

Automaten und Formale Sprachen

5. Vorlesung

Prof. Dr. Dietrich Kuske

FG Automaten und Logik, TU Ilmenau

Wintersemester 2023/24

Reguläre Ausdrücke

reguläre Ausdrücke erlauben es, die regulären Sprachen kompakt in „Textform“ zu beschreiben

Definition

Die Menge $\text{Reg}(\Sigma)$ der **regulären Ausdrücke über dem Alphabet Σ** ist die kleinste Menge mit folgenden Eigenschaften:

- $\emptyset \in \text{Reg}(\Sigma)$, $\lambda \in \text{Reg}(\Sigma)$, $\Sigma \subseteq \text{Reg}(\Sigma)$.
- Wenn $\alpha, \beta \in \text{Reg}(\Sigma)$, dann auch $(\alpha \cdot \beta)$, $(\alpha + \beta)$, $(\alpha^*) \in \text{Reg}(\Sigma)$.

Bemerkungen:

- für $(\alpha \cdot \beta)$ schreibt man oft $(\alpha \beta)$
- für $(\alpha + \beta)$ schreibt man auch $(\alpha \mid \beta)$

Beispiel

 $\Sigma = \{a, b, c, d\}:$

$$\lambda \quad ((ab)b) \quad ((a+d)a)$$

$$\left(\left(((ab)a) \right) + \left((((ba)b))^* \right) \right) \quad \left(\left((ab)((a+b)^*) \right) (ba) \right)$$

$$\Sigma = \{0, 1\}: \left(1((0+1)^*) \right) \quad \left(((0^*)(1^*)) (0^*) \right)$$

Beobachtung: Jeder reguläre Ausdruck über Σ ist ein Wort über dem Alphabet $\Sigma \cup \{\emptyset, \lambda, \cdot, +, *, (,)\}$.

Nach der Festlegung der Syntax regulärer Ausdrücke müssen wir auch deren Bedeutung (= Semantik) festlegen.

Dabei betrachten wir einen regulären Ausdruck α als ein „Wortschema“, $L(\alpha)$ soll die Menge der Wörter sein, die zum „Wortschema“ „passen“.

Definition

Für einen regulären Ausdruck $\alpha \in \text{Reg}(\Sigma)$ ist die Sprache $L(\alpha) \subseteq \Sigma^*$ induktiv definiert:

$$L(\alpha) = \begin{cases} \emptyset & \text{falls } \alpha = \emptyset \\ \{\varepsilon\} & \text{falls } \alpha = \lambda \\ \{a\} & \text{falls } \alpha = a \in \Sigma \\ L(\beta) \cup L(\gamma) & \text{falls } \alpha = (\beta + \gamma) \\ L(\beta)L(\gamma) & \text{falls } \alpha = (\beta \cdot \gamma) \\ (L(\beta))^* & \text{falls } \alpha = (\beta^*) \end{cases}$$

Beispiel (vgl. Folie 5.3)

 $\Sigma = \{a, b, c, d\}$:

$$L(\lambda) = \{\varepsilon\}$$

$$L\left(\left((ab)b\right)\right) = \{abb\}$$

$$L\left(\left((a+d)a\right)\right) = \{aa, da\}$$

$$L\left(\left(\left(\left((ab)a\right) + \left(\left((ba)b\right)^*\right)\right)\right)\right) = \{aba\} \cup \{bab\}^*$$

$$L\left(\left(\left(\left(ab\right)\left(\left(a+b\right)^*\right)\right)\left(ba\right)\right)\right) = \{ab\} \cdot \{a, b\}^* \cdot \{ba\}$$

 $\Sigma = \{0, 1\}$:

$$L\left(\left(\left(1\left(\left(0+1\right)^*\right)\right)\right)\right) = \{1\} \cdot \{0, 1\}^*$$

$$L\left(\left(\left(\left(0^*\right)\left(1^*\right)\right)\left(0^*\right)\right)\right) = \{0^i 1^j 0^k \mid i, j, k \in \mathbb{N}\}$$

Klammern sparen:

- Äußere Klammern weglassen: $(0 + 1)^*$ statt $((0 + 1)^*)$.
- Da $(L_1 L_2) L_3 = L_1 (L_2 L_3)$ und $(L_1 \cup L_2) \cup L_3 = L_1 \cup (L_2 \cup L_3)$, läßt man Klammern bei Operatoren auf gleicher Ebene weg:
 $(0 + 1 + 2 + 3)$ statt $((0 + 1) + (2 + 3))$ oder $(0 + (1 + (2 + 3)))$
- **Präferenzregeln:**
 - * bindet stärker als \cdot .
 - \cdot bindet stärker als $+$.

