

**Липецкий государственный технический университет**

**Факультет автоматизации и информатики**

**Кафедра автоматизированных систем управления**

**ЛАБОРАТОРНАЯ РАБОТА №1**

**по дисциплине «Прикладные интеллектуальные системы и экспертные  
системы»**

**Бинарная классификация фактографических данных**

Студент

Коровайцев А.А.

Группа М-ИАП-23-1

Руководитель

Кургасов В.В.

Доцент

Липецк 2023 г.

## Цель работы

Получить практические навыки решения задачи бинарной классификации данных в среде Jupiter Notebook. Научиться загружать данные, обучать классификаторы и проводить классификацию. Научиться оценивать точность полученных моделей.

### Задание кафедры

- 1) В среде Jupiter Notebook создать новый ноутбук (Notebook)
- 2) Импортировать необходимые для работы библиотеки и модули
- 3) Загрузить данные в соответствие с вариантом
- 4) Вывести первые 15 элементов выборки (координаты точек и метки класса)
- 5) Отобразить на графике сгенерированную выборку. Объекты разных классов должны иметь разные цвета.
- 6) Разбить данные на обучающую (train) и тестовую (test) выборки в пропорции 75% - 25% соответственно.
- 7) Отобразить на графике обучающую и тестовую выборки. Объекты разных классов должны иметь разные цвета.
- 8) Реализовать модели классификаторов, обучить их на обучающем множестве. Применить модели на тестовой выборке, вывести результаты классификации:
  - 9) Истинные и предсказанные метки классов
  - 10) Матрицу ошибок (confusion matrix)
  - 11) Значения полноты, точности, f1-меры и аккуратности
  - 12) Значение площади под кривой ошибок (AUC ROC)
  - 13) Отобразить на графике область принятия решений по каждому классу
- 14) В качестве методов классификации использовать:
- 15) Метод к-ближайших соседей ( $n\_neighbors = \{1, 3, 5, 9\}$ )
- 16) Наивный байесовский метод
- 17) Случайный лес ( $n\_estimators = \{5, 10, 15, 20, 50\}$ )
- 18) По каждому пункту работы занести в отчет программный код и результат вывода.
- 19) По результатам п.8 занести в отчет таблицу с результатами классификации всеми методами и выводы о наиболее подходящем методе классификации ваших данных.

20) Изучить, как изменится качество классификации, если на тестовую часть выделить 10% выборки, 35% выборки. Для этого повторить п.п. 6 – 10.

Вариант №7

Вид классов: moons

Random\_state: 77

noise: 0.25

Ход работы

Подготовка данных

Для генерации данных воспользуемся функцией `make_moons` из пакета `sklearn.datasets`. Результат генерации данных и вывод первых 15 значений представлен на рисунке 1.

### Генерация выборки

```
1 X, y = make_moons(n_samples=1000, shuffle=True, noise=0.25, random_state=77)
```

```
1 print('Координаты точек: ')
2 print(X[:15])
3 print('Метки класса: ')
4 print(y[:15])
```

Координаты точек:

```
[[ 0.3877582  0.83132146]
 [ 0.75917445  0.37546408]
 [ 0.16596943  1.06109846]
 [ 2.23428045  0.2359786 ]
 [-0.89666798  1.0952051 ]
 [ 0.94876632  0.31861216]
 [-0.81661113  0.04043469]
 [ 0.02592078  0.16408361]
 [-0.92023208  0.20859127]
 [ 0.27797801  0.63569972]
 [ 0.83512001  0.62834727]
 [ 0.0828706  -0.14748687]
 [ 0.15418065  0.92032556]
 [-0.39249897  1.18232379]
 [ 1.19731795 -0.30376657]]
```

Метки класса:

```
[0 0 0 1 0 0 0 1 0 0 0 1 0 0 1]
```

Рисунок 1 – Генерация данных и вывод первых 15-ти значений

Отобразим на графике сгенерированную выборку с выделением классов разными цветами. Для этого воспользуемся функцией `scatter` из библиотеки `matplotlib.pyplot`. Результат визуализации представлен на рисунке 2.

```

1 plt.scatter(X[:, 0], X[:, 1], c=y)
2 plt.show()

```

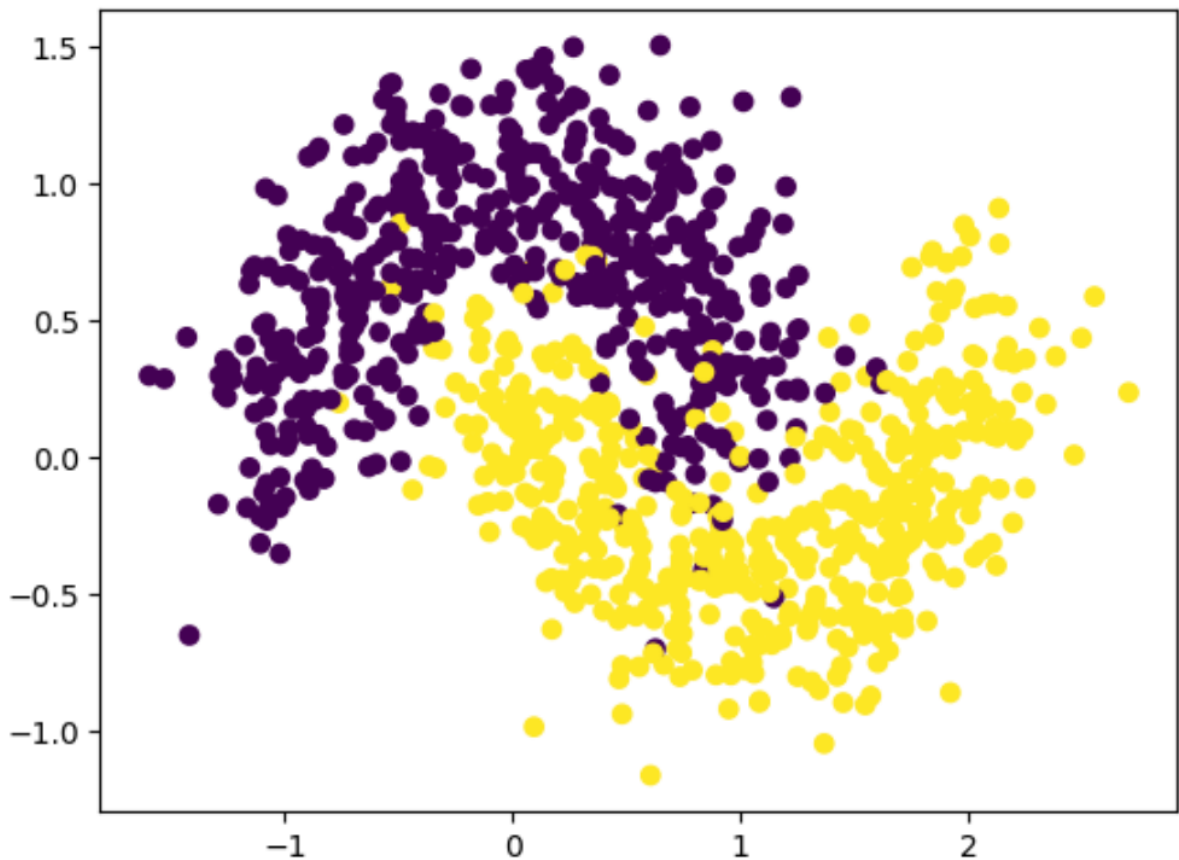


Рисунок 2 – Визуализация выборки

Разделим данных на обучающую и тестовую выборку. Для этого воспользуемся функцией `train_test_split` из пакета `sklearn.model_selection`. Скрипт для разделения данных представлен на рисунке 3. Результат разбиения выборки на тестовую и обучающую с последующей их визуализацией представлены на рисунке 4 и 5.

### Разбитие выборки на обучающее и тестовое множество (75/25)

```

1 X_train, X_test, y_train, y_test = train_test_split(X,
2                                                    y,
3                                                    test_size=0.25,
4                                                    random_state=77)

```

Рисунок 3 – Разделение выборки на обучающее и тестовое множество

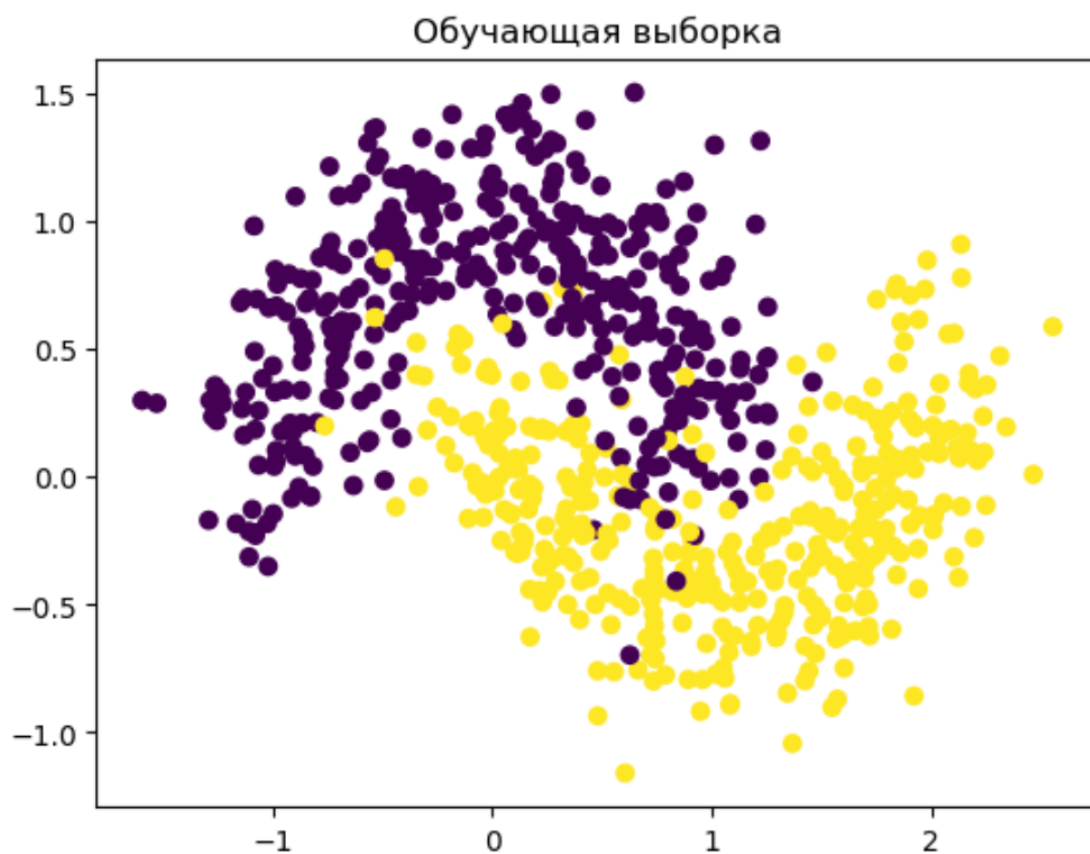


Рисунок 4 – Обучающая выборка

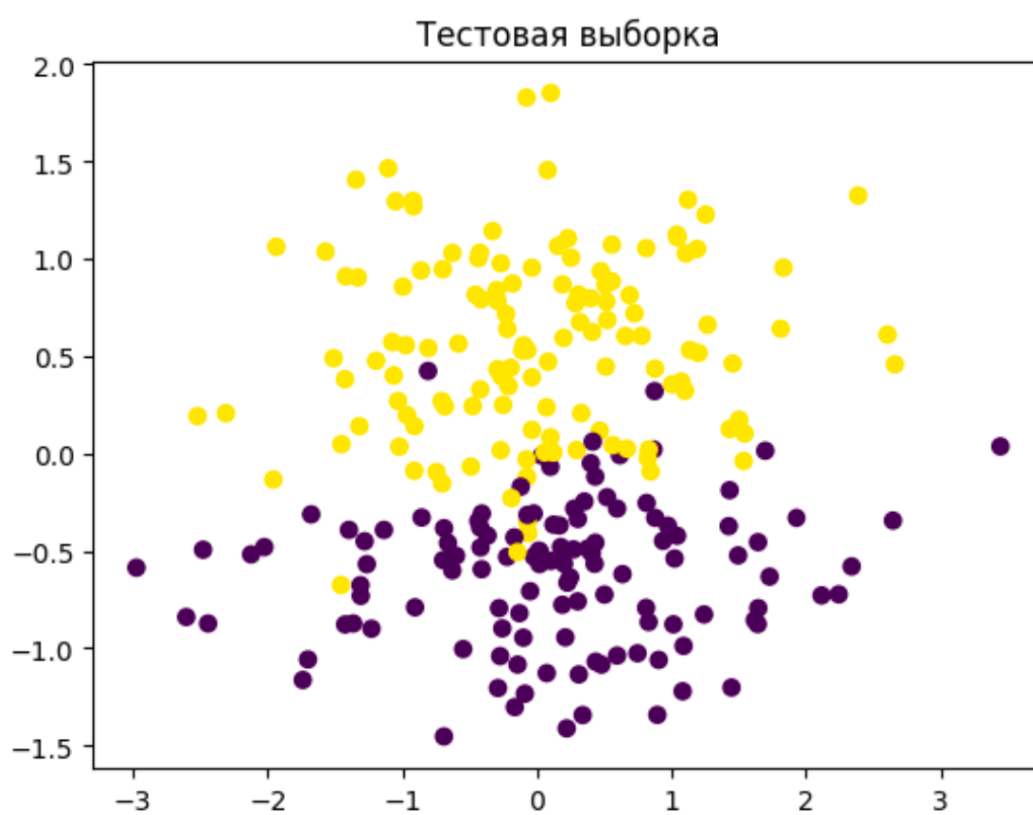


Рисунок 5 – Тестовая выборка

## Классификация с помощью метода к-ближайших соседей

Использование параметра `n_neighbors = 1`. Составленный код для использования данного классификатора с данным параметром представлен на рисунке 6. Полученная информация о точности классификации при использовании данного метода представлена на рисунке 7.

```
knn = KNeighborsClassifier(n_neighbors=1, metric='euclidean')  
  
# Обучаем модель данных  
knn.fit(X_train, y_train)  
  
# Оцениваем качество модели  
prediction = knn.predict(X_test)  
  
# Выводим сводную информацию  
show_info(knn, 'ближайшие соседи (1)', y_test, prediction)
```

