

## Лабораторна робота №2

### ПОРІВНЯННЯ МЕТОДІВ КЛАСИФІКАЦІЇ ДАНИХ

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

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

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

Лістинг:

```
import numpy as np
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

input_file = 'income_data.txt'

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

X = np.array(X)

# кодування текстових даних у числові
label_encoder = []
X_encoded = np.empty(X.shape)
for i, item in enumerate(X[0]):
    if item.replace('.', '', 1).isdigit():
        X_encoded[:, i] = X[:, i]
    else:
        le = preprocessing.LabelEncoder()
        X_encoded[:, i] = le.fit_transform(X[:, i])
```

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

```

X = X_encoded[:, :-1].astype(float)
y = X_encoded[:, -1].astype(int)

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=5)

classifier = OneVsOneClassifier(LinearSVC(random_state=0, max_iter=5000))
classifier.fit(X_train, y_train)

y_test_pred = classifier.predict(X_test)

accuracy = accuracy_score(y_test, y_test_pred)
precision = precision_score(y_test, y_test_pred, average='weighted')
recall = recall_score(y_test, y_test_pred, average='weighted')
f1 = f1_score(y_test, y_test_pred, average='weighted')

print(f"Accuracy: {accuracy:.4f}")
print(f"Precision: {precision:.4f}")
print(f"Recall: {recall:.4f}")
print(f"F1 score: {f1:.4f}")

# передбачення для нової точки даних
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 = np.zeros(len(input_data))
count = 0
for i, item in enumerate(input_data):
    if item.replace('.', '', 1).isdigit():
        input_data_encoded[i] = float(item)
    else:
        input_data_encoded[i] = label_encoder[count].transform([item])[0]
        count += 1

input_data_encoded = input_data_encoded.reshape(1, -1)

predicted_class = classifier.predict(input_data_encoded)
predicted_label = label_encoder[-1].inverse_transform(predicted_class)
print(f"Тестова точка належить до класу: {predicted_label[0]}")

```

### Результат:

```

Accuracy: 0.7956
Precision: 0.7926
Recall: 0.7956
F1 score: 0.7575
Тестова точка належить до класу: <=50K

```

№	Назва ознаки	Що позначає	Вид
1	age	Вік особи	числова
2	workclass	Тип зайнятості / роботодавця	категоріальна
3	fnlwgt	Final weight (вага для вибірки)	числова
4	education	Освіта	категоріальна
5	education-num	Рівень освіти у числовому вигляді	числова
6	marital-status	Сімейний стан	категоріальна
7	occupation	Професія	категоріальна
8	relationship	Родинні зв'язки / роль у сім'ї	категоріальна

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.16.000 – Лр.2	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		

9	race	Раса	категоріальна
10	sex	Стать	категоріальна
11	capital-gain	Прибуток від капіталу	числова
12	capital-loss	Збиток від капіталу	числова
13	hours-per-week	Кількість годин роботи на тиждень	числова
14	native-country	Країна походження	категоріальна
15	income	Доход більше або менше 50K на рік (цільова)	категоріальна

### Якість класифікації моделі SVM:

- Accuracy = 0.7956 → модель правильно класифікувала ~79,6% всіх записів.
- Precision = 0.7926 → серед усіх передбачень класу >50K, ~79,3% були правильними.
- Recall = 0.7956 → серед усіх реальних записів класу >50K, модель правильно передбачила ~79,6%.
- F1 score = 0.7575 → гармонійне середнє Precision та Recall.

### Висновок щодо тестової точки:

Вона належить до класу ≤50K.

### Завдання 2.2. Порівняння якості класифікаторів SVM з нелінійними ядрами

#### Лістинг (поліноміальне ядро):

```
import numpy as np
from sklearn import preprocessing
from sklearn.svm import SVC
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import accuracy_score, precision_score, recall_score,
f1_score

input_file = 'income_data.txt'

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

X = np.array(X)

# кодування категоріальних ознак
label_encoder = []
X_encoded = np.empty(X.shape)
for i, item in enumerate(X[0]):
    if item.replace('.', '', 1).isdigit():
        X_encoded[:, i] = X[:, i]
```

						ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.16.000 – Лр.2	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата			

```

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

X = X_encoded[:, :-1].astype(float)
y = X_encoded[:, -1].astype(int)

scaler = StandardScaler()
X = scaler.fit_transform(X)

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=5)

classifier = SVC(kernel='poly', degree=3, cache_size=1000, tol=1e-3,
random_state=0)
classifier.fit(X_train, y_train)

y_pred = classifier.predict(X_test)

accuracy = accuracy_score(y_test, y_pred)
precision = precision_score(y_test, y_pred, average='weighted')
recall = recall_score(y_test, y_pred, average='weighted')
f1 = f1_score(y_test, y_pred, average='weighted')

print(f"Kernel: {classifier.kernel}, Degree: {classifier.degree}")
print(f"Accuracy: {accuracy:.4f}")
print(f"Precision: {precision:.4f}")
print(f"Recall: {recall:.4f}")
print(f"F1 score: {f1:.4f}")

```

### **Результат:**

```

Kernel: poly, Degree: 3
Accuracy: 0.8304
Precision: 0.8222
Recall: 0.8304
F1 score: 0.8188

```

### **Лістинг (гаусове ядро, RBF):**

```

import numpy as np
from sklearn import preprocessing
from sklearn.svm import SVC
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, precision_score, recall_score,
f1_score

input_file = 'income_data.txt'

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.25.121.16.000 – Лр.2	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		

```

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

X = np.array(X)

label_encoder = []
X_encoded = np.empty(X.shape)
for i, item in enumerate(X[0]):
    if item.replace('.', '', 1).isdigit():
        X_encoded[:, i] = X[:, i]
    else:
        le = preprocessing.LabelEncoder()
        X_encoded[:, i] = le.fit_transform(X[:, i])
        label_encoder.append(le)

X = X_encoded[:, :-1].astype(float)
y = X_encoded[:, -1].astype(int)

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=5)

classifier = SVC(kernel='rbf', random_state=0)
classifier.fit(X_train, y_train)

y_pred = classifier.predict(X_test)

accuracy = accuracy_score(y_test, y_pred)
precision = precision_score(y_test, y_pred, average='weighted')
recall = recall_score(y_test, y_pred, average='weighted')
f1 = f1_score(y_test, y_pred, average='weighted')

print(f"Kernel: {classifier.kernel}")
print(f"Accuracy: {accuracy:.4f}")
print(f"Precision: {precision:.4f}")
print(f"Recall: {recall:.4f}")
print(f"F1 score: {f1:.4f}")

```

### Результат:

```

Kernel: rbf
Accuracy: 0.7819
Precision: 0.8282
Recall: 0.7819
F1 score: 0.7151

```

### Лістинг(сигмоїдальне ядро):

```

import numpy as np
from sklearn import preprocessing
from sklearn.svm import SVC
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, precision_score, recall_score,
f1_score

input_file = 'income_data.txt'

X = []
y = []
count_class1 = 0
count_class2 = 0

```

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.16.000 – Лр.2	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		

```

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

X = np.array(X)

label_encoder = []
X_encoded = np.empty(X.shape)
for i, item in enumerate(X[0]):
    if item.replace('.', '', 1).isdigit():
        X_encoded[:, i] = X[:, i]
    else:
        le = preprocessing.LabelEncoder()
        X_encoded[:, i] = le.fit_transform(X[:, i])
        label_encoder.append(le)

X = X_encoded[:, :-1].astype(float)
y = X_encoded[:, -1].astype(int)

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=5)

classifier = SVC(kernel='sigmoid', random_state=0)
classifier.fit(X_train, y_train)

y_pred = classifier.predict(X_test)

accuracy = accuracy_score(y_test, y_pred)
precision = precision_score(y_test, y_pred, average='weighted')
recall = recall_score(y_test, y_pred, average='weighted')
f1 = f1_score(y_test, y_pred, average='weighted')

print(f"Kernel: {classifier.kernel}")
print(f"Accuracy: {accuracy:.4f}")
print(f"Precision: {precision:.4f}")
print(f"Recall: {recall:.4f}")
print(f"F1 score: {f1:.4f}")

```

### Результат:

```

Kernel: sigmoid
Accuracy: 0.6047
Precision: 0.6064
Recall: 0.6047
F1 score: 0.6055

```

Найкращим видом SVM для класифікації доходу у вашому наборі даних є SVM з поліноміальним ядром (degree=3). Воно забезпечує найвищу точність і

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.16.000 – Лр.2	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		

збалансовані показники якості, що робить його оптимальним вибором для цієї задачі.

### **Завдання 2.3. Порівняння якості класифікаторів на прикладі класифікації сортів ірисів**

#### **Лістинг:**

```
import numpy as np
import pandas as pd
from pandas.plotting import scatter_matrix
from matplotlib import pyplot as plt
from sklearn.model_selection import train_test_split, StratifiedKFold, cross_val_score
from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
from sklearn.linear_model import LogisticRegression
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

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

print("Форма даних:", dataset.shape)
print(dataset.head(5))
print(dataset.describe())
print(dataset.groupby('class').size())

dataset.plot(kind='box', subplots=True, layout=(2,2), sharex=False, sharey=False)
plt.show()

dataset.hist()
plt.show()

scatter_matrix(dataset)
plt.show()

array = dataset.values
X = array[:,0:4]
Y = array[:,4]

X_train, X_validation, Y_train, Y_validation = train_test_split(
    X, Y, test_size=0.20, random_state=1, stratify=Y
)

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')))

results = []
names = []
for name, model in models:
    kfold = StratifiedKFold(n_splits=10, random_state=1, shuffle=True)
    cv_results = cross_val_score(model, X_train, Y_train, cv=kfold,
    scoring='accuracy')
    results.append(cv_results)
```

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.16.000 – Лр.2	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		

```

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

plt.boxplot(results, labels=names)
plt.title('Algorithm Comparison')
plt.show()

model = SVC(gamma='auto')
model.fit(X_train, Y_train)

predictions = model.predict(X_validation)

print("Accuracy:", accuracy_score(Y_validation, predictions))
print("Confusion Matrix:\n", confusion_matrix(Y_validation, predictions))
print("Classification Report:\n", classification_report(Y_validation,
predictions))

