# DETECTING EMAIL-BASED PHISHING WEBSITES USING MACHINE LEARNING

2023-123

# Team Members

J.M.Lindamulage             Mandira Pabasari L.              Yapa S.P.J.                 Perera I.S.S.

# Research Question

How to detect phishing websites and phishing emails using machine learning and deep learning.
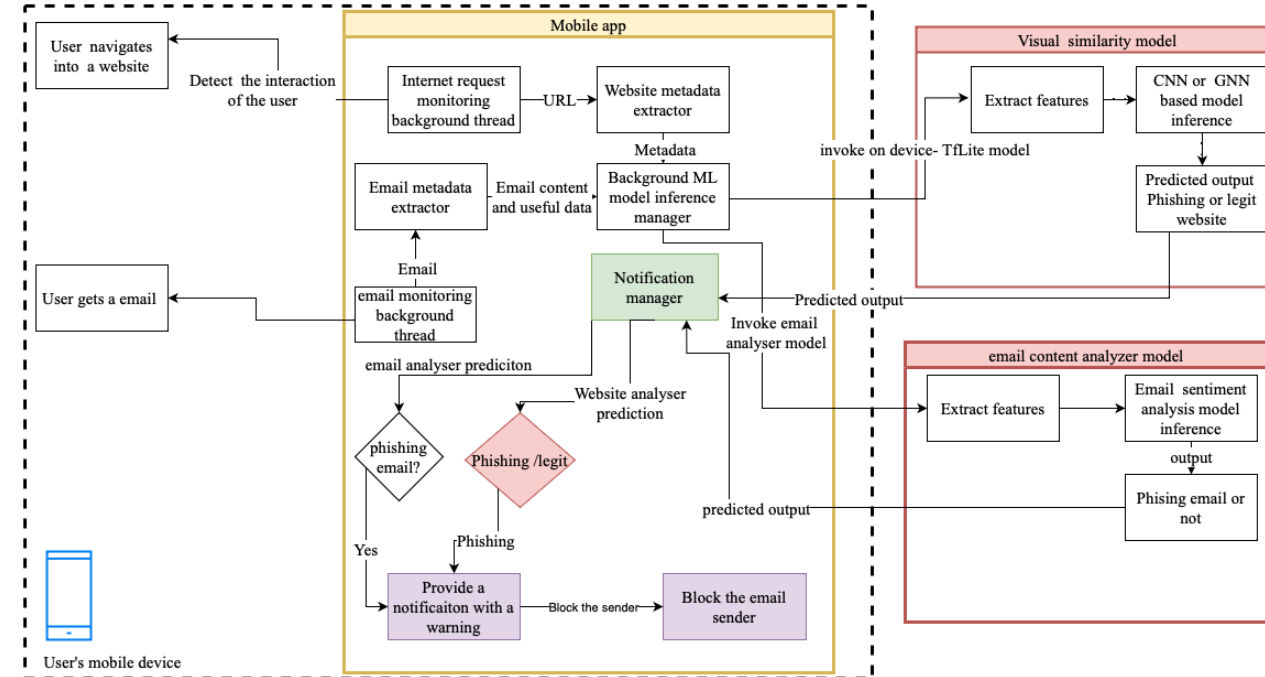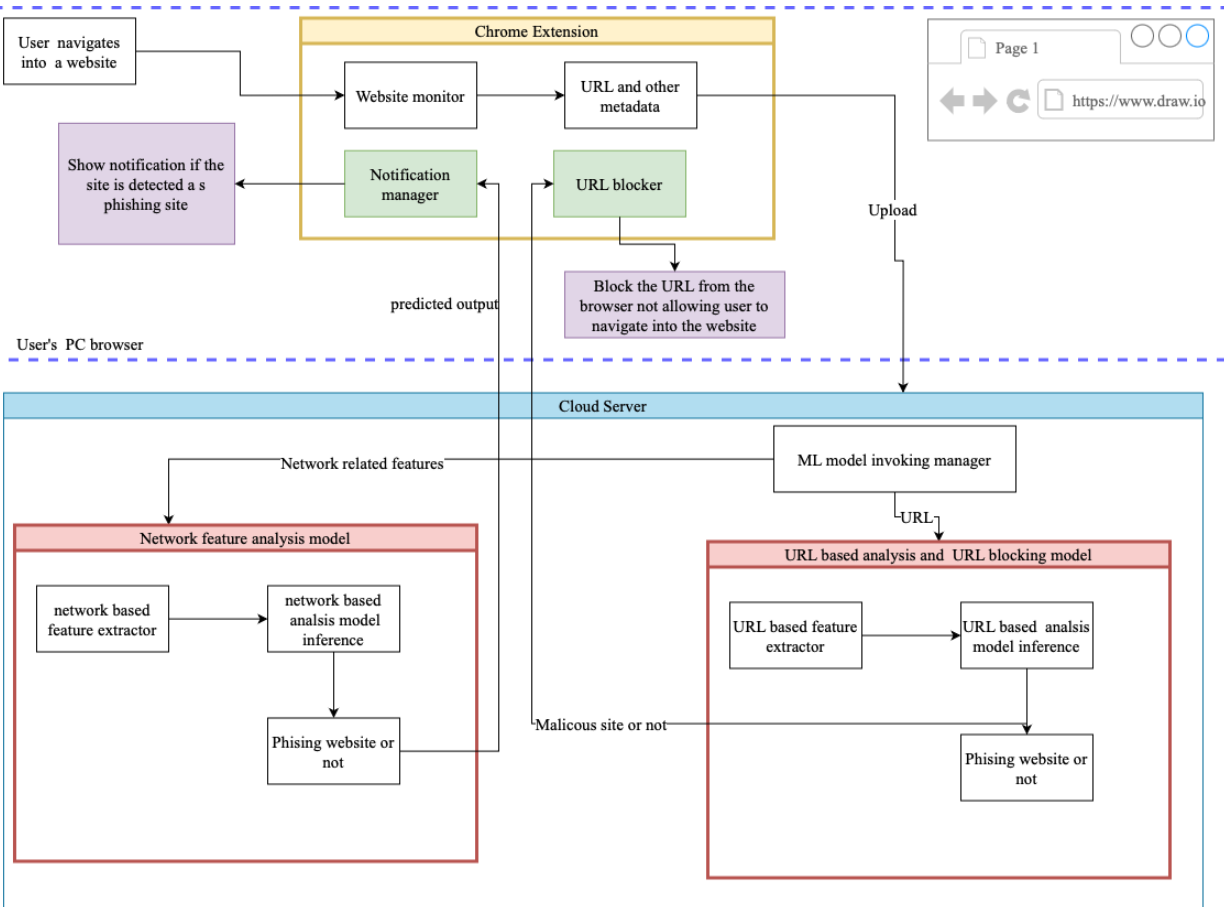
# Main Objective

