

UNISUL - UNIVERSIDADE DO SUL DE SANTA CATARINA
DEPARTAMENTO DE CIÊNCIAS TECNOLÓGICAS
CURSO DE CIÊNCIA DA COMPUTAÇÃO
Prof. Roberto M. Scheffel

Apostila

Compiladores I

Sumário

SUMÁRIO	I
INTRODUÇÃO.....	1
I - TEORIA DA COMPUTAÇÃO.....	2
1. INTRODUÇÃO.....	2
2. ALGORITMOS E DECIDIBILIDADE.....	2
3. INTRODUÇÃO A COMPILADORES	3
3.1 Processadores de Linguagem.....	3
3.2. Modelo Simples de Compilador.....	4
II - TEORIA DE LINGUAGENS.....	6
1. LINGUAGENS FORMAIS.....	6
1.1. Símbolo.....	6
1.2. Cadeia	6
1.3. Linguagens	7
2. GRAMÁTICAS	7
2.1. Definição Formal de Gramática	8
2.2. Derivação	8
2.3. Notação	9
3. LINGUAGENS DEFINIDAS POR GRAMÁTICAS	10
4. TIPOS DE GRAMÁTICAS	10
5. TIPOS DE LINGUAGENS.....	11
EXERCÍCIOS	11
III - LINGUAGENS REGULARES E AUTÔMATOS FINITOS.....	13
1. LINGUAGENS REGULARES	13
1.1. Operações de Concatenação e Fechamento.....	13
1.2. Definição de Expressões Regulares.....	13
2. AUTÔMATOS FINITOS	14
2.1 Autômatos Finitos Determinísticos.....	15
2.2. Autômatos Finitos Não-Determinísticos.....	17
2.3. Comparação entre AFD e AFND	18
2.4. Transformação de AFND para AFD	18
2.5. Autômato Finito com ϵ -Transições.....	19
3. RELAÇÃO ENTRE GR E AF	20
4. MINIMIZAÇÃO DE AUTÔMATOS FINITOS	22
4.1. Algoritmo para Construção das Classes de Equivalência.....	22
4.2. Algoritmo para Construção do Autômato Finito Mínimo	23
5. CONSTRUÇÃO DO ANALISADOR LÉXICO	25
EXERCÍCIOS:	28
IV - LINGUAGENS LIVRES DE CONTEXTO.....	30
1. LINGUAGENS AUTO-EMBEBIDAS.....	30
2. GRAMÁTICAS LIVRES DE CONTEXTO	30
3. ÁRVORE DE DERIVAÇÃO	30
4. DERIVAÇÃO MAIS À ESQUERDA E MAIS À DIREITA	31
5. AMBIGÜIDADE.....	32
6. SIMPLIFICAÇÕES DE GRAMÁTICAS LIVRES DE CONTEXTO	32
6.1. Símbolos Inúteis.....	32
6.2. Eliminação de ϵ -Produções.....	34
6.3. Produções Unitárias.....	36
7. FATORAÇÃO	37
8. ELIMINAÇÃO DE RECURSÃO À ESQUERDA.....	38
9. TIPOS ESPECIAIS DE GLC.....	39
10 PRINCIPAIS NOTAÇÕES DE GLC	40

11. CONJUNTOS FIRST E FOLLOW	41
11.1 Conjunto First	41
11.2. Conjunto Follow	42
EXERCÍCIOS	44
V - AUTÔMATOS DE PILHA	46
1. INTRODUÇÃO.....	46
2. AUTÔMATOS DE PILHA	46
2.1. Definição Formal	47
3. AUTÔMATOS DE PILHA DETERMINÍSTICOS	49
EXERCÍCIOS	50
VI - ANÁLISE SINTÁTICA	51
1. INTRODUÇÃO.....	51
2. CLASSES DE ANALISADORES SINTÁTICOS	51
3. ANALISADORES ASCENDENTES (FAMÍLIA LR)	51
3.1. Estrutura dos Analisadores LR.....	51
3.2. Algoritmo de Análise Sintática LR	52
3.3. A Tabela de Parsing LR.....	52
3.4. A Configuração de um Analisador LR.....	53
3.5. Gramáticas LR(0).....	53
3.6. Itens LR.....	53
3.7. Cálculo dos Conjuntos de Itens Válidos.....	54
3.8. Definição de uma Gramática LR(0)	56
3.9. Construção do Conjunto LR.....	56
3.10. Construção da Tabela de Parsing SLR(1).....	57
4. ANALISADORES DESCENDENTES	58
4.1. Analisadores Descendentes sem Back-Tracking	58
EXERCÍCIOS	61
RESPOSTAS AOS EXERCÍCIOS PROPOSTOS	62
REFERÊNCIAS BIBLIOGRÁFICAS	71

Introdução

Este texto provê os elementos essenciais da disciplina de “**Compiladores I**”. O capítulo I situa a área de linguagens formais no contexto da teoria da computação. Em seguida trata de alguns aspectos relacionados à construção de compiladores.

O capítulo II apresenta algumas definições básicas de linguagens formais e uma breve introdução às gramáticas gerativas, que serão base para o restante da apostila.

O capítulo III apresenta a teoria de linguagens regulares e das máquinas de estados finitos, ou autômatos finitos, que serão utilizados como base para a construção de analisadores léxicos.

O capítulo IV apresenta a teoria específica de gramáticas livres de contexto, que são importantes ferramentas para a construção de analisadores sintáticos de compiladores.

O capítulo V apresenta alguns conceitos relativos a autômatos de pilha, que são utilizados para analisar sentenças geradas pelas gramáticas livres de contexto.

Finalmente, o capítulo VI apresenta o processo de análise sintática propriamente dito.

A partir do capítulo II, ao final de cada capítulo são propostos alguns exercícios de fixação. O exercício cujo enunciado estiver marcado com um asterisco (*) tem sua resolução no final da apostila.

O material apresentado nesta apostila cobre o conteúdo básico da disciplina. Maiores subsídios podem ser encontrados na bibliografia indicada no final da apostila.

I - Teoria da Computação

1. Introdução

Teoria de Linguagens Formais e Teoria de Máquina são tópicos abrangentes que se inserem no estudo da Teoria da Computação em geral. A Teoria da Computação é uma ciência que procura organizar o conhecimento formal relativo aos processo de computação, como complexidade de algoritmos, linguagens formais, problemas intratáveis, etc.

Atualmente, o termo “computar” está associado ao conceito de fazer cálculos ou aritmética, usando para isso máquinas computadoras. Entretanto, a noção de computabilidade pode ser dissociada da implementação física da máquina. Originalmente, a palavra latina *putare* significa “pensar”, mas no caso da teoria da computação o sentido mais adequado seria algo como “manipulação de símbolos”, o que não é, necessariamente, pensamento. Essa manipulação de símbolos envolvendo, principalmente, letras, números e proposições algébricas permitiria às máquinas computadoras, segundo Alan Turing (1939) realizar qualquer operação formal de que o ser humano seria capaz.

Mesmo assim, alguns filósofos sustentam que os computadores não podem computar (fazer aritmética), porque não compreendem a noção de número (apenas manipulam sinais elétricos). Todavia, pode-se falar, sem entrar no mérito da questão, que os computadores realizam uma “computação perceptível”, já que transformam uma entrada em saída como se estivessem realmente computando. Neste caso, o que importa é o efeito final do processo e não a maneira como ele é feito.

2. Algoritmos e Decidibilidade

Os algoritmos surgiram na história da humanidade possivelmente com as tábuas de logaritmos e de números primos, na antigüidade. Certos algoritmos também eram usados para calcular a posição dos planetas e outros auxiliavam a navegação. Pode-se dizer que os computadores surgiram graças aos algoritmos. Por exemplo: Charles Babbage (1820), considerado o “pai da Computação”, irritou-se com a maneira mecânica e repetitiva com que os cálculos da posição de estrelas e planetas eram feitos na Royal Astronomical Society, onde trabalhava, e escreveu o artigo “Observations on the Application of Machinery to the Computatiuon of Mathematical Tables”, onde propunha que estes algoritmos fossem realizados por máquinas.

Mas foi apenas em 1900 que o matemático Frege chegou à conclusão de que se precisava de uma noção precisa e concisa do conceito de algoritmo. Frege procurou usar a lógica como ferramenta para estabelecer os principais algoritmos de matemática. A notação utilizada realmente era precisa mas verificou-se pouco depois que o trabalho continha uma contradição, que ficou conhecida como paradoxo de Russel.

O Paradoxo de Russel

Muitos problemas até hoje não possuem solução algorítmica, isto é, não podem ser resolvidos por máquinas. Destes problemas, alguns não têm solução conhecida, mas existe um grupo de problemas dos quais sabe-se que não existem algoritmos formais que possam resolvê-los. Em muitos casos, o problema é apresentado como uma simples questão, que deve ser respondida com *sim* ou *não*. Posto o problema desta forma, se existe um algoritmo que sempre responda *sim* ou *não* para qualquer entrada, então o problema é decidível. Em alguns casos, porém, provou-se que não existem nem poderão existir estes algoritmos. Tais problemas são chamados, então, de indecidíveis, ou parcialmente decidíveis, se pelo menos uma das respostas *sim* ou *não* puder ser encontrada.

Gödel, em 1931, provou que a aritmética elementar é indecidível, ou seja, não existe um algoritmo para provar a verdade ou falsidade de qualquer proposição da aritmética. Hilbert, em 1925, propôs vários problemas que se supunha fossem indecidíveis, dentre eles, por exemplo, o de

dizer se o polinômio $P(x_1, x_2, \dots, x_n) = 0$ tem solução inteira. Apenas em 1970 Matijacevic provou que não existe nenhum procedimento efetivo para decidir esta questão.

Em termos de linguagem, pode-se dizer que uma linguagem é decidível se dada uma linguagem L e uma sentença w for possível decidir algoritmicamente se w é uma sentença de L . Para as linguagens não-decidíveis, pode-se esperar que para algum w , a pergunta “ $w \in L?$ ” não será respondida, porque o algoritmo que procura pela resposta não parará.

Um algoritmo é uma descrição finita de um processo de computação; é um texto finito, escrito em um linguagem algorítmica na qual cada sentença tem um significado não-ambíguo.

A palavra *algoritmo* é uma corruptela da *Al-Khuwarizmi*, sobrenome de um matemático persa do século IX. O algoritmo deve poder ser resolvido em tempo finito (deve ter parada garantida). Um *procedimento* não termina necessariamente. Assim, existem problemas que possuem procedimentos para serem resolvidos, mas não possuem algoritmos. Em outras palavras, existem um procedimento de procura de solução, que pode chegar a uma resposta em um tempo finito, ou então pode não chegar a uma resposta em momento algum, ao contrário do algoritmo, que sempre chega à uma resposta (afirmativa ou negativa) num tempo finito.

Decidir se um problema possui um algoritmo ou simplesmente um procedimento é o principal problema indecidível da teoria da computação, conhecido como “Problema da Parada”. Outra forma de enunciar o problema da parada é a seguinte: “Dado um programa P e um conjunto de dados d , responda sim se $P(d)$ pára e não caso contrário”.

A Tese de Church

Alguns dos mais importantes modelos formais de computabilidade existentes são a Máquina de Turing, as gramáticas do Tipo 0 e as Funções Recursivas Parciais. A Tese de Church (1936) estabelece que: “Embora todas estas caracterizações sejam diferentes na forma, elas caracterizam o mesmo conjunto de funções”. No caso, Church provou que era possível definir funções de tradução de representações de um formalismo para qualquer um dos outros, donde se concluiu que tudo o que é expressável em um deles pode ser expresso em qualquer um dos outros. Isto garante que todos os modelos formais de computação têm o mesmo poder expressivo.

Como as Funções Recursivas Parciais são suficientemente inclusivas em relação à computação aparente, e a caracterização da Máquina de Turing demonstra o aspecto mecânico do processamento algorítmico, estes formalismos são aceitos como capazes de expressar qualquer função ou algoritmo que possa ser formalmente expresso de alguma maneira.

3. Introdução a Compiladores

3.1 Processadores de Linguagem

Um processador de linguagem é um programa que permite ao computador “entender” os comandos de alto nível fornecidos pelos usuários.

Existem duas espécies principais de processadores de linguagem: os *interpretadores* e os *tradutores*.

Um interpretador é um programa que aceita como entrada um programa escrito em uma linguagem chamada *linguagem fonte* e executa diretamente as instruções dadas nesta linguagem.

Um tradutor é um programa que aceita como entrada um programa escrito em uma linguagem fonte e produz como saída uma outra versão do programa escrito em uma *linguagem objeto*. Muitas vezes a linguagem objeto é a própria linguagem de máquina do computador. Neste caso o programa objeto pode ser diretamente executado pela máquina.

Tradutores são arbitrariamente divididos em montadores e compiladores, os quais traduzem linguagens de baixo nível e de alto nível, respectivamente. Além disso, há os *pré-processadores* que traduzem uma linguagem de alto nível em outra linguagem de alto nível.

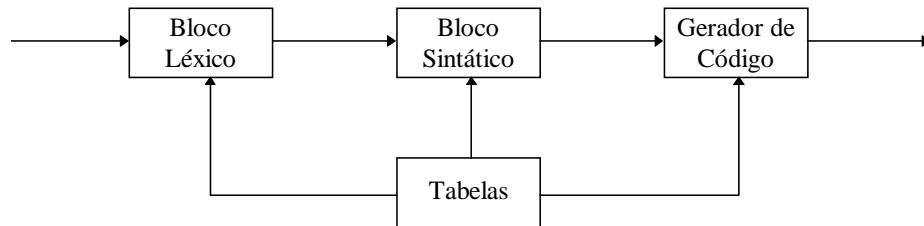
O fundamento matemático da teoria de processamento de linguagens é a teoria de autômatos e linguagens formais.

3.2. Modelo Simples de Compilador

O objetivo de um compilador é traduzir as seqüências de caracteres que representam o programa fonte em instruções de máquina que incorporam a intenção do programador. Esta tarefa é complexa o suficiente para poder ser vista como uma estrutura formada por processos menores interconectados.

Em um modelo simples de compilador o trabalho pode ser feito por uma interconexão serial entre três blocos, chamados: *bloco léxico*, *bloco sintático* e *gerador de código*.

Os três blocos têm acesso a tabelas de informações globais sobre o programa fonte. Uma destas tabelas é, por exemplo, a *tabela de símbolos*, na qual a informação sobre cada variável ou identificador utilizado no programa fonte é armazenada. A conexão entre estes blocos e tabelas é mostrada na figura seguinte:



3.2.1. Bloco Léxico

A entrada para um compilador é a cadeia de caracteres que representa o programa fonte. O objetivo do bloco léxico é quebrar esta cadeia de caracteres separando as palavras que nela estão representadas. Por exemplo, a cadeia de entrada poderia ser:

```
RESULTADO := RESULTADO + 15
```

Neste caso, o bloco léxico deve ser capaz de discernir que esta cadeia contém o nome de uma variável `RESULTADO`, seguido de um comando de atribuição `:=`, seguido novamente de `RESULTADO`, seguido de uma operação aritmética `+` e finalizado por uma constante numérica `15`. Assim, os 23 caracteres da cadeia de entrada são transformados em 5 novas entidades. Estas entidades são freqüentemente denominadas de *tokens*.

Um token consiste de uma par ordenado *Classe x Valor*. A *classe* indica que o token pertence a uma classe de um conjunto finito de classes possíveis e indica a natureza da informação contida em *valor*.

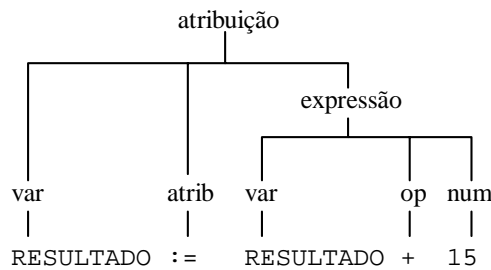
Voltando ao exemplo, a variável `RESULTADO` poderia pertencer à classe *variável* e seu valor seria um ponteiro para a entrada `RESULTADO` na tabela de símbolos. O token `15` poderia pertencer à classe *constante* e ter valor correspondente ao número inteiro `15`. Já os tokens `:=` e `+` poderiam pertencer às classes *símbolos especiais*.

Se a tabela de símbolos for vista como um dicionário, então o processo léxico pode ser visto como o ato de agrupar letras em palavras e determinar suas localizações no dicionário.

Outras funções normalmente atribuídas ao bloco léxico são as de ignorar espaços em branco e comentários, além de detectar os erros léxicos.

3.2.2. Bloco Sintático

O bloco sintático agrupa os tokens fornecidos pelo bloco léxico em estruturas sintáticas, verificando se a sintaxe da linguagem foi obedecida. Tomando novamente o exemplo anterior, os 5 tokens gerados pelo bloco léxico poderiam ser agrupados na estrutura sintática *atribuição*, representada por uma *árvore sintática*, da seguinte forma:



Note que a árvore exige que para a atribuição apareça uma variável, seguida de uma atribuição e de uma expressão. Desta forma, uma entrada do tipo

`15 := RESULTADO + 10`

passaria no bloco léxico, pois as tokens estão escritas da *maneira correta*, porém não passaria no bloco sintático, pois não estão escritas na *ordem correta*.

3.2.3. Gerador de Código

A função do gerador de código é expandir os átomos gerados pelo bloco sintático em uma seqüência de instruções do computador que executam a intenção do programa fonte. Consiste basicamente no *processamento semântico* e na *otimização de código*.

Processamento Semântico

Certas atividades relacionadas ao *significado* dos símbolos são freqüentemente classificadas como processamento semântico. Por exemplo, a semântica de um identificador pode incluir o seu tipo, e, se for um array, por exemplo, seu tamanho. O processamento também pode incluir o preenchimento da tabela de símbolos com as propriedades dos identificadores quando estas se tornam conhecidas (quando são declaradas ou implícitas).

Na maioria dos compiladores, as atividades semânticas são realizadas por um bloco separado, denominado *bloco semântico*, que fica entre o bloco sintático e o gerador de código. Algumas operações semânticas usuais são: verificação de compatibilidade de tipo, análise de escopo, verificação de compatibilidade entre declaração e uso de entidades, verificação de correspondência entre parâmetros atuais e formais e verificação de referência não resolvidas.

Otimização de Código

Certas atividades de um compilador que não são estritamente necessárias, mas que permitem obter melhores programas objeto, são freqüentemente referidas como *otimização*. Para alguns compiladores, a otimização é inserida como um bloco entre o bloco sintático e o semântico (se houver). Este bloco pode resolver problemas como, por exemplo, instruções invariantes em laços, rearranjando os símbolos sintáticos antes de passá-los para a análise semântica.

II - Teoria de Linguagens

1. Linguagens Formais

O estudo dos diversos tipos de linguagens aplicam-se a diversas áreas da informática, sempre que for necessário analisar o formato e o significado de uma seqüência de entrada. Dentre as aplicações possíveis, temos: **editor de estruturas**, que analisa as estruturas de um programa fonte sendo editado, tentando minimizar erros; **pretty printers**, que tenta tornar os programas mais legíveis; **verificadores ortográficos**, **verificadores gramaticais**, **verificadores estáticos**, **interpretadores**, **compiladores**, etc.

Serão apresentados nesta seção conceitos gerais sobre linguagens que servirão para fundamentar o estudo de todos os tipos de linguagens que virão a seguir.

1.1. Símbolo

Um símbolo é uma entidade abstrata que não precisa ser definida formalmente, assim como “ponto” e “linha” não são definidos na geometria. Letras e dígitos são exemplos de símbolos freqüentemente usados.

Símbolos são ordenáveis lexicograficamente e, portanto, podem ser comparados quanto à igualdade ou precedência. Por exemplo, tomando as letras do alfabeto, podemos ter a ordenação $a < b < c < \dots < z$. A principal utilidade dos símbolos está na possibilidade de usá-los como elementos atômicos em definições de linguagens.

1.2. Cadeia

Uma *cadeia* (ou string, ou palavra) é uma seqüência finita de símbolos justapostos (isto é, sem vírgulas separando os caracteres). Por exemplo, se a, b e c são símbolos, então abcb é um exemplo de cadeia usando utilizando estes símbolos.

O *tamanho* de uma cadeia é o comprimento da seqüência de símbolos que a forma. O tamanho de uma cadeia w será denotado por $|w|$. Por exemplo, $|abcb| = 4$. A *cadeia vazia* é denotada por ϵ , e tem tamanho igual a 0; assim, $|\epsilon| = 0$.

Um *prefixo* de uma cadeia é um número qualquer de símbolos tomados de seu início, e um *sufixo* é um número qualquer de símbolos tomados de seu fim. Por exemplo, a cadeia abc tem prefixos ϵ , a, ab, e abc. Seus sufixos são ϵ , c, bc, abc. Um prefixo ou sufixo que não são a própria cadeia são chamados de *prefixo próprio* e *sufixo próprio*, respectivamente. No exemplo anterior, os prefixos próprios são ϵ , a e ab; e os sufixos próprios são ϵ , c e bc.

A *concatenação* de duas cadeias é a cadeia formada pela escrita da primeira cadeia seguida da segunda, sem nenhum espaço no meio. Por exemplo, a concatenação de *compila* e *dores* é *compiladores*. O operador de concatenação é a justaposição. Isto é, se w e x são variáveis que denotam cadeias, então wx é a cadeia formada pela concatenação de w e x . No exemplo acima, se tomarmos $w = \textit{compila}$ e $x = \textit{dores}$, então temos $wx = \textit{compiladores}$.

A *n-ésima potência* de uma cadeia x , denotada por x^n , é a concatenação de x com ela mesma $n-1$ vezes, ou a repetição da cadeia n vezes. Por exemplo, $abc^3 = abcabcabc$.

Um *alfabeto* é definido simplesmente como um conjunto finito de símbolos. Por exemplo, o alfabeto da língua portuguesa é $V = \{a, b, c, \dots, z\}$. O alfabeto utilizado para expressar os números naturais é $V = \{1, 2, 3, 4, 5, 6, 7, 8, 9, 0\}$.

O *fechamento* de um alfabeto V , representado por V^* , é o conjunto de todas as cadeias que podem ser formadas com os símbolos de V , inclusive a cadeia vazia. O *fechamento positivo* de V , denotado por V^+ é definido como $V^* - \{\epsilon\}$. Ou seja, todas as cadeias formadas com os símbolos de V , exceto a cadeia vazia. Por exemplo, dado o alfabeto $V = \{0, 1\}$, então o fechamento de V é dado

por $V^* = \{\epsilon, 0, 1, 00, 01, 10, 11, 000, 001, \dots\}$ e o fechamento positivo de V é dado por $V^+ = V^* - \{\epsilon\} = \{0, 1, 00, 01, 10, 11, 000, 001, \dots\}$

1.3. Linguagens

Uma *linguagem* (formal) é um conjunto de cadeias de símbolos tomados de algum alfabeto. Isto é, uma linguagem sobre o alfabeto V é um subconjunto de V^* . Note-se que esta definição é meramente extensional, isto é, não considera os mecanismos formadores da linguagem, mas apenas a sua extensão. Assim, por exemplo, o conjunto de sentenças válidas da língua portuguesa poderia ser definido extensionalmente como um subconjunto de $\{a, b, c, \dots, z\}^+$.

Uma linguagem é *finita* se suas sentenças formam um conjunto finito. Caso contrário, a linguagem é *infinita*. Uma linguagem infinita precisa ser definida através de uma representação finita. Por exemplo, a linguagem dos números naturais menores que 10 é finita, e pode ser representada, literalmente, como $L = \{1, 2, 3, 4, 5, 6, 7, 8, 9\}$. Já a linguagem dos números naturais como um todo não é uma linguagem finita, já que existem infinitos números naturais. Porém, como veremos adiante, existem mecanismos que permitem expressar esta e outras linguagens infinitas através de representações finitas. É o caso das gramáticas e das expressões regulares.

Um *reconhecedor* para uma linguagem é um dispositivo formal usado para verificar se uma determinada sentença pertence ou não a uma determinada linguagem. São exemplos de reconhecedores de linguagens as Máquinas de Turing, o Autômato de Pilha e o Autômato Finito. Cada um destes mecanismos reconhece um conjunto bastante particular de linguagens. Entretanto, os modelos são inclusivos: todas as linguagens reconhecíveis por um Autômato Finito podem ser reconhecidas por uma Autômato de Pilha e todas as linguagens reconhecíveis por um Autômato de Pilha são reconhecíveis por Máquinas de Turing. As recíprocas, entretanto, não são verdadeiras.

Um *sistema gerador* é um dispositivo formal através do qual as sentenças de uma linguagem podem ser sistematicamente geradas. Exemplos de sistemas geradores são as gramáticas gerativas, definidas pela primeira vez por Noam Chomsky, em seus estudos para sistematizar a gramática da língua inglesa.

Todo reconhecedor e todo sistema gerador pode ser representado por algoritmos e/ou procedimentos. Como será visto mais adiante, pode-se garantir que linguagens reconhecidas por autômatos são reconhecidas algorítmicamente. Já o conjunto das linguagens reconhecidas por Máquinas de Turing, que é mais abrangente, inclui linguagens reconhecíveis algorítmicamente e linguagens indecidíveis, que não possuem algoritmos para reconhecimento.

2. Gramáticas

Uma *gramática gerativa* é um instrumento formal capaz de construir (gerar) conjuntos de cadeias de uma determinada linguagem. As gramáticas são instrumentos que facilitam muito a definição das características sintáticas das linguagens. São instrumentos que permitem definir, de forma formal e sistemática, uma representação finita para linguagens infinitas.

