

企画趣旨説明: 意味論研究の新地平

窪田悠介

(国立国語研究所)

1. はじめに

世間は今、ChatGPT に代表される生成 AI の話題で持ち切りである。これらの生成 AI は(言語学者が考えるような意味での) 文法を組み込んであるわけでもないのに何故か不気味なほどに流暢に日本語や英語などの自然言語を「話し」、少なくとも計算機上で完結するタスクについては、人間と言葉でやり取りしながらかなりの程度まで複雑な作業を行うことができる。このようなものが実現してしまった今、言語の文法や意味の研究は今後どのような方向に向かうのだろうか? 文系の言語研究にはまだ存在意義はあるのか? 本シンポジウムはこのようなことを考えるためのきっかけにしたいと思い企画した。言語を理解することはとりもなおさず意味を理解することであるという立場に立てば、意味の研究が言語の研究の中の最も中心的な課題の一つであることは疑いの余地がない。現在までのいわば「職人技」の集積ともいえる文系の世界での伝統的な言語研究のやり方と、ChatGPT などに代表される(コーパスと計算機資源の)物量にものを言わせた「力技」の言語処理とはどういう関係にあるのか? 意味の研究という領域にひとまず範囲を限定した際に両者に接点があるとすればどこにあるのか? 現在の意味論研究の最前線を眺める作業を通して、このような問題を考えてみたい。

2. 講師とコメンテータの紹介: 意味の研究に対する二つのアプローチ

意味に関する理論的研究としては、大きく分けて、用法基盤的な言語観と親和性が高い認知言語学と、記号处理的なアプローチをとる形式意味論の流れとがある。そこで、それぞれの分野で自然言語処理との関連を明確に意識した研究を行っている研究者である、大谷直輝氏と峯島宏次氏に講演をお願いすることにした。認知言語学の分野からは大谷直輝氏にご登壇いただく。大谷氏は主に英語学の分野で構文文法や用法基盤的なアプローチに基づく言語研究を進めており、最近では生成 AI などの背後にある大規模言語モデルの技術と用法基盤的な言語観との関係性を模索する立場から、自然言語処理の研究者との共同研究も行っている。大谷氏の研究の主な対象言語は英語だが、中心的な研究領域であるコーパスを利用した構文文法研究や語法研究は本学会の会員の研究関心と重なる部分が多く、そのような研究が自然言語処理の研究とどのようにつながりうるかということについて研究の最前線をお話しいただけるはずである。これに対して峯島宏次氏は、哲学と論理学の研究を出発点とし、最近では論理学的手法に基づく計算言語学・自然言語処理研究に精力的に取り組んでいる。特に、生成文法や形式意味論に代表される記号处理的アプローチと、自然言語処理の分野で現在主流の明示的な規則体系なしの "end-to-end" と呼ばれる手法とを組み合わせた研究で興味深い成果を発表している。この二つは一般に対立するものと捉えられがちだが、その対立を止揚し、融合的なアプローチによって自然言語の本質を探る試みの最先端を紹介していただく。このような二名の講師の講演から、自然言語の意味に関する研究

において、認知言語学、形式意味論、自然言語処理の三つのアプローチがどのような関係にあるかが見えてくるはずである。

意味論の理論研究も、自然言語処理の最近の展開も、本学会の会員にとっては必ずしも最も身近なものではないかもしれない。そこで、本企画の立案段階から、確かな見通しのもとに分野の橋渡しをすることができる方、そして、学術的な議論において重要な論点を曖昧なままにしておくことをよしとしない方にぜひコメントをお願いしたいと思っていた。幸いなことに、「この人しか考えられない」と思っていた三宅知宏氏にこの役をご快諾いただいた。三宅氏と講師二名との対話から、そしてフロアとの議論から、今後の意味研究、文法研究の道しるべが見えてくればと思っている。

3. 自然言語処理における言語の意味の扱い

内省に基づく記述研究や理論研究などの「職人技」の言語学は本学会の会員にとっては馴染み深いものであると思われるので、ここで改めて説明する必要はないだろう。そこで、以下では「力技」の言語処理の世界で言語の意味の問題に関してどのようなアプローチが取られてきたかということをごく簡単にまとめてみようと思う(最近の深層学習の手法や大規模言語モデルに関してさらに詳しく知りたい方は、一般書店の棚に並んでいる啓蒙書などを参照していただきたい)。誤解を恐れずに言えば、この手の技術が要するにどのようなものか、そして伝統的な「職人技」の言語研究とどう関連しうるかを大づかみに把握するためには、多少の頭の柔らかさは必要だが、数学や統計の知識はおろか、「理系的/データ・サイエンス的なものの見方」すらほとんど不要である。

