

Системи штучного інтелекту.

Лабораторна робота 5.Федорович Дмитро ІІЗ-21-3

<https://github.com/Dmitrij3/lab5AI>

**Завдання №1.** Створення класифікаторів на основі випадкових та гранично випадкових лісів.

```
import argparse

import numpy as np
import matplotlib.pyplot as plt
from sklearn.metrics import classification_report
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier, ExtraTreesClassifier
from sklearn.metrics import classification_report

from utilities import visualize_classifier

# Argument parser
def build_arg_parser():
    parser = argparse.ArgumentParser(description='Classify data using \
        Ensemble Learning techniques')
    parser.add_argument('--classifier-type', dest='classifier_type',
        required=True, choices=['rf', 'erf'], help="Type of classifier \
        to use; can be either 'rf' or 'erf'")
    return parser

if __name__ == '__main__':
    # Parse the input arguments
    args = build_arg_parser().parse_args()
    classifier_type = args.classifier_type

    # Load input data
    input_file = 'data_random_forests.txt'
    data = np.loadtxt(input_file, delimiter=',')
    X, y = data[:, :-1], data[:, -1]

    # Separate input data into three classes based on labels
    class_0 = np.array(X[y==0])
    class_1 = np.array(X[y==1])
    class_2 = np.array(X[y==2])

    # Visualize input data
    plt.figure()
    plt.scatter(class_0[:, 0], class_0[:, 1], s=75, facecolors='white',
        edgecolors='black', linewidth=1, marker='s')
    plt.scatter(class_1[:, 0], class_1[:, 1], s=75, facecolors='white',
        edgecolors='black', linewidth=1, marker='o')
    plt.scatter(class_2[:, 0], class_2[:, 1], s=75, facecolors='white',
        edgecolors='black', linewidth=1, marker='^')
    plt.title('Input data')

    # Split data into training and testing datasets
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25,
        random_state=5)

    # Ensemble Learning classifier
    params = {'n_estimators': 100, 'max_depth': 4, 'random_state': 0}
    if classifier_type == 'rf':
        classifier = RandomForestClassifier(**params)
    else:
        classifier = ExtraTreesClassifier(**params)

    classifier.fit(X_train, y_train)
    visualize_classifier(classifier, X_train, y_train)
```

```

y_test_pred = classifier.predict(X_test)
visualize_classifier(classifier, X_test, y_test)

# Evaluate classifier performance
class_names = ['Class-0', 'Class-1', 'Class-2']
print("\n" + "#" * 40)
print("\nClassifier performance on training dataset\n")
print(classification_report(y_train, classifier.predict(X_train),
target_names=class_names))
print("#" * 40 + "\n")

print("#" * 40)
print("\nClassifier performance on test dataset\n")
print(classification_report(y_test, y_test_pred, target_names=class_names))
print("#" * 40 + "\n")

# Compute confidence
test_datapoints = np.array([[5, 5], [3, 6], [6, 4], [7, 2], [4, 4], [5, 2]])

print("\nConfidence measure:")
for datapoint in test_datapoints:
    probabilities = classifier.predict_proba([datapoint])[0]
    predicted_class = 'Class-' + str(np.argmax(probabilities))
    print('\nDatapoint:', datapoint)
    print('Predicted class:', predicted_class)

# Visualize the datapoints
visualize_classifier(classifier, test_datapoints, [0]*len(test_datapoints),
                    'Test datapoints')

plt.show()

```

```
PS F:\4 kypc\CWI\lab5> python3 random_forests.py --classifier-type rf
```

```
#####
```

```
Classifier performance on training dataset
```

	precision	recall	f1-score	support
Class-0	0.91	0.86	0.88	221
Class-1	0.84	0.87	0.86	230
Class-2	0.86	0.87	0.86	224
accuracy			0.87	675
macro avg	0.87	0.87	0.87	675
weighted avg	0.87	0.87	0.87	675

```
#####
```

```
#####
```

```
Classifier performance on test dataset
```

	precision	recall	f1-score	support
Class-0	0.92	0.85	0.88	79
Class-1	0.86	0.84	0.85	70
Class-2	0.84	0.92	0.88	76
accuracy			0.87	225
macro avg	0.87	0.87	0.87	225
weighted avg	0.87	0.87	0.87	225

```
#####
```

```
Confidence measure:
```

```
Datapoint: [5 5]
```

```
Predicted class: Class-0
```

```
Datapoint: [3 6]
```

```
Predicted class: Class-0
```

```
Datapoint: [6 4]
```

```
Predicted class: Class-1
```

```
Datapoint: [7 2]
```

```
Predicted class: Class-1
```

```
Datapoint: [4 4]
```

```
Predicted class: Class-2
```

```
Datapoint: [5 2]
```

```
Predicted class: Class-2
```

