

Бинарные деревья поиска

25/02/2023

В предыдущих сериях

- Хотим структуру с *get*, *lower_bound* по значению, добавлением и удалением по значению
- Skip-list подходит

Пара слов про кучи

В куче нет запроса *get*, есть только *get_min*. Такое бывает полезно:

- Например, если вы закинули все элементы массива в кучу, а потом доставали по очереди минимумы и удаляли их, то вы отсортировали массив (это называется heap sort).

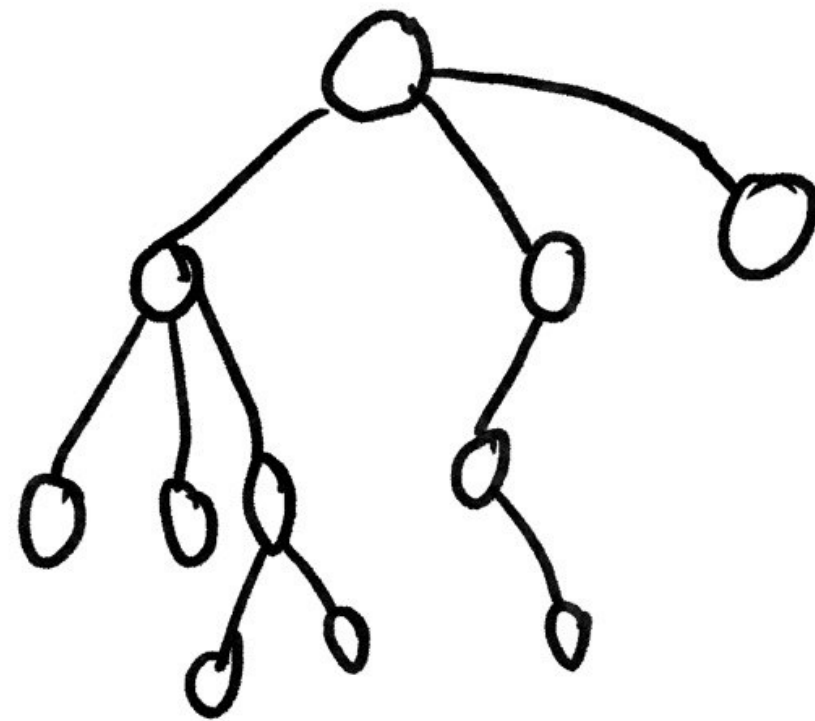
Но нам это не очень интересно, потому что деревья делают все то же, что куча, и даже больше.

Дерево поиска

Мы с вами уже обсуждали списки - это такая структура, которая хранит ссылку на следующий элемент.

В деревьях узлы хранят ссылку на несколько "следующих" узлов, которые называются "сыновьями". Каждый сын является корнем своего поддерева. Все поддерева сыновей независимы и не пересекаются. Точка входа в дерево это корень.

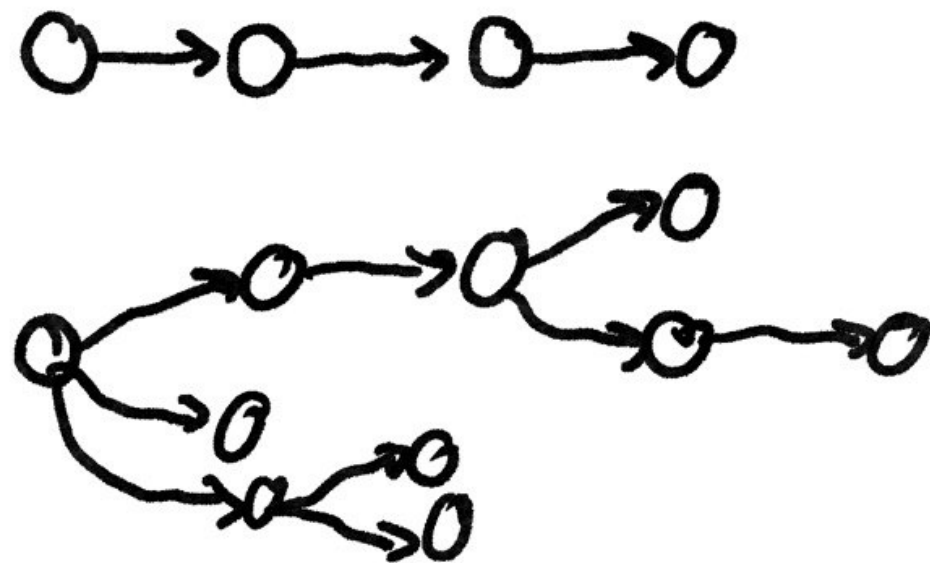
Дерево



Вопрос

Является ли список деревом?

Под другим углом



Поддеревья

Некоторые величины можно пересчитывать через поддерево с помощью рекурсивных определений

