

УДК 004.45(045)

¹Минаев Ю.Н. – д-р техн. наук,
¹Гузий Н.Н. – канд. техн. наук,
²Филимонова О.Ю – канд. техн. наук

ИДЕНТИФИКАЦИЯ АНОМАЛЬНЫХ СОСТОЯНИЙ ТРАФИКА КОМПЬЮТЕРНЫХ СЕТЕЙ НА ОСНОВЕ ПАРАДИГМЫ МНОГОМЕРНЫХ СЕТЕЙ

¹Национальный авиационный университет

²Киевский национальный университет строительства и архитектуры

Рассмотрены вопросы представления трафика компьютерной сети инвариантами тензора 4-го ранга. Трафик рассматривается как поток тензоров 2-го ранга многомерной сети. Показана эффективность применения предложенного подхода для идентификации аномальных состояний компьютерной сети

Введение

Компьютерные сети (КС) относятся к классу сетей с коммутацией пакетов, общей технологией которых является семиуровневая эталонная модель взаимодействия открытых систем. Вместе с тем, анализ трафика выполняется в большинстве случаев на уровне обобщенного трафика без учета его компонент[1, 2].

Одним из актуальных направлений исследования *NGN, FGN* является моделирование на основе многомерных сетей. Использование многомерной структуры сети создает комплекс принципиально новых задач, в частности выявления аномалий трафика[3,4].

Для идентификации аномального трафика компоненты трафика представляются отдельным измерением многомерного массива (тензора). В силу этого требуется разработка принципиально новых алгоритмов идентификации аномальных состояний трафика. Процессы, связанные с передачей информации в многомерной интегрированной сети достаточно сложны, одним из способов анализа таких сетей являются многомерные тензоры и их декомпозиции, в частности, сечения тензоров.

Цель исследования

Создание виртуальных многомерно-матричных моделей сетевых структур в качестве средства для исследования трафика КС.

Современное состояние исследований. В работах Г.Крона [5] показано, что представление объекта исследования (измерения) в виде тензора есть более адекватным. Тензорная модель, которая рассматривается как матричная проекция, позволяет анализировать объект в разных системах координат.

Анализ трафика широко использует понятие «нормального» и «нестандартного» (необычного) трафика. Ожидаемый (нормальный) трафик это трафик, который передается чаще всего. Известно, что пакетный трафик – это дискретный положительный процесс с сингулярной внутренней структурой. Наиболее полно он может быть представлен на уровне иерархических структур.

Проблема выявления аномалий в работе КС на основе анализа трафика относится к т.н. трудно формализуемым проблемам, решение которых современная наука видит в использовании интеллектуальных технологий [6]. Параметры трафика авторами в общем виде определены следующим образом: $x = \{x_i\}$,

$i = 1,9$: x_1 – Protocol ID протокол, связанный с событием (TCP=0, UDP=1, ICMP=2, unknown=3); x_2 – номер порта источника; x_3 – номер порта хоста назначения; x_4 – IP-адрес источника; x_5 – IP-адрес приемника; x_6 – ICMP Type тип ICMP-пакета (Echo Request or Null); x_7 – ICMP Code кодовое поле из ICMP-пакета

(None or Null); x8 – Raw Data Length длина данных в пакете; x9 – Raw Data порция данных в пакете.

Одним из способов выявления аномальных состояний трафика является идентификация при помощи стандартных перцептронов. Для заданных состояний КС -аномального - $X^{(a)} = \{x_i^{(a)}\}$, $i=1, 9$ и нормального $X^{(n)} = \{x_i^{(n)}\}$, $i=1, 9$ сформированы множества $\{^R X^{(a)}\}_{j=1}^J$ и $\{^R X^{(n)}\}_{j=1}^J$, $J= 1000 \div 10000$, у которых каждая компонента трафика лежит в интервале $\pm 15 \%$

от заданных, т.е. $x_i^{(a)} - 0.15 x_i^{(a)} \leq ^R X_i^{(a)} \leq x_i^{(a)} + 0.15 x_i^{(a)} (\forall i)$, $^R X_i^{(a)} \in \{^R X^{(a)}\}_{j=1}^J$ и соответственно $x_i^{(n)} - 0.15 x_i^{(n)} \leq ^R X_i^{(n)} \leq x_i^{(n)} + 0.15 x_i^{(n)} (\forall i)$, $^R X_i^{(n)} \in \{^R X^{(n)}\}_{j=1}^J$. Нейронная сеть (рис.1), обученная на одном из наборов, например, наиболее часто встречающимся- $\{^R X^{(n)}\}_{j=1}^J \rightarrow 1$, $\{^R X^{(a)}\}_{j=1}^J \rightarrow 0$, позволяет практически однозначно идентифицировать состояние.

