# A Report on


## "Web Phishing Detection using Machine Learning Algorithms"


SUBMITTED BY


**Ansh Sharma**
**22070127011**

**Priyanshu Lathi**
**22070127048**

**Under the Guidance of**
**Dr Sachit T.S**


**SYMBIOSIS INSTITUTE OF TECHNOLOGY**


**(A CONSTITUTENT OF SYMBIOSIS INTERNATIONAL**
**UNIVERSITY)**
**Pune – 412115**

# Certificate

The report entitled Tool Wear Classification using Deep Learning submitted to the Symbiosis Institute of Technology, Pune is based on my original work carried out under the guidance of Dr. Sachit T S. The dissertation has not been submitted elsewhere for award of any degree.

The material borrowed from other source and incorporated in the dissertation has been duly acknowledged and/or referenced.

I understand that I myself could be held responsible and accountable for plagiarism, if any, detected later on.

**Signature of the candidate**

Ansh Sharma (22070127011)

Priyanshu Lathi (22070127048)

**Project Supervisor**                                                    **H.O.D.**

Dr. Sachit T S                                                    Dr. Arunkumar Bongale

# Abstract

This research paper investigates the efficacy of machine learning algorithms in detecting web phishing attempts, a prevalent cybersecurity threat. Utilizing a dataset comprising 11,430 instances and 89 features, diverse aspects of URLs and webpage characteristics were examined to discern potential malicious activities. The features encompassed various attributes such as URL length, presence of IP addresses, frequency of special characters, and webpage structural properties. The target variable, denoted as 'status', distinguished between legitimate URLs and phishing attempts.

Three popular machine learning algorithms, namely K-Nearest Neighbors (KNN), Logistic Regression, and Random Forest, were employed to train predictive models on the dataset. The models exhibited commendable accuracies, with KNN achieving 94%, Logistic Regression attaining 92%, and Random Forest yielding 93%. These accuracies underscore the effectiveness of employing machine learning techniques in combating web phishing attacks.

The study contributes to the ongoing efforts in cybersecurity by demonstrating the feasibility of employing machine learning algorithms for automated detection of phishing attempts. The findings highlight the importance of feature-rich datasets in training robust models capable of accurately distinguishing between legitimate web content and malicious phishing schemes. Future research endeavors could focus on exploring ensemble methods or deep learning architectures to further enhance the detection capabilities and resilience against evolving phishing techniques in cyberspace.

# INDEX

# Chapter-1

# Introduction

## 1.1 Introduction to Artificial Intelligence

Artificial intelligence (AI), also known as machine learning or deep learning, is a broad field of computer science that focuses on the development of systems that can perform tasks that usually require human intelligence. It encompasses a wide range of cognitive abilities, such as learning, reasoning, problem solving, natural language, perception and sensory processing, creativity, emotional intelligence, and much more. The field of artificial intelligence has its roots in the early 20th century, with Alan Turing proposing the famous Turing Test (a measure of a machine's intelligence) and John McCarthy coining the term "artificial intelligence" and organizing the Dartmouth Conference (1956), widely regarded as the birthplace of AI as a discipline.

Artificial intelligence (AI) can be implemented in Different approaches:

**Symbolic AI:** The first type of AI is symbolic AI, also known as good old-fashioned AI (GOFAI). In this type of AI, human knowledge and rules are encoded into computer systems so that they can reason and make decisions using symbolic representations.

**Machine Learning:** The second type of AI is machine learning (ML). ML involves creating algorithms and statistical models to help computers perform tasks without explicitly programming them. Instead, ML systems learn from the data, recognize patterns, and use that learning to make predictions or decisions.

## 1.1.1 AI Applications:

1. **Healthcare:** AI plays a crucial role in analyzing medical images for disease detection, analyzing patient records for diagnosis and treatment, and forecasting patient outcomes. This technological advancement is reshaping healthcare delivery, enhancing precision, and optimizing patient care.
2. **Finance:** AI drives fraud detection, risk evaluation, automated trading, and tailored customer service, streamlining financial activities, bolstering security, and refining investment approaches for both individuals and organizations.
3. **Autonomous** Vehicles: AI empowers self-driving vehicles to perceive their surroundings, make split-second decisions, and navigate safely, revolutionizing transportation, reducing accidents, and enhancing the movement of people and goods.

4. **E-commerce**: AI-driven recommendation engines and virtual assistants customize the shopping experience, increase sales, and elevate customer satisfaction, fostering engagement and loyalty in online retail settings.

5. **Manufacturing:** AI-powered predictive maintenance and robotic systems streamline production processes, minimize operational downtime, and boost efficiency, transforming manufacturing operations and elevating productivity levels.
6. **Education:** AI-based platforms for personalized learning and intelligent tutoring systems adapt to the unique needs of each student, fostering engagement, mastery, and positive outcomes, reshaping the delivery of education and promoting continuous learning.

