

# トラフィック特徴量の時系列データにおける 相関特性を用いた変化点からの異常検出

## 1 はじめに

近年、コンピューターネットワークシステムが社会に普及していくにつれ、ネットワークセキュリティの確保が必要不可欠となっている。そこで、インターネット上で起こるセキュリティインシデントを分析する「インシデント分析システム」が注目されている。

インシデント分析に関する重要な要件に、各種ログデータからのインシデント候補のリアルタイム検知があり、時系列データからの変化点検出を用いた手法が提案されている。これまで、この種の研究は統計量の分野で広く扱われてきている。

情報源の外れ値や変化点の特定は、ログデータから異常行動や不正行為につながるデータや、新しいトレンドを示す重要なデータを発見することができるとされている。そのため、データマイニング [1][2][3] や統計量に関する研究において、最も注目されている問題である。

文献[4]では、データマイニングの観点から、外れ値と変化点に明確な関係を与え、外れ値と変化点を区別して検出する変化点検出エンジンを用いて、異常状態を検出している。この手法は、時系列データを当てはめる確率モデルが過去の統計量を次第に忘れていくことによって、リアルタイムで時系列データの特徴をうまく抽出し追跡できる学習アルゴリズムを用いている。

ネットワークにおける異常状態は、トラフィック量やアクセス頻度時系列データの変化に特徴が現れることが多い。これまでの研究では、そのような特徴量を単一で利用することにより異常を検出している。しかし、大量の通信量やアクセス頻度を伴うネットワークにおいては、単一の特徴量の変化だけで正常異常を区別することは困難である。また、全体の通信量に対して、異常な通信量が少ない場合においても、統計的手法では検出することができない。

そこで本稿では、それらの問題点を改善する手法を提案する。現在、インターネットやイントラネットで標準的に使われているプロトコルは TCP/IP

(Transmission Control Protocol/Internet Protocol) であり、TCP/IP に基づいて通信を行っている。そのため、トラフィック量の特徴を抽出した時系列データは、強い相関を持つデータが含まれていると考えられる。この相関特性を利用して、統計的手法[4]による変化点検出エンジンを用いることで、相関が大きく変化する時点を異常と検出することができる。

以降 2 章においてトラフィック特徴の相関特性について述べ、3 章で変化点検出エンジンについて説明する。4 章で提案手法を説明して、5 章で実験と考察、6 章でむすびとする。

## 2 トラフィック特徴量の相関特性

### 2.1 トラフィック特徴の概要

インターネット通信は、TCP/IP というプロトコルに基づいて、ある一定のルールが決められている。そのため、ルールから逸脱した通信や、通常の通信では起こりにくいトラフィック特徴が現れた場合、それはしばしば異常状態である可能性が高いと考えられる。以下では、TCP/IP で定められている通信方法と相関特性について説明する。

#### 2.1.1 TCP/IP

現在、標準的に使われているプロトコルは TCP/IP である。国際標準化機構 (ISO) に制定された、「OSI 基本参照モデル」(異機種間のデータ通信を実現するためのネットワーク構造の設計方針) には階層的に、物理層、データリンク層、ネットワーク層、トランスポート層、セッション層、プレゼンテーション層、アプリケーション層が存在する。OSI 基本参照モデルにおいてネットワーク層に属する IP は、ネットワークに接続している機器の住所付け(アドレッシング)や、相互に接続された複数のネットワーク間での通信経路の選択(ルーティング)を行っている。また、トランスポート層に属する TCP に関しては、コネクション型のエンドシステム間における通信の信頼性を高めている。

#### 2.1.2 3way-handshake

IP パケット(ネットワーク層やトランスポート層

を流れる分割されたデータの単位) に付与されるヘッダのフィールドは数多くあるが、TCP ヘッダの「TCP Flag」は、通信に信頼性を持たせるためのコネクション型の通信を行う際に重要な機能を担っている。TCP における通信の流れを図 1 に示す。まず、送信元からコネクション確立のために SYN フラグがセットされた IP パケットを送出し、送信先が SYN と ACK のフラグが立った IP パケットを返送してくると、最後にもう一度 ACK フラグがセットされた IP パケットを送出して通信を開始する。このコネクション確立のための手順を「3way handshake」と呼び、TCP の信頼性が高い大きな理由となっている。また、データの通信を終える際には、送受信するホスト間で相互に FIN フラグのセットされた IP パケットを送り、片方ずつ通信を切断していく。このコネクション切断のための手順を「ハーフクローズ」と呼ぶ。図 1 に 3way handshake の手順を示す。

