

## Лабораторна робота № 2

### ПОРІВНЯННЯ МЕТОДІВ КЛАСИФІКАЦІЇ ДАНИХ

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

Github: <https://github.com/igogol-hie/myrepoz987/tree/main/SAI-LR-2>

#### Хід роботи

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

Створіть класифікатор у вигляді машини опорних векторів, призначений для прогнозування меж доходу заданої фізичної особи на основі 14 ознак (атрибутів). Метою є з'ясування умов, за яких щорічний прибуток людини перевищує \$50000 або менше цієї величини за допомогою бінарної класифікації.

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

Обчисліть значення інших показників якості класифікації (акуратність, повнота, точність) та разом з F1 занесіть їх у звіт. (Див. ЛР-1).

Зробіть висновок до якого класу належить тестова точка.

14 ознак набору даних:

age – позначає вік, числові;

workclass – позначає робочий клас, категоріальні;

fnlwgt – числові;

education – категоріальні;

education-num – числові;

marital-status – категоріальні;

occupation – категоріальні;

relationship – категоріальні;

race – категоріальні;

sex – категоріальні;

capital-gain – числові;

capital-loss – числові;

hours-per-week – числові;

native-country – категоріальні.

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

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[:-1].split(',')
        income_class = data[-1]
        if income_class == '<=50K' and count_class1 < max_datapoints:
            X.append(data)
            count_class1 += 1

        if income_class == '>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.isdigit():
        X_encoded[:, i] = X[:, i]
    else:
        current_label_encoder = preprocessing.LabelEncoder()
        label_encoder.append(current_label_encoder)
        X_encoded[:, i] = current_label_encoder.fit_transform(X[:, i])

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

classifier = OneVsOneClassifier(LinearSVC(random_state=0))
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
                                                    random_state=5)
classifier.fit(X_train, y_train)
y_test_pred = classifier.predict(X_test)

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

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

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

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

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

input_data_encoded = [-1] * len(input_data)
count = 0
for i, item in enumerate(input_data):
    if item.isdigit():
        input_data_encoded[i] = int(input_data[i])
    else:
        encoder = label_encoder[count]
        input_data_encoded[i] = int(encoder.transform([input_data[i]])[-1])
        count += 1

input_data_encoded = np.array(input_data_encoded)
predicted_class = classifier.predict([input_data_encoded])
print(label_encoder[-1].inverse_transform(predicted_class)[0])

```

```
Accuracy score: 62.64%
C:\Users\ziraf\PycharmProjects\pythonProject\venv\lib\site-packages\sklearn\svm\_base.py:1244:
iterations.
warnings.warn(
C:\Users\ziraf\PycharmProjects\pythonProject\venv\lib\site-packages\sklearn\svm\_base.py:1244:
iterations.
warnings.warn(
C:\Users\ziraf\PycharmProjects\pythonProject\venv\lib\site-packages\sklearn\svm\_base.py:1244:
iterations.
warnings.warn(
Precision score: 69.18%
C:\Users\ziraf\PycharmProjects\pythonProject\venv\lib\site-packages\sklearn\svm\_base.py:1244:
iterations.
warnings.warn(
C:\Users\ziraf\PycharmProjects\pythonProject\venv\lib\site-packages\sklearn\svm\_base.py:1244:
iterations.
warnings.warn(
C:\Users\ziraf\PycharmProjects\pythonProject\venv\lib\site-packages\sklearn\svm\_base.py:1244:
iterations.
warnings.warn(
Recall score: 38.24%
C:\Users\ziraf\PycharmProjects\pythonProject\venv\lib\site-packages\sklearn\svm\_base.py:1244:
iterations.
warnings.warn(
C:\Users\ziraf\PycharmProjects\pythonProject\venv\lib\site-packages\sklearn\svm\_base.py:1244:
iterations.
warnings.warn(
C:\Users\ziraf\PycharmProjects\pythonProject\venv\lib\site-packages\sklearn\svm\_base.py:1244:
iterations.
warnings.warn(
F1 score: 56.15%
<=50K
```

Рис. 1 Результат виконання програми

Висновок: тестова точка належить до класу  $\leq 50k$ , тобто людина заробляє менше або рівно 50к.

## 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.

Лістинг програми:

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn import preprocessing
from sklearn.svm import SVC
from sklearn.multiclass import OneVsOneClassifier
from sklearn.model_selection import train_test_split, cross_val_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[:-1].split(', ')
        income_class = data[-1]
        if income_class == '<=50K' and count_class1 < max_datapoints:
            X.append(data)
            count_class1 += 1

        if income_class == '>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.isdigit():
        X_encoded[:, i] = X[:, i]
    else:
        current_label_encoder = preprocessing.LabelEncoder()
        label_encoder.append(current_label_encoder)
        X_encoded[:, i] = current_label_encoder.fit_transform(X[:, i])

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

classifier = OneVsOneClassifier(SVC(kernel='poly', degree=8, max_iter=5000))
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
                                                    random_state=5)
classifier.fit(X_train, y_train)
y_test_pred = classifier.predict(X_test)

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

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

