

ルールベース機械翻訳によって常体⇔敬体へと変換するプログラム

情報班：小山 翼

1. はじめに

今日、外国人観光客や在留外国人が増え、外国人が日本語を使う機会が増えた。その時、彼らを困らせるのが敬語の使い方である。日本語の敬語はとても複雑であり、日本人でさえ、敬語を正しく使えない人もいる。しかし、日本で暮らすには敬語の使用は避けられない。以上の理由から、常体の文を敬体の文に翻訳するプログラムを作成すれば、日本語が得意でない人にとって、敬語の学習をする時や目上の人と話す時等、様々な場面で助けになるのではないかと考え、この研究を進めたいと思った。機械翻訳には、コーパスと呼ばれる言語データベースを用いて翻訳するルールベース機械翻訳と、大量の対訳を AI に学習させて翻訳する統計的機械翻訳がある。統計的機械翻訳については、敬語の対訳が少なく、プログラムを作成するのが困難と考えた。よってルールベース機械翻訳を用いて研究を進める。

2. 作成方法

プログラミング言語は python、実行環境は Google Colaboratory で行った。形態素解析には MeCab という形態素解析エンジンを用いた。また、翻訳に必要な辞書は IPA 辞書を用いた。

《方法》

(1)python の辞書型を用いてよく使われる尊敬語の表(sonkeigo)と謙譲語の表(kenjogo)を作成した(図 1 参照)。この表を用いて単語の変換を行う。

(2)MeCab を用いて、翻訳したい文について形態素解析をする関数(mecab)を作成した(図 2 参照)。この関数では翻訳したい文を引数(task)として、原文を単語に区切ったリスト(sentence)、それぞれの単語の原型のリスト(genkei)、それぞれの単語の活用形のリスト(katuyo)、それぞれの単語の品詞のリスト(hinshi)の 4 つの戻り値を返すようにした。また、尊敬語と謙譲語の両方に訳せるようにするため、変数 sonkei が 1 ならば尊敬表現に、変数 kenjo が 1 ならば謙譲表現に翻訳するようにした。

(3)genkei の中に sonkeigo もしくは kenjogo と共通した単語を探し、その単語に対応する敬語に変換するプログラムを作成した(図 3 参照)。しかし、変換先の敬語を正しく活用させなければ不自然な文になってしまうことが判明した。例えば、「食べた」の語幹である「食べ」を、そのまま「召し上がる」に変換すると、「召し上がるた」になってしまう。よって、目的単語の原型と原文における目的単語の品詞と活用形から、目的単語を指定された活用形に変換するプログラム(kensaku)を作成した(図 4 参照)。この関数では、IPA 辞書の動詞をまとめたファイル(Verb.csv)を読み込み、その中から入力した単語の原型と活用形に該当する単語を探し出し、戻り値として返すようにした。しかし、変換する単語とその敬語の活用形が一致しない場合があった。例えば、「食べた」における「食べる」の活用形は連用形であるが、「召し上がった」

における「召し上がる」の活用形は連用タ接続である。もし、変換先の敬語の活用形を連用形のまま訳してしまうと、「召し上がりた」になってしまう。そこで、変換する単語に続く助詞が「う」もしくは「た」、「て」であるか判定し、「う」であれば活用形を「〇〇ウ接続」に、「た」「て」であれば「〇〇タ接続」にする関数(henkan)を作成した(図5参照)。関数 henkan には関数 kensaku を組み込み、一つの処理にまとめた。

(4) 丁寧語を付け加える処理をするため、(3)までの処理をした文を変数 semifinal に格納し、それを関数 mecab の引数として形態素解析を行い、(2)と同じように、単語のリストを sentence に、原型のリストを genkei に、活用形のリストを katuyo に、品詞のリストを hinshi に代入するようにした(図6参照)。新しく得たこれらのデータを用いて、丁寧語を付け加える処理を行う。

(5) 文中の助動詞「だ」を「です」に変換するプログラムを作成した(図7参照)。(3)で作成した関数 henkan を用いて、活用形を合わせて変換できるようにした。

(6) 助動詞「ます」を文末に付け加えるプログラムを作成した(図8参照)。このプログラムでは、助動詞のリスト(jodoshi)(図9参照)を用いて文末の単語が、「ます」が接続できる動詞型の活用をする助動詞であるか、または動詞であるかを判定し、その単語を a とする。そして「ます」は連用形に接続するので、a を連用形にして、「ます」を a の活用形に合わせた上で、a の直後に活用させた「ます」を付け加えるようにした。また、原文に「ます」が含まれている場合は、break を用いてすぐに処理を終えるようにした。

(7) 敬体の文を常体の文に翻訳するプログラムも作成するために、(1)で作成した辞書型の尊敬語と謙譲語の表の key と value を逆にした新しい表を作成した(図10参照)。これを用いて(1)から(3)の処理を行えば、敬体の文を常体の文にすることができた。また、丁寧語の処理については、助動詞の「です」を「だ」に変換し、助動詞の「ます」を取り除くプログラムを作成した(図11参照)。