## 1.2 **Introduction to Machine Learning**

**Machine Learning (ML)** is a branch of Artificial Intelligence (AI) concentrated on creating algorithms and statistical models. These models empower computers to learn from data and formulate predictions or decisions. Unlike conventional programming, which relies on specific instructions, ML systems discern patterns and connections from data without direct programming.

Machine Learning (ML) represents a specialized subset of Artificial Intelligence (AI), honing in on the refinement of algorithms and statistical models. These tools facilitate computers in assimilating insights from data, thereby enabling them to formulate predictions or decisions. Unlike traditional programming paradigms reliant on explicit instructions, ML systems autonomously discern patterns and correlations within data, independent of direct programming input. This approach fosters adaptability and scalability, empowering machines to continuously refine their understanding and performance. ML's emphasis on data-driven learning distinguishes it from conventional programming, positioning it as a dynamic and versatile toolset for addressing complex challenges across various domains.

**Machine learning is Divided into two subfields on the basis of training data:**

1. **Supervised Learning:** When the model is trained on the labelled data or data which has Output features corresponding to its input features.
2. **Unsupervised Learning:** When the model is trained on the unlabelled data or data which does not have output features corresponding to its input features.

## 1.2.1 **Supervised Machine Learning**

Supervised learning, a branch of machine learning, involves training algorithms using labeled data, which consists of input features and corresponding output labels. Through this process, the algorithm learns to map input data to output labels, facilitating predictions or decisions on unseen data. The primary objective of supervised learning is to develop a model capable of accurately predicting output labels based on input data. This approach enables machines to generalize patterns and relationships within the data, enhancing their predictive capabilities in various applications.
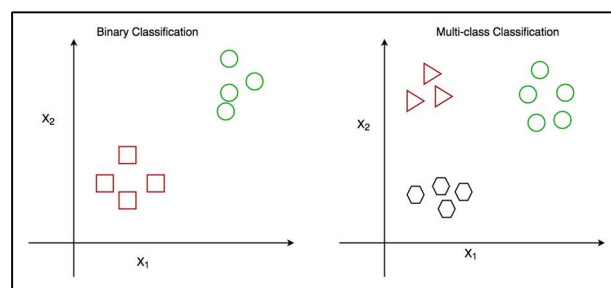
**Supervised Learning is Categorized into two categories:**

### 1. Classification:

Classification, a fundamental aspect of supervised learning within machine learning, entails categorizing input data into predefined classes or categories. This process involves training algorithms using labeled datasets, where input features are associated with discrete output labels. The ultimate aim is to establish a mapping from input features to specific output labels, enabling accurate classification of new, unseen data instances. Through this approach, machines learn to recognize patterns and characteristics within the data, enhancing their ability to categorize and make decisions across a range of applications.

**There are two Types of Classification in Machine Learning:**

i.  **Binary Classification:** Binary classification is a machine learning task where data is categorized into two exclusive classes, such as spam or non-spam in email filtering, disease or non-disease in medical diagnosis, and positive or negative sentiment in sentiment analysis.



ii. **Multiclass Classification:** multiclass classification extends beyond binary classification, involving the classification of data into three or more distinct categories. Examples include identifying objects in images, categorizing documents into topics, and classifying animals into different species.

The classification can be done using various Algorithms such as Random Forest, Decision Support Vector machine, KNN (K nearest neighbours)

### 2. Regression:

Regression in machine learning serves as a cornerstone of supervised learning endeavors, dedicated to predicting continuous numerical outcomes through the utilization of input features. The essence of this task lies in grasping the intricate mapping between input variables and a continuous output, denoting a numerical value of substantial import. By delving into regression analysis, algorithms dissect intricate patterns inherent within input data, thus facilitating the prediction of continuous outcomes. This predictive capability extends its reach across diverse domains, encompassing realms such as sales projections, stock market prognostication, and medical diagnosis grounded in comprehensive patient data analysis.

There are many types of Regression in ML:

1. Linear Regression
2. Logistic Regression
3. Ridge Regression
4. Lasso Regression

## 1.2.2 Unsupervised Machine Learning

Unsupervised learning, a pivotal facet of machine learning, involves algorithms discerning patterns and structures from unlabelled data. This denotes that the training dataset lacks corresponding output labels, distinguishing it from supervised learning, which relies on labelled examples. In contrast to supervised learning, where algorithms are guided by labelled data, unsupervised learning mandates algorithms to autonomously uncover inherent structures within input data. Through this process, algorithms delve into the intricate web of data, identifying patterns and organizing information without external guidance. This inherent adaptability positions unsupervised learning as a powerful tool for uncovering latent insights across various domains and applications.

