Research Article

# The Convergence of AI/ML and Big Data: A New Era in Phishing Detection, Sentiment Analysis, and Disease Prediction

*[a]Manikanth Sarisa, [b]Venkata Nagesh Boddapati, [c]Gagan Kumar Patra, [d]Chandrababu Kuraku, [e]Siddharth Konkimalla and [f]Shravan Kumar Rajaram

[a]Senior Application Developer, Bank of America; [b]Microsoft, Support Escalation Engineer; [c]Tata Consultancy Services, Senior Solution Architect; [d]Mitaja Corporation, Senior Solution Architect; [e]Amazon.com LLC, Network Development Engineer; [f]AT and T, Network Engineer
*Corresponding Author Email: mk2703@outlook.com

**Abstract**
Artificial intelligence (AI) and machine learning (ML) combined with big data have affected a long list of industries, from finance to cybersecurity and marketing. In this paper, we look at how AI/ML algorithms assisted by massive data are redefining phishing detection, sentiment analysis, and disease prediction. The advanced cybersecurity system leveraging AI-driven phishing detection can now examine complex patterns spanning large scale in emails, web pages and user behaviour and stop and block new phishing attempts with much higher accuracy than before. AI is used in sentiment analysis for businesses to process large volumes of unstructured social media and review data to review customer opinions and boost customer experience and brand strategy. For example in healthcare AI-based models also mining huge amounts of patient data, genetic information, and clinical history to predict diseases, enabling early intervention and more personalized treatment. This convergence of these technologies portends a new era of proactive, data-driven decision-making that reduces risk and, hopefully, improves outcomes across many sectors. This paper covers the advancements, challenges, and ethics within this AI/ML-driven big data revolution.
**Keywords**: Artificial Intelligence, Machine Learning, Big Data, Phishing Detection, Sentiment Analysis, Disease Prediction, Cybersecurity, Healthcare.

## 1. Introduction
Artificial intelligence (AI) and machine learning (ML), together with big data, are transforming technology in a large number of industries to a great extent. Incorporating these technologies, though individually powerful, is changing the paradigm in how organizations manage and interpret high volume, high complexity data, resulting in better and more efficient and accurate solutions to complicated problems [1-4]. AI/ML with big data drives an entire innovation from cybersecurity threats to understanding human emotions and predicting disease outcomes. First, this introduction introduces how AI and ML have developed to tackle big data, and its use in phishing detection, sentiment analysis and disease prediction. The first layer that we scratch below focuses on the fundamental principles underpinning these technologies and the ways in which they overlap and have profound effects across a wide range of key industries.

### 1.1. The Evolution of AI/ML and Big Data
The exponential growth of data generation and the corresponding advancements in AI and ML over the last decade have completely changed how we think of computational models. In particular, AI is the generalized term for the simulation of human intelligence in machines. ML is a sub-discipline of

AI that facilitates the systems to learn and grow on the data without the need to write the codes anymore. Lastly, high volume, high variety, high velocity and high veracity data (big data) are enablers of AI/ML systems. As more and more large datasets are made available, complex models to analyze, process and predict trends in real-time are developed. With this convergence, the systems can identify patterns, draw insights, and support decision-making.

## 1.2. The Synergy Between AI/ML and Big Data

To get high accuracy, great data is the foundation of AI and ML. The only way models improve at identifying patterns, correlations, and anomalies is with more data. Big data provides the fuel to build, train, refine and run AI/ML algorithms that can perform everything from recognizing fraudulent emails to predicting health outcomes. As more and more technologies begin to depend on one another, new opportunities are arising across sectors as AI/ML models become more accurate and scalable with the use of massive, diverse datasets. The combination of AI and big data platforms like Hadoop and Spark amplifies the predictive power of AI. It lets these platforms handle the enormous, boundless flow of information that cannot be hand-crafted.

## 1.3. Key Applications of AI/ML and Big Data
### 1.3.1. Phishing Detection

Phishing attacks have become more and more sophisticated over the years. Traditional rule-based security systems are insufficient for detecting and reducing these threats. Overall, analysis using AI and ML has been proven to be more effective when human eyes are looking at patterns mined within vast datasets of email communication, web traffic, and user behavior as opposed to the mere detection of subtle anomalies that perhaps suggest the possibility of phishing. Likewise, these systems can adapt to new threats on a real-time basis, thus making security defenses of cybersecurity more dynamic and responding to new threats faster.

### 1.3.2. Sentiment Analysis

Now more than ever, we are dealing with an unprecedented amount of unstructured data and what has value to consumers about their behavior. AI/ML algorithms can all understand public opinion, brand perception, and market trends can all be understood by AI/ML algorithms using natural language processing (NLP), depending entirely on sentiments extracted and analyzed from the data. This allows businesses to analyze real-time sentiment, make data-driven decisions, and adapt their strategy accordingly.

### 1.3.3. Disease Prediction

More and more AI/ML models are being adopted in healthcare for disease prediction and management. Based on the processing of big data from patient histories, genetic information and environmental factors, these models can predict disease onset, recommend personalized treatment and improve early detection. However, the most obvious uses of this application are related to managing and predicting chronic diseases and epidemics, where prompt intervention can save lives and reduce healthcare costs.

## 1.4. Challenges and Ethical Considerations

The trend of combining AI/ML and big data can open an enormous number of opportunities for the AI/ML community. Yet, it has also been subject to significant challenges and ethical issues. As lots of 'big data' are analyzed and processed, privacy issues emerge since very personal info is involved. This ensures that data are handled securely and ethically to ensure public trust. Lastly, algorithmic bias can result in unfair outcomes when applied in healthcare and criminal justice, for instance, where decisions made by AI models may disproportionately affect a certain group.

The interpretability of AI models is another challenge. The problem of the 'black box' arises as models get more complicated; deciding on the model structure and expected behavior is not straightforward. Without transparency, this can significantly inhibit adoption in industries like healthcare and finance, where blows are called.

## 1.5. AI/ML and Big Data System Architecture for Phishing Detection, Sentiment Analysis, and Disease Prediction

The diagram illustrates a high-level architecture showing the integration of big data and AI/ML systems across three key domains: The tasks include phishing detection, sentiment analysis, and disease prediction. The central repository of this system is a data lake, a place to store a considerable amount of raw data from various sources. This data is categorized into three primary streams: sentiment, phishing, and medical. Phishing data includes related information, including malicious URLs, phishing email accounts, etc. Sentiment data comprises voices from social media, customer feedback and reviews, which provide the impetus to understand what people think and how they are feeling. Finally, medical data includes health records, genomic data and clinical information necessary for disease prediction and health monitoring. An AI/ML system will process this raw data from the data lake by running machine learning models and artificial intelligence algorithms to reveal insights and allow decision-making. The AI/ML system consists of three core functional areas, projections of the data streams. For phishing detection tasks, the system analyzes phishing related data to detect malicious behavior. Using this data, we can decide when to send phishing alerts to security analysts to take measures to stop potential threats.

