



OPPORTUNITY
CUP

Разработка алгоритма для определения мошеннических банковских операций

ЛИНГВИСТОЧКИ



Резюме команды

Лобанова Алина

НИУ ВШЭ, 3 курс

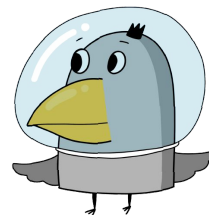
ОП “Фундаментальная и компьютерная лингвистика”



Ревак Ксения

НИУ ВШЭ, 3 курс

ОП “Фундаментальная и компьютерная лингвистика”



Смирнова Варвара

НИУ ВШЭ, 3 курс

ОП “Фундаментальная и компьютерная лингвистика”



Ткач Анна

НИУ ВШЭ, 3 курс

ОП “Фундаментальная и компьютерная лингвистика”



Для всех первый опыт настоящего кейс-чемпионата, участвовали только в хакатоне в рамках курса по программированию.

Этапы решения кейса

Этап 1

Погружение в тему

Обзор данных и литературы по теме, выбор технического средства работы с данными

Этап 2

ETL-1

Обработка данных и применение unsupervised machine learning для промежуточных результатов

Этап 3

ETL-2

Совмещение результатов с этапов 1 и 2
=
Выделение паттернов

Этап 4

Написание алгоритма

Перекладывание теоретически выведенных паттернов на данные

Этап 5

Итоги

Оценка результатов, создание презентации

Этап 1

1. Данные

Для работы с данными выбирается библиотека pandas языка Python, как удобный инструмент анализа данных (в частности формата json).

По датафрейму транзакций строятся графики и сводные таблицы для проверки и понимания данных, в частности выясняется, что:

- Аккаунты и карты соотносятся 1 к 1;
- не бывает такого, что 1 клиент - 2 человека;
- а также открываются разные интересные факты, как сменивший фамилию человек или сработавшая после окончания срока действия аккаунта карта;
- и другое.

2. Литература по теме

Об обнаружении мошенничестве в целом:

Bolton, R. J., & Hand, D. J. (2002).
Statistical Fraud Detection: A Review.
Statistical Science, 17(3), 235–249.
Transaction aggregation as a strategy for
credit card fraud detection

Об автоматическом анализе:

Clark and Niblett, 1989
Quinlan, 1990
Cohen, 1995
Breiman, Friedman, Olshen and Stone,
1984
Quinlan, 1993

О паттернах, основывающихся на клиентах:

Hand and Blunt (2001)

Об обучении без учителя

Fawcett and Provost (1997)
Bolton and Hand (2001)

Этап 2

ETL-1

Аномальные явления умеют выделять машинным обучением, из прочитанных статей узнаем, что аномальность отдельной транзакции выделить тяжело => определяем подозрительность клиента.

Extract

Json данные транзакций переводятся в формат pandas датафрейма.

Дополнительно вычисляется и добавляется столбец кол-во аккаунтов у пользователя

Transform

Из изначальных данных убираются смыслово дублирующие (e.g. клиент - вместе имен, дат рождений и др)

Выделяются числовые значимые данные: кол-во удачных/неудачных транзакций по типам и др

Load-1

Выделяется датафрейм из числовых данных

Модель машинного обучения без учителя

Проводится обучение по выделенным данным

Что использовалось? Модуль anomaly библиотеки PyCaret.

На выходе получаются бинарная оценка не/аномальности на каждого клиента + уверенность модели в решении

Этап 3

ETL-2

Объединение данных для
выделения итоговых паттернов

Extract

Изначальные данные транзакций объединяются по клиентам с результатами определения степени “подозрительности” клиентов



Transform

На основании полученных представлений об области задачи и новых данных проводится проверка и объединение идей в цельные паттерны



Load

Результатом являются 4 вербализованных повторяющихся паттерна потенциально мошеннических операций

Выделенные паттерны

По результатам кумуляции профессиональной информации из прочитанных статей о классификациях fraud-a, анализа выделенных аномальных клиентов и общечеловеческой логики выделяются 4 паттерна подозрительных транзакций:

- ◆ транзакции клиентов, у которых более 5 транзакций за сутки;
- ◆ транзакции больших сумм в ночное время (суммы более 30 тыс. Рублей в период между 10 pm. и 6 am.);
- ◆ транзакции клиентов, у которых более 5 аккаунтов;
- ◆ транзакции, сделанные одним клиентом в разных городах в пределах суток.

Стоит отметить,

что степень уверенности во “fraud-овости” транзакций по некоторым из паттернов невысока, но из соображений цены ошибки исследуемой задачи приоритет отдаётся полноте, а не точности решения.

Заключительный этап

Прописываются алгоритмы по выделению реальных транзакций, соответствующих найденным паттернам

+

Описывается процесс решения в презентации

Перспективы улучшения:

- Выделение “более уверенных” паттернов ошибки на основе комбинации выделенных “простых” паттернов
- Исследование иных подходов и библиотек по выделению аномальных явлений в данных
- Выведение большего количество потенциально значимых числовых параметров