Figure 2

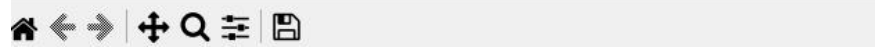
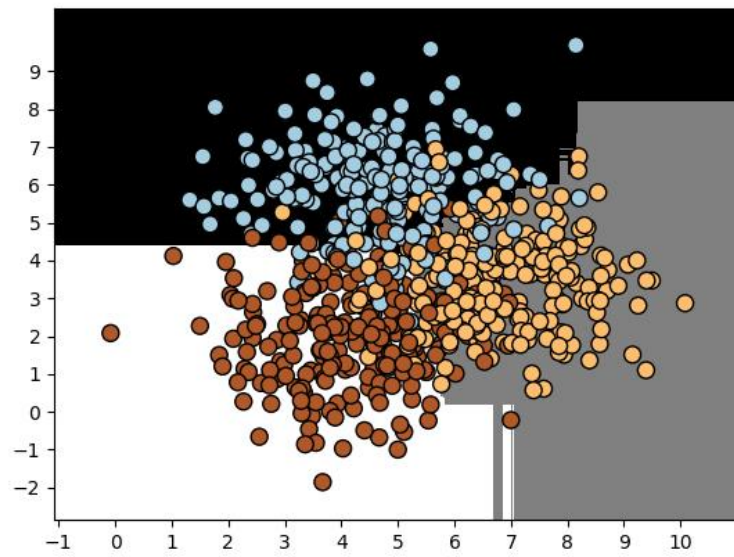
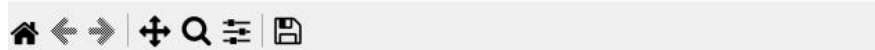
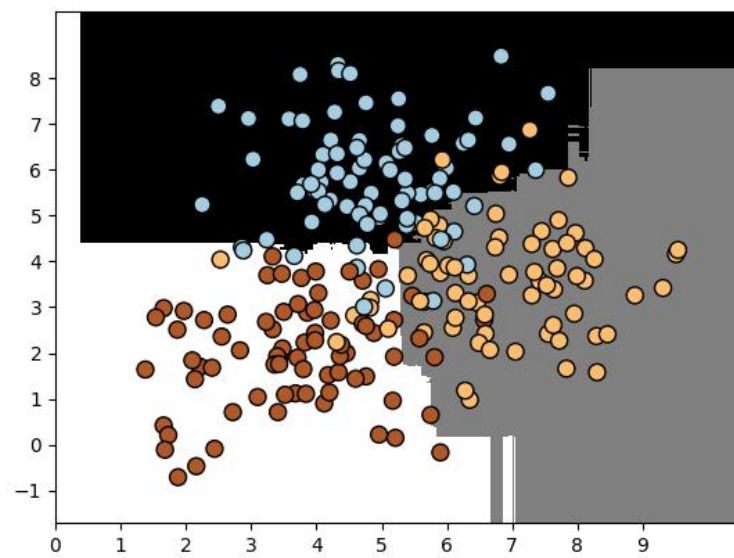
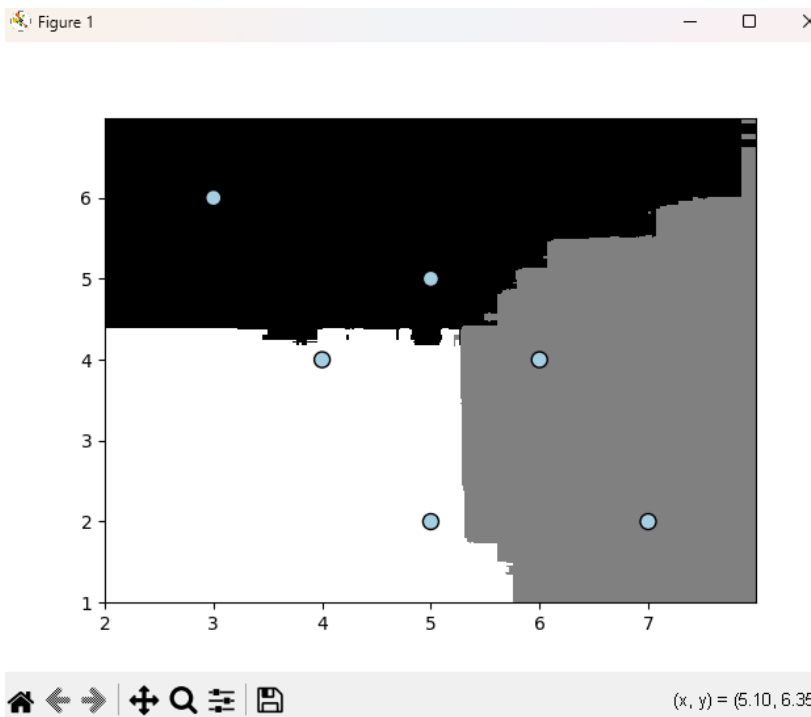
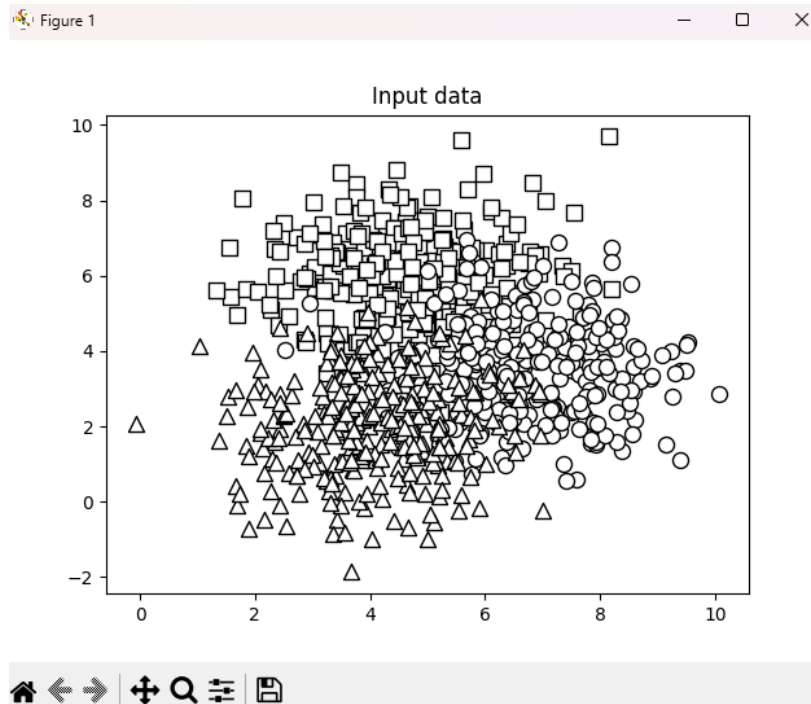


