

Softwareparadigmen

Dieses Skriptum basiert auf der Softwareparadigmen Übung im
Sommersemester 2011 und dem Vorlesungsskriptum 2007.
Vorlesung von Alexander Felfernig

Übungsskriptum verfasst von Daniel Gruß. Version 24. Dezember 2011.
Fehlerfunde bitte melden an gruss@student.tugraz.at.

Inhaltsverzeichnis

1	Einleitung	1
1.1	Syntax	1
1	Definitionen	1
2	Grammatiken und Sprachen	2
3	Compilerbau	6
1.2	Semantik von Programmiersprachen	7

Kapitel 1

Einleitung

Dieses Skriptum versucht einen kompakten und übersichtlichen Zugang zu den Themen dieser Lehrveranstaltung zu schaffen. Der Vollständigkeit halber verweisen wir in der Fußnote auf die entsprechenden Kapitel des Vorlesungsskriptums.

1.1 Syntax

In der Informatik stehen wir oft vor dem Problem, dass wir eine Art Text auf Fehler überprüfen wollen und in eine andere Darstellungsform übersetzen wollen. Ein Webbrowser prüft beispielsweise HTML-Dokumente auf Fehler und übersetzt diese sofern möglich in eine praktische visuelle Darstellung.

Dieses Problem wollen wir auch lösen können. Dazu brauchen wir einige Definitionen.

1 Definitionen

Definition 1.1 (Alphabet): Ein Alphabet Σ ist eine endliche Menge von Symbolen.

Binärzahlen haben das Alphabet $\Sigma = \{0, 1\}$. Eine einfache Variante der Markup-Sprache HTML hat das Alphabet $\Sigma = \{<, /, >, \dots, a, b, c, \dots\}$.

Definition 1.2: Σ^* ist die Menge aller beliebigen Konkatenationen von Symbolen aus Σ . Ein Element aus Σ^* nennen wir Wort.

Für $\Sigma = \{0, 1\}$ ist $\Sigma^* = \{\varepsilon, 0, 1, 00, 01, 10, 11, 000, 001, \dots\}$.

Definition 1.3 (Sprache): Eine Sprache \mathcal{L} ist eine Teilmenge (\subseteq) von Σ^* .

Sei $\Sigma = \{0, 1\}$. Definieren wir die Sprache der Binärzahlen \mathcal{B} , müssen wir zu jedem Wort aus Σ^* entscheiden ob dieses Wort in \mathcal{B} enthalten ist. Im Fall der Binärzahlen könnten

wir z.B. ε herausnehmen, dann ist die Sprache $\mathcal{B} = \Sigma^* \setminus \{\varepsilon\}$.

Würden wir die Programmiersprache \mathcal{C} als formale Sprache definieren, so wäre ein gesamtes (gültiges) \mathcal{C} -Programm ein Wort der Sprache, also ein Element der Menge \mathcal{C} .

Definition 1.4 (Compiler): Seien \mathcal{A} und \mathcal{B} Programmiersprachen. Ein Compiler ist ein Programm, welches Programme von \mathcal{A} nach \mathcal{B} übersetzt.

Ein einfacher Compiler würde beispielsweise Binärzahlen zu Dezimalzahlen übersetzen.

2 Grammatiken und Sprachen

Oft möchte man nur prüfen ob ein Wort in einer Sprache enthalten ist. Es ist umständlich die Menge aller Wörter einer Sprache aufzuschreiben und darin zu suchen. Daher wollen wir eine kürzere und Methode finden um Sprachen zu beschreiben und diese Überprüfung durchzuführen.

Über Automaten (FSM) ist dies möglich. Für Binärzahlen könnte man dies wie in Abbildung 1.1 machen.

Wir haben eine Eingabe w (einem potentiellen Wort der Sprache). Wir beginnen im Zustand S . Die einzigen gültigen Übergänge sind die zu Zustand A . Es muss also entweder eine 0 oder eine 1 gelesen werden (wir schreiben auch gematcht bzw. geparkt). Andernfalls gibt es keinen gültigen Übergang und da S kein Endzustand ist wäre die Eingabe nicht gültig. Im Zustand A wurde bereits eine 0 oder 1 gelesen, also ist $w \in \mathcal{B}$. Die Übergänge von A zu A erlauben weitere Zeichen zu lesen und so längere Wörter zu erhalten, die in der Sprache enthalten sind.