このように、TCP の通信方式である 3way-handshake を使って通信を行っているため、TCP/IP 通信においては、SYN フラグと FIN フラグの相関は非常に高いと考えることができる。

## 2.2 異常事象の特徴

本稿では、不正アクセスとして大きく取り上げられる問題のうち、「ポートスキャン」と「DoS 攻撃」に注目し、それらの特徴を抽出することで、トラフィック特徴量の相関特性を利用することを考える。

### 2.2.1 ポートスキャン

ポートスキャンとは、サーバのサービスポートに順番にアクセスし、動作しているサービスや OS の種類を調べ、侵入口となりうる脆弱なサービスポートがないかどうかを調べる行為のことである。一般的なポートスキャンには一定の特徴が見られる。まず、悪意のあるユーザが外部から進入を試みる際に最初に行われることが多い。また、侵入口を探すため、広範囲のサービスポートへアクセスをする。さらに、アクセスに対するサーバの反応を調べるため、送信元 IP アドレスは偽装されない、などが挙げられる。

ポートスキャンにおけるトラフィック特徴量の相関特性の例として、しばしば悪意のあるユーザが用いる TCP SYN スキャンを以下で述べる。これは、ネットワークログにポートスキャンの痕跡を残すことを回避するために行われ、通常の TCP 接続を確立するプロセスを途中まで行うことで、ネットワークサービ

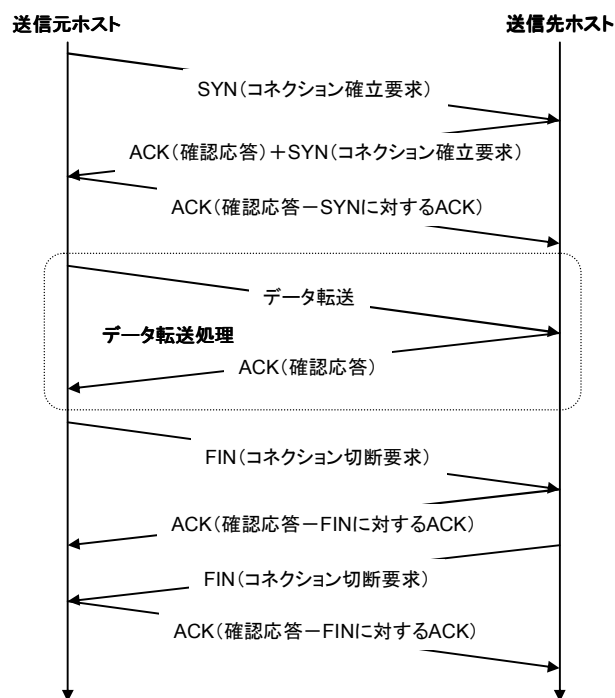


図 1 : 3way handshake

スの存在を確認するために用いられる。TCP 接続を確立するためには、2.1.2 で述べた 3way handshake と呼ばれるプロセスが用いられ、正常なデータ転送を行うための通信の確立の流れは、図 1 のようになる。しかし、TCP SYN スキャンでは、接続先ホストが SYN と ACK のフラグが立った IP パケットを返送してくると、接続の痕跡を残さないように、接続の強制終了である RST のフラグを立てて IP パケットを送出する。このような流れにより、TCP SYN スキャンは接続先に接続してきたことを知られることなく、接続先がどのようなサービスを提供しているかを知ることができる。

ポートスキャンは、不正アクセスの準備段階として用いられることが多いが、途中までは通常の TCP コネクションを成立させる手順であるため、正常なアクセスと判別が難しく、防ぐことが困難であるとされている。

しかし、通常のアクセスでは、3way handshake を用いて接続を確立しデータ転送が完了した後は、ハーフクローズという切断方法を用いて、通信を終了させることとなっている。その際に、ポートスキャンでは RST フラグによる強制切断を行うため、FIN フラグと SYN フラグの相関が崩れると考えられる。この特徴を利用することで、相関が大きく変化した時点で、異常状態である可能性が高くなると考えられる。

