

Домашнее задание №6

1. Для реализованной в методичке модели градиентного бустинга построить графики зависимости ошибки от количества деревьев в ансамбле и от максимальной глубины деревьев. Сделать выводы о зависимости ошибки от этих параметров.

In [1]:

```
1 from sklearn.tree import DecisionTreeRegressor
2 from sklearn import model_selection
3 import numpy as np
4 import random
5 from sklearn.datasets import load_diabetes
6 import matplotlib.pyplot as plt
7 X, y = load_diabetes(return_X_y=True)
```

Разделим выборку на обучающую и тестовую в соотношении 75/25.

In [2]:

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

В качестве функционала ошибки будем использовать среднеквадратичную ошибку. Реализуем соответствующую функцию.

In [3]:

```
1 def mean_squared_error(y_real, prediction):
2     return (sum((y_real - prediction) ** 2)) / len(y_real)
```

Используем L_2 loss $L(y, z) = (y - z)^2$, ее производная по z примет вид $L'(y, z) = 2(z - y)$. Реализуем ее также в виде функции.

In [4]:

```
1 def bias(z, y):
2     return 2 * (z - y)
```

Напишем функцию, реализующую предсказание в градиентном бустинге.

In [5]:

```
1 def gb_predict(X, trees_list, coef_list, eta):
2     # Реализуемый алгоритм градиентного бустинга будет инициализироваться нулевыми значениями
3     # поэтому все деревья из списка trees_list уже являются дополнительными и при предсказании
4     # прибавляются с шагом eta
5     return np.array([sum([eta * coef * alg.predict([x])[0] for alg, coef in zip(trees_list, coef_list)]) for x in X])
```

Реализуем функцию обучения градиентного бустинга.

In [6]:

```
1 def gb_fit(n_trees, max_depth, X_train, X_test, y_train, y_test, coefs, eta):
2
3     # eta - скорость обучения
4     # Деревья будем записывать в список
5     trees = []
6
7     # Будем записывать ошибки на обучающей и тестовой выборке на каждой итерации в список
8     train_errors = []
9     test_errors = []
10
11     for i in range(n_trees):
12         tree = DecisionTreeRegressor(max_depth=max_depth, random_state=42)
13
14         # инициализируем бустинг начальным алгоритмом, возвращающим ноль,
15         # поэтому первый алгоритм просто обучаем на выборке и добавляем в список
16         if len(trees) == 0:
17             # обучаем первое дерево на обучающей выборке
18             tree.fit(X_train, y_train)
19
20             train_errors.append(mean_squared_error(y_train, gb_predict(X_train, trees,
21             test_errors.append(mean_squared_error(y_test, gb_predict(X_test, trees, coefs, eta)
22         else:
23             # Получим ответы на текущей композиции
24             target = gb_predict(X_train, trees, coefs, eta)
25
26             # алгоритмы начиная со второго обучаем на сдвиг
27             tree.fit(X_train, bias(y_train, target))
28
29             train_errors.append(mean_squared_error(y_train, gb_predict(X_train, trees,
30             test_errors.append(mean_squared_error(y_test, gb_predict(X_test, trees, coefs, eta)
31
32     trees.append(tree)
33
34     return trees, train_errors, test_errors
```

Обучим несколько моделей с разным количеством деревьев и посмотрим как оно влияет на ошибку

In [7]:

```
1 max_depth = 3
2 eta = 0.05
3 list_mse = []
4 for n_trees in range(1, 81, 10):
5     coefs = [1] * n_trees
6     trees, train_errors, test_errors = gb_fit(n_trees, max_depth, X_train, X_test, y_train, y_test, coefs, eta)
7     train_prediction = gb_predict(X_train, trees, coefs, eta)
8     test_prediction = gb_predict(X_test, trees, coefs, eta)
9     round(mean_squared_error(y_test, test_prediction))
10    list_mse.append([n_trees, round(mean_squared_error(y_train, train_prediction)), round(mean_squared_error(y_test, test_prediction))])
11
```

