


Application of OPTICS and ensemble learning for Database Intrusion

Поляков Степан, Грицаенко Александра, Екатерина Степура, Слипченко Вера,
Хугаева Дана, Рублева Вита, Тихомирова Алина

Актуальность

- Базы данных — ядро информационных систем
- Хранят критически важные сведения (финансы, клиенты, наука)
- Потери организаций $\approx 5\%$ дохода/год (ACFE, 2016)
- Атаки бывают внешние и внутренние

Проблема

Внешний атакующий	База данных	Внутренний атакующий
уязвимости, слабые пароли		злоупотребление правами

Стандартная защита:

- аутентификация
- контроль доступа
- пароли
- firewall (межсетевые экраны)
- шифрование



НЕДОСТАТОЧНО

Системы DIDS

Ключевые требования:

- высокая точность
- низкий уровень ложных срабатываний
- надежность и скорость
- работа в реальном времени

Существующие подходы

Год	Метод	Описание	Методы
2015	ANFIS (нейро-нечёткие правила), <i>Brahma & Panigrahi</i>	Система на основе нейро-нечёткой логики. Строит правила поведения пользователей, транзакции вне правил считаются подозрительными.	ANFIS – Adaptive Neuro-Fuzzy Inference System; Нечёткая логика – работа с неопределённостью; Нейронные сети.
2016	Ассоциативные правила + кластеризация, <i>Singh et al.</i>	Используются ассоциативные правила для профиля пользователей и кластеризация для объединения схожих транзакций.	Ассоциативные правила – закономерности вида 'если А, то В'; Кластеризация – группировка объектов по схожести.
2016	PCA + Random Forest, <i>Ronao & Cho</i>	PCA уменьшает размерность данных, Random Forest классифицирует транзакции, выявляя аномалии.	PCA – метод главных компонент; Random Forest – ансамбль деревьев решений.

Существующие подходы

Год	Метод	Описание	Методы
2017	CNN + система правил LCS, <i>Bu & Cho</i>	Гибридный метод: система правил LCS выделяет аномалии, CNN уточняет классификацию.	CNN – свёрточная нейросеть; LCS – система обучающихся классификаторов.
2018	Квантовые протоколы, <i>Wei et al.</i>	Применение квантовой криптографии для защиты транзакций и запросов.	Квантовая криптография – методы защиты на основе квантовой механики.
2018	Контроль доступа по схеме БД, <i>Yesin et al.</i>	Использование структуры базы данных для построения правил доступа.	Контроль доступа – система разрешений для пользователей.
2018	Наивный Байес для SQL-запросов, <i>Jayaprakash & Kandasamy</i>	Применение Наивного Байеса для классификации SQL-запросов как нормальных или атакующих.	Наивный Байес – вероятностный метод классификации, предполагающий независимость признаков.

Основная проблема

- IDS часто генерирует слишком много тревог
- Эффект базовой частоты (Axelsson, 2000)

Даже редкие атаки -> низкая эффективность IDS, если высокий уровень ложноположительных результатов

Зачем нужен новый подход?

Необходим метод, который:

- сохраняет высокую точность
- уменьшает количество ложных тревог

Предложение авторов:

- кластеризация OPTICS (поиск аномалий)
- ансамблевое обучение (устойчивая классификация)

Обзор используемых методов - Кластеризация с использованием OPTICS

Вход:

ε — радиус окрестности

MinPts — минимальное количество точек,
необходимое для образования кластера.

Выход:

Упорядоченная последовательность точек,
отражающая структуру кластеризации

Основная идея алгоритма OPTICS: Для каждого объекта k в кластере существует как минимум MinPts объектов в его ε -окрестности.

Ядровая точка — объект, в ε -окрестности которого находится как минимум MinPts объектов ($\geq \text{MinPts}$).

$$cd_{\varepsilon, \text{MinPts}}(k) = \begin{cases} \text{undefined} & \text{if } |N_{\varepsilon}(k)| < \text{MinPts} \\ \text{Minpts_distance}(k) & \text{otherwise} \end{cases}$$

$$rd_{\varepsilon, \text{MinPts}}(k, q) = \begin{cases} \text{undefined} & \text{if } |N_{\varepsilon}(k)| < \text{MinPts} \\ \max(cd(q), \text{dist}(q, k)) & \text{otherwise} \end{cases}$$

$$lrd_{\text{MinPts}}(k) = 1 / \frac{\sum_{o \in N_{\text{MinPts}}(k)} rd_{\infty, \text{MinPts}}(k, o)}{|N_{\text{MinPts}}(k)|}$$

$$OF_{\text{MinPts}}(k) = \frac{\sum_{o \in N_{\text{MinPts}}(k)} \frac{lrd_{\text{MinPts}}(o)}{lrd_{\text{MinPts}}(k)}}{|N_{\text{MinPts}}(k)|}$$

