

# ЛАБОРАТОРНА РОБОТА № 1

## ПОПЕРЕДНЯ ОБРОБКА ТА КОНТРОЛЬОВАНА КЛАСИФІКАЦІЯ ДАНИХ

**Мета роботи:** використовуючи спеціалізовані бібліотеки та мову програмування Python дослідити попередню обробку та класифікацію даних.

**Хід роботи:**

**Завдання 2.1-2.1.4. Попередня обробка даних.**

```

1  import numpy as np
2  from sklearn import preprocessing
3
4  input_data = np.array([[5.1, -2.9, 3.3],
5                        [-1.2, 7.8, -6.1],
6                        [3.9, 0.4, 2.1],
7                        [7.3, -9.9, -4.5]])
8
9  # Бінаризація даних
10 data_binarized = preprocessing.Binarizer(threshold=2.1).transform(input_data)
11 print("\n Binarized data:\n", data_binarized)
12 # Виведення середнього значення та стандартного відхилення
13 print("\nBEFORE: ")
14 print("Mean =", input_data.mean(axis=0))
15 print("Std deviation =", input_data.std(axis=0))
16 # Виключення середнього
17 data_scaled = preprocessing.scale(input_data)
18 print("\nAFTER: ")
19 print("Mean =", data_scaled.mean(axis=0))
20 print("Std deviation =", data_scaled.std(axis=0))
21 # Масштабування MinMax
22 data_scaler_minmax = preprocessing.MinMaxScaler(feature_range=(0, 1))
23 data_scaled_minmax = data_scaler_minmax.fit_transform(input_data)
24 print("\nMin max scaled data:\n", data_scaled_minmax)
25 # Нормалізація даних
26 data_normalized_l1 = preprocessing.normalize(input_data, norm='l1')
27 data_normalized_l2 = preprocessing.normalize(input_data, norm='l2')
28 print("\nL1 normalized data:\n", data_normalized_l1)
29 print("\nL2 normalized data:\n", data_normalized_l2)

```

Рис. 1. Відображення коду програм.

					ЖИТОМИРСЬКА ПОЛІТЕХНІКА.23.125.08.000 - Лр1		
Зм	Арк.	№ докум.	Підпис	Дата			
Розроб.		Куліш М.В.			Звіт з лабораторної роботи №1		
Перевір.		Пулеко І.В.					
Реценз.							
Н. Контр.							
Зав.каф.							
					Літ.	Арк.	Акрушів
						1	15
					ФІКТ, гр. КБм-22-1		

```

LR_1_task_1 x
/home/xtr99/labs/ai/lab01/venv/bin/python /home/xtr99/labs/ai/lab01/LR_1_task_1.py

Binarized data:
[[1. 0. 1.]
 [0. 1. 0.]
 [1. 0. 0.]
 [1. 0. 0.]]

BEFORE:
Mean = [ 3.775 -1.15 -1.3 ]
Std deviation = [3.12039661 6.36651396 4.0620192 ]

AFTER:
Mean = [1.11022302e-16 0.00000000e+00 2.77555756e-17]
Std deviation = [1. 1. 1.]

Min max scaled data:
[[0.74117647 0.39548023 1.
 [0. 1. 0.
 [0.6 0.5819209 0.87234043]
 [1. 0. 0.17021277]]

L1 normalized data:
[[ 0.45132743 -0.25663717 0.2920354 ]
 [-0.0794702 0.51655629 -0.40397351]
 [ 0.609375 0.0625 0.328125 ]
 [ 0.33640553 -0.4562212 -0.20737327]]

L2 normalized data:
[[ 0.75765788 -0.43082507 0.49024922]
 [-0.12030718 0.78199664 -0.61156148]
 [ 0.87690281 0.08993875 0.47217844]
 [ 0.55734935 -0.75585734 -0.34357152]]

```

Рис. 2. Виконання коду програм.

L1-нормалізація використовує метод найменших абсолютних відхилень для забезпечення того, що сума абсолютних значень вектора буде дорівнювати 1. За своєю суттю, L1-нормалізація вважається більш надійною порівняно з L2-нормалізацією, оскільки вона менш чутлива до викидів або значень, що відхиляються від середнього значення.

L2-нормалізація, натомість, використовує метод найменших квадратів для забезпечення того, що сума квадратів значень вектора буде дорівнювати 1. Хоча L2-нормалізація також є ефективним методом нормалізації, вона може бути більш чутливою до викидів, що можуть впливати на результуючий вектор.

Отже, L1-нормалізація вважається більш надійною технікою порівняно з L2-нормалізацією через її меншу чутливість до викидів або відхилень від середнього значення.

		Куліш М.В.			ЖИТОМИРСЬКА ПОЛІТЕХНІКА.23.125.08.000 - Лр1	Арк.
		Пулеко І.В.				2
Змн.	Арк.	№ докум.	Підпис	Дата		

Для прикладу L1-нормалізації та L2-нормалізації скористаємося даними з рис. 2. рядками 4 масиву.