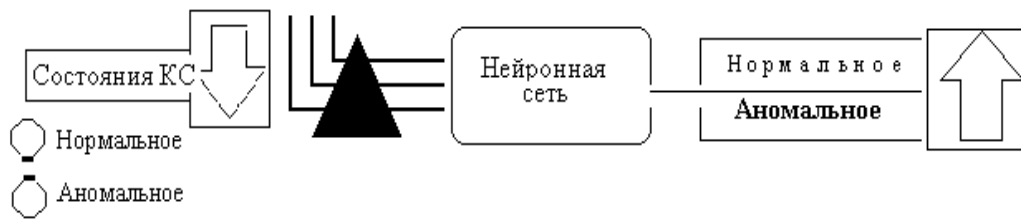


Рис. 1. Идентификация состояний КС при помощи НС

Недостаток данного подхода состоит в том, что предварительное определение базовых состояний $X^{(a)}$ и $X^{(n)}$ является экспертной оценкой и по своей сложности адекватно исходной задаче.

Современные приложения (трафик Internet, телекоммуникационные записи и др.) генерируют большие массивы данных. Тензоры (мультимерные массивы) обеспечивают естественное представление для таких данных. Тензорные декомпозиции стали инструментальными средствами для обобщения и анализа временных рядов [6–8].

На основе тензорных декомпозиций можно обнаруживать общности групп объектов, предсказывать потерю связей, ранжировать объекты и многие другие задачи [9]. На основе обобщения низкоранговой матричной аппроксимации. Rank-1 тензор определяется как внешнее произведение r векторов $x(1), \dots, x(r-1), x(r)$ заданных через r-мерный массив у которого (i_1, \dots, i_r) -тый вход есть $x_{i_1 i_2 \dots i_r}(1), x_{i_1 i_2 \dots i_r}(2), \dots, x_{i_1 i_2 \dots i_r}(r)$; это обозначено как $x(1) \otimes x(2) \otimes \dots \otimes x(r)$. Для любого r-мерного массива A существует аппроксимация

через сумму малого числа rank-1 тензоров.

Основная часть исследования

Многомерный трафик КС представляется многомерным массивом с достаточно большим количеством компонент, для анализа которого целесообразно использовать аппарат тензорного анализа и мультилинейной алгебры. Уникальное свойство тензоров состоит в том, что индексы могут иметь произвольную топологию [10]. Тензор-декомпозиции – модели тензора в смысле наименьших квадратов: если T- тензор, то \tilde{T} -модель тензора (в общем случае меньшей размерности), полученная, если минимизировать $\delta = \|T - \tilde{T}\|_2$ или точнее

$$\delta = \sum_{i_1, i_2, \dots} |T(i_1, i_2, \dots) - \tilde{T}(i_1, i_2, \dots)|^2$$

Тензор может храниться в т.н. отложенной форме как сумма rank-1 тензоров [11]. На рис.2 показан многомерный тензор как 3-мерный массив данных размера $I1 \times I2 \times I3$.

Определение 1. Пусть A - тензор размерности $I1 \times I2 \times \dots \times IN$. Порядок A есть N, n-ая размерность (способ или

путь) A есть размер In . Нотация двоеточия использована, чтобы обозначить полную область данного индекса (i -тая колонка матрицы дается $A(i, :)$ и j -тый столбец – $A(:, j)$).

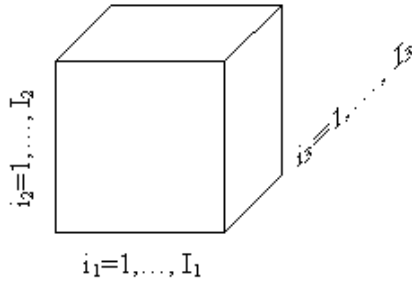


Рис. 2. Представление 3-х мерного массива
 Рассмотрим тензор 3-го-порядка (рис.2). В этом случае определение единственного индекса связано с вырезкой (*slice*), которая является матрицей в специфической ориентации. Так $A(i, :, :)$ формируют i -тую горизонтальную вырезку, $A(:, j, :)$ - j -тая боковая вырезка и $A(:, :, k)$ k -тая фронтальная вырезка.

