ЛАБОРАТОРНА РОБОТА №2 ПОРІВНЯННЯ МЕТОДІВ КЛАСИФІКАЦІЇ ДАНИХ

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

Хід роботи:

Завдання 2.1.

Класифікація за допомогою машин опорних векторів (SVM) Код програми:

```
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
from sklearn.model_selection import cross_val_score
input_file = "income_data.txt"
count_class1 = 0
count_class2 = 0
max_datapoints = 25000
with open(input_file, "r") as f:
       if count_class1 >= max_datapoints and count_class2 >= max_datapoints:
            break
        if data[-1] == '<=50K' and count_class1 < max_datapoints:</pre>
            count_class1 += 1
            X.append(data)
            count_class2 += 1
X = np.array(X)
label_encoder = []
X_encoded = np.empty(X.shape)
   if item.isdigit():
       X_encoded[:, i] = X[:, i]
        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))
classifier = OneVsOneClassifier(LinearSVC(random_state=0))
```

					ДУ «Житомирська політехніка».20.121.12.		21.12.	
Змн.	Арк.	№ докум.	Підпис	Дата	· · · · · · · · · · · · · · · · · · ·			
Розр	0 б.	Анкудевич Д.Р.				Літ.	Арк.	Аркушів
Пере	евір.	Філіпов В.О.					1	19
Керіє	зник							
Н. контр.						ФІКТ Гр. ІПЗк-20-1		13к-20-1
Зав.	каф.						•	

```
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))
40 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))
44 X_train = scaller.fit_transform(X_train)
classifier.fit(X=X_train, y=y_train)
    ত C:/Users/ankud/Desktop/myAi/lab2/LR_2_task_1.py
    classifier: OneVsOneClassifier = OneVsOneClassifier(LinearSVC(random_state=0))
49 print("Accuracy: " + str(round(100 * accuracy_values.mean(), 2)) + "%")
50 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)
55 print("F1: " + str(round(100 * f1_values.mean(), 2)) + "%")
56 print("F1 score: " + str(round(100 * f1.mean(), 2)) + "%")
input_data = ['37', 'Private', '215646', 'HS-grad', '9', 'Never-married', 'Handlers-cleaners',
input_data_encoded = np.array([-1] * len(input_data))
61 for i, item in enumerate(input_data):
   if item.isdigit():
            input_data_encoded[i] = item
            input_data_encoded[i] = int(label_encoder[count].transform([item]))
input_data_encoded = input_data_encoded.astype(int)
68 input_data_encoded = [input_data_encoded]
69 predicate_class = classifier.predict(input_data_encoded)
70 print(label_encoder[-1].inverse_transform(predicate_class)[0])
```

Рис 2.1 Код файлу LR_2_task_1.py

```
Python 3.11.3 (tags/v3.11.3:
Accuracy: 81.95%
Precision: 80.94%
Recall: 81.95%
F1: 80.13%
F1 score: 80.13%
>50K
```

Рис 2.2 Результат програми файлу LR_2_task_1.py

 $Ap\kappa$.

		Анкудеувич Д.Р			
		Філіпов В.О.			ДУ «Житомирська політехніка».20.121.12 — Лр2
Змн.	Арк.	№ докум.	Підпис	Дата	

Завдання 2.2.

Порівняння якості класифікаторів SVM з нелінійними ядрами

Код програми:

```
from sklearn import preprocessing
from sklearn.model_selection import train_test_split
from sklearn.model_selection import cross_val_score
input_file = "income_data.txt"
count_class1 = 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
        data = line[:-1].split(', ')
        if data[-1] == '<=50K' and count_class1 < max_datapoints:</pre>
            X.append(data)
            count_class1 += 1
        if data[-1] == '>50K' and count_class2 < max_datapoints:</pre>
            X.append(data)
            count_class2 += 1
X = np.array(X)
label_encoder = []
X_encoded = np.empty(X.shape)
       X_{encoded}[:, i] = X[:, i]
        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)
```