Damit erhalten die Beispiel-Ausdrücke die lesbarere Form:

$$\begin{array}{ccc}
 \lambda & abb & (a + d)a \\
 aba + (bab)^* & ab(a + b)^* ba & \\
 1(0 + 1)^* & 0^* 1^* 0^* &
 \end{array}$$

Beispiel

$\Sigma = \{a, b\}$.

- $L(a(a + b)^*bb)$ ist die Menge der Wörter, die ???
- $L((a + b)^*aba(a + b)^*)$ ist die Menge der Wörter, die ???
- $L((b^*ab^*a)^*b^*)$ ist die Menge der Wörter, die ???

Wo treten reguläre Ausdrücke in der Praxis auf?

- **Suchen und Ersetzen** in Editoren
(z.B. vi, emacs, ...)
- **Pattern-Matching** und Verarbeitung großer Texte und Datenmengen,
z.B. beim Data-Mining
(z.B. sed, awk, ...)
- **Übersetzung** von Programmiersprachen:
Lexikalische Analyse – Umwandlung einer Folge von Zeichen (das Programm) in eine Folge von Tokens, in der bereits die Schlüsselwörter, Bezeichner, Daten, etc. identifiziert sind.
(z.B. lex, flex, ...)

Diese Anwendungen beruhen auf dem Zusammenhang zu endlichen Automaten, den wir jetzt untersuchen werden.

Proposition

Ist γ ein regulärer Ausdruck, so ist $L(\gamma)$ eine reguläre Sprache.

Beweis: per Induktion über den Aufbau von γ .

IA Da die Sprachen \emptyset , $\{\varepsilon\}$ und $\{a\}$ für $a \in \Sigma$ regulär sind, gilt die Aussage für die regulären Ausdrücke $\gamma \in \{\emptyset, \lambda\} \cup \Sigma$.

IS Die Klasse der regulären Sprachen ist abgeschlossen unter Vereinigung, Konkatenation und Kleene-Iteration (Sätze auf Folien 4.14, 4.16 und 4.21). □

Proposition

Zu jedem DFA M gibt es einen regulären Ausdruck γ mit $L(M) = L(\gamma)$.

Beweis:

Sei $M = (\{1, \dots, n\}, \Sigma, 1, \delta, E)$ ein DFA.

Wir konstruieren einen regulären Ausdruck γ mit $L(M) = L(\gamma)$.

Für ein Wort $w \in \Sigma^*$ sei

$$\text{Pref}(w) = \{u \in \Sigma^* \mid \exists v : w = uv, \varepsilon \neq u \neq w\}$$

die Menge aller nicht-leeren echten Präfixe von w .

Für $i, j \in \{1, \dots, n\}$ und $k \in \{0, \dots, n\}$ sei

$$L_{i,j}^k = \{w \in \Sigma^* \mid \widehat{\delta}(i, w) = j, \forall u \in \text{Pref}(w): 1 \leq \widehat{\delta}(i, u) \leq k\}.$$

Intuitiv: Ein Wort w gehört zu $L_{i,j}^k$ genau dann, wenn w den Zustand i in den Zustand j überführt und wenn dabei kein Zwischenzustand (außer vielleicht ganz am Anfang und ganz am Ende) $> k$ vorkommt.

Behauptung: Für alle $i, j \in \{1, \dots, n\}$ und $k \in \{0, \dots, n\}$ existieren reguläre Ausdrücke $\gamma_{i,j}^k$ mit $L(\gamma_{i,j}^k) = L_{i,j}^k$.

Bemerkung: Falls $E = \{i_1, i_2, \dots, i_\ell\}$, ergibt sich dann

$$L(\gamma_{1,i_1}^n + \gamma_{1,i_2}^n + \dots + \gamma_{1,i_\ell}^n) = L(M),$$

womit der Beweis abgeschlossen sein wird.

Beweis der Behauptung, d.h. Konstruktion von $\gamma_{i,j}^k$ durch Induktion über $k \in \{0, \dots, n\}$.

IA $k = 0$. Es gilt:

$$L_{i,j}^0 = \begin{cases} \{a \in \Sigma \mid \delta(i, a) = j\} & \text{falls } i \neq j \\ \{a \in \Sigma \mid \delta(i, a) = j\} \cup \{\varepsilon\} & \text{falls } i = j \end{cases}$$

Gilt $\{a_1, \dots, a_m\} = \{a \in \Sigma \mid \delta(i, a) = j\}$, so setzen wir:

$$\gamma_{i,j}^0 = \begin{cases} a_1 + a_2 + \dots + a_m + \emptyset & \text{falls } i \neq j \\ a_1 + a_2 + \dots + a_m + \lambda & \text{falls } i = j \end{cases}$$

In jedem Fall gilt dann $L(\gamma_{i,j}^0) = L_{i,j}^0$.

