

## コスト最小法形態素解析のコストルールの作成方法

小松 英二                      安原 宏  
日本電子化辞書研究所      沖電気工業株式会社

コスト最小法形態素解析のコストルールの作成方法について述べる。コストルールの条件部は、既存の経験的規則や解析の失敗例を基に作成する。コストルールが与えるコストの値は、形態素解析の実行結果の誤り部分と正解とのペアを制約条件として決定する。本稿の方法を実際の形態素解析プログラムに適用した例も示す。

### A method to construct cost evaluation rules for morphological analyzers

Eiji Komatsu  
Japan Electronic Dictionary  
Research Institute,LTD.

Hiroshi Yasuhara  
Oki Electric Industry,Co.,LTD.

C/O system laboratory,Oki electric industry co.,LTD. 11-22,Shibaura 4-Chome,Minato-ku,Tokyo,Japan  
11-22,Shibaura 4-Chome,Minato-ku,Tokyo,Japan Yasuhara@okilab.oki.co.jp  
komatsu@edr6r.edr.co.jp

In this paper ,we propose a method to construct cost evaluation rules for Japanese morphological analyzers . As for condition parts of rules,our method makes them from existing inductive facts or examples of failure of a morphological analyzer. As for values of cost that each rule gives, our method decides values by using pairs of an error part of an result by a morphological analyzer and a right solution, as constraints.We also show examples that we applied this method to our own morphological analyzer.

## 1. はじめに

日本語形態素解析の曖昧性の解消方法には、最長一致法、二文節最長一致法、ヒューリスティック的方法、文節数最小法、コスト最小法等がある[1]。これらの方法は、処理上の違いはあるが、論理的にはすべて最後のコスト最小法の特別な場合である。

コスト最小法とは、全ての解を候補として生成し、単語及び単語の接続等に対して一定のコストを与え、コストの和が最も小さくなる単語列を解として選ぶ方法である。コスト最小法では、全解が必要なため、全解を効率よく生成するグラフスタック等の圧縮表現が研究されている[2,3,4,5]。また、コストの与え方[6,7]、形態素解析の曖昧性調査結果[8,9]についても多くの研究がある。

筆者らは、コスト最小法形態素解析プログラムを作成したが、コストを与える規則（以下、コストルールと呼ぶ）の作成にあたって、次のような2つの問題が生じた。

- (1) 単語のどのような属性に着目してコストを与えるか
- (2) 各コストルールの与えるコストの値をいくつにするか

(1)については、以下で述べるようにコストの値の決定方法を工夫することにより、逐次必要に応じて新しい属性を追加することで解決した。

(2)については、プログラム作成上必須であり、なんらかの定式的な方法を考えなければならなかった。各ルールの与えるコストの値の決定は、大量のデータによって統計的に決定することもできるが、筆者らのプログラムでは、単語同士の接続にもコストを与えるため、属性の組み合わせのデータも必要となり、すべての組み合わせの頻度を調べることは非常に労力が大きく、得られた頻度の意味付けもむずかしいことが予想される。さらに形態素解析プログラムのインプリメント上も記憶容量の点から効率的でない。

このような理由から、筆者らの作成した形態素解析プログラムでは、個々の実例を制約条件としてコストの値を決定する方法を検討した。

本論文は、沖電気工業(株)在職中の成果に基づいて書かれたものである。

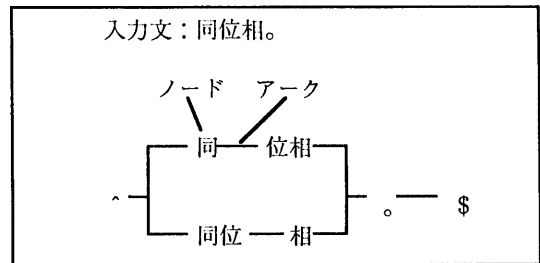


図1 グラフスタックを用いた形態素解析結果の例

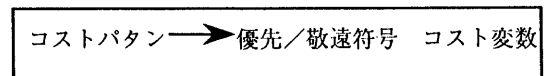


図2 コストルールの論理的なフォーマット

## 2. コスト最小法形態素解析プログラム

コストルール作成の対象として想定しているコスト最小法形態素解析プログラムについて述べる。

全解の表現形式としては、グラフスタック[1]を用いた。グラフスタックは、複数解の共通部分を共有することにより記憶の節約を行うデータ構造である。図1にグラフスタックを用いた形態素解析結果の例を示す。以下、グラフスタックの単語をノード、単語同士の接続をアークと呼ぶことにする。「^」と「\$」は、それぞれ文頭と文末を表すために便宜上つけ加えられるノードである。

コストは、ノード及びアークに対して与え、最もコストの和が小さくなる経路を選び、その経路の単語列を形態素解析結果として出力する。

辞書は、沖電気の機械翻訳用辞書に3文字から6文字のひらがな表記を加えたものを用いている。品詞数は約40である。

コストルールは、手続きとして実現したが、論理的には図2のようなプロダクションルールと等価である。以下の説明では、この論理的な形式のルールを扱う。ルールの左辺は、見出し、品詞、字種、文字数等必要に応じてあらゆる情報を用いて記述した条件部であり、ルールの右辺は、各ルールの与えるコストに対応する変数である。左辺をコストパターン、右辺をコスト変数と呼ぶことにする。コスト変数は、3節で述べる方法により、形態素解析実行時には確定した値に置き換えて用

いる。コスト変数は、ルール数を減らす目的のため、0を基準としており、優先したいボタンにはコスト変数にマイナスの符号を与え、敬遠したいボタンにはプラスの符号を与えている。この符号を、優先／敬遠符号と呼ぶことにする。

コストルールは、ノード用ルールとアーク用ルールを作成した。図3にコストルールの例を示す。例えば、同じ見出しの語について固有名詞より普通名詞を積極的に用いたい場合には、ノード用ルールとして、品詞を用いて、図3のルール1のようなルールを作る。また、一文字漢字列へ分割する防ぎたい場合には、アーク用ルールとして、文字数と字種を用いて、図3のルール2のようなルールを作る。

### 3. コストルールの作成方法

2節のような形態素解析プログラムのコストルールの作成方法について述べる。上記のコストルールの作成には、次の2段階の作業が必要である。

- (1) コストパターン、優先／敬遠符号の決定
- (2) コスト変数の値の決定

以下、それぞれの作業について述べる。

#### (1) コストパターン・優先／敬遠符号の選択

まず、従来の調査により得られている規則・内観に基づく規則を集め、コストパターンを作成する。また、解析結果の誤りと正解を弁別するのに必要なルールを逐次加えていく。コストパターンの選択は、作成者の主観・知識に依存する面が大きいため、コストパターンの妥当性の裏付けは別途行っておく必要が有る。ただし、不適切なボタンを登録しても、(2)で述べるコスト変数決定過程において、0に近い値が与えられるため、不適切なボタンは自然に排除できる。

#### (2) コスト変数の値の決定

(1)で作成したコストパターンに対応するコスト変数の値を決定する方法について述べる。

形態素解析の結果に誤りがある場合に、誤った解析結果を取り出し、正しい解析結果を加えたペアを作成する。これを、制約データと呼ぶことにする。図4に制約データの例を示す。

制約データは、さらに、コストルール作成用の

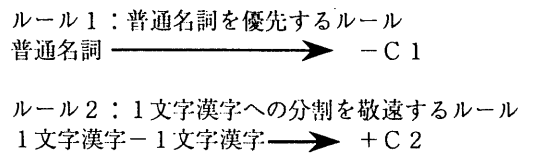


図3 コストルールの例

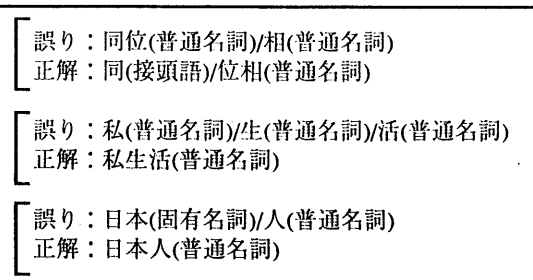


図4 制約データの例

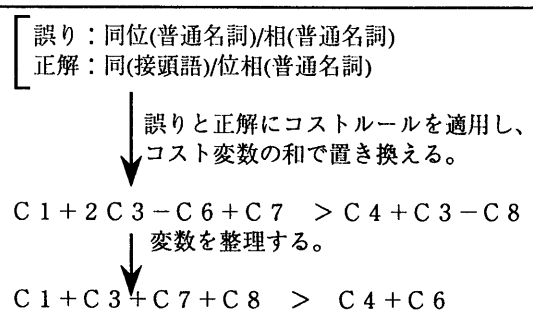


図5 制約データから制約不等式への変換の例

プログラムによりコスト評価を行い、誤りと正解の両方をコスト変数の和に置き換える。正解のコストは誤りのコストより小さくなければならないから、前者が後者より小さくなるように不等号で結び、コスト変数の不等式に変換する。この変換された不等式を制約不等式と呼ぶことにする。図5に制約データから制約不等式への変換の例を示す。

また、制約データを蓄積したものを制約データベースと呼ぶことにする。

コスト変数の決定は、まず、コスト変数に適当な初期値を設定し、形態素解析を行い、制約データを収集し、制約データベースに追加する。制約データベースの制約データをすべて制約不等式に変換し、これらを満たすコスト変数の値を決定し、得られたコスト変数を用いて、形態素解析を

行う。この作業を繰り返し、制約データを収集しながら、コスト変数値を修正していく。もし、制約データベースに対してコスト変数の値が存在しなくなった場合には、制約データを削除するか、(1)に戻って、コストルールを変更(追加・削除・修正)する。コストルールの変更があった場合には、制約不等式を新たに生成する必要がある。図6に上記の処理の概略フローを示す。また、図7にコスト変数の決定の各処理の内容を示す。

#### 4. コストルールの作成例

3節の方法で、コスト変数の値を決定するようすを実例で示す。図8に初期状態を示す。コストルールは、文節数最小、単語数最小の規則に、字種による未知語の規則を加えた簡単な規則である。この時点でできるだけ多くのルールを用意しておいた方が望ましい。例文1として、図9のような文を選ぶ。コスト変数の初期値として全て1を設定する。

プログラムのコスト変数を修正し、1回目の解析を行う。図10に例文1-1回目の実行結果を示す。図11のような過程を経て、コスト変数の値を修正する。

プログラムのコスト変数の値を新しい値に変更し、形態素解析を行い2回目の解析を行う。図12に例文1-2回目の実行結果を示す。図13のように、新たに得られた制約データを制約データベースに加え、コスト変数の値を決定する。コスト変数の値を変更して形態素解析を実行すると、図14のように正しい結果が得られる。

次に例文2として図15のような文を選ぶ。図16は例文2-1回目の実行結果である。今度は、図17に示すように、新たに得られた制約不等式が解けないため、図18のようなルールを加えると、図19のようにコスト変数の値が決定でき、図20のように形態素解析の結果も正しくなる。この場合、例文1の結果が正しいことは保障するためには再度チェックもする必要がある。図21にコストルールの最終状態を示す。また、図22に制約データベースの最終状態を示す。

上記の例では、優先/敬遠符号は+しか現われていないが「-」の場合もある。また、コストルールの修正や制約データの削除などが必要になる場合もある。制約データ間の矛盾は制約不等式においてチェックする。

実際のコスト作成に当たっては、いくつかの例文をまとめて選んで制約データを作成した。また、上記の例では、誤りを直すために主観的にもっともらしいと思われるルールを逐次作成したが、実際のルール作成時には、最初に用意するルールと同様に十分に調査をしたうえで、確実度の高いルールを追加していくことがコスト変数を早く収束させるために必要である。

#### 5. おわりに

コスト最小法形態素解析のコストルールの使われ方はルール同士の兼合により結果的に非決定的にな利、誤りに対応してプログラムを修正することとは個となる。また、逆にルールの妥当性をコストの値として得ることができるため、統計的データ収集の前段階として、必要な属性の洗いだしや基礎的なルールの抽出に有効であると考えている。

本稿の方法により、約60ルールが得られ、新聞記事で95%以上の解析精度を得ている。

本稿の方法の応用として、構文解析のコストルールの作成、あるいは、文書要約における重要文の抽出のルール作成などを考えている。

#### [参考文献]

- [1]長尾真監修：日本語情報処理、(社)電子通信学会、1984
- [2]杉村：グラフスタックを用いた日本語形態素解析、情報処理学会第37回全国大会、1988
- [3]杉村他：論理型形態素解析LAX、Proceedings of the logic programming conference '88、1988
- [4]久保他：LTB形態素解析LAXの開発環境、情報処理学会第37回全国大会、1988
- [5]吉村他：未知語を含む日本語文の形態素解析、情報処理学会論文誌、Vol.30 No.3、1989
- [6]久光他：接続コスト最小法による日本語形態素解析、情報処理学会第42回全国大会、1991
- [7]N.Maruyama, etc : A Japanese sentence analyzer、IBM J. RES. DEVELOP. Vol.32 No.2、1989
- [8]芦沢他：日英機械翻訳用前編集支援システム(2)形態素の曖昧性の検出方式、情報処理学会第36回全国大会、1989
- [9]坂本：文節の認定、日本語情報処理シンポジウム報告書、1978

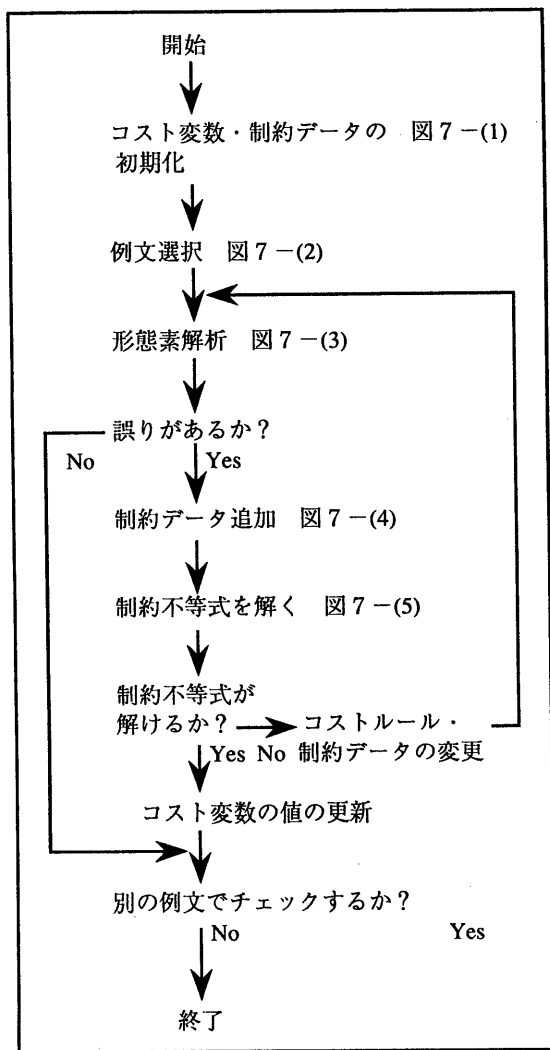


図6 コスト変数の決定の処理フロー

- 1)コストルールのコスト変数に初期値を設定する。初期値は、既にコスト変数の値が決定されている場合はその値を用い、新たに決定する場合には適当な正值を初期値として設定する。制約データベースとしては、それまでにコスト変数の値を決定するために用いた全ての制約データを用いる。
- (2)例文は1文ずつ選んでも、まとめて数文選んでもよい。
- (3)既に得られているコスト変数の値を用いて、形態素解析を行う。
- (4)形態素解析の結果に誤りがある場合は、誤りの前後を含む部分を取り出し、制約データベースに追加する。
- (5)制約データベースの制約データをすべて制約不等式に変換し、連立不等式を満たす値の組を1つ決定する。制約不等式を満たす解がない場合は、互いに矛盾する制約データを発見し、制約データの変更を行う。変更の方法は次の3つが考えられる。
  - コストルールの追加：  
誤りと正解を弁別する要因に関するコストルールがないと判断される場合に、新たな要因に関するコストルールを追加する。
  - コストルールの削除：  
確信度の低いルールは、コスト変数の収束を遅らすため、削除する。
  - コストルールの修正：  
コストルールのパターンが例外を含んだり、優先/敬遠の指定が間違っていると判断される場合に、ルールのパターンが指定する条件を厳しくしたり、優先/敬遠の指定を変更したりする。
- 制約データの削除：  
他の制約データと矛盾するがどちらの制約データも同じ位の頻度で出現し、かつ、コストルールの変更でも対処できないと思われる場合には、制約データを削除する。

図7 コスト変数の決定の各処理の内容

コストルール：

ルール1：文節数の少ない解を選ぶルール

単語-単語， {文節の切れ目} →+C1

ルール2：単語数の少ない解を選ぶルール

単語-単語 →+C2

ルール3：未知語の範囲を決定するルール

未知語-未知語， {同一字種} →+C3

ルール4：未知語の範囲を決定するルール

未知語-未知語， {異なる字種} →+C4

ルール5：未知語の範囲を決定するルール

登録語-未知語， {同一字種} →+C5

ルール6：未知語の範囲を決定するルール

登録語-未知語， {異なる字種} →+C6

ルール7：未知語の範囲を決定するルール

未知語-登録語， {同一字種} →+C7

ルール8：未知語の範囲を決定するルール

未知語-登録語， {異なる字種} →+C8

コスト変数：

C1=1, C2=1, C3=1, C4=1,

C5=1, C6=1, C7=1, C8=1

制約データベース：空

図8 初期状態

例文1：為替相場をいかにして安定させるか。

図9 例文の選択 (例文1)

形態素解析結果：

為替相場(普通名詞)/をいかにして(未知語)/安定(サ変名詞)//させ(下一段動詞)/る(一段動詞語尾)/か(終助詞)。(記号)

図10 実行結果 (例文1-1回目)

追加した制約データ：

誤り：為替相場(普通名詞)/をいかにして(未知語)/安定(サ変名詞)

正解：為替相場(普通名詞)/を(格助詞)/いかに(副詞)//し(サ変動詞)/て(接続助詞)/安定(サ変名詞)

誤り：安定(サ変名詞)//させ(下一段動詞)/る(一段動詞語尾)

正解：安定(サ変名詞)/さ(サ変名詞語尾)/せ(助動詞)/る(一段動詞語尾)

制約不等式：

$$2C2+C6+C8 > C1+5C2 \implies C6+C8 > C1+3C2$$

$$C1+2C2 > 3C2 \implies C1 > C2$$

コスト変数の値：

C1=2, C2=1, C3=1, C4=1,

C5=1, C6=3, C7=1, C8=3

図11 コスト変数決定 (例文1-1回目)

形態素解析結果：

為替相場(普通名詞)/を(格助詞)/いかに(普通名詞)/にして(未知語)/安定(未知語)/させる(未知語)/か(終助詞)。(記号)

図12 実行結果 (例文1-2回目)

追加した制約データ：

誤り：いかに(普通名詞)/にして(未知語)/安定(未知語)/させる(未知語)/か(終助詞)

正解：いかに(副詞)//し(サ変名詞)/て(接続助詞)//安定(サ変名詞)/ さ(サ変名詞語尾)/せ(助動詞)/る(一段動詞語尾)/か(終助詞)

制約不等式：

$$2C2+C6+C8 > C1+5C2 \implies C6+C8 > C1+3C2$$

$$C1+2C2 > 3C2 \implies C1 > C2$$

・・・以上前の制約データからの制約不等式

$$4C2+2C3+2C4+C5+C7 > 2C1+7C2$$

$$\implies 2C3+2C4+C5+C7 > 2C1+3C2$$

・・・追加した制約データからの制約不等式

コスト変数の値：

C1=1, C2=1, C3=1, C4=1,

C5=2, C6=3, C7=2, C8=3

図13 コスト変数決定 (例文1-2回目)

形態素解析結果：  
 為替相場(普通名詞)/を(格助詞)//いかに(副詞)//し  
 (サ変動詞)/て(接続助詞)//安定(サ変名詞)//させ(下  
 一段動詞)/る(一段動詞語尾)/か(終助詞)。(記号)  
 誤りなし。

図 1 4 実行結果 (例文 1 - 3 回目)

例文 2：そして累積債務を抱えた開発途上国対  
 策をどうするか。

図 1 5 例文の選択 (例文 2)

形態素解析結果：  
 そして(順接接続詞)//累積(サ変名詞)/債務(普通名  
 詞)/を(格助詞)//抱え(普通名詞)/た(接頭語)/開発途  
 上国(普通名詞)/対策(普通名詞)/を(格助詞)//どう  
 (副詞)//する(サ変動詞)/か(終助詞)。(記号)

図 1 6 実行結果 (例文 2 - 1 回目)

追加した制約データ：

誤り：抱え(普通名詞)/た(接頭語)/開発途上国(普  
 通名詞)  
 正解：抱え(下一段動詞)/た(助動詞)/開発途上国  
 (普通名詞)

制約不等式：  
 $2C2+C6+C8 > C1+5C2 \implies C6+C8 > C1+3C2$   
 $C1+2C2 > 3C2 \implies C1 > C2$   
 $4C2+2C3+2C4+C5+C7 > 2C1+7C2$   
 $\implies 2C3+2C4+C5+C7 > 2C1+3C2$   
 ・・・・以上は前の制約データからの制約不等式  
 $2C2 > 2C2 \implies 0 > 0$   
 ・・・・追加した制約データからの制約不等式

コスト変数は決定できない。

図 1 7 コスト変数決定 (例文 2 - 1 回目)

ルール 9：/\*普通名詞と接頭語の接続を敬遠する  
 ルール\*/  
 @名詞-@接頭語 → +C9

注意) @は、複数の品詞のグループを表わす。

図 1 8 追加したコストルール

制約不等式：  
 $2C2+C6+C8 > C1+5C2 \implies C6+C8 > C1+3C2$   
 $C1+2C2 > 3C2 \implies C1 > C2$   
 $4C2+2C3+2C4+C5+C7 > 2C1+7C2$   
 $\implies 2C3+2C4+C5+C7 > 2C1+3C2$   
 ・・・・以上は前の制約データからの制約不等式  
 $2C2 + C9 > 2C2 \implies C9 > 0$   
 ・・・・追加した制約データからの制約不等式

$C1 = 1, C2 = 1, C3 = 1, C4 = 1,$   
 $C5 = 2, C6 = 3, C7 = 2, C8 = 3,$   
 $C9 = 1$

図 1 9 コスト変数決定  
 (例文 2 - 1 回目、ルール 9 追加後)

形態素解析結果：  
 そして(接続詞)//累積(サ変動詞)/債務(普通名詞)/  
 を(格助詞)//抱え(下一段動詞)/た(助動詞)//開発途  
 上国(普通名詞)/対策(普通名詞)/を(格助詞)//どう  
 (副詞)//する(サ変動詞)/か(終助詞)。(記号)

誤りなし。

図 2 0 実行結果 (例文 2 - 2 回目)

コストルール：

ルール 1：文節数の少ない解を選ぶルール

単語-単語, {文節の切れ目} → +C 1

ルール 2：単語数の少ない解を選ぶルール

単語-単語 → +C 2

ルール 3：未知語の範囲を決定するルール

未知語-未知語, {同一字種} → +C 3

ルール 4：未知語の範囲を決定するルール

未知語-未知語, {異なる字種} → +C 4

ルール 5：未知語の範囲を決定するルール

登録語-未知語, {同一字種} → +C 5

ルール 6：未知語の範囲を決定するルール

登録語-未知語, {異なる字種} → +C 6

ルール 7：未知語の範囲を決定するルール

未知語-登録語, {同一字種} → +C 7

ルール 8：未知語の範囲を決定するルール

未知語-登録語, {異なる字種} → +C 8

ルール 9：名詞と接頭語の接続のルール

@名詞-@接頭語 → +C 9

コスト変数：

C 1 = 1, C 2 = 1, C 3 = 1, C 4 = 1,

C 5 = 2, C 6 = 3, C 7 = 2, C 8 = 3,

C 9 = 1

図 2 1 コストルール・コスト変数の最終状態

誤り：為替相場(普通名詞)/をいかにして(未知語)  
/安定(サ変名詞)

正解：為替相場(普通名詞)/を(格助詞)いかに(副  
詞)//し(サ変動詞)/て(接続助詞)/安定(サ変  
名詞)

誤り：安定(サ変名詞)//させ(一段動詞)/る(一段  
動詞語尾)

正解：安定(サ変名詞)/さ(サ変名詞語尾)/せ(助動  
詞)/る(一段動詞語尾)

誤り：いか(普通名詞)/にして(未知語)/安定(未知  
語)/させる(未知語)/か(終助詞)

正解：いかに(副詞)//し(サ変名詞)/て(接続助詞)//  
安定(サ変名詞)/ さ(サ変名詞語尾)/せ(助  
動詞)/る(一段動詞語尾)/か(終助詞)

誤り：抱え(普通名詞)/た(接頭語)/開発途上国(普  
通名詞)

正解：抱え(一段動詞)/た(助動詞)/開発途上国  
(普通名詞)

図 2 2 制約データベースの最終状態