**Algorithms in Unsupervised Learning:**

1. **Clustering Algorithms:**

   K-means Clustering methodically partitions data into predetermined clusters, wherein each cluster embodies a set of data points sharing similar traits.

   Hierarchical Clustering constructs a hierarchical arrangement of clusters, progressively amalgamating data points into clusters grounded on their likenesses.

   DBSCAN (Density-Based Spatial Clustering of Applications with Noise) adeptly pinpoints clusters of varying configurations and densities in data, meticulously distinguishing core points, border points, and noise points.

2. **Dimensionality Reduction Techniques:**

   Principal Component Analysis (PCA) adeptly diminishes data dimensionality by transmuting it into a lower-dimensional realm while conserving the lion's share of variability.

   t-Distributed Stochastic Neighbour Embedding (t-SNE) adeptly portrays high-dimensional data in a lower-dimensional space, accentuating localized relationships among data points.

Autoencoders, sophisticated neural network architectures, adeptly assimilate and reconstruct data, encapsulating its underlying structure in a lower-dimensional latent space.

Association Rule Mining: The Apriori Algorithm adroitly ferrets out frequent item sets in transactional data and extrapolates association rules to elucidate item interrelations.

3. **Anomaly Detection:**

The Isolation Forest proficiently identifies anomalies or aberrations in data by segregating them in a configuration akin to a random forest, where anomalies are segregated with fewer partitions.

4. **Generative Models:**

Variational Autoencoders (VAEs) proficiently acquire the knack of generating fresh data samples by encoding input data into a latent space and subsequently decoding it into the original data realm.

Generative Adversarial Networks (GANs) seamlessly integrate two neural networks, a generator and a discriminator, engaged in a competition to fashion lifelike data samples.

5. **Self-Organizing Maps (SOMs):**

SOMs ingeniously organize high-dimensional data into a low-dimensional grid, thereby conserving topological relationships between data points.

## 1.3 **Introduction to Web phishing:**

Web phishing refers to the fraudulent practice of attempting to acquire sensitive data like usernames, passwords, or credit card information by posing as a trustworthy source in electronic communications. This deceitful tactic is commonly executed through emails, instant messages, or deceptive websites that appear legitimate. The primary aim is to trick recipients into divulging confidential details, often by directing them to counterfeit websites that imitate reputable organizations or services. Phishing schemes can be quite sophisticated, employing various strategies such as crafting fake login portals, leveraging psychological manipulation techniques, and instilling a sense of urgency or fear to compel immediate responses. To guard against phishing, it's crucial for individuals and businesses to exercise vigilance and caution when dealing with unsolicited messages or sharing personal information online. Furthermore, implementing security measures like email filters, anti-phishing software, and educational initiatives can help reduce the likelihood of falling prey to phishing scams.

**Why to Detect Web Phishing:**

Detecting web phishing holds significant importance for several reasons:

1. **Preservation of Personal Data:** Phishing attacks are designed to pilfer sensitive information such as login credentials, financial particulars, and personal data. Detecting these nefarious attempts acts as a barrier against unauthorized access to such information.
2. **Mitigation of Identity Theft**: Phishing endeavours frequently culminate in identity theft, wherein perpetrators exploit pilfered information to assume the identities of victims or engage in fraudulent activities. Early identification serves as a deterrent against these malevolent actors from capitalizing on personal identities.
3. **Upholding Trust**: For enterprises and establishments, the detection and prevention of phishing attacks are pivotal in upholding trust with clientele and stakeholders. Succumbing to phishing can inflict reputational harm upon a company and diminish trust in its online services.
4. **Prevention of Financial Setbacks:** Phishing assaults can yield financial setbacks for individuals and entities alike, encompassing deceitful transactions, unauthorized account access, and ransomware payouts. Timely detection of phishing endeavors aids in circumventing these financial setbacks.
5. **Protection of Data:** Phishing schemes may also serve as conduits for disseminating malware, ransomware, or other malicious software onto victims' devices. Detecting phishing attempts serves as a bulwark against these perils and shields sensitive data.
6. **Adherence to Regulatory Mandates:** Across numerous sectors, regulatory mandates exist concerning the safeguarding of personal and financial information. Identifying and thwarting phishing attacks aids organizations in meeting these mandates and evading potential repercussions.

In summation, the detection of web phishing is indispensable for safeguarding personal information, thwarting identity theft and financial losses, upholding trust, preserving data integrity, and ensuring compliance with regulations.

## 1.3.1 **Machine Learning in Web phishing Detection:**

The realm of web phishing detection has witnessed a significant evolution thanks to the emergence of artificial intelligence (AI) and machine learning (ML) technologies. These advancements have ushered in a new era of more precise and efficient prediction and identification of phishing attempts across diverse online platforms.