Например,

$$sum(tree) = root_{val} + \sum_{sons} sum(subtree(sons))$$

$$sz(tree) = 1 + \sum_{sons} sz(subtree(sons))$$

Бинарное дерево поиска

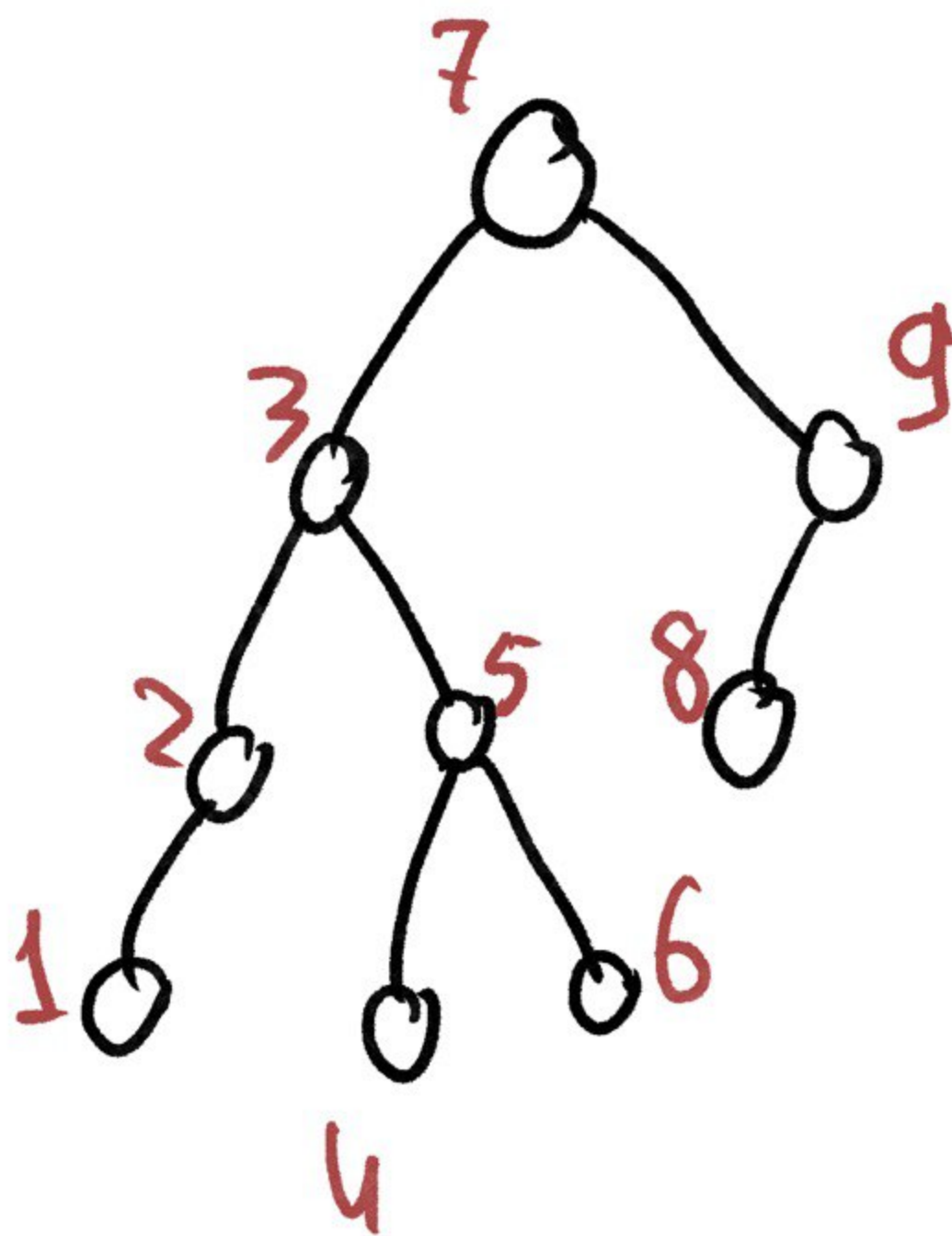
Дерево называется *бинарным*, если у каждой вершины не больше 2 сыновей. Тогда можно их назвать левым и правым

```
class Node:  
    left: None  
    right: None  
    val: 0
```

Бинарное дерево становится деревом поиска, если для каждой вершины *node*:

$$\max(subtree(node_l)) \leq node_{val} \leq \min(subtree(node_r))$$

BST

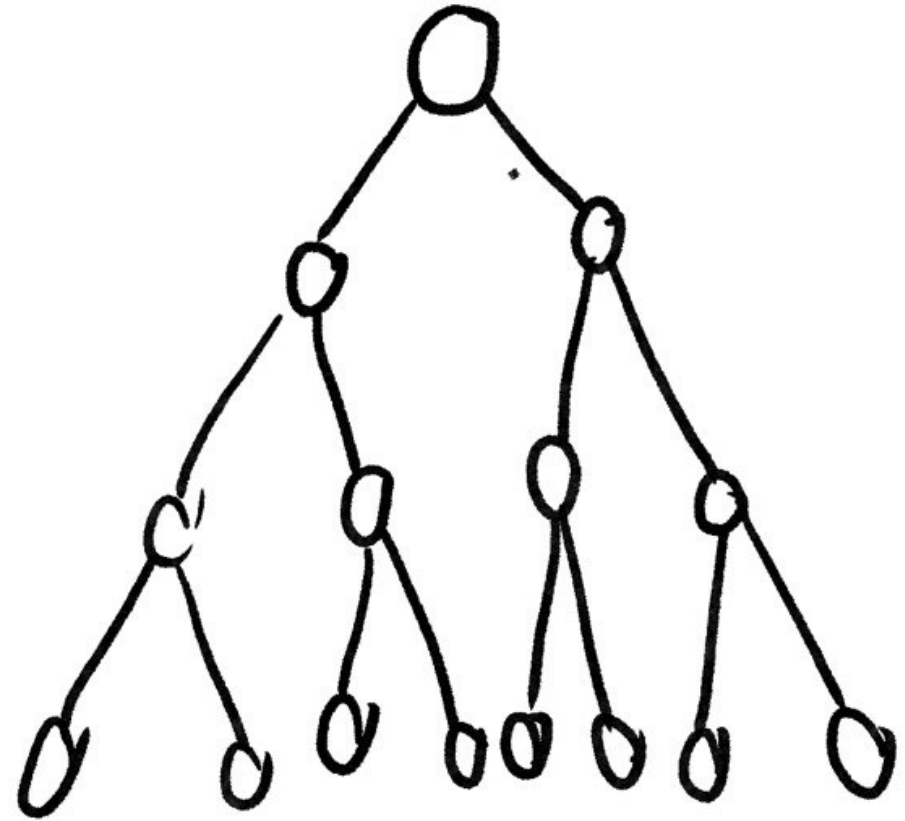


Алгоритм поиска

```
def find(node, x):  
    if node.val == x:  
        return node  
    if node.val < x:  
        return find(node.r, x)  
    return find(node.l, x)
```