Результат для L1-нормалізації:  $|0.33640553| + |-0.4562212| + |-0.20737327| = 1$

Результат для L2-нормалізації:  $0.55734935^2 + (-0.75585734)^2 + (-0.34357152)^2 = 1$

### Завдання 2.1.5. Кодування міток.

```
# Надання позначок вхідних даних
input_labels = ['red', 'black', 'red', 'green', 'black', 'yellow', 'white']
# Створення кодувальника та встановлення відповідності між мітками та числами
encoder = preprocessing.LabelEncoder()
encoder.fit(input_labels)
# Виведення відображення
print("\nLabel mapping:")
for i, item in enumerate(encoder.classes_):
    print(item, '-->', i)
# перетворення міток за допомогою кодувальника
test_labels = ['green', 'red', 'black']
encoded_values = encoder.transform(test_labels)
print("\nLabels =", test_labels)
print("Encoded values =", list(encoded_values))
# Декодування набору чисел за допомогою декодера
encoded_values = [3, 0, 4, 1]
decoded_list = encoder.inverse_transform(encoded_values)
print("\nEncoded values =", encoded_values)
print("Decoded labels =", list(decoded_list))
```

Рис. 3. Відображення коду програми.

```
Label mapping:
black --> 0
green --> 1
red --> 2
white --> 3
yellow --> 4

Labels = ['green', 'red', 'black']
Encoded values = [1, 2, 0]

Encoded values = [3, 0, 4, 1]
Decoded labels = ['white', 'black', 'yellow', 'green']
```

Рис. 4. Результат виконання програми.

		Куліш М.В.			ЖИТОМИРСЬКА ПОЛІТЕХНІКА.23.125.08.000 - Лр1	Арк.
		Пулеко І.В.				3
Змн.	Арк.	№ докум.	Підпис	Дата		

Як можна побачити на рис.4 кожній мітці присвоєне число за допомогою якого можна закодувати та декодувати значення мітки.

## Завдання 2.2. Попередня обробка нових даних.

Таблиця 1

№ варіанту	Значення змінної input_data												Поріг бінаризації
8.	4.6	9.9	-3.5	-2.9	4.1	3.3	-2.2	8.8	-6.1	3.9	1.4	2.2	2.2

```

LR_1_task_2.py x
1 import numpy as np
2 from sklearn import preprocessing
3
4 input_data = np.array([[4.6, 9.9, -3.5],
5                        [-2.9, 4.1, 3.3],
6                        [-2.2, 8.8, -6.1],
7                        [3.9, 1.4, 2.2]])
8
9 # Бінаризація даних
10 data_binarized = preprocessing.Binarizer(threshold=2.1).transform(input_data)
11 print("\n Binarized data:\n", data_binarized)
12
13 # Виключення середнього
14 data_scaled = preprocessing.scale(input_data)
15 print("Mean =", data_scaled.mean(axis=0))
16 print("Std deviation =", data_scaled.std(axis=0))
17
18 # Масштабування MinMax
19 data_scaler_minmax = preprocessing.MinMaxScaler(feature_range=(0, 1))
20 data_scaled_minmax = data_scaler_minmax.fit_transform(input_data)
21 print("\nMin max scaled data:\n", data_scaled_minmax)
22
23 # Нормалізація даних
24 data_normalized_l1 = preprocessing.normalize(input_data, norm='l1')
25 data_normalized_l2 = preprocessing.normalize(input_data, norm='l2')
26 print("\nL1 normalized data:\n", data_normalized_l1)
27 print("\nL2 normalized data:\n", data_normalized_l2)

```

Рис. 5. Відображення коду програм.

		Куліш М.В.			ЖИТОМИРСЬКА ПОЛІТЕХНІКА.23.125.08.000 - Лр1	Арк.
		Пулеко І.В.				4
Змн.	Арк.	№ докум.	Підпис	Дата		

```

Binarized data:
[[1. 1. 0.]
 [0. 1. 1.]
 [0. 1. 0.]
 [1. 0. 1.]]

```

Рис. 6. Результат бінарізації даних.

```

Mean = [0.00000000e+00 1.11022302e-16 2.77555756e-17]
Std deviation = [1. 1. 1.]

```

Рис. 7. Результат виключення середнього.

```

Min max scaled data:
[[1.          1.          0.27659574]
 [0.          0.31764706 1.          ]
 [0.09333333 0.87058824 0.          ]
 [0.90666667 0.          0.88297872]]

```

Рис. 8. Результат масштабування MinMax.

```

L1 normalized data:
[[ 0.25555556  0.55      -0.19444444]
 [-0.2815534   0.39805825  0.32038835]
 [-0.12865497  0.51461988 -0.35672515]
 [ 0.52         0.18666667  0.29333333]]

L2 normalized data:
[[ 0.40126114  0.86358375 -0.30530739]
 [-0.4825966   0.68229174  0.54916164]
 [-0.20125974  0.80503895 -0.55803836]
 [ 0.83129388  0.29841319  0.46893501]]

```

Рис. 9. Результат нормалізації L1 та L2.

