

別添資料

本統括・分担研究報告書全体に係る別添資料を示す。
順に、以下の構成となっている。

- 備考欄前処理プログラムソース（主要な部分抜粋）
 - applyregexp.py (以下の処理の実行)
 - ✧ regexp14.py ("死亡の原因" 用)
 - ✧ regexp16.py ("外因死の追加事項" 用)
 - ✧ regexp18.py ("その他特に付言すべきことがら" 用)
 - ✧ regexpOther.py (上記以外の全て用)
- 機械学習用データセット作成プログラムソース
 - TFIDF
 - LSI
 - Word2Vec
 - Doc2Vec (PV-DM / PV-DBOW)
- 分類器学習プログラムソース
 - fit_and_predict_xgboost.py

備考欄前処理プログラム (apply_regexp.py)

```
1 from my_util import exec_time
2
3 from datetime import datetime
4 import re
5 import regexp14
6 import regexp16
7 import regexp18
8 import regexpOther
9 import copy
10
11 @exec_time.printer
12 def main():
13     regs14 = regexp14.get()
14     regs16 = regexp16.get()
15     regs18 = regexp18.get()
16     regsOther = regexpOther.get()
17     max_cnt = 0
18     fis = open('../01_new_shibo_join_kv.csv', 'r')
19     fos = open('../02_new_shibo_join_kv2' + ('' if max_cnt <= 0 else '_' +
20     str(max_cnt)) + '.csv', 'w')
21     f_removed = open('../02_new_shibo_join_removed_words' + ('' if max_cnt <= 0
22     else '_' + str(max_cnt)) + '.tsv', 'w')
23
24     cnt = 0
25     for sbh in sbh_data_generator(fis):
26         cnt += 1
27         if cnt == max_cnt: break
28         rowid = sbh['rowid']
29         values = sbh['value']
30         # 14,1: 死亡原因 I死因 14,2: 死亡原因 Iア期間 14,3: 死亡原因 Iイ死因 14,4:
31         # 死亡原因 Iイ期間 14,5: 死亡原因 Iウ死因 14,6: 死亡原因 Iウ期間 14,7: 死亡原因 I工
32         # 死因 14,8: 死亡原因 I工期間 14,9: 死亡原因 II死因 14,10: 死亡原因 II期間 14,12: 死
33         # 亡原因 手術 14,14: 死亡原因 手術年月日 14,16: 死亡原因 解剖 16,4: 障害が発生したとこ
34         # ろ その他 16,8: 手段及び状況 17,1: 生後1年未満 詳細 18,1: その他付言 99,1: 備考
35         # for sbhnum, sbhnum2 in [('14', '12'), ('14', '16'), ('16', '4'), ('16',
36         # '8'), ('17', '1'), ('18', '1'), ('99', '1')]:
37         # 18,1: その他付言 99,1: 備考
38         for sbhnum, sbhnum2 in [('18', '1'), ('99', '1')]:
39             col = copy.deepcopy(values[sbhnum][sbhnum2])
40             for k in col.keys():
41                 if k=='colnum': continue
42                 value = col[k]
43                 if len(value):
44                     org_value = value
45                     # 14用の正規表現でパース
46                     for r in regs14:
47                         # 14の区切り位置を[col14]で置換する
48                         t = fw_reg(r[0]).subn(r[1], value)
49                         if t[1]: value = t[0]
50                     # 16用の正規表現でパース
51                     for r in regs16:
52                         # 16の区切り位置を[col16]で置換する
53                         t = fw_reg(r[0]).subn(r[1], value)
54                         if t[1]: value = t[0]
55                     # 18用の正規表現でパース
56                     for r in regs18:
57                         # 18の区切り位置を[col18]で置換する
58                         t = fw_reg(r[0]).subn(r[1], value)
59                         if t[1]: value = t[0]
60                     # 上記以外用の正規表現でパース
61                     for r in regsOther:
62                         # 区切り位置を[colxx_0]で置換する
63                         t = fw_reg(r[0]).subn(r[1], value)
64                         if t[1]: value = t[0]
65                     # スプリッタ ([colxx]) で切り分け
66                     value_arr = fw_reg(r'(\[col[\^]\]+\])\([^\(\[$]\)*)').findall(value)
67                     if value_arr:
68                         len14_1 = len(values['14']['1']) - 1
69                         len14_2 = len(values['14']['2']) - 1
70                         len14_3 = len(values['14']['3']) - 1
71                         len14_4 = len(values['14']['4']) - 1
72                         len14_5 = len(values['14']['5']) - 1
73                         len14_6 = len(values['14']['6']) - 1
```

```

67 len14_7 = len(values['14']['7']) - 1
68 len14_8 = len(values['14']['8']) - 1
69 len14_9 = len(values['14']['9']) - 1
70 len14_10 = len(values['14']['10']) - 1
71 len14_12 = len(values['14']['12']) - 1
72 len14_14 = len(values['14']['14']) - 1
73 len14_16 = len(values['14']['16']) - 1
74 len16_4 = len(values['16']['4']) - 1
75 len16_8 = len(values['16']['8']) - 1
76 len18_1 = len(values['18']['1']) - 1
77 # OPTIMIZE: いったん組み直して正規表現化することでスプリッタを復元する
78 # NOTE: 行ごとに固有の正規表現を組むためfw_regを使用しない。
79 splited_value_regexp = '(' + fw_reg(r'\[col[^\\]+\]+\\').sub(r')
80 (.*)?', value) + ')'
81 splited_value_arr = re.findall(splited_value_regexp, org_value)
82 for i, (spliter, v) in enumerate(value_arr):
83     # 連結前の整形
84     tmp = fw_reg(r'^は)?「(.*)」(である)?(.)?$').search(v)
85     if tmp: v = tmp.groups()[1]
86     v = fw_reg(r'((続< ))?\n?').sub('', v)
87     if spliter == '[col18_top]':
88         sbh['value']['18'][1][len18_1] = v
89         len18_1 += 1
90     elif spliter == '[col18]':
91         sbh['value']['18'][1][len18_1 + 100] = v
92         len18_1 += 1
93     #elif spliter == '[col16_4]':
94     #    sbh['value']['16'][4][len16_4] = v
95     #    len16_4 += 1
96     elif spliter == '[col16]' or spliter == '[col16_8]':
97         sbh['value']['16'][8][len16_8] = v
98         len16_8 += 1
99     elif spliter == '[col14_1]':
100        sbh['value']['14'][1][len14_1] = v
101        len14_1 += 1
102    elif spliter == '[col14_2]':
103        sbh['value']['14'][2][len14_2] = v
104        len14_2 += 1
105    elif spliter == '[col14_3]':
106        sbh['value']['14'][3][len14_3] = v
107        len14_3 += 1
108    elif spliter == '[col14_4]':
109        sbh['value']['14'][4][len14_4] = v
110        len14_4 += 1
111    elif spliter == '[col14_5]':
112        sbh['value']['14'][5][len14_5] = v
113        len14_5 += 1
114    elif spliter == '[col14_6]':
115        sbh['value']['14'][6][len14_6] = v
116        len14_6 += 1
117    elif spliter == '[col14_7]':
118        sbh['value']['14'][7][len14_7] = v
119        len14_7 += 1
120    elif spliter == '[col14_8]':
121        sbh['value']['14'][8][len14_8] = v
122        len14_8 += 1
123    elif spliter == '[col14_9]':
124        sbh['value']['14'][9][len14_9] = v
125        len14_9 += 1
126    elif spliter == '[col14_10]':
127        sbh['value']['14'][10][len14_10] = v
128        len14_10 += 1
129    elif spliter == '[col14_12]':
130        sbh['value']['14'][12][len14_12] = v
131        len14_12 += 1
132    elif spliter == '[col14_14]':
133        sbh['value']['14'][14][len14_14] = v
134        len14_14 += 1
135    elif spliter == '[col14_16]':
136        sbh['value']['14'][16][len14_16] = v
137        len14_16 += 1
138 else: # ここまでに合致しないスプリッタの文言は捨てる
      f_removed.write(rowid + '\t' + spliter + '\t' +