		Анкудеувич Д.Р		
		Філіпов В.О.		
Змн.	Арк.	№ докум.	Підпис	Дата

```
classifier.fit(X=X, y=Y)
    X_train, X_test, y_train, y_test \
   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)
   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)) + "%")
   input_data = ['37', 'Private', '215646', 'HS-grad', '9', 'Never-married', 'Handlers-cleaners',
input_data_encoded = np.array([-1] * len(input_data))
   for i, item in enumerate(input_data):
       if item.isdigit():
           input_data_encoded[i] = item
            input_data_encoded[i] = int(label_encoder[count].transform([item]))
    input_data_encoded = input_data_encoded.astype(int)
    input_data_encoded = [input_data_encoded]
    predicate_class = classifier.predict(input_data_encoded)
    print(label_encoder[-1].inverse_transform(predicate_class)[0])
```

Рис 2.3 Код файлу LR_2_task_2_1.py

Accuracy: 83.96%
Precision: 83.18%
Recall: 83.96%
F1: 82.95%
F1 score: 82.95%
<=50K

Рис 2.4 Результат Поліномінального ядра

		Анкудеувич Д.Р		
		Філіпов В.О.		
Змн.	Арк.	№ докум.	Підпис	Дата

Accuracy: 83.96% Precision: 83.18%

Рис 2.7 Результат гаусового ядра

Accuracy: 83.96% Precision: 83.18%

Рис 2.8 Результат сигмоїдального ядра

Висновок: в даній ситуації краще за всього справляється RBF, має кращу точність та швидкість. Сигмоїдне ядро не так добре, так як відстає по швидкості.

Завдання 2.3

Порівняння якості класифікаторів на прикладі класифікації сортів ірисів Код програми:

		Анкудеувич Д.Р		
		Філіпов В.О.		
Змн.	Арк.	№ докум.	Підпис	Дата

```
from sklearn.datasets import load_iris
                                                                            A3 ±50
i⊯port 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
iris_dataset = load_iris()
print("Ключі iris dataset : \n{}".format(iris_dataset.keys()))
print(iris_dataset["DESCR"][:193] + "\n...")
print("Назви відповідей: {}".format(iris_dataset["target_names"]))
print("Назви ознак: \n{}".format(iris_dataset["feature_names"]))
print("Тип масиву date: {}".format(type(iris_dataset["data"])))
print("Форма масиву data: {}".format(iris_dataset["data"].shape))
print("Тип масиву target: {}".format(type(iris_dataset['target'])))
print("Відповіді:\n{}".format(iris_dataset['target']))
```

		Анкудеувич Д.Р		
		Філіпов В.О.		
Змн.	Арк.	№ докум.	Підпис	Дата

```
print("Ключі iris dataset : \n{}".format(iris_dataset.keys()))
    print(iris_dataset["DESCR"][:193] + "\n...")
   print("Назви відповідей: {}".format(iris_dataset["target_names"]))
   print("Назви ознак: \n{}".format(iris_dataset["feature_names"]))
   print("Тип масиву date: {}".format(type(iris_dataset["data"])))
   print("Форма масиву data: {}".format(iris_dataset["data"].shape))
   print("Tun macusy target: {}".format(type(iris_dataset['target'])))
   print("Відповіді:\n{}".format(iris_dataset['target']))
   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)
32 print(dataset.shape)
34 print(dataset.head(20))
36 print(dataset.describe())
38 print(dataset.groupby('class').size())
   dataset.plot(kind='box', subplots=True, layout=(2, 2), sharex=False, sharey=False)
41 pyplot.show()
42 # Гістограма розподілу атрибутів датасета
```

		Анкудеувич Д.Р		
		Філіпов В.О.		
Змн.	Арк.	№ докум.	Підпис	Дата

```
dataset.plot(kind='box', subplots=True, layout=(2, 2), sharex=False, sharey=False)
42 # Гістограма розподілу атрибутів датасета
43 dataset.hist()
44 pyplot.show()
45 # Матриця діаграм розсіювання
46 scatter_matrix(dataset)
47 pyplot.show()
49 array = dataset.values
51 X = array[:, 0:4]
53 y = array[:, 4]
55 X_train, X_validation, Y_train, Y_validation = train_test_split(X, y, test_size=0.20,
   random_state=1)
   models = []
58 | models.append(('LR', LogisticRegression(solver='liblinear', multi_class='ovr')))
59 models.append(('LDA', LinearDiscriminantAnalysis()))
60 models.append(('KNN', KNeighborsClassifier()))
61 models.append(('CART', DecisionTreeClassifier()))
62 models.append(('NB', GaussianNB()))
63 models.append(('SVM', SVC(gamma='auto')))
64 results = []
```