Figure 1





Модель Random Forest демонструє стабільну та надійну продуктивність як на навчальній, так і на тестовій вибірках. Середні значення точності, відгуку та F1-міри на навчальному наборі становлять 0.87 для всіх класів. Найкраща точність спостерігається для класу Class-0 (0.91), тоді як відгук і F1-міра є збалансованими між усіма класами. На тестовому наборі загальна точність моделі також дорівнює 0.87, а значення F1-міри для окремих класів знаходяться в межах від 0.85 до 0.88, що свідчить про хорошу узгодженість у роботі моделі. Передбачення для окремих тестових точок відповідають результатам навчання, що підтверджує стабільність моделі. Графіки ілюструють розподіл вхідних даних, межі між класами, а також передбачення та рівень впевненості моделі. У підсумку, модель демонструє високу точність і ефективність у виконанні

класифікаційного завдання, забезпечуючи збалансовані показники точності, відгуку та F1-міри.

## Erf

```
PS F:\4 курс\СШІ\lab5> python3 random_forests.py --classifier-type erf
```

```
#####
```

Classifier performance on training dataset

	precision	recall	f1-score	support
Class-0	0.89	0.83	0.86	221
Class-1	0.82	0.84	0.83	230
Class-2	0.83	0.86	0.85	224
accuracy			0.85	675
macro avg	0.85	0.85	0.85	675
weighted avg	0.85	0.85	0.85	675

```
#####
```

```
#####
```

Classifier performance on test dataset

	precision	recall	f1-score	support
Class-0	0.92	0.85	0.88	79
Class-1	0.84	0.84	0.84	70
Class-2	0.85	0.92	0.89	76
accuracy			0.87	225
macro avg	0.87	0.87	0.87	225
weighted avg	0.87	0.87	0.87	225

```
#####
```

```
Confidence measure:  
  
Datapoint: [5 5]  
Predicted class: Class-0  
  
Datapoint: [3 6]  
Predicted class: Class-0  
  
Datapoint: [6 4]  
Predicted class: Class-1  
  
Datapoint: [7 2]  
Predicted class: Class-1  
  
Datapoint: [4 4]  
Predicted class: Class-2  
  
Datapoint: [5 2]  
Predicted class: Class-2
```

