

Лабораторна робота №2
Тема: Порівняння методів класифікації даних

Завдання 1: Класифікація за допомогою машин опорних векторів SVM
Посилання на гіт хаб: <https://github.com/ViMIMercurysMight/python-ai.git>

На початку роботи випишімо атрибути що надаються з набору даних та визначимо їх тип (таб. 1).

Таблиця 1. Опис атрибутів

Атрибут	Опис	Тип
Вік	Вік	Числова
Тип зайнятості	Тип зайнятості (ФОП, працює на неволу ставку тощо)	Категоріальна
Вага	Коефіцієнт вани у файлах поточного опитування населення	Числова
Освіта	Рівень освіти (бакалавр, неповне середнє тощо)	Категоріальна
Освітній номер	Номер рівня освіти	Числова
Подружній стан	Одружений, розлучений тощо (поточний статус)	Категоріальна
Професійна діяльність	Рід занять	Категоріальна
Сімейний стан	Має дитину, розлучений тощо	Категоріальна
Раса	Расова приналежність	Категоріальна
Стать	Статева приналежність	Категоріальна
Крїїна походження	Країна проживання (поточна)	Категоріальна
Дохід	Середній дохід в рік	Категоріальна
Годин в неденю працює	Кількість робочих годин в неділю	Числова
Капітал зростання	Рівень росту доходу	Числова

Після описання атрибутів напишімо програму для опису значення показника якості F1 за зразком доповнивши програму пошуком інших показників якості. Обрана точка позначення доходу не дивлячись на наявність чисел в нашому випадку скоріше виступає в якості категоріальної (визначає групи чий дохід вище або нижче за деяке значення)

Лістинг програми:

```
import numpy as np
import matplotlib.pyplot as plt
import sklearn.svm
from sklearn import preprocessing
from sklearn.svm import LinearSVC
from sklearn.svm import SVC

from sklearn.multiclass import OneVsOneClassifier
from sklearn.model_selection import train_test_split
from sklearn.model_selection import cross_val_score
```

					<i>Житомирська політехніка</i> 22.121.06.000 – Лр2		
Змн.	Арк.	№ докум.	Підпис	Дата			
Розроб.	Медведєв В.В.				Звіт з лабораторної роботи	Літ.	Арк.
Перевір.	Філіпов В.О						1
Керівник						ФІКТ Гр. ПІ-61	
Н. контр.							
Зав. каф.							

```

from sklearn.metrics import accuracy_score
from sklearn.metrics import recall_score
from sklearn.metrics import precision_score

#Input Data
input_file = "income_data.txt"

X = []
Y = []
count_class1 = 0
count_class2 = 0
max_datapoints = 500

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

X = np.array(X)

print("-----[X After file reading]-----\n")
print(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:
        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)
scaller = preprocessing.MinMaxScaler(feature_range=(0,1))
X = scaller.fit_transform(X)

classifier = OneVsOneClassifier(LinearSVC(random_state=0))

classifier.fit(X=X, y=Y)

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

scaller = preprocessing.MinMaxScaler(feature_range=(0,1))
X_train = scaller.fit_transform(X_train)

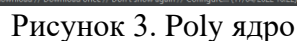
classifier.fit(X=X_train, y=y_train)
y_test_pred = classifier.predict(X_test)

f1 = cross_val_score(classifier, X, Y, scoring="f1_weighted", cv=3)

```

		Медведєв В.В			Житомирська політехніка 22.121.06.000 – Лр2	Арк.
		. Філіпов В.О				
Змн.	Арк.	№ докум.	Підпис	Дата		2

Наступним кроком обчислимо значення показників з використанням різних ядер класифікатора. Коди програм наведено окремими файлами в репозиторії.