To implement a mobile application and a web extension capable of detecting phishing emails and websites utilizing machine learning and deep learning models.

SLIIT
FACULTY OF COMPUTING

# Sub-objectives

- To employ the visual similarity features to classify phishing websites out of legit websites.

- Detecting Phishing sites using website feature analysis

- To identify  phishing emails using the heading and the textual content in the email.

- To discover phishing websites using the URL.

# system diagram

# IT20222840 |Judith Malshini L.

BSc (Hons) Degree in Information Technology (specialization in Cyber Security )

SLIIT
FACULTY OF COMPUTING

# Research Problem

How to use visual similarity features of a website to detect if a website is a phishing website or a legitimate website

# Objectives

- To utilize Vision GNN for the first time for classifying phishing websites.
  - Vision GNN: An Image is Worth Graph of Nodes(2022)
  - Was trained on ImageNet dataset.
- To optimize the developed model for mobile devices.
- To deploy the deep learning model on android mobile app.
- How to run the prediction on the edge device itself without running it on cloud server

# Contributions

Introduced VisionGNN architecture based on Graph neural networks into phishing website classification for the first time

Utilized visual features alone with graph neural network representations for the first time.

Tuned hyperparameters to obtain the best accuracy.

Converted implemented PyTorch model into PyTorch mobile version to be deployed on Android mobile app.

Implemented a mobile app to detect phishing websites using a screenshot of the page.

Integrated deep learning model into the mobile app as a PyTorch mobile model to done the prediction on the edge device itself.

SLIIT
FACULTY OF COMPUTING

# Graph Neural Network Concept and Vision GNN

Propose the image as a graph structure and introduce a new Vision GNN (ViG) architecture to extract graph level feature for visual tasks.

- Split the image to several patches which are viewed as nodes.
- Patches are transferred into the feature vector.



Fig 1. Image to graph construction (source-VisionGNN)

[20]    Kai Han1,2∗ Yunhe Wang2∗ Jianyuan Guo2 Yehui Tang2,3 Enhua Wu1,4, "Vision GNN: An Image is Worth Graph of Nodes".

# Dataset

- Visual phish dataset
  - Contain 9363 screenshots of PhishTank phishing pages that target 155 websites and 1195 phishing pages.

