

# Projekt ZUM - dokumentacja wstępna

Michał Bartnicki, Michał Piotrak

Temat zadania: **Analityczny - detekcja anomalii (kod - DAz11)**

Wybrany zbiór danych: [Credit Card Fraud Detection \(Kaggle\)](#)

## Interpretacja tematu projektu

Realizowany projekt będzie miał charakter szczegółowego przeglądu rozwiązań problemu detekcji anomalii. Jest to obszar analizy danych, którego celem jest identyfikowanie obserwacji odbiegających znacznie od oczekiwanego wzorca. Takie podejście jest wymagane w m.in. obszarze ujawniania wszelkiego rodzaju oszustw, w tym przypadku dotyczących transakcji dokonanych za pomocą kart kredytowych.

Do rozwiązania problemu detekcji anomalii można wykorzystać różne metody, które możemy podzielić na dwie kategorie:

- metody nadzorowane - wykorzystują dane etykietowane, gdzie możemy wyróżnić podział na przykłady standardowe oraz anomalie.
- metody nienadzorowane - identyfikują anomalie na podstawie samej struktury zbioru danych, odległości między obserwacjami lub gęstości rozkładu danych. Te metody są użyteczne, gdy nie posiadamy pełnej wiedzy o występujących anomaliiach, czyli w zbiorach bez etykiet.

Przeprowadzone badania w ramach projektu będą dotyczyły zarówno podejścia nadzorowanego, jak i nienadzorowanego. Dogłębna analiza tematu będzie wymagała również poruszenia aspektu niezrównoważenia klas, który wpływa na skuteczność klasycznych metod klasyfikacji. Poza tym, trzeba będzie przeprowadzić rozważania dotyczące wyboru miar oceny jakości modeli do przypadku silnie niezbalansowanych klas.

## Lista algorytmów, które będą wykorzystane w eksperymentach

Metody nadzorowane:

Metoda	Klasa	Biblioteka
SVM z wagami klas	<a href="#">sklearn.svm.SVC</a> * * w razie problemów z wydajnością i w zależności od rozkładu anomalii - <a href="#">sklearn.svm.LinearSVC</a>	scikit-learn
Drzewo decyzyjne z wagami klas	<a href="#">sklearn.tree.DecisionTreeClassifier</a>	scikit-learn
Las losowy	<a href="#">sklearn.ensemble.RandomForestClassifier</a>	scikit-learn
XGBoost	<a href="#">xgboost.XGBClassifier</a>	XGBoost

Metody nienadzorowane:

Metoda	Klasa	Biblioteka
OC - SVM	<a href="#">sklearn.svm.OneClassSVM</a>	scikit-learn
Las izolacyjny	<a href="#">sklearn.ensemble.IsolationForest</a>	scikit-learn
Lokalny czynnik odstający	<a href="#">sklearn.neighbors.LocalOutlierFactor</a>	scikit-learn
DBSCAN	<a href="#">sklearn.cluster.DBSCAN</a>	scikit-learn

Pozostałe algorytmy, które planujemy zastosować z krótkim opisem celu:

Metoda	Klasa	Biblioteka	Cel
Oversampling (SMOTE)	<a href="#">imblearn.over_sampling.SMOTE</a>	Imbalanced-learn	zbalansowanie klas

Undersampling	<a href="#">imblearn.under_sampling.RandomUnderSampler</a>	Imbalanced-learn	zbalansowanie klas
Walidacja krzyżowa	<a href="#">sklearn.model_selection.StratifiedKFold</a>	scikit-learn	procedury walidacyjne

## Plan badań

### Cel eksperymentów

Nasza praca ma na celu porównanie skuteczności różnych algorytmów klasyfikacji w przypadku dwuklasowym, działając na silnie niezbalansowanym zbiorze danych. Ponadto dla powyższych rozwiązań zweryfikujemy pozytywny wpływ metod zapewniających zbalansowanie klas, gdzie w szczególności znaczenie to powinno mieć dla podejścia nadzorowanego. Chcemy także potwierdzić tezę, że algorytmy nadzorowane gwarantują wyższą skuteczność detekcji znanych wzorców anomalii niż algorytmy nienadzorowane, które z kolei powinny być lepsze w przypadkach gdy charakter anomalii jest zmienny.

### Charakterystyka wybranego zbioru danych

Wybrany zbiór zawiera zanonimizowane informacje o transakcjach dokonanych za pomocą kart kredytowych przez Europejczyków we wrześniu 2013 roku przez okres dwóch dni. Liczba transakcji: 284807, w tym 492 oznaczone jako fałszywe.

Transakcje są opisane przy pomocy następujących cech:

Cecha	Opis
Wektor cech [V1, ..., V28]	Cechy transakcji po transformacji PCA
Time	Upływ czasu w sekundach pomiędzy pierwszą a określoną transakcją (bez transformacji PCA)
Amount	Kwota transakcji (bez transformacji PCA)
Class	Klasa transakcji, wartość "1" - transakcja fałszywa (oszustwo), wartość "0" - transakcja uczciwa

Opisywany zbiór jest wysoce niezbalansowany - przykłady określające transakcje fałszywe (które uznajemy za egzemplarze klasy reprezentującej przypadki nieprawidłowe) stanowią jedynie 0,172% wszystkich transakcji.

