

АЛГОРИТМЫ УЛУЧШЕНИЯ КАЧЕСТВА R-ДЕРЕВЬЕВ

Рассматривается новый подход к работе со структурой для пространственного индексирования неточечных объектов в виде R-дерева, заключающийся в первоначальном глобальном построении эффективной структуры R-дерева и последующей работы с ней посредством обычных динамических алгоритмов. Предлагается глобальная стратегия построения R-дерева, сводимая к задаче разбиения множества прямоугольных объектов на K частей с минимальным взаимным пересечением. Предлагается три алгоритма разбиения: базовый, клеточный и «Разделяй и властвуй». Обсуждаются результаты экспериментального моделирования работы различных алгоритмов.

Введение

Задача графического поиска является одной из базовых задач вычислительной геометрии. Наиболее часто она возникает в машинной графике, геоинформатике, при решении различных задач, учитывающих влияние смежных объектов, и т.д. Так как форма различных объектов может значительно различаться, то для упрощения алгоритмов регионального поиска обычно для каждого из объектов вычисляется минимальный объемлющий прямоугольник со сторонами, параллельными осям координат. Тогда задача регионального поиска определяется как задача нахождения всех прямоугольных объектов, имеющих непустое перекрытие с прямоугольным регионом, являющимся запросом.

В настоящее время существует множество структур для организации данных, позволяющих решать данную задачу регионального поиска в той или иной мере. При этом наиболее эффективной структурой, позволяющей обрабатывать неточечные объекты при эффективном размещении данных во вторичной памяти, считаются R-деревья [1] и ряд их модификаций, в частности R*-деревья [2]. Практическое сравнение эффективности различных вариантов R-деревьев проведено в [3]. Наиболее быстродействующими из них являются R*-деревья, поэтому в дальнейшем в данной работе будет использоваться именно эта разновидность R-деревьев.

Индексная структура R-дерева является сбалансированным по высоте деревом с индексными записями в листьях, содержащими указатели на объекты данных. Узлам дерева соответствуют дисковые страницы, если индекс располагается на диске. Структура является полностью динамической, вставка и удаление объектов может свободно перемежаться с запросами на поиск, периодической реорганизации данных при этом не требуется.

Данная индексная структура предназначена для эффективного выполнения регионального поиска, поэтому основная идея дерева заключается в разбиении пространства в каждом узле дерева на минимально пересекающиеся группы, что позволяет при выполнении поиска эффективно отсекал неверные ветви работы алгоритма. Задача разбиения в общем случае, видимо, по меньшей мере является NP-полной, так как в настоящее время даже для двух групп не известно полиномиального алгоритма [1]. Поэтому во всех алгоритмах построения R-деревьев используются эвристические подходы, дающие на реальных данных далеко не оптимальные разбиения.

Высокая степень перекрытия входов в узлах дерева происходит из-за динамичности алгоритмов его построения, когда при вставке очередного элемента его необходимо за

минимальное время разместить и при необходимости перестроить дерево. При этом глобальная оптимизация дерева на каждом шаге не производится.

На практике взаимодействие с индексными структурами состоит из двух основных этапов: начального построения дерева и последующей оперативной работы. При этом во многих случаях во время второго этапа производится только поиск объектов без вставки новых и удаления старых.

В данной работе предлагается учесть специфику начального построения структуры, построив вначале более эффективное R-дерево, а на последующих этапах использовать обычные алгоритмы для работы с R-деревьями. В частном случае, когда вначале ничего не строится, получается обычное R-дерево.

1. Глобальные алгоритмы

Во всех существующих структурах R-деревьев основным параметром, определяющим скорость поиска, считается степень взаимного пересечения потомков в узлах дерева [1,2]. Именно поэтому в алгоритмах построения R-деревьев одной из главных целей является минимизация такого пересечения. Это позволяет при выполнении регионального поиска на ранних стадиях эффективно отсекал тупиковые ветви работы алгоритма.

В связи с этим можно ввести следующее определение:

Определение. Построенная для заданного набора объектов структура R-дерева называется оптимальной, если сумма площадей попарных перекрытий входов во всех нелистовых узлах дерева является минимальной среди всех возможных структур R-деревьев:

$$\min_{R\text{-деревья}} \sum_{i \neq j} S(R_i \cap R_j),$$

где S – функция площади.

