

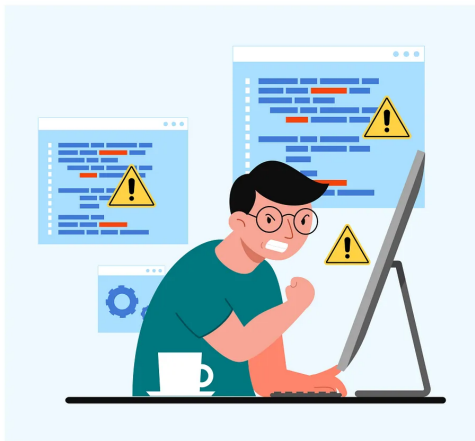
Методы обнаружения мошеннических операций

Егор Горюнов
yegor_goryunov@vk.com

МТУСИ

13.12.2023

Проблема мошенничества в финансовой сфере



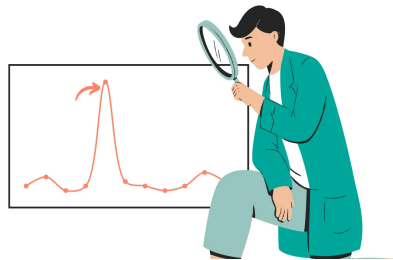
Подходы к детектированию мошеннических транзакций

1. Традиционная ручная проверка
2. Правила и эвристические методы
3. Методы машинного обучения



Машинное обучение для обнаружения мошенничества

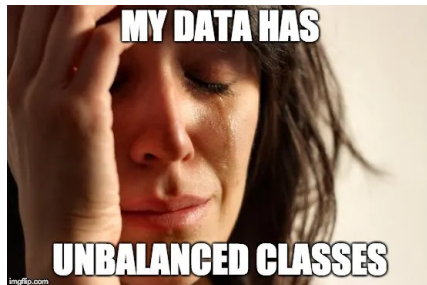
1. Supervised подход: решаем задачу классификации
 - 1.1 NaiveBayes
 - 1.2 Logistic Regression
 - 1.3 SVM
 - 1.4 Gradient Boosting
 - 1.5 Multilayer perceptron (MLP)
2. Unsupervised подход: решаем задачу обнаружения аномалий
 - 2.1 Метрические методы (PCA, LOF, etc.)
 - 2.2 Кластеризация (DBSCAN, GMM, etc.)
 - 2.3 OneClassSVM
 - 2.4 Isolation Forest
 - 2.5 Autoencoder
3. Гибридный подход



