

● 例文を基にした文章の生成 ●

文章は、いくつかの単語が順に並ぶことで生み出されます。

猿 も 木 から 落ちる 。

今回は次のような例文をもとに、日本語として意味の通る別の文章を自動的に生成することを目指します。

```
富山県はブラックラーメンが有名です。石川県は人が減少中。その高校はブラック校則が有名ですか。
```

ステップ1 「分かち書き」による例文からの単語の抽出

日本語は英語と異なり、語句の間をスペースで区切ったりしないので、どこまでが単語なのかがわかりにくい。

```
すもももももものうち
```

↓ これを、単語ごとにスペースで区切ることを「分かち書き」という。

```
すもも も もも も もものうち
```

Python では、**Mecab** や **janome** などのモジュールをインストールすることによって日本語の文章の分かち書きを（多少の不完全さはあるにせよ）自動で行うことができる。（詳細は後ほど。）

例文を **janome** をつかって分かち書きしてみます。

```
富山 県 は ブラック ラーメン が 有名 です 。
石川 県 は 人 が 減少 中 。
その 高校 は ブラック 校則 が 有名 です か 。
```

すると、3つの例文は、17種類¹の単語で構成されていることがわかります。（**gosu = 17**）

```
富山 県 は ブラック ラーメン が 有名 です 。 石川 人 減少 中 その 高校 校則 か
```

これらの語からいくつか選び出して適当に並べてみることで、新しい文章の生成を試みる。と…

```
ブラック が 県 校則 人 。 その
```

```
有名 石川 は です ラーメン。
```

…これだと「カタコトの日本語」とも言えないような文章になってしまいますね。さて。どうしたものか。

【実習①】 Anaconda Prompt を用いて、**janome** をインストールします。

```
> pip install janome
```

すると、spyder で次のように書くことで分かち書きができる。

```
1 # -*- coding: utf-8 -*-
2 from janome.tokenizer import Tokenizer
3 t = Tokenizer()
4 txt_data = '富山県はブラックラーメンが有名です。石川県は人が減少中。その高校はブラック
  校則が有名ですか。'
5 print("*****原文*****")
6 print(txt_data)
7 print("*****分かち書き*****")
8 Word_list = list(t.tokenize(txt_data,wakati=True))
9 print(Word_list)
```

ステップ2 マルコフ連鎖についての基礎知識

話は変わりますが、数学に「マルコフ連鎖」という考え方があります。名前からして難しそうですが、実は単純な話で、次に起こる出来事が、「今の状態」から確率分岐する（過去によらない）という考え方です。

数学で習ったさいころの独立試行に近いものがあります。例えば、1つのさいころを繰り返し投げるような「実験」を考えてみると、「過去に1の目がたくさん出たから、次にさいころを振ると1の目は出にくい」ということはありません。さいころが特に細工のない普通のもので、1が出る確率が $1/6$ なら、そのさいころで過去1がたくさん出たか否かに関わらず、次にさいころをふれば $1/6$ の確率で1が出ます。

今この「実験」を、[今の状態:[次にとりうる状態]]という記号で、次のように表現することにします。

1:	[1, 2, 3, 4, 5, 6]
2:	[1, 2, 3, 4, 5, 6]
3:	[1, 2, 3, 4, 5, 6]
4:	[1, 2, 3, 4, 5, 6]
5:	[1, 2, 3, 4, 5, 6]
6:	[1, 2, 3, 4, 5, 6]

さいころの場合は1～6の6つの状態をとり、各数字が出る割合が等しいので上記のようになります。

ちょっと一般化してみましよう。たとえば、A,B,Cの3つの状態があり、いまAの状態にあるものが状態Bに移る確率が $1/4$ 、状態Cに移る確率が $3/4$ であるときは、次のように表現することができます。

A:	[B,C,C,C]
----	-----------

つまり、[B,C,C,C]のなかからランダムでどれか一つが選ばれることで、どの遷移が起こるのかが確率的に選択される（状態Cへの遷移は、Bへの遷移の3倍起こりやすい。）ように表現します。

問1：いま状態Cにあり、そこから状態Aに遷移する可能性が $2/5$ 、Bに遷移する可能性が $3/5$ である。それを記号であらわすとどのように書けるか。空欄を補いましょう。[C [

--

]]

問2：

[A:[B,C,C,C]]

[B:[A,C]]

[C:[A,B,B]]

