

Wyszukiwanie z wildcardami w tekście i wzorcu

Kamil Rajtar

Czerwiec 2020

Opis problemu

Na wejściu do algorytmu dostajemy tekst i wzorec w których oprócz liter mogą występować wildcardy (?) pasujące do dowolnego znaku. Na wyjściu ma się znaleźć lista pozycji w których wzorec pasuje do tekstu.

Wersja bez wildcardów

W rozwiązaniu problemu użyty jest spłot obliczany za pomocą szybkiej transformaty Fouriera (FFT) zdefiniowany następująco:

$$p \otimes t \stackrel{\text{def}}{=} \left(\sum_{j=0}^{m-1} p_j t_{i+j}, 0 \leq i \leq n-m \right)$$

Normalne (bez wildcardów) dopasowanie wzorca p do tekstu t na pozycjach od i do $i+m$ możemy obliczyć za pomocą spłotu. Korzystamy z własności wzoru:

$$\sum_{j=0}^{m-1} (p_j - t_{i+j})^2 = \sum_{j=0}^{m-1} (p_j^2 - 2p_j t_{i+j} + t_{i+j}^2)$$

Zobaczmy że lewa część równania równa 0 to znaleźliśmy dopasowanie. Tekst nie różni się od wzorca na żadnej z m kolejnych pozycji. Prawą stronę potrafimy szybko obliczyć ponieważ podniesienie każdej pozycji do kwadratu wykonywane jest w czasie stałym a środkowy składnik liczymy za pomocą FFT.

Wersja z wildcardami

Zdefiniujmy ciągi:

$$p'_j = \begin{cases} 0 & p_j = '?' \\ 1 & \text{wpp} \end{cases}$$
$$t'_j = \begin{cases} 0 & t_j = '?' \\ 1 & \text{wpp} \end{cases}$$

Wtedy łatwo jest widzieć że następujące równanie jest naturalnym rozszerzeniem rozwiązania poprzedniego problemu.

$$\sum_{j=0}^{m-1} p'_j t'_{i+j} (p_j - t_{i+j})^2 = 0$$

Zobaczmy że tekst część pod sumą jest zawsze nieujemna p'_j oraz t'_{i+j} przyjmują wartości $\{0,1\}$ oraz kwadrat różnicy jest zawsze nieujemny. Więc zastanówmy się kiedy suma jest równa zero. Wtedy kiedy wszystkie składniki są równe zero. Kiedy składnik jest równy zero? W jednym z 3 przypadków. $p'_j = 0$ wtedy mamy wildcard we wzorcu, $t'_{i+j} = 0$ wtedy mamy wildcard w tekście, $(p_j - t_{i+j})^2$ wtedy litery wzorca i tekstu są równe. Zobaczmy że to zachodzi wtedy i tylko wtedy gdy mamy do czynienia z dopasowaniem.

Po rozpisaniu nawiasu dostajemy formę:

$$\sum_{j=0}^{m-1} (p'_j p_j^2 t'_{i+j} - 2p'_j p_j t_{i+j} t'_{i+j} + p'_j t_{i+j}^2 t'_{i+j})$$

Takie rozwinięcie łatwo jest obliczyć wykorzystując 3xFFT oraz podstawowe operacje na ciągach.

Złożoność

Trywialnie jest pokazać że algorytm działa $O(n \log n)$ gdy użyjemy FFT dla tekstu i patternu rozszerzonego do długości tekstu. Można jednak osiągnąć $O(n \log m)$ gdy podzielimy tekst na n/m zachodzących na siebie kawałków wielkości $2m$ i dla każdego policzymy matching ze wzorcem a następnie wyniki złączymy. Zobaczmy że jest to poprawne ponieważ dla każdego miejsca początkowego istnieje kawałek, który zawiera część od długości co najmniej m od tego miejsca. Złożoność: $n/m * O(m \log m) = O(n \log m)$.