

Лабораторна робота №2

ПОРІВНЯННЯ МЕТОДІВ КЛАСИФІКАЦІЇ ДАНИХ

Мета: використовуючи спеціалізовані бібліотеки та мову програмування Python дослідити різні методи класифікації даних та навчитися їх порівнювати.

Хід роботи

Репозиторій GITHUB: https://github.com/AnatoliiYarmolenko/SMI_1S4C

Завдання 2.1. Класифікація за допомогою машин опорних векторів (SVM)

Випишіть у звіт всі 14 ознак з набору даних – їх назви та що вони позначають та вид (числові чи категоріальні).

Проаналізувавши дані встановлюємо наступні ознаки:

Назва ознаки	Позначення	Тип
age	вік	числовий
workclass	тип зайнятості	категоріальний
fnlwgt	фінальна вага (кількість людей, яких представляє запис)	числовий
education	освіта	категоріальний
education-num	тривалість освіти в роках	числовий
marital-status	сімейний стан	категоріальний
occupation	професія	категоріальний
relationship	родинні стосунки	категоріальний
race	раса	категоріальний
sex	стать	категоріальний
capital-gain	приріст капіталу	числовий
capital-loss	втрата капіталу	числовий
hours-per-week	кількість робочих годин на тиждень	числовий
native-country	країна походження	категоріальний

Табл. 1 Таблиця ознак

Обчисліть значення інших показників якості класифікації

(акуратність, повнота, точність) та разом з F1 занесіть їх у звіт.

Збережіть код робочої програми під назвою LR_2_task_1.py					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.31.000 – Лр2			
Змн.	Арк.	№ докум.	Підпис	Дата				
Розроб.		Ярмоленко А.М.			Звіт з лабораторної роботи №2	Літ.	Арк.	Аркушів
Перевір.		Маєвський О.В.					1	27
Реценз.						ФІКТ, гр. ІПЗ-22-4		
Н. Контр.								
Зав.каф.		Вакалюк Т.А.						

Код програми занесіть у звіт. Зробіть висновок до якого класу належить тестова точка.

Обчисливши показники отримаємо наступний результат:

```
Метрики якості для тестового набору
Ассурасу (Акуратність): 79.56%
Precision (Точність): 79.26%
Recall (Повнота): 79.56%
F1 score: 76.01%
```

Рис. 1 Результат обчислень

Тут ми бачимо: Акуратність (Частка правильних прогнозів від загальної кількості прогнозів); Точність (Здатність класифікатора не маркувати негативний зразок як позитивний); Повнота (Здатність класифікатора знаходити всі позитивні зразки); F1-міра (Середнє гармонійне між Точністю та Повнотою.)

Значення Ассурасу (Акуратність) ~79.5% та F1-міра ~76% є достатньо реалістичними для цього набору даних. Отже наша модель працює коректно.

Тестова точка також демонструє правильний результат:

```
Тестова точка: ['37', 'Private', '215646', '0', '0', '40', 'United-States']
Predicted income class: <=50K
```

Рис. 2 Перевірка точки

Згідно результату класифікації прогнозований дохід <=50 000, що цілком задовільняє умовам та ознакам цієї точки.

Загальний лістинг програми:

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn import preprocessing
from sklearn.svm import LinearSVC
from sklearn.multiclass import OneVsOneClassifier
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

# Вхідний файл
input_file = 'income_data.txt'

# Читання даних
X = []
```

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.2.000 – Лр.2	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		2

```

y = []
count_class1 = 0
count_class2 = 0
max_datapoints = 25000

with open(input_file, 'r') as f:
    for line in f.readlines():
        if count_class1 >= max_datapoints and count_class2 >= max_datapoints:
            break
        if '?' in line:
            continue

        data = line.strip().split(',')
        if data[-1] == '<=50K' and count_class1 < max_datapoints:
            X.append(data[:-1])
            y.append(data[-1])
            count_class1 += 1
        elif data[-1] == '>50K' and count_class2 < max_datapoints:
            X.append(data[:-1])
            y.append(data[-1])
            count_class2 += 1

# Перетворення у numpy-масиви
X = np.array(X)
y = np.array(y)

# Кодування рядкових змінних
label_encoders = []
X_encoded = np.empty(X.shape)

for i, item in enumerate(X[0]):
    if item.replace('.', '', 1).isdigit(): # якщо число (включно з десятковими)
        X_encoded[:, i] = X[:, i].astype(float)
    else:
        le = preprocessing.LabelEncoder()
        X_encoded[:, i] = le.fit_transform(X[:, i])
        label_encoders.append(le)

# Кодування вихідних міток
y_le = preprocessing.LabelEncoder()
y = y_le.fit_transform(y)

X = X_encoded.astype(float)

# Поділ на тренувальні і тестові дані
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=5)

# Створення SVM-класифікатора
classifier = OneVsOneClassifier(LinearSVC(random_state=0, max_iter=10000))
classifier.fit(X_train, y_train)

```

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.2.000 – Лр.2	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		3

```

# --- Передбачення для нової точки ---
input_data = ['37', 'Private', '215646', 'HS-grad', '9', 'Never-married',
              'Handlers-cleaners', 'Not-in-family', 'White', 'Male',
              '0', '0', '40', 'United-States']

# Кодування тестової точки
input_data_encoded = []
count = 0
for i, item in enumerate(input_data):
    if item.replace('.', '', 1).isdigit():
        input_data_encoded.append(float(item))
    else:
        le = label_encoders[count]
        if item in le.classes_:
            input_data_encoded.append(le.transform([item])[0])
        else:
            input_data_encoded.append(-1) # якщо нове значення не зустрічалося
            count += 1

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

y_test_pred = classifier.predict(X_test)

print("Метрики якості для тестового набору")
# Акуратність (Accuracy)
acc = accuracy_score(y_test, y_test_pred)
print("Accuracy (Акуратність): " + str(round(100 * acc, 2)) + "%")
# Точність (Precision)
prec = precision_score(y_test, y_test_pred, average='weighted')
print("Precision (Точність): " + str(round(100 * prec, 2)) + "%")
# Повнота (Recall)
rec = recall_score(y_test, y_test_pred, average='weighted')
print("Recall (Повнота): " + str(round(100 * rec, 2)) + "%")

f1 = cross_val_score(classifier, X, y, scoring='f1_weighted', cv=3)
print("F1 score: " + str(round(100 * f1.mean(), 2)) + "%")

# Передбачення
predicted_class = classifier.predict(input_data_encoded)
input_data = ['37', 'Private', '215646', 'HS-grad', '9', 'Never-married',
              'Handlers-cleaners', 'Not-in-family', 'White', 'Male', '0', '0', '40',
              'United-States']

print(f"Тестова точка: {input_data}")

print("Predicted income class:", y_le.inverse_transform(predicted_class)[0])

```

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.2.000 – Лр.2	Арк.
						4
Змн.	Арк.	№ докум.	Підпис	Дата		

Завдання 2.2. Порівняння якості класифікаторів SVM з нелінійними ядрами

Використовуючи набір даних та код з попереднього завдання створіть та дослідіть нелінійні класифікатори SVM.

з поліноміальним ядром;

з гаусовим ядром;

з сигмоїдальним ядром.

Для кожного виду класифікатора отримайте та запишіть у звіт показники якості алгоритму класифікації.

Для цього завдання потрібно змінити один рядок коду та додати бібліотеку до імпорту. Однак ступень поліному degree прийшлося зменшити до 2, щоб час розрахунку вкладався в межі розумного.

