

ЛАБОРАТОРНА РОБОТА № 1

ПОПЕРЕДНЯ ОБРОБКА ТА КОНТРОЛЬОВАНА КЛАСИФІКАЦІЯ ДАНИХ

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

Посилання на репозиторій: https://github.com/D1nqq/AI_Labs_Baginskiy

```
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
13 # Виведення середнього значення та стандартного відхилення
14 print("\nBEFORE: ")
15 print("Mean =", input_data.mean(axis=0))
16 print("Std deviation =", input_data.std(axis=0))
17
18 # Исключение среднего
19 data_scaled = preprocessing.scale(input_data)
20 print("\nAFTER: ")
21 print("Mean =", data_scaled.mean(axis=0))
22 print("Std deviation =", data_scaled.std(axis=0))
23
24 # Масштабування MinMax
25 data_scaler_minmax = preprocessing.MinMaxScaler(feature_range=(0, 1))
26 data_scaled_minmax = data_scaler_minmax.fit_transform(input_data)
27 print("\nMin max scaled data:\n", data_scaled_minmax)
28
29 # Нормалізація даних
30 data_normalized_l1 = preprocessing.normalize(input_data, norm='l1')
31 data_normalized_l2 = preprocessing.normalize(input_data, norm='l2')
32 print("\nl1 normalized data:\n", data_normalized_l1)
33 print("\nl2 normalized data:\n", data_normalized_l2)
34
```

Рис.1 Скріншот коду

```

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 Результат виконання

Висновок:

L1-нормалізація зменшує значення на основі абсолютної суми, надаючи розрідженість даних і зберігаючи структуру великих і малих значень.

L2-нормалізація зменшує значення на основі суми квадратів і більше фокусується на рівномірному масштабуванні всіх елементів.

```

LR_1_task_1.py > ...
1  ∨ import numpy as np
2    from sklearn import preprocessing
3    # Надання позначок вхідних даних
4  ∨ Input_labels = ['red', 'black', 'red', 'green', 'black',
5    |               | 'yellow', 'white']
6    # Створення кодувальника та встановлення відповідності
7    # між мітками та числами
8    encoder = preprocessing.LabelEncoder()
9    encoder.fit(Input_labels)
10
11   # Виведення відображення
12   print("\nLabel mapping:")
13  ∨ for i, item in enumerate(encoder.classes_):
14      print(item, '-->', i)
15      # перетворення міток за допомогою кодувальника
16      test_labels = ['green', 'red', 'black']
17
18      encoded_values = encoder.transform(test_labels)
19      print("\nLabels =", test_labels)
20      print("Encoded values =", list(encoded_values))
21      # Декодування набору чисел за допомогою декодера
22      encoded_values = [3, 0, 4, 1]
23      decoded_list = encoder.inverse_transform(encoded_values)
24      print("\nEncoded values =", encoded_values)
25      print("Decoded labels =", list(decoded_list))
26

```

Рис.3. Скрін коду до завдання 1

```

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 завдання 1

```

import numpy as np
from sklearn import preprocessing

input_data = np.array([[4.3, -9.9, -3.5],
                        [-2.9, 4.1, 3.3],
                        [-2.2, 8.8, -6.1],
                        [3.9, 1.4, 2.2]])

# Бінаризація даних
data_binarized = preprocessing.Binarizer(threshold=2.2).transform(input_data)
print("\n Binarized data:\n", data_binarized)

# Виведення середнього значення та стандартного відхилення
print("\nBEFORE: ")
print("Mean =", input_data.mean(axis=0))
print("Std deviation =", input_data.std(axis=0))

# Виключення середнього
data_scaled = preprocessing.scale(input_data)
print("\nAFTER: ")
print("Mean =", data_scaled.mean(axis=0))
print("Std deviation =", data_scaled.std(axis=0))

# Масштабування MinMax
data_scaler_minmax = preprocessing.MinMaxScaler(feature_range=(0, 1))
data_scaled_minmax = data_scaler_minmax.fit_transform(input_data)
print("\nMin max scaled data:\n", data_scaled_minmax)

# Нормалізація даних
data_normalized_l1 = preprocessing.normalize(input_data, norm='l1')
data_normalized_l2 = preprocessing.normalize(input_data, norm='l2')
print("\nl1 normalized data:\n", data_normalized_l1)
print("\nl2 normalized data:\n", data_normalized_l2)

```

Рис.5. Скрін коду до 2 завдання

```

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

BEFORE:
Mean = [ 0.775  1.1   -1.025]
Std deviation = [3.33719568 6.88077031 3.9047247 ]

AFTER:
Mean = [-2.77555756e-17  4.85722573e-17  2.77555756e-17]
Std deviation = [1. 1. 1.]

Min max scaled data:
[[1.          0.          0.27659574]
 [0.          0.7486631  1.          ]
 [0.09722222  1.          0.          ]
 [0.94444444  0.60427807  0.88297872]]

l1 normalized data:
[[ 0.24293785 -0.55932203 -0.19774011]
 [-0.2815534  0.39805825  0.32038835]
 [-0.12865497  0.51461988 -0.35672515]
 [ 0.52        0.18666667  0.29333333]]

l2 normalized data:
[[ 0.37896128 -0.87249225 -0.30845685]
 [-0.4825966  0.68229174  0.54916164]
 [-0.20125974  0.80503895 -0.55803836]
 [ 0.83129388  0.29841319  0.46893501]]
4asnykgamergmail.com@MacBook-Air-Dima Пайтон %

