

発見的グラフ探索法を用いた日本語形態素解析

長尾健司 菅野祐司
 松下電器産業(株) 情報通信東京研究所

グラフ探索のアルゴリズム A^* を、探索コストが、最小コスト + α (α は任意の正定数) 以下の範囲で、 m 番目 (m は任意整数) までのパスを取り出すことができるように拡張し、これをベースに、高速な形態素解析を実現した。この形態素解析アルゴリズムの時間計算量は、入力文字列長を L とすると、 $O(L^2)$ となり (α や m に無関係ゆえ) 求める解の範囲を大きくしても、処理の負荷は大きくならないという特徴がある。本稿では、合わせて、本アルゴリズムをもとに、文節数最小法を実現し、最小文節数の解析をただ一つ取り出す場合と、最小文節数 + 1 以内の解析を全て取り出す場合の実験を行い、処理時間がさほど変わらないことを確認し、アルゴリズムの有効性を立証する。

A JAPANESE MORPHOLOGICAL ANALYSIS
 USING HEURISTIC SEARCH ALGORITHM

Kenji NAGAO Yuji KANNO

Tokyo Information and Communications
 Research Laboratory,
 Matsushita Electric Industrial Co., Ltd.
 3-10-1 Higashimita, Tama-Ku,
 Kawasaki, Kanagawa 214, Japan

An efficient algorithm for Japanese morphological analysis has been developed based on a heuristic graph search procedure which is a new extension of the A^* . This extension enables to find m paths of which costs are the least and less than or equal to the minimal cost + α , where m and α are a given integer and an arbitrary positive value, respectively. This algorithm runs in $O(L^2)$ steps for an input Japanese sentence of which length is L . The advantage of this algorithm is its speed since the running time to find m paths is comparable to the time for finding the minimal one. This is expected from the fact that the function, $O(L^2)$, is independent of m or α . The experiments have been carried out on a system in which the costs of graphs are determined by Minimum Pause Group Method.

1. はじめに

筆者らは、種々の日本語処理応用システムの中枢部に組み込めるような汎用の日本語処理基本システムの開発を進めている¹⁾。システムの基本的構成(解析システムのみ)を図1に示す。

モジュール構成は、一文単位に、形態素解析が終了した後で結果を受け取って構文解析を行うというもっともよく見られる処理手順に基づいたものとなっている。特徴は、この間に受け渡される結果のデータ構造と、これを前提とした処理のアルゴリズムにある。

即ち、日本語の様々なレベルでの曖昧性を強く意識して、受渡しのデータである形態素解析結果をグラフにより表現している。さらにこのグラフを出力する形態素解析は発見的なグラフ探索のアルゴリズムを用い、グラフを入力として受け取る構文解析では、構文解析固有の曖昧性と形態素解析結果の曖昧性を統一的に扱える独自のアルゴリズムを開発し³⁾、それぞれ処理の効率化をはかっている。

本稿では、このうちの日本語形態素解析システムの探索アルゴリズムとその処理効率について報告する。

一般に、漢字かな混じりの入力文字列に対する辞書検索の結果は複数あり、通常形態素解析で行なうような接続ルールを用いた絞り込みでは、これを一通りに決めることはできない。これが形態素解析レベルでの曖昧性である。

ここで、接続ルールを絶対的なものとして受け入れてしまうと、形態素解析は、入力文字列と日本語辞書及び接続ルールによって決まる探索グラフの中で、経験的な知識等に照らしてもっともらしい幾つかの解析に対応するサブグラフを探索するという問題として捉えることができる。

この捉え方に素直に基づくならば、この探索にはグラフ探索のアルゴリズムを用いることができる。さらに、ここで形態素解析という問題領域を考えると、探索においては発見的な知識を利用することが可能であり、より効率化をはかることができることに気付く。

そこで、本形態素解析システムでは、発見的なグラフ探索の手法A*をもっともらしい解パスよりなるサブグラフを取り出せるように拡張し、これをベースにしたアルゴリズムを用いる。