Рисунок 6 – Код для классификатора с помощью метода к-ближайших соседей с параметром `n_neighbors = 1`



Метод классификации: ближайшие соседи (1)

Предсказанные и реальные значения:

```
[1 1 0 1 0 0 0 1 0 0 1 1 0 1 0 1 0 1 0 0 1 0 1 1 0 0 1 0 0 1 0 1 1
0 1 0 0 0 1 0 0 0 0 1 0 0 1 1 0 1 1 1 1 0 0 1 0 1 1 0 1 1 0 0 0 0 0 1
0 0 0 0 1 0 1 0 0 1 0 1 0 0 0 0 1 0 0 0 0 1 1 1 0 1 1 0 1 0 0 1 0 0 1
1 1 0 1 0 1 1 0 1 0 0 0 1 1 0 0 1 0 1 0 0 0 1 0 0 1 1 1 0 0 1 0 0 1 0 1 0
1 1 1 1 1 1 0 1 1 1 0 1 0 0 0 0 1 1 1 0 1 1 0 0 1 1 1 0 0 0 1 0 1 1 1 0 1
0 1 0 0 0 1 0 1 0 1 0 0 0 0 1 0 1 1 1 0 1 1 1 1 1 0 0 0 0 1 1 0 1 0 1 0 1
0 1 1 1 1 1 1 1 0 1 1 1 1 0 0 0 1 0 0 1 1 1 1 1 1 1 1 0 1]
[1 1 0 1 0 0 0 1 0 0 1 1 0 0 0 1 0 1 0 1 0 1 1 0 1 1 0 0 1 0 0 1 0 0 1
0 1 0 0 0 1 0 0 0 1 1 0 0 1 1 0 1 1 1 0 0 0 1 0 1 1 0 1 0 1 1 0 0 0 0 1
1 0 0 0 1 0 1 0 0 1 0 1 0 1 0 0 1 0 0 0 1 1 1 0 1 1 0 1 1 0 1 0 0 1 0 0 1
1 1 0 1 0 1 1 0 1 0 0 0 1 1 0 0 1 0 1 0 0 0 1 0 0 1 1 1 0 0 1 1 0 1 0 1 0
1 1 1 1 0 1 0 1 1 1 0 1 0 0 0 0 1 1 1 0 1 1 0 0 1 1 1 0 0 0 1 0 1 1 1 0 1
0 1 0 0 0 1 0 1 0 1 0 0 0 1 0 1 1 0 1 1 1 1 1 1 0 0 0 0 1 1 0 0 0 1 1 1
0 1 1 1 0 1 1 1 0 0 1 1 1 0 0 0 1 0 0 0 1 1 1 1 1 1 0 1 1 0 1]
```

Матрица неточностей

```
[[116  8]
 [  8 118]]
```

Точность классификации: 0.936

Полнота:

	precision	recall	f1-score	support
0	0.94	0.94	0.94	124
1	0.94	0.94	0.94	126
accuracy			0.94	250
macro avg	0.94	0.94	0.94	250
weighted avg	0.94	0.94	0.94	250

Площадь под кривой: 0.9359959037378393

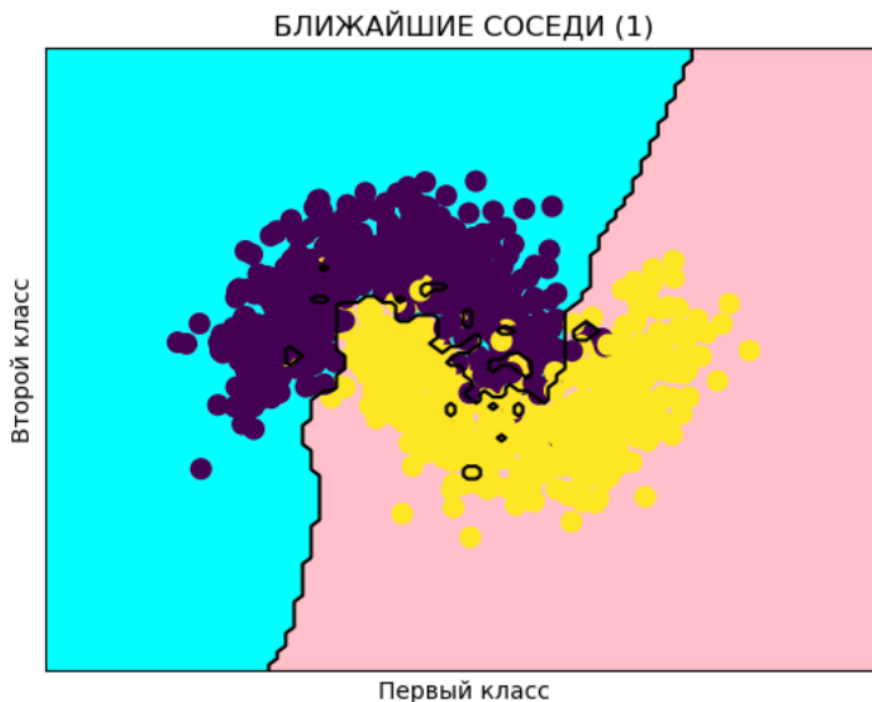


Рисунок 7 – Результат классификации с помощью метода к-ближайших соседей с параметром  $n\_neighbors = 1$

Запустим классификацию с использованием параметра  $n\_neighbors = 3$ , 5 и 9. Результаты классификации представлены на рисунках 8, 9 и 10 соответственно.

Метод классификации: ближайшие соседи (3)

Предсказанные и реальные значения:

```
[1 1 0 1 0 0 0 1 0 0 1 1 0 1 0 1 0 1 0 1 0 1 0 0 1 0 0 1 0 0 1 0 1 1
0 1 0 0 0 1 0 0 0 0 1 0 0 1 1 0 1 1 1 1 0 0 1 0 1 1 0 1 0 1 1 0 0 0 0 1
0 0 0 0 1 0 1 0 0 1 0 1 1 0 0 0 1 0 0 0 1 1 1 0 1 1 0 1 0 0 1 0 0 1 0 0 1
1 1 0 1 0 1 1 0 1 0 0 0 1 1 0 0 1 0 1 0 0 0 1 0 0 1 1 1 0 0 1 1 0 1 0 1 0
1 1 1 1 1 1 0 1 1 1 0 1 1 0 0 0 1 1 1 0 1 1 0 0 1 1 1 0 0 0 1 0 1 1 1 0 1
0 1 0 0 0 1 0 1 0 1 0 0 0 0 1 0 1 1 1 0 1 1 1 1 1 0 0 0 0 1 1 0 0 0 1 1 1
0 1 1 1 1 1 1 1 0 0 1 1 1 0 0 0 1 0 0 1 1 1 1 1 1 1 1 0 1 1]
[1 1 0 1 0 0 0 1 0 0 1 1 0 0 0 1 0 1 0 1 0 1 0 1 0 1 1 0 0 1 0 0 1 0 0 1
0 1 0 0 0 1 0 0 0 1 1 0 0 1 1 0 1 1 1 0 0 0 1 0 1 1 0 1 1 0 0 0 0 0 0 1
1 0 0 0 1 0 1 0 0 1 0 1 0 1 0 0 1 0 0 0 1 1 1 0 1 1 0 1 1 0 1 0 0 1 0 0 1
1 1 0 1 0 1 1 0 1 0 0 0 1 1 0 0 1 0 0 0 1 0 0 0 1 0 0 1 1 1 0 0 1 1 0 1 0
1 1 1 1 0 1 0 1 1 1 0 1 0 0 0 0 1 1 1 0 1 1 0 0 1 1 1 0 0 0 1 0 1 1 0 1
0 1 0 0 0 1 0 1 0 1 1 0 0 0 1 0 1 1 1 0 1 1 1 1 0 0 0 0 1 1 0 0 0 1 1 1
0 1 1 1 0 1 1 1 0 0 1 1 1 0 0 0 1 0 0 0 1 1 1 1 1 1 1 1 0 1 1]
```

Матрица неточностей

```
[[116  8]
 [  5 121]]
```

Точность классификации: 0.948

Полнота:

	precision	recall	f1-score	support
0	0.96	0.94	0.95	124
1	0.94	0.96	0.95	126
accuracy			0.95	250
macro avg	0.95	0.95	0.95	250
weighted avg	0.95	0.95	0.95	250

Площадь под кривой: 0.9479006656426012

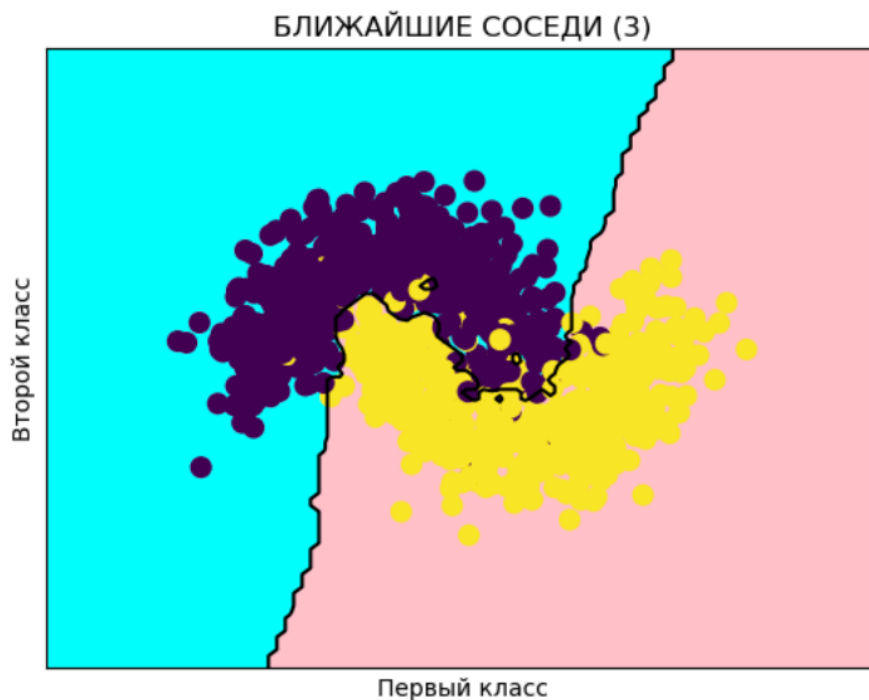


Рисунок 8 – Результат классификации с помощью метода к-ближайших соседей с параметром  $n\_neighbors = 3$

Метод классификации: ближайшие соседи (5)

