

「意味が通る」や「面白い」に迫れるか？

An Approach to Understandable and Interesting Texts



名古屋大学大学院工学研究科教授

佐藤 理史

京都大学大学院工学研究科電気工学第二専攻博士課程研究指導認定退学。博士(工学)。北陸先端科学技術大学院大学、京都大学を経て、2005年より現職。

1 はじめに

適当に語をいくつか並べても、ほとんどの場合、まっとうな日本語表現にはならない。まっとうな日本語表現とするためには、文法的な制約を守るとともに、その表現が意味するところが理解されるものになっていなければならない。「意味の通る」日本語表現とは、そのような表現である。

コンピュータによる日本語文章生成では、意味の通る日本語文章を作る必要があるが、与えられた任意の文字列が意味の通る文章となっているかどうかを判定する技術は、まだまだ発展途上である。さらにその先には、わかりやすい文章になっているか、面白い文章になっているかを判定するという、より高度な課題が控えている。

私がことば遊びを題材とした研究を行うのは、これらの課題に対して、なんらかの突破口を見つけたいからである。これまで、クロスワードパズルを解く研究^[1,2]や回文やアナグラムを生成する研究^[3,4,5]を断続的に行ってきたが、今回、「名大 MIRAI GSC」^[6]というプログラムで高校生と一緒に研究する機会があり、その研究テーマとしてアナグラムの自動生成を選択した。

アナグラムとは、「言葉の綴りの順番を変えて別の語や文を作る遊び」(広辞苑)である。日本語のアナグラムは読み(ひらがな文字列)の順番を変えて作るのが一般的であり、たとえば、「大根(だいこん)」から「コインだ」を作ることがこれに相当する。

前回、アナグラムの自動生成に挑戦したときは、「意味が通る」を判定する有効な方法がなく^[5]、大量のア

ナグラム候補を前に途方に暮れたが、近年、利用可能になった大規模言語モデルによって、この状況を打開できるのであろうか。研究テーマを選んだ時点での私の興味は、この点にあった。

2 アナグラムの自動生成システム

今回作成したアナグラムの自動生成システム^[7]の構成を図1に示す。自動生成は、候補リスト生成とランキングの2ステップで構成されている。

前半の候補リスト生成のステップでは、まず、読み(ひらがな列)の並び替えをすべて作成し、そのそれぞれをかな漢字変換する。たとえば、「名古屋大学(なごやだいがく)」の読みは7文字なので、読みの並び替えは $7! = 5,040$ 個存在する。候補生成の最初の段階では、すべての候補を作り、かな漢字変換する。次に、か

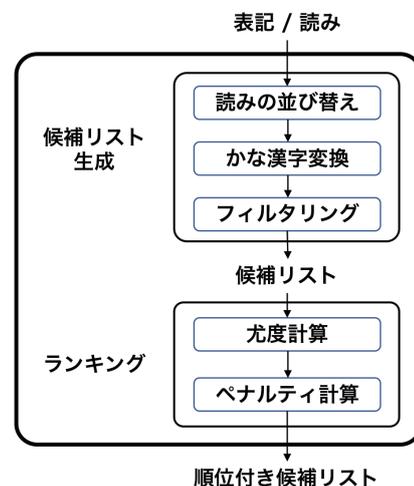


図1 アナグラム自動生成システムの構成

な漢字変換結果のうち、日本語の文字列としてほとんど可能性がないものをフィルタリングする。このフィルタリングには、国立国語研究所が公開している日本語文字 n-gram データ^[8] を利用し、頻度 10 回以下の文字 3-gram を含む場合は、その候補を削除する。このフィルタはかなり強力で、「名古屋大学」の場合は候補は 435 個に絞られる。