Определение 2. Обобщенная матрицизация. Пусть A есть $I_1 \times I_2 \times \dots \times I_N$ - тензор r и предположим, что необходимо перестроить его в матрицу размером $J_1 \times J_2$. Число входов в матрицу должно быть таким же, что и число входов в тензор $\prod_{n=1}^N I_n = J_1 \cdot J_2$. Данные J_1 и J_2 удовлетворяют следующему свойству: отображение может быть задано любым количеством измерений такой длины, что имеем отображение π один в один, так что

$$\pi: \{1, \dots, I_1\} \times \{1, \dots, I_2\} \times \dots \times \{1, \dots, I_N\} \rightarrow \{1, \dots, J_1\} \times \{1, \dots, J_2\}.$$

Класс *tensor_as_matrix* (*tensor toolbox Mat-Lab*) поддерживает конверсию тензора в матрицу. Рассмотрим n -мерную матрицизацию. Обычно тензор матрицизирован так, что все его волокна связаны с отдельными измерениями, выровненными как колонки результирующей матрицы. Это специальный случай обобщенной матрицизации, где только одна размерность отображается в строке, так что $K = 1$ и $\{r_1\} = \{n\}$. Результирующая матрица типично обозначена как

$A(n)$. Колонки могут быть упорядочены по любому направлению.

Одно из преимуществ матрицизации состоит в том, что тензор, хранимый в матричной форме, может быть манипулирован как матрицы, сокращая n -мерное умножение. Если $B = A \times_n M$, тогда $B(n) = MA(n)$. Последовательность n -мерных произведений, которая записана в матричной формулировке, может быть представлена как Кронекерова последовательность произведений, включающих U матрицы. Последовательность n -мерных произведений имеет вид

$$B_{(n)} = U^{(n)} A_{(n)} (U^{(c_{n-1})} \otimes U^{(c_{n-2})} \otimes \dots \otimes U^{(c_1)})$$

Можно также создать два дополнительных класса для поддержки представления тензоров в декомпозированной форме в виде суммы *rank-1* тензоров. *rank-1* тензор – это тензор, который может быть записан как внешнее произведение векторов, т. е

$$A = \lambda u(1) \circ u(2) \circ \dots \circ u(N),$$

где λ - скаляр и каждый $u(n)$ – In -вектор для $n = 1 \dots N$. Символ \circ обозначает внешнее произведение; в случае $(i_1 i_2 \dots i_N)$ входа A представляется в виде

$$A(i_1, i_2, \dots, i_N) = \lambda u_{i_1}^{(1)} u_{i_2}^{(2)} \dots u_{i_N}^{(N)},$$

где u_i обозначает i -ый вход вектора u . Существует несколько декомпозиций тензора. Декомпозиция *SVD* (показана на рис.3), применяемая к тензору, дает

$$A = \sum_{k=1}^r u_1 \times u_2 \times u_3 \times w_k,$$

$$A = \sum_{i=1}^{r_1} \sum_{j=1}^{r_2} \sum_{k=1}^{r_3} \sigma_{ijk} u_i \circ v_j \circ w_k,$$

где σ – корневой (сердцевинный) тензор, U, V, W – компоненты или факторы, корневой тензор – диагональный, колонки в компонентах ортонормальны.

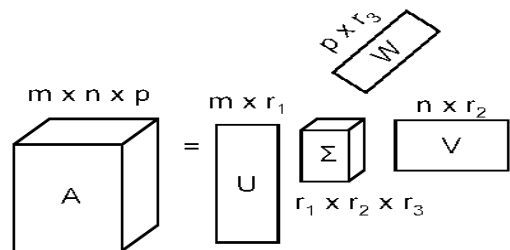


Рис.3. Декомпозиция *SVD*

Определение аномалий трафика рассмотрено в [11]. Трафик сети рассматривается на уровне существования тензоров порядка $M=3$ с направлениями "IP-источник", "IP-приемник", "порт".

На рис. 4 представлен сетевой поток данных как тензор 3-го ранга, разложенный на корневой тензор и 3 матрицы.

