

# K-d 树数据结构

杨子凡

Jul 29, 2025

在高维数据查询领域，传统二叉树结构面临显著局限性。当数据维度升高时，二叉树无法有效组织空间关系，导致查询效率急剧下降。K-d 树「K-dimensional Tree」正是为解决这一挑战而生的数据结构。其核心思想是通过递归的轴对齐分割「axis-aligned splits」，将多维空间逐层划分。这种结构在最近邻搜索「KNN」、范围查询、空间数据库索引及计算机图形学「特别是光线追踪」等场景有广泛应用价值。

## 1 K-d 树基础理论

K-d 树中的维度参数  $k$  表示数据空间的维度数。每个节点包含四个关键属性：存储的数据点坐标、左子树指针、右子树指针及当前划分维度  $axis$ 。空间划分遵循特定规则：划分维度通常采用轮转策略「 $depth \bmod k$ 」或基于最大方差选择；划分点则选择当前维度上数据的中位数值，这是保证树平衡性的关键。例如在 2D 空间中，根节点按  $x$  轴分割，第二层节点按  $y$  轴分割，第三层再次回到  $x$  轴，如此递归形成空间划分。

## 2 K-d 树的构建算法

构建过程采用递归分割策略。以下 Python 实现展示了核心逻辑：

```

1 def build_kdtree(points, depth=0):
2     if not points:
3         return None
4     k = len(points[0]) # 获取数据维度
5     axis = depth % k # 轮转选择划分轴
6     points.sort(key=lambda x: x[axis]) # 按当前轴排序
7     median = len(points) // 2 # 确定中位索引
8
9     # 递归构建子树
10    return Node(
11        point=points[median],
12        left=build_kdtree(points[:median], depth+1),
13        right=build_kdtree(points[median+1:], depth+1),
14        axis=axis
15    )

```

此代码中，`depth` 参数控制维度的轮转切换。关键优化在于中位数选择：当数据量较大时，应采用快速选择算法「quickselect」将时间复杂度优化至  $O(n)$ 。对于重复点处理，可在排序时添加次要比较维度。在理想平衡状态下，树高为  $O(\log n)$ ，这是高效查询的基础。

### 3 K-d 树的查询操作

#### 3.1 范围搜索

范围搜索的核心是递归剪枝策略。从根节点开始，判断查询区域与当前节点划分平面的位置关系：若查询区域完全在当前点某一侧，则只需搜索对应子树；若跨越划分平面，则需搜索两侧子树。时间复杂度平均为  $O(\log n)$ ，最坏情况  $O(n)$ 。

#### 3.2 最近邻搜索

最近邻搜索采用递归回溯与剪枝策略：

---

```

1 def nn_search(root, target, best=None):
2     if root is None:
3         return best
4
5     axis = root.axis
6     # 选择初始搜索分支
7     next_branch = root.left if target[axis] < root.point[axis] else root.right
8     best = nn_search(next_branch, target, best)
9
10    # 更新最近点
11    curr_dist = distance(root.point, target)
12    if best is None or curr_dist < distance(best, target):
13        best = root.point
14
15    # 超球体剪枝判断
16    axis_dist = abs(target[axis] - root.point[axis])
17    if axis_dist < distance(best, target):
18        other_branch = root.right if next_branch == root.left else root.left
19        best = nn_search(other_branch, target, best)
20
21    return best

```

---

算法首先向下递归到叶节点「第 7 行」，回溯时更新最近邻点「第 11 行」。关键优化在于超球体剪枝「第 16 行」：若目标点到分割平面的距离小于当前最近距离，则另一侧子树可能存在更近邻点，需进行搜索。距离函数通常采用欧氏距离  $\sqrt{\sum_{i=1}^k (p_i - q_i)^2}$  或曼哈顿距离  $\sum_{i=1}^k |p_i - q_i|$ 。扩展 K 近邻搜索时，需使用优先队列维护候选集。

## 4 复杂度分析与性能考量

构建阶段时间复杂度取决于中位数选择策略：采用排序时为  $O(n \log n)$ ，使用快速选择可优化至  $O(n \log n)$ ；最坏未优化情况达  $O(n^2)$ 。查询操作在低维空间平均为  $O(\log n)$ ，但维度升高时出现「维度灾难」现象：当  $k > 10$ ，剪枝效率显著降低，最坏情况退化为  $O(n)$ 。这是因为高维空间中，超球体半径趋近于最远点距离，导致剪枝失效。此时可考虑 Ball Tree 或局部敏感哈希「LSH」等替代方案。

## 5 Python 完整实现示例

以下展示关键类的实现框架：

```
1 class KDNode:
2     __slots__ = ('point', 'left', 'right', 'axis')
3
4     def __init__(self, point, left=None, right=None, axis=0):
5         self.point = point # 数据点坐标
6         self.left = left # 左子树
7         self.right = right # 右子树
8         self.axis = axis # 划分维度索引
9
10    class KDTree:
11        def __init__(self, points):
12            self.root = self._build(points)
13
14        def _build(self, points, depth=0):
15            # 构建代码同前文
16
17        def range_search(self, bounds):
18            results = []
19            self._range_search(self.root, bounds, results)
20            return results
21
22        def _range_search(self, node, bounds, results):
23            if not node: return
24            # 检查当前节点是否在边界内
25            if all(bounds[i][0] <= node.point[i] <= bounds[i][1] for i in range(len(bounds))
26                  → )):
27                results.append(node.point)
28            # 递归剪枝逻辑
29            axis = node.axis
30            if bounds[axis][0] <= node.point[axis]:
```

```
29         self._range_search(node.left, bounds, results)
30     if bounds[axis][1] >= node.point[axis]:
31         self._range_search(node.right, bounds, results)
```

范围搜索通过边界框「bounds」进行区域过滤。可视化虽无法展示，但可通过 Matplotlib 绘制 2D 树的递归分割线及查询区域，直观展示空间划分过程。

## 6 优化与扩展方向

工程实践中，可通过节点缓存「caching」存储搜索路径，加速后续查询。对于高维数据，近似最近邻搜索「Approximate Nearest Neighbor」可牺牲少量精度换取显著性能提升。变种结构中，VP-Tree 采用球面空间划分，更适合非欧几里得空间；R\* 树则更适合动态数据场景。选择依据在于：当维度  $k < 10$  且数据静态时，K-d 树是理想选择；动态数据或超高维场景则需其他结构。

K-d 树通过「递归空间划分 + 回溯剪枝」机制，在低维空间实现了高效的范围查询与最近邻搜索。其优势在于结构简单、易于实现，但在高维场景存在退化风险。核心在于平衡构建与查询效率，理解轴对齐分割的几何意义。建议读者参考 GitHub 的完整实现进行实验，通过修改维度参数  $k$  直观观察维度灾难现象。