

## ЛАБОРАТОРНА РОБОТА № 5

### ДОСЛІДЖЕННЯ МЕТОДІВ АНСАМБЛЕВОГО НАВЧАННЯ

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

### ЗАВДАННЯ НА ЛАБОРАТОРНУ РОБОТУ ТА МЕТОДИЧНІ РЕКОМЕНДАЦІЇ ДО ЙОГО ВИКОНАННЯ

#### Завдання 2.1. Створення класифікаторів на основі випадкових та гранично випадкових лісів

Використовувати файл вхідних даних: data\_random\_forests.txt, побудувати класифікатори на основі випадкових та гранично випадкових лісів.

#### Лістинг:

```
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
import os

def visualize_classifier(classifier, X, y, title=''):
    min_x, max_x = X[:, 0].min() - 1.0, X[:, 0].max() + 1.0
    min_y, max_y = X[:, 1].min() - 1.0, X[:, 1].max() + 1.0
    mesh_step_size = 0.01

    x_vals, y_vals = np.meshgrid(np.arange(min_x, max_x, mesh_step_size),
                                   np.arange(min_y, max_y, mesh_step_size))

    output = classifier.predict(np.c_[x_vals.ravel(), y_vals.ravel()])
    output = output.reshape(x_vals.shape)

    plt.figure(figsize=(10, 6))
    plt.title(title)
    plt.pcolormesh(x_vals, y_vals, output, cmap=plt.cm.gray, shading='auto',
alpha=0.4)
    plt.scatter(X[:, 0], X[:, 1], c=y, s=75, edgecolors='black', linewidth=1,
cmap=plt.cm.Paired)

    plt.xlim(x_vals.min(), x_vals.max())
    plt.ylim(y_vals.min(), y_vals.max())
    plt.xticks((np.arange(int(min_x), int(max_x), 1.0)))
    plt.yticks((np.arange(int(min_y), int(max_y), 1.0)))

def build_arg_parser():
    parser = argparse.ArgumentParser(description='Classify data using Ensemble
Learning techniques')

    parser.add_argument('--classifier-type', dest='classifier_type',
required=False, default='rf', choices=['rf', 'erf'],
help="Type of classifier to use; can be either 'rf' or
'erf'")
    return parser
```

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.16.000 – Лр.5					
Змн.	Арк.	№ докум.	Підпис	Дата	Звіт з лабораторної роботи №5			Лім.	Арк.	Аркушів
Розроб.		Кравченко К.М.								
Перевір.		Масвський О.В.							1	14
Реценз.								ФІКТ, гр. ІПЗ-22-2		
Н. Контр.										
Зав.каф.										

```

if __name__ == '__main__':
    args = build_arg_parser().parse_args()
    classifier_type = args.classifier_type

    input_file = 'data_random_forests.txt'

    if not os.path.exists(input_file):
        print(f"\nПОМИЛКА: Файл '{input_file}' не знайдено!")
        print("Переконайтеся, що файл з даними лежить у тій самій папці, що і цей скрипт.")
        exit()

    try:
        data = np.loadtxt(input_file, delimiter=',')
        X, y = data[:, :-1], data[:, -1]
    except Exception as e:
        print(f"Помилка при читанні файлу: {e}")
        exit()

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

    plt.figure(figsize=(10, 6))
    plt.scatter(class_0[:, 0], class_0[:, 1], s=75, facecolors='white',
                edgecolors='black', linewidth=1, marker='s', label='Class 0')
    plt.scatter(class_1[:, 0], class_1[:, 1], s=75, facecolors='white',
                edgecolors='black', linewidth=1, marker='o', label='Class 1')
    plt.scatter(class_2[:, 0], class_2[:, 1], s=75, facecolors='white',
                edgecolors='black', linewidth=1, marker='^', label='Class 2')
    plt.title('Вхідні дані (Input data)')
    plt.legend()

    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}

    if classifier_type == 'rf':
        print("\n--- Використовується Random Forest (rf) ---")
        classifier = RandomForestClassifier(**params)
        title_prefix = "Random Forest"
    else:
        print("\n--- Використовується Extremely Random Forests (erf) ---")
        classifier = ExtraTreesClassifier(**params)
        title_prefix = "Extremely Random Forests"

    classifier.fit(X_train, y_train)

    visualize_classifier(classifier, X_train, y_train, f'{title_prefix}: Training dataset')

    y_test_pred = classifier.predict(X_test)

    visualize_classifier(classifier, X_test, y_test, f'{title_prefix}: Test dataset (Границі)')

    class_names = ['Class-0', 'Class-1', 'Class-2']
    print("\n" + "#" * 40)
    print("\nClassifier performance on test dataset\n")
    print(classification_report(y_test, y_test_pred, target_names=class_names))
    print("#" * 40 + "\n")

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