		Куліш М.В.			ЖИТОМИРСЬКА ПОЛІТЕХНІКА.23.125.08.000 - Лр1	Арк.
		Пулеко І.В.				5
Змн.	Арк.	№ докум.	Підпис	Дата		



### Завдання 2.3. Класифікація логістичною регресією або логістичний класифікатор

```

1  import numpy as np
2  from sklearn import linear_model
3  from utilities import visualize_classifier
4
5  # Визначення зразка вхідних даних
6  X = np.array([[3.1, 7.2], [4, 6.7], [2.9, 8], [5.1, 4.5],
7               [6, 5], [5.6, 5], [3.3, 0.4],
8               [3.9, 0.9], [2.8, 1],
9               [0.5, 3.4], [1, 4], [0.6, 4.9]])
10 y = np.array([0, 0, 0, 1, 1, 1, 2, 2, 2, 3, 3, 3])
11
12 # Створення логістичного класифікатора
13 classifier = linear_model.LogisticRegression(solver='liblinear', C=1)
14
15 # Тренування класифікатора
16 classifier.fit(X, y)
17 visualize_classifier(classifier, X, y)

```

Рис. 10. Відображення коду програми.

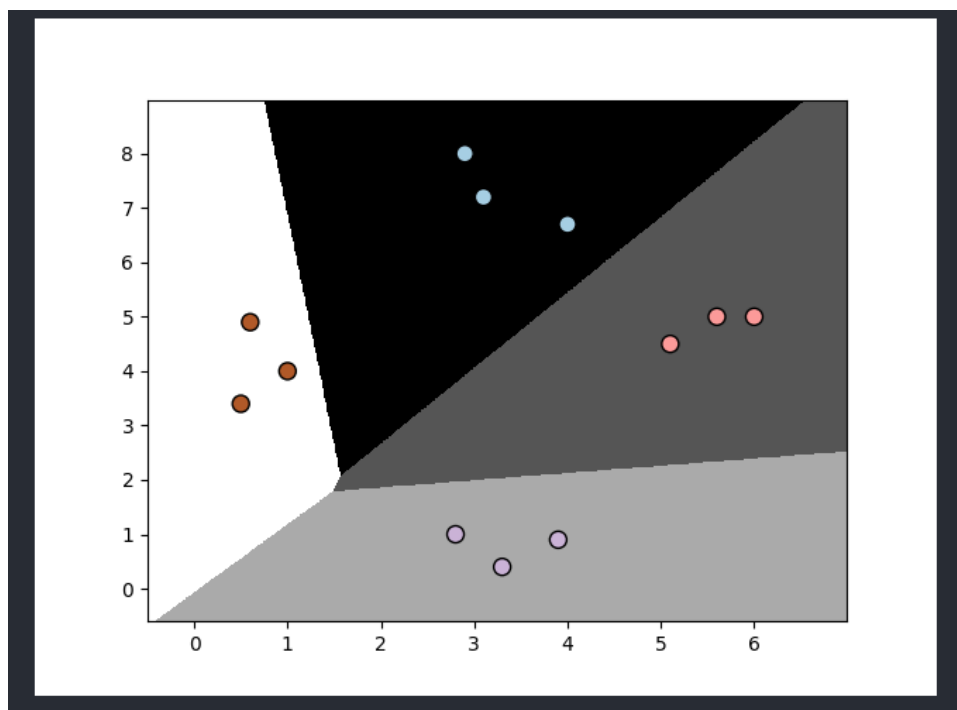


Рис. 11. Результат виконання програми.

		Куліш М.В.			ЖИТОМИРСЬКА ПОЛІТЕХНІКА.23.125.08.000 - Лр1	Арк.
		Пулеко І.В.				6
Змн.	Арк.	№ докум.	Підпис	Дата		

## Завдання 2.4. Класифікація наївним байєсовським класифікатором

```

1  import numpy as np
2      from sklearn.naive_bayes import GaussianNB
3      from sklearn.model_selection import train_test_split, cross_val_score
4      from utilities import visualize_classifier
5
6      # Вхідний файл, який містить дані
7      input_file = 'data_multivar_nb.txt'
8
9      # Завантаження даних із вхідного файлу
10     data = np.loadtxt(input_file, delimiter=',')
11     X, y = data[:, :-1], data[:, -1]
12
13     # Створення наївного байєсовського класифікатора
14     classifier = GaussianNB()
15
16     # Тренування класифікатора
17     classifier.fit(X, y)
18
19     # Прогнозування значень для тренувальних даних
20     y_pred = classifier.predict(X)
21
22     # Обчислення якості класифікатора
23     accuracy = 100.0 * (y == y_pred).sum() / X.shape[0]
24     print("Accuracy of Naive Bayes classifier =", round(accuracy, 2), "%")
25
26     # Візуалізація результатів роботи класифікатора
27     visualize_classifier(classifier, X, y)

```

Рис. 12. Відображення коду програми.

```

LR_1_task_4 x
/home/xtr99/labs/ai/lab01/venv/bin/python /home/xtr99/labs/ai/lab01/LR_1_task_4.py
Accuracy of Naive Bayes classifier = 99.75 %

```

Рис. 13. Результат якості класифікатора.