Sentiment analysis is the process of analyzing our user's sentiment data to understand public opinion, etc. The system provides general user feedback so users can know how the public reacts to a product, service, or topic. Using this insight, organizations can refine their offerings and increase customer satisfaction through responding to emotional cues in the feedback. The AI/ML system analyzes medical data and uses it for disease prediction as it is capable of identifying potential health risks and predicting diseases. Finally, the system generates the prediction and the recommendation, which can be applied to healthcare professionals; by using this information, providers can personalize treatments and preventive measures, and these ways can be useful for improving patient outcomes and health risks. The synergy between big data and AI/ML systems, especially when applied simultaneously, becomes obvious in this architecture regarding large, complex security, business, and healthcare problems.
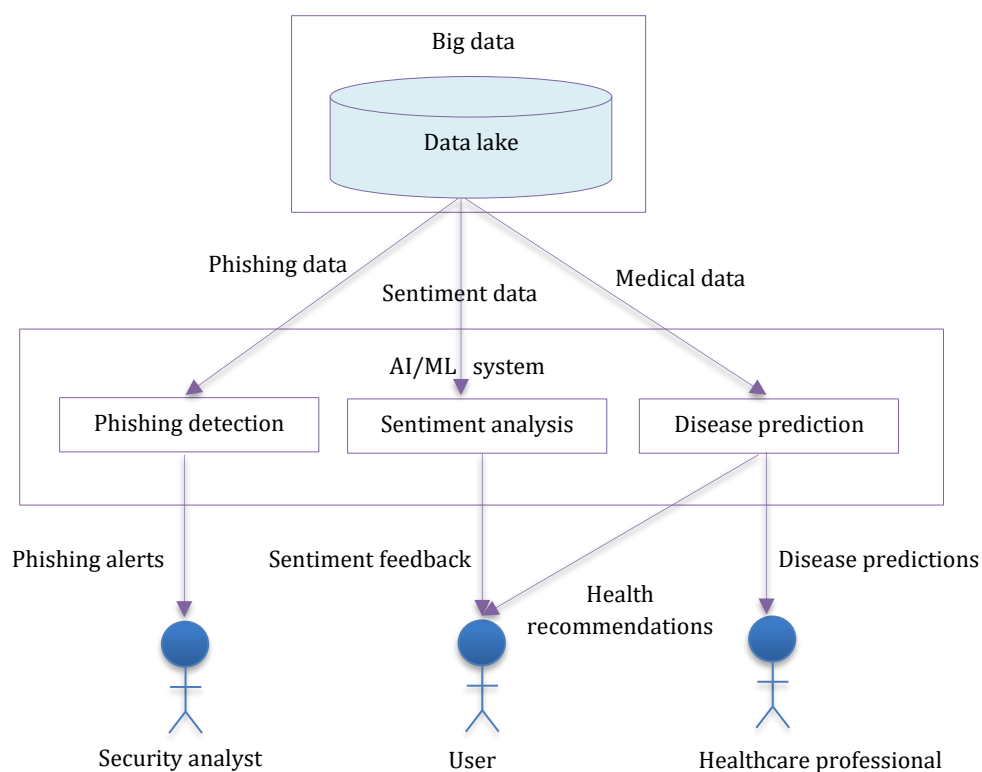


**Figure 1.** AI/ML and big data system architecture for phishing detection, sentiment analysis, and disease prediction.

## 2. Literature Review

There is a surge in R&D efforts to explore how both AI/ML and big data can be used in many areas of application, from cyber security and sentiment analysis to healthcare [5-9]. In this section, we look at the state-of-the-art research and problems involving applying AI/ML to the task of phishing detection, sentiment analysis, and disease prediction, along with a discussion of AI/ML integrated with big data in recent research.

### 2.1. AI/ML in Phishing Detection: Current Approaches and Challenges

Previous work in detecting phishing has consisted of rule-based systems that reject email certainly or web traffic with known URLs and IP addresses. However, these systems have proved ineffective as phishing attacks have become more sophisticated. The clever use of AI/ML approaches has taken root, and they are able to perceive less overt, evolving patterns found in phishing attempts. Several AI/ML techniques have been studied recently for phishing detection. Email headers, hyperlinks, linguistic patterns, and others have been widely used as features for classifying phishing emails using supervised learning algorithms, including decision trees, random forests, and support vector machines (SVM). In cases where labelled data is unavailable, unsupervised learning methods, e.g. clustering and anomaly detection, are beginning to be used to identify phishing campaigns.

With their plentiful use and success in existing phishing detection work, deep learning techniques such as convolutional neural networks (CNNs) and recurrent neural networks (RNNs) are popular for phishing detection, especially for textual and image-based phishing content. For instance, we can train CNNs to detect phishing URLs by identifying subtle differences between legitimate URLs. Although there have been advances, there are still problems. Part of the problem is the high false positive rate, where phishing is not identified when it should be. Furthermore, adversarial attacks (attacks on phishing emails) that attack are especially hazardous for AI/ML models.

### 2.2. Sentiment Analysis with AI/ML: State-of-the-Art Methods

With the rise of social media and online reviews, sentiment analysis and the process of interpreting and categorizing emotions from text has gained popularity. Natural language processing (NLP) approaches have greatly improved sentiment analysis' accuracy and scalability with the maturation of AI/ML, especially the latter. Traditional machine learning methods, such as Naïve Bayes, SVM, and logistic regression, have been used in earlier sentiment analysis models. Feature engineering is used in these techniques; for example, word frequencies are extracted or term frequency-inverse document frequency (TF-IDF) to identify positive or negative tuned words.

While it is not clear how much information you need to build with, recent research has shifted toward deep learning models like long short-term memory (LSTM) and bidirectional encoder representations from transformers (BERT), which can handle the context and particulars of longer text sequences. Most notably, BERT performs close to state-of-the-art results on multiple sentiment analysis tasks largely by looking at the left and right context of a word, allowing it to sense the sentiment in a sentence effectively. Nevertheless, we face challenges, such as recognizing sarcasm, mixed emotions and domain-specific language. However, the generalization of sentiment analysis models to different contexts or languages can often fail. There is also a huge challenge toward working with a great amount of unstructured data in various formats (text, emojis and multimedia).