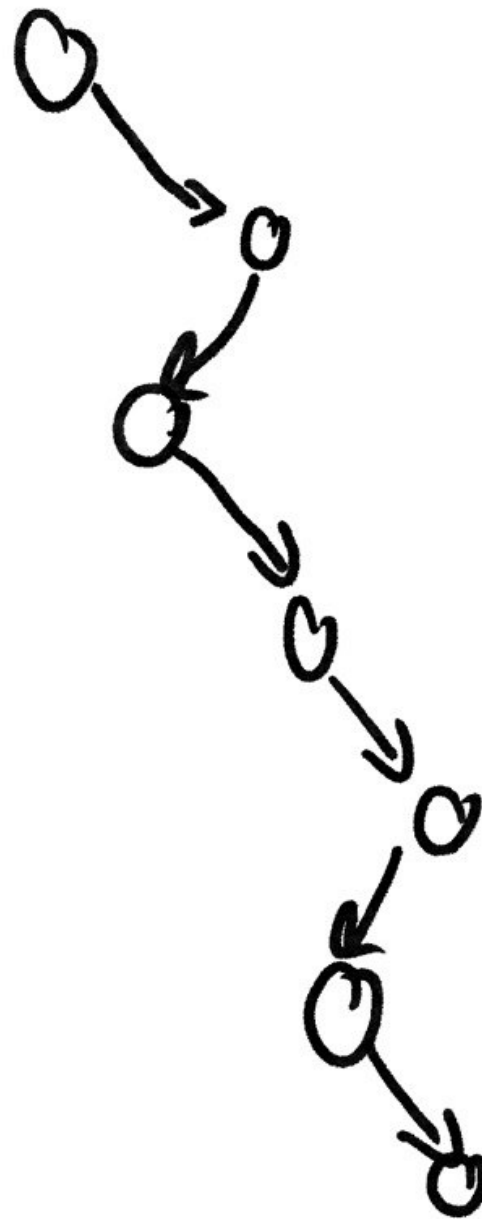
Вопрос

Какая асимптотика у `find` в таком дереве?



Вопрос

А в таком?



Проблема

Дерево должно иметь небольшую глубину. При наивном `insert` (`find` + добавление в конце, если `None`) дерево может иметь глубину $O(n)$.

Принято делать перебалансировку после `insert` вдоль соответствующего пути в дереве.

Детерминированно:

- AVL-tree
- RB-tree

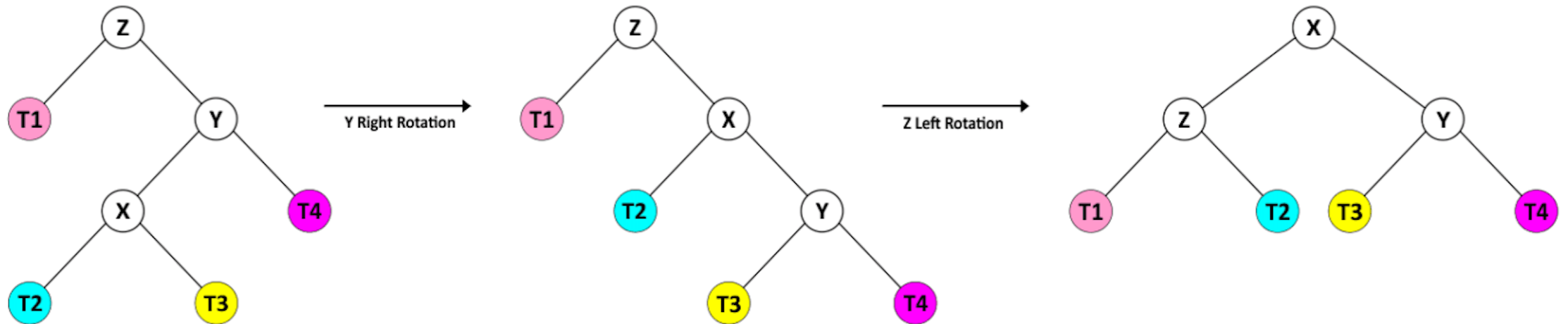
Недетерминированно:

- Декартово дерево (Treap)

AVL

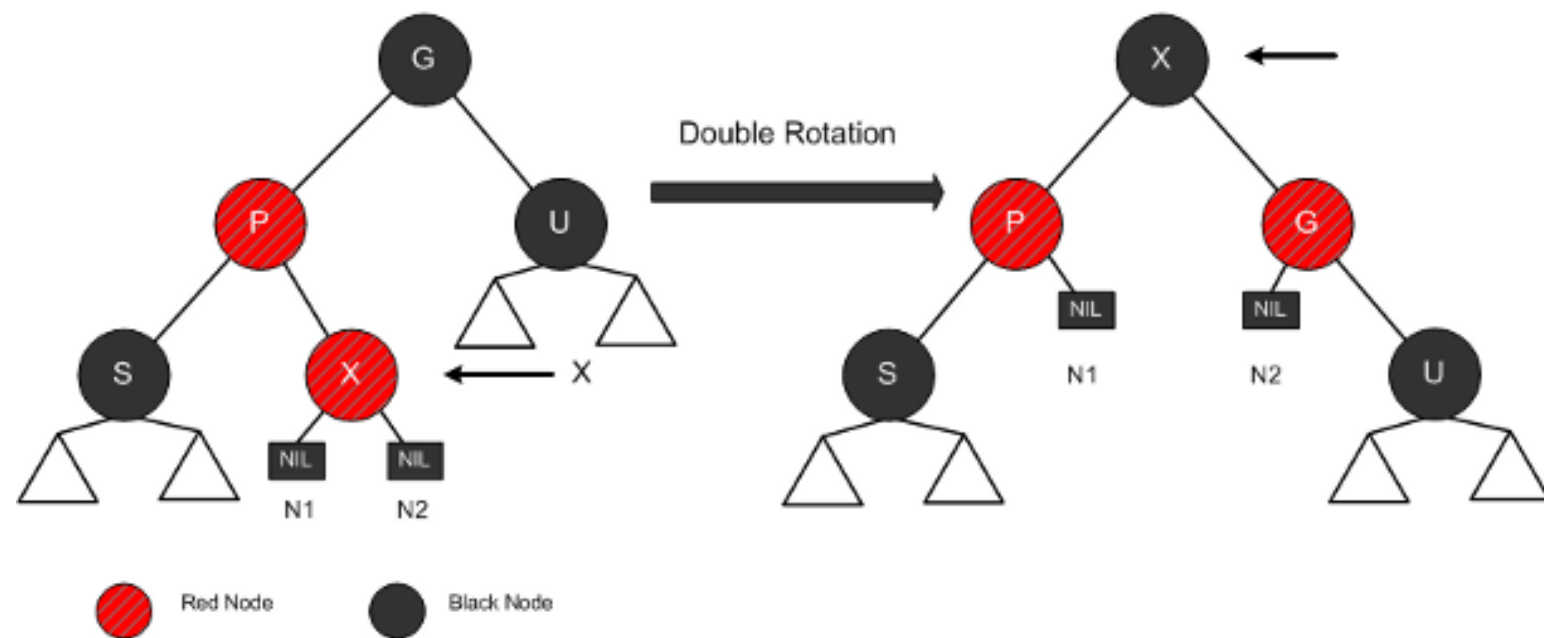
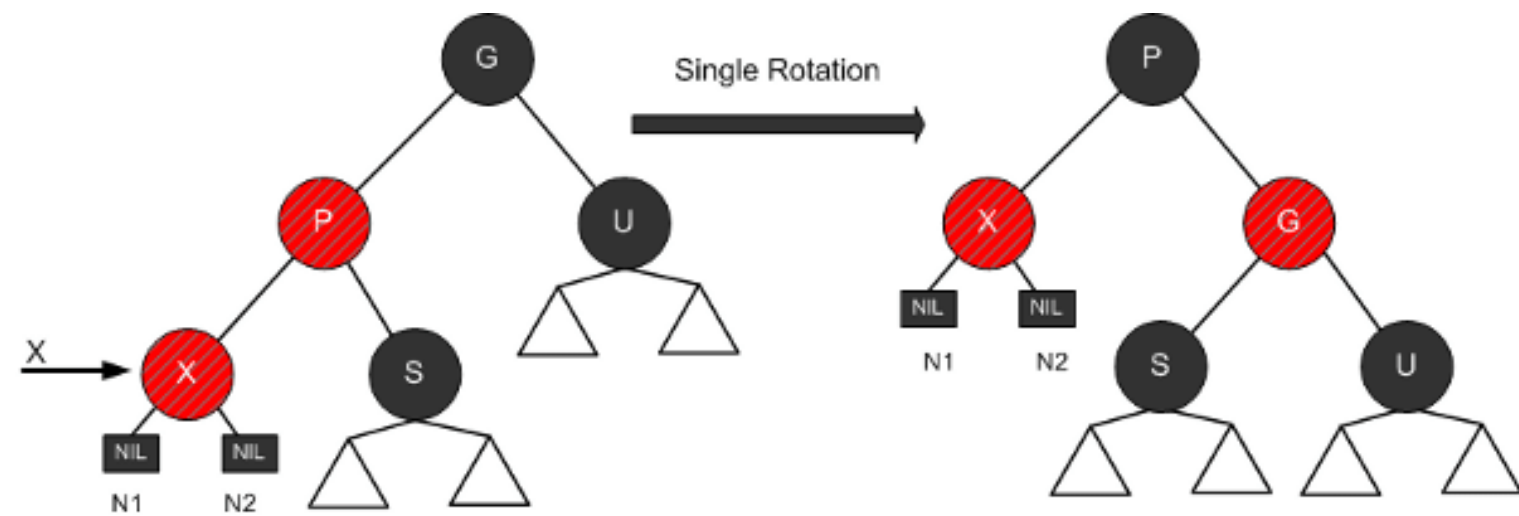
Основное свойство - у любой вершины глубина левого и правого поддеревьев отличается не более, чем на 1. Если это не так, то можно сделать некороткие переподвешивания и починить это свойство. Глубина $O(\log n)$.