Предсказанные и реальные значения:

```
[1 1 0 1 0 0 0 1 0 0 1 1 0 1 0 1 0 1 0 1 1 0 1 0 1 1 0 0 1 0 0 1 0 1 1
0 1 0 0 0 1 0 0 0 0 1 0 0 1 1 0 1 1 1 1 0 0 1 0 1 1 0 1 1 0 0 0 0 0 1
0 0 0 0 1 0 1 0 0 1 0 1 0 0 0 0 1 0 0 0 1 1 1 0 1 1 0 1 0 0 1 0 0 1
1 1 0 1 0 1 1 0 1 0 0 0 1 1 0 0 1 0 1 0 0 0 1 0 0 1 1 1 0 0 1 1 0 1 0
1 1 1 1 1 1 0 1 1 1 0 1 1 0 0 0 1 1 1 0 1 1 0 0 1 1 1 0 0 0 1 0 1 1 1 0 1
0 1 0 0 0 1 0 1 0 1 1 0 0 0 1 0 1 1 1 0 1 1 1 1 1 0 0 0 0 1 1 0 0 0 1 1 1
0 1 1 1 1 1 1 1 0 0 1 1 1 0 0 0 1 0 0 1 1 1 1 1 1 1 1 0 1]
[1 1 0 1 0 0 0 1 0 0 1 1 0 0 0 1 0 1 0 1 1 0 1 0 1 1 0 0 1 0 0 1 0 0 1
0 1 0 0 0 1 0 0 0 1 1 0 0 1 1 0 1 1 1 0 0 0 1 0 1 1 0 1 0 1 1 0 0 0 0 1
1 0 0 0 1 0 1 0 0 1 0 1 0 1 0 0 1 0 0 0 1 1 1 0 1 1 0 1 1 0 1 0 0 1 0 0 1
1 1 0 1 0 1 1 0 1 0 0 0 1 1 0 0 1 0 1 0 0 0 1 0 0 1 1 1 0 0 1 1 0 1 0 1 0
1 1 1 1 0 1 0 1 1 1 0 1 0 0 0 0 1 1 1 0 1 1 0 0 1 1 1 0 0 0 1 0 1 1 1 0 1
0 1 0 0 0 1 0 1 0 1 1 0 0 0 1 0 1 1 1 0 1 1 1 1 1 0 0 0 0 1 1 0 0 0 1 1 1
0 1 1 1 0 1 1 1 0 0 1 1 1 0 0 0 1 0 0 0 1 1 1 1 1 1 1 0 1]
```

Матрица неточностей

```
[[117  7]
 [ 4 122]]
```

Точность классификации: 0.956

Полнота:

	precision	recall	f1-score	support
0	0.97	0.94	0.96	124
1	0.95	0.97	0.96	126
accuracy			0.96	250
macro avg	0.96	0.96	0.96	250
weighted avg	0.96	0.96	0.96	250

Площадь под кривой: 0.9559011776753714

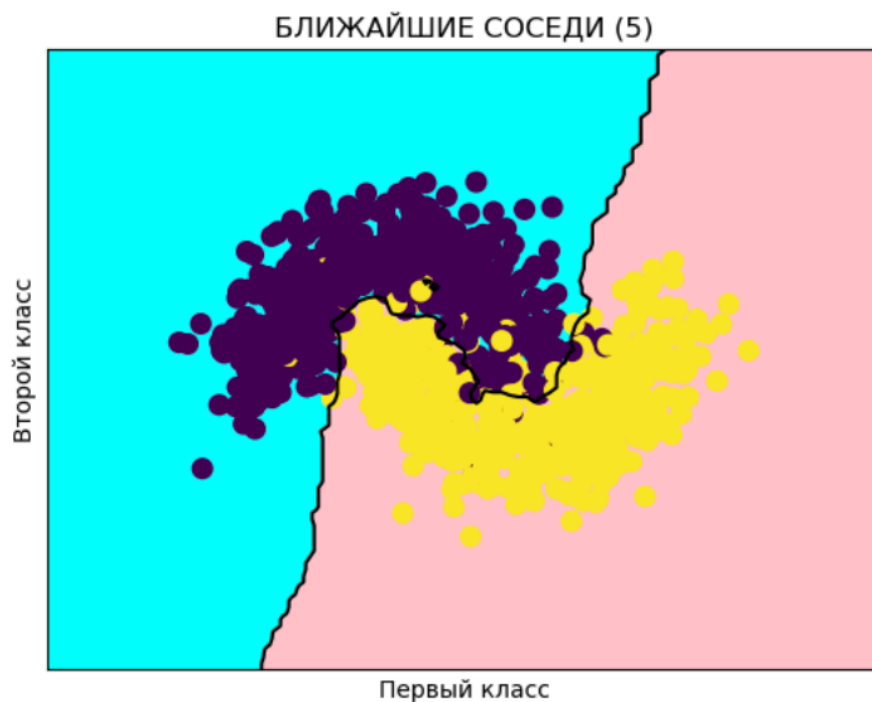


Рисунок 9 – Результат классификации с помощью метода к-ближайших соседей с параметром  $n\_neighbors = 5$

Метод классификации: ближайшие соседи (9)

Предсказанные и реальные значения:

```
[1 1 0 1 0 0 0 1 0 0 1 1 0 1 0 1 0 1 0 1 1 0 1 0 1 1 0 0 1 0 0 1 0 1 1
0 1 0 0 0 1 0 0 0 0 1 0 0 1 1 0 1 1 1 0 0 0 1 0 1 1 0 1 1 0 0 0 0 0 1
0 0 0 0 1 0 1 0 0 1 0 1 0 0 0 0 1 0 0 0 1 1 1 0 1 1 0 1 0 0 1 0 0 1 0 0 1
1 1 0 1 0 1 1 0 1 0 0 0 1 1 0 0 1 0 1 0 0 0 1 0 0 1 1 1 0 0 1 1 0 1 0 1 0
1 1 1 1 1 1 0 1 1 0 1 0 0 0 0 1 1 1 0 1 1 0 0 1 1 1 0 0 0 1 0 1 1 1 0 1
0 1 0 0 0 1 0 1 0 1 0 0 0 0 1 0 1 1 1 0 1 1 1 1 0 0 0 0 1 1 0 0 0 1 1 1
0 1 1 1 0 1 1 1 0 0 1 1 1 0 0 0 1 0 0 1 1 1 1 1 1 1 1 0 1]
[1 1 0 1 0 0 0 1 0 0 1 1 0 0 0 1 0 1 0 1 0 1 1 0 1 0 1 1 0 0 1 0 0 1 0 0 1
0 1 0 0 0 1 0 0 0 1 1 0 0 1 1 0 1 1 1 0 0 0 1 0 1 1 0 1 1 1 0 0 0 0 0 1
1 0 0 0 1 0 1 0 0 1 0 1 0 1 0 0 1 0 0 0 1 1 1 0 1 1 0 1 1 0 1 0 0 1 0 0 1
1 1 0 1 0 1 1 0 1 0 0 0 1 1 0 0 1 0 1 0 0 0 1 0 0 1 1 1 0 0 1 1 0 1 0 1 0
1 1 1 1 0 1 0 1 1 1 0 1 0 0 0 0 1 1 1 0 1 1 0 0 1 1 1 0 0 0 1 0 1 1 1 0 1
0 1 0 0 0 1 0 1 0 1 1 0 0 0 1 0 1 1 1 0 1 1 1 1 0 0 0 0 1 1 0 0 0 1 1 1
0 1 1 1 0 1 1 1 0 0 1 1 1 0 0 0 1 0 0 0 1 1 1 1 1 1 1 0 1]
```

Матрица неточностей

```
[[120  4]
 [ 5 121]]
```

Точность классификации: 0.964

Полнота:

	precision	recall	f1-score	support
0	0.96	0.97	0.96	124
1	0.97	0.96	0.96	126
accuracy			0.96	250
macro avg	0.96	0.96	0.96	250
weighted avg	0.96	0.96	0.96	250

Площадь под кривой: 0.9640296979006657

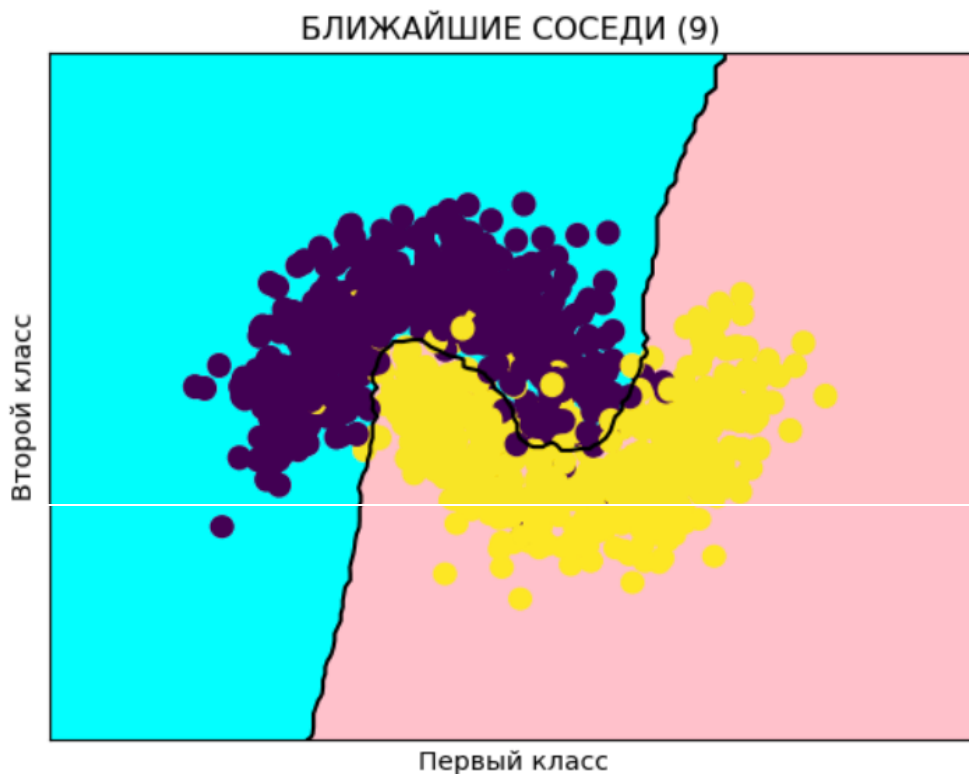


Рисунок 10 – Результат классификации с помощью метода к-ближайших соседей с параметром  $n\_neighbors = 9$

## Классификация с помощью наивного байесовского классификатора

Составленный код для использования данного классификатора представлен на рисунке 11. Полученная информация о точности классификации при использовании данного метода представлена на рисунке 12.

```
: 1 from sklearn.naive_bayes import GaussianNB
  2
  3 nb = GaussianNB()
  4
  5 # Обучаем модель данных
  6 nb.fit(X_train, y_train)
  7
  8 # Оцениваем качество модели
  9 prediction = nb.predict(X_test)
 10
 11 # Выводим сводную информацию
 12 show_info(nb, 'Наивный байесовский классификатор', y_test, prediction)
```

Рисунок 11 – Код для классификатора с помощью наивного байесовского классификатора

Метод классификации: Наивный байесовский классификатор

Предсказанные и реальные значения:

```
[1 1 0 1 0 0 0 1 0 0 1 1 0 1 0 1 0 1 1 1 0 0 1 0 1 0 1 0 0 0 1 0 0 1 0 1 1
0 1 0 0 1 1 0 0 0 0 1 0 0 0 1 0 1 1 0 1 0 0 1 0 0 1 0 1 1 0 0 0 0 0 1
1 0 0 0 1 0 1 0 0 1 0 1 0 1 0 1 0 0 0 1 1 1 0 0 1 0 0 0 0 1 0 0 1 0 1 1
1 0 0 1 0 1 1 0 1 0 1 0 1 0 0 1 0 1 0 1 0 1 0 0 1 1 1 0 0 0 1 1 0 1 0 1 1
1 0 1 0 1 1 0 1 1 0 1 1 0 0 0 0 1 1 1 0 1 1 0 0 1 1 0 0 0 0 1 0 1 1 1 0 1
0 1 0 0 0 1 0 1 0 0 1 0 0 0 1 0 1 0 1 0 1 1 0 1 1 0 0 0 0 1 1 0 0 0 1 0 1
0 1 1 1 0 1 1 1 0 0 0 1 1 0 0 0 1 0 0 1 1 1 1 1 1 1 1 1 0 1]
[1 1 0 1 0 0 0 1 0 0 1 1 0 0 0 1 0 1 0 1 0 1 1 0 1 1 0 1 1 0 0 1 0 0 1 0 0 1
0 1 0 0 0 1 0 0 0 1 1 0 0 1 1 0 1 1 1 0 0 0 1 0 1 1 0 1 1 0 1 1 0 0 0 0 1
1 0 0 0 1 0 1 0 0 1 0 1 0 1 0 0 1 0 0 0 1 1 1 0 1 1 0 1 1 0 1 0 0 1 0 0 1
1 1 0 1 0 1 1 0 1 0 0 0 1 1 0 0 1 0 1 0 0 0 1 0 0 1 1 1 0 0 0 1 1 0 1 0 1 0
1 1 1 1 0 1 0 1 1 1 0 1 0 0 0 0 1 1 1 0 1 1 0 0 1 1 1 0 0 0 1 0 1 1 1 0 1
0 1 0 0 0 1 0 1 0 1 1 0 0 0 1 0 1 1 1 0 1 1 1 1 0 0 0 0 1 1 0 0 0 1 1 1
0 1 1 1 0 1 1 1 0 0 1 1 1 0 0 0 1 0 0 0 1 1 1 1 1 1 1 1 0 1]
```

Матрица неточностей

```
[[112 12]
 [ 19 107]]
```

Точность классификации: 0.876

Полнота:

	precision	recall	f1-score	support
0	0.85	0.90	0.88	124
1	0.90	0.85	0.87	126
accuracy			0.88	250
macro avg	0.88	0.88	0.88	250
weighted avg	0.88	0.88	0.88	250

