

機械学習型侵入検知システムにおける 能動学習を用いた誤検知低減手法の検討

小池 泰輔† 梅澤 猛† 大澤 範高†

千葉大学大学院融合科学研究科†

1 はじめに

ネットワークセキュリティの手法の1つとして、侵入検知システムが提案されている。検知手法は、パターンマッチングにより攻撃を検知するシグネチャ型と、正常時の振る舞いととの相違から異常を検知するアノマリ型の2種類に大別できる。後者は解析が済んでいない攻撃にも有効であり、機械学習を用いた手法が期待されている。しかしアノマリ型はシグネチャ型と比較して、多くの偽陽性を発生する傾向がある。

機械学習で効率良く学習モデルを構築する手法の1つに能動学習[1]がある。能動学習とは機械学習を効率良く行う為の教師データ選択に関する研究分野であり、分類が最も曖昧なデータを選択することで、少ない教師データで高精度なモデルの構築を目指す手法である。そこで本研究では能動学習を用いて選択した教師データによるモデル再構築を検討し、真陽性と偽陽性に対する効果を調べた。

2 能動学習

構築したモデルに対し、学習していないデータが入力される可能性がある。そのような入力に対して、学習済みのデータから規則性を導いて正しい入力を得る必要があるが、偏った分布を示すような学習データを用いた場合に過学習によって判別能力が低下する可能性が生じる。

能動学習とは、分類が難しいデータを学習データとして選択することで、より汎化能力の高いモデルの構築を目指す手法である。例えば2値分類問題において、分類が難しいデータとは境界周辺に分布するデータを指す。

3 提案手法

本研究ではモデル更新用データ $T_{K,n}$ が与えられた時、能動学習を用いることで、より少ないデー

タ数で高精度のモデル構築する手法を提案する。能動学習を用いたモデル構築の流れを図1に示す。

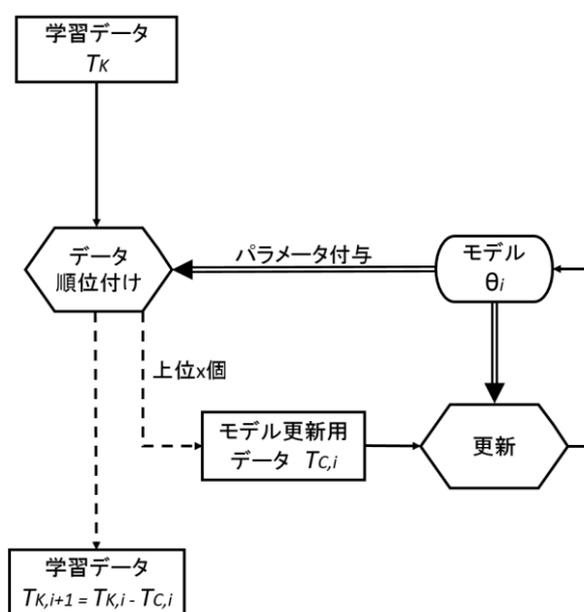


図 1:モデル θ_i による能動学習

まず、モデル θ_i で学習用データ $T_{K,i}$ を評価し、分類の難しさの順位をデータに付与する。順位が上位のデータから構成されたモデル更新用データ $T_{C,i}$ を抽出し、そのデータに基づいてモデルを更新し、モデル θ_{i+1} を構築する。次からモデルを更新する時はモデル更新用データ $T_{K,i+1} = T_{K,i} - T_{C,i}$ を用いる。ここで x は追加データ数、 n はデータのサブセット数、そして i はモデルの更新回数を表す。

能動学習には Margin Sampling を用いる。Margin Samplingとは、事前に構築したモデルを用いてデータの中から分離超平面とのマージンの大きさを比較し、よりマージンが小さいデータを選択する手法である。これにより、構築したモデルにとって分類が難しいデータの選択が可能となる。また、通常と異常の2値分類を対象とし、分離平面からの距離を出力できるSVMを採用した。

False positives reduction using active learning in intrusion detection system

†Daisuke Koike, Takeshi Umezawa, Noritaka Osawa, Graduate School of Advanced Integration Science, Chiba University