2. A*アルゴリズムの拡張

A*は、与えられた探索グラフからコストが最小である任意のパスをただ一つ探索して停止するものであり、現在でも頻繁に利用さ

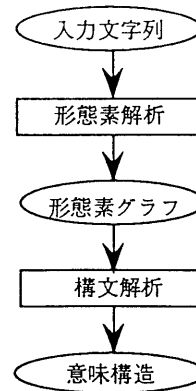


図1 日本語処理基本システム構成

れている。A*アルゴリズムを表1に示すが、これはオリジナルのもの⁴⁾に対して、Albert Martelliが改良を施し⁵⁾、効率アップをはかったものである。

アルゴリズムの時間計算量として、ノードの選択-展開の回数をとるならば、次のような集合Vの要素数をN、即ち、 $N = \dim V$ とすると、オリジナルのA*が、最悪の場合EXP(N)のオーダーになるのに対して、このアルゴリズムは、 N^2 のオーダーに抑えることができる。

(I) $S \in V$

(II) $f(n) \leq f^*(S)$

かつ $\exists n_0 \in V, n_0 \in E(P)$

ならば $n \in V$ ■

以下では、この改良されたアルゴリズムをA*と呼ぶ。

さて、拡張の主目的は、

(1) もっともらしいパスよりなるサブグラフを取り出せるようにする。

であるが、サブグラフを取り出すとなると、当然処理の負荷の増大が予想されるので、

(2) (1)に伴う処理の負荷の増大を緩和する。

も要求される。以下では、上の(1)(2)を実現するため、2段階に分けてA*の変更を試みる。第一の変更は、主として、(1)を確実に実現しようとするものである。このアルゴリズムをA*₊₁と呼ぶ。第2の変更は、A*₊₁を処理効率改善のために、A*₊₁をさらに改良するものである。このアルゴリズムをA*₊₂と呼ぶ。

2-1. 第一の変更: A*₊₁の設計

上述のように、ここでは(1)に主眼を置くが、(2)についても考慮する。改良をそれぞれ次のようにリファインする。

表1 A*アルゴリズム

記号の定義	アルゴリズム
S =スタートノード G =ゴールノード (解ノード) $g(n)$ =スタートノード S からノード n までのコスト $h(n)$ =ノード n から G までのコストの予測値 $h^*(n)$ =ノード n から G までのコストの正確な最小値 $f(n)=g(n)+h(n)$ $f^*(n)=g(n)+h^*(n)$ $c(m, n)$ =ノード m, n 間のパスのコスト $CLOSE$ =すでに探索されたノードの集合 $OPEN$ = $CLOSE$ 中のノードの継続ノードで、 $CLOSE$ に属していないノードの集合 (次の段階で探索される候補となっているノードの集合) $\Gamma(n)$ =ノード n の継続のノードの集合 $P(n)$ =ノード n の親のノードの集合	step 0: F を十分小さい値にする。 step 1: スタートノード S を集合 $CLOSE$ に入れ、すべての $\Gamma(S)$ をつくり、これらを $OPEN$ に入れる。 step 2: もし、 $OPEN$ が空ならば、失敗として外に出る。 step 3: 集合 $OPEN$ の中で最も $f(x)$ の小さいノード n_{rmin} を見つける。もし、 $f(n_{rmin}) < F$ であれば、 $f(n) < F$ であるような n の中で、 $g(n)$ が最小であるようなノード n を候補として選択する。そうでなければ、 n_{rmin} を候補とし、 $F=f(n_{rmin})$ とする。 step 4: 候補ノード n がゴールノードであれば、親ノードとの関係を示すポイントによって、スタートノードからゴールノードまでのパスを見出し、探索を終了する。さもないければ、次にいく。 step 5: ノード n が展開可能であれば、すべての継続ノード $x \in \Gamma(n)$ を生成し、 $f(x)_{new}$ を計算する。ノード n を集合 $CLOSE$ に移す。 step 6: $CLOSE$ にも $OPEN$ にも入っていない継続ノード、及び $CLOSE$ あるいは $OPEN$ に入っているが以前の評価関数の値より新しく計算した $f(x)_{new}$ の方が小さいような継続ノードについては、親ノードを n とするポイントをつけ、評価関数の値を $f(x)=f(x)_{new}$ として $OPEN$ に移す。それ以外の継続ノードについては棄却する。 step 7: step 2 に行く。