Площадь под кривой: 0.876216077828981

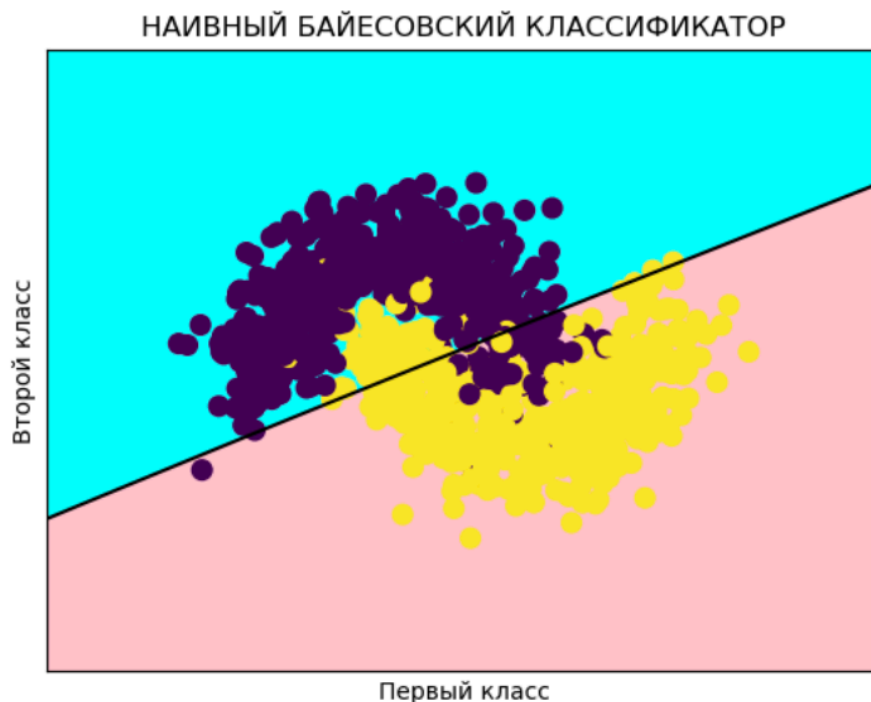


Рисунок 12 – Результат классификации с помощью наивного байесовского классификатора

## Классификация с помощью случайного леса

Использование параметра `n_estimators = 5`. Составленный код для использования данного классификатора с данным параметром представлен на рисунке 13. Полученная информация о точности классификации при использовании данного метода представлена на рисунке 14.

```
1 rfc = RandomForestClassifier(n_estimators=5)
2
3 # Обучаем модель данных
4 rfc.fit(X_train, y_train)
5
6 # Оцениваем качество модели
7 prediction = rfc.predict(X_test)
8
9 # Выводим сводную информацию
10 show_info(rfc, 'случайный лес (5)', y_test, prediction)
```

Рисунок 13 – Код для классификатора с помощью случайного леса с параметром `n_estimators = 5`



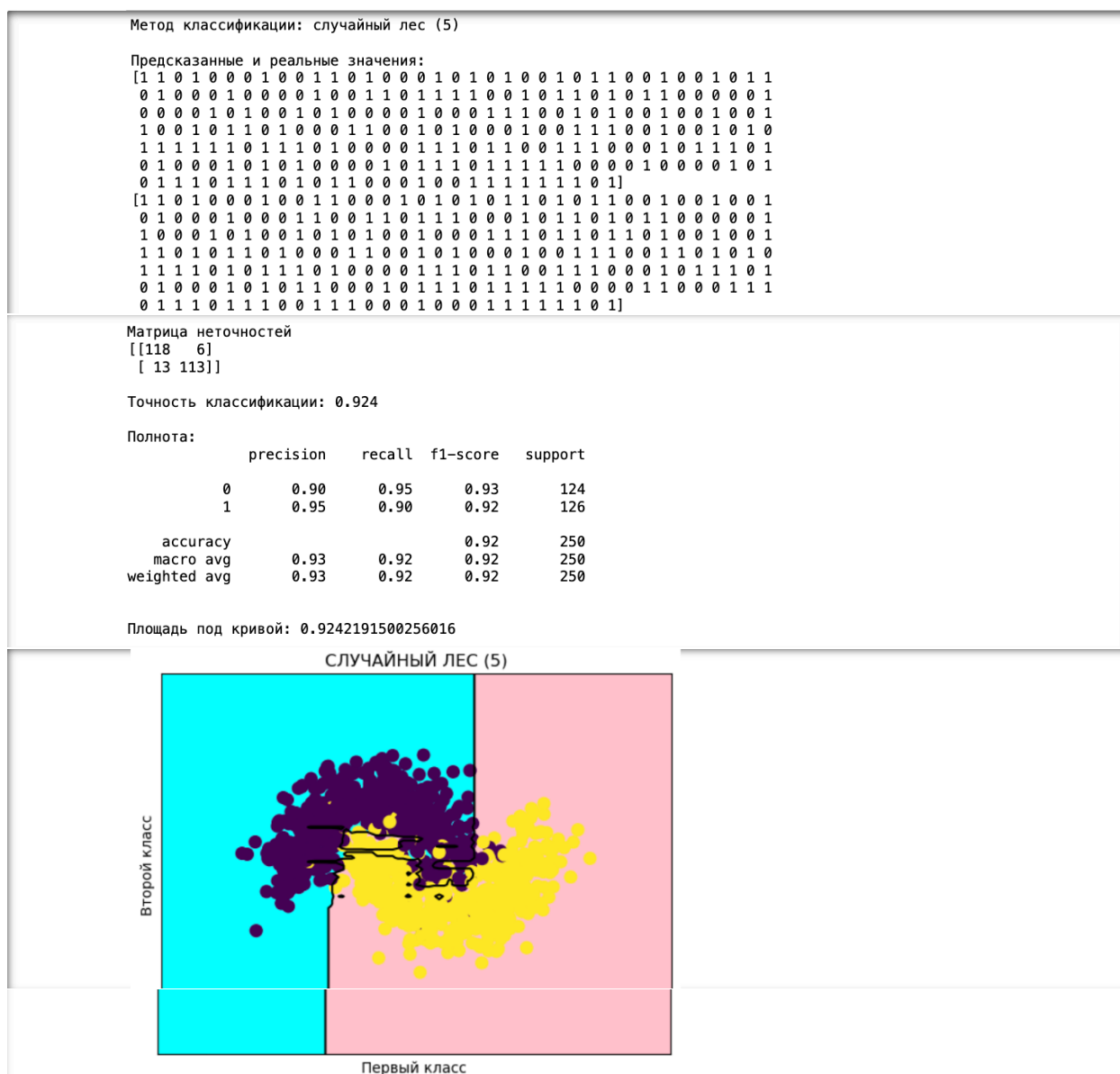


Рисунок 14 – Результат классификации с помощью случайного леса с параметром  $n\_estimators = 5$

Запустим классификацию с использованием параметра  $n\_estimators = 10, 15, 20$  и  $50$ . Результаты классификации представлены на рисунках 15, 16, 17 и 18 соответственно.

Метод классификации: случайный лес (10)

Предсказанные и реальные значения:

```
[1 1 0 1 0 0 0 1 0 0 1 1 0 1 0 1 0 1 0 1 0 0 1 0 1 1 0 0 1 0 0 1 0 1 1
0 1 0 0 0 1 0 0 0 0 1 0 0 1 1 0 1 1 1 1 0 0 1 0 1 0 0 1 0 1 1 0 0 0 0 0 1
0 0 0 0 1 0 1 0 0 1 0 1 0 0 0 0 1 0 0 0 1 1 1 0 1 1 0 1 0 0 1 0 0 1 0 0 1
1 1 0 1 0 1 1 0 1 0 0 0 1 1 0 0 1 0 1 0 0 0 1 0 0 1 1 1 0 0 1 0 0 1 0 1 0
1 1 1 1 1 1 0 1 1 1 0 1 0 0 0 0 1 1 1 0 1 1 0 0 1 1 1 0 0 0 1 0 1 1 1 0 1
0 1 0 0 0 1 0 1 0 1 0 0 0 0 1 0 1 1 1 0 1 1 1 1 1 0 0 0 0 1 1 0 0 0 1 0 1
0 1 1 1 1 1 1 1 0 0 1 1 1 0 0 0 1 0 0 1 1 1 1 1 1 1 1 1 0 1]
[1 1 0 1 0 0 0 1 0 0 1 1 0 0 0 1 0 1 0 1 1 0 1 0 1 1 0 1 0 1 1 0 0 1 0 0 1
0 1 0 0 0 1 0 0 0 1 1 0 0 1 1 0 1 1 1 0 0 0 1 0 1 1 0 1 1 0 1 1 0 0 0 0 1
1 0 0 0 1 0 1 0 0 1 0 1 0 1 0 0 1 0 0 0 1 1 1 0 1 1 0 1 1 0 1 0 0 1 0 0 1
1 1 0 1 0 1 1 0 1 0 0 0 1 1 0 0 1 0 1 0 0 0 1 0 0 1 1 1 0 0 1 1 0 1 0 1 0
1 1 1 1 0 1 0 1 1 1 0 1 0 0 0 0 1 1 1 0 1 1 0 0 1 1 1 0 0 0 1 0 1 1 1 0 1
0 1 0 0 0 1 0 1 0 1 1 0 0 0 1 0 1 1 1 0 1 1 1 1 1 0 0 0 0 1 1 0 0 0 1 1 1
0 1 1 1 0 1 1 1 0 0 1 1 1 0 0 0 1 0 0 0 1 1 1 1 1 1 1 1 0 1]
```

Матрица неточностей

```
[[118  6]
 [ 9 117]]
```

Точность классификации: 0.94

Полнота:

	precision	recall	f1-score	support
0	0.93	0.95	0.94	124
1	0.95	0.93	0.94	126
accuracy			0.94	250
macro avg	0.94	0.94	0.94	250
weighted avg	0.94	0.94	0.94	250

Площадь под кривой: 0.9400921658986175

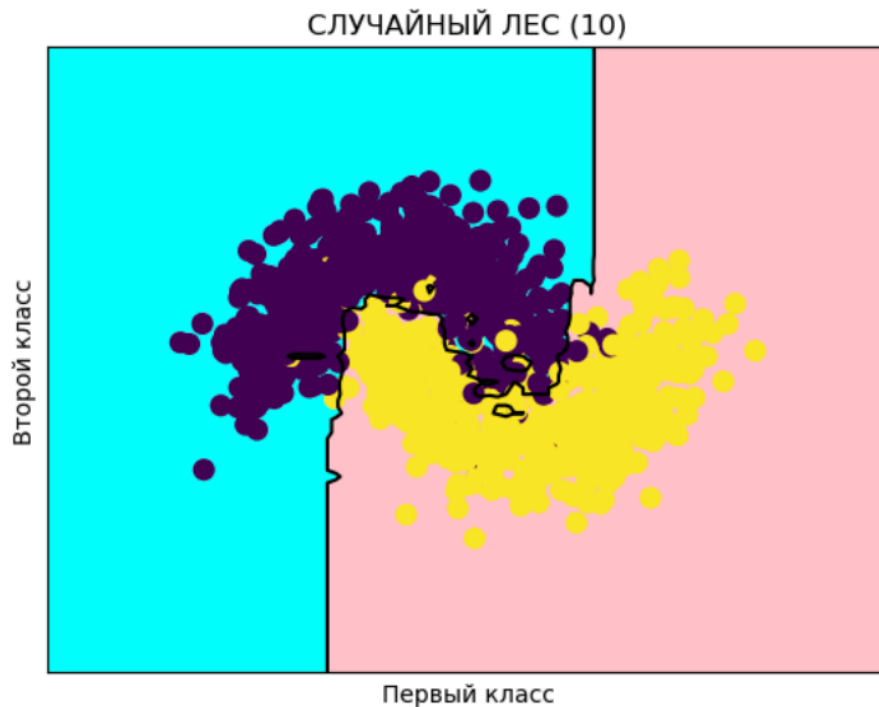


Рисунок 15 – Результат классификации с помощью случайного леса с параметром  $n\_estimators = 10$

Метод классификации: случайный лес (15)

