Support Vector Classifier – SVC)

تم استخدام هذه الخوارزميات لتدريب النموذج واختبار كفاءته في تصنيف الروابط، مما يسمح بتحليل الأنماط المختلفة والتمييز بين الروابط الضارة والأمنة بفعالية.

نحو نموذج هجين أكثر دقة وكفاءة

لتحقيق تحسينات كبيرة في الأداء، تم تطوير نموذج هجين يجمع بين هذه الخوارزميات الخمس، حيث يستفيد من نقاط القوة في كل خوارزمية ويقلل من تأثير نقاط الضعف، مما يؤدي إلى زيادة دقة التصنيف مقارنة باستخدام كل خوارزمية على حدة. يعتمد النموذج على تحليل شامل للروابط المنشورة في محتوى منصات التواصل الاجتماعي وتصنيفها وفقاً لعدة معايير، من أبرزها:

1. التحليل الهيكلي للروابط (Structural Analysis): يتم فحص بنية الروابط والتعرف على الأنماط المشبوهة التي قد تشير إلى محاولات تصيد أو احتيال.
 2. تحليل ميزات النص باستخدام TF-IDF: يتم استخراج وتحليل الكلمات المفتاحية الأكثر شيوعاً داخل الروابط لتحديد ما إذا كانت تحمل دلالات تشير إلى مخاطر أمنية أو لا.
- يساهم هذا النهج المتكامل في بناء نظام أكثر موثوقية وفعالية، مما يعزز قدرات الذكاء الاصطناعي في التصدي للمحتوى الضار والروابط المشبوهة على منصات التواصل الاجتماعي.

**Algorithm(3.1) : The Proposed System algorithm**

Input : Data (Dataset) consist of most relevant features.

ML– Algo:K Nearest Neighbors', 'Decision Tree', 'Random Forest', 'Logistic Regression', 'Support Vector Classifier' and hybrid ML

Output : malicious url (M) , Normal (N)

Begin {

Step 1: Input Given: Read Dataset (Data)

Step 2: Data Pre-processing

Step 3: Features extraction

Step 4: split s-df for 80 for training (x_sdf) and 20 for testing(t_sdf)

Step 5: creat training ML model

Step 6: Training using (ML– Algo) (s_df,b)

Step 7: Testing model.

Step 8: Evaluation the result Step 9 Return N or M

}

End Algorithm

توضح الخوارزمية (1) المنهجية المعتمدة في هذه الأطروحة لتصنيف المحتوى، حيث تتضمن الخطوات التفصيلية لمعالجة البيانات، استخراج الميزات، وتدريب النموذج باستخدام تقنيات التعلم الآلي. تبدأ العملية بجمع وتحليل البيانات من منصات التواصل الاجتماعي، يليها تطبيق تقنيات المعالجة المسبقة لإزالة الضوضاء وتحسين جودة البيانات. بعد ذلك، يتم استخراج الميزات الهيكلية والنصية باستخدام تقنيات مثل تحليل TF-IDF. وأخيراً، يتم تدريب النموذج باستخدام مجموعة من خوارزميات التعلم الآلي لضمان دقة تصنيف المحتوى الضار والحميد، مما يعزز من كفاءة النظام في الكشف عن التهديدات المحتملة.

4 النتائج ومناقشتها

بعد تدريب النموذج المقترح باستخدام 80% من البيانات المتاحة، يتم اختبار النموذج على الـ 20% المتبقية من البيانات التي لم تُستخدم في التدريب. يتم تقسيم هذه العملية إلى مراحل يتم خلالها تحليل أداء النموذج بشكل دقيق: المرحلة الأولى : يتم فيها اختبار الخوارزميات المختلفة للتعلم الآلي التي تم اختيارها. يشمل ذلك مراقبة سلوك كل خوارزمية في معالجة البيانات، حيث يتم تقييم استجابة الخوارزميات للبيانات التي لم تُرها أثناء التدريب. تُجمع البيانات الناتجة من هذه الخوارزميات في هذه المرحلة لفهم مدى كفاءتها في التنبؤ بالنتائج الصحيحة.

