

文章の「行間」を読むコンピュータの実現を目指して

東北大学大学院情報科学研究科 教授

乾 健 太 郎



「文章を理解する」という自然言語処理最大の未解決問題

自然言語処理技術は、この20年で発展した機械学習ベースの統計的手法の成功によって長足の進歩をとげた。しかし、実用的な水準に達している技術はいずれも、バラバラな文やフレーズから表面的な情報を抽出し、処理しているに過ぎない。自然言語処理のルーツであり、今なお毎年おびただしい数の論文が生産されている機械翻訳（自動翻訳）にしても、個別の文をバラバラに翻訳する方式が支配的であり、前後の文脈との繋がりを考慮して翻訳する試みは研究レベルでも極めて少ない。文章の文脈的な繋がりを認識する問題を談話解析あるいは談話理解と呼ぶが、文をバラバラに処理せざるを得ないのは、この談話理解がいまだに非常に手強く、大きな未解決問題として残っているためである。

自然言語の文章は、その中の個々の文がそれぞれ互いに何らかの意味で繋がっており、全体として一つのストーリーを伝える。いくつかの文を単に脈絡なく並べただけでは文章にはならない。文と文が互いに「繋がっている」とは例えば次のようなことである。

Ed shouted at Tim. He crashed the car.

第1文で Ed が Tim を怒鳴りつけている (shout at) のは、おそらく第2文にあるように Ed が車を壊してしまったからだろう。つまり、第2文は第1文と「理由」の

関係で繋がっている。この繋がりがわかるからこそ、第2文の He が第1文の Ed でなく Tim を指していることもわかるし、さらには第2文の the car がおそらく Ed の車であることも想像できる。このように文章の中の文と文は互いに「繋がって」おり、その繋がりがわかることが文章を理解する上で本質的に重要である。

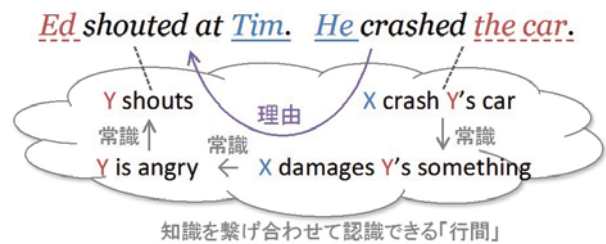


図1 文章の理解では文と文の意味的な繋がりを認識することが重要

ところが、これをコンピュータにやらせるとなると大問題である。第2文が第1文の理由になっていることを認識するためには、「怒鳴りつける」と「車を壊す」とこの間の因果関係がわからないといけない。われわれ人間であれば、他人のものを「壊す」のはその人に「迷惑をかける」ことであり、「迷惑」を被った人は怒って「怒鳴りつける」かもしれない、といった常識のようなものを持っていて、それを使って図1のような繋がりを

導くことができる(人間が毎回このような推論を本当にやっているかは議論の余地があるが、少なくともこうした繋がりを想起することは可能である)。しかしながら、コンピュータはそもそもこうした常識を持っていないし、常識を使いこなして推論する仕組みも脆弱である。上のような推論を頑健に行えるコンピュータはまだまったく実現していない。

これまでの談話解析の研究は、一文ごとの構文解析や機械翻訳でうまく行った機械学習ベースの方法を転用するアプローチが主流であった。単語や構文などの表層的な手がかりを特徴量として、正解付きの訓練データから教師あり学習によって解析モデルを訓練するアプローチである。しかし、上の例のように明示的に書かれていない事態や意味の関係の認識、いわば「行間」の解析まで考えると、表層的な手がかりだけで機械学習する方法の限界は明らかである。この問題は、従来の古典的な機械学習を現在大ブームのディープラーニングに置き換えたとしても解決しない。実際、文の境界を越えた共参照解析(上の例の *He* のような代名詞や省略などの参照先の解析)の精度は 50%にも届いておらず、接続詞が省略された文の間の談話関係(理由や例示等の意味的關係)を解析する問題も 60%程度の精度に留まっている。「行間」の存在は、何も「言わなくても伝わる」や「阿吽の呼吸」といったことがしばしば言われる日本語に限った話ではない。「行間」は上の例のように英語にも他の言語にもあまねく存在し、我々の言語コミュニケーションを効率的で心地良いものにしてくれる反面、コンピュータにとっては非常にやっかいなハードルである。

「行間を読む」コンピュータへの道のり

こうした自然言語の「行間」を読む深い理解をコンピュータ上に実現するにはどうすればよいか。第1に、単語や構文などの言語知識の他に、広範な常識的知識や背景知識をコンピュータに与える必要がある。上の例で言えば、「他人のものを壊す」ことと「その人に迷惑をかける」こと間の因果関係の知識だ。人間はこうした知識を無数に持っている。人手で構築した小規模な知識ベースでは実際の文章の理解にまったく歯が立たない。これは「知識構築のボトルネック」と呼ばれる問題で、自然言語処理がなかなか深い言語の理解に進めない最大の理由の一つとなってきた。また第2に、テキストに明示的に書かれた内容と関連する背景知識を組み合わせ

て、「行間」を含めたテキストの解釈を合成する推論機構を計算的に実現する必要がある。これは平たく言えば、先の図1のように知識の断片を繋げ合わせて、文章の部分どうしの意味的な繋がりを見つける処理と考えることができる。これらの課題は自然言語理解の本質的かつ究極的な目標として 1980年代から 90年代にかけて盛んに論じられた。しかし、当時は解決への道筋を見いだすことができず、その後の機械学習による統計的アプローチの隆盛とともに次第に顧みられることもなくなり、いわば「忘れられた課題」として現在に至っている。

ところが今、こうした状況が大きく変わりつつある。変革のトリガーになっているのは、言語ビッグデータの出現と計算資源の爆発的拡大、そしてクラウドソーシングだ。大量の言語データがネットから手に入るようになり、これまで決定的に欠けていた常識的知識をコンピュータ自身がそこから自動的に獲得できる可能性が出てきた。そうやって自動獲得した知識を実際に談話理解やそのための推論に使うことができるようになってくれば、大きなブレイクスルーになる。もちろん、「行間を読む」コンピュータへの道のりはまだ長く険しい。しかし、それに向けて技術は多方面で確実に進歩を重ねており、この大きな課題に再挑戦する素地が整いつつある。好機到来と考えるべきだろう。

東北大学・自然言語処理研究グループの挑戦

上述のような背景のなか、東北大学の我々の研究グループでも、大規模言語データからの知識獲得を発展させることによって知識構築のボトルネックを解消し、知識と推論によって文章の行間を解析する計算モデルを構築する研究を進めている。図2にその全体像を示す。

知識獲得については、世界最大規模の知識獲得基盤を構築し、とくに行間解析に必要な因果関係等の事態間関係知識の獲得に注力してきた。既存の大規模知識ベースには Freebase や DBpedia などがあるが、これらが所蔵する知識は固有の具体物(特定の人物、組織、施設など)に関する関係知識(「Barack Obama の出身地は Honolulu である」といった知識)が支配的であり、因果関係などの知識はほとんど入っていない。我々のねらいは、既存の知識ベースがカバーしていない常識的な知識を大規模に獲得する新しい仕組みを構築することである。個々の記述から抽出した断片的な事態間関係をどのように汎化し知識に集約するか、獲得した知識を言語理解に柔軟に適用するための知識表現(とくに言語の知