Повороты в дереве



Red-black tree

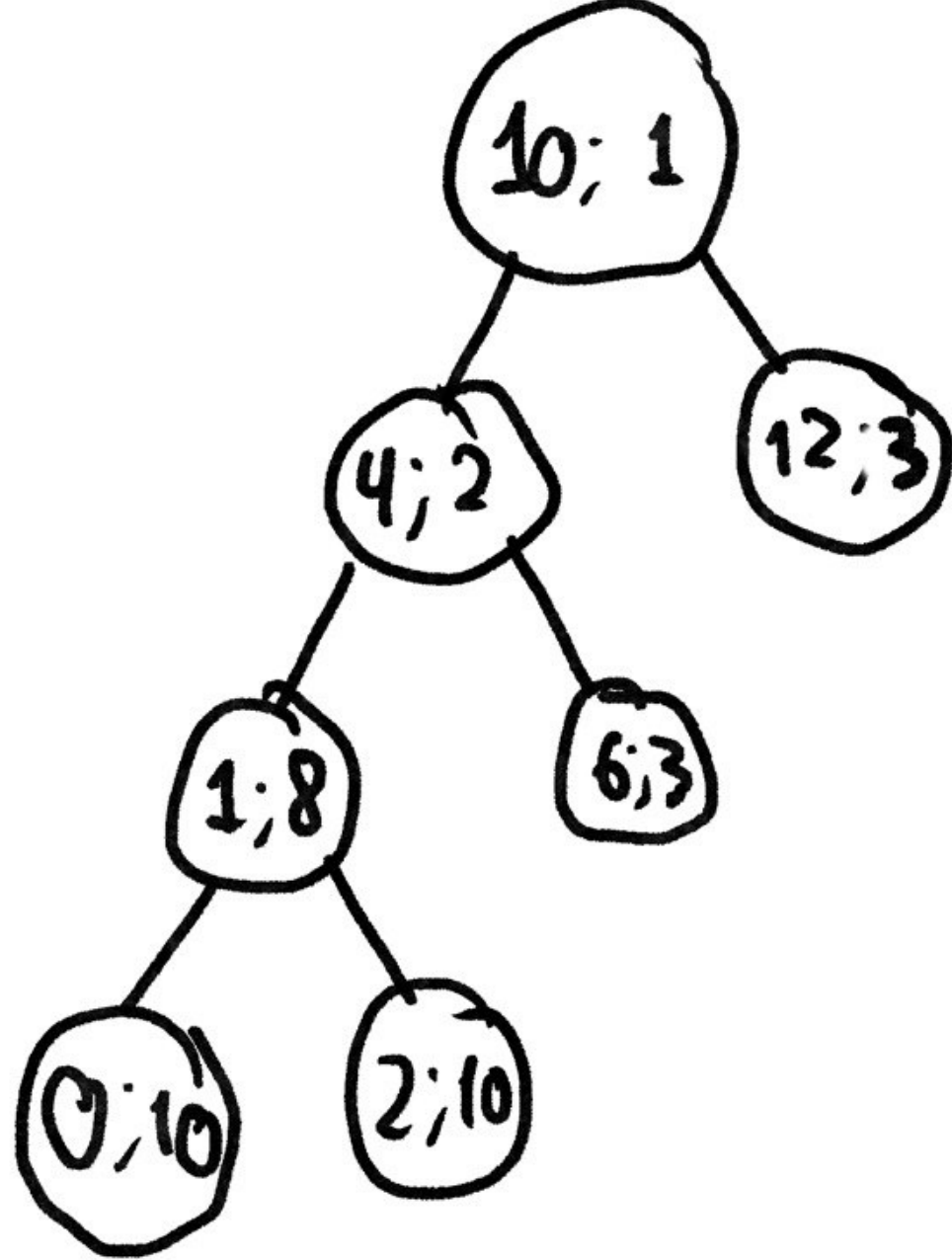
Вершина имеет цвет, и на цвета наложены некоторые правила. Если правило нарушено, можно починить его поворотом. Структура дерева менее четкая, но зато нужно меньше изменений структуры. Глубина $O(\log n)$ все равно гарантируется



Декартово дерево

Вершинам выдается *приоритет* y . Декартово дерево является бинарным деревом поиска по x и кучей по y . Таким образом, оно строится однозначно.

