

Syntaktische Analyse

Contents

Begriffe.....	2
First(A).....	2
Follow(A)	2
Indizmenge	2
SLL	2
Allgemein.....	2
Spezialfall.....	2
Linksfaktorisierung.....	3
FAQ	3
Wofür braucht man die Indizmenge?	3
Warum dürfen die Indizmengen von verschiedenen Produktionen von einem Nichtterminal sich nicht überlappen?	3
Warum soll eine Grammatik SLL(1) sein?.....	4
Wie erstellt man ein Parser für die Grammatik (Rekursiver Abstieg)?	4
Was macht lexer?	4
Was ist lexer.current und lexer.lex()?	4
Was ist expect()?	4
Wie beweist man Formal, dass Grammatik nicht SLL1 ist?	4
Aufgaben	5
First und Follow	5
ohne eps	5
mit eps.....	5
Lambda Ausdrücke	5
Grammatik.....	5
Indizmenge	6
Erklärung.....	6
Vorgehensweise für die Indizberechnung.....	6
Parser	6
JSON	8
Grammatik.....	8
Parser	8

Linksfaktorisierung.....	9
Prolog Parser	9
Grammatik.....	9
Parser	10
Relationen zwischen First und Follow (Teilmengen)	11
SGML	12
Aufgabe	12
Parser	13

Begriffe

First(A)

First₁(A) enthält alle Terminale, die bei irgendeiner Ableitung von A an erster Stelle stehen.

! Wenn es ein Regel (A → ε) gibt ----> ε ist immer in First(A)

Follow(A)

Follow₁(A) enthält alle Terminale, die bei irgendeiner Ableitung von S (dem Startsymbol) direkt hinter A stehen. Sie müssen sich also überlegen, welche Ableitungsschritte zu einem Vorkommen von A führen, und was dann direkt dahinter stehen kann.

Indizmenge

Indizmenge I von A → a ist mit k Token Lookahead genau

First_k(a Follow_k(A)).

Diese komische Konstruktion mit Follow ist nur dazu da, um die Fälle a=ε zu kompensieren.

Wenn es nur eine Produktion für diesen Nichtterminal gibt, dann kann man als Indizmenge für diesen Nichtterminal einfach „egal“ schreiben.

SLL

a und b sind die verschiedene Produktionen von einem Nichtterminal. Hier werden nicht die verschiedene Nichtterminale verglichen!

Allgemein

Eine KFG ist genau dann eine SLL(k)-Grammatik, wenn für alle Paare von Produktionen (von einem Nichtterminal) A → a | b, (a≠b) gilt:

First_k(a Follow_k(A)) SCHNITT First_k(b Follow_k(A)) = Leer

(== verschiedene Produktionen von einem Nichtterminal haben vollständig disjunkte Indizmengen)

Sonst wäre die Unterscheidung mit k-Lookahead nicht möglich.

Spezialfall

(oft in Klausuren!)

k=1, a !->* ε und b !->* ε (keine ε möglich):

SLL <=> First(a) SCHNITT First(b) = leer

$k=1$, $a \rightarrow^* \epsilon$ und $b \not\rightarrow^* \epsilon$ (Eps nur bei a möglich):
 $SLL \leftrightarrow Follow(A) \cap First(b) = \emptyset$

Linksfaktorisierung

Idee: Gemeinsame Anfang nach eine Variable verschieben, und eps | Rest nach andere Regel

Grammar Engineering:

Linksfaktorisierung



Grammatik mit Produktionen $X \rightarrow \gamma\alpha \mid \gamma\beta$ mit gemeinsamen Anfang γ ist nicht SLL, wenn Länge $|\omega|$, $\gamma \Rightarrow^* \omega$, unbeschränkt.⁹

Lösung: vertage Entscheidung zwischen α und β bis nachdem γ konsumiert wurde:

$$\begin{aligned} X &\rightarrow \gamma X' \\ X' &\rightarrow \alpha \mid \beta \end{aligned}$$

Wobei X' ein neues Nichtterminal ist.

Beispiel zur Linksfaktorisierung



If-Produktionen von SIMPLE:

$$\begin{aligned} If &\rightarrow \text{if } BraceExp \text{ Statement} \\ &\quad \mid \text{if } BraceExp \text{ Statement } \text{else } Statement \end{aligned}$$

haben gemeinsamen Anfang $\text{if } BraceExp \text{ Statement}$.
 Linksfaktorisierung ergibt:

$$\begin{aligned} If &\rightarrow \text{if } BraceExp \text{ Statement } If' \\ If' &\rightarrow \epsilon \mid \text{else } Statement \end{aligned}$$

Nach Linksfaktorisierung sind die Produktionen SLL.