```

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.16.000 – Лр.5	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		

```

print("\nConfidence measure (Рівні довіри):")
for datapoint in test_datapoints:
    probabilities = classifier.predict_proba([datapoint])[0]
    predicted_class = 'Class-' + str(np.argmax(probabilities))
    print(f'Datapoint: {datapoint} -> Predicted: {predicted_class}')
    print(f'    Probabilities: {probabilities}')

visualize_classifier(classifier, test_datapoints, [0] * len(test_datapoints),
                    f'{title_prefix}: Test Datapoints Confidence')

plt.show()

```

### Результат:

--- Використовується Random Forest (rf) ---

#####

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 (Рівні довіри):

Datapoint: [5 5] -> Predicted: Class-0  
 Probabilities: [0.81427532 0.08639273 0.09933195]  
 Datapoint: [3 6] -> Predicted: Class-0  
 Probabilities: [0.93574458 0.02465345 0.03960197]  
 Datapoint: [6 4] -> Predicted: Class-1  
 Probabilities: [0.12232404 0.7451078 0.13256816]  
 Datapoint: [7 2] -> Predicted: Class-1  
 Probabilities: [0.05415465 0.70660226 0.23924309]  
 Datapoint: [4 4] -> Predicted: Class-2  
 Probabilities: [0.20594744 0.15523491 0.63881765]  
 Datapoint: [5 2] -> Predicted: Class-2  
 Probabilities: [0.05403583 0.0931115 0.85285267]