In [8]:

```
1 list_mse
```

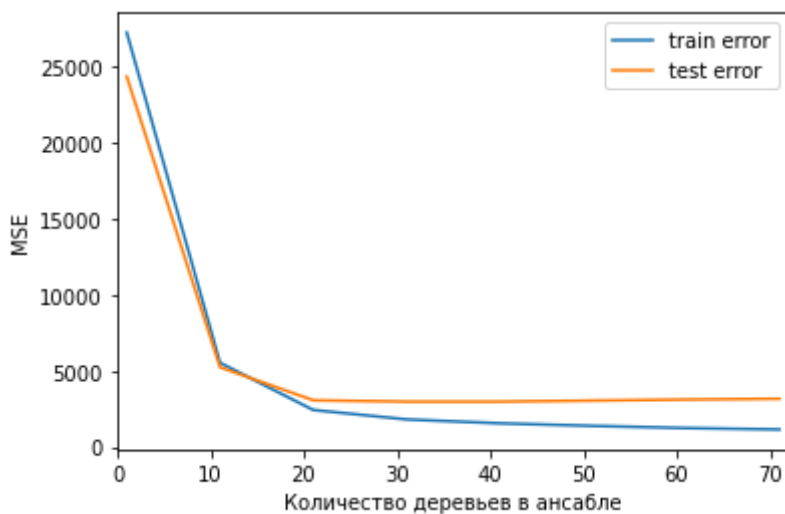
Out[8]:

```
[[1, 27244.0, 24347.0],  
 [11, 5543.0, 5247.0],  
 [21, 2435.0, 3074.0],  
 [31, 1815.0, 2982.0],  
 [41, 1558.0, 2982.0],  
 [51, 1389.0, 3053.0],  
 [61, 1244.0, 3119.0],  
 [71, 1154.0, 3172.0]]
```

Построим графики зависимости ошибки на обучающей и тестовой выборках от количества деревьев в ансамбле.

In [9]:

```
1 plt.xlabel('Количество деревьев в ансамбле')  
2 plt.ylabel('MSE')  
3 plt.xlim(0, list(zip(*list_mse))[0][-1] + 1)  
4 plt.plot(list(zip(*list_mse))[0], list(zip(*list_mse))[1], label='train error')  
5 plt.plot(list(zip(*list_mse))[0], list(zip(*list_mse))[2], label='test error')  
6 plt.legend(loc='best')  
7 plt.show()
```



Теперь обучим несколько моделей с разной глубиной деревьев и посмотрим как глубина влияет на ошибку

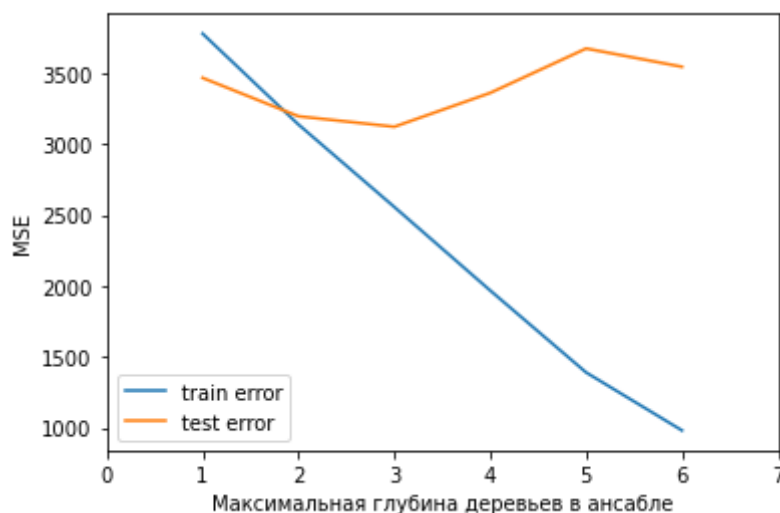
In [10]:

```
1 n_trees = 20
2 eta = 0.05
3 list_mse = []
4 for max_depth in range(1, 7):
5     coefs = [1] * n_trees
6     trees, train_errors, test_errors = gb_fit(n_trees, max_depth, X_train, X_test, y_train, y_test, eta)
7     train_prediction = gb_predict(X_train, trees, coefs, eta)
8     test_prediction = gb_predict(X_test, trees, coefs, eta)
9     round(mean_squared_error(y_test, test_prediction))
10    list_mse.append([max_depth, round(mean_squared_error(y_train, train_prediction)), round(mean_squared_error(y_test, test_prediction))])
11
```