Обзор используемых методов - Ансамблевое обучение

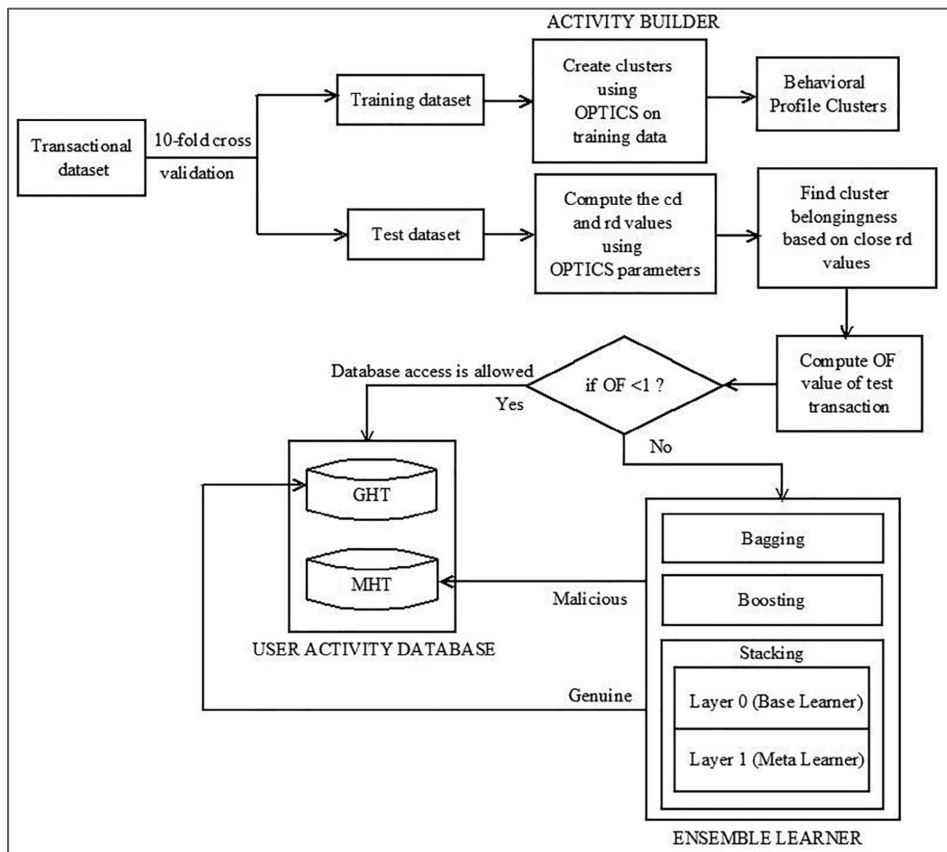
Ансамбль — комбинация или агрегирование нескольких классификаторов.

—————→ Снижение риска увеличения числа ложноположительных результатов за счёт выбора удачного сочетания разных классификаторов.

Типы ансамблевых методов:

- **Bagging** — метод построения ансамбля, при котором на разных подвыборках обучаются одинаковые базовые модели. Предсказания комбинируются путём усреднения или голосования.
- **Boosting** — итеративное построение ансамбля, где каждая следующая модель обучается на ошибках предыдущих. Ошибочно классифицированные объекты получают больший вес. Предсказания объединяются с учётом весов моделей.
- **Stacking** — ансамбль, состоящий из разнородных базовых моделей. Их предсказания используются в качестве признаков для обучения метамоделей («модели второго уровня»).

Предложенный подход



Activity Builder — модуль построения активности

Построение профилей пользователей базы данных на основе их прошлых транзакций.