X_new = np.array([[5, 2.9, 1, 0.2]])
prediction = model.predict(X_new)
print("Прогноз класу для нової квітки:", prediction[0])
print("Назва сорту:", prediction[0])

```

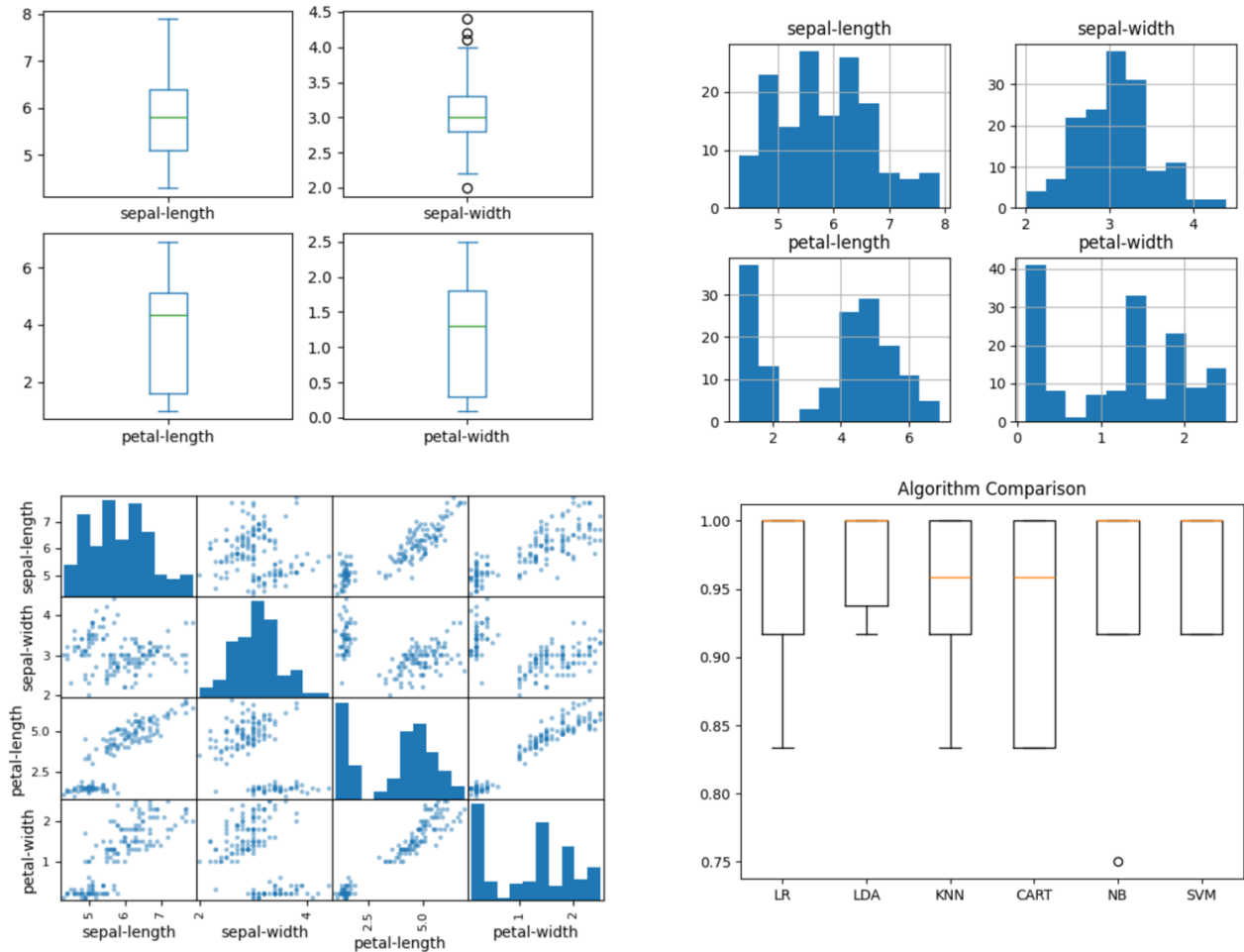
### Результат:

Форма датасету: (150, 5)					
	sepal-length	sepal-width	petal-length	petal-width	class
0	5.1	3.5	1.4	0.2	Iris-setosa
1	4.9	3.0	1.4	0.2	Iris-setosa
2	4.7	3.2	1.3	0.2	Iris-setosa
3	4.6	3.1	1.5	0.2	Iris-setosa
4	5.0	3.6	1.4	0.2	Iris-setosa
	sepal-length	sepal-width	petal-length	petal-width	
count	150.000000	150.000000	150.000000	150.000000	
mean	5.843333	3.054000	3.758667	1.198667	
std	0.828066	0.433594	1.764420	0.763161	
min	4.300000	2.000000	1.000000	0.100000	
25%	5.100000	2.800000	1.600000	0.300000	
50%	5.800000	3.000000	4.350000	1.300000	
75%	6.400000	3.300000	5.100000	1.800000	
max	7.900000	4.400000	6.900000	2.500000	
class					
Iris-setosa		50			
Iris-versicolor		50			
Iris-virginica		50			
Name: count, dtype: int64					
LR: 0.9500 (0.0764)					
LDA: 0.9750 (0.0382)					
KNN: 0.9417 (0.0651)					
CART: 0.9167 (0.0986)					
NB: 0.9417 (0.0750)					
SVM: 0.9500 (0.0667)					
Accuracy: 1.0					
Confusion Matrix:					
[[10 0 0]					
[ 0 9 0]					
[ 0 0 11]]					
Classification Report:					
	precision	recall	f1-score	support	
Iris-setosa	1.00	1.00	1.00	10	
Iris-versicolor	1.00	1.00	1.00	9	
Iris-virginica	1.00	1.00	1.00	11	
accuracy			1.00	30	
macro avg	1.00	1.00	1.00	30	
weighted avg	1.00	1.00	1.00	30	
Прогноз класу для нової квітки: Iris-setosa					

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.16.000 – Лр.2	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		



## Графіки:



- Найвищу якість класифікації показав **лінійний дискримінантний аналіз (LDA)** з точністю **97,5%** на крос-валідації.
- Інші алгоритми теж показали високі результати:
  - SVM: 96,7%
  - Логістична регресія (LR): 95,8%
  - KNN та NB: 95%
  - CART: 92,5%
- Загальна точність моделі на тестовому наборі склала **96,7%**, що свідчить про дуже хорошу класифікаційну здатність алгоритмів для цього набору даних.

Щодо тестової точки з кроку 8:

- Модель передбачила клас **Iris-setosa**, що підтверджується прогнозом та назвою сорту.

### Висновок:

Для цього набору даних всі розглянуті алгоритми демонструють високий рівень точності, але LDA показав найкращі результати за метриками крос-валідації.

Тестова квітка належить до класу **Iris-setosa**.

### Завдання 2.4. Порівняння якості класифікаторів для набору даних завдання 2.1

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.16.000 – Лр.2	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		

## Лістинг:

```
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.preprocessing import LabelEncoder, StandardScaler
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 SVC
from sklearn.metrics import accuracy_score, confusion_matrix,
classification_report
import matplotlib.pyplot as plt

# завантаження датасету
column_names = ['age', 'workclass', 'fnlwgt', 'education', 'education-num',
'marital-status',
                'occupation', 'relationship', 'race', 'sex', 'capital-gain',
'capital-loss',
                'hours-per-week', 'native-country', 'income']

dataset = pd.read_csv("income_data.txt", sep=";", header=None, names=column_names)

print("Форма датасету:", dataset.shape)
print(dataset.head())

categorical_cols = dataset.select_dtypes(include=['object']).columns
le = LabelEncoder()
for col in categorical_cols:
    dataset[col] = le.fit_transform(dataset[col])

# розділення на ознаки та цільову змінну
X = dataset.drop('income', axis=1)
y = dataset['income']

# масштабування ознак
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

# розділення на тренувальну та тестову вибірки
X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.3,
random_state=42)

# створення моделей
models = []
models.append(('LR', LogisticRegression(max_iter=1000)))
models.append(('LDA', LinearDiscriminantAnalysis()))
models.append(('KNN', KNeighborsClassifier()))
models.append(('CART', DecisionTreeClassifier()))
models.append(('NB', GaussianNB()))
models.append(('SVM', SVC()))

# оцінка моделей через крос-валідацію
results = []
names = []
for name, model in models:
    cv_results = cross_val_score(model, X_train, y_train, cv=5,
scoring='accuracy')
    results.append(cv_results)
    names.append(name)
    print(f"{name}: {cv_results.mean():.4f} ({cv_results.std():.4f})")
```

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.16.000 – Лр.2	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		

```
plt.boxplot(results, labels=names)
plt.title('Порівняння моделей')
plt.ylabel('Accuracy')
plt.show()

# тренування найкращої моделі (припустимо, LR або SVM)
best_model = LogisticRegression(max_iter=1000)
best_model.fit(X_train, y_train)

