

МОДЕЛЬ СИСТЕМИ ВИЯВЛЕННЯ ВТОРГНЕНЬ З ВИКОРИСТАННЯМ ДВОСТУПЕНЕВОГО КРИТЕРІЮ ВИЯВЛЕННЯ МЕРЕЖНИХ АНОМАЛІЙ

У статті розглянуто методи, що впливають на роботу систем виявлення вторгнень. Застосовуючи «задачу про розладнання» сформовано двоступеневий критерій виявлення аномалій в комп'ютерних мережах, що забезпечить проведення аналізу характеристик мережевої інфраструктури та їх ототожнення з певними комп'ютерними атаками та надасть можливість реагування на можливі атаки в реальному масштабі часу.

Ключові слова: системи виявлення вторгнень, аномалії, сигнатурний метод аналізу, статистичний метод аналізу, нейронні мережі.

Вступ і постановка задачі

Виходячи з даних про щорічні темпи зростання кількості інцидентів у кібернетичному просторі можна зробити висновок про необхідність обов'язкового включення до складу комплексних систем інформаційної безпеки критично-важливих об'єктів інфраструктури автоматизованих засобів виявлення комп'ютерних атак та інших небезпечних подій, включаючи загрози техногенного та природного (випадкового) характеру.

Сучасні системи виявлення вторгнень (СВВ, англ. *Intrusion Detection Systems - IDS*) для досягнення цілей інформаційної безпеки здійснюють постійний моніторинг стану функціонування апаратних і програмних платформ мережевої інфраструктури та реєструють певну множину кількісних і якісних показників їх роботи:

$$\{X_1(t), \dots, X_N(t)\}, \quad (1)$$

де t - момент часу у який здійснюється вимірювання показника $X_k(t)$.

Такими показниками, зокрема у системі *NIDES* [1], є:

- здатність використання CPU окремо системою та користувачем;
- час на виконання процесу;
- загальний обсяг пам'яті що використана під час виконання процесу та його максимальний розмір під час виконання;
- кількість відкритих файлів під час виконання;
- кількість збоїв сторінок;
- обсяг зчитаної з диска інформації;
- кількість символів вводу/виводу під час виконання додатка;
- чи змінювалось ім'я користувача під час виконання додатка;
- час початку виконання додатка;
- кількість сигналів що отримані під час виконання додатка;
- чи виконувався додаток на віддаленій станції та ім'я цієї станції;
- ім'я додатка що був використаний на віддаленій станції;
- чи виконувався додаток на локальній станції та ім'я цієї станції, ім'я додатка використаного на локальній станції тощо.

Зважаючи на це, метою даної статті є побудова такої моделі системи виявлення вторгнень, яка б дозволила виявляти мережні аномалії в реальному масштабі часу та з найменшою похибкою, а також знаходити відповіді на наступні питання [2]:

- що відбулось у мережі?;
- що зазнало нападу та наскільки небезпечна атака?;
- коли та звідки почалася атака?;
- хто зловмисник?;
- яким чином та внаслідок чого відбулося вторгнення?

Виклад основного матеріалу досліджень

Відомо, що ефективність роботи системи СВВ суттєво залежить від застосованого

методу аналізу вихідних даних. Нині розрізняють наступні основні методи [1]:

- 1) сигнатурні методи аналізу;
- 2) статистичні методи аналізу;
- 3) гібридні методи аналізу з функцією самонавчання.

Сигнатурні методи аналізу засновані на тому, що більшість атак та їх сценаріїв у загальних рисах відомі. У даному підході сигнатури вторгнень визначають характерні особливості та умови функціонування об'єктів, виникнення подій, що є ознаками спроб атак (вторгнення), та їх взаємозв'язку. Звичайно, сигнатурні методи аналізу використовують бази даних сигнатур вторгнень, яка підтримується системою безпеки. В цьому випадку порядок (послідовність) дій, що виконуються або ініціюються користувачем або інформаційним процесом (програмою), порівнюється з відомими сигнатурами. Ознакою спроби порушення безпеки може служити часткова відповідність послідовності подій сигнатурі.

Типовими представниками, що реалізують дану ідею, є антивірусні сканери, що працюють з базою даних сигнатур вірусів. Їх перевагою є достатньо висока швидкодія проведення аналізу. Ефективність сигнатурних методів СВВ може бути підвищена за рахунок застосування методів штучного інтелекту.

До переваг статистичних методів аналізу в СВВ слід віднести:

- відсутність потреби у великому об'ємі пам'яті для зберігання змінних, що контролюються;
- простоту виявлення відхилень у даних що характеризують поведінку користувачів та процесів;
- можливість аналізу кількісних та якісних даних різної природи походження у якості параметру при аналізі.

До недоліків статистичних методів відносять труднощі з формуванням статистики звичайної поведінки користувачів та процесів.

До гібридних методів аналізу з функцією самонавчання відносяться:

- методи навчання на класифікації прикладів;
- нейромережі;
- генетичні алгоритми.

