



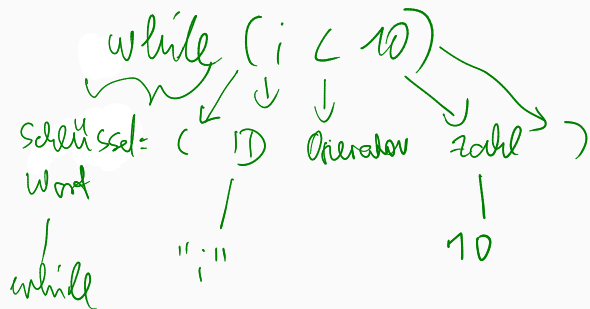
Reguläre Sprachen, Ausdrucksstärke

BC George (FH Bielefeld)

Unless otherwise noted, this work is licensed under CC BY-SA 4.0.

Motivation

Was muss ein Compiler wohl als erstes tun?



`while i < 10:`

\Leftarrow **Symbol, Token**

1. Phase: Zerlegung
der Eingabe:

Lexer

Scanner

Tokenizer

Tokenerkennung:

regex, DFA,

Themen für heute

- Endliche Automaten
- Reguläre Ausdrücke

Endliche Automaten

Alphabete

Sigma

Def.: Ein *Alphabet* Σ ist eine endliche, nicht-leere Menge von Symbolen. Die Symbole eines Alphabets heißen *Buchstaben*.

Epsilon

Def.: Ein *Wort* w über einem *Alphabet* Σ ist eine endliche Folge von Symbolen aus Σ . ϵ ist das leere Wort. Die *Länge* $|w|$ eines Wortes w ist die Anzahl von Buchstaben, die es enthält (Kardinalität).

Def.: $\Sigma^k = \{w \text{ über } \Sigma \mid |w| = k\}$

$$\Sigma^0 = \{\epsilon\}$$

$\Sigma^* = \bigcup_{i \in \mathbb{N}_0} \Sigma^i$ (die Kleene-Hülle von Σ)

$\Sigma^+ = \bigcup_{i \in \mathbb{N}} \Sigma^i$

Sprachen über Alphabete


$$x \cdot y = xy = x \circ y$$

Def.: Seien $x = a_1 a_2 \dots a_n$ und $y = b_1 b_2 \dots b_m$ Wörter. Wir nennen $xy = x \circ y = a_1 \dots a_n b_1 \dots b_m$ die Konkatenation von x und y .

Def.: Eine Sprache L über einem Alphabet Σ ist eine Teilmenge von Σ^* : $L \subseteq \Sigma^*$

\subseteq ; L kann auch
 $= \Sigma^*$ sein

Deterministische endliche Automaten



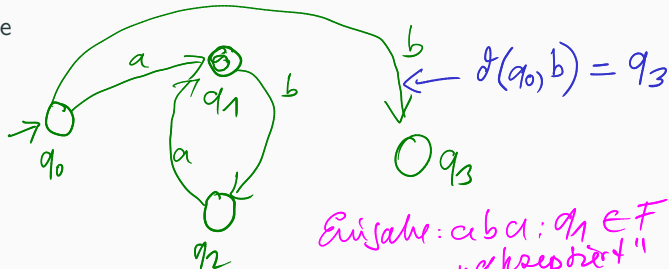
State Machine

Def.: Ein *deterministischer endlicher Automat* (DFA) ist ein 5-Tupel $A = (Q, \Sigma, \delta, q_0, F)$ mit

- Q : eine endliche Menge von Zuständen
- Σ : ein Alphabet von Eingabesymbolen
- δ : die Übergangsfunktion $(Q \times \Sigma) \rightarrow Q$, δ kann partiell sein
- $q_0 \in Q$: der Startzustand
- $F \subseteq Q$: die Menge der Endzustände

z. B. $Q = \{q_1, q_2, q_3, q_0\}$
 $\Sigma = \{a, b\}$

q	a	b
$\rightarrow q_0$	q_1	q_3
$* q_1$	—	q_2
q_2	q_1	—
q_3	—	—



Eingabe: $ab a$; $q_1 \in F$
 "akzeptiert"

Die Übergangsfunktion

alle Wörter in Σ^* sind von
endlicher Länge

Def.: Wir definieren $\delta^* : (Q \times \Sigma^*) \rightarrow Q$: induktiv wie folgt:

- Basis: $\delta^*(q, \epsilon) = q \quad \forall q \in Q$
- Induktion: $\delta^*(q, a_1, \dots, a_n) = \delta(\delta^*(q, a_1, \dots, a_{n-1}), a_n)$

Def.: Ein DFA akzeptiert ein Wort $w \in \Sigma^*$ genau dann, wenn $\delta^*(q_0, w) \in F$.