3.1 分布主義に基づく意味観

まず、「意味」というものをどう捉えるかというのが、自然言語処理と(一般的な)言語学とで大きく異なる。もちろん言語学内部でも、「意味とは何か」という問題に関しては様々な見解がある。しかしながら、多くの場合、言語の「意味」というものは、何らかの離散的な記号体系により表記される知識の抽象的な表現として捉えられているように思われる。たとえば形式意味論における高階論理式での意味表現、生成文法における LF の「表示」、ジャッケンドフの概念意味論の意味表記などのようなものを考えてみていただきたい。(もちろん形式意味論の伝統を墨守し、真理条件的・モデル理論的意味論の立場を貫くならば、意味は真理条件でしかないので、高階論理の表記は単なる便法に過ぎないのだが、今ではこの立場は言語学内部における意味の概念の捉え方としても少数派だろう。)

自然言語処理で一般的な意味に対するアプローチはこのような記号主義的な立場と大きく異なり、徹頭徹尾分布主義的である。これは一言でいうと、単語の意味とはその分布にほかならない、という立場のことであり、実は決して新しい考え方ではない(少なくとも構造主義言語学の Firth が 1950 年代に主張したスローガン 'A word is characterized by the company it keeps' にまで遡ると言われている)。もう少し直感的に言うると、これは要するに、単語の意味を知っていることはその単語がどのような前後環境に出現するかを知っていることにほかならないということである。たとえば日本語の「リンゴ」という単語について考えると、「赤いリンゴが木になっている」「太郎は

昨日リンゴを二個食べた」「リンゴが八百屋に売っていた」などの単語列は文法的かつ意味的に整合性のある文であり現実の言語使用で出現する可能性が高いが、「#リンゴが太郎を食べた」は意味的に不適格なので、(小説や詩などの特殊なジャンルにおける逸脱的な事例でない限り) 現実のコーパスにはまず現れない。とすると、「リンゴ」という単語の共起情報の集積、ないし、その情報を何らかの形に効率よく圧縮した数学的表現を「リンゴ」という単語の意味(の近似)として扱ってよい、と考えることができる。このような共起情報は数学的には高次元のベクトル空間に表現することができ、(次元の圧縮などの工学的な処理により)それを効率よく扱う手法が発展している。十分に大規模なコーパスからそのような共起情報を抽出することで、任意の二つの単語の間の意味的類似度をベクトル空間における距離の計算として扱うことができ、これによってたとえば類義語の判定などが可能となる。この手の意味に対するアプローチは「分布意味論」や「(意味の)分散表現」などと呼ばれる。工学的には、言語の「意味」というとらえどころのない概念を、コーパスから機械的に抽出できる分布特徴であると割り切ってしまうことで、扱いやすいデータとして各種応用タスクに即座に利用可能となる点が最大の利点である。特に、言語の意味を、言語学で一般的な離散的な記号ではなく完全に連続量として扱えることに、実用面で大きな利点がある。

3.2 深層学習以降の動向

分布意味論的手法による自然言語の意味の扱いは内容語の意味に対しては非常に有効に機能し、自然言語処理研究においては、キーワード抽出や情報検索などのタスクにおいて深層学習のような複雑な手法が出てくる前から広く使われている。これに対して、機能語の意味などの、より抽象的な文法知識はこの手のアプローチではそのままの形ではうまく扱えない。たとえば形式意味論の研究で扱われるような、自然言語を用いた論理的な推論(「すべての学生が試験を受けた。太郎は学生である」=>「太郎は試験を受けた」)を正しく行うためには、否定、モダリティ、量化詞などの要素の解釈を踏まえた文の意味の正確な理解が必要となる。このようなタスクを正確に遂行するためには、記号处理的なアプローチ(つまり「職人技」の叡智の結晶)を何らかの形で組み込むことが必要であると考えられる。そして、機械翻訳や対話システムなど、少なくとも素朴に考えると文の意味の正しい理解なしには十分に正確に遂行できないように思われるタスクについても、当然記号处理的アプローチによる「深い」意味解析が正確にできることが前提とされるように思われる。実際、深層学習が現れるまでは、これらのタスクは実用レベルに達しておらず、またその理由もコーパスから機械的に分布的特徴を抽出するだけではこのような高次のタスクを正確に実行することは不可能であるからだと考えられていた。

