

# Regulární výrazy

## Regulární výrazy

Regulární výrazy jsou jedním z prostředků, které můžeme používat k popisu tzv. regulárních jazyků. Nejenom s nimi se seznámíte v průběhu letního semestru v rámci předmětu Automaty a gramatiky. Abychom s nimi ale dokázali pracovat po praktické stránce, bez většího teoretického zázemí se určitě obejdeme.

S regulárními výrazy se setkáváme v řadě kontextů, nejenom v rámci např. programovacích jazyků nebo dotazovacích jazyků v databázových systémech. Obvykle pracujeme s variantami, které nabízí celou řadu uživatelsky přívětivých zjednodušení a zkratk, bez nich se však jsme schopni obejít. Podíváme se tedy na regulární výrazy v podobě, která je nezbytně nutná pro dosažené očekávané vyjadřovací síly.

Regulární výraz  $v$  nad nějakou abecedou  $\Sigma$  (např. symbolů anglické abecedy) definujeme induktivním způsobem takto. Nejprve zavedeme následující jednoduché regulární výrazy:

- $a$  pro každé  $a \in \Sigma$  odpovídající jazyku  $h(a) = \{a\}$ , tedy jazyku, který obsahuje jeden jednosymbolový řetězec  $a$ ,
- $\varepsilon$  (prázdný řetězec) odpovídající jazyku  $h(\varepsilon) = \{\varepsilon\}$ , tedy jazyku, který obsahuje jen prázdný řetězec, který značíme jako  $\varepsilon$  (někdy také  $\lambda$ ) a
- $\emptyset$  (prázdný jazyk) odpovídající prázdnému jazyku  $h(\emptyset) = \{\}$ , tedy jazyku, který neobsahuje žádný řetězec.

Jsou-li  $r$  a  $s$  již definované regulární výrazy (jakkoli složité), jsme pomocí nich a následujících operací schopni zkonstruovat následující složené výrazy takto:

- $r \cdot s$  (operace  $\cdot$  zřetězení) odpovídající jazyku  $h(r \cdot s) = h(r) \cdot h(s) = \{uv \mid u \in h(r), v \in h(s)\}$ , tedy jazyku, který obsahuje všechny řetězce vzniklé tak, že je umíme rozdělit na dva podřetězce takové, že první z nich se dá vygenerovat prvním výrazem  $r$  a druhý druhým  $s$ ,
- $r + s$  (operace  $+$  alternativa) odpovídající jazyku  $h(r + s) = h(r) \cup h(s)$ , tedy jazyku, který obsahuje všechny řetězce, které umí vygenerovat první výraz  $r$  nebo druhý  $s$  a
- $r^*$  (operace  $*$  iterace) odpovídající jazyku  $h(r^*) = \bigcup_{i \in \mathbb{N}_0, i \geq 0} L^i$ , kde  $L^0 = \{\varepsilon\}$  a  $L^i = L \cdot L^{i-1}$  pro  $\forall i \in \mathbb{N}, i \geq 1$ , tedy jazyku, který obsahuje všechny řetězce vzniklé tak, že je umíme rozdělit na libovolný počet podřetězců, z nichž každý se dá vygenerovat výrazem  $r$ .

Pro získání jistoty správného vyhodnocení regulárního výrazu je potřeba vhodně používat pomocné kulaté závorky. To však často není úplně nezbytně nutné, stejně tak si můžeme dovolit i další zjednodušení. My budeme konkrétně uvažovat následující dvě:

- Pokud to není nutné, závorky psát nemusíme, takže např.  $((a + b) + c)$  odpovídá  $a + b + c$  a  $((a \cdot b) \cdot c)$  odpovídá  $a \cdot b \cdot c$ .
- Pokud zřetěžueme několik symbolů naší abecedy, nebudeme psát dokonce ani operaci  $\cdot$  jako takovou, takže např.  $abba$  odpovídá  $a \cdot b \cdot b \cdot a$ . Tuto myšlenku můžeme zobecnit tak, že po jakémkoli symbolu  $a \in \Sigma$  nebo zavírací závorce  $)$  nebo operátoru iterace  $*$  může bezprostředně následovat další libovolný symbol  $a \in \Sigma$  nebo otevírací závorka  $($ , což ve všech těchto případech opět vede na implicitní operaci zřetězení. Např.  $(a + b)^*aa(a + b)$  odpovídá  $(a + b)^* \cdot a \cdot a \cdot (a + b)$ .

