

ЛАБОРАТОРНА РОБОТА № 5

ДОСЛІДЖЕННЯ МЕТОДІВ АНСАМБЛЕВОГО НАВЧАННЯ

Мета роботи: використовуючи спеціалізовані бібліотеки та мову програмування Python дослідити методи ансамблів у машинному навчанні.

Хід роботи

Завдання 1. Створення класифікаторів на основі випадкових та гранично випадкових лісів

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

# Парсер аргументів
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__':
    # Вилучення вхідних аргументів
    args = build_arg_parser().parse_args()
    classifier_type = args.classifier_type

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

    # Розбиття вхідних даних на три класи
    class_0 = np.array(X[y == 0])
    class_1 = np.array(X[y == 1])
    class_2 = np.array(X[y == 2])

    # Візуалізація вхідних даних
    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')
```

					ДУ «Житомирська політехніка».22.121.07.806 – ІПЗк								
Змн.	Арк.	№ докум.	Підпис	Дата									
Розроб.		Кияшенко А.С.			Системи штучного інтелекту Лабораторна №5				Літ.	Арк.	Аркушів		
Перевір.											1	19	
Реценз.									ФІКТ Гр. ІПЗк-19-1				
Н. Контр.													
Затверд.													

```

plt.scatter(class_2[:, 0], class_2[:, 1], s=75, facecolors='white',
            edgecolors='black', linewidth=1, marker='^')
plt.title('Input data')

# Розбивка даних на навчальний та тестовий набори
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.25, random_state=5)

# Класифікатор на основі ансамблевого навчання
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, 'Training dataset')

y_test_pred = classifier.predict(X_test)
visualize_classifier(classifier, X_test, y_test, 'Тестовий набір даних')

# Перевірка роботи класифікатора
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")

# Обчислення параметрів довірливості
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)

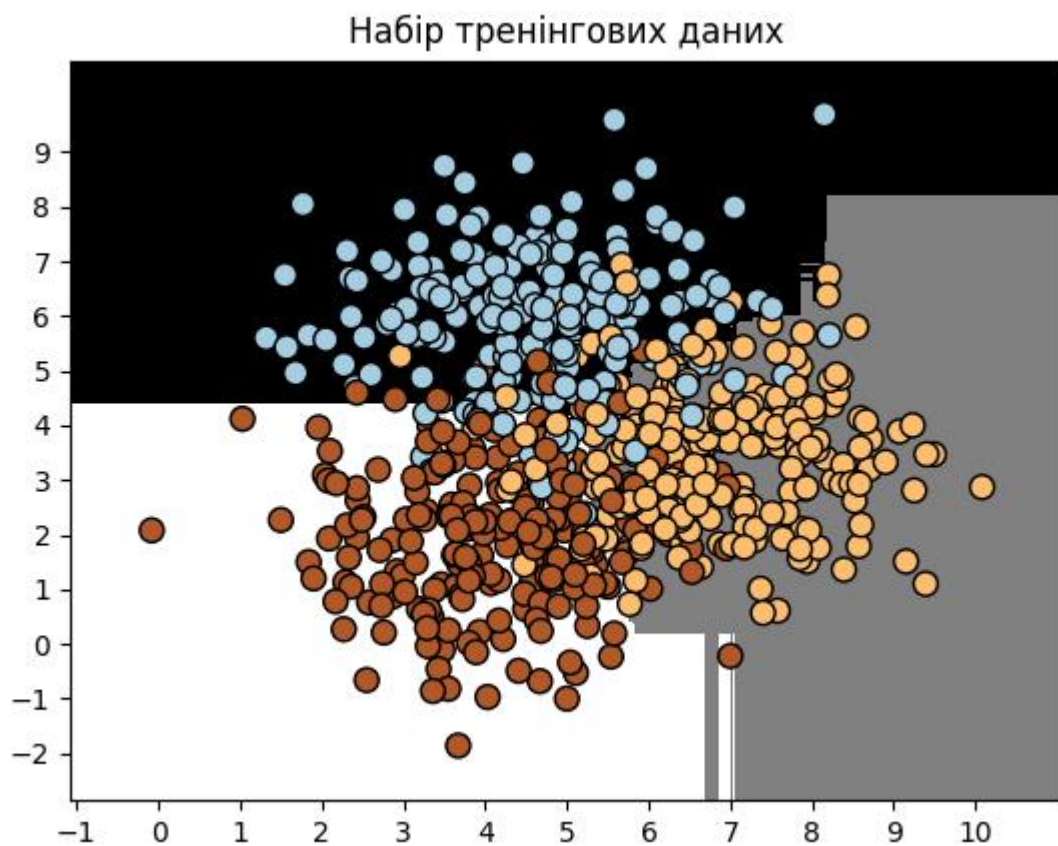
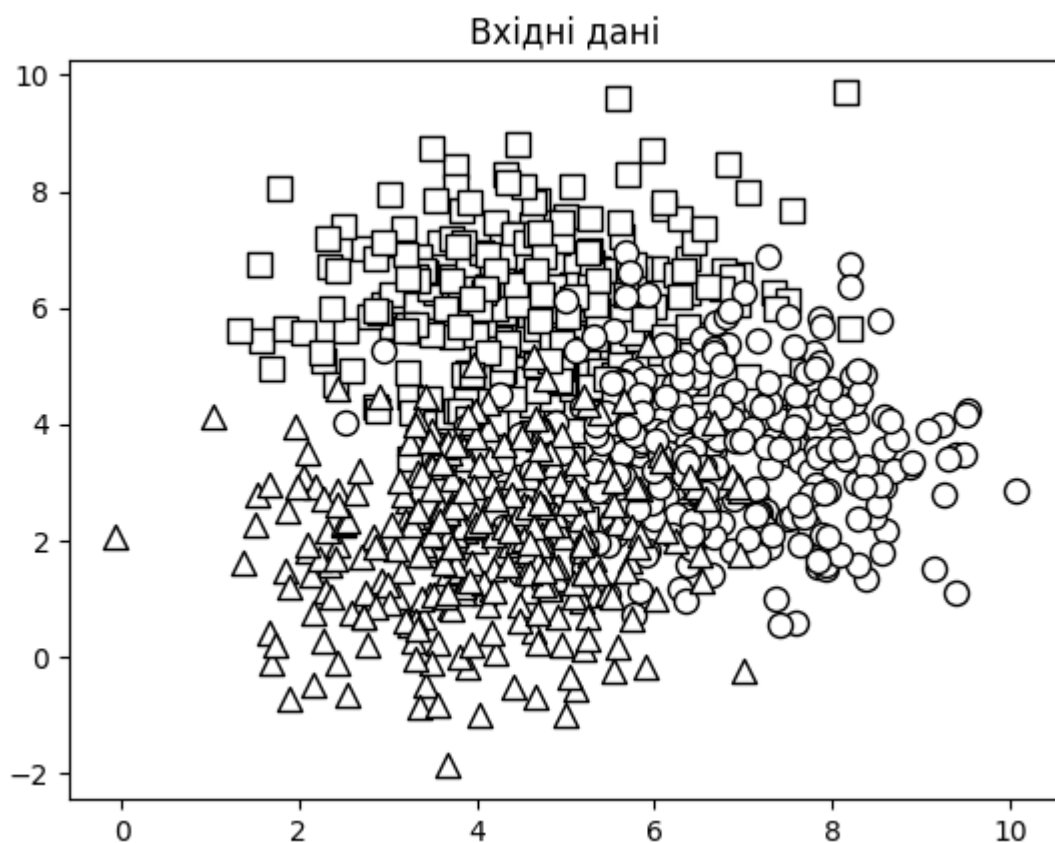
# Візуалізація точок даних
visualize_classifier(classifier, test_datapoints, [0] *
len(test_datapoints),
                    'Тестові точки даних')