AI/ML algorithms are pivotal in the fight against web phishing, leveraging historical data to construct predictive models aimed at spotting potential phishing attacks. Through real-time analysis of vast datasets, these algorithms swiftly pinpoint patterns and irregularities that may indicate ongoing phishing endeavors.

The importance of AI/ML in web phishing detection is underscored by its capacity to:

1. Process enormous volumes of data in real-time, facilitating prompt detection and mitigation of phishing threats.
2. Accurately discern and categorize phishing attempts, thereby refining the accuracy of detection mechanisms.
3. Streamline the detection process through automation, diminishing the need for manual scrutiny and enhancing overall operational efficiency.
4. Interpret intricate data and visuals, fostering deeper comprehension and recognition of phishing patterns.
5. Enhance precision by integrating diverse data sources, encompassing elements like website attributes, user actions, and network traffic, thereby yielding more dependable detection outcomes.

In essence, AI/ML algorithms are reshaping the landscape of web phishing detection by furnishing advanced capabilities for scrutinizing, interpreting, and prognosticating phishing attempts.

Using Algorithms for classification Such as:

1. Random Forest
2. KNN
3. Decision Trees
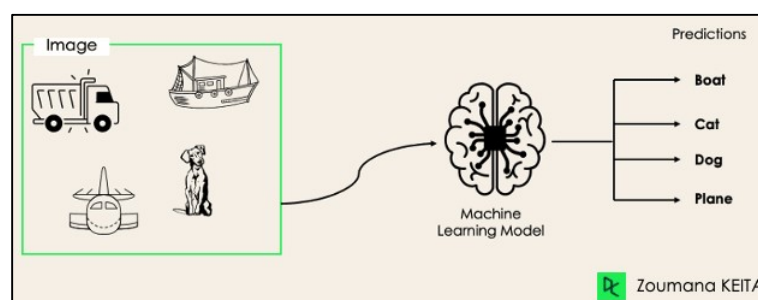
# Chapter-2
# Literature Review

## 2.1 Machine Learning

Machine learning is a subfield of artificial intelligence (AI) that allows computers to perform better without explicit programming by learning from data. It entails the creation of statistical models and algorithms that analyze and comprehend sizable datasets in order to spot trends and base judgments and predictions on the information. supervised learning, unsupervised learning, and reinforcement learning are a few of the different kinds of machine learning techniques. In supervised learning, algorithms are trained on labeled data; in unsupervised learning, patterns are found in unlabeled data; and in reinforcement learning, algorithms are taught to interact with an environment in order to gradually reach a desired objective.

Machine learning finds applications in many different industries, such as autonomous cars, finance, healthcare, natural language processing, recommendation systems, and image and speech recognition. Its ability to gather information, draw conclusions from it, and make rational decisions has revolutionized many industries, fostering productivity and innovation and providing the means for the complex problem-solving of the modern world.

## 2.2 Classification

In machine learning, classification is the process of grouping data points according to their characteristics or attributes into predefined classes or categories. Using labeled training data, the algorithm is trained to classify new, unseen instances into one of the predefined classes using supervised learning.



In order to teach the algorithm to understand the relationship between the input features and the relevant classes, the algorithm is trained on a dataset in which each data point is assigned a class label. Neural networks, k-nearest neighbours (KNN), logistic regression, decision trees, and support vector machines (SVM) are examples of common categorization techniques.

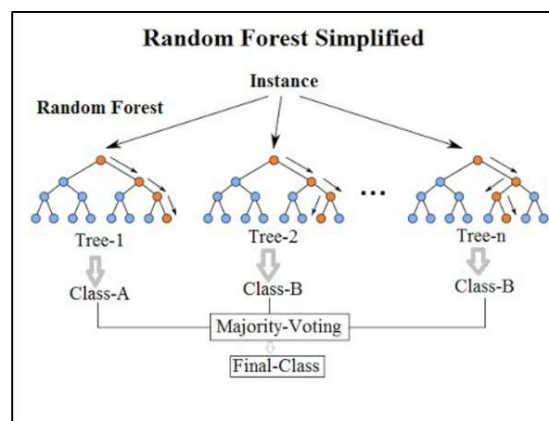Many Algorithms can be used for classification purposes in Machine Learning:

1. KNN(K nearest neighbours)

2. Random Forest
3. Decision Trees
4. Support Vector machine
And many more

Classification is widely used in many fields, such as speech recognition, image identification, credit scoring, medical diagnosis, sentiment analysis in social media, and spam detection in emails. In many real-world contexts, it facilitates activities like prediction, identification, and categorization by enabling automated decision-making based on incoming data. Classification algorithms are extremely useful tools for resolving classification issues in a variety of fields and sectors due to their accuracy and efficiency.