Abychom se ujistili, že konstrukci a významu regulárních výrazů rozumíme, pojďme se podívat na několik příkladů nad abecedou  $\Sigma = \{a, b\}$ :

- $(a + b)^*abba(a + b)^*$  popisuje jazyk  $\{w \mid w = uabba v, u, v \in \{a, b\}^*\}$  aneb řetězce obsahující  $abba$
- $((a + b)(a + b)(a + b))^*(a + b)$  popisuje jazyk  $\{w \mid w \in \{a, b\}^*, |w| = 3k + 1, k \in \mathbb{N}_0\}$ , tedy jazyk obsahující řetězce s délkou dělitelnou 3 se zbytkem 1
- $a(a + b)^*a + b(a + b)^*b + a + b$  popisuje jazyk  $\{w \mid w \in \{a, b\}^*, w \text{ začíná a končí na stejný symbol}\}$

## Metoda sousedů

Častým problémem, který ve spojitosti s regulárními výrazy potřebujeme řešit, je situace, kdy pro daný řetězec máme rozhodnout, jestli odpovídá danému regulárnímu výrazu, tedy jestli jej dokáže vygenerovat. Např. řetězec **aba** našemu předchozímu poslednímu uvedenému regulárnímu výrazu (řetězce začínající a končící na stejný symbol) odpovídá, zatímco řetězec **abb** nikoli.

Uvedený problém je v praxi řešen pomocí tzv. konečných automatů. Takový automat však nejprve pro zadaný výraz musíme umět zkonstruovat. A k tomu slouží řada metod, např. metoda sousedů navržená Viktorem Michajlovičem Gluškovem. Abychom ji mohli použít, musíme nejprve pro zadaný regulární výraz induktivně spočítat čtveřici pomocných funkcí, na jejichž základě takový automat následně umíme zkonstruovat. To však není naším cílem (naučíte se to až v letním semestru), tím je jen vyčíslení těchto funkcí.

Nejprve však potřebujeme ke každému výskytu nějakého symbolu abecedy  $a \in \Sigma$  v zadaném regulárním výrazu přiřadit pomocné číslo, označme jej jako index, pomocí kterého dokážeme navzájem rozlišit jednotlivé výskyty stejného symbolu. To se dá udělat různými způsoby, např. přidělováním postupných čísel. Aneb pro regulární výraz  $(a + b)^*ab$  získáme jeho oindexovanou verzi jako  $(a_1 + b_2)^*a_3b_4$ . Pozor,  $\varepsilon$  ani  $\emptyset$  nejsou symboly abecedy, a tedy indexy jim nepřisuzujeme.

Teď už jen musíme popsat, jaké pomocné funkce uvažujeme a jak je dokážeme pro libovolně komplikovaný oindexovaný regulární výraz  $r'$  spočítat:

- **Starting**( $r'$ ) označuje množinu všech oindexovaných symbolů abecedy, na které může začínat nějaký řetězec odpovídající  $r'$ .
- **Neighbors**( $r'$ ) označuje množinu všech dvojic oindexovaných symbolů, které se mohou vyskytovat bezprostředně za sebou v nějakém řetězci odpovídajícímu  $r'$ .
- **Ending**( $r'$ ) označuje množinu všech oindexovaných symbolů abecedy, na které může končit nějaký řetězec odpovídající  $r'$ .
- **Epsilon**( $r'$ )  $\in \{\text{true}, \text{false}\}$  označuje příznak, zda  $r'$  dokáže vygenerovat prázdný řetězec  $\varepsilon$ .