		Медведєв. В.В			Житомирська політехніка 22.121.06.000 – Лр2	Арк.
		. Філіпов В.О				4
Змн.	Арк.	№ докум.	Підпис	Дата		

```

18 print("-----[X] ---- [Y]-----\n")
19 print(X)
20 print(Y)
21
22 print("-----NORMALIZED-----\n")
23 scaler = preprocessing.MinMaxScaler(feature_range=(0,1))
24 X = scaler.fit_transform(X)
25 print(X)
26
27 #Classifier = OneVsOneClassifier(LinearSVC(random_state=0))
28 #Classifier = OneVsOneClassifier(SVC(kernel="sigmoid"))
29 #Classifier = OneVsOneClassifier(SVC(kernel="poly", degree=9))
30 Classifier = OneVsOneClassifier(SVC(kernel="rbf"))
31 print("-----CLASSIFIER One Vs One-----")
32 Classifier.fit(X=X, y=Y)
33
34 X_train, X_test, y_train, y_test \
35     = train_test_split(X, Y, test_size=0.2, random_state=5)
36
37 print("-----X_train, X_test, y_train, y_test -----")
38 print(X_train)
39 print(X_test)

```

```

1 1 0 0 1 1 0 1 0 1 1 0 0 0 0 1 0 1 1 1 0 0 1 0 0 0 1 1 0 1 0 1 1 1 1 1 0
0 0 0 1 1 0 1 1 1 1 1 0 0 1 1 1 1 0 0 0 0 1 1 0 0 1 0 0 1 1 0 0 1 0 0 1
0 1 1 0 0 1 1 0 0 1 1 1 1 1 1 1 1 1 0 1 0 1 1 1 0 1 1 1 1 1 1 1 1 1 0 1 0
0 0 1 0 0 0 1 0 1 1 1 1 1 1 1 1 1 1 0 0 0 1 1 1 0 1 1 1 1 1 0 1 0 1 0 1 0
0 0 1 0 0 0 0 0 0 1 1 1 1 1 0
Accuracy: 76.9%
Precision: 77.82%
Recall: 76.9%
F1 score: 76.71%
LABEL ENCODERS
[-1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1]
array([[ 37,  2, 215646, 11,  9,  4,  5,  1,
         4,  1,  0,  0, 40, 27]])
<=50K
Process finished with exit code 0

```

Рисунок 4. Rbf ядро

Як можна бачити rbf дало гарні результати дещо поступившись полі ядру, але набагато вигравши в швидкодії. Сигмоїдне ядро дало більш низькі результати, що дає нам основи вважати що rbf ядро, для нашого випадку, надає найкращі поєднання швидкодії та точності.

Завдання 3: Порівняння якості класифікаторів на прикладі класифікації сортів ірисів

Для виконання даного завдання використаємо набір даних що надається бібліотекою sklearn. Нижче наведено повний лістинг програми.

Лістинг програми:

```

from sklearn.datasets import load_iris
iris_dataset = load_iris()

print("Ключи Iris Dataset : \n{}".format(iris_dataset.keys()))
print(iris_dataset["DESCR"][:193] + "\n...")
print("Names answers: {}".format(iris_dataset["target_names"]))
print("Names for description: \n{}".format(iris_dataset["feature_names"]))
print("Type of data array: {}".format(type(iris_dataset["data"])))
print("Form of array data: {}".format(iris_dataset["data"].shape))
print("Тип масиву target: {}".format(type(iris_dataset['target'])))
print("Відповіді:\n{}".format(iris_dataset['target']))

# Завантаження бібліотек
import numpy as np
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

```

		Медведєв. В.В			Житомирська політехніка 22.121.06.000 – Лр2	Арк.
		. Філіпов В.О				5
Змн.	Арк.	№ докум.	Підпис	Дата		

```

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 = 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()

# 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]

# Разделение X и y на обучающую и контрольную выборки
X_train, X_validation, Y_train, Y_validation = train_test_split(X, y,
test_size=0.20, random_state=1)

#LOAD MODELS
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')))

#Quality models
results = []