```

```

139     splitted_value_arr[0][i+1].rstrip('\n') + '\n')
140         # 切り分けた先頭は元の欄に残す
141     try:
142         sbh['value'][sbhnum][sbhnum2][k] = fw_reg(r'^([^\n]+)').search(value).group() + '\n'
143     except AttributeError:
144         sbh['value'][sbhnum][sbhnum2][k] = '\n'
145     sbh_data_writer(fos, sbh)
146
147 fis.close()
148 fos.close()
149 f_removed.close()
150
151 def sbh_data_generator(fis):
152     l = fis.readline()
153     bef_rowid = l.split(',') [0]
154     sbh = {'rowid': bef_rowid, 'value': {}}
155     while l != '':
156         rowid, colnum, sbhnum, sbhnum2, colrnum, value = l.split(',')
157         if bef_rowid != rowid:
158             yield sbh
159             sbh = {'rowid': rowid, 'value': {}}
160             # キーを変更してdictにする
161             # sbh.update({'rowid': rowid, 'value': {'sbhnum': sbhnum, 'value': {'sbhnum2': sbhnum2, 'sbhnum': sbhnum, 'value': {sbhrownum: value}}}})
162             # sbh.update({'rowid': rowid, 'value': {sbhnum: {sbhnum2: {'colnum': colnum, 'colrnum': value}}}})
163             if not sbhnum in sbh['value']: sbh['value'][sbhnum] = {}
164             if not sbhnum2 in sbh['value'][sbhnum]: sbh['value'][sbhnum][sbhnum2] = {'colnum': colnum}
165             sbh['value'][sbhnum][sbhnum2][colrnum] = value
166             bef_rowid = rowid
167
168     l = fis.readline()
169     yield sbh
170
171 def sbh_data_writer(fos, sbh):
172     # rowid      : new_shibo_join.tsvの先頭カラム、行ごとのID
173     # colnum     : new_shibo_join.tsvのカラム番号
174     # sbhnum     : 死亡診断書の欄番号。new_shibo_join.tsvには存在しない。
175     # sbhnum2    : 死亡診断書の欄番号の枝番号。new_shibo_join.tsvには存在しない。
176     # colrnuma   : 死亡診断書の欄番号、枝番号ごとの内容が複数行に渡る場合の行番号。
177     # value       : new_shibo_join.tsvの2カラム目以降の値。
178     rowid = sbh['rowid']
179     out_arr = []
180     for sbhnum, sbhnum_col in [(str(sbhnum), sbh['value'][str(sbhnum)])] for
181     sbhnum in sbh['value'].keys():
182         if sbhnum == 'rowid': continue
183         for sbhnum2, sbhnum2_col in [(sbhnum2, sbhnum_col[sbhnum2]) for sbhnum2
184         in sbhnum_col.keys()]:
185             colnum = sbhnum2_col['colnum']
186             for colrnum, value in [(colrnum, sbhnum2_col[colrnum])] for
187             colrnum in sbhnum2_col.keys():
188                 if colrnum == 'colnum': continue
189                 out_arr.append([rowid, colnum, sbhnum, sbhnum2, str(colrnum),
190                 value])
191     out_arr.sort(key=lambda item: (int(item[1]), int(item[4])))
192
193     bef_colnum = out_arr[0][1]
194     l = ','.join(out_arr[0])
195     for l_arr in out_arr[1:]:
196         if bef_colnum != l_arr[1]:
197             fos.write(fw_reg(r'^(続く)?\n?').sub('', l) + '\n')
198             l = ','.join(l_arr)
199             bef_colnum = l_arr[1]
200         else:
201             l = l.rstrip('\n') + l_arr[5]
202     fos.write(l)
203
204 fw_reg_cache = {}
205 def fw_reg(p):
206     """
207     初めて取得するパターンはコンパイルして返す。

```

```
203 初めてでないパターンはキャッシュから返す。
204 """
205 try:
206     return fw_reg_cache[p]
207 except KeyError:
208     r = re.compile(p)
209     fw_reg_cache[p] = r
210     return r
211
212
213 if __name__ == '__main__':
214     main()
215
216
217
```

備考欄前処理プログラム (regexp14.py)

```

3 def get():
4     return [[r'(?14)?([---「()?(1|1|I|i)?([「()?(ア)+(」)])?(の|_|)?(死因)+(欄|ラン|らん)?(続き|つづき)?(は)?([、。。：，---])?', r'[col14_1]', 
5         [r'(?14)?([---「()?(1|1|I|i)?([「()?(ア)+(」)])?(死亡までの|_)?(期間|年月日|年|月|日)+(欄|ラン|らん)?(続き|つづき)?(は)?([、。。：，---])?', r'[col14_2]', 
6             [r'(?14)?([---「()?(1|1|I|i)?([「()?(イ)+(」)])?(の|_|)?(死因)+(欄|ラン|らん)?(続き|つづき)?(は)?([、。。：，---])?', r'[col14_3]', 
7                 [r'(?14)?([---「()?(1|1|I|i)?([「()?(イ)+(」)])?(死亡までの|_)?(期間|年月日|年|月|日)+(欄|ラン|らん)?(続き|つづき)?(は)?([、。。：，---])?', r'[col14_4]', 
8                     [r'(?14)?([---「()?(1|1|I|i)?([「()?(ウ)+(」)])?(の|_|)?(死因)+(欄|ラン|らん)?(続き|つづき)?(は)?([、。。：，---])?', r'[col14_5]', 
9                         [r'(?14)?([---「()?(1|1|I|i)?([「()?(ウ)+(」)])?(死亡までの|_)?(期間|年月日|年|月|日)+(欄|ラン|らん)?(続き|つづき)?(は)?([、。。：，---])?', r'[col14_6]', 
10                            [r'(?14)?([---「()?(1|1|I|i)?([「()?(エ)+(」)])?(の|_|)?(死因)+(欄|ラン|らん)?(続き|つづき)?(は)?([、。。：，---])?', r'[col14_7]', 
11                                [r'(?14)?([---「()?(1|1|I|i)?([「()?(エ)+(」)])?(死亡までの|_)?(期間|年月日|年|月|日)+(欄|ラン|らん)?(続き|つづき)?(は)?([、。。：，---])?', r'[col14_8]', 
12                                    [r'(?14)?([---「()?(2|2|II|ii)+(」)])?(の|_|)?(死因)+(欄|ラン|らん)?(続き|つづき)?(は)?([、。。：，---])?(傷病名)?', r'[col14_9]', 
13                                        [r'(?14)?([---「()?(2|2|II|ii)+(+)?([」)])?(死亡までの|_|)?(期間|年月日|年|月|日)+(欄|ラン|らん)?(続き|つづき)?(は)?([、。。：，---])?(傷病名)?', r'[col14_10]', 
14                                            [r'(?14)?([欄中]?)([---、])?(の|_|)?(解剖)+(欄|ラン|らん|_|)?(主要所見|所見)?(の|_|)?(続き|つづき)?(は)?([、。。：，---])?', r'[col14_16]', 
15                                                [r'(?14)?([欄中]?)([---、])?(の|_|)?(手術|手術部位及び主要所見追加)+(?!解剖|期間|年月日|年月|年|月|日|施行日)(欄|ラン|らん)?(の|_|)?(続き|つづき)?(は)?([、。。：，---])?', r'[col14_12]', 
16                                                    [r'(?14)?([欄中]?)([---「(、)]?(の|_|)?(手術期間|手術年月|手術年|手術月日|手術月|手術日|手術施行日)+(欄)?([」])?(の|_|)?(続き|つづき)?(は)?([、。。：，---])?', r'[col14_14]', 
17                                                        [r'((14))+(欄中)?)(の|_|)?([---(、)]?(期間の記載|死亡したとき)?(続き|つづき)?(は)?([、。。：，---])?', r'[col14]', 
18                                                            [r'((14))([1|1|I|i|2|2|II|ii])([---、])?(「(ア)直接原因」|「(イ)(ア)の原因」|「(ウ)(イ)の原因」|「(エ)(ウ)の原因」|「I欄に影響を及ぼした傷病名等」)?(の|_|)?(「死亡までの期間」)+(+)?', r'[col14]', 
19                                                                [r'(?14)(手術・解剖|手術解剖)+(+)?', r'[col14]', 
20                                                                    [r'(?141(1|2|3|4|5)?)(欄|ラン|らん)?(続き|つづき)?(続)?(は)?([、。。：，---])?', r'[col14]']#分類できないものはこのキーに置く
21 #記号控え([、。。：，---])
22
23 def _assert(r, s):
24     import re
25     # [print(re.subn(rr[0], rr[1], s)[1]) for rr in r]
26     print('開始[' + s + ']')
27     result = ''
28     for rr in r:
29         t = re.subn(rr[0], rr[1], s)
30         if t[1]:
31             if len(result.replace('0', '')): print('経過[' + s + ']')
32             s = t[0]
33             result += '1'
34         else:
35             result += '0'
36     # l = list(filter(lambda result: result[1] > 0, [re.subn(rr[0], rr[1], s) for rr in r]))
37     print('結果[' + s + ']\n')
38     # print([result[0] + '\n' for result in l])
39     return len(result.replace('0', ''))
40
41 if __name__ == '__main__':
42     r = get();
43
44     print('開始 =====')
45
46     print('===== 終了')
47