المرحلة الثانية: يتم فيها تحليل أداء الخوارزميات بشكل أعمق باستخدام مقاييس التقييم المختلفة. يتم تقييم كل خوارزمية بناءً على عدة معايير:

- الدقة (Accuracy): مقياس يُستخدم لقياس النسبة المئوية للتنبؤات الصحيحة مقارنة بالتنبؤات الإجمالية.
- الاسترجاع (Recall): يُظهر قدرة الخوارزمية على اكتشاف العناصر الإيجابية الحقيقية، ويُستخدم بشكل خاص عندما تكون البيانات غير متوازنة.
- درجة: (F1 Score) مزيج بين الاسترجاع والدقة لتوفير مقياس متوازن، خاصة عندما يكون هناك اهتمام متساوٍ بالدقة والاسترجاع.
- الدقة (Precision): مقياس يُظهر مدى دقة التنبؤات الإيجابية التي تم إجراؤها من قبل الخوارزمية.
- مصفوفة الارتباك (Confusion Matrix): أداة تستخدم لفحص كيفية تصنيف الخوارزمية للبيانات إلى فئات صحيحة وخاطئة، مما يساعد في تحديد الأخطاء النوعية.

4.1 محاكاة خوارزميات التعلم الآلي

يقدم هذا القسم النتائج التجريبية للنماذج المقترحة للتنبؤ بعنوان URL الضارة في المحتوى على منصات وسائل التواصل الاجتماعي. أظهرت نتائج المحاكاة الموضحة بالجدول 1 والأشكال من 3-15 تفوقاً واضحاً لخوارزميات التعلم الآلي في تصنيف عناوين URL الضارة عبر منصات وسائل التواصل الاجتماعي. حيث حققت خوارزميات مثل شجرة القرار، والغابة العشوائية، وآلة الدعم الناقل (SVM)، وأقرب الجيران (KNN) دقة تصنيف تصل إلى 100%، مما يعني أنه تم تصنيف جميع العينات بشكل صحيح، سواء كانت غير ضارة أو ضارة. هذه النتائج تدل على قدرة هذه النماذج العالية في التمييز بين المحتوى الضار وغير الضار بشكل مثالي، وهو أمر بالغ الأهمية في مجال الأمن السيبراني حيث تتطلب منصات التواصل الاجتماعي الكشف المبكر عن المحتوى الضار. أما بالنسبة لخوارزمية الانحدار اللوجستي، فقد أظهرت أداءً أقل بكثير، حيث حققت دقة 81.7%، مما يشير إلى أن النموذج كان أقل قدرة في تصنيف المحتوى الضار بشكل دقيق. تحديدًا، كان لديه مشكلة في تصنيف عينات من الفئة الضارة بشكل صحيح، حيث تم تصنيف العديد من العينات الضارة على أنها غير ضارة، مما يعكس خطورة وجود تهديدات أمنية فعلية غير مكتشفة.

على الجانب الآخر، أثبت النموذج الهجين الذي يجمع بين خمسة خوارزميات (الانحدار اللوجستي، الغابة العشوائية، SVM، KNN، وشجرة القرار) فعاليته الكبيرة. حيث تم تحقيق دقة تصنيف 100%، مما يعكس قدرته على دمج نقاط القوة في كل خوارزمية فردية وتحسين الأداء العام. النموذج الهجين أظهر تحسناً كبيراً مقارنة بالنماذج الفردية، حيث لم يتم تصنيف أي عينات بشكل خاطئ، سواء كانت ضارة أو غير ضارة، مما يضمن تمييزاً دقيقاً بين المحتوى الموثوق وغير الموثوق.



Training Decision Tree...

Accuracy for Decision Tree: 1.0

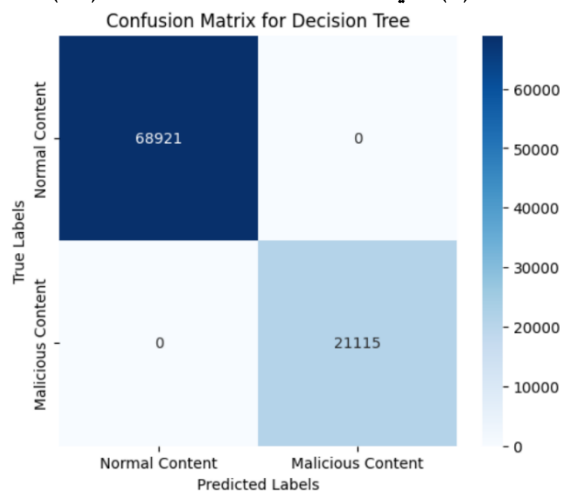
Classification Report for Decision Tree:

	precision	recall	f1-score	support
benign	1.00	1.00	1.00	68921
malicious	1.00	1.00	1.00	21115
accuracy			1.00	90036
macro avg	1.00	1.00	1.00	90036
weighted avg	1.00	1.00	1.00	90036