また、通常のアクセスでは、サーバで提供されてい

るサービスはある程度決まっているため、短い時間間隔で、提供しているサービス以上のポートに特定の IP から複数回接続してきた場合も、異常状態である可能性が高くなると考えられる。つまり、単位時間における、特定 IP からのサービス提供ポート以外へのポートに対する通信は、全体の通信に対して、相関関係はないと考えられる。

### 2.2.2 DoS 攻撃

DoS (Denial of Services) 攻撃とは、サーバに大量の接続要求を送って回線速度を低下させたり、過負荷でサーバを停止させる攻撃のことである。送出するパケットサイズは非常に小さいものから、ネットワークが受け入れられる限界までの大きなサイズまで多岐に渡る。一般的に、サーバのマシンスペックは性能があるので、サーバに過負荷を与えるために攻撃側は単独ホストからの攻撃ではなく、複数ホストから DoS 攻撃を同時に行う DDoS (Distributed DoS) 攻撃を実行する。

DoS 攻撃では、サーバからの反応を逐一確認する必要がないため、送信元 IP アドレスは偽装されることが多い。また、掲示板などによって人手を集め、単純に Web サイトを再読み込み (リロード) し続けることで DoS 攻撃を実行することができる。このため、一つ一つのアクセスは正常な通信であることが多く、DoS 攻撃は防ぐことが難しい攻撃といえる。しかし、DoS 攻撃を行う場合は、送信先 IP と送信先ポートは一定であり、さらにはパケットサイズが等しいことが多い。つまり、通常の通信において、これらの特徴を持つパケットは、全体の通信において相関はあるものの、急激な相関の変化や、相関が崩れることはないと考えられる。

## 3 変化点検出エンジン

ここでは、提案手法において用いる変化点検出エンジン (ChangeFinder) の理論を述べる。ChangeFinder の特筆すべき点は、2 段階の学習過程を繰り返すところにある。最初に、第 1 段階で学習したモデルを利用して外れ値を検出する。次に、第 2 段階で学習したモデルを用いて変化点検出を行う。ChangeFinder による変化点検出の流れを図 2 に示す。

### 3.1 ChangeFinder

第一段階学習では、まず各時刻  $t$  において、AR モ

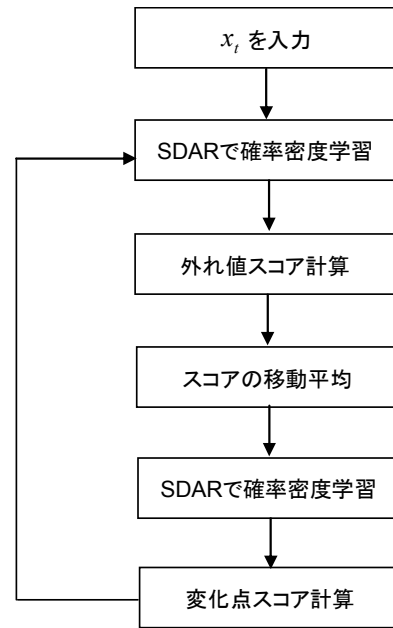


図 2 : ChangeFinder の流れ

デルを SDAR アルゴリズムによって学習する。そして、時系列データに対する外れ値らしさを示す、外れ値スコアを計算する。

第二段階学習では、外れ値スコアに対して、再度 AR モデルをあてはめ、これを学習して、変化点スコアを計算する。変化点スコアが大きいほど、 $t$  が変化点である度合いが高い。

ChangeFinder の特徴は、第一段階学習では時系列中の外れ値しか検出できないところを、外れ値スコアの平滑化を通じて、本質的なモデルの変動を検出しているところにある。計算量に関しても、データ数  $n$  に対して、統計的検定に基づく方式が  $O(n^2)$  であるのに対して、ChangeFinder の計算量は  $O(n)$  で済むため、明らかに効率がよいことがわかる。

さらに、ChangeFinder は平均値の変化だけでなく、AR モデルのパラメータ (AR 係数や分散) の変化も原理的には検出できる。実際に、分散が突然変化する場合でも、十分な効果が得られる。

