

Compression Methods

neděle 30. listopadu 2025 17:56

How to compress a string? We will look at three measures of compression.

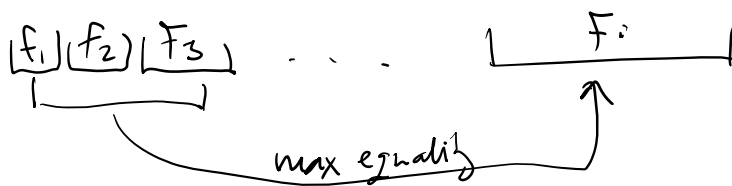
1) Lempel-Ziv (LZ) type compression:

$$w \in \Sigma^N \quad w = f_1 f_2 \dots f_z \quad f_i \in \Sigma^{l_i}$$

$$f_i = w_{b_i \dots b_i + l_i - 1} \quad \text{where } l_i \text{ is maximal s.t.}$$

$$\exists a_i \in \{1, \dots, \sum_{j < i} |f_j| \}$$

$$w_{a_i \dots a_i + l_i - 1} = w_{b_i \dots b_i + l_i - 1}$$



$f_i \dots$ Fragment

$z_w = z \dots$ # of fragments

2) grammar representation

$$w \in \Sigma^N$$

Context-free grammar with rules

$$A \rightarrow BC$$

$A, B, C \in \text{variables}$

$$B \rightarrow b$$

$b \in \Sigma$

(Chomsky normal form)

Consider the smallest grammar (# of rules) that generates w

& nothing else $\dots g_w$

3) Substring diversity: $w \in \Sigma^N$ $d_\ell(w) = |\{w_{i, i+\ell-1}, 1 \leq i \leq n-\ell+1\}|$

$$\mathcal{D}_w = \max_{\ell \in [n]} \frac{d_\ell(w)}{\ell}$$

Claim: $\tilde{O}(z_w) = \tilde{O}(g_w) = \tilde{O}(\mathcal{D}_w)$

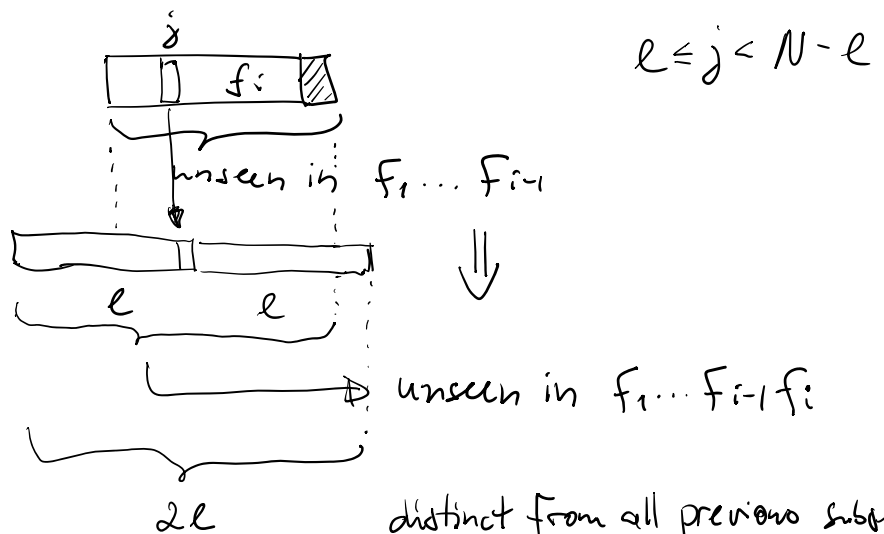
Pf: 1) $z_w \leq 4(\mathcal{D}_w + 1) \lg N$

Claim: $\forall \ell \leq n \quad \sum_{k=1}^{\ell} k n_k \leq 2\ell(\mathcal{D}_w + 1)$ where n_k is the # of fragments f_i of length k .

[Raskhodnikova
- Ron-Rubinfeld
- Smith '13]

Pf: $d_{2\ell} \leq \mathcal{D}_w \cdot 2\ell$

Consider strings of length 2ℓ . We will assign each position $j \in \{\ell, \dots, N-\ell\}$ that falls within a fragment f_i of size $\leq \ell$ to a distinct substring of length 2ℓ .



distinct from all previous substrings of length 2ℓ which have their middle point $< j$.

middle point $< j$.

