

Алгоритмы. Теор.листок 2

Артемий Клячин

28 октября 2021 г.

Задача 1.

Построим общее сжатое суффиксное дерево *sufftree* для строк S, S^{-1} за $O(n)$.

$LCP(S, T)$ - наибольшая общий префикс строк S и T .

$LCA(S, T)$ - наименьший общий предок вершин соответствующих строкам S и T .

Пусть $LCE(S, T, i, j) = LCP(S[i:], T[j:])$

$LCE(S, T, i, j)$ определяется меткой LCA на *sufftree* для $S[i:]$ и $T[j:]$. Т.е. для листов суффиксов $S[i:]$ и $T[j:]$ надо найти ближайшего общего предка, а по глубине предка узнает длину наибольшего общего префикса.

Есть способ находить LCA за $O(1)$ с препроцессингом за $O(N)$ (алгоритм Фарах-Колтона и Бендера). Таким образом мы можем делать LCE за $O(1)$.

$LCE(S, S^{-1}, i, (n-1)-i)$ - это длина наибольшего палиндрома нечётной длины с центром в i .

$LCE(S, S^{-1}, i+1, (n-1)-i)$ - это длина наибольшего палиндрома чётной длины с центром в i и $i+1$.

Таким образом за $O(n)$ мы получим ответ на задачу.

Итоговая асимптотика $O(n)$.

Задача 2.

Введём понятия:

Строка y *сильно* продолжает строку x , если после каждого вхождения x в y следует вхождение y .

Строка y *очень слабо* продолжает строку x , если y слабо и несильно продолжает строку x .

Построим сжатое суффиксное дерево.

Возьмём произвольный x - ему соответствует какая-то позиция на каком-то ребре.

Количество сильных продолжений - это расстояние в символах до ближайшей вершины дерева. (далее будет разветвление и сильного продолжения не существует)

Далее посчитает очень слабые продолжения. Если мы подсчитаем их для вершины v , то их количество будет таким же для всех позиций на ребре (p, v) (кроме самой вершины p), где p - вершина-предок вершины v .

Так как любые два очень слабых продолжения x либо совпадают, либо один является префиксом другого, то очень слабые продолжения находятся на одном пути в некоторый лист u .

Другие продолжения не могут быть такой же длины, что и очень слабое продолжение, следовательно, посимвольная глубина позиции, в которую попадает очень слабое продолжение должна быть больше, чем у любого другого продолжения не на пути к u .

Любое подобное продолжение на пути к u будет очень слабым.

Следовательно, количество очень слабых продолжений x будет $h_1 - h_2$, где h_1 и h_2 - длины двух самых глубоких невложенных продолжения ($h_1 \geq h_2$).

Алгоритм:

за $O(N)$ строим дерево, за $O(N)$ при помощи DFS ищем 2 самых больших продолжения для каждой вершины, за $O(N)$ обходим все ребра и для каждого x , за $O(1)$ вычисляем количество слабых (сильных и очень слабых) продолжений.

e - ребро сжатого суффиксного дерева.

$len(e)$ - количество x на ребре e .

v_e - более глубокая вершина e

Ответ: $\sum_e \frac{len(e) \cdot (len(e)-1)}{2} + m \cdot (h_1(v_e) - h_2(v_e))$

Задача 3.

Построить суффиксное дерево $s_1\$1s_1\$2...s_ns_n\$n$, где $\$, ..., \n - разные окончания.

Далее каждой строкой пройдем по дереву. Если для строки s_i есть продолжение отличное от $\$,$ то для s_i ответ YES, иначе NO.

Оценка асимптотики:

1) Построение дерева - $O(\sum_{i=1}^n |s_i|)$

2) Проход всех строк по дереву происходит суммарно за $O(\sum_{i=1}^n |s_i|)$.

Итоговая асимптотика: $O(\sum_{i=1}^n |s_i|)$

Задача 4.

Построим общее сжатое суффиксное дерево для $s_1, s_1, ..., s_n$. В каждой вершине дерева будем хранить список строк, для которых данная вершина является терминальной. Такие списки будем называть списками терминальных состояний.

Далее каждой строкой s_i пройдемся по дереву и обработаем для каждой вершины списки терминальных состояний. Если для строки s_j , вершина, в которую мы попали - терминальная, то существует суффикс s_j , который является префиксом s_i . И наоборот, если p строка является префиксом s_i и суффиксом s_j , то она обязательно попадет при обходе. Так как нас интересуют максимальные такие p , то будем искать последнее вхождение строки s_j в списки терминальных состояний.

По итогу мы получили по $n - 1$ значений для каждой строки s_i . Которые и будут ответом на задачу.

Оценка асимптотики:

1) Построение дерева - $O(n)$

2) Пусть $m = \sum_{i=1}^n |s_i|$.