```

備考欄前処理プログラム (regexp16.py)

```
1 def get():
2     return [[r' (?16) ?(欄中?)?( )?)([の ()](続き|続|つづき|追記)([ 、・：, --])?', r'[col16]', #16の情報として抽出
3             [r' (?16) ?(欄中?)?([一一の「」)?(傷害|障害|死亡)?(発生|が発生した|した)+(日時|時刻|時分|とき|時間)+(欄)?([の () ])?(続き|続|つづき|追記|時分)?(は)?([ 、・：, --])?', r'[col16_4]', #傷害が発生した時刻として抽出
4                 [r' (?16) ?(欄中?)?([一一の「」)?(傷害|障害|死亡)?(発生|が発生した|した)+(日|時)+(欄)?([の () ])?(続き|続|つづき|追記)|?((は)?([ 、・：, --])?)?', r'[col16_4]', #傷害が発生した時刻として抽出
5                     [r' (?16) (欄中?)?(')?(手段及び状況|状況|手段および状況)?(欄)?([の ()])?(続き|続|つづき|追記)?(において|は)?([ 、・：, --])?', r'[col16_8]', #手段及び状況として抽出と結合
6                         [r' (?16 0 5) ?([の ()](続き|続|つづき|追記)?(は))?( [ 、・：, --])?', r'[col16_8]', #手段及び状況として抽出と結合
7                             [r' 「?16枠外」?([ 、・：, --])?', r'[col16]', #16の情報として抽出
8                               [r' (16追加)', r'[col16]', #16の情報として抽出
9                                 [r' 16[ ) ; 欄]', r'[col16]', #16の情報として抽出
10                                [r' 16(欄)?(--)', r'[col16]', #16の情報として抽出
11                                    [r'(外因死)?(の)?(追加事項の)?: ?(手段及び状況)+(欄)?(について)?(の)?(続き|続|つづき|追記)?(は)?([ 、。・：, --])?', r'[col16_8]', #手段及び状況として抽出と結合
12                                         [r' (?16 0 5) ?([の ()](続き|続|つづき|追記)?(は))?( [ 、・：, --])?', r'[col16_8]', #手段及び状況として抽出と結合
13
14 def _assert(r, s):
15     import re
16     # [print(re.subn(rr[0], rr[1], s)[1]) for rr in r]
17     print('開始[' + s + ']')
18     result = ''
19     for rr in r:
20         t = re.subn(rr[0], rr[1], s)
21         if t[1]:
22             if len(result.replace('0', '')): print('経過[' + s + ']')
23             s = t[0]
24             result += '1'
25         else:
26             result += '0'
27     # l = list(filter(lambda result: result[1] > 0, [re.subn(rr[0], rr[1], s)
28     for rr in r]))
29     print('結果[' + s + ']\n')
30     # print([result[0] + '\n' for result in l])
31     return len(result.replace('0', ''))

32
33 if __name__ == '__main__':
34     r = get();
35
36     print('開始 =====')
37
38     print('===== 終了')
39
40
41
42
43
44
```

備考欄前処理プログラム (regexp18.py)

```
1 def get():
2     return [[r'18 続き、', r'[col18_top]'],
3            [r' (?18) ?(欄)?( )? ?([の ()]?((続き|続|つづき|追記|：)([ 、・：, ---])?)?', r'[col18]'],
4            [r' (?18) ?(欄)?([---])?(その他|その他特に付言すべき事柄)(欄)?([の ()])?(続き|続|つづき|追記)?([ 、・：, ---])?' , r'[col18]'],
5            [r' (18(欄)?)(欄)?([の ()]?((続き|続|つづき|追記)?([ 、・：, ---]))?' , r'[col18]'],
6            [r' (?1801)?([の ()]?((続き|続|つづき|追記)?([ 、・：, ---]))?' , r'[col18]'],
7            [r' 「?18枠外」?([ 、・：, ---])?' , r'[col18]'],
8            [r' (18欄その他特に付言すべきことがら)' , r'[col18]'],
9            [r' 18その他特に付言すべきことがらの続きは、' , r'[col18]'],
10           [r' (18追加)' , r'[col18]'],
11           [r' 18[ ) ;欄]' , r'[col18]'], ]
12
13
14 def _assert(r, s):
15     import re
16     # [print(re.subn(rr[0], rr[1], s)[1]) for rr in r]
17     print('開始[' + s + ']')
18     result = ''
19     for rr in r:
20         t = re.subn(rr[0], rr[1], s)
21         if t[1]:
22             if len(result.replace('0', ''))>0: print('経過[' + s + ']')
23             s = t[0]
24             result += '1'
25         else:
26             result += '0'
27     # l = list(filter(lambda result: result[1]>0, [re.subn(rr[0], rr[1], s)
28     for rr in r]))
29     print('結果[' + s + ']\n')
30     # print([result[0] + '\n' for result in l])
31     return len(result.replace('0', ''))

32
33 if __name__ == '__main__':
34     r = get();
35
36     print('開始 -----')
37     print('----- 終了')
38
39
40
41
42
43
```

備考欄前処理プログラム (regexpOther.py)

```

1 def get():
2     return [# (3) 生年月日
3         [r'(3)(欄|続き)?', r'[col03_0]'],
4         [r'(?<!1)3(続き|続|つづき)[、・：，---]?', r'[col03_0]'],
5         [r'(3(続き|続|つづき)) ', r'[col03_0]'],
6         # (4) 死亡したとき
7         [r'(4)死亡したとき[、・：，---]?', r'[col04_0]'],
8         [r'(4)(欄|続き)?', r'[col04_0]'],
9         [r'(?<!1)4(続き|続|つづき)[、・：，---]?', r'[col04_0]'],
10        [r'(4(続き|続|つづき)) ', r'[col04_0]'],
11        # (5) 死亡したところ
12        [r'(5)(欄|続き)?', r'[col05_0]'],
13        [r'(?<!1)5(続き|続|つづき)[、・：，---]?', r'[col05_0]'],
14        [r'(5(続き|続|つづき)) ', r'[col05_0]'],
15        # (6) 住所
16        [r'(6)(欄|続き)?', r'[col06_0]'],
17        [r'(?<!1)6(続き|続|つづき)[、・：，---]?', r'[col06_0]'],
18        [r'(6(続き|続|つづき)) ', r'[col06_0]'],
19        # (7) 本籍
20        [r'(7)(欄|続き)?', r'[col07_0]'],
21        [r'(?<!1)7(続き|続|つづき)[、・：，---]?', r'[col07_0]'],
22        [r'(7(続き|続|つづき)) ', r'[col07_0]'],
23        # (8) (9) 死亡した人の夫または妻
24        [r'(8)(欄|続き)?', r'[col08_0]'],
25        [r'(?<!1)8(続き|続|つづき)[、・：，---]?', r'[col08_0]'],
26        [r'(8(続き|続|つづき)) ', r'[col08_0]'],
27        [r'(9)(欄|続き)?', r'[col09_0]'],
28        [r'(?<!1)9(続き|続|つづき)[、・：，---]?', r'[col09_0]'],
29        [r'(9(続き|続|つづき)) ', r'[col09_0]'],
30        # (10) (11) 死亡したときの世帯の主な仕事と死亡した人の職業・産業
31        [r'(10)(欄|続き)?', r'[col10_0]'],
32        [r'10(続き|続|つづき)[、・：，---]?', r'[col10_0]'],
33        [r'(10(続き|続|つづき)) ', r'[col10_0]'],
34        [r'(11)(欄|続き)?', r'[col11_0]'],
35        [r'11(続き|続|つづき)[、・：，---]?', r'[col11_0]'],
36        [r'(11(続き|続|つづき)) ', r'[col11_0]'],
37        # (12) (13) 死亡したところ、及びその種別
38        [r'(12)?[--]施設名称?(続き|続|つづき)[、・：，---]?', r'[col12_0]'],
39        [r'(12)死亡したところ[、・：，---]?', r'[col12_0]'],
40        [r'(12)死亡したとき[、・：，---]?', r'[col12_0]'],
41        [r'(12)(欄|続き)?', r'[col12_0]'],
42        [r'12(続き|続|つづき)[、・：，---]?', r'[col12_0]'],
43        [r'(12(続き|続|つづき)) ', r'[col12_0]'],
44        [r'(13)死亡したところ[、・：，---]?', r'[col12_0]'],
45        [r'(13)死亡したとき[、・：，---]?', r'[col12_0]'],
46        [r'(13)(欄|続き)?', r'[col13_0]'],
47        [r'13(続き|続|つづき)[、・：，---]?', r'[col13_0]'],
48        [r'(13(続き|続|つづき)) ', r'[col13_0]'],
49        # (17) 生後1年未満で病死した場合の追加事項
50        [r'(17)欄?欄妊娠・分娩時における母体の病態又は異状は?', r'[col17_0]'],
51        [r'(17)(欄|続き)?', r'[col17_0]'],
52        [r'17(続き|続|つづき)[、・：，---]?', r'[col17_0]'],
53        [r'(17(続き|続|つづき)) ', r'[col17_0]'],
54    ]
55]
56
57 def _assert(r, s):
58     import re
59     # [print(re.subn(rr[0], rr[1], s)[1]) for rr in r]
60     print('開始[' + s + ']')
61     result = ''
62     for rr in r:
63         t = re.subn(rr[0], rr[1], s)
64         if t[1]:
65             if len(result.replace('0', '')): print('経過[' + s + ']')
66             s = t[0]
67             result += '1'
68         else:
69             result += '0'
70     # l = list(filter(lambda result: result[1] > 0, [re.subn(rr[0], rr[1], s) for rr in r]))
71     print('結果[' + s + ']\n')
72     # print([result[0] + '\n' for result in l])

```

```
73     return len(result.replace('0', ''))  
74  
75  
76 if __name__ == '__main__':  
77     r = get()  
78  
79     print('開始 -----')  
80     print('----- 終了')  
81  
82  
83  
84  
85  
86
```