- 4072 data were divided into train and test with 20% split.
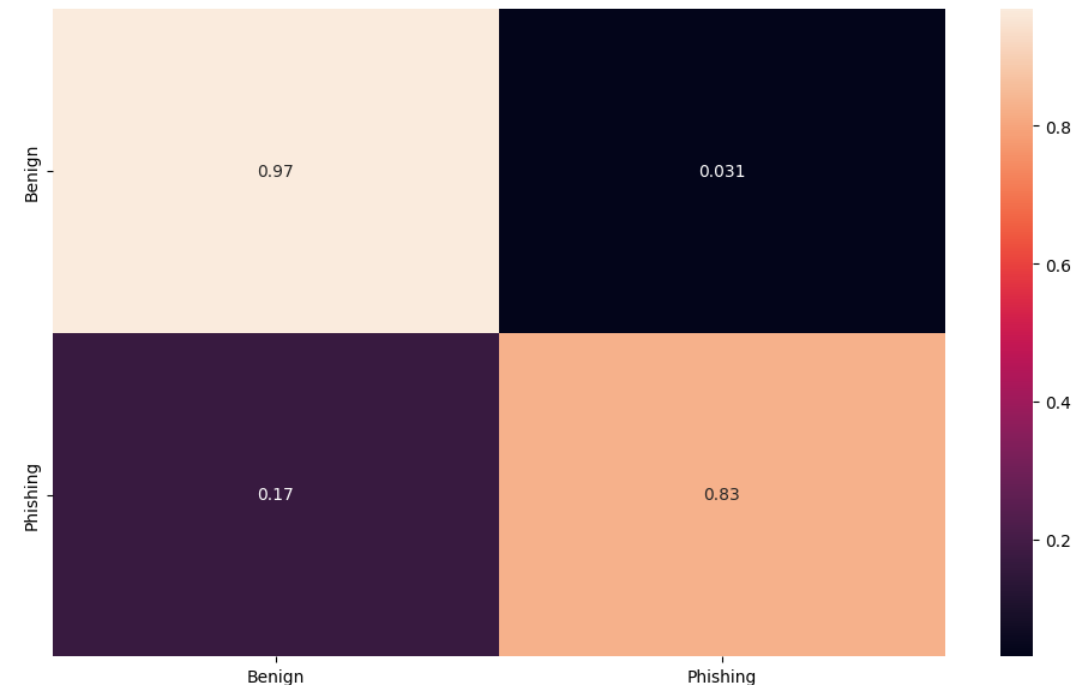
- Train data-3054

- Test data-1018

# Model results

| Model | No of parameters($10^6$) | Accuracy(100) |
|---|---|---|
| Tiny | 9.69 | 93.5 |
| Small | 26.23 | 97.4 |
| Medium | 48.50 | 91.8 |
| Large | 91.96 | 93.5 |

Maxium accuracy 97.4% in small model
with 26.23 x $10^6$ parameters

83% phishing sites detected
97% beign sites are correctly classified

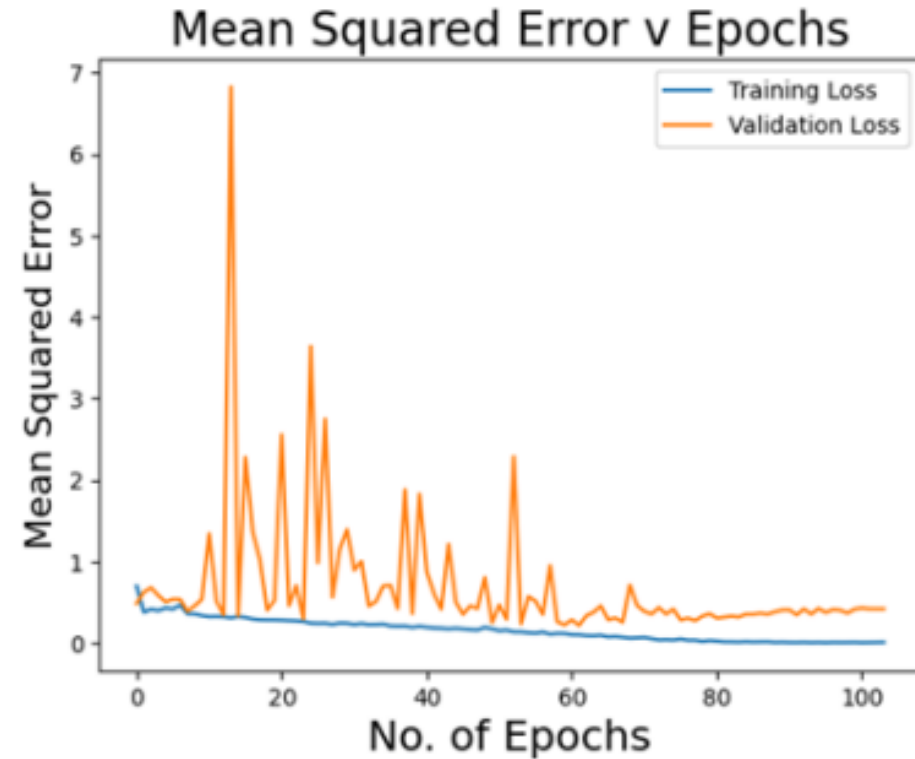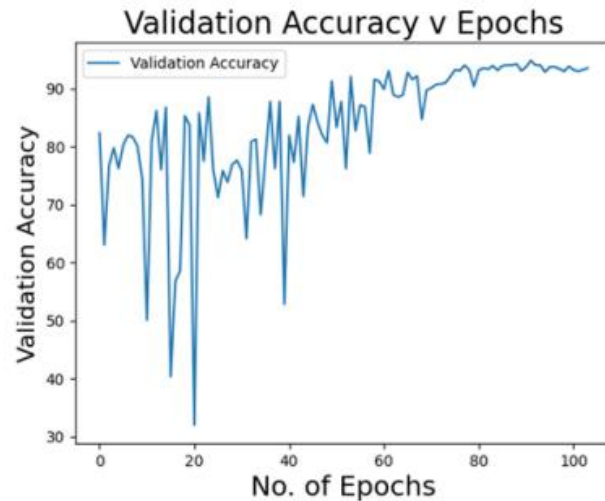# VisionGNN for phishing website detection model training