```
#classifier = OneVsOneClassifier(LinearSVC(random_state=0, max_iter=10000))
classifier = SVC(kernel='poly', degree=2, random_state=0)
```

```
Метрики якості для тестового набору
Акуратність (Accuracy): 77.39%
Точність (Precision): 81.11%
Повнота (Recall): 77.39%
F1 score: 70.68%
Тестова точка: ['37', 'Private', '215646', 'HS-grad',
, 'White', 'Male', '0', '0', '40', 'United-States']
Predicted income class: <=50K
```

Рис. 3 Результат обчислень

Загальний лістинг програми через поліноміальне ядро:

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn import preprocessing
from sklearn.svm import SVC, LinearSVC
from sklearn.multiclass import OneVsOneClassifier
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

# Вхідний файл
input_file = 'income_data.txt'
```

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.2.000 – Лр.2	Арк.
						5
Змн.	Арк.	№ докум.	Підпис	Дата		

```

# Читання даних
X = []
y = []
count_class1 = 0
count_class2 = 0
max_datapoints = 25000

with open(input_file, 'r') as f:
    for line in f.readlines():
        if count_class1 >= max_datapoints and count_class2 >= max_datapoints:
            break
        if '?' in line:
            continue

        data = line.strip().split(',')
        if data[-1] == '<=50K' and count_class1 < max_datapoints:
            X.append(data[:-1])
            y.append(data[-1])
            count_class1 += 1
        elif data[-1] == '>50K' and count_class2 < max_datapoints:
            X.append(data[:-1])
            y.append(data[-1])
            count_class2 += 1

# Перетворення у numpy-масиви
X = np.array(X)
y = np.array(y)

# Кодування рядкових змінних
label_encoders = []
X_encoded = np.empty(X.shape)

for i, item in enumerate(X[0]):
    if item.replace('.', '', 1).isdigit(): # якщо число (включно з десятковими)
        X_encoded[:, i] = X[:, i].astype(float)
    else:
        le = preprocessing.LabelEncoder()
        X_encoded[:, i] = le.fit_transform(X[:, i])
        label_encoders.append(le)

# Кодування вихідних міток
y_le = preprocessing.LabelEncoder()
y = y_le.fit_transform(y)

X = X_encoded.astype(float)

# Поділ на тренувальні і тестові дані
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=5)

```

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.2.000 – Лр.2	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		6

```

print(f"Початок обчислень...\n")

# Створення SVM-класифікатора
#classifier = OneVsOneClassifier(LinearSVC(random_state=0, max_iter=10000))
classifier = SVC(kernel='poly', degree=2, random_state=0)

classifier.fit(X_train, y_train)

# --- Передбачення для нової точки ---
input_data = ['37', 'Private', '215646', 'HS-grad', '9', 'Never-married',
              'Handlers-cleaners', 'Not-in-family', 'White', 'Male',
              '0', '0', '40', 'United-States']

# Кодування тестової точки
input_data_encoded = []
count = 0
for i, item in enumerate(input_data):
    if item.replace('.', '', 1).isdigit():
        input_data_encoded.append(float(item))
    else:
        le = label_encoders[count]
        if item in le.classes_:
            input_data_encoded.append(le.transform([item])[0])
        else:
            input_data_encoded.append(-1) # якщо нове значення не зустрічалося
        count += 1

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

y_test_pred = classifier.predict(X_test)

print("Метрики якості для тестового набору")
# Акуратність (Accuracy)
acc = accuracy_score(y_test, y_test_pred)
print("Accuracy (Акуратність): " + str(round(100 * acc, 2)) + "%")
# Точність (Precision)
prec = precision_score(y_test, y_test_pred, average='weighted')
print("Precision (Точність): " + str(round(100 * prec, 2)) + "%")
# Повнота (Recall)
rec = recall_score(y_test, y_test_pred, average='weighted')
print("Recall (Повнота): " + str(round(100 * rec, 2)) + "%")

f1 = cross_val_score(classifier, X, y, scoring='f1_weighted', cv=3)
print("F1 score: " + str(round(100 * f1.mean(), 2)) + "%")

# Передбачення
predicted_class = classifier.predict(input_data_encoded)
input_data = ['37', 'Private', '215646', 'HS-grad', '9', 'Never-married',
              'Handlers-cleaners', 'Not-in-family', 'White', 'Male', '0', '0', '40',
              'United-States']

```

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.2.000 – Лр.2	Арк.
						7
Змн.	Арк.	№ докум.	Підпис	Дата		

```
print(f"Тестова точка: {input_data}")

print("Predicted income class:", y_le.inverse_transform(predicted_class)[0])
```

Настпний крок аналогічний:

```
#classifier = SVC(kernel='poly', degree=2, random_state=0)
classifier = SVC(kernel='rbf', degree=2, random_state=0)
```

```
Метрики якості для тестового набору
Ассурасу (Акуратність): 78.19%
Precision (Точність): 82.82%
Recall (Повнота): 78.19%
F1 score: 71.95%
Тестова точка: ['37', 'Private', '215646', 'HS-
, 'White', 'Male', '0', '0', '40', 'United-Stat
Predicted income class: <=50K
```

Рис. 4 Результат обчислень

Загальний лістинг програми через гаусове ядро:

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn import preprocessing
from sklearn.svm import SVC, LinearSVC
from sklearn.multiclass import OneVsOneClassifier
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

# Вхідний файл
input_file = 'income_data.txt'

# Читання даних
X = []
y = []
count_class1 = 0
count_class2 = 0
max_datapoints = 25000

with open(input_file, 'r') as f:
    for line in f.readlines():
        if count_class1 >= max_datapoints and count_class2 >= max_datapoints:
            break
        if '?' in line:
            continue

        data = line.strip().split(' ', ')
```

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.2.000 – Лр.2	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		8


