

Tik-79.148
Tietojenkäsittelyteorian perusteet
Laskuharjoitus 3
Ratkaisut

Kevät 2001

Säännölliset lausekkeet määritellään induktiivisesti:

- \emptyset ja kaikki $a \in \Sigma$ ovat säännöllisiä lausekkeita.
- Mikäli α ja β ovat säännöllisiä lausekkeita, niin myös $(\alpha\beta)$, $(\alpha \cup \beta)$ ja α^* ovat säännöllisiä lausekkeita.
- Mitkään muut lausekkeet eivät ole säännöllisiä.

Säännöllisestä lausekkeesta \emptyset^* käytetään lyhennysmerkintää e . Yleensä jätetään ylimääräiset sulut pois: $((a(bc)) \cup (cd^*)) = abc \cup cd^*$.

Säännöllinen lauseke α määrittelee kielen $L(\alpha)$, joka on joukon Σ^* osajoukko. Esimerkiksi:

$$L(ab \cup ac^*d) = \{ab, ad, acd, accd, accd, acccd, acccd, \dots\} .$$

Mikäli sekaantumisen vaaraa ei ole¹, jätetään $L()$ melkein aina pois säännöllisen lausekkeen ympäriltä, ja käytetään lauseketta itseään tarkoittamaan niiden sanojen joukkoa, jotka kuuluvat sen kuvaamaan kieleen.

Säännölliset lausekkeet on ehkä helpointa ajatella muottina. Sana kuuluu kieleen, jos sen saa sovitettua lausekkeen antamaan muottiin:

- $ab \Rightarrow$ otetaan ensin a , sitten b
- $a \cup b \Rightarrow$ otetaan joko a tai b
- $a^* \Rightarrow$ otetaan kuinka monta a :ta tahansa (myös 0).

Oppikirja esittää hieman epäselvästi lausekkeiden \emptyset ja \emptyset^* välisen eron ja siirtyy käyttämään jälkimmäisestä merkintää e määrittelemättä sitä sen tarkemmin. Määritelmän mukaan $L(\emptyset) = \emptyset$, eli tyhjä kieli, johon ei kuulu yhtään sanaa. Vastaavasti $L(\emptyset^*) = L(\emptyset)^*$, eli kieli, joka muodostetaan seuraavasti:

$$L(\emptyset^*) = \{w \mid w = w_1 \circ w_2 \circ \dots \circ w_k, \text{ jollekin } k \geq 0 \text{ ja } w_1, \dots, w_k \in \emptyset\}$$

Koska tyhjään joukkoon ei kuulu yhtään sanaa, on ainoa mahdollinen k :n arvo yllä on 0. Koska määritelmän perusteella nollan merkkijonon katenaatio on tyhjä merkkijono e , saadaan kieleksi

$$L(\emptyset^*) = \{e\}$$

¹Ja valitettavan usein silloinkin, kun sekaantumisen vaara on olemassa, kuten tehtävän 4. b- ja c-kohdissa.

Toisin sanoen, säännöllinen lauseke \emptyset tarkoittaa kieltä, johon ei kuulu yhtään sanaa, ja lauseke \emptyset^* tarkoittaa kieltä, johon kuuluu ainoastaan tyhjä sana e .

Säännöllisten lausekkeiden unionia voidaan merkitä monilla eri tavoilla. Yleisimmät ovat \cup , $+$ sekä $|$. Useimmissa käytännön toteutuksissa määritellään lisäksi ylimääräisiä operaattoreita yksinkertaistamaan lausekkeitä. Lausekkeiden ilmaisuvoima pysyy kuitenkin yleensä samana.

Säännöllisiä lausekkeitä käytetään mm. ohjelmointikielten kääntäjien syntaktiseen analyysiin sekä määrämuotoisen tekstin etsimiseen tekstitiedoista (grep).

- 4 a) Väite $baa \in a^*b^*a^*b^*$ on tosi, sillä valitsemalla yksi ensimmäisistä b -kirjaimista ja kaksi toisista a -kirjaimista saadaan sana baa .
 b) Väite $b^*a^* \cap a^*b^* = a^* \cup b^*$ on tosi, sillä

$$\begin{aligned} b^*a^* \cap a^*b^* &= \\ \{e, b^+, a^+, b^+a^+\} \cap \{e, a^+, b^+, a^+b^+\} &= \\ \{e, a^+, b^+\} &= a^* \cup b^* \end{aligned}$$

Merkintä $a^+ = aa^*$

Tässä kannattaa huomata, että yhtälön vasemmalla puolella oleva merkintä **ei** ole säännöllinen lauseke, sillä säännöllisessä lausekkeessa ei voi esiintyä leikkausta, vaan kyseessä on kahden säännöllisen kielen leikkaus.

