

計算理論 第9回
第5章: 文脈自由文法と言語(2/2)

基礎工学部情報科学科
中川 博之

本日の概要

- 第5章: 文脈自由文法と言語 (の後半)
 - テキスト: p.216～
 - 5.3 文脈自由文法の応用
 - 5.4 文法と言語のあいまいさ
- 重要概念
 - 構文木, パーサー, あいまいな文法

(復習) 文脈自由言語, 文脈自由文法

- 文脈自由言語 (CFL)
 - 正則言語よりも広い言語クラス
 - 再帰的な構造を扱うことが可能
 - 正則言語では再帰的構造を扱うことができない
- 文脈自由文法 (CFG)
 - 文脈自由言語の記述に用いる文法
 - 以下の4つ組で定義: $G = (V, T, P, S)$
 - V: 変数(非終端記号)集合
 - T: 終端記号集合
 - P: 生成規則集合
 - S: 出発記号

5.3 文脈自由文法の応用

文脈自由文法の歴史

- Chomskyが考案
 - 言語学者
 - 元々は、自然言語を記述するために考案
 - 文脈依存文法、句構造文法なども考案
- 計算機科学にて再帰的構造の有用性が評価
 - 構造データ処理の自動化に使える
 - 応用: パーサー, マークアップ言語

パーサー (Parser)

- 構文解析器とも呼ばれる
 - ソースコードに記載された構文を解析する
- ソースコードでは再帰構造がよく出現
 - 数式における括弧: “(” と “)”
 - ブロックの範囲記述: “{” と “}” や begin-end
- 再帰構造の開始と終了は正しく対応する必要がある
 - ○: {{{}}}{}}
 - ×: {{{{}}}}

正則言語の限界

- バランスの取れた括弧の一例を表す記述
 - 言語 $L = \{ (^n)^n : n \geq 1 \}$
- 言語 L は正則言語ではない
 - $L' = \{ a^n b^n : n \geq 1 \}$ が正則言語でないのと同じ
 - 言語 L を表現した正則表現は存在しない
 - 正則表現では再帰構造を表現できない
 - プログラミング言語は正則表現では十分に表現できない
- 文脈自由文法であれば再帰構造を表現可能

バランスが取れた括弧を表す文法

- バランスが取れた括弧の列の例
 - ε , $(())$, $()()$, $((())(()))$
- 文法 $G_{\text{bal}} = (\{B\}, \{ (,) \}, P, B)$
- 生成規則の集合 P
 - $B \rightarrow BB$
 - $B \rightarrow (B)$
 - $B \rightarrow \varepsilon$

IF文を表す文法

- IF文にはelse節が無くて良い
- 文法 $G_{if} = (\{S\}, \{i, e\}, P, S)$
 - i はif節を, e はelse節を表す
- 生成規則の集合 P
 - $S \rightarrow \varepsilon \mid SS \mid iS \mid iSeS$
- 導出の例
 - $S \Rightarrow iS \Rightarrow iiSeS \Rightarrow iieS \Rightarrow iie$
 - $S \Rightarrow SS \Rightarrow iSeSS \Rightarrow iiSeeSS \Rightarrow iieeSS \Rightarrow iieeS$
 $\Rightarrow iieeiS \Rightarrow iieei$

yacc

- yacc: Parser-Generatorのひとつ
- Parser-Generator
 - 入力: 文脈自由文法
 - 出力: パーサー(構文解析器)のプログラムコード
 - 文法に合致するパーサープログラムを生成
- 代表的なParser-Generator
 - yacc, bison (yaccの上位互換), JavaCC

yaccの一例

Exp: Id {...}	Id : 'a' {...}
Exp '+' Exp {...}	'b' {...}
Exp '*' Exp {...}	Id 'a' {...}
'(' Exp ')' {...}	Id 'b' {...}
;	Id '0' {...}
	Id '1' {...}
	;

- 生成記号の “→” は “:” で表現
- 終端記号は引用符 (single quote (')) で囲む
- 引用符無しの文字は変数
- {...} 内には実行して欲しいコードを記述 (アクション)