Обучение с учителем для Fraud Detection



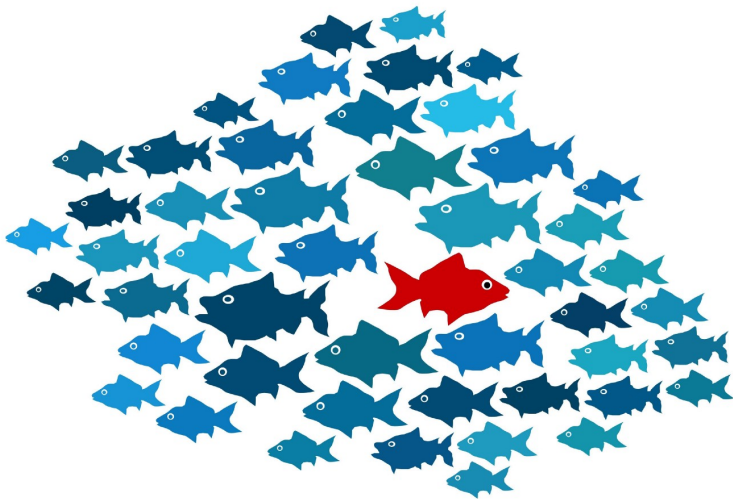
Проблема: экстремально редкий класс



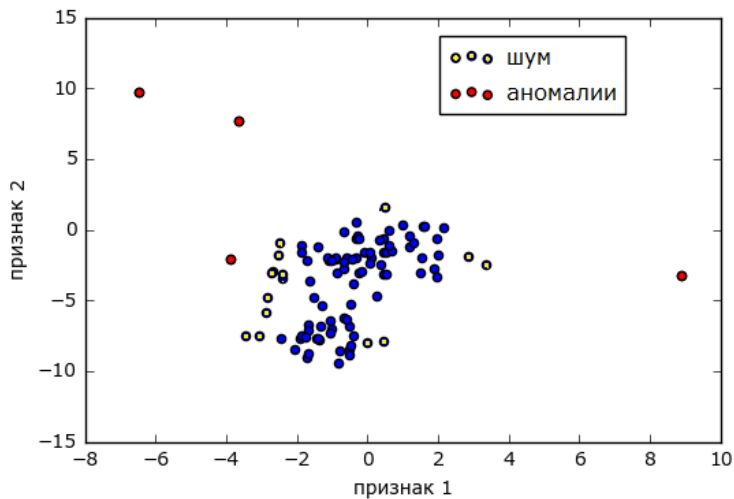
Сравнение различных алгоритмов

Метод	Точность	Время тренировки, с	Время предсказания, с	Количество ложноположительных результатов	Количество ложноотрицательных результатов
Наивный Байес	0.692	0.1479	0.0749	566	44
k-ближайших соседей	0.349	4.265	5.2474	17	95
Логистическая регрессия	0.881	1.2275	0.006	3730	11
Метод опорных векторов	0.9185	3.8093	0.008	1694	9
Случайный лес	0.8081	11.6932	0.128	11	28
XGBoost – CPU	0.8149	0.8328	0.1971	17	27
XGBoost – GPU	0.8081	0.5453	0.1899	14	28
LightGBM	0.8012	0.8809	0.2443	13	29
CatBoost – CPU	0.9041	3.8849	0.1366	589	13
CatBoost – GPU	0.913	1.754	0.1099	413	12

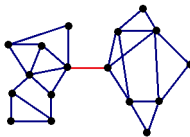
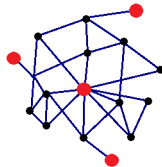
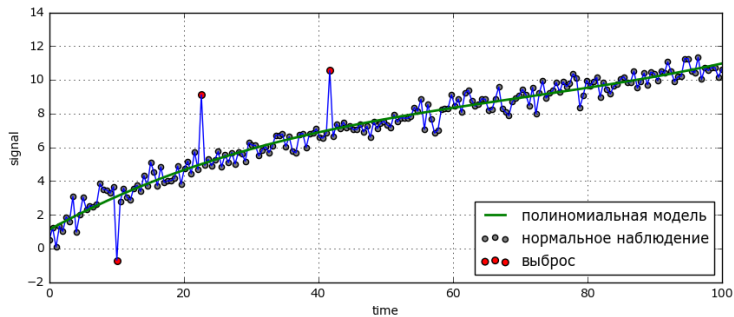
Unsupervised подход: поиск аномалий