Предсказанные и реальные значения:

```
[1 1 0 1 0 0 0 1 0 0 1 1 0 1 0 1 0 1 0 1 1 0 1 0 1 1 0 0 1 0 0 1 0 1 1
0 1 0 0 0 1 0 0 0 0 1 0 0 1 1 0 1 1 1 1 0 0 1 0 1 1 0 1 1 0 0 0 0 0 1
0 0 0 0 1 0 1 0 0 1 0 1 0 0 0 0 1 0 0 0 1 1 1 0 0 1 0 1 0 0 1 0 0 1
1 0 0 1 0 1 1 0 1 0 0 0 1 1 0 0 1 0 1 0 0 0 1 0 0 1 1 1 0 0 1 1 0 1 0
1 1 1 1 1 1 0 1 1 1 0 1 0 0 0 0 1 1 1 0 1 1 0 0 1 1 1 0 0 0 1 0 1 1 1 0 1
0 1 0 0 0 1 0 1 0 1 0 0 0 0 1 0 1 1 1 0 1 1 1 1 1 0 0 0 0 1 1 0 0 0 1 0 1
0 1 1 1 1 1 1 1 0 0 0 1 1 0 0 0 1 0 0 1 1 1 1 1 1 1 1 0 1]
[1 1 0 1 0 0 0 1 0 0 1 1 0 0 0 1 0 1 0 1 0 1 1 0 1 1 0 0 1 0 0 1 0 0 1
0 1 0 0 0 1 0 0 0 1 1 0 0 1 1 0 1 1 1 0 0 0 1 0 1 1 0 1 1 0 0 0 0 0 1
1 0 0 0 1 0 1 0 0 1 0 1 0 1 0 0 1 0 0 0 1 1 1 0 1 1 0 1 1 0 1 0 0 1 0 0 1
1 1 0 1 0 1 1 0 1 0 0 0 1 1 0 0 1 0 1 0 0 0 1 0 0 1 1 1 0 0 1 1 0 1 0 1 0
1 1 1 1 0 1 0 1 1 1 0 1 0 0 0 0 1 1 1 0 1 1 0 0 1 1 1 0 0 0 1 0 1 1 1 0 1
0 1 0 0 0 1 0 1 0 1 1 0 0 0 1 0 1 1 1 0 1 1 1 1 1 0 0 0 0 1 1 0 0 0 1 1 1
0 1 1 1 0 1 1 1 0 0 1 1 1 0 0 0 1 0 0 0 1 1 1 1 1 1 1 0 1]
```

Матрица неточностей

```
[[118   6]
 [   9 117]]
```

Точность классификации: 0.94

Полнота:

	precision	recall	f1-score	support
0	0.93	0.95	0.94	124
1	0.95	0.93	0.94	126
accuracy			0.94	250
macro avg	0.94	0.94	0.94	250
weighted avg	0.94	0.94	0.94	250

Площадь под кривой: 0.9400921658986175

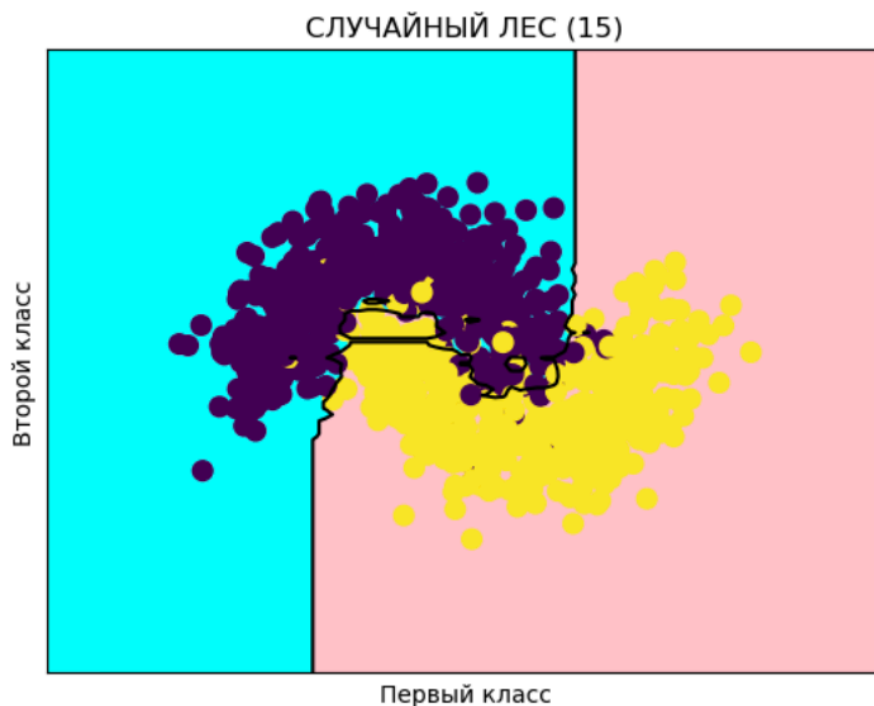


Рисунок 16 – Результат классификации с помощью случайного леса с параметром  $n\_estimators = 15$

Метод классификации: случайный лес (20)

Предсказанные и реальные значения:

```
[1 1 0 1 0 0 0 1 0 0 1 1 0 1 0 1 0 1 0 1 1 0 0 1 0 0 1 0 1 1
 0 1 0 0 0 1 0 0 0 0 1 0 0 1 1 0 1 1 1 0 0 1 0 1 1 0 0 0 0 0 1
 0 0 0 0 1 0 1 0 0 1 0 1 0 0 0 0 1 0 0 0 1 1 1 0 1 1 0 1 0 0 1
 1 1 0 1 0 1 1 0 1 0 0 0 1 1 0 0 1 0 1 0 0 0 1 0 0 1 1 1 0 0 1
 1 1 1 1 1 1 0 1 1 1 0 1 0 0 0 0 1 1 1 0 1 1 0 0 1 1 1 0 0 1 0
 0 1 0 0 0 1 0 1 0 1 0 0 0 0 1 0 1 1 1 0 1 1 1 1 0 0 0 0 1 1 0
 0 1 1 1 1 1 1 1 0 0 0 1 1 1 0 0 0 1 0 0 1 1 1 1 1 1 1 0 1]
[1 1 0 1 0 0 0 1 0 0 1 1 0 0 0 1 0 1 0 1 0 1 1 0 1 0 0 1 0 0 1
 0 1 0 0 0 1 0 0 0 0 1 1 0 0 1 1 0 1 1 1 0 0 0 1 0 1 1 0 0 0 1
 1 0 0 0 1 0 1 0 0 1 0 1 0 1 0 0 1 0 0 0 1 1 1 0 1 1 0 1 1 0 0 1
 1 1 0 1 0 1 1 0 1 0 0 0 1 1 0 0 1 0 1 0 0 0 1 0 0 1 1 1 0 0 1
 1 1 1 1 0 1 0 1 1 1 0 1 0 0 0 0 1 1 1 0 1 1 0 0 1 1 1 0 0 1 0
 0 1 0 0 0 1 0 1 0 1 1 0 0 0 1 0 1 1 1 0 1 1 1 1 0 0 0 0 1 1 0
 0 1 1 1 0 1 1 1 0 0 1 1 1 0 0 0 1 0 0 0 1 1 1 1 1 1 1 0 1]
```

Матрица неточностей

```
[[118   6]
 [   7 119]]
```

Точность классификации: 0.948

Полнота:

	precision	recall	f1-score	support
0	0.94	0.95	0.95	124
1	0.95	0.94	0.95	126
accuracy			0.95	250
macro avg	0.95	0.95	0.95	250
weighted avg	0.95	0.95	0.95	250

Площадь под кривой: 0.9480286738351255

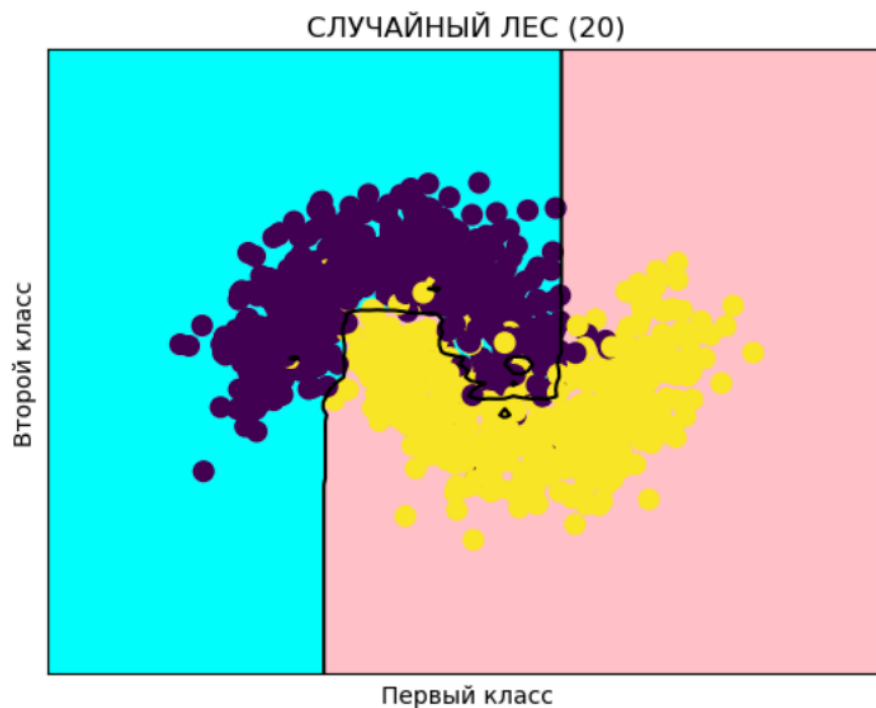


Рисунок 17 – Результат классификации с помощью случайного леса с параметром  $n\_estimators = 20$

Метод классификации: случайный лес (50)

Предсказанные и реальные значения:

```
[1 1 0 1 0 0 0 1 0 0 1 1 0 1 0 1 0 1 0 1 1 0 0 1 0 0 1 0 1 1
0 1 0 0 0 1 0 0 0 0 1 0 0 1 1 0 1 1 1 0 0 1 0 1 1 0 1 0 0 0 0 1
0 0 0 0 1 0 1 0 0 1 0 1 0 0 0 0 1 0 0 0 1 1 1 0 1 1 0 1 0 0 1
1 0 0 1 0 1 1 0 1 0 0 0 1 1 0 0 1 0 1 0 0 0 1 0 0 1 1 1 0 0 1 0
1 1 1 1 1 1 0 1 1 1 0 1 0 0 0 0 1 1 1 0 1 1 0 0 1 1 0 0 0 1 0 1
0 1 0 0 0 1 0 1 0 1 0 0 0 0 1 0 1 1 1 0 1 1 1 1 0 0 0 0 1 1 0 1
0 1 1 1 1 1 1 1 0 0 1 1 1 0 0 0 1 0 0 1 1 1 1 1 1 1 0 1]
[1 1 0 1 0 0 0 1 0 0 1 1 0 0 0 1 0 1 0 1 0 1 1 0 1 0 0 1 0 0 1
0 1 0 0 0 1 0 0 0 1 1 0 0 1 1 0 1 1 1 0 0 0 1 0 1 1 0 1 0 0 0 1
1 0 0 0 1 0 1 0 0 1 0 1 0 1 0 0 1 0 0 0 1 1 1 0 1 1 0 1 0 0 1 0
1 1 0 1 0 1 1 0 1 0 0 0 1 1 0 0 1 0 1 0 0 0 1 0 0 1 1 1 0 0 1 0
1 1 1 1 0 1 0 1 1 1 0 1 0 0 0 0 1 1 1 0 1 1 0 0 1 1 0 0 1 0 1 1
0 1 0 0 0 1 0 1 0 1 1 0 0 0 1 0 1 1 1 0 1 1 1 1 0 0 0 0 1 1 0 1
0 1 1 1 0 1 1 1 0 0 1 1 1 0 0 0 1 0 0 0 1 1 1 1 1 1 0 0 1 1 1]
```

Матрица неточностей

```
[[118   6]
 [   8 118]]
```

Точность классификации: 0.944

Полнота:

	precision	recall	f1-score	support
0	0.94	0.95	0.94	124
1	0.95	0.94	0.94	126
accuracy			0.94	250
macro avg	0.94	0.94	0.94	250
weighted avg	0.94	0.94	0.94	250

Площадь под кривой: 0.9440604198668715

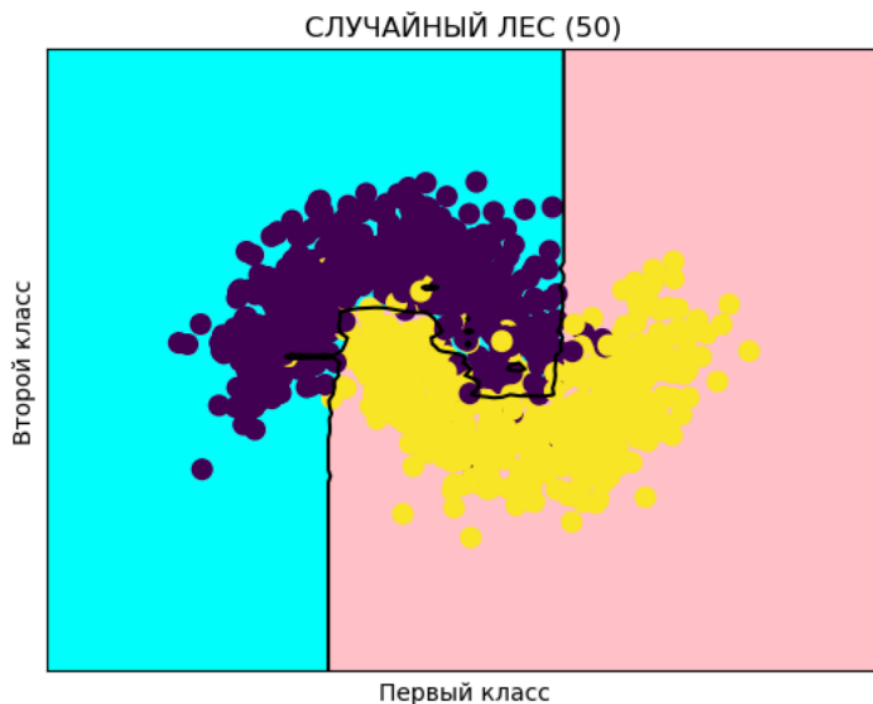


Рисунок 18 – Результат классификации с помощью случайного леса с параметром  $n\_estimators = 50$

## Анализ результатов

Сведем полученные данные в таблицу 1 и сделаем вывод.