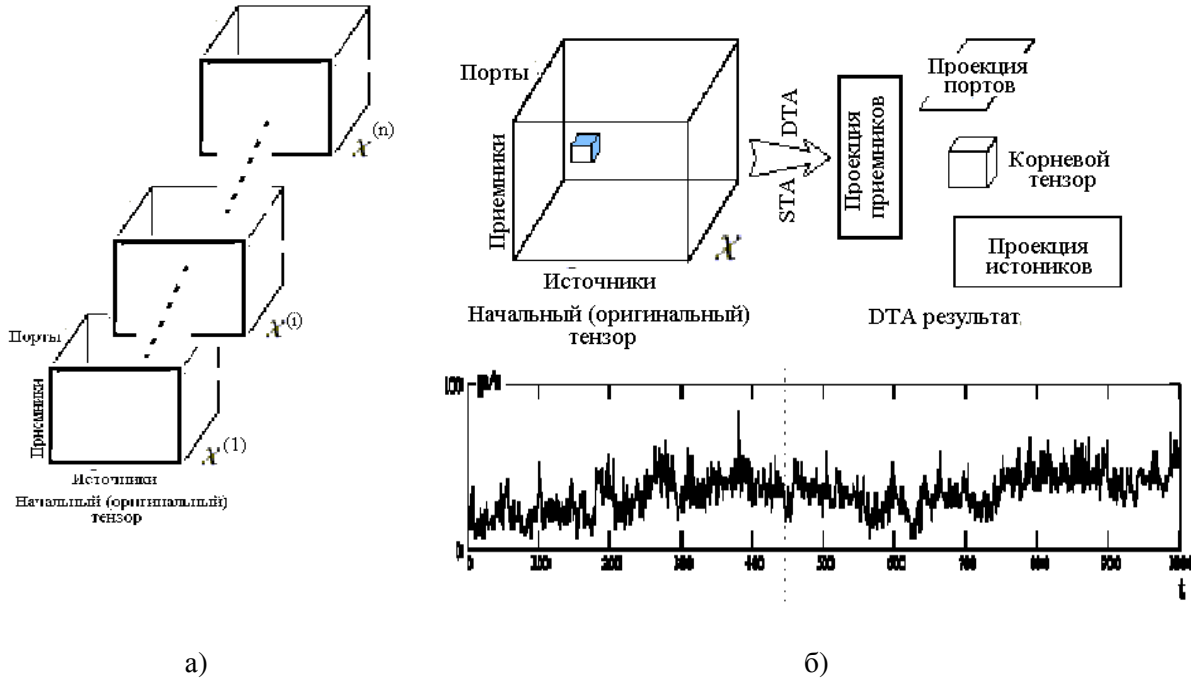


Рис.4 Тензор 3-го порядка сетевого потока данных: а) сетевой поток, б) разложение начального тензора и кривая количества пакетов трафика., пак/с.

Для примера сетевой поток как тензор 3-го порядка, имеющий для рассматриваемого периода времени три измерения: источник, приемник и порт, которые могут быть рассмотрены как 3D куб данных (рис.2). Вход (i, j, k) в этом тензоре имеет количество пакетов от соответствующего источника (i) в приемник j через порт k , в течение данного периода времени. Динамический аспект исходит из того, что новые тензоры прибывают непрерывно во времени.

В тестовой сети система определяет наличие аномалии, анализируя образцы трафика на уровне множеств – источники – приемники - порты соответственно. Для 1-го уровня (тензор уровень) моделируем аномальность тензора \mathfrak{X} через реконструкцию ошибки:

$$e_i = \|\mathfrak{X}_i - \Psi_1 \prod_{i=1}^M U_i^T I_F\|_F = \|\mathfrak{X}_i - \mathfrak{X}_i \prod_{i=1}^M U_i U_i^T I_F^2\|_F$$

Для 2-го уровня (уровень направления – отдельного измерения), реконструкция ошибки 1-го направления включает только одну проекцию матрицы U_l для данного тензора \mathfrak{X} :

$$e_d = \|\mathfrak{X} - \mathfrak{X} \times_l U_l U_l^T\|_F^2$$

Для 3-го уровня (уровень размерности), ошибка размерности d на l -ом направлении есть только реконструкция слайса размерности d тензора вдоль l -го направления. Формально уровень ошибки определяется из выражения

$$e_{T+1} \geq \text{mean}(e_{i|i=1}^{T+1}) + \alpha \cdot \text{std}(e_{i|i=1}^{T+1})$$

Условие $e_{T+1} \geq e_{T+1}$ есть признаком аномалии.

