

Objektorientierte Programmierung

Kapitel 3: Syntaxdiagramme und Grammatikregeln

Stefan Brass

Martin-Luther-Universität Halle-Wittenberg

Wintersemester 2012/13

<http://www.informatik.uni-halle.de/~brass/oop12/>

Inhalt

- 1 Syntaxformalismen
 - Warum Syntaxformalismen wichtig sind
- 2 Syntaxdiagramme
 - Einführung in Syntaxdiagramme
- 3 Grammatik-Regeln
 - Grammatik-Regeln: Erklärung mit Syntax-Diagrammen
 - Beispiel: Zahlkonstanten
 - Hinweis zur Java-Spezifikation

Syntax-Formalismen (1)

- Ein Java-Programm ist eine Folge von Zeichen, aber nicht jede Folge von Zeichen ist ein Java-Programm.
- Ein Alphabet ist eine endliche, nicht-leere Menge, deren Elemente Zeichen heißen.

Java basiert auf dem Unicode-Zeichensatz. Dies ist eine Obermenge des ASCII-Zeichensatzes, der in Kapitel 1 gezeigt wurde.

- Ein Wort über einem Alphabet ist eine endliche Folge von Zeichen des Alphabets.

Nach dieser Definition ist ein Java-Programm ein Wort über dem Unicode-Zeichensatz.

- Eine formale Sprache (über einem Alphabet) ist eine (meist unendliche) Menge von Worten (über dem Alphabet).

Syntax-Formalismen (2)

- Es gibt verschiedene Formalismen, um klar und eindeutig eine formale Sprache (Menge von Zeichenfolgen) zu definieren, z.B.
 - Syntax-Diagramme
 - (Kontextfreie) Grammatik-Regeln
 - BNF (Backus-Naur-Form) ist eine spezielle Syntax für kontextfreie Grammatikregeln (recht verbreitet).
- Syntaxdiagramme und kontextfreie Grammatiken haben die gleiche Ausdruckskraft, d.h. können die gleichen Mengen von Zeichenfolgen beschreiben.

Syntax-Formalismen (3)

- Als professioneller Programmierer sollte man die Definition der Programmiersprache lesen können.
- Die Grammatik ist eine Referenz in unklaren Fällen, hilft aber auch zum Verständnis der Sprache.

Die syntaktischen Kategorien sind oft nützliche Konzepte.

- Es könnte ja sein, dass der Compiler einen Fehler enthält.

Man sollte sich nicht zum “Sklaven” eines einzelnen Compilers machen. Der Compiler hat nicht immer recht (aber meistens schon). Es gibt ein unabhängiges Gesetz (die Java-Spezifikation). Zumindest hilft die Aufklärung eines Fehlers, ihn zukünftig zu vermeiden. Herumprobieren (“wie hätte er es denn gern?”) ist alleine keine richtige Lösung.

Syntax-Formalismen (4)

- Mit kontextfreien Grammatiken definiert man nur eine Obermenge der Java-Programme, und stellt dann weitere einschränkende Forderungen, z.B.
 - der Typ des zugewiesenen Wertes muss zum Typ der Variablen passen.
- Kontextfreie Grammatiken (und damit auch Syntaxdiagramme) können solche Bedingungen nicht ausdrücken.

Dennoch bildet die kontextfreie Grammatik sozusagen das Rückgrat der Sprachspezifikation: Die syntaktischen Kategorien, die in der Grammatik eingeführt sind, werden gebraucht, um die semantischen Zusatzbedingungen überhaupt ausdrücken zu können. Ein allgemeines Prinzip der Programmierung ist, schwierige Aufgaben in Teilaufgaben zu zerlegen. Daher ist man zunächst damit zufrieden, dass die kontextfreie Grammatik eine Obermenge der gültigen Java-Programme beschreibt, in einem zweiten Schritt schränkt man die Menge dann durch Zusatzbedingungen weiter ein.

Inhalt

- 1 Syntaxformalismen
 - Warum Syntaxformalismen wichtig sind
- 2 Syntaxdiagramme
 - Einführung in Syntaxdiagramme
- 3 Grammatik-Regeln
 - Grammatik-Regeln: Erklärung mit Syntax-Diagrammen
 - Beispiel: Zahlkonstanten
 - Hinweis zur Java-Spezifikation

Syntaxdiagramme (3)

- Solche Diagramme können auf zwei Arten verwendet werden:
 - Um ein gültiges Wort der Sprache zu erzeugen, folgt man einem Pfad durch das Diagramm vom Start zum Ziel und druckt jedes Symbol in einem Kreis/Oval aus, den man durchläuft.
 - Um festzustellen, ob eine gegebene Eingabe syntaktisch korrekt ist, muss man einen Pfad durch das Diagramm finden, so dass beim Durchlaufen eines Kreises/Ovals das Zeichen darin das nächste Eingabezeichen ist.

