

自然言語処理と知識表現

辻井 潤一
(京大・工学部)

§ 1. はじめに

自然言語の解析は、与えられた單語列としての文を受け取り、その文の持つ「構造」を明示的に表現する何らかの表現を作り出す過程である。人工知能における自然言語理解の研究においては、この自然言語の解析において、現実世界に関する人間の持っている様々な知識が関与していることを明らかにしてきた。

しかしながら、自然言語表現そのものについての知識、すなわち、文法的な知識も、当然のことながら、自然言語理解には大きな役割を果している。ここでは、文法的知識とは、必ずしも、名詞・動詞・名詞句といった品詞情報にもとづく統語的な知識だけではなく、ある言語表現がどのようなメッセージを聞き手に伝達しようとしているかに関する総ての知識を指すことにする。事実として同じ内容を表現する文であっても、文の焦点・話題等が異なる場合がある。これまでの人工知能からの自然言語理解の研究は、あまりにも「外界に関する知識」に焦点をあてすぎていたために作成されたシステムは非常に domain specific なものになりすぎていた。

本報告では、「外界に関する知識」は直接対象とはせず、一般的言語現象に即した知識をどのように表現するか、また、言語表現と外界知識がどのように関連を持っているかについて整理し、これを計算機処理するためにはどのような機能が必要かについて考えることにする⁽¹⁾。

§ 2. 言語現象に関する知識

我々が「文法」という言葉から想起するのは、「名詞」・「動詞」・「名詞句」といった「品詞」レベルでの規則である。ここでは、まず、議論に入る前に、自然言語に関する知識として、このような品詞レベルでの規則性以外にどのようなものがあるかをいくつか考えてみよう。

(a) 形態素的規則性：單語の諸形変化は、單語内での現象であり、例えば、言語学における文法理論では、品詞レベルでの規則性を抱えた「句構造規則部門」・「变形規則部門」とは別に取扱われる。しかしながら、文の解析し、次元の單語列としての文から、その構造を抽出する解析の過程においては、この形態素的規則性も品詞レベルの規則性と一緒に取扱わなければならぬ。例えば、

(1) the gentlemen in the room {
 (1-1) which
 (1-2) who
 (1-3) that were

では、(1-3)の關係代名詞節の先行詞が room ではなく the gentlemen であることは、この形態素的規則性から決定できる。

(b) 簡単な意味標識：(1)の例文で、gentlemen が單数の gentleman であっても、(1-1)の關係節の先行詞は、the room, (1-2)の關係節の先行詞は the gentleman と決定することができます。(a)と(b)の規則を同時にとらえようとするとき、すく

なくとも、名詞句 the gentleman や the room には、「数」の情報と HUMAN の標識が必要となる。

(c) 統語構造から意味構造への写像：能動と受動、使役等の構文構造からわかるように、表層上の言語表現は、さまざまな変形をうけることになる。言語学では変形部門を設けて、一般的な規則で、この表層表現と深層上の意味を関連づけていた。しかしながら、最近の傾向として、変形規則の取扱う範囲をできるだけ少なくし、意味解釈規則でこれを行なうようになつてている。意味解釈規則は、個々の単語個別の規則であるので、辞書中に、多くの規則性が書かれることになる。このことは、今まで、比較的少數の規則によって一般化できることと考えられていてものが、各単語個別の規則として定式化されることであり、Frame 等の A.I. 研究の動向と一致している。

'Finally, I assume that it is easier for us to look something up than it is to compute it. It does in fact appear that our lexical capacity -- the long-term capability to remember lexical information -- is very large.'⁽²⁾

日本語についても、Case Linking 規則⁽³⁾と呼ばれる単語個別の規則で、これを行なおうとする試みがみられる。

計算機処理の観点からすると、言語処理のための規則の多くの部分をこれまでデータと考えられてきた辞書の記述中に埋め込むことになり、ここでも、人工知能の分野で開発されてきた手法が有効になる。例えば、