		Куліш М.В.			ЖИТОМИРСЬКА ПОЛІТЕХНІКА.23.125.08.000 - Лр1	Арк.
		Пулеко І.В.				7
Змн.	Арк.	№ докум.	Підпис	Дата		

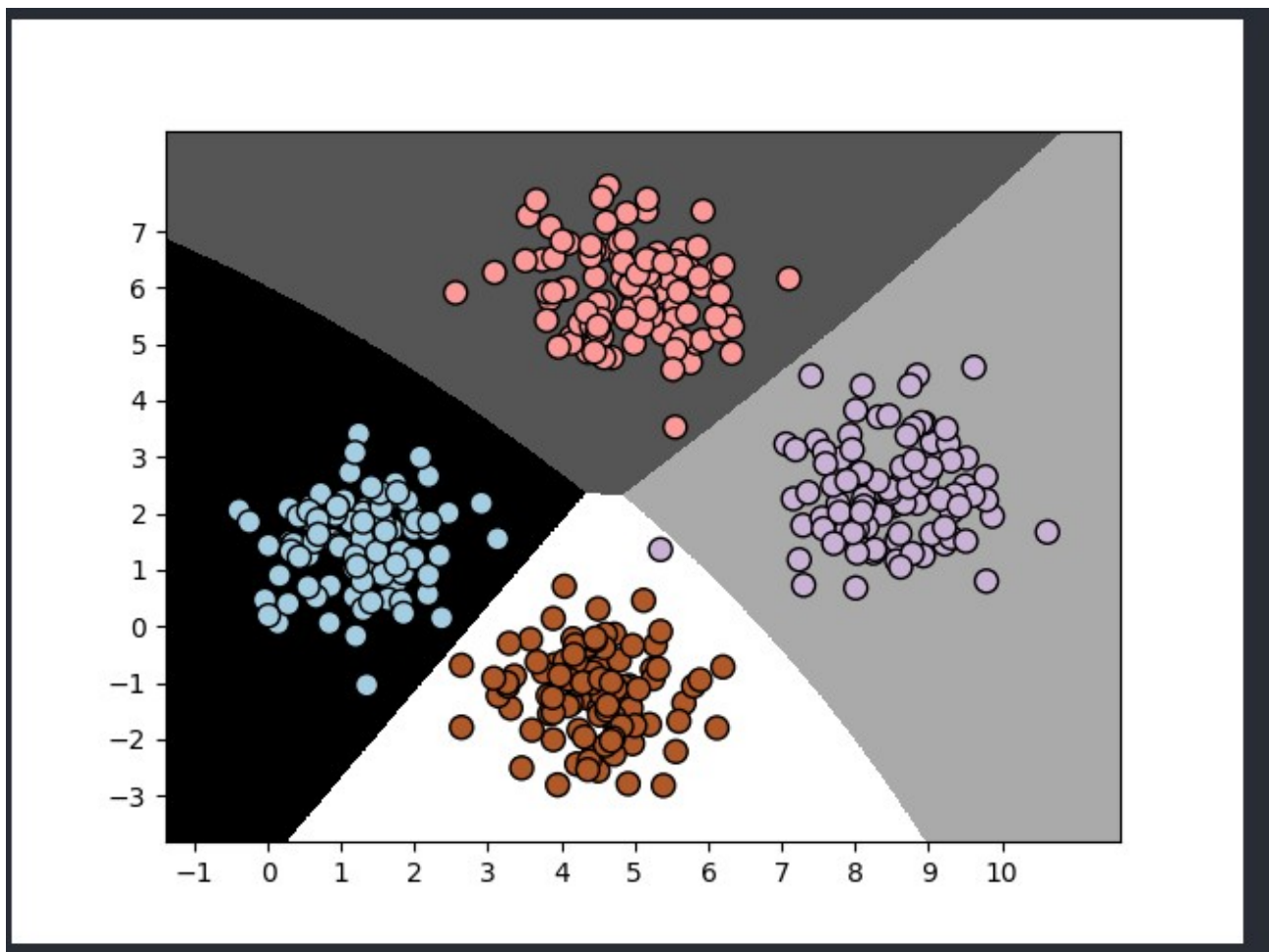


Рис. 14. Результат виконання програми.

```
# Обчислення якості класифікатора
accuracy = 100.0 * (y_test == y_test_pred).sum() / X_test.shape[0]
print("Accuracy of the new classifier =", round(accuracy, 2), "%")
# Візуалізація роботи класифікатора
visualize_classifier(classifier_new, X_test, y_test)
```

Рис. 15. Додавання до програми обчислення якості та візуалізації роботи класифікатора.

```
num_folds = 3
accuracy_values = cross_val_score(classifier, X, y, scoring='accuracy', cv=num_folds)
print("Accuracy: " + str(round(100 * accuracy_values.mean(), 2)) + "%")
precision_values = cross_val_score(classifier, X, y, scoring='precision_weighted', cv=num_folds)
print("Precision: " + str(round(100 * precision_values.mean(), 2)) + "%")
recall_values = cross_val_score(classifier, X, y, scoring='recall_weighted', cv=num_folds)
print("Recall: " + str(round(100 * recall_values.mean(), 2)) + "%")
f1_values = cross_val_score(classifier, X, y, scoring='f1_weighted', cv=num_folds)
print("F1: " + str(round(100 * f1_values.mean(), 2)) + "%")
```

Рис. 16. Додавання до програми вбудовані функції для обчислення якості (accuracy), точності (precision) 2 та повноти (recall) 3 класифікатора на підставі потрійний перехрестної перевірки.

