

# 単語分散表現における女性標示語のステレオタイプの定量化

陳 蕾思<sup>†</sup> 杉本 徹<sup>‡</sup>

芝浦工業大学大学院理工学研究科<sup>†</sup> 芝浦工業大学工学部<sup>‡</sup>

## 1. はじめに

人工知能における差別と公平性の問題は重要な課題の一つである。Bolukbasi らの研究における、単語ベクトル計算で  $\overrightarrow{man} - \overrightarrow{woman} \approx \overrightarrow{programmer} - \overrightarrow{homemaker}$  という結果[1]は、人間社会のジェンダーバイアスが機械学習の結果の公平性に影響を与えていることを証明した。

一方で、伝統文化の影響を受け、日本語では男性視点から女性を描く言葉が多い。男性中心社会の価値基準で、女性の希少性を強調した「女性標示語」である[2]。例えば、「大学生」という言葉は性別が特定されていないのに対し、女性の大学生を表す「女子大生」という特定の単語がある。コーパスに含まれるこのような女性標示語に女性へのステレオタイプがあるか、人工知能の公平性に影響を与えているかに関する研究はまだ少ないのが現状である。

そこで、本研究では、単語分散表現を用いた日本語の女性標示語のベクトル化により、女性標示語と性格用語の関係を分析し、女性標示語に隠されている女性へのネガティブなステレオタイプを分析する。

## 2. 単語分散表現を用いた女性標示語のステレオタイプの定量化

### 2.1 単語分散表現

ニューラルネットワーク言語モデルの word2vec は、単語のベクトルを生成することで、非構造化テキストを構造化された数学モデルに変換することを可能にする。単語ベクトルにより、単語の意味だけでなく、コサイン類似度計算を通して単語間の関係も把握することができる。そこで、本研究では、女性標示語の語彙にステレオタイプがあるかどうかを明らかにすることを目的として、女性標示語とポジティブ・ネガティブな性格形容詞の関係を計算するため、word2vec による単語分散表現を用いる。

Quantifying stereotypes of gender-specific words in word embedding

<sup>†</sup>Leisi Chen <sup>‡</sup>Toru Sugimoto

<sup>†</sup>Graduate School of Engineering and Science, Shibaura Institute of Technology

<sup>‡</sup>Faculty of Engineering, Shibaura Institute of Technology

### 2.2 コーパス

本研究では、日本語 Wikipedia の記事データから作成したコーパスを用いた。Wikipedia は誰でも匿名で編集ができ、大量の記事があるという特徴を持っている。そのため、社会で特定の身分の女性がどのように認識されているかを一定のレベルで反映していると考えられる。

### 2.3 女性標示語と性格特性用語の選定

Google 1gram, 2gram, 3gram の中で「女」という文字から始まるもので出現回数が 1 万回以上のものから「女～」「女性～」「女子～」のような 40 個の女性標示語を選び本研究の実験対象とした。また、これらの女性標示語に対応する性別中立語 40 語（例：「女子大生」に対する「大学生」）も実験対象として選定した。

本研究では女性標示語に含まれるネガティブな印象を直感的に捉えるために、性格特性を表す単語を用いる。性格特性用語リスト[3]における、男女いずれかの選択率が 70% 以上の性格特性用語と、男女とも選択率が 10% 以下の性格特性用語の中から、ポジティブな意味の性格用語とネガティブな意味の性格用語をそれぞれ 15 個抽出した。

選定した女性標示語と性別中立語、性格特性用語の分散表現をコーパスから word2vec を用いて生成する。コーパスの分かち書きには MeCab を使用した。例えば「女子大生」という単語が「女子大」と「生」に分割されないように、あらかじめ MeCab の単語辞書に女性標示語をすべて追加してから分かち書きを実行した。

### 2.4 女性標示語のステレオタイプの可視化

本研究では、Bolukbasi らの研究[1]を参考にし、性格特性用語と女性標示語の関係を 2 次元平面における単語の分布として表す。ここで平面の X 軸は性別方向への射影を表し、X 軸の正方向は女性標示語、負方向は性別中立語の方向を表すとする。そのために、ある女性標示語とそれに対応する性別中立語の分散表現の差ベクトル  $v_{gender}$  を求めて、性格特性用語の分散表現  $v$  と  $v_{gender}$  の内積の値をこの性格特性用語の X 座標と