Построим графики зависимости ошибки на обучающей и тестовой выборках от максимальной глубины деревьев в ансамбле.

In [11]:

```
1 plt.xlabel('Максимальная глубина деревьев в ансамбле')
2 plt.ylabel('MSE')
3 plt.xlim(0, list(zip(*list_mse))[0][-1] + 1)
4 plt.plot(list(zip(*list_mse))[0], list(zip(*list_mse))[1], label='train error')
5 plt.plot(list(zip(*list_mse))[0], list(zip(*list_mse))[2], label='test error')
6 plt.legend(loc='best')
7 plt.show()
```



Выводы:

1. На графиках отчетливо видно, что на обучающей выборке при увеличении количества деревьев и максимальной глубины деревьев ошибка постоянно уменьшается
2. На тестовой выборке до определенного момента при увеличении количества деревьев и максимальной глубины деревьев ошибка также уменьшается
3. После некоторых значений данных гиперпараметров ошибка на тестовой выборке становится все больше, чем на обучающей, показывая тем самым переобучение модели
4. При дальнейшем увеличении значений данных гиперпараметров наступает такой момент, после которого ошибка на тестовой выборке вообще начинает увеличиваться

2. (*) Модифицировать реализованный алгоритм, чтобы получился стохастический градиентный бустинг. Размер подвыборки принять равным 0.5. Сравнить на одном графике кривые изменения ошибки на тестовой выборке в зависимости от числа итераций.

In [12]:

```
1  # Создадим функцию для формирования случайных наборов обучающих признаков
2  def get_stochastic_data(data, labels, n_trees):
3      n_samples = int(X_train.shape[0]*0.5)
4      stochastic_data = []
5
6      for i in range(n_trees):
7          b_data = np.zeros(data[:n_samples].shape)
8          b_labels = np.zeros(labels[:n_samples].shape)
9
10         for j in range(n_samples):
11             sample_index = random.randint(0, n_samples-1)
12             b_data[j] = data[sample_index]
13             b_labels[j] = labels[sample_index]
14         stochastic_data.append((b_data, b_labels))
15
16     return stochastic_data
```

In [13]:

```
1 def gb_fit(n_trees, max_depth, X_train, X_test, y_train, y_test, coefs, eta):
2
3     # eta - скорость обучения
4     # Деревья будем записывать в список
5     trees = []
6
7     # Будем записывать ошибки на обучающей и тестовой выборке на каждой итерации в список
8     train_errors = []
9     test_errors = []
10    stochastic_data = get_stochastic_data(X_train, y_train, n_trees)
11
12    for i in range(n_trees):
13        tree = DecisionTreeRegressor(max_depth=max_depth, random_state=42)
14
15        X_train, y_train = stochastic_data[i][0], stochastic_data[i][1]
16
17
18        # инициализируем бустинг начальным алгоритмом, возвращающим ноль,
19        # поэтому первый алгоритм просто обучаем на выборке и добавляем в список
20        if len(trees) == 0:
21            # обучаем первое дерево на обучающей выборке
22            tree.fit(X_train, y_train)
23
24            train_errors.append(mean_squared_error(y_train, gb_predict(X_train, trees,
25            test_errors.append(mean_squared_error(y_test, gb_predict(X_test, trees, coefs, eta)
26        else:
27            # Получим ответы на текущей композиции
28            target = gb_predict(X_train, trees, coefs, eta)
29
30            # алгоритмы начиная со второго обучаем на сдвиг
31            tree.fit(X_train, bias(y_train, target))
32
33            train_errors.append(mean_squared_error(y_train, gb_predict(X_train, trees,
34            test_errors.append(mean_squared_error(y_test, gb_predict(X_test, trees, coefs, eta)
35
36        trees.append(tree)
37
38    return trees, train_errors, test_errors
```

Теперь обучим несколько моделей с разными параметрами и исследуем их поведение.

