

付録A スコア算出プログラムについて

3. 1 (3) 項で述べたスコアを算出することができるプログラムを図A. 1に示す。なお、図A. 1において、左側の列は説明のための行番号であり、プログラムの一部ではないことに注意する。本プログラムは、近年最もよく用いられるコンピュータ言語の1つであるPythonで記述されている。また、本プログラムはDWPI抄録を処理するものである。その他のデータを処理する場合については、一部を書き換える必要がある。書き換えるべき内容については詳細を後述する。

以下、本プログラムを実行するための環境整備と、プログラムの使い方、最後に、DWPI以外のデータへの対応等を含めプログラムの詳細な解説を行う。

A. 1 環境整備

図A. 1のプログラムを実行するためには、以下の手順で、Pythonプログラムを実行する環境と各種ライブラリを事前にインストールする必要がある。

1. anaconda のインストール
2. nltk の形態素解析機能をインストール
3. nagisa のインストール

インストール手順の詳細を以下に示す。

(1) anaconda のインストール

anaconda は、Python でデータサイエンスや機械学習、自然言語処理を行うために必要なソフトウェアが一括でインストールできるソフトウェア（ディストリビューション）である。

anaconda のWebサイト¹⁾からこれをダウンロ

ードし、インストールする。ここで、Python 3. X版と Python 2. 7版が選択できるが、3. X版(2019年5月26日時点ではXは7)を選択する。

Anaconda の使い方については説明を省略する。

(2) nltk のインストール

次に英語の形態素解析（文章を単語に切り分けて、単語の品詞を推定する処理）や単語の組み合わせを抽出する処理を実行するための機能を追加する。コンソール(anaconda コンソール)を開き図A. 2のaのコマンドを実行する。

(3) nagisa のインストール

最後に、日本語の形態素解析を実行するためのライブラリを追加する。先と同様にコンソールを開き図A. 2のbのコマンドを実行する。ここで、Windows 環境では、エラーが発生することがある。その場合、Build Tools for Visual Studio 2019²⁾をインストールする。

以上の手順でプログラムを実行するための準備が整う。

(4) 補足：プロキシが存在する場合

インターネット接続時にプロキシサーバを通す必要がある場合、図A. 2のコマンドではインストールできない場合がある。その場合は、図A. 3のコマンドを用いてインストールを行う。なお、b' において、パスワードに記号が含まれる場合は URI エンコーディング（記号を%nn (nは数字) に置き換える）する。

A. 2 使い方

上記環境が整備されたコンピュータならば、
図A. 1のプログラムを実行可能である。具体的には、プログラムを任意のディレクトリに保存後、環境整備時と同様にコンソールを開き、
図A. 2のcのコマンドを実行する。

コマンドにおいて、データファイルとは、処理対象とする csv ファイル名である。処理対象の csv ファイルは、1行目にヘッダ、2行目以降に1行1公報の形でテキストが格納された csv である。入力ファイルの例を図A. 4に示す。図A. 4のように、本プログラムでは、処理対象以外の列（公報番号やタイトルなど）の列が含まれていてもよい。すなわち、特許データベースからエクスポートした csv ファイルをそのまま入力可能である。なお、図A. 1のプログラムでは、入力 csv ファイルの文字コードとして Shift_JIS（正確には Windows が用いている Shift_JIS）を想定している。

また、出力も csv 形式であり、各行にキーワードとそのスコアが出力される。

A. 2 プログラム解説

次に、DWPI 以外の形式のデータへの対応等を行う際に重要な部分についてプログラムの詳細を解説する。図A. 1のプログラムは以下のよう構成される。

1. 必要なライブラリの準備 (1~6 行目)
2. 設定 (8~19 行目)
3. キーワードのカウント (24~47 行目)
4. スコアの算出 (49~54 行目)
5. 出力 (56~58 行目)

上記構成のうち、ポイントとなる部分の詳細を述べる。

- (1) 各種設定 (8~19 行目)