Транзакция: $\langle u_id, transaction_id, query_type, table_list, att_list, time_slot, loc, time_gap \rangle$

Пользователь с $u_id = 10$:

Q1: SELECT x, y FROM T1 WHERE z = 1

Q2: DELETE FROM T2 WHERE w = 1

query_type = {SELECT, DELETE}

att_list = {z, x, y, w}

table_list = {T1, T2}

transaction_id = 1

query_type {SELECT, DELETE} = {1, 4}

att_list {z, x, y, w} = {40, 23, 12, 6}

table_list {T1, T2} = {3, 6}

loc = 1 (из офиса)

time_slot = 37 (в интервале с 18:00 до 18:30)

time_gap = 21

u_id	transaction_id	query_type	table_list	att_list	time_slot	loc	time_gap
1	10	select	3, 1, 14	1, 6, 2, 3	9	1	4
	11	alter	17, 33	29, 30, 31, 32, 25	12	2	1
2	78	update	10, 4	4, 6, 50	33	1	5

{10, 1, {1,4}, {3,6}, {40,23,12,6}, 37, 1, 21}

User Activity Database — база активности пользователей

u_id	transaction_id	query_type	table_list	att_list	time_slot	loc	time_gap	label
1	10	select	3, 1, 14	1, 6, 2, 3	9	1	4	genuine
	11	alter	17, 33	29, 30, 31, 32, 25	12	2	1	genuine
2	78	update	10, 4	4, 6, 50	33	1	5	genuine

(a) Genuine History Table (GHT)

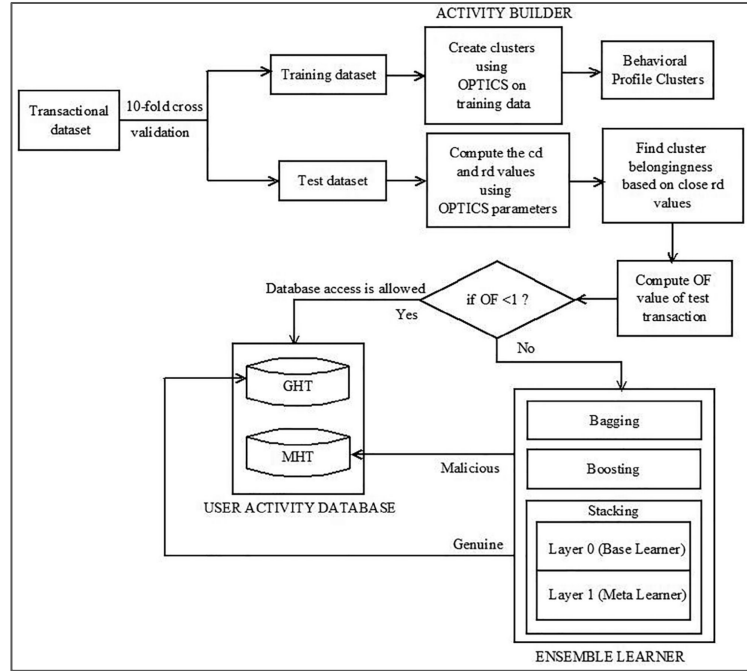
query_type	table_list	att_list	time_slot	loc	time_gap	label
delete	8, 11	16, 24, 5, 16	11	0	2	malicious
alter	55	55, 57, 75, 67	9	2	4	malicious
insert	66	1, 82, 40	27	1	1	malicious

(b) Malicious History Table (MHT)

Activity Builder — модуль построения активности

Базовые классификаторы:

- Наивный Байесовский классификатор (Naïve Bayes)
- Дерево решений (Decision Tree, DT)
- Индукция правил (Rule Induction, RI)
- Метод к-ближайших соседей (k-NN)
- Сети с радиальными базисными функциями (Radial Basis Function Network, RBFN)