```

		Медведєв В.В			Житомирська політехніка 22.121.06.000 – Лр2	Арк.
		. Філіпов В.О				
Змн.	Арк.	№ докум.	Підпис	Дата		6

```

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()

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

# Оцінюємо прогноз
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]])

for name, model in models:
    model.fit(X_train, Y_train)
    prediction = model.predict(X_new)
    print("Прогноз___: {}".format(prediction))
    print(accuracy_score(Y_validation, predictions))
    print(confusion_matrix(Y_validation, predictions))
    print(classification_report(Y_validation, predictions))

```

При виконанні роботи обробивши модель ми отримали наступну діаграму розміху ірисів(рис.5)

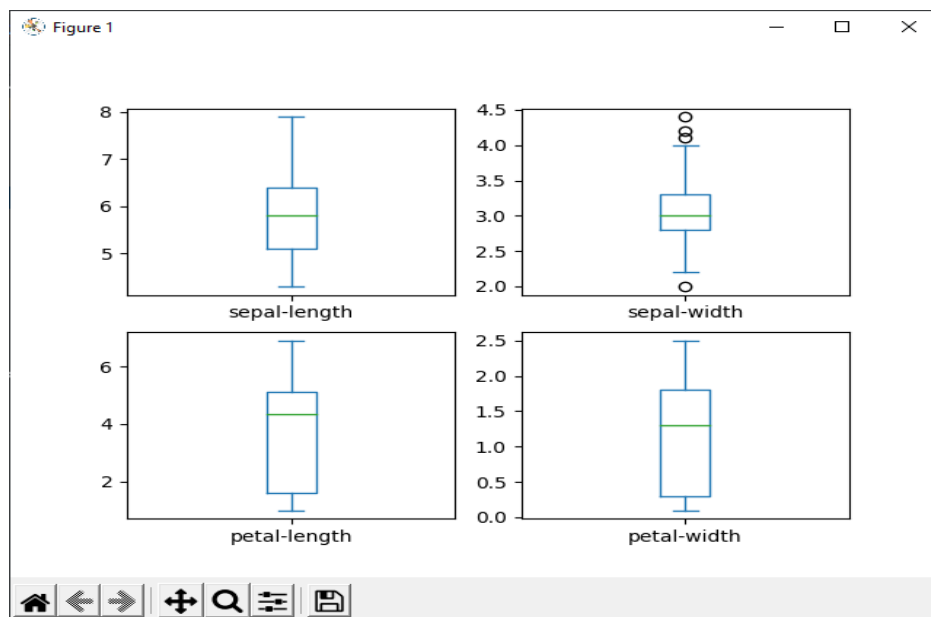


Рисунок 5. Діаграма розмаху

Гістограму розподілу(рис.6)

		Медведєв В.В			Житомирська політехніка 22.121.06.000 – Лр2	Арк.
		Філіпов В.О				7
Змн.	Арк.	№ докум.	Підпис	Дата		

