

ソフトウェア工学最終レポート

201811319 永崎遼太¹

2020.12.28

¹筑波大学情報学群情報科学類 3 年

目 次

1	対象ソフトウェアプロダクト (プロセス) の説明	2
2	対象ソフトウェアプロダクト (プロセス) の問題	4
3	プロダクト (プロセス) の改善	5
4	改善のまとめと今後の展望	6

1 対象ソフトウェアプロダクト (プロセス) の説明

ソフトウェアサイエンス専攻実験「Web 文書分類」にて作成した、ナイーブベイズ (NB) 法により入力文書のクラスを決定するソフトウェアプロダクトを対象とする。

NB 法では、あらかじめ単語の生起確率を求めておく訓練ステップと、求めた生起確率を利用して、入力文書のクラスを決定する分類ステップに分けて考えることができる。

手続き

訓練ステップ 教師文書を利用して各クラス、各単語について $P(w_j, c_i)$ を求める

分類ステップ 入力文書 x に対して、分類クラスを決定する。

1. 入力文書の各単語に対して $N(w_j, x)$ を求める
2. 入力文書の各単語に対して、訓練フェーズの結果を参照して $P(w_j, c_j)$ を得る
3. $N(w_j, x)$ と $P(w_j | c_i)$ の情報から分類クラスを決定する

訓練ステップ、分類ステップをそれぞれ以下の nb_train.py、nb_classification.py で実装している。

```
1 import MeCab as mc
2 from collections import Counter
3 import ast
4 import math
5 ##### 訓練ステップ
6 ##### 教師文書を利用して、各クラス、各単語について  $p(w_j, c_i)$  を求める
7 # ----- # 関数・の定義 delta
8 # 特徴抽出結果を受け取り辞書に変換
9 def str_dic(x):
10     d = x
11     d1 = '{'+d.replace(':', ':').replace(' ', ',') + '}'
12     d2 = ast.literal_eval(d1)
13     return d2
14 # smoothing
15 delta = 1
16 # ----- # クラスについて cleanerP(wj | ci) を求める
17 # 教師文書を読み込み
18 tid = 0
19 space = {}
20 for line in open('train_cleaner_tf.txt', 'r'):
21     line = line.replace(" \n", "\n")
22     line_ = str_dic(line)
23     space[tid] = line_
24     tid += 1
25 # 各単語の生起頻度を登録
26 pwj_cleaner = {}
27 for id in space.keys():
28     for val in space[id].keys():
29         if (val not in pwj_cleaner):
30             pwj_cleaner[val] = space[id][val]
```

```

31         else:
32             pwj_cleaner[val] = pwj_cleaner[val]+space[id][val]
33 # クラスの全単語数を計算cleaner
34 all_words_cleaner = 0
35 for id in pwj_cleaner.keys():
36     tmp = pwj_cleaner[id] + delta
37     all_words_cleaner = all_words_cleaner + tmp
38 # 最終的なP(wj|ci)を計算
39 for id in pwj_cleaner.keys():
40     tmp = pwj_cleaner[id] + delta
41     pwj_cleaner[id] = tmp / (all_words_cleaner)
42 # P(wj|ci)をファイルに出力
43 with open('pwj_cleaner.txt', 'w') as f:
44     print(pwj_cleaner, file=f)
45 # ----- # クラスについてmp3playerP(wj|ci)を求める
46 # 教師文書を読み込み
47 tid = 0
48 space = {}
49 for line in open('train.mp3player_tf.txt', 'r'):
50     line = line.replace("\n", "\n")
51     line_ = str_dic(line)
52     space[tid] = line_
53     tid += 1
54 # 各単語の生起頻度を登録
55 pwj_mp3player = {}
56 for id in space.keys():
57     for val in space[id].keys():
58         if (val not in pwj_mp3player):
59             pwj_mp3player[val] = space[id][val]
60         else:
61             pwj_mp3player[val]=pwj_mp3player[val]+space[id][val]
62 # クラスの全単語数を計算mp3player
63 all_words_mp3player = 0
64 for id in pwj_mp3player.keys():
65     tmp = pwj_mp3player[id] + delta
66     all_words_mp3player = all_words_mp3player + tmp
67 # 最終的なP(wj|ci)を計算
68 for id in pwj_mp3player.keys():
69     tmp = pwj_mp3player[id] + delta
70     pwj_mp3player[id] = tmp / (all_words_mp3player)
71 # P(wj|ci)をファイルに出力
72 with open('pwj_mp3player.txt', 'w') as f:
73     print(pwj_mp3player, file=f)