(2) 太郎は、花子を美しいと思う。

という文では、「思う」は、「NP が S を思う」という格パターン ところが、「花子」のように本来は、S 中で格助詞「が」でマークされる要素が、格助詞「を」を伴って、S から外へ出されることがある。格助詞「を」のこのような用法は、一連の動詞に個別のものであり、こういった動詞の辞書中に、[NP が NP を S を思う] という格パターンとともに、「NP を」が S の「が」格に入るという写像関係を記述しておく必要がある。しかも、この現象は、S 中の述部のアスペクト (d) の項参照) が「状態的」な場合に限られており、このような制御条件も「思う」の辞書記述中に書いておく必要がある。

(d) アスペクト：日本語の述部には、「へている」・「へてある」・「へはじめる」・「～つつある」といったアスペクト形式素がつく。これらのアスペクト形式素に関してすくなくとも、次のような規則性が記述される必要がある。

① 動詞の格パターンの選択

(3) 利用者が パラメータを 指定する

(4) 利用者によって パラメータが 指定してある

②主述部のアスペクト素性(瞬間, 繼続 etc)とアスペクト形式素によって決まる述部アスペクト

(5) 彼は死んでいる。 (結果状態)

(6) 彼は、本を読んでいる。 (進行, 経験, 習慣, 結果状態)

①を取扱うにめには、アスペクト形式素の辞書記述によって、動詞の辞書記述(これ自身がまた処理手順を記述する規則になつてゐる)の一部を修正したり、とり出していなければならぬ。

②の規則化は、主として、主述部のアスペクト素性と形式素間の規則性として扱えられるが、実際には、(6)の例文のように一意に決定できない場合も多い。しかしながら、

(7) 彼は、今本を読んでいる。 (進行)

(8) 彼は、よく本を読んでいる。 (習慣)

(9) 彼は、いつも本を読んでいる。 (習慣)

といった副詞との共起関係によって、その可能性が減少すること、また、

(10) 彼は、いつも本を読んでいる。 (進行)

(11) 彼は、朝日新聞を読んでいる。 (習慣)

(12) 彼は、ジャン・クリストフを読んでいる。 (習慣)

といった名詞の属性によつても、その主な読みが決まる。(10)・(11)・(12)のような名詞句と述部アスペクトの関係は、現実世界の知識との接点として重要である。しかしながら、副詞の共起と述部のアスペクトの決定は、一般的な言語知識として処理できらる必要がある。

述部のアスペクトは、文の間の関係を決定する際にも参照される。

(13) プログラムを実行していきため、----- (原因)

(14) プログラムを実行するため、----- (目的, 原因)

これらアスペクトに関する現象は、日本語において、アスペクトを決定する場合には、単に局所的な単語のまとまりだけではなく、かなり離れた構成要素間の相互関係をうまく扱える必要のあること、また、現実世界の知識をどのように言語処理に反映すべきかの示唆を与えてくれる。

(e) 文脈処理と構文：日本語における照応現象については、これまでに言語学の分野で比較的よく研究されてゐる。一般に、自然言語理解の研究においては、文脈処理は、文の統括的・意味的な処理が完了した後に、文の意味内容を中心に、現実世界の知識を参考にして行われることが多かつた。多くの文脈処理が、現実世界の知識を必要とするることは確かであるが、一方で、文中での語順や統語構

造も重要な役割を果している。例えば、言語学の側から提案されている proceed と command の原則を考えてみよう⁽⁴⁾。この原則は、代名詞は、「その先行詞よりも先行し、同時に、先行詞を含む構成素を統御する」ことにはない、という原則である。

(15) 花子は、彼から太郎がかって住んでいた家を譲り受けた。
X ↑

(16) 花子は、太郎がかって住んでいた家を彼から譲り受けた。
X ↑ OK

(17) 花子は、彼がかって住んでいた家を太郎から譲り受けた。
OK

(15) は、図 1 に示すように、代名詞「彼」が、「太郎」よりも先行し、かつ、これを統御する位置にあるために、「彼」と「太郎」とは同一指示とはなり得ない。これに対して、(16)・(17) は、このいずれかの条件が破られているために、同一指示であり得る。

このような条件だけでは、必ずしも、「彼」と「太郎」が同一指示であると断定することはできないが、少なくとも代名詞の指示対象の範囲を限定する役割は果す。このことは、文脈処理においても、「表層上の諸順」（先行の条件）、及び「統語上の構造」（統御の条件）を参照する必要のあること、したがって、文脈処理のように、外界に関する知識に強く依存するような処理においても、言語構造を cue として使う必要のあることを示唆している。