Примеры аномалий

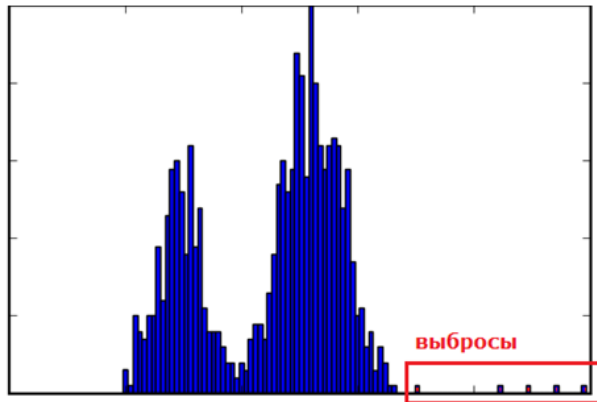


Примеры аномалий

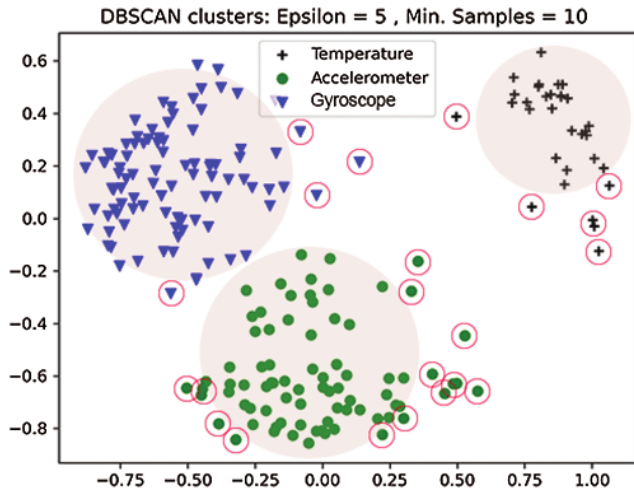


AAABVCCAABBBVCAAABVCAVBBCCABAAABVCCAABBBVCAVBBCC

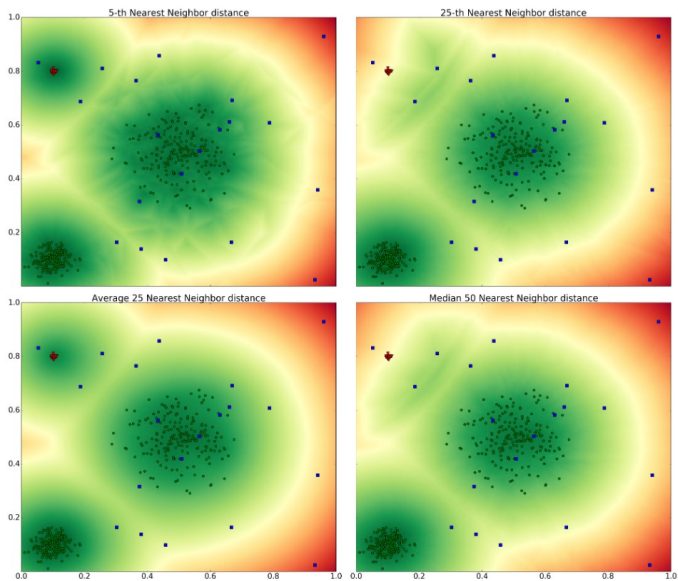
Статистический подход



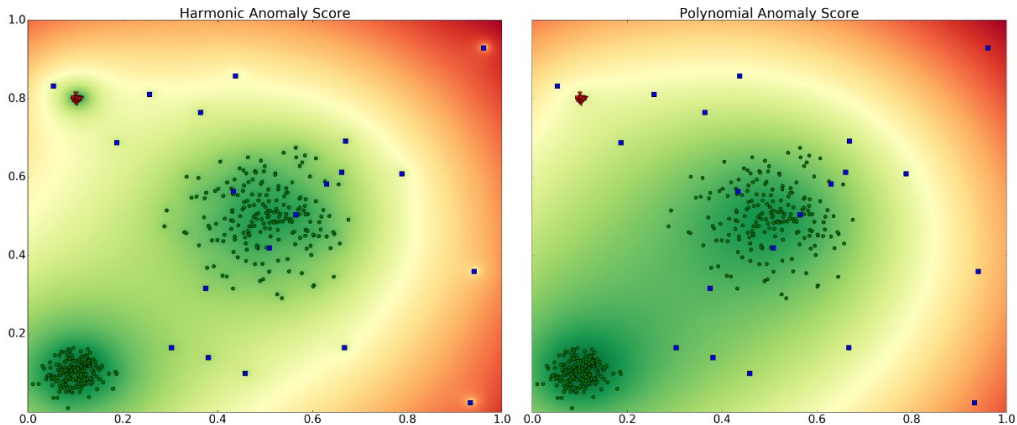
Кластеризация



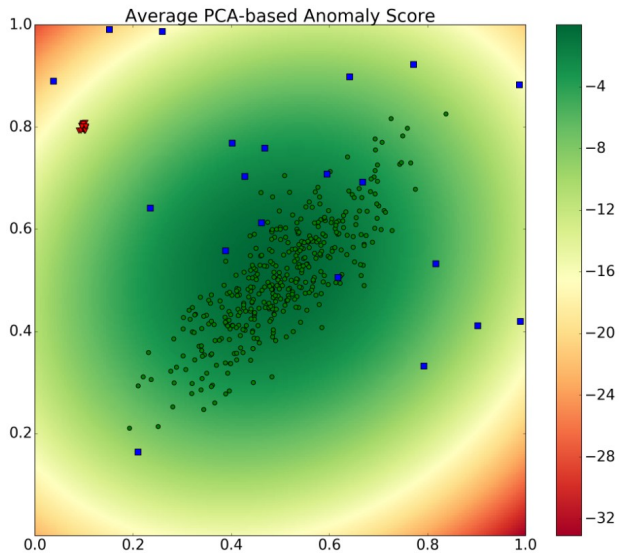
Метрические методы



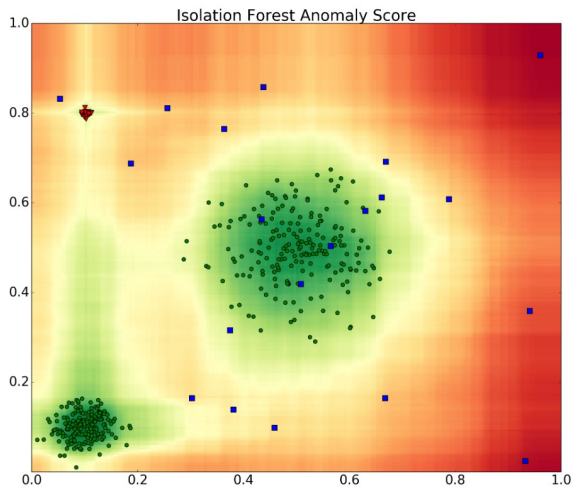
Метрические методы



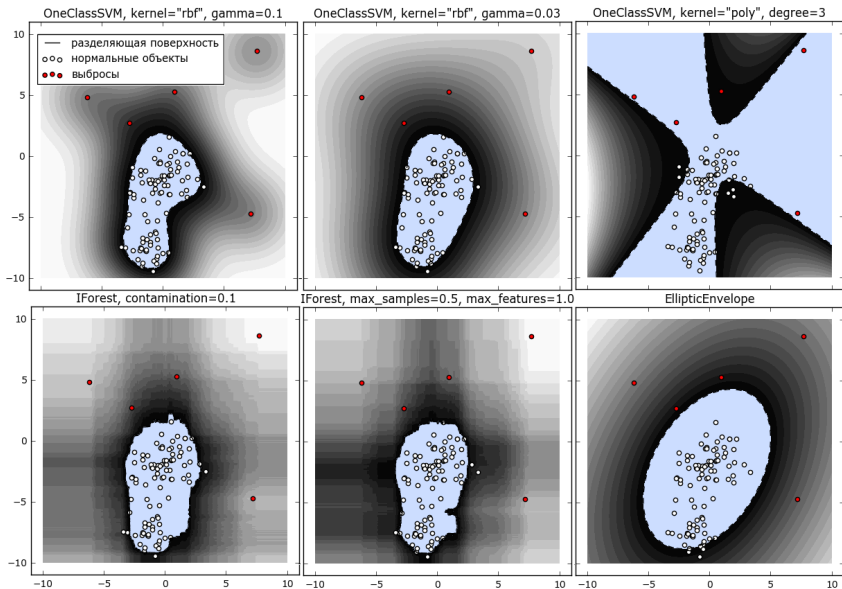
PCA



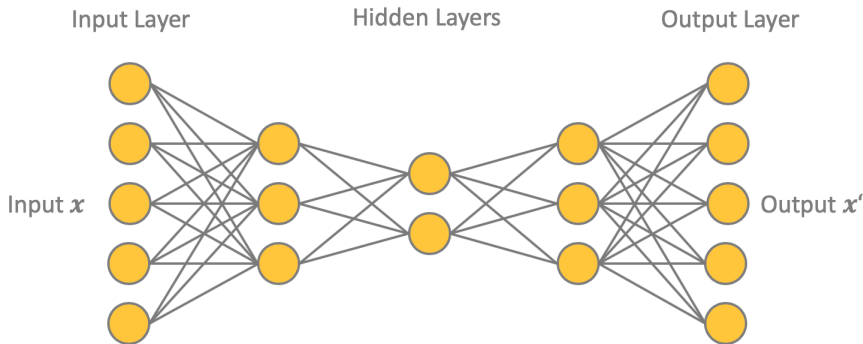
Изолирующий лес



OneClassSVM



Автокодировщики: обучаем восстанавливать неаномальные данные



Execution of the Network:

$$\sqrt{(x_{new} - x'_{new})^2} > \delta \Rightarrow \text{anomaly}$$

Датасет



MACHINE LEARNING GROUP - ULB · UPDATED 6 YEARS AGO



10885

New Notebook



Download (89 MB)



Credit Card Fraud Detection

Anonymized credit card transactions labeled as fraudulent or genuine

[Data Card](#)[Code \(4497\)](#)[Discussion \(103\)](#)

About Dataset

Context

It is important that credit card companies are able to recognize fraudulent credit card transactions so that customers are not charged for items that they did not purchase.