マークアップ言語

- 各種マーク(タグ)を文書内に埋め込むための規則
 - plain text上で構造を記述する手法
- 例: HTML (Hypertext Markup Language)
 - Webページ用の文書ファイル記述言語
- 例: XML (eXtensible Markup Language)
 - 文書構造を記述するための言語

HTMLの例 (図5.12)

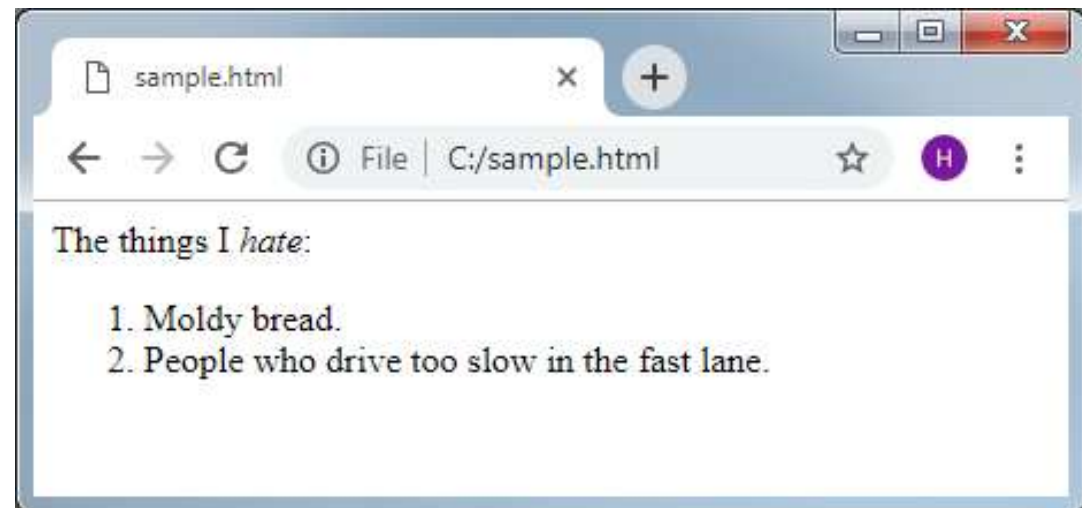
<P>The things I hate:

Moldy bread.

People who drive too slow in the fast lane.

HTML文法の一部

- Doc $\rightarrow \epsilon$ | Element Doc
- Element \rightarrow Text |
 - Doc |
 - <P> Doc |
 - List | ...
- List $\rightarrow \epsilon$ | ListItem List
- ListItem \rightarrow Doc
- ...



XML

- XML (eXtensible Markup Language)
 - DTD (Document-Type Definition)により文書構造を定義
 - DTDは本質的に文脈自由文法
 - DTDに従うXML文書のみが受理される

パソコン仕様記述用DTD (図5.14)

```
<!DOCTYPE PcSpecs[
  <!ELEMENT PCS (PC*)>
  <!ELEMENT PC (MODEL, PRICE, PROCESSOR, RAM, DISK+)>
  <!ELEMENT MODEL (\#PCDATA)>
  <!ELEMENT PRICE (\#PCDATA)>
  <!ELEMENT PROCESSOR (MANF, MODEL, SPEED)>
  <!ELEMENT MANF (\#PCDATA)>
  ...
  <!ELEMENT DISK (HARDDISK | CD | DVD)>
  ...
]>
```

CFG PcSpecsの生成規則:

- PCS → PC*
- PC → MODEL PRICE PROCESSOR RAM DISK+
- MODEL → 文字列
- ...
- DISK → HARDDISK | CD | DVD
- ...