In [14]:

```
1 # Число деревьев в ансамбле
2 n_trees = 50
3
4 # для простоты примем коэффициенты равными 1
5 coefs = [1] * n_trees
6
7 # Максимальная глубина деревьев
8 max_depth = 3
9
10 # Шаг
11 eta = 0.05
12
13 trees, train_errors, test_errors = gb_fit(n_trees, max_depth, X_train, X_test, y_train,
```

In [52]:

```
1 test_errors
```

Out[52]:

```
[24352.018177889797,  
24352.018177889797,  
20174.233067110566,  
16827.461929944642,  
14145.118203572792,  
11908.7435124343,  
10168.287085396425,  
8761.74231287536,  
7679.266283527152,  
6760.926782609632,  
5988.3735653758295,  
5383.001916241956,  
4946.9729081678415,  
4546.938262690729,  
4257.657517742874,  
3992.3083255373303,  
3775.1943248690895,  
3628.0911838071092,  
3486.608877049317,  
3372.9550691014306,  
3298.9625511157856,  
3238.686892283594,  
3172.049699766943,  
3114.413746776039,  
3076.1748047457722,  
3030.5784070748746,  
3012.9294179289514,  
2982.9826875993704,  
2953.3212177232085,  
2929.760092443027,  
2924.415337374638,  
2917.1335662535666,  
2911.922522741922,  
2895.874233096451,  
2883.127261201743,  
2870.7669517273016,  
2860.828241180026,  
2864.533704865219,  
2850.845765221888,  
2842.1135277359876,  
2840.64410335889,  
2847.1149889090834,  
2852.435367509558,  
2847.30251900122,  
2836.8390451614573,  
2826.4483351904464,  
2832.040591118581,  
2832.9465945741417,  
2857.9821956806395,  
2855.7192764103265,  
2851.8365538934377,  
2854.0356844110656,  
2859.0966433390167,  
2872.963682116981,  
2865.905092588673,
```

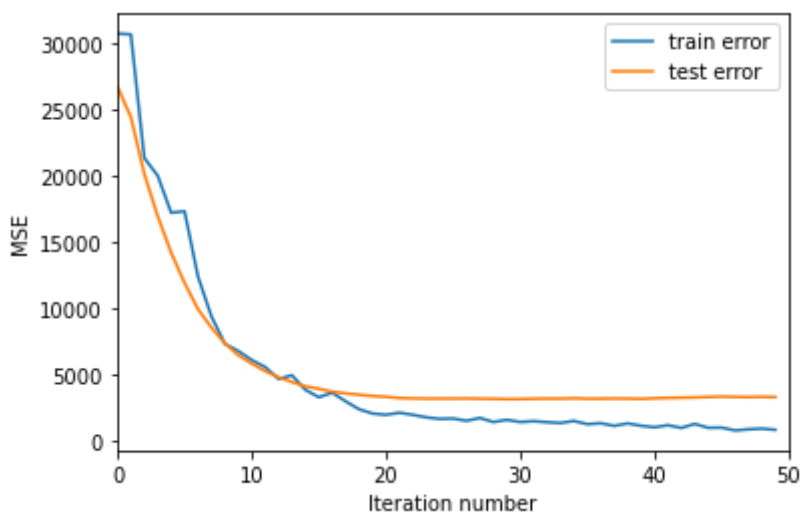
```
2859.18079259134,  
2848.046489329593,  
2867.4461893694606,  
2871.5230103370022,  
2881.3111774842073,  
2901.0572193541693,  
2898.5469577452213,  
2897.9895588351956,  
2902.521962413254,  
2891.2638168866642,  
2891.0760864236167,  
2887.9272534611932,  
2901.6990554064087,  
2925.9822490816573,  
2915.8997896634864,  
2899.4967756851315]
```

