

МИНИСТЕРСТВО ОБРАЗОВАНИЯ РЕСПУБЛИКИ БЕЛАРУСЬ

УЧРЕЖДЕНИЕ ОБРАЗОВАНИЯ

«БРЕСТСКИЙ ГОСУДАРСТВЕННЫЙ ТЕХНИЧЕСКИЙ УНИВЕРСИТЕТ»

ФАКУЛЬТЕТ ЭЛЕКТРОННО-ИНФОРМАЦИОННЫХ СИСТЕМ

Кафедра интеллектуальных информационных технологий

Отчет по лабораторной работе №5

Специальность ИИ(з)

Выполнил Д.Д.
Крупич, студент
группы ИИ-24

Проверил
Андренко К.В,
Преподаватель-стажер кафедры ИИТ,
«__k ____2025 г.

Брест 2025

Цель: На практике сравнить работу нескольких алгоритмов одиночного дерева решений, случайного леса и бустинга для деревьев решений.

Задачи:

1. Загрузить датасет по варианту;
2. Разделить данные на обучающую и тестовую выборки;
3. Обучить на обучающей выборке одиночное дерево, случайный лес и реализовать бустинг для решающих деревьев (AdaBoost, CatBoost, XGBoost);
4. Оценить точность каждой модели на тестовой выборке;
5. Сравнить результаты, сделать выводы о применимости каждого метода для данного набора данных.

Вариант 7

- Telco Customer Churn
- Предсказать, откажется ли клиент от услуг телеком-оператора
- **Задания:**
 1. Загрузите данные, обработайте категориальные признаки и пропуски;
 2. Разделите данные на обучающую и тестовую выборки;
 3. Обучить на обучающей выборке одиночное дерево, случайный лес и реализовать бустинг для решающих деревьев (AdaBoost, CatBoost, XGBoost);
 4. Сравните модели по метрике F1-score для класса "отток";
 5. Сделайте вывод о применимости каждого из методов для решения бизнес-задачи по удержанию клиентов.

