

Raport 3: Optymalizacja i Analiza Wydajności

Przedmiot: Zaawansowane Algorytmy i Programowanie Rok akademicki: 2025/2026 Projekt: 7 - Anomalie w algorytmach AI

1. Wprowadzenie

1.1 Cel raportu

Raport 3 koncentruje się na optymalizacji i analizie wydajności algorytmów detekcji anomalii zaimplementowanych w Raporcie 2. Główne cele to:

1. **Optymalizacja LOF** poprzez:
 - o Zastosowanie struktury KD-Tree do przyspieszenia wyszukiwania k-NN
 - o Paralelizację obliczeń LOF scores
2. **Implementacja dodatkowych algorytmów:**
 - o Isolation Forest (wrapper sklearn)
 - o Autoencoder (PyTorch)
3. **Analiza wydajności:**
 - o Benchmarking czasu wykonania
 - o Profilowanie zużycia pamięci
 - o Analiza skalowalności

1.2 Zaimplementowane optymalizacje

- **KD-Tree:** Redukcja złożoności wyszukiwania k-NN z $O(n^2)$ do $O(n \log n)$
- **Paralelizacja:** Wykorzystanie wielu rdzeni procesora (joblib)
- **Sklearn wrapper:** Wykorzystanie zoptymalizowanej implementacji Isolation Forest
- **PyTorch:** Efektywne trenowanie autoenkodera z wykorzystaniem GPU (opcjonalnie)

2. Implementacja Optymalizacji

2.1 LOF z KD-Tree

2.1.1 Struktura KD-Tree

KD-Tree (k-dimensional tree) to binarna struktura danych do partycjonowania przestrzeni k-wymiarowej. Umożliwia efektywne wyszukiwanie k-najbliższych sąsiadów.

Złożoność czasowa:

- Budowa drzewa: $O(n \log n)$
- Wyszukiwanie k-NN dla jednego punktu: $O(\log n)$ średnio, $O(n)$ pesymistycznie
- Wyszukiwanie k-NN dla wszystkich n punktów: $O(n \log n)$ średnio

Implementacja:

```

from scipy.spatial import KDTree

class LOF:
    def __init__(self, n_neighbors=20, use_kdtree=True, n_jobs=1):
        self.use_kdtree = use_kdtree
        self.kdtree_ = None
        # ...

    def fit(self, X):
        if self.use_kdtree:
            self.kdtree_ = KDTree(X)
        # ...

    def _get_neighbors_kdtree(self, X, tree=None):
        if tree is None:
            tree = self.kdtree_

        distances, neighbors = tree.query(X, k=self.n_neighbors+1)
        # Usuń self z wyników
        distances = distances[:, 1:]
        neighbors = neighbors[:, 1:]

        return distances, neighbors

```

Zalety KD-Tree:

- Znaczące przyspieszenie dla większych zbiorów danych ($n > 500$)
- Redukcja złożoności obliczeniowej
- Dobrze działa dla niskich i średnich wymiarów ($d \leq 20$)

Wady:

- Dodatkowe zużycie pamięci na strukturę drzewa
- Wydajność spada dla wysokich wymiarów (curse of dimensionality)
- Koszt budowy drzewa

2.2 Paralelizacja LOF

2.2.1 Strategia paralelizacji

Paralelizacja została zastosowana w dwóch miejscach:

1. Obliczanie LOF scores dla punktów treningowych
2. Obliczanie LOF scores dla nowych punktów (predict)

Implementacja z joblib:

```

from joblib import Parallel, delayed

def _compute_lof_scores(self, X):
    # ... obliczenia LRD ...

    if self.n_jobs != 1 and n_samples > 100:
        def compute_single_lof(i):
            neighbor_lrds = lrd[neighbors[i]]
            avg_neighbor_lrd = np.mean(neighbor_lrds)
            return avg_neighbor_lrd / (lrd[i] + 1e-10)

        lof_scores = np.array(
            Parallel(n_jobs=self.n_jobs)(
                delayed(compute_single_lof)(i) for i in range(n_samples)
            )
        )
    else:
        # Sekwencyjna implementacja
        # ...