		Анкудеувич Д.Р		
		Філіпов В.О.		
Змн.	Арк.	№ докум.	Підпис	Дата

```
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)
74 pyplot.title('Algorithm Comparison')
    pyplot.show()
   model = SVC(gamma='auto')
78 model.fit(X_train, Y_train)
    predictions = model.predict(X_validation)
   print(accuracy_score(Y_validation, predictions))
82 print(confusion_matrix(Y_validation, predictions))
   print(classification_report(Y_validation, predictions))
   X_{new} = np.array([[5, 2.9, 1, 0.2]])
85 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))
```

Рис 2.9 Код файлу LR_2_task_3.py

		Анкудеувич Д.Р		
		Філіпов В.О.		
Змн.	Арк.	№ докум.	Підпис	Дата

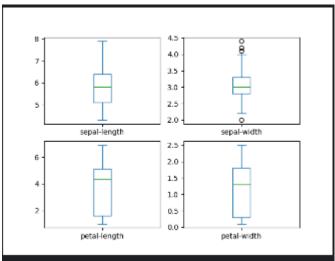


Рис. 2.10 Результат діаграми розмаху

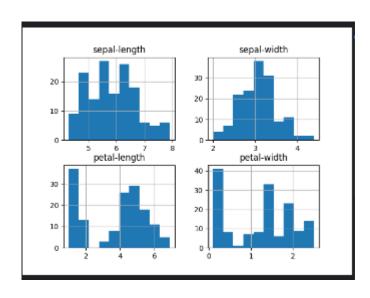
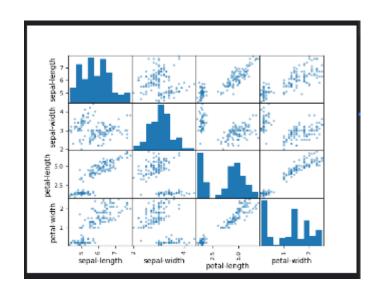


Рис. 2.11 Гістрограма розподілу атрибутів



		Анкудеувич Д.Р		
		Філіпов В.О.		
Змн.	Арк.	№ докум.	Підпис	Дата

Рис. 2.12 Матриця діаграми розсіювання

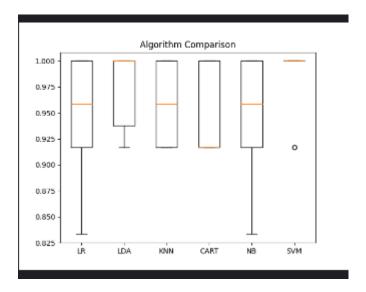


Рис. 2.13 Рисунок порівняння алгоритмів

```
Ключі iris dataset :
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, pre
Назви відповідей: ['setosa' 'versicolor' 'virginica']
Назви ознак:
Тип масиву date: <class 'numpy.ndarray'>
Форма масиву data: (150, 4)
Тип масиву target: <class 'numpy.ndarray'>
2 2]
(150, 5)
    sepal-length sepal-width petal-length petal-width
                                            0.2 Iris-setosa
```

		Анкудеувич Д.Р		
		Філіпов В.О.		
Змн.	Арк.	№ докум.	Підпис	Дата

se	pal-length se	epal-width pe	etal-length	petal-width	class
Θ	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
	sepal-length	sepal-width			
count	150.000000	150.000000	150.00000	150.0000	000
mean	5.843333	3.054000	3.75866	7 1.1986	67
std	0.828066	0.433594	1.76442	0.7631	.61
min	4.300000	2.000000			000
25%	5.100000	2.800000	1.60000		
50%	5.800000	3.000000	4.35000		
75%	6.400000	3.300000	5.10000		
max	7.900000	4.400000	6.90000	0 2.5000	000
class					
Iris-s	etosa :	50			

Рис. 2.14 Результат програми

Висновок: найкраще показала себе модель лінійного дискримінантного аналізу. Квітка належала до класу Iris-setosa.

		Анкудеувич Д.Р		
		Філіпов В.О.		
Змн.	Арк.	№ докум.	Підпис	Дата

Завдання 2.4.