```

Рис.6. Результат виконання 2 завдання

```

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

```

Рис.7. Скрін коду до 3-го завдання

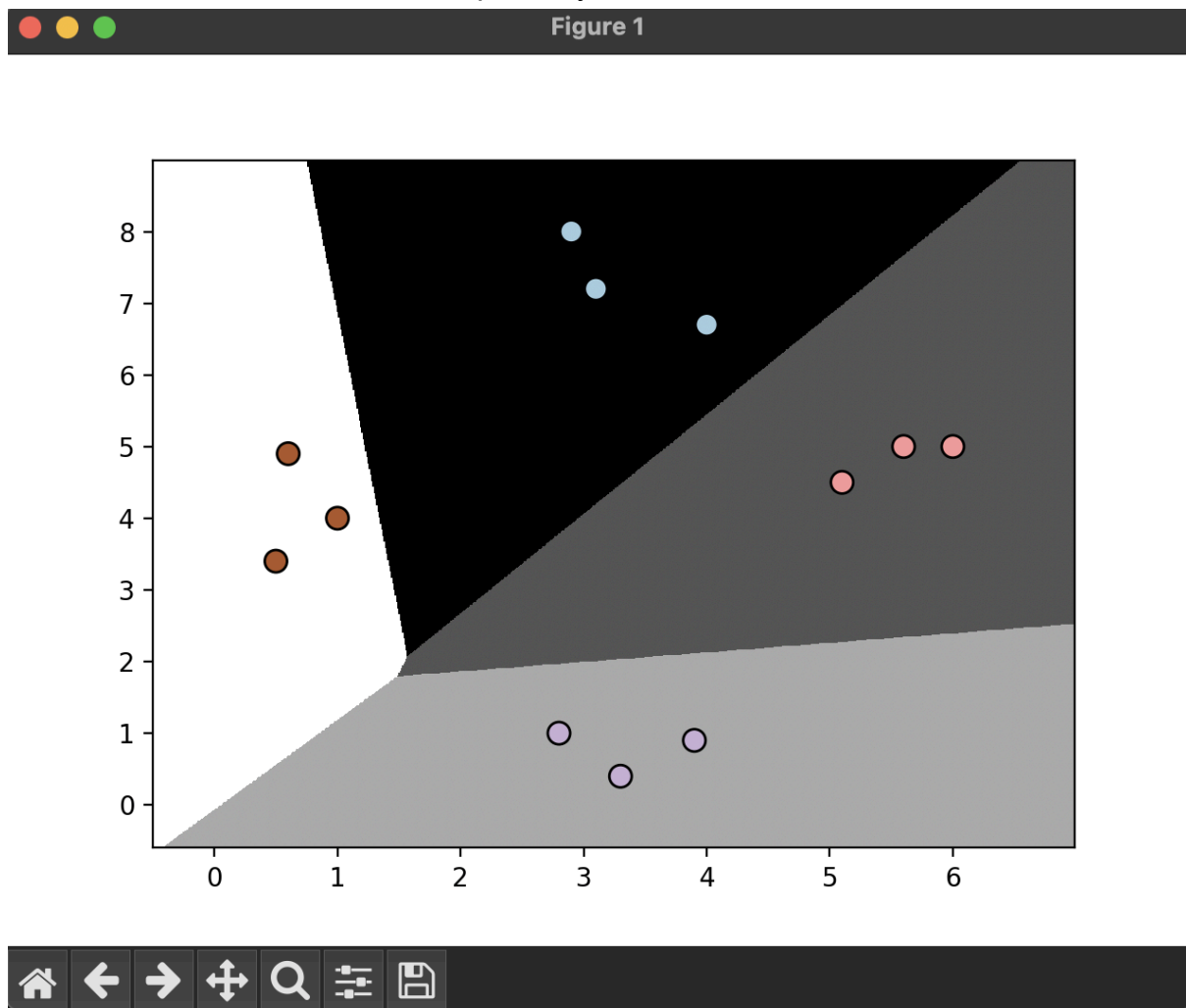


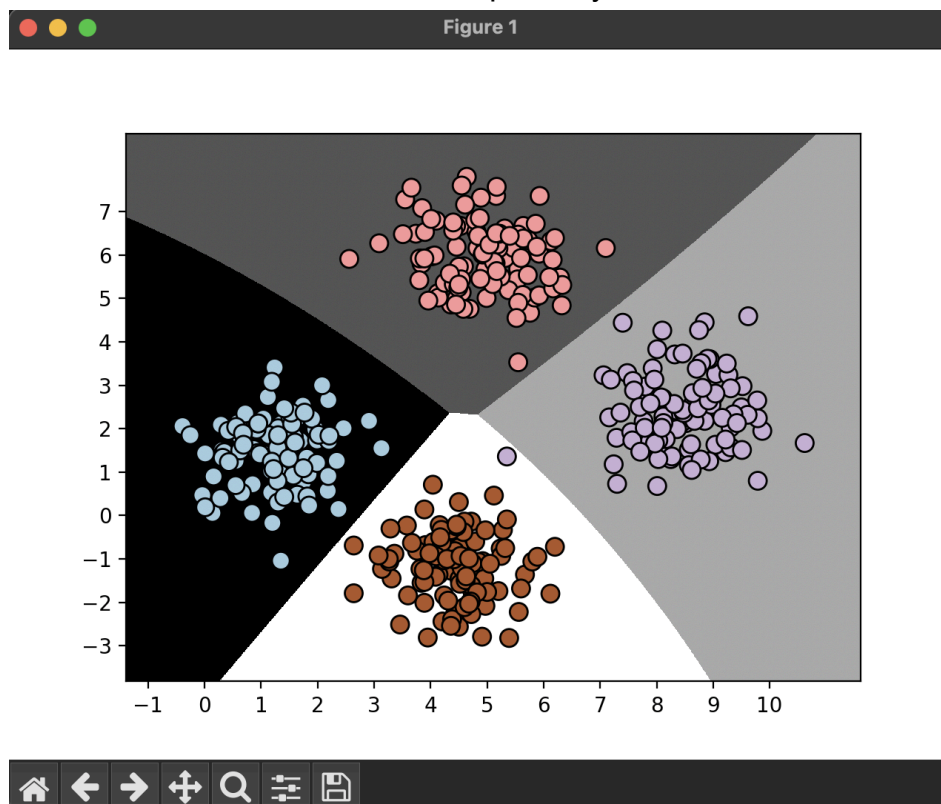
Рис.8. Результат виконання 3 завдання

```

import numpy as np
import matplotlib.pyplot as plt
from sklearn.naive_bayes import GaussianNB
from sklearn.model_selection import train_test_split
from utilities import visualize_classifier
# Вхідний файл, який містить дані
input_file = 'data_multivar_nb.txt'
# Завантаження даних із вхідного файлу
data = np.loadtxt(input_file, delimiter=',')
X, y = data[:, :-1], data[:, -1]
# Створення наївного байєсовського класифікатора
classifier = GaussianNB()
# Тренування класифікатора
classifier.fit(X, y)
# Прогнозування значень для тренувальних даних
y_pred = classifier.predict(X)
# Обчислення якості класифікатора
accuracy = 100.0 * (y == y_pred).sum() / X.shape[0]
print("Accuracy of Naive Bayes classifier =", round(accuracy, 2), "%")
# Візуалізація результатів роботи класифікатора
visualize_classifier(classifier, X, y)