Syntaxdiagramme (5)

- Es gibt Verzweigungsknoten, an denen man verschiedene ausgehende Kanten benutzen kann.
 - Z.B., nach der Eingangskante könnte man nach unten gehen, und die Ziffer 0 ausgeben/einlesen, oder man kann nach rechts gehen, um eine der Ziffern 1 bis 9 zu erhalten.
- Um eine gegebene Eingabe zu prüfen, hilft es meist, das nächste Eingabesymbol anzuschauen, um den richtigen Pfad auszuwählen.
 - Dies ist, was der Compiler auch macht. Man versucht sicherzustellen, dass an jeder Verzweigung, die Kreise/Ovale, die man in verschiedenen Richtungen erreichen kann, disjunkt sind.
- Gibt es keinen passenden Pfad: Syntaxfehler.

Syntaxdiagramme (7)

- Man kann an anderer Stelle definierte Diagramme als Module benutzen.
- Dies wird als Rechteck dargestellt, das den Namen des anderen Diagramms enthält:

Digit Sequence:



- Ein Rechteck steht also für eine syntaktische Kategorie (wie “Subjekt”, “Prädikat”, “Objekt”).

Syntaxdiagramme (9)

- Natürlich muss man nicht explizit das Diagramm der benutzten syntaktischen Kategorie einfügen:
 - Man merkt sich einfach, wo man im übergeordneten Diagramm war (bei welchem Rechteck),
 - durchläuft dann das Diagramm für die syntaktische Kategorie, die in dem Rechteck steht,
 - anschließend (wenn man mit diesem Diagramm fertig ist) kehrt man zum übergeordneten Diagramm zurück, und folgt dort dem Pfeil, der das Rechteck verläßt.

Syntaxdiagramme (10)

- Natürlich weiß man nach einiger Zeit, wofür z.B. die syntaktische Kategorie “Digit” steht, und braucht das zugehörige Diagramm dann nicht mehr nachzuschlagen.
- Jedes Diagramm definiert eine formale Sprache.
 - D.h. eine Menge von Zeichenketten.
- Wenn man ein Rechteck durchläuft, kann man jedes Element der zugehörigen Sprache ausdrucken bzw. einlesen.

Syntaxdiagramme (12)

- Syntaxdiagramme können “rekursiv” definiert sein, d.h. das Diagramm für eine syntaktische Kategorie X kann selbst einen mit X beschrifteten Kasten enthalten.

Entsprechend könnten sich auch zwei syntaktische Kategorien X und Y gegenseitig aufrufen. Es ist nicht verlangt, dass eine syntaktische Kategorie vor ihrer Verwendung definiert ist. Es müssen nur am Ende alle verwendeten syntaktischen Kategorien auch wirklich definiert sein.

- Der Pfad, den man für ein Wort der Sprache durch die Diagramme geht, muss immer endlich sein.

Man kann jetzt zwar nicht mehr vorab die Diagramme für die Kästen vollständig einsetzen, aber man könnte noch die tatsächlich durchlaufenen Kästen nach Bedarf “expandieren”.

Inhalt

- 1 Syntaxformalismen
 - Warum Syntaxformalismen wichtig sind
- 2 Syntaxdiagramme
 - Einführung in Syntaxdiagramme
- 3 **Grammatik-Regeln**
 - **Grammatik-Regeln: Erklärung mit Syntax-Diagrammen**
 - **Beispiel: Zahlkonstanten**
 - **Hinweis zur Java-Spezifikation**

Grammatik-Regeln (1)

- Man kann die formale Syntax einer Sprache auch mithilfe kontextfreier Grammatik-Regeln definieren.
- Dies geschieht z.B. in der Java-Sprachspezifikation.

Die Referenz-Handbücher zu den Datenbanksystemen “Oracle” und “DB2” enthalten dagegen Syntax-Diagramme. Beides wird also in der Praxis verwendet. Syntaxdiagramme sind am Anfang intuitiver, aber Grammatik-Regeln sind kompakter.

- Auch in dieser Notation wird unterschieden zwischen
 - syntaktischen Kategorien (im Diagramm: Rechteck, jetzt in *kursiver* Schrift) und
 - den letztendlich erzeugten Zeichen der definierten Sprache (im Diagramm: Kreis/Oval, jetzt in Teletype Schrift).