① 最小コスト + α 以下の範囲内で、最小のパスから m 番目までのパスを取り出せるようにする。

② 探索の過程に於て、発見的関数が徐々に高精度になっていく場合でもサブグラフを完全に取り出せるようにする。

形態素解析の場合は、文節数最小法等の非常に有効なヒューリスティックがあるため、 α は最小コストに比べて小さくてよいが、形態素解析応用システムのことを考慮すると、 m 番目までの解という指定の仕方も可能である方がよいと考える。

以下にそれぞれの実現について説明する。

【①の実現】

既に述べたように、 A^* はコストが最小であるパスをひとつだけ探索するアルゴリズムである。この点は、step 4 と step 6 で実現されている。step 4 では、選択された候補ノードがゴールノードであれば、探索を終了する。また、step 6 では、親ノードの展開によって生成されたノードが既に他の親ノードから生成されているノードである場合は、常に S からのコストが最小であるパスのみを残すようにしている。即ち S から

のコストが最小でないパス (パス上のノード) は全て棄却するようにしている。

すると、①を実現するためには、以下の2点の変更を行えばよい。

(i) step 4 で、選択されたノードがゴールノードであっても停止せず、探索を続ける。

(ii) step 6 で、新たに生成されたノードが既に生成され $OPEN$ 或は $CLOSE$ 内に存在しても、それが最終的に求めるパス上のノードである可能性があるならば繰り返して有効とする ($OPEN$ に入れる)。

(ii) のためには、同一のノードであっても、それを生成した親ノードや、 S からのコスト (g に相当する) が異なればあたかも別々のノードのように扱われなければならない。この結果、実質的に、探索はノードを f や g の値で評価して、選択 (探索) するのではなく、ノードと親ノードと g の値の3つ組を選択することを繰り返すという形を取ることになる。

【②の実現】

①の実現のために (i) を実現すると、ゴールノードが選択される度に、親ノードを逆に

表2 A*-1

記号の定義	拡張したアルゴリズム
<p> $n_p = \langle \text{node, ancestor, cost} \rangle$ = 拡張ノード = A*のノードに変わって 探索の単位となるもの $n_p.\text{node} = \text{ノード}$ $n_p.\text{ancestor} = n_p.\text{node}$の 祖先ノード $n_p.\text{cost} = S$から$n_p.\text{node}$までのコスト $S_p = \langle S, \text{NIL}, 0 \rangle =$ 拡張スタートノード $G_p =$ 拡張ゴールノード $g_p(n_p) = n_p.\text{cost}$ $h_p(n_p) = h(n_p.\text{node})$ h_{sup}, C_{min} $h_p^*(n_p) = h(n_p.\text{node})$ $f_p(n_p) = g_p(n_p) + h_p(n_p)$ $f_p^*(n_p) = g_p(n_p) + h_p^*(n_p)$ $c_p(n_1, n_2) = c(n_1.\text{node}, n_2.\text{node})$ CLOSEP=すでに探索された 拡張ノードの集合 OPENP=CLOSEP中の拡張ノードの 継続拡張ノードで、 CLOSEに属していない拡張ノードの 集合(次の段階で探索される候補と なっている拡張ノードの集合) $\Gamma_p(n_p) =$ 拡張ノードn_pの 継続拡張ノードの集合 $P_p(n_p) =$ 拡張ノードn_pの 親拡張ノードの集合 $\alpha = 1 - \beta$ が指定するパラメータで、求める ノードのコストの範囲を定める $m = 1 - \beta$ が指定するパラメータで、最小コスト からm番目までの解ノードを全て求める </p>	<p> step 0: F_pを十分小さい値にkを0に C_{min}を未設定とする。 step 1: S_pを集合CLOSEPに入れ、すべての $\Gamma_p(S_p)$をつくり、これらを集合OPENPに入れる。 step 2: もしOPENPが空ならば、kが0のままなら 失敗として外に出る、そうでない場合は探索を終了 する。OPENPが空でないならstep3に行く。 step 3: 集合OPENPの中で最も$f_p(x)$の小さい拡張ノ ード$n_{p, n}$を見つける。もし、$f(n_{p, n}) < F_p$であれば、 $f_p(n_p) < F_p$であるようなnの中で、$g_p(n_p)$が最小であ るような拡張ノードn_pを候補として選択する。そうで なければ$n_{p, min}$を候補とし、$F_p = f_p(n_{p, min})$とする。 step 4: 候補拡張ノードn_pがG_pであれば、 $n_{p1}.\text{ancestor} = n_{p2}.\text{node}$ $n_{p1}.\text{cost}$ $= n_{p2}.\text{cost} + c_p(n_{p1}.\text{node}, n_{p2}.\text{node})$の関係に よって、S_pからG_pまでのノードを見出し、ノード上 のノードにマーカーを$k = k + 1$とする。 $k = m$の場合は探索を終了する。 C_{min}が未設定のままなら$C_{min} = f_p(n_{p, min})$とする。 $k < m$の場合及びn_pがG_pでない場合はstep5に行く。 step 5: 拡張ノードn_pが展開可能ならば、すべての継続 拡張ノード$n_{p, neu} \in \Gamma_p(n_p)$を生成し、これらについて、 $g_p(n_{p, neu})$を計算する。C_{min}設定済みの場合は、$f_p(n_{p, neu}) \leq C_{min} + \alpha$をみだす$n_{p, neu}$に対してのみstep6を 行なう。 step 6: OPENP或はCLOSEPの拡張ノードの中に、node $= n_{p, neu}.\text{node}$であるようなものが既に存在する場合、 それらのg_pの中で($g_p(n_{p, neu})$を含めて)、最小の値を $g_{p, min}$とすると、$g_p(n_{p, neu}) \leq g_{p, min} + \alpha$となる場合に 限り、同一の値も数えて、m番目までに$g_p(n_{p, neu})$が 入っている場合は、$n_{p, neu}$をopenpに入れ、この結果 m番目からあふれた拡張ノードは、棄却する。そうで ない場合は、$n_{p, neu}$を棄却する。 これ以外の場合は、$n_{p, neu}$をOPENPに入れる。 n_pをCLOSEPに移す。 step 7: step 2に行く。 </p>