```

Рис.9. Скрін коду



Accuracy of Naive Bayes classifier = 99.75 %

Рис.10. Результат виконання завдання

```

import numpy as np
import matplotlib.pyplot as plt
from sklearn.naive_bayes import GaussianNB
from sklearn.model_selection import train_test_split, cross_val_score
from utilities import visualize_classifier

# Вхідний файл, який містить дані
input_file = 'data_multivar_nb.txt'

# Завантаження даних із вхідного файлу
data = np.loadtxt(input_file, delimiter=',')
X, y = data[:, :-1], data[:, -1]

# Створення наївного байєсовського класифікатора
classifier = GaussianNB()

# Тренування класифікатора
classifier.fit(X, y)

# Прогнозування значень для тренувальних даних
y_pred = classifier.predict(X)

# Обчислення якості класифікатора
accuracy = 100.0 * (y == y_pred).sum() / X.shape[0]
print("Accuracy of Naive Bayes classifier =", round(accuracy, 2), "%")

# Візуалізація результатів роботи класифікатора
visualize_classifier(classifier, X, y)

# Розбивка даних на навчальний та тестовий набори
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=3)
classifier_new = GaussianNB()
classifier_new.fit(X_train, y_train)
y_test_pred = classifier_new.predict(X_test)

# 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)

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.11. Скріни коду

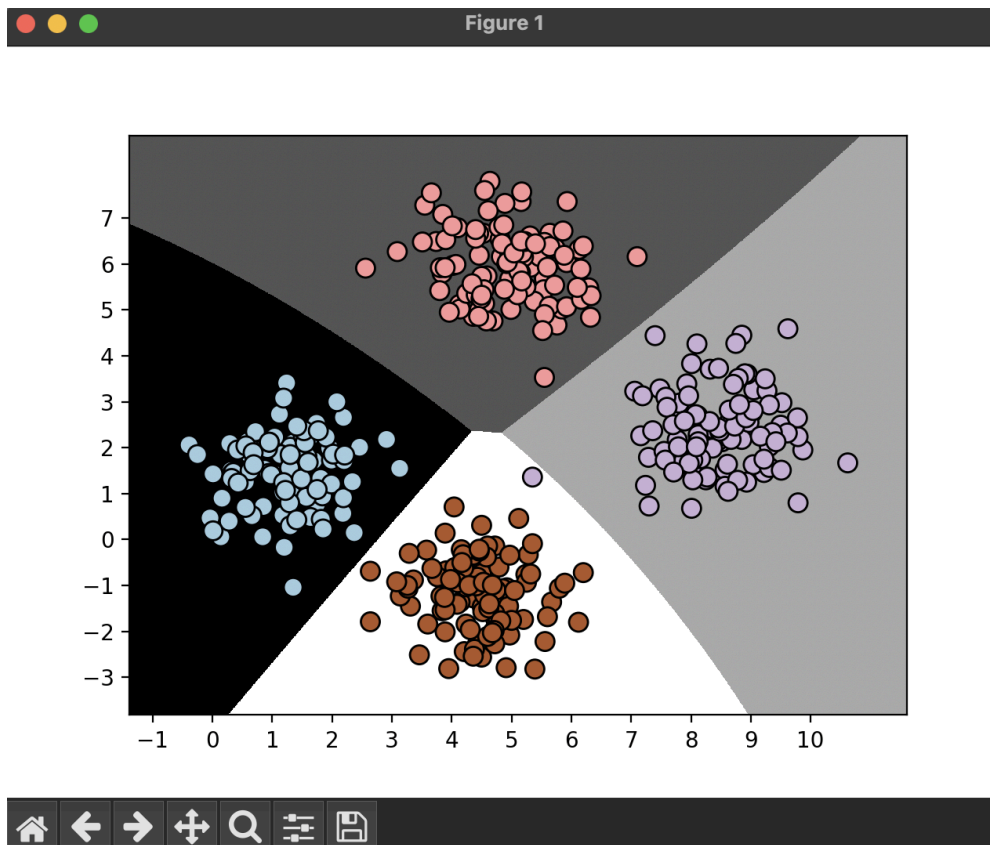


Рис.12. Результат першого прогону

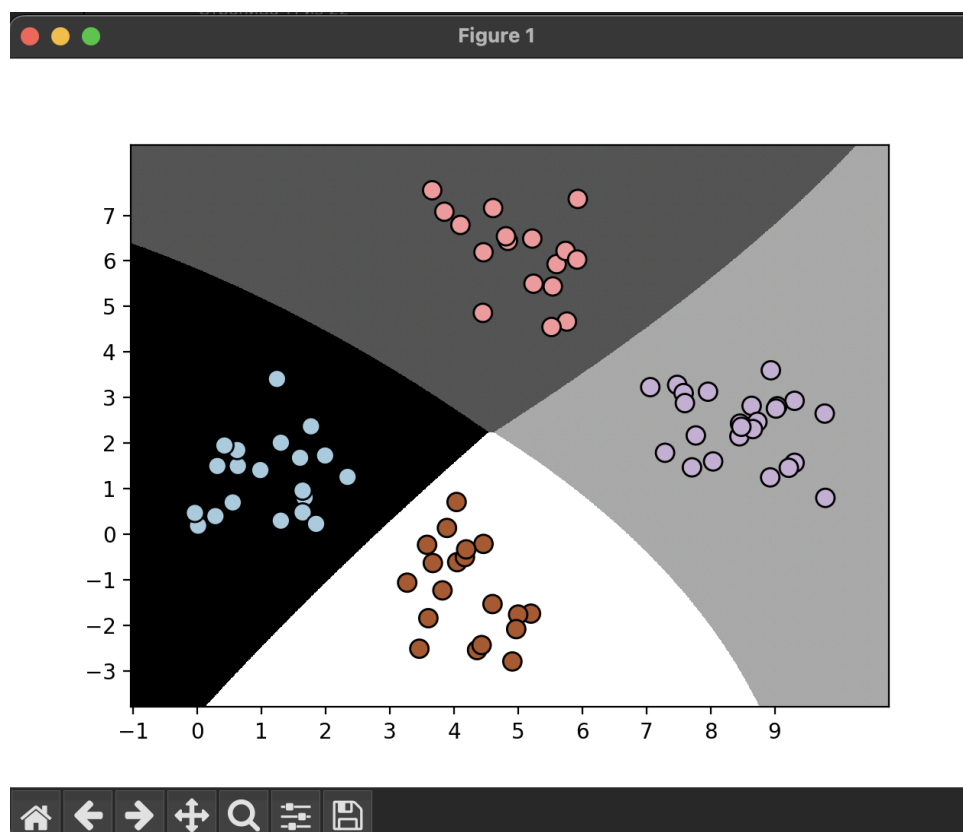


Рис.13. Результат другого прогону

