

日本語文係り受けの分析・解析・学習

白井克彦 林 良彦 平田裕一

(早稲田大学 理工学部)

1. はじめに

複雑でかつ多様な自然言語の現象に対処可能な処理システムを実現しようとする際、それに必要な多様なレベルの知識をどのようにして獲得し、利用するかについて、多くの有効な処理モデルが提案されているが、なお大きな問題である。〔1〕

我々は、人手により加工された形で蓄積された実テキストデータに対する分析により直接的に初期的な知識構造を構成し、これを文解析に使用することにより、この知識構造の構成手法を評価するとともに、この過程を通して得られる情報を、知識構造に反映させるというアプローチを日本語文の係り受け構造に適応してきた。〔2〕～〔4〕本アプローチの概要を Fig. 1 に示す。

ここでいう知識とは、対象世界における言語要素の使われ方に関するものであり、辞書データベースとして記述的な形で構造化される。このように、使われ方

に注目して、知識の獲得・構造化を行なうため、対象世界における拘束を緩やかに含み、かつインクリメンタルに成長しうる知識構造の構成が可能となる。

2. 初期辞書データベースの構成

初期辞書データベースは、実テキストデータの分析により構成されるものであり、以後の成長に対してその核となるものである。

本辞書中では、分析対象の文中に存在する係り受け関係を構成する単語間の依存関係は、単語の集合であるクラスタ間の係り受け可能関係として構造化される。この知識構造を Fig. 2 に示す。以下、実データの分析により、この知識構造を構成する手法について述べるが、その概要を Fig. 3 に示す。