4 実験

本研究では能動学習の有無による影響を確認するため、2種類のモデルを構築し、比較評価する。学習・評価データには KDDCuP1999[2]を用いた。データには特徴ベクトルとして 41次元の属性が付加されている。まず、用意したデータから、初期モデル構築用データ T_I と学習用データ T_K を抽出し、それに基づいて比較評価する。 T_I と T_K は母集団 T から重複しないようにランダム抽出する。また、データ選択による依存の影響を小さくするために、 $n = 100$ 回のデータ抽出とそれに基づいた比較評価を行う。学習データ T_K に含まれる通常と異常のデータ数を同一とする。

能動学習有りのモデルでは、 T_I に基づいて構築した初期モデル θ_0 が構築されているとする。また、 $T_{K,0} = T_K$ とする。 θ_i を用いてデータ $T_{K,i}$ に Margin Sampling を適用し、順位付けに基づいた更新用データ $T_{C,i}$ を抽出し、 $T_{C,i}$ を用いてモデル θ_i を更新して、 θ_{i+1} を構築するという工程を $i = 0$ から $i = 9$ まで 10 回繰り返す。またここでは更新用データ数 $x = |T_{C,i}| = 100$ 件とする。

対する能動学習無しのモデルは、用意したデータ T_K から無作為に 100 組のデータを抽出し、それぞれのデータを 10 分割したデータを用いてモデルの更新を 10 回繰り返す。

こうして得られた 2 つのモデルに基づいて、評価データに対する検知率と誤検知率を求め、能動学習の効果を検証した。

初期モデル構築用データは $|T_I| = 20$ 件、学習用データは $|T_K| = 1000$ 件とした。これらのデータセットはランダムに $n = 100$ 種類用意した。学習データ T_K とは別に用意した評価用データには通常 5000 件、異常 5000 件の計 10000 件を用いた。図 2 は各モデルの真陽性率を、図 3 は偽陽性率を示す。図のエラーバーは標準誤差を表す。

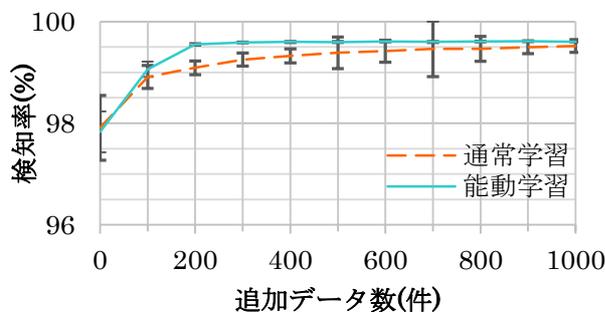


図 2 : 能動学習の有無による真陽性率の推移

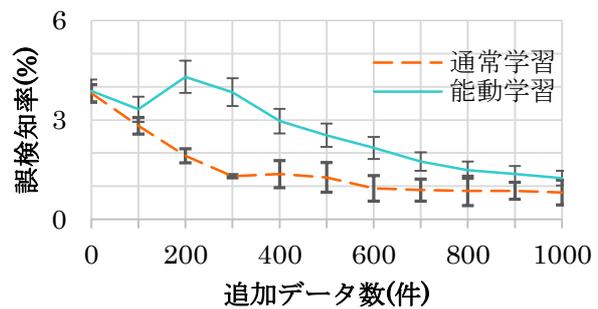


図 3 : 能動学習の有無による偽陽性率の推移

図 2 は異常を正しく検知した割合を示す。能動学習有りのモデルの方が少ないデータ数で高い検知率となった。図 3 は通常パケットを異常として誤検知した割合を示す。能動学習無しのモデルの方が全体的に低い誤検知率を表す結果となった。

5 考察

能動学習有りの偽陽性率が悪化した原因として、初期モデルの性能が関係している可能性がある。能動学習はモデルの判別能力などに依存するため、追加データ数が少ない内はモデルの判別精度が不十分で不適切なデータが選択された可能性がある。初期モデルの性能向上によって能動学習の偽陽性を低減させることができるか検証の必要がある。

6 まとめ

能動学習を用いて選択した教師データによるモデル再構築を検討し、真陽性と偽陽性に対する効果を検証した。今回は真陽性率向上に対しては効果があることがわかったが、偽陽性率低減に対しては課題が残る結果となった。今後は偽陽性低減のため、初期モデルの性能の違いによる能動学習への影響を評価する予定である。

参考文献

- [1]Settles, B. Active learning literature survey.” University of Wisconsin, Madison,” Computer Sciences Technical Report 1648, (2010)
- [2]UCI KDD Archive: KDD Cup 1999 Data, UCI (online), available from <<http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html>> (accessed 2014-12-14).