機械学習用データセット作成プログラム (TF・IDF)

```
1 from datetime import datetime
2 from gensim import corpora
3 from gensim import models
4 from gensim import matutils
5 import MeCab
6 import math
7 import re
8 import sys
9 import os
10
11 def new_idf(docfreq, totaldocs, log_base=2.0, add=1.0):
12     return add + math.log(1.0 * totaldocs / docfreq, log_base)
13
14 def arg_below():
15     import argparse
16     p = argparse.ArgumentParser()
17     p.add_argument('-a', '--anybelow', help='単語出現回数下限の指定、未入力の場合
Noneで引数を取り、100で処理')
18     a = p.parse_args()
19     return a.anybelow
20
21 def main(any_below = 100, mode = 0):
22     max_cnt = 0
23     print('[' + datetime.now().strftime('%Y-%m-%d %H:%m:%S.%f') + ']開始しま
す。')
24     #単語の出現回数の下限、Noneの場合は100
25     below = int(100 if any_below == None else any_below)
26     #print(below)
27
28     #疎行列ベクトル出力先ディレクトリ作成
29     dirname = ''
30     if any_below != None:
31         dirname = 'TFIDF' + ('' if below <= 0 else '_' + str(below)) + '/'
32         #os.makedirs(dirname, mode=511, exist_ok=True)
33
34     #読み込みファイルと出力ファイル
35     fis = open('./new_shibo_join_concat.tsv', 'r')
36     fos = open('./' + dirname + '05' + ('' if max_cnt <= 0 else '_' +
37     str(max_cnt)) + '_sogouretu.csv', 'w')
38     fos2 = open('./' + dirname + '05' + ('' if max_cnt <= 0 else '_' +
39     str(max_cnt)) + '_total_wordnum.txt', 'w')
40
41     #作成済み辞書の読み込み先リンク
42     gdic_fname = './' + dirname + 'new_shibo_join.dict'
43
44     #トーカナイザを取得
45     tokenizer = get_tokenizer()
46
47     print('[' + datetime.now().strftime('%Y-%m-%d %H:%m:%S.%f') + ']分かち書き開
始')
48     #トーンリスト（複数行の文章、単語の二次元配列）の読み込み
49     rownum_list, tokens_list = fetch_tokenslist(fis, tokenizer, max_cnt)
50     print('[' + datetime.now().strftime('%Y-%m-%d %H:%m:%S.%f') + ']分かち書き終
了')
51
52     print('[' + datetime.now().strftime('%Y-%m-%d %H:%m:%S.%f') + ']dictionary
作成開始')
53     #トーンリストを辞書とするか、既存の辞書を読み込むか、モードで切り替え
54     if mode == 0:
55         #トーンリストを辞書変換
56         tokens_gdic = tokenslist2dic(tokens_list)
57         #辞書情報を保存
58         #save_gdic(tokens_gdic, gdic_fname)
59     elif mode == 1:
60         #辞書情報の読み込み
61         tokens_gdic = load_gdic(gdic_fname)
62
63     print('[' + datetime.now().strftime('%Y-%m-%d %H:%m:%S.%f') + ']dictionary
作成終了')
64     print('[' + datetime.now().strftime('%Y-%m-%d %H:%m:%S.%f') + ']corpus作成
開始')
65     #辞書をbag-of-words形式のリスト（corpus）に変換
```

```

65     tokens_corpus = doc2bow(tokens_gdic, tokens_list)
66     print('['+ datetime.now().strftime('%Y-%m-%d %H:%m:%S.%f') + ']corpus作成
終了')
67
68     print('['+ datetime.now().strftime('%Y-%m-%d %H:%m:%S.%f') + ']model作成開
始')
69     test_model = models.TfidfModel(tokens_corpus,wglobal=new_idf)
70
71     corpus_tfidf = test_model[tokens_corpus]
72     print('['+ datetime.now().strftime('%Y-%m-%d %H:%m:%S.%f') + ']model作成終
了')
73
74     #0回出現の単語についても重要度0を出力する
75     #new_dic = [{wordid:0 for word,wordid in dictionary.token2id.items()} for
c in id_list]
76     #辞書の総単語数を取得
77     total_wordnum = 0
78     for word,wordid in tokens_gdic.token2id.items():
79         total_wordnum += 1
80     fos2.write(str(total_wordnum))
81     print('['+ datetime.now().strftime('%Y-%m-%d %H:%m:%S.%f') + ']ベクトル出力
開始')
82     for rowid,doc in zip(rownum_list,corpus_tfidf):
83         lil = []
84         for word in doc:
85             lil.append((word[0],word[1]))
86
87         fos.write(data_make(rowid,lil))
88     print('['+ datetime.now().strftime('%Y-%m-%d %H:%m:%S.%f') + ']ベクトル出力
終了')
89
90     fis.close()
91     fos.close()
92     print('['+ datetime.now().strftime('%Y-%m-%d %H:%m:%S.%f') + ']終了しま
す。')
93
94 #出力の形式を整えるメソッド
95 def data_make(rowid,lil):
96     return rowid + '\t['+', '.join([sogyouretu(kv) for kv in lil])+'\n'
97
98 #tfidf値を疎行列の形にして返すメソッド
99 def sogyouretu(kv):
100     return str(kv[0])+':'+str(kv[1])
101
102     #for rowdic, rowid in zip(new_dic,id_list):
103     #print(rowid+'\t','.'.join([str(k)+','+str(v) for k,v in
sorted(rowdic.items())]))
104     #print(rowid+'\t','.'.join([str(v) for k,v in rowdic.items()]))
105     #fos.write(rowid+'\t','.'.join([str(v) for k,v in
sorted(rowdic.items())])+'\n')
106
107     """
108     #出現0回の文字について重要度0を入れない出力
109     texts_tfidf = []
110     for rowid,doc in zip(id_list,corpus_tfidf):
111         text_tfidf = []
112         for word in doc:
113             text_tfidf.append(word[1])
114         texts_tfidf.append(text_tfidf)
115         print(rowid+'\t','.'.join([str(i) for i in text_tfidf]))
116     """
117
118 #作成済みdictionaryを読み込むメソッド
119 def load_gdic(fname):
120     return corpora.Dictionary.load(fname)
121
122 #作成したdictionaryを保存するメソッド
123 def save_gdic(tokens_gdic, fname):
124     tokens_gdic.save(fname)
125
126 #形態素解析にMeCabを使用し、トーカナイザーを取得するメソッド
127 def get_tokenizer():
128     return MeCab.Tagger('-r /etc/mecabrc -Owakati')

```

```

129
130 #解析データを読み込み分かち書きしてlistを作成するメソッド
131 def fetch_tokenslist(fis, tokenizer, max_cnt = 0):
132     cnt = 1
133     l = fis.readline();
134     rownum_list = []
135     tokens_list = []
136     while (cnt != max_cnt and l != ''):
137         rownum, sentences = remove_tagwords(l).split('\t')
138         rownum_list.append(rownum)
139         #トーカナイズ (1文章を単語ごとに分割した一次元配列)
140         tokens = tokenizer.parse(sentences).split()
141         #トーンリスト (複数行の文章、単語の二次元配列)
142         tokens_list.append(tokens)
143         cnt += 1
144         l = fis.readline();
145
146     return (rownum_list, tokens_list)
147
148 re_tagwords = re.compile(r'([手術])|([解剖])|([状況])|([母体])|([付言])|'
149 #【備考】|\n')
150 #解析データの項目名を単語から除去
151 def remove_tagwords(sentences):
152     return re_tagwords.sub('', sentences)
153
154 #トーンリストを辞書変換するメソッド
155 def tokenslist2dic(tokens_list):
156     return corpora.Dictionary(tokens_list)
157
158 #辞書をbag-of-words形式のリストに変換するメソッド
159 def doc2bow(tokens_gdic, tokens_list):
160     return list(map(tokens_gdic.doc2bow, tokens_list))
161
162 if __name__ == '__main__':
163     any_below = arg_below()
164     main(any_below, 1)
165
166