Usar um método de definição de conjuntos particular para definir uma linguagem pode ser um passo importante no projeto de um reconhecedor para a linguagem, principalmente se existirem métodos sistemáticos para converter a descrição do conjunto em um programa que processe o conjunto. Como veremos adiante neste curso, certos tipos de gramáticas, que possuem algoritmos para definir se uma sentença pertence ou não à linguagem que geram, são utilizadas para implementar compiladores. Já gramáticas que não permitem a definição de algoritmos não podem ser utilizadas para tal.

Mais adiante serão vistos métodos para transformar gramáticas em máquinas abstratas capazes de reconhecer e traduzir os conjuntos definidos pelas gramáticas. Antes, porém, é necessário entender como conjuntos de sentenças são definidos por gramáticas.

2.1. Definição Formal de Gramática

Uma gramática possui um conjunto finito de *variáveis* (também chamadas de *não-terminais*). A linguagem gerada por uma gramática é definida recursivamente em termos das variáveis e de *símbolos primitivos* chamados *terminais*, que pertencem ao alfabeto da linguagem.

As regras que relacionam variáveis e terminais são chamadas de *produções*, que sempre iniciam num ponto, o *símbolo inicial*, pertencente ao conjunto de variáveis (não-terminais). Uma regra de produção típica estabelece como uma determinada configuração de variáveis e terminais pode ser reescrita, gerando uma nova configuração.

Formalmente, uma gramática é especificada por uma quádrupla (N, T, P, S) , onde:

- N é o conjunto de não-terminais, ou variáveis.
- T é o conjunto de terminais, ou alfabeto, sendo que N e T são conjuntos disjuntos.
- P é um conjunto finito de produções, ou regras, sendo que $P \subseteq (N \cup T)^+ \times (N \cup T)^*$.
- S é um símbolo não-terminal inicial, sendo $S \in N$. A partir dele são geradas todas as sentenças da linguagem.

Uma produção (n, p) pode ser escrita $n ::= p$, para facilitar a leitura. No caso de existir mais de uma alternativa para uma mesma configuração e símbolos do lado esquerdo da regra, como por exemplo:

$\langle \text{substantivo} \rangle ::= \text{coisa}$
 $\langle \text{substantivo} \rangle ::= \text{cão}$
 $\langle \text{substantivo} \rangle ::= \text{pedra}$

pode-se escrever as opções em uma única produção, separada por $|$. O exemplo acima ficaria:

$\langle \text{substantivo} \rangle ::= \text{coisa} \mid \text{cão} \mid \text{pedra}$.

O símbolo $::=$ pode ser lido como “é definido por” e o símbolo $|$ pode ser lido como “ou”.

2.2. Derivação

Uma gramática pode ser usada para gerar uma linguagem através de *reescrita* ou *derivação*. A derivação pode ser definida da seguinte maneira: dada uma sentença a_1ba_2 , com a_1 e a_2 sendo cadeias de terminais e não-terminais, se existir uma regra da forma $b ::= c$, então a sentença inicial pode ser reescrita como a_1ca_2 . A derivação, ou reescrita, é denotada da seguinte forma:

$a_1ba_2 \rightarrow a_1ca_2$

Por exemplo, digamos que exista a seguinte seqüência de terminais (letras minúsculas) e não-terminais (letras maiúsculas): “aaaDbbbb”, e que exista a produção “ $D ::= ab$ ” na gramática. Assim, a derivação “ $aaaDbbbb \rightarrow aaaabbbb$ ” é uma derivação válida na gramática dada.

Se $a \rightarrow b$, então se diz que a produz diretamente b . O fecho transitivo e reflexivo da relação \rightarrow é \rightarrow^* . Portanto, se $a \rightarrow b \rightarrow \dots \rightarrow z$, então se diz que a produz z , simplesmente, anotando a produção por $a \rightarrow^* z$. O fecho transitivo pode ser definido da seguinte forma:

- $a \rightarrow^0 a$
- $a \rightarrow^1 b$ se $a \rightarrow b$
- $a \rightarrow^{n+1} z$ se $a \rightarrow^n y$ e $y \rightarrow^1 z$
- $a \rightarrow^* z$ se $a \rightarrow^n z$, para algum $n \geq 0$.

Seja, por exemplo, a seguinte gramática, onde as regras são numeradas para acompanhar o desenvolvimento da derivação:

- 1: $\langle S \rangle ::= \langle X \rangle + \langle X \rangle =$
- 2,3: $\langle X \rangle ::= 1 \langle X \rangle \mid 1$
- 4: $\langle X \rangle + 1 ::= \langle X \rangle 1 +$
- 5: $\langle X \rangle 1 ::= 1 \langle X \rangle$
- 6: $\langle X \rangle + = ::= \langle X \rangle$
- 7: $1 \langle X \rangle ::= \langle X \rangle 1$

A derivação de uma sentença a partir desta gramática consiste em uma seqüência de aplicações das produções a partir do símbolo inicial. Por exemplo:

$\langle S \rangle \rightarrow \langle X \rangle + \langle X \rangle =$ [1]
 $\rightarrow 1 \langle X \rangle + \langle X \rangle =$ [2 no primeiro $\langle X \rangle$]
 $\rightarrow 1 \langle X \rangle + 1 =$ [3 no segundo $\langle X \rangle$]
 $\rightarrow 1 \langle X \rangle 1 + =$ [4]
 $\rightarrow 11 \langle X \rangle + =$ [5]
 $\rightarrow 11 \langle X \rangle$ [6]
 $\rightarrow 111$ [3]

O fato das regras terem sido aplicadas praticamente em seqüência foi mera coincidência. A princípio, qualquer regra de derivação pode ser aplicada a qualquer ponto de sentença desde que a configuração de símbolos do lado esquerdo da regra apareça na sentença sendo derivada.

Outro exemplo é a seguinte gramática:

$\langle S \rangle ::= 1 \langle A \rangle \mid 0 \langle S \rangle \mid 1$
 $\langle A \rangle ::= 0 \langle A \rangle \mid 1 \langle S \rangle$

A sentença 100101 é gerada pela seguinte seqüência de derivações:

$\langle S \rangle \rightarrow 1 \langle A \rangle \rightarrow 10 \langle A \rangle \rightarrow 100 \langle A \rangle \rightarrow 1001 \langle S \rangle \rightarrow 10010 \langle S \rangle \rightarrow 100101$

Originalmente, a motivação do uso de gramáticas foi o de descrever estruturas sintáticas da língua natural. Poderia-se escrever regras (produções) como:

$\langle \text{frase} \rangle ::= \langle \text{artigo} \rangle \langle \text{substantivo} \rangle \langle \text{verbo} \rangle \langle \text{artigo} \rangle \langle \text{substantivo} \rangle. \mid$
 $\langle \text{artigo} \rangle \langle \text{substantivo} \rangle \langle \text{verbo} \rangle.$
 $\langle \text{verbo} \rangle ::= \text{dorme} \mid \text{escuta} \mid \text{ama}$
 $\langle \text{substantivo} \rangle ::= \text{gato} \mid \text{rádio} \mid \text{rato}$
 $\langle \text{artigo} \rangle ::= \text{o}$

Desta gramática simples pode-se obter derivações como:

$\langle \text{frase} \rangle \rightarrow \text{o} \langle \text{substantivo} \rangle \langle \text{verbo} \rangle \langle \text{artigo} \rangle \langle \text{substantivo} \rangle.$
 $\rightarrow \text{o gato} \langle \text{verbo} \rangle \langle \text{artigo} \rangle \langle \text{substantivo} \rangle.$
 $\rightarrow \text{o gato ama} \langle \text{artigo} \rangle \langle \text{substantivo} \rangle.$
 $\rightarrow \text{o gato ama o} \langle \text{substantivo} \rangle.$
 $\rightarrow \text{o gato ama o rato.}$

Entretanto a mesma gramática também permite derivar “o gato dorme o rato”, frase que não faz sentido. Para evitar a geração de tais frases, pode-se usar uma gramática sensível ao contexto, estabelecendo que $\langle \text{verbo} \rangle$ só pode ser reescrito para “dorme” se for sucedido pelo ponto final da frase:

$\langle \text{frase} \rangle ::= \langle \text{artigo} \rangle \langle \text{substantivo} \rangle \langle \text{verbo} \rangle \langle \text{artigo} \rangle \langle \text{substantivo} \rangle. \mid$
 $\langle \text{artigo} \rangle \langle \text{substantivo} \rangle \langle \text{verbo} \rangle.$
 $\langle \text{verbo} \rangle. ::= \text{dorme.}$
 $\langle \text{verbo} \rangle \langle \text{artigo} \rangle ::= \text{escuta} \langle \text{artigo} \rangle \mid \text{ama} \langle \text{artigo} \rangle$
 $\langle \text{substantivo} \rangle ::= \text{gato} \mid \text{rádio} \mid \text{rato}$
 $\langle \text{artigo} \rangle ::= \text{o}$

Para esta gramática, é impossível derivar as sentenças “o gato dorme o rato”, “o gato ama” e “o rato escuta”. O verbo dorme só pode ser derivado se for sucedido por um ponto final e os verbos escuta e ama só podem ser derivados se forem sucedidos por um artigo. Porém ainda é possível derivar uma sentença do tipo “o rádio ama o rato”. O leitor é convidado a modificar a gramática acima para não aceitar sentenças desse tipo (dica: permitir que o substantivo “rádio” apareça somente no final da sentença).

2.3. Notação

Como convenção notacional, pode-se admitir que símbolos não terminais serão sempre representados por letras maiúsculas, e terminais por minúsculas. Assim, uma regra como

$\langle a \rangle ::= \langle a \rangle b \mid b$ poderia ser escrita como $A ::= Ab \mid b$. Além disso, os conjuntos T e N podem ficar subentendidos e não precisam ser expressos sempre.

Pode-se ainda arbitrar que o símbolo inicial S será sempre o não-terminal que aparecer primeiro do lado esquerdo da primeira produção da gramática. Assim, bastaria listar as regras de produção para se ter a definição completa da gramática. Por exemplo, a gramática

S ::= ABSc | ϵ
 BA ::= AB
 Bc ::= bc
 Ab ::= ab
 Bb ::= bb
 Aa ::= aa

Tem como símbolo inicial S, o conjunto de terminais é {a, b, c} e o conjunto de não-terminais (variáveis) é {A, B, S}. As produções estão detalhadas. Tem-se assim a definição completa da gramática.

3. Linguagens Definidas por Gramáticas

A linguagem definida por uma gramática consiste no conjunto de cadeias de terminais que esta gramática pode potencialmente gerar. Este conjunto será denotado por L(G). A linguagem definida por uma gramática $G=(N,T,P,S)$ é dada por:

$$L(G) = \{ \alpha \mid \alpha \in T^* \wedge S \rightarrow^* \alpha \}$$

Em outras palavras, uma cadeia pertence a L(G) se e somente se ela consiste somente de terminais (pode ser a cadeia vazia ϵ), e pode ser produzida a partir de S em 0 ou mais derivações.

Uma cadeia em $(T \cup N)^*$ é chamada de *forma sentencial* se ela pode ser produzida a partir de S (mas ainda contém não-terminais). Assim, pode-se caracterizar L(G) como o conjunto de todas as formas sentenciais que contém apenas terminais, e foram geradas a partir do símbolo inicial de G. Desta forma, “aaaBbbC” é uma forma sentencial, pois contém variáveis (B e C). Já “aabb” é uma sentença, pois contém apenas terminais.

4. Tipos de Gramáticas

Impondo restrições na forma das produções, pode-se identificar quatro tipos diferentes de gramáticas. Esta classificação foi feita por Chomsky e é conhecida como *hierarquia de Chomsky*. Os tipos de gramática definidos nesta hierarquia são:

a) Tipo 0: Não há restrição na forma das produções. Este é o tipo mais geral. Exemplo:

S ::= ABS | ab Ba ::= aB | a Aa ::= aa | a | ϵ

b) Tipo 1: Seja $G=(N,T,P,S)$. Se:

- O símbolo inicial S não aparece no lado direito de nenhuma produção, e
- para cada produção $\sigma_1 ::= \sigma_2$, é verdade que $|\sigma_1| \leq |\sigma_2|$ (com exceção da regra $S ::= \epsilon$),

então se diz que G é uma gramática do tipo 1, ou *Gramática Sensível ao Contexto*. Exemplo:

S ::= aBC | aABC A ::= aABC | aBC CB ::= BC aB ::= ab
 bB ::= bb bC ::= bc cC ::= cc

c) Tipo 2: Uma gramática $G=(N,T,P,S)$ é do tipo 2, ou uma *Gramática Livre de Contexto*, se cada produção é livre de contexto, ou seja, cada produção é da forma

$A ::= \sigma$, com $A \in N$ e $\sigma \in (T \cup N)^*$.

Exemplo:

S ::= aSb | A A ::= cAd | e

O tipo 2 não é definido como restrição do tipo 1, porque se permite produções da forma $A ::= \epsilon$, com $A \neq S$. Também se permite que o símbolo inicial S apareça no lado direito das produções.

Entretanto, existe um teorema que prova que a partir de uma gramática de tipo 2 se pode construir outra gramática equivalente, de tipo 2, que satisfaz as restrições do tipo 1.

d) Tipo 3. Uma gramática G é de tipo 3, ou regular, se cada produção é da forma

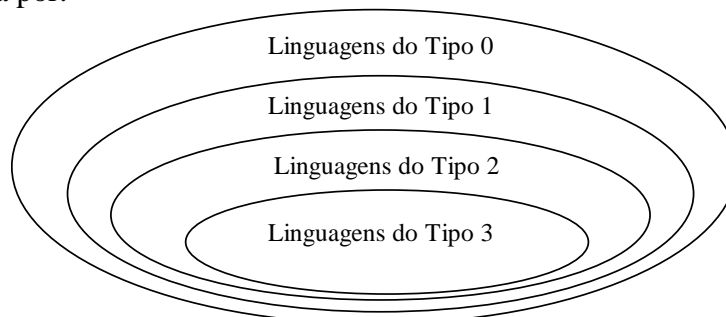
$A ::= aB$, $A ::= a$ ou $A ::= \epsilon$, onde A e B são não-terminais e a é um terminal.

Exemplo:

$S ::= aS \mid bB$ $B ::= bB \mid bC$ $C ::= cC \mid \epsilon$

5. Tipos de Linguagens

Uma linguagem L é de tipo i se existe uma gramática G de tipo i tal que $L = L(G)$, para i igual a 0, 1, 2 ou 3. Pode-se ver que cada linguagem de tipo 3 é também de tipo 2, cada linguagem de tipo 2 é também de tipo 1, e cada linguagem de tipo 1 é também de tipo 0. Estas inclusões são estritas, isto é, existem linguagens de tipo 0 que não são de tipo 1, existem linguagens de tipo 1 que não são do tipo 2 e existem linguagens do tipo 2 que não são do tipo 3. Graficamente esta inclusão pode ser representada por:



Pode-se relacionar cada tipo de gramática com uma máquina reconhecedora da seguinte maneira: a gramática de tipo 0 gera linguagens reconhecíveis por máquinas de Turing; a gramática de tipo 2 gera linguagens reconhecidas por autômatos de pilha; a gramática de tipo 3 gera linguagens reconhecidas por autômatos finitos.

Não há uma classe de reconhecedores para apenas linguagens de tipo 1, porque qualquer máquina capaz de reconhecer linguagens de tipo 1 é poderosa o suficiente para reconhecer linguagens do tipo 0. A principal razão para identificar as linguagens de tipo 1 como um tipo separado é porque toda a linguagem do tipo 1 é decidível, isto é, se $G=(N,T,P,S)$ é do tipo 1, então existe um algoritmo tal que, para qualquer sentença w , responda “sim” se $w \in L(G)$ e “não” caso contrário.

Assim, pode-se dizer que as linguagens do tipo 1 são reconhecíveis por Máquinas de Turing, sendo decidíveis, ou seja, a máquina sempre vai parar num tempo finito com uma resposta, afirmativa ou negativa, sobre a sentença analisada. Já as linguagens do tipo 0 também podem ser analisadas por Máquinas de Turing. O problema é que podem existir linguagens, ou sentenças dentro de linguagens, para as quais a máquina não pare. Ou seja, a máquina ficará analisando a sentença indefinidamente, sem chegar a nenhuma resposta.

EXERCÍCIOS

1. Uma palíndrome é uma cadeia que pode ser lida da mesma maneira da direita para a esquerda e da esquerda para a direita. Exemplos: a, ama, osso, asdds. Defina uma gramática para a linguagem das palíndromes (conjunto de todas as palíndromes sobre o alfabeto $\{a, b, c, \dots, z\}$).

* 2. Dada a gramática (livre de contexto):

$S ::= aSba \mid \epsilon$

mostre três cadeias de $L(G)$ e a derivação de cada uma delas.

3. Dado o conjunto de terminais (alfabeto) $T=\{a, b, c\}$, defina as gramáticas para as seguintes linguagens:

* a) Sentenças que começam e terminam com a.

- * b) Sentenças que tenham tamanho ímpar.
- c) Sentenças em que todos os a's apareçam consecutivos.
- * d) Sentenças na forma $1^n 0^n 2^n$, com $n > 0$.
- e) Sentenças na forma $a^n b^m c^n a^n$, $n > 0$ e $m > 0$.
- * f) Sentenças na forma $a^n b^m c^n$, com $n > 0$ e $m \geq 0$.
- g) Sentenças na forma $a^m b^n c^n$, com $m \geq 0$ e $n > 0$.
- h) Sentenças em que o número de a's seja par.
- i) Sentenças na forma $a^n b^{2n} c^m$, $n > 0$ e $m > 2$.

4. Defina uma gramática que gere as expressões numéricas com o seguinte alfabeto: dígitos (0 a 9), sinais (+, -, *, /) e parênteses (e). Teste sua definição mostrando a derivação das seguintes sentenças:

- a) $(5+5)/(4+1)$
- b) $4*(7/9)$
- c) $((5)/((4+5)*(1+3)))$

5. Escreva uma gramática para cada uma das linguagens abaixo:

- * a) $\{ a^{2i+1} b^{i+3} \mid i \geq 0 \} \cup \{ a^{i+4} b^{3i} \mid i \geq 0 \}$
- * b) $\{ a^i b^k \mid k > 0, i > k \}$
- c) $\{ a^i b^j c^j d^i e^3 \mid i, j \geq 0 \}$
- d) $\{ a^i b^j c^k d^l \mid i, j, k, l > 0, i \neq l, j > k \}$

6. Indique o tipo de cada gramática:

- | | | | |
|-----|------------------------------------|---------------------------------|-------------------|
| () | $S ::= aS \mid b \mid \varepsilon$ | | |
| () | $S ::= AAA \mid aA$ | $AA ::= AbcAA \mid \varepsilon$ | $A ::= a \mid b$ |
| () | $S ::= aB \mid b$ | $B ::= aB \mid b$ | |
| () | $S ::= (S) \mid S + S$ | $(S+S) ::= [id + id]$ | |
| () | $S ::= ABc \mid abc$ | $A ::= ab \mid \varepsilon$ | $B ::= bc \mid b$ |

7. Construa expressões regulares para as cadeias de 0s e 1s descritas abaixo:

- * a) com número par de 1s e 0s.
- b) com número ímpar de ocorrências do padrão 00.
- * c) com pelo menos duas ocorrências do padrão 101.
- d) com qualquer cadeia de entrada.
- e) com nenhuma cadeia de entrada.
- f) as cadeias 0110 e 1001.
- e) todas as cadeias que começam com 01 e terminam com 10.
- g) todas as cadeias que contenham exatamente quatro 1s.
- h) a cadeia nula e 001.

III - Linguagens Regulares e Autômatos Finitos

1. Linguagens Regulares

As linguagens regulares constituem um conjunto de linguagens decidíveis bastante simples e com propriedades bem definidas e compreendidas. Essas linguagens podem ser reconhecidas por autômatos finitos e são facilmente descritas por expressões simples, chamadas *expressões regulares* (E.R.) Antes de introduzir o mecanismo das E.R.s, será feita uma revisão das operações de concatenação e fechamento de conjuntos, já que são estas as operações básicas para a definição de linguagens regulares.

1.1. Operações de Concatenação e Fechamento

Seja C um conjunto finito de símbolos (alfabeto), e sejam L, L_1 e L_2 subconjuntos de C^* (linguagens sobre o alfabeto C). A concatenação de L_1 e L_2 , denotada por L_1L_2 é o conjunto $\{xy \mid x \in L_1 \text{ e } y \in L_2\}$. Em outras palavras, as cadeias da linguagem L_1L_2 são formadas por uma cadeia de L_1 concatenada a uma cadeia em L_2 , nesta ordem, incluindo aí todas as combinações possíveis.

Define-se $L^0 = \{\varepsilon\}$ e $L^n = LL^{n-1}$ para $n \geq 1$. O *fecho de Kleene* (ou simplesmente *fecho*) de uma linguagem L , denotado por L^* , é o conjunto:

$$L^* = \bigcup_{i=0}^{\infty} L^i = L^0 \cup L^1 \cup L^2 \cup L^3 \cup \dots$$

e o *fecho positivo* da linguagem L , denotado por L^+ , é o conjunto:

$$L^+ = \bigcup_{i=1}^{\infty} L^i = L^1 \cup L^2 \cup L^3 \cup L^4 \cup \dots$$

Em outras palavras, L^* denota as cadeias construídas pela concatenação de qualquer número de cadeias tomadas de L . O conjunto L^+ é semelhante, mas neste caso, as cadeias de zero palavras, cuja concatenação é definida como ε , são excluídas. Note-se, porém, que L^+ contém ε se e somente se L contém ε . Esta definição difere da definição do fechamento de alfabetos, onde A^+ era definido como $A^* - \{\varepsilon\}$. Note-se que no caso de linguagens, podem ocorrer dois casos:

- a) Se $\varepsilon \in L$, então $L^+ = L^*$.
- b) Se $\varepsilon \notin L$, então $L^+ = L^* - \{\varepsilon\}$.

Exemplos:

1. Seja $L_1 = \{10, 1\}$ e $L_2 = \{0011, 11\}$. Então $L_1L_2 = \{100011, 1011, 10011, 111\}$
2. $\{10, 11\}^* = \{\varepsilon, 10, 11, 1010, 1011, 1110, 1111, \dots\}$

1.2. Definição de Expressões Regulares

Seja A um alfabeto. As expressões regulares sobre o alfabeto A e os conjuntos (linguagens) que elas denotam são definidas indutivamente como segue:

- a) \emptyset é uma expressão que denota o conjunto vazio.
- b) ε é uma expressão regular que denota o conjunto $\{\varepsilon\}$.
- c) Para cada símbolo a em A , a é uma expressão regular que denota o conjunto $\{a\}$.
- d) Se r e s são expressões regulares que denotam os conjuntos R e S , respectivamente, então $(r+s)$, (rs) , e (r^*) são expressões regulares que denotam os conjuntos $R \cup S$, RS e R^* , respectivamente.

Quando se escreve expressões regulares, pode-se omitir muitos parênteses se for assumido que o símbolo * tem precedência maior que a concatenação e +, e que a concatenação tem precedência maior que +. Por exemplo, $((0(1^*)))+0$ pode ser escrita como 01^*+0 . Além disso, pode-se aplicar simplificações abreviando expressões como rr^* para r^+ .

As expressões regulares possuem as seguintes propriedades algébricas:

AXIOMA	DESCRIÇÃO
$r + s = s + r$	+ é comutativo.
$r + (s + t) = (r + s) + t$	+ é associativo.
$(rs)t = r(st)$	a concatenação é associativa.
$r(s + t) = rs + rt$ $(s + t)r = sr + tr$	a concatenação é distributiva sobre +.
$\epsilon r = r$ $r\epsilon = r$	ϵ é o elemento neutro (identidade) da concatenação.
$r^* = (r + \epsilon)^*$	relação entre ϵ e $*$.
$r^{**} = r^*$	$*$ é idempotente.

Exemplos:

- 00 é uma expressão regular que denota a linguagem $\{ 00 \}$.
- A expressão $(0+1)^*$ denota todas as cadeias de 0s e 1s.
- $(0+1)^*00(0+1)^*$ denota as cadeias de 0s e 1s com pelo menos dois 0s consecutivos.
- $(1+10)^*$ denota as cadeias de 0s e 1s que começam com 1 e não tem 0s consecutivos.
- $(0+1)^*001$ denota as cadeias de 0s e 1s que terminam por 001.
- a^+b^* denota as cadeias que tem qualquer quantidade de a's (mínimo 1) seguidos de qualquer quantidade de b's (possivelmente 0). Exemplo = $\{ a, aaa, aab, aabbbbb, ab, abbbb, \dots \}$

É importante notar que as linguagens regulares podem ser expressas tanto por expressões regulares como por gramáticas regulares. Assim, qualquer linguagem descrita na forma de uma expressão regular pode ser descrita por uma gramática regular. A seguir é apresentada uma tabela que mostra linguagens regulares representadas por expressões regulares e por gramáticas regulares.

Expressão Regular	Gramática Regular
$(0 + 1)^*$	$S ::= 0S \mid 1S \mid \epsilon$
001^*	$S ::= 0A \quad A ::= 0B \quad B ::= 1B \mid \epsilon$
$(00)^*(11)^+$	$S ::= 0A \mid 1B \quad A ::= 0S \quad B ::= 1 \mid 1C \quad C ::= 1B$
$a(aa)^*(b+c)^*$	$S ::= aA \quad A ::= \epsilon \mid aB \mid bC \mid cC \quad B ::= aA$ $C ::= bC \mid cC \mid \epsilon$

2. Autômatos Finitos

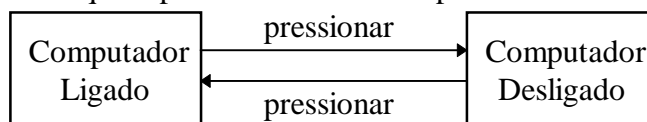
O Autômato Finito (AF) é o tipo mais simples de reconhecedor de linguagens. Ele é usado como reconhecedor de padrões em processamento de textos e também como analisador léxico de linguagens.