Для досягнення поставленої у роботі мети доцільно використовувати саме статистичні методи. Це дозволить, застосовуючи відому «задачу про розладнання» [3]:

- сформуванню двоступеневий критерій виявлення аномалій в комп'ютерних мережах,
- подолати прив'язку статистичних методів до моделі звичайної (нормальної) поведінки користувачів.

На першому кроці має бути досліджено статистику аномальної поведінки (стану) мережі. При цьому виявлення змін станів (з будь-яким ступенем точності) може бути зведено до виявлення зміни математичного очікування (МОЧ) в деякій новій випадковій послідовності, сформованій з вихідної.

Пояснимо це положення. Нехай стосовно випадкової послідовності, що аналізується, $X = \{x_t, x_t \in R, t = \overline{1, N}\}$ розглядаються дві складних гіпотези: H_0 : X є стаціонарною послідовністю з єдиною функцією розподілу ймовірностей, H_1 : X є конкатенацією (результатом "склеювання") двох стаціонарних випадкових послідовностей з різними функціями розподілу:

$$X = X_1 || X_2, \text{ де } X_1 = \{x_t, t = \overline{1, n^*}\}, X_2 = \{x_t, t = \overline{n^* + 1, N}\}, n^* = [\theta N], 0 < \theta < 1. \quad (2)$$

Потрібно оцінити точку «склейки» n^* .

Вважається, що послідовності X_1 і X_2 відрізняються між собою однією з двовимірних функцій розподілу, а саме, розподілу ймовірностей вектору (x_t, x_{t+2}) :

$$F(u_0, u_1) = P\{x_t \leq u_0, x_{t+2} \leq u_1\} \quad (3)$$

до моменту $t_1^* = n^* - 2$ включно дорівнює $F_1(u)$, а при $t \geq t_2^* = n^* + 1$ дорівнює $F_2(u)$, причому

$$\|F_1(u) - F_2(u)\| \geq \varepsilon > 0, \text{ де } \|\dots\| - \text{звичайна sup-норма.} \quad (4)$$

Відомо [3], що функція розподілу кінцевомірного випадкового вектора може бути наближена рівномірно з будь-якою точністю функцією розподілу ймовірностей випадкового вектора з кінцевим числом значень. Звідси випливає, що якщо подати множину R у вигляді поєднання досить великої кількості областей $\{A_j, j = \overline{1, r}\}$, що не перетинаються $A_i \cap A_j = \emptyset$ для $i \neq j$, то вектор (x_t, x_{t+2}) можна апроксимувати по розподілу вектором з кінцевим числом значень.

Тому, якщо ввести нові випадкові послідовності

$$V_t^{ij} = I(x_t \in A_i, x_{t+2} \in A_j), \text{ де } 1 \leq i \leq r, 1 \leq j \leq r, \quad (5)$$

де $I(A)$ - індикатор множини A , то хоча б в одній з них відбувається зміна МОЧ.

Отже, якщо скористатися алгоритмом, який виявляє зміну математичного очікування, то цей же алгоритм виявить і зміну функції розподілу. Ця обставина дозволила в роботі [3] обмежитися розробкою тільки одного, базового алгоритму, який може виявляти зміну МОЧ. Для цього з метою виявлення моментів "розладнань" запропоновано сімейство статистик виду:

$$Y_N(n, \delta) = \left[\frac{n}{N} \left(1 - \frac{n}{N} \right) \right]^\delta \cdot \left[\frac{1}{n} \sum_{k=1}^n x_k - \frac{1}{N-n} \sum_{k=n+1}^N x_k \right], \quad (6)$$

де $0 < \delta \leq 1, 1 \leq n \leq N - 1, X = \{x_k, k = \overline{1, N}\}$ - послідовність, що досліджується.

Наведене сімейство статистик у випадку фіксованого n є узагальненим варіантом статистики Колмогорова - Смірнова, що використовується для перевірки гіпотез збігу або відмінності функцій розподілу у двох вибірках. В роботі [3] також доведено, що статистика виду (6) в разі $\delta = 1$ при $N \rightarrow \infty$ та збереженні співвідношення між обсягами "склеєних" реалізацій мінімізує максимально можливу ймовірність помилки оцінювання моменту "розладнання" (мінімаксна по порядку).

При цьому:

$$P\{\max_{1 \leq n \leq N-1} \sqrt{N} |Y_N(n, 1)| > C^{(1)}\} \rightarrow 2 \sum_{k=1}^{\infty} (-1)^{k+1} \exp(-2k^2 \left(\frac{C^{(1)}}{\sigma_*}\right)^2) \equiv f(C^{(1)}), \quad (7)$$

де параметр σ_* є стандартним відхиленням, $C^{(1)}$ - границя критерію, перевищення якої буде сприйматися, як виникнення «розладнання», значення n_* , для якого це відбулося, є шуканим моментом «розладнання».