Pro náš ukázkový regulární výraz  $r = (a + b)^*ab$  a jeho oindexovanou verzi  $r' = (a_1 + b_2)^*a_3b_4$  pak konkrétně získáme tyto hodnoty:

- **Starting**( $r'$ ) =  $\{a_1, b_2, a_3\}$
- **Neighbors**( $r'$ ) =  $\{a_1a_1, a_1b_2, a_1a_3, b_2a_1, b_2b_2, b_2a_3, a_3b_4\}$
- **Ending**( $r'$ ) =  $\{b_4\}$
- **Epsilon**( $r'$ ) = **false**

Chceme-li umět tyto čtyři funkce spočítat obecně, stačí postupovat podle induktivní struktury našeho regulárního výrazu. Aneb pro jednoduché výrazy vše vyjádříme triviálně, u výrazů vzniklých jednotlivými operacemi vyjdeme ze znalosti těchto funkcí u jednotlivých operandů. Ty jsou jednodušší, a tedy je opravdu vyčíslit umíme. Začneme nejprve již zmíněnými jednoduchými výrazy:

- Oindexovaný symbol  $a_i$ ,  $a \in \Sigma$ ,  $i \in \mathbb{N}$ 
  - **Starting**( $a_i$ ) =  $\{a_i\}$ , **Neighbors**( $a_i$ ) =  $\{\}$ , **Ending**( $a_i$ ) =  $\{a_i\}$ , **Epsilon**( $a_i$ ) = **false**
- Prázdný řetězec  $\varepsilon$ 
  - **Starting**( $\varepsilon$ ) =  $\{\}$ , **Neighbors**( $\varepsilon$ ) =  $\{\}$ , **Ending**( $\varepsilon$ ) =  $\{\}$ , **Epsilon**( $\varepsilon$ ) = **true**
- Prázdný jazyk  $\emptyset$ 
  - **Starting**( $\emptyset$ ) =  $\{\}$ , **Neighbors**( $\emptyset$ ) =  $\{\}$ , **Ending**( $\emptyset$ ) =  $\{\}$ , **Epsilon**( $\emptyset$ ) = **false**

Jsou-li  $r$  a  $s$  libovolně složité oindexované výrazy, pak pro jednotlivé operace spočítáme naše pomocné funkce takto:

- Zřetězení  $r \cdot s$

- $\text{Starting}(r \cdot s) = \begin{cases} \text{Starting}(r) & \text{pokud } \text{Epsilon}(r) = \text{false} \\ \text{Starting}(r) \cup \text{Starting}(s) & \text{v opačném případě} \end{cases}$
- $\text{Neighbors}(r \cdot s) = \text{Neighbors}(r) \cup \text{Neighbors}(s) \cup \{xy \mid x \in \text{Ending}(r), y \in \text{Starting}(s)\}$
- $\text{Ending}(r \cdot s) = \begin{cases} \text{Ending}(s) & \text{pokud } \text{Epsilon}(s) = \text{false} \\ \text{Ending}(s) \cup \text{Ending}(r) & \text{v opačném případě} \end{cases}$
- $\text{Epsilon}(r \cdot s) = \text{Epsilon}(r) \wedge \text{Epsilon}(s)$
- Alternativa  $r + s$ 
  - $\text{Starting}(r + s) = \text{Starting}(r) \cup \text{Starting}(s)$
  - $\text{Neighbors}(r + s) = \text{Neighbors}(r) \cup \text{Neighbors}(s)$
  - $\text{Ending}(r + s) = \text{Ending}(r) \cup \text{Ending}(s)$
  - $\text{Epsilon}(r + s) = \text{Epsilon}(r) \vee \text{Epsilon}(s)$
- Iterace  $r^*$ 
  - $\text{Starting}(r^*) = \text{Starting}(r)$
  - $\text{Neighbors}(r^*) = \text{Neighbors}(r) \cup \{xy \mid x \in \text{Ending}(r), y \in \text{Starting}(r)\}$
  - $\text{Ending}(r^*) = \text{Ending}(r)$
  - $\text{Epsilon}(r^*) = \text{true}$

## Zadání úkolu

Cílem úkolu je naprogramovat jednoduchou aplikaci, která pro množinu regulárních výrazů na vstupu spočítá uvedené pomocné funkce používané v metodě sousedů a vypíše je v určeném formátu na standardní výstup.

Vstupní regulární výrazy mohou být zadány dvojím způsobem: buď je najdeme přímo jako jednotlivé argumenty předané z příkazové řádky nebo je najdeme uložené v textových souborech, jejichž jména jsou opět předána formou argumentů. Pro rozlišení obou situací budeme používat dva přepínače, a to **-a** pro první způsob chování (argumenty) a **-f** pro druhý (soubory).