3. 成果

ルールベース機械翻訳によって、特定の動詞を尊敬語もしくは謙譲語に変換した上で丁寧語を付け加える、またはその逆方向に変換するプログラムを作成できた。

例 ・「先生はリンゴを食べた」→「先生はリンゴを召し上がりました」

・「私はコロッケを召し上がりました」→「私はコロッケを食べた」

しかし、「聞く」のように「伺う」「拝聴する」といった複数の敬語表現がある単語について、訳し分けができなかった。

4. 今後の展望

常体から敬体、敬体から常体の双方向に翻訳するプログラムを作成できた。しかし、ニュアンスの違いによる訳し分けはできなかった。そこで、常体と敬体の対訳データを集め、統計的機械翻訳のプログラムを作成すれば訳し分けも可能ではないか、と考察した。また、今回作成したこのプログラムを誰でも利用できるように、Web フレームワークを用いて、Google Colaboratory においてだけでなく、Google などのブラウザ上でも利用できるようにしたい。

5. 謝辞

本研究を進めるに当たり、適切な助言を賜り、丁寧に指導して下さいました大阪工業大学の小林裕之教授に深く感謝申し上げます。

6. 参考文献ならびに参考 Web ページ

@ menon(2019-11-08), 『Python と MeCab で形態素解析 (on Windows)』, <<https://qiita.com/menon/items/f041b7c46543f38f78f7>>, (参照 2019-8-14)

@ R-Yoshi(2018-12-05), 『【機械学習】Google 翻訳(みたいなもの)を自作してみました。』, <<https://qiita.com/R-Yoshi/items/9a809c0a03e02874fabb>>, (参照 2019-11-6)

```

5 sonkeigo = {'なさる': 'する',
6             'いらっしゃる': '来る',
7             'ご覧になる': '見る',
8             'お耳に入る': '聞く',
9             'おっしゃる': '言う',
10            'お思いになる': '思う',
11            'ご存じです': '知る',
12            'くださる': 'あげる',
13            'お受けになる': 'もらう',
14            '召し上がる': '食べる',
15            'お亡くなりになる': '死ぬ',
16            'お気に召す': '気に入る'}
17
18
19 kenjogo = {'おる': 'いる',
20            '参る': '来る',
21            '拝見する': '見る',
22            'ご覧に入れる': '見せる',
23            '伺う': '聞く', '申し上げる':
24            '言う', '存じ上げる': '知る',
25            '差し上げる': 'あげる',
26            '頂戴する': '貰う',
27            'いただく': '食べる',
28            '亡くなる': '死ぬ',
29            'お目にかかる': '会う',
30            '拝借する': '借りる'}

```

図 1 尊敬語と謙譲語の表

```

43 def mecab(task):
44     tagger = MeCab.Tagger("-Ochasen")
45     node = tagger.parseToNode(task)
46     sentence = []
47     genkei = []
48     katuyo = []
49     hinshi = []
50
51     while node:
52         word = node.surface
53         wclass = node.feature.split(",")
54
55         if not wclass[0]=="BOS/EOS":
56
57             sentence.append(word)
58
59             genkei.append(wclass[6])
60
61             katuyo.append(wclass[5])
62
63             hinshi.append(wclass[0])
64
65         node = node.next
66
67     return sentence,genkei,katuyo,hinshi

```

図 2 翻訳する文を形態素解析する関数

```

140 if sonkei==1:
141     #print("尊敬にする")
142     num=-1
143     for a in genkei:
144         num+=1
145         if a in sonkeigo:
146             katuyokei=katuyo[num]
147             keigo=sonkeigo[a]
148             hinsi=hinshi[num]
149
150             result=henkan(keigo,hinsi,katuyokei,num,sentence)
151
152             sentence[num]=result
153
154
155
156 if kenjo==1:
157     #print("謙讓にする")
158     num=-1
159     for a in genkei:
160         num+=1
161         if a in kenjogo:
162             katuyokei=katuyo[num]
163             keigo=kenjogo[a]
164             hinsi=hinshi[num]
165
166             result=henkan(keigo,hinsi,katuyokei,num,sentence)
167
168             sentence[num]=result
169
170 semifinal="".join(sentence)

```

図 3 単語を尊敬語または謙讓語に変換する関数

```

69 def kensaku(part, word, form):
70     if part == "動詞":
71         file_name = "Verb.csv"
72     elif part == "助動詞":
73         file_name = "Auxil.csv"
74     else :
75         sys.exit()
76
77     f = open(file_name,'r')
78     dataReader = csv.reader(f)
79     for row in dataReader:
80         if word == row[10]:
81             if form in row[9]:
82                 return row[0]