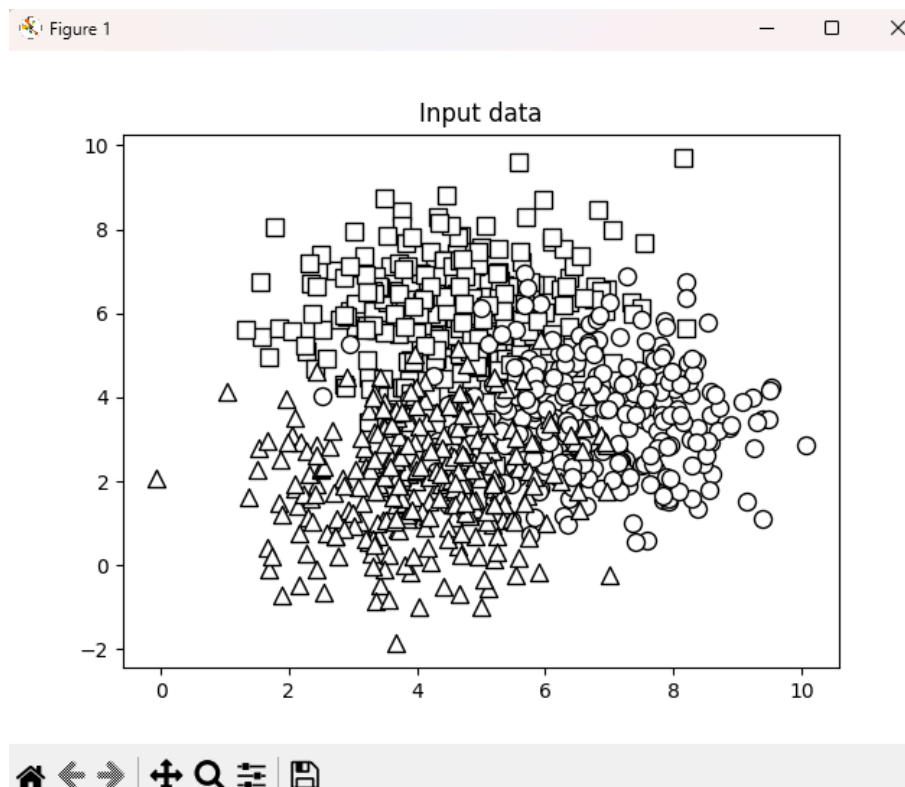


Figure 2

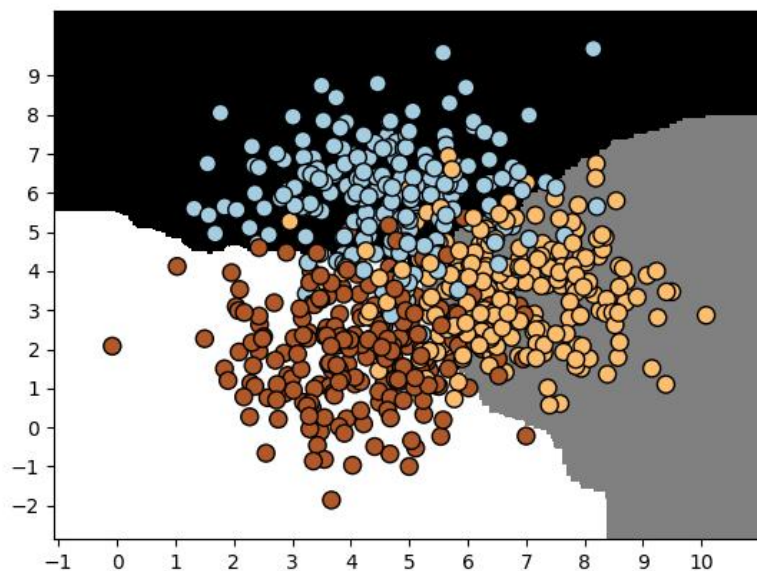
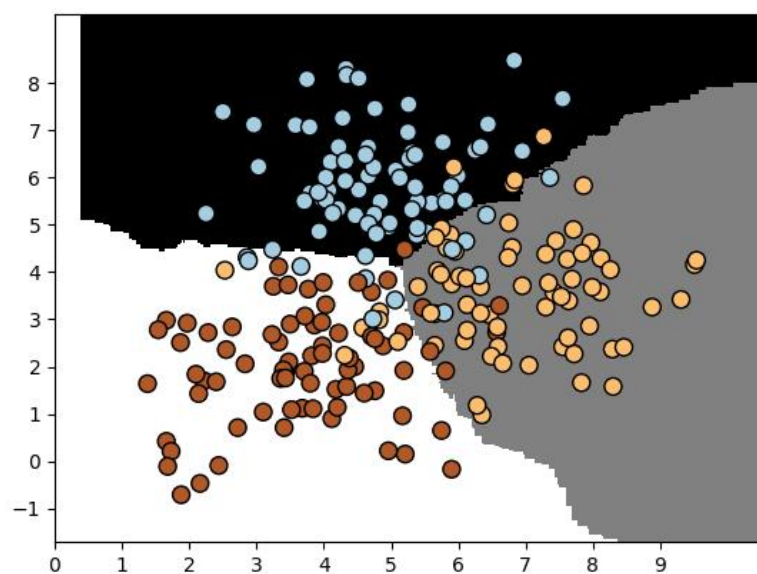
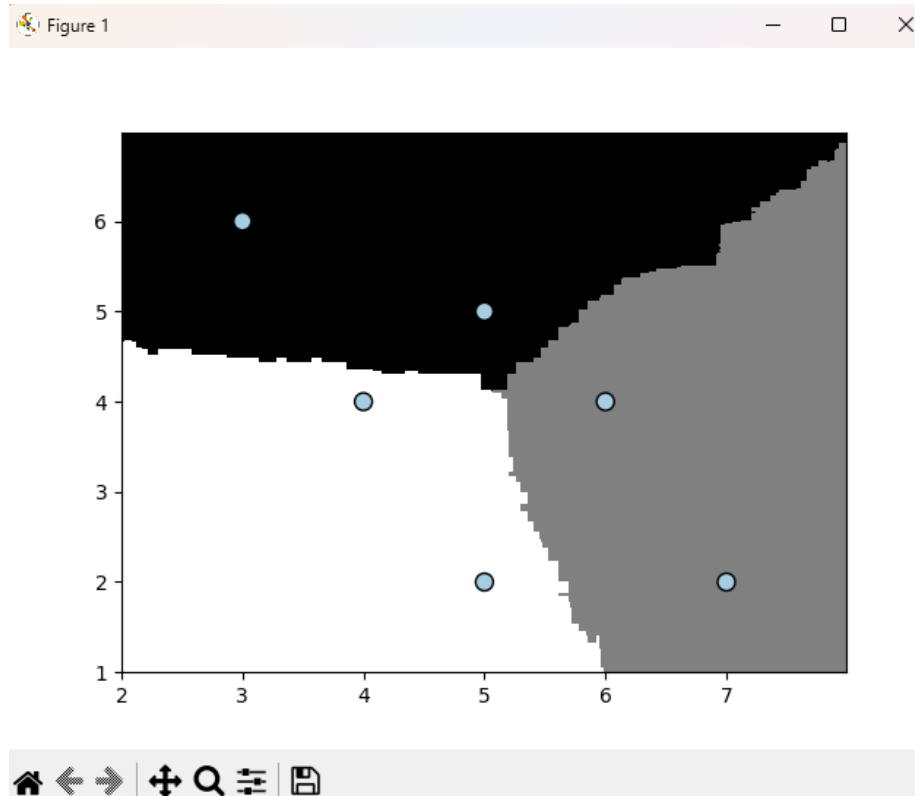