Autômatos finitos são usados com vantagens na compilação porquê:

- Eles têm a capacidade de realizar várias tarefas simples de compilação. Em particular, o projeto do bloco léxico pode ser feito baseado em um autômato finito.
- Uma vez que a simulação de um AF por um computador requer apenas um pequeno número de operações para processar um símbolo de entrada individual, ela opera rapidamente.

- c) A simulação de um AF requer uma quantidade finita de memória.
- d) Existem vários teoremas e algoritmos para construir e simplificar autômatos finitos para vários propósitos.

O conceito fundamental do AF é o conceito de *estado*. O estado representa uma condição sobre o modelo representado. Por exemplo, um computador pode ser representado por um modelo com dois estados: *ligado* e *desligado*. Pode-se ainda dizer que o ato de pressionar um botão ocasiona uma mudança de estado, dependendo do estado atual. Esta mudança de estado é denominada *transição*. O AF que representa este modelo pode ser:



Um Autômato Finito é um reconhecedor de Linguagens Regulares. Entende-se por *reconhecedor* de uma linguagem L um dispositivo que tomando uma seqüência w como entrada, responde “sim” se $w \in L$ e “não” caso contrário.

Os Autômatos Finitos classificam-se em:

- Autômatos Finitos Determinísticos (AFD)
- Autômatos Finitos Não-Determinísticos (AFND)

2.1 Autômatos Finitos Determinísticos

Formalmente definimos um AFD como sendo um sistema formal $M = (K, \Sigma, \delta, q_0, F)$, onde:

- $K \rightarrow$ É um conjunto finitos não-vazio de ESTADOS;
- $\Sigma \rightarrow$ É um ALFABETO, finito, de entrada;
- $\delta \rightarrow$ FUNÇÃO DE MAPEAMENTO (ou função de transição), definida em $K \times \Sigma \rightarrow K$;
- $q_0 \rightarrow \in K$, é o ESTADO INICIAL;
- $F \rightarrow \subseteq K$, é o conjunto de ESTADOS FINAIS.

2.1.1. Interpretação de δ

A interpretação de $\delta(q, a) = p$, onde $\{q, p\} \in K$ e $a \in \Sigma$, é de que se o “controle de M” está no estado “q” e o próximo símbolo de entrada é “a”, então “a” deve ser reconhecido e o “controle” passa para o próximo estado (no caso, “p”). Note que a função de mapeamento é definida por $K \times \Sigma \rightarrow K$, ou seja, de um estado, através do reconhecimento de um símbolo, o controle passa para UM ÚNICO outro estado. Por isso este tipo de autômato é chamado de DETERMINÍSTICO.

2.1.2. Significado Lógico de um Estado

Logicamente um estado é uma situação particular no processo de reconhecimento de uma sentença. Mais genericamente, um estado representa uma condição sobre a parte da sentença que já foi reconhecida.

2.1.3. Sentenças Aceitas por M

Uma sentença x é aceita (reconhecida) por um AF $M = (K, \Sigma, \delta, q_0, F)$ se, e somente se, $\delta(q_0, x) = p \wedge p \in F$. Em outras palavras, a sentença x é reconhecida se, partindo do estado inicial q_0 , através das transições para cada símbolo de x, um estado final é alcançado.

2.1.4. Linguagem Aceita por M

É um conjunto de todas as sentenças aceitas por M. Formalmente, definimos por:

$$L(M) = \{x \mid \delta(q_0, x) = p \wedge p \in F\}.$$

Ou seja, é o conjunto de todas as sentenças que, partindo do estado inicial, fazem com que o autômato alcance um estado final quando termina de reconhecer a sentença.

Obs.: Todo conjunto aceito por um Autômato Finito é um Conjunto Regular.

2.1.5. Diagrama de Transição

Um diagrama de transição para um AF M é um grafo direcionado e rotulado. Os vértices representam os estados e fisicamente são representados por círculos, sendo que o estado inicial é diferenciado por uma seta com rótulo “início” e os estados finais são representados por círculos duplos. As arestas representam as transições, sendo que entre dois estados “ p ” e “ q ” existirá uma aresta direcionada de “ p ” para “ q ” com rótulo “ a ” ($a \in \Sigma$) $\leftrightarrow \exists \delta(p, a) = q$ em M .

2.1.6. Tabela de Transições

É uma representação tabular de um AF

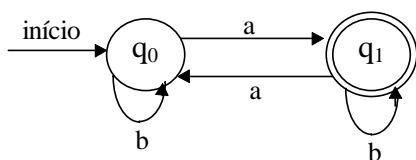
Nesta tabela as linhas representam os estados (o inicial é indicado por uma seta e os finais por asteriscos), as colunas representam os símbolos de entrada e o conteúdo da posição (q, a) será igual a “ p ” se existir $\delta(q, a) = p$, senão será indefinida.

Exemplos

Nesta seção daremos alguns exemplos de AFD nas formas gráficas e tabulares.

Exemplo 1: AFD que reconheça as sentenças da linguagem $L = \{(a,b)^+, \text{onde o número de } a\text{'s é ímpar}\}$.

Forma Gráfica:



Forma Tabular:

δ	a	b
\rightarrow q0	q1	q0
* q1	q0	q1

É possível interpretar o estado q_0 como a condição lógica “há um número par de a 's na sentença” e o estado q_1 como a condição “há um número ímpar de a 's na sentença”.

Para reconhecer a sentença “abbaba” o autômato realiza as seguintes transições, partindo do estado inicial:

q0 abbaba
a q1 bbaba
ab q1 baba
abb q1 aba
abba q0 ba
abbab q0 a
abbaba q1

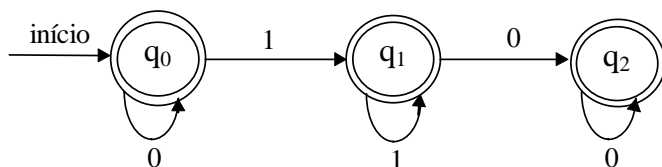
Como a sentença terminou com o autômato num estado final, a mesma é reconhecida, ou seja, faz parte da linguagem definida. Já a sentença “ababaab” leva às seguintes transições do AFD acima:

q0 ababaab
a q1 babaab
ab q1 abaab
aba q0 baab
abab q0 aab
ababa q1 ab
ababaa q0 b
ababaab q0

Ao final da sentença o AFD está num estado não-final. Assim, a sentença não faz parte da linguagem.

Exemplo 2: AFD que reconheça as sentenças da linguagem $L = \{(0,1)^*\}$, onde todos os 1's apareçam consecutivos}.

Forma Gráfica:



Forma Tabular:

δ	0	1
\rightarrow^* q0	q0	q1
* q1	q2	q1
* q2	q2	-

Note que neste caso, o não reconhecimento de sentenças com 1's não consecutivos é impossível pois não há transição no estado q2 para o símbolo 1. A sentença "0101" leva às seguintes transições:

- q0 0101
- 0 q0 101
- 01 q1 01
- 010 q2 1 (?)

Quando o AFD, no estado atual, não apresentar transição para um símbolo da entrada, a sentença também é rejeitada.

2.2. Autômatos Finitos Não-Determinísticos

Um AFND é um sistemas formal $M = M = (K, \Sigma, \delta, q_0, F)$ onde:

$K, \Sigma, q_0, F \rightarrow$ possuem a mesma definição de AFD

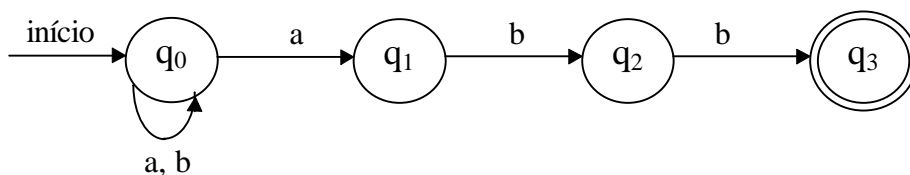
$\delta \rightarrow$ é uma função de mapeamento, definida em $K \times \Sigma = \rho(K)$ (um conjunto de estados), sendo que $\rho(K)$ é um subconjunto de K . Isto equivale a dizer que $\delta(q, a) = p_1, p_2, \dots, p_n$. A interpretação de δ é que M no estado "q", com o símbolo "a" na entrada pode ir tanto para o estado p_1 , como para o estado p_2, \dots , como para o estado p_n .

Uma sentença x é aceita por um AFND se $\delta(q_0, x) = P$, e $\{P \cap F\} \neq \emptyset$. Em outras palavras, uma sentença x é aceita se, dentre os estados alcançados, pelo menos um é final.

Na análise de uma sentença, quando um símbolo leva à mais de um estado, pode-se dizer que o AFND se divide em vários, sendo que cada um deles continua o reconhecimento da sentença por uma das possibilidades. Se pelo menos um deles chegar a um estado final no final da sentença, então ela é válida. Se nenhum chegar a um estado final, então a sentença é inválida.

Exemplo: Seja a linguagem $L = \{(a,b)^* abb\}$. O AFND seria:

Forma Gráfica:



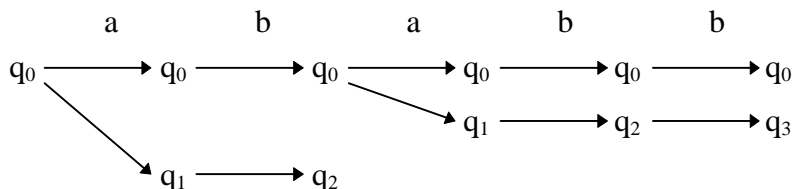
Forma

Tabular:

δ	a	b
\rightarrow q0	q0, q1	q0

	q ₁	-	q ₂
	q ₂	-	q ₃
*	q ₃	-	-

Para reconhecer a sentença “ababb”, o autômato poderia seguir as seguintes alternativas:



Há três possibilidades no reconhecimento da sentença:

- um caminho termina num estado não-final (q₀);
- um caminho chega a um estado (q₂) onde não há transição definida para o próximo símbolo (a);
- um caminho chega a um estado final (q₃) no final da sentença.

Como há a possibilidade de chegar a um estado final com a sentença, então a mesma é válida, ou seja, deve ser reconhecida como parte da linguagem definida.

2.3. Comparação entre AFD e AFND

Os AFDs e os AFNDs representam a mesma classe de linguagens, ou seja, as linguagens regulares. A tabela abaixo mostra as principais diferenças entre AFD e AFND:

	Vantagens	Desvantagens
AFD	Implementação trivial	Não representa claramente algumas L. R.
AFND	Representa mais claramente algumas L. R.	Implementação complexa

2.4. Transformação de AFND para AFD

Como os autômatos finitos representam a mesma classe de linguagens, sejam determinísticos ou não-determinísticos, é de se esperar que haja uma equivalência entre eles. Realmente esta equivalência existe. Assim, para cada AFND é possível construir um AFD equivalente. Nesta seção é apresentado um algoritmo que converte um AFND num AFD equivalente. Assim, é possível utilizar a clareza da representação de um AFND, e para fins de facilidade de implementação, transformá-lo num AFD.

Teorema 3.1: “Seja L um conjunto aceito por um AFND., então \exists um AFD que aceita L.”

Prova: Seja $M = (K, \Sigma, \delta, q_0, F)$ um AFND Construa um A. F. D. $M' = (K', \Sigma', \delta', q_0', F')$ como segue:

- 1 - $K' = \{\rho(K)\} \rightarrow$ isto é, cada estado de M' é formado por um conjunto de estados de M .
- 2 - $q_0' = [q_0] \rightarrow$ obs.: representaremos um estado $q \in K'$ por $[q]$.
- 3 - $F' = \{\rho(K) \mid \rho(K) \cap F \neq \emptyset\}$.
- 4 - $\Sigma' = \Sigma$.
- 5 - Para cada $\rho(K) \subset K', \delta'(\rho(K), a) = \rho'(K)$, onde $\rho'(K) = \{p \mid \text{para algum } q \in \rho(K), \delta(q, a) = p\}$, ou seja:

Se $\rho(K) = [q_1, q_2, \dots, q_r] \in K'$ e

se $\delta(q_1, a) = p_1, p_2, \dots, p_j$

$\delta(q_2, a) = p_{j+1}, p_{j+2}, \dots, p_k$

...

$\delta(q_r, a) = p_i, p_{i+1}, \dots, p_n$ são as transições de M ,

então $\rho'(K) = [p_1, p_2, \dots, p_j, p_{j+1}, \dots, p_k, \dots, p_i, p_{i+1}, \dots, p_n]$ será um estado de M' e M' conterá a transição: $\delta'(\rho(K), a) = \rho'(K)$.

Para concluir a prova do teorema, basta provar que a linguagem de M' é igual à linguagem de M ($L(M) = L(M')$), prova esta que não é apresentada neste texto.

Exemplo: Tomemos por exemplo o AFND da seção 3, representado na tabela abaixo:

	δ	a	b
→	q_0	q_0, q_1	q_0
	q_1	-	q_2
	q_2	-	q_3
*	q_3	-	-

Para determinizar o AFND acima, basta definir o AFD $M' = (K', \Sigma', \delta', q_0', F')$, onde

	δ'	a	b
→	$[q_0]$	$[q_0, q_1]$	$[q_0]$
	$[q_0, q_1]$	$[q_0, q_1]$	$[q_0, q_2]$
	$[q_0, q_2]$	$[q_0, q_1]$	$[q_0, q_3]$
*	$[q_0, q_3]$	$[q_0, q_1]$	$[q_0]$

$K' = \{[q_0], [q_0, q_1], [q_0, q_2], [q_0, q_3]\}$

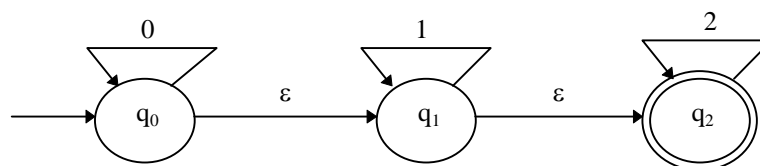
$\Sigma' = \{a, b\}$

$q_0' = [q_0]$

$F' = \{[q_0, q_3]\}$

2.5. Autômato Finito com ϵ -Transições

Pode-se estender o modelo de autômato finito não-determinístico para incluir transições com entrada vazia: ϵ . O AFND da figura seguinte aceita uma linguagem constituída por qualquer número de 0's, seguidos por qualquer número de 1's, seguido por qualquer número de 2's (incluindo a quantidade 0, para cada um dos casos).



Definição:

Um *autômato finito não-determinístico com ϵ -transições* é uma quintupla $(K, \Sigma, \delta, q_0, F)$, com todos os componentes de um AFND, mas com a relação de transição δ definida sobre $K \times \Sigma \cup \{\epsilon\} \rightarrow K$.

Na representação da tabela de transições do AF, é utilizada uma coluna para a ϵ -transição. No caso do exemplo, a tabela seria a seguinte:

	δ	0	1	2	ϵ
→	q_0	q_0			q_1
	q_1		q_1		q_2
*	q_2			q_2	

2.5.1. Equivalência Entre AFs Com e Sem ϵ -Transições

Para um autômato finito com ϵ -transições, define-se o ϵ -fechamento de um estado q como sendo o conjunto de estados alcançáveis a partir de q utilizando-se somente ϵ -transições.

Para o AF do exemplo anterior, têm-se os seguintes ϵ -fechamento dos estados:

ε -fecho(q_0) = { q_0, q_1, q_2 }

ε -fecho(q_1) = { q_1, q_2 }

ε -fecho(q_2) = { q_2 }

Seja $M = (K, \Sigma, \delta, q_0, F)$ um AF com ε -transições, pode-se construir $M' = (K, \Sigma, \delta', q_0, F')$ sem ε -transições, utilizando-se o seguinte algoritmo:

Algoritmo: Eliminação de ε -transições

Entrada: Um AFND com ε -transições $M = (K, \Sigma, \delta, q_0, F)$

Saída: Um AFND com ε -transições $M = (K, \Sigma, \delta', q_0, F')$

F' será definido por $F \cup \{q \in K \text{ e } \varepsilon\text{-fecho}(q) \text{ contém um estado de } F\}$.

δ' deve conter todas as transições de t que não são ε -transições, mais aquelas que são obtidas pela composição de transições de δ' com as ε -transições de δ , da seguinte maneira:

Se $(q_1, a) = q_2 \in \delta'$, e $q_3 \in \varepsilon\text{-fecho}(q_2)$ então

coloque $(q_1, a) = q_3$ em δ'

Fim Se;

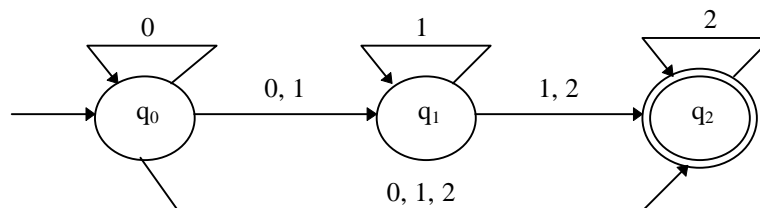
Se $(q_1, a) = q_2 \in \delta'$, e $q_3 \in \varepsilon\text{-fecho}(q_3)$ então

coloque $(q_3, a) = q_2$ em δ'

Fim Se;

Fim;

Para o autômato do exemplo anterior, têm-se o seguinte AFND sem ε -transições:



Claramente, para todo AFND com ε -transições, há um AFD equivalente, uma vez que basta eliminar as ε -transições, e em seguida determinar o AFND resultante, chegando-se ao AFD equivalente.

3. Relação Entre GR e AF

As gramáticas regulares são sistemas geradores das linguagens regulares, enquanto os AFs são sistemas reconhedores do mesmo conjunto de linguagens. Assim, há uma correspondência entre as GR e os AF. Nesta seção são apresentados algoritmos para determinar o AF para reconhecer a linguagem de um GR, bem como para determinar a GR que gera a linguagem de um dado AF.

Teorema 3.2: “Seja $G = (N, T, P, S)$ uma Gramática Regular, então \exists um AF $M = (K, \Sigma, \delta, q_0, F) \mid L(M) = L(G)$ ”.

Prova: a - Mostrar que M existe;

b- Mostrar que $L(M) = L(G)$.

a) Defina M como segue:

1 - $K = N \cup \{A\}$, onde A é um novo símbolo não-terminal;

2 - $\Sigma = T$;

3 - $q_0 = S$;

4 - $F = \{B \mid B \rightarrow \varepsilon \in P\} \cup \{A\}$, ou seja, todos os estados que produzem ε ,

juntamente com o novo estado.

5 - Construa δ de acordo com as regras a, b e c:

a) Para cada produção da forma $B \rightarrow a \in P$, crie a transição $\delta(B,a) = A$;

- b) Para cada produção da forma $B \rightarrow aC \in P$, crie a transição $\delta(B,a) = C$;
- c) Para produções da forma $B \rightarrow \varepsilon$, não é criada nenhuma transição.
- d) Para todo $a \in T$, $\delta(A,a) = -$ (indefinida)

b) Para mostrar que $L(M) = L(G)$, devemos mostrar que (1) $L(G) \subseteq L(M)$ e (2) $L(M) \subseteq L(G)$. Esta demonstração é fácil considerando-se as produções da GR e a definição das transições de M, sendo omitida neste texto.

Exemplo: Dada a GR abaixo, que gera as sentenças da linguagem $\{(a^n b^m)\}$, com n par e $m > 0$, iremos definir um AF que reconheça as sentenças desta linguagem.

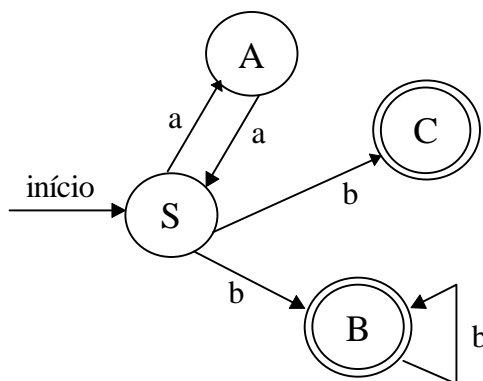
Gramática: $S ::= aA \mid bB \mid b$
 $A ::= aS$
 $B ::= bB \mid \varepsilon$

Autômato:

Primeiramente devemos definir os estados do autômato, com um estado novo, denominado aqui C. Assim, temos que $K = \{S, A, B, C\}$. O alfabeto é o mesmo, assim $\Sigma = (a, b)$, o estado inicial é S, e $F = \{B, C\}$, pela definição acima. Na forma tabular o autômato fica o seguinte:

δ	a	b
→ S	A	B, C
A	S	-
* B	-	B
* C	-	-

Na forma gráfica o autômato fica o seguinte:



Teorema 3.3: “Seja um AF $M = (K, \Sigma, \delta, q_0, F)$. Então \exists uma GR $G = (N, T, P, S) \mid L(G) = L(M)$ ”.

Prova: a - Mostra que G existe;

b - Mostrar que $L(G) = L(M)$.

a) Seja $M = (K, \Sigma, \delta, q_0, F)$ um AF. Construa uma G. R. $G = (N, T, P, S)$ como segue:

1 - $N = K$;

2 - $T = \Sigma$;

3 - $S = q_0$;

4 - Defina P como segue:

4.a) Se $\delta(B,a) = C$, então adicione $B ::= aC$ em P;

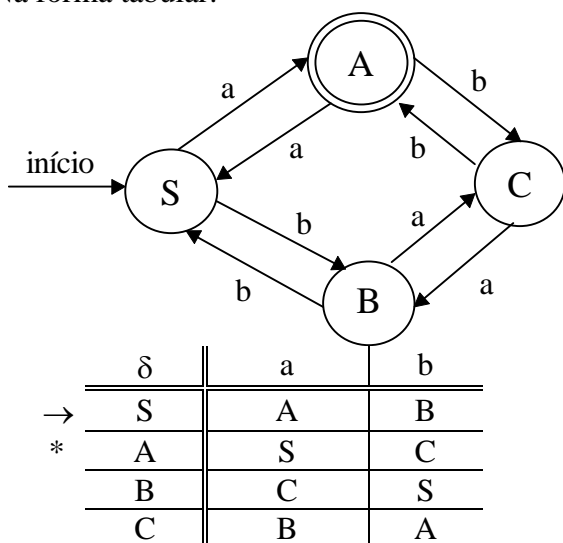
4.b) Se $\delta(B,a) = C$ e $C \in F$, então adicione $B ::= a$ em P;

4.c) Se $q_0 \in F$, então $\varepsilon \in L(M)$. Assim a gramática deve ser transformada encontrar uma outra G.R. G1, tal que $L(G1) = L(G) \cup \{\varepsilon\}$, e $L(G1) = L(M)$. Senão $\varepsilon \notin L(M)$, e $L(G) = L(M)$.

b) Para mostrar que $L(G) = L(M)$, devemos mostrar que (1) $L(M) \subseteq L(G)$ e (2) $L(G) \subseteq L(M)$. Esta demonstração é análoga à parte (b) da prova do teorema 5.1, sendo omitida neste texto.

Exemplo: Tomemos o seguinte A. F.:

Na forma tabular:



A GR definida pelo algoritmo acima resulta em:

$S ::= aA \mid bB \mid a$

$A ::= aS \mid bC$

$B ::= aC \mid bS$

$C ::= aB \mid bA \mid b$

Pelo algoritmo acima, a GR gerada a partir de um AF nunca apresentará produções do tipo $A ::= \epsilon$.

4. Minimização de Autômatos Finitos

Definição: Um A. F. D. $M = (K, \Sigma, \delta, q_0, F)$ é MÍNIMO se:

- 1 - Não possui estados INACESSÍVEIS;
- 2 - Não possui estados MORTOS;
- 3 - Não possui estados EQUIVALENTES.

Estados Inacessíveis:

Um estado $q \in K$ é inacessível quando não existe w tal que a partir de q_0 , q seja alcançado, ou seja, $\sim \exists w \mid \delta(q_0, w) = q$, onde w é uma sentença ou parte dela.

Estados Mortos:

Um estado $q \in K$ é morto se ele $\notin F$ e $\sim \exists w \mid \delta(q, w) = p$, onde $p \in F$ e w é uma sentença ou parte dela. Ou seja, q é morto se ele não é final e a partir dele nenhum estado final pode ser alcançado.

Estados Equivalentes:

Um conjunto de estados q_1, q_2, \dots, q_j são equivalentes entre si, se eles pertencem a uma mesma CLASSE DE EQUIVALÊNCIA.

Classe de Equivalência (CE):

Um conjunto de estados q_1, q_2, \dots, q_j estão em uma mesma CE se $\delta(q_1, a), \delta(q_2, a), \dots, \delta(q_j, a)$, para cada $a \in \Sigma$, resultam respectivamente nos estados q_i, q_{i+1}, \dots, q_n e estes pertencem a uma mesma CE.

4.1. Algoritmo para Construção das Classes de Equivalência

1. Se necessário, crie um estado ϕ para representar as indefinições.
2. Devida K em duas CE, uma contendo F e outra contendo $K-F$.

3. Divida as CE existentes, formando novas CE (de acordo com a definição - lei de formação das CE), até que nenhuma nova CE seja formada.

4.2. Algoritmo para Construção do Autômato Finito Mínimo

Entrada: Um AFD $M = (K, \Sigma, \delta, q_0, F)$;

Saída: Um AFD Mínimo $M' = (K', \Sigma, \delta', q_0', F') \mid M' \equiv M$;

Método:

- 1 - Elimine os estados inacessíveis.
- 2 - Elimine os estados mortos.
- 3 - Construa todas as possíveis classes de equivalência.
- 4 - Construa M' como segue:
 - a) $K' \rightarrow$ é o conjunto de CE obtidas.
 - b) $q_0' \rightarrow$ será a CE que contiver q_0 .
 - c) $F' \rightarrow$ será o conjunto das CE que contenham pelo menos um elemento $\in F$;

isto é, $\{[q] \mid \exists p \in [q] \text{ e } p \in F, \text{ onde } [q] \text{ é um CE}\}$.

d) $\delta' \rightarrow \delta'([p], a) = [q] \leftrightarrow \delta(p, a) = q$ é uma transição de M e p e q são elementos de $[p]$ e $[q]$, respectivamente.