```

図 4 指定された活用形の動詞を検索する関数

```

85 def henkan(keigo,hinsi,katuyokei,num,sentence):
86
87     result=mecab(keigo)
88     tango=result[0]
89     a=tango[-1]
90
91     new_result=kensaku(hinsi,a,katuyokei)
92
93     if len(sentence)>num+1:
94         if sentence[num+1]=="う":
95             if "ウ接続" in katuyokei:
96                 if new_result==None:
97                     katuyokei=katuyokei.strip("ウ接続")
98
99             else:
100                katuyokei=katuyokei.strip("形")+"ウ接続"
101                if kensaku(hinsi,a,katuyokei)==None:
102                    katuyokei=katuyokei.strip("ウ接続")
103
104                if sentence[num+1]=="た" or sentence[num+1]=="て":
105                    if "タ接続" in katuyokei:
106                        if new_result==None:
107                            katuyokei=katuyokei.strip("タ接続")
108                    else:
109                        katuyokei=katuyokei.strip("形")+"タ接続"
110
111                    if kensaku(hinsi,a,katuyokei)==None:
112                        katuyokei=katuyokei.strip("タ接続")
113
114
115     final_result=(kensaku(hinsi,a,katuyokei))
116
117     if not final_result==None:
118         tango[-1]=final_result
119
120     semifinal=""
121     for a in tango:
122         semifinal+=a
123
124     return semifinal

```

図5 活用形をタ接続もしくはウ接続に変換する関数

```

167 semifinal="" .join(sentence)
168 print("↓")
169 print(semifinal)
170
171 #ここから丁寧語
172 result=mecab(semifinal)
173 sentence=result[0]
174 genkei=result[1]
175 katuyo=result[2]
176 hinshi=result[3]

```

図6 (3)までの処理を行った文を再び形態素解析するプログラム

```

178 num=-1
179 for a in genkei:
180     num+=1
181     if a=="だ":
182         if hinshi[num]=="助動詞":
183
184             katuyokei=katuyo[num]
185
186             result=henkan("です","助動詞",katuyokei,num,sentence)
187             sentence[num]=result
188

```

図7 助動詞の「だ」を「です」に変換するプログラム

```

190 num=len(genkei)
191 for a in reversed(genkei):
192     num-=1
193     if a=="ます":
194         num-=1
195
196     if genkei[num] in jodoshi:
197         result=kensaku("動詞",a,"連用形")
198         sentence[num]=result
199
200     katuyokei=katuyo[num]
201     result=henkan("ます","助動詞",katuyokei,num,sentence)
202     sentence.insert(num+1,result)
203     break
204
205     if hinshi[num]=="動詞":
206         result=kensaku("動詞",a,"連用形")
207         sentence[num]=result
208
209     katuyokei=katuyo[num]
210     result=henkan("ます","助動詞",katuyokei,num,sentence)
211     sentence.insert(num+1,result)
212     break

```

図8 助動詞「ます」を文末に付け加えるプログラム

```

36 jodoshi=["せる","させる","れる","られる","がる"]

```

図9 動詞型の活用をする助動詞のリスト


```

5 sonkeigo = ['なさる': 'する',
6             'いらっしゃる': '来る',
7             'ご覧になる': '見る',
8             'お耳に入る': '聞く',
9             'おっしゃる': '言う',
10            'お思いになる': '思う',
11            'ご存じです': '知る',
12            'くださる': 'あげる',
13            'お受けになる': 'もらう',
14            '召し上がる': '食べる',
15            'お亡くなりになる': '死ぬ',
16            'お気に召す': '気に入る']
17
18 kenjogo = ['おる': 'いる',
19            '参る': '来る',
20            '拝見する': '見る',
21            'ご覧に入れる': '見せる',
22            '伺う': '聞く',
23            '申し上げる': '言う',
24            '存じ上げる': '知る',
25            '差し上げる': 'あげる',
26            '頂戴する': '貰う',
27            'いただく': '食べる',
28            '亡くなる': '死ぬ',
29            'お目にかかる': '会う',
30            '拝借する': '借りる']

```

図 10 尊敬語と謙譲語の表

```

206 result=mecab(semifinal)
207
208 sentence=result[0]
209 genkei=result[1]
210 katuyo=result[2]
211 hinshi=result[3]
212
213 for a in reversed(genkei):
214
215     if a=="ます":
216         num=genkei.index(a)
217         katuyokei=katuyo[num]
218         tango=genkei[num-1]
219         hinsi=hinshi[num-1]
220         result=henkan(tango,hinsi,katuyokei,num,sentence)
221         sentence[num-1]=result
222         sentence.pop(num)
223
224     if a=="です":
225         num=genkei.index(a)
226         katuyokei=katuyo[num]
227         result=henkan("だ","助動詞",katuyokei,num,sentence)
228         sentence[num]=result

```

図 11 丁寧語を省くプログラム