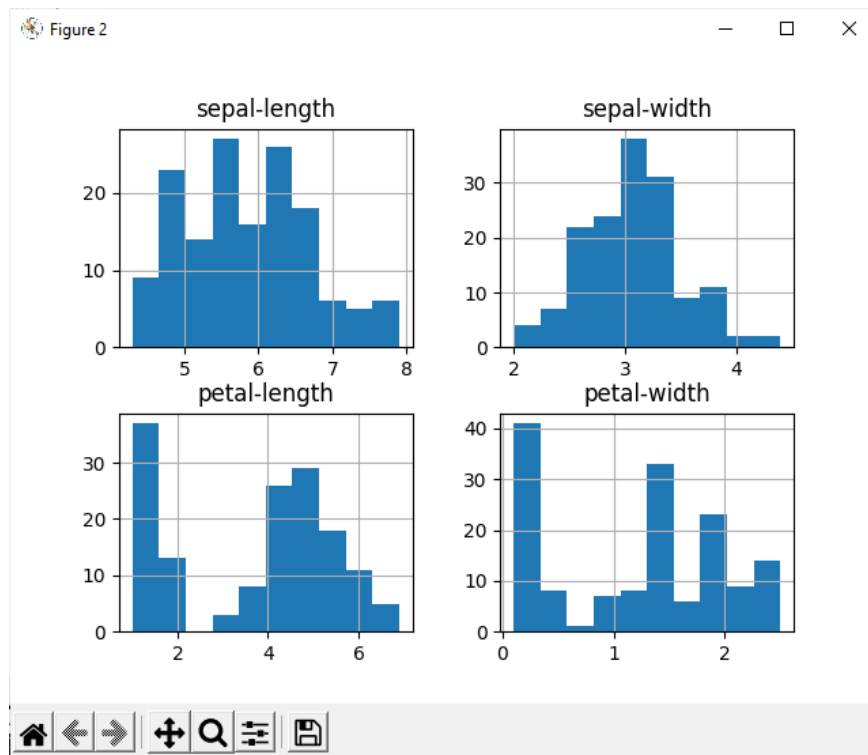


Рисунок 6. Гістограма розподілу атрибутів датасет

Та матрицю діаграми розсіювання (рис.7)

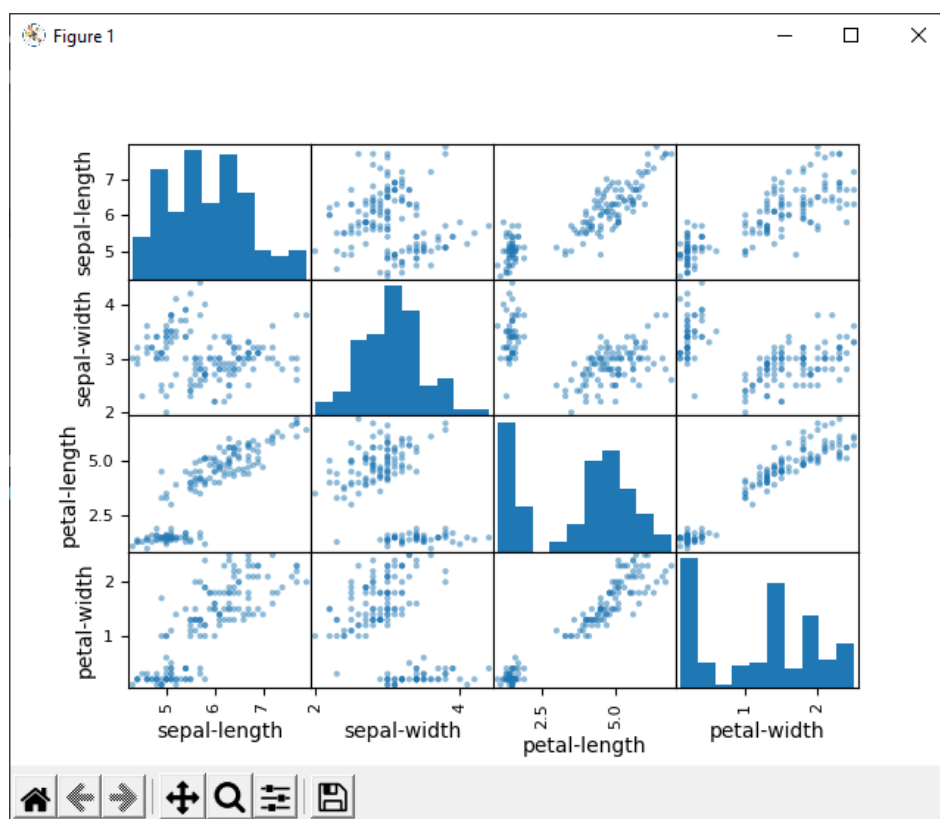


Рисунок 7. Матриця діаграми розсіювання

Сформувавши оцінки алгоритмів було отримано наступну діаграму для поточних конфігурації

		Медведєв. В.В			Житомирська політехніка 22.121.06.000 – Лр2	Арк.
		. Філіпов В.О				8
Змн.	Арк.	№ докум.	Підпис	Дата		