Grammatik-Regeln (3)

- Man wendet Grammatik-Regeln an, indem man ausgehend vom Startsymbol der Grammatik (z.B. *CompilationUnit*)
 - immer eine syntaktische Kategorie durch einen der darunter eingerückten Texte ersetzt,
 - bis der so erzeugte Text schließlich keine syntaktischen Kategorien mehr enthält, sondern nur noch Zeichen der zu definierenden Sprache.
- Weil die Ersetzung bei den Zeichen der erzeugten Sprache (Teletype-Schrift) endet, heißen diese auch Terminalsymbole der Grammatik.
- Die syntaktischen Kategorien (Worte in *kursiver* Schrift) heißen entsprechend Nichtterminalsymbole.

Das ganze Wort, d.h. die Kategorie, ist ein Nichtterminalzeichen.

Grammatik-Regeln (4)

- Es gibt viele Notationen für Grammatikregeln:
 - Z.B. kann man syntaktischen Kategorien auch in spitze Klammern $\langle \dots \rangle$ schreiben.
 - Manche Autoren unterstreichen auch die Zeichen der erzeugten Sprache, oder schreiben sie in $\underline{\dots}$.
 - Statt Einrückung und eine Alternative pro Zeile kann man die Alternativen auch durch “|” trennen, und den jeweils letzten Ersetzungstext durch “;” beenden.
 - In der Theorie ist jede Alternative eine eigene Grammatikregel, mit der zu ersetzenden syntaktischen Kategorie (ein Nichtterminalzeichen) auf der linken Seite, und dem Ersetzungstext (Folge von Terminal- und Nichtterminalzeichen) auf der rechten Seite, dazwischen meist “ \rightarrow ”.

Grammatik-Regeln (5)

- Im Java-Standard werden noch zwei Abkürzungen benutzt:
 - Wenn die alternativen Ersetzungstexte sehr kurz sind, z.B. einzelne Zeichen, wird nach der zu definierenden syntaktischen Kategorie “one of” geschrieben, und die Alternativen darunter in einer Zeile.

NonZeroDigit: one of

1 2 3 4 5 6 7 8 9

- Im Ersetzungstext können syntaktische Kategorien mit dem Index “*opt*” markiert werden: Dann kann man sie an dieser Stelle auch durch die leere Zeichenfolge ersetzen.
Falls mehrere Zeichen mit “*opt*” markiert sind, ist jedes für sich optional, d.h. man kann eine beliebige Teilmenge weglassen.

DecimalIntegerLiteral:

DecimalNumeral IntegerTypeSuffix_{opt}

Beispiel: Zahlkonstanten (1)

- Zahlkonstanten in Dezimalschreibweise (z.B. 123) sind in der Java-Spezifikation durch folgende Grammatikregeln definiert:

DecimalNumeral:

0

NonZeroDigit Digits_{opt}

Digits:

Digit

Digits Digit

Digit:

0

NonZeroDigit

NonZeroDigit: one of

1 2 3 4 5 6 7 8 9

Beispiel: Zahlkonstanten (2)

- Die Regel

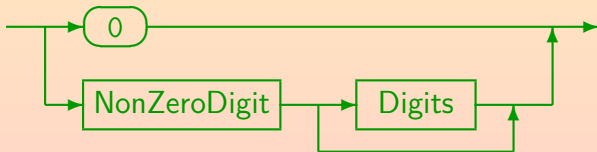
DecimalNumeral:

0

NonZeroDigit Digits_{opt}

entspricht folgendem Syntaxgraphen:

DecimalNumeral:



Diese etwas komplizierte Lösung ist gewählt, um führende Nullen auszuschließen, z.B. wäre 009 nicht gültig. Führende Nullen markieren die Oktalschreibweise, siehe nächstes Kapitel.

Hinweis zur Java-Spezifikation

- Die Java-Spezifikation enthält im Anhang eine weitere Grammatik für die Sprache Java.
- Dies ist nicht einfach eine Zusammenstellung der im Hauptteil angegebenen Grammatik-Regeln, sondern eine Variante, die besser für die Programmierung der Syntaxanalyse geeignet sein soll.

Die Regeln im Hauptteil sollen dagegen für die Erklärung verständlicher sein.

- Im Anhang wird eine etwas andere Grammatik-Syntax benutzt, in der mit
 - [...] optionale Teile,
 - {...} wiederholbare Teile, und mit
 - | Alternativen markiert sind.