

## 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.

#### Хід роботи:

#### Завдання 1. Класифікація за допомогою машин опорних векторів (SVM)

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

Назва	Позначення	Вид
Вік	age	Числове
Працевлаштування	workclass	Категоріальне
Особи, які мають такі ж ознаки	fnlwgt	Числове
Освіта	education	Категоріальне
Роки навчання	education-num	Числове
Сімейне положення	marital-status	Категоріальне
Вид діяльності	occupation	Категоріальне
Взаємовідносини	relationship	Категоріальне
Раса	race	Категоріальне
Стать	sex	Категоріальне
Прирість капіталу	capital-gain	Категоріальне
Втрата капіталу	capital-loss	Числове
Кількість робочих годин на тиждень	hours-per-week	Числове
Рідна країна	native-country	Категоріальне

					Житомирська Політехніка.22.121.15.000 – Лр02					
Змн.	Арк.	№ докум.	Підпис	Дата	Звіт з лабораторної роботи			Літ.	Арк.	Аркушів
Розроб.		Сівченко О. О.								
Перевір.		Філіпов В. О.							1	23
Керівник								ФІКТ Гр. ПІ-60[2]		
Н. контр.										
Зав. каф.										

## 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

# Вхідний файл, який містить дані
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[:-1].split(',')
        if data[-1] == '<=50K' and count_class1 < max_datapoints:
            X.append(data)
            count_class1 += 1
        if data[-1] == '>50K' and count_class2 < max_datapoints:
            X.append(data)
            count_class2 += 1

# Перетворення на масив numpy
X = np.array(X)

# Перетворення рядкових даних на числові
label_encoder = []
X_encoded = np.empty(X.shape)
for i, item in enumerate(X[0]):
    if item.isdigit():
        X_encoded[:, i] = X[:, i]
    else:
        label_encoder.append(preprocessing.LabelEncoder())
        X_encoded[:, i] = label_encoder[-1].fit_transform(X[:, i])
X = X_encoded[:, :-1].astype(int)
y = X_encoded[:, -1].astype(int)

# Створення SVM-класифікатора
classifier = OneVsOneClassifier(LinearSVC(random_state=0))

# Навчання класифікатора
classifier.fit(X, y)

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

# Обчислення F-міри для SVM-класифікатора
```

		Сівченко О. О.			Житомирська Політехніка. 22.121.15.000 – Лр02	Арк.
		Філіпов В. О.				2
Змн.	Арк.	№ докум.	Підпис	Дата		

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

# Передбачення результату для тестової 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 даних
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)

# Використання класифікатора для кодованої точки даних
# та виведення результату
predicted_class = classifier.predict([input_data_encoded])
print(label_encoder[-1].inverse_transform(predicted_class)[0])

accuracy_values = cross_val_score(classifier, X, y, scoring='accuracy', cv=3)
print("Accuracy: " + str(round(100 * accuracy_values.mean(), 2)) + "%")
precision_values = cross_val_score(classifier, X, y, scoring='precision_weighted',
cv=3)
print("Precision: " + str(round(100 * precision_values.mean(), 2)) + "%")
recall_values = cross_val_score(classifier, X, y, scoring='recall_weighted', cv=3)
print("Recall: " + str(round(100 * recall_values.mean(), 2)) + "%")
```

F1 score: 56.15%  
 <=50K

Accuracy: 62.64%

Precision: 75.88%

Recall: 62.64%

Рисунок 1.1 – Результат роботи програми

Тестова точка належить до класу - <=50K

## Завдання 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;

		Сівченко О. О.			Житомирська Політехніка.22.121.15.000 – Лр02	Арк.
		Філіпов В. О.				
Змн.	Арк.	№ докум.	Підпис	Дата		3

- 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

# Вхідний файл, який містить дані
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[:-1].split(',')
        if data[-1] == '<=50K' and count_class1 < max_datapoints:
            X.append(data)
            count_class1 += 1
        if data[-1] == '>50K' and count_class2 < max_datapoints:
            X.append(data)
            count_class2 += 1

# Перетворення на масив numpy
X = np.array(X)

# Перетворення рядкових даних на числові
label_encoder = []
X_encoded = np.empty(X.shape)
for i, item in enumerate(X[0]):
    if item.isdigit():
        X_encoded[:, i] = X[:, i]
    else:
        label_encoder.append(preprocessing.LabelEncoder())
        X_encoded[:, i] = label_encoder[-1].fit_transform(X[:, i])
X = X_encoded[:, :-1].astype(int)
y = X_encoded[:, -1].astype(int)

# Створення SVM-класифікатора
classifier = OneVsOneClassifier(SVC(kernel='poly', degree=8, max_iter=10000))

# Навчання класифікатора
classifier.fit(X, y)
```

		Сівченко О. О.			Житомирська Політехніка. 22.121.15.000 – Лр02	Арк.
		Філіпов В. О.				4
Змн.	Арк.	№ докум.	Підпис	Дата		