8~19 行目には様々な設定（定数）が記述されている。各設定項目を表A. 1に示す。表A. 1の各項目を書き換えることで、DWPI 抄録以外のデータへ対応が可能となる。

表A. 1において、TARGET とはスコアの対象とする列の名前であり、COMPARISON は、TARGET と比較される列の名前である。図A. 1の例では、「抄録 - DWPI 用途」と、「抄録 - DWPI 新規性」または「抄録 - DWPI 優位性」を比較して、用途らしさを表すスコアを算出する。

また、N_GRAM はキーワードを構成する語の最大数を表す。

また、POSTAG はキーワードを構成する語の品詞を表し、STOP_WORD はスコア算出の対象外とする語を表す。ここで、POSTAG と STOP_WORD については、正規表現で記載されている。詳細は省略するが、正規表現とは、文字列パターンの記述方法である。図A. 1の場合、POSTAG は CD、NN、または FW で始まる品詞を表し、STOP_WORD は、数字のみで構成されるキーワード、または「e.g.」、「i.e.」、「etc.」、「claimed」のいずれかを含むキーワードを除外している。なお、POSTAG と STOP_WORD の書き換え例では、re.compile を省略しており、実際は re.compile(?)の「?」部分を表の文字列に置き換える。

- (2) キーワードのカウント (24~47 行目)

プログラムの大半を占める 24 行目から 47 行目は、各列に現れるキーワードの頻度をカウントする処理である。25~30 行目は「csv ファイルの各行の各処理対象の列（スコア算出対象+比較対象）に対して処理を行う」ことを意味し、31 行目~32 行目で言語の判定を行っている。本プログラムでは簡便に「アルファベット、数字、カンマ、ピリオドがテキストの 9 割以上を占めたら英語であり、それ以外は日本語」という判

定方法を採用している。その後、事前にインストールした各ライブラリを利用してテキストの形態素解析処理を行っている（英語は 34 行目、日本語は 37, 38 行目）。形態素解析を行うことで、文字列を語の配列に変換している。最後に、語の配列から 1~N_GRAM 個で構成されるキーワードの候補を抽出し（41 行目）、その候補がすべて対象品詞（POSTAG）に該当するかをチェックし（42 行目）、そのチェックを通過した候補をキーワードとして追加している。なお、本プログラムでは、キーワードをすべて小文字に変換し、かつ、語の境界の左右の文字がアルファベットであるときのみ境界にスペースを挿入している。

なお、先にも述べたように、本プログラムでは Shift_JIS の csv ファイルを入力として想定している。もし、Shift_JIS 以外の csv ファイルを扱いたい場合は、25 行目の cp932 の部分を csv ファイルの文字コードに書き換える。

（2）スコアの算出（49~54 行目）

キーワードの頻度のカウントが終了した後、49 行目から 54 行目でキーワードのスコアを算出している。この際、長さが 1 文字のキーワードや、キーワードが STOP_WORD とマッチした場合は、そのキーワードをスコア算出の対象から除外している。

最後に、スコアが高い順にキーワードを並び替え、キーワードとスコアを出力するのが図 A. 1 のプログラムの流れとなる。

注 記

1) Anaconda

<https://www.anaconda.com/distribution/>

（参照日 2019 年 5 月 26 日）

2) Build Tools for Visual Studio 2019

<https://visualstudio.microsoft.com/ja/>

[downloads/](#)