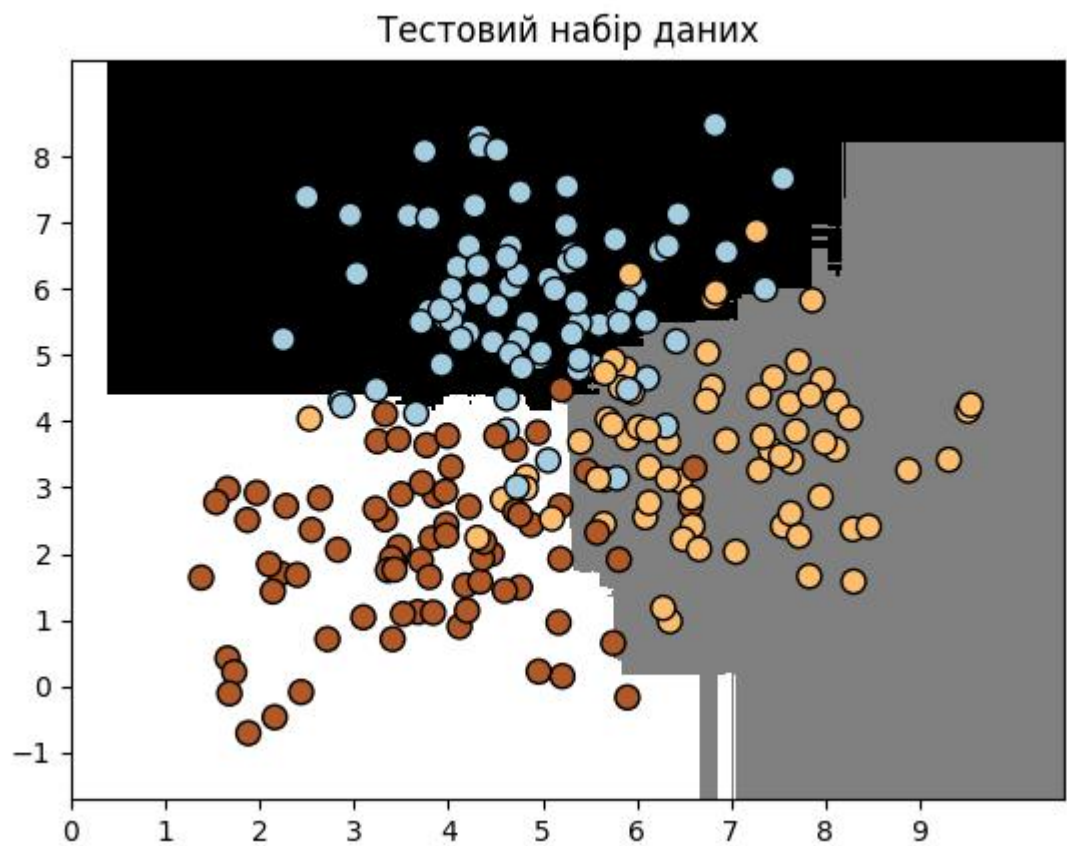
plt.show()

```

Створення класифікатора на основі випадкового лісу за допомогою прапорця **rf**.

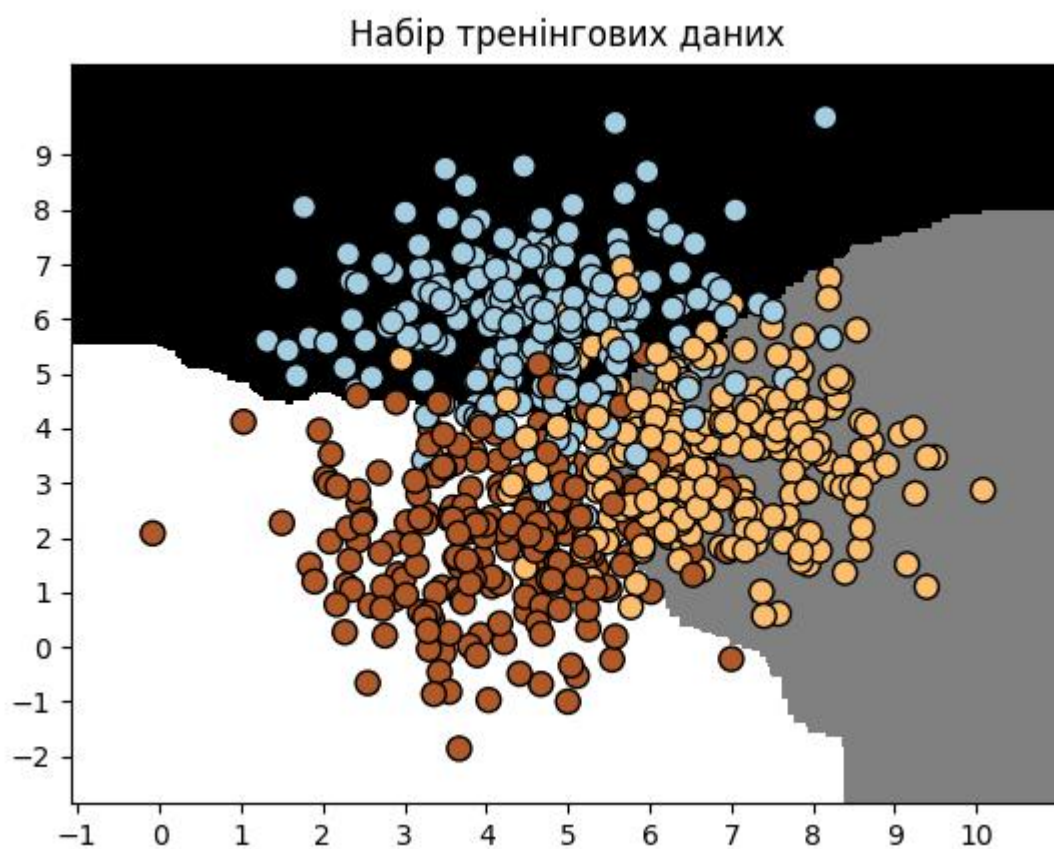
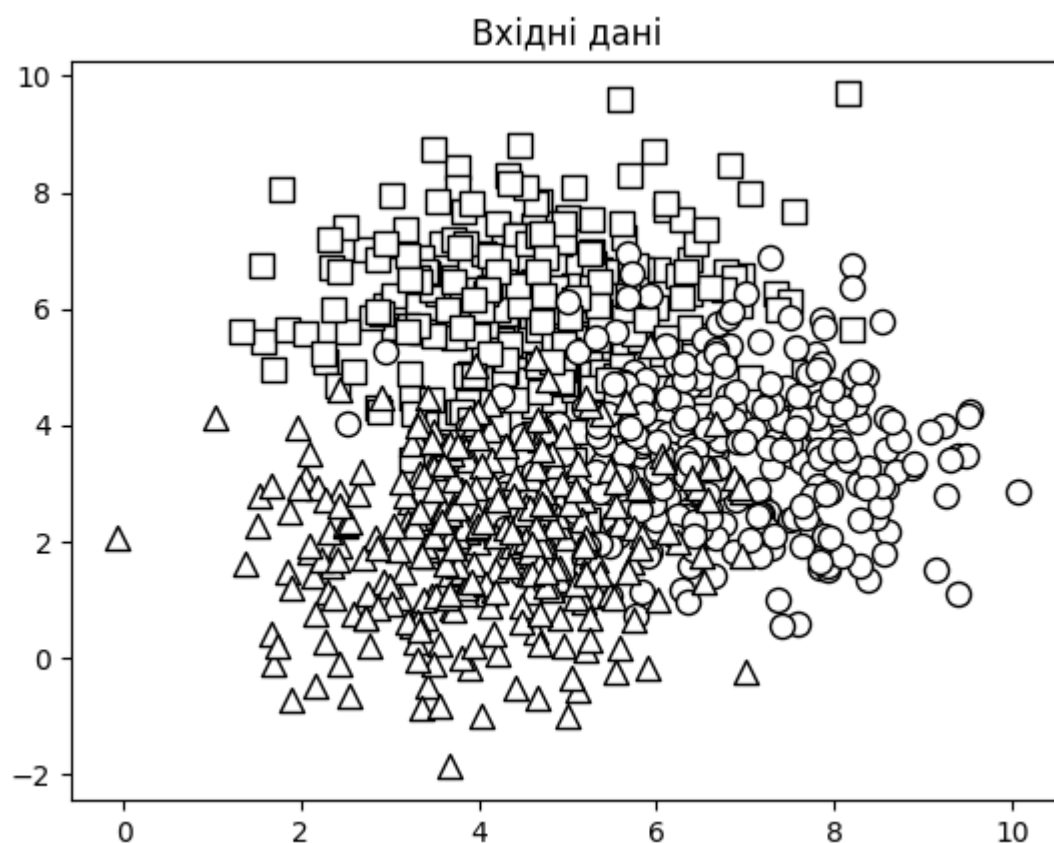
		Кияшенко А.С.			ДУ «Житомирська політехніка».22.121.07. 806– ІПЗк	Арк.
		.				2
Змн.	Арк.	№ докум.	Підпис	Дата		

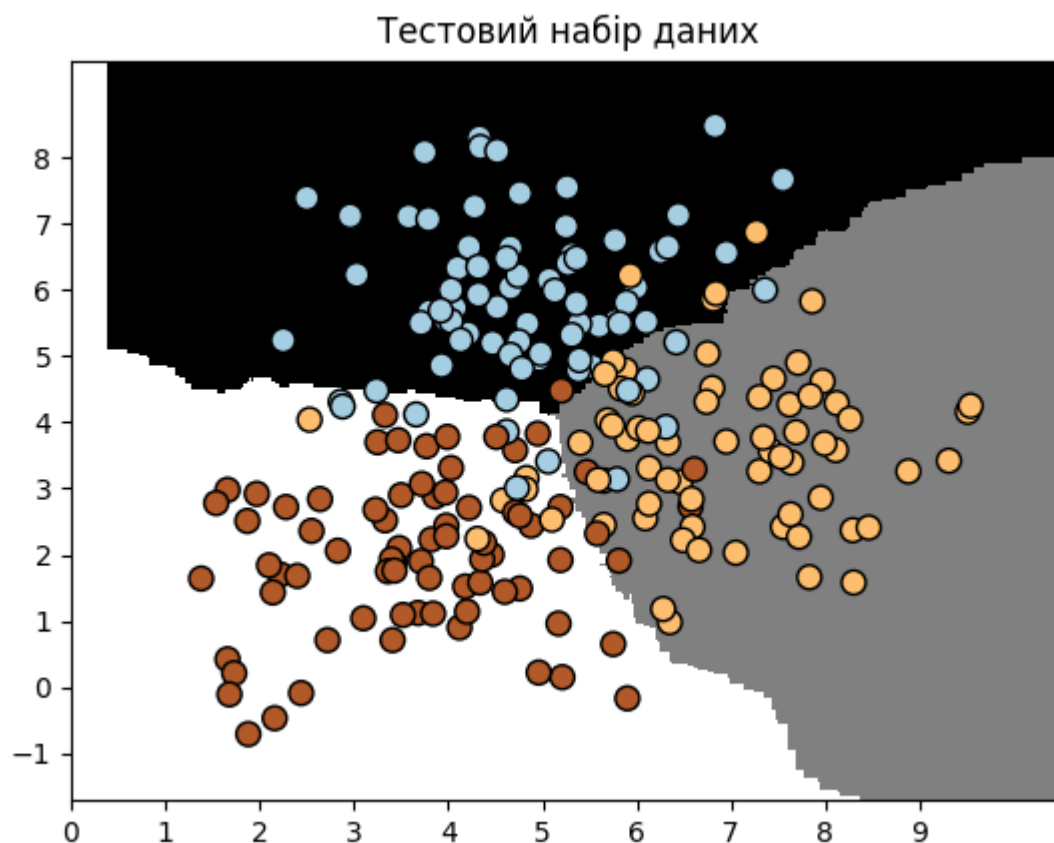




Створення класифікатора на основі гранично випадкового лісу за допомогою прапорця **erf**.

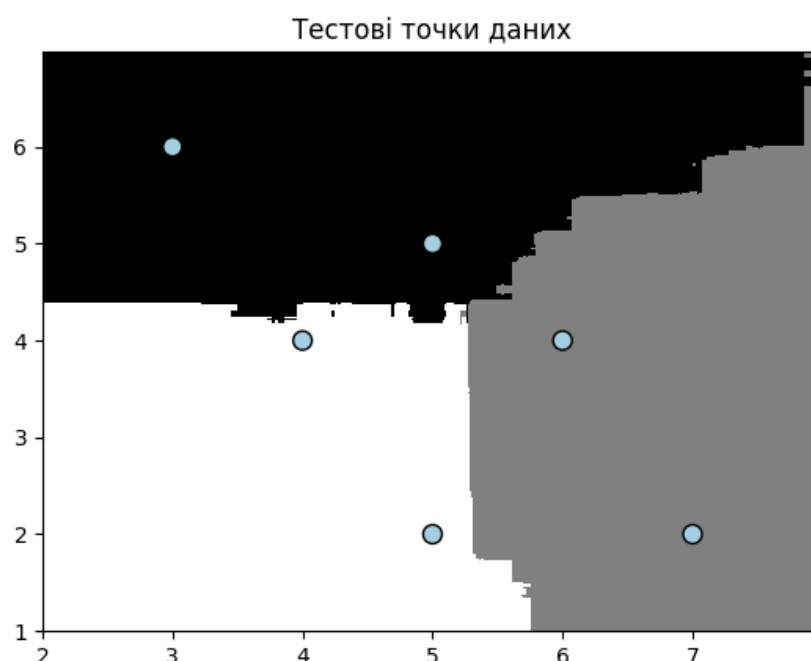
		Кияшенко А.С.			ДУ «Житомирська політехніка».22.121.07. 806– ІПЗк	Арк.
						4
Змн.	Арк.	№ докум.	Підпис	Дата		





В цьому випадку були отримані більш лагідні піки. Це обумовлено тим, що в процесі навчання гранично випадкові ліси мають більше можливостей для вибору оптимальних дерев рішень, тому, як правило, вони забезпечують отримання кращих границь.

Оцінка мір достовірності прогнозів із прапором **rf**.



		Кияшенко А.С.			ДУ «Житомирська політехніка».22.121.07. 806– ІПЗк	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		6


