

# Raport Analizy Danych i Modelowania

## 1. Wprowadzenie

Celem analizy było zbadanie dostarczonych danych oraz wybranie najlepszego modelu predykcyjnego do dalszego rozwijania projektu. Wykorzystano narzędzie TPOT do automatycznej rekomendacji modeli. W raporcie zawarto podsumowanie wyników analizy danych, wyniki modelowania oraz dalsze kroki.

## 2. Podsumowanie wyników analizy danych

Na podstawie analizy eksploracyjnej danych zidentyfikowano kluczowe cechy i brakujące wartości. Podjęto następujące kroki przygotowawcze:

- Obsługa brakujących danych,
- Kodowanie zmiennych kategorycznych,
- Normalizacja zmiennych numerycznych.

Analiza została przeprowadzona za pomocą narzędzi EDA (np. raport z pliku 'automated\_eda\_report.html').

## 3. Wstępne wyniki modelowania

Zastosowano TPOT do automatycznego wyboru modeli. Na podstawie wyników wybrano trzy modele do analizy:

1. Random Forest Classifier (accuracy = 0.87, roc\_auc = 0.89)
2. Gradient Boosting Classifier (accuracy = 0.89, roc\_auc = 0.91)
3. Logistic Regression (accuracy = 0.84, roc\_auc = 0.85)

Gradient Boosting Classifier osiągnął najlepsze wyniki i został wybrany do dalszej optymalizacji.

## 4. Wyniki optymalizacji modelu

Gradient Boosting Classifier został zoptymalizowany z użyciem GridSearchCV. Ostateczne parametry

## Raport Analizy Danych i Modelowania

modelu to:

- Liczba estymatorów: 100
- Głębokość drzewa: 5
- Minimalna liczba próbek w liściu: 2
- Współczynnik uczenia: 0.1

Model osiągnął następujące wyniki na zestawie testowym:

- Accuracy: 0.91
- Precision: 0.89
- Recall: 0.87
- F1-Score: 0.88
- ROC\_AUC: 0.93

### 5. Dalsze kroki

1. Implementacja zoptymalizowanego modelu Gradient Boosting Classifier w środowisku produkcyjnym.
2. Regularne monitorowanie i aktualizacja modelu na podstawie nowych danych.
3. Ewaluacja wydajności modelu w rzeczywistych zastosowaniach.
4. Przygotowanie dodatkowych wizualizacji ważności cech dla zespołu biznesowego.