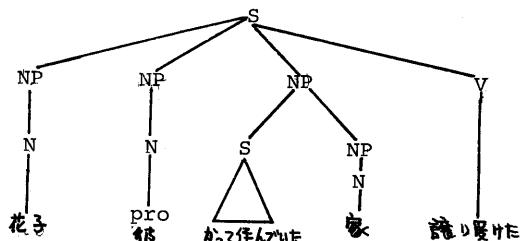


図 1. 「花子は、彼から……」の構造

(f) 意味処理と統語処理：自然言語理解における意味処理の重要性は、例えれば、「の」によって接続される名詞間の意味的関係の決定等を考えれば明らかである。しかしながら、意味処理が重要であるということと、すべての処理が意味によって行なわれるということとは別のことである。重要なことは、統語処理でどの程度のこと方が行なわれるか、また、どのような統語構造があらわれたときに、どのような意味処理を起動すべきかをきめ細く観察する必要がある。簡単な例を示す。

(18) 日立の本社が東京にある。

(19) 日立が東京に本社がある。

この例は、「象は鼻が長い」と並行的な例であるが、可もなくとも、格助詞「が」でマークされる名詞句が、述部と直接係り受け関係をとらず、文中の名詞句と意味的に関係することを認めなければならない。このような現象は、格助詞「が」あるいは副助詞「は」について生じ、しかも、述部のアスペクトが「状態的」である場合に生じる。また、このように宿にういて「が」名詞句が意味的に関係するのは、述部と「が」格でつながる名詞句である。したがって、このような状

況が文中に生じた場合に、名詞間の意味的関係を調べる適切な処理を起動すればよいことになる。

§3. 計算機処理に必要な機能

§2で議論したことは、必ずしも日本語や英語の言語現象のすべて網羅したものではもちろんないし、その定式化にも多くの問題が残されている。実際、これらの現象のための適切な文法モデル（どのような情報を参照すれば良いか etc）を作成することは、それ自体で一つの研究となり得るものであろう。ここで、議論していけるは、このような言語現象を処理するには、少なくとも §2で述べたような各種の情報を参照する必要があるということであり、そのためには、我々はどのようなソフトウェア体系を用意することができるか、ということである。以下、いくつかの観点から、これを考えてみよう。

[保守の容易性] §2で述べたのは、言語にみられる規則性の非常に狭い部分を記述したものである。これと同じ程度の複雑な規則性は、いたるところにみられる。これらをすべて把握した後、はじめて自然言語処理が可能になると考えることはできない（このように考えることは、機械処理を全く放棄することにつながる。）すくなくとも、現在把えられている規則性の範囲内で、計算機処理のためのプログラムを作成する必要がある（このプログラム作成は、作成された文法自体の、気のつかなかつて欠陥を明らかにしてくれるだろう）。

もし、言語のモデルが未完成のままで、計算機処理システムを作成してゆくとしたら、そのシステムは、以後の実験・検討によって頻繁に修正や追加ができるなければならないであろう。Production System にもとづく知識ベースシステムが、evolutionalなシステム開発を容易にするためのプログラミング手法であるとしたら、その手法は、自然言語処理プログラムを開発するための手法ともなり得よう。とくに、自然言語の文法記述によく使われる書き換え規則は、Production rule の考え方の基本となったものである。Production System あるいは、書き換え規則の考え方とは、個々の知識単位（自然言語処理では文法規則）の独立性を保ち、システムの保守容易性を向上させるという意味において有効であろう。

[記述の汎用性と多層性]これまでの書き換え規則の記述の欠陥は、自然言語の持つ規則性を品詞という单一のレベルで把えようとしてきたことである。§2で述べたように、言語の規則性は、品詞レベルだけでなく、

- ① 単・複の区別のような形態素的な情報
- ② HUMAN や Aspects のような意味的な情報
- ③ 特定の格助詞や表層上の単語並びに関する情報
- ④ 統語的な構造・意味的な構造に関する情報

といった、レベルのちがった情報を参照することによって記述される。したがって、文法規則を記述する体系、および処理結果を表現するためのデータ構造は、品詞とそれを基本にした統語構造といった单一の構造ではなく、1つの表現形の中に上記のような各種の情報が有機的に関連づけられて、記述できなければなら