```

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

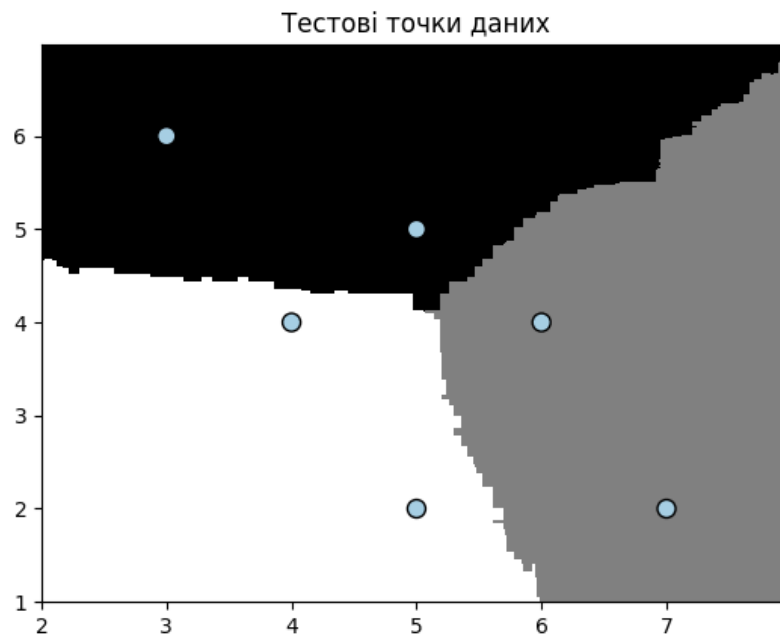
Datapoint: [4 4]
Predicted class: Class-2

