

Липецкий государственный технический университет

Факультет автоматизации и информатики

Кафедра автоматизированных систем управления

ЛАБОРАТОРНАЯ РАБОТА №2

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

Бинарная классификация

Студент

Сухоруков К.О.

Группа М-ИАП-22

Руководитель

Кургасов В.В.

Липецк 2022 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

Вариант 3.

Вариант	3
Вид классов	blobs
Random_state	41
cluster_std	3
noise	-
Centers	2

- 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:
 - Истинные и предсказанные метки классов
 - Матрицу ошибок (confusion matrix)
 - Значения полноты, точности, f1-меры и аккуратности
 - Значение площади под кривой ошибок (AUC ROC)
 - Отобразить на графике область принятия решений по каждому классу
- В качестве методов классификации использовать:
- a) Метод к-ближайших соседей ($n_neighbors = \{1, 3, 5, 9\}$)
 - b) Наивный байесовский метод
 - c) Случайный лес ($n_estimators = \{5, 10, 15, 20, 50\}$)

9) по каждому пункту работы занести в отчет программный код и результат вывода;

10) по результатам п.8 занести в отчет таблицу с результатами классификации всеми методами и выводы о наиболее подходящем методе классификации ваших данных.

Импорт необходимых библиотек

```
] import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import make_blobs
from sklearn.metrics import confusion_matrix
from sklearn.metrics import classification_report
from sklearn.metrics import accuracy_score
from sklearn.metrics import roc_auc_score
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.naive_bayes import GaussianNB
```

Рисунок 1 – Импорт библиотек

Генерация выборки

Вид класса	random_state	cluster_std	noise	centers
blobs	41	3	-	2

```
] X, y = make_blobs(centers=2, random_state=41, cluster_std=3)
```

```
] print("Координаты точек:\n", X[:15])
print("Метки класса: ", y[:15])
```

```
Координаты точек:
[[ 2.79914151 -14.93242644]
 [ 1.28649023 -8.07969146]
 [-7.35362913 -9.59045013]
 [ 0.22481978 -5.91330839]
 [ 5.73201763 -12.02990087]
 [-3.63420157 -11.90547127]
 [-4.60802251 -8.7525575 ]
 [-4.50479918 -7.91385674]
 [-1.48357432 -8.93044011]
 [ 5.81697962 -7.64675616]
 [ 3.46154729 -12.64729445]
 [ 1.63394511 -10.8150589 ]
 [-2.31842133 -8.90512868]
 [-9.20951536 -12.37297525]
 [ 5.16437178 -8.14431504]]
Метки класса: [1 1 0 1 1 0 0 0 0 1 1 0 0 1]
```

Рисунок 2 – Генерация выборки

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

<matplotlib.collections.PathCollection at 0x7f5209b854c0>

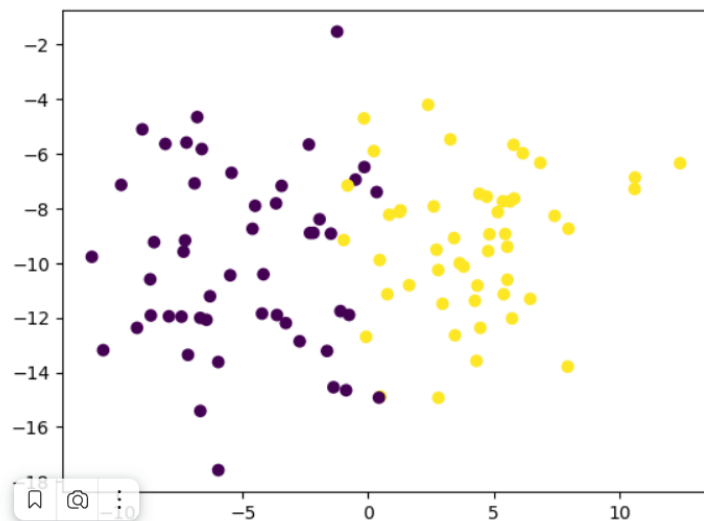


Рисунок 3 – График сгенерированной выборки

Разобьем данные на обучающие (train) и тестовые (test) выборки в пропорции 90% - 10% соответственно.

```
x_train, x_test, y_train, y_test = train_test_split(X, y, test_size = 0.10, random_state=1)
```

Рисунок 4 – Разбиение данных на обучающие и тестовые

Обучающая выборка

```
plt.scatter(X_train[:,0], X_train[:,1], c=y_train)
```

<matplotlib.collections.PathCollection at 0x7f5207a15f70>

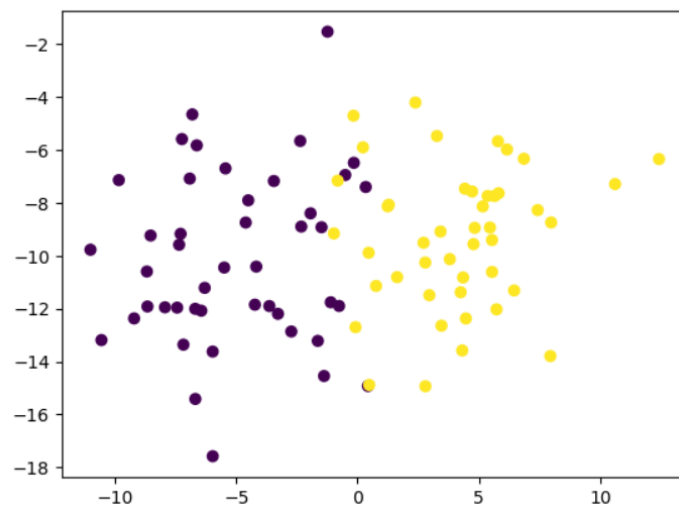


Рисунок 5 – График обучающей выборки

Тестовая выборка

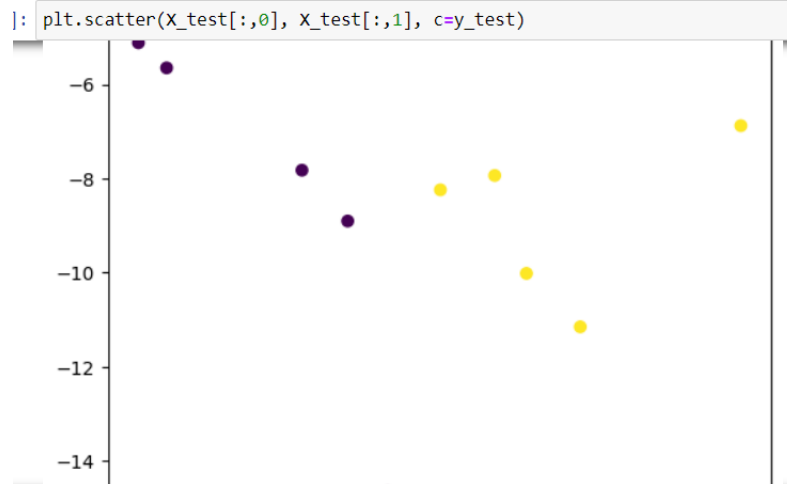


Рисунок 6 – График тестовой выборки

Метод k-ближайших соседей

```
for i in [1, 3, 5, 9]:
    knn = KNeighborsClassifier(n_neighbors=i, metric='euclidean')
    knn.fit(X_train, y_train)
    prediction = knn.predict(X_test)
    print("n_neighbors = ", i)
    print_classification_metrics(knn, X, y, prediction, y_test)
```

n_neighbors = 1
Предсказанные и истинные значения
[0 1 1 0 0 1 1 1 0 0]
[0 1 1 0 0 1 1 1 0 0]
Матрица ошибок
[[5 0]
[0 5]]
Точность классификации: 1.0
Значения полноты, точности, f1-меры и аккуратности

	precision	recall	f1-score	support
0	1.00	1.00	1.00	5
1	1.00	1.00	1.00	5
accuracy			1.00	10
macro avg	1.00	1.00	1.00	10
weighted avg	1.00	1.00	1.00	10

Значение площади под кривой ошибок (AUC ROC)
1.0

Рисунок 7 – Метод k-ближайших соседей (n=1)



Рисунок 8 – Метод k-ближайших соседей (n=1)

```
n_neighbors = 3
Предсказанные и истинные значения
[0 1 1 0 0 1 1 1 0 0]
[0 1 1 0 0 1 1 1 0 0]
Матрица ошибок
[[5 0]
 [0 5]]
Точность классификации: 1.0
Значения полноты, точности, f1-меры и аккуратности
precision    recall  f1-score   support

      0       1.00      1.00      1.00         5
      1       1.00      1.00      1.00         5

 accuracy          1.00         10
 macro avg          1.00         10
 weighted avg        1.00         10

Значение площади под кривой ошибок (AUC ROC)
1.0
```

Рисунок 9 – Метод k-ближайших соседей (n=3)

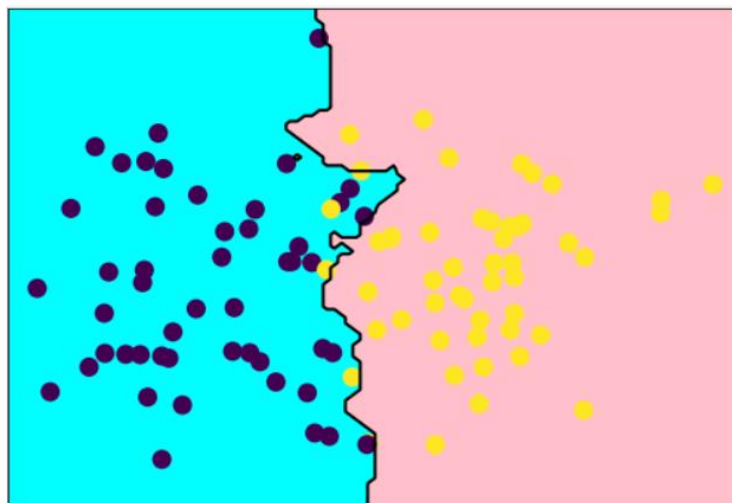


Рисунок 10 – Метод k-ближайших соседей (n=3)


```

n_neighbors = 5
Предсказанные и истинные значения
[0 1 1 0 0 1 1 1 0 0]
[0 1 1 0 0 1 1 1 0 0]
Матрица ошибок
[[5 0]
 [0 5]]
Точность классификации: 1.0
Значения полноты, точности, f1-меры и аккуратности
precision    recall  f1-score   support

      0         1.00      1.00      1.00         5
      1         1.00      1.00      1.00         5

 accuracy          1.00         10
 macro avg         1.00      1.00      1.00         10
weighted avg         1.00      1.00      1.00         10

Значение площади под кривой ошибок (AUC ROC)
1.0

```

Рисунок 11 – Метод k-ближайших соседей (n=5)

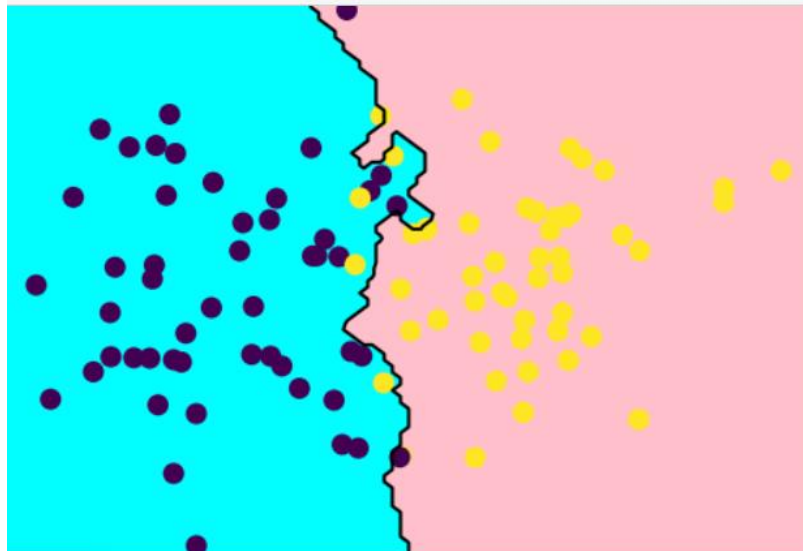


Рисунок 12 – Метод k-ближайших соседей (n=5)

```

n_neighbors = 9
Предсказанные и истинные значения
[0 1 1 0 0 1 1 1 0 0]
[0 1 1 0 0 1 1 1 0 0]
Матрица ошибок
[[5 0]
 [0 5]]
Точность классификации: 1.0
Значения полноты, точности, f1-меры и аккуратности
precision    recall  f1-score   support

      0         1.00      1.00      1.00         5
      1         1.00      1.00      1.00         5

 accuracy          1.00         10
 macro avg         1.00      1.00      1.00         10
weighted avg         1.00      1.00      1.00         10

Значение площади под кривой ошибок (AUC ROC)
1.0