Автором предлагается следующая стратегия построения R-дерева, близкого к оптимальному.

А л г о р и т м *Построить Дерево.* По заданному множеству из N объектов E_i строит R-дерево с параметрами m, M [1].

Шаг 1. *Начало.* Построить корень дерева и определить количество уровней дерева по формуле $L(N) = \lceil \log_m N \rceil$.

Шаг 2. *Построение дерева.* Вызвать алгоритм *Построить_Узел*, передав ему в качестве параметра корень дерева и всё множество объектов.

А л г о р и т м *Построить Узел.* По заданному множеству из N объектов E_i и узлу T строит потомков узла, имеющего уровень L ($L=1$ для листьев).

Шаг 1. *Построение листьев.* Если $L=1$, то заполнить узел объектами E_i и закончить.

Шаг 2. *Выбор количества потомков узла.* Вычислить количество потомков у данного узла $K = \lceil N^{1/L} \rceil$.

Шаг 3. *Построение узлов.* Вызвать алгоритм *Разделить_На_Группы* для разбиения всего множества объектов на K групп, минимизируя сумму попарных пересечений групп. При этом в алгоритм необходимо передать номер уровня L для правильной оценки минимального и максимального допустимого количества объектов в группах. Для каждой полученной группы построить пустой узел – потомок текущего узла и вызвать для него рекурсивно алгоритм *Построить_Узел* с этой группой в качестве множества.

Наиболее сложной частью предложенной автором стратегии является алгоритм

Разделить_На_Группы, от трудоёмкости которого зависит общая сложность алгоритма. В работе предлагается три варианта данного алгоритма.

2. Алгоритмы разбиения множества объектов на минимально пересекающиеся группы

Во все существующие алгоритмы построения R-деревьев входят алгоритмы разбиения множества объектов на две части с минимальным пересечением. При этом выбирается два базовых объекта, последовательными присоединениями к которым всех остальных строятся необходимые группы.

Обобщая данный алгоритм на случай нескольких групп, можно выбрать по одному базовому объекту для каждой группы и, используя аналогичные итеративные алгоритмы добавления оставшихся объектов, произвести разбиение. На этом принципе построен следующий базовый алгоритм разбиения.

А л г о р и т м *Разделить На Группы* (базовый алгоритм). Заданное множество из N объектов E_i делит на K групп с соблюдением ограничения на количество объектов в группах снизу m^{L-1} и сверху M^{L-1} , где L – текущий строящийся уровень дерева.

Шаг 1. Выбор базовых объектов для групп. Разбить всё множество объектов вертикальными и горизонтальными равноотстоящими линиями на K групп. После этого в качестве базовых для каждой группы выбрать по одному любому из объектов, попадающих в построенные клетки. Если в некоторые клетки не попадает ни одного объекта, то необходимо разделить любую из клеток с более чем двумя объектами пополам и взять из получившихся частей по объекту. Если базовых объектов по-прежнему недостаточно, то процесс деления клеток повторить. В качестве текущего размера групп взять размер соответствующих им клеток (на рис. 1,а приведён пример разбиения и выбора базовых объектов при $K=10$).

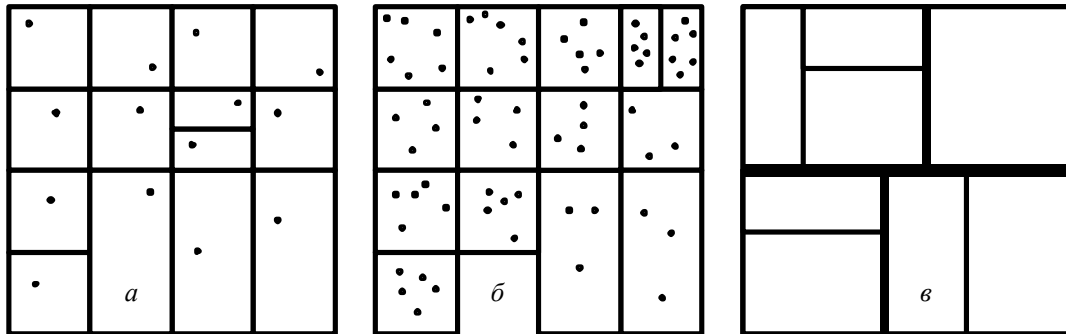


Рис. 1. Разбиение на клетки и выбор базовых объектов (а), клеточное разбиение (б), разбиение «Разделяй и властвуй» (в)