Datapoint: [5 2]
Predicted class: Class-2
PS D:\LabsPoli\AI\Lab5>

```

Оцінка мір достовірності прогнозів із прапором **erf**.

		Кияшенко А.С.			ДУ «Житомирська політехніка».22.121.07. 806– ІПЗк	Арк.
		.				7
Змн.	Арк.	№ докум.	Підпис	Дата		



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

#####

#####

Datapoint: [5 5]
Predicted class: Class-0

Datapoint: [3 6]
Predicted class: Class-0

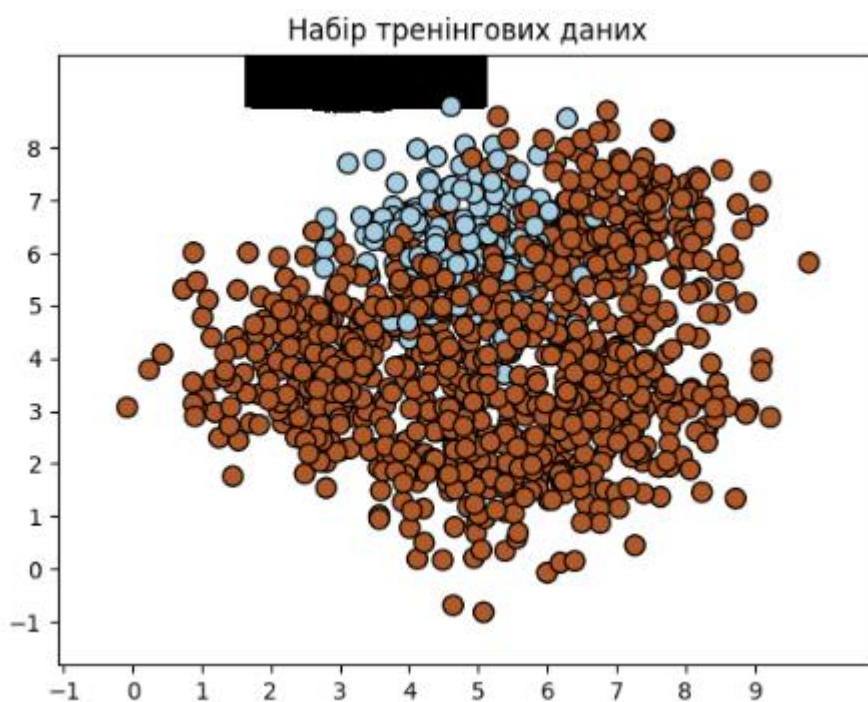
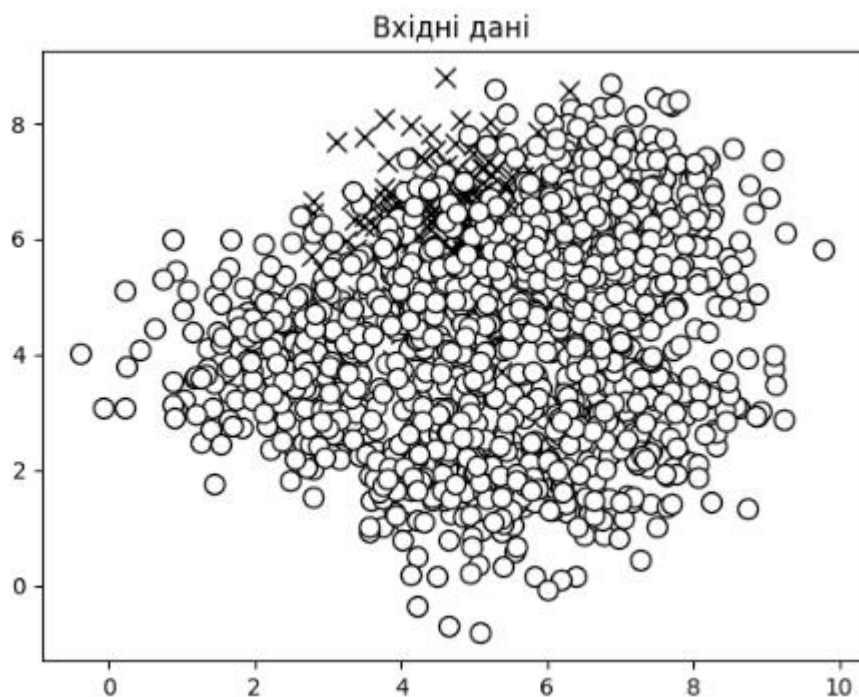
Datapoint: [6 4]
Predicted class: Class-1

Datapoint: [7 2]
Predicted class: Class-1

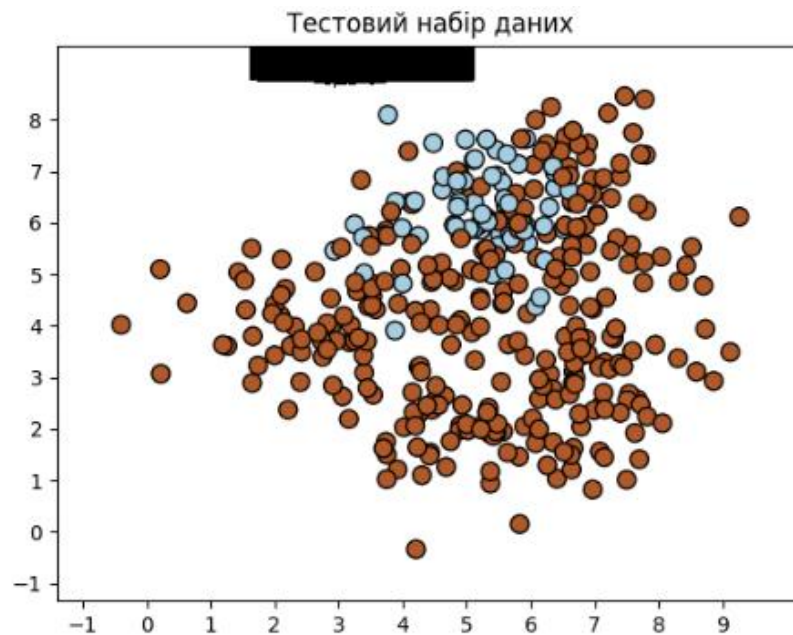
Datapoint: [4 4]
Predicted class: Class-2

Datapoint: [5 2]
Predicted class: Class-2

Завдання 2. Обробка дисбалансу класів



		Кияшенко А.С.			ДУ «Житомирська політехніка».22.121.07. 806– ІПЗк	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		9