FAQ

Wofür braucht man die Indizmenge?

Mit Hilfe der Indizmenge kann Compiler eine Entscheidung machen, welche Produktionen gewählt sein müssen.

Warum dürfen die Indizmengen von verschiedenen Produktionen von einem Nichtterminal sich nicht überlappen?

Sonst gäbe es mehrere korrekte Entscheidungen für Compiler -> Codegeneration nicht eindeutig.

Warum soll eine Grammatik SLL(1) sein?

Sonst wäre die Unterscheidung mit 1-Lookahead nicht möglich -> Kein effizienten Compiler möglich.

Wie erstellt man ein Parser für die Grammatik (Rekursiver Abstieg)?

Es muss ein Parser für ein Nichtterminal implementiert werden.

Wichtigste Bausteine: `lexer.current`, `lexer.lex()`, `switch(lexer.current)`, Schleifen, default: error, manchmal Schleifen mit `while(true)` und `break` nach dem `switch` mit `break` und `continue` in verschiedenen Cases, `return new <Classname>`.

Jeder Parser wird als eine Methode implementiert, die entweder ein Objekt von passendem Typ zurückliefert oder `error()` aufruft.

Name: „**Rekursiver Abstieg**“

Was macht lexer?

Lexer ist ein „Lexical-Analyzer“-Objekt .

Was ist lexer.current und lexer.lex()?

`lexer.current` ist current Token Type.

`lexer.lex()` setzt `lexer.current` auf das nächste Token.

Was ist expect()?

Gute Sache: `expect(Type)` prüft Typ und verschiebt Token-Pointer

Vereinfachung: Gemeinsame Parse-Funktion für (fast) alle Terminale:

```
void expect (TokenType e)
{ if lexer.current == e then lexer.lex() else error() }
```

Wie beweist man formal, dass Grammatik nicht SLL1 ist?

Gegeben sei die folgende Grammatik, die eine Untermenge von SGML beschreibt:

$$\begin{aligned} SGML &\rightarrow < \text{ident} > Children < / > \\ Children &\rightarrow \varepsilon \mid SGML Children \end{aligned}$$

Das Startsymbol dieser Grammatik ist *SGML*.

(a) Begründen Sie formal, warum die obige Grammatik nicht in SLL(1)-Form ist. [3 Punkte]

Children: a | b , a -> eps, b !-> eps --> 2.Spezialfall für SLL

(a) Es gilt $\text{Follow}_1(Children) = \{<\}$. Wir prüfen die SLL(1)-Bedingung anhand der Produktionen von *Children*:

$$\begin{aligned} &\text{First}_1(\varepsilon \cdot \text{Follow}_1(Children)) \cap \text{First}_1(SGML \cdot Children \cdot \text{Follow}_1(Children)) \\ &= \text{Follow}_1(Children) \cap \text{First}_1(SGML) \\ &= \{<\} \cap \{<\} \\ &\neq \emptyset \end{aligned}$$

Aufgaben

First und Follow

ohne eps

Grammatik ($=$, id , $*$ sind Terminale) :

$S \rightarrow L = R \mid R$

$L \rightarrow *R \mid id$

$R \rightarrow L$

$First1(L) = \{*, id\}$

$First1(R) = \{*, id\}$

$First1(S) = First1(L) \cup First1(R) = \{*, id\}$

$Follow1(S) = \{\#\}$

$Follow1(S) = \{=, \#\}$

$Follow1(S) = \{=, \#\}$

mit eps

Grammatik:

$A \rightarrow BC$

$B \rightarrow A \mid \epsilon$

$C \rightarrow id$

$First(A) = \{id\}$: $A \rightarrow BC \rightarrow \epsilon$ $C \rightarrow id$

$First(B) = \{id, \#\}$: $\#$ wegen ϵ , id : $B \rightarrow A \rightarrow BC \rightarrow \epsilon$ $C \rightarrow C \rightarrow id$

$First(C) = \{id\}$: trivial

$Follow(A) = \{id, \#\}$: $A \rightarrow BC \rightarrow AC \rightarrow \dots \rightarrow Follow(A) + \#$ wegen Startsymbol