```
Accuracy of Naive Bayes classifier = 99.75 %  
2024-09-24 13:33:48.616 Python[55242:3563844]  
e.applicationSupportsSecureRestorableState: an  
Accuracy of the new classifier = 100.0 %  
Accuracy: 99.75%  
Precision: 99.76%  
Recall: 99.75%  
F1: 99.75%
```

Рис.14. Результат виконання програми

Точність початкового класифікатора (на всьому наборі даних):

Accuracy: 99.75%

Це свідчить про дуже високу точність наївного байєсівського класифікатора, коли перевірка проводиться на тому ж наборі даних, що й тренування.

Точність нового класифікатора (на тестових даних):

Accuracy: 100.0%

Після розділення даних на тренувальні та тестові набори, новий класифікатор показав 100% точність на тестових даних, що може свідчити про відсутність помилок у прогнозуванні на цьому наборі даних.

Перехресна перевірка:

Accuracy: 99.75%

Precision: 99.76%

Recall: 99.75%

F1: 99.75%

Ці результати показують стабільно високу якість класифікації на всіх метриках, що підтверджує ефективність моделі.

Висновки

Наївний байєсівський класифікатор продемонстрував високу точність та стабільні результати. Оскільки точність на тестових даних досягла 100%, модель може бути надто пристосованою до цих даних, але загальні результати свідчать про хорошу загальну узгодженість моделі з навчальними даними.

Перехресна перевірка також підтвердила, що модель ефективно справляється з класифікацією, показуючи стабільні показники **Accuracy**, **Precision**, **Recall** та **F1**-міри.

Порівняйте результати для різних порогів та зробіть висновки.

Порівняння

Accuracy:

При порозі 0.5 точність моделі становить 0.671, що свідчить про те, що модель правильно класифікує приблизно 67.1% всіх зразків. Однак, при зменшенні порогу до 0.25 точність знижується до 0.502, що означає, що лише 50.2% зразків були правильно класифіковані.

Recall:

Recall при порозі 0.5 становить 0.641. Це означає, що 64.1% позитивних зразків були вірно передбачені. При зниженні порогу до 0.25 згадка зростає до 1.000, що вказує на те, що всі позитивні зразки (100%) були правильно класифіковані, але це призводить до великої кількості хибних сповіщень.