```

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

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

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

input_data_encoded = [-1] * len(input_data)
count = 0
for i, item in enumerate(input_data):
    if item.isdigit():
        input_data_encoded[i] = int(input_data[i])
    else:
        encoder = label_encoder[count]
        input_data_encoded[i] = int(encoder.transform([input_data[i]])[-1])
        count += 1

input_data_encoded = np.array(input_data_encoded)
predicted_class = classifier.predict([input_data_encoded])
print(label_encoder[-1].inverse_transform(predicted_class)[0])

```

```

Accuracy score: 58.41%
C:\Users\ziraf\PycharmProjects\pythonProject\venv\lib\site-packages\sklearn\svm\_base.py:299:
pre-processing your data with StandardScaler or MinMaxScaler.
  warnings.warn(
C:\Users\ziraf\PycharmProjects\pythonProject\venv\lib\site-packages\sklearn\svm\_base.py:299:
pre-processing your data with StandardScaler or MinMaxScaler.
  warnings.warn(
C:\Users\ziraf\PycharmProjects\pythonProject\venv\lib\site-packages\sklearn\svm\_base.py:299:
pre-processing your data with StandardScaler or MinMaxScaler.
  warnings.warn(
Precision score: 41.6%
C:\Users\ziraf\PycharmProjects\pythonProject\venv\lib\site-packages\sklearn\svm\_base.py:299:
pre-processing your data with StandardScaler or MinMaxScaler.
  warnings.warn(
C:\Users\ziraf\PycharmProjects\pythonProject\venv\lib\site-packages\sklearn\svm\_base.py:299:
pre-processing your data with StandardScaler or MinMaxScaler.
  warnings.warn(
C:\Users\ziraf\PycharmProjects\pythonProject\venv\lib\site-packages\sklearn\svm\_base.py:299:
pre-processing your data with StandardScaler or MinMaxScaler.
  warnings.warn(
Recall score: 33.05%
C:\Users\ziraf\PycharmProjects\pythonProject\venv\lib\site-packages\sklearn\svm\_base.py:299:
pre-processing your data with StandardScaler or MinMaxScaler.
  warnings.warn(
C:\Users\ziraf\PycharmProjects\pythonProject\venv\lib\site-packages\sklearn\svm\_base.py:299:
pre-processing your data with StandardScaler or MinMaxScaler.
  warnings.warn(
C:\Users\ziraf\PycharmProjects\pythonProject\venv\lib\site-packages\sklearn\svm\_base.py:299:
pre-processing your data with StandardScaler or MinMaxScaler.
  warnings.warn(
F1 score: 46.5%
>50K

```

Рис. 2 Результат виконання програми

Лістинг програми:

```

import numpy as np
import matplotlib.pyplot as plt
from sklearn import preprocessing
from sklearn.svm import SVC
from sklearn.multiclass import OneVsOneClassifier
from sklearn.model_selection import train_test_split, cross_val_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[:-1].split(' ', 1)
        income_class = data[-1]
        if income_class == '<=50K' and count_class1 < max_datapoints:
            X.append(data)
            count_class1 += 1

        if income_class == '>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.isdigit():
        X_encoded[:, i] = X[:, i]
    else:
        current_label_encoder = preprocessing.LabelEncoder()
        label_encoder.append(current_label_encoder)
        X_encoded[:, i] = current_label_encoder.fit_transform(X[:, i])

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

classifier = OneVsOneClassifier(SVC(kernel='rbf', max_iter=10000))
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
                                                    random_state=5)
classifier.fit(X_train, y_train)
y_test_pred = classifier.predict(X_test)

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

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

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

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

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

input_data_encoded = [-1] * len(input_data)
count = 0
for i, item in enumerate(input_data):
    if item.isdigit():
        input_data_encoded[i] = int(input_data[i])
    else:
        encoder = label_encoder[count]
        input_data_encoded[i] = int(encoder.transform([input_data[i]])[-1])
        count += 1

input_data_encoded = np.array(input_data_encoded)
predicted_class = classifier.predict([input_data_encoded])
print(label_encoder[-1].inverse_transform(predicted_class)[0])

```



```
C:\Users\ziraf\PycharmProjects\pythonProject\venv\Scripts\python.exe
Accuracy score: 78.61%
Precision score: 98.72%
Recall score: 14.26%
F1 score: 71.95%
<=50K