Шаг 2. Построение групп. Все оставшиеся объекты (не базовые) поместить в образовавшиеся группы по принципу минимального увеличения их площадей. При этом добавление объекта в группу не должно привести к числу объектов в группе больше чем M^{L-1} , а также к ситуации, когда оставшихся объектов недостаточно для обеспечения минимально допустимого числа объектов в других группах m^{L-1} .

Одним из недостатков работы такого алгоритма является не очень высокое качество разбиения на неравномерных распределениях, а также в случае, когда N незначительно превосходит K (при построении нижних уровней дерева). Это проявляется в ситуациях,

когда на втором шаге работы алгоритма некоторые группы уже заполнены, но при этом в соответствии с критерием выбора в них необходимо добавить очередной объект. Но так как это невозможно, то приходится добавлять объекты в некоторые другие группы, что в конечном итоге приводит к сильному перекрытию образовавшихся групп.

В качестве одного из вариантов решения данной проблемы автором предлагается временно удалять объекты, которые нельзя поместить в требуемые группы, из списка участвующих в глобальном алгоритме. При этом после окончания его работы все эти объекты необходимо поместить в дерево, используя уже обычные динамические алгоритмы построения R-дерева.

Как показало проведённое автором экспериментальное исследование работы данного алгоритма, количество объектов, которые нельзя разместить в требуемые группы, чаще всего незначительно и они эффективно размещаются по окончании работы глобального алгоритма динамическим способом.

В следующей теореме устанавливается трудоёмкость работы глобального алгоритма построения R-дерева, использующего предложенный базовый алгоритм разбиения.

Теорема 1 (трудоёмкость базового глобального алгоритма построения R-дерева). Пусть даны N объектов, на которых необходимо построить R-дерево. Тогда трудоёмкость глобального алгоритма построения R-дерева, использующего базовый алгоритм разбиения, составляет $O(N \log N)$.

Доказательство. Вначале установим трудоёмкость работы базового алгоритма разбиения. Сложность работы первого шага этого алгоритма состоит из сложности разбиения на клетки, которая составляет линейное время, а также сложности деления клеток, когда базовых объектов недостаточно. Если при этом в качестве базовых объектов случайно выбрать любые другие объекты, не являющиеся базовыми, то трудоёмкость всего первого шага будет $O(1)$.

Трудоёмкость работы второго шага алгоритма является не более чем линейной, так как в цикле по всем оставшимся объектам (не базовым) происходит сравнение с не более чем K группами, а $K \leq M = \text{const}$.

Таким образом, общая трудоёмкость базового алгоритма разбиения составляет $O(N)$.

Для глобального алгоритма построения R-дерева трудоёмкость работы всего алгоритма можно записать следующим образом:

$$S(N) = T(N) + K \cdot S(N / K), \tag{1}$$

где $S(N)$ – общая трудоёмкость глобального построения R-дерева; $T(N)$ – трудоёмкость алгоритма разбиения множества объектов.

Если расписать формулу общей трудоёмкости, то получится

$$S(N) = T(N) + K \cdot T(N / K) + K^2 \cdot T(N / K^2) + \dots + K^{\lceil \log_K N \rceil - 1} = \sum_{i=0}^{\lceil \log_K N \rceil - 1} K^i \cdot T(N / K^i). \tag{2}$$

Так как в случае базового алгоритма $O(T(N)) = O(N)$, то

$$O(S(N)) = O\left(\sum_{i=0}^{\lceil \log_K N \rceil - 1} K^i \cdot N / K^i\right) = O(N \cdot \lceil \log_K N \rceil) = O(N \cdot \log N), \tag{3}$$

что и требовалось доказать.

