

Домашнее задание 5

In [1]:

```
1 import matplotlib.pyplot as plt
2 import random
3
4 from matplotlib.colors import ListedColormap
5 from sklearn import datasets
6
7 import numpy as np
```

1. Задание:

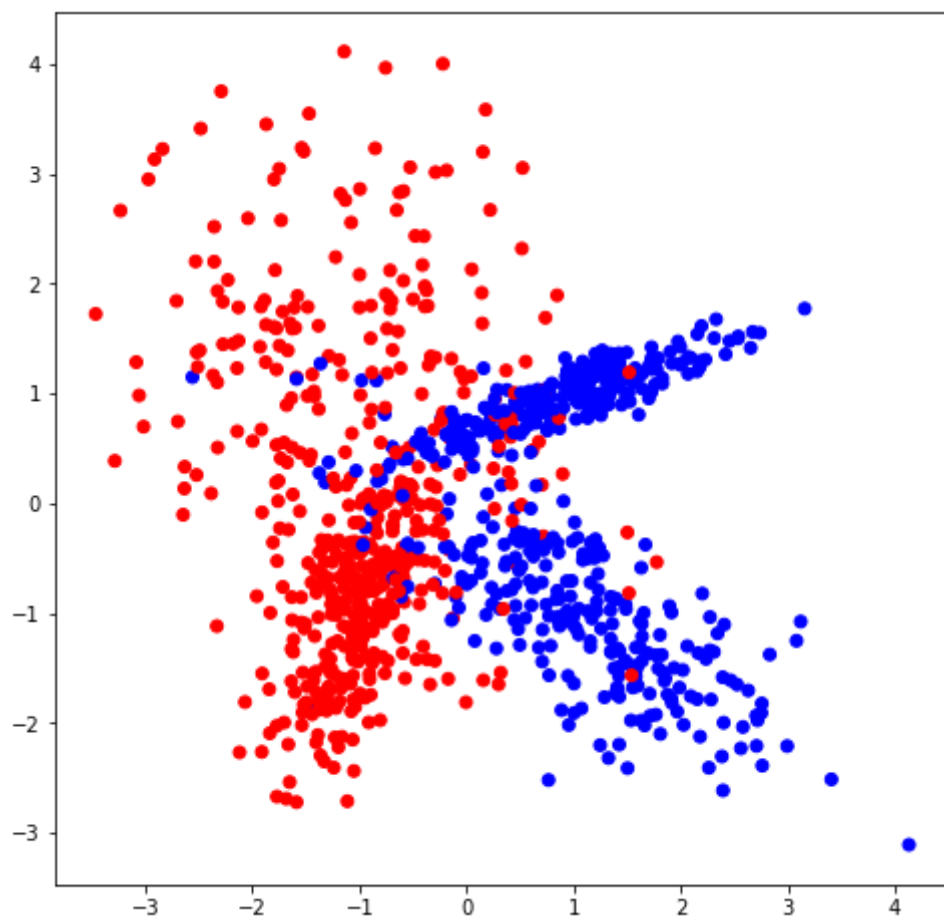
- Сформировать с помощью `sklearn.make_classification` датасет из 1000 объектов с двумя признаками.
- Обучить случайный лес из 1, 3, 10, 50, 100 и 200 деревьев (за основу взять реализацию построения этого алгоритма из урока).
- Визуализировать их разделяющие линии на графиках (по подобию визуализации 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 [2]:

```
1 # сгенерируем данные, представляющие собой 500 объектов с 5-ю признаками
2 classification_data, classification_labels = datasets.make_classification(n_samples=1000,
3
4
5
6
                                     n_features=2, n_informative=2,
                                     n_classes=2, n_redundant=0,
                                     n_clusters_per_class=2,
                                     random_state=6)
```

In [3]:

```
1 # визуализируем сгенерированные данные
2
3 colors = ListedColormap(['red', 'blue'])
4 light_colors = ListedColormap(['lightcoral', 'lightblue'])
5
6 plt.figure(figsize=(8,8))
7 plt.scatter(list(map(lambda x: x[0], classification_data)), list(map(lambda x: x[1], c
8                               c=classification_labels, cmap=colors);
9
10
```