## 2.3 Random Forest

For both classification and regression applications, Random Forest is a flexible machine learning technique that is frequently employed. When it is being trained, it builds a large number of decision trees and outputs the mean prediction (for regression) or mode (for classification) of each tree.



Because each Random Forest decision tree is constructed using a random collection of features and a portion of training data, generalization performance is improved and overfitting is decreased. Because the algorithm aggregates the forecasts of each individual tree, it is able to make predictions; hence, the word "forest."

Random forests can handle huge datasets with high complexity and noisy input because they are very adaptable and resilient. When compared to individual decision trees, they are less likely to overfit and usually provide excellent accuracy and stability for a variety of tasks.
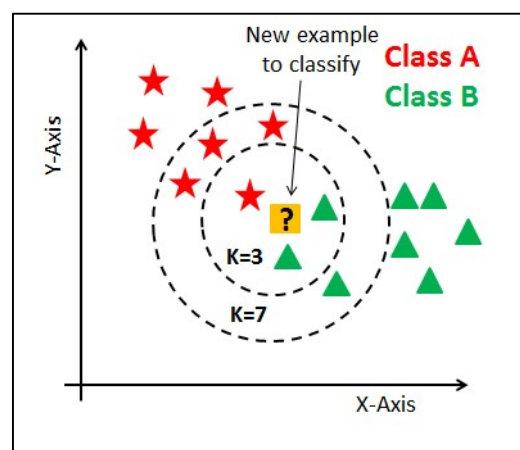
Random Forests are useful in a variety of industries, including bioinformatics, marketing, healthcare, and finance, because of their simplicity and efficacy. When predicted accuracy is crucial and interpretability is not the main goal, they work especially well in classification or regression tasks.

## 2.3 K Nearest Neighbours (KNN)

For problems involving regression and classification, K-Nearest Neighbors (KNN) is a straightforward yet effective machine learning technique. In a KNN, a new data point's prediction is based on how close it is to nearby data points in the feature space.

Without a formal model training phase, the algorithm operates by keeping the whole training dataset in memory. KNN computes the distance between each new data point and every other point in the training dataset before generating predictions for that point. Then, using a distance metric—such as the Manhattan distance or the Euclidean distance—it chooses the K closest neighbors.

KNN selects the class label that is most prevalent among its K nearest neighbors for classification tasks. It calculates the average (or weighted average) of the K nearest neighbor target values in regression tasks.
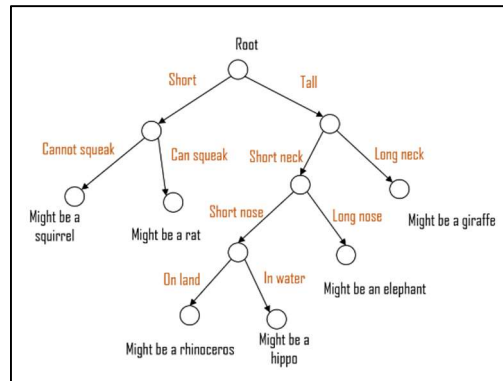


KNN is simple to comprehend and apply thanks to its intuitive approach and simplicity. Nevertheless, the distance metric and K selection may have an impact on how well it performs. Furthermore, KNN can be computationally costly, particularly when working with big datasets.

Applications for KNN can be found in many different fields, such as anomaly detection, picture classification, recommendation systems, and pattern recognition. It is especially helpful when working with non-linear decision limits or when the underlying data distribution is poorly understood.

## 2.4 Decision Trees

For both classification and regression tasks, a decision tree is a widely used machine learning technique. In order to create a tree-like model that can predict the target variable, it operates by recursively partitioning the data into subsets based on the values of input attributes.

Nodes, which stand for features, and edges, which stand for choices or results based on those choices, make up the tree structure. Based on factors like information gain (for classification) or variance reduction (for regression), the algorithm chooses the characteristic at each node that divides the data into more homogeneous subsets the best.



Choice Trees are helpful for understanding how the model generates predictions since they are simple to view and interpret. They can, however, overfit, particularly when dealing with deep trees that pick up noise in the data. Methods such as pruning, restricting the maximum depth of the tree, or establishing a minimum number of samples per leaf node can be used to reduce overfitting.

Decision trees are used in many different fields, such as customer relationship management, marketing, finance, and healthcare. where working with categorical data or where interpretability is crucial to comprehending the decision-making process, they are especially helpful. Furthermore, Decision Trees can be used as the building blocks of more intricate ensemble techniques such as Gradient Boosting Machines and Random Forests.

