

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?

Token in Eingabe finden

⇒ Eingabestrom in Schlüsselwörter,
Namens Operatoren, ...
einteilen.

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*.

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

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

transitiver Abschluss
Hülle

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

Sprachen über Alphabete

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

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^* : $L \subseteq \Sigma^*$

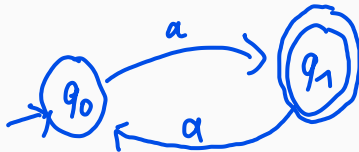


Das leere Wort
 ϵ enthalten

Deterministische endliche Automaten

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

- Q : eine endliche Menge von Zuständen $\{q_0, q_1\}$
- Σ : ein Alphabet von Eingabesymbolen $\{a\}$
- δ : die Übergangsfunktion $(Q \times \Sigma) \rightarrow Q$, δ kann partiell sein *Pfeile*
- ▪ $q_0 \in Q$: der Startzustand $\delta(q_0, a) = q_1$
- ▪ $F \subseteq Q$: die Menge der Endzustände



Die Übergangsfunktion

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

- Basis: $\delta^*(q, \epsilon) = q \ \forall q \in Q$
- 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

Durch 3 Hilbare Dualzahlen
Zustände entsprechen Restklassen:

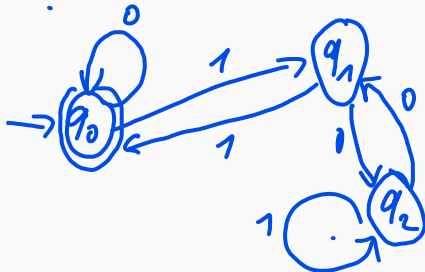
q_0 : die bisher gelesene Zahl einen Rest 0 (bei Div d. 3)

q_1 : "

q_2 : "

1 "

2 "



Nichtdeterministische endliche Automaten

Def.: Ein *nichtdeterministischer endlicher Automat* (NFA) 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 $(Q \times \Sigma) \rightarrow \underline{\underline{\mathcal{P}(Q)}}$
- $q_0 \in Q$: der Startzustand
- $F \subseteq Q$: die Menge der Endzustände

Potenzmenge von Q
= Menge der Teilmengen
von Q

Die Übergangsfunktion eines NFAs

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

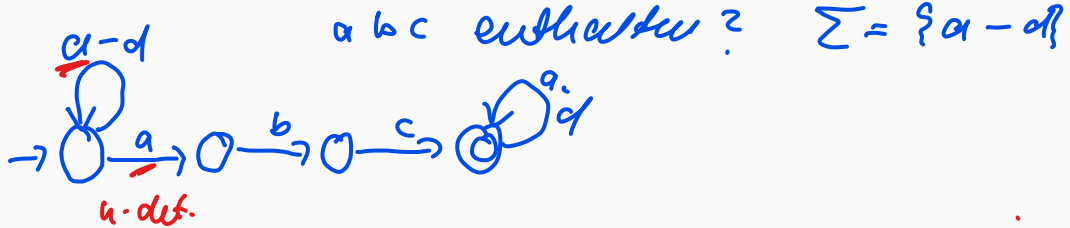
- Basis: $\delta^*(q, \epsilon) = q \ \forall q \in Q$
- Induktion: Sei $w \in \Sigma^*$, $w = xa$, $x \in \Sigma^*$, $a \in \Sigma$ mit

$\delta^*(q, x) = \{p_1, \dots, p_k\}$, $p_i \in Q$, sei

$$A = \bigcup_{i=1}^k \delta(p_i, a) = \{r_1, \dots, r_m\}, r_j \in Q.$$

Dann ist $\delta^*(q, w) = \{r_1, \dots, r_m\}$.

Wozu NFAs im Compilerbau?



Pattern Matching geht mit NFAs.

NFAs sind so nicht zu programmieren, aber:

Satz: Eine Sprache L wird von einem NFA akzeptiert $\Leftrightarrow L$ wird von einem DFA akzeptiert.

Konvertierung eines NFAs in einen DFA

Gegeben: Ein NFA $A = (Q, \Sigma, \delta, q_0, F)$

Wir konstruieren einen DFA $A' = (Q', \Sigma, \delta', q_0, F')$ wie folgt:

```
Q' = {{q0}}  
for each q in Q':  
  for each a ∈ Σ:  
    n = {p | δ(r, a) = p, r ∈ q}  
    Q' = Q' ∪ {n}  
    δ'(q, a) = n  
F' = {q ∈ Q' | q ∩ F ≠ ∅}  
{p, q} -> pq // Umbenennung
```

Abbildung 1: Konvertierung NFA in DFA

Beispiel

δ	a	b
$\rightarrow q_0$	$\{q_0\}$	$\{q_1, q_2\}$
q_1	$\{q_2\}$	$\{q_1\}$
$*q_2$	-	$\{q_0, q_2\}$

NFA

δ'	a	b
$\rightarrow \{q_0\}$	$\{q_0\}$	$\{q_1, q_2\}$
$*\{q_1, q_2\}$	$\{q_2\}$	$\{q_0, q_1, q_2\}$
$*\{q_2\}$	-	$\{q_0, q_2\}$
$*\{q_0, q_2\}$	$\{q_0\}$	$\{q_0, q_1, q_2\}$
$*\{q_0, q_1, q_2\}$	$\{q_0, q_2\}$	$\{q_0, q_1, q_2\}$

ein
zustand

\hookrightarrow umbenennen: $\{q_0, q_1, q_2\} \rightarrow q_0, q_1, q_2$

Minimierung eines DFAs

Ist ist der DFA A nicht vollständig, wird ein Fehlerzustand q_e , der kein Endzustand ist, hinzugefügt und in alle leeren Tabellenfelder eingetragen.