Train epochs -100

Batch size=64

Optimizer – AdamW
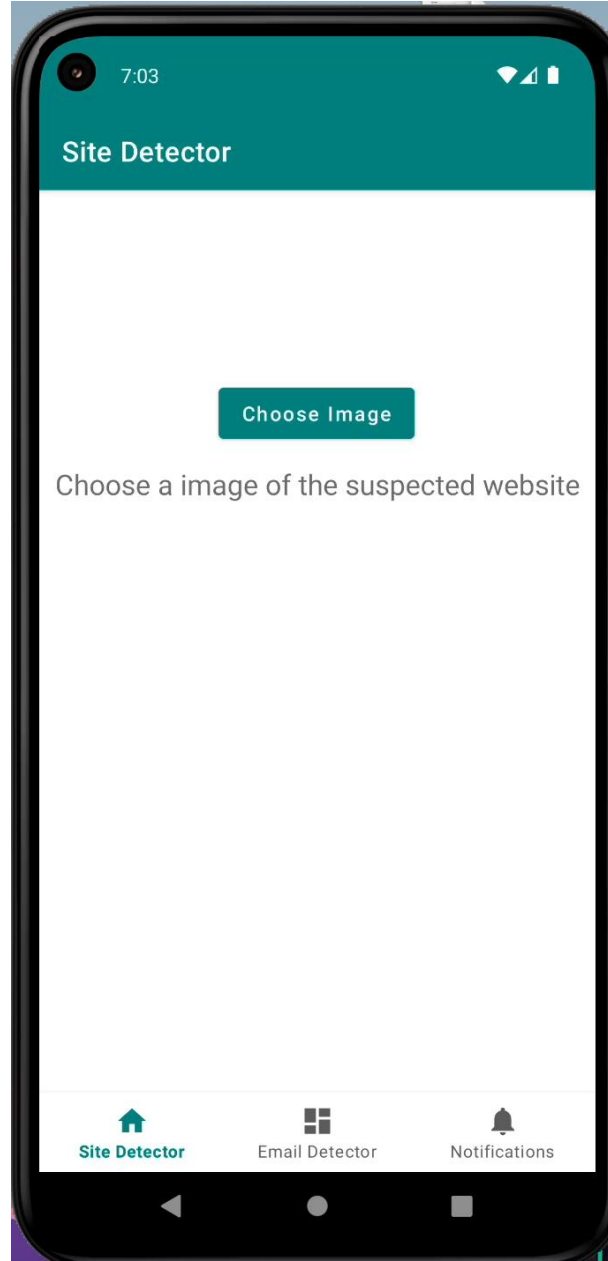
Learning rate – 0.02

Learning rate decay – cosine

Image augmentations – horizontal and verital flips

# Model mobile app integration

- Optimized model for edge computing using Pytorch mobile

- Converted Pytorch model into Torchscript

- Model works on Andorid (works on Java)

- Model runs on GPU if the mobile device has a  compatible GPU

- Advantages:
  - Preserves privacy of the user
  - Saves network bandwidth
  - Low latency