Input: $D_u[m][n]$, ε , $MinPts$

Output: Genuine or Fraud

```

1:  $Norm\_D_u[m][n] = \text{normalize}(D_u[m][n])$   $\triangleright$  Data Normalization
2:  $[Train[m_1][n], Test[m_2][n]] = \text{divide}(Norm\_D_u[m][n])$   $\triangleright$  Train and Test sets are extracted
    $\triangleright m = m_1 + m_2$ 
3:  $Prof_u = \text{OPTICS\_cluster}(Train[m_1][n_1], \varepsilon, MinPts)$   $\triangleright$  Using Eq. (1)- Eq. (2)
    $\triangleright Prof_u = \text{User Profile of User } u$ 
4:  $rd[m_2] = \text{OPTICS\_cluster}(Test[m_2][n], \varepsilon, MinPts)$   $\triangleright$  Generate rd values for Incoming transaction
5:  $OF[m_2] = \text{OPTICS\_cluster}(rd[m_2])$   $\triangleright$  Using Eq.4
6: if  $OF[m_2] < 1$  then
7:   Output ("Database Access is allowed")
8:   Insert\_GHT ( $Test_x$ )  $\triangleright$  Insert  $Test_x$  in GHT
    $\triangleright Test_x \in Test[m_2][n], x = 1, \dots, m_2$ 
9: else
10:   goto EL  $\triangleright$  EL: Ensemble Learner
11: end if
   Collect History Information  $D(U_x)$  about user  $U_x$  from User Activity Database
   EL:
12:  $Bag = \text{Bagging}(Test_x, Prof_u)$   $\triangleright$  Bag: Bagging model
13:  $Boost = \text{Boosting}(Test_x, Prof_u)$   $\triangleright$  Boost: Boosting model
14:  $Stack = \text{Stacking}(Test_x, Prof_u)$   $\triangleright$  Stack: Stacking model
15: if  $(Test_x \in Bag) \vee (Test_x \in Boost) \vee (Test_x \in Stack)$  then
16:   Insert\_GHT ( $Test_x$ )  $\triangleright$  Insert  $Test_x$  in GHT
17:   Output ("Genuine")
18: else
19:   Insert\_MHT ( $Test_x$ )  $\triangleright$  Insert  $Test_x$  in MHT
20:   Output ("Malicious")
21: end if
    
```

Экспериментальные данные

- Использован симулятор транзакций с тремя модулями:
 - GTGM – генерация нормальных действий.
 - MTGM – моделирование атакующих транзакций.
 - MMPPM – управление частотой и структурой потоков.
- Итоговый объём данных: 41 390 транзакций (нормальные + атакующие).

Настройка OPTICS

- OPTICS применяется для кластеризации транзакций и выделения профилей нормального поведения.
- Для подбора оптимальных параметров исследовались комбинации MinPts и ϵ .
- Наилучшие результаты получены при MinPts = 10, ϵ = 0.
- Обеспечили минимальное число ложных тревог.

Performance of OPTICS with Different Parameter Values.

Performance Metrics (in %)	Parameter Values					
	Minpts = 10	Minpts = 10	Minpts = 50	Minpts = 50	Minpts = 100	Minpts = 100
	$\epsilon = 0$	$\epsilon = 5$	$\epsilon = 0$	$\epsilon = 5$	$\epsilon = 0$	$\epsilon = 5$
Accuracy	56.34	56.24	58.37	58.27	62.58	62.48
Precision	69.90	69.80	68.40	68.30	66.43	66.33
F1_Score	50.73	50.60	53.63	53.53	62.30	62.20
TPR	58.79	58.69	60.12	60.02	64.30	64.20
FPR	34.73	35.09	34.94	34.84	37.11	38.46

Сравнение классификаторов и ансамблей

- Проверялись пять алгоритмов: Naïve Bayes, Decision Tree, Rule Induction, k-NN, RBFN.
- При индивидуальном применении модели показывали средние результаты.
- Использование ансамблевых методов (Bagging, Boosting) повысило эффективность:
 - увеличилась точность классификации;
 - снизился уровень ложных тревог.
- Наиболее стабильные показатели зафиксированы у k-NN и RBFN.

Stacking

- Реализована архитектура stacking.
- Этот подход обеспечил наилучшее качество среди всех методов.
- Stacking показал преимущество в точности и устойчивости по сравнению с одиночными классификаторами и другими ансамблями.

Сравнение с существующими подходами

- Проведено сравнение с ANFIS-DIDS и PCA-WRF.
- Новый метод показал преимущество по всем ключевым критериям.
- Дополнительное преимущество — меньшие вычислительные затраты за счёт фильтрации транзакций на этапе OPTICS.

Заключение

- Разработанная система объединяет кластеризацию OPTICS и ансамблевые методы обучения.
- Достоинства подхода:
 - высокая точность обнаружения атак,
 - снижение числа ложных срабатываний,
 - меньшая вычислительная сложность по сравнению с существующими решениями.
- Предложенный метод можно рассматривать как эффективный и перспективный инструмент для построения систем обнаружения вторжений в базы данных.