$Follow(B) = \{id\}$: nur $C \rightarrow id$

$Follow(C) = \{id, \#\}$: $A \rightarrow BC \rightarrow AC \rightarrow BCC \rightarrow BC id$, $A \rightarrow BC \#$

Lambda Ausdrücke

Grammatik

$L \rightarrow Abstr$

$Abstr \rightarrow \lambda id . Abstr \mid App$

$App \rightarrow App Atom \mid Atom$

$Atom \rightarrow Var \mid (Abstr)$

Das Startsymbol der Grammatik soll L heißen, entweder nutzen wir diesen Namen für das niedrigstprioritäre Syntaxelement oder lassen es wie hier auf das Element für Abstraktion ableiten.

Die Abstraktion hat die niedrigste Priorität, damit kommt die dazugehörige Regel in der Grammatik am Anfang.

Ganz allgemein gilt: In der Regel für das Syntaxelement der Priorität n nutzen wir nur die Nichtterminale der Syntaxelemente mit Priorität n oder $n + 1$.

Indizmenge

```
App --> Abstr AppList [egal]
AppList --> eps[ ], # | Abstr AppList [lambda, id, (]
Abstr --> Atom [id, (] | lambda id . Abstr [lambda]
Atom --> Var [id] | ( App ) [ (]
Var --> id [egal]
Indizmenge in [ ]
```

Erklärung

App: egal, nur 1 Produktion (alternativ [id, (, lambda])

AppList: eps: schaue, was nach dem App kommen kann (Follow(App)), da AppList nur links bei App vorkommt. **Abstr -> alle Indizes von Abstr**

Abstr: **Atom -> alle Indizes von Atom**, zweite Produktion beginnt mit Terminalsymbol

Atom: in Var kommt immer id, zweite Produktion beginnt mit Terminalsymbol

Var: egal, nur 1 Produktion (alternativ [id])

Vorgehensweise für die Indizberechnung

1. Falls es nur eine Produktion gibt: „egal“
2. Falls es nur Terminalsymbole gibt: diese Terminalsymbole
3. Gehe von unten nach oben, berechne einfache Indizmengen, die Terminalsymbole als 1. Zeichen haben (hier wäre es Var, Atom, lambda bei Abstr)
4. Indiz von einem Nichtterminal ist die Vereinigung von allen Indizes seinen Produktionen
5. Epsilon bei X: berechne First(Follow(X))

Parser

```
App  → Abstr AppList  egal / {id, (, λ}
AppList → ε { }, # | Abstr AppList {id, (, λ}
Abstr → Atom {id, (} | λ id . Abstr {λ}
Atom  → Var {id} | ( App ) { ( }
Var   → id  egal / {id}
```

Diese Grammatik ist SLL(1). Rechnen Sie nach, indem Sie die nötigen Indizmengen explizit aufstellen und in die dafür vorgesehenen Kästchen eintragen.

- (c) Vervollständigen Sie den unten angegebenen rekursiven Abstiegsparser so, dass er [12 Punkte] die Variante des λ -Kalküls aus (b) in folgende Datenstruktur parst:

```
1 abstract class Expr {} // expression
2
3 class Var extends Expr {
4     public Var(String name) ...
5 }
6 class Abstr extends Expr { // abstraction
7     public Abstr(String varname, Expr body) ...
8 }
9 class App extends Expr { // application
10    public App(Expr left, Expr right) ...
11 }
```

```

1 Expr parseAtom() {
2     // gegeben
3     ...
4 }
5
6 String expect(Token tok) {
7     if (lexer.current != tok)
8         error();
9     lexer.lex();
10    return lexer.string_value;
11 }
12
13 Expr parseApp() {
14     Expr l = parseAbstr();
15     return parseAppList(l);
16 }

```

Beispiellösung:

```

1 Expr parseAbstr() {
2     switch (lexer.current) {
3     case IDENT: case LPAREN:
4         return parseAtom();
5     case LAMBDA:
6         expect(LAMBDA);
7         String varname = expect(IDENT);
8         expect(DOT);
9         Expr body = parseAbstr();
10        return new Abstr(varname, body);
11    default:
12        error();
13    }
14 }
15
16 // Hier müssen die linken Hälften der Applikationen durchgefädelt werden
17 Expr parseAppList(Expr left) {
18     switch(lexer.current) {
19     case RPAREN: case EOF: // epsilon, da war gar keine Applikation
20         return left;
21     case IDENT: case LAMBDA: case LPAREN:
22         Expr right = parseAbstr();
23         left = new App(left, right);
24
25         return parseAppList(left);
26     default:
27         error();
28     }
29 }