Порівняння якості класифікаторів для набору даних завдання 2.1 Код програми:

	·	Анкудеувич Д.Р		
		Філіпов В.О.		
Змн.	Арк.	№ докум.	Підпис	Дата

		Анкудеувич Д.Р		
		Філіпов В.О.		
Змн.	Арк.	№ докум.	Підпис	Дата

Рис 2.15 Код файлу LR_2_task_4.py

Accuracy: 81.82% Precision: 80.69% Recall: 81.82% F1: 80.25%

F1 score: 80.25%

>50K

Рис.2.16 Точність класифікатора LR

Accuracy: 81.14%
Precision: 79.86%
Recall: 81.14%
F1: 79.35%
F1 score: 79.35%
>50K

Рис. 2.17 Точність класифікатора LDA

		Анкудеувич Д.Р		
		Філіпов В.О.		
Змн.	Арк.	№ докум.	Підпис	Дата

Accuracy: 82.16% Precision: 81.53% Recall: 82.16%

F1: 81.75%

F1 score: 81.75%

<=50K

Рис. 2.18 Точність класифікатора KNN

Accuracy: 80.55% Precision: 80.76% Recall: 80.66%

F1: 80.84%

F1 score: 80.77%

>50K

Рис. 2.19 Точність класифікатора CART

Accuracy: 79.76% Precision: 78.2% Recall: 79.76%

F1: 77.13%

F1 score: 77.13%

<=50K

Рис. 2.20 Точність класифікатора NB

Accuracy: 82.38% Precision: 81.51%

Recall: 82.38%

F1: 80.6%

F1 score: 80.6%

>50K

Рис. 2.21 Точність класифікатора SVM

Завдання 2.5.

Класифікація даних лінійним класифікатором Ridge

Код програми:

		Анкудеувич Д.Р			
		Філіпов В.О.			ДУ «Житомирська політехніка».20.121.12 — Лр2
Змн.	Арк.	№ докум.	Підпис	Дата	

```
import numpy as np
from sklearn.datasets import load_iris
from sklearn.linear_model import RidgeClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import confusion_matrix
from io import BytesIO
import matplotlib.pyplot as plt
from sklearn import metrics
iris = load_iris()
X, y = iris.data, iris.target
clf = RidgeClassifier(tol=1e-2, solver="sag")
clf.fit(Xtrain, ytrain)
ypred = clf.predict(Xtest)
print('Accuracy:', np.round(metrics.accuracy_score(ytest, ypred), 4))
    ytest, ypred, average='weighted'), 4))
print('F1 Score:', np.round(metrics.f1_score(ytest, ypred, average='weighted'), 4))
    metrics.cohen_kappa_score(ytest, ypred), 4))
    metrics.matthews_corrcoef(ytest, ypred), 4))
```

```
metrics.classification_report(ypred, ytest))
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")

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

Рис 2.22 Код файлу LR_2_task_5.py

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

		Анкудеувич Д.Р		
		Філіпов В.О.		
Змн.	Арк.	№ докум.	Підпис	Дата

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			

Рис. 2.23 Результат програми

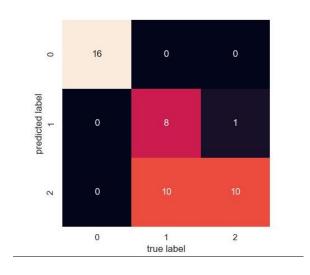


Рис. 2.24 Матриця невідповідності

Висновок: Матриця невідповідності — це таблиця особливого компонування, що дає можливість унаочнювати продуктивність алгоритму. Кожен з рядків цієї матриці представляє зразки прогнозованого класу, тоді як кожен зі стовпців представляє зразки справжнього класу (або навпаки).

Коефіцієнт каппа Коена статистика, яка використовується для вимірювання надійності між експертами для якісних пунктів.

Кореляції Метьюза — використовується в машинному навчанні, як міра якості бінарних мультикласних класифікацій.

		Анкудеувич Д.Р		
		Філіпов В.О.		
Змн.	Арк.	№ докум.	Підпис	Дата

Висновок: в ході виконання лабораторної роботи використовуючи спеціалізовані бібліотеки та мову програмування Python дослідив різні методи класифікації даних та навчився їх порівнювати.

GitHub: https://github.com/Max1648/Artificial-Intelligence2 - новий гіт, з першим виникли проблеми.

Старий гіт - https://github.com/Max1648/Artificial-Intelligence

		Анкудеувич Д.Р		
		Філіпов В.О.		
Змн.	Арк.	№ докум.	Підпис	Дата