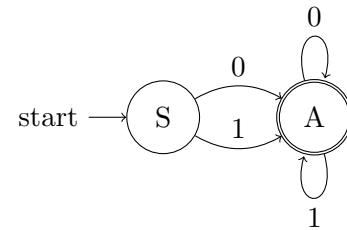


Abb. 1.1: Binärzahlen FSM

Eine ähnliche Herangehensweise stellen Grammatiken dar:

Definition 1.5 (Grammatik): Eine Grammatik ist ein 4-Tupel (V_N, V_T, S, Φ) .

- V_N ist eine endliche Menge von Nonterminalen (entsprechen Zuständen der FSM),
- V_T ist eine endliche Menge von Terminalen (entsprechen Alphabet der Sprache),
- $S \in V_N$ das Startsymbol (wie Initialzustand der FSM),
- $\Phi = \{\alpha \rightarrow \beta\}$ eine endliche Menge von Produktionsregeln (Zustandsübergänge).

α ist hierbei eine beliebige Aneinanderreihung von Terminalen und Nonterminalen, die zumindest ein Nonterminal enthält, also $(V_N \cup V_T)^* V_N (V_N \cup V_T)^*$.

β ist eine beliebige Aneinanderreihung von Terminalen und Nonterminalen, einschließlich der leeren Menge, also $(V_N \cup V_T)^*$.

Anhand der Binärzahlen wollen wir nun genauer betrachten wie die Produktionsregeln funktionieren. Bei den Zustandsübergängen des Automaten hatten wir beispielsweise ein Wort $w = AB$ mit $A \in \{0, 1\}$ und $B \in \{0, 1\}^*$. Mit dem Zustandsübergang würden wir nun das A weglassen und hätten für das verbleibende zu lesende Wort $w' = B$. Wir ersetzen also AB durch B . Bei den Produktionsregeln der Grammatik halten wir ganz konkret fest was wir ersetzen.

Um die Grammatik der Binärzahlen G_B zu definieren müssen wir die Elemente des oben definierten 4-Tupels beschreiben. Die Menge der Terminale ist das Alphabet $\{0, 1\}$. Den Startzustand nennen wir S . In Abbildung 1.2 versuchen wir Produktionsregeln anzugeben um B zu definieren.

$$\begin{aligned} S &\rightarrow 0 \\ S &\rightarrow 1 \\ S &\rightarrow 0S \\ S &\rightarrow 1S \end{aligned}$$

Abb. 1.2: Binärzahlen Grammatik

Um zu verifizieren wann ein Wort von einer Grammatik erzeugt werden kann müssen wir nun zuerst definieren was eine Ableitung ist.

Definition 1.6: Sei (V_N, V_T, S, Φ) eine Grammatik und $\alpha, \beta \in (V_N \cup V_T)^*$. Wenn es 2 Zeichenfolgen τ_1, τ_2 gibt, so dass $\alpha = \tau_1 A \tau_2$, $\beta = \tau_1 B \tau_2$ und $A \rightarrow B \in \Phi$, dann kann β direkt (in einem Schritt) von α abgeleitet werden ($\alpha \rightarrow \beta$).

Diese Definition ist nur geeignet um eine Aussage darüber zu treffen was in genau einem Schritt abgeleitet werden kann. Wir definieren daher die reflexive Hülle dieses Operators.

Definition 1.7: Sei (V_N, V_T, S, Φ) eine Grammatik und $\alpha, \beta \in (V_N \cup V_T)^*$. Wenn es $n \in \mathbb{N}$ Zeichenfolgen τ_1, τ_n gibt, so dass $\alpha \rightarrow \tau_1, \tau_1 \rightarrow \tau_2, \dots, \tau_{n-1} \rightarrow \tau_n, \tau_n \rightarrow \beta$, dann kann β von α (in n Schritten) abgeleitet werden ($\alpha \xrightarrow{+} \beta$).

Diese Definition ist für $n \geq 1$ geeignet. Wenn $\alpha = \beta$ ist, wäre $n = 0$. Wir definieren die reflexive, transitive Hülle durch eine Verknüpfung dieser beiden Fälle.