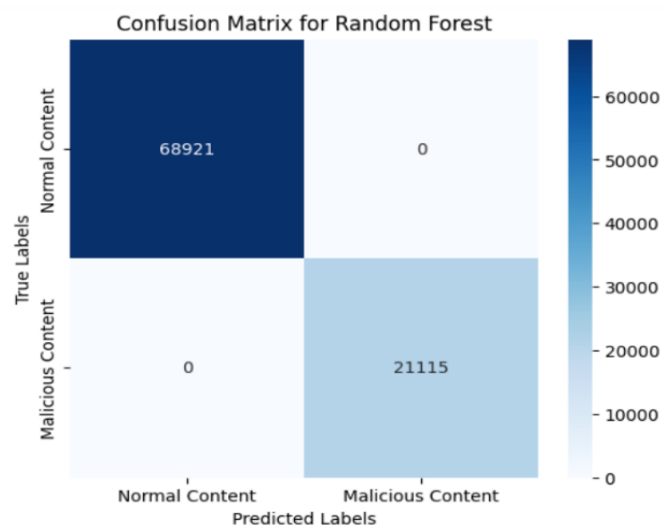
Confusion Matrix for Decision Tree

```
[[68921  0]
 [  0 21115]]
```

الشكل (3) تقرير التصنيف لخوارزمية شجرة القرار. (DT)



الشكل (4) مصفوفة الارتباك لخوارزمية شجرة القرار. (DT)



الشكل (6) مصفوفة الارتباك لخوارزمية الغابات العشوائية. (RF)

Training Random Forest...

Accuracy for Random Forest: 1.0

Classification Report for Random Forest:

	precision	recall	f1-score	support
benign	1.00	1.00	1.00	68921
malicious	1.00	1.00	1.00	21115
accuracy			1.00	90036
macro avg	1.00	1.00	1.00	90036
weighted avg	1.00	1.00	1.00	90036

Confusion Matrix for Random Forest

```
[[68921  0]
 [  0 21115]]
```

الشكل (7) تقرير التصنيف لخوارزمية الغابات العشوائية. (RF)

Training SVM...

Accuracy for SVM: 1.0

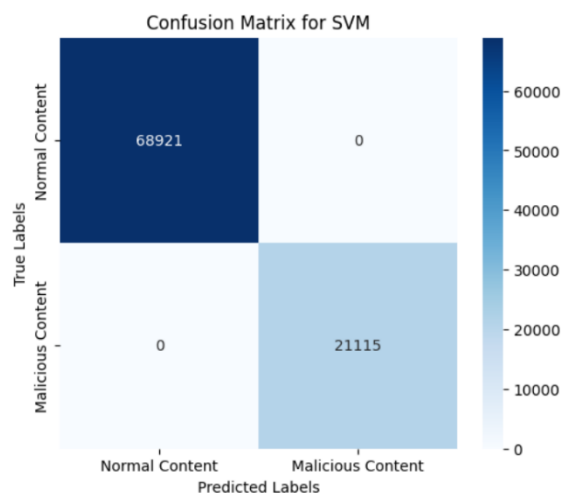
Classification Report for SVM:

	precision	recall	f1-score	support
benign	1.00	1.00	1.00	68921
malicious	1.00	1.00	1.00	21115
accuracy			1.00	90036
macro avg	1.00	1.00	1.00	90036
weighted avg	1.00	1.00	1.00	90036

Confusion Matrix for SVM

```
[[68921  0]
 [  0 21115]]
```

الشكل (8) تقرير التصنيف لخوارزمية آلة المتجهات الداعمة. (SVM)



الشكل (9) مصفوفة الارتباك لخوارزمية آلة المتجهات الداعمة. (SVM)

Training KNN...