Exemplo: Minimize o AFD definido na seguinte tabela de transições:

	δ'	a	b
→*	A	G	B
	B	F	E
	C	C	G
*	D	A	H
	E	E	A
	F	B	C
*	G	G	F
	H	H	D

Temos que os estados acessíveis são: {A, G, B, F, E, C}. Portanto podemos eliminar os estados D e H. Assim o novo AFD é:

	δ'	a	b
→*	A	G	B
	B	F	E
	C	C	G
*	E	E	A
	F	B	C
	G	G	F

Nenhum dos estados do AFD acima é morto. Assim, podemos construir as classes de equivalência. No primeiro passo, duas classes de equivalência são construídas: F e K-F. Cada passo do algoritmo é representado numa linha da tabela abaixo:

	F	K-F
1	{A, G}	{B, C, E, F}
2	{A, G}	{B, F} {C, E}
3	{A, G}	{B, F} {C, E}

A primeira linha da tabela é a divisão dos estados em F e K-F. Para $i > 1$, a linha i é formada pela separação das CEs da linha $i-1$, até que uma linha seja repetida.

Algoritmicamente, a separação das CEs pode ser dada por:

SE houver mais de um estado na CE ENTÃO
 Crie uma nova CE aberta na linha $i+1$ para o primeiro estado da CE;
 PARA todos os outros estados da CE FAÇA
 Procure encaixar o estado numa CE aberta
 SE não for possível ENTÃO
 crie uma nova CE aberta na linha $i+1$ para o estado;
 FIMSE
 FIMPARA
 Feche todas as CEs abertas
 SENÃO
 Repita a CE na linha $i+1$

Na tabela acima, na linha 2, foi aberta uma classe para B, e a seguir os outros estados foram testados. C não pôde ficar na mesma classe de B (pois $\delta(B, b) = E$ e $\delta(C, b) = G$, e E e G pertencem à CEs separadas na primeira linha. Assim, uma nova CE foi aberta para C. O estado E não pôde ficar na CE de B, mas pôde ficar na CE de C, assim não foi aberta uma nova CE para E. Já o estado

F pôde ficar na CE de B. Como não havia mais estados na CE da primeira linha, as duas CEs da segunda linha foram fechadas.

A repetição do algoritmo para a criação da linha 3 fez com que a linha 2 se repetisse. Assim, todas as CEs foram determinadas.

Denominando $\{A, G\} = q_0$, $\{B, F\} = q_1$ e $\{C, E\} = q_2$, temos o seguinte AFD Mínimo:

δ'	a	b
\rightarrow^* q ₀	q ₀	q ₁
q ₁	q ₁	q ₂
q ₂	q ₂	q ₀

5. Construção do Analisador Léxico

No processo de compilação, o Analisador Léxico é responsável pela identificação dos *tokens*, ou seja, das menores unidades de informação repassadas para o Analisar Sintático.

Assim, pode-se dizer que o analisador léxico é responsável pela leitura dos caracteres da entrada, agrupando-os em palavras, que são classificadas em categorias. Estas categorias podem ser, basicamente, as seguintes:

Palavras reservadas: palavras que devem aparecer literalmente na linguagem, sem variações. Por exemplo, na linguagem PASCAL, algumas palavras reservadas são: BEGIN, END, IF, ELSE.

Identificadores: palavras que seguem algumas regras de escrita, porém podem assumir diversos valores. São definidos de forma genérica. Por exemplo, identificadores que devem iniciar com uma letra, seguida de letras ou dígitos. Pode-se representar estes identificadores pela seguinte expressão regular: $(a + b + \dots + c)(a + b + \dots + z + 0 + 1 + \dots + 9)^*$

Símbolos especiais: seqüências de um ou mais símbolos que não podem aparecer em identificadores nem palavras reservadas. São utilizados para composição de expressões aritméticas ou lógicas, por exemplo. Exemplos de símbolos especiais são: “;” (ponto-e-vírgula), “:” (dois pontos), “:=” (atribuição).

Constantes: também denominadas *literais*, são valores que são inseridos diretamente no programa, ao invés de serem uma referência para um endereço de memória. Por exemplo, números (904, 8695, -9, 12E-9), strings (“Hello world”), booleanos (TRUE, FALSE), etc.

Para a construção de um analisador léxico, são necessários 4 passos:

- 1 - definir as gramáticas e autômatos finitos determinísticos mínimos para cada token;
- 2 - juntar os autômatos num único autômato e determinizar (não minimizar);
- 3 - acrescentar um estado de erro, para tratar tokens inválidos;
- 4 - implementar o autômato numa linguagem de programação.

A construção de uma analisador léxico é demonstrada a seguir, com um exemplo:

Sejam os tokens desejados os seguintes:

PALAVRAS RESERVADAS: goto, do.

IDENTIFICADORES: começam com uma letra, seguidas de letras e “_”. Expressão

Regular: $(a + \dots + z)(a + \dots + z + _)^*$

1º passo: definir as gramáticas e autômatos para cada token.

Token “GOTO”:

$C ::= o$

$S ::= gA$

$A ::= oB$

$B ::= tC$

	g	t	o
$\rightarrow S$	A	-	-

III - Linguagens Regulares e Autômatos Finitos

A	-	-	B
B	-	C	-

Token "DO":

$S ::= dD$

$D ::= o$

C	-	-	X
* X	-	-	-
	d	o	
$\rightarrow S$	D	-	
D	-	Y	
* Y	-	-	

Token "Identificador":

$S ::= aE \mid bE \mid \dots \mid zE$

$E ::= aE \mid bE \mid \dots \mid zE \mid _E \mid \varepsilon$

	a	b	...	z	_
$\rightarrow S$	E	E	E	E	-
* E	E	E	E	E	E

Note que os autômatos acima já foram minimizados, sendo que o estado novo do autômato de identificadores foi eliminado por ser inalcançável.

2º passo: juntar os autômatos num só e determinar:

Os autômatos acima serão agora colocados juntos num só autômato. Note que já foi tomado o cuidado de somente o nome do estado inicial coincidir. Todos os outros estados têm nomes diferentes. O autômato não-determinístico geral fica:

	g	t	o	d	a	b	...	z	_
$\rightarrow S$	A, E	E	E	D, E	E	E	E	E	-
A	-	-	B	-	-	-	-	-	-
B	-	C	-	-	-	-	-	-	-
C	-	-	X	-	-	-	-	-	-
D	-	-	Y	-	-	-	-	-	-
* X	-	-	-	-	-	-	-	-	-
* Y	-	-	-	-	-	-	-	-	-
* E	E	E	E	E	E	E	E	E	E

O autômato acima, determinado (e não minimizado) fica:

	g	t	o	d	a	b	...	z	_
→ S [S]	A	E	E	D	E	E	E	E	-
* A [A,E]	E	E	B	E	E	E	E	E	E
* B [B, E]	E	C	E	E	E	E	E	E	E
* C [C, E]	E	E	X	E	E	E	E	E	E
* D [D, E]	E	E	Y	E	E	E	E	E	E
* X [X, E]	E	E	E	E	E	E	E	E	E
* Y [Y, E]	E	E	E	E	E	E	E	E	E
* E [E]	E	E	E	E	E	E	E	E	E

O autômato acima poderia ser minimizado, mas isto faria com que todos os estados finais se juntassem num só. Isto faz com que informações importantes sejam perdidas. Note que o estado final do token **GOTO** está no estado que contém o estado final X, e o estado final do token **DO** está no estado final que contém o estado final Y. Os estados que contém o estado final E são os que reconhecem os identificadores. Onde há conflito entre identificadores e palavras reservadas a prioridade é das palavras reservadas.

Assim, se um token terminar no estado X do autômato acima ele será um GOTO, no estado Y será um DO, e nos demais estados finais será um Identificador.

3º passo: acrescentar o estado de erro:

Note que no autômato acima são definidas as transições somente para letras e o caracter “_”. Para os demais caracteres (que são inválidos na nossa linguagem) é acrescentado um estado de erro. Este estado será alcançado quando houver um caracter inválido, bem como para as transições indefinidas para os caracteres válidos. Este estado de erro será um estado final, que reportará erro. O autômato acima ficará:

	g	t	o	d	a	b	...	z	_	etc.
→ S	A	E	E	D	E	E	E	E	F	F
* A	E	E	B	E	E	E	E	E	E	F
* B	E	C	E	E	E	E	E	E	E	F
* C	E	E	X	E	E	E	E	E	E	F
* D	E	E	Y	E	E	E	E	E	E	F
* X	E	E	E	E	E	E	E	E	E	F
* Y	E	E	E	E	E	E	E	E	E	F
* E	E	E	E	E	E	E	E	E	E	F
* F	F	F	F	F	F	F	F	F	F	F

(etc. denota os caracteres não válidos da linguagem).

4º passo: implementação:

O autômato acima é facilmente implementado numa estrutura de dados do tipo tabela. Ou seja, no reconhecimento de um novo token, parte-se do estado inicial, determinando-se o novo estado pela coluna correspondente ao caracter lido do arquivo. Repete-se este procedimento até achar um separador de token (espaço em branco, quebra de linha ou final de arquivo).

Quando um separador de token é encontrado, verifica-se qual o estado atual do autômato. Se for um estado não-final, houve um erro, que deve ser acusado. Se for um estado final, verifica-se o que o estado representa, e retorna-se o valor desejado.

No autômato acima, se um separador de tokens é encontrado, verifica-se o estado corrente. Cada estado representa um token ou erro, da seguinte forma:

S → erro;

X → token GOTO

Y → token DO

F → erro - caracter inválido

outro estado → Identificador.

EXERCÍCIOS:

1) Construa um A. F. D. para cada linguagem abaixo.

a) $\{w \mid w \in (0, 1)^+ \text{ e } |w| > 3\}$.

* b) $\{w \mid w \in (0, 1)^+ \text{ e } w \text{ comece com } 0 \text{ e termine com } 1\}$

* c) $\{w \mid w \in (0, 1, 2)^+ \text{ e a quantidade de } 1\text{'s em } w \text{ seja ímpar}\}$.

d) $\{w \mid w \in (a, b, c)^* \text{ e } w \text{ não possua letras iguais consecutivas}\}$.

e) $\{w \mid w \text{ é um comentário PASCAL não aninhado, na forma } (*x^*), \text{ e } x \in (a, b, \dots, z)^*\}$.

f) $\{\text{begin, end}\}$. Caso ocorra outra entrada o autômato deve ir para um estado de erro e continuar lá.

2) Construa as G. R. dos autômatos abaixo. A seguir determine-os e minimize-os. Construa a G. R. dos autômatos resultantes.

* a)

δ'	a	b
→ A	A, B	A
B	C	-
C	-	D
* D	D	D

b)

δ'	a	b	c	d
→ A	B, C	D	C	C
B	-	D	C	C
* C	-	-	-	D
D	-	B, C	-	-

c)

	δ'	a	b	c
\rightarrow^*	S	A	B, F	S, F
	A	S, F	C	A
	B	A	-	B, S, F
	C	S, F	-	A, C
*	F	-	-	-

d)

	δ'	0	1
\rightarrow	A	A, B	A
	B	C	A
	C	D	A
*	D	D	D

3) Construa os A. F. que reconhecem as linguagens geradas pelas seguintes G. R. Caso o A. F. seja não-determinístico, determine-o. Minimize todos os A. F.

* a) $S ::= 0S \mid 1S \mid 0A \mid 0C \mid 1B$
 $A ::= 0A \mid 0C \mid 0$
 $B ::= 1B \mid 1$
 $C ::= 0C \mid 0A \mid 0$

b) $S ::= aA \mid aC \mid bB \mid bC$
 $A ::= aF \mid a$
 $B ::= bF \mid b$
 $C ::= aA \mid aC \mid bB \mid bC$
 $F ::= aF \mid bF \mid a \mid b$

c) $S ::= aA \mid bB$
 $A ::= aS \mid aC \mid a$
 $B ::= bS \mid bD \mid b$
 $C ::= aB$
 $D ::= bA$

d) $S ::= 0B \mid 1A \mid 1 \mid \epsilon$
 $A ::= 0B \mid \epsilon$
 $B ::= 0C \mid 0 \mid 1D$
 $C ::= 0B \mid 1A \mid 1$
 $D ::= 1C \mid 1$

IV - Linguagens Livres de Contexto

1. Linguagens Auto-Embebidas

A importância das LLC's reside no fato de que praticamente todas as linguagens de programação podem ser descritas por este formalismo, e que um conjunto bastante significativo destas linguagens pode ser analisado por algoritmos muito eficientes. A maioria das linguagens de programação pertence a um subconjunto das LLC's que pode ser analisado por algoritmos que executam $n \cdot \log(n)$ passos para cada símbolo da sentença de entrada. Outro subconjunto que corresponde às linguagens não ambíguas pode ser analisado em tempo n^2 . No pior caso, uma LLC ambígua pode ser analisada em tempo n^3 . O que ainda é bastante eficiente, se comparado às linguagens sensíveis ao contexto, que podem requerer complexidade exponencial (2^n), sendo, portanto, computacionalmente intratáveis, e às linguagens de tipo 0, que podem ser, inclusive, indecidíveis, isto é, o algoritmo de análise pode não parar nunca.

2. Gramáticas Livres de Contexto

Uma *gramática livre de contexto* (GLC) é denotada por $G=(N,T,P,S)$, onde N e T são conjuntos disjuntos de variáveis e terminais, respectivamente. P é um conjunto finito de produções, cada uma da forma $A ::= \alpha$, onde A é uma variável do conjunto N e α é uma cadeia de símbolos $(N \cup T)^*$. Finalmente, S é uma variável especial denominada símbolo inicial.

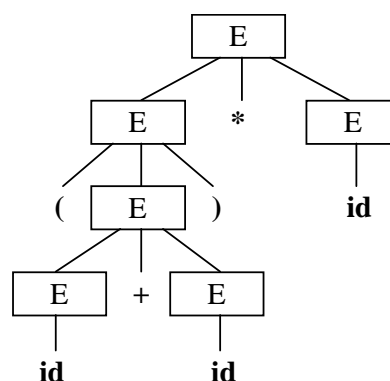
3. Árvore de Derivação

Muitas vezes é útil mostrar as derivações de uma GLC através de árvores. Estas figuras, chamadas árvores de derivação ou árvores sintáticas, impõe certas estruturas úteis às sentenças das linguagens geradas, principalmente as de programação.

Os vértices de uma árvore de derivação são rotulados com terminais ou variáveis, ou mesmo com a cadeia vazia. Se um vértice interior n é rotulado com A, e os filhos de n são rotulados com X_1, X_2, \dots, X_k , da esquerda para a direita, então $A ::= X_1 X_2 \dots X_k$ deve ser uma produção da gramática. Considere, como exemplo, a gramática:

$$\begin{aligned} E & ::= E + E \\ & \quad | E * E \\ & \quad | (E) \\ & \quad | id \end{aligned}$$

Uma árvore de derivação para a sentença “(id+id)*id” gerada por esta gramática poderia ser:



Mais formalmente, seja $G=(N,T,P,S)$ uma GLC. Então uma árvore é uma árvore de derivação para G, se:

- cada vértice tem um rótulo, que é um símbolo de $N \cup T \cup \{\epsilon\}$;
- o rótulo da raiz é o símbolo S;
- os vértices interiores são rotulados apenas com símbolos de N;

IV - Linguagens Livres de Contexto

- d) se o vértice n tem rótulo A , e os vértices n_1, n_2, \dots, n_k são filhos de n , da esquerda para a direita, com rótulos X_1, X_2, \dots, X_k , respectivamente, então $A ::= X_1 X_2 \dots X_k$ deve ser uma produção em P ;
- e) se o vértice n tem rótulo ε , então n é uma folha ε e é o único filho de seu pai.

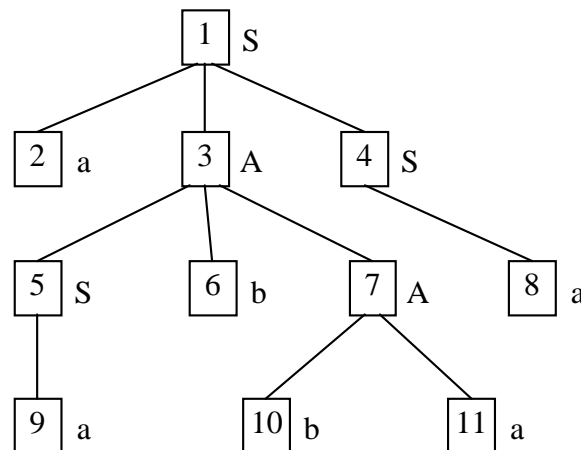
Exemplo:

Considere a gramática $G = (\{S, A\}, \{a, b\}, P, S)$, onde P consiste de:

$S ::= aAS \mid a$

$A ::= SbA \mid SS \mid ba$

A figura seguinte é uma árvore de derivação para uma sentença desta gramática, na qual os nodos foram numerados para facilitar a explicação:



Os vértices 7. O vértice 1 tem o da esquerda para a direita

interiores são 1, 3, 4, 5 e rótulo S, e seus filhos, têm rótulos a, A e S.

Note que $S ::= aAS$ é uma produção. Da mesma forma, o vértice 3 tem rótulo A, e os rótulos de seus filhos são S, b e A, da esquerda para a direita, sendo que $A ::= SbA$ é também uma produção. Os vértices 4 e 5 têm ambos o rótulo S. O único filho deles tem rótulo a, e $S ::= a$ também é uma produção. Finalmente, o vértice 7 tem rótulo A, e seus filhos, da esquerda para a direita, são b e a, e $A ::= ba$ também é produção. Assim, as condições para que esta árvore seja uma árvore de derivação para G foram cumpridas.

Pode-se estender a ordem “da esquerda para a direita” dos filhos para produzir uma ordem da esquerda para a direita de todas as folhas. No exemplo acima, o caminhar da esquerda para a direita nas folhas da árvore produziria a seguinte seqüência: 2, 9, 6, 10, 11 e 8.

Pode-se ver que a árvore de derivação é uma descrição natural para a derivação de uma forma sentencial particular da gramática G .

Uma sub-árvore de uma árvore de derivação é composta por um vértice particular da árvore, juntamente com seus descendentes.

4. Derivação mais à Esquerda e mais à Direita

Se cada passo na produção de uma derivação é aplicado na variável mais à esquerda, então a derivação é chamada derivação mais à esquerda. Similarmente, uma derivação onde a cada passo a variável mais à direita é substituída, é chamada de derivação mais à direita.

Se w está em $L(G)$, então w tem ao menos uma árvore de derivação. Além disso, em relação a uma árvore de derivação particular, w tem uma única derivação mais à esquerda e uma única derivação mais à direita.

Evidentemente, w pode ter várias derivações mais à esquerda e várias derivações mais à direita, já que pode haver mais de uma árvore de derivação para w . Entretanto, é fácil mostrar que, para cada árvore de derivação, apenas uma derivação mais à esquerda e uma derivação mais à direita podem ser obtidas.

Exemplo:

A derivação mais à esquerda correspondente à árvore do exemplo anterior é:

$S \rightarrow aAS \rightarrow aSbAS \rightarrow aabAS \rightarrow aabbaS \rightarrow aabbbaa$

IV - Linguagens Livres de Contexto

A derivação mais à direita correspondente é:
 $S \rightarrow aAS \rightarrow aAa \rightarrow aSbAa \rightarrow aSbbaa \rightarrow aabbaa$

5. Ambigüidade

Se uma gramática G tem mais de uma árvore de derivação para uma mesma sentença, então G é chamada de *gramática ambígua*. Uma linguagem livre de contexto para a qual toda gramática livre de contexto é ambígua é denominada *linguagem livre de contexto inerentemente ambígua*.

6. Simplificações de Gramáticas Livres de Contexto

Existem várias maneiras de restringir as produções de uma gramática livre de contexto sem reduzir seu poder expressivo. Se L é uma linguagem livre de contexto não-vazia, então L pode ser gerada por uma gramática livre de contexto G com as seguintes propriedades:

- Cada variável e cada terminal de G aparecem na derivação de alguma palavra de L .
- Não há produções da forma $A ::= B$, onde A e B são variáveis.

Além disso, se $\varepsilon \notin L$, então não há necessidade de produções da forma $A ::= \varepsilon$.

6.1. Símbolos Inúteis

Pode-se eliminar os símbolos inúteis de uma gramática sem prejudicar seu potencial expressivo. Um símbolo X é útil se existe uma derivação $S \rightarrow^* \alpha X \beta \rightarrow^* w$, para algum w , α e β , onde w é uma cadeia de T^* e α e β são cadeias quaisquer de variáveis e terminais. Caso contrário, o símbolo X é inútil. Tanto terminais quanto não-terminais podem ser úteis ou inúteis. Se o símbolo inicial for inútil, então a linguagem definida pela gramática é vazia.

Há dois tipos de símbolos inúteis: aqueles que não geram nenhuma cadeia de terminais e aqueles que jamais são gerados a partir de S . O primeiro caso corresponde aos *símbolos improdutivo*s (ou *mortos*), e o o segundo caso corresponde aos *símbolos inalcançáveis*. Um terminal sempre é produtivo, mas pode ser inalcançável. Já o não-terminal pode ser tanto improdutivo quanto inalcançável, mas o símbolo inicial sempre é alcançável. A seguir serão vistos algoritmos para eliminar tanto o primeiro quanto o segundo tipo de símbolos inúteis.

O algoritmo para eliminação dos símbolos improdutivo é baseado na idéia de que se um não-terminal A tem uma produção consistindo apenas de símbolos produtivos, então o próprio A é produtivo.

Algoritmo: Eliminação de Símbolos Improdutivo

Entrada: Uma GLC $G=(N,T,P,S)$

Saída: Uma GLC $G'=(N',T,P',S)$, sem símbolos improdutivo.

$SP := T \cup \{\varepsilon\}$

Repita

$Q := \{X \mid X \in N \text{ e } X \notin SP \text{ e (existe pelo menos uma produção } X ::= X_1X_2\dots X_n \text{ tal que } X_1 \in SP, X_2 \in SP, \dots, X_n \in SP)\};$

$SP := SP \cup Q;$

Até $Q = \emptyset;$

$N' = SP \cap N;$

Se $S \in SP$ Então

$P' := \{p \mid p \in P \text{ e todos os símbolos de } p \text{ pertencem a } SP\}$

Senão " $L(G) = \emptyset$ ";

$P' := \emptyset;$

Fim Se

Fim.

No seguinte exemplo, para facilitar o acompanhamento do algoritmo, os símbolos do conjunto SP serão sublinhados a cada passo.

IV - Linguagens Livres de Contexto

Seja G a seguinte gramática:

$S ::= ASB \mid BSA \mid SS \mid aS \mid \varepsilon$

$A ::= ABS \mid B$

$B ::= BSSA \mid A$

Inicialmente $SP = \{a\}$, e são marcadas as produções:

$S ::= \underline{a}S$

$S ::= \underline{\varepsilon}$

Uma vez que o lado direito da produção $S ::= \varepsilon$ está completamente marcado, deve-se marcar todas as ocorrências de S . Como S é o único novo símbolo a ser marcado, então $Q := \{S\}$. As marcações resultam:

$\underline{S} ::= \underline{A}\underline{S}\underline{B} \mid \underline{B}\underline{S}\underline{A} \mid \underline{S}\underline{S} \mid \underline{a}\underline{S} \mid \underline{\varepsilon}$

$A ::= \underline{A}\underline{B}\underline{S} \mid \underline{B}$

$B ::= \underline{B}\underline{S}\underline{S}\underline{A} \mid \underline{A}$

Neste momento $SP = \{a, S\}$. Agora $Q := \emptyset$, e assim a repetição do algoritmo termina com $SP = \{a, S\}$, os únicos símbolos produtivos.

Para simplificar a gramática, os símbolos improdutivos podem ser removidos, bem como todas as produções de P que contenham estes símbolos. Assim, a gramática simplificada para o exemplo seria:

$S ::= SS \mid aS \mid \varepsilon$

O segundo tipo de símbolos inúteis são aqueles que jamais são gerados a partir de S . Estes símbolos, chamados de inalcançáveis, são eliminados pelo seguinte algoritmo:

Algoritmo: Eliminação de Símbolos Inalcançáveis

Entrada: Uma GLC $G=(N,T,P,S)$

Saída: Uma GLC $G'=(N',T,P',S)$, sem símbolos inalcançáveis.

$SA := \{S\}$

Repita

$M := \{X \mid X \in N \cup T \text{ e } X \notin SA \text{ e existe uma produção } Y ::= \alpha X \beta \text{ e } Y \in SA\};$

$SA := SA \cup M;$

Até $M = \emptyset;$

$N' = SA \cap N;$

$T' = SA \cap T;$

$P' := \{p \mid p \in P \text{ e todos os símbolos de } p \text{ pertencem a } SA\};$

Fim.