## 2.4 Previous Studies on Web Phishing detection in Machine Learning

- On October 2018, a study called **Phishing Website Detection using Machine Learning Algorithms** by **Rishikesh Mahajan** and **Irfan Siddavatam** on Web phishing detection by their urls. The study focuses on The paper focuses on using machine learning algorithms to detect phishing websites by analyzing features extracted from legitimate and phishing URLs. 16 different features were extracted from the URLs such as presence of IP address, number of dots, dashes, redirects, sensitive words, URL length, etc. Random Forest gave the best performance with 97.14% accuracy and lowest false positive/negative rates when trained on 90% of the data.
- On August 2022, s study called **Phishing URL detection using machine learning methods** by **SK Hasane Ahammad , Sunil D. Kale , Gopal D. Upadhye , Sandeep Dwarkanath Pande , E Venkatesh Babu , Amol V. Dhumane , Mr. Dilip Kumar Jang Bahadur.** This paper focuses on detecting phishing URLs (malicious websites designed to steal user data) using machine learning techniques. The authors collected a dataset of 3000 URLs, with 1500 malicious and 1500 benign URLs. They extracted 15 different features from the URLs related to the address bar, domain, and other characteristics. The study gave following results: Decision Tree= 85%, Logistic Regression=86%, SVM = 84%.

- On November 2021, a study called **Detecting Phishing Websites Using Machine Learning** **by Aniket Garje, Namrata Tanwani, Sammed Kandale , Twinkle Zope , Prof. Sandeep Gore.** The paper discusses using machine learning techniques to detect phishing websites, which are a major cybersecurity threat aimed at stealing personal information like passwords and credit card numbers. It provides an overview of phishing and cybercrime, and does a literature survey of existing work on using machine learning for phishing detection. The results show that the Decision Tree algorithm performs best, with an F1-score of 0.94 and good precision-recall balance.

- **A Survey of Machine Learning-Based Solutions for Phishing Website Detection L. Tang, Q. Mahmoud [22]:** The proposed approach in the current study uses URLs collected from a variety of platforms, including Kaggle, Phish Storm, Phish Tank, and ISCX-UR, to identify phishing websites. The researchers made a big contribution since they created a browser plug-in that can quickly recognize phishing risks and offer warnings. Various datasets and machine learning techniques were investigated, and the proposed RNN-GRU model outperformed SVM, Random Forest (RF), and Logistic Regression with a maximum accuracy rate of 99.18%. On the other hand, the suggested method is not always accurate in identifying if short links are phishing risks.

# Chapter-3
# Methodology

**Problem Statement:** To classify the malicious websites using their URLs by training from the given data set.

## 3.1 Dataset Description

### Overview

The dataset utilized in this research project on "Web Phishing Detection Using Machine Learning" comprises 11,430 instances with 89 features each. The features encompass a diverse range of attributes associated with URLs and webpage characteristics, aimed at discerning potential phishing attempts.

### Features

**length_url:** Represents the length of the URL.

**length_hostname:** Indicates the length of the hostname within the URL.

**ip:** Binary feature denoting whether the URL contains an IP address.

**nb_dots:** Quantifies the number of dots present in the URL.

**nb_qm:** Records the count of question marks within the URL.

**nb_eq:** Reflects the number of equal signs present in the URL.

**nb_slash:** Measures the frequency of slashes in the URL.

**nb_www:** Binary indicator representing the presence of 'www' in the URL.

**ratio_digits_url:** Proportion of digits within the URL.

**ratio_digits_host:** Proportion of digits within the hostname.

**tld_in_subdomain:** Binary flag indicating whether the top-level domain is within the subdomain.

**prefix_suffix:** Binary feature signaling the presence of a prefix or suffix in the URL.

**shortest_word_host:** Length of the shortest word present in the hostname.

**longest_words_raw:** Length of the longest word within the URL.

**longest_word_path:** Length of the longest word within the URL path.

**phish_hints:** Binary indicator of phishing hints within the URL.

**nb_hyperlinks:** Count of hyperlinks present in the webpage.

**ratio_intHyperlinks:** Ratio of internal hyperlinks within the webpage.

**empty_title:** Binary flag indicating the absence of a title in the webpage.

**domain_in_title:** Binary indicator of whether the domain is included in the webpage title.

**domain_age:** Age of the domain hosting the webpage.

**google_index:** Binary feature representing whether the webpage is indexed by Google.

**page_rank:** Page rank of the webpage.

## 3.2 Data Preprocessing

Data preprocessing is a critical step in preparing the dataset for machine learning model training. In this research, the primary focus of data preprocessing was on handling missing values. The dataset underwent a straightforward approach of null point removal to ensure the integrity and quality of the data.

1. **Null Point Removal:**
- Missing values, if any, were identified across all features of the dataset.
- Instances containing missing values were removed from the dataset to maintain data consistency.
- This approach ensured that only complete and reliable data points were used for model training, thus minimizing the risk of bias or inaccuracies introduced by missing values.