ない。

このようば「記述の多重性」の考え方には、人工知能研究の他の分野にもあらわれている。例えば、画像理解における Marr の Primal Sketch の考え方も、外部知識という、画像データとは無関係なものに結びつく以前に、画像データがそれ自身で持つ情報をだけ落さないように、忠実に記述できること「多重の記述を持つべきだ」という主張であろう。例えば、画像データにおいて、Waltz の線分に関する規則性を計算機処理に反映しようとすれば、まず、画像データから線分の情報を抽出し、それから、線分に関する規則性を適用することになる。しかしながら、画像の持つ情報は、線分だけでなく、画に関する情報、濃度や濃度勾配に関する情報 etc があり、しかも、これらが一枚の画面の中に相互に関係し合って混在している。これら別の情報に関する規則性を表現しようとすれば、これらの情報を画面から抽出し、記号化あるいは数量化する必要がある。一枚の画面の持つ情報を、忠実に記号化や数量化しようとすると、いくつもの視点から行なわなければならぬ。また、線分に関する規則の記述が、その線分の周囲をとりまく面の情報を参照して記述されるとすると、この 2 つの記号化的結果は、相互に有機的に関係づけられている必要がある。自然言語処理においても同様である。我々は、まず、外部世界の知識という言語表現とは独立したものと結びつく以前に、言語表現のもの上記のような様々な情報を有機的に関連づけて表現しなければならない。品詞情報は、あたかも、画像における線分情報のような役割を果すものであり、重要な手掛りとはなるが、それですべての規則性を表現することはできない。これまでの書き換え規則は、品詞という单一の種類の情報に関する規則性を記述するものであつて、これを拡張する必要がある。

【アロダクション・システムと言語処理】これまで、アロダクション・システムは、知識工学における主要な手法として、各種の応用分野に適用されてきた。とくに、診断システム的な分野への応用が顕著である。何故アロダクション・システムが診断システム的な分野で成功してきたかを考えてみると（あるいは、診断システム以外には適用し難いか、とくに、その今まで言語処理に適用するものが何故困難であるか、を考えてみると）、アロダクション・システムには、次のような特徴があることがわかる。

診断システムにおける処理は、基本的にもとのデータから判明した情報が additive につけ加わる処理であり、つけ加わる順序や、つけ加わった情報間の相互関係は、それほど重要ではない。したがって、データが処理過程全体を通じて、加工されてゆくという通常の計算機処理とは異なる。

診断的システムでは、あるデータ集合から含意されるクラスタが決定される。パターン認識の framework と同じである。これに対して、航空写真のような画像から、その画像にうつっている対象物についての記述、あるいは、対象物相互間の関係に関する記述を得ようとすると、単に情報が additive に、既知の情報の集合の中につけ加わってゆくだけでなく、つけ加わる先の「既知の情報集合」 자체が構造（画像処理の場合には、空間的位置関係 etc）を持っているために、この

構造内に適切に新たな情報を加える必要がある。また、規則も、「既知の情報集合」内に、Aという情報が、他の情報とRという関係にあるれば、…」という形式で定式化されることが多い。このような定式化を可能にするためには、プロダクション・システムの規則間のコミュニケーションを支配するWM(Working Memory)にある種の構造を持たせる必要がある。もちろん、述語を使うことによって、通常のWMをそのまま使うこともできようが、この方式には煩雑な操作が必要となる。このような応用においては、WMが、Blackboardという別称が与えられる所以である。

言語処理の場合には、この「既知の情報集合」中の要素間の関係が、出発時点では単語の並びであるが、処理が進行するに伴って、木構造となり、さらに複雑になる。とくに、自然言語のもつ構造的な特徴のために、「既知の情報集合」の要素間に成立する関係が複数個生じることになり、問題がさらに複雑になる。これまでの書き換え規則によるシステムが使っていた、Recognition Matrixのようなもので、WMに相当する部分を augment する必要がある。

[データとプログラム] 人工知能研究で問題とされた「データとプログラム」の区別が、自然言語処理のプログラムを開発する際、すなわち、自然言語に関する我々の知識を記述する際に、どのような意味をもつのかについて考えてみよう。「データかプログラムか」の議論を我々なりに整理すると、次のようになる。