Para exemplificar, seja G :

$S ::= aS \mid SB \mid SS \mid \varepsilon$

$A ::= ASB \mid c$

$B ::= b$

Inicialmente, $SA = \{S\}$ e $M = \{S\}$. Assim, se obtém:

$S ::= \underline{aS} \mid \underline{SB} \mid \underline{SS} \mid \varepsilon$

e a e B são os únicos símbolos recém marcados que não estão em SA . Assim, a M é atribuído o conjunto $\{a, B\}$ e a SA é atribuído $\{S, a, B\}$. Em seguida se obtém:

$B ::= \underline{b}$

e assim M passa a ser igual a $\{b\}$ e SA igual a $\{S, a, B, b\}$. Já que M não contém não-terminais, a execução do laço deixará $M = \emptyset$ e SA não será alterado. Assim, S, a, B e b são alcançáveis e A e c são inalcançáveis e suas produções podem ser eliminadas da gramática, resultando:

$S ::= aS \mid SB \mid SS \mid \varepsilon$

$B ::= b$

6.2. Eliminação de ε -Produções

Conforme foi dito anteriormente, uma GLC pode ter produções do tipo $A ::= \varepsilon$. Mas toda GLC pode ser transformada em uma GLC equivalente sem este tipo de produções (chamadas ε -produções), com exceção da produção $S ::= \varepsilon$, se esta existir (ou seja, a cadeia vazia pertence à linguagem). Assim procedendo, é possível mostrar que toda GLC pode obedecer à restrição das GSC (tipo 1). É o método da eliminação de ε -produções, da qual se tratará agora.

O método consiste em determinar, para cada variável A em N , se $A \rightarrow^* \varepsilon$. Se isto for verdade, se diz que a variável A é anulável. Pode-se assim substituir cada produção da forma $B ::= X_1 X_2 \dots X_n$ por todas as produções formadas pela retirada de uma ou mais variáveis X_i anuláveis.

Exemplo:

$A ::= BCDe$

$B ::= \varepsilon \mid e$

$C ::= \varepsilon \mid a$

$D ::= b \mid cC$

Neste caso, as variáveis anuláveis são B e C , assim a gramática pode ser transformada na seguinte:

$A ::= BCDe \mid CDe \mid BDe \mid De$

$B ::= e$

IV - Linguagens Livres de Contexto

$C ::= a$

$D ::= b \mid c \mid cC$

Os não-terminais que derivam a sentença ε são chamados de ε -não-terminais. Um não-terminal A é um ε -não-terminal se existir uma derivação $A \rightarrow^* \varepsilon$ em G . Note que ε está em $L(G)$ se e somente se S é um ε -não-terminal. Se G não tem ε -não-terminais, ela é dita ε -livre.

Se G tem ε -produções, então em uma derivação sentencial da forma:

$S \rightarrow \alpha_1 \rightarrow \alpha_2 \rightarrow \dots \rightarrow \alpha_n$, os tamanhos das sentenças irão variar não-monotonicamente um em relação ao outro, isto é, ora aumentam, ora diminuem. Entretanto, se G não tem ε -produções, então:

$$|S| \leq |\alpha_1| \leq |\alpha_2| \leq \dots \leq |\alpha_n|$$

isto é, os tamanhos das sentenças são monotonicamente crescentes. Esta propriedade é útil quando se quer testar se uma dada palavra é ou não gerada por uma GLC.

Antes de apresentar um algoritmo para eliminar as ε -produções, será apresentado um algoritmo para encontrar os ε -não-terminais.

Algoritmo: Encontrar o Conjunto dos ε -não-terminais

Entrada: Uma GLC $G=(N,T,P,S)$

Saída: O conjunto E dos ε -não-terminais.

$E := \{\varepsilon\}$

Repita

$Q := \{X \mid X \in N \text{ e } X \notin E \text{ e existe pelo menos uma produção } X ::= Y_1 Y_2 \dots Y_n \text{ tal que } Y_1 \in E, Y_2 \in E, \dots, Y_n \in E\};$

$E := E \cup Q;$

Até $Q = \emptyset;$

Fim.

Seja G a seguinte gramática:

$S ::= aS \mid SS \mid bA$

$A ::= BC$

$B ::= CC \mid ab \mid aAbC$

$C ::= \varepsilon$

Inicialmente temos $E = \{\varepsilon\}$.

Quando o laço é executado pela primeira vez, $Q = \{C\}$, e se obtém

$S ::= aS \mid SS \mid bA \quad A ::= \underline{BC} \quad B ::= \underline{CC} \mid ab \mid aAb\underline{C} \quad C ::= \underline{\varepsilon}$

Como Q não é vazio, temos $E = \{\varepsilon, C\}$ e uma segunda iteração, que deixa $Q = \{B\}$, e

obtem-se:

$S ::= aS \mid SS \mid bA \quad A ::= \underline{BC} \quad B ::= \underline{CC} \mid ab \mid aAb\underline{C} \quad C ::= \underline{\varepsilon}$

Novamente Q não é vazio, temos $E = \{\varepsilon, B, C\}$, e numa nova iteração, que faz $Q = \{A\}$

$S ::= aS \mid SS \mid b\underline{A} \quad A ::= \underline{BC} \quad B ::= \underline{CC} \mid ab \mid a\underline{A}b\underline{C} \quad C ::= \underline{\varepsilon}$

O símbolo A é acrescentado a E , que passa a ser $E = \{\varepsilon, A, B, C\}$. Na próxima iteração não são acrescentados símbolos a Q , terminando assim o algoritmo, e temos que os ε -não-terminais em G são A, B e C .

Para eliminar os ε -não-terminais de uma GLC, pode-se usar o seguinte algoritmo, que elimina ε -não-terminais sem introduzir novos. A estratégia é baseada na seguinte idéia: Seja A um ε -não-terminal em G . Então ele é dividido conceitualmente em dois não-terminais A' e A'' , tal que A' gera todas as cadeias não-vazias e A'' gera apenas ε . Agora, do lado direito de cada produção onde A aparece uma vez, por exemplo $B ::= \alpha A \beta$, é trocado por duas produções $B ::= \alpha A' \beta$ e $B ::= \alpha A'' \beta$. Já que A'' só gera ε , ele pode ser trocado por ε (ou seja, eliminado), deixando $B ::= \alpha \beta$. Depois disso, pode-se usar A em lugar de A' . A produção final fica:

IV - Linguagens Livres de Contexto

$B ::= \alpha\beta \mid \alpha A\beta$, onde A não é mais um ε -não-terminal. Se $B ::= \alpha A\beta A\gamma$, então são obtidas quatro combinações possíveis pelas trocas de A por ε ($\alpha A\beta A\gamma$, $\alpha A\beta\gamma$, $\alpha\beta A\gamma$, $\alpha\beta\gamma$) e assim por diante.

Algoritmo: Eliminar Todos os ε -não-terminais

Entrada: Uma GLC $G=(N,T,P,S)$

Saída: Uma GLC $G=(N',T,P',S')$ ε -livre.

Construa o conjunto E

$P' := \{p \mid p \in P \text{ e } p \text{ não é } \varepsilon\text{-produção}\}$

Repita

Se P' tem uma produção da forma $A ::= \alpha B\beta$, tal que $B \in E$, $\alpha\beta (N \cup T)^*$ e $\alpha\beta \neq \varepsilon$, então inclua a produção $A ::= \alpha\beta$ em P' ;

Até que nenhuma nova produção possa ser adicionada a P' ;

Se $S \in E$, então

Adicione a P' as produções $S' ::= S \mid \varepsilon$;

$N' := N \cup \{S'\}$;

Senão

$S' := S$;

$N' := N$;

Fim Se

Fim.

Exemplo:

Dada a seguinte GLC:

$S ::= SaBC \mid ABb \mid BC$

$A ::= bC \mid aBC \mid a$

$B ::= BCA \mid bC \mid \varepsilon$

$C ::= ab \mid \varepsilon$

O conjunto dos ε -não-terminais é $E = \{S, B, C\}$. A gramática ε -livre equivalente é:

$S' ::= S \mid \varepsilon$

$S ::= SaBC \mid ABb \mid BC \mid aBC \mid SaC \mid SaB \mid Ab \mid C \mid B \mid aC \mid aB \mid Sa \mid a$

$A ::= bC \mid aBC \mid a \mid b \mid aB \mid aC$

$B ::= BCA \mid bC \mid CA \mid BA \mid b \mid A$

$C ::= ab$

6.3. Produções Unitárias

Uma produção da forma $A ::= \alpha$ em uma GLC $G=(N,T,P,S)$ é chamada de *produção unitária* se α é um não-terminal. Um caso especial de produção unitária é a produção $A ::= A$, também chamada *produção circular*. Tal produção pode ser removida imediatamente sem afetar a capacidade de geração da gramática. O algoritmo para eliminar produções unitárias assume que as produções circulares já tenham sido eliminadas anteriormente. Essa eliminação é trivial, pois a identificação de produções circulares é bastante simples.

Algoritmo: Eliminar Produções Unitárias

Entrada: Uma GLC $G=(N,T,P,S)$ sem produções circulares.

Saída: Uma GLC $G=(N,T,P',S)$ sem produções unitárias.

Para toda $A \in N$ faça:

$N_A := \{B \mid A \rightarrow^* B, \text{ com } B \in N\}$;

Fim Para

IV - Linguagens Livres de Contexto

$P' := \emptyset$;
 Para toda produção $B ::= \alpha \in P$ faça:
 Se $B ::= \alpha$ não é uma produção unitária, então:
 $P' := P' \cup \{A ::= \alpha \mid B \in N_A\}$;
 Fim Se
 Fim Para
 Fim.

Podem surgir símbolos inalcançáveis depois da aplicação deste algoritmo. Isso se deve ao fato de que o símbolo é substituído por suas produções, e se o mesmo aparecer somente em produções unitárias, o mesmo irá desaparecer do lado direito das produções. Para exemplificar, o algoritmo é aplicado na seguinte gramática:

$S ::= aB \mid aA \mid A$ $B ::= aA \mid B \mid A$
 $A ::= C \mid aaa$ $C ::= \varepsilon \mid cBc$

Primeiramente é necessário determinar o conjunto de N_K para cada não terminal K :

$N_S = \{S, A, B, C\}$ $N_B = \{B, A, C\}$ $N_A = \{A, C\}$ $N_C = \{C\}$

As produções de cada símbolo K são as produções $R ::= \alpha$, onde R é um elemento de N_K e α é uma produção não-unitária. No exemplo, as produções de S são as produções não unitárias de S , A , B e C . Assim, a nova GLC sem produções unitárias é

$S ::= aB \mid aA \mid aaa \mid cBc \mid \varepsilon$ $B ::= aA \mid aaa \mid cBc \mid \varepsilon$
 $A ::= aaa \mid \varepsilon$ $C ::= \varepsilon \mid cBc$

Nesta nova GLC, o não-terminal C é inalcançável.

7. Fatoração

Uma GLC está *fatorada* se ela é determinística, ou seja, se ela não possui produções para um mesmo não-terminal no lado esquerdo cujo lado direito inicie com a mesma cadeia de símbolos ou com símbolos que derivam seqüências que iniciem com a mesma cadeia de símbolos. De acordo com esta definição, a seguinte gramática é não-determinística:

$S ::= aSB \mid aSA$ fonte do
 $A ::= a$ não-determinismo
 $B ::= b$

Para fatorar uma GLC, deve-se alterar as produções envolvidas no não-determinismo da seguinte forma:

a) As produções com não-determinismo direto da forma $A ::= \alpha\beta \mid \alpha\gamma$ serão substituídas por:

$A ::= \alpha A'$
 $A' ::= \beta \mid \gamma$

Na gramática acima, temos $\alpha = aS$, assim, com a eliminação do não-determinismo, teríamos a seguinte gramática:

$S ::= aSS'$
 $S' ::= B \mid A$
 $A ::= a$
 $B ::= b$

b) o não-determinismo indireto é retirado através de sua transformação em não-determinismo direto (através de derivações sucessivas) e posterior eliminação segundo o item (a). Tomemos por exemplo a seguinte gramática:

$S ::= AC \mid BC$ $C ::= eC \mid eA$
 $A ::= aD \mid cC$ $D ::= fD \mid AB$
 $B ::= aB \mid dD$

IV - Linguagens Livres de Contexto

onde o não-determinismo indireto se dá pelo fato de que A e B podem iniciar com o terminal a. Assim, transformaremos primeiro o não-determinismo indireto em não-determinismo direto, substituindo os não-terminais A e B pelas suas produções em S:

$$\begin{aligned} S &::= aCC \mid cBC \mid aBC \mid dDC & C &::= eC \mid eA \\ A &::= aD \mid cB & D &::= fD \mid AB \\ B &::= aB \mid dD \end{aligned}$$

O próximo passo é simples, bastando agora eliminar o não-determinismo direto, segundo a regra do item a, ou seja:

$$\begin{aligned} S &::= aS' \mid cBC \mid dDC & C &::= eC' \\ S' &::= CC \mid BC & C' &::= C \mid A \\ A &::= aD \mid cB & D &::= fD \mid AB \\ B &::= aB \mid dD \end{aligned}$$

8. Eliminação de Recursão à Esquerda

Um não-terminal A, em uma GLC $G=(N,T,P,S)$ é *recursivo* se $A \rightarrow \alpha A \beta$, para α e $\beta \in (N \cup T)^*$.

Se $\alpha = \varepsilon$, então A é *recursivo à esquerda*; se $\beta = \varepsilon$, então A é *recursivo à direita*. Esta recursividade pode ser direta ou indireta.

Uma gramática com pelo menos um não-terminal recursivo à esquerda ou à direita é uma *gramática recursiva à esquerda* ou *à direita*, respectivamente.

Uma GLC $G=(N,T,P,S)$ possui *recursão à esquerda direta* se P contém pelo menos uma produção da forma $A ::= A\alpha$.

Uma GLC $G=(N,T,P,S)$ possui *recursão à esquerda indireta* se existe em G uma derivação da forma $A \rightarrow^n A\beta$, para algum $n \geq 2$.

Para eliminar as recursões diretas à esquerda nas produções:

$$A ::= A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_n \mid \dots \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_m$$

onde nenhum β_i começa com A, deve-se substituir estas produções pelas seguintes:

$$\begin{aligned} A &::= \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_m A' \\ A' &::= \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_n A' \mid \varepsilon \end{aligned}$$

onde A' é um novo não-terminal.

A seguir é apresentado um algoritmo para eliminar recursões diretas e indiretas à esquerda.

Algoritmo: Eliminar Recursões à Esquerda

Entrada: Uma GLC $G=(N,T,P,S)$.

Saída: Uma GLC $G=(N',T,P',S)$ sem recursão à esquerda.

$N' := N$;

$P' := P$;

Ordene os não-terminais de N' em uma ordem qualquer (por exemplo, $A_1, A_2, A_3, \dots, A_n$);

Para i de 1 até n, faça:

 Para j de 1 até i-1 faça

 Substitua as produções de P' da forma $A_i ::= A_j \gamma$ por produções da forma $A_i ::= \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$, onde $\delta_1, \delta_2, \dots, \delta_k$ são os lados direitos das produções com lado esquerdo A_j , ou seja, $A_j ::= \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$.

 Fim Para;

 Elimine as recursões diretas das produções de P' com lado esquerdo A_i ;

Fim Para;

Fim.

IV - Linguagens Livres de Contexto

Note que este algoritmo já elimina as recursões à esquerda diretas e indiretas. Para exemplificar, considere a seguinte gramática:

$$\begin{aligned} S &::= Sb \mid Cc \mid Ab \\ A &::= Sc \mid ab \\ C &::= Scd \mid Cba \mid b \end{aligned}$$

Primeiramente, ($i = 1$), a recursão a esquerda de S é eliminada:

$$\begin{aligned} S &::= CcS' \mid AbS' \\ A &::= Sc \mid ab \\ C &::= Scd \mid Cba \mid b \\ S' &::= bS' \mid \varepsilon \end{aligned}$$

No segundo passo ($i = 2$), temos que substituir S nas produções de A que começam por S :

$$\begin{aligned} S &::= CcS' \mid AbS' \\ A &::= CcS'c \mid AbS'c \mid ab \\ C &::= Scd \mid Cba \mid b \\ S' &::= bS' \mid \varepsilon \end{aligned}$$

Em seguida, eliminar a recursão a esquerda direta:

$$\begin{aligned} S &::= CcS' \mid AbS' \\ A &::= CcS'cA' \mid abA' \\ C &::= Scd \mid Cba \mid b \\ S' &::= bS' \mid \varepsilon \\ A' &::= bS'cA' \mid \varepsilon \end{aligned}$$

No terceiro passo ($i = 3$), temos que substituir S ($j = 1$) nas produções de C que começam por S :

$$\begin{aligned} S &::= CcS' \mid AbS' \\ A &::= CcS'cA' \mid abA' \\ C &::= CcS'cd \mid AbS'cd \mid Cba \mid b \\ S' &::= bS' \mid \varepsilon \\ A' &::= bS'cA' \mid \varepsilon \end{aligned}$$

Em seguida substituir A ($j = 2$) nas novas produções de C que começam por A :

$$\begin{aligned} S &::= CcS' \mid AbS' \\ A &::= CcS'cA' \mid abA' \\ C &::= CcS'cd \mid CcS'cA'bS'cd \mid abA'bS'cd \mid Cba \mid b \\ S' &::= bS' \mid \varepsilon \\ A' &::= bS'cA' \mid \varepsilon \end{aligned}$$

Por fim, eliminar as recursões a esquerda diretas:

$$\begin{aligned} S &::= CcS' \mid AbS' \\ A &::= CcS'cA' \mid abA' \\ C &::= abA'bS'cdC' \mid bC' \\ S' &::= bS' \mid \varepsilon \\ A' &::= bS'cA' \mid \varepsilon \\ C' &::= cS'cdC' \mid cS'cA'bS'cd \mid baC' \mid \varepsilon \end{aligned}$$

9. Tipos Especiais de GLC

A seguir serão identificados alguns tipos especiais de GLC

a) Gramática Própria: Uma GLC é própria se:

- a.1) Não possui produções cíclicas (ver definição seguinte);
- a.2) É ε -livre;
- a.3) Não possui símbolos inúteis.

IV - Linguagens Livres de Contexto

b) Gramática Sem Ciclos: $G(N,T,P,S)$ é uma GLC sem ciclos (ou livre de ciclos) se não existe em G nenhuma derivação da forma $A \rightarrow^+ A$, para todo $A \in N$.

c) Gramática Reduzida: Uma GLC $G=(N,T,P,S)$ é uma GLC reduzida se:

c.1) $L(G)$ não é vazia;

c.2) Se $A ::= \alpha \in P$ então $A \neq \alpha$

c.3) G não possui símbolos inúteis

d) Gramática de Operadores: Uma GLC $G=(N,T,P,S)$ é de operadores se ela não possui produções cujo lado direito contenha não-terminais consecutivos.

e) Gramática Unicamente Inversível: Uma GLC reduzida é unicamente inversível se ela não possui produções com lados direitos iguais.

f) Gramática Linear: Uma GLC $G=(N,T,P,S)$ é linear se todas as suas produções forem da forma $A ::= xBw \mid x$, onde A e B pertencem a N , e x e w pertencem a T^* .

g) Forma Normal de Chomsky: Uma GLC está na forma normal de Chomsky se ela é ϵ -livre, e todas as suas produções (exceto, possivelmente, $S ::= \epsilon$) são da forma:

g.1) $A ::= BC$, com A, B e $C \in N$, ou

g.2) $A ::= a$; com $A \in N$ e $a \in T$.

h) Forma Normal de Greinbach: Uma GLC está na forma normal de Greinbach se ela é ϵ -livre e todas as suas produções (exceto, possivelmente, $S ::= \epsilon$) são da forma $A ::= a\alpha$ tal que $a \in T, \alpha \in N^*$ e $A \in N$.

10 Principais Notações de GLC

A BNF (*Backus Naur Form*) é uma notação utilizada na especificação formal da sintaxe de linguagens de programação. Esta é a forma de especificação de gramáticas que tem sido utilizada neste contexto. As gramáticas a seguir são exemplos de gramáticas descritas na notação BNF:

$\langle S \rangle ::= a \langle S \rangle \mid \epsilon$

$\langle E \rangle ::= \langle E \rangle + id \mid id$

A BNFE (Backus Naur Form Extended) é equivalente à BNF, mas permite uma especificação mais compacta da sintaxe de uma linguagem de programação. Os símbolos entre chaves abreviam os fechados das cadeias que eles representam. As mesmas gramáticas descritas acima em BNF podem ser representadas em BNFE da seguinte forma:

$\langle S \rangle ::= \{a\}$

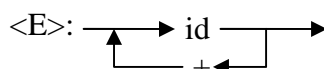
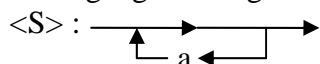
$\langle E \rangle ::= id\{+id\}$

A RRP é uma notação de GLC onde o lado direito das produções é especificado através de expressões regulares, envolvendo os símbolos de N e T de uma gramática. As gramáticas dos exemplos acima podem ser descritas em RRP da seguinte forma:

$\langle S \rangle ::= a^*$

$\langle E \rangle ::= id(+id)^*$

O diagrama sintático é uma notação gráfica para GLC's bastante utilizada em análise de linguagem natural e em manuais de linguagens de programação, para descrever graficamente a sintaxe da linguagem. As gramáticas exemplo acima teriam os seguintes diagramas sintáticos:



11. Conjuntos First e Follow

11.1 Conjunto First

Seja α uma forma sentencial qualquer gerada por G . $FIRST(\alpha)$ será o conjunto de símbolos terminais que podem iniciar (que podem aparecer na posição mais à esquerda das sentenças derivadas desta forma sentencial) α ou seqüências derivadas (direta ou indiretamente) de α . Além disso, se $\alpha = \epsilon$ ou $\alpha \rightarrow^+ \epsilon$ então $\epsilon \in FIRST(\alpha)$.

Para calcular $FIRST(X)$ para todo $X \in (N \cup T)^*$, vamos considerar os seguintes casos especiais:

- a) $X = \epsilon$. Neste caso, $FIRST(X) = \{\epsilon\}$.
- b) $X \in T$. Neste caso, $FIRST(X) = \{X\}$.
- c) $X \in N$. Neste caso, se $X ::= \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$ são todas as produções com lado esquerdo X , então $FIRST(X) = FIRST(\alpha_1) \cup FIRST(\alpha_2) \cup \dots \cup FIRST(\alpha_n)$.

O algoritmo genérico do cálculo do conjunto $FIRST$ para uma cadeia qualquer é recursivo e utiliza as definições para os casos especiais acima.

Algoritmo: Cálculo do Conjunto First de uma Cadeia

Entrada: Uma cadeia de símbolos $X \in (N \cup T)^*$ de uma gramática $G=(N,T,P,S)$.

Uma tabela $TABFIRST$ com o conjunto $FIRST$ de cada não-terminal de G .

Saída: O conjunto $FIRST(X)$, denotado por F .

Se $X = \epsilon$ então:

$F := \{\epsilon\}$.

Senão se $X \in T$ então:

$F = \{X\}$

Senão se $X \in N$ então:

Sejam $X ::= \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$ as produções com lado esquerdo X , então

$F := FIRST(\alpha_1) \cup FIRST(\alpha_2) \cup \dots \cup FIRST(\alpha_n)$;

Senão

$i := 0$;

$F := \emptyset$;

Repita:

$i := i+1$;

$X_i :=$ o i -ésimo símbolo da cadeia X ;

$F := F \cup TABFIRST(X_i) - \{\epsilon\}$;

Se $X = \epsilon$ então:

$F := F \cup \{\epsilon\}$;

Fim Se

Até $X \in T$ ou $X = \epsilon$ ou $(X \in N \text{ e } \epsilon \notin TABFIRST(X))$

Fim Se

Para determinar o conjunto $FIRST$ de cada um dos não-terminais de uma gramática, procede-se como segue:

1. Crie uma tabela $TABFIRST$ cujas linhas são rotuladas com os não-terminais da gramática e com uma coluna rotulada como “ $FIRST$ ”, que irá conter os símbolos pertencentes ao $FIRST$ de cada não-terminal. Inicialize cada uma das posições desta tabela com \emptyset (conjunto vazio).
2. Para cada um dos não-terminais, calcule o seu conjunto $FIRST$ utilizando o algoritmo acima.

IV - Linguagens Livres de Contexto

3. Repita o passo 2 até que não ocorra mais nenhuma modificação na tabela.

Para exemplificar, considere a gramática:

$S ::= ABS \mid aA$

$A ::= \varepsilon \mid a$

$B ::= Bb \mid cd$

Inicialmente constrói-se a tabela:

N	FIRST
S	\emptyset
A	\emptyset
B	\emptyset

Calculando o FIRST de S pelo algoritmo acima, faz-se $FIRST(S) = FIRST(ABS) \cup FIRST(aA)$. Como $FIRST(A) = \emptyset$, na tabela, então $FIRST(S)$ inicialmente é $\emptyset \cup \{a\}$. A tabela fica:

N	FIRST
S	{a}
A	\emptyset
B	\emptyset

Calculando o FIRST(A) temos $\{\varepsilon, a\}$. Então a tabela fica:

N	FIRST
S	{a}
A	{ ε , a}
B	\emptyset

Calculando o FIRST(B) temos $FIRST(B) = FIRST(Bb) \cup FIRST(cd) = \emptyset \cup \{c\}$. A tabela fica:

N	FIRST
S	{a}
A	{ ε , a}
B	{c}

Como a tabela foi modificada, o passo 2 deve ser repetido. Desta vez, ao calcular o FIRST de S teremos: $FIRST(S) = FIRST(ABS) \cup FIRST(aA) = \{a\} \cup FIRST(BS) \cup \{a\} = \{a\} \cup \{c\} \cup \{a\} = \{a, c\}$. Logo a tabela fica:

N	FIRST
S	{a, c}
A	{ ε , a}
B	{c}

O cálculo do FIRST de A e B não modifica a tabela. Mas como houve uma modificação no cálculo do FIRST de S, o passo 2 deve ser novamente executado. Como desta vez não haverá modificações, a última tabela será de fato a definitiva.