であるとき、

いま状態Aにあるものが状態A→状態B→状態A→状態Cの順に遷移する確率はいくらか。（分数で回答可。）

ステップ3 例文からのモデルの作成と、マルコフ連鎖による作文

例文をうまく解釈し、そこにマルコフ連鎖の考え方を適用するだけで、それなりに意味の通る日本語の文章を生成することができます。

つまり、先に挙げた17種類の語それぞれの後に、どの語が続く可能性があるかを例文から整理します。この作業を「modelの作成」と呼ぶことにします。(すごく大事な作業です。)

たとえば、例文を分析してみると、「ブラック」という語の後には「校則」または「ラーメン」が続きます。「は」という言葉のあとには1/3の確率で「人」、2/3の確率で「ブラック」が続きます。

17種類の語の後に続く言葉を整理した(つまり、3つの例文から作られた)モデルは次のようになります。(BoSは、Beginning of Sentence:文章の開始を意味します。)

作成されたモデル	
'BoS': ['富山', '石川', 'その'],	'石川': ['県'],
'富山': ['県'],	'人': ['が'],
'県': ['は', 'は'],	'減少': ['中'],
'は': ['ブラック', '人', 'ブラック'],	'中': ['。'],
'ブラック': ['ラーメン', '校則'],	'その': ['高校'],
'ラーメン': ['が'],	'高校': ['は'],
'が': [(a)],	'校則': ['が'],
'有名': ['です', 'です'],	'か': ['。']
'です': ['。', 'か'],	

つまり、「富山」、「石川」、「その」の3語から1語を選ぶところから文章生成がスタートし、もし「富山」が選ばれたらその後には必ず「県」が続く(「現象」とか「校則」は富山の後には選択されない。)
「県」のあとは必ず(2/2の確率で)「は」、「は」という言葉のあとには先ほど示したように1/3の確率で「人」、2/3の確率で「ブラック」...という風にマルコフ連鎖を続けていくことで、意味のある語が選ばれて、「。」が選択されたら終了とすることで文章生成(作文)を行うことが可能...なのではないかということになります。

問3 : (a) にあてはまる語を答えなさい。ただし、同じ語が複数回入ることがある。

問4 : このモデルを使って文章を生成した場合に、3つの例文以外で生成される可能性のある文章を1つ書きなさい。

問5 : 「石川県はブラック校則が減少中です。」という文章が生成される確率はいくらか。(分数で回答可。)

分ち書きをした1~9行目までのプログラミングに続く形で、たとえば次のようにコードを作成することで、マルコフ連鎖による文章生成ができます。

```
10 #*****重複のない word_list の作成*****
11 Word_list_single = []
12 for word in Word_list:
13     if word not in Word_list_single:
14         Word_list_single.append(word)
15 gosu = len(Word_list_single)#重複しない単語の数
16 print(Word_list_single,gosu,"語")
17 #*****model の生成*****
18 Word_list.insert(0, ".")
19 Model = [[]for i in range(gosu)]
20 for i in range(gosu):
21     Model[i].append(Word_list_single[i])
22 print("*****model の生成*****")
23 print("0 番目の要素に続く可能性のある語を 1 番目、2 番目...の要素に追記します")
24 #print(model)
25 for j in range(gosu):
26     word = Word_list_single[j]
27     for i in range(len(Word_list)-1):
28         if Word_list[i] == word:
29             Model[j].append(Word_list[i+1])
30 print(Model)
31 #*****マルコフ連鎖による作文*****
32 print("*****マルコフ連鎖*****")
33 import random
34 sentence = ""
35 select_word = ".">#BoS (文章の始まり) の代わりに読点を使います。
36 for j in range(40):#最大で 40 語の文章を生成
37     for i in range(gosu):
38         if Model[i][0] == select_word:
39             num = random.randint(1,len(Model[i])-1)
40             select_word = Model[i][num]#後に続く語をランダムに選択
41             break
42     sentence = sentence + select_word
43     print(Model[i][1:len(Model[i])]) #作文過程を表示する記述。省略可。
44     print(sentence) #作文過程を表示する記述。省略可。
45     if select_word == ".":#「。」が来ると作文を停止
46         break
```

【実習②】 実行できることが確認出来たら、4行目の `txt_data` をさまざまな文章に変えて実行してみましょう。作家の作品の一部とすることで、その作家風の作文をする…なんてこともできるかもしれませんね。