

Задание. Дерево решений

В нижеприведенном коде

- 1) Реализовать функцию расчет критерия Джини
- 2) Реализовать расчет критерия прироста информации
- 3) Реализовать критерии останова (не менее двух)
- 4) Реализовать функцию подсчета метрики модели
- 5) Проверить работу самописного дерева для решения задачи классификации (датасет можно сгенерировать самостоятельно при помощи функции `make_classification`, либо взять любой набор данных из задания к ЛР4), сравнив полученные результаты с результатами дерева решений из `sklearn`. Результаты должны совпасть полностью.

```
# 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 узле
        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 его
        prediction = max(classes, key=classes.get)
        return prediction
```

1)

```
# Расчет критерия Джини
def gini(labels):
```

2)

```
# Расчет прироста информации

def gain(left_labels, right_labels, root_gini):

# Разбиение датасета в узле

def split(data, labels, column_index, t):
    left = np.where(data[:, column_index] <= t)
    right = np.where(data[:, column_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
```

3)

```
# Нахождение наилучшего разбиения

def find_best_split(data, labels):
    # обозначим минимальное количество объектов в узле
    min_samples_leaf = 3

    root_gini = gini(labels)

    best_gain = 0
    best_t = None
    best_index = None

    n_features = data.shape[1]

    for index in range(n_features):
        t_values = np.unique(data[:, index])

        for t in t_values:
            true_data, false_data, true_labels, false_labels = split(data,
labels, index, t)
            current_gain = gain(true_labels, false_labels, root_gini)

            if current_gain > best_gain:
                best_gain, best_t, best_index = current_gain, t, index

    return best_gain, best_t, best_index
```

```
import time
# Построение дерева с помощью рекурсивной функции

def build_tree(data, labels):
    gain, t, index = find_best_split(data, labels)
```

```

# Базовый случай - прекращаем рекурсию, когда нет прироста в качества
if gain == 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)

# print(time.time(), true_branch)
false_branch = build_tree(false_data, false_labels)

# print(time.time(), false_branch)

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

```

```

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)
    else:
        return classify_object(obj, node.false_branch)

```

```

def predict(data, tree):

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

```

```

# Напечатаем ход нашего дерева
def print_tree(node, spacing=""):

    # Если лист, то выводим его прогноз
    if isinstance(node, Leaf):
        print(spacing + "Прогноз:", node.prediction)
        return

    # Выведем значение индекса и порога на этом узле
    print(spacing + 'Индекс', str(node.index), '<=', str(node.t))

    # Рекурсионный вызов функции на положительном поддереве
    print(spacing + '--> True:')
    print_tree(node.true_branch, spacing + "  ")

```

```
# Рекурсионный вызов функции на отрицательном поддереве
print(spacing + '--> False:')
print_tree(node.false_branch, spacing + "  ")

print_tree(my_tree)
```

4)

```
# Введем функцию подсчета точности как доли правильных ответов
def accuracy_metric(actual, predicted):
```