

Алгоритмы и структуры данных

Строки

Z-функция

CS Center, Новосибирск

Понятия и обозначения

Строка α – последовательность $\alpha_0\alpha_1 \dots \alpha_{n-1}$

$$\alpha_i \in \Sigma = \{0, 1, \dots, |\Sigma|-1\}$$

Σ^* – множество всех строк

« $\alpha \in \Sigma^*$ »

$|\alpha|$ – длина строки

ε – пустая строка

« $|\varepsilon|=0$ »

$\alpha\beta$ – конкатенация двух строк

$\alpha[i : j)$ – подстрока с позиции i (включительно) по j (не включительно)

$$|\alpha[i : j)| = j - i$$

$\alpha[0 : j)$ – префикс строки α

также обозначается $\alpha[:j)$

$\alpha[i : |\alpha|)$ – суффикс строки α

также обозначается $\alpha[i:]$

Pattern matching

Хотим найти шаблон P в тексте T

- $P, T \in \Sigma^*$
- $P \text{ is_sub } T?$
 - Это точный поиск одного шаблона
- Более сложные задачи:
 - нахождение неточных совпадений
 - поиск сразу нескольких шаблонов в тексте
 - online-задачи:
 - сначала получаем текст, потом разные шаблоны по очереди
 - сначала получаем шаблон, потом разные тексты по очереди

Z-функция; LCP

Z_α – массив целых чисел длины $|\alpha|$

- $Z_\alpha[i] = | \text{LCP}(\alpha[i :), \alpha) |$
 - $\text{LCP}(\alpha, \beta)$ – longest common prefix
 - $j, \alpha[: j) = \beta[: j), \alpha[: j + 1) \neq \beta[: j + 1)$
 - $\text{LCP}(\alpha, \beta) = \alpha[: j)$
- Иначе говоря, $Z_\alpha[i] = \max\{ j, \alpha[i : i+j) = \alpha[: j) \}$

Z-функция – пример

Z_α – массив целых чисел длины $|\alpha|$

- $Z_\alpha[i] = |\text{LCP}(\alpha[i:], \alpha)|$

$\alpha = \text{abacabacaba}$

$\alpha[i]$	a	b	a	c	a	b	a	c	a	b	a
$Z_\alpha[i]$	11	0	1	0	7	0	1	0	3	0	1

Z-функция для поиска шаблона

Хотим найти шаблон P в тексте T

$$P, T \in \Sigma^*, P \text{ is_sub } T?$$

- $\alpha = P \$ T$, $\$$ - sentinel (не содержится в P , в T)
- $Z_\alpha = [\dots, 0, \dots]$

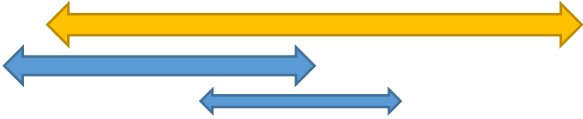
Вхождение с позиции i : $T[i : i+|P|) = P$

$$Z_\alpha[|P|+1+i] = |P|$$

Z-функция за линейное время

- Вычисляем Z_α слева направо
 - добрались до позиции i

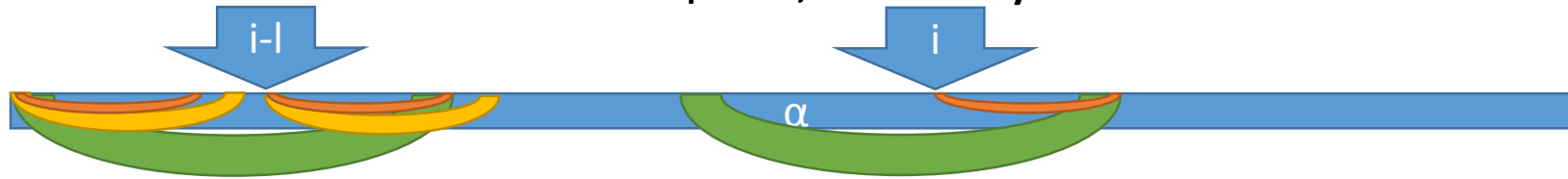
$\alpha[i]$	α_0	α_1	...	α_{i-1}	α_i	...	?
$Z_\alpha[i]$	$ \alpha $	x	...	y	?	???	?