Definition 1.8: Es gilt $\alpha \xrightarrow{*} \beta$, genau dann wenn $\alpha \xrightarrow{+} \beta$ oder $\alpha = \beta$ (reflexive, transitive Hülle).

Mit dieser Definition können wir alle Wörter ableiten die diese Grammatik produziert.

Definition 1.9: Sei (V_N, V_T, S, Φ) eine Grammatik G . G akzeptiert die Sprache $L(G) = \{w \mid S \xrightarrow{*} w, w \in V_T^*\}$, d.h. die Menge aller Wörter w die in beliebig vielen Schritten aus dem Startsymbol S ableitbar sind und in der Menge aller beliebigen Konkatenationen von Terminalsymbolen V_T enthalten sind.

Die Äquivalenz von zwei Sprachen zu zeigen ist im Allgemeinen nicht trivial. Wenn wir versuchen eine Sprache \mathcal{L} durch eine Grammatik G zu beschreiben ist es im Allgemeinen nicht trivial zu zeigen dass $L(G) = \mathcal{L}$.

An dieser Stelle möchten wir festzuhalten: Um die Gleichheit zweier Mengen (Sprachen) M, N zu zeigen muss gezeigt werden, dass jedes Element (Wort) aus M in N enthalten ist und jedes Element (Wort) aus N in M enthalten ist. Ungleichheit ist daher viel leichter zu zeigen, da es genügt ein Element (Wort) zu finden welches nicht in beiden Mengen (Sprachen) enthalten ist. Der geneigte Leser kann probieren die Gleichheit oder Ungleichheit der Sprache unserer oben definierten Grammatik und der Sprache der Binärzahlen zu zeigen.

Definition 1.10: Ein Programm $P_{\mathcal{L}}$ welches für ein Wort w entscheidet ob es in der Sprache \mathcal{L} enthalten ist (d.h. **true** dann und nur dann zurückliefert wenn es enthalten ist), nennen wir **Parser**.

Chomsky-Sprachhierarchie

Durch die Grammatik können wir entscheiden ob ein Wort in einer Sprache enthalten ist. Im Falle der Binärzahlen haben wir eine kurze und einfache Grammatik und können einen Parser schreiben welches entscheidet ob ein Eingabewort in der Sprache enthalten ist. Definieren wir eine Programmiersprache wie \mathcal{C} formal, so wird nicht nur die Anzahl der Produktionsregeln höher sein sondern auch die Komplexität der Produktionsregeln ($\alpha \rightarrow \beta$ mit komplexen Ausdrücken für α und β). Es ist dann erheblich schwieriger einen Parser zu schreiben. Wir haben also Interesse daran möglichst einfache Grammatiken zu finden. Dazu müssen wir Grammatiken nach ihrer Komplexität vergleichen können.

Wir haben Produktionsregeln definiert durch $\alpha \rightarrow \beta$, wobei α zumindest ein Nonterminal enthält.

Definition 1.11: Eine Grammatik ist nach der Chomsky-Sprachhierarchie:

- allgemein/uneingeschränkt (unrestricted)

Keine Restriktionen

- Kontext-sensitiv (context sensitive): $|\alpha| \leq |\beta|$

Es werden nicht mehr Symbole gelöscht als produziert werden)

- Kontext-frei (context free): $|\alpha| \leq |\beta|$, $\alpha \in V_N$

Wie Kontext-sensitiv; außerdem muss α genau **ein** Non-Terminal sein

- regulär (regular): $|\alpha| \leq |\beta|$, $\alpha \in V_N$, $\beta = aA$, $a \in V_T \cup \{\varepsilon\}$, $A \in V_N \cup \{\varepsilon\}$

Wie Kontext-frei; außerdem ist $\beta = aA$ wobei a ein Terminal oder ε ist und A ein Nonterminal oder ε . (Anmerkung: $\varepsilon\varepsilon = \varepsilon$)

Es gilt $\mathbb{L}_{\text{regulär}} \subset \mathbb{L}_{\text{context free}} \subset \mathbb{L}_{\text{context sensitive}} \subset \mathbb{L}_{\text{unrestricted}}$ (\mathbb{L}_x Menge aller Sprachen der Stufe x).