11.2. Conjunto Follow

$FOLLOW(A)$ é definido para todo $A \in N$ como sendo o conjunto de símbolos terminais que podem aparecer imediatamente após A em alguma forma sentencial de G. A seguir, é definido um algoritmo que constrói o FOLLOW de todos os não-terminais simultaneamente.

Algoritmo: Cálculo do Conjunto Follow de um Não-Terminal

Entrada: Uma gramática $G=(N,T,P,S)$.

Saída: Uma tabela TABFOLLOW com o conjunto FOLLOW de cada não-terminal de G.

IV - Linguagens Livres de Contexto

Para todo $X \in N$, exceto S :

$TABFOLLOW(X) := \emptyset$;

Fim Para

$TABFOLLOW(S) := \{\$ \}$; (marca de final de sentença)

Para toda ocorrência de um não-terminal X em uma produção $Y ::= \alpha X \beta \in P$ com $\beta \neq \epsilon$:

$TABFOLLOW(X) := TABFOLLOW(X) \cup FIRST(\beta) - \{\epsilon\}$

Fim Para;

Repita:

Para cada produção da forma $X ::= Y_1 Y_2 \dots Y_n$ faça:

$i ::= n+1$;

Repita:

$i ::= i-1$;

Se $Y_i \in N$ então:

$TABFOLLOW(Y_i) := TABFOLLOW(X) \cup TABFOLLOW(Y_i)$

Fim Se

Até $Y_i \notin N$ ou $\epsilon \notin TABFIRST(Y_i)$

Fim Para

Até que não haja mais modificações na tabela $TABFOLLOW$

Seguem duas observações importantes:

a) A função $FIRST$ é definida para cadeias de símbolos, $FOLLOW$ só é definida para não-terminais.

b) Os elementos de um conjunto $FIRST$ pertencem ao conjunto $T \cup \{\epsilon\}$; os elementos de um conjunto $FOLLOW$ pertencem ao conjunto $T \cup \{\$ \}$.

Como exemplo, será calculado o conjunto $FOLLOW$ da gramática utilizada como exemplo no cálculo do $FIRST$, que é:

$S ::= ABS \mid aA$

$A ::= \epsilon \mid a$

$B ::= Bb \mid cd$

cuja tabela $FIRST$ é:

N	FIRST
S	{a, c}
A	{ ϵ , a}
B	{c}

O tabela $FOLLOW$ dos não-terminais da gramática é inicializada como:

N	FOLLOW
S	{ $\$$ }
A	\emptyset
B	\emptyset

Para o primeiro passo, temos as ocorrências de $Y ::= \alpha X \beta$ em $S ::= ABS$ e $B ::= Bb$, assim:

$FOLLOW(A) = FOLLOW(A) \cup FIRST(BS) = \emptyset \cup \{c\} = \{c\}$

$FOLLOW(B) = FOLLOW(B) \cup FIRST(S) = \emptyset \cup \{a, c\} = \{a, c\}$

$FOLLOW(b) = FOLLOW(B) \cup FIRST(b) = \{a, c\} \cup \{b\} = \{a, b, c\}$

Temos então a tabela $FOLLOW$ intermediária:

N	FOLLOW
S	{ $\$$ }
A	{c}
B	{a, b, c}

IV - Linguagens Livres de Contexto

No segundo passo, as produções são analisadas da esquerda para a direita. A produção $S ::= ABS$ faz com que:

$$\text{FOLLOW}(S) = \text{FOLLOW}(S) \cup \text{FOLLOW}(S) = \{\$\}$$

Como $\varepsilon \notin \text{FIRST}(S)$, a segunda regra termina para esta produção. Para a produção $S ::= aA$, temos:

$$\text{FOLLOW}(A) ::= \text{FOLLOW}(A) \cup \text{FOLLOW}(S) = \{c, \$\}$$

Como $\varepsilon \in \text{FIRST}(A)$, o símbolo a esquerda de A é analisado. Como é um terminal (a), o algoritmo pára. Já que não há mais produções que tenham não-terminais a direita, o cálculo da tabela FOLLOW termina. A tabela final fica:

N	FOLLOW
S	{\\$}
A	{c, \\$}
B	{a, b, c}

EXERCÍCIOS

1. Seja G a seguinte GLC:

$$S ::= aScSa \mid \varepsilon$$

- Mostre três palavras (sentenças) em $L(G)$ e alguma árvore de derivação para cada uma delas.
- Mostre todas as árvores sintáticas para a palavra “aacacaacacaa”.

2. Seja G a seguinte GLC:

$$S ::= SS \mid aSf \mid af \mid bSe \mid be \mid cSd \mid cd$$

mostre todas as árvores de derivação para a palavra “abafcdefbcbdde”.

3. Seja G uma GLC:

$$S ::= ABB \mid CAC$$

$$A ::= a$$

$$B ::= Bc \mid ABB$$

$$C ::= bB \mid a$$

encontre uma GLC reduzida G' equivalente a G retirando os símbolos improdutivos e inalcançáveis.

4. Dada uma GLC:

$$S ::= aSASb \mid Saa \mid AA$$

$$A ::= caA \mid Ac \mid bca \mid \varepsilon$$

- Quais são os ε -não-terminais de G?
- $\varepsilon \in L(G)$
- Construa uma GLC ε -livre equivalente a G (aplicando o algoritmo).

5. Considere a seguinte GLC:

$$S ::= B \mid aB$$

$$A ::= B \mid aB$$

$$B ::= C \mid aC$$

$$C ::= b$$

- Quais são as produções unitárias de G?
- Construa uma GLC equivalente livre de produções unitárias.

6. A linguagem $L(G) = \{ a^i b^j c^m d^n \mid (i=j \text{ e } n=m) \text{ ou } (i=n \text{ e } j=m) \}$ é inerentemente ambígua? Caso positivo, construa G ambígua que a representa. Caso negativo, construa G não-ambígua que a representa. Em qualquer caso, construa todas as árvores de derivação e todas as derivações das seguintes sentenças: abcd e bbcc.

7. Elimine os símbolos inúteis da seguinte GLC:

IV - Linguagens Livres de Contexto

$S ::= aSa \mid FbD$

$D ::= Dd \mid fF \mid c$

$A ::= aA \mid CA \mid \varepsilon$

$E ::= BC \mid eE \mid EB$

$B ::= bB \mid FE$

$F ::= fF \mid Dd$

$C ::= cCb \mid AcA$

8. Transforme em ε -livre e elimine as produções unitárias das seguintes GLC's:

* a) $S ::= AB \mid aS$

$A ::= bA \mid BCD$

$B ::= dB \mid C \mid \varepsilon$

$C ::= cCc \mid BD$

$D ::= CD \mid d \mid \varepsilon$

b) $P ::= KL \mid bKLe$

$K ::= cK \mid TV$

$T ::= tT \mid \varepsilon$

$V ::= vV \mid \varepsilon$

$L ::= LC \mid C$

$C ::= P \mid \text{com} \mid \varepsilon$

9. Elimine as recursões à esquerda das seguintes GLC's:

a) $E ::= E+T \mid E-T \mid T$

$T ::= T*F \mid T/F \mid F$

$F ::= F**P \mid P$

$P ::= (E) \mid \text{id} \mid \text{cte}$

* b) $S ::= BaS \mid \varepsilon$

$B ::= SAa \mid Bb$

$A ::= Sa \mid \varepsilon$

10. Fatore as seguintes GLC's:

* a) $S ::= bcD \mid Bcd$

$B ::= bB \mid b$

$D ::= dD \mid d$

b) $P ::= DL \mid L$

$D ::= dD \mid \varepsilon$

$L ::= L;C \mid C$

c) $C ::= V=\text{exp} \mid \text{id}(E)$

$V ::= \text{id}[E] \mid \text{id}$

$E ::= \text{exp},E \mid \text{exp}$

11. Procure responder as seguintes questões justificando a resposta:

a) Toda linguagem livre de contexto pode ser gerada por uma GLC própria?

b) Toda linguagem livre de contexto pode ser representada por uma GLC na forma normal de Chomsky?

c) Toda GLC pode ser fatorada?

d) Ambigüidade e recursão à esquerda implicam em não-fatoração?

12. A união, concatenação e intercessão das linguagens de duas GLC quaisquer são também linguagens livres de contexto?

13. Calcule o FIRST e FOLLOW das gramáticas resultantes dos exercícios anteriores.

V - Autômatos de Pilha

1. Introdução

Assim como as linguagens regulares têm nos autômatos finitos seus sistemas reconhecedores, também as linguagens livres de contexto têm seus sistemas reconhecedores: os *autômatos de pilha* (AP).

Porém, neste caso, a relação entre GLCs e APs não é tão satisfatória como a relação entre GRs e AFs. O autômato de pilha é um dispositivo não-determinístico, e a versão determinística aceita apenas um subconjunto de todas as linguagens livres de contexto. Felizmente, este subconjunto contém a sintaxe da maioria das linguagens de programação.

2. Autômatos de Pilha

O autômato de pilha é, essencialmente, um autômato finito com controle tanto da entrada quanto de uma pilha. Esta pilha, acoplada a um controle finito, pode ser usada para reconhecer conjuntos não-regulares. Considere a linguagem livre de contexto gerada pela gramática $S ::= OS0 \mid 1S1 \mid c$. Não é difícil mostrar que L não pode ser reconhecida por um autômato finito. Para aceitar L , usaremos um controle finito com dois estados e_1 e e_2 e uma pilha na qual colocamos pratos azuis, verdes e vermelhos. O dispositivo vai operar de acordo com as seguintes regras:

1. A máquina inicia com um prato vermelho no topo e com o controle finito no estado e_1 .
2. Se a entrada é 0 e o controle está no estado e_1 , um prato azul é colocado na pilha. Se a entrada é 1 e o dispositivo está no estado e_1 , um prato verde é colocado na pilha. Em ambos os casos, o dispositivo permanece no estado e_1 .
3. Se a entrada é c e o dispositivo está no estado e_1 , ele muda para o estado e_2 e não modifica a pilha.
4. Se a entrada é 0 e o dispositivo está no estado e_2 com um prato azul, que representa 0, no topo da pilha, o prato é removido. Se a entrada é 1 e o dispositivo está no estado e_2 com um prato verde, que representa 1, no topo da pilha, o prato é removido. Em ambos os casos, o controle finito permanece no estado e_2 .
5. Se o dispositivo está no estado e_2 e um prato vermelho está no topo da pilha, o prato é removido sem esperar pela próxima entrada.
6. Para todos os casos não descritos acima, o dispositivo não pode fazer nenhuma transição ou movimento na pilha.

Dizemos que o dispositivo descrito acima aceita uma cadeia de entrada se, quando processar o último símbolo da cadeia, a pilha de pratos fica completamente vazia. Note que quando a pilha está vazia, mais nenhum movimento é possível.

As regras do dispositivo são resumidas na seguinte tabela:

Prato no topo	Estado	Entrada		
		0	1	c
Azul	e ₁	Empilhe azul Fique em e ₁	Empilhe verde Fique em e ₁	Vá para e ₂
	e ₂	Desempilhe Fique em e ₂	-	-
Verde	e ₁	Empilhe azul Fique em e ₁	Empilhe verde Fique em e ₁	Vá para e ₂
	e ₂	-	Desempilhe Fique em e ₂	-
Vermelho	e ₁	Empilhe azul Fique em e ₁	Empilhe verde Fique em e ₁	Vá para e ₂
	e ₂	Sem consumir a entrada, desempilhe		

Essencialmente, o dispositivo opera da seguinte forma: no estado e₁ o dispositivo faz uma imagem de sua entrada, colocando um prato azul no topo da pilha cada vez que um 0 aparece na entrada e um prato verde cada vez que um 1 aparece na entrada. Quando a entrada é “c”, o dispositivo vai para o estado e₂. A seguir, o restante da entrada é comparado com a pilha, removendo um prato azul cada vez que o símbolo da entrada é um 0 e um prato verde cada vez que o símbolo é um 1. Se o prato no topo da pilha for da cor errada, o dispositivo pára e nenhum processamento posterior na entrada é possível. Se todos os pratos correspondem às entradas, eventualmente o prato vermelho do fundo da pilha será exposto. O prato vermelho é imediatamente removido e o dispositivo aceitou a entrada. Todos os pratos podem ser removidos apenas quando a cadeia após o “c” é a reversa da que precedia “c”.

2.1. Definição Formal

O autômato de pilha (AP) é composto por uma entrada, um autômato finito e uma pilha. A pilha é uma cadeia de símbolos de algum alfabeto. O símbolo mais à esquerda da pilha é considerado como seu topo. A máquina será não-determinística, tendo um número finito de escolhas de movimentos a efetuar em cada situação. Os movimentos podem ser de dois tipos.

No primeiro tipo de movimento o símbolo da entrada é usado. Dependendo do símbolo da entrada, do topo da pilha e do estado do autômato finito, algumas escolhas são possíveis. Cada escolha consiste de um estado seguinte para o autômato finito e uma cadeia de símbolos (possivelmente vazia) para ser trocada pelo topo da pilha. Após selecionar uma escolha, um símbolo é consumido da entrada.

O segundo tipo de movimento (chamado ϵ -movimento) é similar ao primeiro, exceto que ele não consome símbolos da entrada. Este tipo de movimento permite que o AP manipule a pilha sem ler símbolos da entrada.

Finalmente, devemos definir a linguagem aceita por um autômato de pilha. Há duas formas naturais para fazer isto. A primeira, que já foi vista, é definir a linguagem aceita como sendo o conjunto de todas as entradas para as quais alguma seqüência de movimentos faz com que a pilha do autômato fique vazia. Esta linguagem é conhecida como a linguagem aceita por pilha vazia.

A segunda forma de definir a linguagem reconhecida é similar à forma utilizadas para autômatos finitos. Isto é, designamos alguns estados como estados finais e definimos a linguagem aceita como o conjunto de todas as entradas para as quais alguma escolha de movimentos leve a um estado final.

Como podemos ver, as duas definições são equivalentes no sentido de que um conjunto que pode ser aceito por pilha vazia em um AP pode ser aceito por um estado final em outro AP, e vice-versa.

Reconhecimento por estado final é a noção mais comum, mas é fácil provar o teorema básico dos autômatos de pilha usando reconhecimento por pilha vazia. Este teorema estabelece que uma linguagem é aceita por um AP se e somente se ela é uma linguagem livre de contexto.

Um *autômato de pilha* M é uma tupla $(E, A, P, t, e_0, p_0, F)$, onde:

- E é um conjunto finito de estados;
- A é um alfabeto de entrada;
- P é um alfabeto da pilha;
- $e_0 \in E$ é o estado inicial;
- $p_0 \in P$ é o símbolo inicial da pilha;
- $F \subseteq E$ é o conjunto de estados finais;
- t é um mapeamento de $E \times \{A \cup \{\varepsilon\}\} \times P$ em subconjuntos finitos de $E \times P^*$.

A menos que seja explicitamente dito, usaremos letras minúsculas do início do alfabeto para denotar símbolos da entrada e letras do fim do alfabeto para denotar cadeias de símbolos de entrada. Letras maiúsculas denotarão símbolos da pilha e letras gregas indicarão cadeias de símbolos da pilha.

2.1.1. Movimentos

A interpretação de:

$$t(e, a, p) = \{(d_1, \gamma_1), (d_2, \gamma_2), \dots, (d_m, \gamma_m)\}$$

onde “ e ” e d_i , $1 \leq i \leq m$, são estados, a está em A , p é um símbolo da pilha, e γ_i está em P^* para $1 \leq i \leq m$, é que o AP no estado e , com o símbolo de entrada a e p no topo da pilha, pode, para algum i , entrar no estado d_i , trocar p pela cadeia γ_i e consumir um símbolo da entrada. É adotada a convenção de que o símbolo mais à esquerda de γ_i será colocado no topo da pilha.

A interpretação de:

$$t(e, \varepsilon, p) = \{(d_1, \gamma_1), (d_2, \gamma_2), \dots, (d_m, \gamma_m)\}$$

é que o AP no estado e , independente do símbolo da entrada, e com o valor p no topo da pilha, pode entrar no estado d_i e trocar p por γ_i , para algum i , $1 \leq i \leq m$. Neste caso, nenhum símbolo da entrada é consumido.

Considere o seguinte exemplo:

$$M = (\{e_1, e_2\}, \{0, 1, c\}, \{R, B, G\}, t, e_1, R, \emptyset)$$

$$t(e_1, 0, R) = \{(e_1, BR)\} \quad t(e_1, 0, B) = \{(e_1, BB)\} \quad t(e_1, 0, G) = \{(e_1, BG)\}$$

$$t(e_1, 1, R) = \{(e_1, GR)\} \quad t(e_1, 1, B) = \{(e_1, GB)\} \quad t(e_1, 1, G) = \{(e_1, GG)\}$$

$$t(e_1, c, R) = \{(e_2, R)\} \quad t(e_1, c, B) = \{(e_2, B)\} \quad t(e_1, c, G) = \{(e_2, G)\}$$

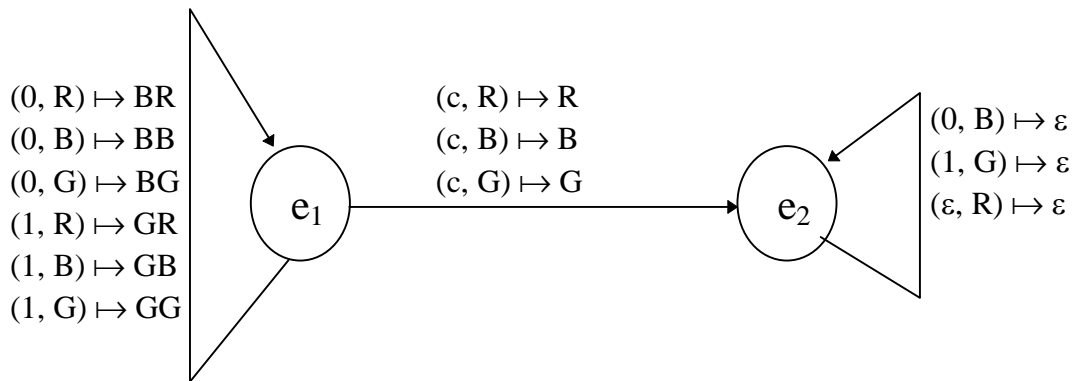
$$t(e_2, \varepsilon, R) = \{(e_2, \varepsilon)\} \quad t(e_2, 0, B) = \{(e_2, \varepsilon)\} \quad t(e_2, 1, G) = \{(e_2, \varepsilon)\}$$

O AP acima reconhece a linguagem $\{wcw^{\text{reversa}} \mid w \in (0+1)^*\}$ por pilha vazia. Note que para um movimento no qual o AP escreve um símbolo no topo da pilha, t tem o valor (e, γ) , onde $|\gamma| = 2$. Por exemplo, $t(e_1, 0, R) = \{(e_1, BR)\}$. Se γ fosse de comprimento 1, o AP simplesmente trocaria o símbolo no topo da pilha. Isto nos permite usar γ igual a ε quando queremos desempilhar.

2.1.2. Notação Gráfica

A notação gráfica dos AP é parecida com a notação dos autômatos finitos. O que muda são os rótulos das transições. Ao invés de apenas indicar o símbolo da entrada, é necessário indicar também o símbolo que deve estar no topo da pilha e a transformação a se fazer na pilha. Por exemplo, uma transição do estado e_1 para e_2 , rotulada com $(c, R) \mapsto R$ representa a transição $t(e_1, c, R) = \{(e_2, R)\}$

Para o AP acima, a representação gráfica é:



3. Autômatos de Pilha Determinísticos

O autômato de pilha mostrado na seção anterior é determinístico no sentido de que no máximo um movimento é possível para cada situação. Formalmente, dizemos que um AP $M = (E, A, P, t, e_0, p_0, F)$ é *determinístico* se:

1. Para cada $e \in E$ e $p \in P$, sempre que $t(e, \varepsilon, p)$ não for vazio, então $t(e, a, p)$ é vazio para todo $a \in A$.
2. Para nenhum $e \in E, p \in P$ e $a \in A \cup \{\varepsilon\}$, acontece que $t(e, a, p)$ contenha mais de um elemento.

A primeira regra descreve o não determinismo entre o consumo ou não do símbolo da entrada. A segunda regra é análoga à regra dos AF determinísticos. Ou seja, se o controle estiver num determinado estado, para as combinações de símbolo de entrada e símbolo no topo da pilha, somente um novo estado e uma substituição na pilha devem ser possíveis.

Sejam as seguintes transições num AP:

- a) $t(e_1, \varepsilon, F) = (e_2, GF)$
- b) $t(e_1, a, F) = (e_2, GF)$
- c) $t(e_2, a, F) = (e_3, GF)$
- d) $t(e_2, a, F) = (e_3, HF)$
- e) $t(e_3, b, G) = (e_4, \varepsilon)$
- f) $t(e_3, b, G) = (e_5, \varepsilon)$

Os seguintes não-determinismos estão representados nas transições acima:

1. O AP, no estado e_1 , pode chegar a e_2 consumindo ou não o símbolo da entrada a (transições (a) e (b));
2. O AP, no estado e_2 , consumindo o símbolo de entrada a , pode chegar a e_3 com símbolos diferentes no topo da pilha (transições (c) e (d));
3. No estado e_3 , o AP tem que consumir o mesmo símbolo de entrada, e fazer as mesmas alterações na pilha, porém pode ir para estados diferentes, e_4 ou e_5 (transições (e) e (f)).

Para autômatos finitos, os modelos determinísticos e não determinísticos são equivalentes com respeito às linguagens reconhecidas. Isto não é verdade, porém, para os autômatos de pilha. De fato, ww^{reversa} pode ser reconhecida por um autômato de pilha não-determinístico, mas não por um AP determinístico.

Para ilustrar, seja o AP para reconhecer a linguagem $\{ ww^{\text{reversa}} \mid w \in (a + b)^* \}$:

$M = (\{e_1, e_2\}, \{a, b\}, \{R, B, G\}, t, e_1, R, \emptyset)$

- | | |
|---|--|
| $t(e_1, a, R) = \{(e_1, BR)\}$ | $t(e_1, a, B) = \{(e_1, BB), (e_2, \varepsilon)\}$ |
| $t(e_1, a, G) = \{(e_1, BG)\}$ | $t(e_1, b, R) = \{(e_1, GR)\}$ |
| $t(e_1, b, B) = \{(e_1, GB)\}$ | $t(e_1, b, G) = \{(e_1, GG), (e_2, \varepsilon)\}$ |
| $t(e_1, \varepsilon, R) = \{(e_2, \varepsilon)\}$ | $t(e_2, a, B) = \{(e_2, \varepsilon)\}$ |
| $t(e_2, b, G) = \{(e_2, \varepsilon)\}$ | $t(e_2, \varepsilon, R) = \{(e_2, \varepsilon)\}$ |

Note que o AP acima é não-determinístico, porque no estado e_1 , toda vez que o símbolo de entrada corresponde ao símbolo no topo da pilha, há duas possibilidades: continuar a empilhar ou começar a desempilhar. Também há a possibilidade de desempilhar o símbolo inicial da pilha, sem consumir nenhuma entrada. Isto é necessário para reconhecer a sentença vazia, que faz parte da linguagem.

Isto porque não há um símbolo que marca a metade da cadeia, como na linguagem do primeiro exemplo deste capítulo. Desta forma, não há possibilidade de se construir um AP determinístico para a linguagem aqui proposta.

Exercícios

1) Construa Autômatos de Pilha para reconhecer as seguintes linguagens:

a) $\{a^{i+3}b^{2i+1} \mid i \geq 0\} \cup \{a^{2i+1}b^{3i} \mid i \geq 0\}$

* b) $\{a^i b^j c^j d^j e^3 \mid i, j \geq 0\}$

* c) $\{a^i b^j \mid i, j \geq 0, \text{ com } i = 2j \text{ ou } 2i = j\}$

d) $\{a^i b^j c^k d^l \mid i, j, k, l \geq 0, i < l \text{ e } j \neq k\}$

e) $\{a^i b^j b^i a^j \mid i, j \geq 0\}$

VI - Análise Sintática

1. Introdução

Um analisador sintático, ou *parser*, é um sistema capaz de construir uma derivação para qualquer sentença em alguma linguagem $L(G)$ baseado em uma gramática G .

Um parser também pode ser visto como um mecanismo para a construção de árvores de derivação.

Seja $G = (N, T, P, S)$ uma GLC com as produções de P numeradas de 1 a p e uma derivação $S \rightarrow^+ x$. A seqüência formada pelo número das produções utilizadas na derivação

$S \rightarrow^+ x$ constitui o *parse* de x em G .

O *parse ascendente* é constituído pela seqüência invertida dos números das produções utilizadas em $S \rightarrow_{\text{dir}}^+ x$, onde $\rightarrow_{\text{dir}}^+$ denota a derivação mais a direita.

O *parse descendente* é constituído pela seqüência dos números das produções utilizadas em $S \rightarrow_{\text{esq}}^+ x$, onde $\rightarrow_{\text{esq}}^+$ denota a derivação mais a esquerda.

Seja a seguinte GLC:

$$S ::= aAB \mid aB \quad (1, 2) \quad A ::= bc \mid bAc \quad (3, 4) \quad B ::= eBf \mid d \quad (5, 6)$$

onde os números entre parênteses representam os números das produções.

Para a sentença *abcceedff* o parse ascendente é $\{3, 4, 6, 5, 5, 1\}$ (que é seqüência de derivações mais a direita $\{1, 5, 5, 6, 4, 3\}$ invertida). Já o parse descendente é $\{1, 4, 3, 5, 5, 6\}$ (seqüência de derivações mais a esquerda).

2. Classes de Analisadores Sintáticos

Existem duas classes fundamentais de analisadores sintáticos, definidas em função da estratégia utilizada na análise.

Os analisadores *ascendentes* procuram chegar ao símbolo inicial da gramática a partir da sentença a ser analisada, olhando a sentença, ou parte dela, para decidir qual produção será utilizada na *redução*.