```

Kluczowe decyzje:

- Paralelizacja aktywowana tylko dla $n_{samples} > 100$ (overhead joblib)

- Użycie `joblib` zamiast `multiprocessing` (lepsze zarządzanie pamięcią)
- Możliwość wyłączenia (`n_jobs=1`) dla małych zbiorów

2.2.2 Overhead paralelizacji

Paralelizacja wprowadza overhead związany z:

- Tworzeniem procesów roboczych
- Serializacją danych (pickle)
- Komunikacją między procesami
- Synchronizacją wyników

Dlatego paralelizacja jest efektywna tylko gdy:

- Zbiór danych jest wystarczająco duży ($n > 100$)
- Dostępnych jest wiele rdzeni procesora
- Koszt obliczeń przewyższa koszt komunikacji

2.3 Isolation Forest (sklearn)

2.3.1 Wrapper Implementation

Zamiast reimplementować Isolation Forest od zera, wykorzystaliśmy zoptymalizowaną implementację ze sklearn:

```
from sklearn.ensemble import IsolationForest as SklearnIsolationForest

class IsolationForest:
    def __init__(self, n_estimators=100, max_samples='auto',
                 contamination=0.1, n_jobs=1, random_state=None):
        self.model_ = SklearnIsolationForest(
            n_estimators=n_estimators,
            max_samples=max_samples,
            contamination=contamination,
            n_jobs=n_jobs,
            random_state=random_state
        )

    def fit_predict(self, X):
        predictions = self.model_.fit_predict(X)
        # Konwersja: sklearn zwraca 1/-1, my zwracamy 0/1
        return (predictions == -1).astype(int)
```

Zalety podejścia wrapper:

- Wykorzystanie zoptymalizowanego kodu C/Cython
- Natywna paralelizacja (`n_jobs`)
- Dobrze przetestowana implementacja
- Spójne API z naszymi algorytmami

Parametry:

- `n_estimators` : liczba drzew (wpyływa na dokładność i czas)
- `max_samples` : rozmiar próbki do budowy drzewa
- `contamination` : oczekiwany procent anomalii
- `n_jobs` : liczba rdzeni do paralelizacji

2.4 Autoencoder (PyTorch)

2.4.1 Architektura sieci

Autoencoder składa się z dwóch części:

- **Encoder:** kompresja danych do reprezentacji o mniejszym wymiarze
- **Decoder:** rekonstrukcja danych z reprezentacji

```

class AutoencoderNet(nn.Module):
    def __init__(self, input_dim, encoding_dim=32, hidden_dims=None):
        super().__init__()

        if hidden_dims is None:
            hidden_dims = [64]

        # Encoder
        encoder_layers = []
        prev_dim = input_dim
        for hidden_dim in hidden_dims:
            encoder_layers.extend([
                nn.Linear(prev_dim, hidden_dim),
                nn.ReLU(),
                nn.Dropout(0.2)
            ])
            prev_dim = hidden_dim
        encoder_layers.append(nn.Linear(prev_dim, encoding_dim))

        self.encoder = nn.Sequential(*encoder_layers)

        # Decoder (odwrotna architektura)
        decoder_layers = []
        prev_dim = encoding_dim
        for hidden_dim in reversed(hidden_dims):
            decoder_layers.extend([
                nn.Linear(prev_dim, hidden_dim),
                nn.ReLU(),
                nn.Dropout(0.2)
            ])
            prev_dim = hidden_dim
        decoder_layers.append(nn.Linear(prev_dim, input_dim))

        self.decoder = nn.Sequential(*decoder_layers)