Dann wird eine Matrix generiert, die für alle Zustandspaare sagt, ob die beiden Zustände zu einem verschmelzen können.

```
for each (p,q) with  $p,q \in Q \times Q$ ,  $p \neq q$ :
  if ( $p \in F$  and  $q \notin F$ ) or ( $p \notin F$  and  $q \in F$ ):
     $D(p,q) = \text{"-"}$ 
  else
     $D(p,q) = \epsilon$ 

repeat
  for each (q,p)  $\in Q \times Q$  with  $p \neq q$ :
    for each  $a \in \Sigma$ :
      if  $D(p,q) = \epsilon$  and
         $D(\delta(p,a), \delta(q,a)) \neq \epsilon$ :
         $D(p,q) = a$ 
until there are no changes

for each entry in  $D$  with  $D(p,q) = \epsilon$ :
  combine  $p$  and  $q$  to one state
```

Abbildung 2: DFA Minimierung

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: *leere Wort*
 - 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:
 - $E + F$ ist ein regex und bezeichnet die Vereinigung $L(E + F) = L(E) \cup L(F)$
 - EF ist ein regex und bezeichnet die Konkatenation $L(EF) = L(E)L(F)$
 - 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, +

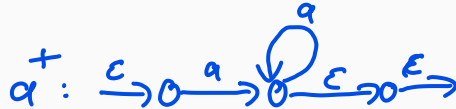
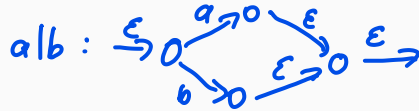
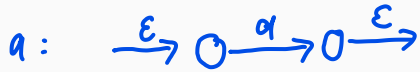
$$ab^* \neq (ab)^*$$

Wichtige Identitäten

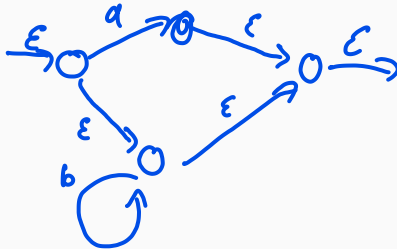
Satz: Sei A ein DFA $\Rightarrow \exists$ regex R mit $L(A) = L(R)$.

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

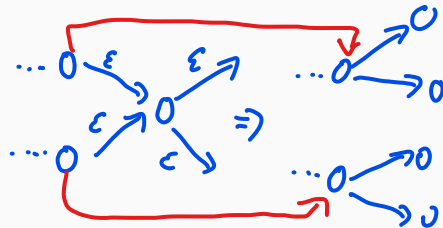
Beispiel: Umwandlung eines regex in einen NFA



NFA
 $a|b^*$



ϵ entfernen:

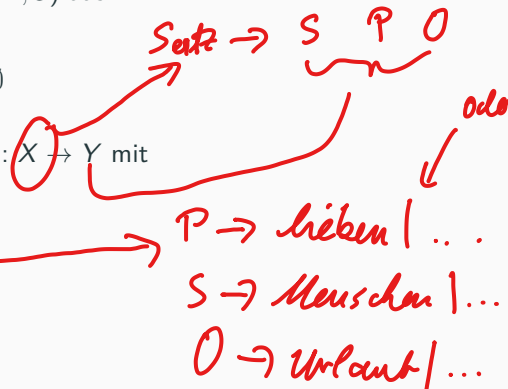


Formale Grammatiken

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

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


mind. 1 Nichtterminal


 $Satz \rightarrow S \ P \ O$
oder
 $P \rightarrow lieben \mid \dots$
 $S \rightarrow Menschen \mid \dots$
 $O \rightarrow Urlaub \mid \dots$

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:

- 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.

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.

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.

Satz: Die von DFA's akzeptierten Sprachen
(+ NFA's)
sind die regulären Sprachen.

Das Pumping Lemma für reguläre Sprachen

Satz: Das *Pumping Lemma* für reguläre Sprachen:

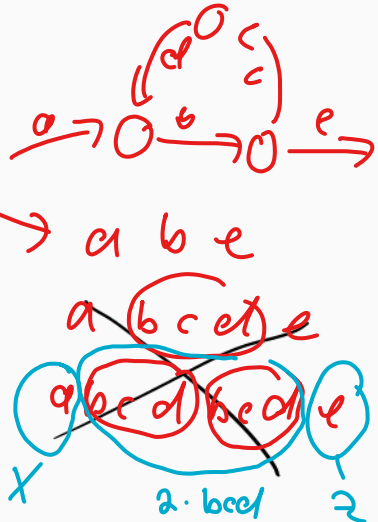
Sei L eine reguläre Sprache.

$\Rightarrow \exists$ Konstante $n \in \mathbb{N}$:

$$\forall \substack{w \in L \\ |w| \geq n} \exists x, y, z \in \Sigma^* \text{ mit } w = xyz, y \neq \epsilon, |xy| \leq n:$$
$$\forall_{k \geq 0} xy^kz \in L$$

↓
an for new pen

abcde
 abcde
 ab...e



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

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.

$a^n b^n$

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

- einfach zu programmieren
 - unübersichtlich \Rightarrow Fehler
- DFA's zu programmieren
- einen Generator benutzen

"langsam"

Wrap-Up

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

LICENSE



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