Def.: Die Sprache eines DFA A $L(A)$ ist definiert durch:

$$L(A) = \{w \mid \delta^*(q_0, w) \in F\}$$

Beispiel

Bsp: durch 3 Hilbbare Dualzahlen

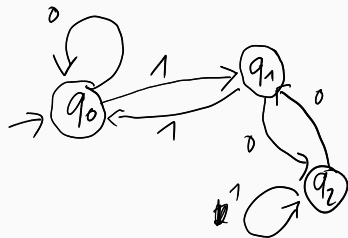
Zustände unterscheiden Partitionen

q_0 : die bisher gelesene Zahl ist mit Rest 0 durch 3 teilbar

q_1 : " " " " " 1

q_2 : " " " " " 2

von links nach rechts lesen z.B. $\overset{4}{1}\overset{2}{0}\overset{1}{1}\overset{1}{0}\overset{1}{1}$ 0: * 2
 $\rightarrow \uparrow$ 1: * 2 + 1



Reguläre Ausdrücke

Def.: Seien L und M Sprachen.

- $L \cup M := \{w \mid w \in L \vee w \in M\}$
- $LM = L \cdot M = L \circ M = \{vw \mid v \in L \wedge w \in M\}$
- Die Kleene-Hülle einer Sprache:
 - Basis: $L^0 = \{\epsilon\}$
 - Induktion: $L^i = \{xw \mid x \in L^{i-1}, w \in L, i > 0\}$,
$$L^* = \bigcup_{i \geq 0} L^i,$$
$$L^+ = \bigcup_{i > 0} L^i$$

Reguläre Ausdrücke

Def.: Induktive Definition von regulären Ausdrücken (*regex*) und der von ihnen repräsentierten Sprache:

▪ Basis:

- ϵ und \emptyset sind reguläre Ausdrücke mit $L(\epsilon) = \{\epsilon\}$, $L(\emptyset) = \emptyset$
- Sei a ein Symbol $\Rightarrow a$ ist ein regex mit $L(a) = \{a\}$

▪ Induktion: Seien E, F reguläre Ausdrücke. Dann gilt:

- $E | F =$
- $E + F$ ist ein regex und bezeichnet die Vereinigung $L(E + F) = L(E) \cup L(F)$
 - EF ist ein regex und bezeichnet die Konkatenation $L(EF) = L(E)L(F)$ *Konkatenation*
 - E^* ist ein regex und bezeichnet die Kleene-Hülle $L(E^*) = (L(E))^*$
 - (E) ist ein regex mit $L((E)) = L(E)$

Vorrangregeln der Operatoren für reguläre Ausdrücke: $*$, Konkatenation, $+$

$$ab^* \neq (ab)^*$$

$$a^4 b^4 \neq \text{Regex}$$
$$a^5 = aaaaa$$

Wichtige Identitäten

Satz: Sei A ein DFA $\Rightarrow \exists$ regex R mit $L(A) = L(R)$.

Satz: Sei E ein regex $\Rightarrow \exists$ DFA A mit $L(E) = L(A)$.

Formale Grammatiken

DFA: Sprachen erkennen

Regex: " " beschreiben

Gram: " " erzeugen

* : \emptyset ~~oder~~ lokal oder endlich oft

Def.: Eine *formale Grammatik* ist ein 4-Tupel $G = (N, T, P, S)$ aus

- N : einer endlichen Menge von *Nichtterminalen* (Variable)
- T : einer endlichen Menge von *Terminalen*, $N \cap T = \emptyset$ zeichnen in den generierten Wörtern
- $S \in N$: dem *Startsymbol*
- P : einer endlichen Menge von Produktionen der Form: $X \rightarrow Y$ mit $X \in (N \cup T)^* N (N \cup T)^*$, $Y \in (N \cup T)^*$

$\in N$ \swarrow $S \rightarrow a A b \mid B$
 \swarrow $A \rightarrow a \mid \epsilon$
 \searrow $B \rightarrow \dots$
 $\epsilon \in T$

kleinbuchstaben: Terminal
großbuchstaben: Nichtterminal

Ableitungen

$$\gamma, \alpha, \beta \in (N \cup T)^* \\ A \in N$$

Def.: Sei $G = (N, T, P, S)$ eine Grammatik, sei $\alpha A \beta$ eine Zeichenkette über $(N \cup T)^*$ und sei $A \rightarrow \gamma$ eine Produktion von G .

Wir sagen: $\alpha A \beta \Rightarrow \alpha \gamma \beta$ ($\alpha A \beta$ leitet $\alpha \gamma \beta$ ab).

Def.: Wir definieren die Relation \Rightarrow^* induktiv wie folgt:

→ endlich viele Schritte

- Basis: $\forall \alpha \in (N \cup T)^* \alpha \Rightarrow^* \alpha$ (Jede Zeichenkette leitet sich selbst ab.)
- Induktion: Wenn $\alpha \Rightarrow^* \beta$ und $\beta \Rightarrow \gamma$ dann $\alpha \Rightarrow^* \gamma$

Def.: {Sei $G = (N, T, P, S)$ eine formale Grammatik. Dann ist $L(G) = \{w \in T^* \mid S \Rightarrow^* w\}$ die von G erzeugte Sprache.