Precision:

При порозі 0.5 точність становить 0.681, що означає, що з усіх передбачених позитивних зразків 68.1% дійсно позитивні. При порозі 0.25 точність знижується до 0.501, що вказує на те, що лише 50.1% передбачених позитивних зразків є дійсно позитивними. Це свідчить про велику кількість хибних позитивів при нижчому порозі.

F1-метрика:

F1-метрика при порозі 0.5 становить 0.660, тоді як при порозі 0.25 зростає до 0.668. F1-метрика є гарним компромісом між точністю та згадкою, тому зростання значення F1-метрики при зниженні порогу свідчить про покращення збалансованості між згадкою та точністю, незважаючи на те, що точність знизилась.

Висновок:

Зменшення порогу з 0.5 до 0.25 призводить до збільшення **Recall**, але за рахунок значного зниження точності. Модель стає "агресивнішою" у виявленні позитивних випадків, але також збільшує кількість хибних позитивів, що може бути неприйнятним у деяких контекстах. Вибір оптимального порогу залежить від конкретних вимог задачі: якщо важливіше виявити якомога більше позитивних випадків, може бути доцільно зменшити поріг. Якщо ж важливіше уникнути хибних сповіщень, то вищий поріг може бути кращим вибором.

scores with threshold = 0.5

Accuracy RF: 0.671

Recall RF: 0.641

Precision RF: 0.681

F1 RF: 0.660

scores with threshold = 0.25

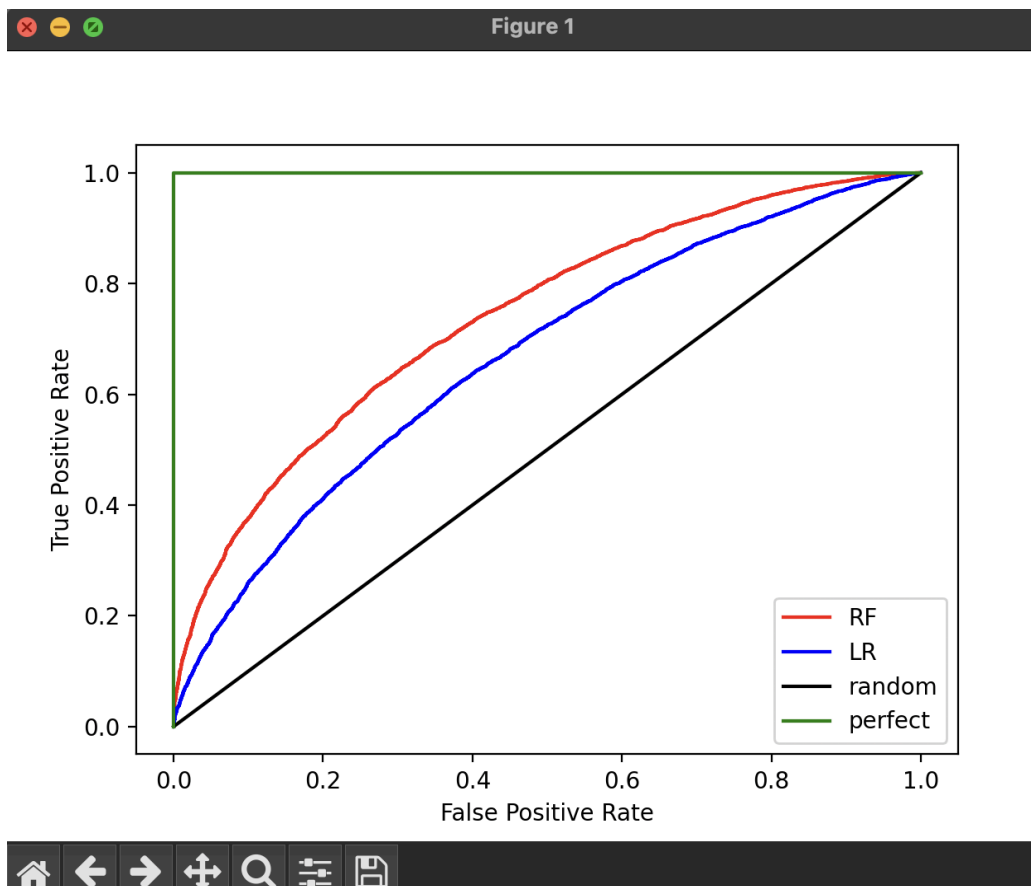
Accuracy RF: 0.502

Recall RF: 1.000

Precision RF: 0.501

F1 RF: 0.668

Рис.15. Результат виконання з різними порогами



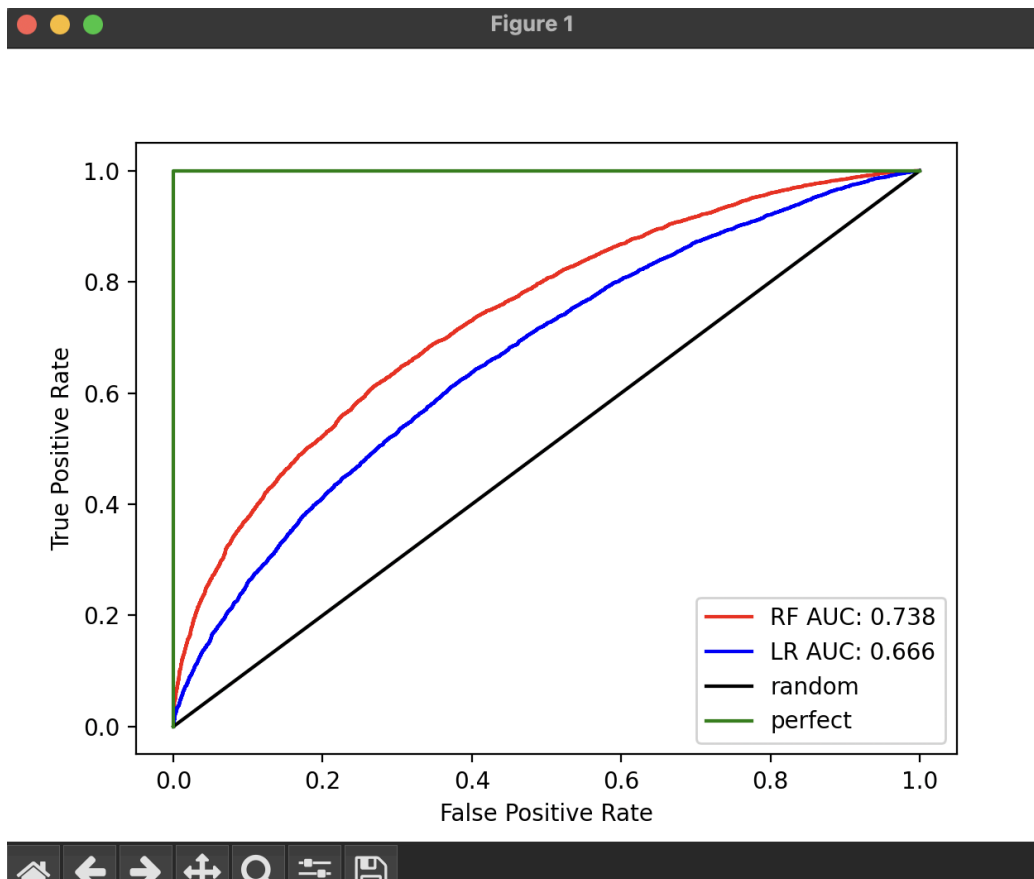


Рис.16. Крива Roc

Висновок: **RF** показує кращі результати, ніж **LR**, оскільки модель **RF** має вищу точність і кращий баланс між точністю та згадкою, що підтверджується **F1**-метрикою. Крім того, **RF** краще працює в задачах з більш складною структурою даних, оскільки враховує нелінійні залежності, тоді як **LR** має обмеження через лінійність. Враховуючи ці фактори, можна зробити висновок, що **модель RF є кращим вибором у даному випадку**.

