

Moderné regulárne výrazy

Tatiana Tóthová*

Školiteľ: Michal Forišek†

Katedra informatiky, FMFI UK, Mlynská Dolina 842 48 Bratislava

Abstrakt: Regulárne výrazy implementované v súčasných programovacích jazykoch ponúkajú omnoho viac operácií ako pôvodný model z teórie jazykov. Už konštrukciou spätných referencií bola prekročená hranica regulárnych jazykov. Náš model obsahuje navyše konštrukcie lookahead a lookbehind. V článku uvidíme zaradenie modelu zodpovedajúcej triedy jazykov do Chomského hierarchie, vlastností tejto triedy a výsledky z oblasti priestorovej zložitosti.

Keľúčové slová: regulárny výraz, regex, lookahead, lookbehind, spätné referencie

1 Úvod

Regulárne výrazy vznikli v 60tych rokoch v teórii jazykov ako ďalší model na vyjadrenie regulárnych jazykov. Z takéhoto popisu ľudský mozog rýchlejšie pochopil o aký jazyk sa jedná, než zo zápisu konečného automatu, či regulárnej gramatiky. Ďalšou výhodou bol kratší a kompaktný zápis.

Vďaka týmto vlastnostiam boli implementované ako vyhľadávacie nástroje. Postupom času sa iniciatívou používateľov s vyššími nárokmi pridávali nové konštrukcie na uľahčenie práce. Nástroj takto rozvíjali až do dnešnej podoby. My sa budeme opierať o špecifikáciu regulárnych výrazov v jazyku Python [Foundation, 2012].

Ako čoskoro zistíme, nové regulárne výrazy vedú reprezentovať zložitejšie jazyky ako regulárne, preto je dobré ich nejako odlíšiť. V literatúre sa zaužíval výraz „regex“ z anglického *regular expression*, ktorý budeme používať aj my.

1.1 Základná definícia

Regulárne výrazy sú zložené zo znakov a metaznakov. Znak a predstavuje jazyk $L(a) = \{a\}$. Metaznak alebo skupina metaznakov určuje, aká operácia sa so znakmi udeje. Základné operácie sú zret'azenie (je definované tým, že regulárne výrazy idú po sebe, bez metaznaku), Kleeneho uzáver ($(0 - \infty)$ -krát zopakuj

výraz, metaznak $*$) a alternácia (vyber výraz naľavo alebo napravo, metaznak $|$). Navyše sa využíva metaznak \backslash , ktorý robí z metaznakov obyčajné znaky a okružle zátvorky na logické oddelenie regulárnych výrazov.

Pre regulárny výraz α a slovo $w \in L(\alpha)$ hovoríme, že α vyhovuje slovu w resp. α matchuje slovo w . Tiež budeme hovoriť, že α generuje jazyk $L(\alpha)$.

1.2 Nové jednoduché konštrukcie

- $+$ – Kleeneho uzáver opakujúci $(1 - \infty)$ -krát
- $\{n, m\}$ ($\{n\}$) – opakuj regulárny výraz aspoň n a najviac m -krát (opakuj n -krát)
- $[a_1 a_2 \dots a_n]$ – predstavuje ľubovoľný znak z množiny $\{a_1, \dots, a_n\}$
- $[\wedge a_1 a_2 \dots a_n]$ – predstavuje ľubovoľný znak, ktorý nepatrí do množiny $\{a_1, \dots, a_n\}$
- $.$ – predstavuje ľubovoľný znak
- $?$ – ak samostatne: opakuj 0 alebo 1-krát
ak za operáciou: namiesto greedy implementácie použi minimalistickú, t.j. zober čo najmenej znakov (platí pre $*, +, ?, \{n, m\}$)¹
- $^$ – metaznak označujúci začiatok slova
- $\$$ – metaznak označujúci koniec slova
- $(?# \text{komentár})$ – komentár sa pri vykonávaní regexu úplne ignoruje

Všetky tieto konštrukcie sú len „kozmetickou“ úpravou pôvodných regexov – to isté vieme popísať pôvodnými regulárnymi výrazmi, akurát je to dlhšie a menej prehľadné

Rozdiely medzi minimalistickou a greedy verziou operácií vníma iba používateľ, pretože ak existuje zhoda regexu so slovom, v oboch prípadoch sa nájde. Viditeľné sú až pri výstupnej informácii pre používateľa, ktorú môže použiť ďalej.