Os analisadores *descendentes* procuram chegar à sentença a partir do símbolo inicial de G , olhando a sentença ou parte dela para decidir que produção deverá ser usada na *derivação*.

3. Analisadores Ascendentes (Família LR)

Sob o ponto de vista lógico, um analisador LR consiste de duas partes: um algoritmo de análise sintática (padrão para todas as técnicas da família) e uma tabela de parsing (específica para cada técnica e para cada gramática).

As principais razões da grande importância desta família de analisadores na teoria de parsing e, conseqüentemente, no desenvolvimento de compiladores, são:

- analisam praticamente todas as construções sintáticas de linguagens de programação que podem ser representadas por GLC;
- são mais gerais que os outros analisadores ascendentes e que a maioria dos descendentes e que a maioria dos descendentes sem back-tracking;
- possuem a propriedade de detecção imediata de erros sintáticos;
- o tempo de análise é proporcional ao tamanho da sentença a ser analisada.

Em contrapartida, pode-se destacar como ponto negativo dos analisadores LR a complexidade da construção da tabela de parsing, bem como o espaço requerido para seu armazenamento (especialmente as técnicas mais abrangentes).

3.1. Estrutura dos Analisadores LR

Um analisador LR se compõe de:

- um algoritmo de análise sintática;

- uma tabela de análise sintática;
- uma pilha de estados (ou pilha sintática) que conterà um histórico da análise efetuada. Ela é inicializada com o estado inicial da análise sintática;
- uma entrada que conterà a sentença a ser analisada, seguida por uma marca de final de sentença (por exemplo, o símbolo “\$”).
- uma GLC com as produções numeradas de 1 a p.

3.2. Algoritmo de Análise Sintática LR

O algoritmo de análise sintática LR consiste em comparar o estado do topo da pilha e o próximo símbolo da entrada e, com base nessas informações, consultar a tabela de parsing para decidir a próxima ação a ser efetuada e efetuá-la. Esta ação pode ser:

- PARAR
- ERRO
- TRANSFERIR(S)
- REDUZIR(R)

Os procedimentos associados às ações têm os seguintes significados:

- PARAR - Fim de análise. O procedimento consiste em encerrar a análise sintática.
- ERRO - Indica a ocorrência de um erro sintático. O procedimento é o seguinte: as rotinas de recuperação, se existirem, deverão ser ativadas para que a situação de erro seja contornada e a análise possa continuar; senão a análise deve ser encerrada.
- TRANSFERIR(S) - Significa o reconhecimento sintático do símbolo da entrada. O procedimento consiste em retirar o símbolo da entrada, e o estado S indicado na tabela de parsing deverá ser empilhado.
- REDUZIR(R) - Significa que uma redução pela produção n^o R deve ser efetuada. O procedimento consiste em retirar da pilha sintática tantos estados quantos forem os símbolos do lado direito da produção R, e o símbolo do lado esquerdo dessa produção deverá ser colocado na entrada antes da próxima ação do analisador.

Observação: o processo de determinar e efetuar as ações sintáticas deverá ser repetido até que a ação PARAR seja encontrada, ou quando um erro for encontrado (nas implementações sem recuperação de erro).

3.3. A Tabela de Parsing LR

Para implementar uma tabela de parsing pode-se usar uma matriz ou uma lista. No caso de matriz, como normalmente as tabelas de parsing são bastante esparsas, utilizam-se técnicas de compactação de matrizes.

Por exemplo, considere a seguinte gramática LR(0):

- $E' ::= E\$$
- $E ::= E+T \mid T$
- $T ::= T * F \mid F$
- $F ::= (E) \mid id$

A tabela SLR(1) correspondente seria:

	id	()	+	*	\$	E	T	F
0	S ₅	S ₄					S ₁	S ₂	S ₃
1				S ₆		PARE			
2			R ₂	R ₂	S ₇	R ₂			
3			R ₄	R ₄	R ₄	R ₄			
4	S ₅	S ₄					S ₈	S ₂	S ₃
5			R ₆	R ₆	R ₆	R ₆			
6	S ₅	S ₄						S ₉	S ₃
7	S ₅	S ₄							S ₁₀
8			S ₁₁	S ₆					

9			R ₁	R ₁	S ₇	R ₁			
10			R ₃	R ₃	R ₃	R ₃			
11			R ₅	R ₅	R ₅	R ₅			

3.4. A Configuração de um Analisador LR

A configuração de um analisador LR é um par cujo primeiro elemento é o conteúdo da pilha sintática e o segundo é a entrada a ser analisada. Por exemplo, a configuração:

$$(S_0 S_1 S_2 \dots S_m, a_i a_{i+1} a_{i+2} \dots a_n \$)$$

tem o seguinte significado: os primeiros $i-1$ símbolos da entrada já foram analisados e a próxima ação a ser efetuada será determinada pelo estado S_m (topo da pilha) e por a_i (próximo símbolo da entrada).

3.5. Gramáticas LR(0)

Uma das motivações para o estudo de linguagens livres de contexto determinísticas é a sua habilidade em descrever a sintaxe de linguagens de programação. Vários sistemas de construção de compiladores requerem especificação sintática na forma de GLC's restritas que permitem apenas a representação de GLC's determinísticas. Além disso, o parser produzido por tais geradores de compiladores é essencialmente um autômato de pilha determinístico.

Nesta seção será mostrado um tipo restrito de GLC chamado gramática LR(0). Esta classe de gramáticas é a primeira de uma família chamada de *família de gramáticas LR*. O nome LR(0) significa "left-to-right scan on the input producing a rightmost derivation and using 0 symbols of lookahead on the input" (busca de símbolos da esquerda para a direita na entrada produzindo uma derivação mais à direita e usando 0 símbolos de contexto de entrada).

As gramáticas LR(0) definem exatamente as linguagens livres de contexto determinísticas que têm a *propriedade do prefixo*. Esta propriedade estabelece que se w está em L , nenhum prefixo próprio de w está em L . Note que a propriedade do prefixo não é muito severa porque a introdução de uma marca de final de sentença (\$) converte toda linguagem livre de contexto determinística em uma linguagem livre de contexto determinística com a propriedade do prefixo. Assim, $L\$ = \{w\$ \mid w \in L\}$ é uma linguagem livre de contexto determinística com a propriedade do prefixo sempre que L for uma linguagem livre de contexto determinística.

3.6. Itens LR

Para introduzir as gramáticas LR(0) são necessárias algumas definições preliminares. Em primeiro lugar, um *item*, para uma dada GLC, é uma produção com um ponto em alguma posição do lado direito, incluindo o início e o fim. No caso da produção vazia, $B ::= \epsilon$, o único item possível é $B ::= \bullet$.

A seguinte gramática será usada em uma série de exemplos subseqüentes:

$$S' ::= Sc \quad S ::= SA \mid A \quad A ::= aSb \mid ab$$

Esta gramática, cujo símbolo inicial é S' , gera as cadeias de parênteses corretamente aninhados, tratando a e b como abre e fecha parênteses, respectivamente, e c como uma marca de final de sentença.

Os itens para a gramática são

$S' ::= \bullet Sc$	$S ::= \bullet SA$	$A ::= \bullet aSb$
$S' ::= S \bullet c$	$S ::= S \bullet A$	$A ::= a \bullet Sb$
$S' ::= Sc \bullet$	$S ::= SA \bullet$	$A ::= aS \bullet b$
	$S ::= \bullet A$	$A ::= aSb \bullet$
	$S ::= A \bullet$	$A ::= \bullet ab$
		$A ::= a \bullet b$
		$A ::= ab \bullet$

A seguir, serão usados os símbolos \rightarrow_D^* e \rightarrow_D para denotar as derivações mais a direita e passos de derivações mais a direita, respectivamente. Uma *forma-sentencial-direita* é uma forma sentencial que pode ser derivada por uma derivação mais a direita. Um *controlador* de uma forma sentencial direita $x\beta y$ para uma GLC G é uma sub-cadeia β , tal que:

$$S \rightarrow_D^* xAy \rightarrow_D x\beta y,$$

com $A \in N$. Em outras palavras, um controlador de $x\beta y$ é uma sub-cadeia que pode ser introduzida no último passo de uma derivação mais a direita de $x\beta y$. Note que neste contexto, a posição de β em $x\beta y$ é importante.

Um *prefixo viável* de uma forma sentencial direita γ é qualquer prefixo de γ que termine antes do (ou no) último símbolo do controlador de γ . Se o controlador ocupa as posições i até j , na cadeia, então os prefixos viáveis são todos os k -prefixos com $k \leq j$.

Por exemplo, na gramática acima há uma derivação mais a direita:

$$S' \rightarrow Sc \rightarrow SAc \rightarrow SaSbc$$

Assim, $SaSbc$ é uma forma sentencial direita, e seu controlador é aSb . Em qualquer gramática não ambígua sem símbolos inúteis, que a derivação mais a direita de uma dada forma sentencial é única, assim, seu controlador é único. Assim, pode-se falar de “o controlador”, ao invés de “um controlador”. Os prefixos viáveis de $SaSbc$ são $\{\varepsilon, S, Sa, SaS, SaSb\}$ (ou seja, do início da forma sentencial até o último símbolo do controlador).

Dizemos que um item $A ::= \alpha \bullet \beta$ é *válido* para um prefixo viável $x\alpha$ se existe uma derivação mais à direita:

$$S \rightarrow_D^* xAy \rightarrow_D x\alpha\beta y$$

Saber quais itens são válidos para um dado prefixo viável é importante para encontrar uma derivação mais a direita reversa, como segue. Um item é denominado *completo* se o ponto está na posição mais a direita do item. Se $A ::= \alpha \bullet$ é um item completo válido para γ , então pode ser que $A ::= \alpha$ tenha sido usada no último passo e que a forma sentencial direita precedente na derivação de γy era xAy .

Porém isto é apenas uma suspeita, já que $A ::= \alpha \bullet$ pode ser válido para γ devido a uma derivação mais a direita $S \rightarrow^* xAy' \rightarrow \gamma y'$. Certamente pode haver dois ou mais itens válidos para γ , ou pode haver um controlador para γy que inclua símbolos de y . Intuitivamente, uma gramática é definida como LR(0) se em cada situação deste tipo xAy é realmente a forma sentencial direita prévia par γy . Neste caso é possível começar com uma cadeia de terminais w que está em $L(G)$, e assim é uma forma sentencial direita de G , e trabalhar para trás, no sentido das formas sentenciais direitas anteriores, até chegarmos a S . Tem-se assim uma derivação mais à direita de w .

Por exemplo, considerando a gramática anterior e a forma sentencial direita abc . Já que $S \rightarrow_D^* Ac \rightarrow abc$ vemos que $A ::= ab \bullet$ é válido para o prefixo viável ab . Também vemos que $A ::= a \bullet b$ é válido para o prefixo viável a , e $A ::= \bullet ab$ é válido para o prefixo viável ε . Como $A ::= ab \bullet$ é um item completo, podemos ser capazes de deduzir que Ac foi a forma sentencial direita prévia de abc .

3.7. Cálculo dos Conjuntos de Itens Válidos

A definição de gramáticas LR(0) e o método para aceitar $L(G)$ para uma gramática LR(0) por um autômato de pilha determinístico dependem de conhecer o conjunto de itens válidos para cada prefixo viável γ . Segue que para toda GLC G , qualquer que seja, o conjunto de prefixos é um conjunto regular, e este conjunto regular é aceito por um autômato finito não-determinístico cujos estados são os itens de G . Aplicando a determinização a este autômato, obtemos um autômato finito determinístico cujo estado, em resposta ao prefixo viável γ é o conjunto de itens válidos para γ .

O AFND M que reconhece os prefixos viáveis para a GLC $G = (N, T, P, S)$ é definido como segue. Seja $A = (E, N \cup T, t, e_0, E)$, onde E é o conjunto de itens para G mais o estado e_0 , o qual não é um item. Considerando $\alpha, \beta, \gamma \in (N \cup T)^*$ e $A, B \in N$, vamos definir:

1. $t(e_0, \varepsilon) = \{S ::= \bullet\alpha \mid S ::= \alpha \text{ é uma produção}\}$
2. $t(A, \alpha\bullet B\beta, \varepsilon) = \{B ::= \bullet\gamma \mid B ::= \gamma \text{ é uma produção}\}$
3. $t(A ::= \alpha\bullet X\beta, X) = \{A ::= \alpha X\bullet\beta\}$

A regra 2 permite a expansão de uma variável B que aparece imediatamente a direita do ponto. A regra 3 permite mover o ponto sobre qualquer símbolo X da gramática, se X é o próximo símbolo da entrada.

O AFND para a gramática do exemplo anterior é mostrado a seguir

Regra 1	$t(q_0, \varepsilon) = S' ::= \bullet Sc$
Regra 2	$t(S' ::= \bullet Sc, \varepsilon) = S ::= \bullet SA$ $t(S' ::= \bullet Sc, \varepsilon) = S ::= \bullet A$ $t(S ::= \bullet SA, \varepsilon) = S ::= \bullet SA$ $t(S ::= \bullet SA, \varepsilon) = S ::= \bullet A$ $t(S ::= \bullet A, \varepsilon) = A ::= \bullet aSb$ $t(S ::= \bullet A, \varepsilon) = A ::= \bullet ab$ $t(S ::= S\bullet A, \varepsilon) = A ::= \bullet aSb$ $t(S ::= S\bullet A, \varepsilon) = A ::= \bullet ab$ $t(S ::= a\bullet Sb, \varepsilon) = S ::= \bullet SA$ $t(S ::= a\bullet Sb, \varepsilon) = S ::= \bullet A$

Regra 3	$t(S' ::= \bullet Sc, S) = S' ::= S \bullet c$ $t(S' ::= S \bullet c, c) = S' ::= Sc \bullet$ $t(S ::= \bullet A, A) = S ::= A \bullet$ $t(S ::= \bullet SA, S) = S ::= S \bullet A$ $t(S ::= S \bullet A, A) = S ::= SA \bullet$ $t(A ::= \bullet aSb, a) = A ::= a \bullet Sb$ $t(A ::= a \bullet Sb, S) = A ::= aS \bullet b$ $t(A ::= aS \bullet b, b) = A ::= aSb \bullet$ $t(A ::= \bullet ab, a) = A ::= a \bullet b$ $t(A ::= a \bullet b, b) = A ::= ab \bullet$
---------	--

3.8. Definição de uma Gramática LR(0)

Dizemos que uma gramática G é uma gramática LR(0) se

- 1) Seu símbolo inicial não aparece no lado direito de nenhuma produção, e
- 2) Para todo prefixo viável γ de G , sempre que $A ::= \alpha \bullet$ é um item completo válido para γ , então nenhum outro item completo e nenhum item com um terminal no lado direito do ponto é válido para γ .

Não há proibição a que vários itens incompletos sejam válidos para γ .

Para calcular o conjunto de itens válidos para cada prefixo viável, converta o AFND cujos estados são itens em um AFD. No AFD, o caminho do estado inicial rotulado com a seqüência γ leva ao estado que é o conjunto de itens válidos para γ . Assim, deve-se construir o AFD e inspecionar cada estado par ver se uma violação das condições LR(0) acontece.

Por exemplo, o AFD mostrado a seguir foi construído a partir do AFND anterior. Todos os estados, exceto I_0 , I_1 , I_3 e I_6 consistem de um único item completo. Os estados com mais de um item não têm itens completos, e o símbolo inicial S' não aparece no lado direito de nenhuma produção. Assim, a gramática que viemos usando como exemplo é, de fato, LR(0).

Sejam os estados do AFD:

$$I_0 = \{e_0, S' ::= \bullet Sc, S ::= \bullet SA, S ::= \bullet A, A ::= \bullet aSb, A ::= \bullet ab\}$$

$$I_1 = \{S' ::= S \bullet c, S ::= S \bullet A, A ::= \bullet aSb, A ::= \bullet ab\}$$

$$I_2 = \{S ::= A \bullet\}$$

$$I_3 = \{A ::= a \bullet Sb, A ::= a \bullet b, S ::= \bullet SA, S ::= \bullet A, A ::= \bullet aSb, A ::= \bullet ab\}$$

$$I_4 = \{S' ::= Sc \bullet\}$$

$$I_5 = \{S ::= SA \bullet\}$$

$$I_6 = \{A ::= aS \bullet b, S ::= S \bullet A, A ::= \bullet aSb, A ::= \bullet ab\}$$

$$I_7 = \{A ::= ab \bullet\}$$

$$I_8 = \{A ::= aSb \bullet\}$$

As transições são:

$$t(I_0, S) = I_1$$

$$t(I_0, A) = I_2$$

$$t(I_0, a) = I_3$$

$$t(I_1, c) = I_4$$

$$t(I_1, a) = I_5$$

$$t(I_3, a) = I_3$$

$$t(I_3, S) = I_6$$

$$t(I_3, A) = I_2$$

$$t(I_3, b) = I_7$$

$$t(I_6, a) = I_3$$

$$t(I_6, b) = I_8$$

3.9. Construção do Conjunto LR

Há um método opcional para calcular o conjunto de itens para uma gramática LR. Este método não envolve determinização de autômatos pois gera o autômato determinizado diretamente a partir da gramática. Porém a aplicação do método é mais complexa.

Serão definidos, inicialmente, dois procedimentos auxiliares. O algoritmo 6.3 define o método propriamente dito.

Algoritmo 6.1: Fechamento de um estado e_i
 Entrada: Um estado e_i e uma GLC $G=(N, T, P, S)$

Saída: O estado e_i fechado

Repita

Se $X ::= \alpha \bullet Y \beta$ pertence ao estado e_i e $Y \in N$ então
 Para todas as produções da forma $Y ::= \gamma$ faça
 Adicione o item $Y ::= \bullet \gamma$ ao estado e_i

FimPara

FimSe

Até não haver mais modificações a fazer em e_i .

Algoritmo 6.2: Cálculo dos estados sucessores de e_i

Entrada: Um estado e_i e uma GLC $G=(N, T, P, S)$

Saída: Um conjunto de núcleos U e um conjunto de transições R

Faça $U = R = \{ \}$

Seja Δ o conjunto dos símbolos terminais e não terminais que ocorrem imediatamente à direita do ponto nos itens pertencentes a e_i .

Para todo símbolo $\delta \in \Delta$ faça

Seja $I = \{ X ::= \alpha \bullet \delta \beta \}$ o conjunto dos itens de e_i onde δ ocorre imediatamente à direita do ponto.

Coloque em U um estado aberto e_δ com os itens de I no núcleo

Coloque em R uma transição de e_i rotulada com δ para e_δ .

Fim Para

Retorne U e R .

Algoritmo 6.3: Construção do autômato de itens válidos

Entrada: Uma GLC $G=(N, T, P, S)$

Saída: Um autômato de Itens LR $A = (E, A, t, e_0, F)$

Crie um estado inicial e_0 aberto com o núcleo $S' ::= \bullet S \$$

Coloque e_0 em E

Coloque $N \cup T$ em A

Repita

Para todo estado aberto $e_i \in E$ faça:

Faça fechamento(e_i)

Adicione a E os estados de sucessores(e_i)

Adicione a t as transições de sucessores(e_i)

Fim Para

Até que E não contenha mais estados abertos.

3.10. Construção da Tabela de Parsing SLR(1)

O algoritmo para a construção da tabela de parsing SLR(1), usando o autômato finito determinístico da seção anterior, é o seguinte:

Algoritmo 6.4: Construção da Tabela de Parsing SLR(1)

Entrada: Um AFD com os itens de uma GLC LR(0)

Saída: Uma tabela de parsing para a GLC da entrada.

Crie uma linha para cada estado do AFD e uma coluna para cada símbolo de $T \cup N \cup \{ \$ \}$.

Coloque TRANSFERE(e_i) nos estados e_j que tem transições para estados e_i nas colunas referentes aos símbolos que rotulam a transição.

Para cada estado e_i com itens completos faça

Para cada item completo faça

Coloque REDUZIR(N) na interseção da linha correspondente ao estado e_i com as colunas correspondentes a símbolos que pertencem a FOLLOW(A), onde N é o número da produção no item completo em questão e A é o símbolo do lado esquerdo desta produção.

Fim Para

Fim Para

Coloque PARE na interseção da linha referente ao estado que contenha o item $S' ::= S \bullet \$$, onde S' é o novo símbolo inicial, na coluna do símbolo \$.

As posições da tabela que ficarem vazias representarão as situações de erro.

A condição *SLR(1)*, que é a condição para a tabela construída pelo algoritmo acima seja *SLR(1)*, é que não existam estados inadequados (estados com conflitos TRANSFERE-REDUZ e/ou conflitos REDUZ-REDUZ). Isto equivale a dizer que cada posição da tabela deverá conter no máximo uma ação.

Para que um analisador LR aceite qualquer sentença correta e detecte pelo menos um erro em cada sentença incorreta, a tabela de parsing deve possuir as seguintes propriedades:

- A condição de erro nunca será encontrada para as sentenças sintaticamente corretas.
- Cada TRANSFERE deverá especificar um único estado, o qual dependerá do símbolo da entrada.
- Cada REDUZ só é realizado quando os estados do topo da pilha (os símbolos que eles representam) forem exatamente os símbolos do lado direito de alguma produção de G.
- O PARE só será encontrado quando a análise estiver completa.

4. Analisadores Descendentes

A idéia geral desta classe de analisadores resume-se a uma tentativa de construir a derivação mais à esquerda da sentença de entrada, ou, de maneira equivalente, a uma tentativa de construir a árvore de derivação da sentença a partir do símbolo inicial da gramática em questão. Estes analisadores podem ser implementados com ou sem “back-tracking”.

A utilização de back-tracking permite que um conjunto maior de GLCs possa ser analisado. Entretanto, apresenta várias desvantagens, dentre as quais pode-se citar:

- Maior tempo necessário para a análise;
- Dificuldade na recuperação de erros;
- Problemas na análise semântica e geração de código.

Todos estes problemas são decorrentes do não-determinismo.

As implementações sem back-tracking, apesar de limitarem a classe de gramáticas que podem ser analisadas (como será visto adiante), tornam-se mais vantajosas pelo fato de serem técnicas determinísticas e superarem as deficiências práticas das implementações sem back-tracking.

4.1. Analisadores Descendentes sem Back-Tracking

Para uma GLC poder ser analisada por estes analisadores, ela deve satisfazer as seguintes condições:

- não possuir recursão à esquerda;
- Estar fatorada, isto é, se $A ::= \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$ são as produções para o não-terminal A de uma determinada GLC, então $FIRST(\alpha_i) \cap FIRST(\alpha_j) = \emptyset$, $1 \leq i, j \leq n$, $i \neq j$.
- Para todo $A \in N$, tal que $A \rightarrow^* \epsilon$, $FIRST(A) \cap FOLLOW(A) = \emptyset$.

A seguir serão apresentadas duas técnicas de análise sintática descendente sem back-tracking.

4.1.1. Técnica de Implementação Descendente Recursivo

A técnica “descendente recursivo” consiste na construção de um conjunto de procedimentos (normalmente recursivos), um para cada símbolo não terminal da gramática em questão.

A principal desvantagem desta técnica é que ela não é geral, ou seja, os procedimentos são específicos para cada gramática. Além disso, o tempo para análise é maior se comparado com a técnica que será vista a seguir. Também existe a necessidade de uma linguagem que permita recursividade para sua implementação.

Por outro lado, a principal vantagem desta técnica é a simplicidade. Além disso, ela aceita uma classe maior de linguagens pelo fato de que em alguns casos a condição (c) pode ser relaxada.

Exemplo:

O analisador sintático descendente recursivo da gramática:

$$E ::= TE' \qquad E' ::= +TE' \mid \varepsilon \qquad T ::= (E) \mid id$$

seria composto pelos seguintes procedimentos:

(* programa principal *)

início

 analisa(símbolo);
 E(símbolo);

fim

procedimento E(simb);

início

 T(simb);
 Elinha(simb);

fim

procedimento Elinha(simb);

início

 se simb = “+” então
 analisa(simb);
 T(simb);
 Elinha(simb);

 fim se;

fim;

procedimento T(simb);

início

 se simb = “(“ então
 analisa(simb);
 E(simb);
 analisa(simb);
 se simb = “)” então
 analisa(simb);

 senão

 erro;

 senão

 se simb = “id” então
 analisa(simb);

 senão

 erro;

 fim se;

fim;

4.1.2. Parser Preditivo (LL)

O parser preditivo é uma maneira eficiente de implementar um parser descendente recursivo. O parser preditivo consiste de:

- Entrada: contendo a sentença (programa) a ser analisada.
- Pilha: usada para simular a recursividade, ela prevê a parte da entrada que está para ser analisada. Ela é inicializada com \$ e o símbolo inicial da gramática em questão.
- Tabela de parsing: ou tabela de análise sintática. Ela contém as ações a serem efetuadas. É uma matriz $M(A, a)$, onde $A \in N$ e $a \in T$.

O algoritmo do parser preditivo tem como função determinar, a partir de X (o elemento do topo da pilha) e de “a” (o próximo elemento da entrada), a ação a ser executada, a qual poderá ser:

- Se $X = a = \$$, o algoritmo anuncia o final da análise.
- Se $X = a \neq \$$, o analisador retira S do topo da pilha e a da entrada.
- Se $X \neq a$ e $X \in T$, então a situação é de erro.
- Se $X \in N$, o analisador consulta a tabela de parsing $M(X, a)$, a qual poderá conter o número de uma produção ou um indicativo de erro.

Se $M(X, a)$ contém o número de uma produção, e esta produção é, por exemplo, $X ::= uVw$, então X, que está no topo da pilha, deve ser substituído por wVu (com u no topo).

Se $M(X, a)$ contém um indicativo de erro, então o analisador deve ativar os procedimentos de recuperação de erros. Em implementações sem recuperação de erros a análise é encerrada.