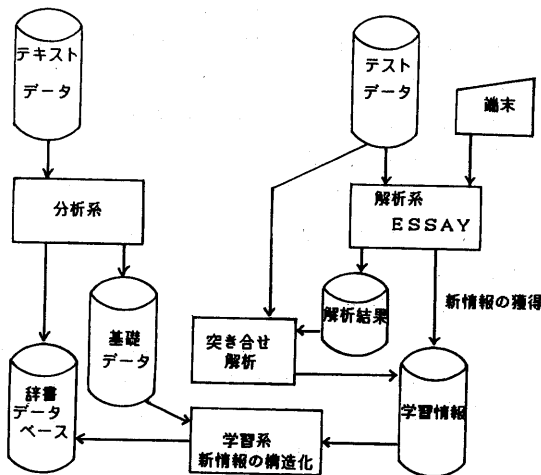


Fig. 1 アプローチの概要

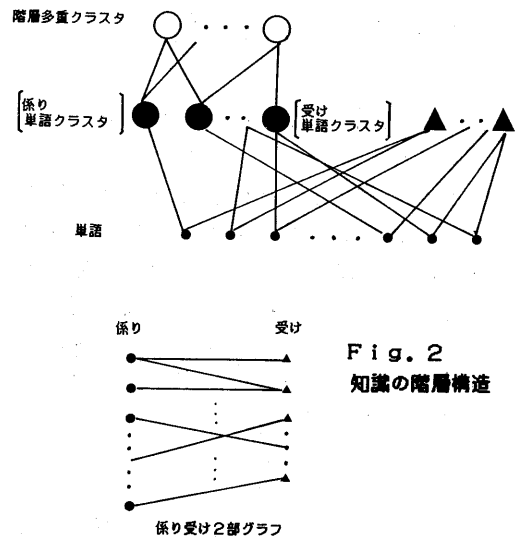


Fig. 2 知識の階層構造

2.1 対象テキスト

含まれる単語数が比較的少ない、複雑な文が少ない等の理由により、小学校の教科書を対象とし、係り受け構造等を付加し約2500文を登録した。

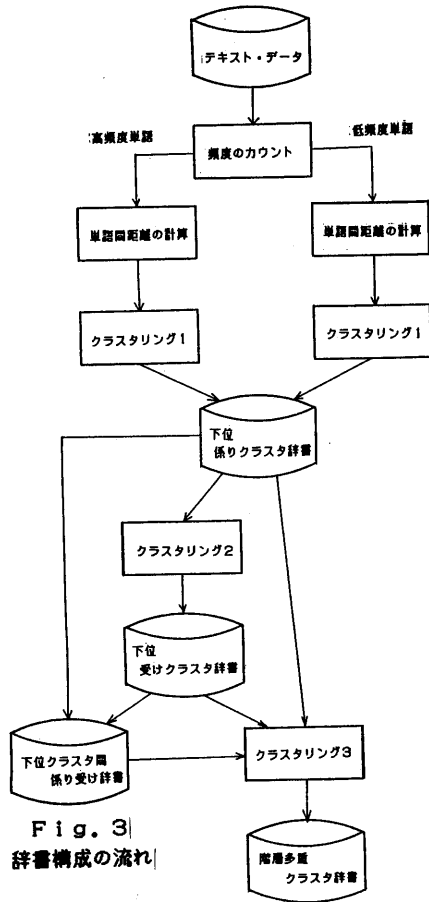


Fig. 3 辞書構成の流れ

2. 2 単語クラスタ辞書の構成

テキストデータ中に存在する係り受け関係を、Fig. 2に示したように構造化することにより、テキストデータ中には実際には存在しない係り受け関係をも処理しようという抽象化の効果が期待できる。単語は、1つの係りクラスタ、受けクラスタ及び複数の階層多重クラスタに属することになる。

2. 2. 1 単語間距離行列の計算

クラスタリングにより単語分類を行なうためには、単語間の距離を計算しておくことが必要となる。そこで係り単語に対しては、同じ付属語を介して同じ単語により多くの頻度で係っている単語間の距離が近いと考える事により、単語A, B間の距離を次の様に定義し、計算を行

なう。

$$D(A,B) = 1 - \frac{\sum [M(A,f,a) + M(B,f,a)]}{\sum F(A,c) + \sum F(B,d)} \quad (1)$$

ここで

$M(\alpha, \beta, \gamma)$: 単語 α が付属語 β を介して単語 γ に係る頻度

$f \in K(A) \cap K(B)$

$K(a)$: 単語 a につく付属語の集合

$F(a,b)$: 単語 a から単語 b への係り受け頻度

である。

2. 2. 2 単語の出現頻度の考慮

現在テキスト中出现する係り単語数は約1500、受け単語数は約1000であるが、テキストに1度しか現れない単語は、係り単語で約40%、受け単語で約44%にもぼっている。このような出現頻度の低い単語、つまり情報量の少ない単語は、クラスタリングの際の雑音的存在となってしまう。そこでクラスタリングを低頻度単語と高頻度単語とで別々に行なった。これにより、クラスタリングにおける雑音を除去し、低頻度単語の少ない情報を有効に使用する事が可能となる。

2. 2. 3 しきい値の定め方

上で求めた距離行列を用いてクラスタリングを行なうが、問題となるのがクラスタリングの際のしきい値である。そこで、このしきい値を定める目安として、品詞分散と距離歪の2つの尺度を用いる。

品詞分散 E_c は、クラスタ内の単語の品詞的なまとまりを表す値で次のように定義する。

単語品詞ベクトル H_T

$$H_T = (x_1, x_2, \dots, x_7) \quad (2)$$

x_i : 7つの品詞カテゴリー(名詞, 形容詞, 連体詞, 代名詞, 副詞, 動詞, 接続詞)に対応し、単語の品詞に該当する要素が1で他は0となる。

1つのクラスタ内の品詞分散 E_{cc} は

$$E_{cc} = \frac{N_{CL}}{T=1} |H_{CL} - H_T|^2 \quad (3)$$

但し

$$H_{CL} = \frac{N_{CL}}{T=1} H_T / N_{CL} \quad (4)$$

$$|H_{CL} - H_T|^2 = \sum_{i=1}^7 (x_i - x_i)^2 \quad (5)$$

N_{CL} : クラスタ内に含まれる単語数
全クラスタに対して E_{cc} を求め、その和を
求める事により品詞分散 E_c が得られる。

$$E_c = \sum_{c=1}^Q E_{cc} \quad (6)$$

距離歪 E_d は、クラスタリングを行なっ
た後に同一のクラスタに含まれた単語の
距離は0と見なされ、単語間距離情報と
の誤差が生じるが、この誤差は元の情報
の保存の指標とも考えられる。