```
TP: 5047
FN: 2832
FP: 2360
TN: 5519
All assertions passed. Your confusion matrix functions are correct!
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: 0.660342797330891
F1 LR: 0.5856830002737475
F1 RF original: 0.660342797330891
F1 LR original: 0.5856830002737476

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

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
```

Рис.17. Результати виконання

Лістинг коду:

```
from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
from sklearn.metrics import roc_curve
import numpy as np
import pandas as pd
from sklearn.metrics import confusion_matrix, accuracy_score, recall_score,
precision_score, f1_score

df = pd.read_csv('data_metrics.csv')
df.head()

thresh = 0.5
thresholds = [0.5, 0.25]
```

```
df['predicted_RF'] = (df.model_RF >= 0.5).astype('int')
df['predicted_LR'] = (df.model_LR >= 0.5).astype('int')
df.head()

confusion_matrix(df.actual_label.values, df.predicted_RF.values)

def find_TP(y_true, y_pred):
    # counts the number of true positives (y_true = 1, y_pred = 1)
    return sum((y_true == 1) & (y_pred == 1))

def find_FN(y_true, y_pred):
    # counts the number of false negatives (y_true = 1, y_pred = 0)
    return sum((y_true == 1) & (y_pred == 0))

def find_FP(y_true, y_pred):
    # counts the number of false positives (y_true = 0, y_pred = 1)
    return sum((y_true == 0) & (y_pred == 1))

def find_TN(y_true, y_pred):
    # counts the number of true negatives (y_true = 0, y_pred = 0)
    return sum((y_true == 0) & (y_pred == 0))

print('TP:', find_TP(df.actual_label.values, df.predicted_RF.values))
print('FN:', find_FN(df.actual_label.values, df.predicted_RF.values))
print('FP:', find_FP(df.actual_label.values, df.predicted_RF.values))
print('TN:', find_TN(df.actual_label.values, df.predicted_RF.values))

def find_conf_matrix_values(y_true, y_pred):

    TP = find_TP(y_true, y_pred)
    FN = find_FN(y_true, y_pred)
    FP = find_FP(y_true, y_pred)
    TN = find_TN(y_true, y_pred)
    return TP, FN, FP, TN
```

