

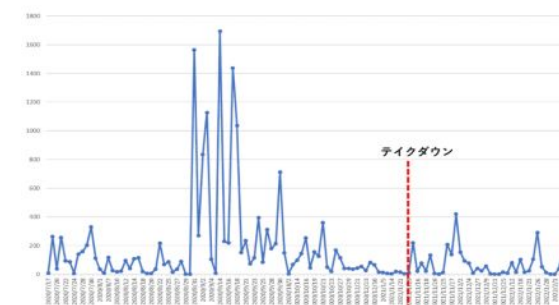
ALBERTを用いたアノマリ検知手法の検討

Consideration of NLP based anomaly detection using ALBERT

脇谷 峡平・ネットワーク分科会・情報セキュリティ大学院大学

研究背景・目的

近年,収束されていたと思われたマルウェアが急速に拡大している。Emotetを例に出すと、2021年初頭に一度テイクダウンが確認されたが、2022年3月にも確認され、各種機関は引き続き注意喚起を行なっている。(右図)そこで、Emotetの侵入経路を調査したところ、Emotetを含むマルウェアは、メール経由で感染拡大を狙っていることが判明した。実際に、IPA情報セキュリティ10大脅威2023内では、昨年に引き続きメールを経由しているとの報告がされている。そこで本研究では、メール経由で侵入するEmotetを含むマルウェアファミリー感染拡大を防ぐため、メールの文面に着目をした。この理由としては、メール内の文章から高精度の分類・検知を行えた場合、メールに添付されたマルウェア全体の検知及びフィルタリングが可能となるためである。



月別Emotet感染数

提案手法

本研究では、近年注目されている自然言語処理の一種ALBERTを用いる。ALBERTの特徴として、事前学習にBERTより困難なタスクであるSOP(Sentence-order-prediction)を実施しており、より複雑なタスクで精度向上を狙っている。更に、BERTの問題点であったメモリを大量に使用する箇所をTransformer block間のパラメータの共有により解消している点が特徴的である。更に本研究では、1ステップ先のフェーズとして、学習済みのALBERTが予測した値と検知用文章の双方を比較し、類似度を測ることによって閾値設定を可能とする。

現在までの進捗

他の自然言語処理でLSTM,Bi-LSTM,BERT,ALBERTでの比較実験を行った。結果を下図に示す。

ネットワークモデル	精度	f1テスト時
LSTM	40%	82.89%
Bi-LSTM	60%	85.53%
BERT	40%	87.23%
ALBERT	50%	91.48%

自然言語処理を用いた各学習環境下での比較結果

今後の展望

今後の展望としては2点存在する。
1.フィルタリングに最適な閾値の設定
2.データセットの大量確保による精度向上

1の解決方法は、全ての学習環境下において手動で最も良い値を模索し、提案を行う。
2の解決方法として、データセットの大量確保が重要となる機械学習内であるにも関わらず、現在の環境ではデータセットが不足しているため、より多くのデータセットを収集する計画を策定中である。