通りパスを見いだすことになるが、この際、
 次の計算によって解のパス上の各ノード n から
 ゴールまでの正確な最小コスト h^* がわかる。
 $h^*(n) = g(G) - g(n)$
 但しこれは、同一のノードについては、最初
 に計算された値のみについて言えることであ
 る。従って、探索が進むに連れノードからゴ
 ールまでの正確な最小コストがわかっていく
 わけで、以降の探索にはこの正確な値を利用
 するのがよいのではないかと考えられる。
 これは、発見的関数が真の値に近付いてい

くことが、無駄な探索の減少につながることを
 期待するものであり、発見的関数が(スタ
 ティックに)高精度であるほど、無駄な探索
 が少ないというA*の性質からの類推である。
 このような探索の途中段階で得られたゴール
 までの正確なコストを利用することも含め
 て、次の条件を定める。
【条件T】
 時刻 t に於けるノード n の発見的関数値を
 $h(n, t)$ とすると、 h は次の条件を充さ
 なければならない。

$$h(n, t_1) \leq h(n, t_2) \leq h^*(n) \quad (t_1 \leq t_2) \blacksquare$$

条件Tの右辺の不等式により、アルゴリズムは、求めるサブグラフを完全に取り出すことができる。このことは、A*に於ける、実行可能性と同様にして証明することができる。

以上のような変更により得られるA*₊₁をその記述の中で用いる記号の定義と共に表2に示す。A*の記述で用いた記号は、その定義に従う。

【A*₊₁の時間計算量】

A*₊₁の時間計算量について考察する。

(定義)

OPENPからCLOSEPに移された拡張ノードのf_p値の列を順に、

$$S_0: f_1, f_2, \dots, f_r$$

S₀の部分列S₁を

S₁: f_{i1}, f_{i2}, ..., f_{ip} を f_{i1} = f₁ とし、f_{ik+1} (1 ≤ k ≤ i_{p-1}) を f_{ik} のノードがCLOSEPに移されて後、最初に、f_{ik+1} ≥ f_{ik} を満足するものとする。■

(系1)

発見的関数を任意の拡張ノードn_pについて、時間的変化無しとした場合、即ち、h_p(n_p, t) = h_{p'}(n_p)とした場合、A*₊₁の時間計算量は、m² · N²以下である。