Figure 1







Модель ExtraTreesClassifier (ERF) демонструє стабільну ефективність як на навчальному, так і на тестовому наборах даних. На навчальній вибірці середні показники точності, відгуку та F1-міри дорівнюють 0.85, що свідчить про узгоджену роботу моделі. Найкращий рівень точності серед класів спостерігається для Class-0 і складає 0.89. На тестовому наборі загальна точність моделі становить 0.87, а F1-міра для всіх класів знаходиться на аналогічному рівні, що підтверджує її надійність.

Передбачення для конкретних точок демонструють здатність моделі ефективно класифікувати приклади: точки [5; 5] та [3; 6] належать до Class-0, а точки [6; 4] і [7; 2] — до Class-1.

Візуалізації, що показують розподіл вхідних даних і зони рішень моделі, підтверджують її збалансовану продуктивність, відображаючи межі між класами та рівень впевненості в передбаченнях.

## Завдання №2. Обробка дисбалансу класів

```
import sys
import numpy as np
import matplotlib.pyplot as plt
from sklearn.ensemble import ExtraTreesClassifier
from sklearn.metrics import classification_report
from sklearn.model_selection import train_test_split

from utilities import visualize_classifier

input_file = 'data_imbalance.txt'
data = np.loadtxt(input_file, delimiter=',')
X, y = data[:, :-1], data[:, -1]
class_0 = np.array(X[y==0])
class_1 = np.array(X[y==1])
```

```

plt.figure()
plt.scatter(class_0[:, 0], class_0[:, 1], s=75, facecolors='black',
            edgecolors='black', linewidth=1, marker='x')
plt.scatter(class_1[:, 0], class_1[:, 1], s=75, facecolors='white',
            edgecolors='black', linewidth=1, marker='o')
plt.title('Input data')

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25,
                                                    random_state=5)

# Класифікатор на основі гранично випадкових лісів
params = {'n_estimators': 100, 'max_depth': 4, 'random_state': 0}

if len(sys.argv) > 1:
    if sys.argv[1] == 'balance':
        params = {'n_estimators': 100, 'max_depth': 4, 'random_state': 0,
                  'class_weight': 'balanced'}
    else:
        raise TypeError("Invalid input argument; should be 'balance'")

classifier = ExtraTreesClassifier(**params)
classifier.fit(X_train, y_train)
visualize_classifier(classifier, X_train, y_train)

# Передбачимо та візуалізуємо результат для тестового набору даних
y_test_pred = classifier.predict(X_test)
visualize_classifier(classifier, X_test, y_test)

class_names = ['Class-0', 'Class-1']
print("\n" + "#" * 40)
print("\nClassifier performance on training dataset\n")
print(classification_report(y_train, classifier.predict(X_train),
                             target_names=class_names))
print("#" * 40 + "\n")

print("#" * 40)
print("\nClassifier performance on test dataset\n")
print(classification_report(y_test, y_test_pred, target_names=class_names))
print("#" * 40 + "\n")

plt.show()

```

```

#####

Classifier performance on training dataset

              precision    recall  f1-score   support

   Class-0       0.44        0.93        0.60        181
   Class-1       0.98        0.77        0.86        944

 accuracy              0.80        1125
 macro avg           0.71        0.85        0.73        1125
weighted avg           0.89        0.80        0.82        1125

#####

```

```
#####
```

Classifier performance on test dataset

	precision	recall	f1-score	support
Class-0	0.45	0.94	0.61	69
Class-1	0.98	0.74	0.84	306
accuracy			0.78	375
macro avg	0.72	0.84	0.73	375
weighted avg	0.88	0.78	0.80	375

```
#####
```

