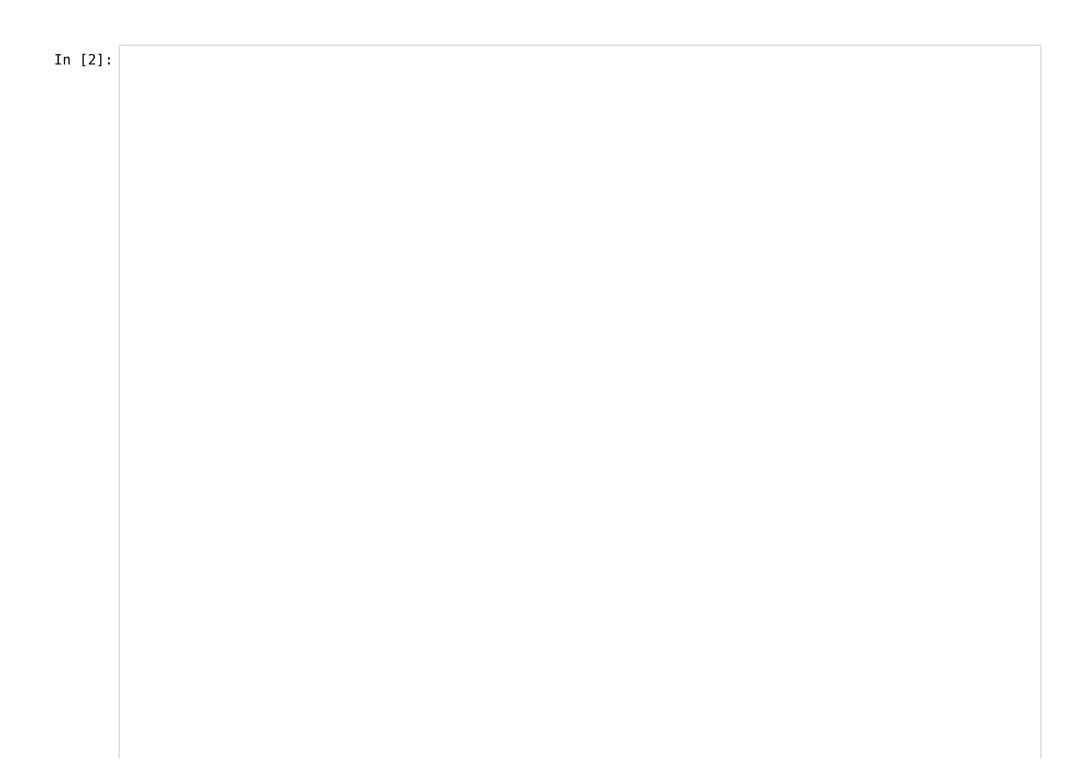
```
In [1]: from sklearn.model_selection import cross_val_score
    from IPython.display import Image
    from sklearn import tree
    from sklearn.metrics import accuracy_score, classification_report
    import numpy as np
    import pandas as pd
    import matplotlib.pyplot as plt
    import math
    from sklearn.cross_validation import train_test_split
    from sklearn.metrics import accuracy_score, classification_report
    from sklearn.metrics import mean squared error
```

/home/nimloth/anaconda3/envs/python27/lib/python2.7/site-packages/sklearn/cross_validation.py:44: Deprecatio nWarning: This module was deprecated in version 0.18 in favor of the model_selection module into which all the refactored classes and functions are moved. Also note that the interface of the new CV iterators are different from that of this module. This module will be removed in 0.20.

"This module will be removed in 0.20.", DeprecationWarning)

