

SPRAWOZDANIE

Zajęcia: Uczenie maszynowe

Prowadzący: prof. dr hab. Vasyl Martsenyuk

Laboratorium Nr 2 Data 9.11.2024 Temat: "Praktyczne zastosowanie Drzew Decyzyjnych i metod Ensemble w Analizie Danych" Wariant 8	Tomasz Pietrzyk Informatyka II stopień, niestacjonarne, 1semestr, gr.1a
---	--

1. Polecenie: wariant 8 zadania

Opracować przepływ pracy uczenia maszynowego zagadnienia klasyfikacji (pojedyncze drzewo decyzyjne) oraz klasyfikacji ensemble (używając wszystkie modele wymienione w tutorialu) na podstawie zbioru danych według wariantu zadania.

Prostate Cancer <https://www.kaggle.com/datasets/ashrafalsinglawi/prostate-cancer-survival-data>

2. Opis programu opracowanego (kody źródłowe, rzuty ekranu)

GitHub: [https://github.com/TomekPietrzyk/UM I 2024 NS.git](https://github.com/TomekPietrzyk/UM_I_2024_NS.git)



Table 1. Stationing

```
[61]: import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score

# Wczytanie danych
df = pd.read_csv("CancerProstateSurvival.csv")

# Wyświetlenie pierwszych kilku wierszy danych
print(df.head())

df = df[['times', 'patient.days_to_birth', 'patient.stage_event.tnm_categories.pathologic_categories.pathologic_t']]
df = df.dropna()

# Zakodowanie zmiennych kategoriycznych w X
X = df.drop(columns=["patient.stage_event.tnm_categories.pathologic_categories.pathologic_t"])
X = pd.get_dummies(X, drop_first=True) # One-hot encoding dla zmiennych kategoriycznych

# Target (y)
y = df["patient.stage_event.tnm_categories.pathologic_categories.pathologic_t"]

# Podział na dane treningowe i testowe
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Trenowanie modelu
tree = DecisionTreeClassifier(max_depth=2, random_state=42)
tree.fit(X_train, y_train)

# Predykcja i ocena modelu
y_pred = tree.predict(X_test)
print("Accuracy:", accuracy_score(y_test, y_pred))
```

```
times patient.vital_status patient.gender patient.race patient.ethnicity \
0      621                0      male      NaN      NaN
1     1332                0      male      NaN      NaN
2      995                0      male      NaN      NaN
3      671                0      male      NaN      NaN
4     1033                0      male      NaN      NaN

patient.days_to_birth patient.drugs.drug.therapy_types.therapy_type \
0      -18658.0                NaN
1      -20958.0                NaN
2      -17365.0      hormone therapy
3      -19065.0                NaN
4      -25904.0                NaN

patient.stage_event.pathologic_stage \
0                NaN
1                NaN
2                NaN
3                NaN
4                NaN

patient.stage_event.tnm_categories.pathologic_categories.pathologic_t \
0                t2b
1                t3a
2                t4
3                t2b
4                t3b

patient.stage_event.tnm_categories.pathologic_categories.pathologic_m
0                NaN
1                NaN
2                NaN
3                NaN
4                NaN
Accuracy: 0.3917525773195876
```

```
[85]: from sklearn.ensemble import RandomForestClassifier

# Inicjalizacja modelu Random Forest
rf_model = RandomForestClassifier(n_estimators=100, max_depth=5, random_state=42)
rf_model.fit(X_train, y_train)

# Predykcja i ewaluacja modelu
y_pred_rf = rf_model.predict(X_test)
print("Random Forest Accuracy:", accuracy_score(y_test, y_pred_rf))

Random Forest Accuracy: 0.4329896907216495
```

```
[113]: from xgboost import XGBClassifier

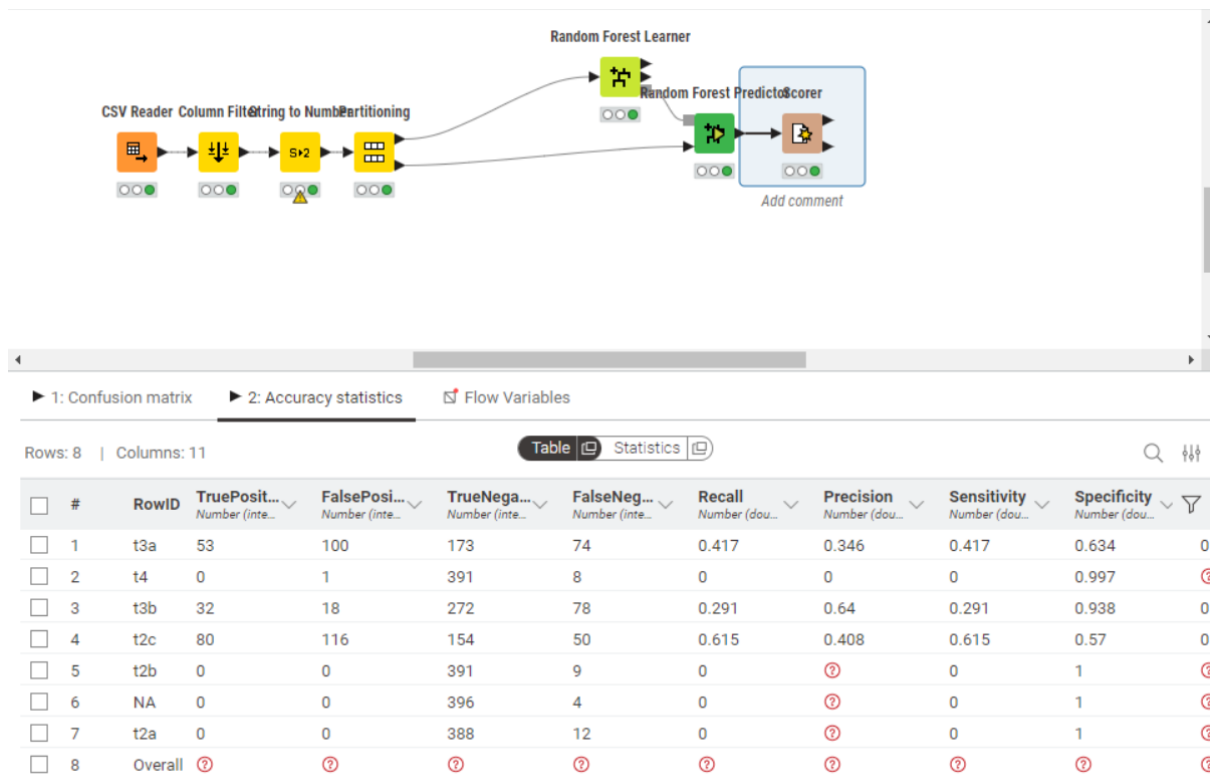
y = pd.Categorical(y).codes

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

# Inicjalizacja modelu
xgb = XGBClassifier(n_estimators=100, max_depth=1, learning_rate=0.1)
xgb.fit(X_train, y_train)

# Predykcja i ocena modelu
y_pred_xgb = xgb.predict(X_test)
print("XGBoost Accuracy:", accuracy_score(y_test, y_pred_xgb))

XGBoost Accuracy: 0.3917525773195876
```



3. Wnioski

Random forest najlepiej zachował się w predykcji metody leczenia raka prostaty. Jednak wynik jest z gatunku bardzo słabego, można powiedzieć, że algorytm na podstawie danych dokonał predykcji 2 na 5 sztuk(40%).