		Куліш М.В.			ЖИТОМИРСЬКА ПОЛІТЕХНІКА.23.125.08.000 - Лр1	Арк.
		Пулеко І.В.				8
Змн.	Арк.	№ докум.	Підпис	Дата		



```

Accuracy of Naive Bayes classifier = 99.75 %
Accuracy of the new classifier = 100.0 %
Accuracy: 99.75%
Precision: 99.76%
Recall: 99.75%
F1: 99.75%

```

Рис. 17. Результат обчислення якості (ассурасу), точності (precision) 2 та повноти (recall) 3 класифікатора на підставі потрійний перехрестної перевірки.

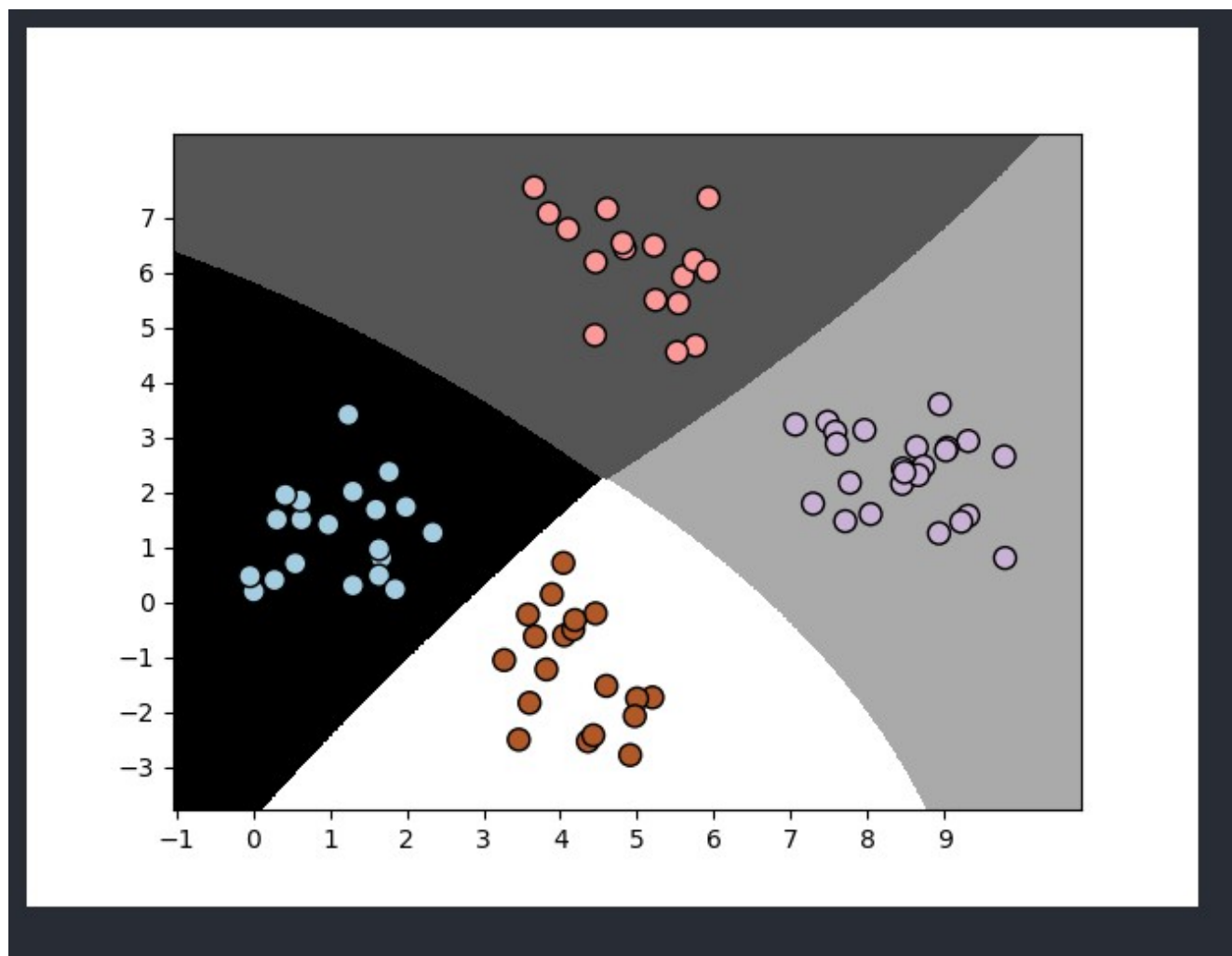


Рис. 18. Результат виконання програми.

