

Алгоритъм на Lempel-Ziv. Построяване на суфиксни масиви.

08.01.2021 г.

Една от мотивациите да може да строим суфиксно дърво on-line е, че това позволява компресия, при това ефективна такава (линейно време и теоритично смачкване на данните, - възможно най-много - алгоритъма на Lempel-Ziv).

Връзка между суфиксни масиви, суфиксни дървета и компресия.

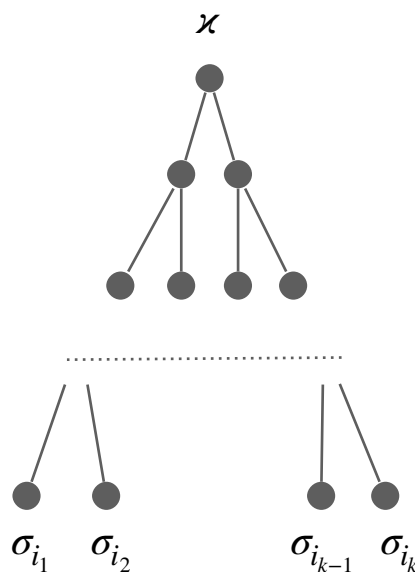
Нека имаме текст $T = t_1 t_2 \dots t_n$, $t_i \in \Sigma$, $i = \overline{1, n}$. Искаме да представим този текст T с възможно най-малко битове (най-малко памет). В тази връзка, когато Шенън е започнал да разсъждава върху този въпрос през 60-те години на миналия век, той заключил следното: ако искаме да представим символите независимо един от друг, т.е. искаме някаква компресия $\kappa : \Sigma \rightarrow \{0,1\}^*$, така че ако имаме две букви

$\sigma_1, \sigma_2 \in \Sigma$ и $\sigma_1 \neq \sigma_2$, $\kappa(\sigma_1) \not\prec_{pref} \kappa(\sigma_2)$. Идеята е следната: ако имаме такава

(*)

функция $\kappa(\sigma)$, то тогава може да компресиране текста $T \mapsto \kappa(t_1)\kappa(t_2)\dots\kappa(t_n)$.

Условието (*) означава, че няма многозначности. Може да си представим κ като едно двоично дърво с нули и единици:



и някъде по листата на това дърво седи някаква пермутация на символите σ_i .

Разсъждавайки по този начин стигаме до следната парадигма - искаме да

минимизираме $\sum_{i=1}^n |\kappa(t_i)|$, т.е. $\min \sum_{i=1}^n |\kappa(t_i)| = M$ което означава, че като

преброим колко пъти се среща всеки символ в текста,

$M = \min_{\sigma \in \Sigma} \underbrace{\text{осс}_T(\sigma)}_{\text{брой срещания на } \sigma} \underbrace{|\kappa(\sigma)|}_{\text{дължина на } \sigma}$. Оптимумът се достига при

$|x(\sigma)| = -\log_{(2)} \frac{\text{occ}_T(\sigma)}{|T|}$ - честота на символа σ в текста T . Крайният резултат

казва, че честосрещаните символи в текстовете искаме да ги компресируем с възможно най-малко бита, така че те да оказват най-малко влияние върху цялата сума. И обратното - рядко срещаните символи съответно ще бъдат кодирани с по дълго представяне. Това е мястото, от където се появява тази ентропия, за която става въпрос през цялото време.

Ако означим $p(\sigma) = \mathbb{P}(\sigma) = \frac{\text{occ}_T(\sigma)}{|T|}$ да е честотата на срещане на символа σ в T ,

то $H(T) \stackrel{\text{def.}}{=} - \sum_{\sigma} p(\sigma) \ln p(\sigma)$ е ентропията на това колко добре може да

компресируем един текст. Алгоритъм на Хъфман (Huffman). Това е най-простата схема, ако не искаме да вземаме под внимание контекста, в който символите се срещат, с която човек може да реши проблема за компресията. Обаче обикновено има зависимост между символите и ние бихме искали да отчетем тази зависимост. Алгоритъмът на Lempel-Ziv постига точно това.