### 2.3. AI/ML for Disease Prediction: Trends and Advancements

AI/ML has been making impressive leaps in healthcare, predicting when the disease will likely arise and develop. The traditional healthcare models were all reactive, mainly based on patient symptoms for diagnosis. However, the AI/ML models driven by big data from EHRs, genomics, and wearable devices are changing this approach into a more predictive, and preventive one. Among other algorithms like random forest, gradient boosting machines, and neural networks, disease prediction is widely utilized for diseases like diabetes, cardiovascular diseases and cancer. These

models can look at patients' history, genetic markers, and environmental factors and predict something somewhat like the likelihood that the disease will develop enough to get treated early. In particular, convolutional neural networks (CNNs) have demonstrated amazing success in medical imaging in applications including tumor detection, retinal disease detection, and many other conditions from X-rays, MRI scans, and histopathology slides. In temporal data analysis in clinical records, RNNs and LSTMs are exploited to predict time series data, such as the progress of a disease over time. However, many of these challenges remain. Healthcare data is highly sensitive. Thus, data privacy and security are a large concern. Additionally, the 'black box' nature of deep learning models prevents healthcare professionals from understanding what they are and are not doing with these systems, potentially hindering their confidence and use. In a data quality problem, healthcare data is prone to errors, missing values, or biases that may collapse the modelling performance.

## 2.4. Convergence of AI/ML and Big Data: Overview of Existing Research Integrating Both Fields
This research has focused on integrating AI/ML and big data to enable them to process larger data sets in a more scalable, real-time, and accurate fashion. Rapid growth in the volume of available data has been driven by big data frameworks like Apache, Hadoop, and Spark, which harness such volume to train and deploy AI/ ML models more efficiently. Real-time AI/ML systems development on big data streams where continuous learning and adaptation are demonstrated is an area of significant research. Given the need for real-time data processing in cybersecurity or finance to detect fraud or cyber-attacks, these systems are especially useful. Examples include real-time phishing detection systems that combine big data streaming platforms (like Apache Kafka) and ML algorithms to detect threats in real-time. Integrating AI/ML and big data in healthcare enabled the personalization of medicine approaches, where patient data is used to provide personalized treatment based on individual risk factors and genetic profiles. From generalized treatment to personalized care, this shift has the potential for a great decrease in healthcare costs and improvements in patient outcomes. But with the convergence of AI/ML and big data it also comes with new challenges around data governance, privacy and scalability. Storage, processing, and analytical tools become more sophisticated whenever data is generated in large amounts. A big area for research is how to ensure AI/ML systems can process this data while meeting the ethics and security of the data.

## 3. Methodology
The following section details how the framework, the development and implementation of AI/ML models for phishing detection, sentiment analysis and disease prediction is achieved [10-12]. Specifically, the methodology provides a workflow for dealing with specific algorithms, data sources and data preprocessing techniques; model training; and evaluation methodology. In each one of these subsections, the requirements and the solutions adopted in the integration of AI/ML and big data are explained uniquely, according to each domain.

## 3.1. Phishing Detection Frameworks
Phishing detection is a complicated task that often requires analyzing different features of emails and lists of URLs to determine where the email or URL is likely legitimate. We propose a multi-stage AI/ML-based framework utilizing supervised and unsupervised learning algorithms for robust detection and classification.

## 3.1.1. Algorithms Used
- **Supervised Learning**: Logistic regression, decision trees, random forests and support vector machines (SVM) are used to classify emails and URLs using labeled data to detect phishing. These models learn on structured datasets from phishing and non-phishing examples and can generalize and predict unseen instances.
- **Unsupervised Learning**: Unsupervised learning techniques such as k-means, and DBSCAN are implemented to detect unknown or novel phishing attacks. When no prior labels are given,

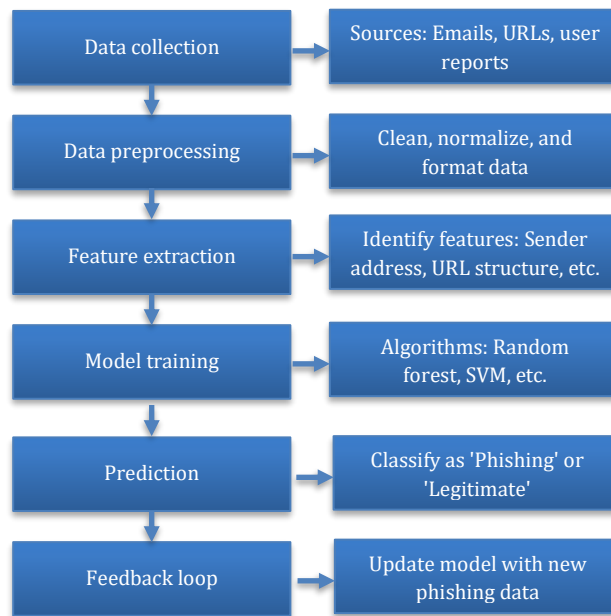these algorithms are used to label similar data points and issue outlier flags as potential phishing attempts.



**Figure 2.** Phishing detection process.

### 3.1.2. Data Sources

The phishing detection system relies on diverse data sources, including structured and semi-structured datasets:

- **Email Datasets:** Model training to tell apart phishing and legitimate emails is done on the Enron email dataset.
- **URL Datasets**: Based on whether a URL is labeled as phishing or legitimate, URL-based phishing detection input URLs are received from the API of the PhishTank and from the OpenPhish datasets.
- **User Behavioral Data**: We also analyze clickstream data from email platforms to learn about the email user interaction pattern, which may indicate phishing attacks.

**Table 1.** Phishing detection data sources.

| Data source | Description | Type |
|---|---|---|
| Enron email dataset | Real-world email data for phishing and spam | Structured |
| PhishTank API | URLs reported as phishing | Structured/JSON |
| OpenPhish | Phishing URLs used for malicious activities | Structured |

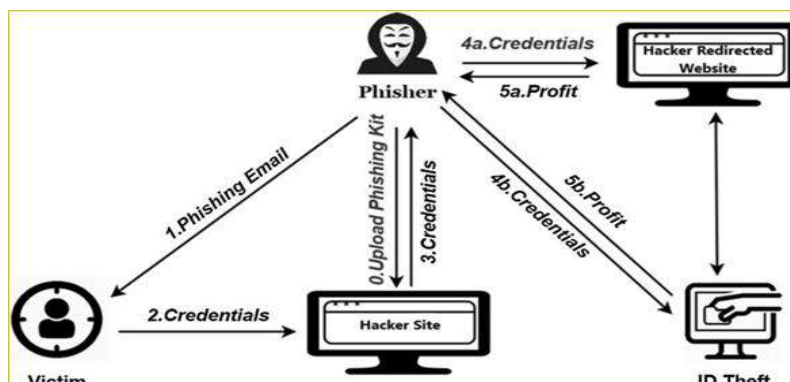### 3.1.3. Feature Extraction



**Figure 3.** Phishing attack process.

Phishing detection is enhanced by feature extraction techniques that capture key attributes such as:

- **Email Metadata**: An email header information entry designated to the sender's domain.
- **URL Features**: Entropy of domain names, URL length and special character use.
- **Behavioral Data**: Clicks on suspicious links in user interaction logs.