Как показало экспериментальное исследование, такой алгоритм хорошо ведёт себя на

равномерных распределениях непересекающихся объектов. Но на неравномерных распределениях часто происходит переполнение строящихся групп, что приводит к необходимости помещения очередных объектов в неоптимальные группы. В итоге получаются недопустимо большие перекрытия построенных групп. Поэтому в работе предлагается другой алгоритм, в котором вначале на основе клеточного разбиения плоскости строятся группы объектов и только затем анализируются их количество и наполнение. В случае необходимости полученные группы делятся пополам или объединяются с другими.

А л г о р и т м *Разделить На Группы* (клеточный алгоритм). Заданное множество из N объектов E_i делит на K групп с соблюдением ограничения на количество объектов в группах снизу m^{L-1} и сверху M^{L-1} , где L – текущий строящийся уровень дерева.

Шаг 1. *Построение базовых групп.* Разбить всё множество объектов вертикальными и горизонтальными линиями на $\lceil \sqrt{K} \rceil^2$ одинаковых клеток. После этого все объекты распределить по группам, соответствующим полученным клеткам, по принципу попадания центра объектов в построенные прямоугольники (рис. 1,б).

Шаг 2. *Уничтожение неполных групп.* Для каждой группы G_0 , в которой количество объектов меньше m^{L-1} , выполнить слияние с некоторой другой группой G_i . При этом необходимо выбрать такую группу G_i для слияния, чтобы площадь минимального объемлющего группы G_0 и G_i прямоугольника была минимальной среди всех возможных групп G_i .

Шаг 3. *Уничтожение лишних групп.* Пока общее количество групп больше K , выполнять следующее. Группу с минимальной площадью уничтожить, передав её объекты в некоторую другую группу по критерию, описанному в шаге 2.

Шаг 4. *Создание недостающих групп.* Пока общее количество групп меньше K , выполнять следующее. Группу с максимальным количеством объектов необходимо уничтожить, распределив её объекты по двум новым группам с минимальным пересечением полученных групп с помощью алгоритма *Расщепить_Узел*, входящего в состав классических алгоритмов для манипуляции с R-деревьями [1,2].

Шаг 5. *Распределение переполненных групп.* Для каждой группы G_0 , в которой количество объектов превышает M^{L-1} , выполнять шаг 6.

Шаг 6. *Перемещение объектов по группам.* Выбрать ближайшую (например, в смысле расстояний между центрами) к группе G_0 группу G_i , в которой количество объектов меньше M^{L-1} . После этого выбрать из группы G_0 такой объект g , что перемещение его в группу G_i вызовет минимальное увеличение площади последней группы. Переместить найденный объект g из группы G_i в G_0 .

Одним из недостатков работы такого алгоритма является увеличение времени работы алгоритма на неравномерных распределениях, так как при этом часто происходит недополнение и переполнение групп, а также приходится выполнять перемещение объектов по группам. В худшем случае алгоритм может иметь квадратичную сложность в случае многократного перемещения объектов по группам в пятом шаге.

В следующей теореме устанавливается трудоёмкость работы глобального алгоритма построения R-дерева, использующего предложенный автором клеточный алгоритм разбиения.

Теорема 2 (трудоёмкость клеточного глобального алгоритма построения R-дерева).

Пусть даны N объектов, на которых необходимо построить R-дерево. Тогда трудоёмкость глобального алгоритма построения R-дерева, использующего клеточный алгоритм разбиения, составляет $O(N^2)$.

Доказательство. Вначале установим трудоёмкость работы алгоритма клеточного разбиения. Его трудоёмкость состоит из трудоёмкости работы шагов 1-4 алгоритма, которая линейна ($T_{1-4}(N) = O(N)$), что следует из описания алгоритма при условии использования линейного по трудоёмкости алгоритма *Расщепить_Узел* [1], а также из трудоёмкости работы шагов 5-6. При выполнении шага 5 возможна ситуация, когда во всех группах, кроме одной, находится минимально допустимое количество объектов m^{L-1} . Если при этом необходимо разбить всё множество объектов на группы по M^{L-1} объектов (в случае, если $N = K \cdot M^{L-1}$), то в одной из групп находится $K \cdot M^{L-1} - (K-1) \cdot m^{L-1}$ объектов. Поэтому на шестом шаге алгоритма будет перемещено всего $R(N) = (K-1) \times (M^{L-1} - m^{L-1})$ объектов. Так как $m \leq M/2$ и $N \leq K \cdot M^{L-1}$, то общее количество перемещаемых объектов становится в худшем случае

