

SILA A ZLOŽITOSŤ MODERNÝCH REGULÁRNYCH VÝRAZOV

Tatiana Tóthová

školiteľ: RNDr. Michal Forišek, PhD.

Regulárne výrazy vznikli v 60tych rokoch ako ďalší model reprezentujúci regulárne jazyky. Vďaka jednoduchosti a kompaktnosti boli implementované ako vyhľadávacie nástroj. V praxi boli postupne rozširované o ďalšie konštrukcie. My sme tento model skúmali z hľadiska teórie formálnych jazykov. Niektoré nové konštrukcie sa dajú prepísať pomocou základných operácií, iné model rozšírili.

Moderné regulárne výrazy (regexy) obsahujú základné operácie regulárnych výrazov – zretazenie, alternácia, Kleeneho *. Navyše majú metaznak pre ľubovoľný znak ., začiatok slova ^ a koniec slova \$. Zátvorky regexov sú číslované, zľava doprava podľa otváracej zátvorky.

Zložitejšie konštrukcie rozširujúce model regulárnych výrazov sú:

- **Spätné referencie** ($\backslash k$) sa odkazujú na k -te zátvorky. Pre regex $\alpha \left(\beta \right) \gamma \backslash k \delta$ výpočet na slove w vyzerá takto:

$$w = \underbrace{x_1 \dots x_{i-1}}_{\alpha} \underbrace{\overbrace{x_i \dots x_{j-1}}^{w_k}}_{\left(\beta \right)} \underbrace{x_j \dots x_{l-1}}_{\gamma} \underbrace{\overbrace{x_l \dots x_{m-1}}^{w_k}}_{\backslash k} \underbrace{x_m \dots x_n}_{\delta}$$

pričom platí: $w_k = x_i \dots x_{j-1} = x_l \dots x_{m-1}$.

- **Lookahead** (nazeranie dopredu ($?= \dots$)) musí matchovať nejaký prefix od aktuálneho pracovného miesta. Pre regex $\alpha(?=\beta)\gamma$:

$$w = \underbrace{x_1 \dots x_{i-1}}_{\alpha} \underbrace{\overbrace{x_i \dots x_j}^{\beta} x_{j+1} \dots x_n}_{\gamma}$$

- **Lookbehind** (nazeranie dozadu ($?<= \dots$)) musí matchovať nejaký sufix od aktuálneho pracovného miesta.
- Operácie lookahead a lookbehind majú aj **negatívne** verzie (pridáva sa !), kde regex vnútri konštrukcie nesmie matchovať žiaden prefix/sufix.

Vlastnosti lookaheadu a lookbehindu

Už len uzavretosť na zretazenie nie je triviálna, pretože treba zabrániť presahu lookaheadu a lookbehindu vo výpočte na slove z $L(\alpha)L(\beta)$:

$$\left(\left(\alpha \right) \left(\beta \right) \right) \alpha' \backslash k+2 \left(?<= \wedge \backslash 1 \beta' \right)$$

Zaviedli sme nasledujúce triedy regexov nad množinami operácií:

Regex – základné operácie. *Eregex* – má navyše spätné referencie (\mathcal{L}_{ERE}). $LEregex = Eregex + \text{lookahead, lookbehind}$ (\mathcal{L}_{LERE}). $nLEregex = LEregex + \text{neg. lookahead, neg. lookbehind}$ (\mathcal{L}_{nLEERE}).

Ukázali sme, že trieda *Regex* s pozitívnym a negatívnym lookaheadom a lookbehindom je ekvivalentná triede \mathcal{R} . Pre zadané triedy zároveň platí nasledovná hierarchia:

$$\mathcal{R} \subsetneq \mathcal{L}_{ERE} \subsetneq \mathcal{L}_{LERE} \subseteq \mathcal{L}_{nLEERE} \subsetneq \mathcal{L}_{CS}$$

Vidieť, že pokiaľ sú v modeli aj spätné referencie, lookahead a lookbehind ho rozširujú.

Trieda *nLEregex* nie je porovnateľná s \mathcal{L}_{CF} .

Priestorová zložitosť

Ukázali sme, že $\mathcal{L}_{LERE} \subseteq NSPACE(\log n)$.

Zo Savitchovej vety vyplýva: $\mathcal{L}_{LERE} \subseteq DSPACE(\log^2 n)$ a výsledok sme rozšírili: $\mathcal{L}_{nLEERE} \subseteq DSPACE(\log^2 n)$

Pre problém, kde dostávame na vstup aj regex aj vstupné slovo sme ukázali nasledovné výsledky:

$$L(\text{regex}\#\text{word}|\text{regex} \in LEregex) \in NSPACE(n \log n)$$

Nech regex r je taký, že dĺžka jeho akceptačného výpočtu na slove w sa dá zhora ohraničiť hodnotou $f(|w|)$, potom platí:

$L(r\#w) \in DSPACE(w \cdot \log n \cdot \log(f(n)))$, pričom n je dĺžka vstupu