```

JSON

Grammatik

$Value \rightarrow \underline{string} \mid \underline{number} \mid Object$
 $Object \rightarrow \{ Members \}$
 $Members \rightarrow Pair\ Members'$
 $Members' \rightarrow _ Members \mid \epsilon$
 $Pair \rightarrow \underline{string} _ Value$

Parser

```
JSONValue parseValue() {
    String s;
    switch (lexer.current.type) {
        case STRING: s = lexer.current.text; lexer.lex(); return new JSONString(s);
        case NUMBER: s = lexer.current.text; lexer.lex(); return new JSONNumber(s);
        case LCURLY: return parseObject();
        default: error();
    }
}

JSONObject parseObject() {
    switch (lexer.current.type) {
        case LCURLY: lexer.lex();
            JSONObject o = new JSONObject(parseMembers());
            if (lexer.current.type != RCURLY) error();
            lexer.lex();
            return o;
        default: error();
    }
}

List<Pair> parseMembers() {
    List<Pair> l = new LinkedList<Pair>();
    l.add(parsePair());
    while (true) {
        switch (lexer.current.type) {
            case COMMA: lexer.lex();
                l.add(parsePair());
                continue;
            case RCURLY: break;
            default: error();
        }
        break;
    }
    return l;
}

Pair parsePair() {
    if (lexer.current.type != TokenType.STRING) error();
    JSONString s = new JSONString(lexer.current.text);
    lexer.lex();

    if (lexer.current.type != TokenType.COLON) error();
    lexer.lex();

    JSONValue v = parseValue();
    return new Pair(s, v);
}
```


Linksfaktorisierung

Gegeben sei die folgende Grammatik für Prolog-Fakten mit dem Startsymbol *Fact*:

$$\begin{aligned} Fact &\rightarrow Term . \\ Term &\rightarrow \mathbf{atom} \mid \mathbf{atom} (Termlist) \mid \mathbf{number} \mid \mathbf{variable} \\ Termlist &\rightarrow Term Termlist' \\ Termlist' &\rightarrow \varepsilon \mid , Term Termlist' \end{aligned}$$

- (a) Geben Sie die $First_1$ - und $Follow_1$ -Mengen für die Nichtterminale *Fact* und *Term* an. [4 Punkte]
- (b) Geben Sie eine Linksfaktorisierung an, so dass die Grammatik SLL(1) ist. [2 Punkte]
Es genügt geänderte bzw. neue Zeilen der Grammatik anzugeben.

Nicht SLL(1), da es die gleiche First-Mengen bei verschiedenen Produktionen vom „Term“ gibt.

(a)

$$\begin{aligned} First_1 (Fact) &= \{\mathbf{atom}, \mathbf{number}, \mathbf{variable}\} \\ First_1 (Term) &= \{\mathbf{atom}, \mathbf{number}, \mathbf{variable}\} \\ Follow_1 (Fact) &= \{\#\} \\ Follow_1 (Term) &= \{., ,, \} \end{aligned}$$

(b)

$$\begin{aligned} Term &\rightarrow \mathbf{atom} Atom' \mid \mathbf{number} \mid \mathbf{variable} \\ Atom' &\rightarrow \varepsilon \mid (Termlist) \end{aligned}$$

Prolog Parser

Grammatik

Fact -> Term .

Term -> **atom** Atom' | **number** | **variable**

Atom' -> **eps** | (Termlist)

Termlist -> Term Termlist'

Termlist' -> eps | , Term Termlist'

Parser

```
1 void parseFact() {
2     parseTerm();
3     if (token.getType() != dot) {
4         error();
5     }
6     nextToken();
7 }
8
9 void parseTerm() {
10     switch (token.getType()) {
11         case atom:
12             nextToken();
13             parseAtom';
14             break;
15         case number:
16         case variable:
17             nextToken();
18             break;
19         default:
20             error();
21     }
22 }
23
24 void parseAtom'() {
25     switch (token.getType()) {
26         case dot:
27         case comma:
28         case rp:
29             break;
30         case lp:
31             nextToken();
32             parseTermlist();
33             if (token.getType() != rp) {
34                 error();
35             }
36             nextToken();
37             break;
38         default:
39             error();
40     }
41 }
42
43 void parseTermlist() {
44     parseTerm();
45     parseTermlist';
46 }
47
48 void parseTermlist'() {
49     switch (token.getType()) {
50         case rp:
51             break;
52         case comma:
53             nextToken();
54             parseTerm();
55             parseTermlist';
56             break;
57         default:
58             error();
59     }
60 }
```