### 3.2 外れ値の検出

最初に、時間  $t$  に対して、 $\{x_t : t=1,2,\dots\}$  で表すことができる時系列データを考える。  $t$  が変化するときの  $x_t$  が、本稿における重要な意味を示す  $d$  次元の実数値ベクトルである。ChangeFinder は、第 1 段階で、 $\{p_t : t=1,2,\dots\}$  として表される時系列データの確率密度関数を計算する。データ  $x_t$  が入力されると、この時系列データは  $\{x_t\}$  から徐々に学習していく。一般的

に、それぞれの  $p_t$  が確率過程の密度を表すと考える。確率過程  $p$  に対して、 $x^t = x_1 x_2 \cdots x_t$  を与える  $x_{t-1}$  の条件付き確率密度関数は、 $p(x_{t-1} | x^t)$  で表される。その確率密度関数の系列を  $\{p_t : t=1,2,\dots\}$  とする。学習方法には、確率過程  $p$  を推定するために、各入力  $x_t$  に対して、次式(1)を用いて、 $x_t$  の外れ値スコアを計算する：

$$\text{Score}(x_t) = -\log p_{t-1}(x_t | x^{t-1}) \quad (1)$$

式(1)の左辺は、確率密度関数  $p_{t-1}(\cdot | x^{t-1})$  に対する  $x_t$  の対数予測損失を表し、対数損失スコアと呼ぶことにする。

### 3.3 SDAR アルゴリズム

学習モデルで用いるアルゴリズムについて説明する。 $\theta_t$  は、 $x_t$  が与えられたときの  $\theta$  の推定値を表し、 $p_t = (\cdot | \theta_t)$ 、 $\theta = (A_1, \dots, A_k, \mu, \Sigma)$  とする。 $\theta$  の評価のために、以下の量を最大にする  $\theta$  の値を計算するアルゴリズムを提案する。

$$\sum_{i=1}^t (1-r)^{t-i} \log p(x_i | x^{i-1}, \theta)$$

これは、オンラインで使用するための、最尤推定法の変形である。このとき、重さが時間  $t$  で指数的に減少するところで、尤度が最大になる。これを、SDAR (sequentially discount-ing AR model estimating) アルゴリズムと呼ぶことにする。

SDAR アルゴリズムは、バッチ学習方式の AR モデルを改良した、逐次型学習方式である。SDAR アルゴリズムでは、逐次学習と忘却機能という 2 つのポイントがある。

逐次学習とは、新たなデータを 1 つ読み込むごとにパラメータの推定値を更新する。

忘却機能とは、 $i$  時点前のデータの影響が  $(1-r)^i$  倍に減少するようにパラメータの推定値を更新する。これによって、非定常な情報源に対応できる。アルゴリズムのパラメータ  $r$  を忘却パラメータと呼び、 $1/r$  個程度の過去データの情報を蓄積するようにする。以下に SDAR アルゴリズムを示す。

SDAR アルゴリズム ( $0 < r < 1$ )

STEP 1. 初期化

Set  $\hat{\mu}, C_j, \hat{A}_j (j=1, \dots, k), \hat{\Sigma}$ .

STEP 2. パラメータ更新

For  $t=1, 2, \dots$

$x_t$  を読み込む:

$$\hat{\mu} := (1-r)\hat{\mu} + rx_t$$

$$C_j := (1-r)C_j + r(x_t - \hat{\mu})(x_{t-j} - \hat{\mu})^T$$

以下の連立方程式を  $A_i$  について解く:

$$C_j = \sum_{i=1}^k A_i C_{j-i} \quad (j=1, \dots, k). \quad (2)$$

方程式(2)の解を  $\bar{A}_1, \dots, \bar{A}_k$  とし、以下を計算

$$\hat{x}_t := \sum_{i=1}^k \bar{A}_i (x_{t-k} - \hat{\mu}) + \hat{\mu}$$

$$\hat{\Sigma} := (1-r)\hat{\Sigma} + r(x_t - \hat{x}_t)(x_t - \hat{x}_t)^T$$

このアルゴリズムにおいて  $t$  番目のデータまで用いて得られる確率密度関数を、 $p_t$  とする。

### 3.4 変化点検出

次に変化点スコアを求める。 $T$  を正の整数とする。データ列  $\{x_t\}$  に対して、 $T$  移動平均スコア  $y_t$  を次式(3)で定義する。

$$y_t = \frac{1}{T} \left( \sum_{i=t-T+1}^t \text{Score}(x_i) \right) \quad (3)$$

