

Задание 1.

Сформировать с помощью `sklearn.make_classification` датасет из 100 объектов с двумя признаками, обучить случайный лес из 1, 3, 10 и 50 деревьев и визуализировать их разделяющие гиперплоскости на графиках (по подобию визуализации деревьев из предыдущего урока, необходимо только заменить вызов функции `predict` на `tree_vote`). Сделать выводы о получаемой сложности гиперплоскости и недообучении или переобучении случайного леса в зависимости от количества деревьев в нем.

In [50]:

```

import matplotlib.pyplot as plt
import random

from matplotlib.colors import ListedColormap
from sklearn import datasets

import numpy as np

random.seed(42)

def get_bootstrap(data, labels, N):
    n_samples = data.shape[0]
    bootstrap = []

    for i in range(N):
        b_data = np.zeros(data.shape)
        b_labels = np.zeros(labels.shape)

        for j in range(n_samples):
            sample_index = random.randint(0, n_samples-1)
            b_data[j] = data[sample_index]
            b_labels[j] = labels[sample_index]
        bootstrap.append((b_data, b_labels))

    return bootstrap

def get_subsample(len_sample):
    # будем сохранять не сами признаки, а их индексы
    sample_indexes = [i for i in range(len_sample)]

    len_subsample = int(np.sqrt(len_sample))
    subsample = []

    random.shuffle(sample_indexes)
    for _ in range(len_subsample):
        subsample.append(sample_indexes.pop())

    return subsample

# Реализуем класс узла

class Node:

    def __init__(self, index, t, true_branch, false_branch):
        self.index = index # индекс признака, по которому ведется сравнение с порогом в этом
        self.t = t # значение порога
        self.true_branch = true_branch # поддерево, удовлетворяющее условию в узле
        self.false_branch = false_branch # поддерево, не удовлетворяющее условию в узле

# И класс терминального узла (листа)
class Leaf:

    def __init__(self, data, labels):
        self.data = data
        self.labels = labels
        self.prediction = self.predict()

    def predict(self):

```

```

# подсчет количества объектов разных классов
classes = {} # сформируем словарь "класс: количество объектов"
for label in self.labels:
    if label not in classes:
        classes[label] = 0
    classes[label] += 1
# найдем класс, количество объектов которого будет максимальным в этом листе и вернем
prediction = max(classes, key=classes.get)
return prediction

```

Расчет критерия Джини

```

def gini(labels):
    # подсчет количества объектов разных классов
    classes = {}
    for label in labels:
        if label not in classes:
            classes[label] = 0
        classes[label] += 1

    # расчет критерия
    impurity = 1
    for label in classes:
        p = classes[label] / len(labels)
        impurity -= p ** 2

    return impurity

```

Расчет качества

```

def quality(left_labels, right_labels, current_gini):

    # доля выбоки, ушедшая в левое поддерево
    p = float(left_labels.shape[0]) / (left_labels.shape[0] + right_labels.shape[0])

    return current_gini - p * gini(left_labels) - (1 - p) * gini(right_labels)

```

Разбиение датасета в узле

```

def split(data, labels, index, t):

    left = np.where(data[:, index] <= t)
    right = np.where(data[:, index] > t)

    true_data = data[left]
    false_data = data[right]
    true_labels = labels[left]
    false_labels = labels[right]

    return true_data, false_data, true_labels, false_labels

```

Нахождение наилучшего разбиения

```

def find_best_split(data, labels):

    # обозначим минимальное количество объектов в узле
    min_leaf = 1

    current_gini = gini(labels)

    best_quality = 0
    best_t = None
    best_index = None

    n_features = data.shape[1]

```

```

# выбор индекса из подвыборки длиной sqrt(n_features)
subsample = get_subsample(n_features)

for index in subsample:
    # будем проверять только уникальные значения признака, исключая повторения
    t_values = np.unique([row[index] for row in data])

    for t in t_values:
        true_data, false_data, true_labels, false_labels = split(data, labels, index, t)
        # пропускаем разбиения, в которых в узле остается менее 5 объектов
        if len(true_data) < min_leaf or len(false_data) < min_leaf:
            continue

        current_quality = quality(true_labels, false_labels, current_gini)

        # выбираем порог, на котором получается максимальный прирост качества
        if current_quality > best_quality:
            best_quality, best_t, best_index = current_quality, t, index

    return best_quality, best_t, best_index

# Построение дерева с помощью рекурсивной функции
def build_tree(data, labels):

    quality, t, index = find_best_split(data, labels)

    # Базовый случай - прекращаем рекурсию, когда нет прироста в качества
    if quality == 0:
        return Leaf(data, labels)

    true_data, false_data, true_labels, false_labels = split(data, labels, index, t)

    # Рекурсивно строим два поддерева
    true_branch = build_tree(true_data, true_labels)
    false_branch = build_tree(false_data, false_labels)

    # Возвращаем класс узла со всеми поддеревьями, то есть целого дерева
    return Node(index, t, true_branch, false_branch)

def random_forest(data, labels, n_trees):
    forest = []
    bootstrap = get_bootstrap(data, labels, n_trees)

    for b_data, b_labels in bootstrap:
        forest.append(build_tree(b_data, b_labels))

    return forest

# Функция классификации отдельного объекта
def classify_object(obj, node):

    # Останавливаем рекурсию, если достигли листа
    if isinstance(node, Leaf):
        answer = node.prediction
        return answer

    if obj[node.index] <= node.t:
        return classify_object(obj, node.true_branch)