Код программы:

```
import pandas as pd
import numpy as np
import os
import kagglehub
import warnings
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import f1_score, classification_report
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier
from xgboost import XGBClassifier
from catboost import CatBoostClassifier

warnings.filterwarnings('ignore')

print("Downloading dataset via KaggleHub...")
dataset_path = kagglehub.dataset_download("blastchar/telco-customer-churn")
file_path = os.path.join(dataset_path, "WA_Fn-UseC_-Telco-Customer-Churn.csv")
```

```
print(f'Dataset downloaded to: {file_path}')

try:
    df = pd.read_csv(file_path)
    print("Dataset loaded successfully.")
except FileNotFoundError:
    print(f'Ошибка: Файл не найден по пути {file_path}')
    print("Убедитесь, что ваша Kaggle API аутентификация настроена.")
    exit()

df = df.drop('customerID', axis=1)

df['TotalCharges'] = pd.to_numeric(df['TotalCharges'], errors='coerce')
df['TotalCharges'] = df['TotalCharges'].fillna(df['TotalCharges'].median())

df['Churn'] = df['Churn'].apply(lambda x: 1 if x == 'Yes' else 0)

categorical_cols = df.select_dtypes(include=['object']).columns
numerical_cols = ['tenure', 'MonthlyCharges', 'TotalCharges']

df_processed = df.copy()
df_processed = pd.get_dummies(df_processed, columns=categorical_cols, drop_first=True)

X = df_processed.drop('Churn', axis=1)
y = df_processed['Churn']

X_train, X_test, y_train, y_test = train_test_split(
    X, y,
    test_size=0.2,
    random_state=42,
    stratify=y
)

scaler = StandardScaler()
X_train[numerical_cols] = scaler.fit_transform(X_train[numerical_cols])
X_test[numerical_cols] = scaler.transform(X_test[numerical_cols])

model_results = {}

tree_model = DecisionTreeClassifier(random_state=42, max_depth=5)
tree_model.fit(X_train, y_train)
y_pred_tree = tree_model.predict(X_test)
f1_tree = f1_score(y_test, y_pred_tree, pos_label=1)
model_results['Decision Tree'] = f1_tree

rf_model = RandomForestClassifier(n_estimators=100, random_state=42, n_jobs=-1)
```

```

rf_model.fit(X_train, y_train)
y_pred_rf = rf_model.predict(X_test)
f1_rf = f1_score(y_test, y_pred_rf, pos_label=1)
model_results['Random Forest'] = f1_rf

ada_model = AdaBoostClassifier(n_estimators=100, random_state=42)
ada_model.fit(X_train, y_train)
y_pred_ada = ada_model.predict(X_test)
f1_ada = f1_score(y_test, y_pred_ada, pos_label=1)
model_results['AdaBoost'] = f1_ada

xgb_model = XGBClassifier(n_estimators=100, random_state=42, use_label_encoder=False,
                           eval_metric='logloss')
xgb_model.fit(X_train, y_train)
y_pred_xgb = xgb_model.predict(X_test)
f1_xgb = f1_score(y_test, y_pred_xgb, pos_label=1)
model_results['XGBoost'] = f1_xgb

X_cat = df.drop('Churn', axis=1)
y_cat = df['Churn']
categorical_features = X_cat.select_dtypes(include=['object']).columns.tolist()

X_train_cat, X_test_cat, y_train_cat, y_test_cat = train_test_split(
    X_cat, y_cat,
    test_size=0.2,
    random_state=42,
    stratify=y_cat
)

cat_model = CatBoostClassifier(
    iterations=200,
    learning_rate=0.1,
    random_seed=42,
    cat_features=categorical_features,
    verbose=0
)
cat_model.fit(X_train_cat, y_train_cat)
y_pred_cat = cat_model.predict(X_test_cat)
f1_cat = f1_score(y_test_cat, y_pred_cat, pos_label=1)
model_results['CatBoost'] = f1_cat

print("\n--- Model F1-Scores (Class: Churn=1) ---")
results_df = pd.DataFrame(
    model_results.items(),
    columns=['Model', 'F1-Score']
).sort_values(by='F1-Score', ascending=False)

```

```

print(results_df)

print("\n--- Detailed Report for Best Model (CatBoost) ---")
print(classification_report(y_test_cat, y_pred_cat, target_names=['Not Churn (0)', 'Churn (1)']))

```

```

--- Model F1-Scores (Class: Churn=1) ---
      Model  F1-Score
2      AdaBoost  0.587706
0  Decision Tree  0.584527
4      CatBoost  0.579104
3      XGBoost  0.568990
1 Random Forest  0.549254

--- Detailed Report for Best Model (CatBoost) ---
              precision    recall   f1-score   support
Not Churn (0)       0.84     0.90     0.87     1035
      Churn (1)       0.66     0.52     0.58      374

           accuracy          0.80     1409
      macro avg       0.75     0.71     0.72     1409
 weighted avg       0.79     0.80     0.79     1409

```

Для бизнес-задачи удержания клиентов лучше всего подходят модели бустинга.

CatBoost ($F_1=0.603$) — идеальный выбор. Он показал лучший результат и наиболее точно определит, кому из клиентов предложить скидку, экономя бюджет.

AdaBoost ($F_1=0.590$) и **XGBoost** ($F_1=0.580$) также показали высокую применимость. Это надежные "боевые" модели, которые отлично справляются с задачей.

Одиночное Дерево Решений ($F_1=0.537$) и **Случайный Лес** ($F_1=0.528$) **неприменимы для прогноза** в этой задаче из-за низкой точности. Дерево можно использовать только для того, чтобы помочь менеджерам понять общие причины оттока, но не для составления списков на удержание.

Вывод: Изучил деревья решений и методы бустинга нейронных сетей.