DTDに従うXMLの記述例(図5.15)

<pre><PCS> <PC> <MODEL>4560</MODEL> <PRICE>\$2295</PRICE> <PROCESSOR> <MANF> Intel</MANF> <MODEL>Pentium </MODEL> <SPEED>800MHz</SPEED></pre>	<pre></PROCESSOR> <RAM>256</RAM> <DISK><HARDDISK> ... </PC> <PC> ... </PC> </PCS></pre>
---	---

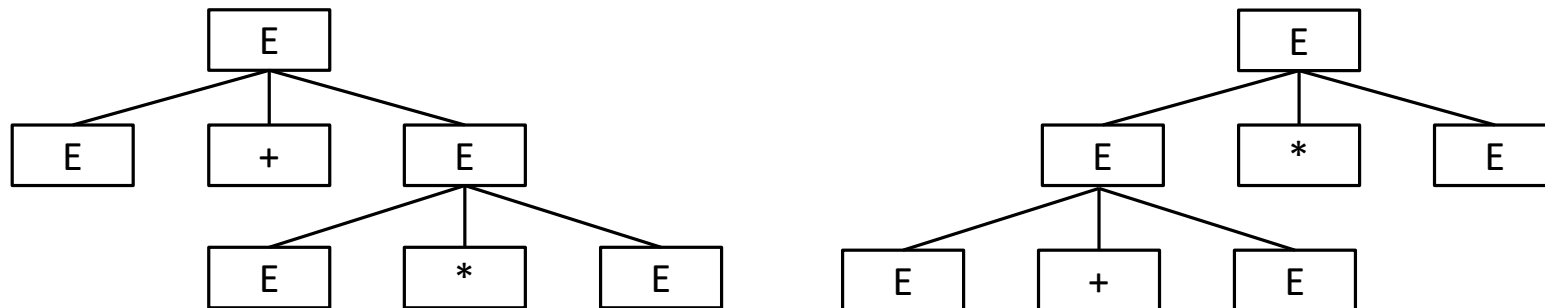
5.4 文法と言語のあいまいさ

あいまいな文法

- 定義: 文法 $G = (V, T, P, S)$ があいまいである
→ ある終端記号列 w に2つ以上の異なる構文木が存在
- 定義: 文法 $G = (V, T, P, S)$ があいまいでない
→ 任意の終端記号列 w の構文木は高々1つ存在

例5.25: 式文法 G_{exp}

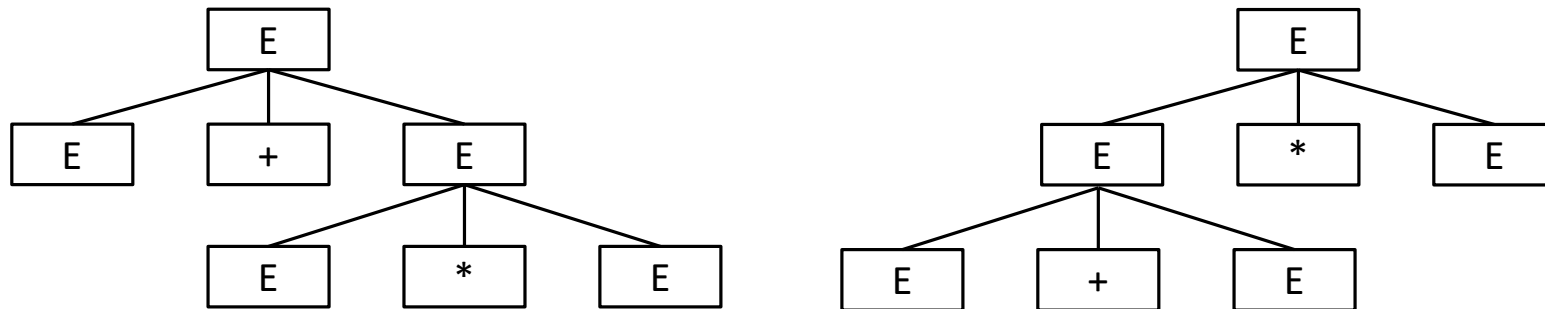
- (再掲) 式文法 G_{exp}
 - $E \rightarrow I \mid E+E \mid E^*E \mid (E)$
 - $I \rightarrow Ia \mid Ib \mid I0 \mid I1 \mid a \mid b$
- 文形式 $E+E^*E$ に対し2通りの構文木が存在