Nicht alle Grammatiken können in eine äquivalente Grammatik einer stärker eingeschränkten Stufe umgewandelt werden. Um zu zeigen dass es sich um echte Teilmengen ($A \subset B$) handelt müssen wir zeigen dass alle Elemente aus A in B enthalten sind und mindestens ein Element aus B nicht in A enthalten ist. Eine Beweisskizze dazu findet sich im Vorlesungsskriptum auf Seite 12. Dort wird gezeigt, dass es für $L = \{a^n b a^n | n \in \mathbb{N}_0\}$ äquivalente reguläre Grammatik G gibt, d.h. $\nexists G \in \mathbb{L}_x : L(G) = L$.

Die Chomsky-Sprachhierarchie unterscheidet Sprachen anhand der Komplexität der produzierten Sprache.

Parser

Wir überspringen an dieser Stelle den BPARSE-Algorithmus (siehe Vorlesungsskriptum) und betrachten stattdessen einen Recursive Descent Parser (RDP).

Bei einem RDP werden alle Nonterminale in Funktionen übersetzt und diese Funktionen behandeln die verschiedenen Produktionsregeln. Die Eingabe wird in die Terminale unterteilt (auch Tokens genannt). Wir verwenden im Pseudo-Code die Variable `token`, die immer das aktuelle Token enthält, sowie die Funktion `nextToken()`, die `token` auf das nächste Token setzt. Der Parser ruft die Startfunktion auf und gibt `true` zurück wenn `token` leer ist (d.h. das Ende der Eingabe erreicht wurde). `ERROR` führt dazu dass der Parser die Eingabe (das Wort) nicht akzeptiert.

Sei $G_{\text{bin1}} = (\{S, T\}, \{0, 1\}, S, \{S \rightarrow 0, S \rightarrow 1T, T \rightarrow 0T, T \rightarrow 1T, T \rightarrow \varepsilon\})$. Der Pseudo-Code des RDP zu dieser Grammatik könnte so aussehen:

<pre> FUNC S() IF token == 0 nextToken() ELSE IF token == 1 nextToken() T() ELSE ERROR </pre>	<pre> FUNC T() IF token == 0 nextToken() T() ELSE IF token == 1 nextToken() T() ELSE IF token != epsilon ERROR </pre>
---	---

Sei $G_{\text{bin2}} = (\{S, T\}, \{0, 1\}, S, \{S \rightarrow 0, S \rightarrow 1T, T \rightarrow T0, T \rightarrow T1, T \rightarrow \varepsilon\})$. Wir wissen $L(G_{\text{bin1}}) = L(G_{\text{bin2}})$ (ohne Beweis) und versuchen auch hier einen einfachen rekursiven Parser zu schreiben.

<pre> FUNC S() IF token == 0 nextToken() ELSE IF token == 1 nextToken() T() ELSE ERROR </pre>	<pre> FUNC T() IF token == 0 T() // Endlos-Rekursion nextToken() ELSE IF token == 1 T() // Endlos-Rekursion nextToken() ELSE IF token != epsilon ERROR </pre>
---	---

Wir erhalten hier im Parser eine Endlos-Rekursion, da in $T \rightarrow T0, T \rightarrow T1$ ein Nonterminal ganz links steht. Wir nennen dies “Linksrekursion”.

Definition 1.12 (Mehrdeutig, Eindeutig): Sei $G = (V_N, V_T, S, \Phi)$ eine Grammatik. Wenn es für ein Wort $w \in L(G)$ mehrere unterschiedliche Ableitungssequenzen ω, ψ , d.h. $S \rightarrow \omega_1, \omega_1 \rightarrow \dots, \dots \rightarrow \omega_n, \omega_n \rightarrow w$, $S \rightarrow \psi_1, \psi_1 \rightarrow \dots, \dots \rightarrow \psi_k, \psi_k \rightarrow w$ wobei $\exists \psi_i : \psi_i \neq \omega_i$, ist es mehrdeutig. Wenn eine Grammatik ist genau dann eindeutig, wenn sie nicht mehrdeutig ist.

