

Алгоритм Решения Задачи № 16

Задача: Даны α - регулярное выражение и слово $u \in \{a, b, c\}^*$. Найти длину самого длинного подслова u , являющегося также подсловом некоторого слова в языке L .

Решение: основным юнитом решения является структура *Language*, которая хранит 4 верхнетреугольные булевы матрицы размера $(N + 1) \times (N + 1)$, где $N = \text{len}(u)$.

1. *has_prefix_(i, j)* - содержит ли *Language* слово, для которого $[i, j]$ (подслово u , начинающееся на i -том индексе и заканчивающееся на j -1-ом) является префиксом.
 2. *has_suffix_(i, j)* - содержит ли *Language* слово, для которого $[i, j]$ является суффиксом.
 3. *has_whole_(i, j)* - содержит ли *Language* слово, которое совпадает с $[i, j]$.
 4. *has_substr_(i, j)* - содержит ли *Language* слово, для которого $[i, j]$ является подсловом.
- Во время парсинга регулярного выражения будем вычислять эти предикатные матрицы.

Когда рег.выражение состоит из одного символа, очевидно как заполнять эти матрицы. Они все единичные, если в рег.выражении только "1". Если регулярка не состоит из "1", то на главной диагонали всех матриц, кроме *has_whole_* (ее диагональ нулевая, т.к. когда рег.выражение состоит не из "1", в L не будет пустого слова) лежат единички, а также если i -тый символ входного слова то $(i, i + 1)$ элементы каждой матрицы -- единички.

Для удобства записи, пусть $(\text{has_substr_}, \text{has_whole_}, \text{has_prefix_}, \text{has_suffix_}) = (ss, w, p, s)$. Итак, мы поняли как выглядят предикатные матрицы для языков задаваемых рег.выражениями длины

1. Опишем, вычислять предикатные матрицы для объединения, произведения и итерации языков. Мы знаем, что:

- 1) $\text{Language}(a + b) = \text{Language}(a) \cup \text{Language}(b)$;
- 2) $\text{Language}(a.b) = \text{Language}(a).\text{Language}(b)$;
- 3) $\text{Language}(a^*) = [\text{Language}(a)]^*$

Достаточно очевидно, что, если подслово $[i, j]$ префикс некоторого слова в языке $\text{Language}(a) \cup \text{Language}(b)$, то оно является префиксом некоторого слова в $\text{Language}(b)$ или в $\text{Language}(a)$. В случае суффикса, совпадения $[i, j]$ с некоторым словом в $\text{Language}(a) \cup \text{Language}(b)$ или подсловости рассуждения аналогичны. Значит, чтобы вычислить предикатные матрицы для языка $\text{Language}(a) \cup \text{Language}(b)$ достаточно просто поэлементно применять операцию ИЛИ для соответствующих матриц языков $\text{Language}(a)$ и $\text{Language}(b)$.

А также, очевидно, что:

$$[i, j] \text{ --- префикс некоторого слова в } \text{Language}(a).\text{Language}(b) \iff (\exists k \in N, i \leq k \leq j, \text{Language}(a).w(i, k) = \text{true} \wedge \text{Language}(b).p(k, j) = \text{true}) \quad (1)$$

Значит:

$$\text{Language}(a.b).p = \text{Language}(a).p + \text{Language}(a).w * \text{Language}(b).p \quad (2)$$

В умножении матриц переопределим оператор $+$ как дизъюнкцию ИЛИ, а умножение двух булевых значений -- конъюнкция. Тогда, формула (2) действительно соответствует утверждению (1). Умножение матриц $\text{Language}(a).w * \text{Language}(b).p$ означает, что мы рассматриваем все k , для всех i и j , что $i \leq k \leq j$ и $\text{Language}(a).w(i, k) = \text{true} \wedge \text{Language}(b).p(k, j) = \text{true}$.

Аналогичными рассуждениями выводим, что:

$$\text{Language}(a.b).s = \text{Language}(b).s + \text{Language}(a).s * \text{Language}(b).w \quad (3)$$

$$\text{Language}(a.b).w = \text{Language}(a).w * \text{Language}(b).w \quad (4)$$

$$\text{Language}(a.b).ss = \text{Language}(a).ss + \text{Language}(b).ss + \text{Language}(a).s * \text{Language}(b).p \quad (5)$$

Итак, научились вычислять предикатные матрицы для объединения и умножения языков.

Осталось вычислить их для итерации языка.

рассмотрим $\text{Language}(a^*).w$

$$\text{Language}(a^0).w = \text{Language}(1).w = E$$

$$\text{Language}(a^2).w = (\text{Language}(a).w)^2$$

.....

$$\text{Language}(a^n).w = (\text{Language}(a).w)^n$$

$$\text{Language}(a^*).w = \text{Language}(a^0).w + \text{Language}(a).w + \text{Language}(a^2).w + \dots +$$

$$+ \text{Language}(a^n).w + \dots = E + \text{Language}(a).w + (\text{Language}(a).w)^2 + \dots +$$

$$+ (\text{Language}(a).w)^n + \dots \quad (6)$$

Если, рассмотреть предикатную матрицу A , как ориентированный граф, то согласно утверждению из теории графов A^n - все пути длины n . Значит матрица в формуле (6) имеет все пути длины от 0, до максимально возможной, т.е. это транзитивное замыкание графа $\text{Language}(a).w$. Это значит, чтобы вычислить матрицу $\text{Language}(a^*).w$ достаточно транзитивно замкнуть $\text{Language}(a).w$. Для вычисления остальных матриц воспользуемся $\text{Language}(a^*).w$. Если, слово - префикс некоторого слова в $\text{Language}(a^*)$, то некоторый префикс этого слова может целиком лежать в языке $\text{Language}(1) \cup \text{Language}(a) \cup \text{Language}(a^2) \cup \dots = \text{Language}(a^*)$, а оставшая часть этого слова является префиксом какого-то слова в $\text{Language}(a)$. Из этого получаем естественную формулу для $\text{Language}(a^*).p$:

$$\text{Language}(a^*).p = \text{Language}(a^*).w * \text{Language}(a^*).p$$

Аналогичными рассуждениями выводим формулы для s и ss :

$$\text{Language}(a^*).s = \text{Language}(a).s * \text{Language}(a^*).w$$

$$\text{Language}(a^*).ss = \text{Language}(a).s * \text{Language}(a^*).w * \text{Language}(a).p + \text{Language}(a).ss$$

Итак, научились вычислять предикатные матрицы для любого рег.выражения. Теперь, используя эти матрицы посчитаем ответ задачи.

После того, как тривиальными операциями push/pop в стек распарсили входное рег.выражение и одновременно вычислили все 4 предикатные матрицы, итерируемся по $L(\alpha).ss$ и находим в ней значение true с максимальной разницей индекса столбца и индекса строки, т.к. эта разница и дает нам длину максимального подслова удовлетворяющего условию задачи.

Оценим работу нашего алгоритма. пусть N -- длина входного слова u и R -- длина входного рег.выражения α . Тогда, каждое каждое перемножение матриц делается за $O((N+1)(N+1)(N+1)) = O(N^3)$. Сложение за $O(N^2)$, последняя итерация по матрице `has_substr` для поиска ответа также за $O(N^2)$, $O(N^3 + N^2) = O(N^3)$, а такие операции выполняются максимально R раз \Rightarrow общая сложность алгоритм составляет $O(R * N^3)$.