これまで「データとみなされてきた知識の多くが、実は、「入力データ」を処理するための手続を規定する」という、プログラムとしての性質を持っている。「知識」というデータがあり、それを処理するための基本的な手続がある、と考えるよりも、「知識」そのものがプログラムだと考えられる。ただし、この議論は、どのレベルで「記述」とみるかによって変化する。例えば、論理式集合は、theorem prover によって処理されるデータとみることもできるし、prolog にみられるように、これをプログラムとみることもできる。システムを作成する側からみると、ある記述を行なったときに、その記述形を解釈し処理するための記述（プログラム）を別途自分で用意する必要があるかどうかによって、その記述が「プログラムであるか、データであるか」が決まることになろう。

2に述べたように、言語表現の規則性の多くは、個々の単語に依存するものが多く、したがって、辞書中に多くの記述を行なう必要がある。システムの evolutional な改造も、この辞書を通じて行なわれることが多い。もし、辞書の記述を解釈し、処理に反映するには、プログラムを別途用意しなければならないとしたら、システムの保守容易性はきわめて悪くなる。すなわち、一般的な規則性を記述する文法規則と、個別単語に依存する辞書記述に、大きな区別があるので、辞書記述を変更すると、それに伴って、その記述を解釈するための一般的な文法規則も修正しなければならないとしたら、システムの保守容易性はきわめて悪くなる。辞書といふ、従来はデータとみなされていて記述を変更することによって処理手順も柔軟に変更されること、すなわち、辞書記述も一般規則と同じように取扱われるこれが、システムの保守容易性を向上させるためには必要である。

「辞書の記述を、一般文法規則の記述と同じレベルで取扱うようにする」という上記の主張は、「辞書記述の中に、使用者定義のプログラムを埋め込む」という

人工知能研究において従来行なわれていて主張とは異なる。一般文法規則も辞書記述も同じ枠組み記述がまき、かつ、(Prologにおいて論理式をプロログラムとみなすことができることの同じ意味)この記述を解釈実行する機能をシステム側が持つことが必要がある、という主張がある。

辞書記述や一般規則に、人工知能において従来行なわれていて「プロログラムの埋め込み」を行なうことには、上記の主張とはまた別の意味が必要がある。すなはち、簡単な意味標識のチェック以上の意味処理や、外部知識の処理を伴う文脈処理を行なうための機能等は、言語表現を処理するための機能とは別の(例えは、推論や演繹の機能)ものがあり、それ専用の記述の体制と処理手順を用意する必要がある。このような機能とのリンクは、言語現象を処理するためのソフトウェアシステムからみて、一種の「プロロジミングの埋め込み」として実現されるにあろう。

34. おわりに

現在我々のグループでは、以上のような目的意識から自然言語処理用のソフトウェアシステムを開発していきます。そのプロトタイプシステムは一応完成し、これによって実験的な文法を記述していきます。このようなソフトウェアシステムは、実際に文法記述を行ない、その有効性の検証と不備な部分の改良・拡充を不斷に行なってゆく必要がある。いずれにしても、自然言語処理は、我々が自然言語表現に關して持つべき知識を、いかにもうまく計算機能にかかる形式に記述するか、の研究である。適切な文法規則の作成と同時に、そのような文法を記述がまき、処理ができるソフトウェアの開発が重要である。

[参考文献]

- 1) 立井・中村：「自然言語処理のためのソフトウェア構造」，京都大学数理解析研究所講究録396, 1980
- 2) J. Bresnan: Realistic Transformational Grammar, in Linguistic Theory and Psychological Reality, MIT Press, 1979
- 3) N. Ostler: Case Linking; A Theory of Case and Verb Diathesis Applied to Classical Sanskrit, Ph.D Thesis, MIT, 1979
- 4) 寺津・福田・山梨：「日本語における照應現象について(その2)」，計算機による日本語談話行動の統合モデル化，昭和54年度研究報告書(代表者:石川敏雄), 1980
- 5) M. Kay: Functional Grammar, Xerox, PARC report (1978)
- 6) 中村：「自然言語解析のための文法記述システム」，京都大学修士論文，1981