In [4]:

```
1 random.seed(42)
2
3 def get_bootstrap(data, labels, N):
4     n_samples = data.shape[0]
5     bootstrap = []
6
7     for i in range(N):
8         b_data = np.zeros(data.shape)
9         b_labels = np.zeros(labels.shape)
10
11         for j in range(n_samples):
12             sample_index = random.randint(0, n_samples-1)
13             b_data[j] = data[sample_index]
14             b_labels[j] = labels[sample_index]
15         bootstrap.append((b_data, b_labels))
16
17     return bootstrap
```

In [5]:

```
1 def get_subsample(len_sample):
2     # будем сохранять не сами признаки, а их индексы
3     sample_indexes = [i for i in range(len_sample)]
4
5     len_subsample = int(np.sqrt(len_sample))
6     subsample = []
7
8     random.shuffle(sample_indexes)
9     for _ in range(len_subsample):
10         subsample.append(sample_indexes.pop())
11
12     return subsample
```

Далее повторим реализацию построения дерева решений из предыдущего урока с некоторыми изменениями

In [6]:

```
1 # Реализуем класс узла
2
3 class Node:
4
5     def __init__(self, index, t, true_branch, false_branch):
6         self.index = index # индекс признака, по которому ведется сравнение с порогом
7         self.t = t # значение порога
8         self.true_branch = true_branch # поддерево, удовлетворяющее условию в узле
9         self.false_branch = false_branch # поддерево, не удовлетворяющее условию в узле
```

In [7]:

```
1  # И класс терминального узла (листа)
2
3  class Leaf:
4
5      def __init__(self, data, labels):
6          self.data = data
7          self.labels = labels
8          self.prediction = self.predict()
9
10     def predict(self):
11         # подсчет количества объектов разных классов
12         classes = {} # сформируем словарь "класс: количество объектов"
13         for label in self.labels:
14             if label not in classes:
15                 classes[label] = 0
16             classes[label] += 1
17         # найдем класс, количество объектов которого будет максимальным в этом листе
18         prediction = max(classes, key=classes.get)
19         return prediction
```

In [8]:

```
1  # Расчет критерия Джини
2
3  def gini(labels):
4      # подсчет количества объектов разных классов
5      classes = {}
6      for label in labels:
7          if label not in classes:
8              classes[label] = 0
9              classes[label] += 1
10
11     # расчет критерия
12     impurity = 1
13     for label in classes:
14         p = classes[label] / len(labels)
15         impurity -= p ** 2
16
17     return impurity
```

In [9]:

```
1  # Расчет качества
2
3  def quality(left_labels, right_labels, current_gini):
4
5      # доля выбоки, ушедшая в левое поддерево
6      p = float(left_labels.shape[0]) / (left_labels.shape[0] + right_labels.shape[0])
7
8      return current_gini - p * gini(left_labels) - (1 - p) * gini(right_labels)
```

In [10]:

```
1  # Разбиение датасета в узле
2
3  def split(data, labels, index, t):
4
5      left = np.where(data[:, index] <= t)
6      right = np.where(data[:, index] > t)
7
8      true_data = data[left]
9      false_data = data[right]
10     true_labels = labels[left]
11     false_labels = labels[right]
12
13     return true_data, false_data, true_labels, false_labels
```

In [11]:

```
1  # Нахождение наилучшего разбиения
2
3  def find_best_split(data, labels):
4
5      # обозначим минимальное количество объектов в узле
6      min_leaf = 1
7
8      current_gini = gini(labels)
9
10     best_quality = 0
11     best_t = None
12     best_index = None
13
14     n_features = data.shape[1]
15
16     # выбор индекса из подвыборки длиной sqrt(n_features)
17     subsample = get_subsample(n_features)
18
19     for index in subsample:
20         # будем проверять только уникальные значения признака, исключая повторения
21         t_values = np.unique([row[index] for row in data])
22
23         for t in t_values:
24             true_data, false_data, true_labels, false_labels = split(data, labels, index, t)
25             # пропускаем разбиения, в которых в узле остается менее 1 объекта
26             if len(true_data) < min_leaf or len(false_data) < min_leaf:
27                 continue
28
29             current_quality = quality(true_labels, false_labels, current_gini)
30
31             # выбираем порог, на котором получается максимальный прирост качества
32             if current_quality > best_quality:
33                 best_quality, best_t, best_index = current_quality, t, index
34
35     return best_quality, best_t, best_index
```

In [12]:

```
1  # Построение дерева с помощью рекурсивной функции
2
3  def build_tree(data, labels):
4
5      quality, t, index = find_best_split(data, labels)
6
7      # Базовый случай - прекращаем рекурсию, когда нет прироста в качества
8      if quality == 0:
9          return Leaf(data, labels)
10
11     true_data, false_data, true_labels, false_labels = split(data, labels, index, t)
12
13     # Рекурсивно строим два поддеревя
14     true_branch = build_tree(true_data, true_labels)
15     false_branch = build_tree(false_data, false_labels)
16
17     # Возвращаем класс узла со всеми поддеревьями, то есть целого дерева
18     return Node(index, t, true_branch, false_branch)
```

Теперь добавим функцию формирования случайного леса.