# Mobile app demo

# REFERENCES

| [1] | **Aya Hashim∗Razan Medani, Dr.Tahani Abdalla Attia, "Defences Against web Application Attacks and Detecting Phishing Links Using Machine Learning," in 2020 International Conference on Computer, Control, Electrical, and Electronics Engineering (ICCCEEE), Sudan, 2020.** |
|---|---|
| [2] | "PHISHING ACTIVITY TRENDS REPORT 3rd Quarter," APWG, 2022. |
| [3] | Dr. Moulana Mohammed,K. Koteswara Prasanth,S. Venkata Sai Subhash , "PHISHING DETECTION USING MACHINE LEARNING ALGORITHMS," in Proceedings of the Fourth International Conference on Smart Systems and Inventive Technology (ICSSIT-2022), India, 2022. |
| [4] | A. Lakshmanarao , P.Surya Prabhakara Rao,M M Bala Krishna, "Phishing website detection using novel machine learning fusion approach," in Proceedings of the International Conference on Artificial Intelligence and Smart Systems (ICAIS-2021), India, 2021. |
| [5] | Surbhi Gupta,Abhishek Singhal ,Akanksha Kapoor, "A Literature Survey on Social Engineering Attacks:Phishing Attack," in International Conference on Computing, Communication and Automation, India, 2016. |
| [6] | Ankit Kumar Jain and B. B. Gupta, "Phishing Detection: Analysis of Visual Similarity Based Approaches," WILEY, p. 2017, 2017. |
| [7] | R. Dhamija, J. D. Tygar, and M. A. Hearst, "Why phishing works," in SIGCHI Conference on Human, 2006. |
| [8] | Z. Fan, "Detecting and Classifying Phishing Websites by Machine Learning," in 2021 3rd International Conference on Applied Machine Learning (ICAML), china, 2021. |
| [9] | Malak Aljabri ,Samiha Mirza , "Phishing Attacks Detection using Machine Learning and Deep Learning Models," in 2022 7th International Conference on Data Science and Machine Learning Applications (CDMA) , Saudi Arabia, 2022. |
| [10] | "Phishing Website Detection Using Random Forest and Support Vector Machine: A Comparison," in 2nd International Conference on Artificial Intelligence and Data Sciences (AiDAS), Malaysia, 2021. |
| [11] | Jaydeep Solanki, Rupesh G. Vaishnav, "Website Phishing Detection using Heuristic Based Approach," International Research Journal of Engineering and Technology (IRJET), vol. 03, no. 05, 2016. |
| [12] | Eric Medvet,Engin Kirda,Christopher Kruegel, "Visual-Similarity-Based Phishing Detection," in SecureComm 2008 , Turkey, 2008. |
| [13] | Igino Corona1,2, Battista Biggio1,2, Matteo Contini2Luca Piras1,2, Roberto Corda2Mauro, Guido Mureddu2, "DeltaPhish: Detecting Phishing Webpages in Compromised Websites∗," in ESORICS 2017., Italy, 2017. |
| [14] | Sahar Abdelnabi,Katharina Krombholz,Mario Fritz, "VisualPhishNet: Zero-Day Phishing Website Detection by Visual Similarity," 2020. |
| [15] | U. Saeed, "Visual similarity-based phishing detection using deep learning," Journal of Electronic Imaging, vol. 31, no. 5, 2022. |
| [16] | Padmalochan Panda 1,†, Alekha Kumar Mishra 1,† and Deepak Puthal , "A Novel Logo Identification Technique for Logo-Based Phishing Detection in Cyber-Physical Systems," Future Internet, vol. 14, no. 8, 2022. |
| [17] | Saad Al-Ahmadi1 Yasser Alharbi, "DEEP LEARNING TECHNIQUE FOR WEB PHISHING DETECTION COMBINED URL FEATURES AND VISUAL SIMILARITY," International jpurnal of computer networks abd comunications(IJCNC), vol. 12, no. 5, 2020. |
| [18] | Rundong Yang,Kangfeng Zheng, Bin Wu,Chunhua Wu,Xiujuan Wang, "Phishing Website Detection Based on Deep Convolutional Neural Network and Random Forest Ensemble Learning," sensors, 2021. |
| [19] | "Phishing Detection: Analysis of Visual Similarity Based Approaches," WILEY, vol. 2017, 2017. |
| [20] | Kai Han1,2∗ Yunhe Wang2∗ Jianyuan Guo2 Yehui Tang2,3 Enhua Wu1,4, "Vision GNN: An Image is Worth Graph of Nodes". |

# IT19966236 | Mandira Pabasari L.

BSc (Hons) Degree in Information Technology (Specialization in Cyber Security)

# WEBSITE FEATURES



- **Attributes and characteristics of a web page.**

- **These attributes provide insights into a webpage's behavior, structure, and content.**

- **Are like pieces of a puzzle that make up a webpage.**

- **They include elements related to URLs, HTML, JavaScript, domains, and more.**

- **Think of website features as clues that help us understand a webpage's intent.**

- **They can reveal whether a webpage is legitimate or potentially malicious (phishing).**
  - Address Bar-based Features
  - Abnormal Behavior Features
  - HTML and JavaScript Features
  - Domain-related Features
  - Statistical Reports-based Features

# Current Progress

- Model Development

- Model Serialization

- Extension Integration

# Connection Flow:

```
┌─────────────────────────┐     ┌─────────────────────────┐     ┌─────────────────────────┐     ┌─────────────────────────┐
│ The training dataset is │ ──► │ The trained models are  │ ──► │ The Chrome browser      │ ──► │ The collected data is   │
│ preprocessed and used to│     │ converted to JSON format│     │ extension's content     │     │ passed to the predict   │
│ train machine learning  │     │ and saved as            │     │ script (content.js)     │     │ function in content.js, │
│ models in training.py.  │     │ 'classifier.json' using │     │ runs on web pages,      │     │ where a prediction is   │
│                         │     │ dump.py.                │     │ collecting data to      │     │ made based on the       │
│                         │     │                         │     │ perform phishing        │     │ trained model's weights.│
│                         │     │                         │     │ detection checks.       │     │                         │
└─────────────────────────┘     └─────────────────────────┘     └─────────────────────────┘     └─────────────────────────┘
```

Depending on the prediction result, a message is logged in the console.

The background script (background.js) listens for requests from content scripts.

When a request is received, it checks the prediction result and triggers an alert message in the Chrome browser based on the result.