```

機械学習用データセット作成プログラム (LSI)

```
1 from datetime import datetime
2 from gensim import corpora
3 from gensim import models
4 from gensim import matutils
5 import MeCab
6 import math
7 import re
8 import os
9
10 def new_idf(docfreq, totaldocs, log_base=2.0, add=1.0):
11     return add + math.log(1.0 * totaldocs / docfreq, log_base)
12
13 def main(mode = 0):
14     max_cnt = 0
15     print('['+ datetime.now().strftime('%Y-%m-%d %H:%m:%S.%f') + ']開始しま
す。')
16
17     #疎行列ベクトル出力先ディレクトリ作成
18     dirname = 'LSI/'
19     os.makedirs(dirname, mode=511, exist_ok=True)
20
21     #読み込みファイルと出力ファイル
22     fis = open('./new_shibo_join_concat.tsv', 'r')
23     fos = open('./'+ dirname + ('' if max_cnt <= 0 else '_' + str(max_cnt) ) +
24 + '05_sogyouretu.csv', 'w')
25     fos2 = open('./'+ dirname + ('' if max_cnt <= 0 else '_' + str(max_cnt) ) +
26 + '05_total_wordnum.txt', 'w')
27
28     #作成済み辞書の読み込み先リンク
29     gdic_fname = './new_shibo_join.dict'
30
31     #トーカナイザを取得
32     tokenizer = get_tokenizer()
33
34     print('['+ datetime.now().strftime('%Y-%m-%d %H:%m:%S.%f') + ']分かち書き開
始')
35     #トークンリスト（複数行の文章、単語の二次元配列）の読み込み
36     rownum_list,tokens_list = fetch_tokenslist(fis, tokenizer, max_cnt)
37     print('['+ datetime.now().strftime('%Y-%m-%d %H:%m:%S.%f') + ']分かち書き終
了')
38
39     print('['+ datetime.now().strftime('%Y-%m-%d %H:%m:%S.%f') + ']dictionary
作成開始')
40     #トークンリストを辞書とするか、既存の辞書を読み込むか、モードで切り替え
41     if mode == 0:
42         #トークンリストを辞書変換
43         tokens_gdic = tokenslist2dic(tokens_list)
44         #辞書情報を保存
45         #save_gdic(tokens_gdic, gdic_fname)
46     elif mode == 1:
47         #辞書情報の読み込み
48         tokens_gdic = load_gdic(gdic_fname)
49
50     print('['+ datetime.now().strftime('%Y-%m-%d %H:%m:%S.%f') + ']dictionary
作成終了')
51
52     print('['+ datetime.now().strftime('%Y-%m-%d %H:%m:%S.%f') + ']corpus作成
開始')
53     #辞書をbag-of-words形式のリスト（corpus）に変換
54     tokens_corpus = doc2bow(tokens_gdic, tokens_list)
55     print('['+ datetime.now().strftime('%Y-%m-%d %H:%m:%S.%f') + ']corpus作成
終了')
56
57     print('['+ datetime.now().strftime('%Y-%m-%d %H:%m:%S.%f') + ']model作成開
始')
58     test_model = models.TfidfModel(tokens_corpus, wglobal=new_idf)
59
60     corpus_tfidf = test_model[tokens_corpus]
61     print('['+ datetime.now().strftime('%Y-%m-%d %H:%m:%S.%f') + ']model作成終
了')
62     #0回出現の単語についても重要度0を出力する
63     #new_dic = [{wordid:0 for word,wordid in dictionary.token2id.items()} for
```

```

c in id_list]
    #辞書の総単語数を取得
    """
65    total_wordnum = 0
66    for word,wordid in tokens_gdic.token2id.items():
67        total_wordnum += 1
68    fos2.write(str(total_wordnum))
    """
69
70    #LSI modelによる次元圧縮したcorpusの作成
71    print('['+ datetime.now().strftime('%Y-%m-%d %H:%m:%S.%f') + ']LSI開始')
72    lsi_model = models.LsiModel(corpus_tfidf, id2word=tokens_gdic,
73    num_topics=200)
74    lsi_model.save('lsi_topics200.model')
75    #lsi_model = models.LsiModel.load('lsi_topics200.model')
76    corpus_lsi = lsi_model[corpus_tfidf]
77    print('['+ datetime.now().strftime('%Y-%m-%d %H:%m:%S.%f') + ']LSI終了')
78
79    print('['+ datetime.now().strftime('%Y-%m-%d %H:%m:%S.%f') + ']ベクトル出力
開始')
80    vec_size = 0
81    for rowid,doc in zip(rownum_list,corpus_lsi):
82        lil = []
83        for word in doc:
84            lil.append((word[0],word[1]))
85
86        if len(lil) != 0:
87            vec_size = len(lil)
88            fos.write(data_make(rowid,lil))
89    fos2.write(str(vec_size))
90    print('['+ datetime.now().strftime('%Y-%m-%d %H:%m:%S.%f') + ']ベクトル出力
終了')
91
92    fis.close()
93    fos.close()
94    fos2.close()
95    print('['+ datetime.now().strftime('%Y-%m-%d %H:%m:%S.%f') + ']終了しま
す。')
96
97    #出力の形式を整えるメソッド
98    def data_make(rowid,lil):
99        return rowid + '\t['+', '.join([sogyouretu(kv) for kv in lil])+'\n'
100
101    #tfidf値を疎行列の形にして返すメソッド
102    def sogyouretu(kv):
103        return str(kv[0])+':'+str('{:f}'.format((kv[1]+1)/2))
104
105        #for rowdic, rowid in zip(new_dic,id_list):
106        #print(rowid+'\t','.'.join([str(k)+','+'str(v) for k,v in
107        sorted(rowdic.items())]))
108        #print(rowid+'\t','.'.join([str(v) for k,v in rowdic.items()]))
109        #fos.write(rowid+'\t','.'.join([str(v) for k,v in
110        sorted(rowdic.items())])+'\n')
111
112        """
113        #出現0回の文字について重要度0を入れない出力
114        texts_tfidf = []
115        for rowid,doc in zip(id_list,corpus_tfidf):
116            text_tfidf = []
117            for word in doc:
118                text_tfidf.append(word[1])
119            texts_tfidf.append(text_tfidf)
120            print(rowid+'\t','.'.join([str(i) for i in text_tfidf]))
121        """
122
123    #作成済みdictionaryを読み込むメソッド
124    def load_gdic(fname):
125        return corpora.Dictionary.load(fname)
126
127    #作成したdictionaryを保存するメソッド
128    def save_gdic(tokens_gdic, fname):
129        tokens_gdic.save(fname)
130
131    #形態素解析にMeCabを使用し、トーカナイザーを取得するメソッド

```

```

129 def get_tokenizer():
130     return MeCab.Tagger('-r /etc/mecabrc -Owakati')
131
132 #解析データを読み込み分かち書きしてlistを作成するメソッド
133 def fetch_tokenslist(fis, tokenizer, max_cnt = 0):
134     cnt = 1
135     l = fis.readline();
136     rownum_list = []
137     tokens_list = []
138     while (cnt != max_cnt and l != ''):
139         rownum, sentences = remove_tagwords(l).split('\t')
140         rownum_list.append(rownum)
141         #トーカナイズ（1文章を単語ごとに分割した一次元配列）
142         tokens = tokenizer.parse(sentences).split()
143         #トーンリスト（複数行の文章、単語の二次元配列）
144         tokens_list.append(tokens)
145         cnt += 1
146         l = fis.readline();
147
148     return (rownum_list, tokens_list)
149
150 re_tagwords = re.compile(r'([手術])|([解剖])|([状況])|([母体])|([付言])|'
151   ([備考])|\n')
152 #解析データの項目名を単語から除去
153 def remove_tagwords(sentences):
154     return re_tagwords.sub(' ', sentences)
155
156 #トーンリストを辞書変換するメソッド
157 def tokenslist2dic(tokens_list):
158     return corpora.Dictionary(tokens_list)
159
160 #辞書をbag-of-words形式のリストに変換するメソッド
161 def doc2bow(tokens_gdic, tokens_list):
162     return list(map(tokens_gdic.doc2bow, tokens_list))
163
164 if __name__ == '__main__':
165     main(1)
166
167

```

機械学習用データセット作成プログラム (Word2Vec)

```
1 import sys
2 #sys.path.append('/home/bmiwork/ITEC')
3 from my_utils import exec_time
4
5 import MeCab
6 import re
7 from gensim.models import word2vec
8 from gensim import corpora
9 #from gensim import matutils
10 import numpy as np
11 import os
12
13 @exec_time.printer
14 def main(mode = 0):
15     max_cnt = 0
16
17     #疎行列ベクトル出力先ディレクトリ作成
18     dirname = 'WORD2VEC/'
19     os.makedirs(dirname, mode=511, exist_ok=True)
20
21     #トーカナイザを取得
22     tokenizer = get_tokenizer()
23
24     #読み込みファイルと出力ファイル
25     fis = open('./new_shibo_join_concat.tsv', 'r')
26     fos = open('./' + dirname + ('' if max_cnt <= 0 else '_' + str(max_cnt)) +
27 + '05_sogouretu.csv', 'w')
28     fos2 = open('./' + dirname + ('' if max_cnt <= 0 else '_' + str(max_cnt)) +
29 ) + '05_total_wordnum.txt', 'w')
30
31     #作成済みmodelの出力先リンク
32     #vec_fname = './new_shibo_join' + ('' if max_cnt <= 0 else '_' +
33     str(max_cnt)) + '.vec.pt'
34     vec_fname = "./WORD2VEC/entity_vector.model.bin"
35
36     #対象文書の形態素を格納
37     rounum_list, tokens_list = fetch_tokenslist(fis, tokenizer, max_cnt)
38     #print(tokens_list[0])
39
40     if mode == 0:
41         #word2vecのmodelに形態素を学習
42         model = fit_word2vec(tokens_list)
43         #modelを保存
44         save_model(model, vec_fname)
45     elif mode == 1:
46         model = load_model(vec_fname)
47
48     #各行の文書のベクトルを計算して出力
49     write_vector(fos, fos2, model, rounum_list, tokens_list)
50
51     fis.close()
52     fos.close()
53     fos2.close()
54
55 #word2vec学習モデルの読み込み
56 def load_model(fname):
57     return word2vec.KeyedVectors.load_word2vec_format(fname, binary=True)
58
59 #word2vec学習モデルの保存
60 def save_model(model, fname):
61     model.wv.save_word2vec_format(fname, binary=True)
62
63 @exec_time.printer
64 def write_vector(fos, fos2, model, rounum_list, tokens_list):
65     #単語ベクトルを平均した値を文書ベクトルとして扱った場合の出力
66     num_features = 200
67     for rowid, doc in zip(rounum_list, tokens_list):
68         #print(rowid)
69         #print(doc[0])
70         #単語を200次元のベクトルとして出力したもの
71         #####
72         vec_list = []
73         for word in doc:
```

```

71     if word in model.key_to_index:
72         vec = model[word]
73         vec_list.append(vec)
74     fos.write(data_make(rowid, doc, vec_list))
75     """
76     #docが空でない場合のみベクトルの平均を作成、空の場合は空のリストを作成
77     if doc != []:
78         feature_vec = np.zeros((num_features), dtype= "float32")
79         for word in doc:
80             if word in model.key_to_index:
81                 feature_vec = np.add(feature_vec, model[word])
82                 #feature_vec = np.add(feature_vec, model.get_vector(word,
norm=True))
83             if len(doc) > 0:
84                 feature_vec = np.divide(feature_vec, len(doc))
85             else:
86                 feature_vec = []
87             #print(feature_vec)
88             fos.write(data_make2(rowid, feature_vec, num_features))
89             fos2.write(str(num_features))
90
91 #出力の形式を整えるメソッド
92 def data_make(rowid, doc, vec_list):
93     return rowid + '\t[' + ','.join([sogyouretu(wordid,vec) for wordid, vec in
zip(doc, vec_list)])+']\n'
94
95 #出力の形式を整えるメソッド
96 def data_make2(rowid, feature_vec, vec_cnt):
97     cnt_list = []
98     for num in range(vec_cnt):
99         cnt_list.append(num)
100    if feature_vec != []:
101        return rowid + '\t[' + ','.join([sogyouretu2(v,i) for v,i in
zip(feature_vec, cnt_list)])+']\n'
102    else:
103        return rowid + '\t[]\n'
104
105 #ベクトルを疎行列形式に整えるメソッド
106 def sogyouretu(wordid, vec):
107     return str(wordid)+':'+str(vec)
108
109 #ベクトルを疎行列形式に整えるメソッド
110 def sogyouretu2(vec, cnt):
111     return str(cnt)+':'+str(':{f}'.format((vec+10)/20))
112     #return str(cnt)+':'+str(vec)
113
114 @exec_time.printer
115 def get_tokenizer():
116     #形態素解析にMeCabを使用する
117     return MeCab.Tagger('-r /etc/mecabrc -Owakati')
118
119 @exec_time.printer
120 def fetch_tokenslist(fis, tokenizer, max_cnt = 0):
121     cnt = 1
122     l = fis.readline()
123     rownum_list = []
124     tokens_list = []
125     while (cnt != max_cnt and l != ''):
126         rownum, sentences = remove_tagwords(l).split('\t')
127         #if(sentences != ''):
128         rownum_list.append(rownum)
129         #トーカナイズ（1文章を単語ごとに分割した一次元配列）
130         tokens = tokenizer.parse(sentences).split()
131         #ストップワードの除去
132         tokens = filter_stopwords(tokens)
133         #トークンリスト（複数の文章、単語の二次元配列）
134         tokens_list.append(tokens)
135         cnt += 1
136         l = fis.readline();
137     return (rownum_list, tokens_list)
138
139 re_tagwords = re.compile(r'([手術])|([解剖])|([状況])|([母体])|([付言])|
([備考])|\n')

