

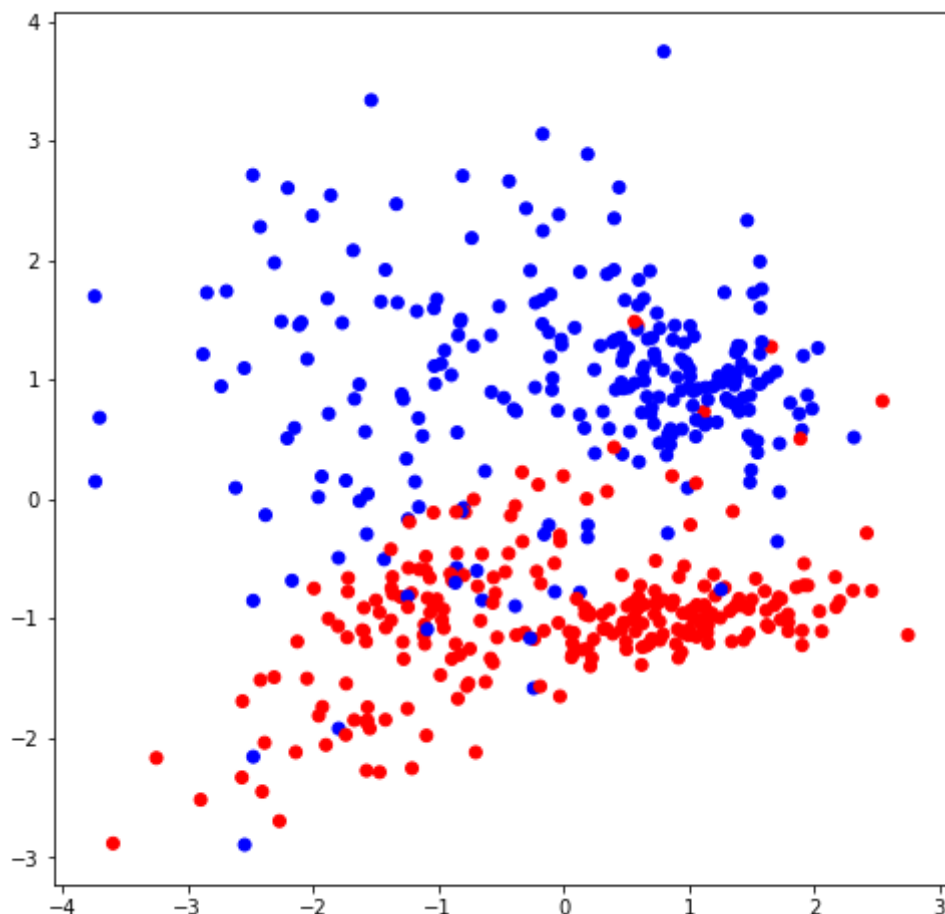
Домашнее задание

1. В коде из методички реализуйте один или несколько критериев останова (количество листьев, количество используемых признаков, глубина дерева и т.д.).

Сгенерируем данные

In [2]:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 # Сгенерируем данные
5 from sklearn import datasets
6 from matplotlib.colors import ListedColormap
7
8 # сгенерируем данные с помощью sklearn.datasets
9 classes = datasets.make_classification(n_samples=500, n_features=2, n_informative=2,
10                                     n_redundant=0, n_classes=2, random_state=6)
11 # datasets.make_blobs(centers = 5, cluster_std = 1, random_state=1)
12
13 # и изобразим их на графике
14 colors = ListedColormap(['red', 'blue'])
15 light_colors = ListedColormap(['lightcoral', 'lightblue'])
16
17 plt.figure(figsize=(8, 8))
18 plt.scatter([x[0] for x in classes[0]], [x[1] for x in classes[0]], c=classes[1], cmap=
```



Далее разделим выборку на обучающую и тестовую.

In [3]:

```
1 # перемешивание датасета
2 np.random.seed(41)
3 shuffle_index = np.random.permutation(classes[0].shape[0])
4 X_shuffled, y_shuffled = classes[0][shuffle_index], classes[1][shuffle_index]
5
6 # разбишка на обучающую и тестовую выборки
7 train_proportion = 0.7
8 train_test_cut = int(len(classes[0]) * train_proportion)
9
10 train_data, test_data, train_labels, test_labels = \
11     X_shuffled[:train_test_cut], \
12     X_shuffled[train_test_cut:], \
13     y_shuffled[:train_test_cut], \
14     y_shuffled[train_test_cut:]
```

In [4]:

```
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 [5]:

```
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
```

Индекс Джини:

$$H(X) = 1 - \sum_{k=1}^K p_k^2.$$

In [6]:

```
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 # коэффициент неопределенности Джини
13     for label in classes:
14         p = classes[label] / len(labels)
15         impurity -= p ** 2
16
17     return impurity
```

Функционал качества:

$$Q = H(X_m) - \frac{|X_l|}{|X_m|} H(X_l) - \frac{|X_r|}{|X_m|} H(X_r)$$

In [7]:

```
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 [8]:

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

In [9]:

```
1  # Нахождение наилучшего разбиения
2
3  def find_best_split(data, labels):
4
5      # обозначим минимальное количество объектов в узле
6      min_leaf = 5
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     for index in range(n_features):
17         # будем проверять только уникальные значения признака, исключая повторения
18         t_values = np.unique([row[index] for row in data])
19
20         for t in t_values:
21             true_data, false_data, true_labels, false_labels = split(data, labels, index, t)
22
23             # пропускаем разбиения, в которых в узле остается менее 5 объектов
24             if len(true_data) < min_leaf or len(false_data) < min_leaf:
25                 continue
26
27             current_quality = quality(true_labels, false_labels, current_gini)
28
29             # выбираем порог, на котором получается максимальный прирост качества
30             if current_quality > best_quality:
31                 best_quality, best_t, best_index = current_quality, t, index
32
33     return best_quality, best_t, best_index
```

In [10]:

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

In [11]:

```
1 def predict(data, tree):
2
3     classes = []
4     for obj in data:
5         prediction = classify_object(obj, tree)
6         classes.append(prediction)
7     return classes
```

Ограничим глубину дерева

In [13]:

```
1 d = int(input("Введите максимальную глубину дерева: "))
2
3 # Построение дерева с помощью рекурсивной функции
4 def build_tree(data, labels):
5     global depth, true_branch, false_branch
6     print("Глубина", depth)
7     quality, t, index = find_best_split(data, labels)
8
9     # Базовый случай - прекращаем рекурсию, когда нет прироста в качества
10    if quality == 0:
11        return Leaf(data, labels)
12        print("quality == 0")
13
14    if depth == d:
15        return Leaf(data, labels)
16
17    # Рекурсивно строим два поддеревья
18    print("Делаем ветвление на глубине ", depth)
19    depth += 1
20    true_data, false_data, true_labels, false_labels = split(data, labels, index, t)
21    true_branch = build_tree(true_data, true_labels)
22    false_branch = build_tree(false_data, false_labels)
23
24    # Возвращаем класс узла со всеми поддеревьями, то есть целого дерева
25    return Node(index, t, true_branch, false_branch)
```

Введите максимальную глубину дерева: 5

In [14]:

```
1 # Построим дерево по обучающей выборке
2 depth = 0
3 true_branch = None
4 false_branch = None
5 my_tree = build_tree(train_data, train_labels)
```

Глубина 0

Делаем ветвление на глубине 0

Глубина 1

Делаем ветвление на глубине 1

Глубина 2

Делаем ветвление на глубине 2

Глубина 3

Делаем ветвление на глубине 3

Глубина 4

Делаем ветвление на глубине 4

Глубина 5

Глубина 5

Глубина 5

Глубина 5

Глубина 5

Глубина 5

In [15]:

```
1 # Напечатаем ход нашего дерева
2 def print_tree(node, spacing=""):
3
4     # Если лист, то выводим его прогноз
5     if isinstance(node, Leaf):
6         print(spacing + "Прогноз:", node.prediction)
7         return
8
9     # Выведем значение индекса и порога на этом узле
10    print(spacing + 'Индекс', str(node.index))
11    print(spacing + 'Порог', str(node.t))
12
13    # Рекурсионный вызов функции на положительном поддереве
14    print (spacing + '--> True:')
15    print_tree(node.true_branch, spacing + " ")
16
17    # Рекурсионный вызов функции на отрицательном поддереве
18    print (spacing + '--> False:')
19    print_tree(node.false_branch, spacing + " ")
20
21    print_tree(my_tree)
```

Индекс 1

Порог -0.107726691940364

--> True:

Индекс 1

Порог -0.3549877301133244

--> True:

Индекс 0

Порог -0.23882214445115557

--> True:

Индекс 0

Порог -0.4381064214473833

--> True:

Индекс 0

Порог -1.793845443769914

--> True:

Прогноз: 0

--> False:

Прогноз: 0

--> False:

Прогноз: 1

--> False:

Прогноз: 0

--> False:

Прогноз: 0

--> False:

Прогноз: 1

In [16]:

```
1 # Получим ответы для обучающей выборки
2 train_answers = predict(train_data, my_tree)
3
4 # Получим ответы для тестовой выборки
5 test_answers = predict(test_data, my_tree)
```

In [17]:

```
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
```

In [18]:

```
1 # Точность на обучающей выборке
2 train_accuracy = accuracy_metric(train_labels, train_answers)
3 train_accuracy
```

Out[18]:

91.71428571428571

In [19]:

```
1 # Точность на тестовой выборке
2 test_accuracy = accuracy_metric(test_labels, test_answers)
3 test_accuracy
```

Out[19]:

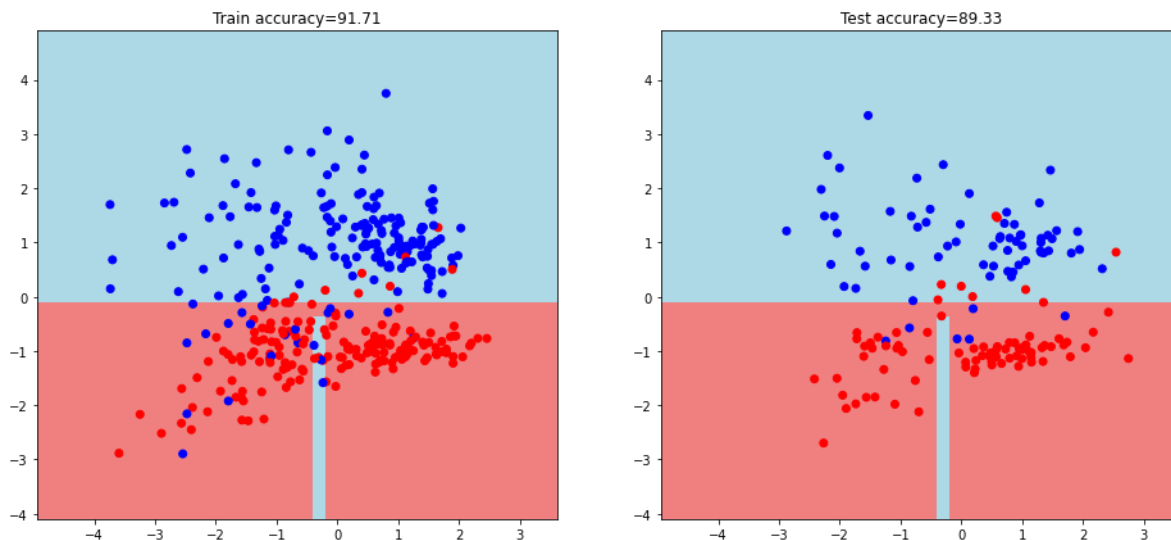
89.33333333333333

In [20]:

```
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
11 plt.subplot(1,2,1)
12 xx, yy = get_meshgrid(train_data)
13 mesh_predictions = np.array(predict(np.c_[xx.ravel(), yy.ravel()], my_tree)).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={train_accuracy:.2f}')
17
18 # график тестовой выборки
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'Test accuracy={test_accuracy:.2f}')
```

Out[20]:

Text(0.5, 1.0, 'Test accuracy=89.33')



2. Для задачи классификации обучите дерево решений с использованием критериев разбиения Джини и Энтропия. Сравните качество классификации, сделайте выводы.

In [21]:

```
1 # Расчет критерия энтропия
2
3 def entropy(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 # коэффициент неопределенности энтропии
13     for label in classes:
14         p = classes[label] / len(labels)
15         impurity -= np.sum(p * np.log2(p))
16
17     return impurity
```

Функционал качества:

$$Q = H(X_m) - \frac{|X_l|}{|X_m|} H(X_l) - \frac{|X_r|}{|X_m|} H(X_r)$$

In [22]:

```
1 # Расчет качества
2
3 def quality(left_labels, right_labels, current_entropy):
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_entropy - p * entropy(left_labels) - (1 - p) * entropy(right_labels)
```

In [23]:

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

In [24]:

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

In [25]:

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

In [26]:

```
1  def predict(data, tree):
2
3      classes = []
4      for obj in data:
5          prediction = classify_object(obj, tree)
6          classes.append(prediction)
7      return classes
```

Органичим глубину дерева

In [27]:

```
1 d = int(input("Введите максимальную глубину дерева: "))
2
3 # Построение дерева с помощью рекурсивной функции
4 def build_tree(data, labels):
5     global depth, true_branch, false_branch
6     print("Глубина", depth)
7     quality, t, index = find_best_split(data, labels)
8
9     # Базовый случай - прекращаем рекурсию, когда нет прироста в качества
10    if quality == 0:
11        return Leaf(data, labels)
12        print("quality == 0")
13
14    if depth == d:
15        return Leaf(data, labels)
16
17    # Рекурсивно строим два поддерева
18    print("Делаем ветвление на глубине ", depth)
19    depth += 1
20    true_data, false_data, true_labels, false_labels = split(data, labels, index, t)
21    true_branch = build_tree(true_data, true_labels)
22    false_branch = build_tree(false_data, false_labels)
23
24    # Возвращаем класс узла со всеми поддеревьями, то есть целого дерева
25    return Node(index, t, true_branch, false_branch)
```

Введите максимальную глубину дерева: 5

In [28]:

```
1 # Построим дерево по обучающей выборке
2 depth = 0
3 true_branch = None
4 false_branch = None
5 my_tree = build_tree(train_data, train_labels)
```

Глубина 0
Делаем ветвление на глубине 0
Глубина 1
Делаем ветвление на глубине 1
Глубина 2
Делаем ветвление на глубине 2
Глубина 3
Делаем ветвление на глубине 3
Глубина 4
Делаем ветвление на глубине 4
Глубина 5
Глубина 5
Глубина 5
Глубина 5
Глубина 5
Глубина 5

In [29]:

```
1 # Напечатаем ход нашего дерева
2 def print_tree(node, spacing=""):
3
4     # Если лист, то выводим его прогноз
5     if isinstance(node, Leaf):
6         print(spacing + "Прогноз:", node.prediction)
7         return
8
9     # Выведем значение индекса и порога на этом узле
10    print(spacing + 'Индекс', str(node.index))
11    print(spacing + 'Порог', str(node.t))
12
13    # Рекурсионный вызов функции на положительном поддереве
14    print (spacing + '--> True:')
15    print_tree(node.true_branch, spacing + " ")
16
17    # Рекурсионный вызов функции на отрицательном поддереве
18    print (spacing + '--> False:')
19    print_tree(node.false_branch, spacing + " ")
20
21    print_tree(my_tree)
```

Индекс 1

Порог -0.107726691940364

--> True:

Индекс 1

Порог -0.5175343065438042

--> True:

Индекс 0

Порог -0.23882214445115557

--> True:

Индекс 0

Порог -0.46469463939970623

--> True:

Индекс 0

Порог -2.165676399540467

--> True:

Прогноз: 0

--> False:

Прогноз: 0

--> False:

Прогноз: 1

--> False:

Прогноз: 0

--> False:

Прогноз: 0

--> False:

Прогноз: 1

In [30]:

```
1 # Получим ответы для обучающей выборки
2 train_answers = predict(train_data, my_tree)
3
4 # Получим ответы для тестовой выборки
5 test_answers = predict(test_data, my_tree)
```

In [31]:

```
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
```

In [32]:

```
1 # Точность на обучающей выборке
2 train_accuracy = accuracy_metric(train_labels, train_answers)
3 train_accuracy
```

Out[32]:

91.71428571428571

In [33]:

```
1 # Точность на тестовой выборке
2 test_accuracy = accuracy_metric(test_labels, test_answers)
3 test_accuracy
```

Out[33]:

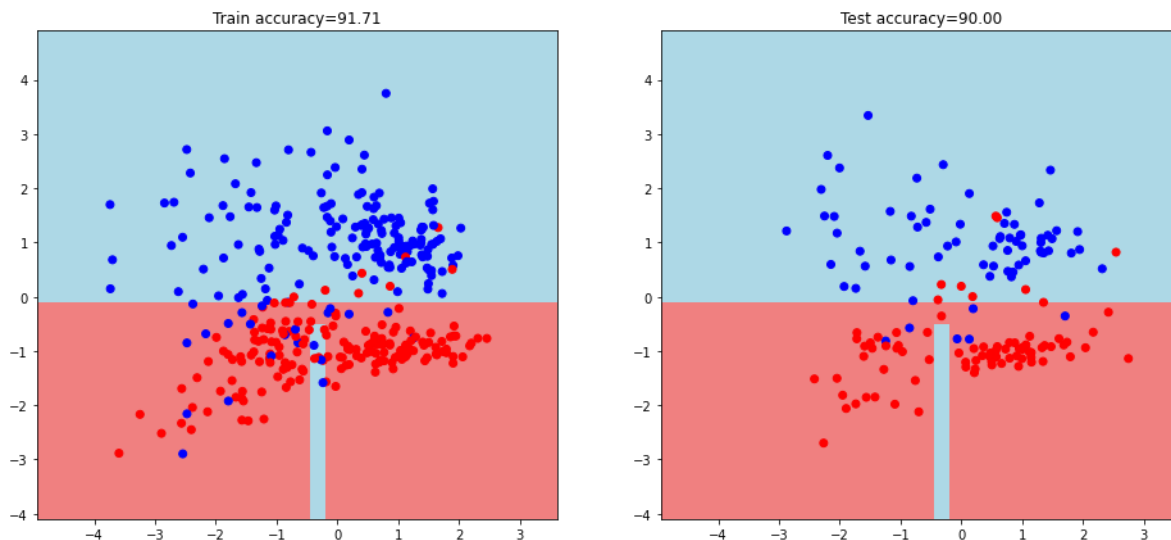
90.0

In [34]:

```
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
11 plt.subplot(1,2,1)
12 xx, yy = get_meshgrid(train_data)
13 mesh_predictions = np.array(predict(np.c_[xx.ravel(), yy.ravel()], my_tree)).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={train_accuracy:.2f}')
17
18 # график тестовой выборки
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'Test accuracy={test_accuracy:.2f}')
```

Out[34]:

Text(0.5, 1.0, 'Test accuracy=90.00')



Вывод: Видим, что качество работы моделей на этих данных идентичное

3. *Реализуйте дерево для задачи регрессии. Возьмите за основу дерево, реализованное в методичке, заменив механизм предсказания в листе на взятие среднего значения по выборке, а критерий Джини на дисперсию значений.

Сгенерируем датасет для регрессии

In [35]:

```
1 from sklearn.datasets import load_boston
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 import numpy as np
6 from sklearn.model_selection import learning_curve
7 from sklearn.metrics import make_scorer
8 %matplotlib inline
9
10 np.random.seed(42)
11
12 boston_data = load_boston()
13 boston_df = pd.DataFrame(boston_data.data, columns=boston_data.feature_names)
14 boston_df_2 = pd.DataFrame(boston_data.data, columns=boston_data.feature_names)
15 # boston_df = np.array(boston_df[['RM', 'LSTAT']])
```

C:\ProgramData\Anaconda3\lib\site-packages\statsmodels\tools_testing.py:19:
FutureWarning: pandas.util.testing is deprecated. Use the functions in the public API at pandas.testing instead.
import pandas.util.testing as tm

In [36]:

```
1 boston_df
```

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT
0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1.0	296.0	15.3	396.90	4.98
1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2.0	242.0	17.8	396.90	9.14
2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2.0	242.0	17.8	392.83	4.03
3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3.0	222.0	18.7	394.63	2.94
4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3.0	222.0	18.7	396.90	5.33
...
501	0.06263	0.0	11.93	0.0	0.573	6.593	69.1	2.4786	1.0	273.0	21.0	391.99	9.67
502	0.04527	0.0	11.93	0.0	0.573	6.120	76.7	2.2875	1.0	273.0	21.0	396.90	9.08
503	0.06076	0.0	11.93	0.0	0.573	6.976	91.0	2.1675	1.0	273.0	21.0	396.90	5.64
504	0.10959	0.0	11.93	0.0	0.573	6.794	89.3	2.3889	1.0	273.0	21.0	393.45	6.48
505	0.04741	0.0	11.93	0.0	0.573	6.030	80.8	2.5050	1.0	273.0	21.0	396.90	7.88

Переведем таблицу в матрицу

In [37]:

```
1 boston_df = np.array(boston_df)
2 target = boston_data.target
```


In [38]:

```
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        prediction = np.mean(self.labels)
12        return prediction
```

In [39]:

```
1 #перемешаем данные
2 np.random.seed(41)
3 shuffle_index = np.random.permutation(boston_df.shape[0])
4 X_shuffled, y_shuffled = boston_df[shuffle_index], np.array(target[shuffle_index])
5
```

In [40]:

```
1 # разбишка на обучающую и тестовую выборки
2 train_proportion = 0.7
3 train_test_cut = int(len(classes[0]) * train_proportion)
4
5 train_data, test_data, train_labels, test_labels = \
6     X_shuffled[:train_test_cut], \
7     X_shuffled[train_test_cut:], \
8     y_shuffled[:train_test_cut], \
9     y_shuffled[train_test_cut:]
```

В случае регрессии разброс будет характеризоваться дисперсией, поэтому критерий информативности будет записан в виде

$$H(X) = \frac{1}{X} \sum_{i \in X} (y_i - \bar{y}(X))^2,$$

где $\bar{y}(X)$ - среднее значение ответа в выборке X :

$$\bar{y}(X) = \frac{1}{|X|} \sum_{i \in X} y_i.$$

In [41]:

```
1 # Расчет дисперсии
2 def variance(labels):
3
4     impurity = (np.sum())/labels
5     return impurity
```

In [42]:

```
1 # Расчет дисперсии
2
3 def variance(labels):
4
5     impurity = np.var(labels)
6     return impurity
```

In [43]:

```
1 # Расчет качества
2
3 def quality(left_labels, right_labels, current_variance):
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_variance - p * variance(left_labels) - (1 - p) * variance(right_labels)
```

In [44]:

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

In [45]:

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

In [46]:

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

In [47]:

```
1  def predict(data, tree):
2
3      classes = []
4      for obj in data:
5          prediction = classify_object(obj, tree)
6          classes.append(prediction)
7      return classes
```

Органичим глубину дерева

In [48]:

```
1 d = int(input("Введите максимальную глубину дерева: "))
2
3 # Построение дерева с помощью рекурсивной функции
4 def build_tree(data, labels):
5     global depth, true_branch, false_branch
6     print("Глубина", depth)
7     quality, t, index = find_best_split(data, labels)
8
9     # Базовый случай - прекращаем рекурсию, когда нет прироста в качества
10    if quality == 0:
11        return Leaf(data, labels)
12        print("quality == 0")
13
14    if depth == d:
15        return Leaf(data, labels)
16
17    # Рекурсивно строим два поддеревя
18    print("Делаем ветвление на глубине ", depth)
19    depth += 1
20    true_data, false_data, true_labels, false_labels = split(data, labels, index, t)
21    true_branch = build_tree(true_data, true_labels)
22    false_branch = build_tree(false_data, false_labels)
23
24    # Возвращаем класс узла со всеми поддеревьями, то есть целого дерева
25    return Node(index, t, true_branch, false_branch)
```

Введите максимальную глубину дерева: 5

In [49]:

```
1 # Построим дерево по обучающей выборке
2 depth = 0
3 true_branch = None
4 false_branch = None
5 my_tree = build_tree(train_data, train_labels)
```

Глубина 0

Делаем ветвление на глубине 0

Глубина 1

Делаем ветвление на глубине 1

Глубина 2

Делаем ветвление на глубине 2

Глубина 3

Глубина 3

Делаем ветвление на глубине 3

Глубина 4

Делаем ветвление на глубине 4

Глубина 5

Глубина 5

Глубина 5

Глубина 5

Глубина 5

In [50]:

```
1  # Напечатаем ход нашего дерева
2  def print_tree(node, spacing=""):
3
4      # Если лист, то выводим его прогноз
5      if isinstance(node, Leaf):
6          print(spacing + "Прогноз:", node.prediction)
7          return
8
9      # Выведем значение индекса и порога на этом узле
10     print(spacing + 'Индекс', str(node.index))
11     print(spacing + 'Порог', str(node.t))
12
13     # Рекурсионный вызов функции на положительном поддереве
14     print (spacing + '--> True:')
15     print_tree(node.true_branch, spacing + " ")
16
17     # Рекурсионный вызов функции на отрицательном поддереве
18     print (spacing + '--> False:')
19     print_tree(node.false_branch, spacing + " ")
20
21 print_tree(my_tree)
```

Индекс 12

Порог 9.53

--> True:

Индекс 5

Порог 7.42

--> True:

Индекс 7

Порог 1.7573

--> True:

Индекс 5

Порог 6.538

--> True:

Прогноз: 22.87068965517241

--> False:

Прогноз: 27.54

--> False:

Индекс 5

Порог 6.794

--> True:

Индекс 5

Порог 6.538

--> True:

Прогноз: 22.87068965517241

--> False:

Прогноз: 27.54

--> False:

Прогноз: 32.768571428571434

--> False:

Прогноз: 44.915

--> False:

Прогноз: 17.74455445544554

In [51]:

```
1 # Получим ответы для обучающей выборки
2 train_answers = predict(train_data, my_tree)
3
4 # Получим ответы для тестовой выборки
5 test_answers = predict(test_data, my_tree)
```

In [52]:

```
1 from sklearn.metrics import mean_absolute_error
```

In [53]:

```
1 # Точность на обучающей выборке
2 mean_absolute_error(train_labels, train_answers)
```

Out[53]:

3.6628201238843094

In [54]:

```
1 # Точность на тестовой выборке
2 mean_absolute_error(test_labels, test_answers)
```

Out[54]:

3.911579368622728