```

Рисунок 13 – Метод k-ближайших соседей (n=9)



Рисунок 14 – Метод k-ближайших соседей (n=9)

Предсказанные и истинные значения

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

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

Матрица ошибок

[[5 0]

[0 5]]

Точность классификации: 1.0

Значения полноты, точности, f1-меры и аккуратности

	precision	recall	f1-score	support
0	1.00	1.00	1.00	5
1	1.00	1.00	1.00	5
accuracy			1.00	10
macro avg	1.00	1.00	1.00	10
weighted avg	1.00	1.00	1.00	10

Значение площади под кривой ошибок (AUC ROC)

1.0

Рисунок 15 – Наивный байесовский метод

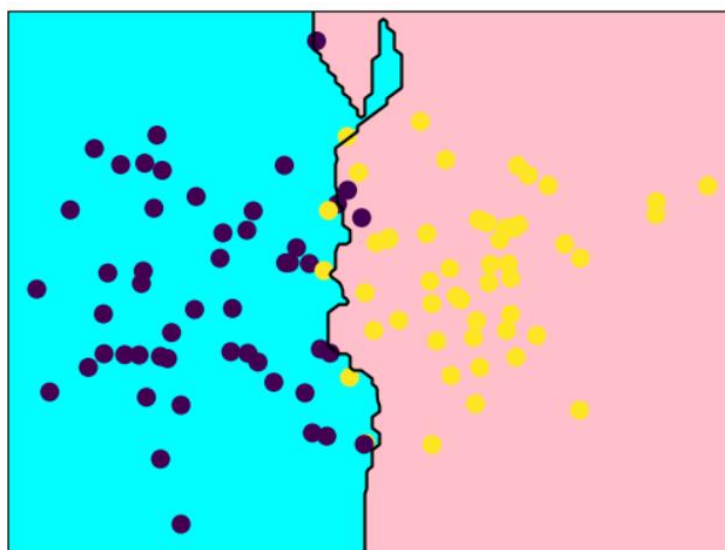


Рисунок 16 – Наивный байесовский метод

```

n_estimators = 5
Предсказанные и истинные значения
[0 1 1 0 0 1 1 1 0 0]
[0 1 1 0 0 1 1 1 0 0]
Матрица ошибок
[[5 0]
 [0 5]]
Точность классификации: 1.0
Значения полноты, точности, f1-меры и аккуратности
      precision    recall  f1-score   support

      0         1.00      1.00      1.00         5
      1         1.00      1.00      1.00         5

   accuracy                   1.00         10
  macro avg         1.00      1.00      1.00         10
 weighted avg         1.00      1.00      1.00         10

Значение площади под кривой ошибок (AUC ROC)
1.0

```

Рисунок 17 – Случайный лес n = 5

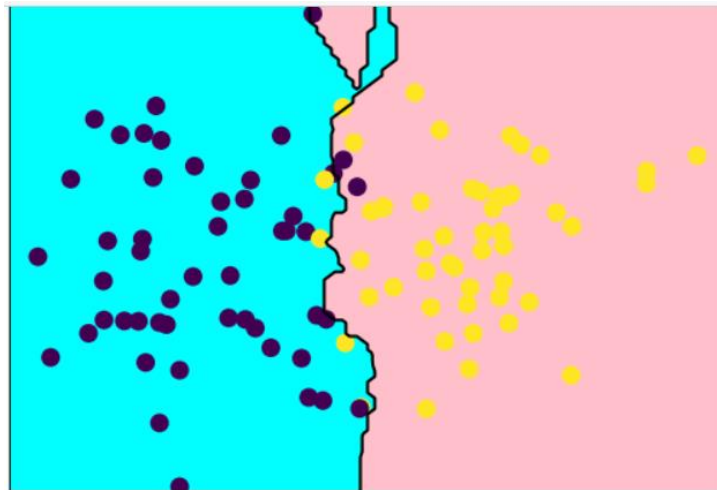


Рисунок 18 – Случайный лес n = 5

```

n_estimators = 10
Предсказанные и истинные значения
[0 1 1 0 0 1 1 1 0 0]
[0 1 1 0 0 1 1 1 0 0]
Матрица ошибок
[[5 0]
 [0 5]]
Точность классификации: 1.0
Значения полноты, точности, f1-меры и аккуратности
      precision    recall  f1-score   support

      0         1.00      1.00      1.00         5
      1         1.00      1.00      1.00         5

   accuracy                   1.00         10
  macro avg         1.00      1.00      1.00         10
 weighted avg         1.00      1.00      1.00         10

Значение площади под кривой ошибок (AUC ROC)
1.0

```

Рисунок 19 – Случайный лес n = 10

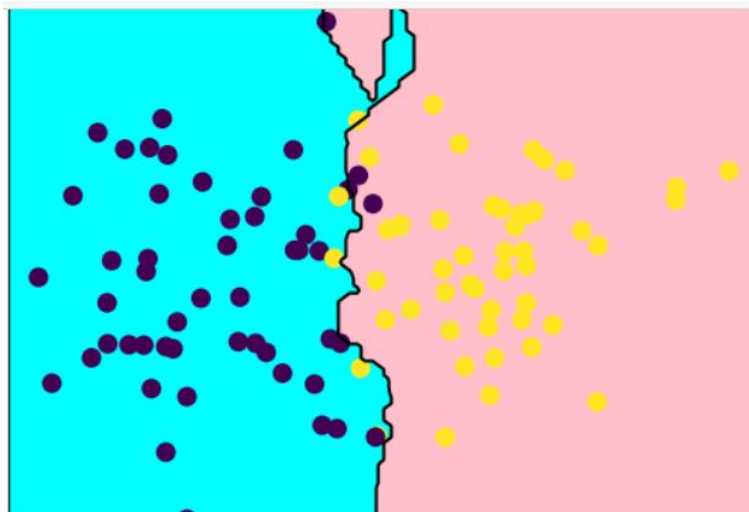


Рисунок 20 – Случайный лес $n = 10$

```
n_estimators = 15
Предсказанные и истинные значения
[0 1 1 0 0 1 1 1 0 0]
[0 1 1 0 0 1 1 1 0 0]
Матрица ошибок
[[5 0]
 [0 5]]
Точность классификации: 1.0
Значения полноты, точности, f1-меры и аккуратности
precision    recall  f1-score   support

      0       1.00      1.00      1.00         5
      1       1.00      1.00      1.00         5

 accuracy          1.00         10
 macro avg          1.00         10
weighted avg          1.00         10

Значение площади под кривой ошибок (AUC ROC)
1.0
```

Рисунок 21 – Случайный лес $n = 15$



Рисунок 22 – Случайный лес $n = 15$

```

n_estimators = 20
Предсказанные и истинные значения
[0 1 1 0 0 1 1 1 0 0]
[0 1 1 0 0 1 1 1 0 0]
Матрица ошибок
[[5 0]
 [0 5]]
Точность классификации: 1.0
Значения полноты, точности, f1-меры и аккуратности
precision    recall  f1-score   support

      0         1.00      1.00      1.00         5
      1         1.00      1.00      1.00         5

 accuracy          1.00         10
 macro avg         1.00      1.00      1.00         10
weighted avg         1.00      1.00      1.00         10

Значение площади под кривой ошибок (AUC ROC)
1.0

```

Рисунок 23 – Случайный лес $n = 20$



Рисунок 24 – Случайный лес $n = 20$

```

n_estimators = 50
Предсказанные и истинные значения
[0 1 1 0 0 1 1 1 0 0]
[0 1 1 0 0 1 1 1 0 0]
Матрица ошибок
[[5 0]
 [0 5]]
Точность классификации: 1.0
Значения полноты, точности, f1-меры и аккуратности
      precision    recall  f1-score   support

      0         1.00      1.00      1.00         5
      1         1.00      1.00      1.00         5

   accuracy                   1.00         10
  macro avg         1.00      1.00      1.00         10
 weighted avg         1.00      1.00      1.00         10

Значение площади под кривой ошибок (AUC ROC)
1.0

```

Рисунок 25 – Случайный лес n = 50



Рисунок 26 – Случайный лес n = 50

Таблица 1 – Результаты работы программы

Метод	Истинные и предсказанные метки классов	Матрица ошибок	Значения полноты, точности, f1-меры и аккуратности	Значение площади под кривой ошибок
k-ближайших соседей	<pre> [0 1 1 0 0 1 1 1 0 0] [0 1 1 0 0 1 1 1 0 0] </pre>	<pre> [5 0] [0 5] </pre>	Precision(0) = 1	1

n = 1			Precision(1) = 1 recall(0)= 1 recall(1)= 1 f1-score(0)= 1 f1-score(1)= 1 accuracy = 1	
k-ближайших соседей n = 3	$\begin{bmatrix} 0 & 1 & 1 & 0 & 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 & 1 & 1 & 0 & 0 \end{bmatrix}$	$\begin{bmatrix} 5 & 0 \\ 0 & 5 \end{bmatrix}$	Precision(0) = 1 Precision(1) = 1 recall(0)= 1 recall(1)= 1 f1-score(0)= 1 f1-score(1)= 1 accuracy = 1	1
k-ближайших соседей n = 5	$\begin{bmatrix} 0 & 1 & 1 & 0 & 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 & 1 & 1 & 0 & 0 \end{bmatrix}$	$\begin{bmatrix} 5 & 0 \\ 0 & 5 \end{bmatrix}$	Precision(0) = 1 Precision(1) = 1 recall(0)= 1 recall(1)= 1 f1-score(0)= 1 f1-score(1)= 1 accuracy = 1	1
k-ближайших соседей n = 9	$\begin{bmatrix} 0 & 1 & 1 & 0 & 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 & 1 & 1 & 0 & 0 \end{bmatrix}$	$\begin{bmatrix} 5 & 0 \\ 0 & 5 \end{bmatrix}$	Precision(0) = 1 Precision(1) = 1 recall(0)= 1 recall(1)= 1 f1-score(0)= 1 f1-score(1)= 1 accuracy = 1	1
Наивный байесовский метод	$\begin{bmatrix} 0 & 1 & 1 & 0 & 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 & 1 & 1 & 0 & 0 \end{bmatrix}$	$\begin{bmatrix} 5 & 0 \\ 0 & 5 \end{bmatrix}$	Precision(0) = 1 Precision(1) = 1 recall(0)= 1 recall(1)= 1 f1-score(0)= 1 f1-score(1)= 1 accuracy = 1	1
Случайный лес n = 5	$\begin{bmatrix} 0 & 1 & 1 & 0 & 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 & 1 & 1 & 0 & 0 \end{bmatrix}$	$\begin{bmatrix} 5 & 0 \\ 0 & 5 \end{bmatrix}$	Precision(0) = 1 Precision(1) = 1 recall(0)= 1 recall(1)= 1 f1-score(0)=	1

			1 f1-score(1)= 1 accuracy = 1	
Случайный лес n = 10	<code>[0 1 1 0 0 1 1 1 0 0]</code> <code>[0 1 1 0 0 1 1 1 0 0]</code>	<code>[5 0]</code> <code>[0 5]</code>	Precision(0) = 1 Precision(1) = 1 recall(0)= 1 recall(1)= 1 f1-score(0)= 1 f1-score(1)= 1 accuracy =	1
Случайный лес n = 15	<code>[0 1 1 0 0 1 1 1 0 0]</code> <code>[0 1 1 0 0 1 1 1 0 0]</code>	<code>[5 0]</code> <code>[0 5]</code>	Precision(0) = 1 Precision(1) = 1 recall(0)= 1 recall(1)= 1 f1-score(0)= 1 f1-score(1)= 1 accuracy = 1	1
Случайный лес n = 20	<code>[0 1 1 0 0 1 1 1 0 0]</code> <code>[0 1 1 0 0 1 1 1 0 0]</code>	<code>[5 0]</code> <code>[0 5]</code>	Precision(0) = 1 Precision(1) = 1 recall(0)= 1 recall(1)=1 f1-score(0)= 1 f1-score(1)= 1 accuracy = 1	1
Случайный лес n = 50	<code>[0 1 1 0 0 1 1 1 0 0]</code> <code>[0 1 1 0 0 1 1 1 0 0]</code>	<code>[5 0]</code> <code>[0 5]</code>	Precision(0) = 1 Precision(1) = 1 recall(0)= 1 recall(1)= 1 f1-score(0)= 1 f1-score(1)= 1 accuracy = 1	1

Аккуратность при данном разбиении выборки одинакова, следовательно все методы подходят для классификации данных.