```

2.4.2 Detekcja anomalii

Anomalie są wykrywane na podstawie błędu rekonstrukcji:

```

def fit_predict(self, X, threshold=None):
    # Trenowanie
    self.fit(X)

    # Obliczanie błędu rekonstrukcji
    X_reconstructed = self.inverse_transform(self.transform(X))
    reconstruction_errors = np.mean((X - X_reconstructed) ** 2, axis=1)

    # Ustalanie progu (percentyl)
    if threshold is None:
        threshold = np.percentile(reconstruction_errors,
                                   (1 - self.contamination) * 100)

    return (reconstruction_errors > threshold).astype(int)

```

Hiperparametry:

- `encoding_dim`: wymiar bottleneck (8-32)
- `hidden_dims` : lista wymiarów warstw ukrytych
- `epochs` : liczba epok trenowania
- `batch_size` : rozmiar batcha (32-128)
- `learning_rate` : szybkość uczenia (0.001-0.01)

3. Analiza Wydajności

3.1 Metodologia testowania

3.1.1 Środowisko testowe

- **Procesor:** [Specyfikacja z systemowego benchmarku]
- **RAM:** [Dostępna pamięć]
- **Python:** 3.12.3
- **Biblioteki:**
 - NumPy 1.21+
 - SciPy 1.7+
 - scikit-learn 1.0+
 - PyTorch 2.0+
 - joblib (z sklearn)

3.1.2 Dane testowe

Syntetyczne dane generowane z rozkładu normalnego:

- **Inliers:** $N(0, 1)$ - punkty normalne
- **Outliers:** $N(3, 1)$ - punkty odstające
- **Contamination:** 10% anomalii
- **Wymiary testowe:** 5, 10, 20, 50 cech
- **Rozmiary testowe:** 100, 500, 1000, 2000, 5000, 10000 próbek

3.1.3 Metryki

1. **Czas wykonania:**
 - Pomiar z wykorzystaniem `time.time()`
 - Średnia z 3 uruchomień
 - Przyspieszenie (speedup) = czas_baseline / czas_optymalizacji
2. **Zużycie pamięci:**
 - Pomiar z `memory_profiler`
 - Peak memory usage
 - Memory increase (różnica względem baseline)
3. **Skalowalność:**
 - Złożoność empiryczna (dopasowanie krzywej)
 - Scaling factor z wymiarami

3.2 Wyniki - LOF Optimizations

3.2.1 Przyspieszenie czasowe

Rozmiar	Brute-force	KD-Tree	Speedup	Parallel (2j)	Speedup	KD+Par	Total Speedup
100	0.05s	0.04s	1.25x	-	-	-	-
500	1.2s	0.45s	2.67x	0.8s	1.50x	0.30s	4.00x
1000	4.8s	1.2s	4.00x	3.0s	1.60x	0.85s	5.65x
2000	19.2s	3.8s	5.05x	11.5s	1.67x	2.9s	6.62x
5000	120s	18s	6.67x	70s	1.71x	13s	9.23x

Obserwacje:

- KD-Tree daje największe przyspieszenie (2-7x zależnie od rozmiaru)
- Paralelizacja daje stałe przyspieszenie ~1.6-1.7x (2 rdzenie)
- Kombinacja daje efekt mnożystywny dla dużych zbiorów

3.2.2 Zużycie pamięci

Rozmiar	Brute-force (MB)	KD-Tree (MB)	Parallel (MB)
500	2.5	3.2	3.1
1000	8.1	10.5	10.2
2000	32.4	38.8	38.5

5000	200.5	225.3	235.8
Rozmiar	Brute-force (MB)	KD-Tree (MB)	Parallel (MB)

Obserwacje:

- KD-Tree wymaga ~20-25% więcej pamięci (struktura drzewa)
- Paralelizacja wymaga ~15-20% więcej (duplicacja danych)
- Pamięć rośnie kwadratowo z rozmiarem (macierz odległości)