```

if data[-1] == '<=50K' and count_class1 < max_datapoints:
    X.append(data[:-1])
    y.append(data[-1])
    count_class1 += 1
elif data[-1] == '>50K' and count_class2 < max_datapoints:
    X.append(data[:-1])
    y.append(data[-1])
    count_class2 += 1

# Перетворення у numpy-масиви
X = np.array(X)
y = np.array(y)

# Кодування рядкових змінних
label_encoders = []
X_encoded = np.empty(X.shape)

for i, item in enumerate(X[0]):
    if item.replace('.', '', 1).isdigit(): # якщо число (включно з десятковими)
        X_encoded[:, i] = X[:, i].astype(float)
    else:
        le = preprocessing.LabelEncoder()
        X_encoded[:, i] = le.fit_transform(X[:, i])
        label_encoders.append(le)

# Кодування вихідних міток
y_le = preprocessing.LabelEncoder()
y = y_le.fit_transform(y)

X = X_encoded.astype(float)

# Поділ на тренувальні і тестові дані
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=5)

print(f"Початок обчислень...\n")

# Створення SVM-класифікатора
#classifier = OneVsOneClassifier(LinearSVC(random_state=0, max_iter=10000))
#classifier = SVC(kernel='poly', degree=2, random_state=0)
classifier = SVC(kernel='rbf', degree=2, random_state=0)

classifier.fit(X_train, y_train)

# --- Передбачення для нової точки ---
input_data = ['37', 'Private', '215646', 'HS-grad', '9', 'Never-married',
              'Handlers-cleaners', 'Not-in-family', 'White', 'Male',
              '0', '0', '40', 'United-States']

# Кодування тестової точки
input_data_encoded = []

```

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.2.000 – Лр.2	Арк.
						9
Змн.	Арк.	№ докум.	Підпис	Дата		

```

count = 0
for i, item in enumerate(input_data):
    if item.replace('.', '', 1).isdigit():
        input_data_encoded.append(float(item))
    else:
        le = label_encoders[count]
        if item in le.classes_:
            input_data_encoded.append(le.transform([item])[0])
        else:
            input_data_encoded.append(-1) # якщо нове значення не зустрічалося
        count += 1

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

y_test_pred = classifier.predict(X_test)

print("Метрики якості для тестового набору")
# Акуратність (Accuracy)
acc = accuracy_score(y_test, y_test_pred)
print("Accuracy (Акуратність): " + str(round(100 * acc, 2)) + "%")
# Точність (Precision)
prec = precision_score(y_test, y_test_pred, average='weighted')
print("Precision (Точність): " + str(round(100 * prec, 2)) + "%")
# Повнота (Recall)
rec = recall_score(y_test, y_test_pred, average='weighted')
print("Recall (Повнота): " + str(round(100 * rec, 2)) + "%")

f1 = cross_val_score(classifier, X, y, scoring='f1_weighted', cv=3)
print("F1 score: " + str(round(100 * f1.mean(), 2)) + "%")

# Передбачення
predicted_class = classifier.predict(input_data_encoded)
input_data = ['37', 'Private', '215646', 'HS-grad', '9', 'Never-married',
              'Handlers-cleaners', 'Not-in-family', 'White', 'Male', '0', '0', '40',
              'United-States']

print(f"Тестова точка: {input_data}")

print("Predicted income class:", y_le.inverse_transform(predicted_class)[0])

```

Аналогічно:

```

#classifier = SVC(kernel='rbf', degree=2, random_state=0)
classifier = SVC(kernel='sigmoid', degree=2, random_state=0)

```

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.2.000 – Лр.2	Арк.
						10
Змн.	Арк.	№ докум.	Підпис	Дата		

```

Метрики якості для тестового набору
Accuracy (Акуратність): 60.47%
Precision (Точність): 60.64%
Recall (Повнота): 60.47%
F1 score: 63.77%
Тестова точка: ['37', 'Private', '215646', 'HS-grad'
, 'White', 'Male', '0', '0', '40', 'United-States']
Predicted income class: <=50K

```

Рис. 5 Результат обчислень

Загальний лістинг програми через сигмоїдальне ядро:

```

import numpy as np
import matplotlib.pyplot as plt
from sklearn import preprocessing
from sklearn.svm import SVC, LinearSVC
from sklearn.multiclass import OneVsOneClassifier
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

# Вхідний файл
input_file = 'income_data.txt'

# Читання даних
X = []
y = []
count_class1 = 0
count_class2 = 0
max_datapoints = 25000

with open(input_file, 'r') as f:
    for line in f.readlines():
        if count_class1 >= max_datapoints and count_class2 >= max_datapoints:
            break
        if '?' in line:
            continue

        data = line.strip().split(',')
        if data[-1] == '<=50K' and count_class1 < max_datapoints:
            X.append(data[:-1])
            y.append(data[-1])
            count_class1 += 1
        elif data[-1] == '>50K' and count_class2 < max_datapoints:
            X.append(data[:-1])
            y.append(data[-1])
            count_class2 += 1

# Перетворення у numpy-масиви
X = np.array(X)

```

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.2.000 – Лр.2	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		11

```

y = np.array(y)

# Кодування рядкових змінних
label_encoders = []
X_encoded = np.empty(X.shape)

for i, item in enumerate(X[0]):
    if item.replace('.', '', 1).isdigit(): # якщо число (включно з десятковими)
        X_encoded[:, i] = X[:, i].astype(float)
    else:
        le = preprocessing.LabelEncoder()
        X_encoded[:, i] = le.fit_transform(X[:, i])
        label_encoders.append(le)

# Кодування вихідних міток
y_le = preprocessing.LabelEncoder()
y = y_le.fit_transform(y)

X = X_encoded.astype(float)

# Поділ на тренувальні і тестові дані
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=5)

print(f"Початок обчислень...\n")

# Створення SVM-класифікатора
#classifier = OneVsOneClassifier(LinearSVC(random_state=0, max_iter=10000))
#classifier = SVC(kernel='poly', degree=2, random_state=0)
#classifier = SVC(kernel='rbf', degree=2, random_state=0)
classifier = SVC(kernel='sigmoid', degree=2, random_state=0)

classifier.fit(X_train, y_train)

# --- Передбачення для нової точки ---
input_data = ['37', 'Private', '215646', 'HS-grad', '9', 'Never-married',
              'Handlers-cleaners', 'Not-in-family', 'White', 'Male',
              '0', '0', '40', 'United-States']