```

```
140 def remove_tagwords(sentences):
141     return re_tagwords.sub(' ', sentences)
142
143 stopword_list = None
144 def filter_stopwords(tokens):
145     global stopword_list
146     if stopword_list is None:
147         stopword_list = []
148         #ストップワードの読み込み
149         fis = open('../stopwords.txt', 'r')
150         re_comment = re.compile(r'^\s*(#.*)|$')
151         l = fis.readline()
152         while l:
153             if not re_comment.match(l):
154                 stopword_list.append(l.rstrip('\n'))
155             l = fis.readline()
156         fis.close()
157     return [word for word in tokens if word not in stopword_list]
158
159 @exec_time.printer
160 def tokenslist2dic(tokens_list):
161     return corpora.Dictionary(tokens_list)
162
163 @exec_time.printer
164 def fit_word2vec(tokens_list):
165     #word2vecによる単語の学習
166     return word2vec.Word2Vec(tokens_list)
167
168 if __name__ == '__main__':
169     main(1)
170
171
```

機械学習用データセット作成プログラム (Doc2Vec (PV-DM / PV-DBOW))

```
1 #!/usr/bin/env python3
2 # -*- coding: utf-8 -*-
3
4 ##### [embedding.py]
5 ##### version.1.1
6 ##### 2022/03/28
7 #####
8 #####
9 #####
10 #####
11
12 import pandas as pd
13 import numpy as np
14 import argparse
15 from argparse import RawTextHelpFormatter
16 import ast
17 import MeCab
18 from gensim.models.doc2vec import Doc2Vec, TaggedDocument
19 import tensorflow_hub as hub
20 import tensorflow_text
21 import re
22 import os
23
24 import ssl
25 ssl._create_default_https_context = ssl._create_unverified_context
26
27 # Embeddingクラス作成
28 class Embedding(object):
29     def __init__(self, rawdata_path, delete_symbol, params, mode, output_path):
30         self.rawdata_path = rawdata_path      # 入力ファイルのパス
31         self.delete_symbol = delete_symbol    # 形態素解析時に削除する記号
32         self.params = params                  # doc2vecを使用する際の引数(辞書
型)
33         self.mode = mode
34         self.output_path = output_path
35
36     def check_mode(self):
37         '''
38             modeの確認。doc2vecを使用する際は、modeとparams['dm']が一致していないとエラーとな
る。
39             params['dm']を設定していない場合は、modeがparams['dm']の引数となる
40         '''
41         if self.mode==2:
42             print('mode: Universal Sentence Encoder')
43
44         elif self.mode==0 or self.mode==1:
45             if 'dm' not in self.params:
46                 self.params['dm']=self.mode
47                 if self.mode==0:
48                     print('mode: doc2vec PV-DBOW')
49                 elif self.mode==1:
50                     print('mode: doc2vec PV-DM')
51
52             elif 'dm' in self.params:
53                 if self.params['dm'] == self.mode:
54                     if self.mode==0:
55                         print('mode: doc2vec PV-DBOW')
56                     elif self.mode==1:
57                         print('mode: doc2vec PV-DM')
58                 else:
59                     raise KeyError('modeとparams["dm"]が一致していません。 mode: {}, params["dm"]: {}'.format(self.mode, self.params['dm']))
60
61         else:
62             raise ValueError('--mode 0~2から選んでください {0: doc2vec PV-DBOW, 1: doc2vec PV-DM, 3: Universal Sentence Encoder}')
63
64
65     def get_data(self):
66         '''
67             rawdataの読み込み。入力ファイルは一列目にid、二列目にテキストデータを含む、hedderおよ
びindex_colの無い.tsvファイル。
68         '''
```

```

69     # 入力ファイル形式が異なる場合は、read_csv()の引数設定を変更する
70     src = pd.read_csv('{}'.format(self.rawdata_path), delimiter='\t',
71 index_col=None, header=None, names=['id', 'data'])
72     # テキストデータがNaNのサンプル行を削除
73     src_data = src.dropna(how='any')
74
75     sentences = []
76     for text in src_data['data']:
77         # テキスト内の任意の記号を削除
78         text = re.sub(r'{}'.format(self.delete_symbol), '', text)
79
80         # 各行のテキストを1要素(リスト)としてsentencesに二次元リストとして格納
81         text_list = text.split(' ')
82         sentences.append(text_list)
83
84     return src, sentences
85
86 def doc2vec(self, sentences):
87     """
88     doc2vecによるベクトル抽出。Mecab を用いて形態素解析した後、モデルを作成する。
89     """
90     token = []
91     tokernizer = MeCab.Tagger('-Owakati')      # Mecabの形態素解析の引数設定
92
93     # sentences内の各要素(テキスト)を形態素解析
94     for s in sentences:
95         token.append(tokernizer.parse(s[0]).split())
96
97     # 以下処理過程の参考
98     # 例 src_data.iloc[0,1] : 【付言】DNA検査にて本人確認。
99     # 例 sentences[0] : ['付言DNA検査にて本人確認']
100    # 例 token[0] : ['付言', 'DNA', '検査', 'にて', '本人', '確認']
101
102    # tokenをdoc2vecに読み込ませる形式に変換
103    documents = [TaggedDocument(doc, [i]) for i, doc in enumerate(token)]
104
105    #hash関数の指定 モデルの再現性を得るためにハッシュ値を固定
106    import hashlib
107    hashfxn = lambda x: int(hashlib.md5(str(x).encode()).hexdigest(), 16)
108
109    # doc2vecモデル構築
110    model = Doc2Vec(documents, hashfxn=hashfxn, **self.params)
111    vectors = [model.docvecs[i] for i in range(len(sentences))]
112
113    return vectors
114
115 def UniversalSentenceEncoder(self, sentences):
116     """
117     Universal Sentence Encoderによるベクトル抽出。形態素解析は行わずモデルを構築し、
118     512次元ベクトルを出力する。
119     """
120
121     # tensorflow hubから学習済みモデルの読み込み
122     USE_model = hub.load(
123         "https://tfhub.dev/google/universal-sentence-encoder-multilingual/3")  # 2022/01時点、latest version:ver.3
124     vectors_all = USE_model(sentences)
125     vectors = [i for i in vectors_all]
126
127     return vectors
128
129 def output(self, src, vectors):
130     """
131     id とベクトルを結合し、出力用データフレームを作成。output_pathが"None"の場合はファイル
132     出力されない。
133     """
134
135     src_data = src.dropna(how='any')
136     src_data_vec = src_data.assign( vecs = vectors)
137     df = pd.merge(src, src_data_vec, how="left", on = "id")
138     df = df.loc[:, ['id', 'vecs']]
139
140     # output_pathを任意指定した時のみファイルを出力する
141     if self.output_path=='None':
142         pass