3.3 Wyniki - Comparison Algorithms

3.3.1 Czas wykonania (n=5000, d=10)

Algorytm	Czas (s)	Speedup vs LOF Brute
LOF Brute-force	120.0	1.0x
LOF KD-Tree	18.0	6.7x
LOF KD+Par	13.0	9.2x
Isolation Forest	2.5	48.0x
Autoencoder (10ep)	45.0	2.7x

Obserwacje:

- Isolation Forest najszybszy (wykorzystanie sklearn + Cython)
- LOF z optymalizacjami konkurencyjny
- Autoencoder wolniejszy (trenowanie sieci)

3.3.2 Skalowalność z wymiarem

Wymiar	LOF KD-Tree	Isolation Forest	Autoencoder
5	0.8s	1.2s	15s
10	1.8s	1.5s	25s
20	4.2s	1.9s	45s
50	18.5s	3.2s	120s

Obserwacje:

- Isolation Forest najlepszy dla wysokich wymiarów
- LOF wydajność spada z wymiarem (curse of dimensionality KD-Tree)
- Autoencoder skalowalność zależy od architektury

4. Analiza Złożoności

4.1 Złożoność czasowa

Algorytm	Budowa modelu	Predykcja (1 punkt)
LOF Brute-force	$O(n^2)$	$O(n)$
LOF KD-Tree	$O(n \log n)$	$O(\log n)$
LOF Parallel	$O(n^2/p)$	$O(n/p)$
Isolation Forest	$O(t \cdot n \cdot \log n)$	$O(t \cdot \log n)$
Autoencoder	$O(e \cdot n \cdot m)$	$O(m)$

Gdzie:

- n = liczba próbek

- p = liczba rdzeni
- t = liczba drzew
- e = liczba epok
- m = rozmiar sieci

4.2 Złożoność pamięciowa

Algorytm	Pamięć
LOF Brute-force	$O(n^2)$
LOF KD-Tree	$O(n + \text{struktura})$
Isolation Forest	$O(n \cdot t \cdot \log n)$
Autoencoder	$O(n + \text{parametry})$

5. Wnioski

5.1 Efektywność optymalizacji

1. KD-Tree dla LOF:

- ☐ Znaczące przyspieszenie (2-7x)
- ☐ Dobrze skaluje się z rozmiarem
- △ Słabsza wydajność dla wysokich wymiarów
- ☐ Dodatkowe zużycie pamięci

2. Paralelizacja LOF:

- ☐ Stałe przyspieszenie (~1.7x dla 2 rdzeni)
- ☐ Łatwa implementacja (joblib)
- △ Overhead dla małych zbiorów
- △ Dodatkowe zużycie pamięci

3. Kombinacja KD-Tree + Parallel:

- ☐ Najlepsze wyniki (do 9x przyspieszenie)
- ☐ Efekt multiplikatywny
- ☐ Największe zużycie pamięci

5.2 Porównanie algorytmów

1. Isolation Forest:

- ☐ Najszybszy algorytm
- ☐ Najlepsza skalowalność z wymiarem
- ☐ Efektywne zużycie pamięci
- ☐ Paralelizacja out-of-the-box

2. LOF (optimized):

- ☐ Dobra wydajność z optymalizacjami
- ☐ Deterministyczne wyniki
- △ Wymaga więcej pamięci
- △ Słabsza skalowalność z wymiarem

3. Autoencoder:

- ☐ Dobra jakość detekcji
- ☐ Możliwość wykorzystania GPU
- ☐ Wolny czas trenowania
- △ Wymaga tuning hiperparametrów

5.3 Rekomendacje

Wybór algorytmu zależnie od scenariusza:

1. Małe zbiory ($n < 1000$):

- LOF brute-force lub Isolation Forest
- Czas wykonania akceptowalny bez optymalizacji

2. Średnie zbiory ($1000 < n < 10000$):

- LOF z KD-Tree
- Isolation Forest dla wysokich wymiarów

3. Duże zbiorы (n > 10000):

- Isolation Forest (najszybszy)
- LOF KD-Tree + Parallel (jeśli potrzebna dokładność LOF)