Разобьем данные на обучающие (train) и тестовые (test) выборки в пропорции 75% - 25% соответственно.

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, random_state=1)
```

Рисунок 27 – Разбиение данных на обучающие и тестовые

Обучающая выборка

```
plt.scatter(X_train[:,0], X_train[:,1], c=y_train)
```

<matplotlib.collections.PathCollection at 0x273cb7980a0>

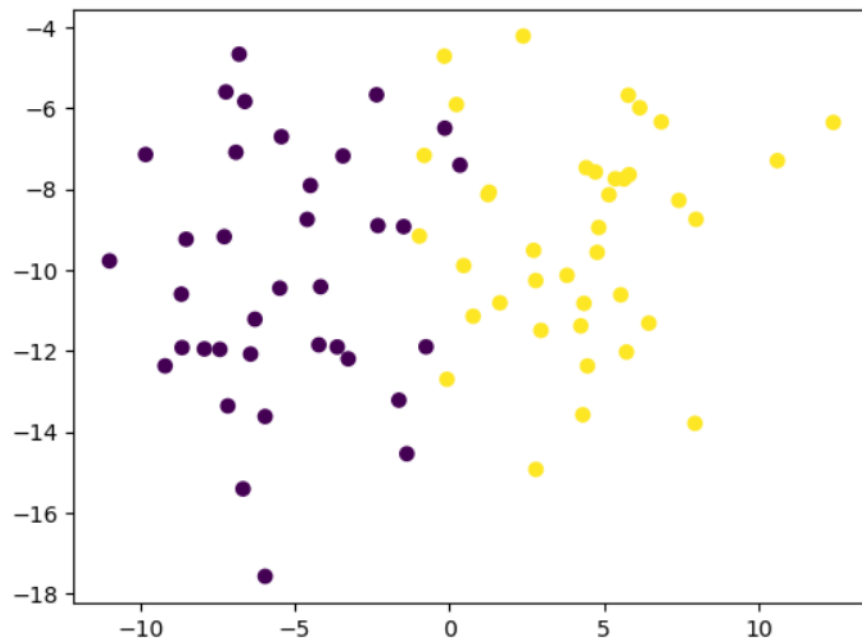


Рисунок 28 – График обучающей выборки

Тестовая выборка

```
plt.scatter(X_test[:,0], X_test[:,1], c=y_test)
```

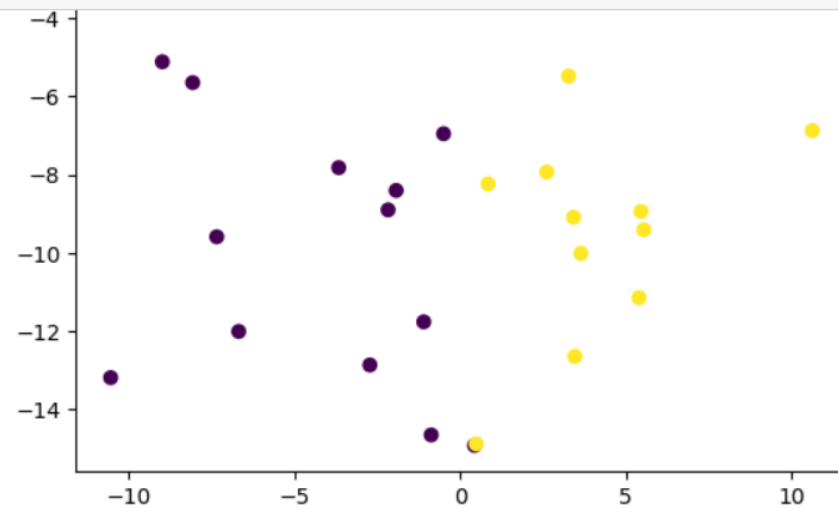


Рисунок 29 – График тестовой выборки

```

n_neighbors = 1
Предсказанные и истинные значения
[0 1 1 0 0 1 1 1 0 0 0 1 0 0 1 0 1 1 0 1 0 0 1 1 0]
[0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 1 0 0]
Матрица ошибок
[[12  2]
 [ 1 10]]
Точность классификации: 0.88
Значения полноты, точности, f1-меры и аккуратности
      precision    recall  f1-score   support

      0         0.92      0.86      0.89         14
      1         0.83      0.91      0.87         11

   accuracy              0.88         25
  macro avg              0.88      0.88      0.88         25
 weighted avg              0.88      0.88      0.88         25

Значение площади под кривой ошибок (AUC ROC)
0.8831168831168832

```

Рисунок 30 – Метод k-ближайших соседей (n=1)



Рисунок 31– Метод k-ближайших соседей (n=1)

```
n_neighbors = 3
Предсказанные и истинные значения
[0 1 1 0 0 1 1 1 0 0 1 1 0 0 1 0 1 1 0 1 1 0 0]
[0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 1 0 0]
Матрица ошибок
[[12  2]
 [ 0 11]]
Точность классификации: 0.92
Значения полноты, точности, f1-меры и аккуратности
precision    recall  f1-score   support

      0         1.00      0.86      0.92         14
      1         0.85      1.00      0.92         11

 accuracy          0.92          0.92          0.92         25
 macro avg          0.92          0.93          0.92         25
 weighted avg        0.93          0.92          0.92         25

Значение площади под кривой ошибок (AUC ROC)
0.9285714285714286
```

Рисунок 32 – Метод k-ближайших соседей (n=3)



Рисунок 33 – Метод k-ближайших соседей (n=3)

```
n_neighbors = 5
Предсказанные и истинные значения
[0 1 1 0 0 1 1 1 0 0 0 1 0 0 1 0 1 1 0 1 0 0 1 1 0]
[0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 1 0 0]
Матрица ошибок
[[12  2]
 [ 1 10]]
Точность классификации: 0.88
Значения полноты, точности, f1-меры и аккуратности
precision    recall  f1-score   support

      0         0.92      0.86      0.89         14
      1         0.83      0.91      0.87         11

 accuracy          0.88          0.88          0.88         25
 macro avg          0.88          0.88          0.88         25
 weighted avg        0.88          0.88          0.88         25

Значение площади под кривой ошибок (AUC ROC)
0.8831168831168832
```

Рисунок 34 – Метод k-ближайших соседей (n=5)

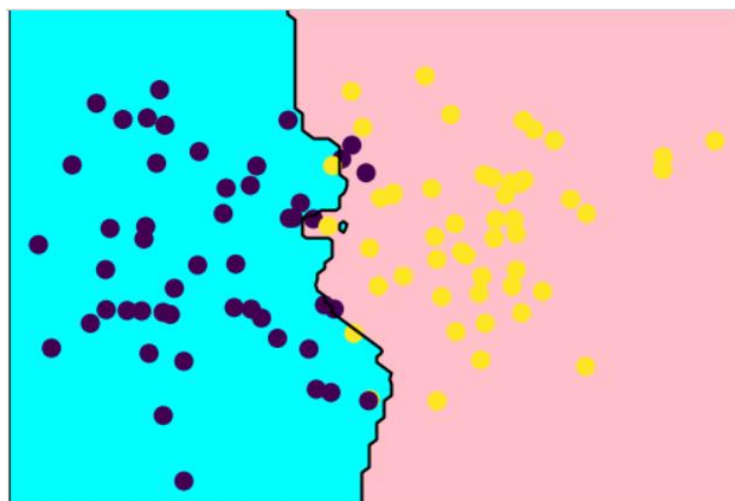


Рисунок 35 – Метод k-ближайших соседей (n=5)

```
n_neighbors = 9
Предсказанные и истинные значения
[0 1 1 0 0 1 1 1 0 0 1 0 0 0 1 0 1 1 0 1 1 1 0]
[0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 0]
Матрица ошибок
[[12  2]
 [ 0 11]]
Точность классификации: 0.92
Значения полноты, точности, f1-меры и аккуратности
precision    recall  f1-score   support

      0         1.00      0.86      0.92         14
      1         0.85      1.00      0.92         11

 accuracy          0.92          0.92         25
 macro avg         0.92         0.93      0.92         25
weighted avg         0.93         0.92      0.92         25

Значение площади под кривой ошибок (AUC ROC)
0.9285714285714286
```

Рисунок 36 – Метод k-ближайших соседей (n=9)

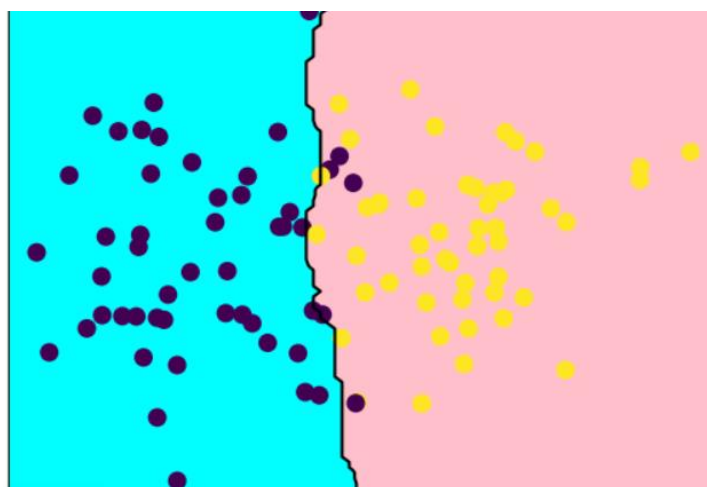


Рисунок 37 – Метод k-ближайших соседей (n=9)

Предсказанные и истинные значения
 [0 1 1 0 0 1 1 1 0 0 1 0 0 0 1 0 1 1 0 1 1 1 0]
 [0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 1 1 0]
 Матрица ошибок
 [[12 2]
 [0 11]]
 Точность классификации: 0.92
 Значения полноты, точности, f1-меры и аккуратности

	precision	recall	f1-score	support
0	1.00	0.86	0.92	14
1	0.85	1.00	0.92	11
accuracy			0.92	25
macro avg	0.92	0.93	0.92	25
weighted avg	0.93	0.92	0.92	25

Значение площади под кривой ошибок (AUC ROC)
 0.9285714285714286

Рисунок 38 – Наивный байесовский метод

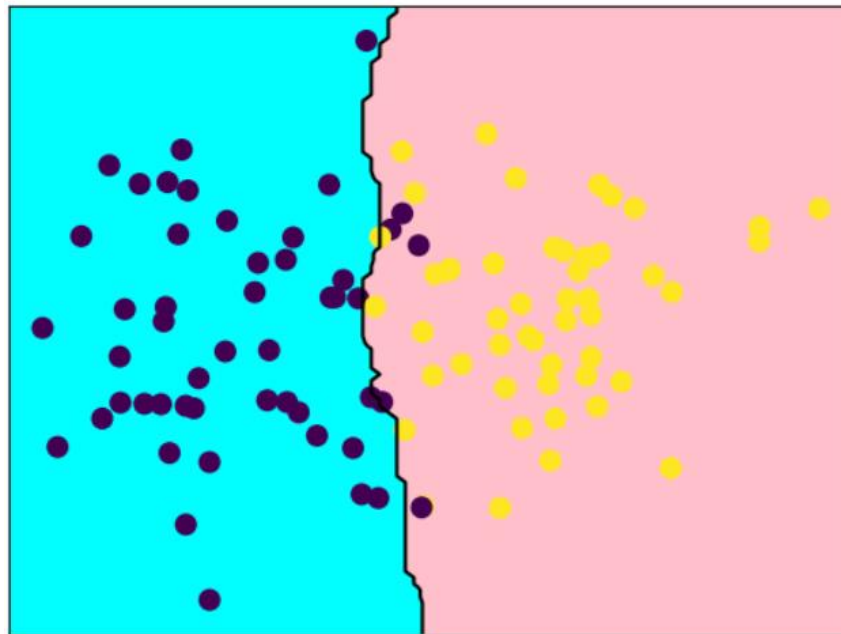


Рисунок 39 – Наивный байесовский метод

```

n_estimators = 5
Предсказанные и истинные значения
[0 1 1 0 0 1 1 1 0 0 1 0 0 0 1 0 1 1 0 1 0 1 1 1 0]
[0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 1 0 0]
Матрица ошибок
[[12  2]
 [ 0 11]]
Точность классификации: 0.92
Значения полноты, точности, f1-меры и аккуратности
      precision    recall  f1-score   support