1. Schritt: $S \Rightarrow aAb$

2. Schritt: $S \Rightarrow aab$

$S \Rightarrow^* aab$

Produktionen

$S \rightarrow aAb \mid B \quad A \rightarrow a \quad B \rightarrow bA$

$S \rightarrow aAb$ Produktion

$S \Rightarrow aAb$, prod. $A \rightarrow a$ anwenden

$S \Rightarrow^* aab$

Def.: Eine *reguläre (oder type-3-) Grammatik* ist eine formale Grammatik mit den folgenden Einschränkungen:

- Alle Produktionen sind entweder von der Form
 - $X \rightarrow aY$ mit $X \in N, a \in T, Y \in N$ (*rechtsreguläre Grammatik*) oder
 - $X \rightarrow Ya$ mit $X \in N, a \in T, Y \in N$ (*linksreguläre Grammatik*)
- $X \rightarrow \epsilon$ ist in beiden Fällen erlaubt.

$$A \rightarrow aB$$



Satz: Die von rechtsregulären Grammatiken erzeugten Sprachen sind genau die von linksregulären Grammatiken erzeugten Sprachen. Beide werden *reguläre* Sprachen genannt.

Satz: Die von regulären Ausdrücken beschriebenen Sprachen sind die regulären Sprachen.

Die Klasse der regulären Sprachen ist abgeschlossen unter

- Vereinigung
- Konkatenation
- Kleene-Stern
- Komplementbildung
- Durchschnitt

Entscheidbarkeit für reguläre Sprachen

Satz: Es ist entscheidbar,

- ob eine gegebene reguläre Sprache leer ist
- ob $w \in \Sigma^*$ in einer gegebenen regulären Sprache enthalten ist (Das “Wort-Problem”)
- ob zwei reguläre Sprachen äquivalent sind

Grenzen der regulären Sprachen

$a^n b c^n \notin \text{regex}$ z.B. $\{ \}$ in Programmen

Reguläre Sprachen sind von ihrer Struktur her einfach. Schon Sprachen, in denen etwas “gematcht” werden muss, lassen sich nicht mehr regulär beschreiben, weil z. B. die fixe Anzahl von Zuständen eines DFAs die Erkennung solcher Sprachen verhindert.

Wozu das Ganze?

Im Compilerbau werden reguläre Ausdrücke benutzt, um die Schlüsselwörter und weitere Symbole der zu erkennenden Sprache anzugeben. Daraus wird mit Hilfe eines Generators, der aus den regulären Ausdrücken DFAs (oder einen großen DFA) macht, der sog. Scanner oder Lexer genannt, generiert. Seine Aufgabe ist es, die Folge von Zeichen in der Quelldatei in eine Folge von sog. Token umzuwandeln. Z. B. wird so aus den Zeichen des Schlüsselwortes *while* im Programmtext das Token für *while* gemacht, das in der Syntaxanalyse weiterverarbeitet wird. Die Tokenfolge eines Programms ist ein Wort einer Sprache, die der Parser erkennt. Jedes vom Lexer erkannte Token ist dort also ein terminales Symbol.

Ein Lexer ist mehr als ein DFA

Was ist zu beachten:

- Man braucht mindestens eine Liste von Paaren aus regulären Ausdrücken und Tokennamen.
- Neben den Schlüsselwörtern und Symbolen wie (,), *, ... müssen auch Namen für Variablen, Funktionen, Klassen, Methoden, ... (sog. Identifier) erkannt werden
- Namen haben meist eine gewisse Struktur, die sich mit regulären Ausdrücken beschreiben lassen.
- Erlaubte Token sind in der Grammatik des Parsers beschrieben, d. h. für literale Namen, Strings, Zahlen liefert der Scanner zwei Werte:
 - z. B. <ID, "radius">, <Integerzahl, 558>
- Kommentare und Strings müssen richtig erkannt werden. (Schachtelungen)

Man kann natürlich auch einen Lexer selbst programmieren, d. h. die DFAs für die regulären Ausdrücke implementieren.

Automatisch oder händisch

- alles ausprogrammieren
- einen DFA ausprogrammieren
- einen Scannergenerator benutzen \Leftarrow

Wrap-Up

- Definition und Aufgaben von Lexern
- DFAs und NFAs
- Reguläre Ausdrücke
- Reguläre Grammatiken
- Zusammenhänge zwischen diesen Mechanismen und Lexern, bzw. Lexergeneratoren



Unless otherwise noted, this work is licensed under CC BY-SA 4.0.