Process finished with exit code 0
```

Рис. 3 Результат виконання програми

Лістинг програми:

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn import preprocessing
from sklearn.svm import SVC
from sklearn.multiclass import OneVsOneClassifier
from sklearn.model_selection import train_test_split, cross_val_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[:-1].split(', ')
        income_class = data[-1]
        if income_class == '<=50K' and count_class1 < max_datapoints:
            X.append(data)
            count_class1 += 1

        if income_class == '>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.isdigit():
        X_encoded[:, i] = X[:, i]
    else:
        current_label_encoder = preprocessing.LabelEncoder()
        label_encoder.append(current_label_encoder)
        X_encoded[:, i] = current_label_encoder.fit_transform(X[:, i])

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

```

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

classifier = OneVsOneClassifier(SVC(kernel='sigmoid', max_iter=10000))
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=5)
classifier.fit(X_train, y_train)
y_test_pred = classifier.predict(X_test)

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

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

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

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

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

input_data_encoded = [-1] * len(input_data)
count = 0
for i, item in enumerate(input_data):
    if item.isdigit():
        input_data_encoded[i] = int(input_data[i])
    else:
        encoder = label_encoder[count]
        input_data_encoded[i] = int(encoder.transform([input_data[i]])[-1])
        count += 1

input_data_encoded = np.array(input_data_encoded)
predicted_class = classifier.predict([input_data_encoded])
print(label_encoder[-1].inverse_transform(predicted_class)[0])

```

```

C:\Users\ziraf\PycharmProjects\pythonProject\venv\Scripts\python.exe
Accuracy score: 63.89%
Precision score: 27.01%
Recall score: 26.48%
F1 score: 63.77%
<=50K

Process finished with exit code 0

```

Рис. 4 Результат виконання програми

Висновок: за акуратністю і точністю найкращий вид класифікатора є нелінійний класифікатор SVM з гаусовим ядром, але за повнотою найкраще справився нелінійний класифікатор SVM з поліноміальним ядром. Загалом, найкраще виконує завдання класифікатор з гаусовим ядром.

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

Необхідно класифікувати сорти ірисів за деякими їх характеристиками: довжина та ширина пелюсток, а також довжина та ширина чашолистків.

Також, в наявності є вимірювання цих же характеристик ірисів, які раніше дозволили досвідченому експерту віднести їх до сортів: *setosa*, *versicolor* і *virginica*.

Використовувати класичний набір даних у машинному навчанні та статистиці - Iris. Він включений у модуль `datasets` бібліотеки `scikit-learn`.

Лістинг програми:

```
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
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
import numpy as np

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)

print(dataset.shape)
print(dataset.head(20))
print(dataset.describe())
print(dataset.groupby('class').size())

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

scatter_matrix(dataset)
pyplot.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=0.2, random_state=1)

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)
    names.append(name)
    print('%s: %f (%f)' % (name, cv_results.mean(), cv_results.std()))

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

model = SVC(gamma='auto')
model.fit(X_train, y_train)
predictions = model.predict(X_validation)

print(accuracy_score(y_validation, predictions))
print(confusion_matrix(y_validation, predictions))
print(classification_report(y_validation, predictions))

X_new = np.array([[5, 2.9, 1, 0.2]])
print("Форма массива X_new: {}".format(X_new.shape))

prediction = model.predict(X_new)
print("Прогноз: {}".format(prediction))
print("Спрогнозована мітка: {}".format(prediction[0]))

```

```

C:\Users\ziraf\PycharmProjects\pythonProject\venv\Scripts\python.exe "F:/nice stuff/lab 4 course/штучний інтелект/lab2/LR_2_task_3.py"
(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
5	5.4	3.9	1.7	0.4	Iris-setosa
6	4.6	3.4	1.4	0.3	Iris-setosa
7	5.0	3.4	1.5	0.2	Iris-setosa
8	4.4	2.9	1.4	0.2	Iris-setosa
9	4.9	3.1	1.5	0.1	Iris-setosa
10	5.4	3.7	1.5	0.2	Iris-setosa
11	4.8	3.4	1.6	0.2	Iris-setosa
12	4.8	3.0	1.4	0.1	Iris-setosa
13	4.3	3.0	1.1	0.1	Iris-setosa
14	5.8	4.0	1.2	0.2	Iris-setosa
15	5.7	4.4	1.5	0.4	Iris-setosa
16	5.4	3.9	1.3	0.4	Iris-setosa
17	5.1	3.5	1.4	0.3	Iris-setosa
18	5.7	3.8	1.7	0.3	Iris-setosa
19	5.1	3.8	1.5	0.3	Iris-setosa

Рис. 5 Результат виконання програми

```

      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
dtype: int64
LR: 0.941667 (0.065085)
LDA: 0.975000 (0.038188)
KNN: 0.958333 (0.041667)
CART: 0.950000 (0.040825)
NB: 0.950000 (0.055277)
SVM: 0.983333 (0.033333)
0.9666666666666667
[[11  0  0]
 [ 0 12  1]
 [ 0  0  6]]