      0         1.00      0.86      0.92         14
      1         0.85      1.00      0.92         11

   accuracy                    0.92         25
  macro avg         0.92      0.93      0.92         25
 weighted avg         0.93      0.92      0.92         25

Значение площади под кривой ошибок (AUC ROC)
0.9285714285714286

```

Рисунок 40 – Случайный лес $n = 5$

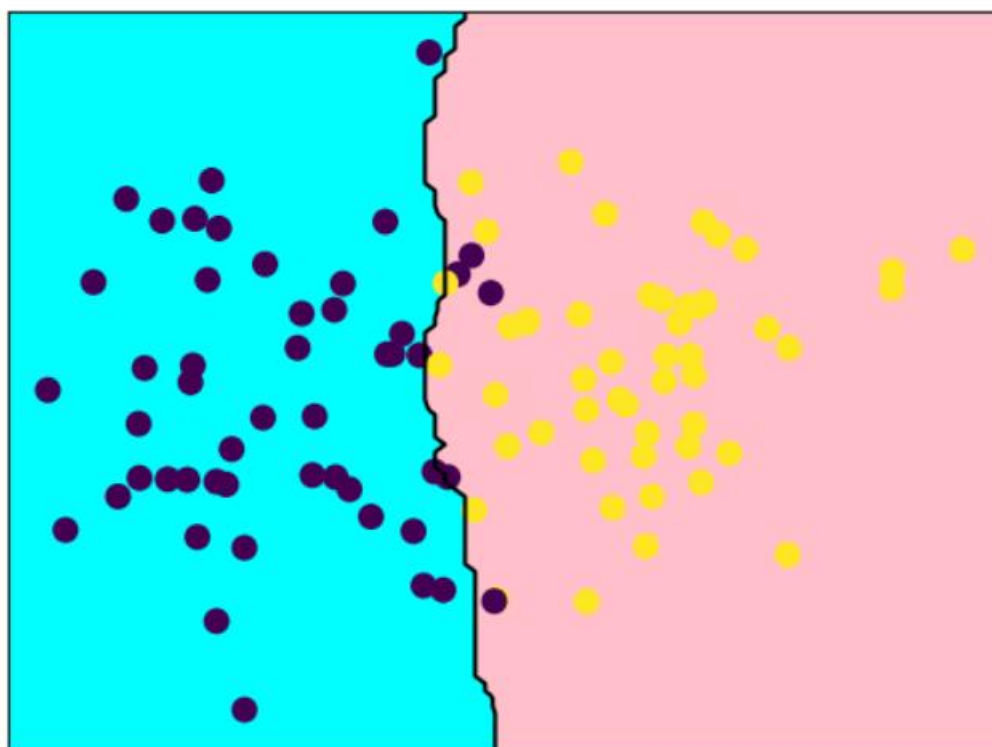


Рисунок 41 – Случайный лес $n = 5$

```

n_estimators = 10
Предсказанные и истинные значения
[0 1 1 0 0 1 1 1 0 0 1 1 0 0 1 0 1 1 0 1 0 1 1 1 0]
[0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 1 0 0]
Матрица ошибок
[[11  3]
 [ 0 11]]
Точность классификации: 0.88
Значения полноты, точности, f1-меры и аккуратности
precision    recall  f1-score   support

      0       1.00      0.79      0.88        14
      1       0.79      1.00      0.88        11

 accuracy          0.88        25
 macro avg         0.89      0.89      0.88        25
weighted avg         0.91      0.88      0.88        25

Значение площади под кривой ошибок (AUC ROC)
0.8928571428571428

```

Рисунок 42 – Случайный лес $n = 10$

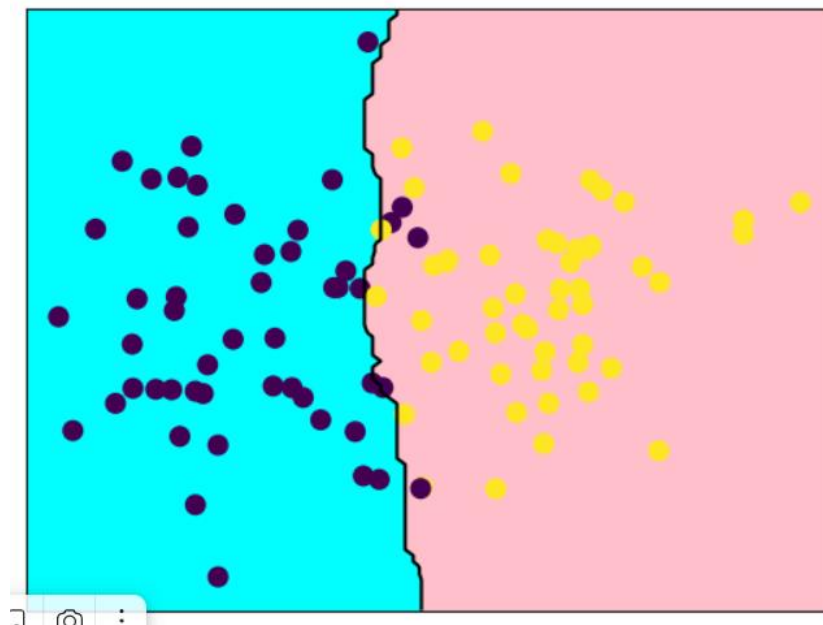


Рисунок 43 – Случайный лес $n = 10$

```

n_estimators = 15
Предсказанные и истинные значения
[0 1 1 0 0 1 1 1 0 0 1 0 0 0 1 0 1 1 0 1 1 1 0]
[0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 0]
Матрица ошибок
[[12  2]
 [ 0 11]]
Точность классификации: 0.92
Значения полноты, точности, f1-меры и аккуратности
      precision    recall  f1-score   support

      0         1.00      0.86      0.92         14
      1         0.85      1.00      0.92         11

   accuracy                   0.92         25
  macro avg       0.92      0.93      0.92         25
 weighted avg     0.93      0.92      0.92         25

Значение площади под кривой ошибок (AUC ROC)
0.9285714285714286

```

Рисунок 44 – Случайный лес n = 15

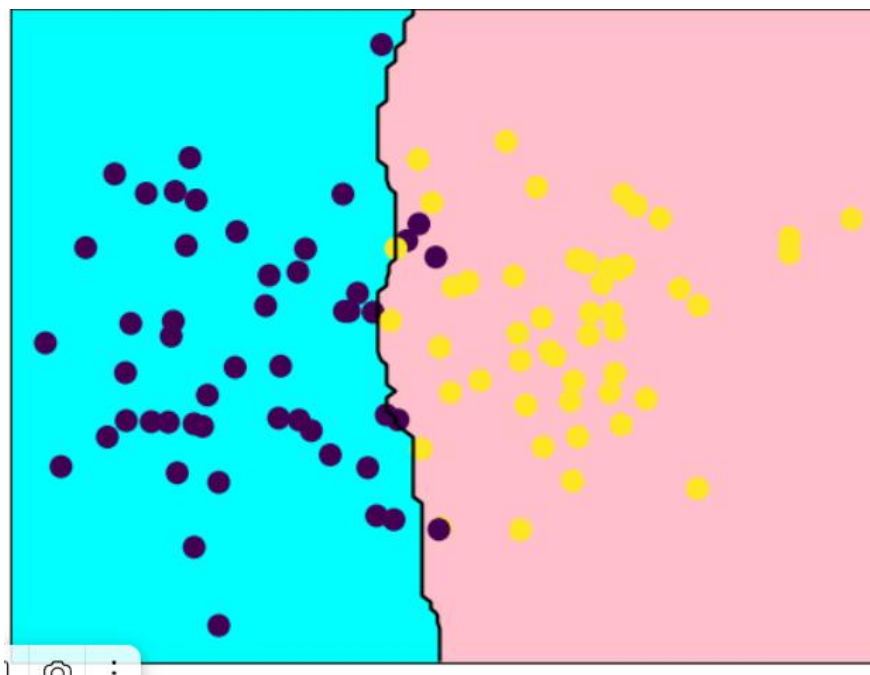


Рисунок 45 – Случайный лес n = 15


```

n_estimators = 20
Предсказанные и истинные значения
[0 1 1 0 0 1 1 1 0 0 1 0 0 0 1 0 1 1 0 1 1 1 0]
[0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 1 1 0]
Матрица ошибок
[[12  2]
 [ 0 11]]
Точность классификации: 0.92
Значения полноты, точности, f1-меры и аккуратности
      precision    recall  f1-score   support

      0         1.00      0.86      0.92         14
      1         0.85      1.00      0.92         11

   accuracy                   0.92         25
  macro avg         0.92      0.93      0.92         25
 weighted avg         0.93      0.92      0.92         25

Значение площади под кривой ошибок (AUC ROC)
0.9285714285714286

```

Рисунок 46 – Случайный лес $n = 20$

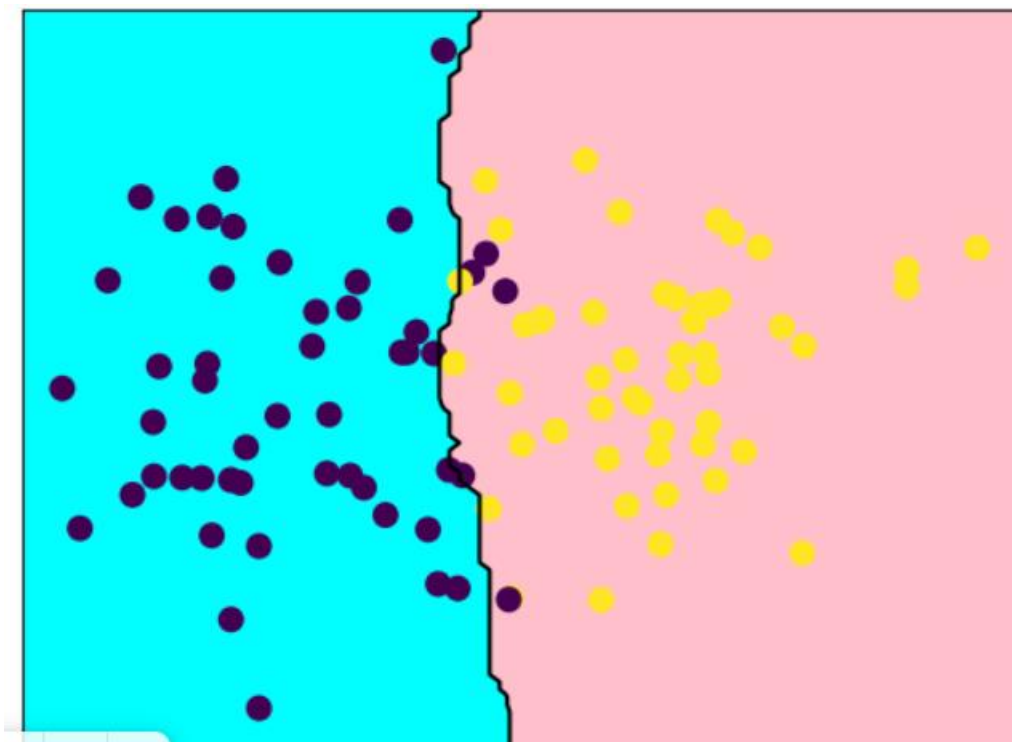


Рисунок 47 – Случайный лес $n = 20$

```

n_estimators = 50
Предсказанные и истинные значения
[0 1 1 0 0 1 1 1 0 0 1 0 0 0 1 0 1 1 0 1 0 1 1 1 0]
[0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 1 0 0]
Матрица ошибок
[[12  2]
 [ 0 11]]
Точность классификации: 0.92
Значения полноты, точности, f1-меры и аккуратности
      precision    recall  f1-score   support