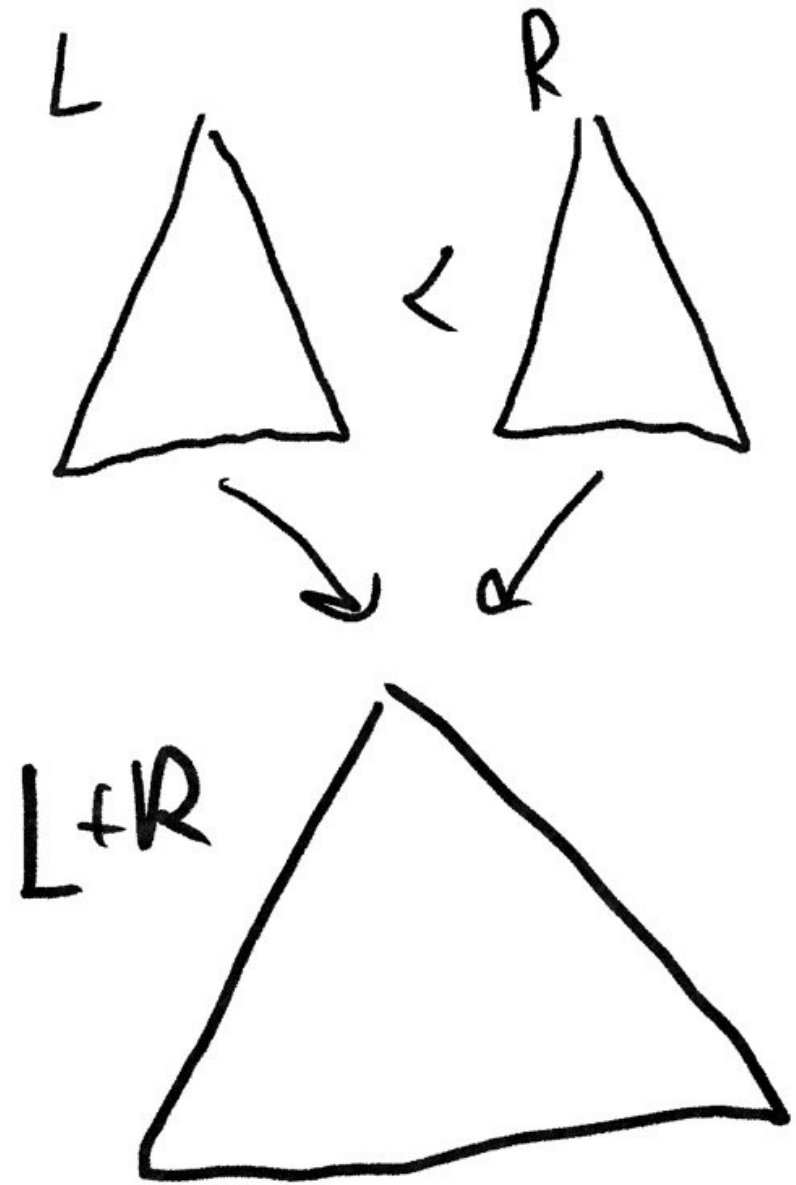
```
class Node:
    def __init__(self, x):
        self.x = x
        self.y = random()
```



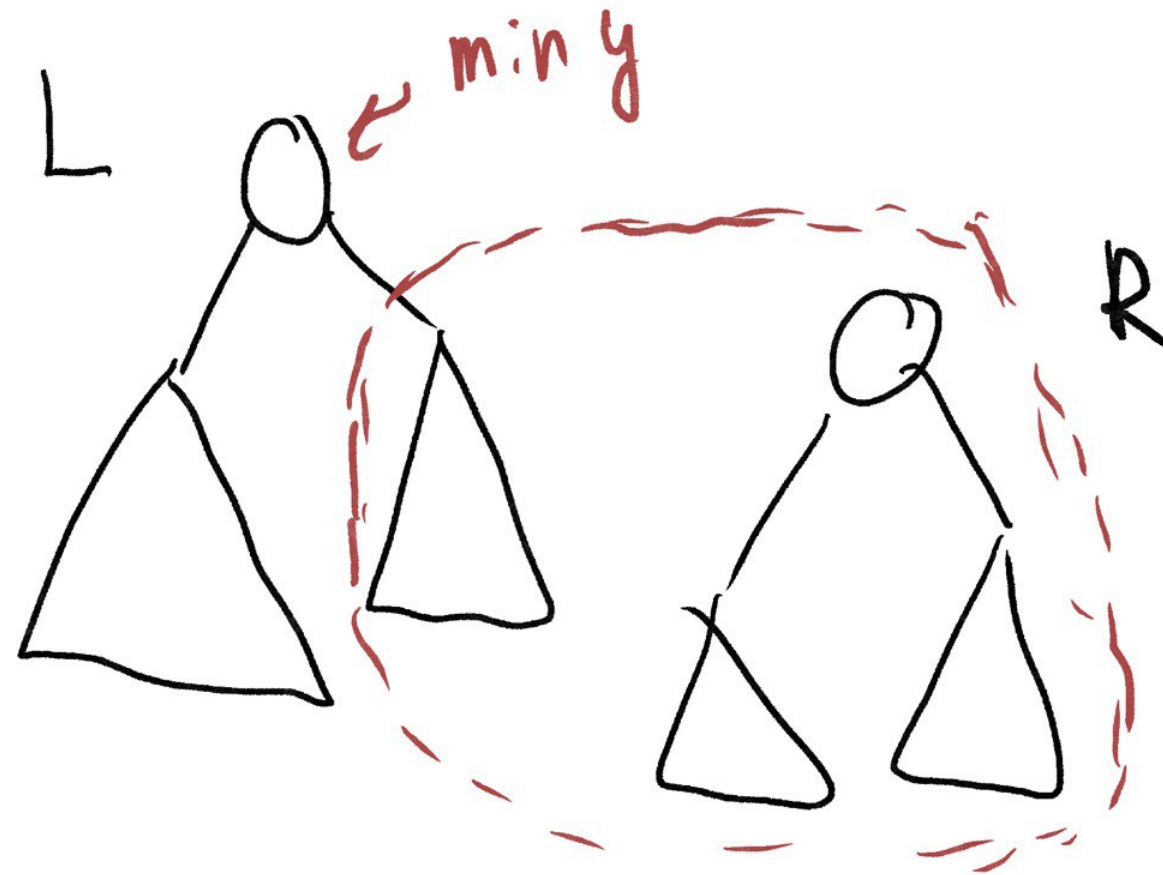
Операции

Декартово дерево строится на двух операциях: `merge` и `split`.

