

機械学習に基づくネットワークトラフィックの分析手法

和田見 一真

ネットワークトラフィックの分類はネットワーク管理やセキュリティ監視、通信の最適化などネットワーク関連の様々な活動に有用である。近年、情報化社会の進展によって、インターネット上のネットワークやアプリケーションのデータ量が急激な増大を続け、トラフィック分類の重要性がより高まっている。2017年の総務省の発表によれば、日本におけるトラフィックは2004年から2016年にかけて約40倍も増えている。しかし、最近ではSSLに代表される暗号化された通信が主流なので、DPI(ディープ・パケット・インスペクション)などのパケット全体を解析する必要がある従来の手法では十分な分析が難しくなりつつある。そこで最新の研究では従来の手法が抱える問題に対処する有望な方法として機械学習が用いられるようになってきた。機械学習はパケットの中身を分析するだけでなく、その統計的な特性をもとに分類することもできる。そこで本研究ではネットワークフローに注目して機械学習によるネットワークトラフィックの分類とその分類に適したネットワークの属性の分析を行う。

本研究では、機械学習ソフトウェア Weka を用いて教師なし学習アルゴリズムの一つである k-means 法で、ネットワークフローをいくつかのクラスターにグループ化した。研究で使ったデータセットはデルカウア大学の学内ネットワークでパケットキャプチャにより、合計 3,577,296 件のネットワークフローを収集したものである。そのデータセットでは、レイヤー7としてSSLを含むHTTPを用いたアプリケーションのラベルがあらかじめ付与されている。そこから5つの主要なアプリケーションごとに約2万件になるようにランダムサンプリングし、計約10万件としたものを4セット準備して実験を行った。クラスターの数は3から7の範囲で指定して、それぞれの数におけるクラスターリングの精度の変化を確かめた。それらのクラスターごとに振り分けられたプロトコルとフローごとのラベルを照合して、その適合率を求めた。トラフィック分類に大きな影響を与えるネットワークフローの因子を特定するために、一回のクラスターリングごとに利用する属性を必要に応じて変えながら、クラスターリングとその適合率の算出を繰り返した。

実験の結果、クラスターの数に関係なく順方向フローを含んだモデルと逆方向フローを含んだモデルでは、明らかに順方向フローの方が優れた精度をもたらすことがわかった。データセットのプロトコル数とクラスター数が一致する場合には、平均して4.1ポイント順方向フローの精度が高かった。より精度を高めるためには異なるクラスターリングアルゴリズムの使用や教師あり学習アルゴリズムとの併用が考えられる。

(指導教員 阪口哲男)