      0         1.00      0.86      0.92         14
      1         0.85      1.00      0.92         11

   accuracy          0.92
  macro avg          0.92
 weighted avg          0.92

Значение площади под кривой ошибок (AUC ROC)
0.9285714285714286

```

Рисунок 48 – Случайный лес $n = 50$

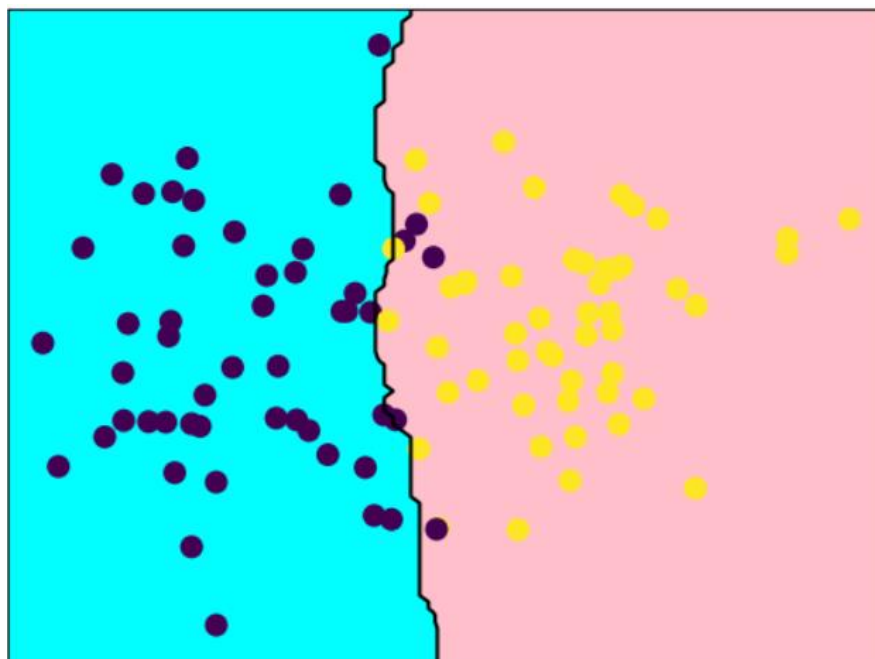


Рисунок 49 – Случайный лес $n = 50$

Таблица 2 – Результаты работы программы

Метод	Истинные и предсказанные метки классов	Матрица ошибок	Значения полноты, точности, f1-меры и	Значени е площад и под

			аккуратность и	кривой ошибок
k- ближайших соседей n = 1	[0 1 1 0 0 1 1 1 0 0 0 1 0 0 1 0 1 1 0 1 0 0 1 1 0] [0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 1 0 0]	[12 2] [1 10]	Precision(0) = 0.92 Precision(1) = 0.83 recall(0)= 0 .86 recall(1)= 0 .91 f1-score(0)= 0.89 f1-score(1)= 0.87 accuracy = 0 .88	0.88312
k- ближайших соседей n = 3	[0 1 1 0 0 1 1 1 0 0 1 1 0 0 1 0 1 1 0 1 0 1 1 0 0] [0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 1 0 0]	[12 2] [0 11]	Precision(0) = 1 Precision(1) = 0.85 recall(0)= 0 .86 recall(1)= 1 f1-score(0)= 0.92 f1-score(1)= 0.92 accuracy = 0 .92	0.92857
k- ближайших соседей n = 5	[0 1 1 0 0 1 1 1 0 0 0 1 0 0 1 0 1 1 0 1 0 0 1 1 0] [0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 1 0 0]	[12 2] [1 10]	Precision(0) = 0.92 Precision(1) = 0.83 recall(0)= 0 .86 recall(1)= 0 .91 f1-score(0)= 0.89 f1-score(1)= 0.87 accuracy = 0 .88	0.88312
k- ближайших соседей n = 9	[0 1 1 0 0 1 1 1 0 0 1 0 0 0 1 0 1 1 0 1 0 1 1 1 0] [0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 1 0 0]	[12 2] [0 11]	Precision(0) = 1 Precision(1) = 0.85 recall(0)= 0 .86 recall(1)= 1 f1-score(0)= 0.92 f1-score(1)= 0.92 accuracy = 0 .92	0.928571

Наивный байесовский метод	<pre> [0 1 1 0 0 1 1 1 0 0 1 0 0 0 1 0 1 1 0 1 0 1 1 1 0] [0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 1 0 0] </pre>	<pre> [[12 2] [0 11]] </pre>	Precision(0) = 1 Precision(1) = 0.85 recall(0) = 0.86 recall(1) = 1 f1-score(0) = 0.92 f1-score(1) = 0.92 accuracy = 0.92	0.92857
Случайный лес n = 5	<pre> [0 1 1 0 0 1 1 1 0 0 1 0 0 0 1 0 1 1 0 1 0 1 1 1 0] [0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 1 0 0] </pre>	<pre> [[12 2] [0 11]] </pre>	Precision(0) = 1 Precision(1) = 0.85 recall(0) = 0.86 recall(1) = 1 f1-score(0) = 0.92 f1-score(1) = 0.92 accuracy = 0.92	0.92857
Случайный лес n = 10	<pre> [0 1 1 0 0 1 1 1 0 0 1 1 0 0 1 0 1 1 0 1 0 1 1 1 0] [0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 1 0 0] </pre>	<pre> [[11 3] [0 11]] </pre>	Precision(0) = 1 Precision(1) = 0.79 recall(0) = 0.79 recall(1) = 1 f1-score(0) = 0.88 f1-score(1) = 0.88 accuracy = 0.88	0.892857
Случайный лес n = 15	<pre> [0 1 1 0 0 1 1 1 0 0 1 0 0 0 1 0 1 1 0 1 0 1 1 1 0] [0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 1 0 0] </pre>	<pre> [[12 2] [0 11]] </pre>	Precision(0) = 1 Precision(1) = 0.85 recall(0) = 0.86 recall(1) = 1 f1-score(0) = 0.92 f1-score(1) = 0.92 accuracy = 0.92	0.92857
Случайный лес n = 20	<pre> [0 1 1 0 0 1 1 1 0 0 1 0 0 0 1 0 1 1 0 1 0 1 1 1 0] [0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 1 0 0] </pre>	<pre> [[12 2] [0 11]] </pre>	Precision(0) = 1 Precision(1) = 0.85 recall(0) = 0.86 recall(1) = 1 f1-score(0) =	0.92857

			0.92 f1-score(1)= 0.92 accuracy = 0 .92	
Случайный лес n = 50	представление в виде матрицы [0 1 1 0 0 1 1 1 0 0 1 0 0 0 1 0 1 1 0 1 0 1 1 1 0] [0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 1 0 0]	матрица [12 2] [0 11]	Precision(0) = 1 Precision(1) = 0.85 recall(0)= 0 .86 recall(1)= 1 f1-score(0)= 0.92 f1-score(1)= 0.92 accuracy = 0 .92	0.92857

Наиболее подходящие методы для классификации данных случайный лес n=5, случайный лес n=15, случайный лес n=20, случайный лес n=50, наивный байесовский метод, k-ближайших соседей n = 3, k-ближайших соседей = 9 так как в этих методах наибольшая аккуратность.

Разобьем данные на обучающие (train) и тестовые (test) выборки в пропорции 65% - 35% соответственно.

```
x_train, x_test, y_train, y_test = train_test_split(X, y, test_size = 0.35, random_state=1)
```

Рисунок 50 – Разбиение данных на обучающие и тестовые

Обучающая выборка

```
plt.scatter(X_train[:,0], X_train[:,1], c=y_train)
<matplotlib.collections.PathCollection at 0x273cca75ac0>
```

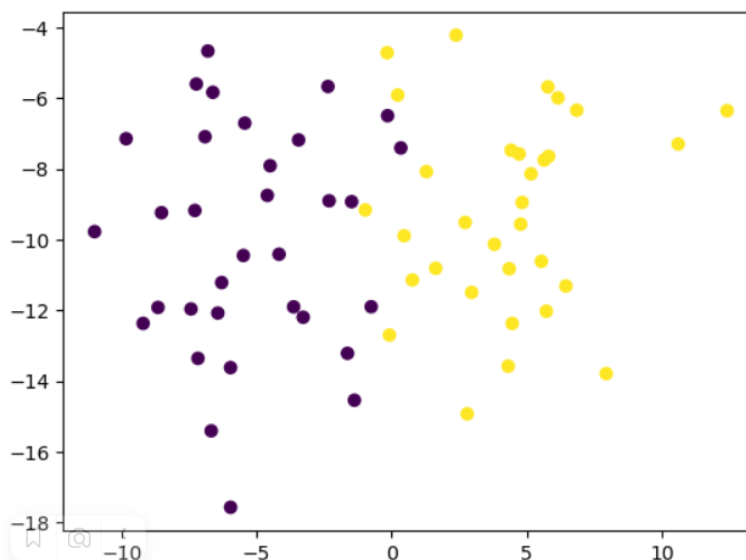


Рисунок 51 – График обучающей выборки

Тестовая выборка

```
plt.scatter(X_test[:,0], X_test[:,1], c=y_test)
```

<matplotlib.collections.PathCollection at 0x273ccadb40>

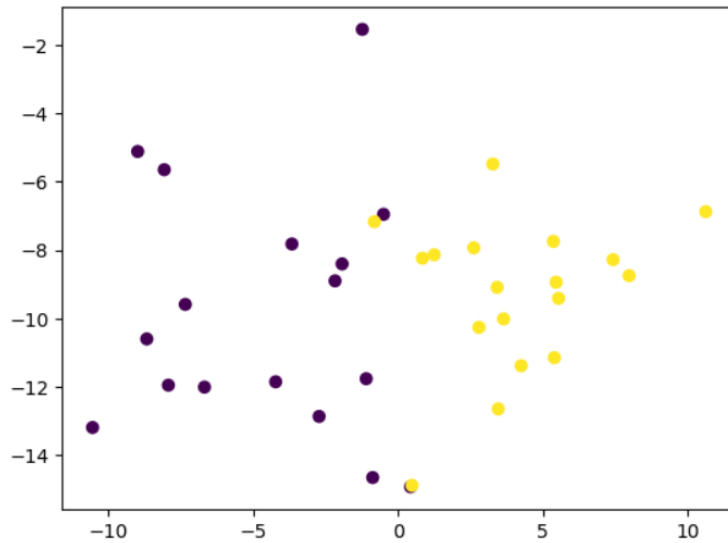


Рисунок 52 – График тестовой выборки

```
n_neighbors = 1
```

Предсказанные и истинные значения

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

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

Матрица ошибок

```
[[16  1]
 [ 2 16]]
```

Точность классификации: 0.9142857142857143

Значения полноты, точности, f1-меры и аккуратности

	precision	recall	f1-score	support
0	0.89	0.94	0.91	17
1	0.94	0.89	0.91	18
accuracy			0.91	35
macro avg	0.92	0.92	0.91	35
weighted avg	0.92	0.91	0.91	35

Значение площади под кривой ошибок (AUC ROC)

0.9150326797385621

Рисунок 53 – Метод k-ближайших соседей (n=1)

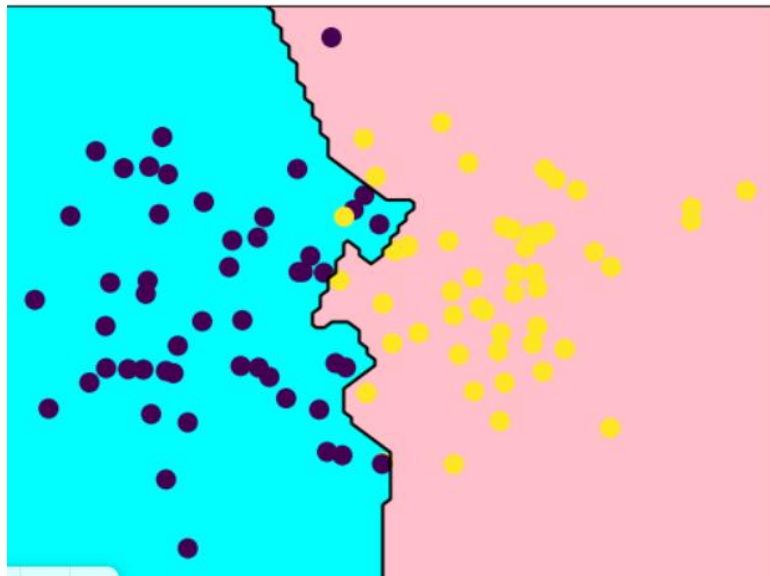


Рисунок 54 – Метод k-ближайших соседей (n=1)

```
n_neighbors = 3
Предсказанные и истинные значения
[0 1 1 0 0 1 1 1 0 0 1 1 0 0 1 0 1 1 0 1 1 0 0 1 1 1 0 1 1 0 0 1 0]
[0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 1 0 0 1 1 1 0 1 1 0 1 0]
Матрица ошибок
[[15  2]
 [ 1 17]]
Точность классификации: 0.9142857142857143
Значения полноты, точности, f1-меры и акkuratности
precision    recall  f1-score   support