O termo “preditivo” deve-se ao fato de que a pilha sempre contém a descrição do restante da sentença (se ela estiver correta); isto é, a pilha prevê a parte da sentença que deve estar na entrada para que a sentença esteja correta.

Algoritmo 6.5: Análise Sintática Descendente

Repita:

```

Faça X ser o topo da pilha.
Faça “a” ser o próximo símbolo da entrada
Se X é terminal ou $ então
    Se X = a então
        retire X do topo da pilha
        retire a da entrada
    Senão
        erro
    Fim Se
Senão (* X é não-terminal *)
    Se  $M(X, a) = X ::= Y_1Y_2\dots Y_k$  então
        retire X da pilha
        coloque  $Y_kY_{k-1}\dots Y_2Y_1$  na pilha (* com  $Y_1$  no topo *)
    Senão
        erro
    Fim Se
Fim Se

```

Até $X = \$$ (* pilha vazia, análise concluída *)

4.1.2.1. Construção da Tabela de Parsing para o Parser Preditivo

A idéia geral é a seguinte: Se $A ::= \alpha \in P$ e $a \in \text{FIRST}(\alpha)$, então se A está no topo da pilha e a é o próximo símbolo da entrada, deve-se expandir (derivar) A , usando a produção

$A ::= \alpha$.

Se $\alpha = \varepsilon$ ou $\alpha \rightarrow^* \varepsilon$ deve-se observar que ε nunca aparecerá na entrada. Neste caso, se $a \in \text{FOLLOW}(A)$ deve-se expandir (derivar) A através da produção $A ::= \alpha$.

Algoritmo 6.6: Construção da Tabela de Parsing

Para cada produção $A ::= \alpha \in P$, faça:

Para todo $a \in \text{FIRST}(\alpha)$, exceto ε , faça:

Coloque o número da produção $A ::= \alpha$ em $M(A, a)$;

Fim Para;

Se $\varepsilon \in \text{FIRST}(\alpha)$ então

Coloque o número da produção $A ::= \alpha$ em $M(A, b)$ para todo $b \in \text{FOLLOW}(A)$;

Fim Se;

Fim Para;

As posições de M que ficarem indefinidas representarão as situações de erro.

A tabela de parsing dos analisadores preditivos deve possuir a propriedade de que, para cada entrada da tabela M exista no máximo uma produção (seu número). Isto viabiliza a análise determinística da sentença de entrada. Para que esta propriedade se verifique, a gramática considerada deverá satisfazer as seguintes condições:

- Não possuir recursão à esquerda;
- Estar fatorada;
- Para todo $a \in N$ tal que $A \rightarrow^* \varepsilon$, $\text{FIRST}(A) \cap \text{FOLLOW}(A) = \emptyset$.

As GLCs que satisfazem estas condições são denominadas GLCs $LL(k)$, isto é, GLCs que podem ser analisadas deterministicamente da esquerda para a direita, e o analisador construirá uma derivação mais a esquerda, sendo necessário a cada passo o conhecimento de k símbolos de “look-ahead” (símbolos da entrada que devem ser vistos para que uma ação seja determinada).

Somente GLCs $LL(k)$ podem ser analisadas pelos analisadores preditivos (as demais causam conflitos na construção da tabela de parsing ou fazem com que o analisador entre em “loop”). Por isso, os analisadores preditivos são também denominados *analisadores $LL(k)$* . Na prática, usa-se $k = 1$, obtendo-se desta forma *analisadores $LL(1)$* para GLC $LL(1)$.

Exercícios

- 1) Efetue a análise sintática das sentenças seguintes usando a técnica $LR(1)$, com base na tabela do texto:
id * id
(id + id) * id
(id +) id * id
- 2) Construa a tabela de parsing $SLR(1)$ para as seguintes GLCs:
 - a) $E ::= E+T \mid T$
 $T ::= T * F \mid F$
 $F ::= (E) \mid id$
 - b) $C ::= \text{if } E \text{ then } C \text{ else } C \mid \text{if } E \text{ then } C \mid \text{com}$
 $E ::= \text{exp}$
 - * c) $B ::= AB'$
 $B' ::= \text{or } AB' \mid \varepsilon$
 $A ::= CA'$
 $A' ::= \text{and } CA' \mid \varepsilon$
 $C ::= \text{exp} \mid (B) \mid \text{not}(C)$
- 3) Construa o parse descendente recursivo para as gramáticas do exercício 2.
- 4) Construa o parse preditivo para as gramáticas do exercício 2.

Respostas aos Exercícios Propostos

Capítulo II

Exercício 2:

Uma sentença possível é: “aababa”, e sua derivação é:

$$S \rightarrow aSba \rightarrow aaSbaba \rightarrow aababa$$

A linguagem da gramática é $L(G) = \{a^n(ba)^n, n \geq 0\}$

Exercício 3:

a) $S ::= aB \mid a$

$$B ::= aB \mid bB \mid cB \mid a$$

ou, alternativamente:

$$S ::= a \mid aBa$$

$$B ::= aB \mid bB \mid cB \mid \varepsilon$$

ou, ainda:

$$S ::= aB$$

$$B ::= \varepsilon \mid bC \mid cC \mid aC \mid a$$

$$C ::= bC \mid cC \mid aC \mid a$$

b) $S ::= aA \mid bA \mid cA \mid \varepsilon$

$$A ::= aS \mid bS \mid cS$$

d) $S ::= 1SBC \mid 1BC$

$$CB ::= BC$$

$$1B ::= 10$$

$$0B ::= 00$$

$$0C ::= 02$$

$$2C ::= 22$$

f) $S ::= aSc \mid aBc$

$$B ::= bB \mid \varepsilon$$

Exercício 5

a) $S ::= A \mid B$

$$A ::= aaAb \mid abbb$$

$$B ::= aBbbb \mid aaaa$$

b) $S ::= aSb \mid aAb$ (estas produções garantem que $i, k > 0, i = k$)

$$A ::= aA \mid a$$
 (garantia de que $i > k$, já que um a deve ser gerado para parar.)

Exercício 7

a) $((0 + 1)(0 + 1))^*$

b) $(0 + 1)^* 101(0 + 1)^* 101(0 + 1)^*$

Capítulo III

Exercício 1:

a)

	δ	0	1
\rightarrow	q_0	q_1	-
	q_1	q_1	q_2
*	q_2	q_1	q_2

b)

	δ	0	1	2
\rightarrow	q_0	q_0	q_1	q_0
*	q_1	q_1	q_0	q_1

Exercício 2:

a) Gramática regular:

$A ::= aA \mid aB \mid bA$

$B ::= aC$

$C ::= bD \mid b$

$D ::= aD \mid bD \mid a \mid b$

Determinização:

	δ'	a	b
\rightarrow	[A]	[A, B]	[A]
	[A, B]	[A, B, C]	[A]
	[A, B, C]	[A, B, C]	[A, D]
*	[A, D]	[A, B, D]	[A, D]
*	[A, B, D]	[A, B, C, D]	[A, D]
*	[A, B, C, D]	[A, B, C, D]	[A, D]

Renomeando os estados:

	δ'	a	b
\rightarrow	A	B	A
	B	C	A
	C	C	D
*	D	E	D
*	E	F	D
*	F	F	D

Classes de Equivalência:

K-F	F
{A, B, C}	{D, E, F}
{A, B} {C}	{D, E, F}
{A} {B} {C}	{D, E, F}
{A} {B} {C}	{D, E, F}

Chamando a CE {D, E, F} de D, temos o seguinte AFD mínimo:

	δ'	a	b
\rightarrow	A	B	A
	B	C	A
	C	C	D
*	D	D	D

A GR do AFD mínimo é a seguinte:

$A ::= aB \mid bA$

$B ::= aC \mid bA$

$C ::= aC \mid bD \mid b$

$D ::= aD \mid bD \mid a \mid b$

Exercício 3:

a) Autômato Finito:

δ'	0	1
→ S	S, A, C	S, B
A	A, C, D	-
B	-	B, D
C	A, C, D	-
* D	-	-

Determinização:

δ'	0	1
→ [S]	[S, A, C]	[S, B]
[S, A C]	[S, A, C, D]	[S, B]
[S, B]	[S, A, C]	[S, B, D]
* [S, A, C, D]	[S, A, C, D]	[S, B]
* [S, B, D]	[S, A, C]	[S, B, D]

Renomeando:

δ'	0	1
→ S	E	F
E	G	F
F	E	H
* G	G	F
* H	E	H

Classes de Equivalência:

K-F	F
{S, E, F}	{G, H}
{S} {E} {F}	{G} {H}

Como as CE's contêm apenas um estado, não há mais o que fazer, e o AFD já é mínimo.

A GR equivalente é:

S ::= 0E | 1F

E ::= 0G | 1F | 0

F ::= 0E | 1H | 1

G ::= 0G | 1F | 0

H ::= 0E | 1H | 1

Capítulo IV

Exercício 8:

a) Conjunto dos ϵ -não-terminais: E = {B, D, C, A, S}

Eliminação das ϵ -produções:

S ::= AB | aS | A | B | a

A ::= bA | BCD | b | CD | BD | BC | C | D | B

B ::= dB | C | d

C ::= cCc | BD | cc | B | D

D ::= CD | d | D | C

Eliminando as produções unitárias:

$$N_S = \{S, A, B, C, D\}$$

$$N_A = \{A, C, D, B\}$$

$$N_B = \{B, C, D\}$$

$$N_C = \{C, B, D\}$$

$$N_D = \{D, C, B\}$$

Gramática sem produções unitárias:

$$S ::= AB \mid aS \mid a \mid bA \mid BCD \mid b \mid CD \mid BD \mid BC \mid dB \mid d \mid cCc \mid cc$$

$$A ::= bA \mid BCD \mid b \mid CD \mid BD \mid BC \mid cCc \mid cc \mid d \mid dB$$

$$B ::= dB \mid d \mid cCc \mid BD \mid cc \mid CD$$

$$C ::= cCc \mid BD \mid cc \mid dB \mid d \mid CD$$

$$D ::= CD \mid d \mid cCc \mid BD \mid cc \mid dB$$

Exercício 9:

b) $S ::= BaS \mid \varepsilon$

$$B ::= AaB'$$

$$A ::= aA' \mid A'$$

$$B' ::= aSAab' \mid bB' \mid \varepsilon$$

$$A' ::= aB'aSaA' \mid \varepsilon$$

Obs: S não gerou S' porque não tinha recursão à esquerda direta. Em B houve a substituição do S mais à esquerda. Em A houve a substituição de S e depois de B.

Exercício 10:

a) $S ::= bcD \mid Bcd$

$$B ::= bB \mid b$$

$$D ::= dD \mid d$$

Substituindo B nas produções de S, temos:

$$S ::= bcD \mid bBcd \mid bcd$$

$$B ::= bB \mid b$$

$$D ::= dD \mid d$$

Eliminando o não-determinismo direto, temos:

$$S ::= bS'$$

$$S' ::= cD \mid Bcd \mid cd$$

$$B ::= bB'$$

$$B' ::= B \mid \varepsilon$$

$$D ::= dD'$$

$$D' ::= D \mid \varepsilon$$

Substituindo D em S', temos:

$$S ::= bS'$$

$$S' ::= cdD' \mid Bcd \mid cd$$

$$B ::= bB'$$

$$B' ::= B \mid \varepsilon$$

$$D ::= dD'$$

$$D' ::= D \mid \varepsilon$$

E, por último, tirar o não-determinismo direto:

$$S ::= bS'$$

$$S' ::= cdS'' \mid Bcd$$

$$S'' ::= D' \mid \varepsilon$$

$$B ::= bB'$$

$$B' ::= B \mid \varepsilon$$

$$D ::= dD'$$

$$D' ::= D \mid \varepsilon$$

Exercício 13:

Cálculo do FIRST e FOLLOW do resultado do exercício 9:

$$S ::= BaS \mid \varepsilon$$

$$B ::= AaB'$$

$$A ::= aA' \mid A'$$

$$B' ::= aSAab' \mid bB' \mid \varepsilon$$

$$A' ::= aB'aSaA' \mid \varepsilon$$

N	FIRST(N)	FOLLOW(N)
S	{a, ε}	{\$, a}
B	{a, ε}	{a}
A	{a, ε}	{a}
B'	{a, b, ε}	{a}
A'	{a, ε}	{a}

FIRST e FOLLOW da GLC:

$$S ::= ABC \mid aB \mid bA \mid \varepsilon$$

$$A ::= ab \mid Cb \mid cB \mid \varepsilon$$

$$B ::= bA \mid b \mid \varepsilon$$

$$C ::= cA \mid ABc \mid \varepsilon$$

N	FIRST(N)	FOLLOW(N)
S	{a, b, ε, c}	{\$}
A	{a, c, ε, b}	{b, c, a, \$}
B	{b, ε}	{c, a, b, \$}
C	{c, a, b, ε}	{b, \$}

CAPÍTULO V

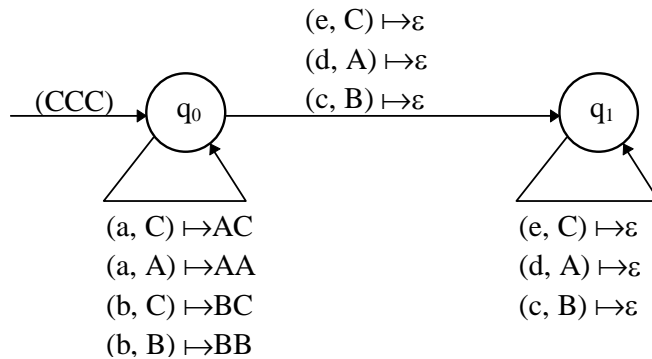
Exercício 1:

Os AP's aqui apresentados reconhecem as sentenças por pilha vazia:

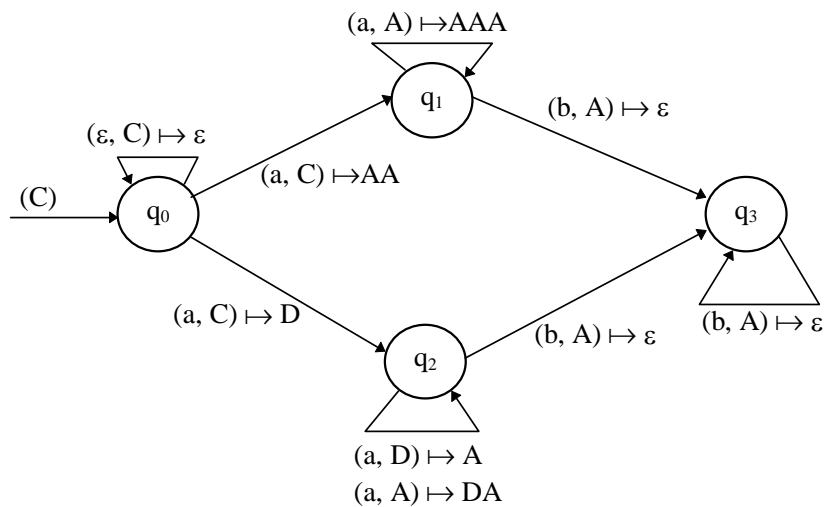
b) Alfabeto da pilha: A, B, C.

A = Empilhado quando reconhece a, para controlar o número de d's

B = Empilhado quando reconhece b, para controlar o número de c's
 C = Controla o número de e's
 Estado inicial da pilha: CCC (três C's empilhados).
 Forma gráfica:



- c) Alfabeto da pilha: A, C, D.
 A = Empilhado quando reconhece a, para controlar o número de b's.
 C = Utilizado para reconhecer a sentença vazia.
 D = usado para o caso de $2i = j$, em que um A deve ser empilhado para cada 2 a's reconhecido. Para o primeiro empilha-se um D, e para o segundo o D é trocado por A.
 Estado inicial da pilha: C.
 O AP é não-determinístico, já que não se sabe em qual dos dois casos a sentença irá se encaixar ($i = 2j$ ou $2i = j$).
 Forma gráfica:



Capítulo VI

Exercício 2:

c) Primeiramente, introduz-se a propriedade do prefixo na gramática, e numera-se as produções:

- 0 $S ::= \underline{B}\$$
- 1 $B ::= \underline{A}B'$
- 2, 3 $B' ::= \text{or } \underline{A}B' \mid \epsilon$
- 4 $A ::= \underline{C}A'$
- 5, 6 $A' ::= \text{and } \underline{C}A' \mid \epsilon$
- 7, 8, 9 $C ::= \text{exp} \mid \underline{(B)} \mid \text{not } C$

O algoritmo utilizado é o fechamento de estados, que já cria o AF determinístico. Assim, e_0 é iniciado com (o núcleo de cada estado é sublinhado):

$$e_0 = \{\underline{S} ::= \bullet \underline{B}\$\}$$

e o fechamento de e_0 resulta em:

$$e_0 = \{ \underline{S ::= \bullet B \$} \quad B ::= \bullet AB' \quad A ::= \bullet CA' \quad C ::= \bullet \text{exp} \quad C ::= \bullet (B) \quad C ::= \bullet \text{not } C \}$$

As transições de e_0 são:

$$t(e_0, B) = e_1 \quad t(e_0, A) = e_2 \quad t(e_0, C) = e_3 \quad t(e_0, \text{exp}) = e_5$$

$$t(e_0, () = e_4 \quad t(e_0, \text{not}) = e_6$$

Os núcleos dos novos estados são:

$$e_1 = \{ \underline{S ::= B \bullet \$} \}$$

$$e_2 = \{ \underline{B ::= A \bullet B'} \}$$

$$e_3 = \{ \underline{A ::= C \bullet A'} \}$$

$$e_4 = \{ \underline{C ::= (\bullet B)} \}$$

$$e_5 = \{ \underline{C ::= \text{exp} \bullet} \}$$

$$e_6 = \{ \underline{C ::= \text{not} \bullet C} \}$$

Os estados do autômato são:

$$e_0 = \{ \underline{S ::= \bullet B \$} \quad B ::= \bullet AB' \quad A ::= \bullet CA' \quad C ::= \bullet \text{exp} \quad C ::= \bullet (B) \quad C ::= \bullet \text{not } C \}$$

$$e_1 = \{ \underline{S ::= B \bullet \$} \}$$

$$e_2 = \{ \underline{B ::= A \bullet B'} \quad B' ::= \bullet \text{or } AB' \quad B' ::= \bullet \}$$

$$e_3 = \{ \underline{A ::= C \bullet A'} \quad A' ::= \bullet \text{and } CA' \quad A' ::= \bullet \}$$

$$e_4 = \{ \underline{C ::= (\bullet B)} \quad B ::= \bullet AB' \quad A ::= \bullet CA' \quad C ::= \bullet \text{exp} \quad C ::= \bullet (B) \quad C ::= \bullet \text{not } C \}$$

$$e_5 = \{ \underline{C ::= \text{exp} \bullet} \}$$

$$e_6 = \{ \underline{C ::= \text{not} \bullet C} \quad C ::= \bullet \text{exp} \quad C ::= \bullet (B) \quad C ::= \bullet \text{not } C \}$$

$$e_7 = \{ \underline{B' ::= AB' \bullet} \}$$

$$e_8 = \{ \underline{B' ::= \text{or} \bullet AB'} \quad A ::= \bullet CA' \quad C ::= \bullet \text{exp} \quad C ::= \bullet (B) \quad C ::= \bullet \text{not } C \}$$

$$e_9 = \{ \underline{A' ::= CA' \bullet} \}$$

$$e_{10} = \{ \underline{A' ::= \text{and} \bullet CA'} \quad C ::= \bullet \text{exp} \quad C ::= \bullet (B) \quad C ::= \bullet \text{not } C \}$$

$$e_{11} = \{ \underline{C ::= (B \bullet)} \}$$

$$e_{12} = \{ \underline{C ::= \text{not} \bullet C} \}$$

$$e_{13} = \{ \underline{B' ::= \text{or} \bullet AB'} \quad B' ::= \bullet \text{or } AB' \quad B' ::= \bullet \}$$

$$e_{14} = \{ \underline{A' ::= \text{and} \bullet CA'} \quad A' ::= \bullet \text{and } CA' \quad A' ::= \bullet \}$$

$$e_{15} = \{ \underline{C ::= (B) \bullet} \}$$

$$e_{16} = \{ \underline{B' ::= \text{or} \bullet AB' \bullet} \}$$

$$e_{17} = \{ \underline{A' ::= \text{and} \bullet CA' \bullet} \}$$

As transições do autômato são:

$$t(e_0, b) = e_1 \quad t(e_0, A) = e_2 \quad t(e_0, C) = e_3 \quad t(e_0, () = e_4$$

$$t(e_0, \text{exp}) = e_5 \quad t(e_0, \text{not}) = e_6 \quad t(e_2, B') = e_7 \quad t(e_2, \text{or}) = e_8$$

$$t(e_3, A') = e_9 \quad t(e_3, \text{and}) = e_{10} \quad t(e_4, B) = e_{11} \quad t(e_4, A) = e_2$$

$$t(e_4, C) = e_3 \quad t(e_4, \text{exp}) = e_5 \quad t(e_4, () = e_4 \quad t(e_4, \text{not}) = e_6$$

$$t(e_6, C) = e_{12} \quad t(e_6, \text{exp}) = e_5 \quad t(e_6, () = e_4 \quad t(e_6, \text{not}) = e_6$$

$$t(e_8, A) = e_{13} \quad t(e_8, C) = e_3 \quad t(e_8, \text{exp}) = e_5 \quad t(e_8, () = e_4 \quad t(e_8, \text{not}) = e_6$$

$$e_6 \quad t(e_{10}, C) = e_{14} \quad t(e_{10}, \text{exp}) = e_5 \quad t(e_{10}, () = e_4 \quad t(e_{10}, \text{not}) = e_6$$

$$t(e_{11}, ()) = e_{15} \quad t(e_{13}, B') = e_{16} \quad t(e_{13}, \text{or}) = e_8 \quad t(e_{14}, A') = e_{17}$$

$$t(e_{14}, \text{and}) = e_{10}$$

Para a montagem da tabela SLR(1), é necessário o FOLLOW de cada não-terminal:

N	FIRST(N)	FOLLOW(N)
B	{ exp, (, not }	{ \$,) }
B'	{ or, ε }	{ \$,) }
A	{ exp, (, not }	{ or, \$,) }
A'	{ and, ε }	{ or, \$,) }
C	{ exp, (, not }	{ and, or, \$ }

A tabela SLR(1) é a seguinte:

	or	and	exp	()	not	\$	B	B'	A	A'	C
0			S5	S4		S6		S1		S2		S3
1							PARE					
2	S8				R3		R3		S7			
3	R6	S10			R6		R6				S9	
4			S5	S4		S6		S11		S2		S3
5	R7	R7			R7		R7					
6			S5	S4		S6						S12
7					R1		R1					
8			S5	S4		S6				S13		S3
9	R4				R4		R4					
10			S5	S4		S6						S14
11					S15							
12	R9	R9			R9		R9					
13	S8				R3		R3		S16			
14	R6	S10			R6		R6				S17	
15	R8	R8			R8		R8					
16					R2		R2					
17	R5				R5		R5					

Para exemplificar, analise a sentença

not(exp) and (exp or exp)

Obs: quando a tabela SLR(1) mandar reduzir por uma produção do tipo $X ::= \epsilon$, nenhum símbolo é retirado da pilha, somente X é escrito na entrada. Por exemplo, a configuração $\{e_0 e_6 e_4 e_3,) \text{ and } (\text{exp or exp})\}$, consultada na tabela, leva a R6, ou reduzir pela produção 6, que é a produção $A' ::= \epsilon$. A nova configuração será $\{e_0 e_6 e_4 e_3, A') \text{ and } (\text{exp or exp})\}$

Exercício 4:

A mesma gramática é considerada, com suas produções numeradas:

- 1 $B ::= AB'$
- 2, 3 $B' ::= \text{or } AB' \mid \epsilon$
- 4 $A ::= CA'$
- 5, 6 $A' ::= \text{and } CA' \mid \epsilon$
- 7, 8, 9 $C ::= \text{exp} \mid (B) \mid \text{not } C$

Esta GLC é LL(1), pois

- a. não possui recursão à esquerda;
- b. está fatorada;
- c. para todo $X \in N$, se $X \rightarrow^* \epsilon$, então $\text{FIRST}(X) \cap \text{FOLLOW}(X) = \emptyset$.

Os conjuntos FIRST e FOLLOW são:

N	FIRST(N)	FOLLOW(N)
B	{ exp, (, not }	{ \$,) }
B'	{ or, ϵ }	{ \$,) }
A	{ exp, (, not }	{ or, \$,) }
A'	{ and, ϵ }	{ or, \$,) }
C	{ exp, (, not }	{ and, or, \$ }

A tabela de parsing é:

	or	and	exp	()	not	\$
B			1	1		1	

Respostas

B'	2				3		3
A			4	4		4	
A'	6	5			6		6
C			7	8		9	

Novamente, para exemplificar, analise a sentença:
not (exp and exp) and (exp or exp).

Referências Bibliográficas

- AHO, A. V., SETHI, R. and ULLMAN, J. D.: *COMPILERS - Principles, Techniques and Tool*. Addison-Wesley Publishing Company, 1986.
- HOPCROFT, J. E. and ULLMAN, J. D.: *Introduction to Automata Theory, Languages and Computation*. Addison-Wesley, 1979.
- MAK, Ronald. *Writing Compilers and Interpreters - An Applied Approach*. John Willey & Sons, Inc., 1991.
- SETZER, Valdermar W e MELO, Inês S. H. de. *A Construção de Um Compilador*. Rio de Janerio : Campus, 1986.
- TREMBLAY, Jean-Paul and SORENSON, Paul G. *The Theory and Practice of Compiler Writing*. McGraw-Hill, 1985.