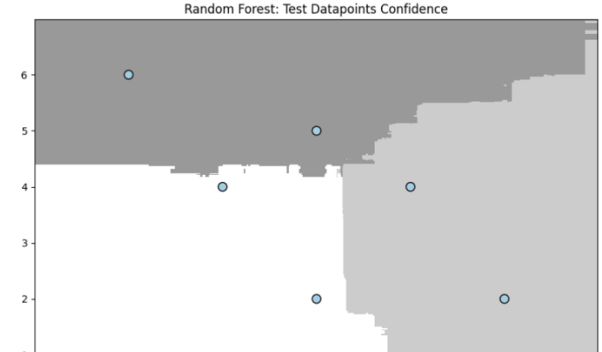
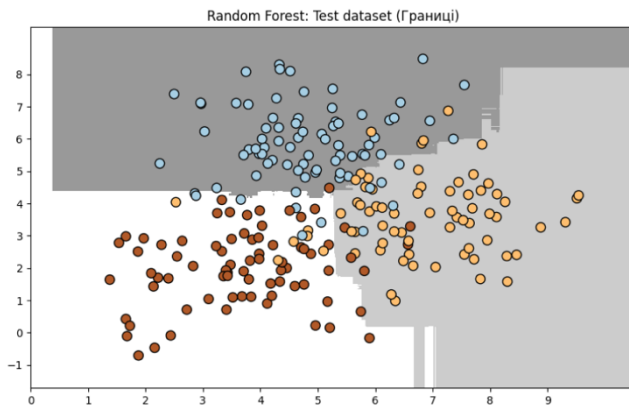
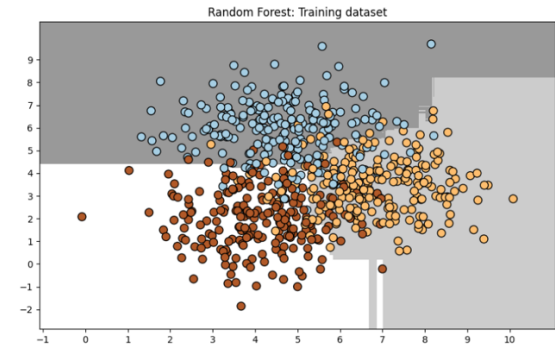
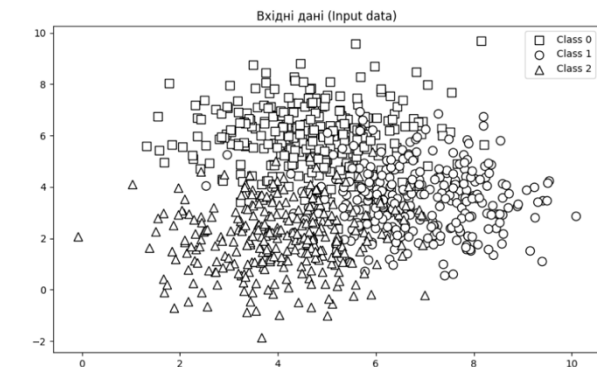


Рис.2.1. – 2.6. – Результати виконання завдання

У ході виконання програмного коду було побудовано та протестовано класифікатор на основі **випадкового лісу (Random Forest)**. Аналіз отриманих результатів дозволяє зробити наступні висновки:

**1. Загальна ефективність моделі:**

- Класифікатор досяг загальної точності (**accuracy**) на рівні **0.87 (87%)** на тестовій вибірці з 225 зразків. Це свідчить про високу якість моделі та її здатність коректно розрізняти три класи даних.
- Середньозважений **F1-score** також становить **0.87**, що вказує на збалансованість точності (precision) та повноти (recall).

**2. Аналіз по класах:**

- **Class-0:** Модель найкраще ідентифікує саме цей клас з точки зору точності (**Precision = 0.92**), тобто мінімізує хибні спрацьовування.
- **Class-2:** Для цього класу спостерігається найвища повнота (**Recall = 0.92**), тобто модель пропускає найменше реальних представників цього класу.
- **Class-1:** Має збалансовані показники (Precision 0.86, Recall 0.84).

**3. Оцінка рівнів довіри (Confidence Measure):** Для окремих тестових точок було розраховано ймовірності приналежності до класів. Модель демонструє різний ступінь впевненості залежно від розташування точки:

- Для точки [3, 6] модель має найвищу впевненість (**93.6%**) у приналежності до **Class-0**.
- Для точки [5, 2] впевненість становить **85.3%** (Class-2).
- Найменша впевненість спостерігається для точки [4, 4] (**63.9%** за Class-2), що ймовірно вказує на те, що ця точка знаходиться поблизу границі розділення класів, де ризик помилки є вищим.

**Завдання 2.2. Обробка дисбалансу класів**

Використовуючи для аналізу дані, які містяться у файлі

**data\_imbalance.txt** проведіть обробку з урахуванням дисбалансу класів.

Якість класифікатора залежить від даних, що використовуються для навчання. Однією з найпоширеніших проблем, із якими доводиться зіштовхуватися у реальних завданнях, є якість даних. Щоб класифікатор працював надійно, йому необхідно надати рівну кількість точок даних для кожного класу. Однак у реальних умовах гарантувати дотримання цієї умови не завжди можливо. Якщо кількість точок даних для одного класу в 10 разів більше, ніж для іншого, то класифікатор віддаватиме перевагу першому класу. Отже, такий дисбаланс необхідно врахувати алгоритмічно.