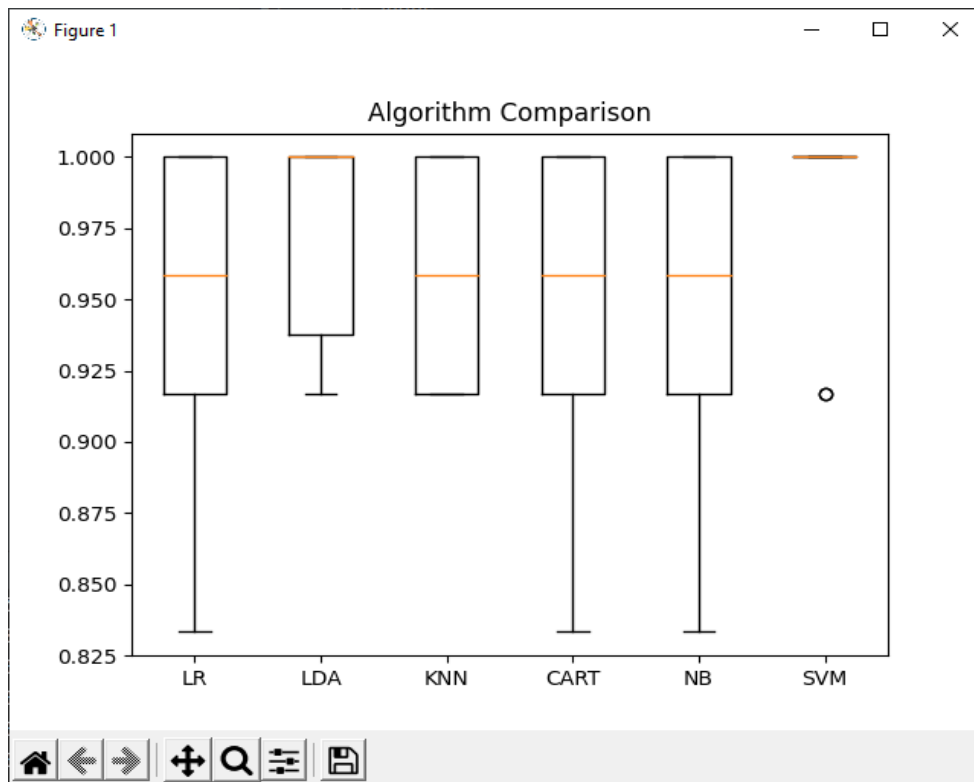


Рисунок 7. Порівняння алгоритмів

Як можна побачити з діаграми порівняння алгоритмів для конфігурації студента Най-краще себе проявила модель лінійного дискримінантного аналізу однак вона не досить стійка й під час тестування займала порівнянно більше часу на своє виконання. Досить непогані результати було отримано з алгоритму к сусідів, але найкраще поєднання передбачиваної точності й швидкої дає алгоритм дерева рішень.

Наступним кроком проведемо оцінки моделі та застосуємо модель для передбачення (код програми засновано на прикладі з методичних рекомендацій з дописом деяких частин)

```

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

# shape
print(dataset.shape)

# 3015 rows head
print(dataset.head(20))

# (Interactive) interactive method describe
print(dataset.describe())

# Formula as a grouped class
print(dataset.groupby('class').size())

# Formula as a grouped class
print(dataset.groupby('class').size())

# Formula as a grouped class
print(dataset.groupby('class').size())

# Formula as a grouped class
print(dataset.groupby('class').size())

# Formula as a grouped class
print(dataset.groupby('class').size())