```

Рис. 6 Результат виконання програми

```

              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

Форма масива X_new: (1, 4)
Прогноз: ['Iris-setosa']
Спрогнозована мітка: Iris-setosa

Process finished with exit code 0

```

Рис. 7 Результат виконання програми

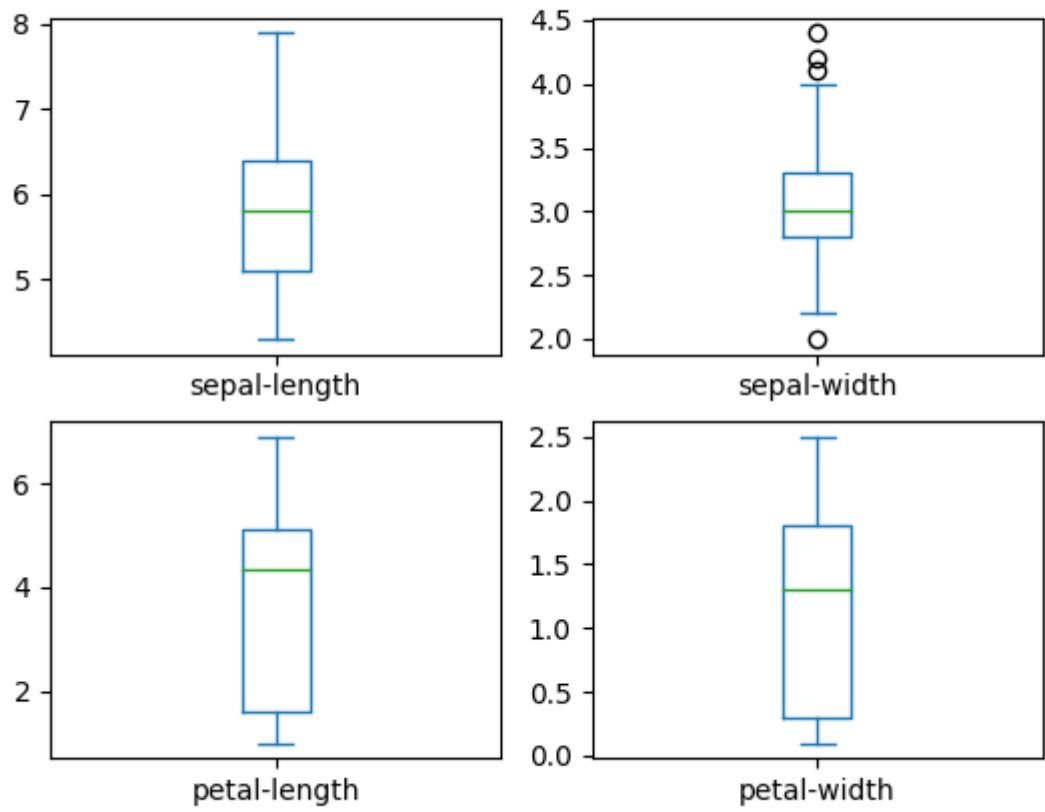


Рис. 8 Діаграма розмаху

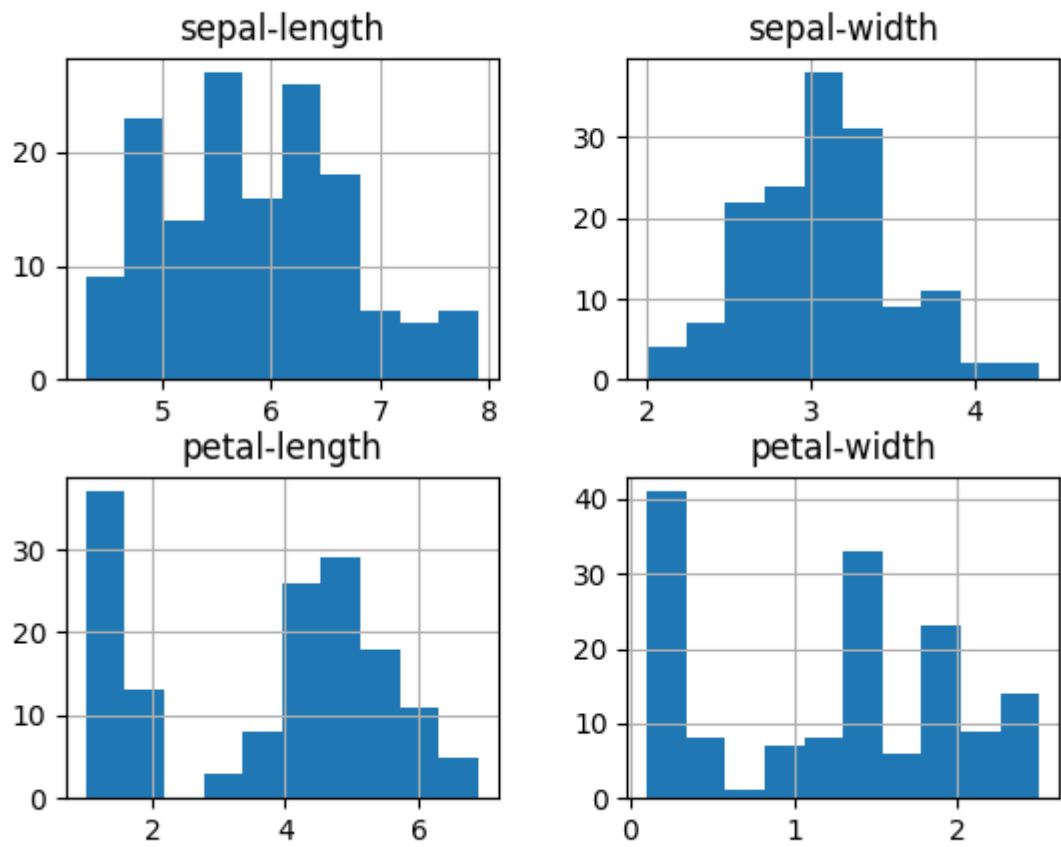


Рис. 9 Гістограма розподілу атрибутів датасета

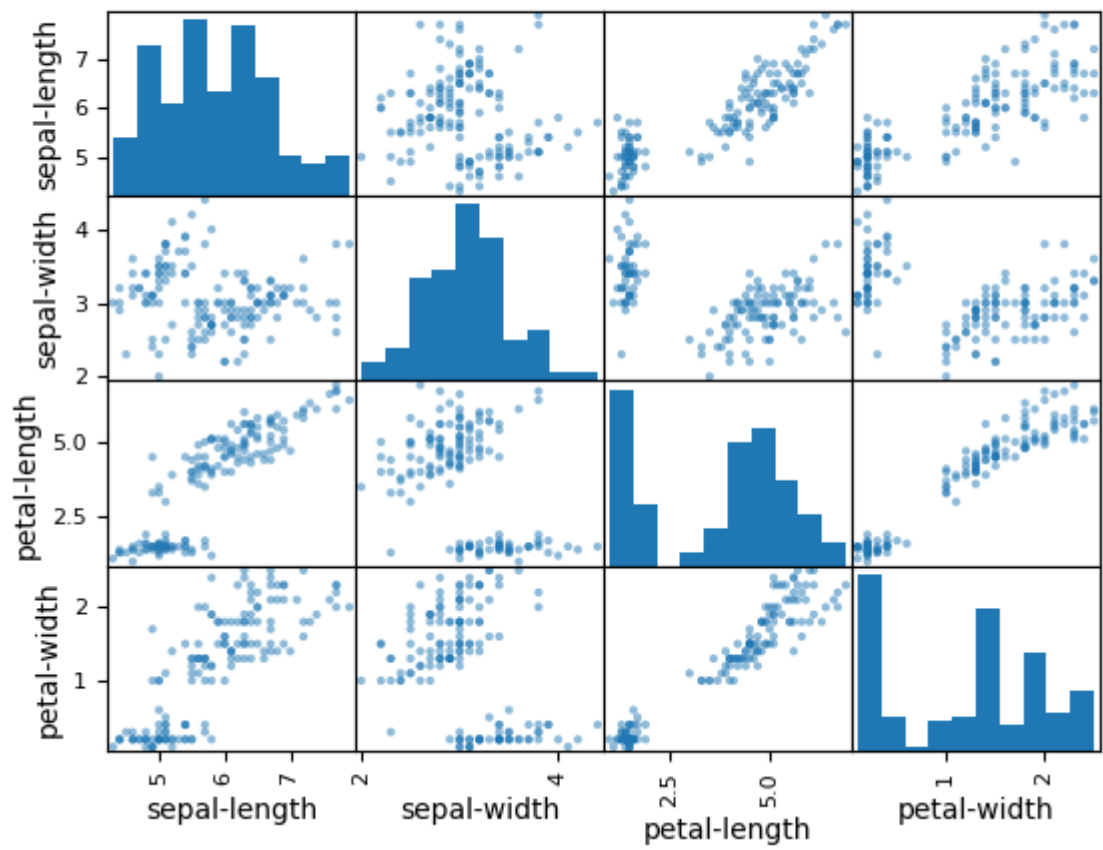


Рис. 10 Матриця діаграм розсіювання



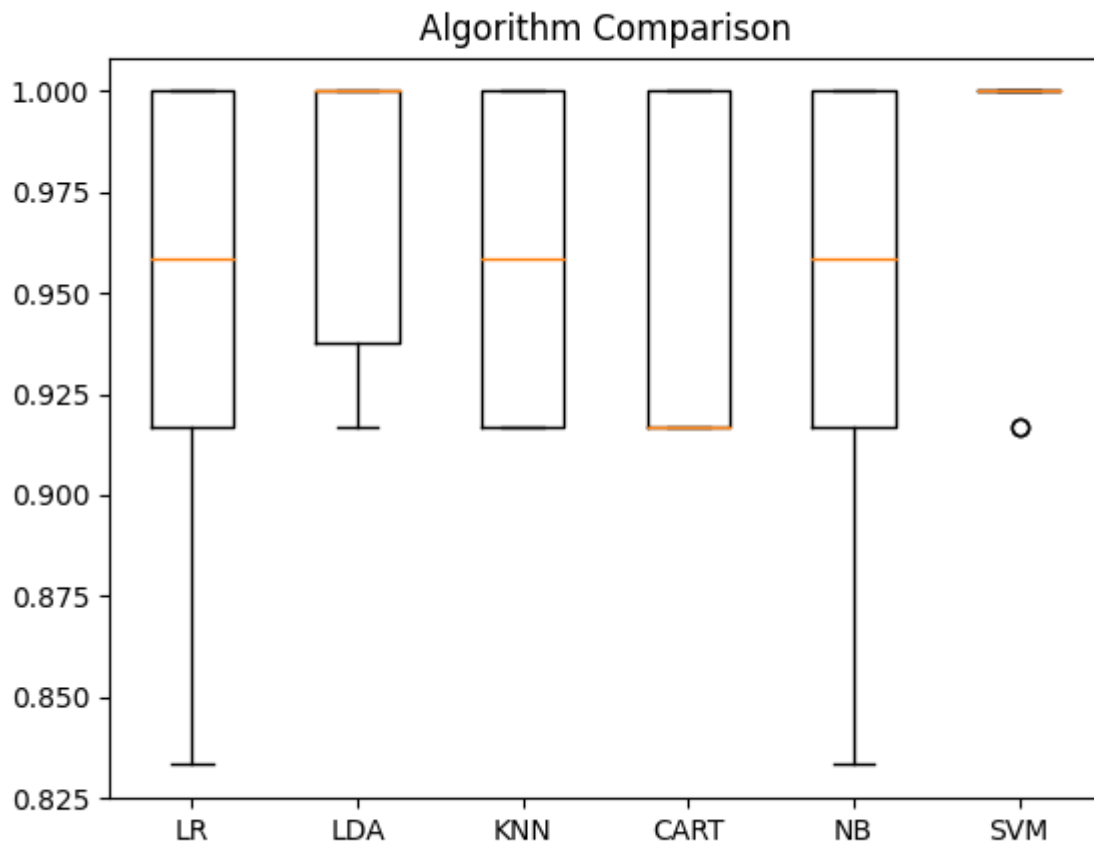


Рис. 11 Алгоритм порівняння

Висновок: Квітка належить до виду *Iris-setosa*. Було вибрано метод опорних векторів (SVM). Вдалося досягти 0.97 показника якості.

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

По аналогії із завданням 2.3 створіть код для порівняння якості класифікації набору даних `income_data.txt` (із завдання 2.1) різними алгоритмами.