```

```

    return classify_object(obj, node.false_branch)

# функция формирования предсказания по выборке на одном дереве
def predict(data, tree):

    classes = []
    for obj in data:
        prediction = classify_object(obj, tree)
        classes.append(prediction)
    return classes

```

In [40]:

```

# сгенерируем данные
X, y = datasets.make_classification(n_samples=100, n_features=2, n_informative=2,
                                   n_classes=2, n_redundant=0, n_clusters_per_class=2)

# визуализируем сгенерированные данные

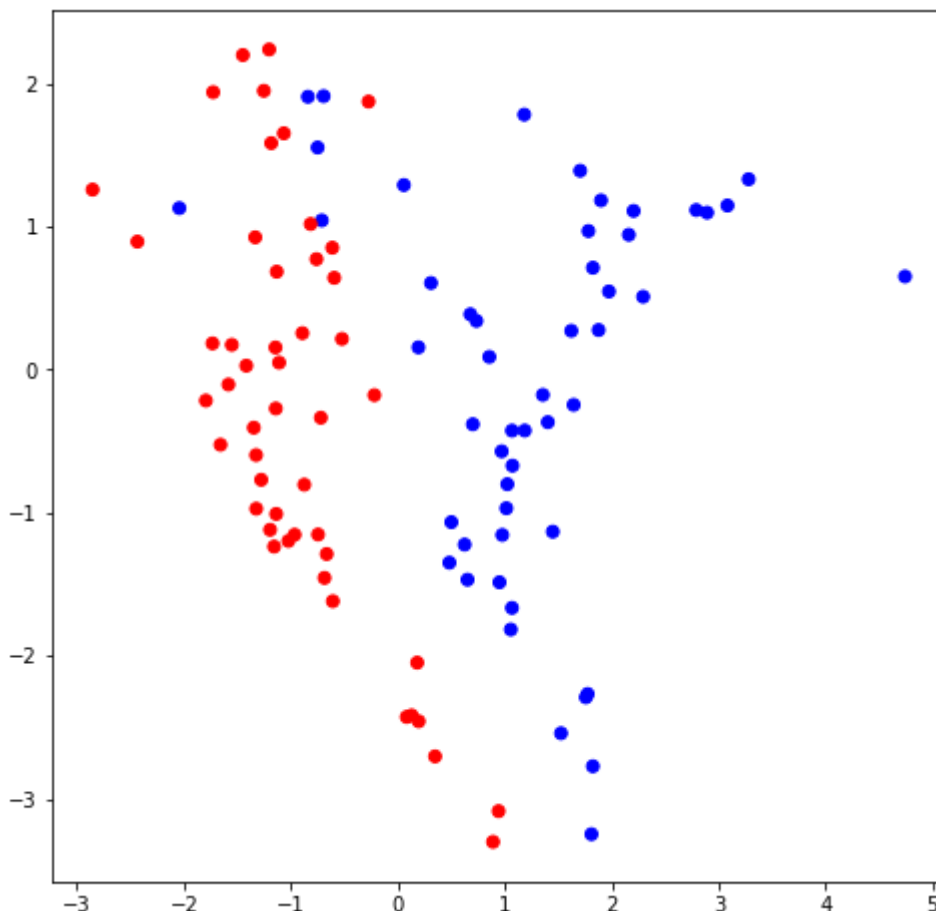
colors = ListedColormap(['red', 'blue'])
light_colors = ListedColormap(['lightcoral', 'lightblue'])

plt.figure(figsize=(8, 8))
plt.scatter(list(map(lambda x: x[0], X)), list(map(lambda x: x[1], X)),
            c=y, cmap=colors)

```

Out[40]:

<matplotlib.collections.PathCollection at 0x7ffe659de280>



In [42]:

```
# предсказание голосованием деревьев

def tree_vote(forest, data):

    # добавим предсказания всех деревьев в список
    predictions = []
    for tree in forest:
        predictions.append(predict(data, tree))

    # сформируем список с предсказаниями для каждого объекта
    predictions_per_object = list(zip(*predictions))

    # выберем в качестве итогового предсказания для каждого объекта то,
    # за которое проголосовало большинство деревьев
    voted_predictions = []
    for obj in predictions_per_object:
        voted_predictions.append(max(set(obj), key=obj.count))

    return voted_predictions
```