2. **Feature Extraction**

In addition to data preprocessing, feature extraction plays a crucial role in enhancing the efficiency and effectiveness of machine learning models. In this research, feature extraction utilizing correlation analysis was employed to identify and select the most relevant features for model training.

- Correlation Analysis:

Correlation analysis was conducted to quantify the relationship between each feature and the target variable, 'status,' indicating whether a URL is classified as phishing or not.

Pearson correlation coefficient, Spearman rank correlation coefficient, or other appropriate measures were calculated to assess the linear or monotonic relationships between features and the target variable.

- Feature Selection Criteria:

Features exhibiting high correlation coefficients with the target variable were considered indicative of significant predictive power for phishing detection.

A threshold correlation value was established to determine the inclusion or exclusion of features in the final feature set. Features surpassing this threshold were retained for subsequent model training, while features below the threshold were discarded.

- Redundant Feature Removal:

In cases where multiple features were highly correlated with each other (multicollinearity), redundant features were identified and removed to prevent model overfitting and enhance interpretability.

Techniques such as variance inflation factor (VIF) analysis or principal component analysis (PCA) may be employed to identify and eliminate redundant features while preserving the essential information captured by the dataset.

- Final Feature Set:

The feature extraction process resulted in a refined set of features characterized by their strong correlation with the target variable, thereby improving the discriminative power of the machine learning models for phishing detection.

## 3.3 Model Training and Testing

Model training and testing are pivotal stages in evaluating the performance of machine learning algorithms for web phishing detection. In this research, four diverse algorithms, namely Random Forest, K-Nearest Neighbors (KNN), Logistic Regression, and Decision Tree, were trained and evaluated using appropriate methodologies.

- **Data Splitting:**

The dataset was divided into two subsets: a training set and a testing set.

Typically, a commonly used split ratio such as 70-30 or 80-20 was employed, where the majority of the data (e.g., 70% or 80%) was used for training, and the remaining portion was reserved for testing.

- **Model Training:**

Each machine learning algorithm (Random Forest, KNN, Logistic Regression, Decision Tree) was trained on the training set using the selected features obtained from the feature extraction process.

The algorithms were provided with labeled data, where the features served as inputs, and the corresponding 'status' label (indicating phishing or legitimate) served as the target variable for supervised learning.

- **Parameter Tuning (Optional):**

Hyperparameter tuning, such as adjusting the number of trees in Random Forest, the number of neighbors in KNN, or the regularization parameter in Logistic Regression, may be performed using techniques like grid search or randomized search to optimize model performance.

## 3.3 Model Evaluation

Once trained, each model was evaluated on the testing set to assess its performance in accurately classifying phishing attempts.

Evaluation metrics such as accuracy, precision, recall, F1-score, and area under the receiver operating characteristic (ROC) curve were computed to quantify the model's performance in detecting phishing attempts.

Cross-Validation (Optional):

To ensure robustness and mitigate overfitting, k-fold cross-validation may be employed, where the dataset is partitioned into k subsets, and the model is trained and evaluated k times, each time using a different subset as the testing set and the remaining data as the training set.

Model Comparison:

The performance of each machine learning algorithm was compared based on the evaluation metrics obtained during testing.

Insights gained from the comparison facilitated the identification of the most effective algorithm for web phishing detection based on the characteristics of the dataset and the requirements of the research objectives.

By following these steps, the research aimed to comprehensively evaluate the efficacy of Random Forest, KNN, Logistic Regression, and Decision Tree algorithms in detecting web phishing attempts, providing valuable insights into the strengths and limitations of each approach.

# Chapter-4
# Results And Discussions

The evaluation of machine learning algorithms for web phishing detection yielded insightful results, shedding light on their effectiveness and performance in discerning between legitimate URLs and phishing attempts.

| Model | Accuracy |
|---|---|
| Random Forest | 97% |
| KNN | 95% |
| Logistic Regression | 92% |
| Decision Tree | 96% |

**Performance Metrics:**

- Random Forest achieved an accuracy of 97%, showcasing its robustness in classifying URLs accurately.
- K-Nearest Neighbors (KNN) exhibited a commendable accuracy of 95%, indicating its efficacy in proximity-based classification.
- Logistic Regression demonstrated a respectable accuracy of 92%, highlighting its simplicity and interpretability in modelling.
- Decision Tree attained an accuracy of 96%, showcasing its ability to capture decision rules for classification.