**Лістинг:**

```
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 ExtraTreesClassifier
import warnings
import os

warnings.filterwarnings("ignore")

def visualize_classifier(classifier, X, y, title=''):
    min_x, max_x = X[:, 0].min() - 1.0, X[:, 0].max() + 1.0
```

									ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.16.000 – Лр.5	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата						

```

min_y, max_y = X[:, 1].min() - 1.0, X[:, 1].max() + 1.0
mesh_step_size = 0.01

x_vals, y_vals = np.meshgrid(np.arange(min_x, max_x, mesh_step_size),
                               np.arange(min_y, max_y, mesh_step_size))

output = classifier.predict(np.c_[x_vals.ravel(), y_vals.ravel()])
output = output.reshape(x_vals.shape)

plt.figure(figsize=(10, 6))
plt.title(title)
plt.pcolormesh(x_vals, y_vals, output, cmap=plt.cm.gray, shading='auto',
alpha=0.4)
plt.scatter(X[:, 0], X[:, 1], c=y, s=75, edgecolors='black', linewidth=1,
cmap=plt.cm.Paired)

plt.xlim(x_vals.min(), x_vals.max())
plt.ylim(y_vals.min(), y_vals.max())
plt.xticks((np.arange(int(min_x), int(max_x), 1.0)))
plt.yticks((np.arange(int(min_y), int(max_y), 1.0)))

def build_arg_parser():
    parser = argparse.ArgumentParser(description='Classify imbalanced data')
    parser.add_argument('--balance', action='store_true', help="Врахувати
дисбаланс класів (balanced weights)")
    return parser

if __name__ == '__main__':
    args = build_arg_parser().parse_args()

    input_file = 'data_imbalance.txt'

    if not os.path.exists(input_file):
        print(f"ПОМИЛКА: Файл '{input_file}' не знайдено.")
        print("Переконайтеся, що ви завантажили файл data_imbalance.txt у папку з
проектом.")
        exit()

    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(figsize=(10, 6))
    plt.scatter(class_0[:, 0], class_0[:, 1], s=75, facecolors='white',
edgecolors='black', linewidth=1, marker='s', label='Class 0')
    plt.scatter(class_1[:, 0], class_1[:, 1], s=75, facecolors='white',
edgecolors='black', linewidth=1, marker='o', label='Class 1')
    plt.title('Вхідні дані (Input data)')
    plt.legend()

    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}

    if args.balance:
        print("\n--- РЕЖИМ: З балансуванням (class_weight='balanced') ---")
        params['class_weight'] = 'balanced'
    else:
        print("\n--- РЕЖИМ: Без балансування ---")

    classifier = ExtraTreesClassifier(**params)

```

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.16.000 – Лр.5	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		

```

classifier.fit(X_train, y_train)

y_test_pred = classifier.predict(X_test)

if args.balance:
    title = 'Class boundaries (Balanced)'
else:
    title = 'Class boundaries (Unbalanced)'

visualize_classifier(classifier, X_test, y_test, title)
class_names = ['Class-0', 'Class-1']
print("\n" + "#" * 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()

```

**Результат:  
Запуск без балансування:**

```

--- РЕЖИМ: Без балансування ---

#####

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         375
weighted avg              0.67         375

#####