ただし、 $\text{Score}(x_t)$  は式(1)より求める。この計算によって、新たな時系列  $\{y_t : t=1,2,\dots\}$  を得る。

次に、 $\{y_t\}$  を入力データとして、再度 SDAR アルゴリズムを用いて AR モデルの学習を行う。 $q_t$  を  $y_t$  まで用いて得られる確率密度関数で表すと、AR モデルによる確率密度関数の列  $\{q_t : t=1,2,\dots\}$  が得られる。

さらに、対数損失式(1)と同様に、 $T$  移動平均スコアを次式(4)で定義する。

$$\text{Score}(y_t) = \frac{1}{T} \sum_{i=t-T+1}^t (-\log q_{i-1}(y_i | y^{i-1})) \quad (4)$$

式(4)は、時間  $t$  の変化点らしさを示す指標となる。すなわち、 $\text{Score}(y_t)$  が大きければ、変化度合いが大きいと解釈することができる。これを、変化点スコアと呼ぶ。

## 4 提案手法

提案手法の仕組みについて説明する。提案手法では、変化点検出エンジンによって変化点を解析する前段階で、トラフィック特徴の相関特性を利用した処理を行う。この前処理を行うことによって、変化がおきる時点で不正アクセスを受けている可能性が高いとする特徴量時系列データを、トラフィック量の時系列データから抽出することができる。今回は、検出したい不正アクセスを、ポートスキャンと DoS 攻撃と定義する。また、以下ではトラフィック量をパケット単位で扱うこととする。

### 4.1 時系列データ処理

一般的に、トラフィック量はパケット単位で取得することができる。それらは、複数の要素 (TCP・IP ヘッダに含まれる情報) から構成されているログデータであるため、変化点検出を行う前処理として、トラフィック量の特徴を反映した数値による時系列データに変換する必要がある。まず、パケットデータを TCP・IP ヘッダ情報別に分割する。次に、分割したパケットデータの単位時間当たりのパケット数をカウントし、ヘッダ情報別特徴量時系列データを作成する。さらに、ヘッダ情報の組み合わせによる特徴量時系列データも作成する。以下では、単位時間を 1 分とし、時系列データ処理の詳細について説明する。図 3 に特徴量抽出の流れを示す。

ここで、TCP・IP ヘッダ情報とは、送信元 IP、送信元ポート、送信先 IP、送信先ポート、TCP フラグ、サービス、パケットサイズである。

#### 4.1.1 ヘッダ情報別特徴量抽出処理

まず、パケットデータをヘッダ情報別に分割する。用いるヘッダ情報は、フラグ単位 (SYN, FIN, ACK)、パケットサイズ単位 (100 バイト毎, 40 バイト)、サービス単位 (SMTP, FTP, HTTP, Telnet, DNS) である。ここで分割したパケットデータを、「ヘッダ情報別パケットデータ」と呼ぶことにする。

これらのヘッダ情報別に分割したパケットデータから、単位時間でのパケット数をカウントし、その数値を一つの時系列要素とする。そのため、時系列データは 60 分間のパケットデータであれば、60 個の時系列要素を含む時系列データが作成できることになる。この処理によって得られるヘッダ情報別パケットデータに対応した時系列データを、「ヘッダ情報

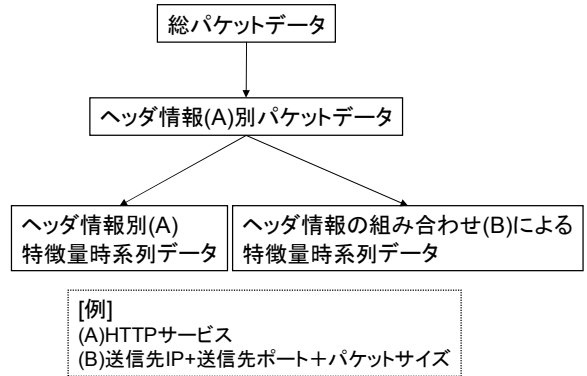


図 3：トラフィック特徴量抽出の流れ

別特徴量時系列データ」と呼ぶことにする。

#### 4.1.2 組み合わせによる特徴量抽出処理

ヘッダ情報別パケットデータに対応した、ヘッダ情報の組み合わせによるトラフィック特徴量の時系列データを作成する。ここで用いるヘッダ情報は、送信元 IP、送信元ポート、送信先 IP、送信元ポート、パケットサイズである。

不正アクセスの特徴として、DoS 攻撃やポートスキャンで用いている特徴を利用したものが多いため、これらの特徴を利用することで、DoS 攻撃やポートスキャン以外の異常状態を検出できる可能性がある。まず、DoS 攻撃で用いられている特徴を抽出する方法について説明する。

DoS 攻撃は、攻撃先ホストの反応を確認する必要がないため、送信元 IP が偽装されている可能性が高い。攻撃特徴として、短時間に大量のアクセスが発生し、アクセス対象は特定ホストの特定ポートである。さらに、大量に発生する通信は同一であることが多いため、パケットサイズが等しい可能性が高い。

提案手法では、これらの特徴に着目し、次式(5)よりヘッダ情報別パケットデータにおいて、単位時間当たりの任意の送信先 IP アドレスと送信先ポート番号とパケットサイズが同一である組み合わせの数をカウントする。次に、同一のヘッダ情報別パケットデータにおいて、組み合わせ内のヘッダ情報を含むパケットの単位時間当たりのパケット数をカウントし、そのパケット数を組み合わせ数で除算する。この処理によって、ヘッダ情報別パケットデータが持っている DoS 攻撃の特徴を数値化する。

また次式(6)より、単位時間当たりのヘッダ情報別パケットデータにおいて、送信元 IP アドレスと送信先 IP アドレスと送信先ポートが同一である組み合わ

せ数をカウントする。次に、その組み合わせのうち、送信元 IP アドレスと送信先 IP アドレスの組み合わせ数をカウントする。そして、送信元 IP アドレスと送信先 IP アドレスと送信先ポートが同一である組み合わせ数を送信元 IP アドレスと送信先 IP アドレスの組み合わせ数で除算する。この処理によって、特定の IP 同士におけるポートの分散状況を数値化することができる。本稿では、式(6)によるポートスキャンの特徴量を用いた実験は行わないこととした。

ヘッダ情報別特徴量

$$= \frac{(\text{送信先IP} + \text{送信先ポート} + \text{パケットサイズ}) \text{ が含まれるパケット}}{(\text{送信先IP} + \text{送信元ポート} + \text{パケットサイズ}) \text{ 組み合わせ数}} \quad (5)$$

ヘッダ情報の組み合わせ特徴量

$$= \frac{(\text{送信元IP} + \text{送信先IP} + \text{送信先ポート}) \text{ 組み合わせ数}}{(\text{送信元IP} + \text{送信先IP}) \text{ 組み合わせ数}} \quad (6)$$

ただし、式 2 においてのヘッダ情報別パケットデータは複数のポートを含む特徴量の場合とする。

## 4.2 相関係数処理

次に、ヘッダ情報別パケットデータに対応した、ヘッダ情報別特徴量時系列データとヘッダ情報の組み合わせによる特徴量時系列データを用いて、トラフィック特徴量の相関係数時系列を作成する処理を説明する。作成された時系列データを変化点検出エンジンに適用することにより、変化点を検出し、異常状態を検出することを目標とする。

### 4.2.1 ヘッダ情報別特徴量同士の相関

ヘッダ情報別特徴量時系列データには、強い相関関係にあるものが含まれているため、それらの相関係数の変化を解析することによって、異常状態を検出することができる。例として、SYN フラグと FIN フラグについて述べる。2.2.1 で述べたポートスキャンの特徴により、SYN フラグと FIN フラグの相関関係が崩れる時点が生じる。この時点において、ポートスキャンと同じ特徴を持つ不正アクセスを検出することができると考えられる。そこで、SYN フラグと FIN フラグのヘッダ情報別特徴量時系列データから、相関係数時系列データを得る。

相関係数を求める際に、窓サイズ  $N$  を決定し、 $N$  を 1 ずつずらすことによって、 $N$  間隔における相関係数の時系列データを得る。2 組の時系列データ  $(x, y) = \{(x_i, y_i)\} (i = 1, 2, \dots, N)$  が与えられた時、相関係

数  $c$  は、ピアソンの積率相関係数によって次式(7)で計算する。

$$c = \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^N (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^N (y_i - \bar{y})^2}} \quad (7)$$

ただし、 $\bar{x}, \bar{y}$  はそれぞれの時系列データの窓サイズ  $N$  内における相加平均である。

### 4.2.2 ヘッダ情報の組み合わせ特徴量による相関

ヘッダ情報別特徴量時系列データと、それに対応するヘッダ情報の組み合わせによる特徴量時系列データとの相関係数を計算する。計算は 4.2.1 と同様の方法を用いる。通常状態では、短時間で送信先 IP アドレスと送信先ポート番号とパケットサイズのパケットが同一である組み合わせが急激に変化することは少ないため、滑らかな変化を持つ相関関係であると考えられる。そのため、相関関係が急激な変化を示した場合、異常状態である可能性が高くなる。

## 5 実験と考察

本稿では、検証実験を行うためのデータとして、MIT の LINCOLN 研究所が作成した IDS 評価用データ [5]の一部を使用した。このデータは同研究所が DARPA (高等研究計画局) の支援によって 1998 年と 1999 年に作成されたもので、IDS の性能を比較するための一般的なデータとして広く利用されている。

今回は、ポートスキャンと、検出が難しいとされている HTTP に対する DoS 攻撃に注目して実験を行った。データには、正常と異常が混在されており、さらにはポートスキャンと HTTP に対する不正アクセスが含まれている Week5 の Tuesday のデータを使用した。具体的な攻撃時間として、211・795 時点において DoS 攻撃、616 時点においてポートスキャンが行われている。データは tcpdump[6]というパケットキャプチャソフトウェアを用いて作成されている。

以下では、ヘッダ情報別特徴量時系列データとヘッダ情報の組み合わせ特徴量時系列データの相関関係を利用する手法を用いることにより、DoS 攻撃検出結果を従来手法と比較した。また、ヘッダ情報別特徴量時系列データ同士の相関関係を用いる手法により、ポートスキャン検出の評価を行った。

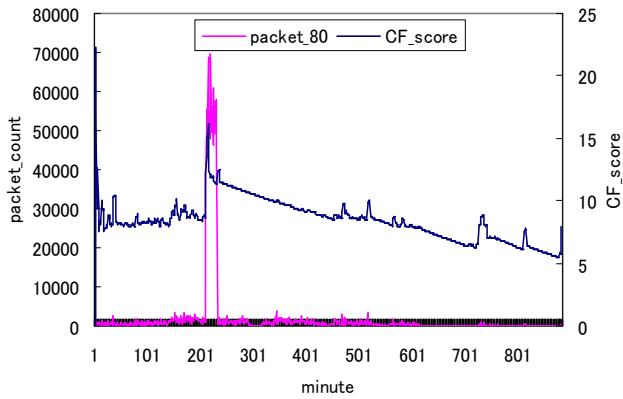


図 4：従来手法による DoS 攻撃検出結果

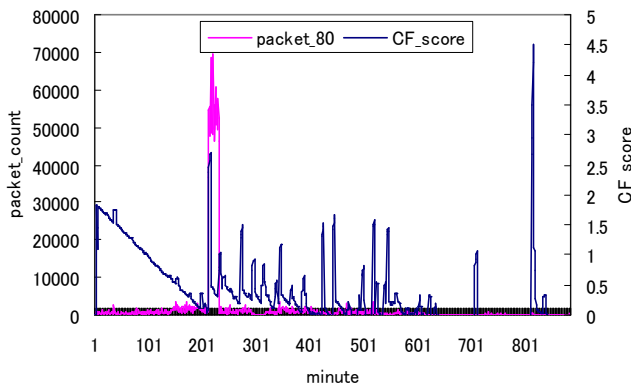


図 5：組み合わせ特徴量による DoS 攻撃検出結果

### 5.1 組み合わせ特徴量時系列データ

従来手法における DoS 攻撃に対する特徴は図 4 となった。ここで、CF\_score とは式(4)における変化点スコアである。従来手法では、異常状態はパケット量の変化に特徴が大きく現れることに注目し、ヘッダ情報別特徴量時系列データに変化点検出エンジンを適応する手法である。その結果、図 4 の 211 時点のようなパケット量に特徴が現れるような DoS 攻撃は検出できているが、795 時点で発生している DoS 攻撃に関しては、特徴を抽出できていない。これは、パケット通信量にあまり変化が現れないため、変化点を検出することができていない。さらに、正常状態であっても、パケット通信量が多く発生する場合があり、その時点を異常と検出してしまうという問題点もあった。提案手法において、ヘッダ情報別特徴量時系列データとヘッダ情報組み合わせ特徴量時系列データの相関を用いる手法で、同様の実験データを用いて実験を行った。結果を図 5 に示す。図 5 より、211 時点の DoS 攻撃の特徴を抽出できており、さらにはパケット通信量にあまり変化が現れない 795 時点の DoS 攻撃に対しても、特徴を抽出するこ

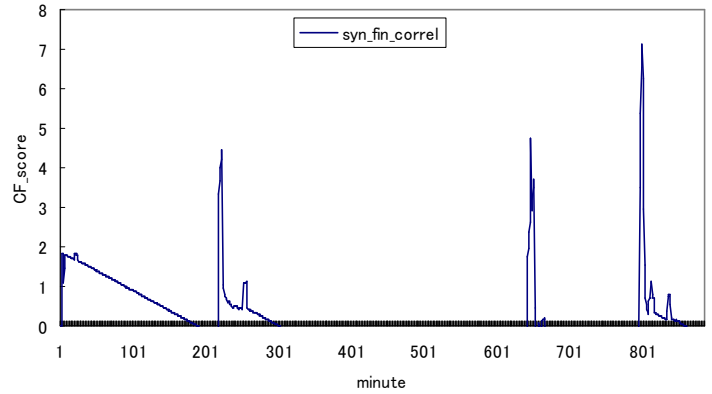


図 6：種類別パケット量による検出結果

とができている。ここでは、実験的に変化点スコアが 2.5 以上の場合に異常を検出することとした。

### 5.2 ヘッダ情報別特徴量時系列データ同士

ポートスキャンなど、サーバログに残らず、パケット量にもあまり変化が見られない攻撃を検出するために、ヘッダ情報別特徴量時系列データ同士の相関関係を用いる手法を用いた。ここでは、2.2.1 の特徴を利用するため、SYN フラグと FIN フラグのヘッダ情報別特徴量時系列データを用いて実験を行い、結果を図 6 に示す。図 6 では、616 時点のポートスキャンの特徴を抽出できており、HTTP サービスに対する DoS 攻撃も検出できていた。

DoS 攻撃の特徴を抽出できていたのは、ポートスキャンのように、ログに記録が残らないようにするために、3way handshake による接続確立をとっていないためと考えられる。

## 6 むすび

本稿では、トラフィック特徴量の相関特性を利用することによって、特徴量時系列データの変化点を異常とする手法を提案した。相関特性を利用することで、従来手法では検出困難であったパケット通信量に特徴が現れない不正アクセスを検出することができた。

今後の課題としては、実験的に決定している相関係数を求める際の窓サイズの適応的決定（標準偏差が 0 になる場合の処理）、不正候補パケットの抽出ルールの改善、変化点検出エンジンの多次元対応が挙げられる。

## 参考文献

- [1] P. Burge and J. Shaw-Taylor, "Detecting Cellular Fraud Using Adaptive Prototypes," Proc. AI Approaches to Fraud Detection and Risk Management, pp. 9-13, 1997.
- [2] V. Guranlink and J. Srivastava, "Event Detection from Time Series Data," Proc. ACM-SIGKDD Int'l Conf. Knowledge Discovery and Data Minig, pp. 33-42, 1999.
- [3] E. M. Knorr and R. T. Ng, "Algorithms for Mining Distance Based Outliers in Large Data Sets," Proc. 24th Very Large Data Bases Conf., pp. 392-403, 1998.
- [4] J. Takeuchi and K. Yamanishi, "A Unifying Framework for Detecting Outliers and Change Points from Time Series," IEEE transactions on Knowledge and Data Engineering, , pp. 482-492, 2006.
- [5] MIT Lincoln Laboratory, "DARPA Intrusion Detection Evaluation",  
<http://www.ll.mit.edu/IST/ideval/>
- [6] "tcpdump", <http://www.tcpdump.org/>