- Блок – это $\alpha[j : j + Z_\alpha[j])$ при $0 < j < i$
- Лучший блок – блок с максимальной правой границей
 - Помним лучший блок

Z-функция за линейное время

- Вычисляем Z_α в позиции i , зная лучший блок



- Пусть лучший блок $[l : r)$ заканчивается правее, чем i ($r > i$)
- Позиция $i-1$ соответствует позиции i
 - $\alpha[: r-1) = \alpha[l : r)$
- Посмотрим на блок, начинающийся в $i-1$; его длина $z[i-1]$
 - $\alpha[: z[i-1]) = \alpha[i-1 : i-1+z[i-1])$
 - $\alpha[i-1 : r-1) = \alpha[i : r)$
 - $\alpha[: k) = \alpha[i-1 : i-1+k) = \alpha[i : i+k)$, $k = \min(z[i-1], r-i)$
 - Известно, что $z[i] \geq \min(z[i-1], r-i)$
 - Будем использовать это значение как начальное для $z[i]$

Z-функция за линейное время

```
(l, r) = (0, 0)
```

```
for i ∈ [ 1 : |α| )
```

```
    z[i] = (r > i) ? min(z[i-1], r-i) : 0
```

```
    while i+z[i] < |α| && α[ z[i] ] == α[ i+z[i] ]
```

```
        ++z[i] // каждая итерация вложенного цикла соответствует увеличению r на единицу
```

```
    if i+z[i] > r
```

```
        (l, r) = (i, i+z[i])
```

- $O(N)$ итераций внешнего цикла
- в сумме $O(N)$ итераций внутреннего цикла – r монотонно растёт до $|\alpha|$

Z-функция для поиска шаблона

- Хотим найти шаблон P в тексте T
 - $\alpha = P \$ T$, $\$$ - sentinel
 - $Z_\alpha[|P|+1+i] = |P| \Rightarrow$ вхождение с позиции i
- Можно уменьшить используемую память
 - $\$$ не даёт блоку быть длиннее $|P|$
 - достаточно хранить только первые $|P|$ значений Z_α
- Можно предобработать образец за $O(|P|)$ и искать в различных текстах за $O(|T|)$

Расстояние редактирования ED

- $ED(\alpha, \beta)$ = сколько нужно изменений чтобы из α получить β
возможные изменения:
 - вставка символа
 - удаление символа
 - замена одного символа на другой
- Неточный поиск шаблона (с одной ошибкой)
 - $P, T \in \Sigma^*$
 - $P' \text{ is_sub } T, ED(P', P) \leq 1?$

Поиск шаблона с одной ошибкой

- Z для $P \ \$ \ T$
- Z_{rev} для $\text{rev}(P) \ \$ \ \text{rev}(T)$
- $\text{cut}(Z)$ – отбросить от Z первые $|P|+1$ элементов
- $Z_1 := \text{cut}(Z), Z_2 := \text{rev}(\text{cut}(Z_{\text{rev}}))$
 - $Z_1[i] + Z_2[i+|P|] = |P| \Rightarrow$ вхождение со вставкой
 - $Z_1[i] + Z_2[i+|P|-1] = |P|-1 \Rightarrow$ вхождение с заменой
 - $Z_1[i] + Z_2[i+|P|-2] = |P|-1 \Rightarrow$ вхождение с удалением

Длиннейшее общее продолжение

- Даны α, β
- Запрос $(i, j) \rightarrow \text{ответ } |\text{LCP}(\alpha[i:], \beta[j:])|$
- z-функция позволяет отвечать на такие запросы за $O(1)$ при фиксированном i
 - подсчёт: Z от $\alpha[i:] \text{ \$ } \beta$