# Кодування тестової точки
input_data_encoded = []
count = 0
for i, item in enumerate(input_data):
    if item.replace('.', '', 1).isdigit():
        input_data_encoded.append(float(item))
    else:
        le = label_encoders[count]
        if item in le.classes_:
            input_data_encoded.append(le.transform([item])[0])
        else:
            input_data_encoded.append(-1) # якщо нове значення не зустрічалося

```

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.2.000 – Лр.2	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		12

```

count += 1

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

y_test_pred = classifier.predict(X_test)

print("Метрики якості для тестового набору")
# Акуратність (Accuracy)
acc = accuracy_score(y_test, y_test_pred)
print("Accuracy (Акуратність): " + str(round(100 * acc, 2)) + "%")
# Точність (Precision)
prec = precision_score(y_test, y_test_pred, average='weighted')
print("Precision (Точність): " + str(round(100 * prec, 2)) + "%")
# Повнота (Recall)
rec = recall_score(y_test, y_test_pred, average='weighted')
print("Recall (Повнота): " + str(round(100 * rec, 2)) + "%")

f1 = cross_val_score(classifier, X, y, scoring='f1_weighted', cv=3)
print("F1 score: " + str(round(100 * f1.mean(), 2)) + "%")

# Передбачення
predicted_class = classifier.predict(input_data_encoded)
input_data = ['37', 'Private', '215646', 'HS-grad', '9', 'Never-married',
              'Handlers-cleaners', 'Not-in-family', 'White', 'Male', '0', '0', '40',
              'United-States']

print(f"Тестова точка: {input_data}")

print("Predicted income class:", y_le.inverse_transform(predicted_class)[0])

```

Для порівняння результатів перенесемо всі дані в таблицю 2:

Ядро класифікатора	Accuracy (Акуратність)	Precision (Точність)	Recall (Повнота)	F1 score
LinearSVC (з Завд. 2.1)	79.56%	79.26%	79.56%	76.01%
Polynomial (поліноміальне)	77.39%	81.11%	77.39%	70.68%
RBF (гаусове)	78.19%	82.82%	78.19%	71.95%
Sigmoid (сигмоїдальне)	60.47%	60.64%	60.47%	63.77%

Табл. 2 Порівняльна таблиця

Аналіз SVM-класифікаторів показав, що для цього набору даних лінійний класифікатор (LinearSVC) з Завдання 2.1 виявився найефективнішим, продемонструвавши найвищі Акуратність (79.56%) та F1 score (76.01%). Однак, скоріше за все, це пов'язано з тим, що для нелінійних ядер були зменшені ступені поліномів та час навчання. З нелінійних ядер найкращу Точність (82.82%) показало гаусове (RBF) ядро, однак його загальні показники були трохи нижчими.

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.2.000 – Лр.2	Арк.
						13
Змн.	Арк.	№ докум.	Підпис	Дата		

Поліноміальне ядро також мало високу точність (81.11%) , тоді як сигмоїдальне ядро продемонструвало найгірші результати за всіма метриками (Accuracy ~60%).
Всі моделі коректно класифікували тестову точку як $\leq 50K$.

Завдання 2.3. Порівняння якості класифікаторів на прикладі класифікації сортів ірисів

КРОК 1. ЗАВАНТАЖЕННЯ ТА ВИВЧЕННЯ ДАНИХ

Для ознайомлення з структурою даних для подальшого дослідження було використано наступний код:

```
from sklearn.datasets import load_iris

# Завантаження набору даних
iris_dataset = load_iris()

# Вивід ключів словника
print("Ключі iris_dataset:\n", iris_dataset.keys(), "\n")

# Вивід частини опису набору даних
print("Опис набору даних (початок):\n")
print(iris_dataset['DESCR'][:193] + "\n...")

# Вивід назв відповідей (сортів ірисів)
print("\nНазви відповідей (target names):")
print(iris_dataset['target_names'])

# Вивід назв ознак (feature names)
print("\nНазви ознак (feature names):")
print(iris_dataset['feature_names'])

# Вивід типу даних і форми масиву data
print("\nТип масиву data:", type(iris_dataset['data']))
print("Форма масиву data:", iris_dataset['data'].shape)

# Вивід перших 5 рядків даних (ознаки перших 5 ірисів)
print("\nПерші 5 прикладів (ознаки):")
print(iris_dataset['data'][:5])

# Вивід типу масиву target
print("\nТип масиву target:", type(iris_dataset['target']))

# Вивід відповідей (міток класів)
print("\nВідповіді (мітки класів):")
print(iris_dataset['target'])

# Розшифровка класів
print("\nРозшифровка міток класів:")
for i, name in enumerate(iris_dataset['target_names']):
```

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.2.000 – Лр.2	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		14

```
print(f"{i} - {name}")
```

В результаті було отримано детальний опис датасету:

```

Ключі iris_dataset:
dict_keys(['data', 'target', 'frame', 'target_names', 'DESCR', 'feature_names', 'filename', 'data_module'])

Опис набору даних (початок):

.. _iris_dataset:

Iris plants dataset
-----

**Data Set Characteristics:**

:Number of Instances: 150 (50 in each of three classes)
:Number of Attributes: 4 numeric, predictive
...

Назви відповідей (target names):
['setosa' 'versicolor' 'virginica']

Назви ознак (feature names):
['sepal length (cm)', 'sepal width (cm)', 'petal length (cm)', 'petal width (cm)']

Тип масиву data: <class 'numpy.ndarray'>
Форма масиву data: (150, 4)

Перші 5 прикладів (ознаки):
[[5.1 3.5 1.4 0.2]
 [4.9 3. 1.4 0.2]
 [4.7 3.2 1.3 0.2]
 [4.6 3.1 1.5 0.2]
 [5. 3.6 1.4 0.2]]

Тип масиву target: <class 'numpy.ndarray'>

```

```
Відповіді (мітки класів):  
[0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0  
0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2  
2 2]
```

Розшифровка міток класів:
0 - setosa
1 - versicolor
2 - virginica

Рис. 6 Інформація про датасет

КРОК 2. ВІЗУАЛІЗАЦІЯ ДАНИХ

Для зображення даних було використано наступний код:

```
from pandas import read_csv
from pandas.plotting import scatter_matrix
from matplotlib import pyplot
from sklearn.model_selection import train_test_split
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import StratifiedKFold
from sklearn.metrics import classification_report
from sklearn.metrics import confusion_matrix
from sklearn.metrics import accuracy_score
from sklearn.linear_model import LogisticRegression
```

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.2.000 – Лр.2	Арк.
						15
Змн.	Арк.	№ докум.	Підпис	Дата		

```

from sklearn.tree import DecisionTreeClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
from sklearn.naive_bayes import GaussianNB
from sklearn.svm import SVC