In [16]:

```
1 def get_error_plot(n_trees, train_err, test_err):  
2     plt.xlabel('Iteration number')  
3     plt.ylabel('MSE')  
4     plt.xlim(0, n_trees)  
5     plt.plot(list(range(n_trees)), train_err, label='train error')  
6     plt.plot(list(range(n_trees)), test_err, label='test error')  
7     plt.legend(loc='upper right')  
8     plt.show()
```

In [17]:

```
1 get_error_plot(n_trees, train_errors, test_errors)
```



Видим характерную ломанную линию отражающую неравномерный процесс обучения при стохастическом градиентном бустинге. Однако ошибка на тестовой выборке снижалась плавно и достигла примерно аналогичных значений с классическим бустингом при одинаковых гиперпараметрах.

3. (*) Модифицировать алгоритм градиентного бустинга, взяв за основу реализацию решающего дерева из ДЗ_4. Сделать выводы о качестве алгоритма по сравнению с реализацией из п.1.

In [20]:

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

In [21]:

```
1 # И класс терминального узла (листа)
2 class Leaf:
3
4     def __init__(self, data, labels):
5         self.data = data
6         self.labels = labels
7         self.prediction = self.predict()
8
9     def predict(self):
10         prediction = np.mean(self.labels)
11         return prediction
```

In [22]:

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

In [23]:

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

In [24]:

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

In [25]:

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

In [26]:

```
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 [27]:

```
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 [35]:

```
1 # Построение дерева с помощью рекурсивной функции
2 def build_tree(data, labels, max_depth):
3     global depth, true_branch, false_branch
4     # print("Глубина", depth)
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    # print("quality == 0")
11
12    if depth == 3:
13        return Leaf(data, labels)
14
15    # Рекурсивно строим два поддерева
16    # print("Делаем ветвление на глубине ", depth)
17    depth += 1
18    true_data, false_data, true_labels, false_labels = split(data, labels, index, t)
19    true_branch = build_tree(true_data, true_labels, depth)
20    false_branch = build_tree(false_data, false_labels, depth)
21
22    # Возвращаем класс узла со всеми поддеревьями, то есть целого дерева
23    return Node(index, t, true_branch, false_branch)
```

In [36]:

```
1 X_train, X_test, y_train, y_test = model_selection.train_test_split(X, y, test_size=0.2)
```

In [37]:

```
1 def gb_predict(X, trees_list, coef_list, eta):
2     # Реализуемый алгоритм градиентного бустинга будет инициализироваться нулевыми значениями
3     # поэтому все деревья из списка trees_list уже являются дополнительными и при предсказании
4     # прибавляются с шагом eta
5     return np.array([sum([eta * coef * predict([x], alg)[0] for alg, coef in zip(trees_list, coef_list)]) for x in X])
```

Обучим несколько моделей с разным количеством деревьев и посмотрим как оно влияет на ошибку

In [42]:

```
1 def gb_fit(n_trees, max_depth, X_train, X_test, y_train, y_test, coefs, eta):
2     global depth
3
4     # eta - скорость обучения
5     # Деревья будем записывать в список
6     trees = []
7
8     # Будем записывать ошибки на обучающей и тестовой выборке на каждой итерации в список
9     train_errors = []
10    test_errors = []
11
12    for i in range(n_trees):
13
14        # инициализируем бустинг начальным алгоритмом, возвращающим ноль,
15        # поэтому первый алгоритм просто обучаем на выборке и добавляем в список
16        if len(trees) == 0:
17            # обучаем первое дерево на обучающей выборке
18
19            depth = 0
20            true_branch = None
21            false_branch = None
22            my_tree = build_tree(X_train, y_train, max_depth)
23            trees.append(my_tree)
24
25            train_errors.append(mean_squared_error(y_train, gb_predict(X_train, trees,
26            test_errors.append(mean_squared_error(y_test, gb_predict(X_test, trees, coefs, eta)
27        else:
28            # Получим ответы на текущей композиции
29            target = gb_predict(X_train, trees, coefs, eta)
30
31            # алгоритмы начиная со второго обучаем на сдвиг
32
33            depth = 0
34            true_branch = None
35            false_branch = None
36            my_tree = build_tree(X_train, bias(y_train, target), max_depth)
37
38
39            train_errors.append(mean_squared_error(y_train, gb_predict(X_train, trees,
40            test_errors.append(mean_squared_error(y_test, gb_predict(X_test, trees, coefs, eta)
41
42            trees.append(my_tree)
43
44    return trees, train_errors, test_errors
```