Z tego powodu modele uzyskane poprzez użycie metod nadzorowanych są narażone na zjawisko nadmiernego dopasowania. Żeby temu zapobiec użyjemy odpowiednich technik balansujących klasy, do których zaliczamy:

- undersampling - ograniczenie liczby przykładów klasy większościowej,
- oversampling - zwiększenie liczby przykładów klasy mniejszościowej, które możemy uzyskać poprzez zastosowanie algorytmu **SMOTE**, który generuje nowe próbki pomiędzy istniejącymi obserwacjami na podstawie ich lokalnej gęstości.

Poza tym, w celu przygotowania danych planujemy zrealizować następujące czynności:

- sprawdzić, czy mamy w następującym zbiorze jakiegolwiek braki danych i duplikaty,
- dokonać skalowania cech *Time* oraz *Amount* (w wyniku zastosowania transformacji PCA, wektor cech [V1, ..., v28] powinien być uprzednio odpowiednio wyskalowany),
- pomimo zastosowania już transformacji PCA na części parametrów, zbadamy jeszcze korelacje pomiędzy nimi, w celu potencjalnej redukcji wymiarowości.

### Parametry algorytmów, których wpływ będziemy badać

Dla algorytmów nadzorowanych zbadamy wpływ przekształceń typu undersampling i oversampling. Przetestujemy także różne wartości parametru *sampling\_strategy*, czyli oczekiwany stosunek ilości transakcji klasy prawdziwej i fałszywej uzyskany po zastosowaniu resamplingu.

Ponadto, będziemy eksperymentować z różnymi parametrami samych algorytmów, wybór poniższych parametrów może ulec zmianie, wraz z postępem pracy i lepszym poznaniem oraz zrozumieniem działania badanych algorytmów.

Metoda	Parametry
SVM z wagami klas	<i>kernel</i> <ul style="list-style-type: none"> <li><i>linear</i>: może skutkować szybszym treningiem, ponieważ opiera się na liniowym rozdzieleniu klas</li> <li><i>RBF</i></li> </ul> <i>C, class_weight (None vs balanced)</i>
Drzewo decyzyjne z wagami klas	<i>max_depth</i> (ustawić, aby zapobiec przerośnięciu drzewa), <i>class_weight (None vs balanced)</i>
Las losowy	<i>n_estimators</i> (ilość drzew, poszukujemy największej wartości z rozsądnym czasem treningu), <i>max_depth</i> , <i>class_weight</i>
XGBoost	<i>n_estimators</i> , <i>scale_pos_weight</i> , <i>learning_rate</i> , <i>max_depth</i> , <i>early_stopping_rounds</i> (może zapobiec przetrenowaniu)
OC - SVM	<i>kernel</i> , <i>nu</i> (większa wartość może skutkować większą ilością fałszywych alarmów, ale przetestujemy wartości inne niż domyślna 0.5), <i>gamma</i>
Las izolacyjny	<i>n_estimators</i> , <i>contamination</i>
Lokalny czynnik odstający	<i>n_neighbours</i> , <i>contamination</i>
DBSCAN	<i>eps</i> , <i>min_samples</i>

## Planowane do wykorzystania procedury oceny modeli oraz miary jakości

Podstawą procedury oceny jakości modeli będzie k - krotna walidacja krzyżowa. Ze względu na silnie niezbalansowany zbiór danych zastosujemy technikę **Stratified k - fold**, która zapewnia, że proporcje klas w zbiorze treningowym i testowym są podobne do tych w całym zbiorze danych. Ze względu na to, że liczba anomalii jest mała, wartość parametru k nie powinna być duża i według naszej początkowej oceny optymalna jej wartość powinna wynosić co najwyżej 5.

**Dokładność** nie jest odpowiednią miarą jakości dla detekcji anomalii, ponieważ jej wysoka wartość może zostać osiągnięta poprzez prawidłowe klasyfikowanie przykładów klasy większościowej. Za główną miarę w tym zadaniu uznamy **odzysk**. Powie nam ona jak wiele, spośród wszystkich oszustw, potrafi wykryć model. Jest to ważne, aby jak najwięcej nielegalnych transakcji zostało odkrytych, nawet kosztem większej liczby transakcji fałszywie pozytywnych - które będą wiązać się z mniejszym kosztem - np. dodatkową pracą człowieka obsługującego alert o prawdopodobnym oszustwie, lub dodatkowym zatwierdzeniem transakcji przez użytkownika.

Jako miary pomocnicze zostaną użyte:

- **Precyzja** - określi, ile transakcji oznaczonych jako oszustwo jest faktycznie oszustwem. Ta miara jest mniej ważna w naszym przypadku, ale należy ją wziąć pod uwagę, aby nie obciążać nadmiernie użytkownika / administratora fałszywymi alarmami.
- **Miara F** - powie o tym, jak zrównoważone względem siebie są precyzja i odzysk oraz pozwoli uniknąć skrajności modelu.
- **Krzywa Precision - Recall** - do ustalenia kompromisu między precyzją a odzyskiem.
- **Pole pod krzywą PR (Area Under the PR Curve, w skrócie AUPRC)** - wskaźnik opisujący krzywą PR, ogólnie im wyższa wartość tym lepsza jakość modelu.