IV Sei $0 \leq k < n$ und seien die regulären Ausdrücke $\gamma_{p,q}^k$ für alle $p, q \in \{1, \dots, n\}$ bereits konstruiert.

IS Sei $i, j \in \{1, \dots, n\}$.

Dann gilt:

$$L_{i,j}^{k+1} = L_{i,j}^k \cup L_{i,k+1}^k (L_{k+1,k+1}^k)^* L_{k+1,j}^k$$

Begründung:

„ \subseteq “: Sei $w \in L_{i,j}^{k+1}$ und sei $\ell \geq 0$ so, daß der Zustand $k + 1$ auf dem eindeutigen mit w beschrifteten Pfad von i nach j genau ℓ mal als echter Zwischenzustand auftaucht.

1.Fall: $\ell = 0$, d.h. $k + 1$ kommt nicht als echter Zwischenzustand vor.

$$\leadsto w \in L_{i,j}^k$$

2.Fall: $\ell > 0$.

$\leadsto w = w_0 w_1 \cdots w_{\ell-1} w_\ell$, wobei:

$$\begin{aligned}\widehat{\delta}(i, w_0) &= k + 1 \\ \widehat{\delta}(k + 1, w_p) &= k + 1 \text{ für } 1 \leq p \leq \ell - 1 \\ \widehat{\delta}(k + 1, w_\ell) &= j\end{aligned}$$

$$\leadsto w_0 \in L_{i,k+1}^k, w_p \in L_{k+1,k+1}^k \text{ (} 1 \leq p \leq \ell - 1 \text{)}, w_\ell \in L_{k+1,j}^k$$

$$\leadsto w = w_0(w_1 \cdots w_{\ell-1})w_\ell \in L_{i,k+1}^k (L_{k+1,k+1}^k)^* L_{k+1,j}^k$$

Damit ist die Inklusion „ \subseteq “ der Behauptung von Folie 5.13 gezeigt.

„ \supseteq “:

Für alle $p, q \in \{1, \dots, n\}$ gilt $L_{p,q}^k \subseteq L_{p,q}^{k+1}$.

Damit erhalten wir

$$\begin{aligned} & L_{i,j}^k \cup L_{i,k+1}^k (L_{k+1,k+1}^k)^* L_{k+1,j}^k \\ \subseteq & L_{i,j}^{k+1} \cup L_{i,k+1}^{k+1} (L_{k+1,k+1}^{k+1})^* L_{k+1,j}^{k+1} \\ \subseteq & L_{i,j}^{k+1} \end{aligned}$$

Damit ist auch die Inklusion „ \supseteq “ der Behauptung von Folie 5.13 gezeigt.

Es gilt also wirklich

$$L_{i,j}^{k+1} = L_{i,j}^k \cup L_{i,k+1}^k (L_{k+1,k+1}^k)^* L_{k+1,j}^k.$$

Außerdem haben wir bereits reguläre Ausdrücke $\gamma_{p,q}^k$ mit $L(\gamma_{p,q}^k) = L_{p,q}^k$ konstruiert.

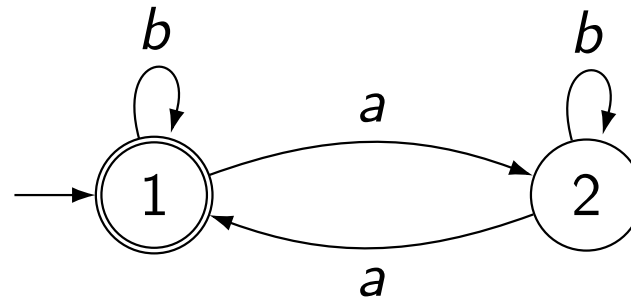
Mit

$$\gamma_{i,j}^{k+1} = \gamma_{i,j}^k + \gamma_{i,k+1}^k (\gamma_{k+1,k+1}^k)^* \gamma_{k+1,j}^k$$

erhalten wir $L(\gamma_{i,j}^{k+1}) = L_{i,j}^{k+1}$.



