

Алгоритмы и структуры данных

Приближённый поиск подстрок

CS Center, Новосибирск

Приближённый поиск подстрок

- $P, T \in \Sigma^*$
 - $P \approx P', P' \text{ is_sub } T?$
- 1. $P \approx P'$:
 - $|P| = |P'|$
 - $\text{distH}(P, P') \leq k$
 - $\text{distH}(P, P') = |\{i, P[i] \neq P'[i]\}|$
- 2. $P \in (\Sigma \cup \{?\})^*, P \approx P'$:
 - $|P| = |P'|,$
 - $P'[i] \in \{P[i], ?\}, i \in [0 : |P|)$
- 3. $P \approx P'$:
 - $\text{ED}(P, P') \leq k$

Приближённый поиск: $P \in (\Sigma \cup \{?\})^*$

- $P = ?^* a_0 P_1 ?^* a_1 P_2 ?^* a_2 \dots ?^*_{k-1} P_k ?^* a_k$
 - Ахо-Корасик: $O(|T| + |P| + |A|)$ – неэффективно
 - Массив счётчиков `cnt`
 - `cnt[i] = k` => вхождение в позиции `i`
 - найден `Pi` в позиции `j` -> увеличить счётчик `cnt[j - shift[i]]`
 - вхождения находятся по порядку
 - достаточно хранить `|P|` счётчиков

#include FFT

- $a = (a_0, \dots, a_n)$
- $b = (b_0, \dots, b_m)$
- $c = (c_0, \dots, c_{n+m})$
 - $c_i = \sum_j a_j b_{i-j}$
- Есть алгоритм со сложностью $O((n+m) \log(n+m))$
- $b'_j = b_{m-j}$
 - $c_i = \sum_j a_j b'_{i-j} = \sum_j a_j b_{m-i+j}$
- $c'_s = c_{m-s}$
 - $c'_i = c_{m-i} = \sum_j a_j b_{i+j}$

Приближённый поиск: $\text{dist}_H(P, P') \leq k$

- $|\Sigma_{abc}| = C$
- C масок для строки α :
 - $\text{mask}_{\alpha x} \in \{0, 1\}$
 - $|\text{mask}_{\alpha x}| = |\alpha|$
 - $\text{mask}_{\alpha x}[i] = \alpha[i]=x ? 1 : 0$
- $\text{match}(i, x) = \sum_j \text{mask}_{Px}[j] \cdot \text{mask}_{Tx}[i+j]$
- $\text{match}(i) = \sum_{x \in \Sigma_{abc}} \text{match}(i, x)$
 - $\text{match}(i) \geq |P| - k \Rightarrow$ вхождение в позиции i
- $N = |P|, M = |T|$
 - Сложность: $O(CM \log M)$
- Делим текст на блоки длины $2N$
 - Блоки перекрываются, чтобы не потерять вхождения
 - $O(M/N)$ блоков
- В каждом блоке ищем шаблон независимо
 - сложность в блоке $O(CN \log N)$
- Общая сложность $O(CM \log N)$

Приближённый поиск: $P \in (\Sigma \cup \{?\})^*$

- Замена Σ на $\{0, 1\}$ – работаем с бинарными строками
 - code: $\Sigma \rightarrow \{0, 1\}^*$
 - $|\text{code}(x)| = \log |\Sigma|$
 - $\text{code}(?) = 00\dots 0$
 - $\text{bin}(\alpha) = \text{code}(\alpha[0]) \text{code}(\alpha[1]) \dots \text{code}(\alpha[|\alpha|-1])$
 - $|\text{bin}(\alpha)| = |\alpha| \log |\Sigma|$
 - Свертка $\text{bin}(P) * \text{bin}(T) \rightarrow \text{res}$
 - $\text{res}[i \log |\Sigma|] = \text{popcnt}(\text{bin}(P)) \Rightarrow$ вхождение в позиции i
 - Сложность $O(M \log C \log(N \log C))$