In [13]:

```
1  def random_forest(data, labels, n_trees):
2      forest = []
3      bootstrap = get_bootstrap(data, labels, n_trees)
4
5      for b_data, b_labels in bootstrap:
6          forest.append(build_tree(b_data, b_labels))
7
8      return forest
```

In [14]:

```
1  # Функция классификации отдельного объекта
2
3  def classify_object(obj, node):
4
5      # Останавливаем рекурсию, если достигли листа
6      if isinstance(node, Leaf):
7          answer = node.prediction
8          return answer
9
10     if obj[node.index] <= node.t:
11         return classify_object(obj, node.true_branch)
12     else:
13         return classify_object(obj, node.false_branch)
```

In [15]:

```
1 # функция формирования предсказания по выборке на одном дереве
2
3 def predict(data, tree):
4
5     classes = []
6     for obj in data:
7         prediction = classify_object(obj, tree)
8         classes.append(prediction)
9     return classes
```

In [16]:

```
1 # предсказание голосованием деревьев
2
3 def tree_vote(forest, data):
4
5     # добавим предсказания всех деревьев в список
6     predictions = []
7     for tree in forest:
8         predictions.append(predict(data, tree))
9
10    # сформируем список с предсказаниями для каждого объекта
11    predictions_per_object = list(zip(*predictions))
12
13    # выберем в качестве итогового предсказания для каждого объекта то,
14    # за которое проголосовало большинство деревьев
15    voted_predictions = []
16    for obj in predictions_per_object:
17        voted_predictions.append(max(set(obj), key=obj.count))
18
19    return voted_predictions
```

Далее мы сделаем обычное разбиение выборки на обучающую и тестовую, как это делалось ранее.

In [17]:

```
1 # Разобьем выборку на обучающую и тестовую
2
3 from sklearn import model_selection
4
5 train_data, test_data, train_labels, test_labels = model_selection.train_test_split(cla
6
7
8
```

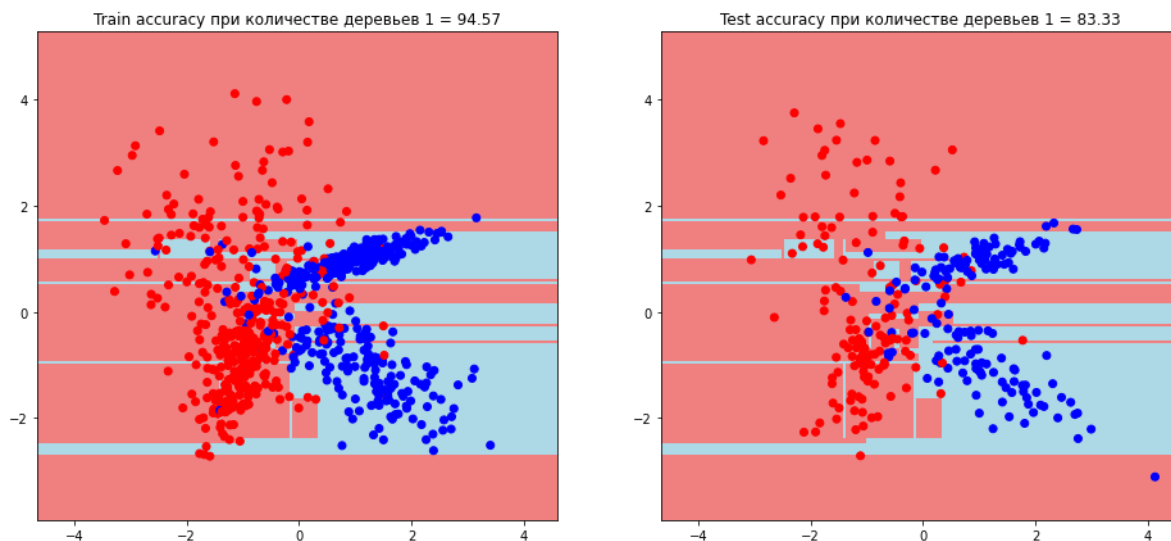
In [18]:

```
1 # Введем функцию подсчета точности как доли правильных ответов
2 def accuracy_metric(actual, predicted):
3     correct = 0
4     for i in range(len(actual)):
5         if actual[i] == predicted[i]:
6             correct += 1
7     return correct / float(len(actual)) * 100.0
```