Сначала мне что-то ударило в голову и я написала классификатор :(



```
class DecisionTreeClassifier:
    def __init__(self, max_depth):
        self.max depth = max depth
        if self.max depth != 1:
            self.L = DecisionTreeClassifier(max_depth=self.max_depth - 1)
            self.R = DecisionTreeClassifier(max depth=self.max depth - 1)
   def advanced log(self, x):
        if x == 0:
            return 0
        else:
            return math.log(x)
   def gini(self, A p 0, A p 1, B p 0, B p 1):
        N = A p 0 + A p 1 + B p 0 + B p 1
       L = A p 0 + A p 1
       R = B p 0 + B p 1
        if L == 0 or R == 0:
            return 1
       A p \theta = A p \theta / L
       A p 1 = A p 1 / L
       B p \theta = B p \theta / R
        B p 1 = B p 1 / R
        g = L / N * (-A p 0 * self.advanced log(A p 0) - A p 1 * self.advanced log(A p 1))
        + R / N * (- B p 0 * self.advanced log(B p 0) - B p 1 * self.advanced log(B p 1))
        return q
    def fit(self, train data, train target):
        class0 = 0
        class1 = 0
        for elem in train target:
            if elem == 0:
                class0 = class0 + 1
            else:
                class1 = class1 + 1
        if class0 >= class1:
            self.ans = 0
        else:
            self.ans = 1
        if self.max depth == 1:
```

```
return
    min g = 1
    for ith arg in range(len(train data[0])):
        for jth threshold in range(len(train data)):
            threshold = train data[jth threshold][ith arg]
            A p 0 = A p 1 = B p 0 = B p 1 = 0
            for jth elem in range(len(train data)):
                if train data[jth elem][ith arg] <= threshold:</pre>
                    if train target[jth elem] == 0:
                         A p 0 = A p 0 + 1
                     else:
                         A p 1 = A p 1 + 1
                else:
                    if train target[jth elem] == 0:
                         B_p_0 = B_p + 1
                     else:
                         B p 1 = B p 1 + 1
            g = self.gini(A_p_0, A_p_1, B_p_0, B_p_1)
            if q < min q:</pre>
                self.threshold = threshold
                self.arg = ith arg
                min g = g
    if min q == 1:
        return
    A train = []
    A target = []
    B train = []
    B target = []
    for jth elem in range(len(train data)):
        if train data[jth elem][self.arg] <= self.threshold:</pre>
            A train.append(train data[jth elem])
            A target.append(train target[jth elem])
        else:
            B train.append(train data[jth elem])
            B target.append(train target[jth elem])
    self.L.fit(A_train, A_target)
    self.R.fit(B train, B target)
def predict_for_one(self, test_data):
    if hasattr(self, 'threshold'):
        if test data[self.arg] <= self.threshold:</pre>
```

```
return self.L.predict_for_one(test_data)
    else:
        return self.R.predict_for_one(test_data)

else:
    if hasattr(self, 'ans'):
        return self.ans
    else:
        print "Node does not have answer"

def predict(self, test_data):
    ans = []
    for elem in test_data:
        ans.append(self.predict_for_one(elem))
    return ans
```

Проверила его на выборке из задания 2.

```
In [3]: samples = pd.read_csv('german_credit.csv')
N = samples.shape[0]
y = samples['Creditability'].reshape(1, N)[0]
x = np.array(samples)
x = np.array([elem[1:] for elem in x]).reshape(N, 20)

/home/nimloth/anaconda3/envs/python27/lib/python2.7/site-packages/ipykernel/__main__.py:3: FutureWarning: re shape is deprecated and will raise in a subsequent release. Please use .values.reshape(...) instead app.launch_new_instance()

In [4]: tree = DecisionTreeClassifier(max_depth=4)

In [5]: train_data, test_data, train_target, test_target = train_test_split(x, y, test_size = 0.2)

In [6]: tree.fit(train_data, train_target)

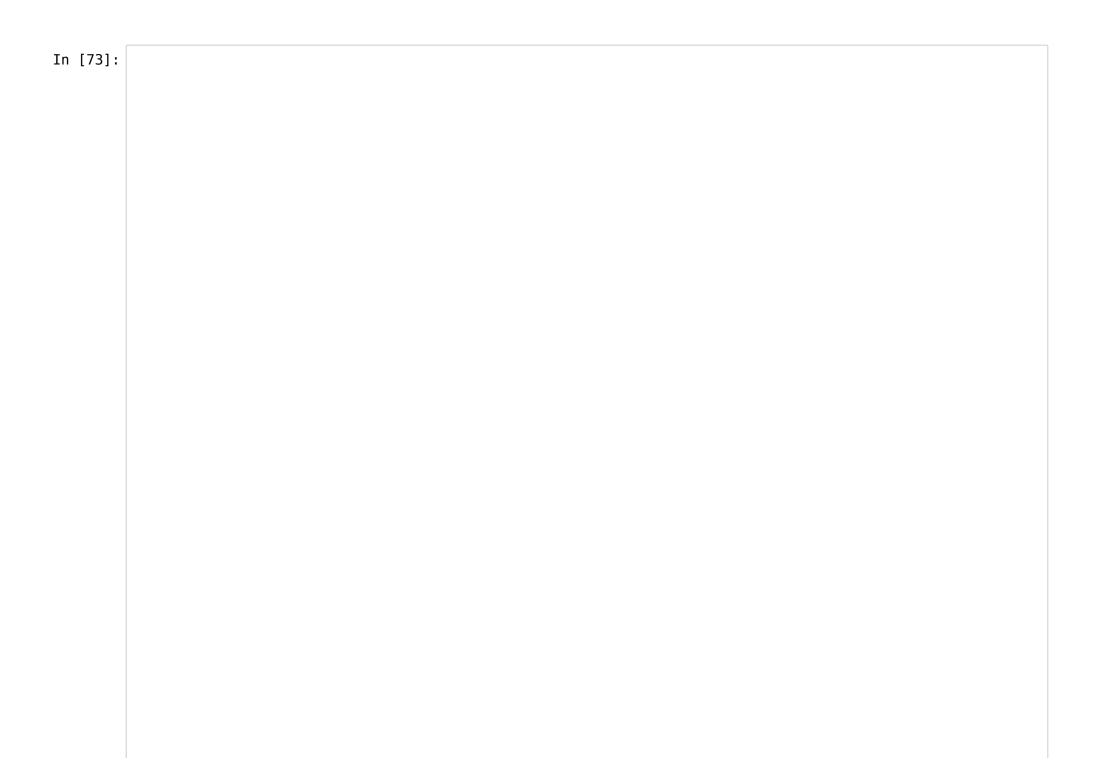
In [7]: test predictions = tree.predict(test data)
```