SLIIT
FACULTY OF COMPUTING

# Future Works

- **User Interface**

- **Enhanced Features**

- **Hosting**

| Works Cited | |  |
|---|---|---|
| [1] | Smriti Dangwal,Arghir-Nicolae Moldovan, "Feature Selection for Machine Learning-based Phishing Websites Detection," in International Conference on Cyber Situational Awareness, Data Analytics and Assessment (CyberSA), 2021. | |
| [2] | Shatha Ghareeb,Mohamed Mahyoub,Jamila Mustafina, "Analysis of Feature Selection and Phishing Website Classification Using Machine Learning," in International Conference on Developments in eSystems Engineering, 2023. | |
| [3] | Smriti Dangwal,Arghir-Nicolae Moldovan, "Feature Selection for Machine Learning-based Phishing Websites Detection," in International Conference on Cyber Situational Awareness, Data Analytics and Assessment (CyberSA), 2021. | |
| [4] | E. Burns, "TechTarget," [Online]. Available: https://www.techtarget.com/searchenterpriseai/definition/machine-learning-ML. | |
| [5] | Rami M. Mohammad, Fadi Thabtah, and Lee McCluskey, "Phishing Websites Features". | |
| [6] | AJAY PRAKASH NAIR, DEVAPRASAD V,VISHNU PRASAD A, "CHROME EXTENSION FOR PHISHING WEBSITE DETECTION".<br> There are no sources in the current document. | |
| | | |
| | | |
| | | |
| | | |

SLIIT
FACULTY OF COMPUTING

# IT20050108 | S. P. J. Yapa

BSc (Hons) Degree in Information Technology (Specialization in Cyber Security)

Component 3 - Phishing email detection with sentiment analysis

# What is Phishing ?

Phishing is a cyberattack technique in which malicious actors impersonate legitimate entities to deceive individuals into revealing sensitive information or performing actions that compromise their security.
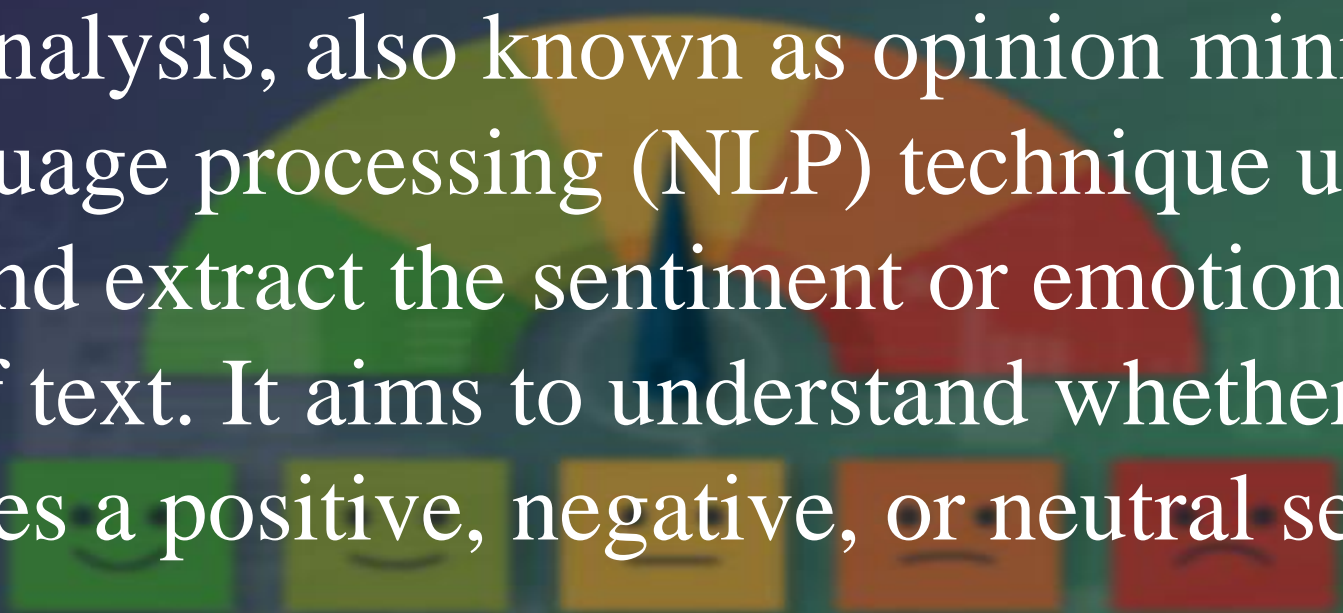
# email Phishing



Email phishing is a deceptive cyberattack method where attackers send fraudulent emails that appear to be from trustworthy sources, aiming to trick recipients into divulging sensitive information or clicking on malicious links.

# Sentiment analysis

Sentiment analysis, also known as opinion mining, is a natural language processing (NLP) technique used to determine and extract the sentiment or emotion expressed in a piece of text. It aims to understand whether a given text expresses a positive, negative, or neutral sentiment.

# Research Problem

How to address the issue of imbalanced datasets, which can affect the accuracy of the model

Using Large dataset to train the model

How to use sentiment analysis for detect phishing emails.