Accuracy for KNN: 0.9999888933315563

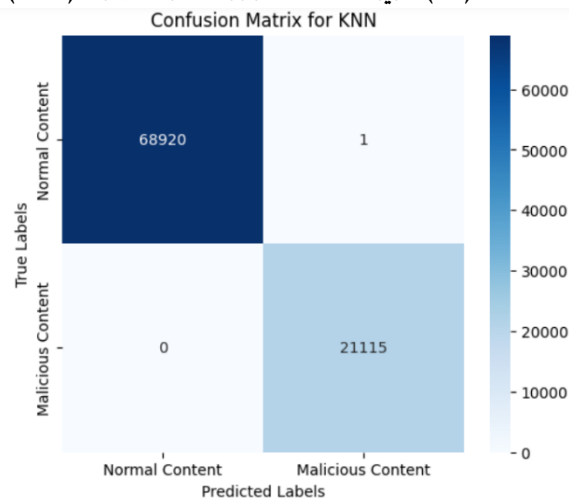
Classification Report for KNN:

	precision	recall	f1-score	support
benign	1.00	1.00	1.00	68921
malicious	1.00	1.00	1.00	21115
accuracy			1.00	90036
macro avg	1.00	1.00	1.00	90036
weighted avg	1.00	1.00	1.00	90036

Confusion Matrix for KNN

[[68920 1]
[0 21115]]

الشكل (10) تقرير التصنيف لخوارزمية أقرب الجيران. (KNN)



الشكل (11) مصفوفة الارتباك لخوارزمية أقرب الجيران. (KNN)

Training Logistic Regression...

Accuracy for Logistic Regression: 0.8170065307210449

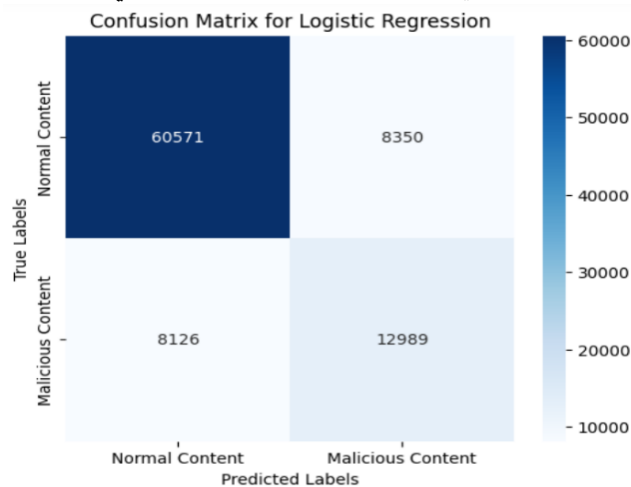
Classification Report for Logistic Regression:

	precision	recall	f1-score	support
benign	0.88	0.88	0.88	68921
malicious	0.61	0.62	0.61	21115
accuracy			0.82	90036
macro avg	0.75	0.75	0.75	90036
weighted avg	0.82	0.82	0.82	90036

Confusion Matrix for Logistic Regression

[[60571 8350]
[8126 12989]]

الشكل (12) تقرير التصنيف لخوارزمية الانحدار اللوجستي (LR)



الشكل (13) مصفوفة الارتباك لخوارزمية الانحدار اللوجستي..

4.2 محاكاة التعلم الآلي الهجين

الاختبار الثاني في هذا البحث يهدف إلى تصميم نموذج هجين يجمع بين خمسة خوارزميات تصنيف مختلفة، وهي: الانحدار اللوجستي، والغابة العشوائية، وآلة المتجهات الداعمة (SVM)، وأقرب الجيران (KNN)، وشجرة القرار. يعتمد هذا النموذج على مصنف التصويت (VotingClassifier) باستخدام التصويت الصارم (Hard Voting)، حيث يتم اتخاذ القرار النهائي بناءً على غالبية توقعات النماذج الفردية. تهدف هذه المقاربة إلى تحسين دقة التصنيف والاستفادة من نقاط القوة في كل خوارزمية، مما يعزز القدرة على اكتشاف المحتوى غير الموثوق به بشكل أكثر فاعلية مقارنةً باستخدام نموذج واحد فقط. تعكس نتائج النموذج التجميعي (Ensemble Model) أداءً مثالياً، حيث حقق دقة تصنيف بلغت 100% في تصنيف المحتوى ضمن فئتيه: الحميد والخبيث.

Training Ensemble Model...