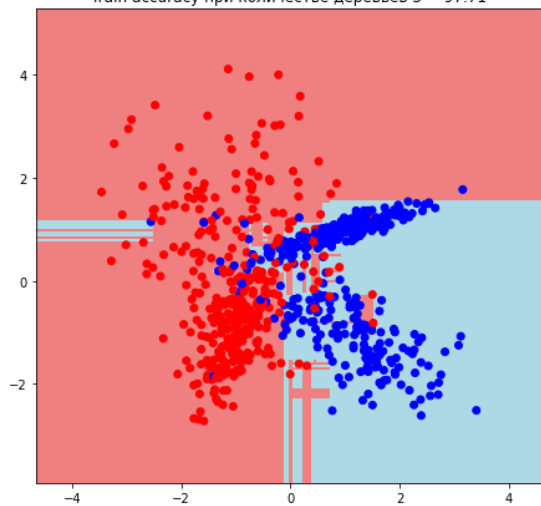
Теперь построим несколько случайных лесов с разным количеством деревьев в них.

In [19]:

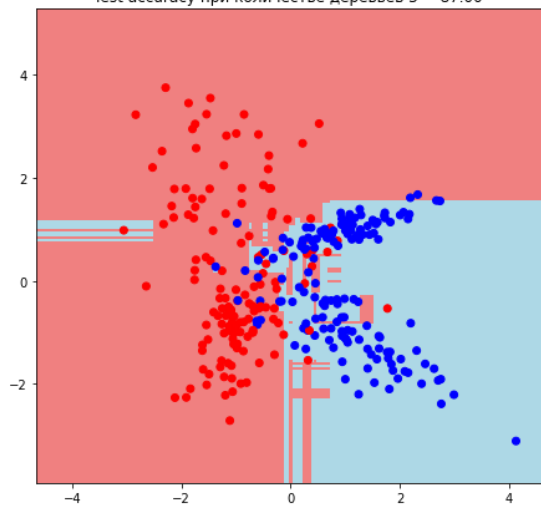
```
1 def get_meshgrid(data, step=.05, border=1.2):
2     x_min, x_max = data[:, 0].min() - border, data[:, 0].max() + border
3     y_min, y_max = data[:, 1].min() - border, data[:, 1].max() + border
4     return np.meshgrid(np.arange(x_min, x_max, step), np.arange(y_min, y_max, step))
5
6 for i in [1, 3, 10, 50, 100, 200]:
7     my_forest = random_forest(train_data, train_labels, i)
8     train_answers = tree_vote(my_forest, train_data)
9     test_answers = tree_vote(my_forest, test_data)
10    train_accuracy = accuracy_metric(train_labels, train_answers)
11    test_accuracy = accuracy_metric(test_labels, test_answers)
12    plt.figure(figsize = (16, 7))
13
14    # график обучающей выборки
15    plt.subplot(1,2,1)
16    xx, yy = get_meshgrid(train_data)
17    mesh_predictions = np.array(tree_vote(my_forest, np.c_[xx.ravel(), yy.ravel()])).ravel()
18    plt.pcolormesh(xx, yy, mesh_predictions, cmap = light_colors)
19    plt.scatter(train_data[:, 0], train_data[:, 1], c = train_labels, cmap = colors)
20    plt.title(f'Train accuracy при количестве деревьев {i} = {train_accuracy:.2f}')
21
22    # график тестовой выборки
23    plt.subplot(1,2,2)
24    plt.pcolormesh(xx, yy, mesh_predictions, cmap = light_colors)
25    plt.scatter(test_data[:, 0], test_data[:, 1], c = test_labels, cmap = colors)
26    plt.title(f'Test accuracy при количестве деревьев {i} = {test_accuracy:.2f}')
27
```