\Rightarrow each fragment of length $k \leq l$ introduces k distinct substrings of length $2l$. (except very end/beginning of w)

$$\Rightarrow \sum_{k=1}^l k n_k - 2l \leq d_{2l}(w) \leq \delta_w \cdot 2l$$

\uparrow
 $l \leq j < N-l$

$$\Rightarrow \sum_{k=1}^l k n_k \leq (\delta_w + 1) \cdot 2l \quad \square$$

$$\Rightarrow \sum_{k=\frac{l}{2}}^l \frac{l}{2} n_k \leq \sum_{k=1}^l k \cdot n_k \leq (\delta_w + 1) \cdot 2l$$

$$\Rightarrow \sum_{k=\frac{l}{2}}^l n_k \leq 4 \cdot (\delta_w + 1)$$

$$\Rightarrow z_w = \sum_{k=1}^N n_k \leq 4(\delta_w + 1) \cdot \lg N \quad \square$$

2) $\delta_w \leq z_w$

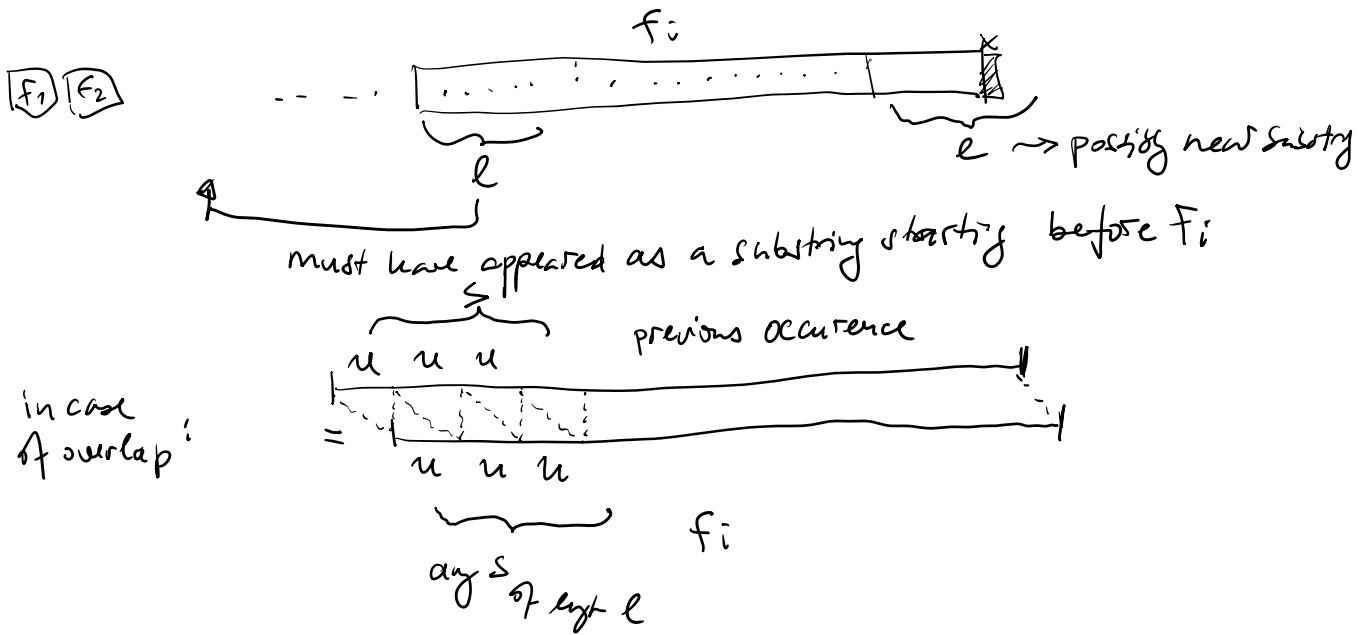
Pf:

For any $l \in [N]$

$$\frac{d_l(w)}{l} \leq \delta_w$$

A fragment F_i of length $k \leq l$ gives only $k \leq l$ substrings of length l that start in it.

A fragment F_i of length $k > l$ gives at most l distinct substrings that start in it:

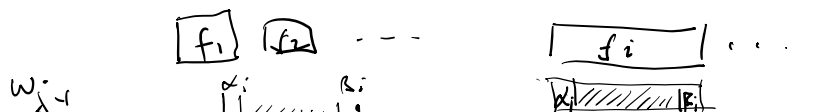


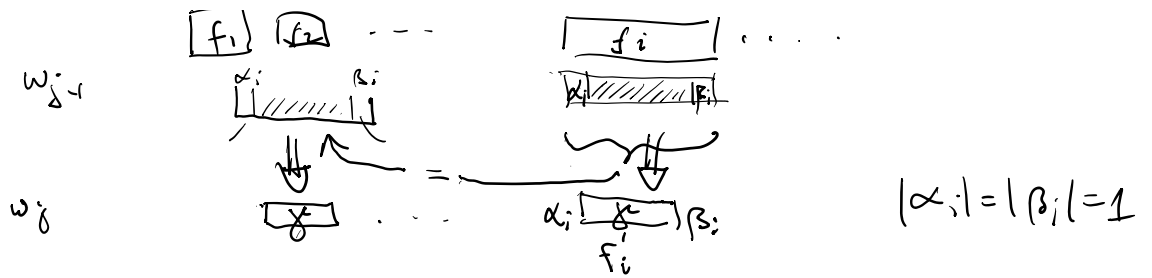
$$\Rightarrow d_c(w) \leq z_w \cdot l \Rightarrow \delta_w \leq z_w \quad \square$$