```

Listing 1: nb_train.py

```

1 import ast
2 import nb_train
3 import math
4 ##### 分類ステップ
5 ##### 入力文書に対して、分類クラスを決定するx
6 # ----- # 関数の定義
7 # 特徴抽出結果を受け取り辞書に変換
8 def str_dic(x):
9     d = x
10    d1 = '{'+x.replace(':', '').replace(' ', ',')+'}'
11    d2 = ast.literal_eval(d1)
12    return d2
13 # ----- # 1
14 # 入力文書の各単語に対して、N(w_j, x)を求める
15 tid = 0
16 nwjx = {}

```

```

17 for line in open('test_tf.txt', 'r'):
18     line = line.replace(" \n", "\n")
19     line_ = str_dic(line)
20     nwjx[tid] = line_
21     tid += 1
22 # ----- # 2
23 # 入力文書の各単語に対して、訓練フェースの結果を参照してp(w_j, c_i)を得る
24 # 訓練フェースの結果を参照
25 pwj_cleaner = nb_train.pwj_cleaner
26 pwj_mp3player = nb_train.pwj_mp3player
27 all_words_cleaner = nb_train.all_words_cleaner
28 all_words_mp3player = nb_train.all_words_mp3player
29 # クラスに関してcleanerP(wj|ci) を求める
30 for id in nwjx.keys():
31     for val in nwjx[id].keys():
32         if (val not in pwj_cleaner):
33             pwj_cleaner[val] = 1 / all_words_cleaner
34 # クラスに関してmp3playerP(wj|ci) を求める
35 for id in nwjx.keys():
36     for val in nwjx[id].keys():
37         if (val not in pwj_mp3player):
38             pwj_mp3player[val] = 1 / all_words_mp3player
39 # ----- # 3
40 # N(w_j, x)とp(w_j, c_i)の情報から分類クラスを決定する
41 for id in nwjx.keys():
42     val_cleaner = 0
43     val_mp3player = 0
44     for val in nwjx[id].keys():
45         tmp = nwjx[id][val]* math.log(pwj_cleaner[val])
46         val_cleaner = val_cleaner + tmp
47     for val in nwjx[id].keys():
48         tmp = nwjx[id][val]* math.log(pwj_mp3player[val])
49         val_mp3player = val_mp3player + tmp
50     output_class='cleaner'if val_cleaner > val_mp3player else '
51     mp3player'
52     print( output_class )

```

Listing 2: nb_classification.py

train.cleaner_tf.txt、train.mp3player_tf.txt という文書を教師文書とし、test_tf.txt というファイルを分類対象のファイルとしている (各ファイルの内容についての説明は割愛)。

2 対象ソフトウェアプロダクト (プロセス) の問題

作成したプログラムでは、まず nb_train.py を実行し、その結果を参照して nb_classification.py を実行するという手順で NB 法を実行する。

交差検定法により実装した NB 法の性能を評価しようとした時、分割されたデータに対して NB 法を繰り返し実行する必要性が生じる。作成したプログラムのままでは、NB 法を繰り返し実行するという処理が困難である。