Merge

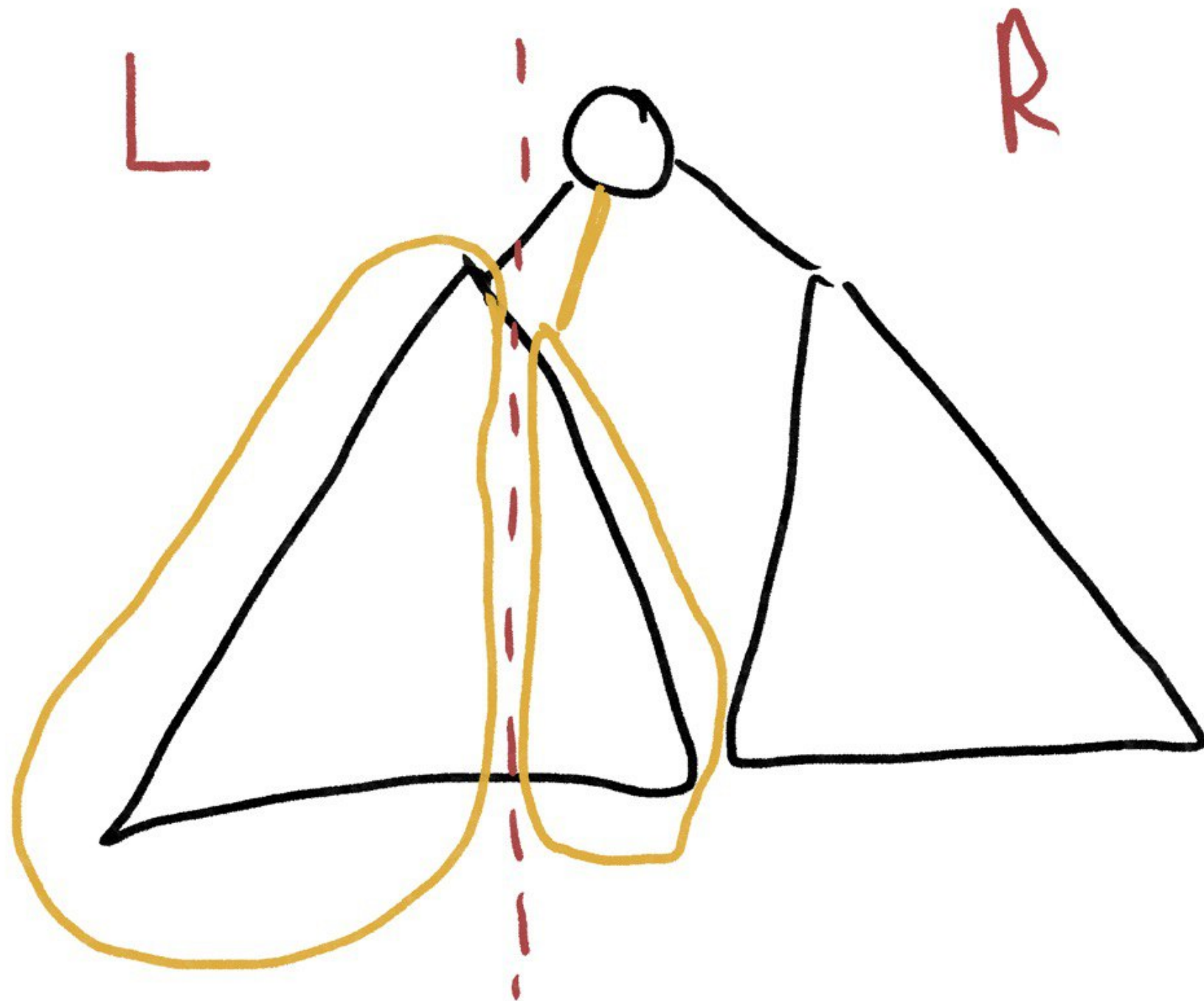


Merge



Merge

```
def merge(l, r):  
    if l is None:  
        return r  
    if r is None:  
        return l  
    if l.y < r.y:  
        l.r = merge(l.r, r)  
        return l  
    else:  
        r.l = merge(l, r.l)  
        return r
```