Relationen zwischen First und Follow (Teilmengen)

Aufgabe 11 (First₁- und Follow₁-Mengen)

[3 Punkte]

Gegeben sei die folgende Produktion einer Grammatik G :

$$S \rightarrow AB$$

Es gilt

$$\text{First}_1(A) \subseteq \text{First}_1(S).$$

Geben Sie zwei weitere gültige Teilmengenbeziehung zwischen den Mengen

$$\begin{array}{ccc} \text{First}_1(S) & \text{First}_1(A) & \text{First}_1(B) \\ \text{Follow}_1(S) & \text{Follow}_1(A) & \text{Follow}_1(B) \end{array}$$

an. Triviale Teilmengenbeziehungen, wie $\text{First}_1(S) \subseteq \text{First}_1(S)$, sind dabei ausgeschlossen.

Beachten Sie, dass die Grammatik G weitere Produktionen mit den Nichtterminalen S , A und B enthalten kann.

Beispiellösung:

$$\begin{array}{l} \text{First}_1(B) \subseteq \text{Follow}_1(A) \\ \text{Follow}_1(S) \subseteq \text{Follow}_1(B) \end{array}$$

SGML

Aufgabe

- (b) Entwickeln Sie für die folgende, linksfaktorierte SGML-Grammatik einen rekursiven Abstiegsparser mit AST-Aufbau in Pseudocode. [16 Punkte]

$$\begin{aligned} SGML &\rightarrow < \mathbf{ident} > ChildrenAndEnd \\ ChildrenAndEnd &\rightarrow < OpenOrClose \\ OpenOrClose &\rightarrow / > \\ &| \mathbf{ident} > ChildrenAndEnd ChildrenAndEnd \end{aligned}$$

Der Parser von *SGML* soll ein Objekt der Klasse *SGML* zurückgeben. Diese ist folgendermaßen definiert:

```
1  class SGML {
2      public String tag;
3      public List<SGML> children;
4      public SGML(String tag, List<SGML> children) { ... }
5  }
```

Hinweis zum AST-Aufbau:

In der Produktion $OpenOrClose \rightarrow \mathbf{ident} > ChildrenAndEnd ChildrenAndEnd$ parst das erste Vorkommen von *ChildrenAndEnd* alle Kinder des gerade geöffneten Tags, das zweite Vorkommen parst die restlichen Kinder des umgebenden Tags.

Lexer-Schnittstelle:

Die globale Variable `token` enthält immer das aktuelle Token. Tokens besitzen die folgenden Methoden:

- `getType()` gibt den Token-Typ zurück
- `getIdent()` gibt einen String zurück, der den Namen eines **ident**-Tokens enthält.

Folgende Token-Typen sind definiert:

IDENT	für einen Bezeichner ident
LT	für das Kleiner-Zeichen <
GT	für das Größer-Zeichen >
SLASH	für den Schrägstrich /

Die globale Methode `nextToken()` setzt `token` auf das nächste Token. Brechen Sie bei Parsefehlern durch Aufruf der globalen `error()`-Methode ohne Fehlermeldung ab.

Parser

```
void expect(TokenType tt) {
    if (token.getType() != tt) {
        error();
    }
    nextToken();
}

SGML parseSGML() {
    expect(LT);

    if (token.getType() != IDENT) {
        error();
    }
    String tag = token.getIdent();
    nextToken();

    expect(GT);

    List<SGML> children = parseChildrenAndEnd();

    return new SGML(tag, children);
}

List<SGML> parseChildrenAndEnd() {
    expect(LT);

    return parseOpenOrClose();
}

List<SGML> parseOpenOrClose() {
    switch (token.getType()) {
        case SLASH:
            nextToken();

            expect(GT);

            return new ArrayList<SGML>();
        case IDENT:
            String tag = token.getIdent();
            nextToken();

            expect(GT);

            List<SGML> children = parseChildrenAndEnd();

            SGML elementHere = new SGML(tag, children);

            List<SGML> siblings = parseChildrenAndEnd();
            siblings.add(0, elementHere);

            return siblings;

        default:
            error();
            return null;
    }
}
```