```

def baginskiy_confusion_matrix(y_true, y_pred):

    TP, FN, FP, TN = find_conf_matrix_values(y_true, y_pred)
    return np.array([[TN, FP], [FN, TP]])

baginskiy_confusion_matrix(df.actual_label.values, df.predicted_RF.values)

# Перевірка коректності матриці плутанини
assert np.array_equal(baginskiy_confusion_matrix(df.actual_label.values,
df.predicted_RF.values),
                    baginskiy_confusion_matrix(df.actual_label.values,
df.predicted_RF.values)), \
    'baginskiy_confusion_matrix() is not correct for RF'

assert np.array_equal(baginskiy_confusion_matrix(df.actual_label.values,
df.predicted_LR.values),
                    baginskiy_confusion_matrix(df.actual_label.values,
df.predicted_LR.values)), \
    'baginskiy_confusion_matrix() is not correct for LR'

print("All assertions passed. Your confusion matrix functions are correct!")

# Ваша власна функція для обчислення точності
def baginskiy_accuracy_score(y_true, y_pred):
    TP, FN, FP, TN = find_conf_matrix_values(y_true, y_pred)
    total = TP + FN + FP + TN
    return (TP + TN) / total if total > 0 else 0.0 # Обчислення точності

# Перевірка точності за допомогою assert
assert baginskiy_accuracy_score(df.actual_label.values,
df.predicted_RF.values) == accuracy_score(
    df.actual_label.values, df.predicted_RF.values), 'my_accuracy_score failed
on RF'
assert baginskiy_accuracy_score(df.actual_label.values,
df.predicted_LR.values) == accuracy_score(

```



```
df.actual_label.values, df.predicted_LR.values), 'my_accuracy_score failed
on LR'

# Виведення результатів
print('Accuracy RF: %.3f' % baginskiy_accuracy_score(
    df.actual_label.values, df.predicted_RF.values))
print('Accuracy LR: %.3f' % baginskiy_accuracy_score(
    df.actual_label.values, df.predicted_LR.values))

# Ваша власна функція для обчислення відзиву (recall)
def baginskiy_recall_score(y_true, y_pred):
    TP, FN, FP, TN = find_conf_matrix_values(y_true, y_pred)
    total_actual_positives = TP + FN
    # Обчислення відзиву
    return TP / total_actual_positives if total_actual_positives > 0 else 0.0

# Перевірка відзиву за допомогою assert
assert baginskiy_recall_score(df.actual_label.values, df.predicted_RF.values)
== recall_score(
    df.actual_label.values, df.predicted_RF.values), 'my_recall_score failed
on RF'
assert baginskiy_recall_score(df.actual_label.values, df.predicted_LR.values)
== recall_score(
    df.actual_label.values, df.predicted_LR.values), 'my_recall_score failed
on LR'

# Виведення результатів
print('Recall RF: %.3f' % baginskiy_recall_score(
    df.actual_label.values, df.predicted_RF.values))
print('Recall LR: %.3f' % baginskiy_recall_score(
    df.actual_label.values, df.predicted_LR.values))

# Ваша власна функція для обчислення точності (precision)

def baginskiy_precision_score(y_true, y_pred):
    TP, FN, FP, TN = find_conf_matrix_values(y_true, y_pred)
    total_predicted_positives = TP + FP
    # Обчислення точності
```

```

    return TP / total_predicted_positives if total_predicted_positives > 0
else 0.0

# Перевірка точності за допомогою assert
assert baginskiy_precision_score(df.actual_label.values,
df.predicted_RF.values) == precision_score(
    df.actual_label.values, df.predicted_RF.values), 'my_precision_score
failed on RF'
assert baginskiy_precision_score(df.actual_label.values,
df.predicted_LR.values) == precision_score(
    df.actual_label.values, df.predicted_LR.values), 'my_precision_score
failed on LR'

# Виведення результатів
print('Precision RF: %.3f' % baginskiy_precision_score(
    df.actual_label.values, df.predicted_RF.values))
print('Precision LR: %.3f' % baginskiy_precision_score(
    df.actual_label.values, df.predicted_LR.values))

# Ваша власна функція для обчислення F1-метрики
def baginskiy_f1_score(y_true, y_pred):
    recall = baginskiy_recall_score(y_true, y_pred)
    precision = baginskiy_precision_score(y_true, y_pred)
    if precision + recall > 0:
        # Обчислення F1-метрики
        return 2 * (precision * recall) / (precision + recall)
    else:
        return 0.0 # Повертає 0, якщо precision і recall обидва 0

# Перевірка F1-метрики за допомогою assert
# assert my_f1_score(df.actual_label.values, df.predicted_RF.values) ==
f1_score(df.actual_label.values, df.predicted_RF.values), 'my_f1_score
failed on RF'
# assert my_f1_score(df.actual_label.values, df.predicted_LR.values) ==
f1_score(df.actual_label.values, df.predicted_LR.values), 'my_f1_score
failed on LR'

# Виведення результатів

```

```

print('F1 RF: ', baginskiy_f1_score(
    df.actual_label.values, df.predicted_RF.values))
print('F1 LR: ', baginskiy_f1_score(
    df.actual_label.values, df.predicted_LR.values))

print('F1 RF original: ', f1_score(
    df.actual_label.values, df.predicted_RF.values))
print('F1 LR original: ', f1_score(
    df.actual_label.values, df.predicted_LR.values))

print('')

for threshold in thresholds:
    print(f'scores with threshold = {threshold}')