Definition 1.13 (Linksrekursiv): Eine Grammatik ist direkt linksrekursiv wenn sie eine Produktion der Form $A\alpha \rightarrow A\beta$ enthält, wobei A ein Nonterminal ist. Eine Grammatik ist indirekt linksrekursiv wenn sie Produktionen der Form $A\alpha \rightarrow A_1\beta_1, A_1\alpha_1 \rightarrow A_2\beta_2, \dots, A_n\alpha_n \rightarrow A\beta_n$ enthält, wobei A, A_i Nonterminale sind. Eine Grammatik ist linksrekursiv wenn sie direkt oder indirekt linksrekursiv ist.

Nun können wir mit Sprachen und Grammatiken umgehen und diese nach ihrer Komplexität einstufen.

3 Compilerbau

Wir wollen uns nun damit beschäftigen wie wir Compiler für Programmiersprachen bauen können. Wir hatten definiert dass ein Compiler eine Wort w einer Sprache \mathcal{L} entgegennimmt und die Übersetzung des Wortes w' in einer Sprache \mathcal{L}' zurückgibt. Konkreter erhält unser Compiler einen beispielsweise einen ASCII-String, wobei unser Alphabet $\Sigma_{\mathcal{L}}$ meist nicht dem ASCII-Alphabet entspricht. Betrachten wir beispielsweise die Programmiersprache \mathcal{C} so können wir auch Schlüsselwörter wie `int` in $\Sigma_{\mathcal{C}}$ haben. Außerdem gibt es eventuell Whitespaces (Leerzeichen, Tabulatoren, Zeilenumbrüche, etc.; je nach dem positionsabhängig), die keinen Einfluss darauf haben ob das Eingabewort ein Wort der Sprache \mathcal{C} ist, d.h. ein \mathcal{C} -Programm ist. Auch Variablen- oder Funktionsnamen werden abgesehen von einer gewissen Form die sie einhalten müssen (z.B. “dürfen nicht nur aus Zahlen bestehen”) keinen Einfluss darauf haben ob das Programm in der Sprache enthalten ist. Wenn wir dies berücksichtigen, verkürzen wir die Grammatik und vereinfachen so unseren Parser sowie eventuelle weitere Rechenschritte.

Lexikalische Analyse

Wir haben bereits in Abschnitt 2 von Tokens geredet. Von Tokens spricht man insbesondere bei einer Sprache die durch Whitespaces getrennt werden. Ein Token ist hierbei die kleinste Sequenz von Zeichen die für die Grammatik eine Bedeutung hat. Der erste Schritt der Kompilervorgangs ist es die Eingabe in eine Sequenz von Tokens umzuwandeln (dies kann natürlich wie bei unserem rekursiven Parser während dem Parsen passieren).

Wir wollen z.B. den C-Code `int main() { printf("helloworld"); return 123; }` in die Sequenz von Tokens `INT ID () { ID (STRING) ; RETURN NUM ; }` umwandeln. Wir nennen dies “Lexikalische Analyse”. Es fällt auf dass in der Token-Sequenz keine Namen oder Zahlen mehr vorkommen außerdem sind nun die Tokens voneinander getrennt. Es ist üblich zur lexikalischen Analyse nur reguläre Grammatiken zu verwenden. Diese sind mächtig genug um beispielsweise nur gewisse Variablennamen zu erlauben und können andererseits sehr schnell berechnet werden. Eine reguläre Grammatiken kann auch über einen regulären Ausdruck (Regular Expression/Regex) kompakt dargestellt werden.

Definition 1.14: Ein regulärer Ausdruck A ist wie folgt rekursiv definiert (mit regulären Ausdrücken Q und R):

$$A = \begin{cases} \varepsilon & \text{Leerstring} \\ t & \text{ein Terminal, d.h. } t \in V_T \\ Q|R & \text{entweder } A = Q \text{ oder } A = R \\ QR & \text{Konkatenation zweier regulärer Ausdrücke} \\ Q? & \text{entweder } A = Q \text{ oder } A = \varepsilon \\ Q^* & \text{beliebige Konkatenation von } Q \text{ mit sich selbst, d.h. } A \in \{\varepsilon, Q, QQ, \dots\} \\ Q+ & \text{Gleiche Bedeutung wie } QQ^*, \text{ d.h. mindestens ein } Q \\ (Q) & \text{Klammerung des Ausdrucks } Q \end{cases}$$

Grammatikalische Analyse

1.2 Semantik von Programmiersprachen