Для каждой строки производим спуск вниз по дереву и обрабатываем терминальные списки. Спуск производится за $O(m)$. Для каждой строки s_i , количество терминальных вершин $|s_i|$, следовательно, длины всех терминальных списков m , поэтому суммарно для всех $s[i]$ их обработка займет $O(m)$.

3) По итогу алгоритма имеется n списков по $n - 1$ значению. Ответ выводим за $O(n^2)$

Итоговая асимптотика: $O(\sum_{i=1}^n |s_i| + n^2)$

Задача 5.

Количество различных подстрок в строке - это сумма длин всех ребер в сжатом суффиксном дереве +1 (пустое подслово).

Решим задачу с помощью алгоритма Укконена. Алгоритм Укконена постепенно строит дерево: в начале строит дерево для $s[0, 0]$, потом $s[0, 1]$, ... $s[0, n - 1]$. Дополнительно, на каждой итерации, считаем сумму длин ребер суфф.дерева для $s[0, i]$, используя результат для предыдущей итерации: $i - 1$

Так как асимптотика Укконена $O(|s|)$, то асимптотика алгоритма $O(|s|)$.

Задача 6.

Построим суффиксный автомат по s . Тогда все подстроки, правые контексты которых совпадают, будут соответствовать только одной вершине. А каждой вершине будут соответствовать только подстроки с одинаковой правой контексткой.

Любому суффиксу s соответствует единственный путь в автомате (с терминальной вершиной на конце), а любому пути в автомате с терминальной вершиной на конце соответствует единственный суффикс s .

Пусть u - некоторый путь в графе. Введём понятие *префикс пути u* - любой путь графа с той же начальной вершиной, что начальная вершина u . Путь в точности совпадает с u , но его конец находится не дальше, чем конечная вершина u .

Проще говоря, префикс пути - это аналог префикса строки.

Пусть $top(x)$ - вершина соответствующая строке x в построенном суффиксном автомате.

Пусть $path(x)$ - путь соответствующий строке x в построенном суффиксном автомате.

Пусть x - произвольный суффикс s . Для любого x_1 - префикса x выполняется: $path(x_1)$ - префикс $path(x)$. А для любой вершины v на пути $path(x)$ найдётся x_1 - префикс x , такой что $top(x_1) = v$.

Рассмотрим x и y - суффиксы s . Исходя из фактов описанных выше, можно сделать вывод, что условие: x и y несравнимы, и условие: $path(x)$ и $path(y)$ не пересекаются по вершинам - эквивалентны.

Следовательно, чтобы решить задачу, необходимо построить как можно больше непересекающихся путей, заканчивающихся в терминальной вершине.

Для этого запустим DFS по суффиксному автомату и найдём все такие терминальные вершины, до которых есть путь без терминальных вершин. Другие терминальные вершины нас не интересуют, так как всегда вместо пути до такой вершины можно взять его префикс, который заканчивается в терминальной вершине.

Теперь у нас есть все терминальные вершины, которые могут быть концами искомых путей. Выкинем все рёбра исходящие из этих терминальных вершин. Понятно, что теперь все тупиковые вершины этого графа будут терминальными. Введём новую вершину (стоковую) и соединим терминальные вершины с ней.

Сведём задачу к поиску максимального потока. В ней потоки величины 1 будут искомыми путями. Для этого зададим пропускную способность всех рёбер - 1. Нам необходимо сделать так, чтобы в каждую вершину (кроме начальной и стоковой) приходил поток величины не более чем 1. (иначе пути пересекутся) Раздвоим все вершины (кроме начальной и стоковой). В одну проведём все входящие рёбра (входившие в прежнюю вершину), а из другой проведём все выходящие рёбра. Обе вершины соединим ребром с потоком 1.

Полученный максимальный поток - как раз и будет состоять из искомых путей. Потоки не будут пересекаться. Если предположить, что возможно было построить больше непересекающихся путей заканчивающихся в терминальных вершинах, то тогда можно было бы сделать больше непересекающихся потоков, чем при решении задачи о максимальном потоке, что противоречит его максимальной. Следовательно, найденное решение - максимально.

Максимальный поток состоит из искомых путей. Величина максимального потока - это искомое количество путей. Найденное множество путей заканчивающихся в терминальных вершинах соответствует искомому множеству суффиксов.

Нашли искомое.

Асимптотика:

Построение суфф.автомата: $O(n)$, где $n = len(s)$

Вершин в графе $O(n)$.

Рёбер в суфф.автомате $O(n^2)$

DFS с поиском терминальных вершин: $O(n^2)$

Видоизменение графа (раздваивание вершин): $O(n^2)$

Поиск максимального потока, когда все рёбра имеют единичную пропускную способность: $O(n^2)$

Итоговая асимптотика: $O(n^2)$