(Lemma 1)

$$\dim S_1 \leq m \cdot N$$

ここで、NはA*の説明のところで定義済み、mはユーザの指定した、求める解の数である。

(証明) f_{ik}の列は非減少列ゆえ、f_p(n_p) ∈ S₁ であるようなn_pで、

$$n_{p1}, node = n_{p2}, node$$

$$n_{p1}, cost > n_{p2}, cost$$

を充す、n_{p1}, n_{p2}がこの順に繰り返して、OPENPからCLOSEPに移されることはない。

従ってf_p(n_p) ∈ S₁である拡張ノードn_pでn_p, nodeが同一のものについては、n_p, costの小さい順に、高々m個選択される。よって、dim S₁ ≤ m · N ■

(Lemma 2)

S₁にf_{ik}が加えられてからf_{ik+1}が加えられるまでの間に、CLOSEPに移される拡張ノードのf_pの値を順に

$$f_{ik} = f_j, f_{j+1}, f_{j+2}, \dots,$$

$$f_{j+1} = f_{ik+1} \text{ とすると、}$$

$$l \leq m \cdot N$$

である。

(証明) f_{j+r} (0 ≤ r < l) に対応する

拡張ノードn_{p1}がCLOSESPに移されたとき、f_p < f_{ik} を満足するf_p値を持つOPENP中の拡張ノードn_pのg_pの中でg_p(n_{p1})が最も小さかったわけで、以後、j+1回のステップに至るまでに、

n_{p2}, node = n_{p1}, node なる拡張ノードn_{p2}がCLOSEPに移されても必ず、g_p(n_{p1}) ≤ g_p(n_{p2}) が成り立つ。

従って、S₁にf_{ik}が加えられてからf_{ik+1}が加えられるまでの間に、n_p, nodeが同一のものについては、n_p, costの小さい順に、高々m個選択される。

よって、l ≤ m · N ■

(系1の証明)

lemma 1 及び lemma 2 より、時間計算量 = dim S₁ · l ≤ m² · N² ■

2-2. 第2の変更: A*₊₂の設計

A*₊₁の設計で行われた2つの変更には、それぞれ次のような問題点がある。

【問題点: Q₁】

A*₊₁の様に3つ組を探索の単位としてしまうと、同一のノードなのに、親ノードやgの値が異なるがために、別々の探索ユニットとして扱われ、以降生成される子孫ノードの系列は、全く同じであるにもかかわらず、(探索ユニットが異なるため) 繰り返して生成されてしまうということが起こりうる。

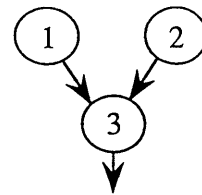


図2 ノードの選択

【問題点: Q₂】

hが条件Tを充しながら正確な値に近づいたとしても必ずしも、選択(探索)回数が減少するとは限らない(実際このことは後の実験により確認される)。むしろ、逆の場合も起こり得る。これは、同じノードを選択する場合でも、選択する順序が異なれば、全体の選択回数が異なってくることもあるからである。例えば、図2に於て、ノード①②③がこの順に大きいf値を持っていて、①→②→③の順に選択されるとすれば、それぞれ一回の選択、合計3回で済むが、①が選択された後で、②のhが条件Tを充す範囲で変化してf値が③のf値より大きくなったとすると、①→③→②→③と合計4回になってしまう。

表3 A**+2 (A**+1の変更点)

step4: 拡張ノード n_0 が G_0 であれば、CLOSEPの中を S_0 の方に廻り、 n_{01} . ancestor= n_{02} . nodeなる親拡張ノードとの関係を基に、 $g_0(n_{02})+c_0(n_{01}, n_{02})+h_0*(n_{01}) \leq C_{min}+\alpha$ を充す拡張ノード全て辿ってマクする。この際、 h^* が、 n_{01} について設定済みで、かつ n_{02} について未設定の場合に、 $g_0(n_{02})+c_0(n_{01}, n_{02})+h_0*(n_{01})=f_0*(n_{01})$ の関係がある場合のみ、 $h_0*(n_{02})=f_0*(n_{01})-g_0(n_{02})$ を計算して保持しておく。そうでない場合は、 $h_{0sup}(n_{02})=f_0*(n_{01})-g_0(n_{02})$ を計算以前の値より小さい場合はこの値を保持しておく。

step6: OPENP或はCLOSEPの中に、 n_{0new} . node= n_{0old} . nodeなる拡張ノード n_{0old} が既に存在する場合、これらの g_0 の中で最小の値を g_{0min} 、その拡張ノードを n_{0min} とする。ここで、以降の記述の簡単のために、つぎの3つの命題をさだめる。