$$\begin{aligned} R(N) &= (K-1) \cdot (M^{L-1} - m^{L-1}) \geq (K-1) \cdot (M^{L-1} - \frac{M^{L-1}}{2}) = \\ &= (K-1) \cdot \frac{M^{L-1}}{2} = \frac{K-1}{2K} K \cdot M^{L-1} \geq \frac{K-1}{2K} N. \end{aligned} \quad (4)$$

С другой стороны, так как $N > (K-1) \cdot M^{L-1}$, то

$$R(N) = (K-1) \cdot (M^{L-1} - m^{L-1}) \leq (K-1) \cdot M^{L-1} < N.$$

Таким образом, $O(R(N)) = O(N)$. А так как в шестом шаге для каждого перемещаемого объекта выполняется цикл по некоторой группе G_i , то общая трудоёмкость работы шагов 5-6 может составлять:

$$O(T_{5-6}(N)) = O(R(N) \cdot M^{L-1}) = O\left(R(N) \cdot \frac{1}{K} K \cdot M^{L-1}\right) = O(N) \cdot O(K \cdot M^{L-1}) = O(N^2). \quad (5)$$

Таким образом, общая трудоёмкость клеточного алгоритма разбиения составляет

$$O(T(N)) = O(T_{1-4}(N)) + O(T_{5-6}(N)) = O(N) + O(N^2) = O(N^2).$$

Для установления общей трудоёмкости клеточного глобального алгоритма построения R-дерева воспользуемся формулой (2):

$$\begin{aligned} O(S(N)) &= O\left(\sum_{i=0}^{\lceil \log_K N \rceil - 1} K^i \cdot T(N/K^i)\right) = O\left(\sum_{i=0}^{\lceil \log_K N \rceil - 1} K^i \cdot O(N^2/K^{2i})\right) = \\ &= \sum_{i=0}^{\lceil \log_K N \rceil - 1} O(N^2/K^i) = O(N^2) \cdot \left(1 + \frac{1}{K} + \frac{1}{K^2} + \dots + \frac{1}{K^{\lceil \log_K N \rceil - 1}}\right) = O(N^2 \cdot F(K)). \end{aligned}$$

Так как $2 \leq m \leq K$, то

$$F(K) = 1 + \frac{1}{K} + \frac{1}{K^2} + \dots + \frac{1}{K^{\lceil \log_K N \rceil - 1}} \leq 1 + \frac{1}{K} + \frac{1}{K^2} + \dots = \frac{1}{1 - \frac{1}{K}} \leq \frac{1}{1 - \frac{1}{2}} = 2.$$

С другой стороны, $F(K) = 1 + \frac{1}{K} + \frac{1}{K^2} + \dots + \frac{1}{K^{\lceil \log_K N \rceil - 1}} \geq 1$. Поэтому

$$O(S(N)) = O(N^2 \cdot F(K)) = O(N^2), \quad (6)$$

что и требовалось доказать.

Теперь рассмотрим другой вариант алгоритма разбиения, основанный на стратегии «Разделяй и властвуй», в котором производится последовательное разбиение всего множества на две части до тех пор, пока всё множество не окажется разбито на требуемое число частей.

А л г о р и т м *Разделить На Группы* (алгоритм «Разделяй и властвуй»). Заданное множество из N объектов E_i делит на K групп с соблюдением ограничения на количество объектов в группах снизу m^{L-1} и сверху M^{L-1} , где L – текущий строящийся уровень дерева.

Шаг 1. *Разделение на две группы*. Вызвать алгоритм *Разделить На Две Группы*, передав в качестве аргументов всё множество объектов E_i и соотношение деления $(\lfloor K/2 \rfloor) : (K - \lfloor K/2 \rfloor)$ для получения групп G_1 и G_2 (рис. 1, в).

Шаг 2. *Рекурсивное разделение*. Если $\lfloor K/2 \rfloor > 1$, то вызвать алгоритм *Разделить На Группы*, передав в качестве аргументов полученную группу G_1 с параметром деления $\lfloor K/2 \rfloor$, иначе выдать в качестве результата G_1 .