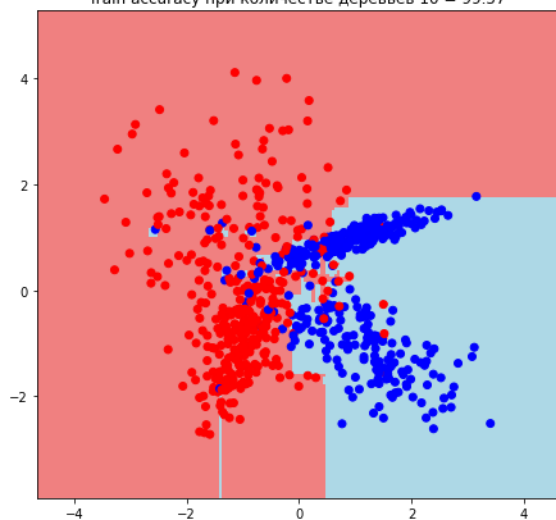
Train accuracy при количестве деревьев 3 = 97.71



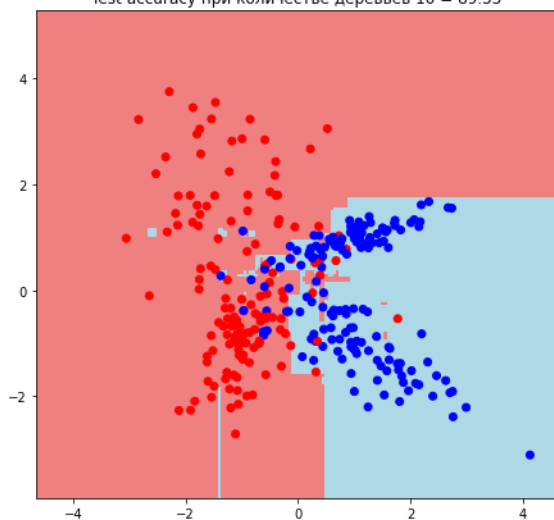
Test accuracy при количестве деревьев 3 = 87.00



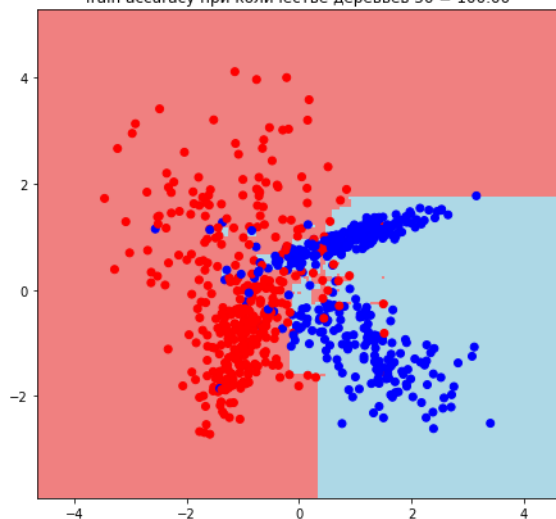
Train accuracy при количестве деревьев 10 = 99.57



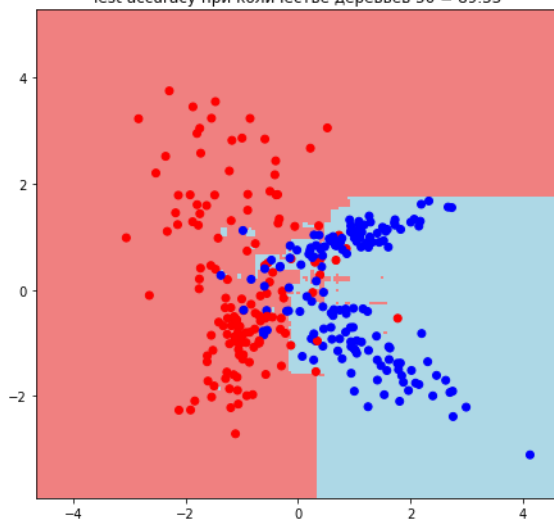
Test accuracy при количестве деревьев 10 = 89.33