(a) $g_0(n_{0new}) \leq (g_{0min} + \alpha)$
 (b) C_{min} 設定済みの場合は、 $f_0(n_{0new}) \leq C_{min} + \alpha$
 (c) 同一ノードを持つものの中で同一の値も数えて、 $g_0(n_{0new})$ が最小のものから m 番目までに入っている。

n_{0min} が CLOSEP の中に在って、既にマクされている場合は step(6-1) を実行する。マクされていない場合は step(6-2) を実行する。そのような n_{0old} が存在しない場合は、step(6-3) を実行する。step(6-1): $h_0*(n_{0old})$ が既にわかっている場合 step(6-1-1) を実行する。わかっていない場合は step(6-1-2) を実行する。

step(6-1-1): $g_0(n_{0new})+h_0*(n_{0old}) \leq m_0+\alpha$ かつ、(c) を充す場合は、 n_{0new} を起点として、 S_p に至るまで、親拡張ノードを廻り、解パスを見いだす。これは step4 で G_0 を起点として廻る場合と全く同じである。上式を充さない場合は、 n_{0new} を棄却する。

step(6-1-2): $g_0(n_{0new}) < g_{0min}$ の場合、step(6-1-2-1) を実行する。そうでない場合は step(6-1-2-2) を実行する。

step(6-1-2-1): 解パスを見いだすと同時にパス上の各拡張ノードの h_{0sup} を計算、保持し、 n_{0new} を OPENP に入れる。 n_{0min} が OPENP 中に存在していた場合は、CLOSEP に移す。この結果、(a) かつ (b) か (c) から漏れた拡張ノードは棄却する。

step(6-1-2-2): $g_0(n_{0new})+h_{0sup}(n_{0old}) \leq m_0+\alpha$ かつ (c) の場合、step(6-1-2-2-1) を実行する。そうでない場合は、step(6-1-2-2-2) を実行する。

step(6-1-2-2-1): 解パスを見いだすと同時にパス上の各拡張ノードの h_{0sup} を計算、保持し、 n_{0new} の CLOSEP に入れる。

step(6-1-2-2-2): n_{0new} を CLOSEP に入れる。

step(6-2): (a) かつ (b) かつ (c) の場合、特に、 $g_0(n_{0new}) < g_{0min}$ の場合は、 n_{0new} を OPENP に入れ、 n_{0min} が OPENP に存在していた場合は n_{0min} を CLOSEP に移す。 $g_0(n_{0new}) \geq g_{0min}$ の場合は、 n_{0new} を即座に CLOSEP に移す。この結果、(a) かつ (b) かつ (c) の条件から漏れた拡張ノードは棄却する。 n_{0new} が (a) かつ (b) かつ (c) を充さない場合は、棄却する。

step(6-3): n_{0new} を OPENP に入れる。

【解決策】

Q1, Q2 のようなことがおこるのは、探索の単位を <ノード、親ノード、g> として、ノード、親ノード、g のどのひとつでも、異なれば全く違うものとしていることが原因であり、結局は Q1 と Q2 の本質は同じである。

そこで、A**+1 に対して、次のような改良を加える。

- (1) 同一のノードを有する拡張ノードが既に生成されている場合は、たとえこれが最終的に解パスに含まれる可能性のある場合でも、必ずしも、新たに探索候補とはしない(拡張ノードを OPENP に入れない)。最も g の小さい

拡張ノードのみを成長させるようにする。そして、解パスを取り出すときには、はじめて g や親ノードの違いを考慮して漏れなく解パス上のノードを拾い出すようにする。

- (2) 既に解パス上のノードだとわかっている拡張ノードは、生成された時点で解パスに含まれることがわかる場合があるので、その場合その子孫の成長を止めて、解パスの取り出しにかかる。

変更は、A**+1 の step 4 と step 6 に対してだけである。こうして得られる A**+2 を表3に示す。

【A+2の時間計算量】**

A_{*+2} の時間計算量について考察する。

【系2】

発見的関数を任意の拡張ノード n_p について、時間的変化無しとした場合、即ち、 $h_p(n_p, t) = h_{p'}(n_p)$ とした場合、 A_{*+2} の時間計算量は、 N^2 以下である。■