Таблица 1 – Результат классификации по методам про 75% обучающей выборки

Метод (параметры)	Точность	Площадь под кривой
Метод к-ближайших соседей (n_neighbors = 1)	0,936	0,936
Метод к-ближайших соседей (n_neighbors = 3)	0,948	0,948
Метод к-ближайших соседей (n_neighbors = 5)	0,956	0,956
Метод к-ближайших соседей (n_neighbors = 9)	0,964	0,964
Наивный байесовский классификатор	0,876	0,876
Случайный лес (n_estimators = 5)	0,924	0,924
Случайный лес (n_estimators = 10)	0,936	0,936
Случайный лес (n_estimators = 15)	0,944	0,944
Случайный лес (n_estimators = 20)	0,948	0,948
Случайный лес (n_estimators = 50)	0,944	0,944

Исходя из таблицы 1, можно сделать вывод о том, что лучше всего себя показал метод к-ближайших соседей (n\_neighbors = 9), хуже всего – наивный байесовский классификатор.

Рассмотрим случай уменьшения тестовой выборки. Установим, что тестовая выборка составляет 10% и построим графики визуализации обучающей и тестовой выборки, данные графики представлены на рисунках 19 и 20 соответственно.

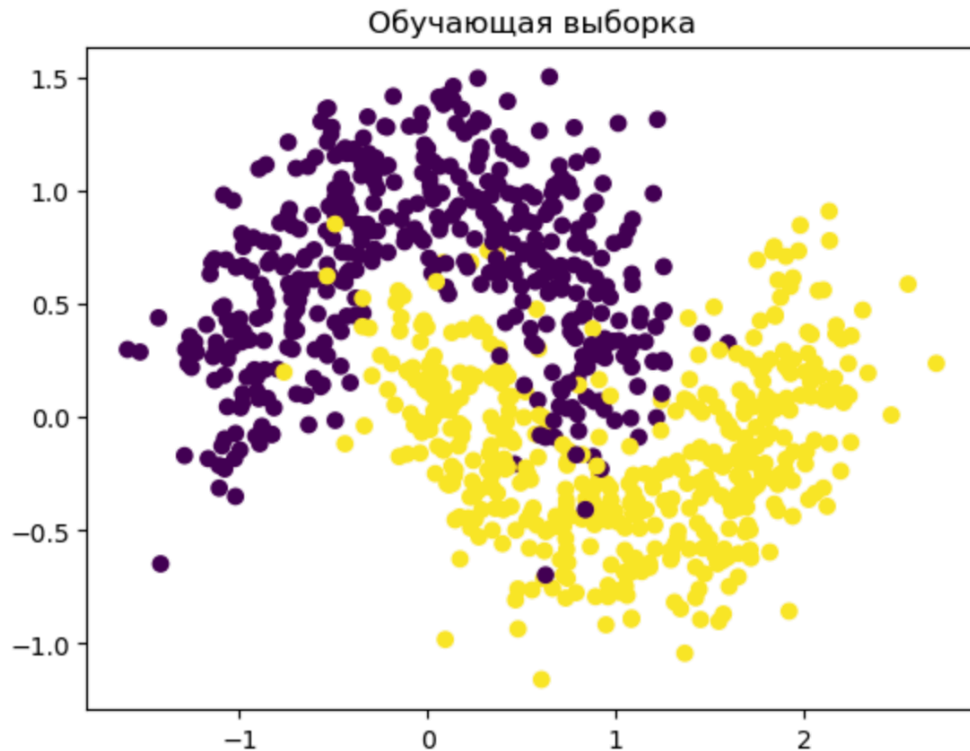


Рисунок 19 – Обучающая выборка при 90% от общего размера

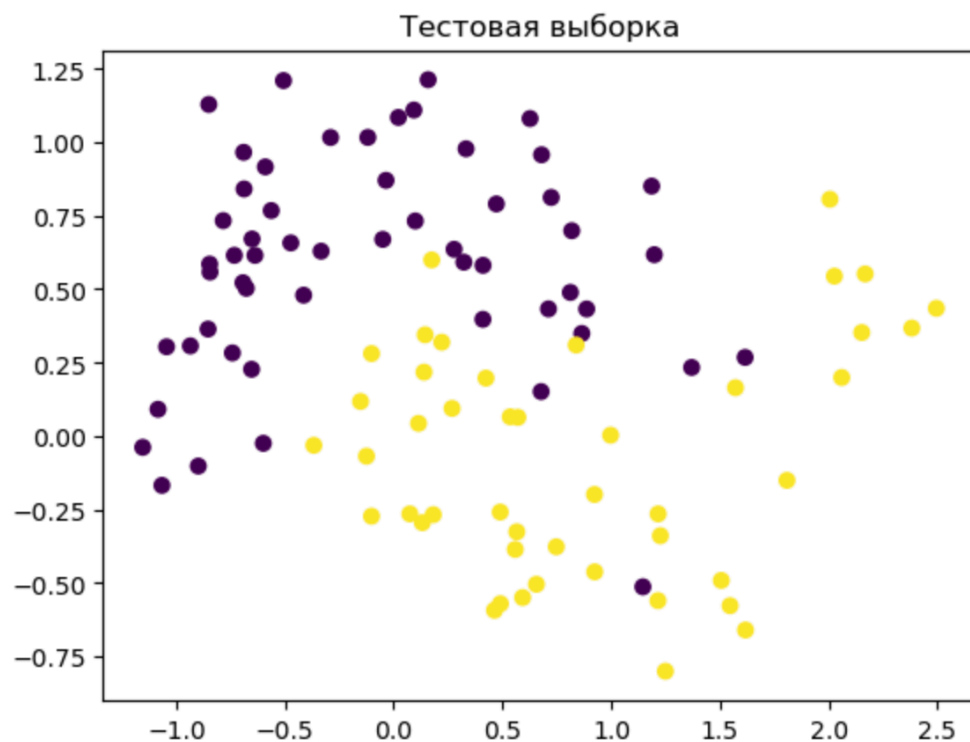


Рисунок 20 – Тестовая выборка при 10% от общего размера

Сведем полученные данные в таблицу 2 и сделаем вывод.

Таблица 2 – Результат классификации по методам про 90% обучающей выборки

Метод (параметры)	Точность	Площадь под кривой
Метод к-ближайших соседей (n_neighbors = 1)	0,93	0,928
Метод к-ближайших соседей (n_neighbors = 3)	0,93	0,928
Метод к-ближайших соседей (n_neighbors = 5)	0,93	0,93
Метод к-ближайших соседей (n_neighbors = 9)	0,94	0,939
Наивный байесовский классификатор	0,88	0,878
Случайный лес (n_estimators = 5)	0,91	0,906
Случайный лес (n_estimators = 10)	0,91	0,906
Случайный лес (n_estimators = 15)	0,92	0,917
Случайный лес (n_estimators = 20)	0,94	0,939
Случайный лес (n_estimators = 50)	0,94	0,939

Исходя из таблицы 2, можно сделать вывод о том, что лучше всего себя показал метод к-ближайших соседей (n\_neighbors = 9) и случайный лес (n\_estimators = 20 и n\_estimators = 50), хуже всего – наивный байесовский классификатор.



Рассмотрим случай уменьшения тестовой выборки. Установим, что тестовая выборка составляет 35% и построим графики визуализации обучающей и тестовой выборки, данные графики представлены на рисунках 21 и 22 соответственно.

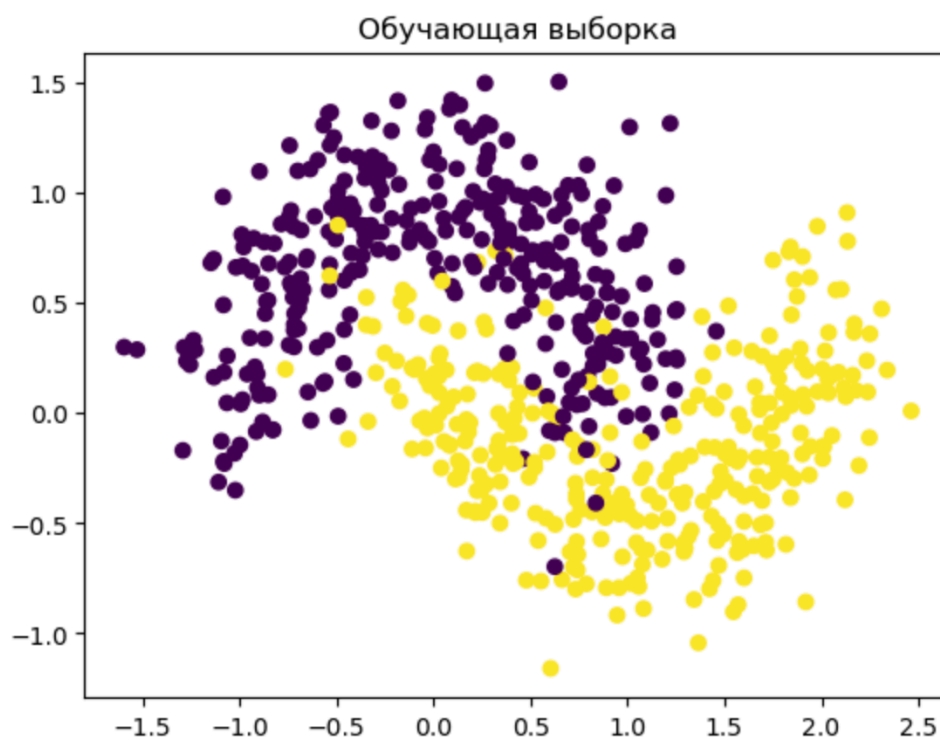


Рисунок 21 – Обучающая выборка при 65% от общего размера

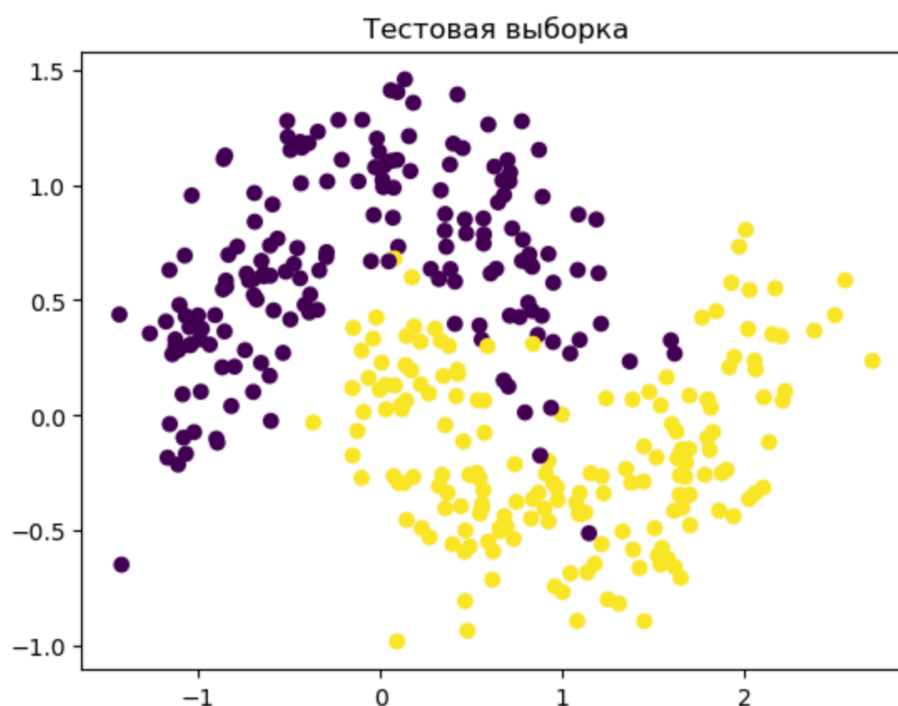


Рисунок 22 – Тестовая выборка при 35% от общего размера

Сведем полученные данные в таблицу 3 и сделаем вывод.