Accuracy for Ensemble Model: 1.0

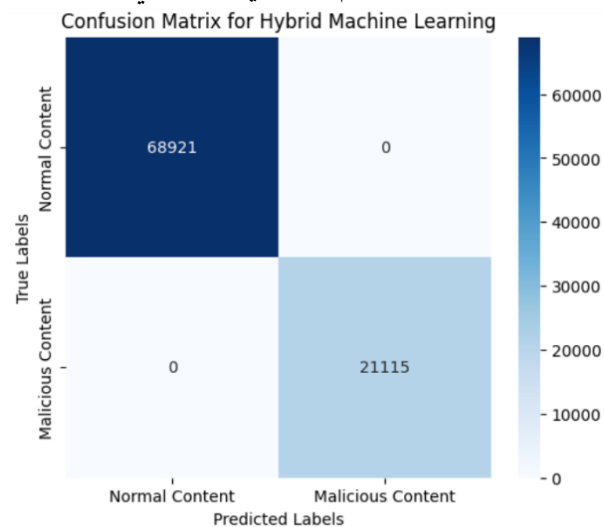
Classification Report for Ensemble Model:

	precision	recall	f1-score	support
benign	1.00	1.00	1.00	68921
malicious	1.00	1.00	1.00	21115
accuracy			1.00	90036
macro avg	1.00	1.00	1.00	90036
weighted avg	1.00	1.00	1.00	90036

```
[[68921  0]
 [  0 21115]]
```

الشكل (14): تقرير التصنيف لخوارزمية الانحدار اللوجستي

يوضح التقرير في الشكل (14) أن جميع المقاييس الأساسية، بما في ذلك الدقة (Precision)، والاسترجاع (Recall)، ومعامل F1-score، قد سجلت 1.00، مما يشير إلى عدم وجود أي أخطاء في التصنيف



الشكل (15) مصفوفة الارتباك لخوارزمية التعلم الآلي الهجين.

تؤكد مصفوفة الارتباك في الشكل 15 أيضاً عدم وجود أي أخطاء في التصنيف، حيث تم تصنيف 68,921 عينة من المحتوى الموثوق به و 21,115 عينة من المحتوى الخبيث بشكل صحيح. يشير هذا الأداء إلى أن النموذج الهجين قد نجح في دمج ميزات النماذج الفردية بطريقة فعالة، مما مكّنه من التمييز بين الفئتين بدقة عالية.

جدول 1. نتائج محاكاة خوارزميات التعلم الآلي لتصنيف عناوين URL الضارة في محتوى وسائل التواصل الاجتماعي

الخوارزمية	الدقة	الاسترجاع	F1-Score
شجرة القرار (DT)	100%	100%	100%
الغابة العشوائية (RF)	100%	100%	100%
آلة الدعم الناقل (SVM)	100%	100%	100%
أقرب الجيران (KNN)	99.99%	100%	1.00
الانحدار اللوجستي (LR)	81.7%	0.88 غير ضار	0.62 ضار
التعلم الآلي الهجين	100%	100%	100%

5 الاستنتاج والاعمال المستقبلية

أظهرت نتائج الدراسة فعالية النموذج التجميعي الهجين في الكشف عن المحتوى الضار بدقة مثالية بلغت 100%، مما يجعله أداة قوية للكشف عن التهديدات الأمنية في منصات التواصل الاجتماعي. على الرغم من أن بعض النماذج الفردية أظهرت بعض القيود، فإن النموذج الهجين جمع بين مزايا خوارزميات متعددة مثل الغابة العشوائية وآلة المتجهات الداعمة (SVM) وأقرب الجيران (KNN)، مما ساعد على تحسين الأداء وتقليل الأخطاء. في المستقبل، يمكن تعزيز النموذج من خلال اختبار أنواع جديدة من الهجمات السيبرانية وتطبيق تقنيات التعلم العميق مثل نماذج BERT وRoBERTa، بالإضافة إلى تحسين كفاءة التنفيذ باستخدام تقنيات ضغط النماذج. كما يمكن دمج النموذج مع أنظمة الأمن السيبراني الفعلية وتطوير حلول لامركزية باستخدام تقنيات البلوكشين لزيادة موثوقيته في التصدي للتهديدات السيبرانية المتطورة.