*tothova166@uniba.sk

†forisek@dcf.fmph.uniba.sk

¹všetky spomenuté operácie sú implementované greedy algoritmom

1.3 Zložitejšie konštrukcie

Spätné referencie

Najprv potrebujeme očíslovať všetky zátvorky v regexe. Číslujú sa všetky, ktoré nie sú tvaru $(? \dots)$. Poradie je určené podľa otváraciej zátvorky.

Spätné referencie umožňujú odkazovať sa na konkrétne zátvorky. Zápis je $\backslash k$ a môže sa nachádzať až za k -tymi zátvorkami. Skutočná hodnota $\backslash k$ sa určí až počas výpočtu – predstavuje posledné podslovo zo vstupu, ktoré matchovali k -te zátvorky.

Lookahead

Zapísaný formou $(?= \dots)$, vnútri je validný regex.

Keď v regexe prideme na pozíciu lookaheadu, zoberieme regex vo vnútri. V slove sa snažíme matchovať ľubovoľný prefix zostávajúcej časti slova. Ak sa to podarí, pokračujeme v regexe ďalej a v slove od pozície, kde lookahead začínal (tzn. ako keby v regexe nikdy nebol).

Má aj negatívnu verziu – negatívny lookahead $(?! \dots)$. Vykondáva sa rovnako ako lookahead, ale má otočnú akceptáciu. Teda ak neexistuje prefix, ktorý by vedel matchovať, akceptuje.

Lookbehind

Zapísaný formou $(?<= \dots)$, vnútri je validný regex.

Pri výpočte zoberieme regex vnútri lookbehindu a snažíme sa vyhovieť ľubovoľnému sufixu už matchovanej časti slova. Ak vyhovíme, pokračujeme v slove a regexe akoby tam lookbehind vôbec nebol.

Aj lookbehind má negatívnu verziu – negatívny lookbehind $(?<! \dots)$ – a pracuje analogicky ako negatívny lookahead.

Lookahead a lookbehind (spolu nazývané jedným slovom lookahead) sú v rôznych implementáciách rôzne obmedzované, aby výpočet netrval príliš dlho. V teórii tieto obmedzenia ignorujeme a prezentujeme model v plnej sile – výsledky tak prezentujú hornú hranicu toho, čo implementácie dokážu.

1.4 Triedy a množiny

Kvôli porovnávaniu a vytvoreniu hierarchie sme rozdelili operácie do niekoľkých množín:

Regex – množina operácií, pomocou ktorých vieme popísať iba regulárne jazyky; presnejšie všetky znaky a metaznaky (bez zložitejších operácií)

Eregex – *Regex* so spätnými referenciami

LEregex – *Eregex* s pozitívnym lookaroundom

nLEregex – *LEregex* s negatívnym lookaroundom

\mathcal{L}_{RE} – trieda jazykov nad *Regex* ($= \mathcal{R}$)

\mathcal{L}_{ERE} – trieda jazykov nad *Eregex*

\mathcal{L}_{LERE} – trieda jazykov nad *LEregex*

\mathcal{L}_{nLERE} – trieda jazykov nad *nLEregex*

Trieda \mathcal{L}_{LERE} už bola hlbšie preskúmaná a výsledky čerpáme z článkov [Câmpeanu et al., 2003] a [Carle and Nadendran, 2009].

2 Vlastnosti lookaroundu

Na začiatok sme zistovali, čo robia samotné operácie lookaroundu.

Veta 1. \mathcal{R} je uzavretá na negatívny a pozitívny lookaround.

Dôkaz. Nech $L_1, L_2, L_3 \in \mathcal{R}$. Chceme ukázať, že $L_1(=L_2)L_3, L_1(?<=L_2)L_3, L_1(?!L_2)L_3, L_1(?<!L_2)L_3 \in \mathcal{R}$. Pre každé $L_i, i \in \{1, 2, 3\}$ existuje determinický konečný automat A_i , ktorý ho akceptuje. Spojíme automaty dohromady nasledujúcim spôsobom:

$L_1(=L_2)L_3$: Začne výpočet A_1 . Ak akceptuje, spustí sa naraz výpočet A_2 a A_3 ... \square

Veta 2. \mathcal{L}_{CF} nie je uzavretá na pozitívny lookaround.

Vieme totiž vygenerovať jazyk $a^n b^n c^n$ prienikom jazykov $a^* b^n c^n$ a $a^n b^n c^*$.