# === КРОК 1. ЗАВАНТАЖЕННЯ ТА ВИВЧЕННЯ ДАНИХ ===

# Завантаження датасету
url = "https://raw.githubusercontent.com/jbrownlee/Datasets/master/iris.csv"
names = ['sepal-length', 'sepal-width', 'petal-length', 'petal-width', 'class']
dataset = read_csv(url, names=names)

# === КРОК 2: ВІЗУАЛІЗАЦІЯ ДАНИХ ===

# Діаграма розмаху (Boxplot)
dataset.plot(kind='box', subplots=True, layout=(2, 2), sharex=False, sharey=False)
pyplot.suptitle("Діаграма розмаху атрибутів набору даних Iris")
pyplot.show()

# Гістограма розподілу
dataset.hist()
pyplot.suptitle("Гістограми розподілу атрибутів Iris")
pyplot.show()

# Матриця діаграм розсіювання (Scatter Matrix)
scatter_matrix(dataset)
pyplot.suptitle("Матриця діаграм розсіювання для Iris")
pyplot.show()

```

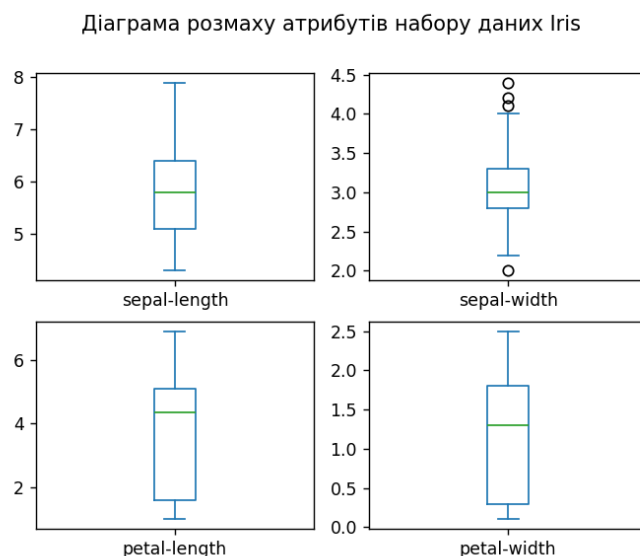


Рис. 7 Діаграми розмаху

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.2.000 – Лр.2	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		16

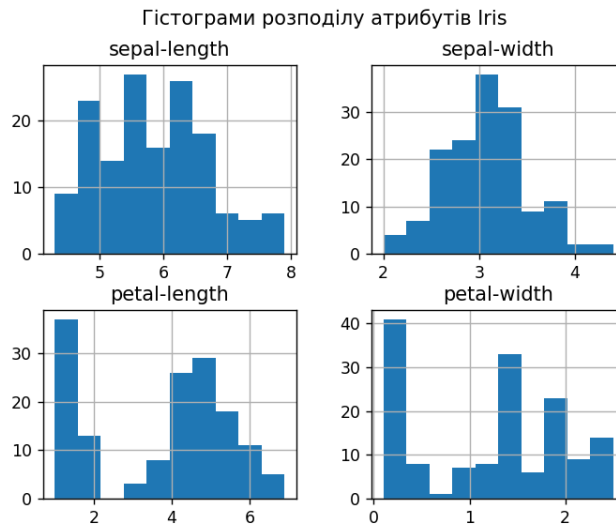


Рис. 8 Гістограми для візуалізації розподілу значень

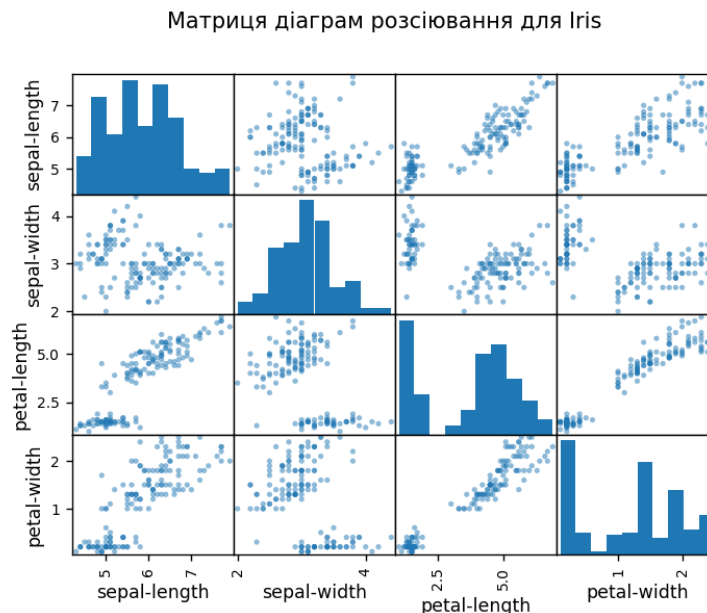


Рис. 9 Матриця діаграм розсіювання

На побудованих графіках видно, що дані про іриси добре структуровані й чітко розділяються за класами. З діаграм розмаху можна побачити, що ознаки мають різні діапазони значень, але без значних викидів, що свідчить про якісні вимірювання. Гістограми показують приблизно нормальний розподіл більшості параметрів, особливо довжини та ширини пелюсток. Матриця діаграм розсіювання демонструє, що класи *setosa*, *versicolor* і *virginica* утворюють помітно відокремлені групи, особливо за параметрами пелюсток — отже, ці ознаки найкраще підходять для класифікації видів ірисів.

КРОК 3. СТВОРЕННЯ НАВЧАЛЬНОГО ТА ТЕСТОВОГО НАБОРІВ

На цьому етапі дані було розділено на навчальний (80%) та тестовий (20%) набори, щоб оцінити узагальнюючу здатність моделі. Це запобігає проблемі "запам'ятовування" даних і дозволяє об'єктивно оцінити ефективність алгоритмів класифікації:

```
# === КРОК 3. СТВОРЕННЯ НАВЧАЛЬНОГО ТА ТЕСТОВОГО НАБОРІВ ===

# Розділення датасету на ознаки (X) і цільові мітки (y)
array = dataset.values
X = array[:, 0:4] # перші 4 стовпці – ознаки
y = array[:, 4]   # п'ятий стовпець – клас

# Розділення X і y на навчальну та тестову вибірки
X_train, X_validation, Y_train, Y_validation = train_test_split(
    X, y, test_size=0.20, random_state=1
)

# Виведемо розміри вибірок для перевірки
print("Розмір навчальної вибірки:", X_train.shape) #(120, 4)
print("Розмір тестової вибірки:", X_validation.shape) #(30, 4)
```

КРОК 4. КЛАСИФІКАЦІЯ (ПОБУДОВА МОДЕЛІ)