```
PS D:\LabsPoli\AI\Lab5> python3 -W ignore .\LR_5_task_2.py
```

```
#####
```

```
Classifier performance on training dataset
```

	precision	recall	f1-score	support
Class-0	1.00	0.01	0.01	181
Class-1	0.84	1.00	0.91	944
accuracy			0.84	1125
macro avg	0.92	0.50	0.46	1125
weighted avg	0.87	0.84	0.77	1125

```
#####
```

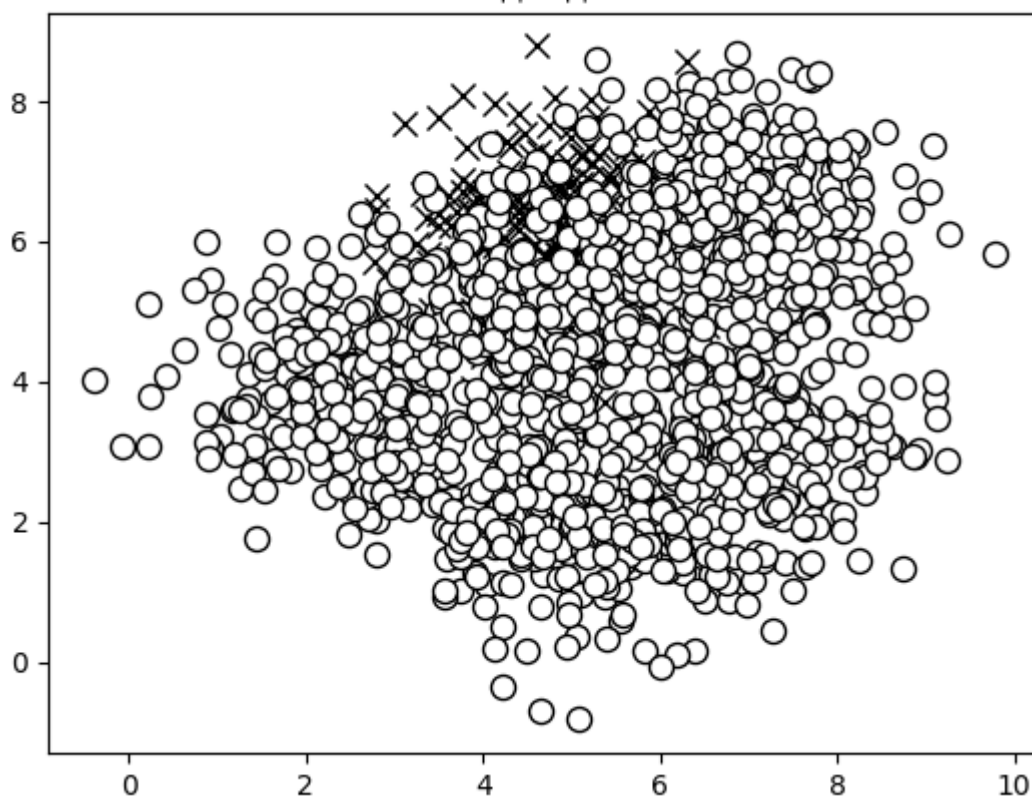
```
#####
```

```
Classifier performance on test dataset
```

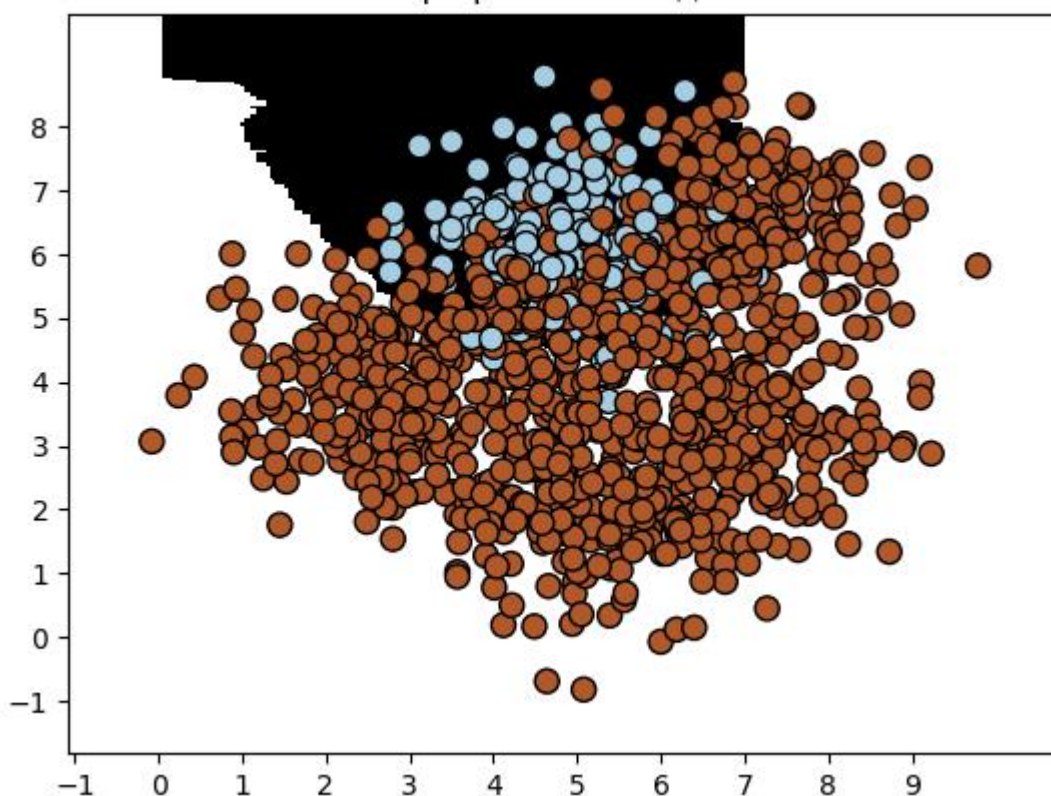
	precision	recall	f1-score	support
Class-0	0.00	0.00	0.00	69
Class-1	0.82	1.00	0.90	306
accuracy			0.82	375
macro avg	0.41	0.50	0.45	375
weighted avg	0.67	0.82	0.73	375

Для врахування дисбалансу використав параметр **balance**.

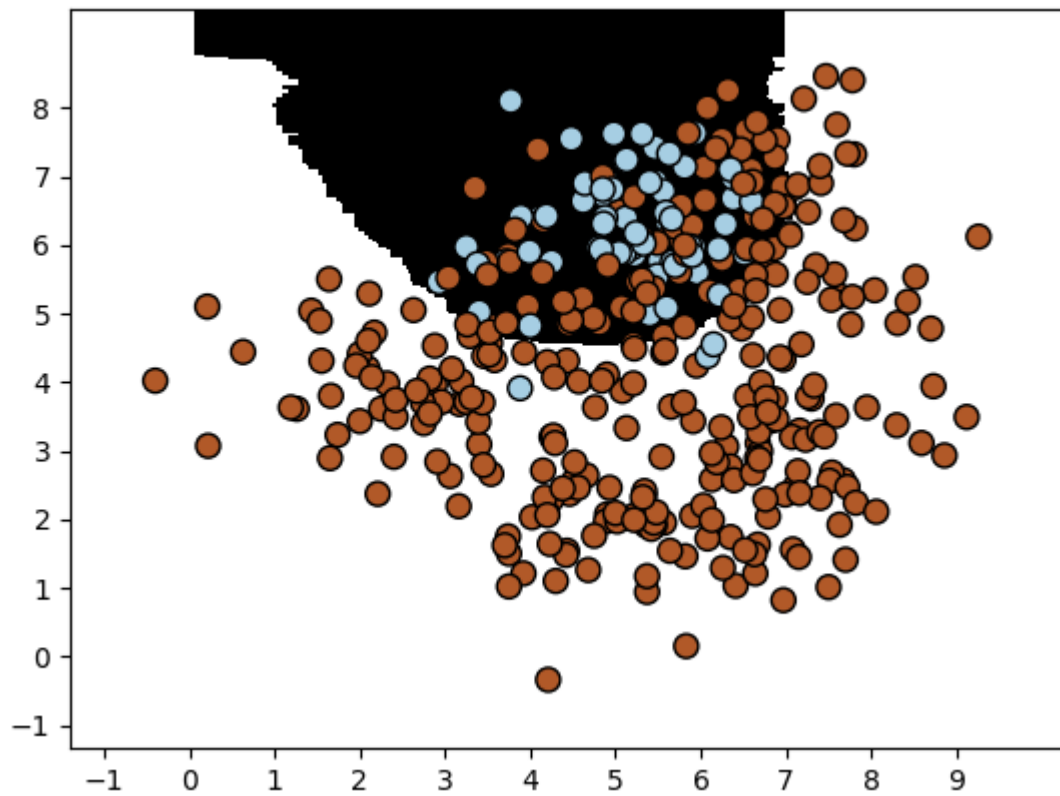
Вхідні дані



Набір тренінгових даних



Тестовий набір даних



```

PS D:\LabsPoli\AI\Lab5> python3 .\LR_5_task_2.py balance

#####

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

#####

```

Використовуючи для аналізу дані, які містяться у файлі data_imbalance.txt я провів обробку з урахуванням дисбалансу класів.

Завдання 3. Знаходження оптимальних навчальних параметрів за допомогою сіткового пошуку

```

import numpy as np
import pandas as pd
from sklearn.ensemble import ExtraTreesClassifier
from sklearn.metrics import classification_report
from sklearn.model_selection import train_test_split, GridSearchCV

```



```

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

# Розбиття даних на три класи на підставі міток
class_0 = np.array(X[y == 0])
class_1 = np.array(X[y == 1])
class_2 = np.array(X[y == 2])

# Розбиття даних на навчальний та тестовий набори
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:")
    print(pd.concat([pd.DataFrame(classifier.cv_results_["params"]),
                     pd.DataFrame(classifier.cv_results_["mean_test_score"],
                                   columns=["Accuracy"])], axis=1))

    print("\nBest parameters:", classifier.best_params_)

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

```

		Кияшенко А.С.			ДУ «Житомирська політехніка».22.121.07. 806– ІПЗк	Арк.
		.				14
Змн.	Арк.	№ докум.	Підпис	Дата		

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

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

	max_depth	n_estimators	Accuracy
0	2	100	0.849757
1	4	100	0.841135
2	7	100	0.843820
3	12	100	0.831955
4	16	100	0.816484
5	4	25	0.845574
6	4	50	0.839851
7	4	100	0.841135
8	4	250	0.844788