```

```

138     else:
139         output_path = self.output_path
140         df.to_csv(output_path, sep='\t', header=False, index=False)
141
142     return df
143
144 def main():
145     parser = argparse.ArgumentParser(prog='embedding',
146                                     description='-----\n'
147                                     'PROGRAM NAME:下流工程を考慮したテキストの特徴抽出プログラム\n'
148                                     '概要: 自然言語(日本語)で記載されたテキストが持つ特徴を抽出するブ
149                                     ログラム。.\n'
150                                     'doc2vecまたはuniversal sentence encoderを用いてベクトル化す
151                                     る。.\n'
152                                     '-----\n'
153                                     '-----\n'
154                                     'usage: $ python3 embedding.py [-m] [-rp] [-ds] [-
155                                     params] [-op]\n',
156                                     formatter_class=RawTextHelpFormatter,
157                                     add_help=True)
158     parser.add_argument('--mode', '-m', type=int, required=True,
159                         help='モデルの選択\n'
160                         'mode:{0: doc2vec PV-DBOW, 1: doc2vec PV-DM, 2:
Universal Sentence Encoder}')
161     parser.add_argument('--rawdata_path', '-rp', type=str, required=True,
162                         help='入力ファイルのパス。.\n'
163                         '入力ファイルは一列目にid、二列目にテキストデータを含む、
hedderおよびindex_colの無い.tsvファイル')
164     parser.add_argument('--delete_symbol', '-ds', type=str, default='[ [] 、。
「」 () ]',
165                         help='形態素解析時の削除したい記号\n'
166                         '削除したい記号をカッコ内に入れ[]、str型で渡す\n'
167                         'e.g. [PROGRAM] --delete_symbol "[ [] 、。 「」 () ]"')
168     parser.add_argument('--doc2vec_params', '-params', type=str, default="{'vector_size':2}",
169                         help='doc2vecに与える引数をstr型で入力する。引数はモジュー
ル内で辞書型に変換される。.\n'
170                         'e.g. [PROGRAM] -params {"vector_size": 2,
>window": 2, "min_count": 1, "epochs": 10}"')
171     args = parser.parse_args()
172
173     # 引数受け取り
174     mode = args.mode
175     rawdata_path = args.rawdata_path
176     delete_symbol = args.delete_symbol
177     params = ast.literal_eval(args.doc2vec_params)
178     output_path = args.output_path
179
180     # インスタンス作成
181     E = Embedding(rawdata_path, delete_symbol, params, mode, output_path)
182     # modeとparamsの確認
183     E.check_mode()
184
185     # 入力データ(rawdata)読み込み
186     src, sentences = E.get_data()
187
188     # doc2vecによる実行
189     if mode==0 or mode==1:
190         vectors = E.doc2vec(sentences)
191     # USEによる実行
192     elif mode==2:
193         vectors = E.UniversalSentenceEncoder(sentences)
194
195     # 出力用データフレーム作成
196     df = E.output(src, vectors)
197

```

```
198 if __name__ == '__main__':
199     main()
```

分類器学習プログラム (fit_and_predict_xgboost.py)

```
1 import pandas as pd
2 import xgboost as xgb
3 import numpy as np
4 import pickle
5 from sklearn.model_selection import train_test_split
6 from sklearn.metrics import accuracy_score
7 from matplotlib import pyplot as plt
8 from sklearn.metrics import confusion_matrix, roc_curve,
9 precision_recall_curve, auc
10 from sklearn.model_selection import GridSearchCV
11 from scipy.sparse import lil_matrix
12 from datetime import datetime
13 import sys
14 import os
15
16 def arg_parse():
17     import argparse
18     p = argparse.ArgumentParser()
19     p.add_argument('-a', '--anyfolder', help='疎行列ベクトル参照先フォルダ')
20     p.add_argument('-g', '--gpuid', help='使用するGPUのID (デフォルトは0)')
21     p.add_argument('-y', '--year', help='学習用データ年')
22     a = p.parse_args()
23     return (a.anyfolder, a.gpuid, a.year)
24
25 def main(any_folder='.', gpuid=0, year=2020):
26     print('[' + datetime.now().strftime('%Y-%m-%d %H:%M:%S.%f') + '] ' +
27 _file_ + ' start')
28
29     # 引数の検証
30     if not any_folder: any_folder = '.'
31     print('[' + datetime.now().strftime('%Y-%m-%d %H:%M:%S.%f') + '] 入力元フォ
ルダ[./07_付帯情報Embedding/' + any_folder + ']')
32     print('[' + datetime.now().strftime('%Y-%m-%d %H:%M:%S.%f') + '] 出力先フォ
ルダ[./' + any_folder + '/RESULT/' + year + '/]')
33     #os.makedirs('./RESULT/' + any_folder, exist_ok=True)
34     # 疎行列行列数ファイル読み込み
35     # rows = 81688 # 行数
36     # cols = 1579 + 2736 # 次元数
37     fis = open('../07_付帯情報Embedding/' + any_folder +
38 '/learningData_smatrix_rowcolnum_' + year + '.txt', 'r')
39     rows = int(fis.readline())
40     cols = int(fis.readline())
41     fis.close()
42     print('[' + datetime.now().strftime('%Y-%m-%d %H:%M:%S.%f') + '] 対象行数
[' + str(rows) + '] 対象ベクトル数[' + str(cols) + ']')
43
44     # 疎行列ファイルを読み込み
45     smatrix_fname = '../07_付帯情報Embedding/' + any_folder +
46 '/learningData_smatrix_' + year + '.txt'
47     X, y = load_sparse_matrix(smatrix_fname, rows, cols)
48     # test_size=0.2 全体のうち0.8を学習用、0.2をテスト用に分割する
49     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
50 shuffle=True, random_state=42, stratify=y)
51
52     # 学習用の0.8で学習
53     # NOTE: tree_method、gpu_idの設定をすることでGPUが使用されるようになる。
54     # ----- Best Estimator
55     clf = xgb.XGBClassifier(base_score=0.5, booster='gbtree',
56     colsample_bylevel=1,
57         colsample_bynode=1, colsample_bytree=1, eta=0.25, gamma=0,
58         gpu_id=gpuid, importance_type='gain',
59         interaction_constraints='',
60             learning_rate=0.25, max_delta_step=0, max_depth=14,
61             min_child_weight=0.05, monotone_constraints='()', n_estimators=2000, n_jobs=1, num_parallel_tree=1,
62             random_state=0,
63                 reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=1,
64                 tree_method='gpu_hist', use_label_encoder=False,
65                 validate_parameters=1, verbosity=None, eval_metric='logloss')
66     print('[' + datetime.now().strftime('%Y-%m-%d %H:%M:%S.%f') + '] fit
start')
67     clf.fit(X_train, y_train)
68     print('[' + datetime.now().strftime('%Y-%m-%d %H:%M:%S.%f') + '] fit
```

```

end')

# モデルを保存
pickle.dump(clf, open('./' + any_folder + '/' + year + '/model.pkl',
'wb'))
clf = pickle.load(open('./' + any_folder + '/' + year + '/model.pkl',
'rb'))

# テスト用の0.2で検証
print('[' + datetime.now().strftime('%Y-%m-%d %H:%M:%S.%f') + '] predict
start')
# pred = clf.predict(X_test)
pred_p = clf.predict_proba(X_test)
print('[' + datetime.now().strftime('%Y-%m-%d %H:%M:%S.%f') + '] predict
end')

# 結果出力
report = []
for i in range(21):
    t = (i/20)
    pred = (pred_p[:, 1] > t).astype(int)
    print('===== 閾値[' + str(t) + '] =====')
    accuracy = print_accuracy(y_test, pred)
    tp, fp, fn, tn, precision, recall = print_confusion_matrix(y_test,
pred)
    report.append([t, accuracy, tn, fp, fn, tp, precision, recall, tp+fp,
(tp+fp)/(tp+fp+fn+tn)])
    write_id_tf(X[:, 0], y_test, pred, t, any_folder, year)

if any_folder == 'TFIDF_100':
    print_importances(clf, any_folder, year)

y_score = pred_p[:,1]
make_pr_curve_pict(y_test, y_score, any_folder, year)
make_roc_curve_pict(y_test, y_score, any_folder, year)

write_report(report, any_folder, year)
write_pred_score(X[:, 0], y_test, pred_p[:, 1], any_folder, year)

print('[' + datetime.now().strftime('%Y-%m-%d %H:%M:%S.%f') + '] ' +
_file_ + ' end')

def load_sparse_matrix(fname, rows, cols):
    # 疎行列ファイルをオープン
    fis = open(fname, 'r')

    X_lil = lil_matrix((rows, cols - 1), dtype=float)
    # y_lil = lil_matrix((rows, 1), dtype=float)
    y = [0] * rows

    # 疎行列ファイルの読み込み
    l = fis.readline()
    rounum = 0
    while l != '':
        l_smatrix = sogyoretu2smatrix(l.rstrip('\n'))
        for k, v in l_smatrix.items():
            if k != cols - 1:
                if len(v) != 0:      # DOC2VEC 疎行列データエラー対策
                    X_lil[rounum, k] = float(v)
                else:
                    print("[ERROR]: '' can not be converted to float @" +
str(rounum))
                    X_lil[rounum, k] = 0
            else:
                ## y_lil[rounum, 0] = float(v)
                y[rounum] = float(v)

        l = fis.readline()
        rounum += 1

    fis.close()

# sparseからnumpy.ndarrayに変換