		Куліш М.В.			ЖИТОМИРСЬКА ПОЛІТЕХНІКА.23.125.08.000 - Лр1	Арк.
		Пулеко І.В.				9
Змн.	Арк.	№ докум.	Підпис	Дата		

```

LR_1_task_4 x
/home/xtr99/labs/ai/lab01/venv/bin/python /home,
Accuracy of Naive Bayes classifier = 99.75 %
Accuracy of the new classifier = 100.0 %
Accuracy: 99.75%
Precision: 99.76%
Recall: 99.75%
F1: 99.75%

Process finished with exit code 0

```

Рис. 19. Повторний результат обчислення якості (ассурасу), точності (precision) 2 та повноти (recall) 3 класифікатора на підставі потрійний перехрестної перевірки.

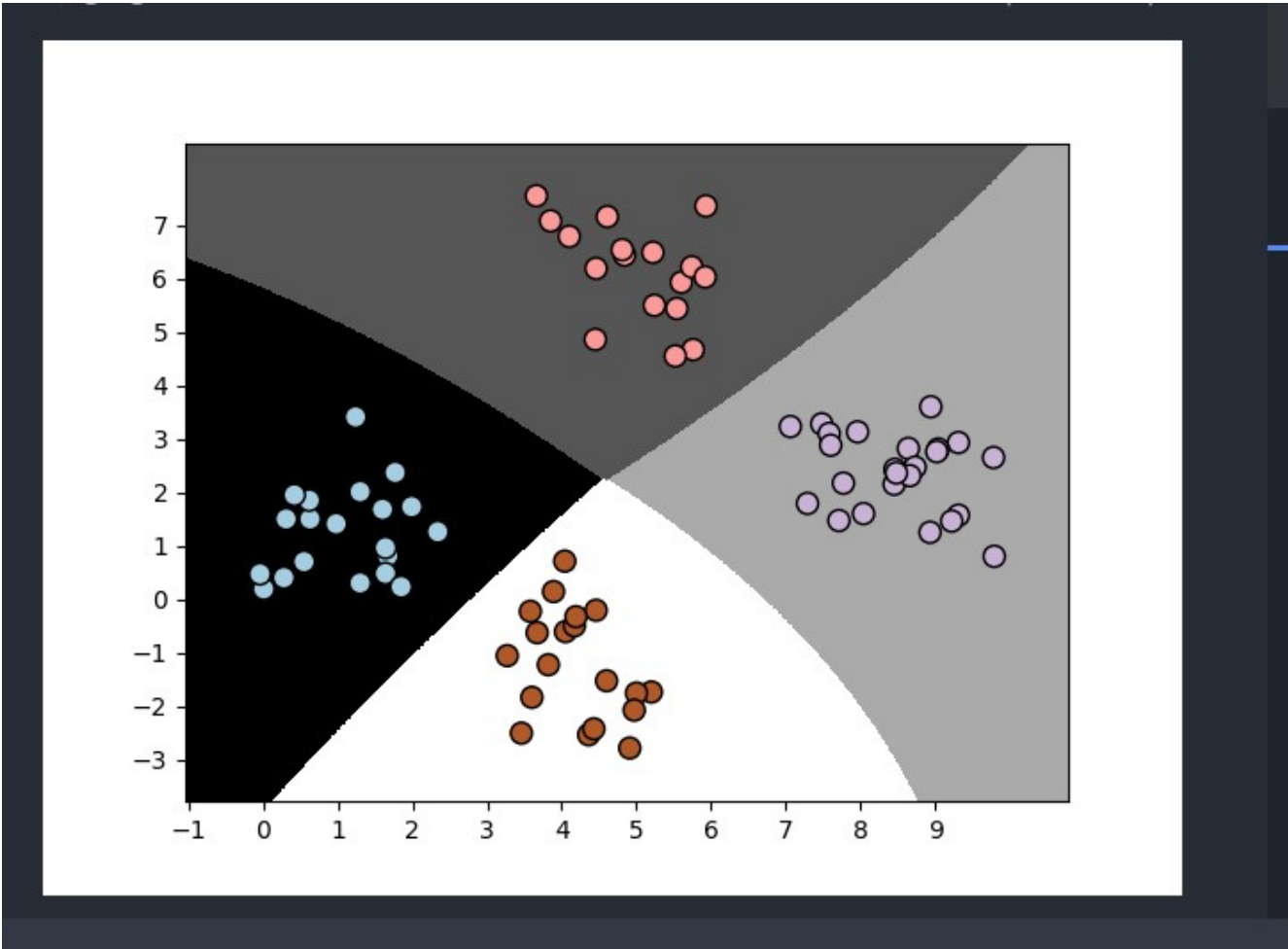


Рис. 20. Повторний результат виконання програми.

Отже, як можна побачити з рис.17-20 результати виконань повністю співпадають, тобто класифікація точна.