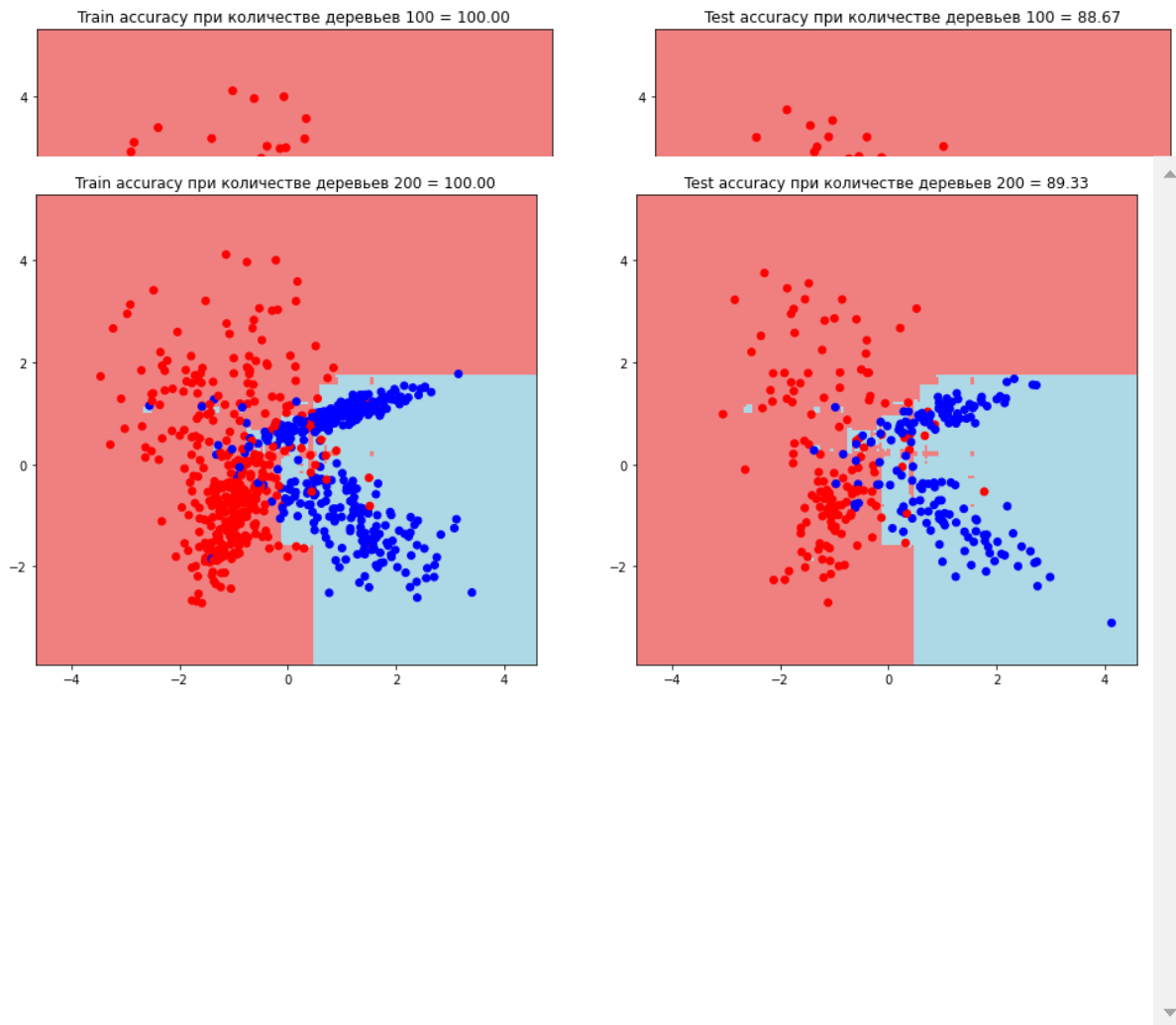


Train accuracy при количестве деревьев 50 = 100.00



Test accuracy при количестве деревьев 50 = 89.33





Выводы:

1. Из графиков видно, что получаемая сложность гиперплоскости при увеличении количества деревьев становится менее рваной и граница между классами более четкой.
2. Начиная с 50 деревьев точность на обучающей выборке становится 100. На 100 деревьях точность тестовой выборки снижается, намекая на переобучение, но уже на 200 деревьях точность тестовой выборки становится максимальной. Что еще раз подтверждает факт того, что случайный лес очень устойчив к переобучению.

2. (*) Заменить в реализованном алгоритме проверку с помощью отложенной выборки на Out-of-Bag.