```

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=5) # ??? cross_validation
classifier = OneVsOneClassifier(SVC(kernel='poly', degree=8, max_iter=10000))
classifier.fit(X_train, y_train)
y_test_pred = classifier.predict(X_test)

# Обчислення F-міри для SVM-класифікатора
f1 = cross_val_score(classifier, X, y, scoring='f1_weighted', cv=3)
print("F1 score: " + str(round(100*f1.mean(), 2)) + "%")

# Передбачення результату для тестової 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 даних
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)

# Використання класифікатора для кодованої точки даних
# та виведення результату
predicted_class = classifier.predict([input_data_encoded])
print(label_encoder[-1].inverse_transform(predicted_class)[0])

accuracy_values = cross_val_score(classifier, X, y, scoring='accuracy', cv=3)
print("Accuracy: " + str(round(100 * accuracy_values.mean(), 2)) + "%")
precision_values = cross_val_score(classifier, X, y, scoring='precision_weighted',
cv=3)
print("Precision: " + str(round(100 * precision_values.mean(), 2)) + "%")
recall_values = cross_val_score(classifier, X, y, scoring='recall_weighted', cv=3)
print("Recall: " + str(round(100 * recall_values.mean(), 2)) + "%")

```

F1 score: 64.51%  
>50K

Accuracy: 75.12%

Precision: 70.07%

Recall: 75.12%

Рисунок 1.2 – Результат роботи програми(з поліноміальним ядром)

		Сівченко О. О.			Житомирська Політехніка.22.121.15.000 – Лр02	Арк.
		Філіпов В. О.				5
Змн.	Арк.	№ докум.	Підпис	Дата		

## 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

# Вхідний файл, який містить дані
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[:-1].split(',')
        if data[-1] == '<=50K' and count_class1 < max_datapoints:
            X.append(data)
            count_class1 += 1
        if data[-1] == '>50K' and count_class2 < max_datapoints:
            X.append(data)
            count_class2 += 1

# Перетворення на масив numpy
X = np.array(X)

# Перетворення рядкових даних на числові
label_encoder = []
X_encoded = np.empty(X.shape)
for i, item in enumerate(X[0]):
    if item.isdigit():
        X_encoded[:, i] = X[:, i]
    else:
        label_encoder.append(preprocessing.LabelEncoder())
        X_encoded[:, i] = label_encoder[-1].fit_transform(X[:, i])
X = X_encoded[:, :-1].astype(int)
y = X_encoded[:, -1].astype(int)

# Створення SVM-класифікатора
classifier = OneVsOneClassifier(SVC(kernel='rbf', max_iter=10000))

# Навчання класифікатора
classifier.fit(X, y)

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

		Сівченко О. О.			Житомирська Політехніка.22.121.15.000 – Пр02	Арк.
		Філіпов В. О.				6
Змн.	Арк.	№ докум.	Підпис	Дата		

```

# Обчислення F-міри для SVM-класифікатора
f1 = cross_val_score(classifier, X, y, scoring='f1_weighted', cv=3)
print("F1 score: " + str(round(100*f1.mean(), 2)) + "%")

# Передбачення результату для тестової 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 даних
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)

# Використання класифікатора для кодованої точки даних
# та виведення результату
predicted_class = classifier.predict([input_data_encoded])
print(label_encoder[-1].inverse_transform(predicted_class)[0])

accuracy_values = cross_val_score(classifier, X, y, scoring='accuracy', cv=3)
print("Accuracy: " + str(round(100 * accuracy_values.mean(), 2)) + "%")
precision_values = cross_val_score(classifier, X, y, scoring='precision_weighted',
cv=3)
print("Precision: " + str(round(100 * precision_values.mean(), 2)) + "%")
recall_values = cross_val_score(classifier, X, y, scoring='recall_weighted', cv=3)
print("Recall: " + str(round(100 * recall_values.mean(), 2)) + "%")

```

```

F1 score: 71.95%
<=50K
Accuracy: 78.61%
Precision: 83.06%
Recall: 78.61%

```

Рисунок 1.3 – Результат роботи програми(з гаусовим ядром)

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

```

		Сівченко О. О.			Житомирська Політехніка.22.121.15.000 – Лр02	Арк.
		Філіпов В. О.				7
Змн.	Арк.	№ докум.	Підпис	Дата		

```

# Вхідний файл, який містить дані
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[:-1].split(',')
        if data[-1] == '<=50K' and count_class1 < max_datapoints:
            X.append(data)
            count_class1 += 1
        if data[-1] == '>50K' and count_class2 < max_datapoints:
            X.append(data)
            count_class2 += 1

# Перетворення на масив numpy
X = np.array(X)

# Перетворення рядкових даних на числові
label_encoder = []
X_encoded = np.empty(X.shape)
for i, item in enumerate(X[0]):
    if item.isdigit():
        X_encoded[:, i] = X[:, i]
    else:
        label_encoder.append(preprocessing.LabelEncoder())
        X_encoded[:, i] = label_encoder[-1].fit_transform(X[:, i])
X = X_encoded[:, :-1].astype(int)
y = X_encoded[:, -1].astype(int)

# Створення SVM-класифікатора
classifier = OneVsOneClassifier(SVC(kernel='sigmoid'))

# Навчання класифікатора
classifier.fit(X, y)

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

# Обчислення F-міри для SVM-класифікатора
f1 = cross_val_score(classifier, X, y, scoring='f1_weighted', cv=3)
print("F1 score: " + str(round(100*f1.mean(), 2)) + "%")

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

```

		Сівченко О. О.			Житомирська Політехніка.22.121.15.000 – Лр02	Арк.
		Філіпов В. О.				8
Змн.	Арк.	№ докум.	Підпис	Дата		



```
# Кодування тестової точки даних
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)

# Використання класифікатора для кодованої точки даних
# та виведення результату
predicted_class = classifier.predict([input_data_encoded])
print(label_encoder[-1].inverse_transform(predicted_class)[0])

accuracy_values = cross_val_score(classifier, X, y, scoring='accuracy', cv=3)
print("Accuracy: " + str(round(100 * accuracy_values.mean(), 2)) + "%")
precision_values = cross_val_score(classifier, X, y, scoring='precision_weighted',
cv=3)
print("Precision: " + str(round(100 * precision_values.mean(), 2)) + "%")
recall_values = cross_val_score(classifier, X, y, scoring='recall_weighted', cv=3)
print("Recall: " + str(round(100 * recall_values.mean(), 2)) + "%")
```

```
F1 score: 63.77%
<=50K
Accuracy: 63.89%
Precision: 63.65%
Recall: 63.89%
```

Рисунок 1.4 – Результат роботи програми(з сигмоїдальним ядром)

За результатами тренування класифікатор SVM з гаусовим ядром найкраще виконує завдання класифікації.

### Завдання 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.

		Сівченко О. О.			Житомирська Політехніка.22.121.15.000 – Пр02	Арк.
		Філіпов В. О.				
Змн.	Арк.	№ докум.	Підпис	Дата		9

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

## КРОК 1. ЗАВАНТАЖЕННЯ ТА ВИВЧЕННЯ ДАНИХ

Лістинг програми:

```
from sklearn.datasets import load_iris
iris_dataset = load_iris()

print("Ключі iris_dataset: \n{}".format(iris_dataset.keys()))
print(iris_dataset['DESCR'][:5000] + "\n..." # 5000

print("Назви відповідей: {}".format(iris_dataset['target_names']))
print("Назва ознак: {}".format(iris_dataset['feature_names']))

print("Тип масиву data: {}".format(type(iris_dataset['data'])))
print("Форма масиву data: {}".format(iris_dataset['data'].shape))

print("Ознаки для 5 прикладів: \n{}".format((iris_dataset['data'][:5])))

print("Тип масиву target:{}".format(type(iris_dataset['target'])))
print("Відповіді:\n{}".format(iris_dataset['target']))
```

```
dict_keys(['data', 'target', 'frame', 'target_names', 'DESCR', 'feature_names', 'filename', 'data_module'])
.. _iris_dataset:

Iris plants dataset
-----

**Data Set Characteristics:**

 :Number of Instances: 150 (50 in each of three classes)
 :Number of Attributes: 4 numeric, predictive attributes and the class
 :Attribute Information:
   - sepal length in cm
   - sepal width in cm
   - petal length in cm
   - petal width in cm
   - class:
     - Iris-Setosa
     - Iris-Versicolour
     - Iris-Virginica
```

Рисунок 1.5 – Результат роботи програми

		Сівченко О. О.			Житомирська Політехніка.22.121.15.000 – Лр02	Арк.
		Філіпов В. О.				10
Змн.	Арк.	№ докум.	Підпис	Дата		

```

:Summary Statistics:

=====
              Min   Max    Mean    SD    Class Correlation
=====
sepal length:  4.3  7.9   5.84   0.83    0.7826
sepal width:   2.0  4.4   3.05   0.43   -0.4194
petal length:  1.0  6.9   3.76   1.76    0.9490 (high!)
petal width:   0.1  2.5   1.20   0.76    0.9565 (high!)
=====

:Missing Attribute Values: None
:Class Distribution: 33.3% for each of 3 classes.
:Creator: R.A. Fisher
:Donor: Michael Marshall (MARSHALL%PLU@io.arc.nasa.gov)
:Date: July, 1988

The famous Iris database, first used by Sir R.A. Fisher. The dataset is taken
from Fisher's paper. Note that it's the same as in R, but not as in the UCI
Machine Learning Repository, which has two wrong data points.

This is perhaps the best known database to be found in the
pattern recognition literature. Fisher's paper is a classic in the field and
is referenced frequently to this day. (See Duda & Hart, for example.) The
data set contains 3 classes of 50 instances each, where each class refers to a
type of iris plant. One class is linearly separable from the other 2; the
latter are NOT linearly separable from each other.