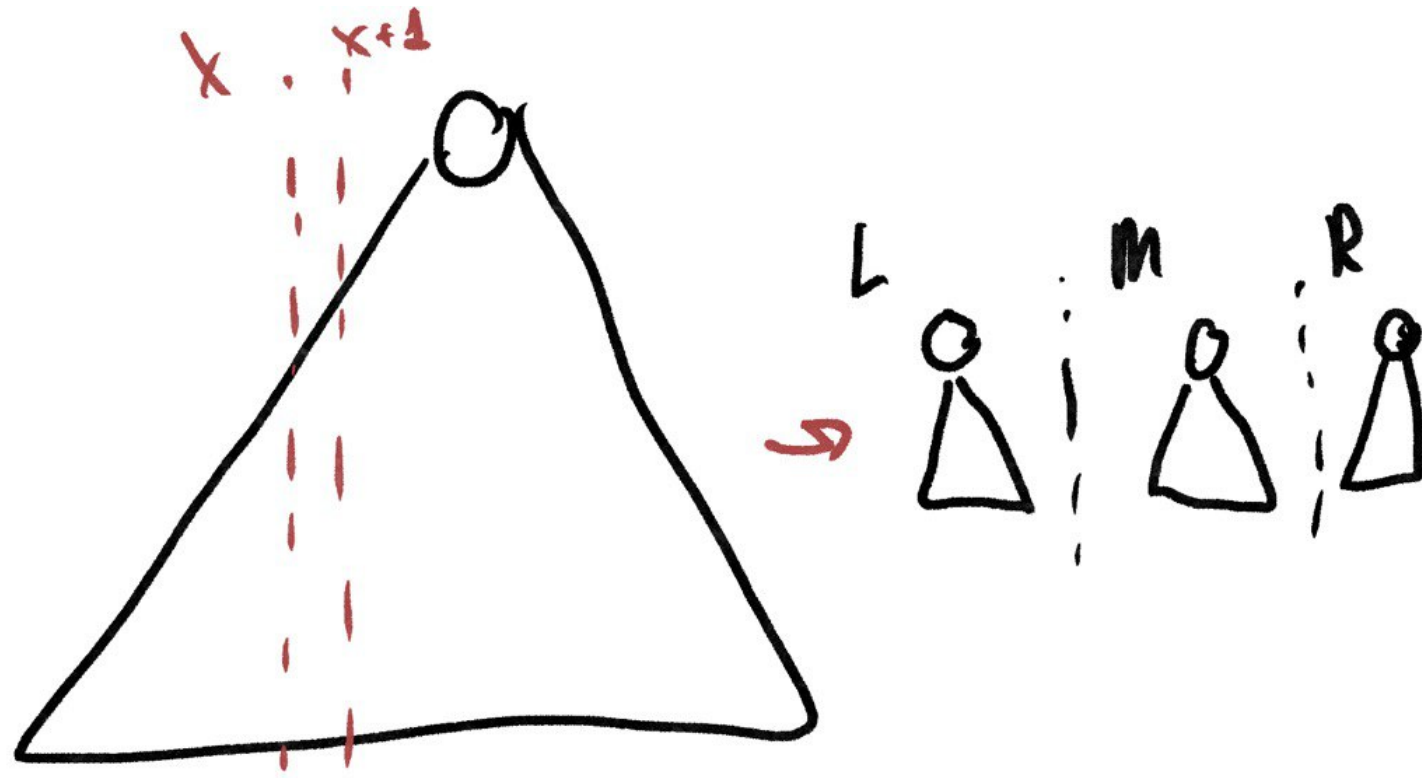
Split

```
def split(node, x):  
    if node is None:  
        return None, None  
    if node.x < x:  
        l, r = split(node.r, x)  
        node.r = l  
        return node, r  
    else:  
        l, r = split(node.l, x)  
        node.l = r  
        return l, node
```

Вопрос

Как реализовать поиск? Добавление? Удаление?

Добавление



Асимптотика

$O(h)$ на операцию. Надо только понять, почему $h \sim \log n$.

Надо понять матожидание количества предков у вершины.

$$P[u \text{ is ancestor of } v] = P[u_y = \max_{w \in [u; v]} w_y] = \frac{1}{|u - v|}$$

$$\sum_u \frac{1}{|u - v|} = O(\log n)$$

B-tree

Давайте сделаем немного другой подход, где у вершины не 2 сына, а произвольное количество. Сын хранит максимума в своем поддереве, поэтому мы можем правильно делать поиск.

```
class Leaf:
    def __init__(self, x):
        self.x = x

class Node:
    def __init__(self):
        self.children = []
        self.max = None
```

B+ tree

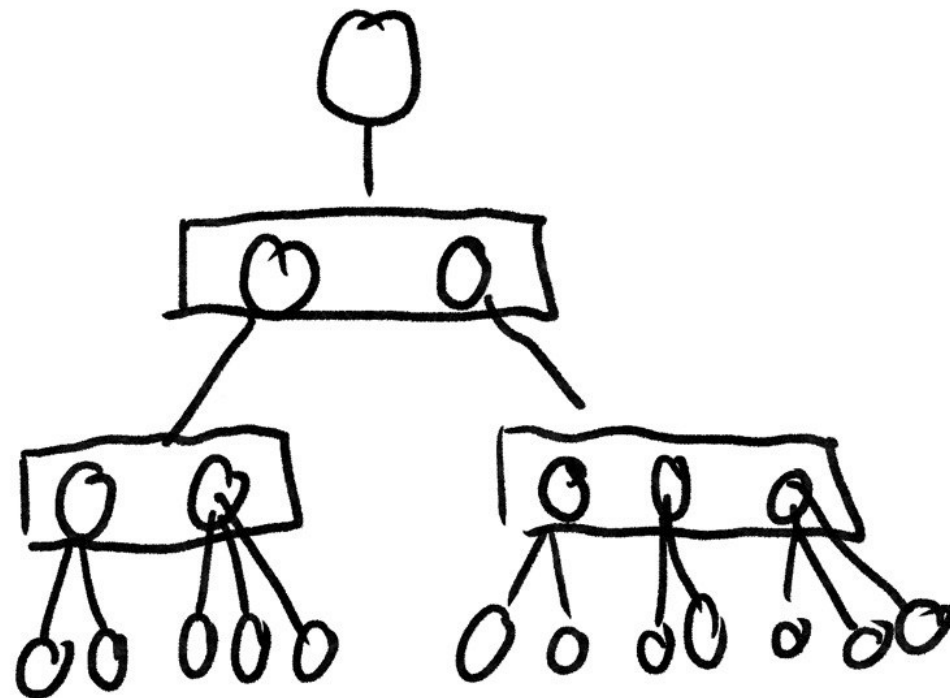
В реализациях B-деревьев есть разные подходы, мы сделаем тот, который называется B+ деревом.

В нашем дереве все значения будут записаны в листьях, а вершины выше будут *поисковыми* вершинами. У каждой из них от b до $b \cdot 2 - 1$ сыновей (отсюда и название.).