Для этого необходимо внести небольшие изменения в код

In [20]:

```
1 random.seed(42)
2
3 def get_bootstrap(data, labels, N):
4     n_samples = data.shape[0]
5     bootstrap = []
6     bootstrap_oobe = []
7     sample_index_com = []
8
9     for i in range(N):
10         b_data = np.zeros(data.shape)
11         b_labels = np.zeros(labels.shape)
12
13         for j in range(n_samples):
14             sample_index = random.randint(0, n_samples-1)
15             # Запоминаем выбранные индексы
16             sample_index_com.append(sample_index)
17             b_data[j] = data[sample_index]
18             b_labels[j] = labels[sample_index]
19             # Формируем список индексов, которые не участвовали в формировании бутстрап выборки
20             oobe_index = list(set([i for i in range(100)]) - set(sample_index_com))
21             # Формируем на основе этих индексов выборки объектов и классов, которые не попали в bootstrap
22             b_data_oobe = data[oobe_index]
23             b_labels_oobe = labels[oobe_index]
24             bootstrap.append((b_data, b_labels))
25             # Формируем наборы oobe выборок, в соответствии с количеством деревьев
26             bootstrap_oobe.append((b_data_oobe, b_labels_oobe))
27
28     return bootstrap, bootstrap_oobe
```

In [21]:

```
1 def random_forest(data, labels, n_trees):
2     forest = []
3     bootstrap, bootstrap_oobe = get_bootstrap(data, labels, n_trees)
4
5     for b_data, b_labels in bootstrap:
6         forest.append(build_tree(b_data, b_labels))
7     # Допишем возвращаемые параметры
8     return forest, bootstrap, bootstrap_oobe
```

Построим лес из одного дерева

In [22]:

```
1 n_trees = 1
2 my_forest_1, bootstrap, bootstrap_oobe = random_forest(classification_data, classification_labels, n_trees)
```

In [23]:

```
1 # Получим ответы для обучающей выборки
2 train_answers_1 = tree_vote(my_forest_1, classification_data)
```

In [24]:

```
1 # Точность на обучающей выборке
2 train_accuracy_1 = accuracy_metric(classification_labels, train_answers_1)
3 print(f'Точность случайного леса из {n_trees} деревьев на обучающей выборке: {train_acc
```

Точность случайного леса из 1 деревьев на обучающей выборке: 93.700

In [25]:

```
1 # Получим ответы для oobe выборки
2 oobe_answers_1 = tree_vote(my_forest_1, bootstrap_oobe[0][0])
```

In [26]:

```
1 # Точность на oobe выборке
2 oobe_accuracy_1 = accuracy_metric(bootstrap_oobe[0][1], oobe_answers_1)
3 print(f'Точность случайного леса из {n_trees} деревьев на oobe выборке: {oobe_accuracy_
```

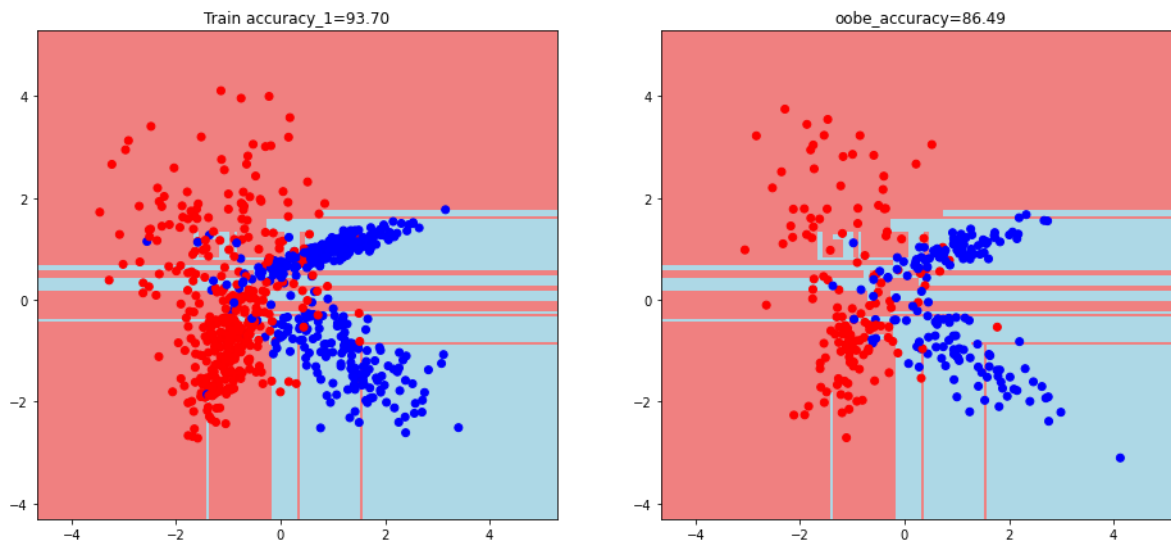
Точность случайного леса из 1 деревьев на oobe выборке: 86.486