Рассмотрим идентификацию реальной аномалии. Главные образцы связей сети должны в быт подобными друг другу на определенных временных интервалах. На рис. 5(а) показана относительная реконструкция ошибки во времени, определения на основе DTA.



Рис.5 Определение аномалии трафика по аппроксимационной точности

Точки аномалии расположены выше строки – порога среднеквадратичного отклонения процента ошибки. Некоторые взрывные ошибки происходят в интервале времени 140-й и 160-й час (рис.5а). Рис. 5b и 5c демонстрируют нормальный и аномальный трафик.

Идентификация аномального состояния КС на основе структурированного тензор-трафика использует специальные математические методы нелинейной и тензорной аппроксимации.

Пусть A – двухуровневая матрица с размерами уровней $n1$ и $n2$. Приближим ее суммой кронекеровых произведений вида $A = \sum_{k=1}^r A_1^k \otimes A_2^k$, где размеры A_1^k, A_2^k соответственно $n1 \times n1$ и $n2 \times n2$. Пусть $A=Ar$ и r - наименьшее возможное число кронекеровых произведений, сумма которых есть A . Тогда r называется тензорным рангом A . Наилучшие приближения, минимизирующие $\|A-Ar\|_F$, могут быть получены алгоритмом SVD, примененным к матрицам, полученных специальным «перемешиванием» элементов.

Проблема трilinearной аппроксимации тензора формулируется следующим образом. Пусть дан трехмерный массив (тензор) $A=aijk$ размером $n1 \times n2 \times n3$. Требуется найти такие матрицы $U=ui\alpha$ $V=vja$, $W=wka$ размеров $n1 \times r$, $n2 \times r$, $n3 \times r$, которые минимизируют функционал

$$\left\| a_{ijk} - \sum_{\alpha=1}^r u_{i\alpha} v_{j\alpha} w_{k\alpha} \right\|,$$

где $\|\Phi_{ijk}\| = \left(\sum_{ijk} \Phi_{ijk}^2 \right)^{1/2}$.

Одно из решений задачи состоит в следующем. Пусть даны некоторые при-

ближения U, V, W к решению U^*, V^*, W^* проблемы. Тогда по V и W находится новое значение U из решения задачи наименьших квадратов

$$U = \underset{U}{\operatorname{argmin}} \left(\left\| a_{ijk} - \sum_{\alpha=1}^r u_{i\alpha} v_{j\alpha} w_{k\alpha} \right\|^2 \right).$$

Задача распадается на $n1$ независимых подзадач: $i=1,2,\dots, n$.

Представление трафика КС тензором позволяет агрегировать все параметры трафика и унифицировать его анализ. Авторами предложен принцип структурирования трафика, представленного в виде тензора четных рангов [7, 13], что позволило создать обобщенную модель трафика. Создание данного тензора реализовано путем тензорного произведения строки (вектора) безразмерных (приведенных) параметров трафика [13] - $X=\{xj\}$ на колонку (вектор) функций псевдопринадлежности

$$M=\{\mu j\}, j=1,2,\dots,9; T=(X=\{xj\}) \bullet (M=\{\mu j\} T), \text{ где } \bullet - \text{знак тензорного произведения, } T=[tij], i,j=1,9; T - \text{символ транспонирования; } x=[x1, x2, x3, x4, x5, x6, x7, x8, x9],$$

$$x^{(1)} = [1 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 1] T, \text{ TRAFIC } T = x \otimes^{(1)} x$$

$$\text{TRAFIC } T = \begin{pmatrix} x1 & x2 & x3 & x4 & x5 & x6 & x7 & x8 & x9 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ x1 & x2 & x3 & x4 & x5 & x6 & x7 & x8 & x9 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ x1 & x2 & x3 & x4 & x5 & x6 & x7 & x8 & x9 \end{pmatrix}$$

В связи с тем, что трафик содержит 3^2 параметров, получаем тензор 4-го ранга.

Объект, который определяется совокупностью коэффициентов $a_{ijk\dots m}$ полилинейной формы $\varphi = \varphi(x, y, z, \dots, w)$, является ортогональным тензором. Совокупность коэффициентов a_{ij} билинейной формы $\varphi = \varphi(x, y)$, образующая матрицу $A = (a_{ij})$, представляет собой тензор валентности 2. Следовательно, компоненты тензора b_k, \dots, m валентности $p - 2$ определяются по формуле $b_k \dots m = a_{11k\dots m} + a_{22k\dots m} + \dots + a_{33k\dots m}$ [12]. Операция получения тензора $b_k \dots m$ из тензора $a_{ijk\dots m}$ называется свертыванием тензора $a_{ijk\dots m}$ по индексам i и j . В результате свертывания тензоров валентностей p и q получается тензор валентности $p + q - 2$. Свертывание тензоров можно производить по любому количеству r таких пар. Валентность результирующего тензора на $2r$

единиц меньше суммы валентностей исходных тензоров.