ところが、深層学習の登場以来この手のより複雑であるはずのタスクの少なくとも一部(典型的には機械翻訳)も実際には end-to-end でかなりの程度解けることが明らかになった。何故このようなことが可能となったのか? 大雑把に言うと、これは計算機の並列計算の処理能力が向上し、web 上に蓄積されたテキスト・データが超大規模コーパスとして利用可能になったことにより、90年代の認知科学におけるコネクショニスト研究のアイデアが実際に計算機上に実装可能となったことによって実現したと考

えることができる。Google 翻訳や ChatGPT などの現代の深層学習に基づく自然言語処理のシステムはすべて、巨大なコーパスからその背後にある抽象的なパターンを暗黙知として保持する「大規模言語モデル」と呼ばれる言語処理に特化した深層学習のモデルを構築し、それを用いることで様々な個別のタスクを解いている。現代の大規模言語モデルが従来の単語分散表現と決定的に異なるのは、単なる表面的な共起情報を超えたより抽象的な言語の分布的特徴を何らかの形でモデルの内部構造の中に保持していると考えられる点である。深層学習のモデルはブラックボックスなので、これが工学的にどう実現しているかの解明はまだ十分に進んでいないが、要するに、従来「職人技」の言語学が品詞や統語構造のような概念で捉えようとしてきたたぐいの、より抽象的なレベルにおいて言語の分布的特徴を支配している法則も、この手のモデルは(おそらく)内在的にある程度「知っている」ということである。

工学的な実用性の観点からは、機械翻訳などのタスクに関してはこのような方法でコーパスの規模と深層学習モデルのアーキテクチャの複雑さに物を言わせて end-to-end の「力技」で解いたほうが、下手に「職人技」の記号处理的なアプローチと組み合わせるよりもうまくいくのである。そしてあっという間に ChatGPT が(少なくとも表面的には)流暢に言語で人間と対話する時代となった。ここまで来ると、この手の技術の背後にある大規模言語モデルが、「現実の抽象的」という自然言語の本質を何らかの形で体現したものになっているのではないかと考えたくなくなる。90年代のコネクションリスト研究が単なる計算機科学でなく認知科学でありえたのは、背後に「人間の脳を模した構成の処理系を電子計算機というハードウェアに実装し、それを用いて人間の脳が行っているような複雑な処理を行うことは可能か」という問いがあったからである。この問いに対して、深層学習技術の進展によりある程度興味深い答えが実際に提示されるに至ったという事実を、我々言語学者はある程度重く受け止める必要があるのではないかと。現実を抽象化する記号体系としての自然言語の有り様は、まさしく今までの「職人技」の言語学において、個々の専門家の鋭い言語直観を突き合わせる連綿たる営みによって我々言語学者が解明しようとしてきたものにほかならない。人間の脳そのものではないにせよ、それと似たような設計を持ち、そのような自然言語の本質を少なくとも部分的には捉えていると考えざるをえないシステムがいまや電子計算機上で実際に動作するのである。人間が言語の意味を「理解」するプロセスは、果たして生成 AI がプロンプトからの指示に対応する際にモデルの内部で行われている計算と何らかの意味で本質的に違うといえるのだろうか? これからの時代の言語研究、そしてより広く人間に関する科学の研究は、このような問題を少なくとも頭の片隅で意識しながら進めていかざるをえないように思われる。

4. まとめ

上のような問いにはもちろんすぐには答えは出ない。だが、やがてはそこに結びついていく具体的な研究課題を探索する営みが、現在意味論研究の最前線で行われている。意味とは何か? 認知言語学の意味観と形式意味論の意味観が自然言語処理における工学的な意味処理の技術を通して結びつくことはありうるのか? 言語学者による「職人技」の意味分析をこれからも続けていくことは今後も有効な研究方略であり続けるのか? 現在自然言語の意味研究は、こうした様々な根本的な問題を再考する時期

に差し掛かっている。本シンポジウムが、現在の意味研究の最前線を眺め、学会全体で、そして本学会の会員ひとりひとりがそれぞれの立場から、これらの問いへの考察を深めていくきっかけとなればと思っている。