In [8]: print accuracy_score(test_predictions, test_target)

0.725

Вау! Он классно работает! Прямо как во втором задании точность.

Немного переписав решающее дерево получим класс для решения задачи регрессии.



```
#основные атрибуты это self.threshold, self.arg - по какому признаку и порогу делим вершину
#атрибут self.ans отвечает за ответ в данной вершине
#aтрибуты self.L, self.R содержат в себе левое и правле поддеревья
class DecisionTreeRegressor:
   #инициалицируем все дерево, используя глубину
   def init (self, max depth):
        self.max depth = max depth
        if self.max depth != 1:
            self.L = DecisionTreeRegressor(max depth=self.max depth - 1)
            self.R = DecisionTreeRegressor(max depth=self.max depth - 1)
   #взвешенная функция нехорошести
   def variance(self, A, B):
       N = len(A) + len(B)
       L = len(A)
       R = len(B)
        if L == 0 or R == 0:
            return 100000000
        g = L / N * np.array(A).var()
       + R / N * np.array(B).var()
        return q
   def fit(self, train data, train target):
        self.ans = np.array(train target).mean()
        if self.max depth == 1:
            return
        #перебираем все признаки и пороги в поиске наименьшей функции нехорошести
       #устанавливаем атрибуты self.arg, self.thresholf
        min q = 100000000
        for ith arg in range(len(train data[0])):
            for jth threshold in range(len(train data)):
                threshold = train data[jth threshold][ith arg]
                A = []
                B = []
                for jth elem in range(len(train data)):
                    if train data[jth elem][ith arg] <= threshold:</pre>
                        A.append(train target[jth elem])
                    else:
                        B.append(train target[jth elem])
                q = self.variance(A, B)
```

```
if q < min q:
                self.threshold = threshold
                self.arg = ith arg
                min g = g
    if min g == 1000000000:
        return
    #делим выборку на 2 вершины и запускаемся от них рекурсивно
    A train = []
    A target = []
    B train = []
    B target = []
    for jth elem in range(len(train data)):
        if train data[jth elem][self.arg] <= self.threshold:</pre>
            A train.append(train_data[jth_elem])
            A target.append(train target[jth elem])
        else:
            B train.append(train data[jth elem])
            B target.append(train target[jth elem])
    self.L.fit(A train, A target)
    self.R.fit(B train, B target)
#пропихивает выборку по вершинам решающего дерева пока может, иначе выдает ответ в вершине
def predict for one(self, test data):
    if hasattr(self, 'threshold'):
        if test data[self.arg] <= self.threshold:</pre>
            return self.L.predict for one(test data)
        else:
            return self.R.predict for one(test data)
    else:
        if hasattr(self, 'ans'):
            return self.ans
        else:
            print "Node does not have answer"
def predict(self, test_data):
    ans = []
    for elem in test data:
        ans.append(self.predict for one(elem))
    return ans
```

Проверим точность на датасете Бостон.

Неплохо... Но могло бы быть и лучше...