## Завдання 2.5. Вивчити метрики якості класифікації

```

LR_1_task_5.py
1  import pandas as pd
2  import numpy as np
3  import matplotlib.pyplot as plt
4  from sklearn.metrics import confusion_matrix, accuracy_score, recall_score, precision_score, f1_score, roc_curve, \
5  roc_auc_score
6
7  df = pd.read_csv('data_metrics.csv')
8  df.head()
9
10 threshold = 0.5
11 df['predicted_RF'] = (df.model_RF >= 0.5).astype('int')
12 df['predicted_LR'] = (df.model_LR >= 0.5).astype('int')
13 df.head()
14
15 confusion_matrix(df.actual_label.values, df.predicted_RF.values)
16
17
18 2 usages
19 def kulish_find_TP(y_true, y_pred):
20     # counts the number of true positives (y_true = 1, y_pred = 1)
21     return sum((y_true == 1) & (y_pred == 1))
22
23 2 usages
24 def kulish_find_FN(y_true, y_pred):
25     # counts the number of false negatives (y_true = 1, y_pred = 0)
26     return sum((y_true == 1) & (y_pred == 0))
27

```

Рис. 21. Відображення частини коду програми.

```

/home/xtr99/labs/ai/lab01/venv/bin/python /home/xtr99/labs/ai/lab01/LR_1_task_5.py
TP: 5047
FN: 2832
FP: 2360
TN: 5519
Accuracy RF: 0.671
Accuracy LR: 0.616
Recall RF: 0.641
Recall LR: 0.543
Precision RF: 0.681
Precision LR: 0.636
F1 RF: 0.660
F1 LR: 0.586

```

Рис. 22. Результати виконання програми.

		Куліш М.В.			ЖИТОМИРСЬКА ПОЛІТЕХНІКА.23.125.08.000 - Лр1	Арк.
		Пулеко І.В.				11
Змн.	Арк.	№ докум.	Підпис	Дата		

```

scores with threshold = 0.25
Accuracy RF: 0.502
Recall RF: 1.000
Precision RF: 0.501
F1 RF: 0.668
AUC RF:0.738
AUC LR:0.666

Process finished with exit code 0

```

Рис. 23. Результат виконання для порогу 0.25

```

F1 LR: 0.586
scores with threshold = 0.5
Accuracy RF: 0.671
Recall RF: 0.641
Precision RF: 0.681
F1 RF: 0.660

```

Рис. 24. Результат виконання для порогу 0.5

Як можна побачити з рис.23-24 при порозі 0.25: акуратність зменшилась, чутливість більша, точність стала меншою і оцінка змінилися трохи.

На графіку(рис. 25) продемонстровано що RF модель є більш зрозумілою, проти LR моделі. Але залежить також і від складності моделі. Тому це не завжди є основним показником.

		Куліш М.В.			ЖИТОМИРСЬКА ПОЛІТЕХНІКА.23.125.08.000 - Лр1	Арк.
		Пулеко І.В.				12
Змн.	Арк.	№ докум.	Підпис	Дата		