```
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


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

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

	max_depth	n_estimators	Accuracy
0	2	100	0.842963
1	4	100	0.837037
2	7	100	0.841481
3	12	100	0.829630
4	16	100	0.814815
5	4	25	0.842963
6	4	50	0.835556
7	4	100	0.837037
8	4	250	0.841481

```
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

Найкращі результати при параметрах max_depth = 2, n_estimators = 100.

Завдання 4. Обчислення відносної важливості ознак

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
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.utils import shuffle

# Завантаження даних із цінами на нерухомість
housing_data = datasets.load_boston()

# Перемішування даних
X, y = shuffle(housing_data.data, housing_data.target, random_state=7)
```

```
# Розбиття даних на навчальний та тестовий набори
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=7)

# Модель на основі регресора AdaBoost
regressor = AdaBoostRegressor(DecisionTreeRegressor(max_depth=4),
                              n_estimators=400, random_state=7)
regressor.fit(X_train, y_train)

# Обчислення показників ефективності регресора AdaBoost
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))

# Вилучення важливості ознак
feature_importances = regressor.feature_importances_
feature_names = housing_data.feature_names

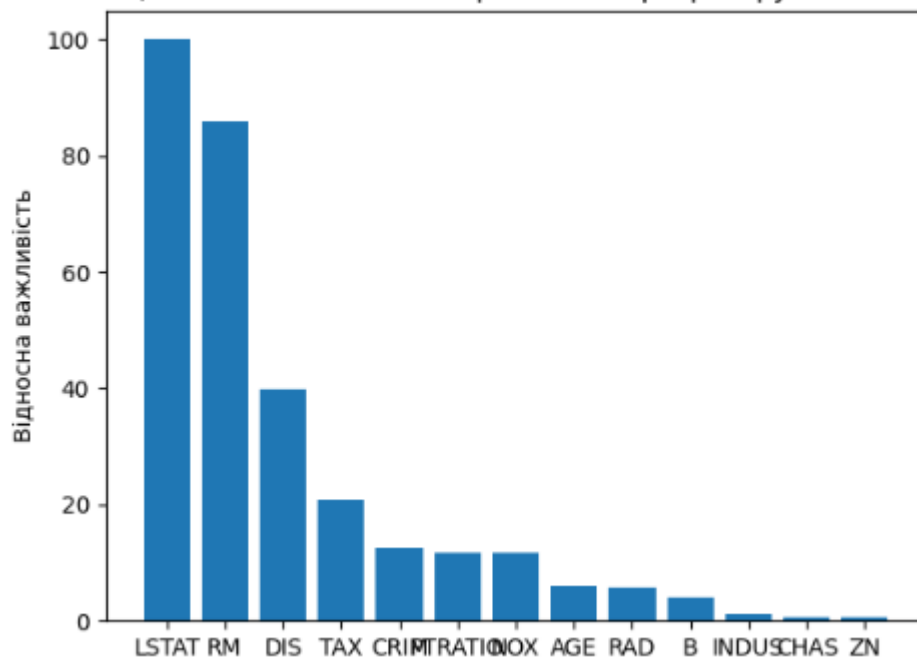
# Нормалізація значень важливості ознак
feature_importances = 100.0 * (feature_importances /
max(feature_importances))

# Сортвання та перестановка значень
index_sorted = np.flipud(np.argsort(feature_importances))

# Розміщення міток уздовж осі X
pos = np.arange(index_sorted.shape[0]) + 0.5

# Побудова стовпчастої діаграми
plt.figure()
plt.bar(pos, feature_importances[index_sorted], align='center')
plt.xticks(pos, feature_names[index_sorted])
plt.ylabel('Відносна важливість')
plt.title('Оцінка важливості з використанням регресору AdaBoost')
plt.show()
```

Оцінка важливості з використанням регресору AdaBoost



		Кияшенко А.С.			ДУ «Житомирська політехніка».22.121.07. 806– ІПЗк	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		17

Завдання 5. Прогнозування інтенсивності дорожнього руху за допомогою класифікатора на основі гранично випадкових лісів

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

# Завантажимо дані із файлу
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)

# Перетворення рядкових даних на числові
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)

# Розбиття даних на навчальний та тестовий набори
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.25, random_state=5)

# Регресор на основі гранично випадкових лісів
params = {'n_estimators': 100, 'max_depth': 4, 'random_state': 0}
regressor = ExtraTreesRegressor(**params)
regressor.fit(X_train, y_train)

# Обчислення характеристик ефективності регресора на тестових даних
y_pred = regressor.predict(X_test)
print("Mean absolute error:", round(mean_absolute_error(y_test, y_pred), 2))

# Тестування кодування на одиночному прикладі
test_datapoint = ['Saturday', '10:20', 'Atlanta', 'no']
test_datapoint_encoded = [-1] * len(test_datapoint)
count = 0
print(test_datapoint)
for i, item in enumerate(test_datapoint):
    print(test_datapoint[i])
    if item.isdigit():
        test_datapoint_encoded[i] = int(test_datapoint[i])
    else:
        test_datapoint_encoded[i] =
int(label_encoder[count].fit_transform(test_datapoint[i]))
        count = count + 1

test_datapoint_encoded = np.array(test_datapoint_encoded)

# Прогнозування результату для тестової точки даних
```

		Кияшенко А.С.			ДУ «Житомирська політехніка».22.121.07. 806– ІПЗк	Арк.
		.				18
Змн.	Арк.	№ докум.	Підпис	Дата		

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

```
C:\Users\Admin\AppData\Local\Microsoft\WindowsApps\python3.10.exe D:/LabsPoli/AI/Lab5/LR_5_task_5.py  
Mean absolute error: 7.42  
['Saturday', '10:20', 'Atlanta', 'no']
```

Посилання на Git: <https://github.com/Grum74/AI>

Висновок

Я, використовуючи спеціалізовані бібліотеки та мову програмування Python дослідив методи ансамблів у машинному навчанні.

		Кияшенко А.С.			ДУ «Житомирська політехніка».22.121.07. 806– ІПЗк	Арк.
		.				19
Змн.	Арк.	№ докум.	Підпис	Дата		