$$3) \quad g_w \leq O(z_w \cdot \lg N) \quad w \in \Sigma^N$$

[Jeř'16]

start with factorization F_1, \dots, F_2 , we will create iteratively a grammar for w . It is created in $O(\lg N)$ stages where in each stage we introduce $O(z)$ new rules. We start with $w_0 = w$ & create w_1, w_2, \dots where in stage j we create w_j from w_{j-1} by "compression" some nonoverlapping pairs of w_{j-1} . The compression (new) is done by replacing a consecutive pair of letters by a \sqrt{z} symbol & adding the appropriate rule. We compress w_{j-1} factor by factor left to right:





except for α_i & β_i which ^{might be} paired outside of f_i , the inner portion of f_i can be compressed as in its corresponding copy to the left.

(Unless the corresponding copy starts at least two symbols to the left of f_i w/ remove the first letter of f_i & leave it singleton. This only happens for $f_i = a^k$.)

Whenever a singleton letter α_i or β_i appears next to another singleton letter, compress it using a new rule.

(In particular, if the inner portion of f_i is a single letter, compress it with neighboring singletons if there is one.)

Since # Factors $\leq z$, # singletons created at each step is at most $3z \Rightarrow$ at most $\frac{3}{2}z$ new rules are created at each iteration.

Since we greedily compress neighboring singletons whenever they are created, the compression cannot leave two singletons next to each other. So each singleton must be flanked by a pair on either side.

\Rightarrow out of every 3 consecutive symbols in w_{i-1} , at least two get compressed into one.

$\Rightarrow |w_i| \leq \frac{2}{3} |w_{i-1}|$

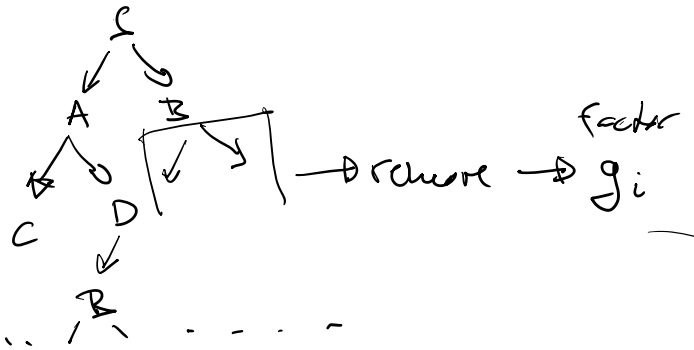
$$\Rightarrow |w_j| \leq \frac{2}{3} |w_{j-1}|$$

\Rightarrow After $\log_{\frac{2}{3}} N$ iterations, w_j consists of only one singleton \Rightarrow done □

4) $z_w \leq 2g_w$

[Rytter '03]

Consider a grammar of size g_w & a derivation tree of w :



$w \rightarrow \dots B \dots$
going left to right, remove a subtree below a non-terminal if B occurs somewhere in the tree to the left

\rightarrow gives a factorization $g_1 \dots g_{\leq 2g_w}$ of w

first occurrence of A might contribute by two factors $g_i \Rightarrow \# \text{fac.} \leq 2g_w$



g_i : each inner node corresponds to a distinct non-terminal. $\# \text{leaves} \leq 2 \cdot \# \text{non-term.}$

where each factor g_i corresponds to some non-terminal or a single letter & g_i appears in $g_1 \dots g_{i-1}$

let $f_1 f_2 \dots f_{z_w}$ be LZ-factorization of w

claim: $\forall i \quad |g_1 \dots g_i| \leq |f_1 \dots f_i|$

Pf: Induct on i : $i = \checkmark$
 $i \rightarrow i+1$:

if $|g_1 \dots g_i| = |f_1 \dots f_i|$

then $|g_1 \dots g_{i+1}| \leq |f_1 \dots f_{i+1}|$

since LZ is greedy & f_{i+1} will contain
at least g_{i+1} .

if $|g_1 \dots g_i| < |f_1 \dots f_i|$

then either g_{i+1} is already contained
in $f_1 \dots f_i$ or some suffix of g_{i+1}
will be covered by f_{i+1} . Again by
greediness $|g_1 \dots g_{i+1}| \leq |f_1 \dots f_{i+1}|$



$\Rightarrow z_w \leq g_w$