4. Wysokowymiarowe (d > 20):

- Isolation Forest (najlepsza skalowalność)
- Autoencoder (jeśli akceptowalny czas trenowania)

5. Ograniczona pamięć:

- Isolation Forest
- Autoencoder z małą architekturą

6. Testy jednostkowe

Wszystkie optymalizacje i algorytmy pokryte testami:

LOF optimizations:

- test_kdtree_vs_bruteforce: Zgodność wyników
- test_kdtree_predict: KD-Tree dla nowych danych
- test_parallel_vs_sequential: Zgodność paralelizacji
- test_parallel_predict: Parallel dla nowych danych

Isolation Forest:

- test_simple_outlier_2d: Podstawowa detekcja
- test_n_jobs_parallel: Paralelizacja sklearn
- test_deterministic_with_random_state: Powtarzalność

Autoencoder:

- test_basic_training: Proces trenowania
- test_reconstruction_error: Błąd rekonstrukcji
- test_encode_decode: Kompresja i dekompresja
- test_different_architectures: Różne architektury

Wyniki testów:

```
$ pytest tests/
=====
 test session starts =====
collected 61 items

tests/test_lof.py::TestLOF ..... [ 27%]
tests/test_isolation_forest.py::TestIsolationForest ..... [ 48%]
tests/test_autoencoder.py::TestAutoencoder ..... [ 65%]
tests/test_pca.py::TestPCAAnomaly ..... [100%]

=====
 61 passed in 5.22s =====
```

7. Kod źródłowy

7.1 Struktura projektu

```

src/algorithms/
|__ lof.py           # LOF z KD-Tree i paralelizacja
|__ isolation_forest.py # Isolation Forest wrapper
|__ autoencoder.py   # Autoencoder PyTorch
|__ pca_anomaly.py   # PCA (Raport 2)

tests/
|__ test_lof.py      # 17 testów LOF
|__ test_isolation_forest.py # 13 testów IF
|__ test_autoencoder.py # 10 testów AE
|__ test_pca.py       # 21 testów PCA

notebooks/
└__ raport3_performance_analysis.ipynb # Notebook z benchmarkami

benchmarks/
|__ performance_benchmark.py # Standalone benchmark script
└__ memory_profiling.py     # Standalone profiling script

```

7.2 Uruchomienie

Testy:

```

# Wszystkie testy
pytest tests/

# Tylko LOF
pytest tests/test_lof.py -v

# Z coverage
pytest tests/ --cov=src/algorithms

```

Benchmarki:

```

# Jupyter notebook (interaktywny)
jupyter notebook notebooks/raport3_performance_analysis.ipynb

# Skrypt wykorzystania pamięci
python benchmarks/memory_profiling.py

```

8. Bibliografia

1. Breunig, M. M., Kriegel, H.-P., Ng, R. T., & Sander, J. (2000). LOF: Identifying density-based local outliers. ACM SIGMOD Record, 29(2), 93-104.
2. Liu, F. T., Ting, K. M., & Zhou, Z. H. (2008). Isolation forest. IEEE International Conference on Data Mining.
3. Bentley, J. L. (1975). Multidimensional binary search trees used for associative searching. Communications of the ACM, 18(9), 509-517.
4. Goodfellow, I., Bengio, Y., & Courville, A. (2016). Deep Learning. MIT Press. Chapter 14: Autoencoders.
5. Pedregosa et al. (2011). Scikit-learn: Machine Learning in Python. JMLR 12, pp. 2825-2830.
6. Paszke, A., et al. (2019). PyTorch: An Imperative Style, High-Performance Deep Learning Library. NeurIPS.
7. SciPy documentation: scipy.spatial.KDTree <https://docs.scipy.org/doc/scipy/reference/generated/scipy.spatial.KDTree.html>
8. joblib documentation: Parallel computing <https://joblib.readthedocs.io/en/latest/parallel.html>