Если $K - \lfloor K/2 \rfloor > 1$, то вызвать алгоритм *Разделить На Группы*, передав в качестве аргументов полученную группу G_2 с параметром деления $K - \lfloor K/2 \rfloor$, иначе выдать в качестве результата G_2 .

А л г о р и т м *Разделить На Две Группы*. Заданное множество из N объектов E_i делит на 2 группы в заданном соотношении площадей $k:l$ с соблюдением ограничения на количество объектов снизу $k \cdot m^{L-1}$ и сверху $k \cdot M^{L-1}$ для первой группы, а также снизу $l \cdot m^{L-1}$ и сверху $l \cdot M^{L-1}$ для второй, где L – текущий строящийся уровень дерева.

Данный алгоритм аналогичен алгоритму *Расщепить Узел* обычного R-дерева, за исключением проверок, при которых сравниваемые значения взвешиваются по величинам k и l .

Шаг 1. *Выбор оси координат для сортировки*. Выбрать такую ось координат i , на которой достигается максимальное значение выражения

$$\max_{\text{оси } i} \left(\frac{\max_{\text{фигуры } j} R_a^{j,i} - \min_{\text{фигуры } j} R_b^{j,i}}{\max_{\text{фигуры } j} R_b^{j,i} - \min_{\text{фигуры } j} R_a^{j,i}} \right),$$

где $R_a^{j,i}$ и $R_b^{j,i}$ – соответственно меньшее и большее значение i -й координаты минимального объемлющего j -ю фигуру прямоугольника.

Шаг 2. *Разделение по группам*. Разделить все объекты на три части по значению $P^i(j) = (R_a^{j,i} + R_b^{j,i})/2$ для выбранного i по квантилям уровня α и $1-\alpha$ ($\alpha \in [0; 0,5]$ – параметр деления по группам): из наименьшей по значениям части образовать первую группу объектов, а из наибольшей – вторую группу. Оставшиеся объекты распределить по группам в шагах 3 и 4.

Шаг 3. *Проверка завершения*. Если все объекты распределены, то закончить. Если одна из групп в результате дальнейших операций будет недополнена (количество объектов в группах $p_1 < k \cdot m^{L-1}$ или $p_2 < l \cdot m^{L-1}$), то вставить в неё все оставшиеся объекты и закончить. Если одна из групп в результате дальнейших операций будет переполнена (количество объектов в группах $p_1 > k \cdot M^{L-1}$ или $p_2 > l \cdot M^{L-1}$), то вставить в другую группу все оставшиеся объекты и закончить.

Шаг 4. *Распределение объектов*. Взять любой нераспределённый объект и поместить его в группу, размер которой увеличится в минимальной степени при взвешивании по k и l . Перейти на шаг 3.

В следующей теореме устанавливается трудоёмкость работы глобального алгоритма

построения R-дерева, использующего предложенный алгоритм разбиения «Разделяй и властвуй».

Теорема 3 (трудоемкость глобального алгоритма построения R-дерева «Разделяй и властвуй»). Пусть даны N объектов, на которых необходимо построить R-дерево. Тогда трудоемкость глобального алгоритма построения R-дерева, использующего алгоритм разбиения «Разделяй и властвуй», составляет $O(N \cdot \log^2 N)$.

Доказательство. Вначале установим трудоемкость работы алгоритма разбиения «Разделяй и властвуй». Его трудоемкость можно записать следующим образом:

$$T(N) = R(N) + 2 \cdot T(N/2), \quad (7)$$

где $R(N)$ – трудоемкость алгоритма разбиения множества объектов на две части.

В алгоритме деления на две части наиболее трудоемким этапом является второй шаг, в котором производится деление множества объектов на три части по заданным квантилям. Если использовать линейные по сложности алгоритм деления по квантилям Хоара или алгоритм цифровой сортировки, то получается $O(R(N)) = O(N)$, так как другие шаги обсуждаемого алгоритма не более чем линейны, что следует из описания алгоритма.