```

Рисунок 9. Обробка якості моделі

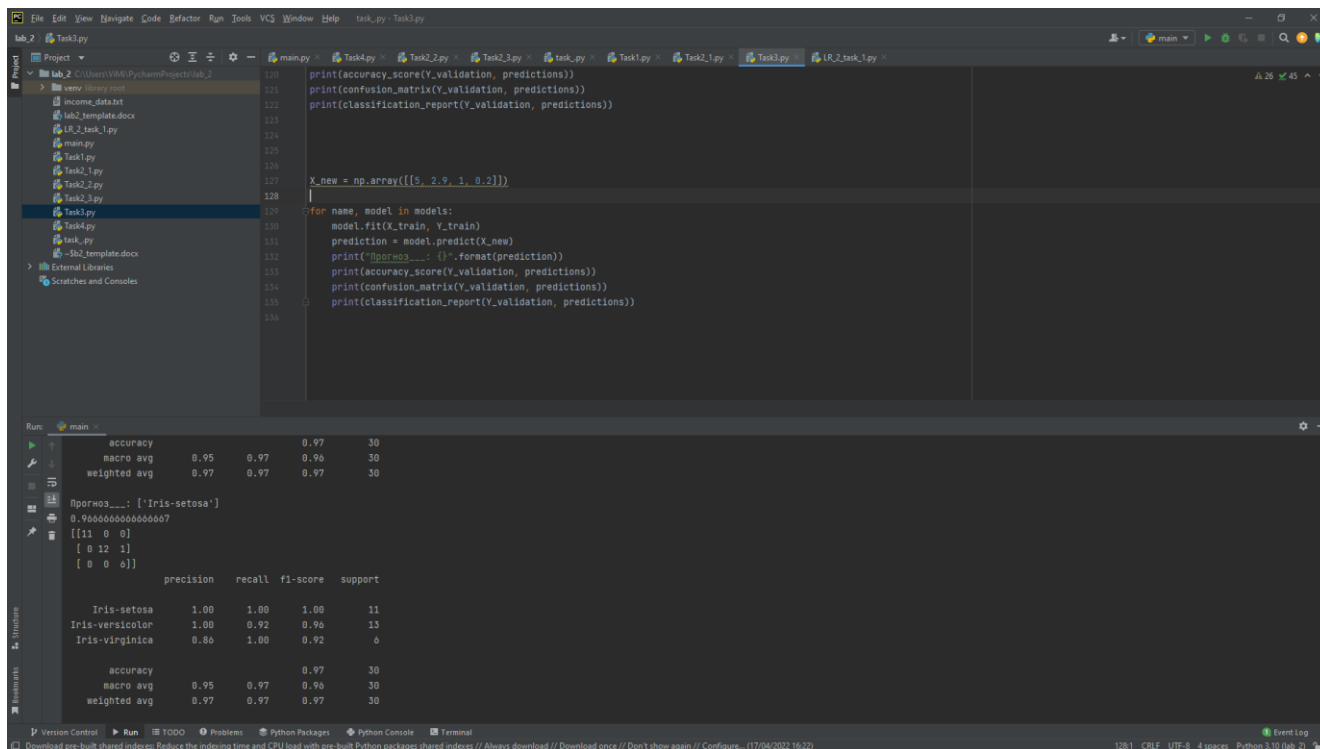


Рисунок 10. Використання моделі для передбачення

Квітка з кроку належить до класу setosa

Завдання 4: Порівняння якості класифікаторів для набору даних

Для виконання цього завдання адаптуємо з виокористанням багатьох моделей до використання на першому наборі даних.

Лістинг програми:

```
import numpy as np
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
import matplotlib.pyplot as plt
import sklearn.svm
from sklearn import preprocessing
from sklearn.svm import LinearSVC
from sklearn.svm import SVC

from sklearn.multiclass import OneVsOneClassifier
from sklearn.model_selection import train_test_split
from sklearn.model_selection import cross_val_score
from sklearn.metrics import accuracy_score
```