In [27]:

```
1 # Визуализируем дерево на графике
2
3 def get_meshgrid(data, step=.05, border=1.2):
4     x_min, x_max = data[:, 0].min() - border, data[:, 0].max() + border
5     y_min, y_max = data[:, 1].min() - border, data[:, 1].max() + border
6     return np.meshgrid(np.arange(x_min, x_max, step), np.arange(y_min, y_max, step))
7
8 plt.figure(figsize = (16, 7))
9
10 # график обучающей выборки на 1 дереве
11 plt.subplot(1,2,1)
12 xx, yy = get_meshgrid(classification_data)
13 mesh_predictions = np.array(tree_vote(my_forest_1, np.c_[xx.ravel(), yy.ravel()])).reshape(xx.shape)
14 plt.pcolormesh(xx, yy, mesh_predictions, cmap = light_colors)
15 plt.scatter(train_data[:, 0], train_data[:, 1], c = train_labels, cmap = colors)
16 plt.title(f'Train accuracy_1={train_accuracy_1:.2f}')
17
18 # график тестовой выборки 1 дереве
19 plt.subplot(1,2,2)
20 plt.pcolormesh(xx, yy, mesh_predictions, cmap = light_colors)
21 plt.scatter(test_data[:, 0], test_data[:, 1], c = test_labels, cmap = colors)
22 plt.title(f'oobe_accuracy={oobe_accuracy_1:.2f}')
```

Out[27]:

Text(0.5, 1.0, 'oobe_accuracy=86.49')



Построим лес из трех деревьев

In [28]:

```
1 n_trees = 3
2 my_forest_3, bootstrap, bootstrap_oobe = random_forest(classification_data, classification_labels, n_trees)
```

In [29]:

```
1 # Получим ответы для обучающей выборки
2 train_answers_3 = tree_vote(my_forest_3, classification_data)
```

In [30]:

```
1 # Точность на обучающей выборке
2 train_accuracy_3 = accuracy_metric(classification_labels, train_answers_3)
3 print(f'Точность случайного леса из {n_trees} деревьев на обучающей выборке: {train_acc
```

Точность случайного леса из 3 деревьев на обучающей выборке: 96.800

In [31]:

```
1 # Получим ответы для oobe выборки
2 oobe_answers_3 = tree_vote(my_forest_3, bootstrap_oobe[0][0])
```

In [32]:

```
1 # Точность на oobe выборке
2 oobe_accuracy_3 = accuracy_metric(bootstrap_oobe[0][1], oobe_answers_3)
3 print(f'Точность случайного леса из {n_trees} деревьев на oobe выборке: {oobe_accuracy_
```

Точность случайного леса из 3 деревьев на oobe выборке: 92.500

In [33]:

```
1  # Визуализируем дерево на графике
2
3  def get_meshgrid(data, step=.05, border=1.2):
4      x_min, x_max = data[:, 0].min() - border, data[:, 0].max() + border
5      y_min, y_max = data[:, 1].min() - border, data[:, 1].max() + border
6      return np.meshgrid(np.arange(x_min, x_max, step), np.arange(y_min, y_max, step))
7
8  plt.figure(figsize = (16, 7))
9
10 # график обучающей выборки на 1 дереве
11 plt.subplot(1,2,1)
12 xx, yy = get_meshgrid(classification_data)
13 mesh_predictions = np.array(tree_vote(my_forest_3, np.c_[xx.ravel(), yy.ravel()])).reshape(xx.shape)
14 plt.pcolormesh(xx, yy, mesh_predictions, cmap = light_colors)
15 plt.scatter(train_data[:, 0], train_data[:, 1], c = train_labels, cmap = colors)
16 plt.title(f'Train accuracy_3={train_accuracy_3:.2f}')
17
18 # график тестовой выборки 1 дереве
19 plt.subplot(1,2,2)
20 plt.pcolormesh(xx, yy, mesh_predictions, cmap = light_colors)
21 plt.scatter(test_data[:, 0], test_data[:, 1], c = test_labels, cmap = colors)
22 plt.title(f'oobe_accuracy_3={oobe_accuracy_3:.2f}')
```

Out[33]:

Text(0.5, 1.0, 'oobe_accuracy_3=92.50')