証明は省略するが、系1と類似の展開によって示される。

系2によれば、 A_{*+2} では求める解の範囲を大きくしても、処理の負荷はさほど大きくならないことが期待できる。

3. 形態素解析アルゴリズム

以下に、 A_{*-1} や A_{*+2} をベースにした形態素解析アルゴリズムを示すが、データ構造の対応と、発見的関数として如何なる情報を用いるかということ、それから時間計算量について示せば十分だろう。

【データ構造の対応】

表4に示す。

表4 データ構造の対応

ノード => (文中左位置、文中右位置、品詞*インク)
文中左位置=入力文字列中の位置
文中右位置=入力文字列中の位置
品詞*インク=文中左位置と文中右位置の間に挟まれた文字列に対応する単語で、品詞を同じくするものの情報を格納するデータ構造への*インク
スタートノード => (0, 0, NIL)
ゴールノード => ノードの文中右位置が入力文字列長と一致するもの。
パスのコスト => パス上の単語列の中の文節の数や、単語の頻度情報等をもとにする。
継続ノード => ノードの単語に接続しうる単語よりなる。

【発見的関数】

既に述べた条件Tを充すようにすることを念頭において一例をあげる。

ヒューリスティックは文節数最小法を用いる。コストはパス上の文節の数とする。初期値の設定には、入力文字列を末尾から先頭に向かってスキャンし、各文字位置で、その文字からはじまる自立語や付属語列の最大長についての予め用意した情報を基に、末尾までの文節数の予測値を計算すればよい。解パスが見いだされ、パス上のノードからゴールまでの正しい最小コストがわかった後は、その値を用いる。

【形態素解析アルゴリズムの時間計算量】

入力文字列の各位置から始まる単語の品詞

の数は有限個、K以下とすると、入力文字列長をLとすれば、 $N \leq K \cdot L$ であるから、発見的関数を一定とすると、系1より A_{*-1} に基づく場合は、 $K^2 \cdot m^2 \cdot L^2$ 以下、系2より A_{*+2} に基づく場合は、 $K^2 \cdot L^2$ 以下となる。

4. 評価実験

A_{*-1} 、 A_{*+2} をベースにした形態素解析アルゴリズムの性能評価を行うため、比較実験を行う。文節数が最小の解析、最小+1以内の解析をすべて取り出すという実験で、実際に解析に要した処理時間、ノードの選択回数を測定する。

【解析対象文】

新聞記事(漢字仮名混じり) 50文

【実験機種・言語】

SUN3-60上のC言語

【比較アルゴリズム】

表5に示す。

表5 比較アルゴリズム

X	A_{*-1} で全てのノードについて $h=0$ とする
Y	A_{*-1} で h を初期値のまま一定とする
Z	A_{*-1} で h^* が求まった後は $h=h^*$ とする
V	A_{*+2} をベースにする
W	唯一つだけ解パスを見い出して停止するように $m=1$ とする。

Xから、発見的探索を用いなくて通常のダイナミック・プログラムにより実現した場合を類推することができる。また、Wは、 A_{*} アルゴリズムをもとに最小コストのパスをただ一つ見いだして終了する場合に等しい。

【発見的関数の初期値】

4. で例として示したものとす。

【実験結果・考察】

実験結果を図3、図4にグラフで示す。

図3、図4は、それぞれ、入力文字列長に対する処理時間 t (1/60 SEC)、close (OPENPからCLOSEPに移すこと)の回数 r 、の常用対数値を示したものである。但し図中の数字はそのまの値を示している。

グラフよりV (A_{*+2} に基づく)が圧倒的に優位であることは明白である。文長Lが大きくなるほど、この傾向は大きくなり、最小文節数+1以下の解析を全て求める場合、 $L=80$ (解パスの数=315)では、処理