		Медведєв. В.В			Житомирська політехніка 22.121.06.000 – Лр2	Арк.
		. Філіпов В.О				10
Змн.	Арк.	№ докум.	Підпис	Дата		

```

from sklearn.metrics import recall_score
from sklearn.metrics import precision_score

#Input Data
input_file = "income_data.txt"

X = []
Y = []
count_class1 = 0
count_class2 = 0
max_datapoints = 500

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

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:
        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)

scaller = preprocessing.MinMaxScaler(feature_range=(0,1))
X = scaller.fit_transform(X)

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')))

for name, model in models:
    classifier = OneVsOneClassifier(model)
    classifier.fit(X=X, y=Y)

X_train, X_test, y_train, y_test \

```

		Медведєв В.В			Житомирська політехніка 22.121.06.000 – Лр2	Арк.
		. Філіпов В.О				
Змн.	Арк.	№ докум.	Підпис	Дата		11

```

    = train_test_split(X, Y, test_size=0.2, random_state=5)

scaller = preprocessing.MinMaxScaler(feature_range=(0,1))
X_train = scaller.fit_transform(X_train)

classifier.fit(X=X_train, y=y_train)
y_test_pred = classifier.predict(X_test)

f1 = cross_val_score(classifier, X, Y, scoring="f1_weighted", cv=3)
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_values = cross_val_score(classifier, X, Y, scoring='f1_weighted', cv=3)
print("F1: " + str(round(100 * f1_values.mean(), 2)) + "%")

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 = np.array([-1] * len(input_data))

print("LABEL ENCODERS")
print(input_data_encoded)
count = 0

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

input_data_encoded = input_data_encoded.astype(int)
input_data_encoded = [ input_data_encoded ]

print("-----")
print(input_data_encoded)

predicate_class = classifier.predict(input_data_encoded)
print(label_encoder[-1].inverse_transform(predicate_class)[0])

```

		Медведєв. В.В			Житомирська політехніка 22.121.06.000 – Лр2	Арк.
		. Філіпов В.О				
Змн.	Арк.	№ докум.	Підпис	Дата		12

LR Accuracy: 73.4% Precision: 73.68% Recall: 73.4% F1: 73.31% F1 score: 73.31% LABEL ENCODERS	LDA Accuracy: 73.4% Precision: 73.61% Recall: 73.4% F1: 73.34% F1 score: 73.34% LABEL ENCODERS	KNN Accuracy: 67.0% Precision: 77.38% Recall: 67.0% F1: 63.48% F1 score: 63.48%
CART Accuracy: 72.2% Precision: 72.05% Recall: 71.8% F1: 72.73% F1 score: 72.25%	SVM Accuracy: 72.1% Precision: 73.28% Recall: 72.1% F1: 71.73% F1 score: 71.73%	NB Accuracy: 50.0% Precision: 25.0% Recall: 50.0% F1: 33.33% F1 score: 33.33%

Рисунок 11. Отримані результати точності для різних класифікаторів

Завдання 5: Класифікація даних лінійним класифікатором Ridge

Наступним кроком проведемо класифікацію за допомогою лінійного класифікатора Ridge

Лістинг програми:

```
import numpy as np
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

#LAST TASK
# =====
# Приклад класифікатора Ridge
# =====

import numpy as np
from sklearn.datasets import load_iris
from sklearn.linear_model import RidgeClassifier

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)

from sklearn import metrics
print('Accuracy:', np.round(metrics.accuracy_score(ytest, ypred), 4))
```

		Медведєв В.В			Житомирська політехніка 22.121.06.000 – Лр2	Арк.
		Філіпов В.О				13
Змн.	Арк.	№ докум.	Підпис	Дата		

```

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

from sklearn.metrics import confusion_matrix
from io import BytesIO

#needed for plot
import seaborn as sns;
sns.set()
import matplotlib.pyplot as plt

mat = confusion_matrix(ytest, ypred)
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")