Veta 3. \mathcal{L}_{CS} je uzavretá na pozitívny lookaround.

Veta 4. Trieda jazykov nad *Regex* s pozitívnym a negatívnym lookaroundom je ekvivalentná \mathcal{R} .

3 Chomského hierarchia

Veta 5. $\mathcal{R} \subsetneq \mathcal{L}_{ERE} \subseteq \mathcal{L}_{LERE} \subsetneq \mathcal{L}_{nLERE} \subsetneq \mathcal{L}_{CS}$

Dôkaz. Všetky \subseteq triviálne platia.

$\mathcal{L}_{ERE} \subsetneq \mathcal{L}_{LERE}$: Nerovnosť dokazuje jazyk $L = \{a^i b a^{i+1} b a^i k \mid k = i(i+1)k' \text{ pre nejaké } k' > 0, i > 0\}$. $L \notin Eregex$ podľa pumpovacej lemy z [Carle and Nadendran, 2009] a tu je regex z *LEregex* pre L : $\alpha = (a^*)b(\backslash 1a)b(=?(\backslash 1)*\$)(\backslash 2)^*$

$\mathcal{L}_{LERE}, \mathcal{L}_{nLERE} \subsetneq \mathcal{L}_{CS}$: Triedy \mathcal{L}_{LERE} a \mathcal{L}_{nLERE} sú neporovnateľné s \mathcal{L}_{CF} . K jazyku $L_1 = \{ww \mid w \in \{a,b\}^*\}$ existuje regex z $LEregex$: $\alpha = ((a|b)^*) \setminus 1$. Ani jedna z tried neobsahuje jazyk $L_2 = \{a^n b^n \mid n \in \mathbb{N}\}$. \square

Intuitívne by malo platiť aj $\mathcal{L}_{ERE} \subsetneq \mathcal{L}_{LERE}$, pretože negatívny lookaround pridáva uzavretosť na komplement. Jazyk dokazujúci nerovnosť by mohol byť napríklad regex $\alpha = (?!(aa+)(\setminus 1)+\$)$, pričom $L(\alpha) = \{a^p \mid p \text{ je prvočíslo}\}$.

4 Vlastnosti triedy \mathcal{L}_{LERE}

Očividne operácia lookahead/lookbehind pridala uzavretosť na prienik. Nech $\alpha, \beta \in LEregex$, potom $L(\alpha) \cap L(\beta) = L(\gamma)$, kde $\gamma = (?=\alpha\$)\beta$ alebo $\beta(?<=\hat{\alpha})$.

Napak ohrozila uzavretosť na základnú operáciu – zret'azenie. Pri zret'azení 2 jazykov, ktorých regexy nutne musia obsahovať lookahead resp. lookbehind nastáva problém. Nemôžeme tieto regexy len tak položiť za seba. Ak sa napríklad v prvom z jazykov nachádza lookahead, počas výpočtu môže zasahovať aj do časti vstupu, ktorú matchuje druhý regex a tým zmeniť výsledok celého výpočtu. Nakoniec sa ukázalo:

Veta 6. \mathcal{L}_{LERE} je uzavretá na zret'azenie.

Dôkaz. Nech $\alpha, \beta \in LEregex$. Jazyku $L(\alpha)L(\beta)$ bude zodpovedať regex

$$\gamma = (?=(\alpha) (\beta) \$) \alpha' \setminus k+2 (?<=\hat{\setminus 1} \beta')$$

$\begin{matrix} 1 & 1 & k+2 & k+2 \end{matrix}$

V α, β treba vhodne prepísať označenie zátvoriek (po poradí). α' je α prepísaný tak, že pre každý lookahead:

- bez $\$$ – na koniec pridáme $. * \setminus k+2 \$$
- s $\$$ – pred $\$$ pridáme $\setminus k+2$

β' je β prepísaný tak, že pre každý lookbehind:

- bez $\hat{}$ – na začiatok pridáme $\hat{\setminus 1} . *$
- s $\hat{}$ – pred $\hat{}$ pridáme $\hat{\setminus 1}$

Slovami vyjadrené, regex γ najprv rozdelí vstupné slovo na 2 podslová w_1, w_2 patriace do príslušných jazykov $L(\alpha), L(\beta)$. Potom spustí ešte raz regex α

upravený tak, že jeho lookaheady sú „skrotené“, pretože ich na konci donúti matchovať w_2 . Rovnako lookbehindy v β' donúti na začiatku matchovať w_1 , až potom normálne pokračuje ich výpočet.