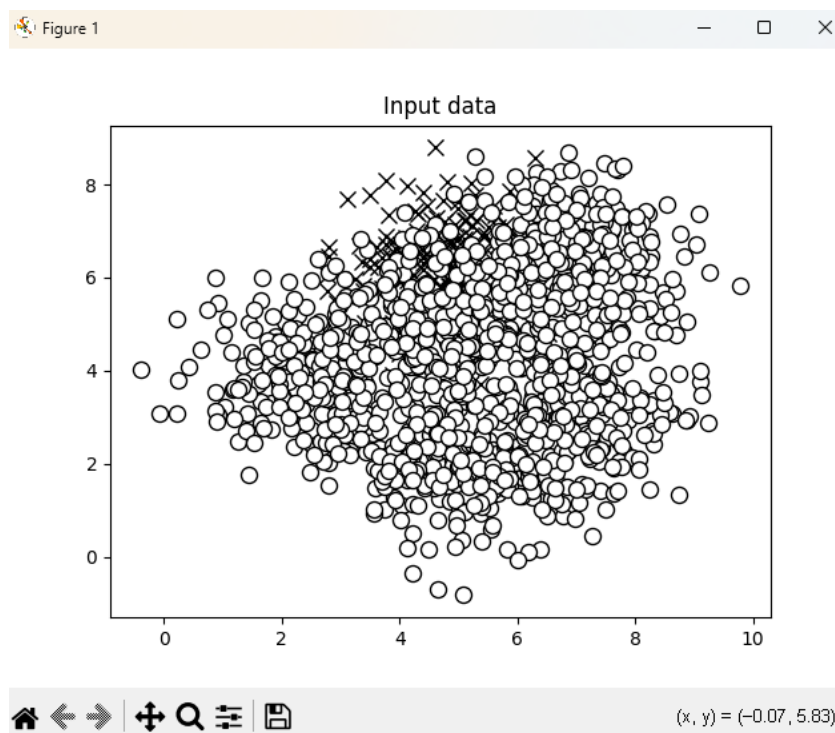


Figure 2

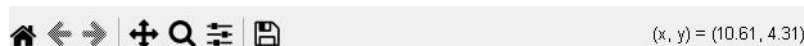
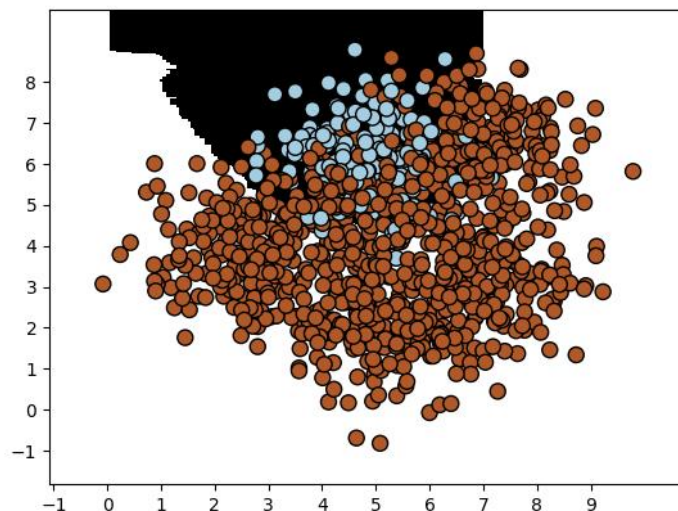
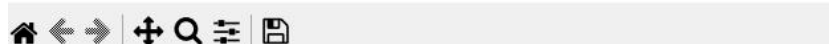
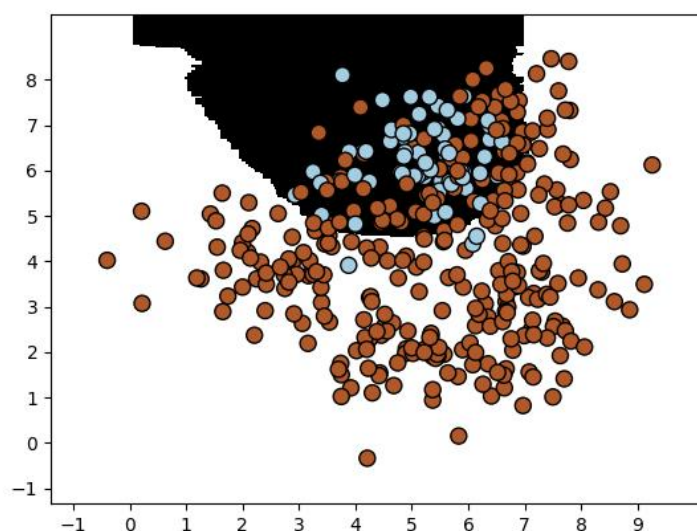


Figure 1



Результати аналізу свідчать про наявність проблеми з дисбалансом класів у моделі. На навчальній вибірці спостерігається низька точність для Class-0 (0.44) при високому відгуку (0.93), що вказує на велику кількість хибнопозитивних прогнозів. У той же час для Class-1 точність дуже висока (0.98), але відгук становить лише 0.77, що свідчить про пропуск значної частини об'єктів цього класу. Загальна точність становить 0.80, а макро F1-міра дорівнює 0.73, що підтверджує суттєвий дисбаланс у показниках для різних класів. Тестова вибірка: точність для Class-0 становить 0.45, відгук — 0.94, тоді як для Class-1 точність сягає 0.98, а відгук — 0.74. Загальна точність — 0.78, макро F1-міра — 0.73.