Using Naive Bayes algorithm to train the model

# System Diagram

# Objectives



❖ Detect emotion of the text using sentiment analysis.
❖ Improve the model to detect the text phishing or not.
❖ Display a warning message when detect a phishing email.
❖ Integrate the developed model to a mobile application.

# REQUIREMENTS

Software Requirements

VS Code
Jupiter Notebook
Google CoLab

Algorithms

Naive Bayes (BernoulliNB)
Random Forest
Natural Language Processing (NLP)

Techniques

Machine Learning
Sentiment Analysis
Model Training

SLIIT
FACULTY OF COMPUTING

# Accuracy
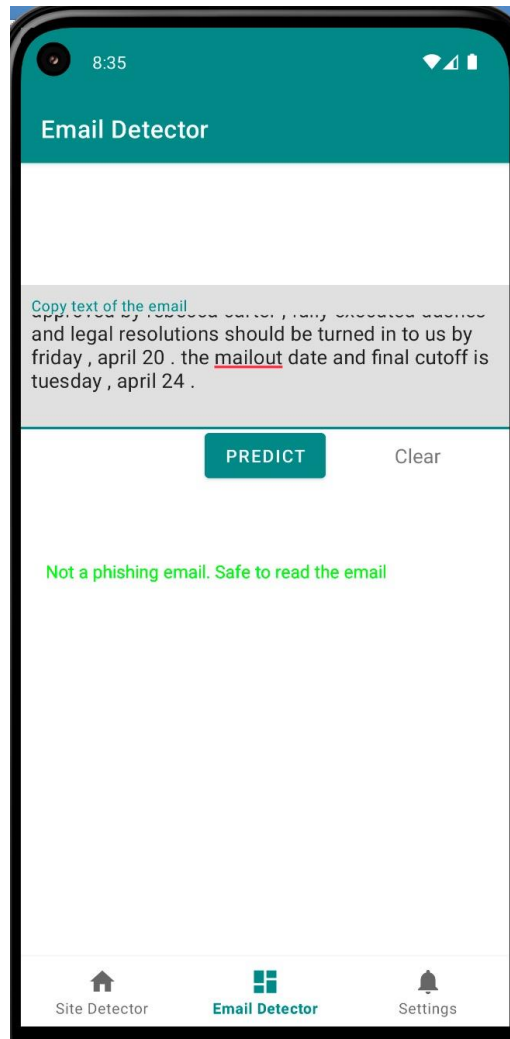
## MultinomialNB

# Accuracy

## BernoulliNB



```python
print('confusion matrics :')
print(confusion_matrix(y_test, y_predict), end='\n\n')

print('accuracy ',accuracy_score(y_test, y_predict))
```
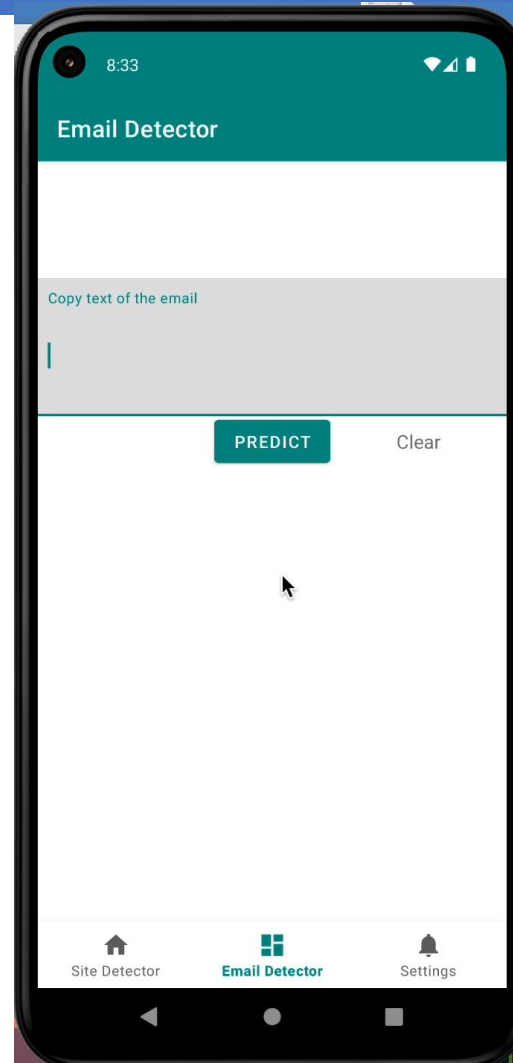```
[47]  ✓  0.2s                                          Python