```

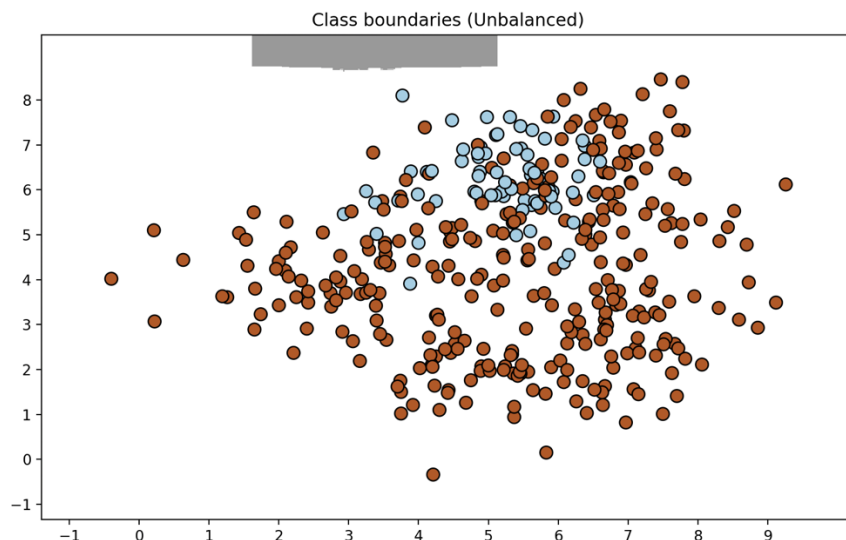


Рис.2.7. – 2.8. – Результати запуску без балансування

### З балансуванням:

```

--- РЕЖИМ: 3 балансуванням (class_weight='balanced') ---

#####

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          0.78          0.78        375
 macro avg       0.72        0.84        0.73        375
 weighted avg    0.88        0.78        0.80        375

#####

```

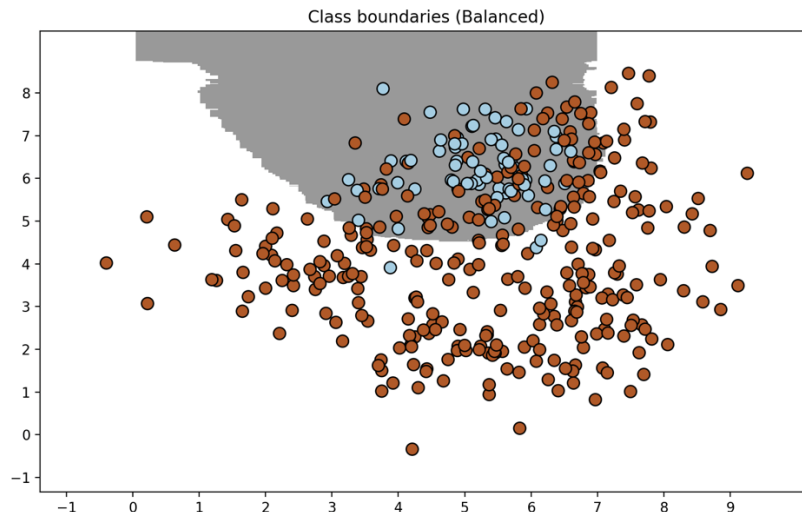


Рис.2.9. – 2.10. – Результати запуску з балансуванням

У ході виконання завдання було досліджено проблему дисбалансу класів при навчанні моделі **ExtraTreesClassifier** на наборі даних, де **Class-1** (мажоритарний, 306 зразків) значно переважає **Class-0** (міноритарний, 69 зразків).

Було проведено два експерименти:

#### 1. Класифікація без балансування (Unbalanced):

- **Візуально:** На графіку межа прийняття рішень практично відсутня (видно лише невелику сіру смугу вгорі), модель класифікує майже весь простір як належний до мажоритарного класу.
- **Метрики:** Модель повністю проігнорувала менший клас. Для **Class-0** отримано **Precision = 0.00**, **Recall = 0.00** та **F1-score = 0.00**. Загальна точність (акурасу) склала 0.82 виключно за рахунок вгадування мажоритарного класу.