Jednotlivé předané argumenty budeme zpracovávat jeden za druhým, zleva doprava. Pokud najdeme nějaký z uvedených přepínačů, všechny další argumenty (žádný nebo více) budeme zpracovávat příslušným způsobem (argument/soubor), a to dokud nenajdeme další přepínač nebo už jsme zpracovali všechny argumenty. Pokud první argument není přepínačem, předpokládáme výchozí režim **-a**. Všechny uvedené soubory budou vypadat tak, že na každém řádku bude jeden regulární výraz (a nic dalšího), případné prázdné řádky budeme přeskakovat.

Pro ilustraci, vstupní argumenty mohou vypadat např. takto: **(a+b) aa\*+\epsilon+\emptyset\* ab\* -f file1.txt file2.txt file3.txt -a (aa)\*bb**.

Všechny regulární výrazy na vstupu budou nad abecedou anglických písmen. Pro operace budeme používat symboly **.** (zřetězení), **+** (alternativa) a **\*** (iterace), pro jednoduché výrazy prázdného řetězce použijeme namísto  $\varepsilon$  řetězec **\epsilon** a analogicky u prázdného jazyka namísto  $\emptyset$  řetězec **\emptyset**. Používat také můžeme pomocné kulaté závorky **( a )**. V souladu s vysvětlenými pravidly tyto závorky mohou být vynechány, stejně jako operace zřetězení. Můžeme předpokládat, že všechny regulární výrazy jsou korektní aneb syntakticky správně formované.

Pro vlastní parsování regulárního výrazu použijte rozšířený shunting-yard algoritmus, s jeho pomocí rovnou vytvořte syntaktický strom odpovídající induktivní struktuře našeho regulárního výrazu. Na základě průchodu tímto stromem spočítejte (a ve vhodné podobě ponechte uložené u každého uzlu) všechny potřebné funkce **Starting**, **Neighbors**, **Ending** a **Epsilon**. Pozor, níže uvedený pseudokód algoritmu neřeší implicitní (chybějící) operátory zřetězení. Z hlediska operátorů očekávejte následující vlastnosti:

- Operátor **\*** je unární postfixový a má nejvyšší precedenci
- Operátor **.** je binární infixový zleva asociativní a má střední precedenci
- Operátor **+** je binární infixový zleva asociativní a má nejnižší precedenci

---

```

1 inicializuj generátor indexů na hodnotu 0
2 foreach token  $t$  ve vstupním výrazu do
3   if  $t$  je symbol abecedy then
4     inkrementuj  $o$  1 hodnotu generátoru indexů na  $i$ 
5     vytvoř pro  $t$  na základě  $i$  nový listový uzel a vlož jej do zásobníku operandů
6   else if  $t$  je otevírací závorka ( then
7     dej ( na zásobník operátorů
8   else if  $t$  je zavírací závorka ) then
9     while na vrcholu zásobníku operátorů je nějaký operátor  $o$  do
10      odeber  $o$  ze zásobníku operátorů
11      odeber dva uzly  $r$  a  $l$  ze zásobníku operandů
12      vytvoř pro  $o$  na základě  $l$  a  $r$  nový vnitřní uzel a vlož jej do zásobníku operandů
13      odeber ( ze zásobníku operátorů
14   else if  $t$  je unární postfixový operátor  $n$  then
15     odeber ze zásobníku operandů jeden uzel  $p$ 
16     vytvoř pro  $n$  na základě  $p$  nový vnitřní uzel a vlož jej do zásobníku operandů
17   else  $t$  je binární infixový operátor  $n$ 
18     while na zásobníku operátorů je operátor  $o$  s precedencí vyšší než  $n$  nebo také stejnou, je-li
19       ovšem  $n$  současně zleva asociativní do
20       odeber  $o$  ze zásobníku operátorů
21       odeber dva uzly  $r$  a  $l$  ze zásobníku operandů
22       vytvoř pro  $o$  na základě  $l$  a  $r$  nový vnitřní uzel a vlož jej do zásobníku operandů
23     přidej  $n$  na zásobník operátorů
24 while zásobník operátorů je neprázdný do
25   odeber  $o$  ze zásobníku operátorů
26   odeber dva uzly  $r$  a  $l$  ze zásobníku operandů
27   vytvoř pro  $o$  na základě  $l$  a  $r$  nový vnitřní uzel a vlož jej do zásobníku operandů
28 zásobník operandů nyní obsahuje jediný uzel, a to kořenový uzel celého stromu