- 左の構文木: $E \Rightarrow E+E \Rightarrow E+E^*E$
- 右の構文木: $E \Rightarrow E^*E \Rightarrow E+E^*E$

あいまいだとななるか

- 複数の解釈が存在

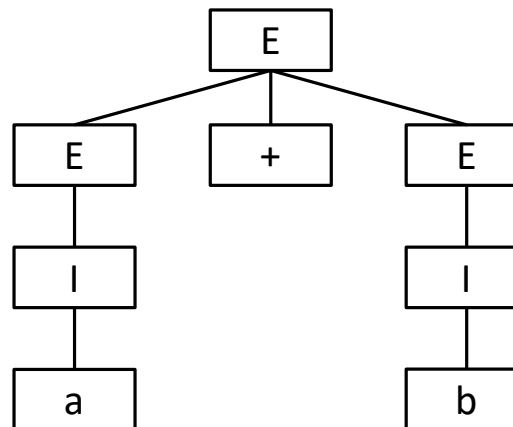


- 解釈1: $E + (E * E)$
- 解釈2: $(E + E) * E$

→ 意図と違う解釈, 計算結果となる可能性があるが,
文形式から正しい解釈(意味)を決定できない

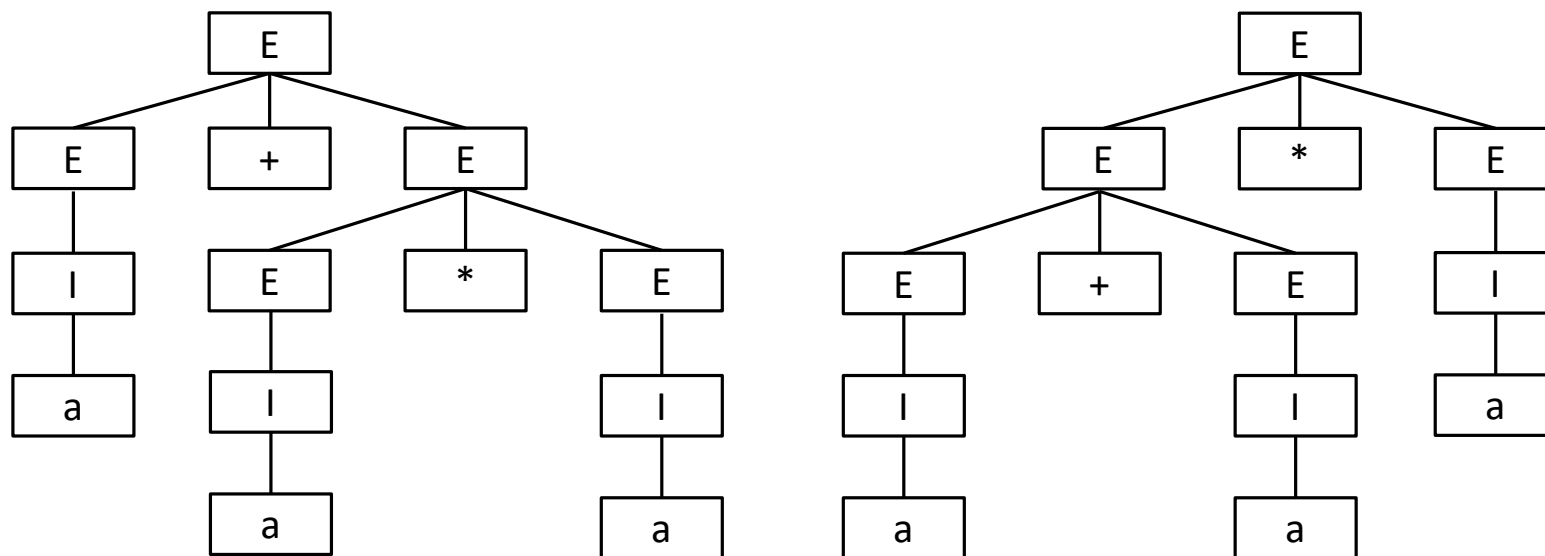
(参考) 複数の導出過程

- 1つの列に対する複数通りの導出
 - あいまいでない文法でもあり得る
- 例: 式文法 G_{exp} における $a+b$ の導出
 - $E \Rightarrow E+E \Rightarrow I+E \Rightarrow a+E \Rightarrow a+I \Rightarrow a+b$
 - $E \Rightarrow E+E \Rightarrow E+I \Rightarrow E+b \Rightarrow I+b \Rightarrow a+b$
 - 構文木は同じ



あいまいであることの証明

- ある終端記号列 w が存在して, 2つ以上の異なる構文木の存在を示せばよい
- 例: $w = a + a * a$



あいまいであることの判別

- 与えられたCFGがあいまいであるかどうかの判別
 - 判定アルゴリズムは存在しない
 - 決定不能であるという
- あいまいなCFGしか存在しないCFLもある
 - 後ほど紹介

あいまいさの除去

- 通常のプログラミング言語で現れるあいまいさは軽微
 - 除去のための良く知られた技法がある
- 先の数式の例: あいまいさの原因は？
 - 原因1: 演算子の優先度が考慮されていない
 - *より+が先に計算されることを許している
 - 原因2: 優先度が同一の演算子に対して結合順序が考慮されていない
 - $E+E+E$ に対し, $E+(E+E)$ と $(E+E)+E$ の両方を許している

あいまいさ除去の方針

- 方針1: 演算子の優先度を導入
 - $*$ は $+$ よりも優先度を高く
- 方針2: 同一優先度の演算子の結合順序を決定
 - 左結合性, 右結合性のいずれかに決めておく
 - $E+E+E$ に対して
 - 左結合性: $(E+E)+E$ と解釈する
 - 右結合性: $E+(E+E)$ と解釈する

具体的な手段

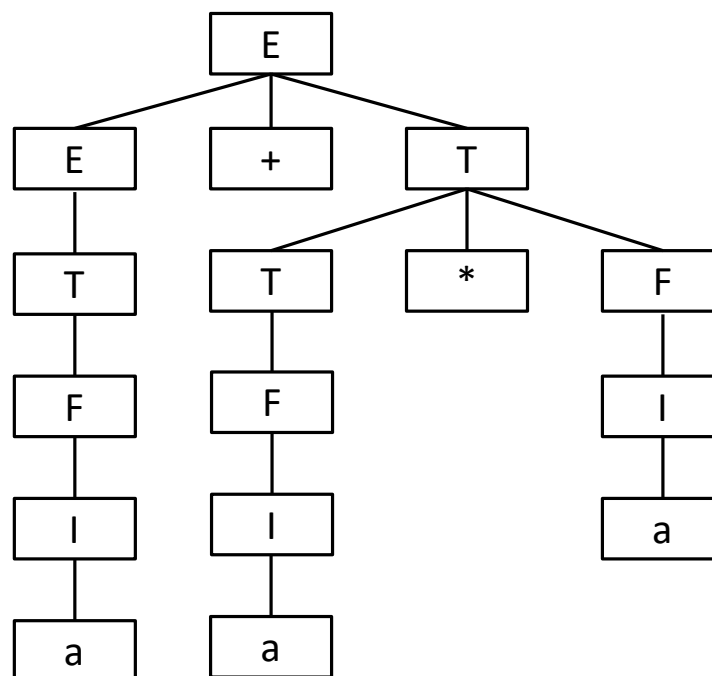
- 結合の強さごとに異なる変数を用意
 - $\langle \text{式} \rangle$ (E): 一つ以上の $\langle \text{項} \rangle$ の和
 - $\langle \text{項} \rangle$ (T): 一つ以上の $\langle \text{因数} \rangle$ の積
 - $\langle \text{因数} \rangle$ (F): $\langle \text{識別子} \rangle$ あるいは $\langle \text{式} \rangle$
 - $\langle \text{式} \rangle$: 括弧で囲まれた式

式に対するあいまいでない文法

- 以下の生成規則を持つ文法
 - $\langle \text{式} \rangle: E \rightarrow T \mid E + T$
 - $\langle \text{項} \rangle: T \rightarrow F \mid T * F$
 - $\langle \text{因数} \rangle: F \rightarrow I \mid (E)$
 - $\langle \text{識別子} \rangle: I \rightarrow a \mid b \mid Ia \mid Ib \mid I0 \mid I1$
 - この文法での演算子は左結合性を持つ

式に対するあいまいでない文法

- $a+a*a$ に対する構文木（これ1つしかない）



(参考)yaccでのあいまいさ対策

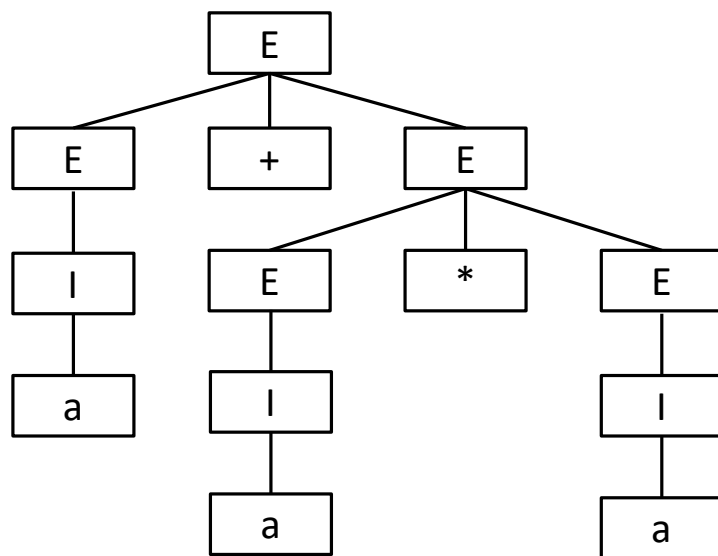
- 演算子の優先順位を指定できる
 - 順位の低いものから高いものへと記述
- 演算子の結合性(左結合or右結合)を指定できる
 - キーワード%left あるいは %right
 - %left ‘+’ *の方が+より優先順位が高く,
– %left ‘*’ +と*はいずれも左結合

あいまいさと最左(最右)導出

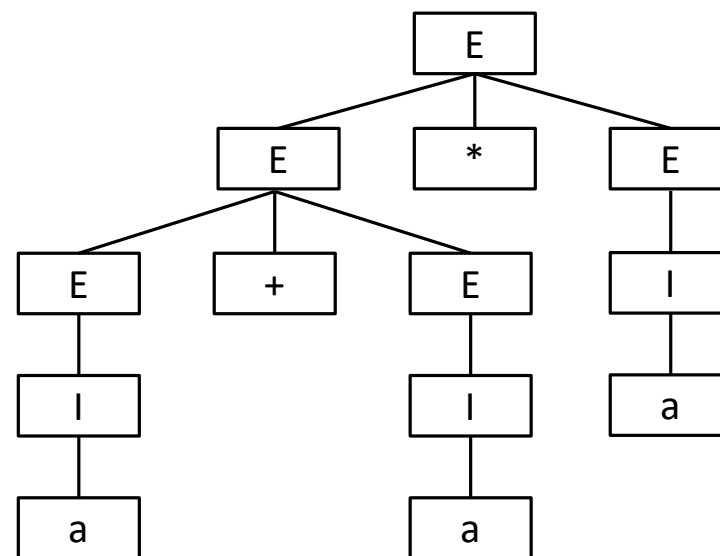
- 文法があいまいでなくても、導出はただ一つとは限らない
- ただし、あいまいでない文法では、最左導出も最右導出もそれぞれただ一つ
- [定理5.29] 文法 $G=(V,T,P,S)$ と終端記号列 w に対して、 w が2つの構文木を持つ
 $\Leftrightarrow S$ から w を導く2つの異なる最左導出が存在
 - 証明は省略
 - 最右導出についても同様

2つの構文木と最左導出

- (a) $E \xRightarrow{\text{左}} E+E \xRightarrow{\text{左}} I+E \xRightarrow{\text{左}} a+E \xRightarrow{\text{左}} a+E^*E \xRightarrow{\text{左}} a+I^*E$
 $\xRightarrow{\text{左}} a+a^*E \xRightarrow{\text{左}} a+a^*I \xRightarrow{\text{左}} a+a^*a$
- (b) $E \xRightarrow{\text{左}} E^*E \xRightarrow{\text{左}} E+E^*E \xRightarrow{\text{左}} I+E^*E \xRightarrow{\text{左}} a+E^*E \xRightarrow{\text{左}} a+I^*E$
 $\xRightarrow{\text{左}} a+a^*E \xRightarrow{\text{左}} a+a^*I \xRightarrow{\text{左}} a+a^*a$



(a)



(b)

本質的なあいまいさ

- [定義] 文脈自由言語Lが本質的にあいまい
 \Leftrightarrow Lを生成する文法のすべてがあいまい
 – Lを生成するあいまいでない文法が1つでも存在
 \Leftrightarrow Lは本質的にあいまいではない
- 本質的にあいまいな言語の例
 – $L = \{a^n b^n c^m d^m \mid n \geq 1, m \geq 1\} \cup \{a^n b^m c^m d^n \mid n \geq 1, m \geq 1\}$

Lの本質的なあいまいさの直観的説明

- Lを生成する文法G -

- $L = \{a^n b^n c^m d^m \mid n \geq 1, m \geq 1\} \cup \{a^n b^m c^m d^n \mid n \geq 1, m \geq 1\}$
- Lを生成する文法G
 - $S \rightarrow AB \mid C$
 - ABが1つ目の集合, Cが2つ目の集合に対応
 - $A \rightarrow aAb \mid ab$
 - $B \rightarrow cBd \mid cd$
 - $C \rightarrow aCd \mid aDd$
 - $D \rightarrow bDc \mid bc$

Lの本質的なあいまいさの直観的説明

- 複数の最左導出の存在 -

- このとき, $a^n b^n c^n d^n$ は以下の2通りの方法で導出できる
 - 導出1: $S \Rightarrow_{\text{左}} AB \Rightarrow_{\text{左}} aAbB \xRightarrow{*}_{\text{左}} a^n b^n B \Rightarrow_{\text{左}} a^n b^n cBd \xRightarrow{*}_{\text{左}} a^n b^n c^n d^n$
 - 導出2: $S \Rightarrow_{\text{左}} C \Rightarrow_{\text{左}} aCd \xRightarrow{*}_{\text{左}} a^n Dd^n \Rightarrow_{\text{左}} a^n bDcd^n \xRightarrow{*}_{\text{左}} a^n b^n c^n d^n$
- 最左導出(構文木)が2つ存在
→ この文法はあいまい

ミニレポート: 9-1

- テキスト p.231 問5.3.2:
 - (概要) 二つの型の括弧, つまり丸括弧() と角括弧[]のバランスの取れた列のすべて, またそれらだけを生成する文法を設計せよ.
- 例:
 - ()
 - []()
 - ([()]([]))[][[([[]])]((()))

ミニレポート: 9-2

- テキストp231 問5.3.5
- 図5.16(下図)のDTDを文脈自由文法に変換せよ.

```
<!DOCTYPE CourseSpecs [  
  <!ELEMENT COURSES (COURSE+)>  
  <!ELEMENT COURSE (CNAME, PROF, STUDENT*, TA?)>  
  <!ELEMENT CNAME (\#PCDATA)>  
  <!ELEMENT STUDENT (\#PCDATA)>  
  <!ELEMENT TA (\#PCDATA)>  
>
```