Content

The dataset contains transactions made by credit cards in September 2013 by European cardholders.

This dataset presents transactions that occurred in two days, where we have 492 frauds out of 284,807 transactions. The dataset is highly unbalanced, the positive class (frauds) account for 0.172% of all transactions.

It contains only numerical input variables which are the result of a PCA transformation. Unfortunately, due to confidentiality issues, we cannot provide the original features and more background information about the data. Features V1, V2, ... V28 are the principal components obtained with PCA, the only features which have not been transformed with PCA are 'Time' and 'Amount'. Feature 'Time' contains the seconds elapsed between each transaction and the first transaction in the dataset. The feature 'Amount' is the transaction Amount, this feature can be used for example-dependant cost-sensitive learning. Feature 'Class' is the response variable and it takes value 1 in case of fraud and 0 otherwise.

Given the class imbalance ratio, we recommend measuring the accuracy using the Area Under the Precision-Recall Curve (AUPRC). Confusion matrix accuracy is not meaningful for unbalanced classification.

Usability ⓘ

8.53

License

[Database: Open Database, Cont...](#)

Expected update frequency

Not specified

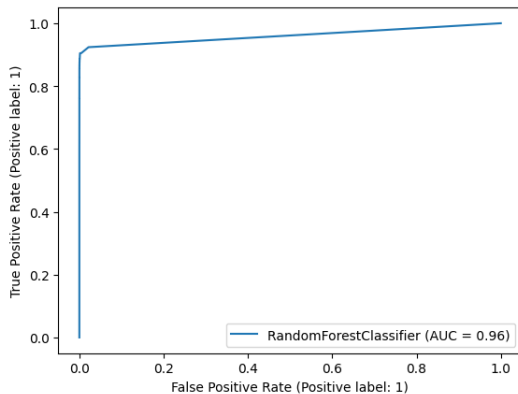
Tags

Finance

Crime

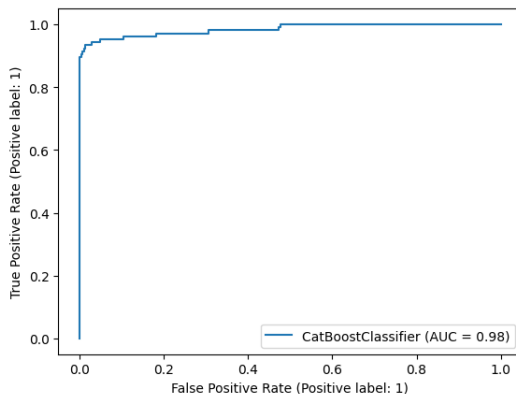
Random Forest

	precision	recall	f1-score	support
0	1.00	1.00	1.00	56856
1	0.52	0.90	0.66	105
accuracy			1.00	56961
macro avg	0.76	0.95	0.83	56961
weighted avg	1.00	1.00	1.00	56961



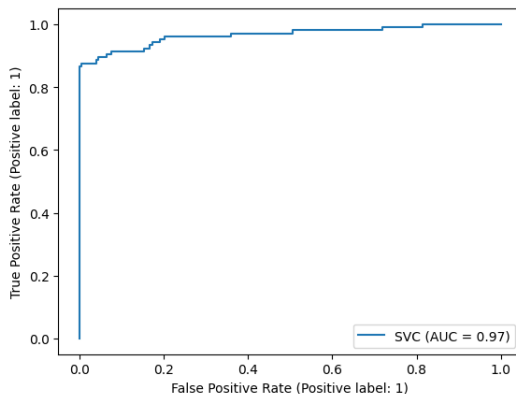
CatBoost

	precision	recall	f1-score	support
0	1.00	1.00	1.00	56856
1	0.54	0.90	0.68	105
accuracy			1.00	56961
macro avg	0.77	0.95	0.84	56961
weighted avg	1.00	1.00	1.00	56961



RBF SVM

	precision	recall	f1-score	support
0	1.00	1.00	1.00	56856
1	0.90	0.85	0.87	105
accuracy			1.00	56961
macro avg	0.95	0.92	0.94	56961
weighted avg	1.00	1.00	1.00	56961



Ссылки



Обсудим?