```

---

Vstupní regulární výrazy zpracujte ve stejném pořadí, v jakém jste je objevili na vstupu, respektujte tedy pořadí předaných argumentů a pořadí výrazů uvnitř souborů. Pro každý výraz vypište výstup v následujícím formátu:

(aa+bb)\*

- Starting: {a:1, b:3}
- Neighbors: {(a:1, a:2), (a:2, a:1), (a:2, b:3), (b:3, b:4), (b:4, a:1), (b:4, b:3)}
- Ending: {a:2, b:4}
- Epsilon: true

ab\*a

- Starting: {a:1}
- Neighbors: {(a:1, b:2), (a:1, a:3), (b:2, b:2), (b:2, a:3)}
- Ending: {a:3}
- Epsilon: false

Každý blok pro jeden konkrétní regulární výraz ukončíte jedním dalším prázdným řádkem (mezi bloky a za posledním tedy bude vždy jeden volný řádek). Regulární výraz jako takový vypíšete přesně v té podobě, v jaké jste jej dostali na vstupu (aneb neprogramujte vlastní funkci, která by výraz dokázala vypsát vhodným průchodem syntaktického stromu, ale opravdu přímo využijte výraz ze zadání). Oindexované symboly u funkcí **Starting** a **Ending** seřadte vzestupně podle indexů těchto symbolů. Dvojice sousedních symbolů u funkce **Neighbors** seřadte primárně podle indexu prvního symbolu a sekundárně podle indexu druhého. Funkce **Epsilon** vypíše **true** nebo **false**.

Pro implementaci funkcí **Starting**, **Ending** i **Neighbors** použijte standardní kontejner pro množinu, tedy kontejner `std::set` dostupný v hlavičkovém souboru `<set>`. Aby to bylo možné, je potřeba pro vkládané objekty implementovat porovnávání, konkrétně operátor `<`. Toho snadno dosáhneme např. prostřednictvím globální funkce `bool operator<(const Symbol& left, const Symbol& right)`. `Symbol` je třída, která naše objekty reprezentuje, funkce samotná vrátí `true` právě tehdy, když první operand je ostře menší než druhý.

Pokud veškerá činnost celého programu proběhne bez problémů, vrátí aplikace kód 0. V případě jakékoli chyby pak kód 1. V takovém případě obsah na standardním výstupu není podstatný. Připomeňme, že vstupní regulární výrazy budou vždy korektní, v jiných případech ale případné chyby (např. během práce se soubory) už ošetřit potřeba je.

Cílem úkolu je ukázat schopnost práce s konstrukty, se kterými jsme se už od začátku semestru setkali. Kromě základních dovedností tedy jde o práci s argumenty programu, textovými soubory a obecně streamy, funkcemi, předávání parametrů, návrh tříd, použití konstruktorů a destruktorů, dědičnost, virtuální metody a stejně tak i dynamickou alokaci.

Předložená implementace pochopitelně musí být korektní a stabilní, kompilace musí proběhnout bez jakýchkoli varování. Hodnotit se ale bude i celková kvalita kódu. Tedy zejména avšak ne výlučně organizace kódu do jednotlivých souborů, tříd a funkcí, použití hlavičkových souborů, pojmenovávání souborů, funkcí nebo i proměnných, celkový vizuální styl kódu a indentace, předávání argumentů hodnotou nebo referencí, kvalita návrhu tříd a použití dědičnosti a virtuálních metod, zbytečné neopakování stejného kódu, deklarace významných konstant, ošetřování chybových situací stejně jako využívání standardních knihoven, kontejnerů nebo funkcí.

Odevzdejte jen zdrojové soubory (`*.cpp`, `*.h`, `*.hpp`).