1つのクラスタ内の距離歪 E_{dc} は

$$E_{dc} = \sum_{i=1}^{N_{CL}-1} d_{ij} \quad (7)$$

d_{ij} : 単語 i がこのクラスタに含ま
れた時の距離

全クラスタに対して E_{dc} を求め、その和を
求める事により距離歪 E_d が得られる。

$$E_d = \sum_{c=1}^Q E_{dc} \quad (8)$$

ここでクラスタリングのしきい値をパ
ラメータとして、ある単語距離行列に対
して、クラスタ数、品詞分散、距離歪の
関係をグラフに表わした例を Fig. 4
に示す。

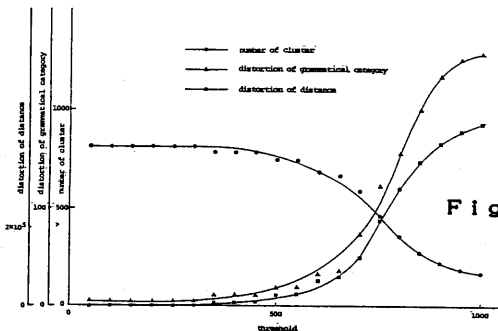


Fig. 4

Fig. 4を参照する事により、クラスタ
リングのしきい値を、品詞分散、距離
歪の急激な変化が現れる以前の約0.3
付近に定めればよいと考えられる。

2. 2. 4 係り単語のクラスタリング

Fig. 3のクラスタリング1では係
り単語を先に分類している。これは、係
り単語数が受け単語数より多く、又受け
単語は同時に多くの単語から係られ、係
り受けの特徴が現われにくいという事を
考えたためである。ここで用いるアルゴ
リズムは、時定のクラスタに多数の単語
が集中することを避けるためCentroid
法[6]とした。

2. 2. 5 受け単語のクラスタリング

Fig. 3のクラスタリング2では受
け単語を分類する。前に述べたように、
受け単語は、多くの意味合いで係られる
ために品詞分散を調べると、非常に高い
値と成りやすい。そこでクラスタリング
の際の時と同じ意味でCentroid
法とした。

2. 3 階層多重クラスタ辞書の構成

日本語の係り受けは非常に複雑で、曖
昧な構造を持っている。このような構造を
捕らえるには1回のクラスタリングでは
不十分だと思われる。そこで前に示した
係りクラスタをもう一度分類し、しかも
1つの下位クラスタは複数の上位クラス
タに含まれるような階層多重クラスタを
考えることにより、係り受け構造をより
大きく捕らえることが可能となる。

2. 3. 1 階層多重クラスタの構成

Fig. 3のクラスタリング3では、先
のクラスタリングで構成した係りクラス
タを再度分類する。

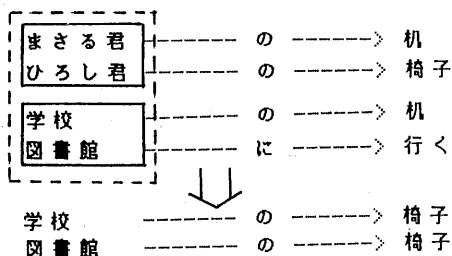
その方法としては、同じ受けクラスタに、
同一の付属語を介して係っているクラス
タは1つの上位クラスタに含まれるよう
にした。これは同じ係り受けが存在する
以上、そこには広い意味でなんらかの共
通点が存在するはずであるという仮定に

よるものである。

下位クラスは、複数の上位クラスに属するために、どの上位クラスとより密接な関係があるかが順位づけされている。

例を Fig. 5 に示す。点線で囲まれているのが、階層多重クラスである。

例 1



例 2

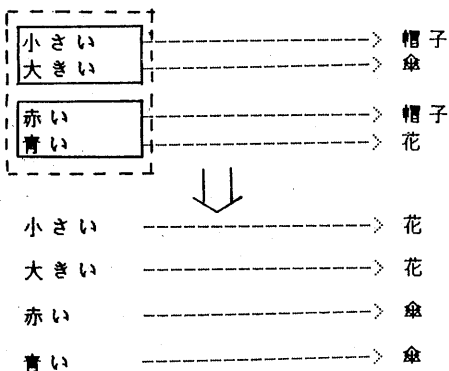


Fig. 5

2. 3. 2 階層多重クラスタの利用

クラスタを階層的にすることにより、曖昧度をそれほど増さずに、係り受けの抽象化が行なえる。つまり、ある係り受け情報を検索する場合、まず下位クラスレベルで検索をし、その情報がそこに存在しなかった時のみ上位クラスを探索すればよいわけで、しきい値を大きくして、大きなクラスタを構成するのにくらべて、かなり曖昧度が押さえられる。