Використати такі алгоритми класифікації:

Логістична регресія або логіт-модель (LR)

Лінійний дискримінантний аналіз (LDA)

Метод k-найближчих сусідів (KNN)

Класифікація та регресія за допомогою дерев (CART)

Наївний баєсовський класифікатор (NB)

Метод опорних векторів (SVM)

Розрахуйте показники якості класифікації для кожного алгоритму.

Порівняйте їх між собою. Оберіть найкращий для рішення задачі. Поясніть чому ви так вирішили у висновках до завдання.

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 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
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

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[:-1].split(', ')
        income_class = data[-1]
        if income_class == '<=50K' and count_class1 < max_datapoints:
            X.append(data)
            count_class1 += 1

        if income_class == '>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.isdigit():
        X_encoded[:, i] = X[:, i]
    else:
        current_label_encoder = preprocessing.LabelEncoder()
        label_encoder.append(current_label_encoder)
```

```

        X_encoded[:, i] = current_label_encoder.fit_transform(X[:, i])

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

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

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', max_iter=10000)))

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)
    names.append(name)
    print('%s: %f (%f)' % (name, cv_results.mean(), cv_results.std()))

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

```

```

C:\Users\ziraf\PycharmProjects\pythonProject\venv\Scripts\python.exe
LR: 0.791993 (0.005400)
LDA: 0.811637 (0.005701)
KNN: 0.767748 (0.003026)
CART: 0.806706 (0.007189)
NB: 0.789133 (0.006934)
SVM: 0.753492 (0.000328)

Process finished with exit code 0

```

Рис. 12 Результат виконання програми

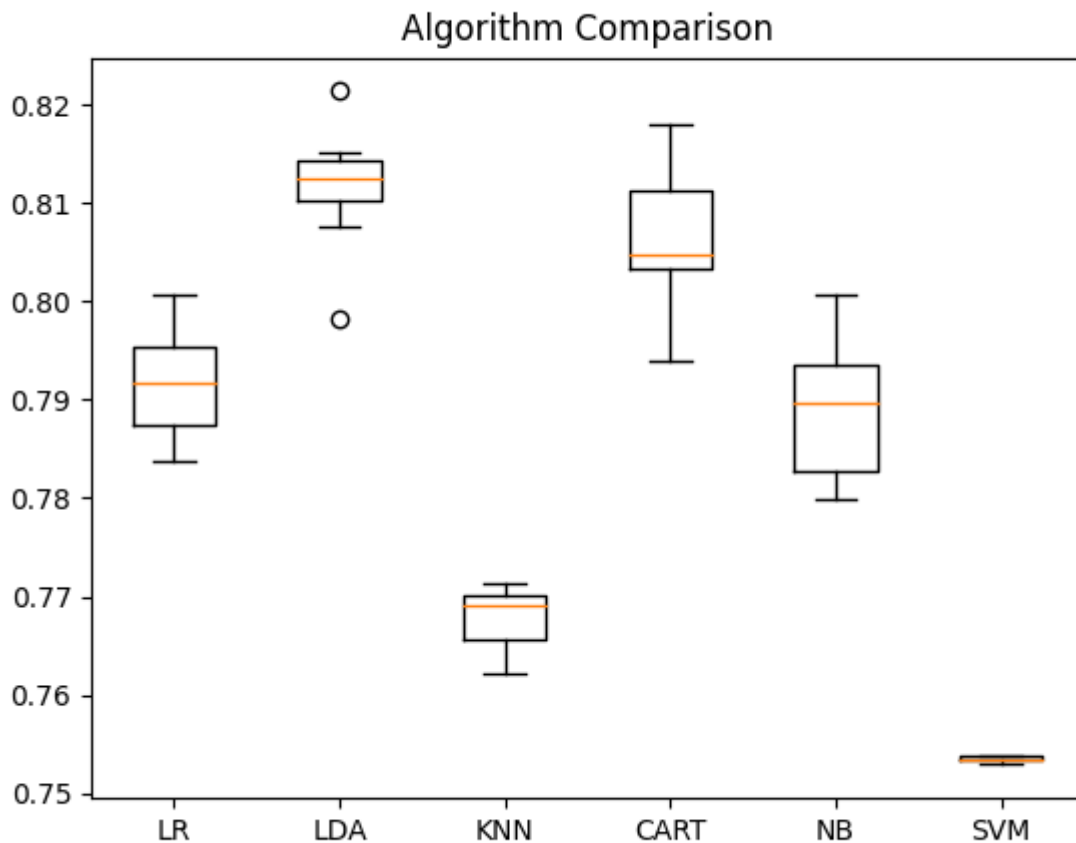


Рис. 13 Результат виконання програми

### Завдання 2.5. Класифікація даних лінійним класифікатором Ridge

Виправте код та виконайте класифікацію. Опишіть які налаштування класифікатора Ridge тут використані та що вони позначають. Опишіть які показники якості використовуються та їх отримані результати. Вставте у звіт та поясніть зображення Confusion.jpg. Опишіть, що таке коефіцієнт Коена Каппа та коефіцієнт кореляції Метьюза. Що вони тут розраховують та що показують.

Лістинг програми:

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

sns.set()
```

```
iris = load_iris()
X, y = iris.data, iris.target
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
random_state=0)
clf = RidgeClassifier(tol=1e-2, solver="sag")
clf.fit(X_train, y_train)
y_pred = clf.predict(X_test)

print('Accuracy:', np.round(metrics.accuracy_score(y_test, y_pred), 4))
print('Precision:', np.round(metrics.precision_score(y_test, y_pred,
average='weighted'), 4))
print('Recall:', np.round(metrics.recall_score(y_test, y_pred,
average='weighted'), 4))
print('F1 Score:', np.round(metrics.f1_score(y_test, y_pred, average='weighted'),
4))
print('Cohen Kappa Score:', np.round(metrics.cohen_kappa_score(y_test, y_pred),
4))
print('Matthews Corrcoef:', np.round(metrics.matthews_corrcoef(y_test, y_pred),
4))
print('\t\tClassification Report:\n', metrics.classification_report(y_pred,
y_test))

mat = confusion_matrix(y_test, y_pred)
sns.heatmap(mat.T, square=True, annot=True, fmt='d', cbar=False)
plt.xlabel('true label')
plt.ylabel('predicted label')
plt.savefig("Confusion.jpg")
# Save SVG in a fake file object.
f = BytesIO()
plt.savefig(f, format="svg")
```

```
C:\Users\ziraf\PycharmProjects\pythonProject\venv\Scripts\python.exe
Accuracy: 0.7556
Precision: 0.8333
Recall: 0.7556
F1 Score: 0.7503
Cohen Kappa Score: 0.6431
Matthews Corrcoef: 0.6831
      Classification Report:
              precision    recall  f1-score   support

         0              1.00      1.00      1.00        16
         1              0.44      0.89      0.59         9
         2              0.91      0.50      0.65        20

   accuracy              0.76              45
  macro avg              0.78      0.80      0.75        45
weighted avg              0.85      0.76      0.76        45

Process finished with exit code 0
```

Рис. 14 Результат виконання програми

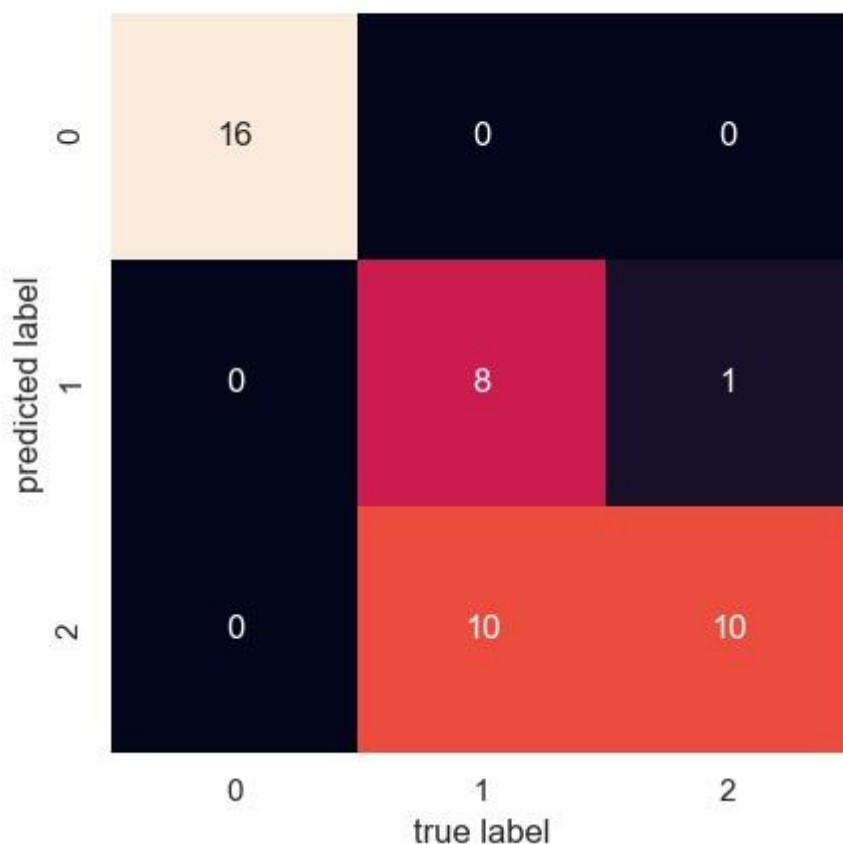


Рис. 15 Confusion.jpg

В класифікаторі Ridge були використані налаштування точності ( $\text{tol}=1\text{e-}2$ ) та розв'язник Stochastic Average Gradient descent ( $\text{solver}=\text{"sag"}$ ).

Показники якості, які використовувались – акуратність, точність, повнота, коефіцієнт Коена Каппа, коефіцієнт кореляції Метьюза.

На рис. 15 (Confusion.jpg) показана матриця confusion, як scikit-learn може навчатися класифікувати.

Коефіцієнт Коена Каппа – це статистика, яка вимірює міжрегіональну згоду якісних (категоріальних) предметів. Зазвичай вважається, що це надійніший захід, ніж простий розрахунок угоди про відсотки, оскільки  $k$  враховує випадкову угоду. Каппа Коена вимірює угоду між двома оцінювачами, кожен із яких класифікує  $N$  предметів на  $C$  взаємовиключних категорій.

В даному випадку коефіцієнт Коена Каппа (0.6431) показує істотну згоду.

Коефіцієнт кореляції Метьюза або коефіцієнт  $\rho_{hi}$  використовується в машинному навчанні як міра якості бінарних (двокласних) класифікацій, запроваджених біохіміком Браяном У. Метьюзом у 1975 році.

Не дивлячись на те, що акуратність точність і повнота в нас доволі високі, коефіцієнт кореляції Метьюза – 0.6831, тому що його оцінка висока в тих випадках, коли класифікатор справляється і з негативними, і з позитивними значеннями.

**Висновки:** на даній лабораторній ми, використовуючи спеціалізовані бібліотеки та мову програмування Python, дослідили різні методи класифікації даних та навчилися їх порівнювати, порівняли різні види класифікатора SVM, проаналізували значення коефіцієнтів Коена Каппа і кореляції Метьюза, порівняли якості класифікаторів на основі класифікації ірисів.