- c) Väite $a^*b^* \cap c^*d^* = \emptyset$ on väärä, sillä

$$\begin{aligned} a^*b^* \cup c^*d^* &= \\ \{e, b^+, a^+, b^+a^+\} \cap \{e, c^+, d^+, c^+d^+\} &= \\ \{e\} &\neq \emptyset \end{aligned}$$

- d) Väite $abcd \in (a(cd)^*b)^*$ on myös väärä, sillä kaikki säännöllisen lausekkeen hyväksymät sanat (tyhjää sanaa lukuunottamatta) päättyvät b -kirjaimen.

- 5 Osoitetaan väite $a(b \cup c) = ab \cup ac$. Määritelmän mukaan lause $b \cup c$ kuvaa joukkoa $\{b, c\}$. Konkatenation määritelmästä saadaan että $a(b \cup c)$ kuvaa joukkoa $\{a\} \circ \{b, c\} = \{ab, ac\}$ mistä väite seuraa.

- 6 Väite: Jos kieli L on säännöllinen, niin myös kieli

$$L' = \{w \mid uw \in L \text{ jollekin merkkijonolle } u\}$$

on säännöllinen.

Kieleen L' kuuluvat kaikki ne merkkijonot, jotka esiintyvät kieleen L kuuluvien sanojen suffiksina. Esimerkiksi: jos

$$\begin{aligned} L &= \{a, aba, abb\}, \text{ niin} \\ L' &= \{e, a, b, ba, bb, aba, abb\} \end{aligned}$$

Tässä konstruoidaan systemaattinen tapa muodostaa L :n määrittelemästä säännöllisestä lausekkeesta R (eli $L(R) = L$) uusi lauseke R' , jonka määrittelemä kieli on L' (eli $L(R') = L'$).

Formaalimmin, merkitään kaikkien säännöllisten lausekkeiden joukkoa symbolilla S . Nyt halutaan löytää jokin funktio $f : S \rightarrow S$, siten, että $f(R) = R'$.

Määritellään funktio f rekursiivisesti:

- i) (perustapaus) $f(\emptyset) = \emptyset$ ja $f(a) = a \cup e, a \in \Sigma \cup \{e\}$. Toisin sanoen, mikäli säännöllinen lauseke R on tyhjä ei sillä ole yhtään suffiksia, ja jos se on pelkästään yksi aakkoston merkki tai e , on lauseke R' saman merkin ja e :n unioni.
- ii) $f(\alpha \cup \beta) = f(\alpha) \cup f(\beta), \alpha, \beta \in S$. Toisin sanoen, mikäli R on kahden lausekkeen α ja β unioni, niin sen määrittelemän kielen sanojen lopussa voi olla mitä tahansa, mitä voi olla jomman kumman alilausekkeen lopussa.
- iii) $f(\alpha^*) = f(\alpha)\alpha^*$, eli sanan lopussa voi olla tähden vaikutusalueella oleva alilauseke kokonaisuutensa mielivaltaisen monta kertaa, ja sen edellä jotain, mihin alilauseke voi päättyä
- iv) $f(\alpha\beta) = f(\alpha)\beta \cup f(\beta)$, eli katenaatiolla muodostetun lausekkeen kuvaaman kielen sanojen lopussa voi joko olla jotain, mitä voi olla jälkimmäisen osan lopussa, tai sitten viimeinen osa on kokonaisuudessaan mukana, ja sen edellä on jotain, mihin ensimmäinen osa voi loppua.

Tämän jälkeen puuttuu enää formaali todistus ylläolevien sääntöjen oikeellisuudesta. Todistus sivuutetaan tässä, mutta sen voi rakentaa osoittamalla erikseen jokaisen säännön oikeellisuuden induktiolla.

Esimerkki: Olkoon $R = ba(aa \cup bab)^*$. Nyt

$$\begin{aligned}
 f(R) &= f(ba(aa \cup bab)^*) \\
 &= f(ba)(aa \cup bab)^* \cup f((aa \cup bab)^*) \\
 &= (f(b)a \cup f(a))(aa \cup bab)^* \cup f(aa \cup bab)(aa \cup bab)^* \\
 &= ((e \cup b)a \cup (e \cup a))(aa \cup bab)^* \cup (f(aa) \cup f(bab))(aa \cup bab)^* \\
 &= (e \cup a \cup ba)(aa \cup bab)^* \cup ((f(a)a \cup f(a)) \cup (f(b)ab \cup f(a)b \cup f(b)))(aa \cup bab)^* \\
 &= (e \cup a \cup ba)(aa \cup bab)^* \cup (e \cup a \cup aa \cup b \cup ab \cup bab)(aa \cup bab)^* \\
 &= (e \cup a \cup aa \cup b \cup ab \cup bab \cup ba)(aa \cup bab)^* \\
 &= (e \cup a \cup b \cup ab \cup ba)(aa \cup bab)^*
 \end{aligned}$$