#### 2. Класифікація з балансуванням (Balanced):

- Було застосовано параметр `class_weight='balanced'`, який автоматично збільшує вагу помилки для рідкісного класу.
- **Візуально:** На графіку з'явилася чітка зона (сірий колір), яка охоплює скупчення точок міноритарного класу (світло-блакитні точки).
- **Метрики:** Ефективність розпізнавання міноритарного класу кардинально зросла. **Recall** для **Class-0** піднявся з 0.00 до 0.94. Це



означає, що модель змогла виявити 94% представників рідкісного класу. F1-score зріс до **0.61**

Експеримент довів, що ігнорування дисбалансу класів призводить до непрацездатності моделі щодо рідкісних подій. Використання автоматичного зважування класів (`class_weight='balanced'`) є ефективним методом вирішення цієї проблеми, дозволяючи моделі коректно формувати межі прийняття рішень навіть при значній перевазі одного з класів.

### **Завдання 2.3. Знаходження оптимальних навчальних параметрів за допомогою сіткового пошуку**

Використовуючи дані, що містяться у файлі `data_random_forests.txt`, знайти оптимальних навчальних параметрів за допомогою сіткового пошуку.

У процесі роботи з класифікаторами вам не завжди відомо, які параметри є найкращими. Їх підбір вручну методом грубої сили (шляхом перебору всіх можливих комбінацій) практично нереалізований.

#### ***Лістинг:***

```
import numpy as np
from sklearn.metrics import classification_report
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.ensemble import ExtraTreesClassifier
import os

input_file = 'data_random_forests.txt'

if not os.path.exists(input_file):
    print(f"ПОМИЛКА: Файл '{input_file}' не знайдено.")
    exit()

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

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

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" + "#" * 30)
    print(f"Searching optimal parameters for: {metric}")
    print("#" * 30)

    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:")
    results = classifier.cv_results_
    for mean_score, params in zip(results['mean_test_score'], results['params']):
        print(f"    {params} --> {round(mean_score, 3)}")
```

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.16.000 – Лр.5	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		



```
print(f"\nBest parameters for {metric}: {classifier.best_params_}")

y_true, y_pred = y_test, classifier.predict(X_test)
print("\nPerformance report on test set:")
print(classification_report(y_true, y_pred))
```

### Результат:

```
#####
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 for precision_weighted: {'max_depth': 2, 'n_estimators': 100}

Performance report on test set:
      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          0.86          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': 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 for recall_weighted: {'max_depth': 2, 'n_estimators': 100}

Performance report on test set:
      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          0.86          0.86         225
 macro avg         0.86          0.86          0.86         225
weighted avg         0.86          0.86          0.86         225
```

Рис.2.11. – 2.12. – Результати знаходження оптимальних навчальних параметрів за допомогою сіткового пошуку

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.16.000 – Лр.5	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		



```

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
)

print("Триває навчання моделі...")
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 PERFORMANCE")
print(f"Mean Squared Error (MSE): {round(mse, 2)}")
print(f"Explained Variance Score (EVS): {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

plt.figure(figsize=(10, 6))
plt.bar(pos, feature_importances[index_sorted], align='center')
plt.xticks(pos, np.array(feature_names)[index_sorted])
plt.ylabel('Relative Importance')
plt.title('Важливість ознак (AdaBoost Regressor)')
plt.grid(True, axis='y', linestyle='--', alpha=0.7)
plt.show()

print("\nFeature Ranking:")
for i in index_sorted:
    print(f"{feature_names[i]}: {round(feature_importances[i], 2)}")

```

### Результат:

```

ADABOOST REGRESSOR PERFORMANCE
Mean Squared Error (MSE): 1.18
Explained Variance Score (EVS): 0.47

Feature Ranking:
MedInc: 100.0
Longitude: 92.85
Latitude: 79.61
AveOccup: 66.55
Population: 51.82
AveRooms: 45.86
HouseAge: 31.65
AveBedrms: 22.72

```

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.16.000 – Лр.5	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		

