# List of Figures

# List of Tables

# List of Abbreviations/Symbols

| Abbreviation/Symbol | Definition |
|---|---|
| $\alpha$ | Symbol for Malicious Page Seek Rate (A value defined for the seek rate of a Malicious Webpage by MalCrawler) or Learning Rate in the Context of DNN |
| Adam | Adaptive Moment Estimation (A Replacement Optimization Algorithm for SGD) |
| ADASYN | Adaptive Synthetic (An Algorithm for Generating Synthetic Samples in ML) |
| AE | Auto Encoder |
| AI | Artificial Intelligence |
| AJAX | Asynchronous JavaScript and XML |
| ANN | Artificial Neural Networks |
| API | Application Programming Interface (An interface that defines interaction between multiple software) |
| APK | Android Package (A Format of Archived Executable File in Android) |
| ASIC | Application Specific Integrated Circuit |
| AUC-ROC | Area Under Curve- Receiver Operating Characteristics |
| **b** | Bias |
| BERT | Bidirectional Encoder Representation from Transformers |
| C4.5 | A Decision Tree Algorithm |
| CNN | Convolutional Neural Network |
| CPU | Central Processing Unit |
| CSS | Cascading Style Sheets (A Style Sheet Language) |
| DBN | Deep Belief Networks |
| DNN | Deep Neural Network |
| DNS | Domain Name System |
| DOM | Document Object Model |
| DP | Differential Privacy |
| DT | Decision Tree |
| **E** | Loss Function |
| ELMo | Embedding from Language Models |
| FFN | Feed Forward Network |

| Abbreviation/Symbol | Definition |
|---|---|
| FIFO | First In First Out (A method of queuing in memory) |
| FL | Federated Learning |
| FN | False Negative |
| FNR | False Negative Rate |
| FP | False Positive |
| FPR | False Positive Rate |
| FWS | Federated Learning based Web Security (An App Developed as Part of this Thesis) |
| GPU | Graphics Processing Unit |
| HFL | Hierarchical Federated Learning |
| HTTP | Hyper Text Transfer Protocol |
| HTTPS | Hyper Text Transfer Protocol Secure (It is a HTTP Protocol that Supports Encryption) |
| IDF | Inverse Document Frequency (A Weighted TF that Increases Weight of Words that Occur Less Frequently) |
| IE | Microsoft Internet Explorer |
| iFrame | An HTML Element |
| IP | Internet Protocol |
| JADX | A JavaScript Disassembler |
| JDBC | Java Database Connectivity (A Java API for Database Connectivity) |
| jQuery | A JavaScript Library |
| JS | JavaScript |
| LSTM | Long Short Term Memory (A Type of RNN) |
| MalCrawler | Name given to a Customized Focused Crawler designed for this Research to Crawl Malicious Websites |
| ML | Machine Learning |
| NLP | Natural Language Processing (A Branch of AI Handling Interpretation of Human Language) |
| NN | Neural Network |
| PCA | Principal Component Analysis |
| RBF | Radial Basis Function |
| RELU | Rectified Linear Unit (An Activation Function for ANN) |
| RNN | Recurrent Neural Network |
| SDK | Software Development Kit |
| SGD | Stochastic Gradient Descent |

| Abbreviation/Symbol | Definition |
|---|---|
| SMOTE | Synthetic Minority Oversampling Technique (An Algorithm for Generating Synthetic Samples in ML) |
| SQL | Structured Query Language |
| SVM | Support Vector Machine |
| TF | Term Frequency (Count of Number of Times a Word Appears in a Document) |
| TFF | Tensor Flow Federated (A Python Library for Federated Learning) |
| TLD | Top Level Domain |
| TN | True Negative |
| TNR | True Negative Rate |
| TP | True Positive |
| TPR | True Positive Rate |
| TPU | Tensor Processing Unit (An ASIC for DNN Learning) |
| UI | User Interface |
| URL | Uniform Resource Locator (Address of a Webpage) |
| **W** | Weights |
| Web 2.0 | Webpage Technology that supports User Generated Content |
| WebView | An Android Class |
| WebView-OAuth | WebView Authentication Method |
| WEKA | Waikato Environment for Knowledge Analysis (A Machine Learning Software) |
| WHOIS | Provides Domain Registration Details |
| **X** | Input to a NN |
| XHR | XML HTTP Request (An API for Transferring Data Between Web Server & Web Browser) |
| XML | eXtensible Markup Language |
| XML | Extensible Markup Language (A Markup Language for Encoding Documents) |
| XSS | Cross-site Scripting (A Type of Web Based Attack) |
| **Y** | Output of a NN |
| $\delta\ model$ | A Small Incremental Model in Federated Learning with Each Participating Device |