Приближённый поиск: $P \in (\Sigma \cup \{?\})^*$

- $R = \{ i, P_i \neq ? \}$
 - $\sum_{j \in R} (P_j - T_{i+j})^2 = 0 \Rightarrow$ вхождение в позиции i
- $\sum_{j \in R} (P_j - T_{i+j})^2 = \sum_j (P_j - T_{i+j})^2 \cdot X_j = \sum_j P_j^2 \cdot X_j - 2 \sum_j P_j \cdot T_{i+j} \cdot X_j + \sum_j T_{i+j}^2 \cdot X_j$
 - $\sum_j P_j^2 \cdot X_j$ не зависит от i
 - $\sum_j P_j \cdot T_{i+j} \cdot X_j$ – свертка T и PX
 - $\sum_j T_{i+j}^2 \cdot X_j$ – свертка T^2 и X
 - сложность $O(M \log N)$, $M = |T|$, $N = |P|$

Приближённый поиск: $ED(P, P') \leq k$

Динамическое программирование

$dp[i][j] = \min \{$

$dp[i-1][j] + (0 < i < |T| ? 1 : 0),$

$dp[i][j-1] + 1,$

$dp[i-1][j-i] + (T[i]=P[j] ? 1 : 0)$

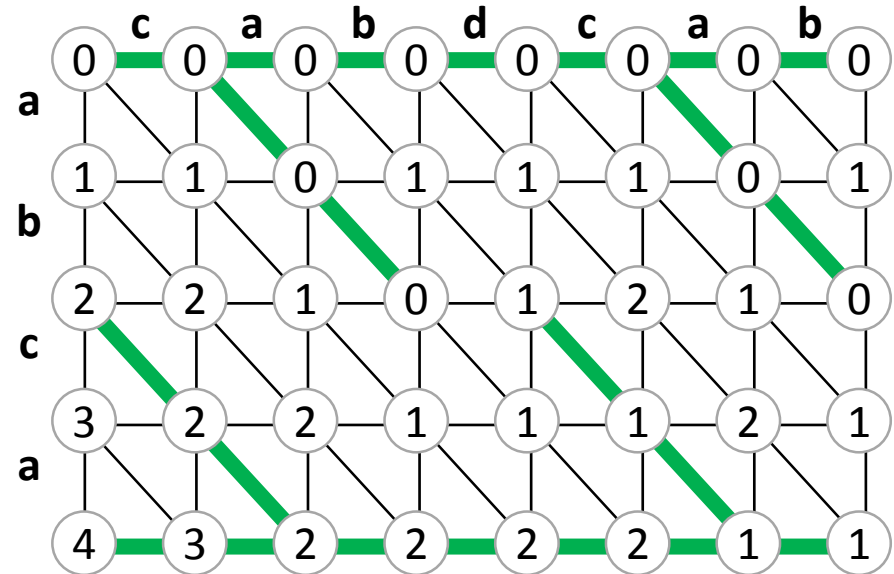
$\}$

$dp[|T|][|P|] \leq k \Rightarrow$ есть вхождение

- Сложность $O(|T||P|)$

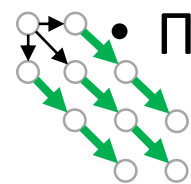
Пример:

$P=abca, T=cabdca$



Приближённый поиск: $ED(P, P') \leq k$

- S – вершины в строке 0, T – вершины в строке $|P|$
- $D_i = \{v \mid \text{dist}(S, v) \leq i\}$
- $D_k \cap T = \emptyset$?
- D_i' вместо D_i – на каждой диагонали оставляем только последнюю вершину
 - $|D_i'| \leq |P| + |T|$
 - $D_k \cap T = \emptyset \Leftrightarrow D_k' \cap T = \emptyset$
- Нужно построить D_0' и научиться переходить $D_i' \rightarrow D_{i+1}'$
 - D_0' строится с помощью $LCP(P, T[i:])$ за линейное время



- Переход $D_i' \rightarrow D_{i+1}'$:

- для каждой вершины 3 запроса LCP
- на диагонали берём наиболее далёкую вершину

- Сложность $O(k|T|)$