（参照日 2019 年 5 月 26 日）

```

1 import sys
2 import csv
3 import re
4 from collections import Counter
5 import nltk
6 import nagisa
7
8 # 計測対象の列名
9 TARGET = '抄録 - DWPI 用途'
10 # 比較対象の列名
11 COMPARISION = ['抄録 - DWPI 新規性', '抄録 - DWPI 優位性']
12 # 処理対象の列名
13 ALL_COLUMNS = [TARGET] + COMPARISION
14 # n-gram の最大値
15 N_GRAM = 3
16 # n-gram の対象とする品詞
17 POSTAG = re.compile(r'^(CD|NN|FW)')
18 # 除外語
19 STOP_WORD = re.compile(r'^\d+$|e\.g|i\.e|etc|claimed')
20
21 # 語の出現頻度カウンター
22 df = {col: Counter() for col in ALL_COLUMNS}
23
24 # キーワードの出現頻度のカウント
25 with open(sys.argv[1], 'r', encoding = 'cp932') as csv_file:
26     for row in csv.DictReader(csv_file):
27         for col in ALL_COLUMNS:
28             if not (col in row) or len(row[col]) == 0:
29                 continue
30             txt = row[col]
31             a_txt = re.sub(r'^\w\s\.\.]', '', txt, flags = re.ASCII)
32             if len(a_txt) / len(txt) >= 0.9:
33                 # 英語の形態素解析
34                 words = nltk.pos_tag(nltk.word_tokenize(txt))
35             else:
36                 # 日本語の形態素解析(nltkの形式に合わせる)
37                 words = nagisa.tagging(txt)
38                 words = list(zip(words.words, words.postags))
39             # 単語 n-gram の取得
40             ngrams = set()
41             for ng in nltk.everygrams(words, 1, min([len(words), N_GRAM])):
42                 if all((POSTAG.search(w[1]) for w in ng)):
43                     # 文字列化
44                     ng_str = ' '.join([w[0] for w in ng]).lower()
45                     ng_str = re.sub(r'([\^a-z])\s([\^a-z])', r'\1\2', ng_str)
46                     ngrams.add(ng_str)
47             df[col].update(ngrams)
48
49 # キーワードのスコア計算
50 r_df = {}
51 for w, f in df[TARGET].items():
52     if len(w) <= 1 or STOP_WORD.search(w):
53         continue
54     r_df[w] = f / (max([df[c][w] for c in COMPARISION]) + 1)
55
56 # 出力
57 for w, f in sorted(r_df.items(), key=lambda x: -x[1]):
58     w = w.replace("'", '""')
60     print(f'"{w}", {f}')

```

図A. 1 スコア算出プログラム (DWPI 版)

```

a. nltk の形態素解析機能のインストール
> python -c "import nltk; nltk.download(['punkt', 'averaged_perceptron_tagger'])"

b. nagisa のインストール
> pip install nagisa

c. プログラムの実行方法
> python analyze.py データファイル > 出力ファイル

```

図A. 2 各種コマンド

```

a' . nltk の形態素解析機能のインストール (改行せずに入力)
> python -c "import nltk; nltk.set_proxy('http://ユーザ名:パスワード@サーバ:ポート');
            nltk.download(['punkt', 'averaged_perceptron_tagger'])"

b' . nagisa のインストール
> pip install nagisa --proxy="http://ユーザ名:パスワード@サーバ:ポート"

```

図A. 3 各種コマンド (プロキシがある場合)

	A	B	C	D	E	F
1	公報番号	タイトル - DWPI	抄録 - DWPI 新規性	抄録 - DWPI 用途	抄録 - DWPI 優位性	...
2	JP2010527805A	Bernoulli gripper ...	The gripper has a clamping ...	Used for contactless admission ...	The damper runs about the
3	JP03138660U	Bipedal movement	A lower drive device moves ...	Bipedal movement robot.	The fall of robot during curve
4	JP04088332B1	Industrial robot ...	The robot has impact sensor ...	Industrial robot for welding.	The control characteristics
5	JP03140558U	Rehabilitation...	The apparatus has a moving ...	Rehabilitation apparatus for ...	The apparatus can be stopped
6	⋮	⋮	⋮	⋮	⋮	⋮

図A. 4 csv ファイルの例 (⋮は省略を表す)

表A. 1 設定 (定数)

行数	定数名	意味	書き換え例
9	TARGET	スコアを計算したい列	'課題'
11	COMPARISION	TARGET と比較する列 (複数可)	['技術分野', '手段']
15	N_GRAM	キーワードを構成する語の最大数	4
17	POSTAG	キーワードを構成する語の品詞	r'^(接頭辞 名詞 接尾辞)'
19	STOP_WORD	スコア算出の対象外とする語	r'特許 出願 特開 特表 '