```

Рисунок 1.6 – Результат роботи програми(продовження)

```

.. topic:: References

- Fisher, R.A. "The use of multiple measurements in taxonomic problems"
  Annual Eugenics, 7, Part II, 179-188 (1936); also in "Contributions to
  Mathematical Statistics" (John Wiley, NY, 1950).
- Duda, R.O., & Hart, P.E. (1973) Pattern Classification and Scene Analysis.
  (Q327.D83) John Wiley & Sons. ISBN 0-471-22361-1. See page 218.
- Dasarathy, B.V. (1980) "Nosing Around the Neighborhood: A New System
  Structure and Classification Rule for Recognition in Partially Exposed
  Environments". IEEE Transactions on Pattern Analysis and Machine
  Intelligence, Vol. PAMI-2, No. 1, 67-71.
- Gates, G.W. (1972) "The Reduced Nearest Neighbor Rule". IEEE Transactions
  on Information Theory, May 1972, 431-433.
- See also: 1988 MLC Proceedings, 54-64. Cheeseman et al's AUTOCLASS II
  conceptual clustering system finds 3 classes in the data.
- Many, many more ...

...

```

Рисунок 1.7 – Результат роботи програми(продовження)

[illegible]

Рисунок 1.8 – Результат роботи програми(продовження)

## КРОК 2. ВІЗУАЛІЗАЦІЯ ДАНИХ

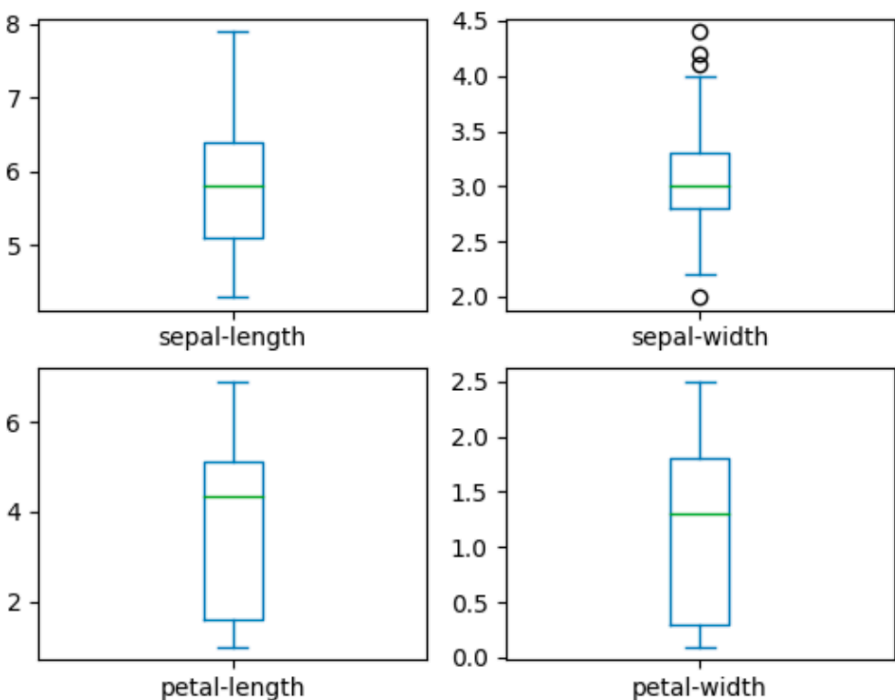


Рисунок 1.9 – Діаграма розмаху атрибутів вхідних даних

		Сівченко О. О.			Житомирська Політехніка.22.121.15.000 – Лр02	Арк.
		Філіпов В. О.				12
Змн.	Арк.	№ докум.	Підпис	Дата		