後半のランキングステップでは、まず、GPT-2^[9,10] を利用して候補文字列の尤度（日本語としてのもっともらしさ）を計算する。「名古屋大学」の例では、一番尤度が高いのは、当然のことながら「名古屋大学」となり、「大学名古屋（だいがくなごや）」、「名護大学や（なごだいがくや）」がそれに続く。これらの候補は、入力と重複する部分が多く、意外性に乏しい。そのため、入力と部分的に重複する候補にはペナルティを与える。具体的には、表記において共通する文字 3-gram と文字 2-gram の数を、読みにおいて共通する文字 4-gram と文字 3-gram の数を数え、それらの数からペナルティを計算して、尤度から差し引く。

「名古屋大学」に対して最終的に得られる結果の上位 10 位までを表 1 に示す。以前の研究で、「長い誤訳だ（ながいごやくだ）」というアナグラムがあることは知っていたが、本システムを動かして、「語学嫌だな（ごがくいやだな）」があることを初めて知った。

表 2 に、本システムで作成した都道府県名のアナグラムをいくつか示したので、元になった都道府県名を考

えてみてほしい。ひらがなに開いてみると、答を見つけやすいかもしれない。すでに述べたように、日本語のアナグラムは読みのアナグラムなので、表面上は似ても似つかないものとなる。この点が、表記を直接並べ替える英語などのアナグラムとは大きく異なる。

表 2 都道府県名のアナグラム例

あんた聞け	なんか毛皮	顔負けやん
焼けた漫画	マトン焼け	火消しロマン
嫌い馬券	和紙関係	クマと試験
また計算	おかず試験	なさげが金

3 大規模言語モデルの性能

先に述べたように、この研究の当初の興味は、大規模言語モデルの性能を身をもって知りたいということにあった。研究を実施した時点（2021 年 8 月）で利用できる日本語の大規模言語モデルはいくつかあったが、テキストの尤度計算に適している GPT-2 の日本語版を用いた。

システムの出力結果を眺めたところ、次のような傾向が観察された。

- ①意味をなす候補が上位に来ることが多い。これは、言語モデルの尤度が、全体としては日本語らしさを反映していることを意味する。
- ②しかし、尤度が高いからといって、必ずしも意味をなすとは限らない。

表 1 「名古屋大学（なごやだいがく）」に対する出力の上位 10 件

スコア	ペナルティ	表記類似度	読み類似度	GPT-2 尤度	表記	読み
-27.053	0.000	0.000	0.000	-27.053	長い 訳語 だ	ながい やくご だ
-27.232	0.000	0.000	0.000	-27.232	訳語 が ない だ	やくご が ない だ
-27.538	-4.500	0.125	0.325	-23.038	名護 大学 や	なご だいがく や
-27.857	0.000	0.000	0.000	-27.857	嫌 だ な 語学	いや だ な ごがく
-28.397	0.000	0.000	0.000	-28.397	大 納屋 が ごく	だいな や が ごく
-28.609	0.000	0.000	0.000	-28.609	嫌 な 語学 だ	い や な ご が く だ
-29.446	0.000	0.000	0.000	-29.446	大 梁 が ごく	だいな や が ごく
-31.420	-9.667	0.542	0.425	-21.753	大学 名古屋	だいがく なごや
-31.650	0.000	0.000	0.000	-31.650	抱く ご 長屋	いだく ご ながや
-31.811	0.000	0.000	0.000	-31.811	語学 嫌 だ な	ごがく いや だ な

③さらに、尤度がそれほど高くなくても、意味をなす候補も存在する。

体感としては、6から9文字のアナグラム候補では、尤度マイナス30がひとつの閾値で、それよりも尤度が高いところには、意味の通る候補が存在することが多い。ただし、その割合は半数程度である。一方、尤度がマイナス30より小さい場合は、意味が通らない場合がほとんどである。しかし、まれに意味が通る面白い候補が見つかることもある。

今回のシステムが有効に機能するのは、大規模言語モデルの尤度計算に負うところが大きい。次節で示すような面白いアナグラムを多数見つけることができるようになったのは、大規模言語モデルのおかげである。その一方で、大規模言語モデルの尤度だけでは、意味の通るアナグラムを十分に絞り込むことはできない。その事実も明らかになった。