3. 文解析実験による辞書の評価

3. 1 文解析システム ESSAY

システム ESSAY は、辞書駆動型のシステムであり、2. に述べた手法により構成される辞書データベースを、その知識源として用いる。システム中に手続的な形で組み込まれている知識は、日本語文の係り受けに関する原則が中心であり、文法的あるいは意味的な知識は考慮されていない。そのため 2. で述べた辞書を用いた文解析実験により辞書データベースの評価を行うことが可能となる。この評価結果は、辞書構成の手法にフィードバックされる。2. において、低頻度単語と高頻度単語の処理を別けて考えるようにしたのはその例である。

3. 1. 1 システムの概要

システムは、文節分かち書きの形で入力された文の係り受け構造を出力する。係り受け解析のさい、基本となるのは、2 文節間の係り受け関係 [5] であり、辞書データベースを参照しつつ組み合わせ処理により係り受け構造を構成する。構成される係り受け構造は順位付け評価され出力される。処理の流れを、Fig. 6 に示す。

3. 1. 2 辞書データベースの利用

(A) 語彙解析：形態素テーブル、付属語テーブルを参照し、入力文の各文節を自立語部と付属語部に切り分ける。このとき、自立語付属語に対する知識(帰属クラス等)も検索され、また同型異義語の存在も検出される。

(B) 係り受け可能性の検索：各文節に対して、係り得る文節を、クラス間係り受け辞書の検索によりすべて求める。このとき、単語辞書、付属語辞書中の統計量を基に Bayes の定理を用いて、係り受けの機能を表わす係り受けラベルの推定を行う。また係り受けの良さをローカルに評価し、係り先候補をソートする。単語クラス情報のみでは係り先

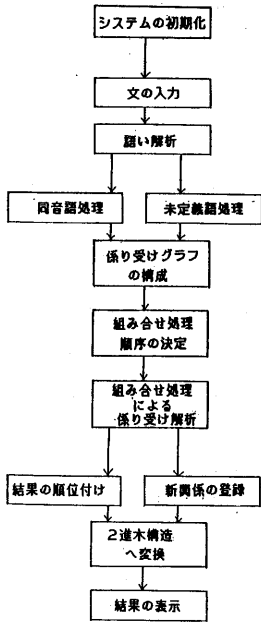


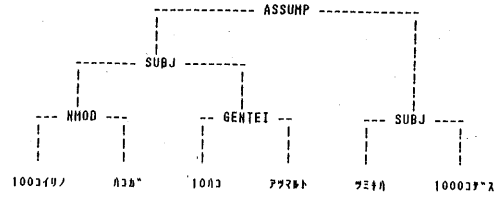
Fig. 6
ESSAY
における処理の流れ

候補文節が求められない場合、階層多重クラスタ情報を用いる。

3. 1. 3 解析結果の順位付け

係り受け可能性をリンクとし、文節をノードとするグラフ構造から解析結果を表わす係り受け構造を組み合わせ処理により抽出する。この組み合わせ処理を行う順序を、係り受けリンクに対するローカルな評価値を基に制御する。この段階で、解析数は評価しうるので、この数が設定した限度数より多い場合は、刈り取り処理を行う。本処理により正しい解析が落ちる率は、限度数を50とした場合4%程度であった。組み合わせ処理により構成された係り受け構造の良さは、次の評価関数により求められ、評価値に従い順位付けを行う。式中、A, B, Cは、0~1の値を取るパラメータであり、その最適近似値は、正しい結果の平均順位、正しい結果が1位に出力されない率を指標とした予備実験によって、(A, B, C) = (0.1, 0.9, 0.9)とした。2進木構造に変換された出力例を、Fig. 7に示す。