In [43]:

```
1 max_depth = 3
2 eta = 0.05
3 list_mse = []
4 for n_trees in range(1, 81, 10):
5     coefs = [1] * n_trees
6     trees, train_errors, test_errors = gb_fit(n_trees, max_depth, X_train, X_test, y_train, y_test)
7     train_prediction = gb_predict(X_train, trees, coefs, eta)
8     test_prediction = gb_predict(X_test, trees, coefs, eta)
9     round(mean_squared_error(y_test, test_prediction))
10    list_mse.append([n_trees, round(mean_squared_error(y_train, train_prediction)), round(mean_squared_error(y_test, test_prediction))])
11
```

In [44]:

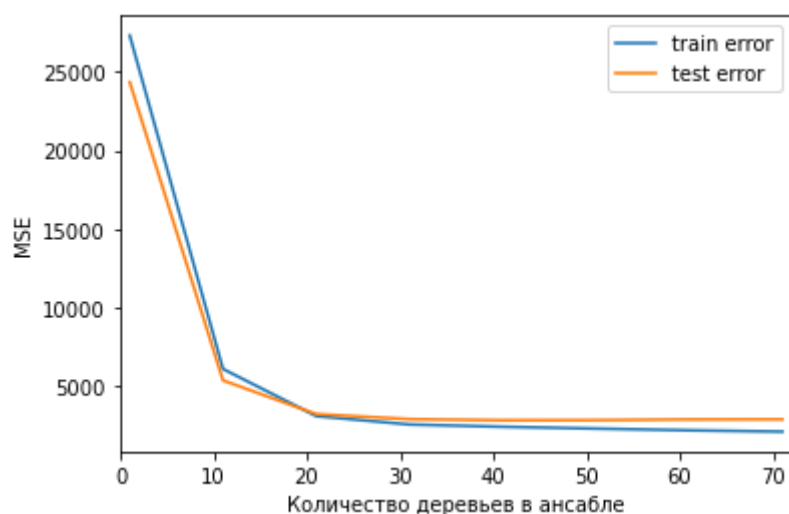
```
1 list_mse
```

Out[44]:

```
[[1, 27317.0, 24352.0],
 [11, 6108.0, 5383.0],
 [21, 3109.0, 3239.0],
 [31, 2576.0, 2917.0],
 [41, 2431.0, 2847.0],
 [51, 2315.0, 2854.0],
 [61, 2205.0, 2899.0],
 [71, 2122.0, 2905.0]]
```

In [45]:

```
1 plt.xlabel('Количество деревьев в ансамбле')
2 plt.ylabel('MSE')
3 plt.xlim(0, list(zip(*list_mse))[0][-1] + 1)
4 plt.plot(list(zip(*list_mse))[0], list(zip(*list_mse))[1], label='train error')
5 plt.plot(list(zip(*list_mse))[0], list(zip(*list_mse))[2], label='test error')
6 plt.legend(loc='best')
7 plt.show()
```



Выводы:

Качество алгоритма на основе решающего дерева из ДЗ_4 оказалось лучше по сравнению с реализацией из п.1. Ошибка, полученная при его работе оказалась меньше. Также алгоритм оказался менее склонен к переобучению, что видно из графиков и таблиц.