4 面白さは多様

アナグラムがそれなりに生成できるようになった時点で、研究目的が「アナグラムを作る」から「面白いアナグラムを作る」に深化し、興味の対象も「アナグラムの面白さ」にシフトした。

4.1 面白いアナグラムを探す

色々な入力に対するシステムの出力を観察したところ、「おお、これは！」というような面白いアナグラムがいくつか見つかった。それが契機となって、「アナグラムの面白さって何なんだろう？」という疑問が湧き上がってきた。

意味が通るアナグラムの中には、面白いものと、そうでもないものが存在する。それを区別するものは何なのか。

当然のところながら、どんなものを面白いと思うかは、個人に依存する。とは言え、比較的多くの人が面白いと同意するものがあるのも事実である。そうであれば、「面白さ」に対して、その要因（理由）を考えることができよう。

そのためには、まずは面白いアナグラムをたくさん集める必要がある。そこで、大量の入力に対して、アナグラム候補を生成し、その出力を調査して面白いアナグラ

ムを探した。入力としては、形態素解析用辞書 IPADic に含まれる読みの長さが6文字から9文字の普通名詞(6,750語)を用いた。

こうして収集できた面白いアナグラムは、500件以上にのぼる。これらを観察して、面白いアナグラムの要因を調査・分析し、以下に示すように、大きく2軸に分けて整理した。