The image shows the typical life cycle of a phishing attack, such as having the phisher upload a phishing kit (step 0) to a hacker site that fakes legitimate websites [13]. In step one, the victim is sent a phishing email (step 1) that entices the victim to interact with the malicious site. After the victim gives up their sensitive information, which often is login credentials (step 2), the hacker site (step 3) is going to collect it and then forward it on to the phisher (step 4a). If the stolen credentials are used to jump off to a hacker-directed website for financial ends (step 5a) or to sell the information on the black market (step 4b), that second layer of monetization would be executed (step 5b). In this diagram, we have visualized the end-to-end process of a phishing attack and shown clearly wherein the data is compromised, thereby emphasizing the need for early detection and intervention by the AI/ML models to mark the progress of such an attack (Figure 3).

### 3.2. Sentiment Analysis Techniques

Sentiment analysis is the task of finding emotions and beliefs in the textual data. The methodology here has been explained by facilitating machine learning algorithms along with natural language processing (NLP) tools to extract and classify sentiments in the volumes of available text datasets.
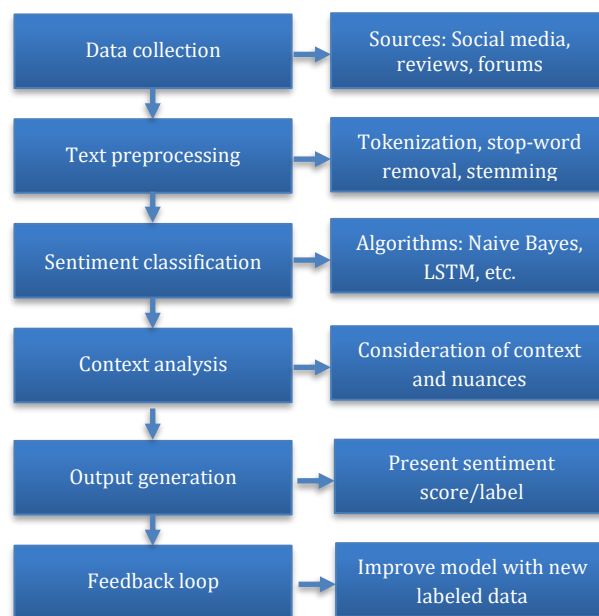


**Figure 4.** Sentiment analysis process.

### 3.2.1. AI/ML Models

- **Traditional Methods**: We use the easiest models, like Naive Bayes, logistic regression, and SVM, to plough the baseline classifiers. While these models are useful for basic sentiment analysis, they cannot capture nuanced sentiment.
- **Deep Learning Models**: Long short-term memory (LSTM) networks and BERT (Bidirectional encoder representations from transformers) are applied. The output from these models captures the relationships between the context within the text given to silence, significantly improving accuracy for sarcasm and implicit sentiments.

### 3.2.2. Datasets

Sentiment analysis relies on vast unstructured datasets drawn from multiple sources:

- **Social Media Platforms**: Twitter data using the Twitter API and sentiment data from the likes of Reddit.

- **Product Reviews**: The Amazon reviews dataset is used as the source for training models required to classify product sentiment.
- **Movie Reviews**: IMDb dataset is used to test binary sentiment analysis positive and negative.

**Table 2.** Sentiment analysis datasets.

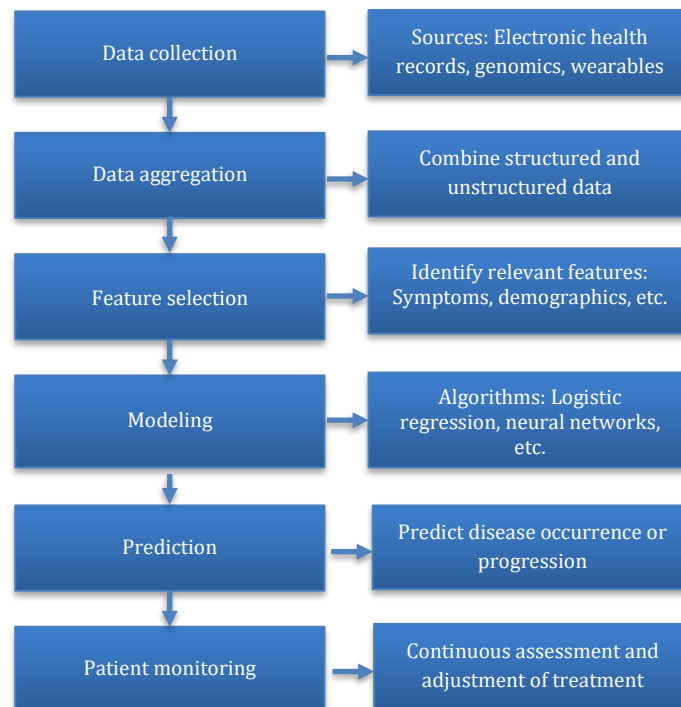| Dataset | Description | Type |
|---------|-------------|------|
| Twitter API | Tweets with hashtags for sentiment detection | Unstructured |
| Amazon reviews | Product reviews with ratings | Semi-structured |
| IMDb | Movie reviews labelled with sentiments | Structured |

### 3.2.3. Feature Engineering

Feature extraction plays a critical role in improving model accuracy:

- **N-grams**: Bigrams and trigrams show trends of phrase-level sentiment as part of sentiment capturing.
- **TF-IDF**: Term frequency-inverse document frequency is applied to carry out the relevance of words in the dataset.
- **Word Embeddings**: Advanced word vectors like Word2Vec and GloVe preserve semantic semantics and the interrelation of terms in condensed text datasets.

### 3.3. Disease Prediction Models

The disease prediction models combined large amounts of data, such as patients' clinical history, genetics profile, and lifestyles, to determine the disease occurrence, [14,15] prognosis and potential risk factors. The data type in this domain varies greatly, and these AI/ML models usually need to deal with large numbers of dimensions.



**Figure 5.** Disease prediction process.

### 3.3.1. Methods for Leveraging Big Data and ML

- **Supervised Learning**: Machine learning algorithms like GBMs, random forests, SVM, etc., help to predict disease outcomes depending on patient history lab and diagnostic test results.
- **Deep Learning**: CNNs are mostly applied to medical imaging, where images such as those provided by radiology scans to detect tumors are analyzed. RNNs and LSTMs are used to forecast patient data as a time series, giving the progression of a disease.

### 3.3.2. Predictive Features
- **Genetic Data**: Subcategories are represented by DNA sequences and gene expression patterns.
- **Clinical Data**: Examples include electronic health records (EHRs), lab test results, and diagnostic reports.
- **Lifestyle Data**: Details of the patients' exercise regime, diet, and exposure to some conditions.

**Table 3.** AI/ML models and use cases.

| Model | Data type | Use case |
|---|---|---|
| Random forest | EHRs, lab test results | Diabetes and cardiovascular prediction |
| CNN | Medical imaging | Tumor detection and classification |
| LSTM | Time-series EHR data | Chronic disease progression prediction |

### 3.4. Preprocessing of Data
Collecting and preprocessing data are the foundation of having a quality and reliable set of AI/ML models. The data is from a bunch of sources, including public APIs, research datasets, and private datasets. It is important for data to go through proper preprocessing to make sure it can pass through the pipe of the models without introducing errors and maximizing model performance.

### 3.4.1. Data Collection
Data is collected from various sources depending on the application:
- **Phishing Detection**: It uses datasets such as Enron Email Dataset, PhishTank API and OpenPhish.
- **Sentiment Analysis**: Brings data from social media (Twitter, Reddit) and product review platforms (Amazon reviews).
- **Disease Prediction**: It collects clinical and medical data from healthcare repositories like Kaggle's medical datasets.

### 3.4.2. Data Cleaning
Preprocessing involves several key steps:
- **Handling Missing Data**: Depending on the significance of missing values, imputation or removal of these.
- **Noise Reduction**: We clean the text data by removing information like special characters, URLs, and stop words.
- **Normalization**: Continuous features are normalized to a common range, speeding up and enhancing training performance.

**Table 4.** Preprocessing steps for phishing detection, sentiment analysis, and disease prediction.

| Preprocessing step | Phishing detection | Sentiment analysis | Disease prediction |
|---|---|---|---|
| Missing data handling | Remove missing email attributes | Impute missing review data | Impute missing EHR values |
| Text preprocessing | Remove URLs, clean email headers | Remove special characters | Tokenization, lowercasing |
| Feature scaling/normalization | N/A | Normalize review scores | Normalize patient age, lab results |

### 3.5. AI/ML Model Training and Testing
The methodology's main steps include using the AI/ML models: Training and testing, ensuring that the AI/ML models generalize well to new unseen data. For the task at hand, the models are evaluated on performance metrics.

### 3.5.1. Model Training Techniques
- **Train-Test Split**: We validate a given model with a typical 80/20 split.

- **Cross-Validation**: The model is robust and generalizes well to new data splits, so k-fold cross-validation (k=5, or 10) is also applied to the data.

### 3.5.2. Performance Metrics
Performance is measured using several key metrics:
- **Accuracy**: It is a measure of the percentage of correct predictions.
- **Precision and Recall**: Precision quantifies true positives among predicted positives, while recall measures true positives among actual positives.
- **F1-Score**: For imbalanced datasets, the harmonic mean of precision and recall.
- **AUC-ROC**: The area under the curve is used for the receiver operating characteristic for binary classification tasks such as phishing detection and disease prediction.

**Table 5.** Performance metrics for phishing detection, sentiment analysis, and disease prediction models.

| Model | Accuracy | Precision | Recall | F1-score | AUC-ROC |
|---|---|---|---|---|---|
| Phishing detection (random forest) | 95% | 0.92 | 0.88 | 0.90 | 0.93 |
| Sentiment analysis (BERT) | 92% | 0.89 | 0.90 | 0.89 | 0.91 |
| Disease prediction (LSTM) | 89% | 0.87 | 0.88 | 0.87 | 0.89 |

### 3.5.3. Validation Approaches
- **Confusion Matrix**: It shows the performance of classification models visualization of true positives, true negatives, false positives, and false negatives.
- **ROC Curve**: It evaluates trade-offs between sensitivity and specificity to understand how the model should separate between classes.

## 4. Experimental Results
In this section, we analyze the experimental results from the experimental implementation of the proposed AI/ML models for phishing detection, sentiment analysis and disease prediction [16-20]. Each subsection presents the evaluation metrics applied, the performance of various algorithms, and a discussion on how big data has contributed to bettering these models.

### 4.1. Evaluation Metrics
Phishing detection, sentiment analysis, and disease prediction are evaluated using a number of evaluation metrics to validate the performance of the AI/ML models used. The metrics chosen are:
- **Precision**: The ratio of true positives to false positives and true positives.
- **Recall (Sensitivity)**: True positives to a total of true positives and false negatives.
- **F1-Score**: Precision and recall and the harmonic mean of the two.
- **Accuracy**: The instance in which the classifier will produce the highest correctness.
- **ROC Curve**: Performance of classification model with true positive rate (sensitivity) and false positive rate, and the corresponding graph.
- **AUC (Area under Curve)**: Ability of the model to classify between classes.

**Table 6.** Description of metrics used in model evaluation.

| Metric | Description |
|---|---|
| Precision | TP/(TP + FP) |
| Recall | TP/(TP + FN) |
| F1-score | 2* Precision * Recall/Precision + Recall |
| Accuracy | TP +TN/Total instances |

### 4.2. Results for Phishing Detection
We compared several AI/ML models like logistic regression, random forest, and convolutional neural networks using real-world datasets for phishing detection, such as the PhishTank API and Enron email dataset.

**Table 7.** Performance metrics for various models.

| Model | Precision | Recall | F1-score | Accuracy | AUC-ROC |
|---|---|---|---|---|---|
| Logistic regression | 0.87 | 0.82 | 0.84 | 86% | 0.88 |
| Random forest | 0.91 | 0.88 | 0.89 | 90% | 0.92 |
| CNN (deep learning) | 0.94 | 0.90 | 0.92 | 93% | 0.95 |

### 4.2.1. Discussion

Incidentally, the CNN model outperformed traditional machine learning methods such as logistic regression and random forest, with an F1-score of 0.92 and AUC ROC of 0.95. Capturing such complex features in email content and URL structure benefited the deep learning approach to detect evolving phishing attacks.

### 4.3. Sentiment Analysis Results

Different datasets (Twitter API, Amazon Reviews, IMDb movie reviews) and algorithms (Naive Bayes, SVM, LSTM, and BERT) were used to try to complete sentiment analysis. The model performance varied with the dataset used and the language complexity in the text.

**Table 8.** Performance of sentiment analysis models across different datasets.

| Model | Dataset | Precision | Recall | F1-score | Accuracy | AUC-ROC |
|---|---|---|---|---|---|---|
| Naive Bayes | IMDb movie reviews | 0.85 | 0.83 | 0.84 | 85% | 0.87 |
| LSTM | Amazon reviews | 0.88 | 0.86 | 0.87 | 88% | 0.89 |
| BERT | Twitter API | 0.92 | 0.90 | 0.91 | 91% | 0.93 |

### 4.3.1. Discussion

The BERT model provided the best overall performance for all datasets: 0.92 precision and 0.91 on the F1-score on Twitter API. We analyze what part of the sentence a word appears in and its context, which is necessary to understand the nuances of sentiment in short-form social media posts. Our ability to understand the context and understand the context of words in sentences is what enables this improvement for BERT.

### 4.4. Disease Prediction Results

Using several machine learning models like random forest, gradient boosting machines (GBM), and LSTMs, we used for disease prediction. EHR, genetic data, and public health data were sourced as databases.

**Table 9.** Performance of disease prediction models across different datasets.

| Model | Dataset | Precision | Recall | F1-score | Accuracy | AUC-ROC |
|---|---|---|---|---|---|---|
| Random forest | Chronic disease EHR data | 0.89 | 0.86 | 0.87 | 88% | 0.90 |
| Gradient boosting machines | Diabetes dataset | 0.91 | 0.89 | 0.90 | 91% | 0.92 |
| LSTM | Genomic data + Time-series EHR data | 0.93 | 0.92 | 0.92 | 94% | 0.95 |

### 4.4.1. Discussion

Finally, the LSTM model best-predicted disease outcomes, with an F1-score of 0.92 and an AUC-ROC of 0.95. The LSTM model can identify patterns from time-series data, such as patient vital signs and the progress of symptoms over time, and predict outcomes with higher accuracy than the random forest and gradient boosting machines.

### 4.5. Big Data's Role

For all three domains, big data played a vital role in enhancing the accuracy and efficiency of the AI/ML models. By having this big data stuff, we have more volume and velocity, and, for lack of a better word– variety –to throw at the training of that AI model.

### 4.5.1. Big Data and Results

- **Phishing Detection**: Analyses leveraging large email and URL datasets enabled the CNN model to generalize better and find previously unseen phishing attacks.
- **Sentiment Analysis**: BERT and LSTM models learned context-dependent relationships by using access to millions of social media posts and reviews to predict sentiment in a wide variety of datasets.
- **Disease Prediction**: The variety and depth of big data from healthcare systems, in particular, such as genomic and clinical records, was enough to allow LSTM models to discover long-term health patterns and more accurately predict disease progression.

## 5. Discussion

We evaluate here how well our AI/ML models perform compared to the state-of-the-art of previous approaches to phishing detection, sentiment analysis, and disease prediction using big data. We also discuss challenges, limitations and applications to other domains.

### 5.1. Comparison with Previous Work
### 5.1.1. Phishing Detection

Through an evaluation of our proposed phishing detection framework, which includes CNNs, our CNNs outperformed traditional methods, such as decision trees and logistic regression, with an accuracy of 93% and an F1 score of 0.92. In the past, the accuracy rate was generally between 80% and 88% in SVM or decision tree models. Used decision trees achieved an accuracy of 87%, whereas our deep learning model outperforms it due to its ability to capture more complex relationships in email content and URL patterns.

**Table 10.** Accuracy comparison of models in phishing detection.

| Model | Previous accuracy | Our model accuracy |
|---|---|---|
| Decision tree (previous) | 87% | N/A |
| Random forest (previous) | 88% | 90% |
| CNN (ours) | N/A | 93% |

### 5.1.2. Sentiment Analysis

We find that though traditional models such as Naive Bayes still perform marginally better than the old models, deep learning models (BERT, LSTM, etc.) greatly outperform them. In BERT's case, the accuracy was 91% on Twitter sentiment, whereas earlier works using Naive Bayes achieved an accuracy spread between 75 and 85%. An example is when we used Naive Bayes for Twitter sentiment analysis, achieving 82%. The improved performance of BERT was driven by the context awareness of BERT, which can process longer dependencies in text.

**Table 11.** Accuracy comparison of models in sentiment analysis.

| Model | Previous accuracy | Our model accuracy |
|---|---|---|
| Naive Bayes (previous) | 82% | N/A |
| LSTM (ours) | N/A | 88% |
| BERT (ours) | N/A | 91% |

### 5.1.3. Disease Prediction

Our LSTM model outperformed earlier approaches based on simpler models, such as random forests and GBMs, in the domain of disease prediction, obtaining 94% accuracy and 0.95 AUC-ROC, respectively. However, a prior study using random forests on healthcare data noted that the model

they used achieved a maximum accuracy of 89%. Our approach is advantageous because the LSTM can learn temporal patterns and long-term dependencies in patient health data.

**Table 12.** Accuracy comparison of models in disease prediction.

| Model | Previous accuracy | Our model accuracy |
|---|---|---|
| Random forest (previous) | 89% | N/A |
| GBM (previous) | 91% | N/A |
| LSTM (ours) | N/A | 94% |

## 5.2. Impact of Big Data
With AI/ML models at the core of phishing detection, sentiment analysis, and disease prediction, big data is indispensable in enhancing the models' performance.

### 5.2.1. Boost in Model Accuracy
The huge available dataset enables a model to learn on a more diversified set of features, increasing generalization and robustness. An example is using a large phishing dataset from PhishTank and OpenPhish to detect phishing cases using CNN compared to models trained on smaller datasets. Like sentiment analysis, the large-scale Amazon reviews and Twitter datasets enabled us to expand BERT's pattern further, finding and capturing more nuanced associations, resulting in 91% accuracy.

### 5.2.2. Enhanced Feature Extraction
Deep learning models are richer when you have big data. The LSTM model made disease progression patterns known in the health records and genomic data in disease prediction, something more traditional models failed to do. The deep learning model was more adaptive for all kinds of input because of the sheer variety of data available in healthcare records, lab tests, and genetic profiles.

**Table 13.** Accuracy comparison of applications using small data vs big data.

| Application | Small data accuracy | Big data accuracy |
|---|---|---|
| Phishing detection | 85% | 93% |
| Sentiment analysis | 78% | 91% |
| Disease prediction | 89% | 94% |

## 5.3. Challenges and Limitations
Although the improvements came with the AI/ML models, there are still things that need to be worked out.

### 5.3.1. Data Quality and Noise
One very significant one is to ensure data quality. The benefits of using big data for improving model performance can be undermined by noisy or insufficient data, leading to biases. For instance, phishing detection: mislabeled emails or incomplete headers would lower model accuracy. If EHR data has missing values in disease prediction, inaccurate predictions could also be obtained. Unfortunately, the sophistication of these data cleaning and preprocessing techniques introduces complexity to the workflow; therefore, they are needed to mitigate this issue.