Таким образом, общая трудоемкость алгоритма разбиения «Разделяй и властвуй» составляет

$$O(T(N)) = O(N) + 2 \cdot O(T(N/2)) = O(N) + 2 \cdot (O(N/2) + 2 \cdot O(T(N/4))) = O(N) + 2 \cdot (O(N/2) + 2 \cdot (O(N/4) + 2 \cdot O(T(N/8)))) = \dots = O(N) \cdot O(\lceil \log_2 N \rceil) = O(N \log N). \quad (8)$$

Для установления общей трудоемкости глобального алгоритма построения R-дерева «Разделяй и властвуй» воспользуемся формулой (2):

$$\begin{aligned} O(S(N)) &= O\left(\sum_{i=0}^{\lceil \log_K N \rceil - 1} K^i \cdot T(N/K^i)\right) = O\left(\sum_{i=0}^{\lceil \log_K N \rceil - 1} K^i \cdot O(N/K^i \cdot \log(N/K^i))\right) = \\ &= \sum_{i=0}^{\lceil \log_K N \rceil - 1} O(N \cdot \log N) = O(N \cdot \log N) \cdot O(\lceil \log_K N \rceil) = O(N \cdot \log^2 N), \end{aligned} \quad (9)$$

что и требовалось доказать.

На втором шаге работы алгоритма деления «Разделяй и властвуй» присутствует параметр α деления по группам. Для установления его влияния на качество получаемого разбиения автором было проведено экспериментальное моделирование работы глобального алгоритма построения R-дерева, использующего алгоритм разбиения «Разделяй и властвуй». При этом было установлено, что уменьшение значения α приводит к некоторому улучшению структуры R-дерева на неравномерных распределениях и распределениях со значительным перекрытием объектов. С другой стороны, это приводит к снижению скорости работы алгоритма и уменьшению качества разбиения на точечных и регулярных наборах данных.

Суммируя достоинства и недостатки поведения алгоритма при разных значениях α на различных распределениях объектов, автор выбрал эмпирическое значение параметра деления для универсального применения, равное 0,45. При этом если на практике имеются некоторые сведения о реальном распределении объектов, то значение параметра α , вероятно, стоит оценить дополнительно в зависимости от требований конкретной задачи.

3. Практический анализ глобальных алгоритмов

Для количественного анализа построенных алгоритмов автором было проведено мо-

делирование работы различных алгоритмов построения R-деревьев на различных распределениях объектов и поисковых регионах. При этом анализировались такие характеристики, как время построения R-дерева, процент попадания при поиске и процент перекрытия потомков [1]. Для сравнения алгоритмов в различных условиях в экспериментах строились наборы данных с 10 000 объектами, располагавшимися в единичном квадрате $[0, 1]^2$, с 8 базовыми типами распределений, приведенными в [2].

Сравнение экспериментальных данных подтверждает правильность выбранной стратегии на минимизацию взаимных перекрытий потомков узлов в R-дереве, так как практически везде меньшему проценту перекрытия потомков соответствует больший процент попадания в нужные узлы при поиске.

По результатам моделирования можно сделать вывод, что предложенные в работе базовый и клеточный алгоритмы наиболее хорошо работают на равномерных распределениях, когда объекты мало пересекаются между собой. Но на неравномерных распределениях заметно ухудшение работы алгоритмов по сравнению с существующими вариантами R-деревьев. В то же время алгоритм «Разделяй и властвуй» одинаково хорош как на равномерных, так и на неравномерных распределениях. Практически во всех тестах он обходит остальные алгоритмы по качеству получающегося построения. Только при построении R-дерева на одном наборе данных (очень большие объекты с сильным перекрытием) этот алгоритм строит дерево, уступающее по проценту перекрытия классическому алгоритму построения R-дерева.

СПИСОК ЛИТЕРАТУРЫ

1. Guttman A. // Proc. ACM SIGMOD Int. Conf. on Management of Data. – 1984. – P. 47-57.
2. Beckmann N., Kriegel H., Schneider R., Seeger B. // ACM. – 1990. – P. 322-331.
3. Kriegel H., Schiwietz M., Schneider R., Seeger B. // Proc. Symp. On the Design and Implementation of Large Spatial Databases, Santa Barbara. – 1989. – P. 413-418.