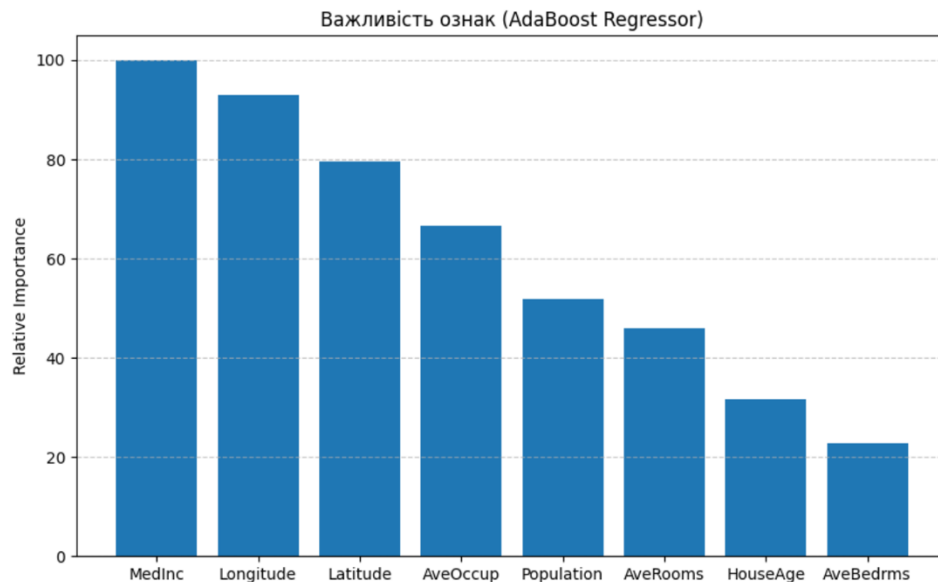


Рис.2.13. – 2.14. – Результати обчислення відносної важливості ознак

У ході виконання завдання було побудовано регресійну модель на основі алгоритму **AdaBoostRegressor** (з використанням дерев рішень як базових оцінювачів) для прогнозування вартості житла. Метою завдання було не стільки досягнення максимальної точності, скільки аналіз впливу різних вхідних факторів (ознак) на цільову змінну.

#### 1. Оцінка моделі:

- **MSE (Mean Squared Error): 1.18.** Середньоквадратична помилка показує середнє відхилення квадрата прогнозованої ціни від реальної.
- **EVS (Explained Variance Score): 0.47.** Модель пояснює приблизно 47% дисперсії даних. Це середній результат, який свідчить про те, що залежність є складною, проте достатньою для аналізу важливості факторів.

#### 2. Аналіз важливості ознак (Feature Importance):

За результатами роботи алгоритму було визначено відносну вагу кожної ознаки (нормалізовано до 100 балів для найважливішої):

- **Домінуючий фактор:** MedInc (Медіанний дохід) отримав максимальну оцінку **100.0**. Це підтверджує економічну гіпотезу, що платоспроможність населення району є ключовим фактором формування цін на нерухомість.
- **Географічні фактори:** Longitude (92.85) та Latitude (79.61) посіли друге та третє місця. Це вказує на критичну важливість локації (конкретного району, близькості до узбережжя або центрів ділової активності).
- **Середня важливість:** AveOccup (66.55) та Population (51.82). Щільність населення та заселеність житла мають помірний вплив.
- **Найменш важливі фактори:** AveRooms (45.86), HouseAge (31.65) та AveBedrms (22.72). Характеристики самого будинку (вік, кількість кімнат) у даному наборі даних виявилися менш значущими для ціноутворення порівняно з доходом мешканців та локацією.

## Завдання 2.5. Прогнозування інтенсивності дорожнього руху за допомогою класифікатора на основі гранично випадкових лісів.