На другому кроці, зафіксувавши рівень ймовірності α "помилкової тривоги" про розладнання, під час статистичної обробки реальних даних визначимо рівень порога першого рівня $C^{(1)}$:

$$\alpha = f(C^{(1)} \sqrt{N} / \bar{\sigma}_*), \quad (8)$$

де $\bar{\sigma}_*$ є оцінкою параметра σ_* (стандартне відхилення), а N є обсягом вибірки - послідовності, що досліджується.

У випадку гіпотези H_1 про подання вихідної послідовності вимірювань у вигляді конкатенації декількох стаціонарних випадкових послідовностей з різними функціями розподілу ймовірностей, застосуємо критерій припустимої кількості спрацювань («розладнань») Z . Тобто, вважаємо що має місце наступне рівняння:

$$n_1 + n_2 + \dots + n_Z = N, \text{ де } 2 \leq Z < N - 2. \quad (9)$$

Для невеликої кількості послідовностей з метою визначення границі критерію можна скористатися нерівністю Чебишева, якщо випадкова величина Z має математичне очікування μ та стандартне відхилення σ для заданого $\varepsilon > 0$:

$$P\{|Z - \mu| \geq \varepsilon\} \leq \frac{\sigma^2}{\varepsilon^2}. \quad (10)$$

Якщо Z є випадковою величиною з одномодовим розподілом ймовірностей з математичним очікуванням μ та стандартним відхиленням $0 < \sigma < \infty$, то для будь якого $\lambda > \sqrt{8/3} \approx 1.63299 \dots$, має місце нерівність Височанського – Петунина [6], що покращує оцінку ймовірності відхилення:

$$P\{|Z - \mu| \geq \lambda\sigma\} \leq \frac{4}{9\lambda^2}. \quad (11)$$

Зокрема, для типового відхилення у 3σ (три сигма) рівень значущості критерію - ймовірність "помилкової тривоги" обчислюється як:

$$\alpha = \frac{4}{9\lambda^2} \approx 0.0494.$$

Більш точний результат можливо отримати для випадку випадкової величини Z , що є сумою достатньо великої кількості незалежних випадкових величин ($m \rightarrow \infty$):

$$Z = z_1 + z_2 + \dots + z_m.$$

Згідно інтегральної граничної теореми [5] ймовірність відхилення можливо апроксимувати за допомогою функції нормального розподілу ймовірностей $N(0,1)$:

$$P\left\{\left|\frac{Z-\mu}{\sigma}\right| \geq t_{1-\alpha}\right\} = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{t_{1-\alpha}} \exp\left(-\frac{t^2}{2}\right) dt. \quad (12)$$

На підставі (3) розраховуємо границю критерію другого рівня:

$$C^{(2)} = \mu + t_{1-\alpha}\sigma.$$

Таким чином, алгоритм обчислення критерію включає наступні кроки (рис. 1):



Рис. 1. Блок – схема алгоритму формування дворівневого критерію виявлення мережних аномалій

Можливо бачити, що складність алгоритму що реалізує запропонований критерій оцінюється величиною $O(N)$. Таким чином, його реалізація за принципом «ковзаючи вікно» не викличе суттєвого навантаження на обчислювальну систему.

Найбільш складними питаннями для реалізації методу виявлення аномалій, як і в інших статистичних методах, є формування штатної поведінки інфраструктури мережі та її користувачів.

Висновки:

1. Запропоновано дворівневий статистичний критерій виявлення аномалій на основі вимірювань характеристик мережевої інфраструктури, що забезпечує можливість їх подальшого аналізу з метою ототожнення з певними комп'ютерними атаками.

2. Застосування критерію не потребує значних обчислювальних ресурсів та надає можливість реагування на можливі атаки практично в реальному масштабі часу.

Література

- 1.Корт С.С. Методы обнаружения нарушителя, <http://www.uran.donetsk.ua/~masters/2011/fknt/brich/library/article6.htm>
- 2.Костров Дм. Системы обнаружения атак. <http://www.bytemag.ru/articles/detail.php?ID=6608>
- 3.Brodsky V. E., Darkhovsky V. S., Non-Parametric Statistical Diagnosis: Problems and Methods, Kluwer, Dordrecht, 2000, 452 pp
- 4.Справочник по прикладной статистике. Под ред. Э. Ллойда, У. Ледермана. Том. 2. Перевод с англ. под ред. С.А. Айвазяна и Ю. Тюрина. М.: Финансы и статистика, 1990. - 526 с.
- 5.Ширяев А.Н. Вероятность: В 2-х кн. – 4-е изд., переработ. и доп. – М.: МЦНМО, 2007, Кн. 1. -552 с.
- 6.Высочанский Д. Ф., Петунин Ю. И. Обоснование правила 3-sigma для одномодальных распределений. — Теория вероятностей и мат. статистика, 1979, вып. 21, с. 23-35.

Надійшла 15.11.2015 р.

Рецензент: д.т.н., проф. Бурячок В.Л.