

Задание (вариант 18):

Даны регулярное выражение в обратной польской записи R и буква x . Найти максимальное k , такое что в языке L есть слова, начинающиеся с x^k . Алфавит: $\{a, b, c, 1, ., +, *\}$

Алгоритм: Для каждого регулярного выражения r , с которым работаем во время перевода регулярного выражения из обратной польской записи в обычную запись, будем хранить два параметра:

$$\begin{cases} \max_pref_len \in [0, inf] & \text{т.е. } \max k : x^k - \text{префикс } w \in r \\ \max_word_len \in [-1, inf] & \text{т.е. } \max m : x^m = w' \in r \end{cases}$$

В написанном в этом репозитории алгоритме $inf = 10^5$. Очевидно, что всегда верно $\max_pref_len \geq \max_word_len$. Поянсим границы для каждого из параметров. Для всех примеров положим $x = a$.

Параметр \max_pref_len принимает значения:

- 0, если r не задает слово, в котором есть префикс, состоящий из буквы x . (Например, $r = ba$)
- inf , если r задает слово, в котором есть префикс из буквы x , размер которого превышает заранее зафиксированное число. (Например, $r = a^*b$)

Параметр \max_word_len принимает значения:

- -1 , если r не задает слово, полностью состоящее из букв x . (Например, $r = ab$)
- 0, если r не задает слово, полностью состоящее из букв x , но при этом задает пустое слово, то есть 1. (Например, $r = ab + 1$)
- inf , если r задает слово, полностью состоящее из букв x , размер которого превышает заранее зафиксированное число. (Например, $r = a^*$)

Тогда база будет вот такой (r – однобуквенное выражение, то есть $r = letter$):

- | | |
|--------------------------|-----------------------------|
| • Если $letter = x$, то | • Если $letter \neq x$, то |
| – $\max_pref_len = 1$ | – $\max_pref_len = 0$ |
| – $\max_word_len = 1$ | – $\max_word_len = -1$ |

Теперь разберем, каким образом определяются значения $r_{\max_pref_len}$ и $r_{\max_word_len}$ для $r = r_1 \# r_2$, где $\#$ – это одна из операций $\{+(2), .(2), *(1)\}$, а r_1, r_2 – корректные регулярные выражения, для которых уже посчитаны $r_{i\max_pref_len}$ и $r_{i\max_word_len}$ $i \in [1, 2]$. Будем называть $w \in r$ – слово с максимальным префиксом или максимальное слово, полностью состоящее из букв x .

- Если это операция $+(2)$, то в терминах регулярных выражений для получения слова w выбирается лучшее слово из r_1 или r_2 . Так как параметры независимые, то

- $r_{\max_pref_len} = \max(r_{1\max_pref_len}, r_{2\max_pref_len})$
- $r_{\max_word_len} = \max(r_{1\max_word_len}, r_{2\max_word_len})$

- Если это операция $.(2)$, то в терминах регулярных выражений для получения слова w склеиваются слова из r_1 и r_2 . Тогда, возможны разные случаи:

- Если r_1 задает слово, полностью состоящее из букв $x \Leftrightarrow r_{1max_word_len} \neq -1$, то $r_{max_pref_len}$ высчитывается как максимум из {склеивания всего слова из букв x из r_1 и префикса из r_2 } или {префикса из r_1 }, следовательно формула такая:

$$r_{max_pref_len} = \max(r_{1max_pref_len}, r_{1max_word_len} + r_{2max_pref_len})$$
- Иначе, префикс r_2 не имеет смысла рассматривать, так как в r_1 после префикса вида $x^{r_{1max_pref_len}}$ идут другие буквы, и тогда $r_{max_pref_len} = r_{1max_word_len}$
- Если r_1 и r_2 содержат слова, состоящие только из букв x , то $r_{max_word_len}$ высчитывается как сумма длин длинейших слов, состоящих только из букв x , из обеих регулярных выражений, то есть $r_{max_word_len} = r_{1max_word_len} + r_{2max_word_len}$
- Иначе, слов, полностью состоящих из букв x , появиться не может и параметр $r_{max_word_len} = -1$

3. Если это операция $*$ ⁽¹⁾, то в терминах регулярных выражений для получения слова w склеиваются слова из r_1 сами с собой необходимое число раз. Пусть для упрощения обозначения $p = r_{1max_word_len}$. Тогда, возможны случаи:

- Если r_1 содержит слово вида x^p (без учета пустого слова, то есть $p \neq 0$), то тогда мы можем получить префикс размера inf , сливая x^p само с собой кучу раз (для $p = 0$, очевидно это неверно). Следовательно, если $r_{1max_word_len} > 0$, то

$$r_{max_pref_len} = r_{max_word_len} = inf$$
- В противном же случае ($p \leq 0$) $r_{max_pref_len} = r_{1max_pref_len}$, так как оператором $*$ мы никак не будем разрастать слово (это бессмысленно), а $r_{max_word_len} = 0$, так как оператор $*$ гарантирует наличие пустого слова в r .

В итоге, обрабатывая с помощью стека исходное регулярное выражение в обратной польской записи, мы приходим к R в обычной записи и у нас будут посчитаны два параметра $R_{max_pref_len}$ и $R_{max_word_len}$. Логично, что ответом на задачу является максимум из этих двух чисел. Однако, как было замечено в начале алгоритма, $max_pref_len \geq max_word_len$, значит максимум всегда будет возвращать первое число. Таким образом, ответом будет являться $R_{max_pref_len}$.

Корректность: Из алгоритма видно, что задача решается индукцией по увеличению приоритета операций над регулярным выражением. Так как исходное регулярное выражение задается нам в обратной польской нотации, то операции поступают к нам сразу в интересующем нас порядке.

- Очевидно, что база индукции (значения параметров для букв) корректна.
- Предполагаем, что параметры были вычислены корректно на всех шагах.
- Переход: делаем последнюю операцию, после выполнения которой, мы получим требуемое регулярное выражение. Так как при парсинге каждой из трех возможных операций наш алгоритм максимизирует ответ по очевидной логике, то в итоге мы получаем величину максимального префикса состоящего из букв x .

Асимптотика: Пусть $len(R)$ – это длина задаваемого регулярного выражения в обратной польской нотации. Заметим, что обработку выражения мы делаем с помощью стека, в котором операции *push*, *top* и *pop* работают за $O(1)$. Мы идем поэлементно по R и взаимодействуем со стеком (кладем в него что-то единожды, либо достаем 1-2 элемента, пересчитываем два параметра за $O(1)$ и кладем это в стек). Значит, один символ из R обрабатывается за $O(1)$. Следовательно, итоговая асимптотика $O(len(R))$.