### Завдання №3. Знаходження оптимальних навчальних параметрів за допомогою сіткового пошуку.

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.metrics import classification_report
from sklearn.ensemble import ExtraTreesClassifier
from sklearn.model_selection import train_test_split, GridSearchCV
from utilities import visualize_classifier

input_file = 'data_random_forests.txt'
data = np.loadtxt(input_file, delimiter=',')
X, y = data[:, :-1], data[:, -1]

class_0 = np.array(X[y==0])
class_1 = np.array(X[y==1])
class_2 = np.array(X[y==2])

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25,
random_state=5)

# Визначення сітки значень параметрів
parameter_grid = [
    {'n_estimators': [100], 'max_depth': [2, 4, 7, 12, 16]},
    {'max_depth': [4], 'n_estimators': [25, 50, 100, 250]}
]

metrics = ['precision_weighted', 'recall_weighted']
for metric in metrics:
    print("\n#### Searching optimal parameters for", metric)

    classifier = GridSearchCV(
        ExtraTreesClassifier(random_state=0),
        parameter_grid, cv=5, scoring=metric)
    classifier.fit(X_train, y_train)

    print("\nGrid scores for the parameter grid:")
    means = classifier.cv_results_['mean_test_score']
    params = classifier.cv_results_['params']
    for mean, param in zip(means, params):
        print(param, '-->', round(mean, 3))
    print("\nBest parameters:", classifier.best_params_)

y_pred = classifier.predict(X_test)
print("\nPerformance report:\n")
print(classification_report(y_test, y_pred))
```

```
C:\Users\dimad\AppData\Local\Programs\Python\Python39\python.exe "F:/4 курс/CWI/lab5/LR_5_task_3.py"
```

```
##### Searching optimal parameters for precision_weighted
```

```
Grid scores for the parameter grid:
```

```
{ 'max_depth': 2, 'n_estimators': 100} --> 0.85  
{ 'max_depth': 4, 'n_estimators': 100} --> 0.841  
{ 'max_depth': 7, 'n_estimators': 100} --> 0.844  
{ 'max_depth': 12, 'n_estimators': 100} --> 0.832  
{ 'max_depth': 16, 'n_estimators': 100} --> 0.816  
{ 'max_depth': 4, 'n_estimators': 25} --> 0.846  
{ 'max_depth': 4, 'n_estimators': 50} --> 0.84  
{ 'max_depth': 4, 'n_estimators': 100} --> 0.841  
{ 'max_depth': 4, 'n_estimators': 250} --> 0.845
```

```
Best parameters: { 'max_depth': 2, 'n_estimators': 100}
```

```
##### Searching optimal parameters for recall_weighted
```

```
Grid scores for the parameter grid:
```

```
{ 'max_depth': 2, 'n_estimators': 100} --> 0.843  
{ 'max_depth': 4, 'n_estimators': 100} --> 0.837  
{ 'max_depth': 7, 'n_estimators': 100} --> 0.841  
{ 'max_depth': 12, 'n_estimators': 100} --> 0.83  
{ 'max_depth': 16, 'n_estimators': 100} --> 0.815  
{ 'max_depth': 4, 'n_estimators': 25} --> 0.843  
{ 'max_depth': 4, 'n_estimators': 50} --> 0.836  
{ 'max_depth': 4, 'n_estimators': 100} --> 0.837  
{ 'max_depth': 4, 'n_estimators': 250} --> 0.841
```

```
Best parameters: { 'max_depth': 2, 'n_estimators': 100}
```

```
Performance report:
```

	precision	recall	f1-score	support
0.0	0.94	0.81	0.87	79
1.0	0.81	0.86	0.83	70
2.0	0.83	0.91	0.87	76
accuracy			0.86	225
macro avg	0.86	0.86	0.86	225
weighted avg	0.86	0.86	0.86	225

## Завдання №4. Обчислення відносної важливості ознак.

```
import numpy as np  
import pandas as pd  
import matplotlib.pyplot as plt  
from sklearn.tree import DecisionTreeRegressor  
from sklearn.ensemble import AdaBoostRegressor  
from sklearn import datasets  
from sklearn.metrics import mean_squared_error, explained_variance_score  
from sklearn.model_selection import train_test_split  
from sklearn.utils import shuffle  
  
housing_data = datasets.fetch_california_housing()
```

```

X, y = shuffle(housing_data.data, housing_data.target, random_state=7)

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=7)

regressor = AdaBoostRegressor(
    DecisionTreeRegressor(max_depth=4),
    n_estimators=400, random_state=7)
regressor.fit(X_train, y_train)

y_pred = regressor.predict(X_test)
mse = mean_squared_error(y_test, y_pred)
evs = explained_variance_score(y_test, y_pred)

print("\nADABOOST REGRESSOR")
print("Mean squared error =", round(mse, 2))
print("Explained variance score =", round(evs, 2))

feature_importances = regressor.feature_importances_
feature_names = housing_data.feature_names

feature_importances = 100.0 * (feature_importances / max(feature_importances))

index_sorted = np.flipud(np.argsort(feature_importances))

pos = np.arange(index_sorted.shape[0]) + 0.5

feature_names = np.array(housing_data.feature_names)

plt.figure()
plt.bar(pos, feature_importances[index_sorted], align='center')
plt.xticks(pos, feature_names[index_sorted])
plt.ylabel('Relative Importance')
plt.title('Оценка важности признаков с использованием регрессора AdaBoost')
plt.show()

```

```

C:\Users\dimad\AppData\Local\Programs\Python\Python39\python.exe "F:/4 кypc/CWI/lab5/LR_4_task_4.py"