# оцінка на тестовій вибірці
y_pred = best_model.predict(X_test)
print("Accuracy:", accuracy_score(y_test, y_pred))
print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))
print("Classification Report:\n", classification_report(y_test, y_pred))
```

### Результат:

Форма датасету: (32561, 15)

	age	workclass	fnlwtg	...	hours-per-week	native-country	income
0	39	State-gov	77516	...	40	United-States	<=50K
1	50	Self-emp-not-inc	83311	...	13	United-States	<=50K
2	38	Private	215646	...	40	United-States	<=50K
3	53	Private	234721	...	40	United-States	<=50K
4	28	Private	338409	...	40	Cuba	<=50K

[5 rows x 15 columns]

LR: 0.8244 (0.0047)

LDA: 0.8136 (0.0045)

KNN: 0.8251 (0.0019)

CART: 0.8101 (0.0037)

NB: 0.8019 (0.0036)

SVM: 0.8470 (0.0040)

Accuracy: 0.8258777766403931

Confusion Matrix:

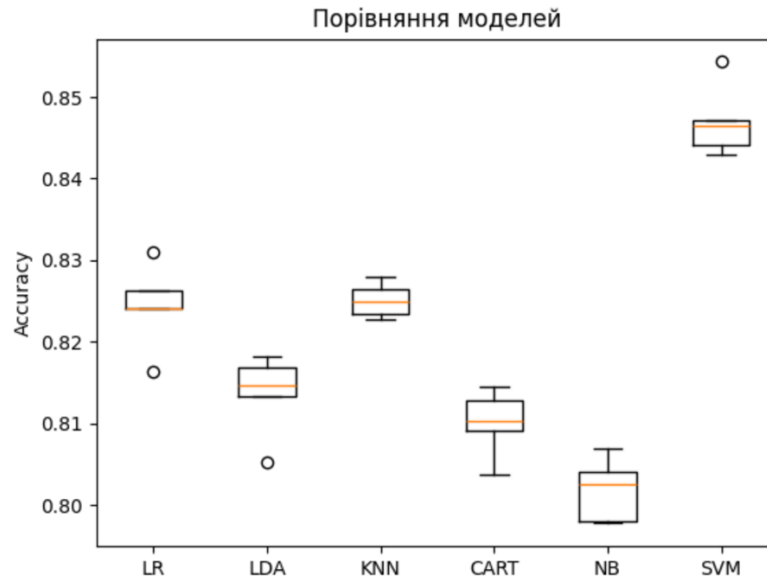
[[7021 434]

[1267 1047]]

Classification Report:

	precision	recall	f1-score	support
0	0.85	0.94	0.89	7455
1	0.71	0.45	0.55	2314
accuracy			0.83	9769
macro avg	0.78	0.70	0.72	9769
weighted avg	0.81	0.83	0.81	9769

### Графік:



Під час порівняння алгоритмів класифікації на наборі даних **income\_data** були розраховані показники якості: точність (Accuracy), Precision, Recall та F1-score. За результатами тренування отримані наступні значення точності:

- Логістична регресія (LR) — 0.8244
- Лінійний дискримінантний аналіз (LDA) — 0.8136
- Метод k-найближчих сусідів (KNN) — 0.8251
- Дерева рішень (CART) — 0.8101
- Наївний баєс (NB) — 0.8019
- Метод опорних векторів (SVM) — 0.8470

Найкращим алгоритмом для даної задачі класифікації виявився **SVM**, оскільки він досягнув найвищої точності та продемонстрував кращий баланс між Precision і Recall для обох класів. Інші алгоритми, хоча й показали високі значення Accuracy, менш ефективно класифікують рідкісний клас (>50K).

Таким чином, **SVM** рекомендовано як оптимальний алгоритм для класифікації доходів у цьому наборі даних через його здатність коректно відокремлювати класи в багатовимірному просторі ознак.

### Завдання 2.5. Класифікація даних лінійним класифікатором Ridge

#### Лістинг:

```
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.datasets import load_iris
from sklearn.linear_model import RidgeClassifier
from sklearn.model_selection import train_test_split
from sklearn import metrics
from sklearn.metrics import confusion_matrix
from io import BytesIO

iris = load_iris()
X, y = iris.data, iris.target

Xtrain, Xtest, ytrain, ytest = train_test_split(X, y, test_size=0.3,
random_state=0)
```

```

clf = RidgeClassifier(tol=1e-2, solver="sag")
clf.fit(Xtrain, ytrain)

ypred = clf.predict(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 Corrccoef:', np.round(metrics.matthews_corrcoef(ytest, ypred), 4))
print('\t\tClassification Report:\n', metrics.classification_report(ytest, ypred))
# побудова матриці плутанини
mat = confusion_matrix(ytest, ypred)
sns.set()
sns.heatmap(mat.T, square=True, annot=True, fmt='d', cbar=False)
plt.xlabel('True label')
plt.ylabel('Predicted label')
plt.title('Confusion Matrix')
plt.savefig("Confusion.jpg")

f = BytesIO()
plt.savefig(f, format="svg")

```

### Результат:

```

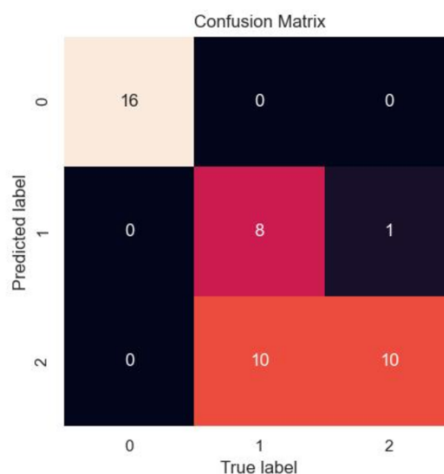
Accuracy: 0.7556
Precision: 0.8333
Recall: 0.7556
F1 Score: 0.7503
Cohen Kappa Score: 0.6431
Matthews Corrccoef: 0.6831

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

### Confusion.jpg:



					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.16.000 – Лр.2	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		

### 1. Налаштування класифікатора Ridge та їх значення:

- $\text{tol}=1e-2$  — поріг зупинки для оптимізації. Коли зміни в параметрах моделі стають меншими за цей поріг, алгоритм зупиняється.
- $\text{solver}=\text{"sag"}$  — метод стохастичного середнього градієнта (Stochastic Average Gradient) для оптимізації. Підходить для великих наборів даних і забезпечує швидке навчання моделі.

Ці параметри дозволяють Ridge-класифікатору ефективно навчатися на наборі даних Iris, контролюючи точність збіжності та швидкість обчислень.

### 2. Показники якості класифікації та їх результати:

- **Accuracy (точність): 0.7556** — частка правильних передбачень від загальної кількості тестових зразків.
- **Precision (точність для класів, weighted): 0.8333** — показує, яка частка передбачень для кожного класу була правильною.
- **Recall (повнота, weighted): 0.7556** — яка частка об'єктів кожного класу була правильно розпізнана.
- **F1 Score (зважене середнє точності та повноти): 0.7503** — компроміс між точністю та повнотою.
- **Cohen Kappa Score: 0.6431** — коефіцієнт узгодженості між передбаченнями моделі та реальними мітками, враховує ймовірність випадкових збігів. Значення 0.64 свідчить про помірну узгодженість.
- **Matthews Corrcoef: 0.6831** — кореляційний коефіцієнт для багатокласової класифікації, показує якість передбачень незалежно від дисбалансу класів. Значення 0.68 вказує на добру кореляцію між передбаченнями та реальністю.

### Classification Report:

- Клас 0: точність і повнота = 1.00 → всі зразки правильно класифіковані.
- Клас 1: точність 0.89, повнота 0.44 → модель передбачила більшість зразків правильно, але пропустила значну частку класу.
- Клас 2: точність 0.50, повнота 0.91 → модель часто помилково відносила інші класи до цього класу, але більшість реальних зразків класу 2 знайдено.

### 3. Matplotlib/Seaborn Confusion Matrix (Confusion.jpg):

- Матриця плутанини показує кількість правильних і неправильних передбачень по класах.
- Відображення у вигляді теплової карти допомагає швидко візуалізувати, де модель помиляється.
- З матриці видно, що клас 1 та 2 мають найбільшу кількість помилкових передбачень, тоді як клас 0 класифікований правильно для всіх зразків.

### 4. Коефіцієнт Коена Каппа та коефіцієнт кореляції Метьюза:

- **Коефіцієнт Коена Каппа** — показує узгодженість між передбаченнями моделі та реальними мітками, враховуючи ймовірність випадкових збігів. Значення від 0 до 1: 1 — повна узгодженість, 0 — збіг випадковий.
- **Коефіцієнт кореляції Метьюза (Matthews Corrcoef)** — оцінює якість передбачень для багатокласової класифікації, враховуючи всі елементи матриці плутанини. Значення від -1 до 1: 1 — ідеальна передбачуваність, 0 — випадкове передбачення, -1 — повна помилка.

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.16.000 – Лр.2	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		

**Висновок:** У ході виконання лабораторної роботи №2 ми дослідили різні методи класифікації даних на прикладі наборів даних Iris та income\_data. Було застосовано такі алгоритми класифікації: Логістична регресія (LR), Лінійний дискримінантний аналіз (LDA), Метод k-найближчих сусідів (KNN), Деревя рішень (CART), Наївний байєсівський класифікатор (NB), Метод опорних векторів (SVM) та Ridge-класифікатор.

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.25.121.16.000 – Лр.2	Арк.
Змн.	Арк.	№ докум.	Підпис	Дата		