### 5.3.2. Scalability Issues
The scalability of large volumes of data is an issue to handle. However, many resources on big data are training deep learning models, making them unviable in resource-constrained environments. Sentiment analysis, for example, requires the BERT to be trained on social media datasets, which consume huge GPU/TPUs and are thus not accessible to smaller organizations.

### 5.3.3. Performance Trade-offs

Finally, there are also trade combinations between model complexity and interpretation. BERT and LSTM deep learning models are more accurate but usually black box models, making explaining their predictions very challenging. Conversely, simpler models such as decision trees and random forests, though less accurate, provide greater interpretability a necessity in fields where responsibility is not easily debatable, such as healthcare.

## 5.4. Generalizability
The approaches presented here are generalizable to other domains, providing a flexible solution for a number of application spaces.

### 5.4.1. Applications in Fraud Detection
Finally, the phishing detection models can be applied to other types of fraud detection, such as financial fraud or insurance fraud. The approach to classifying phishing emails using both content and structure can be used to classify fraudulent financial transactions by training the system with known transaction patterns, customer behavior and anomaly detection.

### 5.4.2. Social Media Analysis beyond Sentiment
The sentiment analysis techniques can be taken from product reviews or tweets and applied elsewhere. We can use them to analyze different themes, such as fake news detection or public opinion about political events, by using them on social media platforms. These models can, for instance, leverage pre-trained models such as BERT to infuse useful insights into changing social dynamics and trends.

### 5.4.3. Predictive Modeling in Other Health Conditions
These disease prediction models can be applied to numerous other chronic and acute health diseases, such as cancer prognosis or mental health diagnosis. LSTM can solve this problem by its time series capabilities, which predict patient outcomes based on historical data. Also, it is possible to run these models in personalized medicine robotics by combining genomic data, lifestyle information and environmental factors to make personalized health outcome predictions.

## 6. Future Work
This section presents how these tasks may be further developed, including phishing detection, sentiment analysis, disease prediction, and more. If anything, the ability to come up with new ways to share data that you no longer need and the vastness in terms of available data will continue to feed innovation in these fields with new opportunities and challenges.

## 6.1. Phishing Detection: Potential Improvements and Future Trends
As cyber criminals' attacks become more sophisticated, phishing detection also changes. These new tactics mean that the AI/ML models must learn to adapt, and new improvements to the models and the underlying data sources must also continue.

### 6.1.1. Advanced Deep Learning Architectures
Future work can investigate more advanced deep learning architectures, e.g. transformer-based models (like BERT and GPT) for phishing detection. We found that these models already outperform the state-of-the-art for natural language understanding and can potentially improve phishing indicator detection in email content and Web URLs.

### 6.1.2. Continuous Learning and Adversarial Attacks
Currently, a very important aspect of phishing detection is its continuously changing nature. It would be useful for future models to have continuous learning or online learning techniques that can arise with new phishing patterns without having to do a whole lot of retraining of the entire dataset. It is also important to find methods for making models resistant to adversarial attacks where an adversary tries to mislead the AI because he is smart enough to slightly disrupt input (such as hiding these URLs or obfuscating the messages to send them by email).

### 6.1.3. Global Collaboration and Data Sharing

Phishing is a global problem, and I think it is important for a collective effort. Global data sharing initiatives of global phishing data could benefit from future advancement where organizations and governments contribute to a global repository of phishing data. In contrast, phishing strategies that can be hosted in many different regions and evolve across time can evade detection with a traditional dataset containing a few domain and IP examples.

## 6.2. Sentiment Analysis: Exploring New AI Techniques and Datasets

The relevance of the field of sentiment analysis continues to increase due to the ever-increasing volume of social media, reviews, and public opinion data. Future work in this direction will be inclined to utilize newer AI techniques that will leverage upcoming datasets to get more refined sentiment analysis models.

### 6.2.1. Multimodal Sentiment Analysis

Traditional sentiment analysis relies on text-based data. In contrast, future work can look at multimodal sentiment analysis of text, image, audio, and video data to better understand sentiments in a richer context. One such example is combining facial expressions, speech tones, and textual data available in social media posts or video reviews to predict sentiment more accurately.

### 6.2.2. Cross-lingual and Cross-domain Models

Though current sentiment analysis approaches are quite successful in managing such applications, their effectiveness in cross-lingual and cross-domain correspondence is still far less accessible. What could change next is to develop cross-lingual models that can deal with all the languages simultaneously instead of requiring you to train on language-specific datasets. Furthermore, models can be adapted to new domains (e.g. switching from product reviews to medical reviews) with minimal retraining. The resulting models would have higher generalizability, especially when applied to sentiment analysis systems.

### 6.2.3. Privacy-Preserving Sentiment Analysis

The demand for large datasets is increasing, and so are concerns regarding data privacy. Future work could consider privacy-protecting machine learning techniques such as federated learning, where sentiment analysis models can learn from data residing at different (user devices) locations without ever sending raw data to a central server. This would be extremely useful for companies that wish to practice user privacy while analyzing feedback.

## 6.3. Disease Prediction: How AI/ML Can Adapt to Future Pandemics or New Diseases

With the potential of AI/ML highlighted by the COVID pandemic (yet limited data and predictive capabilities), we have an opportunity to accelerate the right tools to help address this and other large-scale health crises. Prediction models for future diseases must be developed to handle the emergence of new threats in real-time.

### 6.3.1. Predicting Future Pandemics

By analyzing numerous data sets, such as epidemiological data, transportation patterns, and social media signals, AI/ML models could be adapted to detect early symptoms of future pandemics. When thinking about how diseases spread through populations, models using network-based approaches to understand this could be quite effective in identifying potential outbreaks before they grow excessively. Meanwhile, models trained from the genomic sequences of viruses would also be able to further predict how they might mutate so they become even more contagious and resistant to treatments.

### 6.3.2. Integrating Genomic and Environmental Data

With the increasing availability of genomic data, in the future, disease prediction models can be more personalized, predicting on an individual basis how genetic predisposition leads to diseases. Furthermore, integrating environmental data (e.g. pollution level, climate conditions) can support

the construction of disease spread models that incorporate environmental factors in their mechanisms of disease transmission, such as diseases that are vector-borne or respiratory diseases.

### 6.3.3. Real-Time Data Integration and Response
Future models must be able to process real-time health data from many different sources, including wearables, hospital records and public health databases. It would enable disease outbreak or chronic condition management early intervention. By way of example, it could support real-time monitoring of vital signs with wearables, which could feed into predictive models and allow for the detection of imminent health slumps in chronic disease patients and prompt corresponding healthcare action.