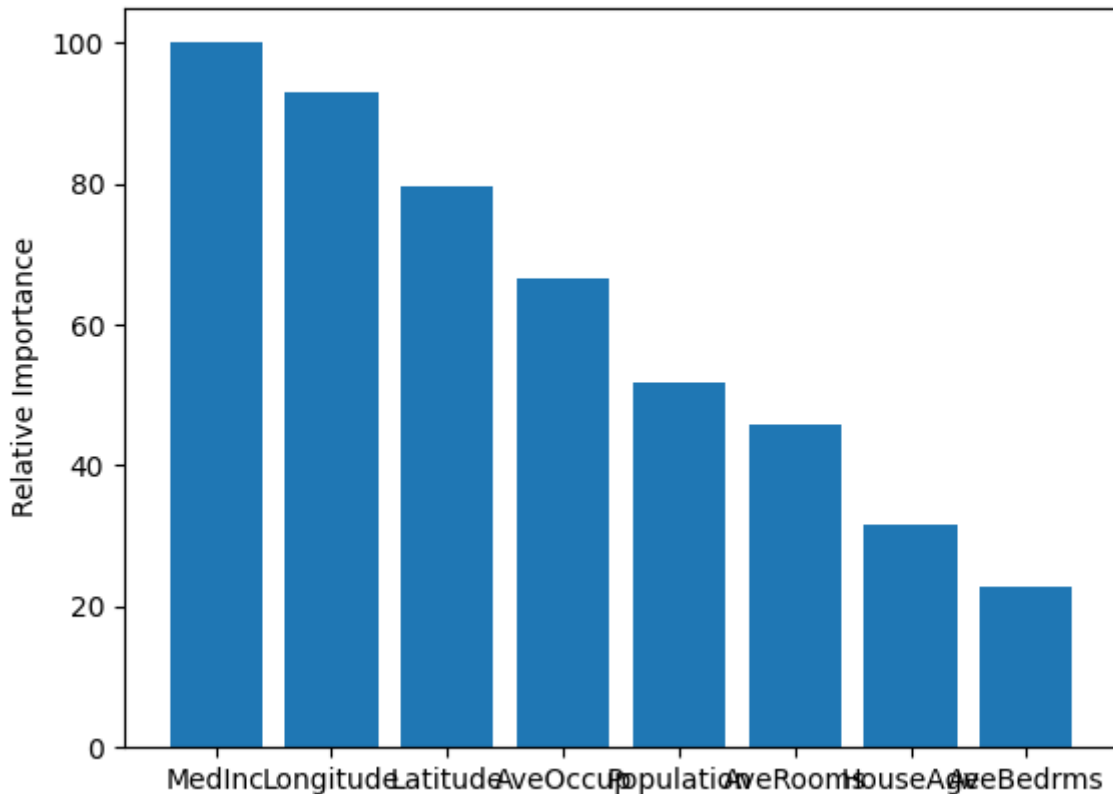
```

```

ADABOOST REGRESSOR
Mean squared error = 1.18
Explained variance score = 0.47

```

## Оценка важности признаков с использованием регрессора AdaBoost



**Завдання №5.** Прогнозування інтенсивності дорожнього руху за допомогою класифікатора на основі гранично випадкових лісів.

```
import numpy as np
from sklearn.metrics import classification_report, mean_absolute_error
from sklearn import preprocessing
from sklearn.model_selection import train_test_split
from sklearn.ensemble import ExtraTreesRegressor

input_file = 'traffic_data.txt'
data = []
with open(input_file, 'r') as f:
    for line in f.readlines():
        items = line[:-1].split(',')
        data.append(items)

data = np.array(data)

label_encoder = []
X_encoded = np.empty(data.shape)
for i, item in enumerate(data[0]):
    if item.isdigit():
        X_encoded[:, i] = data[:, i]
    else:
        label_encoder.append(preprocessing.LabelEncoder())
        X_encoded[:, i] = label_encoder[-1].fit_transform(data[:, i])
X = X_encoded[:, :-1].astype(int)
y = X_encoded[:, -1].astype(int)

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25,
                                                    random_state=5)

params = {'n_estimators' : 100, 'max_depth':4, 'random_state':0}
```



```

regressor = ExtraTreesRegressor(**params)
regressor.fit(X_train, y_train)

y_pred = regressor.predict(X_test)
print("Mean absolute error: ", round(mean_absolute_error(y_test, y_pred), 2))

test_datapoint = ['Saturday', '10:20', 'Atlanta', 'no']
test_datapoint_encoded = [-1]*len(test_datapoint)
count = 0
for i, item in enumerate(test_datapoint):
    if item.isdigit():
        test_datapoint_encoded[i] = int(test_datapoint[i])
    else :
        test_datapoint_encoded[i] =
int(label_encoder[count].transform([test_datapoint[i]])[0])
        count += 1
test_datapoint_encoded = np.array(test_datapoint_encoded)
test_datapoint_encoded = test_datapoint_encoded.reshape(1, -1)

print("Predicted traffic: ", int(regressor.predict(test_datapoint_encoded)[0]))

```

```

C:\Users\dimad\AppData\Local\Programs\Python\Python39\python.exe "F:/4 kypc/CWI/lab5/LR_5_task_5.py"
Mean absolute error: 7.42
Predicted traffic: 26

```