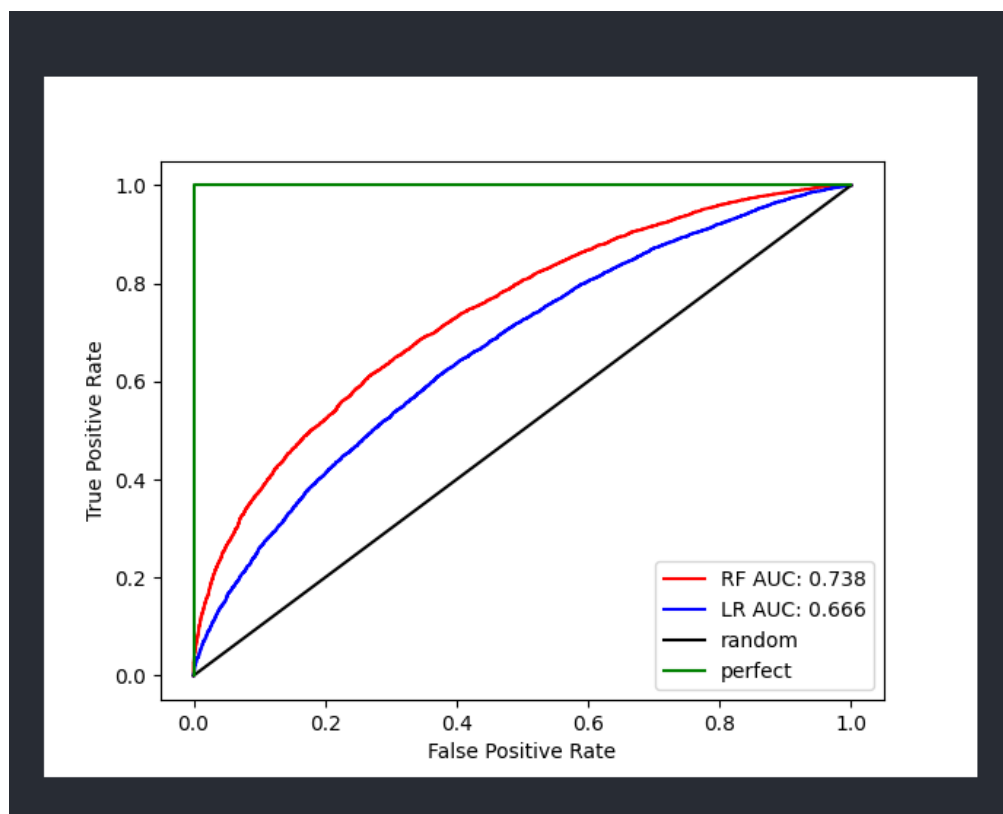


Рис. 25. Графік кривої ROC з урахуванням AUC

**Завдання 2.6. Розробіть програму класифікації даних в файлі data\_multivar\_nb.txt за допомогою машини опорних векторів (Support Vector Machine - SVM). Розрахуйте показники якості класифікації. Порівняйте їх з показниками наївного байєсівського класифікатора. Зробіть висновки яку модель класифікації краще обрати і чому.**

```

1  import numpy as np
2  from sklearn.model_selection import train_test_split
3  from sklearn import svm
4  from sklearn import metrics
5
6  from utilities import visualize_classifier
7  input_file = 'data_multivar_nb.txt'
8  data = np.loadtxt(input_file, delimiter=',')
9  X, y = data[:, :-1], data[:, -1]
10 X_train, X_test, y_train, y_test = train_test_split(X, y.astype(int), test_size=0.2, random_state=3)
11 cls = svm.SVC(kernel='linear')
12 cls.fit(X_train, y_train)
13 pred = cls.predict(X_test)
14 print("Accuracy:", metrics.accuracy_score(y_test, y_pred=pred))
15 print("Precision:", metrics.precision_score(y_test, y_pred=pred, average='macro'))
16 print("Recall", metrics.recall_score(y_test, y_pred=pred, average='macro'))
17 print(metrics.classification_report(y_test, y_pred=pred))
18 visualize_classifier(cls, X_test, y_test)

```

Рис. 26. Відображення коду програми.

		Куліш М.В.			ЖИТОМИРСЬКА ПОЛІТЕХНІКА.23.125.08.000 - Лр1	Арк.
		Пулеко І.В.				13
Змн.	Арк.	№ докум.	Підпис	Дата		

```

/home/xtr99/labs/ai/lab01/venv/bin/python /home/xtr99/labs/ai/lab01/LR_1_task_6.py
Accuracy: 1.0
Precision: 1.0
Recall 1.0

```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	20
1	1.00	1.00	1.00	17
2	1.00	1.00	1.00	24
3	1.00	1.00	1.00	19
accuracy			1.00	80
macro avg	1.00	1.00	1.00	80
weighted avg	1.00	1.00	1.00	80

```

|
Process finished with exit code 0

```

Рис. 27. Результати виконання програми.

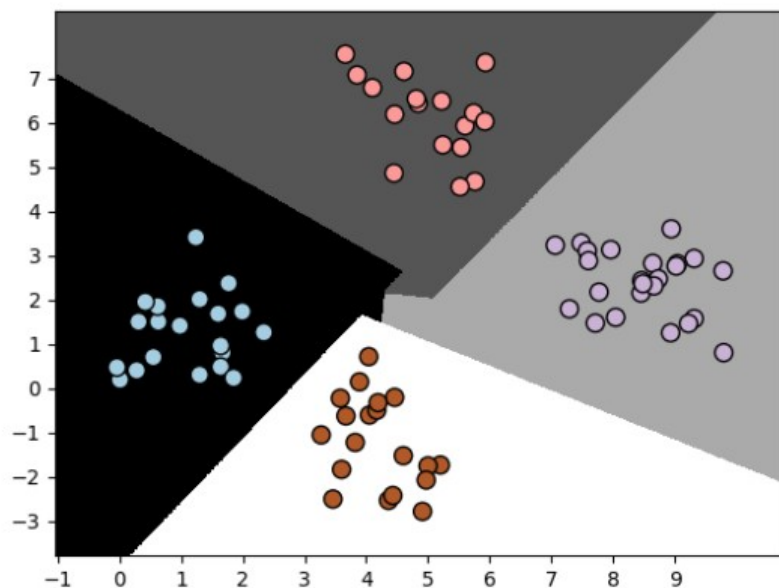


Рис. 28. Результати виконання програми.

Покликання на github: <https://github.com/mrkulish/ai-labs/tree/master/lab01>

**Висновок:** якість класифікації при використанні SVM становить 100%, тоді як при використанні наївного байєсівського класифікатора досягає 99.75%. Однак, враховуючи, що при наївному байєсівському класифікаторі значення порогу відіграють важливу роль і можуть впливати на результати, навіть дуже невелике відхилення від 100% може стати вирішальним. Тому, хоча різниця між значенням 99.75% та 100% є дуже незначною, краще використовувати модель SVM, оскільки вона працює без залежності від певних встановлених значень або порогів. Це робить її більш надійною та стабільною моделлю класифікації.

		Куліш М.В.			ЖИТОМИРСЬКА ПОЛІТЕХНІКА.23.125.08.000 - Лр1	Арк.
		Пулеко І.В.				15
Змн.	Арк.	№ докум.	Підпис	Дата		