

Title: ISmishU: Smishing Detection System for Code-mixed Messages Using Transfer Learning of XLM-RoBERTa

Authors: Capulong, Mark Daniel M.

Digang, Cyrel P.

Morales, Mark Jerico P.

Recio, Paolo Luigi G.

Soriano, Stephanie Elenn V.

Adviser: Prof. Aleta C. Fabregas

I. Introduction

In 2023, phishing and smishing attacks accounted for nearly half of all reported fraud cases in the Philippines, driven by the widespread use of mobile technology. Smishing, which is a combination of short message services or 'SMS' and 'phishing' in which the phishing attack is initiated through text messaging, remains a significant issue despite government efforts like the SIM Registration Act and blocking malicious URLs.

Existing smishing detection systems rely primarily on language models trained in high-resource languages like English, leaving low-resource languages like Taglish underrepresented. This gap poses unique challenges, as Taglish exhibits distinct linguistic patterns that existing systems fail to address.

This study proposes a smishing detection system tailored to Taglish, leveraging XLM-RoBERTa, a multilingual language model known for its ability to analyze message content. By incorporating both text analysis of the fine-tuned XLM-RoBERTa and URL verification through the VirusTotal API, the system aims to detect smishing messages, or "ham" (non-smishing) messages, with improved accuracy.

Training large models like XLM-RoBERTa on low-resource datasets introduces challenges such as the possibility of overfitting, where the model becomes too tailored to training data and performs poorly on new data. To mitigate this, the study explores techniques like layer freezing, where only certain layers of the model are updated, improving generalization and performance. By combining advanced techniques and

addressing the unique linguistic characteristics of Taglish, this research seeks to strengthen cybersecurity efforts and mitigate the threat of smishing in the Philippines.

II. Methodology

This study adopts an experimental quantitative approach to develop a smishing detection system for Taglish SMS using a fine-tuned XLM-RoBERTa model. Data is collected over an eight-week period through a Google Form survey distributed via social media platforms, targeting Filipino residents who have received spam or suspicious messages within the past two years. The collected dataset consists of smishing and legitimate ("ham") messages in Tagalog, English, and Taglish. To ensure accuracy and reliability, all messages are validated by a linguist's expert. Preprocessing steps include standardizing the text by correcting common misspellings (e.g., "anu" to "ano," "22o" to "totoo"), removing special characters and punctuation, and converting all text to lowercase for consistency.

The annotated dataset is used to train and test the XLM-RoBERTa model. The Experiment was conducted using various training dataset ratios and layer freezing configurations are implemented during the training to address possible overfitting by selectively updating the weights of certain layers while freezing others. The model's performance is assessed using accuracy, precision, recall, and F1 score, with testing results validated against evaluations from expert linguists to ensure reliability and robustness. Finally, a one-way Analysis of Variance (ANOVA) test is conducted to assess whether significant differences exist in the model's performance across the different training dataset ratios.

The system will have a two-layer analysis. If a URL is detected in the message, it is sent to VirusTotal through its API for URL analysis. If the URL is identified as malicious, the user interface immediately displays the result indicating that the message is a smishing attempt. If no URL is detected or the URL is not malicious, the text body of the message proceeds to content analysis of the most optimized model configuration, determined through experiments with varying training dataset ratios and layer freezing setups. The combination of text analysis and URL verification ensures a more comprehensive and accurate smishing detection system. of the system's effectiveness.

III. Expected Results/Outcomes

The system is expected to accurately detect whether a message, based on its content or the URL it contains, is classified as smishing or not. When the system performs content analysis using the fine-tuned XLM-RoBERTa model, it will provide a detailed output that includes the classification results, showing the percentage likelihood of the message being "ham" or "smishing" to help explain its decision. Additionally, the system will display a cosine similarity to input embeddings across all layers, offering insights into how the model processes and understands the message. For messages classified as smishing, the system will highlight specific tokens that contributed most to the classification, providing transparency and interpretability of the model's decision-making process.

IV. System