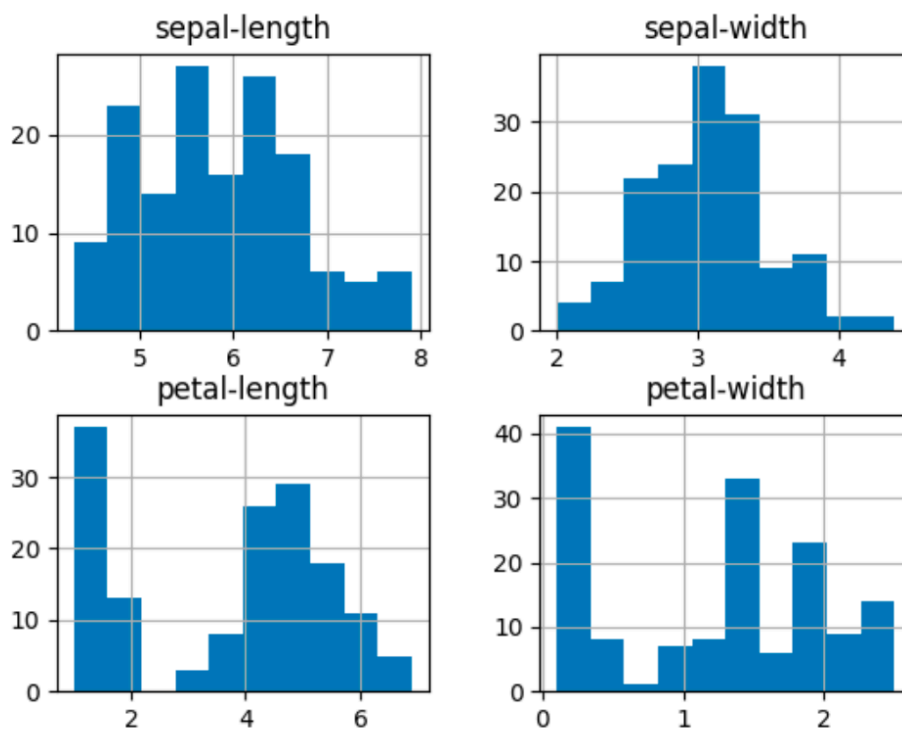


Рисунок 2.1 – Гістограма розподілу атрибутів вхідних даних

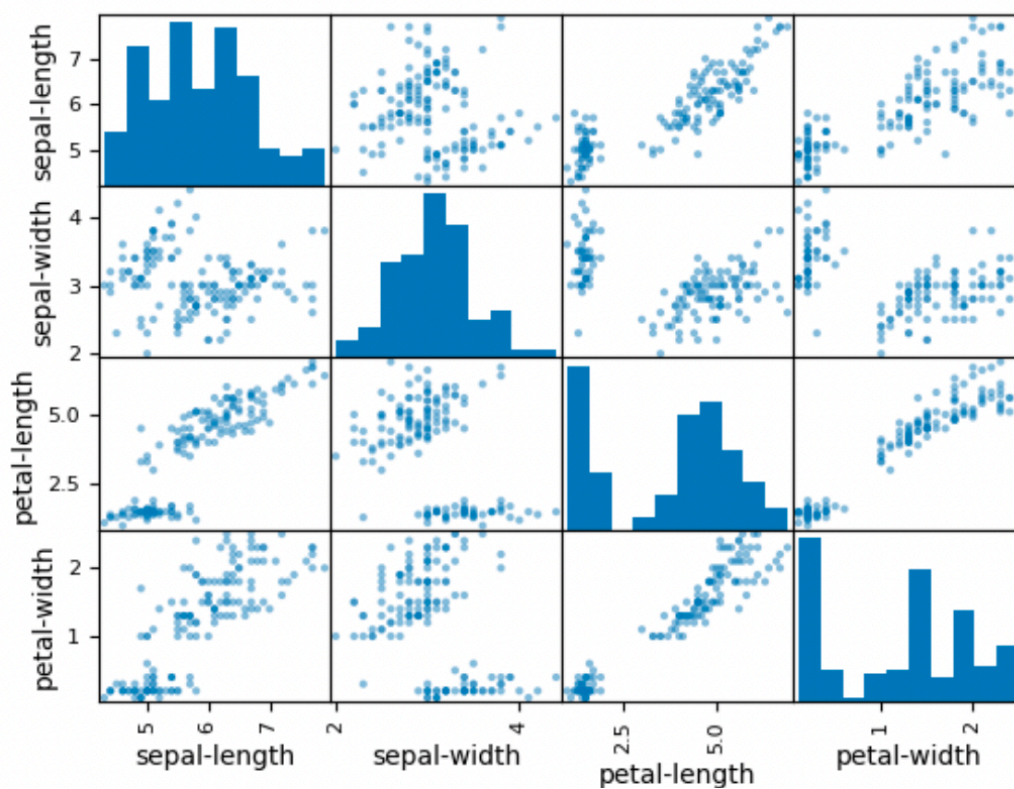


Рисунок 2.2 – Матриця діаграм розсіювання

### Лістинг програми:

```
# Завантаження бібліотек
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

# Завантаження датасету
url = "iris.csv"
names = ['sepal-length', 'sepal-width', 'petal-length', 'petal-width', 'class']
dataset = read_csv(url, names=names)

# shape
print(dataset.shape)

# Зріз даних head
print(dataset.head(20))

# Статистичні зведення методом describe
print(dataset.describe())

# Розподіл за атрибутом class
print(dataset.groupby('class').size())

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

# Гістограма розподілу атрибутів датасета
dataset.hist()
pyplot.show()

# Матриця діаграм розсіювання
scatter_matrix(dataset)
pyplot.show()
```