4.2 出力単体の面白さ

出力単体の面白さとは、元になった語（入力）が何だったかに依存しない面白さである。以下に、下位分類と具体例を示す。

①文形式となるアナグラム

例：頭でっかち → ちまたで悪化
功労賞 → 売ろうコショウ

文は、事態を表す形式であり、読み手は、その事態を想起することができ、具体的なイメージが湧きやすい。そのイメージを我々は面白いと感じるのだろう。

②意外性のある複合語となるアナグラム

例：サウンドトラック → 皿うどん特区
相対性理論 → 総理生熊論

「皿うどん特区」という特区は、おそらく、この世には存在しない。しかし、あっても不思議ではなく、ちゃんと意味が通じる。それが、好奇心を刺激する。

③ナンセンス系アナグラム

例：一日千秋（いちじつせんしゅう）
→ セイウチ出陣（せいうちしゅつじん）
般若心経 → 社運は人魚

このタイプは、思いもよらなかった意外な組み合わせに想像力が刺激される。

④オノマトペを含むアナグラム

例：一昨昨週（いっさくさくしゅう）
→ サクサク一周
ニッポンチャレンジ → レンジにちゃっぼん

オノマトペを含むアナグラムが生成されることは比較的少なく、希少性がある。さらに、オノマトペには臨場感があるので、イメージを想起しやすい。

⑤評価や感情を表す形容詞を含むアナグラム

例：風林火山 → うんざり花粉

ジャイアントパンダ → ジャパンだと安易

このようなアナグラムを面白いと感じる理由は、面白

さの評価が、感覚や感情に基づく評価であることが関係しているのかもしれない。

4.3 入力と出力の関係性の面白さ

一方、入力との関連で面白いアナグラムもある。

①入力と強く関連するアナグラム

例：知識工学 → 講師が鬼畜

方角違い → 北緯が違う

入力と関連するアナグラムを人間が作るのは、かなり難しい。制作者のスキルの高さ（巧みさ）を評価して、面白と思うのだろう。

②入力とひと続きの文になるアナグラム

例：一日中 → 知事に注意

百科事典 → 光ってんじゃ

さらに、入力とひと続きの文になると、「うまい！」と感じる。

③ダジャレ系アナグラム

例：垂成層圏（あせいそうけん） → 麻生政権

交感神経 → 更新関係

これは説明不要であろう。

④入力とギャップがある場合

例：天地神明 → メンチ進呈

全知全能 → 膳所の運賃（ぜぜのうんちん）

逆に、入力との意味的ギャップが面白い場合もある。何だかよくわからないところが、逆に面白いのであろう。

4.4 面白さの要因

これら、「面白い」と感じるアナグラムは、次のように総括できよう。

①具体的なイメージを想起しやすいアナグラム

②驚きや意外性のあるアナグラム、想像力や好奇心を喚起するアナグラム

③入力と強く関連するアナグラム、入力との落差がはげしいアナグラム

この総括からわかるように、面白さの要因はひとつではない。複数の要因が組み合わさると、面白さが際立ってくる。

今回のプロジェクトで見つけた中で、最も面白いと判断したのは、次のアナグラムである。

他力本願 → 本気が足りん

このレベルのアナグラムは、創作とみなしてもかまわ

ないだろう。

5 おわりに

「何の役に立つのですか？」

この研究に対してよく受ける質問はこれである。この質問には、次のように答えることにしている。

「まったく役に立ちません。でも、面白いでしょ？」

本当に役に立たないかどうかは、実はよくわからない。この研究のデモンストレーションでは、見に来てくださった方の名前を入力して、その場でアナグラムを作るデモを行うのだが、これがなかなか好評である。先日も、「英雄の歌」というアナグラムが出力されて、歓声が上がった。研究室の学生のひとは、自分の名前のアナグラムを、就職したときの自己紹介で使えると喜んでいる。

ことばには、人を楽しませる力がある。そんなことばを、コンピュータは作ることができるようになった。人を楽しませることは、「役に立つ」ことには入らないのだろうか。

良い機会なので、「日本特許情報機構（にほんとつきよじょうほうきこう）」のアナグラムを探してみた。読みの長さは15文字なのでなかなか手強かったが、システムをまる2日動かして約11億の並べ替えを調べた結果、上位100件から次のようなアナグラムが見つかった。みなさんは、この中でどれを面白いと思いますか？

1. 発起人と故郷情報
2. 工場基本法に特許
3. 特許法に基本条項
4. 実況に根拠と方法
5. 発起人と候補状況
6. 京都っ子に基本情報
7. 今日本当に国境保持
8. 今日ニコッと基本情報
9. 根拠法と情報日記
10. 情報金庫に特許法



参考文献

- [1] 佐藤理史. 日本語クロスワードパズルを解く. 情報処理学会自然言語処理研究会 NL-147-11, pp.69-76, 2002.
- [2] 内木賢吾, 佐藤理史, 駒谷和範. 日本語クロスワードパズルのカギの解法. 情報処理学会第 74 回全国大会講演論文集, 3R-5, pp.2-267-268, 2012.
- [3] 鈴木啓輔, 佐藤理史, 駒谷和範. 文頭固定法による効率的な回文生成. 言語処理学会第 17 回年次大会発表論文集, pp.826-829, 2011.
- [4] 鈴木啓輔・佐藤理史・駒谷和範. 文節データベースを用いた日本語アナグラムの自動生成. FIT-2011 (第 10 回情報科学技術フォーラム), RF-009, 第 2 分冊, pp.97-102, 2011.
- [5] 鈴木啓輔, 佐藤理史, 駒谷和範. アナグラム生成における文節列の意味的適格性の判定法の検討. 言語処理学会第 18 回年次大会発表論文集, pp.1308-1311, 2012.
- [6] 名大 MIRAI GSC.
<http://nuqa.nagoya-u.ac.jp/miraigsc/>
- [7] 土井遥, 山本優衣奈, 佐藤理史. 面白いアナグラムとはどんなアナグラムか. 言語処理学会第 28 回年次大会発表論文集, pp1472-1477, 2022.
<https://youtu.be/dH45d5wj9EI>
- [8] GSK2020-C 「国語研日本語ウェブコーパス」
n-gram データ・頻度表.
<https://www.gsk.or.jp/catalog/gsk2020-c/>
- [9] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever. Language Models are Unsupervised Multitask Learners.
https://d4mucfpksywv.cloudfront.net/better-language-models/language_models_are_unsupervised_multitask_learners.pdf
- [10] GPT2 日本語モデル.
<https://github.com/tanreinama/gpt2-japanese>