تتمثل أهمية هذا البحث في توفير أداة قوية للكشف المبكر عن المحتوى الضار على منصات التواصل الاجتماعي، وهو أمر حيوي في مكافحة التطرف الرقمي. بالنظر إلى زيادة انتشار الأفكار المتطرفة عبر الإنترنت، يمكن أن يساعد النموذج التجميعي الهجين في التعرف على المحتوى الضار أو المتطرف بدقة عالية، مما يساهم في تقليل انتشار هذه الأفكار الضارة في المجتمع. من خلال استخدام خوارزميات مثل شجرة القرار، وآلة المتجهات الداعمة (SVM)، والغابة العشوائية، يتمكن النموذج من التمييز بفعالية بين المحتوى المضر والمحتوى غير الضار، مما يساهم في حماية المستخدمين من التعرض للأيديولوجيات المتطرفة. إضافة إلى ذلك، فإن التحديث المستمر والتكيف مع أنواع جديدة من المحتوى الضار يعزز من قدرة النموذج على مواجهة تحديات التطرف الرقمي في المستقبل، مما يساهم في خلق بيئة أكثر أماناً وصحة على منصات التواصل الاجتماعي.



Conflict of interest.

There is no conflict of interest

References

- [1] W. Waheeb and R. Ghazali, "Content-based SMS Classification: Statistical Analysis for the Relationship between Number of Features and Classification Performance," *Computacion y Sistemas*, vol. 21, pp. 771–785, 2017.
- [2] A. Tekerek, "Support vector machine based spam SMS detection," *Journal of Polytechnic*, vol. 0900, pp. 779–784, 2018.
- [3] P. Poomka, W. Pongsena, N. Kerdprasop, and K. Kerdprasop, "SMS Spam Detection Based on Long Short-Term Memory and Gated Recurrent Unit," *International Journal of Future Computer and Communication*, vol. 8, pp. 12–15, 2019.
- [4] P. K. Roy, J. P. Singh, and S. Banerjee, "Deep learning to filter SMS Spam," *Future Generation Computer Systems*, vol. 102, pp. 524–533, 2020.
- [5] B. Serkan and K. Onur, "Development of content based SMS classification application by using Word2Vec based feature extraction," *IET Software*, vol. 13, pp. 295–304, 2018.
- [6] A. Barushka and P. Hajek, "Spam filtering using integrated distribution-based balancing approach and regularized deep neural networks," *Applied Intelligence*, vol. 48, pp. 3538–3556, 2018.
- [7] The Apache SpamAssassin Group, "The First Enterprise Open-Source Spam Filter," [Online]. Available: <http://spamassassin.apache.org/>. [Accessed: Jun. 2, 2020].
- [8] D. Ruano-Ordás, J. Fdez-Glez, F. Fdez-Riverola, and J. R. Méndez, "Effective scheduling strategies for boosting performance on rule-based spam filtering frameworks," *Journal of Systems and Software*, vol. 86, pp. 3151–3161, 2013.
- [9] M. Fernandez and H. Alani, "Artificial intelligence and online extremism: Challenges and opportunities," in *Predictive Policing and Artificial Intelligence*, pp. 132–162, 2021.
- [10] I. V. Mashechkin, M. I. Petrovskiy, D. V. Tsarev, and M. N. Chikunov, "Machine learning methods for detecting and monitoring extremist information on the internet," *Programming and Computer Software*, vol. 45, pp. 99–115, 2019.
- [11] A. Rösner et al., "Machine learning-based classification of extremism using explainable artificial intelligence," in *Proc. 2024 IEEE International Conference on Intelligent Systems (IS)*, pp. 1–7, Aug. 2024.
- [12] M. Gaikwad, S. Ahirrao, S. Phansalkar, and K. Kotecha, "Online extremism detection: A systematic literature review with emphasis on datasets, classification techniques, validation methods, and tools," *IEEE Access*, vol. 9, pp. 48364–48400, 2021.
- [13] M. Gaikwad, S. Ahirrao, S. Phansalkar, K. Kotecha, and S. Rani, "Multi-Ideology, Multiclass Online Extremism Dataset, and Its Evaluation Using Machine Learning," *Computational Intelligence and Neuroscience*, vol. 2023, no. 1, Art. no. 4563145, 2023.
- [14] M. Fernandez and H. Alani, "Artificial intelligence and online extremism: Challenges and opportunities," in *Predictive Policing and Artificial Intelligence*, J. McDaniel and K. Pease, Eds. New York, NY, USA: Taylor & Francis, 2020.