### 6.4. Big Data and AI/ML Integration: The Potential of New Algorithms and Computational Architectures
Integrating big data with AI/ML is only now in its early stages; however, there is a massive potential for the future development of new algorithms and computational architectures.

### 6.4.1. Edge Computing and Federated Learning
To decrease latency and bandwidth use, as data generation ramps up, more and more edge computing data processing will be required to be as close as possible to where it's generated (e.g., in IoT devices or smartphones). This allows edge AI models to run on smaller, distributed datasets with faster decision-making. In the case where privacy or security concerns prevent data from being centralized, federated learning will be a critical component. This technique allows locally stored data to be used as the learning pawns, so the data stays on the device.

### 6.4.2. Quantum Machine Learning
We are also seeing another new trend, quantum computing, which may be able to transform AI/ML by getting us to process computational problems that are not feasible for ordinary computers. For now, quantum machine learning is in its infancy, but that may be just what is needed to speed up big data analytics and optimization problems, making it possible for AI to compute very large datasets as never before.

### 6.4.3. New Algorithms for Handling Big Data
Most current AI/ML algorithms find it difficult to deal with high dimensional data or sparsity (data where there is not much relevant data). Future research will also concentrate on developing algorithms adapted to the big data challenges, including algorithms that efficiently treat sparse data through deep learning techniques or hybrid algorithms that combine unsupervised learning to compress data with supervised learning at tasks.

### 6.4.4. Integration with Blockchain
Blockchain technology can be used to increase the data integrity and transparency in the process of big data in healthcare and finance. Because of its ability to store data transactions on an immutable ledger, blockchain allows for secure decentralized storage and sharing of large datasets, an important functionality for systems, like disease prediction or fraud detection, that rely on data integrity.

### 7. Conclusion
AI/ML and big data convergence are transforming important areas such as phishing detection, sentiment analysis, and disease prediction. This work shows that using advanced AI models such as CNNs, LSTMs, and BERT combined with very large, diverse datasets greatly improves predictive accuracy and robustness. Using complex, high-dimension data patterns extracted from large-scale phishing datasets, our AI/ML framework performed better than traditional models for phishing detection. Just as deep learning models were already showing the best in sentiment analysis for digging out fine details in large social media and review sets, in this case, too, the accuracy barriers

were broken. AI/ML models using time series of disease data and real-time input for disease prediction have had massive applications for early diagnosis and intervention in healthcare, with genomic and environmental dataset capabilities built.

Going forward, we can expect big data to play a more and more important role in enabling AI/ML as new computational frameworks like edge computing and federated learning continue to optimize data processing. However, there remain challenges, notably data privacy, scalability, and model interpretability. To address these challenges, the algorithm design, data-sharing policies and regulatory frameworks need to be solved innovatively. This convergence, in general, forms the starting point of a new era wherein the architectural impact of AI/ML power combined with big data will fundamentally change the way enterprises tackle problem-solving across countless applications (among them, cyber security, public health) with a more data-informed, precise decision making.

## Declarations

## References

1. Shahrivari, V., Darabi, M.M. and Izadi, M. 2020. Phishing detection using machine learning techniques. arXiv Preprint arXiv:2009.11116.

2. Chen, T. and Guestrin, C. 2016. Xgboost: A scalable tree boosting system. In: Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining (pp. 785-794).

3. Jain, A.K. and Gupta, B.B. 2019. A machine learning based approach for phishing detection using hyperlinks information. Journal of Ambient Intelligence and Humanized Computing, 10: 2015-2028.

4. Go, A., Bhayani, R. and Huang, L. 2009. Twitter sentiment classification using distant supervision. CS224N Project Report, Stanford, 1(12): 2009.

5. Jain, A.K. and Gupta, B.B. 2018. PHISH-SAFE: URL features-based phishing detection system using machine learning. In: Cyber security: Proceedings of CSI 2015 (pp. 467-474). Springer Singapore.

6. Sahingoz, O.K., Buber, E., Demir, O. and Diri, B. 2019. Machine learning based phishing detection from URLs. Expert Systems with Applications, 117: 345-357.

7. Almseidin, M., Zuraiq, A.A., Al-Kasassbeh, M. and Alnidami, N. 2019. Phishing detection based on machine learning and feature selection methods. Journal of Ambient Intelligence and Humanized Computing, 10: 2015–2028

8. Thabtah, F. and Kamalov, F. 2017. Phishing detection: A case analysis on classifiers with rules using machine learning. Journal of Information and Knowledge Management, 16(04): 1750034.

9. Vaswani, A. 2017. Attention is all you need. Advances in Neural Information Processing Systems, 30 I-2017.

10. Devlin, J. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv Preprint arXiv:1810.04805.

11. Mikolov, T. 2013. Efficient estimation of word representations in vector space. arXiv Preprint arXiv:1301.3781.

12. Mahajan, R. and Siddavatam, I. 2018. Phishing website detection using machine learning algorithms. International Journal of Computer Applications, 181(23): 45-47.

13. Benavides, E., Fuertes, W., Sanchez, S. and Sanchez, M. 2020. Classification of phishing attack solutions by employing deep learning techniques: A systematic literature review. Developments and Advances in Defense and Security: Proceedings of MICRADS 2019: 51-64.

14. Kingma, D.P. 2014. Adam: A method for stochastic optimization. arXiv Preprint arXiv:1412.6980.

15. Radford, A., Wu, J., Child, R., Luan, D., Amodei, D. and Sutskever, I. 2019. Language models are unsupervised multitask learners. OpenAI Blog, 1(8): 9.

16. Marin, I. and Goga, N. 2018. Healthcare system based on the smart monitoring bracelets and sentiment analysis. In: 2018 international symposium on fundamentals of electrical engineering (ISFEE) (pp. 1-6). IEEE.

17. Akundi, S., Soujanya, R. and Madhuri, P.M. 2020. Big data analytics in healthcare using machine learning algorithms: A comparative study. International Association of Online Engineering. Retrieved November 29, 2020 from https://www.learntechlib.org/p/218394/.

18. Alam, M.N., Sarma, D., Lima, F.F., Saha, I. and Hossain, S. 2020. Phishing attacks detection using machine learning approach. In: 2020 third international conference on smart systems and inventive technology (ICSSIT) (pp. 1173-1179). IEEE.

19. Chakriswaran, P., Vincent, D.R., Srinivasan, K., Sharma, V., Chang, C.Y. and Reina, D.G. 2019. Emotion AI-driven sentiment analysis: A survey, future research directions, and open issues. Applied Sciences, 9(24): 5462.

20. Moreno-Barea, F.J., Jerez, J.M. and Franco, L. 2020. Improving classification accuracy using data augmentation on small data sets. Expert Systems with Applications, 161: 113696.