Beispiel: Betrachte den folgenden DFA:



Damit ergibt sich (bei Durchführung offensichtlicher Vereinfachungen):

$$\gamma_{1,1}^0 = \lambda + b \quad \gamma_{1,2}^0 = a \quad \gamma_{2,1}^0 = a \quad \gamma_{2,2}^0 = \lambda + b$$

$$\gamma_{1,1}^1 = \gamma_{1,1}^0 + \gamma_{1,1}^0 (\gamma_{1,1}^0)^* \gamma_{1,1}^0 = \lambda + b + (\lambda + b)(\lambda + b)^* (\lambda + b) \hat{=} b^*$$

$$\gamma_{1,2}^1 = \gamma_{1,2}^0 + \gamma_{1,1}^0 (\gamma_{1,1}^0)^* \gamma_{1,2}^0 = a + (\lambda + b)(\lambda + b)^* a \hat{=} b^* a$$

$$\gamma_{2,1}^1 = \gamma_{2,1}^0 + \gamma_{2,1}^0 (\gamma_{1,1}^0)^* \gamma_{1,1}^0 = a + a(\lambda + b)^* (\lambda + b) \hat{=} ab^*$$

$$\gamma_{2,2}^1 = \gamma_{2,2}^0 + \gamma_{2,1}^0 (\gamma_{1,1}^0)^* \gamma_{1,2}^0 = \lambda + b + a(\lambda + b)^* a \hat{=} \lambda + b + ab^* a$$

$$\gamma_{1,1}^2 = \gamma_{1,1}^1 + \gamma_{1,2}^1 (\gamma_{2,2}^1)^* \gamma_{2,1}^1 = b^* + b^* a (\lambda + b + ab^* a)^* ab^*$$

Zusammenfassender Satz

Sei $L \subseteq \Sigma^*$ eine Sprache. Dann sind äquivalent

- ① L ist regulär, d.h. es gibt einen DFA M mit $L(M) = L$.
- ② Es gibt einen NFA M mit $L(M) = L$.
- ③ L ist rechtslinear, d.h. es gibt eine rechtslineare Grammatik G mit $L(G) = L$.
- ④ Es gibt einen regulären Ausdruck γ mit $L(\gamma) = L$.
- ⑤ L ist erste Komponente einer Lösung eines linearen Gleichungssystems.
- ⑥ L ist in monadischer Logik 2. Stufe definierbar.
- ⑦ L ist abgeschlossen unter einer Kongruenz endlichen Indexes.
- ⑧ Es gibt Homomorphismus $\eta: \Sigma^* \rightarrow S$ in ein endliches Monoid S mit $L = \eta^{-1}\eta(L)$.
- ⑨ L wird von einem alternierenden Automaten akzeptiert.
- ⑩ L wird von einem 2-Weg-Automaten akzeptiert.

Weitere Forschungsfragen

- 1 DFAs mit speziellen Eigenschaften (z.B. „zählerfrei“): Welche Sprachen akzeptieren diese? Kann ich feststellen, ob ein allgemeiner DFA äquivalent zu einem mit spezieller Eigenschaft ist?
- 2 Größenvergleich von äquivalenten Mechanismen (z.B. „DFAs sind notwendigerweise exponentiell größer als NFAs“ oder „Übersetzung von regulären Ausdrücken in NFA ist in polynomieller Zeit möglich“)

Zusammenfassung

verschiedene Modelle zur Beschreibung aller regulären Sprachen:

- **Rechtslineare Grammatiken:**
 - Verbindung zur Chomsky-Hierarchie
 - erzeugen Sprachen
 - nicht geeignet, um zu entscheiden, ob ein geg. Wort zur Sprache gehört
- **NFAs:**
 - erlauben kleine, kompakte Darstellung
 - intuitive graphische Notation
 - nicht geeignet, um zu entscheiden, ob ein geg. Wort zur Sprache gehört
- **DFAs:**
 - für effiziente Beantwortung der Frage, ob ein Wort zur Sprache gehört
 - sind u.U. exponentiell größer als NFA
- **Reguläre Ausdrücke:**
 - erlauben kompakte Darstellung in Textform

Zusammenfassung 5. Vorlesung

in dieser Vorlesung neu

- Abschluß der Klasse der rechtslinearen Sprachen unter positiver Iteration und unter Kleene-Iteration
- reguläre Ausdrücke beschreiben genau die rechtslinearen Sprachen, sind also gleich ausdrucksstark wie Typ-3-Grammatiken, DFAs und NFAs

kommende Vorlesung

- Gibt es Sprachen, die nicht rechtslinear sind?
- Wie kann ich von einer geg. Sprachen zeigen, daß sie nicht rechtslinear ist?