

自然言語処理を用い日本語文章に対応した ピクトグラムデザイン生成システムの提案

廣橋 宣汰[†] 小坂 洋明[†]

奈良工業高等専門学校 システム創成工学専攻電気電子システムコース[†]

1 はじめに

外国人の日本の文化・歴史への興味は高く、コロナ収束後旅行したい国や地域をアジア住居者と欧米豪住居者に聞いたアンケートではどちらも日本が1位になっている¹⁾。そのような事から、外国人の日本語学習人数も年々増加している。一方、日本語教師の人数は過去8年間横ばいになっており、教師不足が問題視されている²⁾。原因として、高い言語力や異文化への受け入れが求められるなどがある。また、日本語学習は世界有数の難易度を有しており、教育方法の工夫が求められている。現在、日本語学習の工夫としてユニバーサルデザイン(UD)を使用した教育方法が提案³⁾されている。UDとは、1980年にロナルド・メイスが提唱した用語で、「身体能力の違いや年齢、性別、国籍に関わらず、すべての人が利用しやすいようにつくられたデザイン」を意味する。日本語学習に使用されるUDの一例にピクトグラムが存在する。ピクトグラムは、文字がわからない人に対し図や形などで意味を表すものである。一方ピクトグラムが示す意味が正しく伝わらないといった問題が生じている。その原因は、デザイナーによって同じ単語から作られているがデザインが異なっていることや、汎用性を重視しているため一つのピクトグラムが複数の意味を表してしまうことにあると著者は考えている。そこで本研究は、ピクトグラムデザインの作成方法を統一化すべく、自然言語処理を用い日本語文章から分かりやすいピクトグラムデザインを自動生成するシステムの構築を目指している。本稿では、試作した、自動生成システムについて記す。

Proposal for pictogram design construction system using natural language processing

[†]Senta Hirohashi, Hiroaki Kosaka · Electrical and Electronic Systems Course, Department of Systems Innovation Engineering, National Institute of Technology, Nara College

2 自然言語処理

自然言語処理とは、人間が日常的に使用している言語(自然言語)を自動的に処理し結果を応用するといった技術分野である。具体的には、機械可読辞書とコーパス、形態素解析、構文解析、意味解析という一連の処理をAIが行うことである。本研究ではオープンソース形態素解析エンジンである「MeCab⁴⁾」を使用する。MeCabは言語や辞書、またデータベース化された言語資料であるコーパスに依存しない汎用的な設計であり条件付き確率場に基づく高い解析精度を保持している。

3 提案システム

3.1 概要

本システムの開発環境はGoogle Colaboratory上で使用言語はpython(Ver. 3.7.15)である。このシステムには主に、形態素解析エンジンであるMeCab(Ver. 3.0.7)、画像生成や処理に使用したライブラリPillow(Ver. 7.1.7)、fastTextの日本語の学習済み単語ベクトルの一種であるWEBクローラーとWikipediaの文章から学習したモデルを使用するためgensim(Ver. 3.6.0)を導入している。その他モジュールやライブラリは省略する。

本システムのピクトグラムの生成フローは以下の通りである。最初に、使用者が日本語文章を一文入力する。この時、主語、動詞、述語はそれぞれ1つずつに制限をする。次に、その文章をMeCabにより形態素解析し意味が分かる最小単位である単語に切り分ける。続いて、ピクトグラム生成に必要な品詞の抽出をする。この時、抽出する品詞は「助動詞(表層形)・形容詞(基本形)・動詞(基本形)・名詞(基本形)」である。ここで表層形とは文章から抽出されたまま単語にし、基本形とは抽出された単語を日本語の学習済み単語ベクトルのモデルに掲載されている言葉へ変換したものである。最後に抽出された単語が示すピクトグラムを一枚の画像に出力する。

3.2 素材と配置テンプレート

上記の方法でピクトグラムを一枚の画像にする際、Fig.1に示すようなテンプレートを設定した。この時、人間がFの字をなぞるように画面を注視するといった根拠をもとに設定している。左上から、天気や周りの環境などの「状態」、その右に否定か肯定かを表す「是非」、左下に人や犬などの「名詞1」、その右に行くや入るとなどの「動詞」、最後に病院や公園などの「名詞2」といった配置テンプレートを複数設けている。

現在、本システムには人や犬などの名詞1を表すピクトグラムが10枚、公園や病院などの名詞2を表すピクトグラムが20枚、その他ピクトグラムを20枚、計50枚搭載している。

3.3 デザイン出力結果

本システムによって自動で生成されたピクトグラムデザインをFig.2に示す。この時使用した文章はaが「犬を公園に入れなくてください」、bが

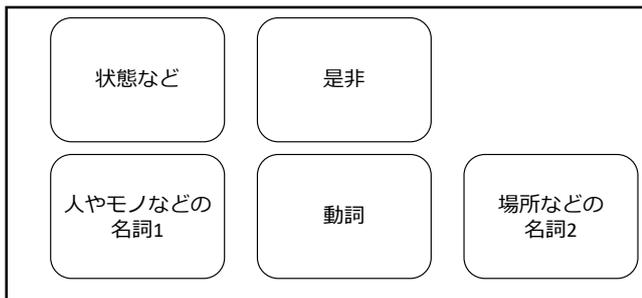


Fig.1 ピクトグラム配置テンプレート



Fig.2 構築システムの出力結果

「火事だ、逃げろ！」である。aにおいて文章は「犬・を・公園・入れ・ない・で・ください」に形態素解析される。次に目的品詞を抽出すると「犬・公園・入れる・くださる・ない」が出力され、「犬・公園・入れる・ない」の意味を表すピクトグラムデザインがFig.1に示すテンプレートに従い一枚の画像に出力されている。bにおいても同様に処理され「火事・逃げろ」からピクトグラムデザインが生成されている。Fig.2において、aは「犬」という主語を設けているがbは主語を設けていないにもかかわらず人を表すピクトグラムを出力している。これは、著者が「逃げろ」という単語に対して人と右矢印を表現しようと考えたからである。このように本システムは、1つの単語に対し複数のピクトグラムを出力することが可能である。

4 まとめ

ピクトグラムの意味が正しく伝わらないといった問題に対し、文章からピクトグラムデザインを生成することが有効だと考えた。手法として、MeCabを使用し、日本語文章から抽出した単語とその単語を意味するピクトグラムを配置テンプレートに配置することでピクトグラムデザインを自動で生成するシステムを構築した。その結果「犬を公園に入れなくてください」などの簡単な文章をピクトグラム化することに成功した。今後複雑な文章のピクトグラム化や2文以上の同時ピクトグラム化が可能なシステムへの改良や、JIS規格ピクトグラムをデータセットに使用したピクトグラム認識システムを構築し、本システムの評価に使用することを予定している。

参考文献

- 1) (株)日本政策投資銀行・(硬材)日本交通公社: DBJ・JTBF・アジア・欧米豪訪日外国人旅行者の意向調査(第3回 新型コロナ影響度特別調査), (2022/2/28)
- 2) 文化庁国語課: 令和3年度 日本語教育実態調査報告書「国内の日本語教育の概要」, (R3/11/1)
- 3) Yokoyama Rieko: 日本語教育における教育のユニバーサルデザインの提案, (基礎教育保障研究 第6号)(2022/8)
- 4) Taku Kudo, Kaoru Yamamoto, Yuji Matsumoto: Applying Conditional Random Fields to Japanese Morphological Analysis, Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing (EMNLP-2004), pp.230-237 (2004.)