### КРОК 3. СТВОРЕННЯ НАВЧАЛЬНОГО ТА ТЕСТОВОГО НАБОРІВ

#### Лістинг програми:

```
# Розділення датасету на навчальну та контрольну вибірки
array = dataset.values
# Вибір перших 4-х стовпців
X = array[:,0:4]
# Вибір 5-го стовпця
y = array[:,4]
# Розділення X и y на навчающую и контрольную выборки
X_train, X_validation, Y_train, Y_validation = train_test_split(X, y,
test_size=0.20, random_state=1)
```

		Сівченко О. О.			Житомирська Політехніка. 22.121.15.000 – Пр02	Арк.
		Філіпов В. О.				14
Змн.	Арк.	№ докум.	Підпис	Дата		

## КРОК 4. КЛАСИФІКАЦІЯ (ПОБУДОВА МОДЕЛІ)

Лістинг програми:

```
# Завантажуємо алгоритми моделі
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')))

# оцінюємо модель на кожній ітерації
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()))

# Порівняння алгоритмів
pyplot.boxplot(results, labels=names)
pyplot.title('Algorithm Comparison')
pyplot.show()
```

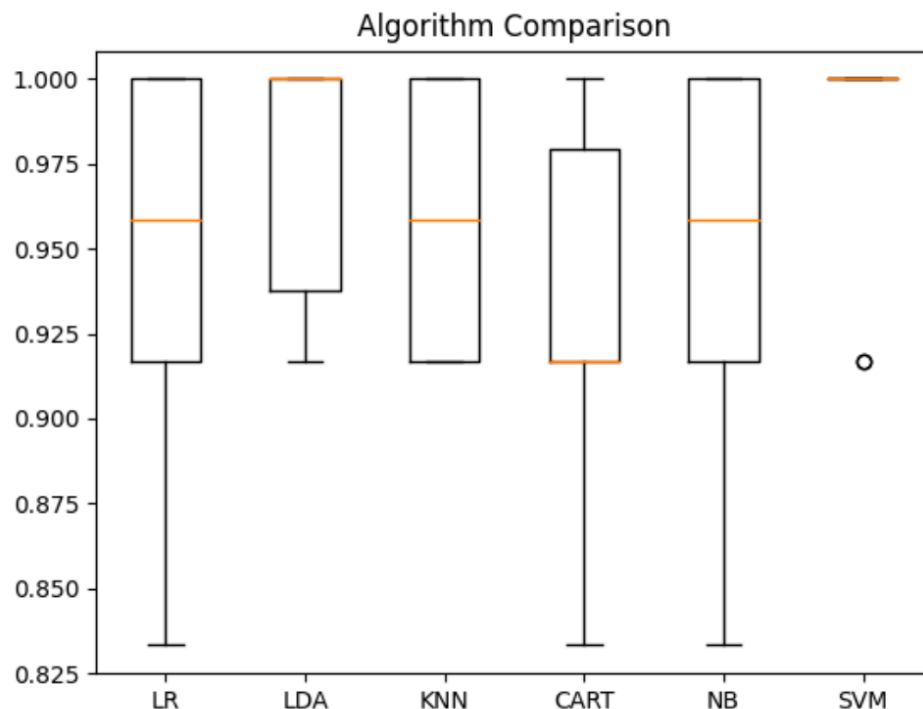


Рисунок 2.3 – Діаграма розмаху атрибутів вихідних даних та їх вусів для кожного розподілу та порівняння розподілів

		Сівченко О. О.			Житомирська Політехніка.22.121.15.000 – Лр02	Арк.
		Філіпов В. О.				15
Змн.	Арк.	№ докум.	Підпис	Дата		

```

LR: 0.941667 (0.065085)
LDA: 0.975000 (0.038188)
KNN: 0.958333 (0.041667)
CART: 0.933333 (0.050000)
NB: 0.950000 (0.055277)
SVM: 0.983333 (0.033333)

```

Рисунок 2.4 – Результат роботи програми

Порівнявши отримані результати кожного з розподілів можна зробити висновок, що найкращий метод це Метод опорних векторів(SVM), тому що згідно графіку його точність краща, а розмах менший.