Таблица 3 – Результат классификации по методам про 65% обучающей выборки

Метод (параметры)	Точность	Площадь под кривой
Метод к-ближайших соседей (n_neighbors = 1)	0,934	0,934
Метод к-ближайших соседей (n_neighbors = 3)	0,954	0,954
Метод к-ближайших соседей (n_neighbors = 5)	0,963	0,963
Метод к-ближайших соседей (n_neighbors = 9)	0,971	0,971
Наивный байесовский классификатор	0,886	0,889
Случайный лес (n_estimators = 5)	0,949	0,949
Случайный лес (n_estimators = 10)	0,929	0,929
Случайный лес (n_estimators = 15)	0,943	0,943
Случайный лес (n_estimators = 20)	0,951	0,951
Случайный лес (n_estimators = 50)	0,946	0,946

Исходя из таблицы 3, можно сделать вывод о том, что лучше всего себя показал метод к-ближайших соседей (n\_neighbors = 9), хуже всего – наивный байесовский классификатор.

## Вывод

В ходе выполнения данной лабораторной работы мною были получены практические навыки решения задачи бинарной классификации данных в среде Jupiter Notebook. Также я научился загружать данные, обучать классификаторы и проводить классификацию и получил опыт в оценивании точности полученных моделей.

В результате анализа полученных результатов выяснилось, что наивысшая точность классификации 0,971 (и наибольшая площадь под кривой) достигается, когда размер обучающей выборки равен 65% и в качестве метода выбран метод к-ближних соседей при  $n\_neighbors = 9$ .

## Приложение А

### Исходный код при 25% тестовой выборки

```
#!/usr/bin/env python
# coding: utf-8

# ## Лабораторная работа №1
# ### Задание
# ### Вариант №7
# ### Вид классов: `moons`
# ### Random state: `77`
# ### noise: `0.25`

# In[25]:

# Модуль numpy (сокращение от "Numerical Python") предоставляет
# функциональность для эффективной работы
# с массивами и математическими операциями на ними.
import numpy as np

# Модуль matplotlib.pyplot используется для создания графиков и визуализации
# данных.
# Он предоставляет множество функций для построения различных типов графиков.
import matplotlib.pyplot as plt

from sklearn.datasets import make_moons
from sklearn.model_selection import train_test_split

# In[26]:

# А для отображения на графике области принятия решения - готовую функцию
plot_2d_separator,
# которой нужно передать на вход объект classifier - модель классификатора и
# X - массив входных данных:
def plot_2d_separator(classifier, X, fill=False, line=True, ax=None,
eps=None):
    if eps is None:
        eps = 1.0
    x_min, x_max = X[:, 0].min() - eps, X[:, 0].max() + eps
    y_min, y_max = X[:, 1].min() - eps, X[:, 1].max() + eps
    xx = np.linspace(x_min, x_max, 100)
    yy = np.linspace(y_min, y_max, 100)
    x1, x2 = np.meshgrid(xx, yy)
    X_grid = np.c_[x1.ravel(), x2.ravel()]
    try:
        decision_values = classifier.decision_function(X_grid)
        levels = [0]
        fill_levels = [decision_values.min(), 0, decision_values.max()]
    except AttributeError:
        decision_values = classifier.predict_proba(X_grid)[:, 1]
        levels = [.5]
        fill_levels = [0, .5, 1]

    if ax is None:
        ax = plt.gca()
    if fill:
        ax.contourf(x1,
                    x2,
                    decision_values.reshape(x1.shape),
```

```

        levels=fill_levels,
        colors=['cyan', 'pink', 'yellow'])
    if line:
        ax.contour(x1,
                   x2,
                   decision_values.reshape(x1.shape),
                   levels=levels,
                   colors='black')

    ax.set_xlim(x_min, x_max)
    ax.set_ylim(y_min, y_max)
    ax.set_xticks(())
    ax.set_yticks(())

# ### Генерация выборки

# In[27]:

X, y = make_moons(n_samples=1000, shuffle=True, noise=0.25, random_state=77)

# In[28]:

print('Координаты точек: ')
print(X[:15])
print('Метки класса: ')
print(y[:15])

# In[29]:

plt.scatter(X[:, 0], X[:, 1], c=y)
plt.show()

# ### Разбитие выборки на обучающее и тестовое множество (75/25)

# In[30]:

X_train, X_test, y_train, y_test = train_test_split(X,
                                                    y,
                                                    test_size=0.25,
                                                    random_state=77)

# In[31]:

plt.title('Обучающая выборка')
plt.scatter(X_train[:, 0], X_train[:, 1], c=y_train)
plt.show()

# In[32]:

plt.title('Тестовая выборка')
plt.scatter(X_test[:, 0], X_test[:, 1], c=y_test)

```

```

plt.show()

# ### Классификация

# In[33]:

from sklearn.metrics import confusion_matrix
from sklearn.metrics import accuracy_score
from sklearn.metrics import classification_report
from sklearn.metrics import roc_auc_score

# In[34]:

def show_info(classifier, classifier_name, real_values, prediction_values):
    print(f'Метод классификации: {classifier_name}\n')

    # Выводим предсказанное и реальное значение
    print('Предсказанные и реальные значения:')
    print(prediction_values)
    print(real_values)

    # Выводим матрицу неточностей
    print('\nМатрица неточностей')
    print(confusion_matrix(real_values, prediction_values))

    # Выводим точность классификации
    print(f'\nТочность классификации: {accuracy_score(prediction_values,
real_values)}')

    # Выводим полноту
    print('\nПолнота: ')
    print(classification_report(real_values, prediction_values))

    # AUC ROC
    print(f'\nПлощадь под кривой: {roc_auc_score(real_values,
prediction_values)}')

    plt.xlabel('Первый класс')
    plt.ylabel('Второй класс')
    plt.title(classifier_name.upper())
    plot_2d_separator(classifier, X, fill=True)
    plt.scatter(X[:, 0], X[:, 1], c=y, s=70)

# ### Метод k-ближайших соседей (1)

# In[35]:

from sklearn.neighbors import KNeighborsClassifier

# In[36]:

knn = KNeighborsClassifier(n_neighbors=1, metric='euclidean')

# Обучаем модель данных
knn.fit(X_train, y_train)

```

```

# Оцениваем качество модели
prediction = knn.predict(X_test)

# Выводим сводную информацию
show_info(knn, 'ближайшие соседи (1)', y_test, prediction)

# ### Метод k-ближайших соседей (3)

# In[37]:

knn = KNeighborsClassifier(n_neighbors=3, metric='euclidean')

# Обучаем модель данных
knn.fit(X_train, y_train)

# Оцениваем качество модели
prediction = knn.predict(X_test)

# Выводим сводную информацию
show_info(knn, 'ближайшие соседи (3)', y_test, prediction)

# ### Метод k-ближайших соседей (5)

# In[38]:

knn = KNeighborsClassifier(n_neighbors=5, metric='euclidean')

# Обучаем модель данных
knn.fit(X_train, y_train)

# Оцениваем качество модели
prediction = knn.predict(X_test)

# Выводим сводную информацию
show_info(knn, 'ближайшие соседи (5)', y_test, prediction)

# ### Метод k-ближайших соседей (9)

# In[39]:

knn = KNeighborsClassifier(n_neighbors=9, metric='euclidean')

# Обучаем модель данных
knn.fit(X_train, y_train)

# Оцениваем качество модели
prediction = knn.predict(X_test)

# In[40]:

# Выводим сводную информацию
show_info(knn, 'ближайшие соседи (9)', y_test, prediction)

```

```

# ## Наивный байесовский классификатор

# In[41]:

from sklearn.naive_bayes import GaussianNB

nb = GaussianNB()

# Обучаем модель данных
nb.fit(X_train, y_train)

# Оцениваем качество модели
prediction = nb.predict(X_test)

# Выводим сводную информацию
show_info(nb, 'Наивный байесовский классификатор', y_test, prediction)

# ### Случайный лес (5)

# In[42]:

from sklearn.ensemble import RandomForestClassifier

# In[43]:

rfc = RandomForestClassifier(n_estimators=5)

# Обучаем модель данных
rfc.fit(X_train, y_train)

# Оцениваем качество модели
prediction = rfc.predict(X_test)

# Выводим сводную информацию
show_info(rfc, 'случайный лес (5)', y_test, prediction)

# ### Случайный лес (10)

# In[44]:

rfc = RandomForestClassifier(n_estimators=10)

# Обучаем модель данных
rfc.fit(X_train, y_train)

# Оцениваем качество модели
prediction = rfc.predict(X_test)

# Выводим сводную информацию
show_info(rfc, 'случайный лес (10)', y_test, prediction)

# ### Случайный лес (15)

# In[45]:

```



```

rfc = RandomForestClassifier(n_estimators=15)

# Обучаем модель данных
rfc.fit(X_train, y_train)

# Оцениваем качество модели
prediction = rfc.predict(X_test)

# Выводим сводную информацию
show_info(rfc, 'случайный лес (15)', y_test, prediction)

# ### Случайный лес (20)

# In[46]:

rfc = RandomForestClassifier(n_estimators=20)

# Обучаем модель данных
rfc.fit(X_train, y_train)

# Оцениваем качество модели
prediction = rfc.predict(X_test)

# Выводим сводную информацию
show_info(rfc, 'случайный лес (20)', y_test, prediction)

# ### Случайный лес (50)

# In[47]:

rfc = RandomForestClassifier(n_estimators=50)

# Обучаем модель данных
rfc.fit(X_train, y_train)

# Оцениваем качество модели
prediction = rfc.predict(X_test)

# Выводим сводную информацию
show_info(rfc, 'случайный лес (50)', y_test, prediction)

```

## Приложение Б

### Исходный код при 10% тестовой выборки

```
#!/usr/bin/env python
# coding: utf-8

# ## Лабораторная работа №1
# ### Задание
# ### Вариант №7
# ### Вид классов: `moons`
# ### Random state: `77`
# ### noise: `0.25`

# In[23]:

# Модуль numpy (сокращение от "Numerical Python") предоставляет
# функциональность для эффективной работы
# с массивами и математическими операциями на ними.
import numpy as np

# Модуль matplotlib.pyplot используется для создания графиков и визуализации
# данных.
# Он предоставляет множество функций для построения различных типов графиков.
import matplotlib.pyplot as plt

from sklearn.datasets import make_moons
from sklearn.model_selection import train_test_split

# In[24]:

# А для отображения на графике области принятия решения - готовую функцию
plot_2d_separator,
# которой нужно передать на вход объект classifier - модель классификатора и
# X - массив входных данных:
def plot_2d_separator(classifier, X, fill=False, line=True, ax=None,
eps=None):
    if eps is None:
        eps = 1.0
    x_min, x_max = X[:, 0].min() - eps, X[:, 0].max() + eps
    y_min, y_max = X[:, 1].min() - eps, X[:, 1].max() + eps
    xx = np.linspace(x_min, x_max, 100)
    yy = np.linspace(y_min, y_max, 100)
    x1, x2 = np.meshgrid(xx, yy)
    X_grid = np.c_[x1.ravel(), x2.ravel()]
    try:
        decision_values = classifier.decision_function(X_grid)
        levels = [0]
        fill_levels = [decision_values.min(), 0, decision_values.max()]
    except AttributeError:
        decision_values = classifier.predict_proba(X_grid)[:, 1]
        levels = [.5]
        fill_levels = [0, .5, 1]

    if ax is None:
        ax = plt.gca()
    if fill:
        ax.contourf(x1,
                    x2,
                    decision_values.reshape(x1.shape),
```

```

        levels=fill_levels,
        colors=['cyan', 'pink', 'yellow'])
    if line:
        ax.contour(x1,
                   x2,
                   decision_values.reshape(x1.shape),
                   levels=levels,
                   colors='black')

    ax.set_xlim(x_min, x_max)
    ax.set_ylim(y_min, y_max)
    ax.set_xticks(())
    ax.set_yticks(())

# ### Генерация выборки

# In[25]:

X, y = make_moons(n_samples=1000, shuffle=True, noise=0.25, random_state=77)

# In[26]:

print('Координаты точек: ')
print(X[:15])
print('Метки класса: ')
print(y[:15])

# In[27]:

plt.scatter(X[:, 0], X[:, 1], c=y)
plt.show()

# ### Разбитие выборки на обучающее и тестовое множество (90/10)

# In[28]:

X_train, X_test, y_train, y_test = train_test_split(X,
                                                    y,
                                                    test_size=0.10,
                                                    random_state=77)

# In[29]:

plt.title('Обучающая выборка')
plt.scatter(X_train[:, 0], X_train[:, 1], c=y_train)
plt.show()

# In[30]:

plt.title('Тестовая выборка')
plt.scatter(X_test[:, 0], X_test[:, 1], c=y_test)