**Model Comparison:**

- Among the tested algorithms, Random Forest Classifier emerged as the top performer, boasting the highest accuracy of 97%. Its ability to handle high-dimensional data and capture complex interactions between features proved highly effective in identifying phishing attempts.
- Decision Tree Classifier followed closely behind with an accuracy of 96%, demonstrating strong performance but with a potential tendency toward overfitting, which may affect generalization to unseen data.
- K-Nearest Neighbors (KNN) achieved an accuracy of 95%, effectively classifying URLs based on their similarity to neighboring instances.
- Logistic Regression, while slightly trailing with an accuracy of 92%, offered transparency in model interpretation, making it suitable for scenarios where interpretability is crucial.

In conclusion, the results of this study provide valuable insights into the effectiveness of machine learning algorithms for web phishing detection, offering guidance for the development of robust and adaptive cybersecurity solutions.

# Chapter-5
# Future Aspects

- **Ensemble Methods Integration:**

Future research could explore the integration of ensemble learning techniques such as bagging, boosting, or stacking to further enhance the performance of phishing detection models. Ensemble methods combine multiple base learners to improve prediction accuracy and generalization capabilities, offering potential advancements in detecting sophisticated phishing attempts.

- **Deep Learning Architectures:**

Investigating deep learning architectures, such as convolutional neural networks (CNNs) or recurrent neural networks (RNNs), presents an avenue for enhancing the detection of complex patterns and subtle indicators of phishing activities. Deep learning models can automatically learn hierarchical representations from raw data, potentially leading to more nuanced and accurate phishing detection systems.

- **Dynamic Feature Engineering:**

Continuous refinement and expansion of feature engineering techniques are essential to adapt to evolving phishing tactics. Incorporating dynamic features derived from real-time web traffic analysis, user behaviour patterns, or domain reputation scores could improve the responsiveness and adaptability of phishing detection models to emerging threats.

- **Adversarial Attack Resilience:**

Addressing the challenge of adversarial attacks against phishing detection systems is crucial for ensuring robust cybersecurity. Future research could focus on developing adversarial robust models capable of withstanding deliberate manipulations or evasions attempted by sophisticated adversaries aiming to circumvent detection mechanisms.

- **Real-Time Deployment and Feedback Loop:**

Implementing real-time deployment of phishing detection systems integrated with a feedback loop mechanism enables continuous learning and refinement based on real-world feedback. This iterative process allows the model to adapt to changing phishing tactics and improve its effectiveness over time, enhancing overall cybersecurity posture.

- **Cross-Domain Generalization:**

Investigating the generalization of phishing detection models across different domains and contexts is essential for scalability and applicability in diverse environments. Future research should explore techniques for cross-domain transfer learning or domain adaptation to ensure the robustness and versatility of detection models across varied scenarios.

Addressing these future aspects will contribute to the advancement of phishing detection technology, ultimately strengthening cybersecurity defences and mitigating the risks posed by malicious actors in the ever-evolving digital landscape.

# Chapter-6
# Conclusion

In conclusion, this research has demonstrated the effectiveness of machine learning algorithms in detecting web phishing attempts using a comprehensive dataset of URL and webpage attributes. The high accuracies achieved by K-Nearest Neighbors, Logistic Regression, and Random Forest models emphasize the potential of these techniques in bolstering cybersecurity measures against phishing attacks.

The study underscores the importance of feature selection and dataset quality in training robust phishing detection models. By leveraging diverse features such as URL structure, webpage content, and domain characteristics, the models were able to discern subtle indicators of phishing attempts with notable accuracy.

Moving forward, further advancements in machine learning techniques, coupled with continuous updates to datasets reflecting evolving phishing tactics, are essential to stay ahead of cyber threats. Additionally, the integration of real-time monitoring and adaptive learning mechanisms could enhance the agility and responsiveness of phishing detection systems.

Overall, this research contributes to the ongoing efforts in cybersecurity by providing insights into the efficacy of machine learning in combating web-based phishing attacks. By leveraging the power of data-driven approaches, organizations can fortify their defences and mitigate the risks posed by malicious actors in the digital landscape.

**References:**

1. Mahajan, Rishikesh & Siddavatam, Irfan. (2018). Phishing Website Detection using Machine Learning Algorithms. International Journal of Computer Applications. 181. 45-47. 10.5120/ijca2018918026.

2. Aniket Garje,Namrata Tanwani,Sammed Kandale,Twinkle Zope,Prof. Sandeep Gore, "DETECTING PHISHING WEBSITES USING MACHINE LEARNING", International Journal of Creative Research Thoughts (IJCRT), ISSN:2320-2882, Volume.9, Issue 11, pp.243-246, November 2021, Available at :http://www.ijcrt.org/papers/IJCRTI020051.pdf

3. Tang, Lizhen & Mahmoud, Qusay. (2021). A Survey of Machine Learning-Based Solutions for Phishing Website Detection. Machine Learning and Knowledge Extraction. 3. 672-694. 10.3390/make3030034.

4. https://www.sciencedirect.com/science/article/abs/pii/S0965997822001892