           0           0.94      0.88      0.91         17
           1           0.89      0.94      0.92         18

 accuracy          0.91         35
 macro avg         0.92      0.91      0.91         35
 weighted avg      0.92      0.91      0.91         35

Значение площади под кривой ошибок (AUC ROC)
0.9133986928104575
```

Рисунок 55 – Метод k-ближайших соседей (n=3)

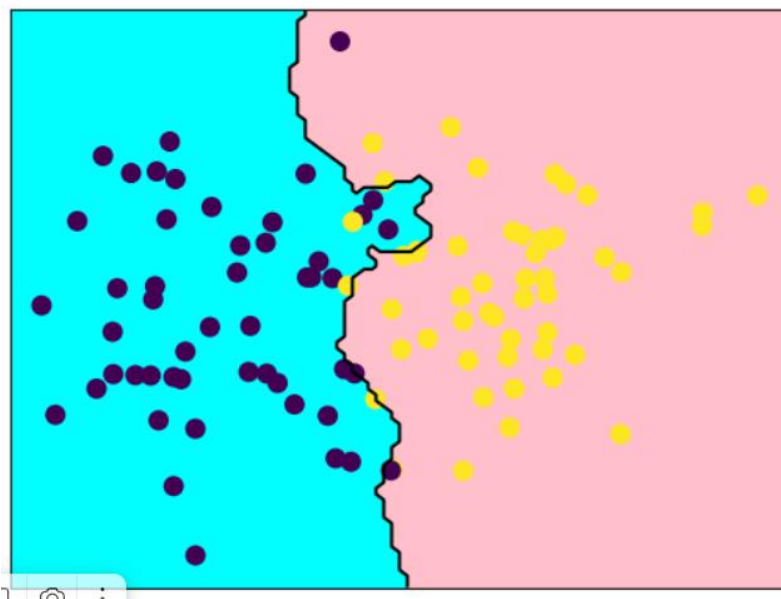


Рисунок 56 – Метод k-ближайших соседей (n=3)

```

n_neighbors = 5
Предсказанные и истинные значения
[0 1 1 0 0 1 1 1 0 0 0 1 0 0 1 0 1 1 0 1 0 0 1 0 0 1 1 1 0 1 1 0 0 1 0]
[0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 1 0 0 1 1 1 0 1 1 0 1 0]
Матрица ошибок
[[16  1]
 [ 2 16]]
Точность классификации: 0.9142857142857143
Значения полноты, точности, f1-меры и аккуратности
precision    recall  f1-score   support

      0       0.89       0.94       0.91        17
      1       0.94       0.89       0.91        18

accuracy          0.91        35
macro avg       0.92       0.92       0.91        35
weighted avg    0.92       0.91       0.91        35

Значение площади под кривой ошибок (AUC ROC)
0.9150326797385621

```

Рисунок 57 – Метод k-ближайших соседей (n=5)

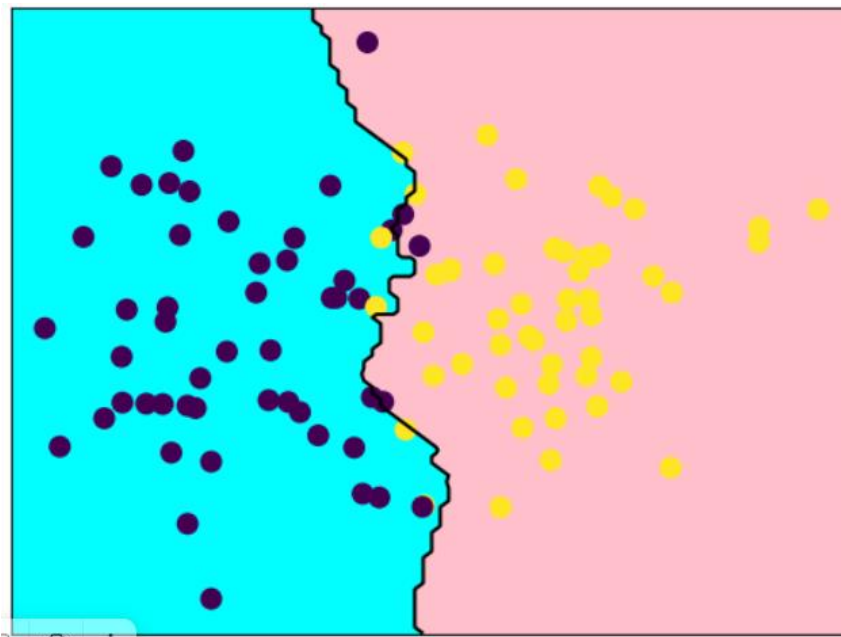


Рисунок 58 – Метод k-ближайших соседей (n=5)

```

n_neighbors = 9
Предсказанные и истинные значения
[0 1 1 0 0 1 1 1 0 0 1 0 0 0 1 0 1 1 0 1 0 1 1 0 0 1 1 1 0 1 1 0 0 1 0]
[0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 1 0 0 1 1 1 0 1 1 1 0 1 0]
Матрица ошибок
[[16  1]
 [ 1 17]]
Точность классификации: 0.9428571428571428
Значения полноты, точности, f1-меры и аккуратности
precision    recall  f1-score   support

      0       0.94       0.94       0.94        17
      1       0.94       0.94       0.94        18

accuracy          0.94        35
macro avg       0.94       0.94       0.94        35
weighted avg    0.94       0.94       0.94        35

Значение площади под кривой ошибок (AUC ROC)
0.9428104575163399

```

Рисунок 59 – Метод k-ближайших соседей (n=9)

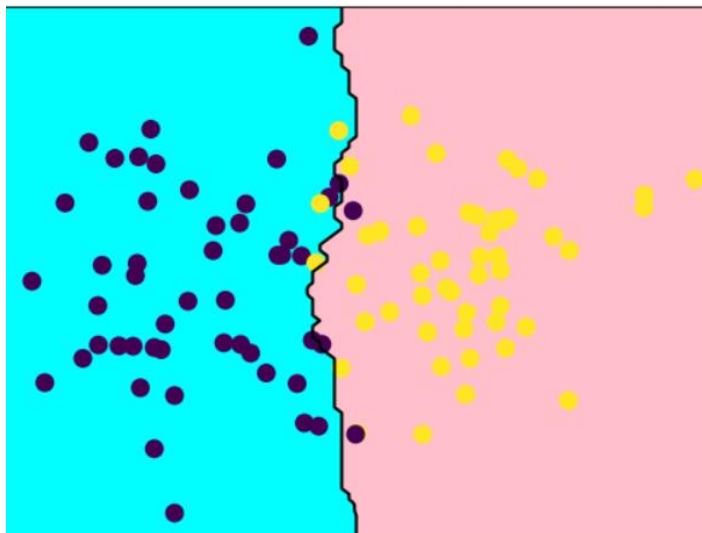


Рисунок 60 – Метод k-ближайших соседей (n=9)

Предсказанные и истинные значения
 [0 1 1 0 0 1 1 1 0 0 1 0 0 0 1 0 1 1 0 1 0 1 1 1 0 1 1 1 0 1 1 0 0 1 0]
 [0 1 1 0 0 1 1 1 0 0 0 0 0 0 0 1 0 1 1 0 1 0 1 1 0 0 1 1 1 0 1 1 0 1 0]
 Матрица ошибок
 [[15 2]
 [1 17]]
 Точность классификации: 0.9142857142857143
 Значения полноты, точности, f1-меры и аккуратности

	precision	recall	f1-score	support
0	0.94	0.88	0.91	17
1	0.89	0.94	0.92	18
accuracy			0.91	35
macro avg	0.92	0.91	0.91	35
weighted avg	0.92	0.91	0.91	35

Значение площади под кривой ошибок (AUC ROC)
 0.9133986928104575

Рисунок 61 – Наивный байесовский метод

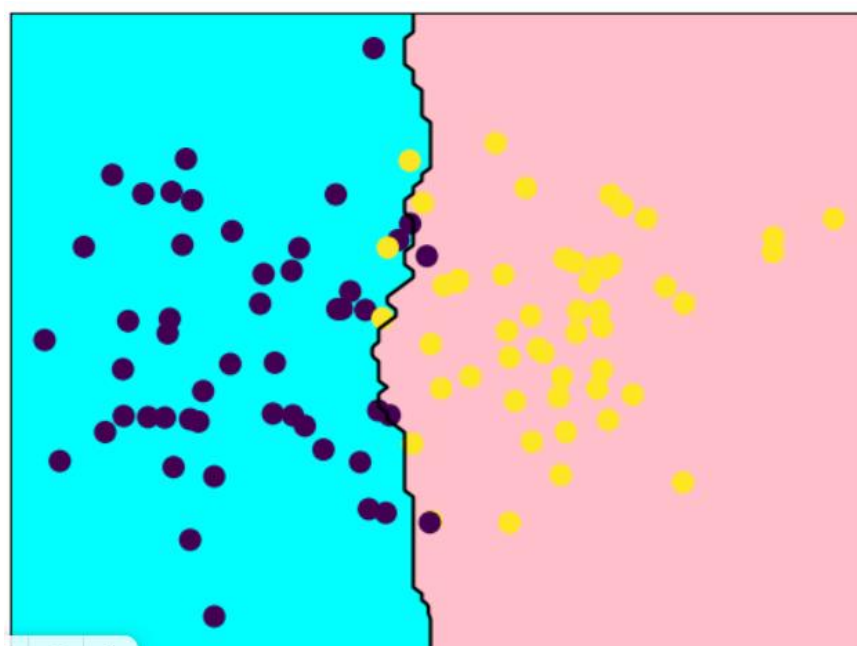


Рисунок 62 – Наивный байесовский метод

```

n_estimators = 5
Предсказанные и истинные значения
[0 1 1 0 0 1 1 1 0 0 1 0 0 0 1 0 1 1 0 1 0 1 1 0 0 1 1 0 0 1 0]
[0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 1 0 0 1 1 0 1 1 0]
Матрица ошибок
[[16  1]
 [ 1 17]]
Точность классификации: 0.9428571428571428
Значения полноты, точности, f1-меры и аккуратности
      precision    recall  f1-score   support

      0       0.94       0.94       0.94        17
      1       0.94       0.94       0.94        18

 accuracy          0.94          35
 macro avg       0.94       0.94       0.94          35
weighted avg       0.94       0.94       0.94          35

Значение площади под кривой ошибок (AUC ROC)
0.9428104575163399

```

Рисунок 63 – Случайный лес $n = 5$

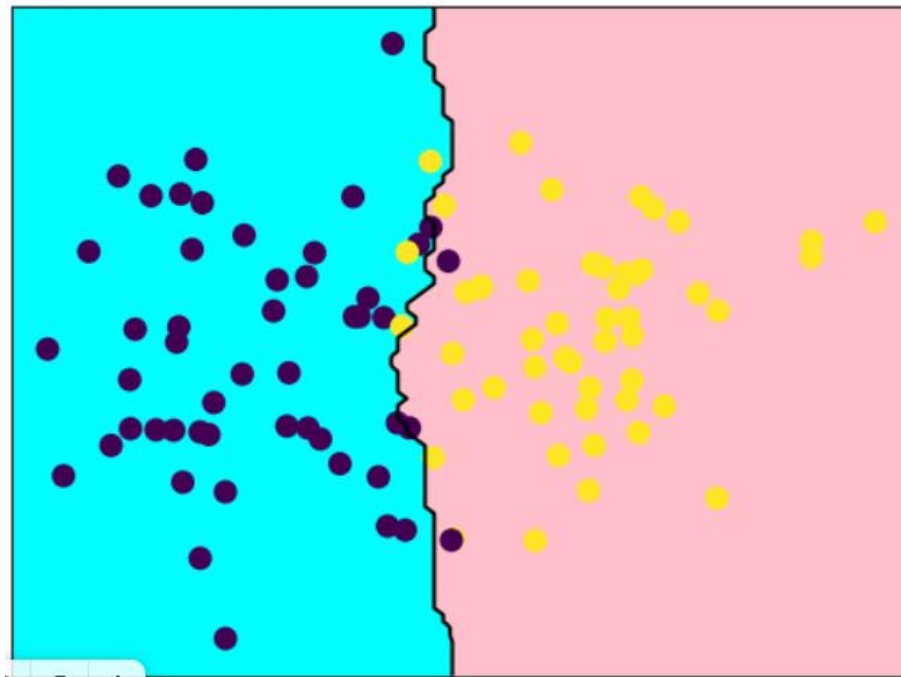


Рисунок 64 – Случайный лес $n = 5$