*** EVALUATIVE POINTS = 81



*** EVALUATIVE POINTS = 77

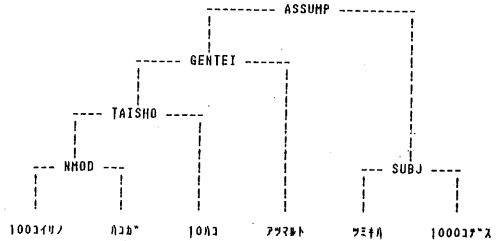


Fig. 7 処理例

$$E = \frac{\sum_{k=2}^n U(k) C(k)}{\sum_{k=2}^n U(k)} \quad (9)$$

$$U(k) = \begin{cases} 1 & \dots k \text{ 文節が受け文節} \\ 0 & \dots \text{それ以外} \end{cases} \quad (10)$$

$$C(k) = \frac{B \left(\sum_{j=1}^{k-1} KU(j,k) (C g(j,k) + (1-C) E_{jk}) \right)}{\sum_{j=1}^{k-1} KU(j,k) + (1-B) L_a(k)} \quad (11)$$

$$KU(j,k) = \begin{cases} 1 & \dots j \text{ 文節} \rightarrow k \text{ 文節} \\ 0 & \dots \text{それ以外} \end{cases} \quad (12)$$

$$g(j,k) = \begin{cases} 1 & \dots \text{品詞レベルでマッチ} \\ 0 & \dots \text{それ以外} \end{cases} \quad (13)$$

$$E_{jk} = A P_{jk} + (1-A) W_{jk} \quad (14)$$

P_{jk} : ラベル推定の尤度

W_{jk} : 辞書中の尤度

E_{jk} : ローカルな評価値

$$L_a(k) = \frac{\text{文中で埋められたスロット数}}{k \text{ 文節中自立語の受けスロット数}} \quad (15)$$

3. 2 文解析実験と評価

初期辞書の構成に用いない文500文に対し実験を行った結果をTable. 1に示す。表中、カバーリンク数とは、処理可能な係り受け単語のペア数である。クラスタリングを行うことによる抽象化の効果が現れているが、得られる解析率は、十分とは言えない。そこで、失敗の場合を考察したところ、係り先を求められない文節をふくむ場合、その文節数の

割合は30%程度であり、学習処理の可能性と必要性が確認された。

Table. 1

	単語レベル	クラスタレベル	階層使用
係り			
クラスタ数	1529	839
受け			
クラスタ数	960	1215
係り受け			
リンク数	6162	4202
カバー			
リンク数	6162	106819
品詞分散	0	35.4
文解析率(%)	32.5	46.7	54.2

4. 学習情報の獲得と構造化

学習処理は、学習情報を文解析処理を通して獲得するフェーズと、獲得した学習情報を既存の構造中に構造化するフェーズの2フェーズにより行われる。学習情報としては、現在の辞書中に存在しない語に関する情報と、現在の辞書記述では処理できない係り受け関係に関する情報の2種類が考えられるが、前者はユーザとの対話により単語情報を得、新たに単語クラスタを割り付けることにより行われる。以下では後者について扱う。これらの処理に対する支援機能はシステムESSAY中に用意されている。

4.1 文解析からの学習情報の獲得

文解析におけるネガティブな状況(解析の失敗、正しい結果が上位に出力されない)は、使用中の辞書の性能が十分でないことを示すものであるが、逆に考えれば、そのような状況を生じさせた原因に関する情報を抽出し、辞書の内容にフィードバックすれば、同じ原因による失敗は以後生じなくすることが可能となる。このような獲得の処理は、文解析をオンラインで行う場合と、オフラインで行う場合とでは現在のところ異なっている。すなわち、オフラインの場合、正しい結果を入力文ファイル中に付加してお

き、出力される解析結果ファイルとの突き合わせ処理により行うが、オンラインの場合、システム側の推定とユーザの対話により行うこととしている。

4.1.1 係り受け構造の推定

入力文が正しい文であり、辞書検索のより得られている構造が正しいとの仮定の下では、推定により係り受け構造が決定できる場合が存在する。簡単な例をFig. 8に示す。この処理は、係り受け構造の性質に基づいてなされるが、一般にはこの場合のように一意に構造を定めることは不可能である。そこでFig. 9に示すようなアルゴリズムに基づきシステムの推定とユーザの応答により正しい構造を決定し、学習情報の獲得を行う。