```

```

plt.show()

# ### Классификация

# In[31]:

from sklearn.metrics import confusion_matrix
from sklearn.metrics import accuracy_score
from sklearn.metrics import classification_report
from sklearn.metrics import roc_auc_score

# In[32]:

def show_info(classifier, classifier_name, real_values, prediction_values):
    print(f'Метод классификации: {classifier_name}\n')

    # Выводим предсказанное и реальное значение
    print('Предсказанные и реальные значения:')
    print(prediction_values)
    print(real_values)

    # Выводим матрицу неточностей
    print('\nМатрица неточностей')
    print(confusion_matrix(real_values, prediction_values))

    # Выводим точность классификации
    print(f'\nТочность классификации: {accuracy_score(prediction_values,
real_values)}')

    # Выводим полноту
    print('\nПолнота: ')
    print(classification_report(real_values, prediction_values))

    # AUC ROC
    print(f'\nПлощадь под кривой: {roc_auc_score(real_values,
prediction_values)}')

    plt.xlabel('Первый класс')
    plt.ylabel('Второй класс')
    plt.title(classifier_name.upper())
    plot_2d_separator(classifier, X, fill=True)
    plt.scatter(X[:, 0], X[:, 1], c=y, s=70)

# ### Метод k-ближайших соседей (1)

# In[33]:

from sklearn.neighbors import KNeighborsClassifier

# In[34]:

knn = KNeighborsClassifier(n_neighbors=1, metric='euclidean')

# Обучаем модель данных
knn.fit(X_train, y_train)

```

```

# Оцениваем качество модели
prediction = knn.predict(X_test)

# Выводим сводную информацию
show_info(knn, 'ближайшие соседи (1)', y_test, prediction)

# ### Метод k-ближайших соседей (3)

# In[35]:

knn = KNeighborsClassifier(n_neighbors=3, metric='euclidean')

# Обучаем модель данных
knn.fit(X_train, y_train)

# Оцениваем качество модели
prediction = knn.predict(X_test)

# Выводим сводную информацию
show_info(knn, 'ближайшие соседи (3)', y_test, prediction)

# ### Метод k-ближайших соседей (5)

# In[36]:

knn = KNeighborsClassifier(n_neighbors=5, metric='euclidean')

# Обучаем модель данных
knn.fit(X_train, y_train)

# Оцениваем качество модели
prediction = knn.predict(X_test)

# Выводим сводную информацию
show_info(knn, 'ближайшие соседи (5)', y_test, prediction)

# ### Метод k-ближайших соседей (9)

# In[37]:

knn = KNeighborsClassifier(n_neighbors=9, metric='euclidean')

# Обучаем модель данных
knn.fit(X_train, y_train)

# Оцениваем качество модели
prediction = knn.predict(X_test)

# Выводим сводную информацию
show_info(knn, 'ближайшие соседи (9)', y_test, prediction)

# ## Наивный байесовский классификатор

# In[38]:

```

```

from sklearn.naive_bayes import GaussianNB

nb = GaussianNB()

# Обучаем модель данных
nb.fit(X_train, y_train)

# Оцениваем качество модели
prediction = nb.predict(X_test)

# Выводим сводную информацию
show_info(nb, 'Наивный байесовский классификатор', y_test, prediction)

# ### Случайный лес (5)

# In[39]:

from sklearn.ensemble import RandomForestClassifier

# In[40]:

rfc = RandomForestClassifier(n_estimators=5)

# Обучаем модель данных
rfc.fit(X_train, y_train)

# Оцениваем качество модели
prediction = rfc.predict(X_test)

# Выводим сводную информацию
show_info(rfc, 'случайный лес (5)', y_test, prediction)

# ### Случайный лес (10)

# In[41]:

rfc = RandomForestClassifier(n_estimators=10)

# Обучаем модель данных
rfc.fit(X_train, y_train)

# Оцениваем качество модели
prediction = rfc.predict(X_test)

# Выводим сводную информацию
show_info(rfc, 'случайный лес (10)', y_test, prediction)

# ### Случайный лес (15)

# In[42]:

rfc = RandomForestClassifier(n_estimators=15)

# Обучаем модель данных

```

```

rfc.fit(X_train, y_train)

# Оцениваем качество модели
prediction = rfc.predict(X_test)

# Выводим сводную информацию
show_info(rfc, 'случайный лес (15)', y_test, prediction)

# ### Случайный лес (20)

# In[43]:

rfc = RandomForestClassifier(n_estimators=20)

# Обучаем модель данных
rfc.fit(X_train, y_train)

# Оцениваем качество модели
prediction = rfc.predict(X_test)

# Выводим сводную информацию
show_info(rfc, 'случайный лес (20)', y_test, prediction)

# ### Случайный лес (50)

# In[44]:

rfc = RandomForestClassifier(n_estimators=50)

# Обучаем модель данных
rfc.fit(X_train, y_train)

# Оцениваем качество модели
prediction = rfc.predict(X_test)

# Выводим сводную информацию
show_info(rfc, 'случайный лес (50)', y_test, prediction)

```

## Приложение В

### Исходный код при 35% тестовой выборки

```
#!/usr/bin/env python
# coding: utf-8

# ## Лабораторная работа №1
# ### Задание
# ### Вариант №7
# ### Вид классов: `moons`
# ### Random state: `77`
# ### noise: `0.25`

# In[1]:

# Модуль numpy (сокращение от "Numerical Python") предоставляет
# функциональность для эффективной работы
# с массивами и математическими операциями на них.
import numpy as np

# Модуль matplotlib.pyplot используется для создания графиков и визуализации
# данных.
# Он предоставляет множество функций для построения различных типов графиков.
import matplotlib.pyplot as plt

from sklearn.datasets import make_moons
from sklearn.model_selection import train_test_split

# In[2]:

# А для отображения на графике области принятия решения - готовую функцию
plot_2d_separator,
# которой нужно передать на вход объект classifier - модель классификатора и
# X - массив входных данных:
def plot_2d_separator(classifier, X, fill=False, line=True, ax=None,
eps=None):
    if eps is None:
        eps = 1.0
    x_min, x_max = X[:, 0].min() - eps, X[:, 0].max() + eps
    y_min, y_max = X[:, 1].min() - eps, X[:, 1].max() + eps
    xx = np.linspace(x_min, x_max, 100)
    yy = np.linspace(y_min, y_max, 100)
    x1, x2 = np.meshgrid(xx, yy)
    X_grid = np.c_[x1.ravel(), x2.ravel()]
    try:
        decision_values = classifier.decision_function(X_grid)
        levels = [0]
        fill_levels = [decision_values.min(), 0, decision_values.max()]
    except AttributeError:
        decision_values = classifier.predict_proba(X_grid)[:, 1]
        levels = [.5]
        fill_levels = [0, .5, 1]

    if ax is None:
        ax = plt.gca()
    if fill:
        ax.contourf(x1,
```



```

        x2,
        decision_values.reshape(x1.shape),
        levels=fill_levels,
        colors=['cyan', 'pink', 'yellow'])
    if line:
        ax.contour(x1,
                   x2,
                   decision_values.reshape(x1.shape),
                   levels=levels,
                   colors='black')

    ax.set_xlim(x_min, x_max)
    ax.set_ylim(y_min, y_max)
    ax.set_xticks(())
    ax.set_yticks(())

# ### Генерация выборки

# In[3]:

X, y = make_moons(n_samples=1000, shuffle=True, noise=0.25, random_state=77)

# In[4]:

print('Координаты точек: ')
print(X[:15])
print('Метки класса: ')
print(y[:15])

# In[5]:

plt.scatter(X[:, 0], X[:, 1], c=y)
plt.show()

# ### Разбитие выборки на обучающее и тестовое множество (65/35)

# In[6]:

X_train, X_test, y_train, y_test = train_test_split(X,
                                                    y,
                                                    test_size=0.35,
                                                    random_state=77)

# In[7]:

plt.title('Обучающая выборка')
plt.scatter(X_train[:, 0], X_train[:, 1], c=y_train)
plt.show()

# In[8]:

```

```

plt.title('Тестовая выборка')
plt.scatter(X_test[:, 0], X_test[:, 1], c=y_test)
plt.show()

# ### Классификация

# In[9]:

from sklearn.metrics import confusion_matrix
from sklearn.metrics import accuracy_score
from sklearn.metrics import classification_report
from sklearn.metrics import roc_auc_score

# In[10]:

def show_info(classifier, classifier_name, real_values, prediction_values):
    print(f'Метод классификации: {classifier_name}\n')

    # Выводим предсказанное и реальное значение
    print('Предсказанные и реальные значения:')
    print(prediction_values)
    print(real_values)

    # Выводим матрицу неточностей
    print('\nМатрица неточностей')
    print(confusion_matrix(real_values, prediction_values))

    # Выводим точность классификации
    print(f'\nТочность классификации: {accuracy_score(prediction_values,
real_values)}')

    # Выводим полноту
    print('\nПолнота: ')
    print(classification_report(real_values, prediction_values))

    # AUC ROC
    print(f'\nПлощадь под кривой: {roc_auc_score(real_values,
prediction_values)}')

    plt.xlabel('Первый класс')
    plt.ylabel('Второй класс')
    plt.title(classifier_name.upper())
    plot_2d_separator(classifier, X, fill=True)
    plt.scatter(X[:, 0], X[:, 1], c=y, s=70)

# ### Метод k-ближайших соседей (1)

# In[11]:

from sklearn.neighbors import KNeighborsClassifier

# In[12]:

knn = KNeighborsClassifier(n_neighbors=1, metric='euclidean')

```

```

# Обучаем модель данных
knn.fit(X_train, y_train)

# Оцениваем качество модели
prediction = knn.predict(X_test)

# Выводим сводную информацию
show_info(knn, 'ближайшие соседи (1)', y_test, prediction)

# ### Метод k-ближайших соседей (3)

# In[13]:

knn = KNeighborsClassifier(n_neighbors=3, metric='euclidean')

# Обучаем модель данных
knn.fit(X_train, y_train)

# Оцениваем качество модели
prediction = knn.predict(X_test)

# Выводим сводную информацию
show_info(knn, 'ближайшие соседи (3)', y_test, prediction)

# ### Метод k-ближайших соседей (5)

# In[14]:

knn = KNeighborsClassifier(n_neighbors=5, metric='euclidean')

# Обучаем модель данных
knn.fit(X_train, y_train)

# Оцениваем качество модели
prediction = knn.predict(X_test)

# Выводим сводную информацию
show_info(knn, 'ближайшие соседи (5)', y_test, prediction)

# ### Метод k-ближайших соседей (9)

# In[15]:

knn = KNeighborsClassifier(n_neighbors=9, metric='euclidean')

# Обучаем модель данных
knn.fit(X_train, y_train)

# Оцениваем качество модели
prediction = knn.predict(X_test)

# Выводим сводную информацию
show_info(knn, 'ближайшие соседи (9)', y_test, prediction)

# ## Наивный байесовский классификатор

```

```

# In[16]:

from sklearn.naive_bayes import GaussianNB

nb = GaussianNB()

# Обучаем модель данных
nb.fit(X_train, y_train)

# Оцениваем качество модели
prediction = nb.predict(X_test)

# Выводим сводную информацию
show_info(nb, 'Наивный байесовский классификатор', y_test, prediction)

# ### Случайный лес (5)

# In[17]:

from sklearn.ensemble import RandomForestClassifier

# In[18]:

rfc = RandomForestClassifier(n_estimators=5)

# Обучаем модель данных
rfc.fit(X_train, y_train)

# Оцениваем качество модели
prediction = rfc.predict(X_test)

# Выводим сводную информацию
show_info(rfc, 'случайный лес (5)', y_test, prediction)

# ### Случайный лес (10)

# In[19]:

rfc = RandomForestClassifier(n_estimators=10)

# Обучаем модель данных
rfc.fit(X_train, y_train)

# Оцениваем качество модели
prediction = rfc.predict(X_test)

# Выводим сводную информацию
show_info(rfc, 'случайный лес (10)', y_test, prediction)

# ### Случайный лес (15)

# In[20]:

rfc = RandomForestClassifier(n_estimators=15)

```

```

# Обучаем модель данных
rfc.fit(X_train, y_train)

# Оцениваем качество модели
prediction = rfc.predict(X_test)

# Выводим сводную информацию
show_info(rfc, 'случайный лес (15)', y_test, prediction)

# ### Случайный лес (20)

# In[21]:

rfc = RandomForestClassifier(n_estimators=20)

# Обучаем модель данных
rfc.fit(X_train, y_train)

# Оцениваем качество модели
prediction = rfc.predict(X_test)

# Выводим сводную информацию
show_info(rfc, 'случайный лес (20)', y_test, prediction)

# ### Случайный лес (50)

# In[22]:

rfc = RandomForestClassifier(n_estimators=50)

# Обучаем модель данных
rfc.fit(X_train, y_train)

# Оцениваем качество модели
prediction = rfc.predict(X_test)

# Выводим сводную информацию
show_info(rfc, 'случайный лес (50)', y_test, prediction)

```