また、文書を読み込むという機能を逐一プログラムを書いて実行していたり、それぞれで str_df() 関数を作成していたりなど、冗長な部分が見られる。

3 プロダクト(プロセス)の改善

両方のプログラムで必要とされる機能を取り出し、関数としてまとめた新たなプログラム my_function.py を作成した。

```
1 import ast
2 num_split = 5
3 # ----- #
4 # 特徴抽出結果を受け取り辞書に変換
5 def str_dic(x):
6     d = x
7     d1 = '{' + '}' + d.replace(':', ':').replace(' ', ',').
8         replace('\"': '\"', '\"': '\"') + '}'
9     d2 = ast.literal_eval(d1)
10    return d2
11 # ----- #
12 # 特徴抽出結果を読み込むtf
13 def read_tf(text):
14     tid = 0
15     space = {}
16     for line in text:
17         line = line.replace(" \n", "\n")
18         line_ = str_dic(line)
19         space[tid] = line_
20         tid += 1
21     return space
```

Listing 3: my_function.py

作成したプログラム内の機能を抽象化して関数として再実装し、呼び出せるようにする。それぞれを改善した結果が以下に示す nb_train2.py、nb_classification2.py である。

```
1 import ast
2 import my_function
3 ##### 訓練ステップ
4 ##### 教師文書を利用して、各クラス、各単語についてp(w_j, c_i)を求める
5 # smoothing
6 delta = 1
7 # ----- # P(w_j|c_i)を求める
8 # あるクラスの教師文書の各単語の生起頻度を計算
9 def train_count_pwj(data):
10     pwj_class = {}
11     for id in data.keys():
12         for val in data[id].keys():
13             if (val not in pwj_class):
14                 pwj_class[val] = data[id][val]
15             else:
16                 pwj_class[val] = pwj_class[val]+data[id][val]
17     return pwj_class
18 # あるクラスの教師文書の全単語数を計算
19 def train_count_word(pwj_class):
20     all_words_class = 0
21     for id in pwj_class.keys():
22         tmp = pwj_class[id] + delta
23         all_words_class = all_words_class + tmp
24     return all_words_class
25 # 最終的なP(pwj|ci)を計算
26 def train_pwj_class(pwj_class, all_words_class):
27     for id in pwj_class.keys():
```

```

28     tmp = pwj_class[id] + delta
29     pwj_class[id] = tmp / (all_words_class)
30     return pwj_class

```

Listing 4: nb_train2.py

```

1  import ast
2  import math
3  import my_function
4  ##### 分類ステップ
5  ##### 入力文書に対して、分類クラスを決定するx
6  # ----- # 1
7  # 入力文書の各単語に対して、 $N(w_j, x)$ を求める
8  # ----- # 2
9  # 入力文書の各単語に対して、訓練フェースの結果を参照して $p(w_j, c_i)$ を得る
10 # あるクラスに関して $P(w_j|c_i)$ を求める
11 def pwj_class(test_data, pwj_class, all_words_class):
12     for id in test_data.keys():
13         for val in test_data[id].keys():
14             if (val not in pwj_class):
15                 pwj_class[val] = 1 / all_words_class
16     return pwj_class
17 # ----- # 3
18 #  $N(w_j, x)$ と $p(w_j, c_i)$ の情報から分類クラスを決定する
19 def decide_class(test_data, pwj_class1, pwj_class2):
20     output_class = {}
21     tid = 1
22     for id in test_data.keys():
23         val_class1 = 0
24         val_class2 = 0
25         for val in test_data[id].keys():
26             tmp = test_data[id][val] * math.log(pwj_class1[val])
27             val_class1 = val_class1 + tmp
28         for val in test_data[id].keys():
29             tmp = test_data[id][val] * math.log(pwj_class2[val])
30             val_class2 = val_class2 + tmp
31         output_class_ = 'pos' if val_class1 > val_class2 else 'neg'
32         output_class[tid] = output_class_
33         tid += 1
34     return output_class

```

Listing 5: nb_classification2.py

4 改善のまとめと今後の展望

機能を抽象化して関数として再実装することで、機能を繰り返し実行することが容易になり、またプログラムの冗長性を改善することが可能になった。今回作成したプログラムを用いて交差検定法を実行するプログラムを実際に作成すると、正しく動作をすることが確認できたので、改善を正しく行うことができたことが確かめられた。