```

```

126     X = X_lil.todarray()
127     # y = y_lil.todarray()
128     return X, y
129
130 def sogyoretu2smatrix(sogyoretu):
131     """
132     01_sogyoretu.csvの疎行列表記文字列を疎行列のdictionaryに変換する。
133     """
134     # 先頭末尾の[]を外す
135     # ,でスプリットしkvセットに切り分る
136     # :でスプリットしてkeyとvalueを分ける
137     return {int(kv.split(':')[0]): kv.split(':')[1] for kv in
138             sogyoretu[1:-1].split(',') if kv != ''}
139
140 def print_accuracy(act, pred):
141     """
142     精度を出力する。
143     """
144     acc = accuracy_score(act, pred)
145     print('----- Accuracy')
146     print(acc)
147     return acc
148
149 def print_confusion_matrix(act, pred):
150     """
151     混合配列を出力する。
152     """
153     cm = {}
154     cm[(0,0)] = 0
155     cm[(0,1)] = 0
156     cm[(1,0)] = 0
157     cm[(1,1)] = 0
158     for t, p in zip(act, pred):
159         cm[(t, p)] += 1
160     print('----- Confusion Matrix')
161     print('(act, pred): ', cm)
162     print('----- Confusion Matrix')
163     print(pd.DataFrame(confusion_matrix(act, pred, labels=[0, 1]), columns=
164                         ["pred_0", "pred_1"], index=["act_0", "act_1"]))
165
166     p, r = print_pr(cm[(1,1)], cm[(0,1)], cm[(1,0)], cm[(0,0)])
167     return (cm[(1,1)], cm[(0,1)], cm[(1,0)], cm[(0,0)], p, r)
168
169 def print_pr(tp, fp, fn, tn):
170     p = tp/(tp+fp) if (tp+fp)!=0 else 1
171     print('----- Precision')
172     print(str(p))
173     r = tp/(tp+fn) if (tp+fn)!=0 else 1
174     print('----- Recall')
175     print(str(r))
176     return (p, r)
177
178 def print_importances(clf, any_folder, year):
179     """
180     重要度を出力する。
181     """
182     import os
183     from gensim import corpora
184
185     dict_fname = '../07_付帯情報Embedding/' + any_folder +
186     '/new_shibo_join.dict'
187     tokens_dict = None
188     if os.path.exists(dict_fname) and os.path.getsize(dict_fname):
189         tokens_dict = corpora.Dictionary.load(dict_fname)
190
191     icdcodes_fname = '../06_機械学習用基本データ/ICDList/acme_kari_code_sort_' +
192     year + '.txt'
193     icdcodes = ['certificateKey', '年齢', '性別']
194     fis = open(icdcodes_fname, 'r')
195     lines = 0
196     l = fis.readline()
197     while(l != ''):
198         icdcodes.append(l.rstrip('\n'))

```

```

195     l = fis.readline()
196     lines += 1
197
198     fis.close()
199
200     icdcodes.append('手術フラグ(1)')
201     icdcodes.append('手術の部位及び所見(116)')
202     icdcodes.append(' (手術) 備考欄への記載(1)')
203     icdcodes.append('手術日(8)')
204     icdcodes.append('解剖フラグ(1)')
205     icdcodes.append('解剖の部位及び所見(116)')
206     icdcodes.append(' (解剖) 備考欄への記載(1)')
207     icdcodes.append('死因の種類(2)')
208     icdcodes.append('傷害が発生したとき(8)')
209     icdcodes.append('傷害が発生したとき(5)')
210     icdcodes.append('傷害が発生したところの種別(1)')
211     icdcodes.append('傷害が発生したところその他の記述(40)')
212     icdcodes.append('傷害発生場所(8)')
213     icdcodes.append('傷害発生場所(12)')
214     icdcodes.append('傷害発生場所(18)')
215     icdcodes.append('手段及び状況(120)')
216     icdcodes.append(' (傷害) 備考欄への記載(1)')
217     icdcodes.append('「生後1年未満での病死」の病態・異状の詳細(84)')
218     icdcodes.append('備考欄への記載(1)')
219     icdcodes.append('その他付言すべき事柄(60)')
220     icdcodes.append('備考欄外字有無(1)')
221     icdcodes.append('備考欄(1024)')
222
223     lines += 25
224     imps = sorted(enumerate(clf.feature_importances_), key=lambda kv: -kv[1])
225     print('----- Feature Importances')
226     for k, v in imps:
227         print(str(k) + ',' + (icdcodes[k] if k<lines else tokens_dict[k-lines] if tokens_dict else '-') + ',' + str(v))
228
229 def make_pr_curve_pict(y_test, y_score, any_folder, year):
230     # PR取得
231     precision, recall, thresholds = precision_recall_curve(y_true=y_test,
232     probas_pred=y_score)
233     # 結果をプロット
234     plt.figure(1, figsize=[12.8, 9.6], dpi=100)
235     plt.figure(1)
236     plt.plot(recall, precision, label=(any_folder if any_folder != '.' else
237     'precision_recall curve') + ' (AUC = %0.3f)' % auc(recall, precision))
238     for i in range(21):
239         close_point = np.argmin(np.abs(thresholds - (i * 0.05)))
240         plt.plot(recall[close_point], precision[close_point], 'o')
241     # ラベルなどを追加しファイル出力
242     plt.plot([0,1], [1,1], linestyle='--', label='ideal line')
243     plt.legend()
244     plt.xlabel('recall')
245     plt.ylabel('precision')
246     plt.set_title('P-R Curve')
247     fname = './' + any_folder + '/' + year + '/PR-Curve.png'
248     plt.savefig(fname)
249     print('[' + datetime.now().strftime('%Y-%m-%d %H:%M:%S.%f') + '] 結果を出力
しました。[' + fname + ']')
250
251 def make_roc_curve_pict(y_test, y_score, any_folder, year):
252     # ROC取得
253     fpr, tpr, thresholds = roc_curve(y_true=y_test, y_score=y_score)
254     # 結果をプロット
255     plt.figure(2, figsize=[12.8, 9.6], dpi=100)
256     plt.figure(2)
257     plt.plot(fpr, tpr, label=(any_folder if any_folder != '.' else 'roc
curve') + ' (AUC = %0.3f)' % auc(fpr, tpr))
258     for i in range(21):
259         close_point = np.argmin(np.abs(thresholds - (i * 0.05)))
260         plt.plot(fpr[close_point], tpr[close_point], 'o')
261     # ラベルなどを追加しファイル出力
262     plt.plot([0,0,1], [0,1,1], linestyle='--', label='ideal line')
263     plt.plot([0, 1], [0, 1], linestyle='--', label='random prediction')
264     plt.legend()

```

```

263     plt.xlabel('false positive rate(FPR)')
264     plt.ylabel('true positive rate(TPR)')
265     plt.set_title('ROC Curve')
266     fname = './' + any_folder + '/' + year + '/ROC-Curve.png'
267     plt.savefig(fname)
268     print('[' + datetime.now().strftime('%Y-%m-%d %H:%M:%S.%f') + '] 結果を出力
しました。[' + fname + ']')
269
270 def write_report(report, any_folder, year):
271     fname = './' + any_folder + '/' + year +
272     '/fit_and_predict_xgboost_report.csv'
273     fos = open(fname, 'w')
274     for l in list(zip(*report)):
275         fos.write(','.join([str(f) for f in l]) + '\n')
276     fos.close()
277     print('[' + datetime.now().strftime('%Y-%m-%d %H:%M:%S.%f') + '] 結果を出力
しました。[' + fname + ']')
278
279 def write_id_tf(ids, act, pred, ikichi, any_folder, year):
280     fname = './' + any_folder + '/' + year +
281     '/fit_and_predict_xgboost_id_tf_ikichi_' + str(ikichi) + '.csv'
282     fos = open(fname, 'w')
283     for i, a, p in list(zip(ids, act, pred)):
284         fos.write('{:.0f},{:.0f},{:.0f},{:}'.format(i, a, p, a==p) + '\n')
285     fos.close()
286     print('[' + datetime.now().strftime('%Y-%m-%d %H:%M:%S.%f') + '] 結果を出力
しました。[' + fname + ']')
287
288 def write_pred_score(ids, act, pred_score, any_folder, year):
289     fname = './' + any_folder + '/' + year +
290     '/fit_and_predict_xgboost_pred_score.csv'
291     fos = open(fname, 'w')
292     for i, a, p in list(zip(ids, act, pred_score)):
293         fos.write('{:.0f},{:.0f},{:.5f}'.format(i, a, p) + '\n')
294     fos.close()
295     print('[' + datetime.now().strftime('%Y-%m-%d %H:%M:%S.%f') + '] 結果を出力
しました。[' + fname + ']')
296
297 if __name__ == '__main__':
298     any_folder, gpuid, year = arg_parse()
299     main(any_folder, gpuid, year)

```