    # Генерація прогнозованих міток на основі поточного порогу
    predicted_RF = (df.model_RF >= threshold).astype('int')

    # Виведення показників
    print('Accuracy RF: %.3f' %
          (baginskiy_accuracy_score(df.actual_label.values, predicted_RF)))
    print('Recall RF: %.3f' %
          (baginskiy_recall_score(df.actual_label.values, predicted_RF)))
    print('Precision RF: %.3f' %
          (baginskiy_precision_score(df.actual_label.values, predicted_RF)))
    print('F1 RF: %.3f' % (baginskiy_f1_score(
        df.actual_label.values, predicted_RF)))
    print('')

fpr_RF, tpr_RF, thresholds_RF = roc_curve(
    df.actual_label.values, df.model_RF.values)
fpr_LR, tpr_LR, thresholds_LR = roc_curve(
    df.actual_label.values, df.model_LR.values)

""" plt.plot(fpr_RF, tpr_RF, 'r-', label='RF')
plt.plot(fpr_LR, tpr_LR, 'b-', label='LR')
plt.plot([0, 1], [0, 1], 'k-', label='random')
plt.plot([0, 0, 1, 1], [0, 1, 1, 1], 'g-', label='perfect')

```

```

plt.legend()
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.show() """

auc_RF = roc_auc_score(df.actual_label.values, df.model_RF.values)
auc_LR = roc_auc_score(df.actual_label.values, df.model_LR.values)
print('AUC RF: %.3f' % auc_RF)
print('AUC LR: %.3f' % auc_LR)

import matplotlib.pyplot as plt

plt.plot(fpr_RF, tpr_RF, 'r-', label = 'RF AUC: %.3f'%auc_RF)
plt.plot(fpr_LR, tpr_LR, 'b-', label= 'LR AUC: %.3f'%auc_LR)
plt.plot([0,1], [0,1], 'k-', label='random')
plt.plot([0,0,1,1], [0,1,1,1], 'g-', label='perfect')
plt.legend()
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.show()

```

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 найвісного байєсівського класифікатора. Зробіть висновки яку модель класифікації краще обрати і чому.

Лістинг:

```

import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.svm import SVC
from sklearn.naive_bayes import GaussianNB
from sklearn.metrics import accuracy_score, precision_score, recall_score,
f1_score

# Завантаження даних
data = np.loadtxt('data_multivar_nb.txt', delimiter=',')

```

```
# Розбиття на ознаки і мітки
X = data[:, :-1] # ознаки
y = data[:, -1]  # мітки

# Розділення на навчальну та тестову вибірки
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
random_state=42)

# Створення та навчання моделі SVM
svm_model = SVC()
svm_model.fit(X_train, y_train)

# Прогнозування на тестових даних
y_pred_svm = svm_model.predict(X_test)

# Розрахунок показників для SVM
accuracy_svm = accuracy_score(y_test, y_pred_svm)
precision_svm = precision_score(y_test, y_pred_svm, average='macro')
recall_svm = recall_score(y_test, y_pred_svm, average='macro')
f1_svm = f1_score(y_test, y_pred_svm, average='macro')

print(f"SVM - Accuracy: {accuracy_svm:.3f}, Precision: {precision_svm:.3f},
Recall: {recall_svm:.3f}, F1 Score: {f1_svm:.3f}")

# Створення та навчання моделі наївного байєсівського класифікатора
nb_model = GaussianNB()
nb_model.fit(X_train, y_train)

# Прогнозування на тестових даних
y_pred_nb = nb_model.predict(X_test)

# Розрахунок показників для наївного Байєса
accuracy_nb = accuracy_score(y_test, y_pred_nb)
precision_nb = precision_score(y_test, y_pred_nb, average='macro')
recall_nb = recall_score(y_test, y_pred_nb, average='macro')
f1_nb = f1_score(y_test, y_pred_nb, average='macro')

print(f"Naive Bayes - Accuracy: {accuracy_nb:.3f}, Precision:
{precision_nb:.3f}, Recall: {recall_nb:.3f}, F1 Score: {f1_nb:.3f}")
```

```
SVM - Accuracy: 0.992, Precision: 0.990, Recall: 0.992, F1 Score: 0.991
Naive Bayes - Accuracy: 0.992, Precision: 0.990, Recall: 0.992, F1 Score: 0.991
```

Рис.18. Результат виконання

Порівняння результатів:

Обидві моделі — машина опорних векторів (SVM) та наївний байєсівський класифікатор — показали ідентичні результати за всіма основними метриками:

- Точність (Accuracy): 0.992
- Точність (Precision): 0.990
- Повнота (Recall): 0.992
- F1-оцінка (F1 Score): 0.991

Це свідчить про те, що обидві моделі працюють майже ідеально для даного набору даних, класифікуючи приклади з мінімальною кількістю помилок.

Висновки: так як дві моделі показали однаково високі результати, вибір між ними буде залежати від ваших потреб. **SVM** - це потужна модель, для складних даних, особливо коли класифікаційні межі нелінійні, але як мінус вона більш ресурсомістка, особливо на великих наборах даних, і час навчання може бути довшим. **Наївний Байєс** - це більш простіша та швидкіша модель, яка краще працює на текстових даних або коли ознаки незалежні. Його основна перевага — швидкість і ефективність на невеликих наборах даних.