```

n_estimators = 10
Предсказанные и истинные значения
[0 1 1 0 0 1 1 1 0 0 1 0 0 0 1 0 1 1 0 1 1 0 0 1 1 1 0 1 1 0 0 1 0]
[0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 1 0 0 1 1 1 0 1 1 1 0 1 0]
Матрица ошибок
[[16  1]
 [ 1 17]]
Точность классификации: 0.9428571428571428
Значения полноты, точности, f1-меры и аккуратности
precision    recall  f1-score   support

      0       0.94      0.94      0.94        17
      1       0.94      0.94      0.94        18

 accuracy          0.94          35
 macro avg       0.94      0.94      0.94          35
 weighted avg    0.94      0.94      0.94          35

Значение площади под кривой ошибок (AUC ROC)
0.9428104575163399

```

Рисунок 65 – Случайный лес n = 10

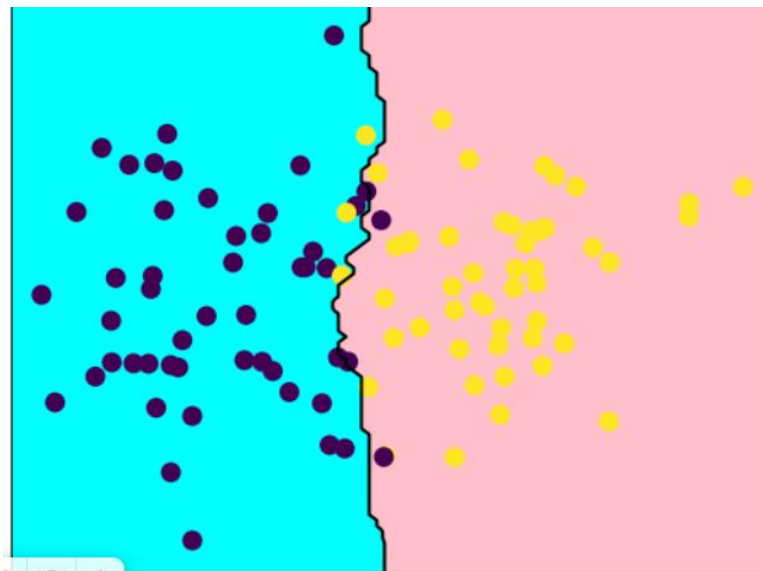


Рисунок 66 – Случайный лес n = 10

```

n_estimators = 15
Предсказанные и истинные значения
[0 1 1 0 0 1 1 1 0 0 1 0 0 0 1 0 1 1 0 1 0 1 1 0 0 1 1 1 0 1 1 0 0 1 0]
[0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 1 0 0 1 1 1 0 1 1 1 0 1 0]
Матрица ошибок
[[16  1]
 [ 1 17]]
Точность классификации: 0.9428571428571428
Значения полноты, точности, f1-меры и аккуратности
precision    recall  f1-score   support

      0       0.94      0.94      0.94        17
      1       0.94      0.94      0.94        18

 accuracy          0.94          35
 macro avg       0.94      0.94      0.94          35
 weighted avg    0.94      0.94      0.94          35

Значение площади под кривой ошибок (AUC ROC)
0.9428104575163399

```

Рисунок 67 – Случайный лес n = 15

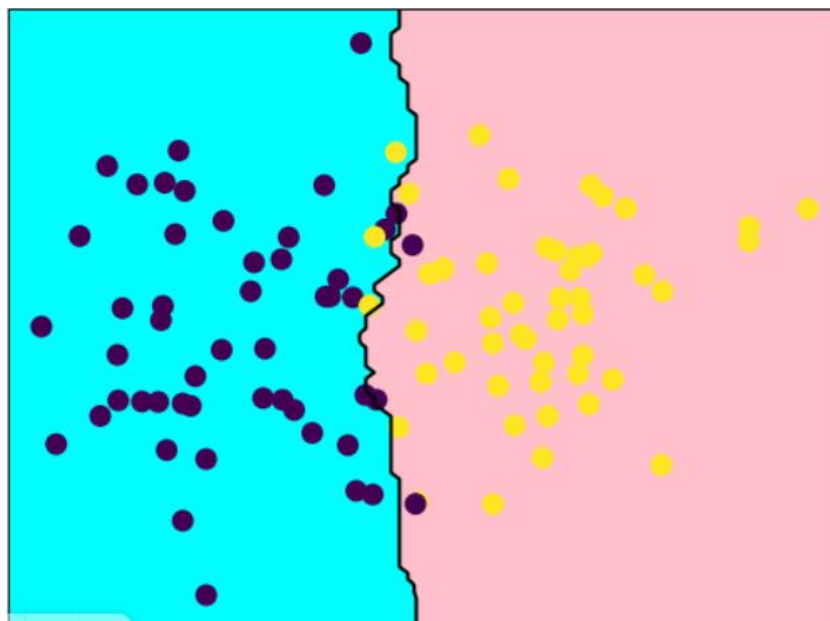


Рисунок 68 – Случайный лес $n = 15$

```
n_estimators = 20
Предсказанные и истинные значения
[0 1 1 0 0 1 1 1 0 0 1 0 0 1 0 1 1 0 1 0 1 1 0 0 1 1 1 0 1 1 0 0 1 0]
[0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 1 0 0 1 1 1 0 1 1 0 1 0]
Матрица ошибок
[[16  1]
 [ 1 17]]
Точность классификации: 0.9428571428571428
Значения полноты, точности, f1-меры и аккуратности
      precision    recall  f1-score   support

     0       0.94       0.94       0.94        17
     1       0.94       0.94       0.94        18

 accuracy          0.94          35
 macro avg       0.94       0.94       0.94          35
 weighted avg    0.94       0.94       0.94          35

Значение площади под кривой ошибок (AUC ROC)
0.9428104575163399
```

Рисунок 69 – Случайный лес $n = 20$

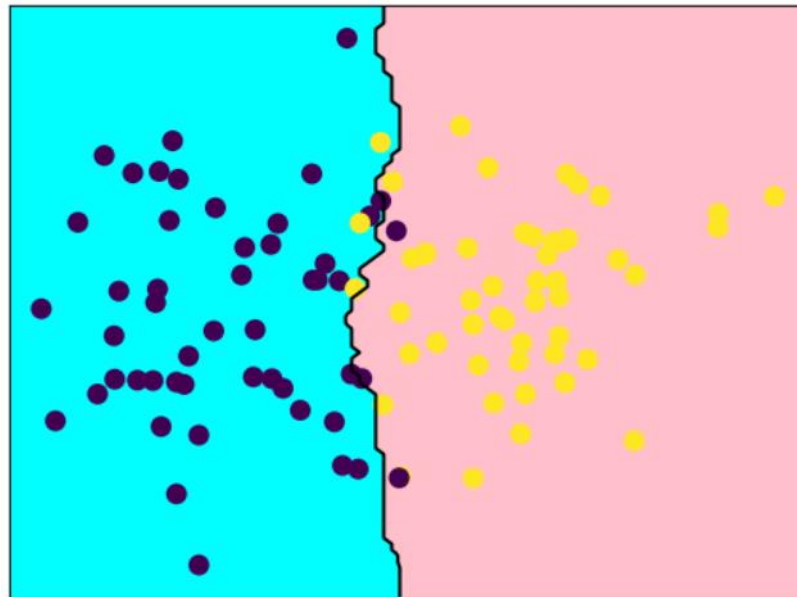


Рисунок 70 – Случайный лес $n = 20$

```
n_estimators = 50
Предсказанные и истинные значения
[0 1 1 0 0 1 1 1 0 0 1 0 0 0 1 0 1 1 0 1 0 1 1 0 0 1 1 1 0 1 1 0 0 1 0]
[0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 1 0 0 1 1 1 0 1 1 1 0 1 0]
Матрица ошибок
[[16  1]
 [ 1 17]]
Точность классификации: 0.9428571428571428
Значения полноты, точности, f1-меры и аккуратности
precision    recall  f1-score   support

      0       0.94       0.94       0.94        17
      1       0.94       0.94       0.94        18

 accuracy          0.94          35
 macro avg       0.94       0.94       0.94          35
weighted avg       0.94       0.94       0.94          35

Значение площади под кривой ошибок (AUC ROC)
0.9428104575163399
```

Рисунок 71 – Случайный лес $n = 50$

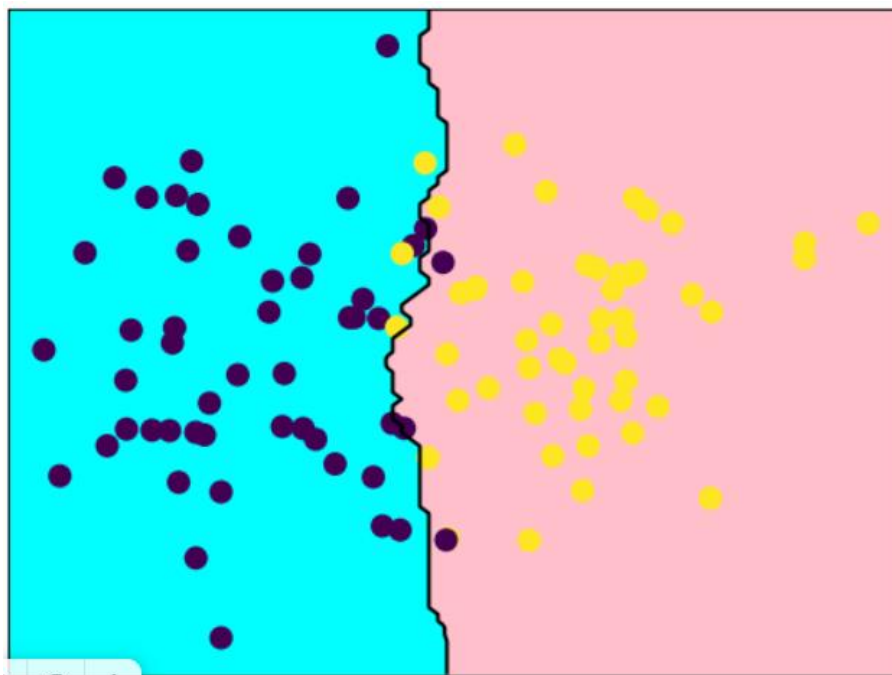


Рисунок 72 – Случайный лес $n = 50$

Таблица 3 – Результаты работы программы

Метод	Истинные и предсказанные метки классов	Матрица ошибок	Значения полноты, точности, f1-меры и аккуратности	Значение площади под кривой ошибок
k-ближайших соседей $n = 1$	<pre> [0 1 1 0 0 1 1 1 0 0 0 1 0 0 1 0 1 1 0 1 0 0 1 0 0 1 1 1 0 1 1 0 0 1 0] [0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 1 0 0 1 1 1 0 1 1 1 0 1 0] </pre>	<pre> [16 1] [2 16] </pre>	Precision(0) = 0.89 Precision(1) = 0.94 recall(0) = 0.94 recall(1) = 0.89 f1-score(0) = 0.91 f1-score(1) = 0.91 accuracy = 0.91	0.915033
k-ближайших соседей $n = 3$	<pre> [0 1 1 0 0 1 1 1 0 0 1 1 0 0 1 0 1 1 0 1 0 1 1 0 0 1 1 1 0 1 1 0 0 1 0] [0 1 1 0 0 1 1 1 0 0 0 0 0 0 1 0 1 1 0 1 0 1 1 0 0 1 1 1 0 1 1 1 0 1 0] </pre>	<pre> [15 2] [1 17] </pre>	Precision(0) = 0.94 Precision(1) = 0.89 recall(0) = 0.88 recall(1) = 0.94	0.913399

			f1-score(0)= 0.91 f1-score(1)= 0.92 accuracy = 0 .91	
k- ближайших соседей n = 5	[01100111000100101101001001110110010] [01100111000000101101011001110111010]	[16 1] [2 16]	Precision(0) = 0.89 Precision(1) = 0.94 recall(0)= 0 .94 recall(1)= 0 .89 f1-score(0)= 0.91 f1-score(1)= 0.91 accuracy = 0 .91	0.915033
k- ближайших соседей n = 9	[01100111001000101101011001110110010] [01100111000000101101011001110111010]	[16 1] [1 17]	Precision(0) = 0.94 Precision(1) = 0.94 recall(0)= 0 .94 recall(1)= 1 f1-score(0)= 0.94 f1-score(1)= 0.94 accuracy = 0 .94	0.94281
Наивный байесовски й метод	[01100111001000101101011101110110010] [01100111000000101101011001110111010]	[15 2] [1 17]	Precision(0) = 0.94 Precision(1) = 0.89 recall(0)= 0 .88 recall(1)= 0 .94 f1-score(0)= 0.91 f1-score(1)= 0.92 accuracy = 0 .91	0.913399
Случайный лес n = 5	[01100111001000101101011001110110010] [01100111000000101101011001110111010]	[16 1] [1 17]	Precision(0) = 0.94 Precision(1) = 0.94 recall(0)= 0 .94 recall(1)= 0 .94 f1-score(0)= 0.94 f1-score(1)= 0.94	0.9428

			accuracy = 0 .94	
Случайный лес n = 10	[01100111001000101101011001110110010] [01100111000000101101011001110111010]	$\begin{bmatrix} 16 & 1 \\ 1 & 17 \end{bmatrix}$	Precision(0) = 0.94 Precision(1) = 0.94 recall(0)= 0 .94 recall(1)= 0 .94 f1-score(0)= 0.94 f1-score(1)= 0.94 accuracy = 0 .94	0.9428
Случайный лес n = 15	[01100111001000101101011001110110010] [01100111000000101101011001110111010]	$\begin{bmatrix} 16 & 1 \\ 1 & 17 \end{bmatrix}$	Precision(0) = 0.94 Precision(1) = 0.94 recall(0)= 0 .94 recall(1)= 0 .94 f1-score(0)= 0.94 f1-score(1)= 0.94 accuracy = 0 .94	0.94281
Случайный лес n = 20	[01100111001000101101011001110110010] [01100111000000101101011001110111010]	$\begin{bmatrix} 16 & 1 \\ 1 & 17 \end{bmatrix}$	Precision(0) = 0.94 Precision(1) = 0.94 recall(0)= 0 .94 recall(1)= 0 .94 f1-score(0)= 0.94 f1-score(1)= 0.94 accuracy = 0 .94	0.94281
Случайный лес n = 50	[01100111001000101101011001110110010] [01100111000000101101011001110111010]	$\begin{bmatrix} 16 & 1 \\ 1 & 17 \end{bmatrix}$	Precision(0) = 0.94 Precision(1) = 0.94 recall(0)= 0 .94 recall(1)= 0 .94 f1-score(0)= 0.94 f1-score(1)= 0.94 accuracy = 0 .94	0.94281

Наиболее подходящие методы для классификации данных k-ближайших соседей $n = 9$, случайный лес $n=5$, случайный лес $n=10$, случайный лес $n=15$, случайный лес $n=20$, случайный лес $n=50$, так как в этих методах наибольшая аккуратность.

Код программы