==辞書検索によるリンク
--推定によるリンク



Fig. 8

文節は文頭側から文末側へ昇順に、1, ..., nと番号付けされているとする。図中で、未定文節とは辞書検索により係り先が見出されなかった文節をいい、限界文節とは、未定文節がこれを越えては係り得ない文節をいう。限界文節及び係り先候補は、非交差条件を考慮して求められる。また、 $EVALC(i, j)$ とは、文節*i*~文節*j*までの部分構造に対する評価関数であり、現在のところ次のように定義している。

$$EVALC(i, j) = \sum_{k=i}^{j-1} |d(k) - k - f(k)| \quad (16)$$

ここで、 $f(k)$ とは文節*k*の付属語が係り文節に存在した場合の、係り先までの隔たりの平均値であり、付属語テーブルに記述されている。また、 $d(k)$ は文節*k*の係り先文節の番号を表わす。

```

for i := n - 1 to 1
  if 文節 i は未定文節 then
    begin
      LM := 限界文節;
      if LM - i = 1 then
        (i -> LM) を登録
      else begin
        NC := 係り先候補数;
        C_LIST := 係り先候補リスト;
        if NC = 1 then
          (i -> NC) を登録
        else begin
          for j := 1 to NC;
            VALUE(j) := EVAL(i, C_LIST(j));
            C_LIST を VALUE に従いソート;
            for k := 1 to NC;
              if 同い合せ(i -> C_LIST(k)) = OK
                then begin
                  (i -> C_LIST(k)) を登録;
                  break;
                end
            end
          end
        end
      end
    end
  end
end

```

Fig. 9 推定アルゴリズム

4. 1. 2 係り受け尤度の適応化

本値は(14)式中のWであり、解析結果の順位付けに際して基本となる量である。正しい結果が1位に出力されなかった場合、本値に対し、適応化処理を行う。

4. 2 学習情報の構造化

4. 1に述べた処理により獲得される新規の係り受けに関する情報は、単にクラスタ間の係り受け関係として既存の構造に新しく付加するよりも、構造化したほうが合理的であり、獲得した学習情報を構造化して保持することも可能となる。

本アプローチにおいては、単語はクラスタ化されるので、構造化処理は単語クラスタ形状を再構成する処理であると言える。しかしながら、現在の初期辞書の構成手法は2. に述べたように、かなり複雑な処理によっているので、文解析された文をすべて原テキストデータに付加し、すべての処理を再度行うことは得策とは言えない。

4. 2. 1 単語クラスタ形状の再構成

そこで、1次的な知識構造である単語クラスタ形状の再構成を、Fig. 10に示す処理により行うこととした。

ここで、特徴的なことは一時的に上位

レベルの情報(一時的階層クラスタとよぶ)を用いることである。

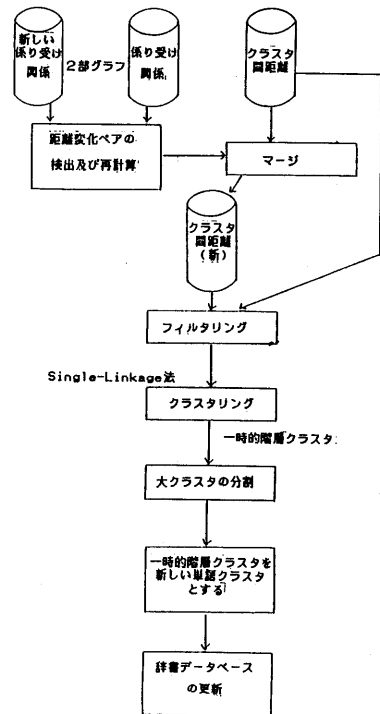


Fig. 10 単語クラスタ形状の再構成

(A) 距離変化ペアの検出と再計算: クラスタ i, j 間の距離を次式で定義する.

$$D_{ij} = 1 - \frac{\sum_r (w(i+r) + w(j+r))}{\sum_p w(i+p) + \sum_q w(i+q)}$$

(17)

$w(*)$ はリンクの頻度

このように、距離はクラスタ間係り受けの2部グラフ上で定義されるため、新しいリンクの付加により距離の変化するペアは容易に検出できる。よって、これらのペアについてのみ、距離の再計算を行えばよい。

