

1.6 Aakkostot, merkkijonot ja kielet

Automaattiteoria ~ diskreetin signaalinkäsittelyn perusmallit ja -menetelmät

(~ diskreettien I/O-kuvausten yleinen teoria)



Automaatin käsite on *matemaattinen abstraktio*. Yleisellä tasolla suunniteltu automaatti voidaan toteuttaa eri tavoin: esim. sähköpiirinä, mekaanisena laitteena tai (tavallisimmin) tietokoneohjelmana.

Tällä kurssilla keskitytään pääosin automaatteihin, joiden:

- (i) syötteet ovat äärellisiä, diskreettejä *merkkijonoja*
- (ii) tulokset ovat muotoa "hyväksy"/"hylkää"
(~ "syöte OK"/"syöte ei kelpaa")

Yleistyksiä:

- (i) äärettömät syötejonot (\rightarrow "reaktiiviset" järjestelmät, Büchi-automaatit)
- (ii) funktioautomaatit (\rightarrow Moore- ja Mealytilakoneet, Turingin funktiokoneet)

Peruskäsitteitä ja merkintöjä

Aakkosto (engl. alphabet, vocabulary): mikä tahansa äärellinen, epätyhjä joukko *alkeismerkkejä* t. *symbolia*. Esim.:

- binääriaakkosto $\{0, 1\}$;
- latinalainen aakkosto $\{A, B, \dots, Z\}$.

Merkkijono (engl. string): äärellinen järjestetty jono jonkin aakkoston merkkejä. Esim.:

- "01001", "0000": binääriaakkoston merkkijonoja;
- "TKTP", "XYZZY": latinalaisen aakkoston merkkijonoja.

Erikoistapaus: *tyhjä merkkijono* (engl. empty string). Tyhjässä merkkijonossa ei ole yhtään merkkiä, mutta havaittavuuden parantamiseksi sen paikka usein osoitetaan erikoismerkillä ε .

Merkkijonon x pituutta, so. siihen sisältyvien merkkien määrää, merkitään $|x|$:llä. Esim.: $|01001| = |XYZZY| = 5$, $|0000| = |TKTP| = 4$, $|\varepsilon| = 0$.

Merkkijonojen välinen perusoperaatio on *katenointi* eli jonojen peräkkäin kirjoittaminen. Katenointi operaatiomerkinä käytetään joskus selkeyden lisäämiseksi symbolia \wedge . Esim.:

- (i) $KALA^\wedge KUKKO = KALAKUKKO$;
- (ii) jos $x = 00$ ja $y = 11$, niin $xy = 0011$ ja $yx = 1100$;
- (iii) kaikilla x on $x\varepsilon = \varepsilon x = x$;
- (iv) kaikilla x, y, z on $(xy)z = x(yz)$;
- (v) kaikilla x, y on $|xy| = |x| + |y|$.

Aakkoston Σ kaikkien merkkijonojen joukko merkitään Σ^* :lla. Esimerkiksi jos $\Sigma = \{0, 1\}$, niin $\Sigma^* = \{\varepsilon, 0, 1, 00, 01, 10, \dots\}$.

Mielivaltaista merkkijonojoukkoa $A \subseteq \Sigma^*$ sanotaan aakkoston Σ (formaaliksi) *kieleksi* (engl. formal language).

Automaatit ja formaalit kielet

Olkoon M automaatti, jonka syötteet ovat jokin aakkoston Σ merkkijonoja, ja tulos on yksinkertaisesti muotoa "syöte hyväksytään"/"syöte hylätään". (Merk. lyhyesti 1/0.)

Merkitään M :n syötteellä x antamaa tulosta $M(x)$:llä ja M :n hyväksymien syötteiden joukkoa A_M :llä, so.

$$A_M = \{x \in \Sigma^* \mid M(x) = 1\}.$$

Sanotaan, että automaatti M tunnistaa (engl. recognizes) kielen $A_M \subseteq \Sigma^*$.

Automaattiteorian (yksi) idea: *automaatin M rakenne heijastuu kielen A_M ominaisuuksissa*.

Kääntäen: olkoon annettuna jokin toivottu I/O-kuvaus $f : \Sigma^* \rightarrow \{0, 1\}$. Tarkastelemalla kielitä

$$A_f = \{x \in \Sigma^* \mid f(x) = 1\}$$

saadaan vihjeitä siitä, millainen automaatti tarvitaan kuvauksen f toteuttamiseen.

35

Vakiintuneita merkintöjä

Em. matemaattisille käsitteille käytetyt merkinnät ovat periaatteessa vapaasti valittavissa, mutta esityksen ymmärrettävyyden parantamiseksi tapana pitää tyypin tietyissä käytännöissä. Seuraavat merkintätavat ovat vakiintuneet:

Aakkostot: Σ, Γ, \dots (isoja kreikkalaisia kirjaimia). *Esim.* binäariaakkosto $\Sigma = \{0, 1\}$.

Aakkoston koko (tai yleisemmin joukon matalvuus): $|\Sigma|$.

Alkeismerkit: a, b, c, \dots (pieniä alkupään latinalaisia kirjaimia). *Esim.*: Olkoon $\Sigma = \{a_1, \dots, a_n\}$ aakkosto; tällöin $|\Sigma| = n$.

Merkkijonot: u, v, w, x, y, \dots (pieniä loppupään latinalaisia kirjaimia).

36

Merkkijonojen ketonaatio: $x \hat{\wedge} y$ tai vain xy .

Merkkijonon pituus: $|x|$. *Esimerkkejä*:

- (i) $|abc| = 3$;
- (ii) olkoon $x = a_1 \dots a_m$, $y = b_1 \dots b_n$; tällöin $|xy| = m + n$.

Tyhjä merkkijono: ε .

Merkkijono, jossa on n kappaletta merkkiä a : a^n . *Esimerkkejä*:

- (i) $a^n = \underbrace{aa \dots a}_{n \text{ kpl}}$;
- (ii) $|a^i b^j c^k| = i + j + k$.

Merkkijonon x toisto k kertaa: x^k . *Esimerkkejä*:

- (i) $(ab)^2 = abab$;
- (ii) $|x^k| = k|x|$.

Aakkoston Σ kaikkien merkkijonojen joukko: Σ^* . *Esim.*:

$$\{a, b\}^* = \{\varepsilon, a, b, aa, ab, ba, bb, aaa, aab, \dots\}.$$

Merkkijonoinduktio

Automaattiteoriassa tehdään usein konstruktioita "induktiolla merkkijonon pituuden suhteen." Tämä tarkoittaa, että määritellään ensin toiminto tyhjän merkkijonon ε (tai joskus yksittäisen aakkosmerkin) tapauksessa. Sitten oletetaan, että toiminto on määritelty kaikilla annetun pituisilla merkkijonoilla u ja esitetään, miten se tällöin määritellään yhtä merkkiä pitemmillä merkkijonoilla $w = ua$.

Esimerkki. Olkoon Σ mielivaltainen aakkosto. Merkkijonon $w \in \Sigma^*$ *käänteisjono* (engl. reversal) w^R määritellään induktiivisesti säännöillä:

- (i) $\varepsilon^R = \varepsilon$;
- (ii) jos $w = ua$, $u \in \Sigma^*$, $a \in \Sigma$, niin $w^R = a^R u^R$.

37

38

Väite. Olkoon Σ aakkosto. Kaikilla $x, y \in \Sigma^*$ on voimassa $(xy)^R = y^R x^R$.

Todistus. Induktio merkkijonon y pituuden suhteen.

(i) *Perustapaus* $y = \varepsilon$.

$$(x\varepsilon)^R = x^R = \varepsilon^R x^R.$$

(ii) *Induktioaskel.* Olkoon y muotoa $y = ua$, $u \in \Sigma^*$, $a \in \Sigma$. Oletetaan, että väite on voimassa merkkijonoilla x, u . Tällöin on:

$$\begin{aligned} (xy)^R &= (xua)^R \\ &= a^R(xu)^R && [R:n määritelmä] \\ &= a^R(u^R x^R) && [\text{induktio-oleetus}] \\ &= (a^R u^R) x^R && [\wedge:n liitännäisyys] \\ &= (ua)^R x^R && [R:n määritelmä] \\ &= y^R x^R. \square \end{aligned}$$

Induktivista (joskus sanotaan myös "rekursiivista") määritelmää voidaan tietenkin käyttää laskujen perustana; esim.:

$$\begin{aligned} (011)^R &= 1^R(01)^R = 1^R(1^R 0^R) \\ &= 11^R(0^R \varepsilon^R) = 110^R \varepsilon^R \\ &= 110^R \varepsilon = 110. \end{aligned}$$

Tärkeämpää on kuitenkin konstruktioiden ominaisuuksien todistaminen määritelmää noudattelevalla induktiolla. Esimerkki seuraa...

1.7 Numeroituvat ja ylinumeroituvat joukot

Määritelmä 1.10 Joukko X on *numeroituvasi* ääretön, jos on olemassa bijektio $f : \mathbb{N} \rightarrow X$. Ääretön joukko, joka ei ole numeroituna on *ylinumeroituna*. Joukko on *numeroituna*, jos se on äärellinen tai numeroituvasti ääretön.

Intuitiivisesti sanoen joukko X on numeroituna, jos sen alkiot voidaan järjestää ja indeksoida luonnollisilla luvuilla:

$$X = \{x_0, x_1, x_2, \dots, x_{n-1}\},$$

jos X on n -alkioinen äärellinen joukko ja

$$X = \{x_0, x_1, x_2, \dots\},$$

jos X on numeroituvasti ääretön.

On helppo osoittaa (HT), että numeroituvan joukon kaikki osajoukot ovat myös numeroituvia, mutta toisaalta ylinumeroituvilla joukoilla on sekä numeroituvia että ylinumeroituvia osajoukkoja. Siten ylinumeroituvat joukot ovat jossain mielessä "isompia" kuin numeroituvat.

Lause 1.11 Minkä tahansa aakkoston Σ merkkijonojen joukko Σ^* on numeroituvasti ääretön.

Todistus. Muodostetaan bijektio $f : \mathbb{N} \rightarrow \Sigma^*$ seuraavasti. Olkoon $\Sigma = \{a_1, a_2, \dots, a_n\}$. Kiinnitetään Σ :n merkeille jokin "aakkosjärjestys"; olkoon se $a_1 < a_2 < \dots < a_n$.

Joukon Σ^* merkkijonot voidaan nyt luetella valitun aakkosjärjestyksen suhteen *kanonisessa* t. *leksikografisessa* *järjestyksessä* (engl. canonical t. lexicographic order) seuraavasti:

- (i) ensin luetellaan 0:n mittaiset merkkijonot ($= \varepsilon$), sitten 1:n mittaiset ($= a_1, a_2, \dots, a_n$), sitten 2:n mittaiset jne.;
- (ii) kunkin pituusryhmän sisällä merkkijonot luetellaan aakkosjärjestyksessä.

Bijektio f on siis:

$$\begin{aligned}
 0 &\mapsto \varepsilon \\
 1 &\mapsto a_1 \\
 2 &\mapsto a_2 \\
 &\vdots \quad \vdots \\
 n &\mapsto a_n \\
 n+1 &\mapsto a_1 a_1 \\
 n+2 &\mapsto a_1 a_2 \\
 &\vdots \quad \vdots \\
 2n &\mapsto a_1 a_n \\
 2n+1 &\mapsto a_2 a_1 \\
 &\vdots \quad \vdots \\
 3n &\mapsto a_2 a_n \\
 &\vdots \quad \vdots \\
 n^2+n &\mapsto a_n a_n \\
 n^2+n+1 &\mapsto a_1 a_1 a_1 \\
 n^2+n+2 &\mapsto a_1 a_1 a_2 \\
 &\vdots \quad \vdots \quad \square
 \end{aligned}$$

Mielikäintoinen huomio on, että millä tahansa ohjelointikielellä kirjoitetut ohjelmat ovat oikeastaan vain kielen perusaakkoston (esim. C-kielessä ASCII-merkistön) merkkijonoja. Lauseen 1.11 mukaan minkä tahansa aakkoston merkkijonojen joukko on numeroituvasti ääretön, joten myös millä tahansa ohjelointikielellä mahdollisten ohjelmien joukko on numeroituva.

Seuraavan lauseen mukaan kuitenkin kaikkien formaalien kielten joukko on ylinumeroituva. Formaaleja kieliä on siis "enemmän" kuin mahdollisia tietokoneohjelmia, ja siksi *millään ohjelointikielellä ei voida laata tunnistusautomatteja kaikille formaaleille kielille*. (Tai toisin sanoen: on olemassa "periaatteessa mahdollisia" I/O-kuvauksia, joita ei voida toteuttaa tietokoneella.)

43

44

Lause 1.12 Minkä tahansa aakkoston Σ kaikkien formaalien kielten perhe on ylinumeroituva.

Todistus (ns. Cantorin diagonaliargumentti). Merkitään aakkoston Σ kaikkien formaalien kielten perhettä $\mathcal{P}(\Sigma^*) = \mathcal{A}$. Oletetaan, että todistettavan väitteen vastaisesti olisi olemassa kaikki Σ :n formaalit kielet kattava numerointi:

$$\mathcal{A} = \{A_0, A_1, A_2, \dots\}.$$

Olkoot Σ^* :n merkkijonot kanonisessa järjestyksessä lueteltuna x_0, x_1, x_2, \dots . Määritellään em. numerointeja käyttäen formaali kieli $\tilde{\mathcal{A}}$:

$$\tilde{\mathcal{A}} = \{x_i \in \Sigma^* \mid x_i \notin A_i\}.$$

Koska $\tilde{\mathcal{A}} \in \mathcal{A}$ ja \mathcal{A} :n numerointi oletettiin kattavaksi, pitäisi olla $\tilde{\mathcal{A}} = A_k$ jollakin $k \in \mathbb{N}$. Mutta tällöin olisi $\tilde{\mathcal{A}}$:n määritelmän mukaan

$$x_k \in \tilde{\mathcal{A}} \Leftrightarrow x_k \notin A_k = \tilde{\mathcal{A}}.$$

Saadun ristiriidan takia oletus, että joukko \mathcal{A} on numeroituva, ei voi pitää paikkaansa. \square

Kuvallisesti todistuksen idea voidaan esittää seuraavasti. Muodostetaan kielten A_0, A_1, A_2, \dots ja merkkijonojen x_0, x_1, x_2, \dots "insidenssimatriisi", jonka rivin i sarakkeessa j on arvo 1 jos $x_i \in A_j$ ja muuten 0. Tällöin kieli $\tilde{\mathcal{A}}$ poikkeaa kustakin kielestä A_k matriisin "diagonaalilla":

$\tilde{\mathcal{A}}$	\searrow	A_0	A_1	A_2	A_3	\dots
x_0		1				
		0	0	0	1	\dots
x_1		0	1	0	0	\dots
				0		
x_2		1	1	1	1	\dots
					1	
x_3		0	0	0	0	\dots
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\ddots

45

46

1.8 *Ekskursio: Turingin pysähtymisongelma

Lauseiden 1.11 ja 1.12 mukaan on siis olemassa formaaleja kieliä (I/O-kuvauksia), joita ei voida toteuttaa esim. C-ohjelmilla. Entä jokin konkreettinen esimerkki tällaisesta?

Tunnetuin esimerkki on ns. *Turingin pysähtymisongelma*. (Alan Turing, 1936). C-ohjelmia käyttäen tulos voidaan muotoilla seuraavasti:

Väite. Ei ole olemassa C-funktiota `halt(p,x)`, joka saa syöteenään mielivaltaisen C-funktion tekstin `p` ja tälle tarkoitettun syöteen `x` ja tuottaa tuloksen 1, jos `p`:n suoritus pysähtyy syötteellä `x`, ja 0 jos `p`:n suoritus `x`:llä jää ikuiseen silmukkaan.

Todistus. Oletetaan väitteen vastaisesti, että tällainen funktio `halt` voitaisiin laatia. Muodostetaan tästä käyttäen toinen funktio `confuse`:

```
void confuse(char *p){  
    int halt(char *p, char *x){  
        ... /* Funktion halt runko. */  
    }  
    if (halt(p,p) == 1) while (1);  
}
```

Merkitään edellä kuvattua funktion `confuse` ohjelmatekstiä c:llä ja tarkastellaan funktion `confuse` laskentaa tällä omalla kuvauksellaan. Saadaan ristiriita:

$$\begin{aligned} \text{confuse}(c) \text{ pysähtyy} &\Leftrightarrow \text{halt}(c,c) == 1 \\ &\Leftrightarrow \text{confuse}(c) \text{ ei pysähdy}. \end{aligned}$$

Ristiriidasta seuraa, että oletettua pysähtymis-testausfunktiota `halt` ei voi olla olemassa. □

Samansukuisia ns. *ratkeamattomia ongelmia* on itse asiassa *paljon*. Asiaan palataan kurssin loppupuolella.