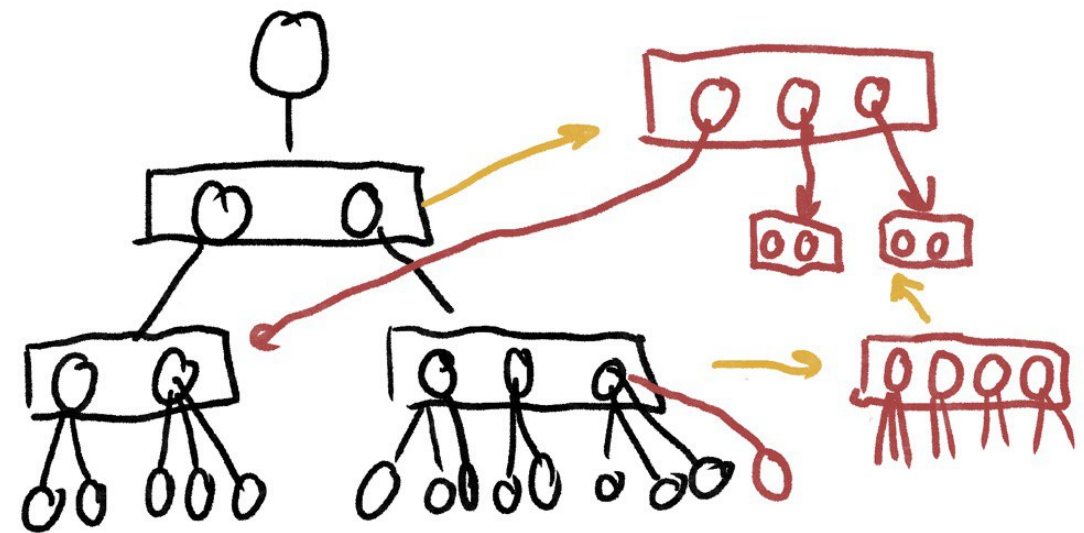
Когда добавляется значение, мы закидываем новый лист в нужное место дерева. Если это ломает инвариант предка, предка нужно разделить на две вершины. Если это ломает инвариант его предка, то рекурсивно решаем дальше.

2-3 Дерево

B+ дерево для $B = 2$.



Добавление



Мотивация для B-дерева

Базы данных обычно хранят данные на диске, куда лучше складывать данные последовательными блоками (условно, 16кб). Поэтому обработать 160 элементов в одном блоке быстрее, чем 10 элементов в разных блоках.

Получается ситуация, в которой выгодно сделать $B \sim 1000$. Тогда глубина дерева становится равна $O(\log_B n)$, а время ответа на запрос получается $O(B \log_B n) + k_{memory} \cdot O(\log_B n)$.

Спуск по дереву

Кроме того, чтобы делать стандартный `get` в структуре, иногда хочется найти первый элемент, который удовлетворяет заданному свойству.

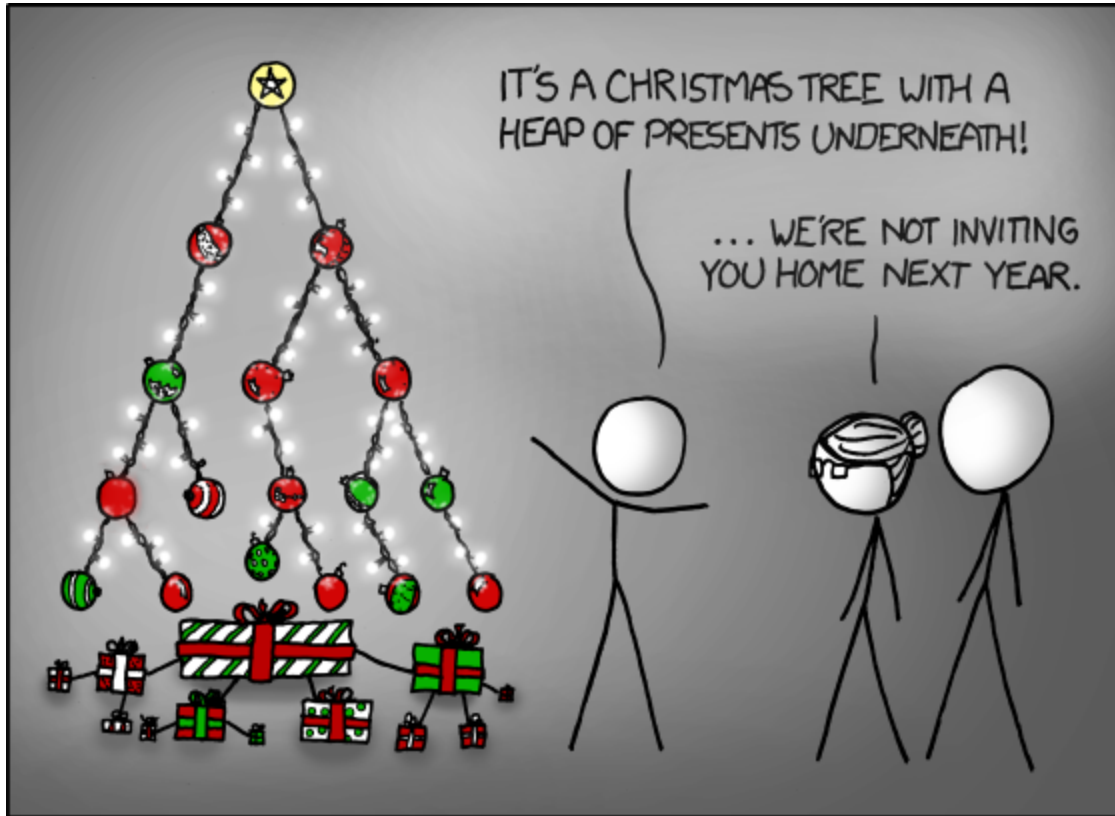
Тогда можно хранить в дереве некоторую дополнительную информацию, и идти в левое поддерево, если в его поддереве есть хотя бы одна подходящая вершина, и в правое поддерево, если такой вершины нет.

Применимо практически в любых деревьях

Спуск по дереву

```
def find_kth(node, k):  
    if node.l.count > k:  
        return find_kth(node.l, k)  
    if node.l.count + 1 == k:  
        return node  
    return find_kth(node.r, k - node.l.count - 1)
```

Mem



Идея пет-проекта

Можно написать с помощью B-дерева свою небольшую файловую систему. Можно, например, выписать все файлы в алфавитном порядке и использовать их полные пути как ключи. Затем можно сделать разбиение файлов на блоки так, чтобы размер блока был $B = 100$. Тогда найти любой файл в вашей файловой системе можно за $\log_{100} n$ обращений к хранилищу.

Поверх вашей файловой системы можно написать свою функцию `open()` и запустить на ней какую-то программу, которая работает с файлами.

```
./0_course_intro  
./0_course_intro/main.pdf  
./1_binary_search  
./1_binary_search/main.pdf  
./2_sorting  
./2_sorting/main.pdf  
./3_containers  
./3_containers/main.pdf
```

Что читать дальше

- AVL
- RB-tree
- Декартово дерево по неявному ключу
- https://en.wikipedia.org/wiki/B-tree#B-tree_usage_in_databases
- <https://ru.algorithmica.org/cs/tree-structures/treap/>