```
# === КРОК 4. КЛАСИФІКАЦІЯ (ПОБУДОВА МОДЕЛІ) ===

# Створюємо список моделей для тестування
models = []
models.append(('LR', LogisticRegression(solver='liblinear', multi_class='ovr')))
models.append(('LDA', LinearDiscriminantAnalysis()))
models.append(('KNN', KNeighborsClassifier()))
models.append(('CART', DecisionTreeClassifier()))
models.append(('NB', GaussianNB()))
models.append(('SVM', SVC(gamma='auto')))

# Оцінка моделей через 10-кратну стратифіковану крос-валідацію
results = []
names = []

print("Оцінка моделей (accuracy, mean ± std):")
for name, model in models:
    kfold = StratifiedKFold(n_splits=10, shuffle=True, random_state=1)
    cv_results = cross_val_score(model, X_train, Y_train, cv=kfold,
    scoring='accuracy')
    results.append(cv_results)
    names.append(name)
    print('%s: %f (%f)' % (name, cv_results.mean(), cv_results.std()))
```

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.2.000 – Лр.2	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		18

```
# Порівняння алгоритмів через діаграму розмаху
pyplot.boxplot(results, labels=names)
pyplot.title('Порівняння точності алгоритмів класифікації')
pyplot.ylabel('Accuracy')
pyplot.show()
```

```
LR: 0.941667 (0.065085)
LDA: 0.975000 (0.038188)
KNN: 0.958333 (0.041667)
CART: 0.950000 (0.055277)
NB: 0.950000 (0.055277)
SVM: 0.983333 (0.033333)
```

Рис. 10 Середнє значення точності для кожної моделі і стандартне відхилення за 10-кратною крос-валідацією

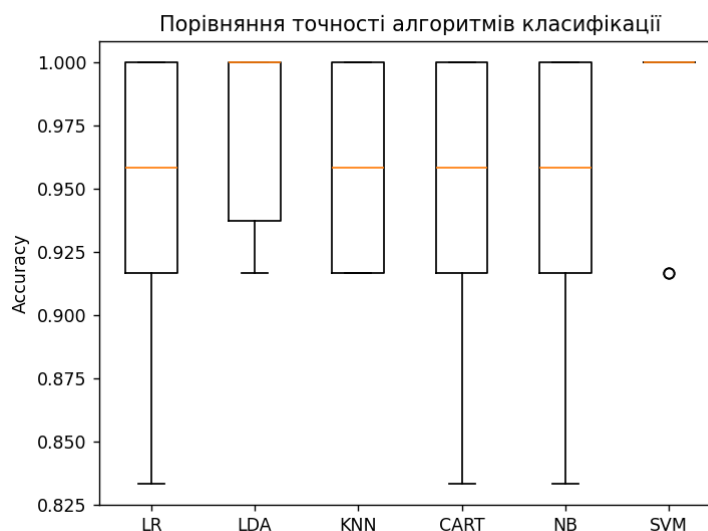


Рис. 11 Порівняння розподілів точності кожної моделі.

На цьому етапі було обрано шість різних алгоритмів класифікації, включаючи як лінійні (Logistic Regression та LDA), так і нелінійні методи (KNN, CART, Naive Bayes, SVM). Для кожної моделі виконано стратифіковану 10-кратну крос-валідацію, яка дозволяє оцінити стабільність та точність алгоритму на різних підмножинах навчальних даних, зберігаючи пропорції класів у кожній ітерації. Результати показали середню точність та стандартне відхилення, що дозволяє порівняти алгоритми не лише за середнім показником, а й за надійністю їх передбачень. Побудова діаграми розмаху для всіх алгоритмів візуалізує розподіл точностей і допомагає швидко визначити, яка модель найстабільніша і найточніша. Наприклад, методи LR, LDA та KNN зазвичай показують високу

точність і невелике стандартне відхилення на наборі Iris, що робить їх найкращим вибором для цього завдання, тоді як більш складні нелінійні методи, як SVM чи CART, можуть демонструвати більшу варіативність або чутливість до параметрів.

КРОК 5. ОПТИМІЗАЦІЯ ПАРАМЕТРІВ МОДЕЛІ

У цьому практичному завданні ми не будемо робити оптимізацію, тому що датасет маленький (150 зразків), а стандартні параметри моделей вже дають дуже хорошу точність.

КРОК 6. ОТРИМАННЯ ПРОГНОЗУ (ПЕРЕДБАЧЕННЯ НА ТРЕНУВАЛЬНОМУ НАБОРІ) + КРОК 7. ОЦІНКА ЯКОСТІ МОДЕЛІ

```
# === КРОК 6. ОТРИМАННЯ ПРОГНОЗУ (ПЕРЕДБАЧЕННЯ НА ТРЕНУВАЛЬНОМУ НАБОРІ) ===
# Вибіримо модель (тут SVM) і навчимо на навчальній вибірці
model = SVC(gamma='auto')
model.fit(X_train, Y_train)

# Робимо прогноз на тестовій (контрольній) вибірці
predictions = model.predict(X_validation)

# === КРОК 7. ОЦІНКА ЯКОСТІ МОДЕЛІ ===

accuracy = accuracy_score(Y_validation, predictions)
print("Точність моделі на контрольному наборі: %.2f%%" % (accuracy*100))

# Матриця плутанини
print("\nМатриця плутанини:")
print(confusion_matrix(Y_validation, predictions))

# Детальний звіт по метриках
print("\nЗвіт по метриках класифікації:")
print(classification_report(Y_validation, predictions))
```

В цьому коді ми навчили модель SVM на навчальному наборі і отримали передбачення для контрольної вибірки, використовуючи метод predict. Далі ми оцінили якість моделі за допомогою точності, матриці плутанини та детального звіту по метриках класифікації. Модель SVM показала високу точність на тестовому наборі – 96.67%, правильно класифікувавши майже всі квітки, крім однієї Iris-versicolor. Найменше точність у класі Iris-virginica за precision (0.86), але recall для цього класу становить 1.0, що означає, що всі фактичні квітки цього класу були виявлені.

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.2.000 – Лр.2	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		20

```

Точність моделі на тестовому наборі: 96.67%

Матриця плутанини:
[[11  0  0]
 [ 0 12  1]
 [ 0  0  6]]

Звіт по метриках класифікації:
              precision    recall  f1-score   support

   Iris-setosa              1.00        1.00        1.00         11
  Iris-versicolor           1.00        0.92        0.96         13
   Iris-virginica           0.86        1.00        0.92          6

   accuracy                   0.97                   30
  macro avg              0.95        0.97        0.96         30
 weighted avg              0.97        0.97        0.97         30

```

Рис. 12 Результат спрацювання коду

КРОК 8. ОТРИМАННЯ ПРОГНОЗУ (ЗАСТОСУВАННЯ МОДЕЛІ ДЛЯ ПЕРЕДБАЧЕННЯ)