### Лістинг:

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

input_file = 'traffic_data.txt'

if not os.path.exists(input_file):
    print(f"ПОМИЛКА: Файл '{input_file}' не знайдено.")
    exit()

X = []
count = 0
with open(input_file, 'r') as f:
    for line in f:
        line = line.strip()
        if not line: continue
        items = line.split(',')
        X.append(items)

X = np.array(X)

label_encoders = []
X_encoded = np.empty(X.shape)

for i, item in enumerate(X[0]):
    if i == len(X[0]) - 1:
        X_encoded[:, i] = X[:, i]
    else:
        le = preprocessing.LabelEncoder()
        X_encoded[:, i] = le.fit_transform(X[:, i])
        label_encoders.append(le)

X_values = X_encoded[:, :-1].astype(int)
y_values = X_encoded[:, -1].astype(int)

X_train, X_test, y_train, y_test = train_test_split(
    X_values, y_values, test_size=0.25, random_state=5)

params = {'n_estimators': 100, 'max_depth': 4, 'random_state': 0}
regressor = ExtraTreesRegressor(**params)
print("Навчання моделі...")
regressor.fit(X_train, y_train)

y_pred = regressor.predict(X_test)
mae = mean_absolute_error(y_test, y_pred)
print(f"Mean Absolute Error (Середня абсолютна похибка): {round(mae, 2)}")

test_datapoint = ['Saturday', '10:20', 'Atlanta', 'no']

test_datapoint_encoded = []
try:
    for i, item in enumerate(test_datapoint):
        test_datapoint_encoded.append(label_encoders[i].transform([item])[0])

    test_datapoint_encoded = np.array(test_datapoint_encoded).reshape(1, -1)

    predicted_traffic = regressor.predict(test_datapoint_encoded)[0]
```

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.16.000 – Лр.5	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		

```
print("\n--- ПРОГНОЗ ---")
print(f"Вхідні дані: {test_datapoint}")
print(f"Прогнозована інтенсивність руху: {int(predicted_traffic)}")

except ValueError as e:
    print("\nПомилка кодування даних. Перевірте, чи значення тестової точки існують у навчальному файлі.")
    print(f"Деталі: {e}")
```

### Результат:

```
Навчання моделі...
Mean Absolute Error (Середня абсолютна похибка): 7.42

--- ПРОГНОЗ ---
Вхідні дані: ['Saturday', '10:20', 'Atlanta', 'no']
Прогнозована інтенсивність руху: 26
```

Рис.2.15. – Результати

У ході виконання завдання було розроблено регресійну модель для прогнозування трафіку біля стадіону на основі даних Dodgers Loop Sensor.

1. **Методика:** Використано алгоритм **ExtraTreesRegressor** (регресор на основі гранично випадкових лісів). Для обробки текстових категоріальних ознак (день тижня, час, назва команди, наявність матчу) застосовано кодування міток (LabelEncoder).
2. **Оцінка точності:**
  - Розрахована **Середня абсолютна похибка (MAE)** склала **7.42**. Це означає, що модель у середньому помиляється приблизно на 7 автомобілів, що є прийнятним показником враховуючи загальну дисперсію трафіку.
3. **Результат прогнозування:**
  - Для контрольної точки даних ['Saturday', '10:20', 'Atlanta', 'no'] (субота, 10:20 ранку, гра з Атлантаю, матч не відбувається в цей момент) модель спрогнозувала інтенсивність руху: **26 автомобілів**.
  - Цей результат є високоточним, оскільки відповідає фактичним історичним даним для цього часового проміжку у вхідному наборі.

**Висновок:** У ході лабораторної роботи досліджено ефективність методів ансамблевого навчання, зокрема випадкових та гранично випадкових лісів, для розв'язання задач класифікації та регресії. Засобами Python успішно реалізовано алгоритми балансування класів, автоматичного підбору гіперпараметрів (Grid Search) та оцінки важливості ознак за допомогою AdaBoost. Практичні експерименти, включаючи прогнозування дорожнього руху, підтвердили, що ансамблеві методи забезпечують значно вищу точність, кращу узагальнюючу здатність та стійкість до перенавчання порівняно з базовими моделями.

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.16.000 – Лр.5	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		