In [43]:

```
# Разобьем выборку на обучающую и тестовую
from sklearn import model_selection

X_train, X_test, y_train, y_test = model_selection.train_test_split(classification_c
```

In [44]:

```
# Введем функцию подсчета точности как доли правильных ответов

def accuracy_metric(actual, predicted):
    correct = 0
    for i in range(len(actual)):
        if actual[i] == predicted[i]:
            correct += 1
    return correct / float(len(actual)) * 100.0
```

In [53]:

визуализируем модель в зависимости от кол-ва деревьев

```
def get_meshgrid(data, step=.05, border=1.2):
    x_min, x_max = data[:, 0].min() - border, data[:, 0].max() + border
    y_min, y_max = data[:, 1].min() - border, data[:, 1].max() + border
    return np.meshgrid(np.arange(x_min, x_max, step), np.arange(y_min, y_max, step))

tree_counts = [1, 3, 10, 50]
fig, axes = plt.subplots(len(tree_counts), 2, figsize=(16, 28))

for i, n_trees in enumerate(tree_counts):

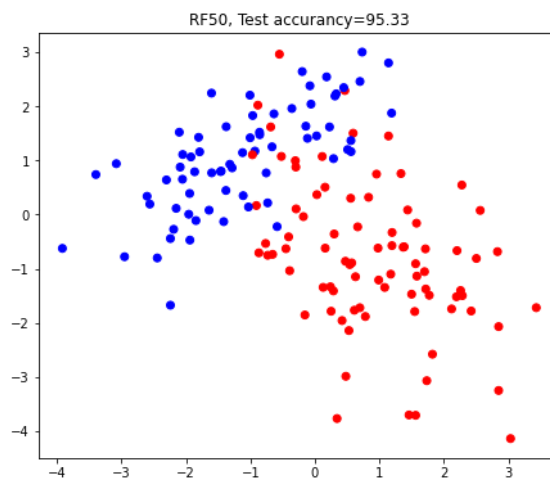
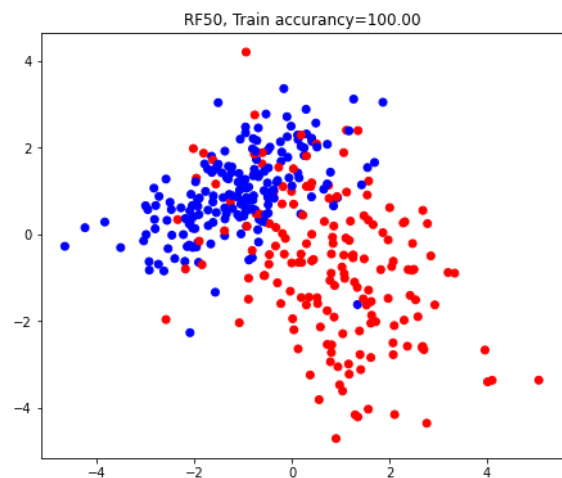
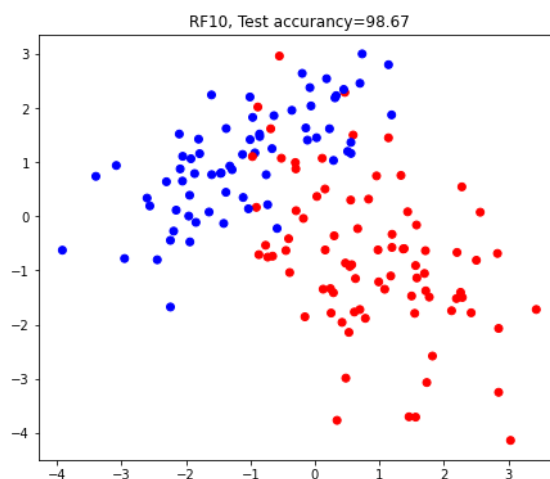
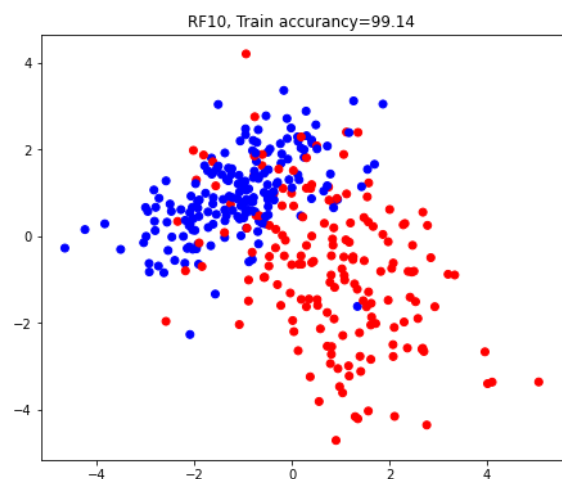
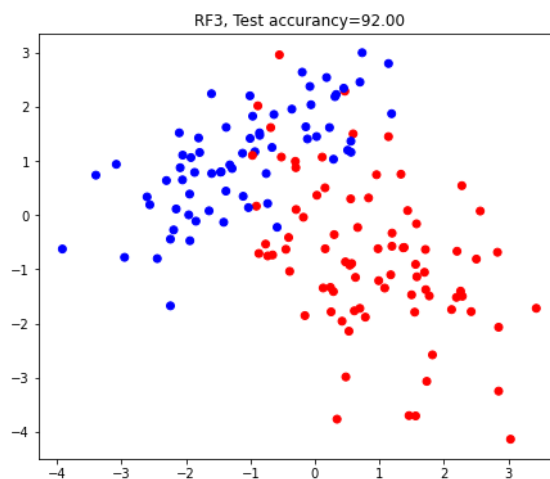
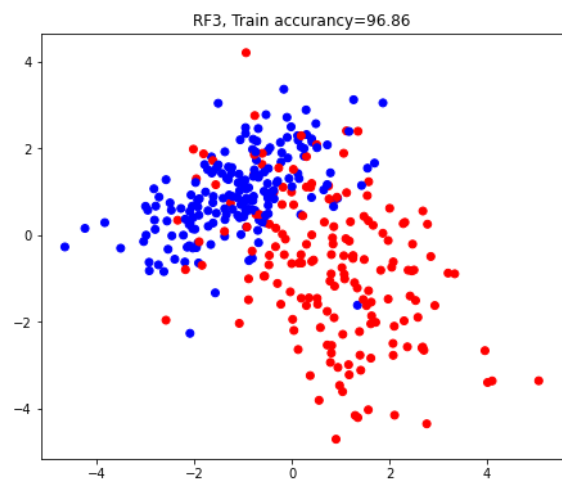
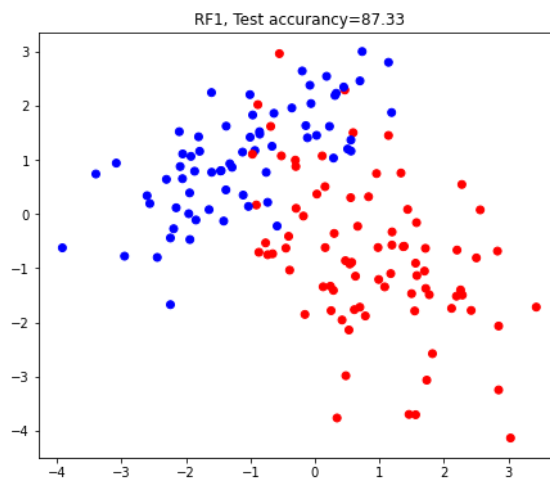
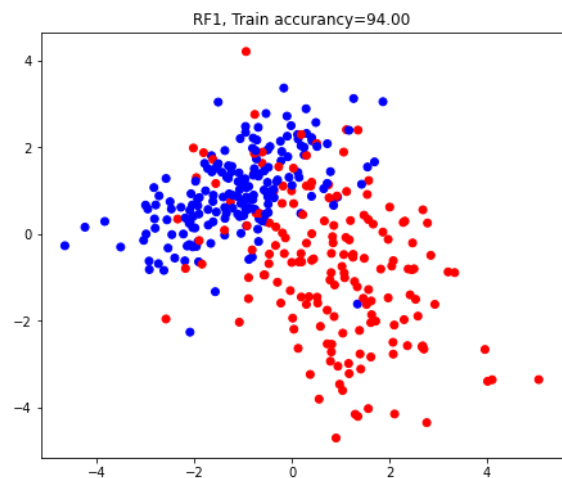
    my_forest = random_forest(X_train, y_train, n_trees)

    train_preds = tree_vote(my_forest, X_train)
    test_preds = tree_vote(my_forest, X_test)

    axes[i][0].scatter(X_train[:, 0], X_train[:, 1], c=y_train, cmap=colors)
    axes[i][0].set_title(f'RF{n_trees}, Train accuracy={accuracy_metric(y_train, tr

    axes[i][1].scatter(X_test[:, 0], X_test[:, 1], c=y_test, cmap=colors)
    axes[i][1].set_title(f'RF{n_trees}, Test accuracy={accuracy_metric(y_test, test

plt.show()
```



присутствует переобучение.

In []: