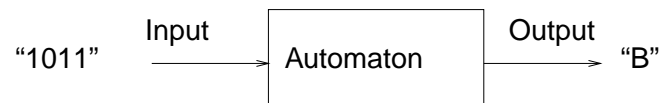


1.6 Aakkostot, merkkijonot ja kielet

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:

1. syötteet ovat äärellisiä, diskreettejä *merkkijonoja*
2. tulokset ovat muotoa "hyväksy"/"hylkää" (~ "syöte OK"/"syöte ei kelpaa")

Yleistyksiä:

- ▶ äärettömät syötejonot (→ "reaktiiviset" järjestelmät, Büchi-automaatit)
- ▶ funktioautomaatit (→ Moore- ja Mealy-tilakoneet, Turingin funktiokoneet)



Peruskäsitteitä ja merkintöjä

Aakkosto (engl. alphabet, vocabulary): äärellinen, epätyhjä joukko *alkeismerkkejä* t. *symboleita*. 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.
- ▶ *tyhjä merkkijono* (engl. empty string). Tyhjässä merkkijonossa ei ole yhtään merkkiä; sitä voidaan merkitä ε :llä.

Merkkijonon x *pituus* $|x|$ on sen merkkien määrä.

Esim.: $|01001| = 5$, $|TKTP| = 4$, $|\varepsilon| = 0$.



Merkkijonojen välinen perusoperaatio on *katenaatio* eli jonojen peräkkäin kirjoittaminen. Katenaation operaatiomerkinä käytetään joskus selkeyden lisäämiseksi symbolia \wedge . Esim.

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



Aakkoston Σ kaikkien merkkijonojen joukkoa 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).



Vakiintuneita merkintöjä

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

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

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

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).



Automaatit ja formaalit kielet

Olkoon M automaatti, jonka syötteet ovat jonkin 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. recognises) 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 kieltä

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

saadaan vihjeitä siitä, millainen automaatti tarvitaan kuvauksen f toteuttamiseen.



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

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

- ▶ $|abc| = 3$;
- ▶ 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:*

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

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

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

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:

- $\varepsilon^R = \varepsilon$;
- jos $w = ua$, $u \in \Sigma^*$, $a \in \Sigma$, niin $w^R = a\hat{u}^R$.



Induktiivista (“rekursiivista”) määritelmää voidaan tietenkin käyttää laskujen perustana; esim.:

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

Tärkeämpää on kuitenkin konstruktoiden ominaisuuksien todistaminen määritelmää noudattelevalla induktiolla. Esimerkki:



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

Todistus. Induktio merkkijonon y pituuden suhteen.

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

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



1.7 Numeroituvat ja ylinumeroituvat joukot

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

Intuitiivisesti sanoen joukko X on numeroituva, 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.

Numeroituvan joukon kaikki osajoukot ovat myös numeroituvia (HT), mutta 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:

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



Bijektio f on siis:

0	\mapsto	ε	
1	\mapsto	a_1	$2n + 1$
2	\mapsto	a_2	\vdots
\vdots	\vdots		$3n$
n	\mapsto	a_n	\vdots
$n + 1$	\mapsto	$a_1 a_1$	$n^2 + n$
$n + 2$	\mapsto	$a_1 a_2$	$n^2 + n + 1$
\vdots	\vdots		$n^2 + n + 2$
$2n$	\mapsto	$a_1 a_n$	\vdots
			\vdots



Itse asiassa millä tahansa ohjelmointikielellä kirjoitetut ohjelmat ovat kielen perusaakkoston (esim. C-kielessä ASCII-merkistön) merkkijonoja. Lauseen 1.11 mukaan minkä tahansa aakkoston merkkijonojen joukko on numeroituvasti ääretön, joten myös millä tahansa ohjelmointikielellä mahdollisten ohjelmien joukko on numeroituva.

Seuraavaksi todistetaan, että kaikkien formaalien kielten joukko on ylinumeroituva. Formaaleja kieliä on siis "enemmän" kuin mahdollisia tietokoneohjelmia, ja siksi *millään ohjelmointikielellä ei voida laatia tunnistusautomaatteja kaikille formaaleille kielille*. (Tai toisin sanoen: on olemassa "periaatteessa mahdollisia" I/O-kuvauksia, joita ei voida toteuttaa tietokoneella.)



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

Todistus (ns. Cantorin diagonaaliargumentti). Merkitään aakkoston Σ kaikkien formaalien kielten perhettä $\mathcal{P}(\Sigma^*) = \mathcal{A}$. Tehdään vastaoletus: oletetaan, että olisi olemassa kaikki Σ :n formaalit kielet kattava numerointi:

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

Olkoot Σ^* :n merkkijonot kanonisessa järjestyksessä 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 $\tilde{\mathcal{A}}$:n määritelmän mukaan

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

Saatu ristiriita osoittaa, että vastaoletus on väärä.



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

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



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 syötteenään mielivaltaisen C-funktion tekstin `p` ja tälle tarkoitetun syötteen `x` ja tuottaa tuloksen 1, jos `p`:n suoritus pysähtyy syötteellä `x`, ja 0 jos `p`:n suoritus `x`:llä jää ikuisen silmukkaan.



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

Merkitään funktion `confuse` ohjelmatekstiä `c`:llä ja tarkastellaan funktion `confuse` laskentaa tällä omalla kuvauksellaan.

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

Saadaan ristiriita:
`confuse(c)` pysähtyy \Leftrightarrow
`halt(c,c) == 1
 \Leftrightarrow
confuse(c) ei pysähdy.`

Ristiriidasta seuraa, että oletettua pysähtymistestausfunktiota `halt` ei voi olla olemassa. \square

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