Zrejme $L(\gamma) = L(\alpha)L(\beta)$. \square

Veta 7. Nech $\alpha \in LEregex$ nad unárnou abecedou $\Sigma = \{a\}$, že neobsahuje lookahead s $\$$ ani lookbehind s $\hat{}$ vnútri iterácie. Existuje konštanta N taká, že ak $w \in L(\alpha)$ a $|w| > N$, potom existuje dekompozícia $w = xy$ s nasledujúcimi vlastnosťami:

- (i) $|y| \geq 1$
- (ii) $\exists k \in \mathbb{N}, k \neq 0; \forall j = 1, 2, \dots : xy^{kj} \in L(\alpha)$

Veta 8. Jazyk všetkých platných výpočtov Turingovho stroja patrí do \mathcal{L}_{LERE} .

Dôkaz. Takýto jazyk pre konkrétny Turingov stroj obsahuje slová, ktoré sú tvorené postupnosťou konfigurácií oddelených oddeľovačom $\#$, ktoré zodpovedajú akceptačným výpočtom na všetkých možných vstupoch.

Turingov stroj má konečný zápis, preto je možné regex pre takýto jazyk vytvoriť. Konštrukcia regexu: $\alpha = \beta(\gamma) * \eta$, kde β predstavuje počiatočnú konfiguráciu² a η akceptačnú konfiguráciu. Ak q_0 je akceptačný stav, potom na koniec α pridáme $|(\#q_0. * \#)$. $\gamma = \gamma_1 \mid \gamma_2 \mid \gamma_3$. Prvok γ_i generuje validnú konfiguráciu a zároveň kontroluje pomocou lookaheadu, či nasledujúca konfigurácia môže podľa δ -funkcie nasledovať. Rozpíšeme si iba jednu možnosť:

$$\gamma_1 = ((\underset{k}{.} \underset{k}{*}) \underset{k}{x} \underset{k}{q} \underset{k+1}{y} (\underset{k+1}{.} \underset{k+1}{*}) \underset{k+1}{\#}) (? = \xi \#)$$

platí pre $\forall q \in K, \forall y \in \Sigma$ a kde $\xi = \xi_1 \mid \xi_2 \mid \dots \mid \xi_n$.

- Ak $(p, z, 0) \in \delta(q, y)$, potom $\xi_i = (\setminus k x p z \setminus k+1)$ pre nejaké i
- Ak $(p, z, 1) \in \delta(q, y)$, potom $\xi_i = (\setminus k x z p \setminus k+1)$ pre nejaké i
- Ak $(p, z, -1) \in \delta(q, y)$, potom $\xi_i = (\setminus k p x z \setminus k+1)$ pre nejaké i

γ_2 a γ_3 sú podobné ako γ_1 , ale matchujú krajné prípady, kedy je hlava Turingovho stroja na ľavom alebo pravom konci pásky. \square

²Musí byť previazaná s nasledujúcou konfiguráciou, aby spĺňala δ -funkciu. Spraví sa to podobne ako v γ_1 .

5 Priestorová zložitosť

Pod'akovanie

Ďakujem školiteľovi za cenné rady a pripomienky.

Literatúra

- [Carle and Nadendran, 2009] Carle, B. and Nadendran, P. (2009). On extended regular expressions. In *Language and Automata Theory and Applications*, volume 3, pages 279–289. Springer.
- [Câmpeanu et al., 2003] Câmpeanu, C., Salomaa, K., and Yu, S. (2003). A formal study of practical regular expressions. *International Journal of Foundations of Computer Science*, 14(06):1007–1018.
- [Ehrenfeucht and Zeiger, 1975] Ehrenfeucht, A. and Zeiger, P. (1975). Complexity measures for regular expressions. *Computer Science Technical Reports*, 64.
- [Ellul et al., 2013] Ellul, K., Krawetz, B., Shallit, J., and wei Wang, M. (2013). Regular expressions: New results and open problems. *Journal of Automata, Languages and Combinatorics* $u(v)w, x-y$.
- [Foundation, 2012] Foundation, P. S. (2012). *Regular expressions operations*.
- [Tóthová, 2013] Tóthová, T. (2013). Moderné regulárne výrazy. Bachelor's thesis, FMFI UK Bratislava.