## КРОК 5. ОПТИМІЗАЦІЯ ПАРАМЕТРІВ МОДЕЛІ

-

## КРОК 6. ОТРИМАННЯ ПРОГНОЗУ (ПЕРЕДБАЧЕННЯ НА ТРЕНУВАЛЬНОМУ НАБОРІ)

Лістинг програми:

```

# Створюємо прогноз на контрольній вибірці
model = SVC(gamma='auto')
model.fit(X_train, Y_train)
predictions = model.predict(X_validation)

```

## КРОК 7. ОЦІНКА ЯКОСТІ МОДЕЛІ

Лістинг програми:

```

# Оцінюємо прогноз
print(accuracy_score(Y_validation, predictions))
print(confusion_matrix(Y_validation, predictions))
print(classification_report(Y_validation, predictions))

```

		Сівченко О. О.			Житомирська Політехніка.22.121.15.000 – Лр02	Арк.
		Філіпов В. О.				16
Змн.	Арк.	№ докум.	Підпис	Дата		



```
0.9666666666666667
[[11  0  0]
 [ 0 12  1]
 [ 0  0  6]]
```

	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

Рисунок 2.5 – Результат роботи програми

## КРОК 8. ОТРИМАННЯ ПРОГНОЗУ (ЗАСТОСУВАННЯ МОДЕЛІ ДЛЯ ПЕРЕДБАЧЕННЯ)

Лістинг програми:

```
import numpy as np

model2 = SVC(gamma='auto')
model2.fit(X_train, Y_train)

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

prediction = model2.predict(X_new)
print("Прогноз: {}".format(prediction))
print("Спрогнозована мітка: {}".format(prediction[0]))
```

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

Рисунок 2.6 – Результат роботи програми

За результатами тренування вдалося досягти таку якість класифікації – 0.9666666666666667. Клас, до якого належить квітка, це Iris-setosa.

		Сівченко О. О.			Житомирська Політехніка.22.121.15.000 – Лр02	Арк.
		Філіпов В. О.				17
Змн.	Арк.	№ докум.	Підпис	Дата		

## Завдання 4. Порівняння якості класифікаторів для набору даних

### завдання 1.

Лістинг програми:

```
# Завантаження бібліотек
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.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
from sklearn import preprocessing

# Вхідний файл, який містить дані
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[:-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

# Перетворення на масив numpy
X = np.array(X)

# Перетворення рядкових даних на числові
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)
```

		Сівченко О. О.			Житомирська Політехніка.22.121.15.000 – Лр02	Арк.
		Філіпов В. О.				18
Змн.	Арк.	№ докум.	Підпис	Дата		

```

# Разделение X и y на обучающую и контрольную выборки
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=5)

# Завантажуємо алгоритми моделі
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')))

# оцінюємо модель на кожній ітерації
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()))

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

```

```

LR: 0.791661 (0.005010)
LDA: 0.811637 (0.005701)
KNN: 0.767707 (0.003104)
CART: 0.807161 (0.007486)
NB: 0.789133 (0.006934)
SVM: 0.753119 (0.000378)

```

Рисунок 2.7 – Результат роботи програми

		Сівченко О. О.			Житомирська Політехніка.22.121.15.000 – Пр02	Арк.
		Філіпов В. О.				19
Змн.	Арк.	№ докум.	Підпис	Дата		

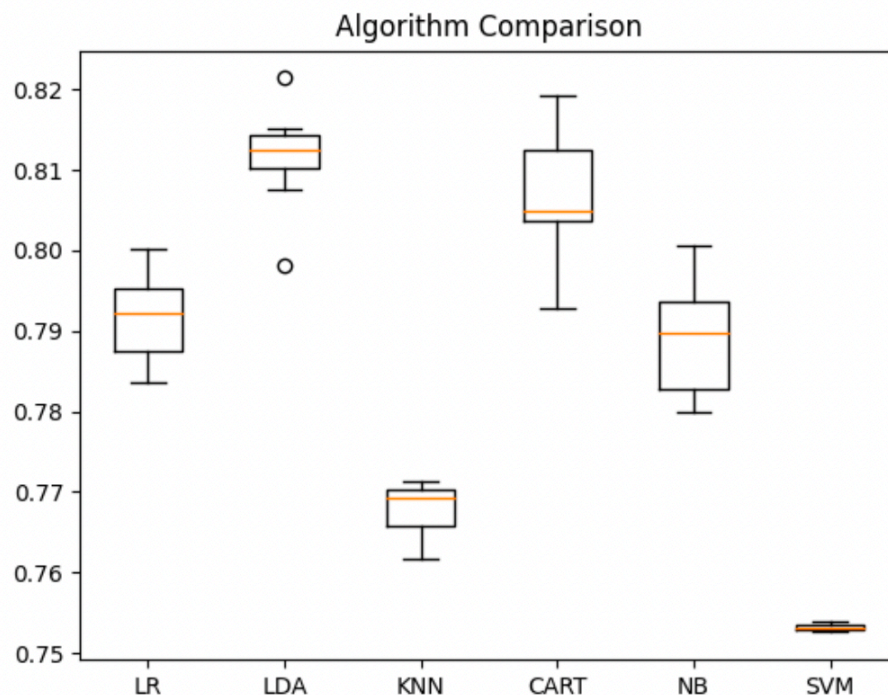


Рисунок 2.8 – Діаграма порівняння розмаху атрибутів вихідних даних та розподілів

Порівнявши отримані результати кожного з розподілів можна зробити висновок, що найкращий метод це Логістична регресія(LDA), тому що згідно графіку його точність більша, а розмах менший.