```

Використаємо два параметри налаштування класифікатора точність рішення (тол) та спосіб рішення що використовується при розрахунках (в нашому випадку це sag (Алгоритм середнього стохастичного градієнту))

```

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 12. Отриманий результат

Як можна побачити з отриманого результату було використано показники якості точність, ф1, recall, коефіцієнт Коена Каппа (статистичне значення, що вимірює міжрегіональну згоду на категоріальні предмети, вважається більш надійним ніж розрахунок у відсотках бо враховує випадковість) та коефіцієнт кореляції Метьюза (використовується в машинному навчанні як міра якості бінарних мультикласних класифікацій, враховує істинні, хибні, позитивні та негативні результати, може бути використаний коли класи мають дуже різні розміри)

Отсанні кроком пояснимо що наведено на малюнку нижче (рис.13)

		Медведєв. В.В			Житомирська політехніка 22.121.06.000 – Лр2	Арк.
		. Філіпов В.О				
Змн.	Арк.	№ докум.	Підпис	Дата		14

Це – Матриця невідповідності, згідно інформації отриманій з мережі інтернет це спеціальна матриця що використовується в галузі машинного навчання, й зокрема в задачі статистичної класифікації, матриця невідповідностей(англ. confusion matrix), також відома як матриця помилок (англ. error matrix) - це таблиця особливого компонування, що дає можливість унаочнювати продуктивність алгоритму, зазвичай керованого навчання (у спонтаннім навчанні її зазвичай називають матрицею допасованості, англ. matching matrix). Кожен з рядків цієї матриці представляє зразки прогнозованого класу, тоді як кожен зі стовпців представляє зразки справжнього класу (або навпаки) Її назва походить від того факту, що вона дає можливість просто бачити, чи допускає система невідповідності між цими двома класами (наприклад, часто помилково маркуючи один як інший).

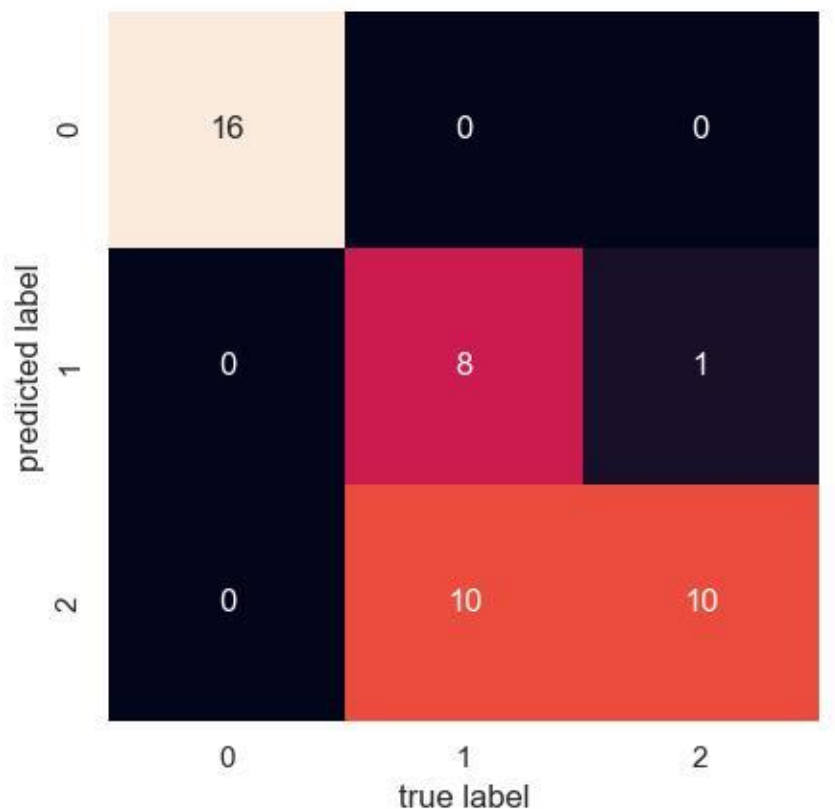


Рисунок 13. Матриця невідповідностей

Висновок: Розглянули особливості роботи з машинами векторів. Провели порівняльний аналіз різних методів класифікації даних.