Свертка тензора ${}^{\text{TRAFIC}}T$ валентности 4 дает тензор валентности 2:

$${}^{\text{TRAFIC}}T = X \otimes^{(1)} X \rightarrow \begin{pmatrix} X_1 & X_2 & X_3 \\ X_4 & X_5 & X_6 \\ X_7 & X_8 & X_9 \end{pmatrix} = (t_{ij}),$$

который в дальнейшем будем называть тензор-трафик (рис.6). Возвращаясь к тензор-трафику, получаем множество матриц 3×3 (рис.7), представленных как $3 \times 3 \times K$ тензор M (где K количество отсчетов трафика).

Для идентификации аномальных режимов трафика можно использовать инварианты тензор-трафика, значения которых для нормального и аномального состояний существенно различаются.

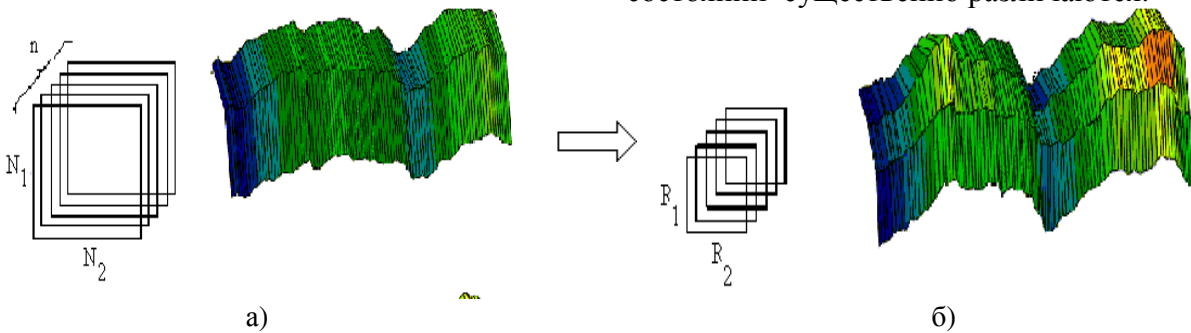


Рис.6. а) реальный тензор-трафик $9 \times 9 \times n$, б) структурированный тензор-трафик $3 \times 3 \times n$

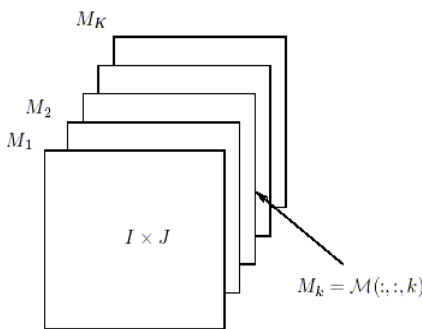


Рис. 7. Матричная модель тензор-трафика

Предположение 1. Аномальное состояние КС определяем по величинам инвариантов тензора и системы правил: если инварианты $I = \{I_1, I_2, \dots, I_9\}$ или собственные значения матрицы тензор-трафика $A = \{\lambda_1, \lambda_2, \dots, \lambda_9\}$ для моментов времени $t^{(1)}$ и $t^{(2)}$ существенно отличаются, т.е. $I^{(1)} \neq I^{(2)}$ или $A^{(1)} \neq A^{(2)}$, где знак \neq

обозначает отсутствие субъективной близости, то имеет место аномальное состояние [13].

Предположение 2. Аномалией может считаться необычное сочетание значений инвариантов тензорного произведения вектора параметров трафика и специального единичного вектора с последующим $I^{(s)} = \{1 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 1\}$, \otimes - символ Кронекера произведения.