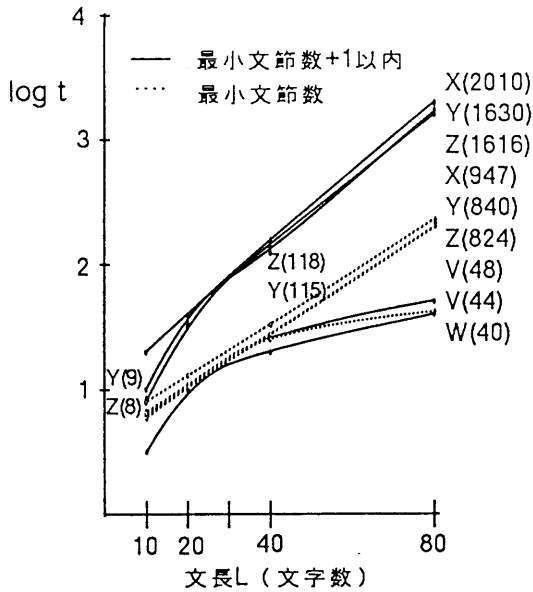


図3 処理時間(t)-文長(L)

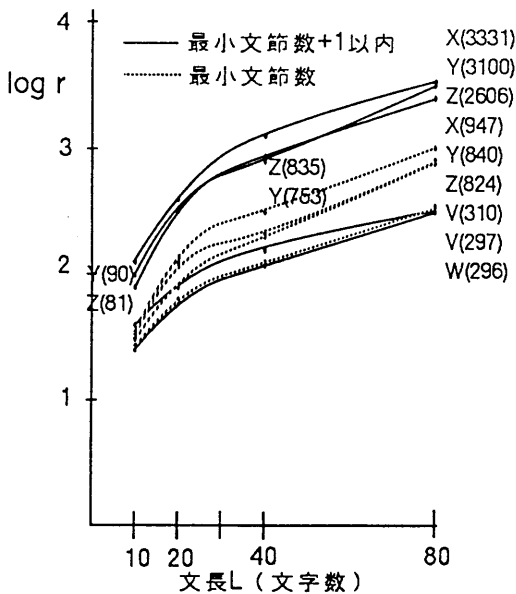


図4 クローズ回数(r)-文長(L)

時間においてVはXの約50倍、Zの40倍、closeの回数でもXの10倍、Zの8倍の効率となっている。又、Vは最小文節数の解析、最小文節数+1以内の解析を取り出す場合で、処理時間、及びcloseの回数において差異がない。しかも、VはWに比べても、処理時間、クローズの回数ともさほど劣ることはない。即ちVは、最小コストの解パ

スをただ一つ求めて停止する場合とほとんど同じ処理時間で終了する。

これは、かなり驚くべき結果であるが、アルゴリズムを比較すると、納得できる部分も多い。即ち、X、Y、Zでは、親ノードから子ノードが生成されるとき、子ノードが既に他の親或いは、異なったg値をもって既に生成されている場合でも、繰り返して子ノードをOPENPに入れることがあるが、Vでは、新しい子ノードのg値が古い子ノードの最小g値よりも小さくないかぎり、子ノードを再びOPENPに入れることがない。従って、求めるコストの範囲(α やm)が大きければ大きいほど、この効果は大きくなる。

また、 $L=40$ のときにYがZを上回る効率を示したのは、既に述べたように、発見的関数が条件Tを充しながら変化するとき、ノードの選択回数が増えてしまうことも有り得るからである。

5. おわりに

1. グラフ探索のアルゴリズムA*を探索コストが、(最小コスト+ α)以下の範囲内で、m番目までの解パスを出力できるように拡張し、これをベースに高速な形態素解析アルゴリズムを設計した。

2. このアルゴリズムの時間計算量は、入力文長をLとすると $O(L^2)$ となり、 α やmに無関係であることを示した。

3. さらに、本アルゴリズムによれば求める解の範囲を大きくしても、処理時間がほとんど増加しないことを実験で確認し、2.を立証すると共に、アルゴリズムの有効性を示した。

参考文献

- 1) 長尾：日本語処理基本システム(1) 情報処理学会第37回全国大会
- 2) 長尾：文節数最小法・形態素解析の実現について 情報処理学会第35回全国大会
- 3) 菅野：曖昧さの効率的処理のための構文解析手法について 本研究会
- 4) PETER E. HART, et al.: A Formal Basis for the Heuristic Determination of Minimum Cost Paths, IEEE TRANSACTIONS ON SYSTEMS SCIENCE AND CYBERNETICS., 7, 1968, pp100-107
- 5) Alberto Martelli: On the Complexity of Admissible Search Algorithms, Artif. Intell., 8, 1977, pp1-13