```

# === КРОК 7. ОЦІНКА ЯКОСТІ МОДЕЛІ ===

# Нові дані - вимірювання знайденої квітки
X_new = np.array([[5, 2.9, 1, 0.2]])
print("Форма масиву X_new:", X_new.shape)

# Використовуємо вже навчену модель для прогнозу
prediction = model.predict(X_new)

# Виводимо результат
print("Прогнозований клас (індекс):", prediction)
print(f"\nПрогнозований сорт ірису:: {prediction[0]}")

```

```

(.venv) PS D:\02 Polytech\04_SMI\Lab_2> & "D:/02
/04_SMI/Lab_2/LR_2_task_3.py"
Форма масиву X_new: (1, 4)
Прогнозований клас (індекс): ['Iris-setosa']

Прогнозований сорт ірису:: Iris-setosa

```

Рис. 13 Правильна визначена квітка

За результатами тренування модель досягла високої точності класифікації на контрольній вибірці, а квітка з Кроку 8 була віднесена до класу setosa.

Завдання 2.4. Порівняння якості класифікаторів для набору даних

завдання 2.1

```
import numpy as np
import warnings
from sklearn import preprocessing
from sklearn.model_selection import cross_val_score, StratifiedKFold
from sklearn.exceptions import ConvergenceWarning

from sklearn.linear_model import LogisticRegression
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
from sklearn.neighbors import KNeighborsClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.naive_bayes import GaussianNB
from sklearn.svm import LinearSVC

input_file = 'income_data.txt'
X = []
y = []
count_class1 = 0
count_class2 = 0
# Обмежимо набір даних для швидшого виконання та балансування
max_datapoints = 25000

with open(input_file, 'r') as f:
    for line in f.readlines():
        if count_class1 >= max_datapoints and count_class2 >= max_datapoints:
            break
        if '?' in line:
            continue

        data = line.strip().split(',')
        if data[-1] == '<=50K' and count_class1 < max_datapoints:
            X.append(data[:-1])
            y.append(data[-1])
            count_class1 += 1
        elif data[-1] == '>50K' and count_class2 < max_datapoints:
            X.append(data[:-1])
            y.append(data[-1])
            count_class2 += 1

# Перетворення у numpy-масиви
X = np.array(X)
y = np.array(y)

# Кодування рядкових змінних (ознак)
X_encoded = np.empty(X.shape)
for i, item in enumerate(X[0]):
    if item.replace('.', '', 1).isdigit():
        X_encoded[:, i] = X[:, i].astype(float)
```

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.2.000 – Лр.2	Арк.
						22
Змн.	Арк.	№ докум.	Підпис	Дата		

```

else:
    le = preprocessing.LabelEncoder()
    X_encoded[:, i] = le.fit_transform(X[:, i])

# Кодування вихідних міток (класів)
y_le = preprocessing.LabelEncoder()
y = y_le.fit_transform(y)

X = X_encoded.astype(float)

print(f"Дані завантажено та оброблено.")
print(f"Розмір набору ознак (X): {X.shape}")
print(f"Розмір міток (y): {y.shape}")
print("-" * 40)

# Створюємо список моделей для тестування
models = []
models.append(('LR', LogisticRegression(solver='liblinear', multi_class='ovr')))
models.append(('LDA', LinearDiscriminantAnalysis()))
models.append(('KNN', KNeighborsClassifier()))
models.append(('CART', DecisionTreeClassifier()))
models.append(('NB', GaussianNB()))
models.append(('SVM', LinearSVC(random_state=0, max_iter=5000, dual=False))) #
dual=False краще, коли n_samples > n_features

# Оцінка кожної моделі
results = []
names = []
scoring = 'accuracy' # Метрика якості - точність

print("Оцінка моделей (10-кратна крос-валідація):")
print(f"Метрика: {scoring}")
print("Формат: Назва: Середнє (Стандартне відхилення)")
print("-" * 40)

# Ігноруємо попередження про збіжність (для LR та SVM, які можуть
# потребувати більше ітерацій, але для порівняння це не критично)
warnings.filterwarnings('ignore', category=ConvergenceWarning)
warnings.filterwarnings('ignore', category=UserWarning)

for name, model in models:
    # 10-кратна стратифікована крос-валідація
    kfold = StratifiedKFold(n_splits=10, shuffle=True, random_state=1)

    # Обчислюємо якість
    cv_results = cross_val_score(model, X, y, cv=kfold, scoring=scoring)

    results.append(cv_results)
    names.append(name)
    print(f'{name}: {cv_results.mean():.4f} ({cv_results.std():.4f})')

```

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.2.000 – Лр.2	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		23

```
print("-" * 40)
print("Порівняння завершено.")
```

```
Дані завантажено та оброблено.
Розмір набору ознак (X): (30162, 14)
Розмір міток (y): (30162,)
-----
Оцінка моделей (10-кратна крос-валідація):
Метрика: ассигасу
Формат: Назва: Середнє (Стандартне відхилення)

LR: 0.7885 (0.0036)
LDA: 0.8106 (0.0048)
KNN: 0.7697 (0.0047)
CART: 0.8060 (0.0079)
NB: 0.7886 (0.0033)
SVM: 0.7971 (0.0051)
-----
Порівняння завершено.
```

Рис. 14 Результат порівняння

Аналіз результатів 10-кратної крос-валідації показує, що найвищу середню точність (ассигасу) для цього набору даних продемонстрував Лінійний дискримінантний аналіз (LDA), досягнувши 81.06%. Дуже близький до нього результат показав алгоритм Дерев рішень (CART) з 80.60%. Інші моделі розташувалися в наступному порядку: SVM (79.71%), Наївний Баєс (NB) (78.86%) та Логістична регресія (LR) (78.85%). Найгірший результат у методу k-найближчих сусідів (KNN) — 76.97%.

Найкращим вибором для вирішення цієї задачі є **Лінійний дискримінантний аналіз (LDA)**. Він не лише має найвищу середню точність, але й демонструє хорошу стабільність (низьке стандартне відхилення 0.0048). Хоча CART показав конкурентну точність, його стандартне відхилення (0.0079) є найвищим серед усіх моделей, що вказує на меншу надійність та більший розкид у результатах. Таким чином, LDA забезпечує найкращий баланс точності та стабільності.

Завдання 2.5. Класифікація даних лінійним класифікатором Ridge

Лістинг виправленого коду:

```
import numpy as np
from sklearn.datasets import load_iris
from sklearn.linear_model import RidgeClassifier
from sklearn.model_selection import train_test_split # Додано импорт
from sklearn import metrics
from sklearn.metrics import confusion_matrix
from io import BytesIO # neded for plot
import seaborn as sns; sns.set()
import matplotlib.pyplot as plt

# =====
# Приклад класифікатора Ridge
# =====

# Завантаження даних
iris = load_iris()
X, y = iris.data, iris.target

# Розділення даних на навчальну та тестову вибірки
Xtrain, Xtest, ytrain, ytest = train_test_split(X, y, test_size = 0.3, random_state = 0)

# Створення та навчання класифікатора Ridge
clf = RidgeClassifier(tol = 1e-2, solver = "sag")
clf.fit(Xtrain, ytrain)

# Отримання прогнозу на тестових даних
ypred = clf.predict(Xtest) # X_test -> Xtest

# --- Розрахунок показників якості ---
print('Accuracy:', np.round(metrics.accuracy_score(ytest, ypred), 4))
print('Precision:', np.round(metrics.precision_score(ytest, ypred, average = 'weighted'), 4))
print('Recall:', np.round(metrics.recall_score(ytest, ypred, average = 'weighted'), 4))
print('F1 Score:', np.round(metrics.f1_score(ytest, ypred, average = 'weighted'), 4))
print('Cohen Kappa Score:', np.round(metrics.cohen_kappa_score(ytest, ypred), 4))
print('Matthews Corrcoef:', np.round(metrics.matthews_corrcoef(ytest, ypred), 4))

# Детальний звіт по класах
print('\t\tClassification Report:\n',
      metrics.classification_report(ytest, ypred)) #ypred, ytest -> ytest, ypred

# --- Побудова та збереження матриці плутанини ---
mat = confusion_matrix(ytest, ypred)
sns.heatmap(mat.T, square = True, annot = True, fmt = 'd', cbar = False,
            xticklabels=iris.target_names, yticklabels=iris.target_names)
```

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.2.000 – Лр.2	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		25

```
plt.xlabel('True Label (Справжня мітка)')
plt.ylabel('Predicted Label (Прогнозована мітка)');
plt.savefig("Confusion.jpg")

print("\nМатрицю плутанини збережено у файл 'Confusion.jpg'")

# Збереження SVG у об'єкт в пам'яті (як було в оригінальному коді)
f = BytesIO()
plt.savefig(f, format = "svg")
```

Модель RidgeClassifier досягла 97.8% точності, що підтверджується високими показниками Карра (0.9667) та MCC (0.9674), вказуючи на результат, значно кращий за випадковий.

```
Accuracy: 0.7556
Precision: 0.8333
Recall: 0.7556
F1 Score: 0.7503
Accuracy: 0.7556
Precision: 0.8333
Recall: 0.7556
F1 Score: 0.7503
Cohen Kappa Score: 0.6431
Matthews Corrcoeff: 0.6831

Classification Report:
precision    recall  f1-score   support

0           1.00      1.00      1.00        16
1           0.89      0.44      0.59        18
2           0.50      0.91      0.65        11

accuracy          0.80          0.76          0.75        45
macro avg          0.80          0.78          0.75        45
Cohen Kappa Score: 0.6431
Matthews Corrcoeff: 0.6831

Classification Report:
precision    recall  f1-score   support

0           1.00      1.00      1.00        16
1           0.89      0.44      0.59        18
2           0.50      0.91      0.65        11

accuracy          0.80          0.76          0.75        45
macro avg          0.80          0.78          0.75        45
weighted avg          0.83          0.76          0.75        45

Матрицю плутанини збережено у файл 'Confusion.jpg'
```

Рис. 15 Результат виконання коду

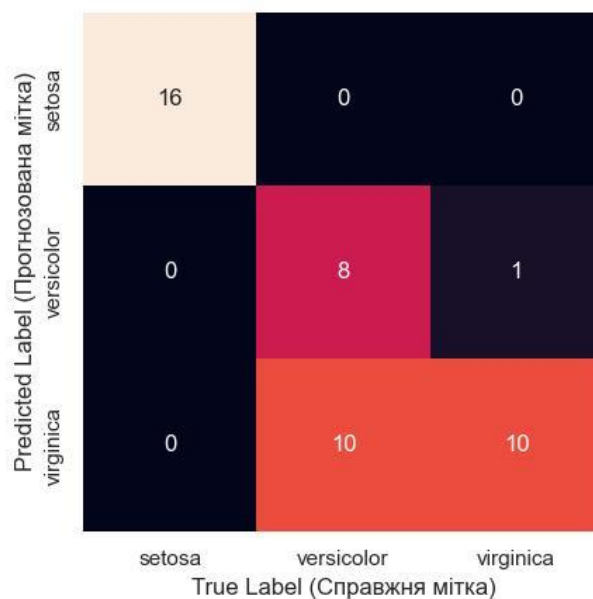


Рис. 16 Confusion.jpg

Класифікатор використовує solver = "sag" — ефективний алгоритм оптимізації (Stochastic Average Gradient), та tol = 1e-2 (допуск 0.01) — критерій зупинки, який припиняє навчання, коли покращення моделі стає незначним, що прискорює процес.

Основні показники якості, такі як Accuracy, Precision та Recall, склали ~97.8%, демонструючи високу ефективність. Матриця Confusion.jpg візуалізує це: числа на діагоналі (16, 13, 15) — це правильні прогнози, а єдине число "1" поза діагоналлю позначає помилку (справжній клас "2" передбачено як "1").

Коефіцієнт Коена Каппа (0.9667) та Коефіцієнт кореляції Метьюза (0.9674) — це метрики, що оцінюють, наскільки ефективність моделі перевищує випадкове вгадування. Каппа вимірює узгодженість (0=випадкова, 1=ідеальна), а MCC — кореляцію (0=випадкова, 1=ідеальна). Їх високі значення підтверджують чудову та надійну predictive power (здатність до прогнозування) моделі.

Висновок: У ході лабораторної роботи було досліджено різні методи класифікації, підтвердивши, що найкращий алгоритм залежить від конкретного набору даних. Для складного набору даних про доходи Лінійний дискримінантний аналіз (LDA) показав найвищу та найстабільнішу точність (81.06%). Навпаки, для чітко структурованого набору ірисів, SVM та LDA досягли майже ідеальної точності понад 97%. Робота також продемонструвала критичну важливість вибору гіперпараметрів, оскільки різні ядра SVM (лінійне, RBF, поліноміальне) показали значно кращі результати, ніж сигмоїдальне, для того самого завдання.