Введем понятия последовательности тензоров и тензорного потока [11]. Последовательность тензоров M -го порядка $\mathfrak{K}^1 \dots \mathfrak{K}^n$, где каждый $\mathfrak{K}^i \in \mathbb{R}^{N_1 \times \dots \times N_M}$ ($1 \leq i \leq n$), есть тензорная последовательность, если n есть фиксированное натуральное число и n есть кардиналом тензорной последовательности. Последовательность тензоров M -го

порядка $\aleph_1 \dots \aleph_n$, где каждый $\aleph_i \in \mathbb{R}^{N_1 \times \dots \times N_M}$ ($1 \leq i \leq n$), названа тензорным потоком, если n есть целое, увеличивающееся во времени.

Тензорный анализ трафика рассматриваем в виде. Дана последовательность тензоров $\aleph_1 \dots \aleph_n$, где $\aleph_i \in \mathbb{R}^{N_1 \times \dots \times N_M}$ ($1 \leq i \leq n$), находим ортогональные матрицы $U_i \in \mathbb{R}^{N_i \times R_i}$ $\prod_{i=1}^M$, для каждого направления, так что ошибка реконструкции есть минимальной:

$$\epsilon = \sum_{i=1}^n \left\| \lambda_i - \lambda_i \prod_{i=1}^M \times_i (U_i U_i^T) \right\|_F^2$$

Входные данные 3D тензора могут быть рассмотрены как $X = \{X_{ijk}\}_{i=1}^{n_1} \{j=1}^{n_2} \{k=1}^{n_3}$, где каждый X_i есть 2D матрица (изображение) размерностью $n_1 \times n_2$.

Выводы

Тензор, представлений в виде мультимерного массива, обеспечивает естественное представление трафика КС.

Многомерные компьютерные сети, один из вариантов которых рассматривает каждую отдельную компоненту трафика в отдельном измерении, дают возможность решения задач анализа трафика.

Предложен метод структурирования многомерного трафика, представленного в виде тензора четных рангов. Показана возможность применения предложенных моделей для выявления аномалий трафика.

Список литературы

1. Минаев Ю.М., Филимонова О.Ю. Тензорная модель трафика компьютерных систем. Сборник трудов конференции Моделирование-2008, том 2. – С. 461–466
2. Коновалов Г.В. Многомерные сети – будущее инфокоммуникационных сетей. «Электросвязь», № 4, 2008. – С.28–34
3. Соколов Н.П. Пространственные матрицы и их приложения / М.: Физматгиз, 1960. – 352 с.
4. Муха В.С. Анализ многомерных данных / Минск.: Технопринт, 2004. – 124 с.

5. Крон Г. Тензорный анализ сетей: Пер.с англ./Под ред. Л.Т.Кузина, П.Г.Кузнецова. – М.: Сов. Радио, 1978. – 720 с.

6. Минаев Ю.Н., Филимонова О.Ю., Гузий Н.Н. Интеллектуальные технологии в системах идентификации и прогнозирования атак на компьютерные сети Электронное моделирование, Т.27. №6. 2005. – С. 37–52.

7. Mark Sears, Brett Bader, Tammy Kolda. Parallel Implementation of Tensor Decompositions for Large Data Analysis. - SIAM AN09 July 8, 2009. – P. 17–25.

8. Brett W. Bader, Tamara G. Kolda. MATLAB Tensor Classes for Fast Algorithm Prototyping. - SANDIA REPORT SAND2004-5187 Unlimited Release Printed October, 2004 <http://www.ntis.gov/ordering.htm>

9. Kolda, T.G. Jimeng Sun. Scalable Tensor Decompositions for Multi-aspect Data Mining Data Mining, 2008. ICDM '08. Eighth IEEE International Conference on. P. 363–372

10. Evrim Acar Daniel M. Dunlavy Tamara G. Kolda Link Prediction on Evolving Data using Matrix and Tensor Factorizations. <http://www.csmr.ca.sandia.gov/~tgkolda/ref/>.

11. Jimeng Sun, Dacheng Tao, Christos Faloutsos. Beyond Streams and Graphs: Dynamic Tensor Analysis. 2006 ACM 1595933395/ 06/0008.- <http://www.cs.cmu.edu/~christos/PROJECTS/GRAPH-MINING/>.

12. Акивис М. А., Гольдберг В. В. Тензорное исчисление: Учеб. пособие. – 3-е изд., перераб. – М.: Физматлит, 2003. – 304 с.

13. Минаев Ю.Н., Филимонова О.Ю. Интеллектуальні технології прогнозування часових рядів на підставі тензорних інваріантів //Зб. наук. праць “Проблеми інформатизації та управління”, №2(26), 2009. –С.104–112

Подано до редакції 01.10.2010