

# Zastosowanie zespołów klasyfikatorów w klasyfikacji niezbalansowanych strumieni danych

Jakub Klikowski

Wrocław 01.12.2021

## Streszczenie

Rozprawa koncentruje się na zadaniach klasyfikacji binarnych niezbalansowanych strumieni danych z dryfem koncepcji. Zaproponowano wykorzystanie hybrydowych technik próbkowania z akumulacją danych i modelami klasyfikatorów jednoklasowych. Przeprowadzone badania i uzyskane wyniki jednoznacznie wskazują, że proponowane podejścia mają znaczny potencjał, a wyniki uzyskane opracowanymi metodami często dorównują lub przewyższają znane z literatury algorytmy. Bazując na powyższych założeniach została sformułowana następująca hipoteza badawcza:

*Można zaprojektować metody zespołów klasyfikatorów wykorzystujące próbkowanie danych oraz jednoklasowe klasyfikatory, które potrafią osiągnąć jakość predykcji lepsza niż znane z literatury klasyfikatory niezbalansowanych strumieni danych.*

Powyższa hipoteza została uprawdopodobniona poprzez realizację następujących celów badawczych:

*Zaprojektowanie metody zespołu klasyfikatorów jednoklasowych w celu klasyfikacji niezbalansowanego strumienia danych.*

Cel ten został osiągnięty poprzez zaproponowanie algorytmu klasyfikacji niezbalansowanych strumieni danych w oparciu o zespół klasyfikatorów jednoklasowy - One Class support vector machine classifier Ensemble for Imbalanced data

Stream (*OCEIS*). Metoda osiąga wyniki na zbliżonym poziomie do porównywanych metod, ale warto zauważyć, że najlepiej sprawdza się na strumieniach rzeczywistych, co jest jej istotną zaletą. Kolejną zaletą jest brak tendencji do nadmiernej klasyfikacji obiektów z jednej z klas. Taka stabilność znacząco przyczynia się do poprawy jakości klasyfikacji i uzyskania satysfakcjonujących wyników.

***Rozszerzenie metody zespołu klasyfikatorów jednoklasowych poprzez wprowadzenie ulepszonej ważonej reguły decyzyjnej dla lepszej adaptacji do niezbalansowanych strumieni danych.***

Cel ten został osiągnięty poprzez zaprojektowanie nowatorskiej metody klasyfikacji niezbalansowanych strumieni danych - One Class support vector machine **W**eighted **E**nsemble (*OCWE*). Jest to metoda oparta na zespołach klasyfikatorów jednoklasowych, rozszerzająca pierwotną ideę z podejścia *OCEIS*. Wprowadzone modyfikacje doprowadziły do wzrostu jakości klasyfikacji niezbalansowanych strumieni danych. W porównaniu z innymi metodami referencyjnymi, algorytm *OCWE* uzyskuje podobną jakość, a w pewnych warunkach *OCWE* osiąga lepszą wydajność. Metoda ta bardzo dobrze sprawdza się w zadaniach klasyfikacji silnie niezbalansowanych danych. Wielość parametrów zapewnia dobrą zdolność adaptacji do aktualnie rozwiązywanego problemu.

***Zaproponowanie klasyfikatora dla niezbalansowanych strumieni danych z hybrydową techniką próbkowania i akumulacji danych.***

Cel ten został osiągnięty poprzez opracowanie metody klasyfikacji strumieni niezbalansowanych. **D**eterministic **S**ampling **C**lassifier (*DSC*) wypadło korzystnie w porównaniu z innymi dedykowanymi algorytmami. Ocene zdolności predykcyjnych technik przeprowadzono na podstawie eksperymentów komputerowych. Algorytm wykorzystuje bufor pamięci w celu propagacji wybranych instancji z poprzednich fragmentów strumieni danych. Ponieważ bufor ma stały rozmiar, po jego wypełnieniu niektóre instancje muszą zostać usunięte. W zaproponowanej wersji algorytmu przyjęto, że usuwane są najstarsze przykłady.

***Ulepszenie klasyfikatora z hybrydowym próbkowaniem akumulacji danych do ważonego zespołu z wykorzystaniem baggingu w celu klasyfikacji niezbalansowanego strumienia danych.***

Cel ten został osiągnięty poprzez zaproponowanie nowej metody klasyfikacji niezbalansowanych strumieni danych - **Deterministic Sampling Classifier with weighted Bagging (DSCB)**. *DSCB* wykorzystuje koncepcje algorytmu *DSC* oraz perturbowane zbiory danych. Proponowana metoda *DSCB* to zespół klasyfikatorów, którego główna idea skupia się na tworzeniu nowych modeli poprzez użycie ważonego *baggingu* i danych, które zostały zgromadzone z poprzednich części strumienia danych. Eksperymenty komputerowe potwierdziły, że zaproponowane podejście uzyskuje dobre wyniki w stosunku do jakości klasyfikacji wybranych metod referencyjnych.

***Poprawa klasyfikatora niezbalansowanych strumieni danych z hybrydowym próbkowaniem akumulacji danych do metody zespołowej z detektorem dryfu koncepcji.***

Cel ten osiągnięto dzięki zaproponowaniu nowatorskiej, skutecznej metody zespołu klasyfikatorów w połączeniu ze strukturą przetwarzania wstępnego dla zadania klasyfikacji dryfującego, niezbalansowanego strumienia danych. Główna idea *DSE* opiera się na akumulacji wybranych próbek z poprzednich fragmentów w celu późniejszego wzmocnienia i zrównoważenia danych. Pozwala to na specyficzny *oversampling*, który nie wymaga generowania sztucznych danych. Badania przeprowadzone na szerokiej gamie strumieni danych rzeczywistych i generowanych komputerowo potwierdziły skuteczność proponowanego rozwiązania. Podkreślono jego mocne strony w porównaniu z metodami referencyjnymi oraz jego wysoka odporność na szum etykiet.

***Zaprojektowanie metody zespołowej do klasyfikacji niezbalansowanych strumieni danych, która łączy hybrydowe próbkowanie akumulacji danych i klasyfikatory jednoklasowe.***

Cel ten osiągnięto dzięki opracowaniu nowej metody klasyfikacji niezbalansowanych strumieni danych — *DSE-OC*, która łączy dwa różne podejścia. Jest to połączenie najlepszych cech algorytmów *DSE* i *OCEIS*. W jednej metodzie istnieje jednocześnie adaptacja problemów binarnych do klasyfikacji przy użyciu jednoklasowych modeli SVM z metodą, która akumuluje dane do półsyntetycznego *oversamplingu* danych niezbalansowanych. Z przeprowadzonych badań wynika, że prezentowana metoda osiąga dość dobrą jakość, a przeprowadzone analiza

statystyczna wyników eksperymentów wykazała, że metoda ta jest statystycznie istotnie lepsza metody referencyjne.

Jakub Klichowski