- [15] P. Fortuna and S. Nunes, "A survey on automatic detection of hate speech in text," *ACM Computing Surveys*, vol. 51, no. 4, pp. 1–30, Sep. 2018, doi: 10.1145/3232676.
- [16] A. Al-Hassan and H. Al-Dossari, "Detection of hate speech in social networks: A survey on multilingual corpus," in *Proc. 6th International Conference on Computer Science and Information Technology*, pp. 1–18, 2019, doi: 10.5121/csit.2019.90208.
- [17] A. S. Rafsanjani et al., "Enhancing malicious URL detection: A novel framework leveraging priority coefficient and feature evaluation," *IEEE Access*, 2024.
- [18] R. Ghanem, H. Erbay, and K. Bakour, "Contents-based spam detection on social networks using RoBERTa embedding and stacked BLSTM," *SN Computer Science*, vol. 4, no. 4, p. 380, 2023.
- [19] Á. Huertas-García, A. Martín, J. Huertas-Tato, and D. Camacho, "Countering malicious content moderation evasion in online social networks: Simulation and detection of word camouflage," *Applied Soft Computing*, vol. 145, Art. no. 110552, 2023.
- [20] Y. J. Chang, K. L. Tsai, W. C. Jiang, and M. K. Liu, "Content-aware malicious webpage detection using convolutional neural network," *Multimedia Tools and Applications*, vol. 83, no. 3, pp. 8145–8163, 2024.
- [21] F. Alkhudair et al., "Detecting malicious URL," in *Proc. 2020 International Conference on Computing and Information Technology (ICCIT-1441)*, pp. 1–5, Sep. 2020.
- [22] V. Abhijith et al., "Detection of Malicious URLs in Twitter," in *Proc. 2021 International Conference on Innovative Computing, Intelligent Communication and Smart Electrical Systems (ICSES)*, pp. 1–7, Sep. 2021.
- [23] S. Zhang, "Cost-sensitive KNN classification," *Neurocomputing*, vol. 391, pp. 234–242, 2020.
- [24] B. W. Yuan et al., "A novel density-based adaptive k nearest neighbor method for dealing with overlapping problem in imbalanced datasets," *Neural Computing and Applications*, vol. 33, no. 9, pp. 4457–4481, 2021.
- [25] M. D. de Lima, J. D. O. R. e Lima, and R. M. Barbosa, "Medical data set classification using a new feature selection algorithm combined with twin-bounded support vector machine," *Medical and Biological Engineering and Computing*, vol. 58, no. 3, pp. 519–528, 2020.
- [26] S. N. B. Jaini et al., "Indirect tool monitoring in drilling based on gap sensor signal and multilayer perceptron feed forward neural network," *Journal of Intelligent Manufacturing*, vol. 32, no. 6, pp. 1605–1619, 2021.
- [27] M. Jafari, Y. Wang, A. Amiryousefi, and J. Tang, "Unsupervised learning and multipartite network models: a promising approach for understanding traditional medicine," *Frontiers in Pharmacology*, vol. 11, Art. no. 1319, 2020.
- [28] F. Yu et al., "Unsupervised domain adaptation for object detection via cross-domain semi-supervised learning," *arXiv preprint arXiv:1911.07158*, 2019.
- [29] T. Reddy et al., "Antlion re-sampling based deep neural network model for classification of imbalanced multimodal stroke dataset," *Multimedia Tools and Applications*, pp. 1–25, 2020.
- [30] G. Forman, "An extensive empirical study of feature selection metrics for text classification," *Journal of Machine Learning Research*, vol. 3, pp. 1289–1305, Mar. 2003.
- [31] R. Solhmirzaei, H. Salehi, V. Kodur, and M. Z. Naser, "Machine learning framework for predicting failure mode and shear capacity of ultra high performance concrete beams," *Engineering Structures*, vol. 224, Art. no. 111221, 2020.



- [32] A. Gumaee et al., "A deep learning-based driver distraction identification framework over edge cloud," *Neural Computing and Applications*, pp. 1–16, 2020.
- [33] D. Chicco and G. Jurman, "The advantages of the Matthews correlation coefficient (MCC) over F1 score and accuracy in binary classification evaluation," *BMC Genomics*, vol. 21, no. 1, pp. 1–13, 2020.
- [34] D. M. Powers, "Evaluation: from precision, recall and F-measure to ROC, informedness, markedness and correlation," *arXiv preprint arXiv:2010.16061*, 2020.