(B) フィルタリング: クラスタリングに対して入力となるのは、クラスタ間の距離行列である(新, 旧をそれぞれ D, D' とする)。このとき、新旧の状況の差に関する情報を用いるため次の変化率行列 R を定義し、

$$R = (R_{ij}) = \left(\frac{D'_{ij}}{D_{ij}} \right) \quad (18)$$

この R_{ij} が閾値条件を満たさない場合は、クラスタリングの対象としないという前処理を行った。

(C) クラスタリング：用いたアルゴリズムは、Single-Linkage 法〔6〕である。(B)の前処理により、本アルゴリズム特有の chaining-effect はある程度さけることができた。

(D) 大クラスタの分割：(C)により構成される一時的階層クラスタには、サイズの大きいものも含まれるため、要素のサイズを考慮し、幾つかの核クラスタを選出し、細分割を行った。

4. 2. 2 階層クラスタ形状の再構成

学習情報を考慮した階層多重クラスタ形状の更新は、単語クラスタ間の係り受けリンクに新規のものを付加し、Fig. 3のクラスタリング3の処理により行う。本処理は、簡易に行えることがその特長である。

なお、4. 2に述べた処理は係りクラスタについてのみ行う。

5. 学習実験と評価

学習処理により再構成された辞書データベースの性能と性質の、学習情報の量に対する変化を調べた。性能については、カバーリンク数と文解析率を指標とした。結果を Fig. 11 に示す。図中、構造化Aとはすべての処理を再度行うものであり、構造化Bとは4. 2に述べた処理によるものである。また非構造化とは、単にリンクを付加したものである。構造化による抽象化の効果が表れていることがわかる。性質については、クラスタ数と品詞分散を指標とした。結果を Fig. 12 に示す。構造化Bの場合、学習情報の加え始めにおいて jump が見られるが、双方の量ともほぼ一定の傾向を示す

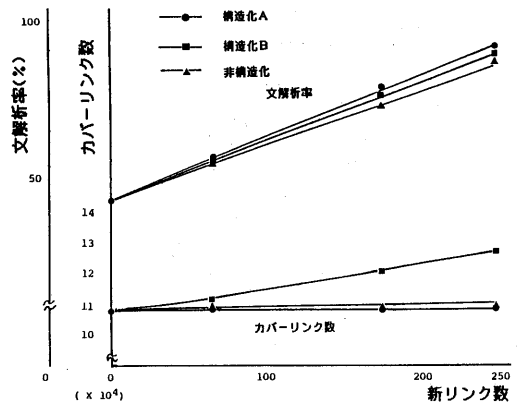


Fig. 11 学習情報の量と辞書の性能

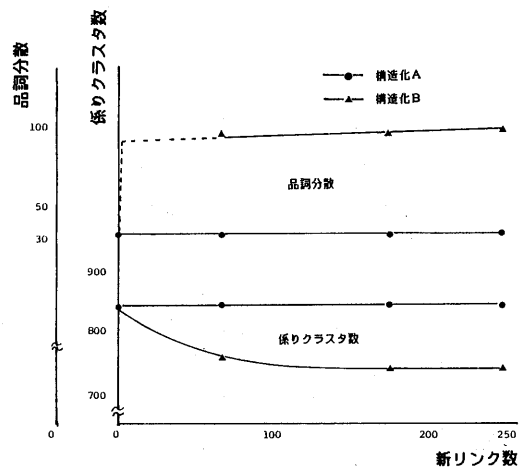


Fig. 12 学習情報の量と辞書の性質

ことがわかり、新しい単語を含まないとする条件下では、新しいリンクを考慮しても大きな変化は生じないことがわかる。

6. まとめ

本アプローチに基づく、知識構造の構成及び成長方法について有効性が確かめられた。今後は、4. 1. 1の考え方を推進し、自律学習の可能性をさぐるとともに階層的な情報をより有効に利用する手法を、解析システムに組み込んでいく予定である。

- 〔参考文献〕
- [1] 辻井・森岡と利用 信学技報 AL-80-27 (1980)
 - [2] Shirai, K., et al.: Japanese Sentence Analysis System ESSAY, COLING 82 (1982)
 - [3][4] 白井, 林, 平田: 情報処理学会第24回全国大会 3K-4, 5 (1982)
 - [5] 吉田: 2文即座の係り受けを基礎とした日本語の構文解析 信学誌D Vol. 55-D No. 4 (1972)
 - [6] Hartigan, J.A.: Clustering Algorithms (1975) John Wiley & Sons