...   confusion matrics :
      [[1584   29]
       [ 105 2667]]

      accuracy  0.9694412770809578
```

# Mobile Application

# Demonstration

# References

[1] Tanusree Sharma, Priscilla Ferronato, and Masooda Bashir, Phishing Email Detection Method: Leveraging Data Across Different Organizations

[2] Duo Pan, Ellen Poplavska, Yichen Yu, Susan Strauss, Shomir Wilson, A Multilingual Comparison of Email Scams, 2020

[3] Enaitz ezpeleta, iñaki velez de mendiz abal, josé maría gómez hidalgo, urko zurutuza , Novel email spam detection method using sentiment analysis and personality recognition, 14 January 2020

[4] Tanusree Sharma and Masooda Bashir, An Analysis of Phishing Emails and How the Human Vulnerabilities are Exploited, 2020

[5] Sikha Bagui, Debarghya Nandi, Subhash Bagui and Robert Jamie White, Machine Learning and Deep Learning for Phishing Email Classification using One-Hot Encoding, 2021

[6] Ala Mughaid, Shadi AlZu'bi, Adnan Hnaif, Salah Taamneh, Asma Alnajjar, Esraa Abu Elsoud, An intelligent cyber security phishing detection system using deep learning techniques, 14 May 2022

[7] C.J. Hutto, Eric Gilbert, VADER: A Parsimonious Rule-based Model for Sentiment Analysis of Social Media Text

[8] Said salloum, tarek gaber, sunil vadera, and khaled shaalan, A systematic literature review on phishing email detection using natural language processing techniques, 2022

[9] Srishti Rawal, Bhuvan Rawal, Aakhila Shaheen, Shubham Malik, Phishing Detection in E-mails using Machine Learning, 2017

[10] Rakesh Verma and Nabil Hossain, Semantic Feature Selection for Text with Application to Phishing Email Detection

[11] Machine Learning and Deep Learning for Phishing Email Classification using One-Hot Encoding - Sikha Bagui, Debarghya Nandi, Subhash Bagui and Robert Jamie White 2021

[12] D. N. S. B. a. R. J. W. Sikha Bagui, "Machine Learning and Deep Learning for Phishing Email," Journal of Computer Science, 2021.

# IT20222468 | I.S.S.Perera

BSc (Hons) Degree in Information Technology (Specialization in Cyber Security)

# Phishing web site detection

# Component 4 – Phishing email detection using URL

# What is phishing

• Nowadays Phishing becomes a main area of concern for security researchers. Because it is not difficult to create the fake website which looks so close to .legitimate website. Experts can identify fake websites but not all the users can .identify the fake website and such users become the victim of phishing attack. .Main aim of the attacker is to steal banks account credentials. Phishing attacks are becoming successful because lack of user awareness. Since phishing attack exploits the weaknesses found in users, it is very difficult to mitigate them, but it is very important to enhance phishing detection techniques.

# Research question

How Effective Are Machine Learning Algorithms in Detecting Phishing Websites Based on URL Features?

# Sub objectives

Creating web site for chcek URLS

URL blocker

Pop up message for user

# Data set

1. URL data sets
   1. 1000-pshing.txt
   2. Legitamate_urls.txt

```
from sklearn.metrics import confusion_matrix,accuracy_score
cpnfusionMatrix = confusion_matrix(labels_test,prediction_label)
print(cpnfusionMatrix)
accuracy_score(labels_test,prediction_label)
```

```
[[267  45]
 [ 67 226]]
0.8148760330578513
```

# Web application

# references

- Malak Aljabri ,Samiha Mirza , "Phishing Attacks Detection using Machine Learning and Deep Learning Models," in 2022 7th International Conference on Data Science and Machine Learning Applications (CDMA) , Saudi Arabia, 2022

- Detecting phishing websites using machine learning technique "Ashit Kumar Dutta"

- Phishing Website Detection using Machine Learning Algorithms "Rishikesh Mahajan MTECH Information Technology K.J. Somaiya College of Engineering, Mumbai - 77"

- Detecting Phishing Domains Using Machine Learning "by Shouq Alnemari *ORCID andMajid Alshammari "

- Phishing-Website-Detection-by-Machine-Learning-Techniques "shreyagopal"

- Model of detection of phishing URLs based on machine learning Kateryna Burbela

- URL-based Phishing Websites Detection via Machine Learning "Qasem Abu Al-Haija; Ahmad Al Badawi"

- Phishing Detection using Machine Learning based URL Analysis: A Survey

- Phishing URLs Detection Using Sequential and Parallel ML Techniques: Comparative Analysis

- Modeling Hybrid Feature-Based Phishing Websites Detection Using Machine Learning Techniques

- Modeling Hybrid Feature-Based Phishing Websites Detection Using Machine Learning Techniques

- Sumitra Das Guptta, Khandaker Tayef Shahriar, Hamed Alqahtani, Dheyaaldin Alsalman & Iqbal H. Sarker

- Phishing URL Detection using Machine Learning Authors: Rutul Patel, Sanjay Kshetry, Sanket Berad, Justin Zirthantlunga

- PHISHING WEBSITE DETECTION USING NOVEL MACHINE LEARNING FUSION APPROACH

- Survey on Phishing Websites Detection using Machine Learning Authors: Mr. B Ravi Raju , S Sai likhitha, N Deepa, S Sushma

- Phishing website detection using machine learning and deep learning techniques (J. Phys.: Conf. Ser. 1916 012169)

- Detectionofphishingwebsitesusingmachinelearning techniques

- Phishing website prediction using base and ensemble classifier techniques with cross-validation Anjaneya Awasthi & Noopur Goel