```
#!/usr/bin/env python
```

```
# coding: utf-8
```

```
# ## Импорт необходимых библиотек
```

```
import numpy as np
```

```
import matplotlib.pyplot as plt
```

```
from sklearn.datasets import make_blobs
```

```
from sklearn.metrics import confusion_matrix
```

```
from sklearn.metrics import classification_report
```

```
from sklearn.metrics import accuracy_score
```

```
from sklearn.metrics import roc_auc_score
```

```
from sklearn.model_selection import train_test_split
```

```
from sklearn.linear_model import LogisticRegression
```

```
from sklearn.neighbors import KNeighborsClassifier
```

```
from sklearn.ensemble import RandomForestClassifier
```

```
from sklearn.naive_bayes import GaussianNB
```

```
# ## Отображение на графике области принятия решения
```

```
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
# | Вид класса | random_state | cluster_std | noise | centers |
# |:-:|:-:|:-:|:-:|:-:|:-:|

```

```
# | blobs | 41 | 3 | - | 2 |
```

```
X, y = make_blobs(centers=2, random_state=41, cluster_std=3)
```

```
print("Координаты точек:\n", X[:15])
```

```
print("Метки класса: ", y[:15])
```

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

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.10,  
random_state=1)
```

```
# ### Обучающая выборка
```

```
plt.scatter(X_train[:,0], X_train[:,1], c=y_train)
```

```
# ### Тестовая выборка
```

```
plt.scatter(X_test[:,0], X_test[:,1], c=y_test)
```

```
# ## Обучение модели и классификация
```

```
def print_classification_metrics(classifier, X, y, prediction, y_test):
```

```
    print("Предсказанные и истинные значения")
```

```
    print(prediction)
```

```
    print(y_test)
```

```
    print("Матрица ошибок")
```

```
    print(confusion_matrix(y_test, prediction))
```

```
    print("Точность классификации: ", accuracy_score(prediction, y_test))
```

```

print("Значения полноты, точности, f1-меры и аккуратности")
print(classification_report(y_test, prediction))
print("Значение площади под кривой ошибок (AUC ROC)")
print(roc_auc_score(y_test, prediction))
print("Область принятия решений")
plt.xlabel("first feature")
plt.ylabel("second feature")
plot_2d_separator(knn, X, fill=True)
plt.scatter(X[:, 0], X[:, 1], c=y, s=70)
plt.show()

# #### Метод k-ближайших соседей

for i in [1, 3, 5, 9]:
    knn = KNeighborsClassifier(n_neighbors=i, metric='euclidean')
    knn.fit(X_train, y_train)
    prediction = knn.predict(X_test)
    print("n_neighbors = ", i)
    print_classification_metrics(knn, X, y, prediction, y_test)

# #### Наивный байесовский метод

naive = GaussianNB()
naive.fit(X_train, y_train)
predict = naive.predict(X_test)
print_classification_metrics(naive, X, y, predict, y_test)

# #### Случайный лес

for i in [5, 10, 15, 20, 50]:

```

```
rand_forest = RandomForestClassifier(n_estimators=i)
rand_forest.fit(X_train, y_train)
prediction = rand_forest.predict(X_test)
print("n_estimators = ", i)
print_classification_metrics(knn, X, y, prediction, y_test)
```

Вывод

В ходе выполнения данной лабораторной работы мы получили базовые навыки работы с языком python и набором функций для анализа и обработки данных. Получили практические навыки решения задачи бинарной классификации данных в среде Jupiter Notebook. Научились загружать данные, обучать классификаторы и проводить классификацию. Научились оценивать точность полученных моделей.

Контрольные вопросы

1) Постановка задачи классификации данных. Что такое бинарная классификация?

Задача классификации — задача, в которой имеется множество объектов (ситуаций), разделённых, некоторым образом, на классы. Задано конечное множество объектов, для которых известно, к каким классам они относятся. Это множество называется выборкой. Классовая принадлежность остальных объектов неизвестна. Требуется построить алгоритм, способный классифицировать (см. ниже) произвольный объект из исходного множества.

Классифицировать объект — значит, указать номер (или наименование) класса, к которому относится данный объект.

Классификация объекта — номер или наименование класса, выдаваемый алгоритмом классификации в результате его применения к данному конкретному объекту.

Бинарная классификация – это один из типов задач классификации в машинном обучении, когда мы должны классифицировать два взаимоисключающих класса. Например, классифицировать сообщения как спам или не спам, классифицировать новости как фальшивые или настоящие.

2) Общий алгоритм решения задачи классификации данных.

1. Конструирование модели: описание множества предопределённых классов.

- Каждый пример набора данных относится к одному предопределённому классу.

- На этом этапе используется обучающее множество, на нем происходит конструирование модели.

- Полученная модель представлена классификационными правилами, деревом решений или математической формулой.

2. Использование модели: классификация новых или неизвестных значений.

- Оценка правильности (точности) модели.

1. Известные значения из тестового примера сравниваются с результатами использования полученной модели.

2. Уровень точности - процент правильно классифицированных примеров в тестовом множестве.

3. Тестовое множество, т.е. множество, на котором тестируется построенная модель, не должно зависеть от обучающего множества.

- Если точность модели допустима, возможно использование модели для классификации новых примеров, класс которых неизвестен

3) Чем отличаются обучающая и тестовая выборки? Какие существуют способы формирования обучающей и тестовой выборок?

Обучающая выборка - это набор, который подается на вход модели в процессе обучения вместе с ответами, с целью научить модель видеть связь между этими признаками и правильным ответом

Тестовая выборка используется для проверки модели. Модель не получает целевой признак на вход и, более того, должна предсказать его величину используя значения остальных признаков. Эти предсказания потом сравниваются с реальными ответами.

Способы формирования выборок:

- метод удерживания
- метод k-кратной перекрёстной проверки
- скользящий экзамен
- стратификация
- самонастройка

4) Как рассчитываются значения полноты и точности классификации?

Точность (precision) и полнота (recall) являются метриками которые используются при оценке большей части алгоритмов извлечения информации. Иногда они используются сами по себе, иногда в качестве базиса для производных метрик, таких как F-мера или R-Precision.

TP — истинно-положительное решение;

TN — истинно-отрицательное решение;

FP — ложно-положительное решение;

FN — ложно-отрицательное решение.

Тогда, точность и полнота определяются следующим образом:

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

5) Как рассчитывается значение площади под кривой ошибок?

AUC-ROC (или ROC AUC) — площадь (Area Under Curve) под кривой ошибок (Receiver Operating Characteristic curve). Данная кривая представляет из себя линию от (0,0) до (1,1) в координатах True Positive Rate (TPR) и False Positive Rate (FPR):

$$TPR = \frac{TP}{TP + FN}$$

$$FPR = \frac{FP}{FP + TN}$$

6) Что показывает и как рассчитывается матрица ошибок?

На практике значения точности и полноты гораздо более удобней рассчитывать с использованием матрицы неточностей (confusion matrix). В случае если количество классов относительно невелико (не более 100-150 классов), этот подход позволяет довольно наглядно представить результаты работы классификатора.

Матрица неточностей — это матрица размера N на N, где N — это количество классов. Столбцы этой матрицы резервируются за экспертными решениями, а строки за решениями классификатора.

Матрица ошибок позволяет оценить эффективность прогноза не только в качественном, но и в количественном выражении

Это таблица с 4 различными комбинациями прогнозируемых и фактических значений.

		Actual Values	
		Positive (1)	Negative (0)
Predicted Values	Positive (1)	TP	FP
	Negative (0)	FN	TN

7) Алгоритм и особенности метода к-ближайших соседей.

Шаг 1 – Загружаем обучающий и тестовый dataset.

Шаг 2 – Выбираем значение К, то есть ближайшие точки данных. Оно может быть любым целым числом.

Шаг 3 - Вычисляем расстояние между тестовыми данными и каждой строкой обучающих данных с помощью любого из методов. Наиболее часто используемый метод вычисления расстояния - евклидов.

Шаг 4– Отсортировываем в порядке возрастания, основываясь на значении расстояния.

Шаг 5 – Алгоритм выбирает верхние К строк из отсортированного массива.

Шаг 6 – Назначаем класс контрольной точке на основе наиболее частого класса этих строк.

Особенности:

- Алгоритм прост и легко реализуем.
- Не чувствителен к выбросам.
- Нет необходимости строить модель, настраивать несколько параметров или делать дополнительные допущения.
- Алгоритм универсален. Его можно использовать для обоих типов задач: классификации и регрессии.

8) Алгоритм и особенности метода случайного леса.

Порядок действий в алгоритме

- Загрузите ваши данные.
- В заданном наборе данных определите случайную выборку.
- Далее алгоритм построит по выборке дерево решений.
- Дерево строится, пока в каждом листе не более n объектов, или пока не будет достигнута определенная высота.
- Затем будет получен результат прогнозирования из каждого дерева решений.

На этом этапе голосование будет проводиться для каждого прогнозируемого результата: мы выбираем лучший признак, делаем разбиение в дереве по нему и повторяем этот пункт до исчерпания выборки.

В конце выбирается результат прогноза с наибольшим количеством голосов. Это и есть окончательный результат прогнозирования.

Особенности:

- имеет высокую точность предсказания, на большинстве задач будет лучше линейных алгоритмов; точность сравнима с точностью бустинга
- практически не чувствителен к выбросам в данных из-за случайного сэмплирования
- не чувствителен к масштабированию значений признаков, связано с выбором случайных подпространств
- способен эффективно обрабатывать данные с большим числом признаков и классов
- одинаково хорошо обрабатывает как непрерывные, так и дискретные признаки
- редко переобучается, на практике добавление деревьев почти всегда только улучшает композицию, но на валидации, после достижения определенного количества деревьев, кривая обучения выходит на асимптоту
- для случайного леса существуют методы оценивания значимости отдельных признаков в модели
- хорошо работает с пропущенными данными; сохраняет хорошую точность, если большая часть данных пропущена

— предполагает возможность сбалансировать вес каждого класса на всей выборке, либо на подвыборке каждого дерева