

# 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 Furiera (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 za pomocą spłotu tekstu ( $t$ ) i odwróconego wzorca ( $p$ ) możemy obliczyć za pomocą 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 w  $O(n \log m)$ .

## 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 widać że następujące równanie jest naturalnym rozszerzeniem rozwiązania poprzedniego problemu.

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

Po rozwinięciu 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**

Analiza złożoności jest trywialna ponieważ algorytm działa w czasie FFT.