### Завдання 5. Класифікація даних лінійним класифікатором Ridge

Лістинг програми:

```
# Приклад класифікатора Ridge
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 # needed for plot
import seaborn as sns; sns.set()
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split

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,
```

		Сівченко О. О.			Житомирська Політехніка.22.121.15.000 – Лр02	Арк.
		Філіпов В. О.				20
Змн.	Арк.	№ докум.	Підпис	Дата		

```

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")

```

```

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

Рисунок 2.9 – Результат роботи програми

		Сівченко О. О.			Житомирська Політехніка.22.121.15.000 – Лр02	Арк.
		Філіпов В. О.				21
Змн.	Арк.	№ докум.	Підпис	Дата		

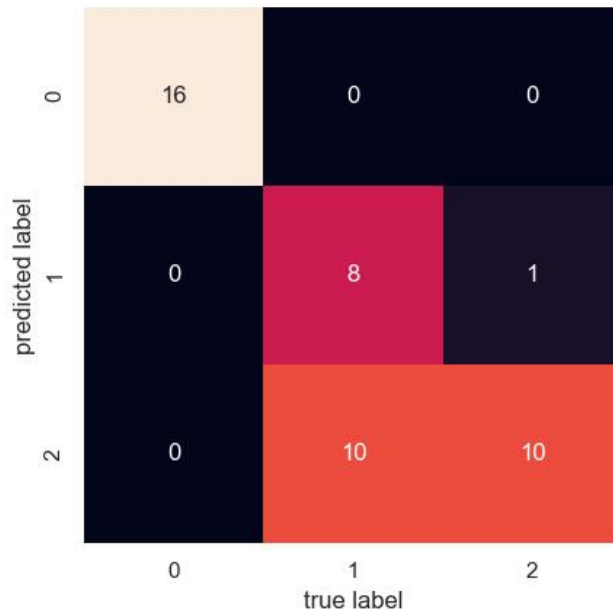


Рисунок 3.1 – Зображення матриці плутанини

У налаштуванні класифікатора Ridge тут використані такі параметри:  $\text{tol}=1\text{e-}2$  – точність рішення,  $\text{solver}=\text{"sag"}$  – вирішувач для використання в обчислювальних програмах, а саме "sag" використовує стохастичний середній градієнтний спуск. Метод використовує ітераційну процедуру і часто є швидшими, ніж інші розв'язувачі, коли  $n\_samples$  і  $n\_features$  великі.

Показники якості, які тут використовуються та їх отримані результати показано на рисунку 2.9. Accuracy – якість, тобто частка вибірок, що правильно спрогнозовані, precision – точність це частка очікуваних позитивних подій, які є позитивними, recall – повнота, що є часткою позитивних подій, які ви правильно передбачили. Оцінка F1 є гармонійним середнім значенням повноти та точності, з більш високою оцінкою як краща модель.

Матриця плутанини обчислюється, щоб оцінити точність класифікації. Ми задаємо основні істинні (правильні) цільові значення та розрахункові цілі, отримані класифікатором.

Коефіцієнт Коєна Каппа обчислює оцінку, яка виражає рівень згоди між двома анотаторами щодо проблеми класифікації.

Коефіцієнт кореляції Метьюза використовується в машинному навчанні як міра якості бінарних і багатокласових класифікацій. Він враховує істинні та хибні,

		Сівченко О. О.			Житомирська Політехніка.22.121.15.000 – Лр02	Арк.
		Філіпов В. О.				22
Змн.	Арк.	№ докум.	Підпис	Дата		

позитивні та негативні показники і, як правило, розглядається як збалансований показник, який можна використовувати, навіть якщо класи дуже різних розмірів.

Посилання на GitHub: <https://github.com/SoylerProfile/SHI>

**Висновки:** в ході виконання лабораторної роботи було досліджено різні методи класифікації даних, використовуючи спеціалізовані бібліотеки та мову програмування Python та навчено їх порівнювати.

		Сівченко О. О.			Житомирська Політехніка.22.121.15.000 – Лр02	Арк.
		Філіпов В. О.				23
Змн.	Арк.	№ докум.	Підпис	Дата		