

МОСКОВСКИЙ ГОСУДАРСТВЕННЫЙ ТЕХНИЧЕСКИЙ
УНИВЕРСИТЕТ ИМЕНИ Н. Э. БАУМАНА
Факультет информатики и систем управления
Кафедра теоретической информатики и компьютерных технологий

Лабораторная работа №17
по курсу «Информационный поиск»
«ЕМ-алгоритм»

Выполнил:
студент группы ИУ9-21М
Беляев А. В.

Проверила:
Лукашевич Н. В.

Москва 2019

1 Цель работы

Подсчитать вероятности отнесения слов и документов к темам за 5 итераций ЕМ-алгоритма (Expectation maximization).

1.1 Ход работы

```
1 from math import sqrt
2 import numpy as np
3
4 np.set_printoptions(precision=3)
5
6 DICT = ['w0', 'w1', 'w2']
7 DOCS = [['w0', 'w1', 'w1'],
8         ['w0', 'w1', 'w2'],
9         ['w0', 'w2', 'w2']]
10 ITERS = 5
11 THEMES = 2
12
13 def rand_init_matrix(cols: int, rows: int) -> list:
14     matrix = np.random.rand(rows, cols)
15     m_sum = sum(np.sum(matrix, 0).tolist())
16     matrix /= m_sum
17
18     matrix = matrix.transpose().tolist()
19     for i in range(len(matrix)):
20         s = sum(matrix[i])
21         for j in range(len(matrix[i])):
22             matrix[i][j] /= s
23     return np.array(matrix).transpose().tolist()
24
25
26 def em_step(m_words: list, m_docs: list, docs: list):
27     def n_dwt(thm_id: int, word_id: int, doc_id: int) -> float:
28         p = m_words[word_id][thm_id] * m_docs[thm_id][doc_id]
29         s = 0
30         for i in range(THEMES):
31             s += m_words[word_id][i] * m_docs[i][doc_id]
32         return docs[doc_id].count(DICT[word_id]) * (p / s)
33
34     def n_wt(thm_id: int, word_id: int) -> float:
35         s = 0
36         for i in range(len(docs)):
37             s += n_dwt(thm_id, word_id, i)
38     return s
39
```

```

40 def n_td(thm_id: int, doc_id: int) -> float:
41     s = 0
42     doc = docs[doc_id]
43     for i in range(len(doc)):
44         s += n_dwt(thm_id, i, doc_id)
45     return s
46
47 def n_t(thm_id: int) -> float:
48     s = 0
49     for i in range(len(DICT)):
50         s += n_wt(thm_id, i)
51     return s
52
53 def n_d(doc_id: int) -> float:
54     s = 0
55     for i in range(THEMES):
56         s += n_td(i, doc_id)
57     return s
58
59 m-fi = np.zeros((len(docs), THEMES))
60 for i in range(len(docs)):
61     for j in range(THEMES):
62         m-fi[i][j] = n-wt(j, i) / n-t(j)
63
64 m-psi = np.zeros((THEMES, len(DICT)))
65 for i in range(THEMES):
66     for j in range(len(DICT)):
67         m-psi[i][j] = n-td(i, j) / n-d(j)
68 return m-fi, m-psi
69
70
71 def main():
72     m-words = rand-init-matrix(THEMES, len(DICT))
73     m-docs = rand-init-matrix(len(DOCS), THEMES)
74
75     print(f'words:\n{np.array(m-words)}')
76     print(f'docs:\n{np.array(m-docs)}')
77
78     for i in range(ITERs):
79         m-words, m-docs = em-step(m-words, m-docs, DOCS)
80         print(f'{i}: words:\n{np.array(m-words)}')
81         print(f'{i}: docs:\n{np.array(m-docs)}')
82
83 if __name__ == '__main__':
84     main()

```

2 Результаты

На листинге ниже приведены исходные матрицы и матрицы, полученные на соответствующих итерациях алгоритма [0, 4]. распределени по двум темам:

```
1 words:
2 [[0.244 0.361]
3  [0.425 0.29 ]
4  [0.331 0.349]]
5 docs:
6 [[0.301 0.596 0.913]
7  [0.699 0.404 0.087]]
8
9 0: words:
10 [[0.293 0.395]
11  [0.267 0.435]
12  [0.44  0.169]]
13 0: docs:
14 [[0.334 0.589 0.898]
15  [0.666 0.411 0.102]]
16
17 1: words:
18 [[0.312 0.363]
19  [0.177 0.557]
20  [0.511 0.08 ]]
21 1: docs:
22 [[0.247 0.59  0.928]
23  [0.753 0.41  0.072]]
24
25 2: words:
26 [[0.333 0.333]
27  [0.099 0.635]
28  [0.568 0.031]]
29 2: docs:
30 [[0.136 0.59  0.964]
31  [0.864 0.41  0.036]]
32
33 3: words:
34 [[0.346 0.318]
35  [0.047 0.672]
36  [0.606 0.01 ]]
37 3: docs:
38 [[0.061 0.579 0.987]
39  [0.939 0.421 0.013]]
40
41 4: words:
```

```

42 [[0.349 0.316]
43    [0.021 0.681]
44    [0.63  0.003]]
45 4: docs:
46 [[0.025 0.559 0.996]
47    [0.975 0.441 0.004]]

```

Распределение по трем темам:

```

1 words:
2 [[0.1987 0.2557 0.2839]
3    [0.4721 0.0767 0.209 ]
4    [0.3292 0.6675 0.5071]]
5 docs:
6 [[0.0703 0.2989 0.3655]
7    [0.7303 0.4927 0.1996]
8    [0.1995 0.2084 0.435 ]]
9
10 0: words:
11 [[0.2438 0.3963 0.329 ]
12    [0.4707 0.2809 0.2849]
13    [0.2855 0.3228 0.3861]]
14 0: docs:
15 [[0.187  0.3539 0.2671]
16    [0.5273 0.4341 0.2562]
17    [0.2857 0.212  0.4767]]
18
19 1: words:
20 [[0.2496 0.3954 0.3246]
21    [0.4328 0.3527 0.2302]
22    [0.3177 0.2519 0.4452]]
23 1: docs:
24 [[0.2284 0.3511 0.2153]
25    [0.5108 0.4359 0.2653]
26    [0.2608 0.213  0.5193]]
27
28 2: words:
29 [[0.2555 0.3934 0.3211]
30    [0.4334 0.411  0.1627]
31    [0.3111 0.1956 0.5162]]
32 2: docs:
33 [[0.2497 0.3491 0.1792]
34    [0.5502 0.4348 0.2285]
35    [0.2001 0.2161 0.5923]]
36

```

```

37 3: words:
38 [[0.2673 0.3905 0.3138]
39  [0.4502 0.4726 0.0876]
40  [0.2825 0.1369 0.5986]]
41 3: docs:
42 [[0.2584 0.3469 0.1383]
43  [0.6203 0.429  0.1653]
44  [0.1214 0.2241 0.6964]]
45
46 4: words:
47 [[0.2882 0.3862 0.3034]
48  [0.469  0.5321 0.031 ]
49  [0.2428 0.0817 0.6656]]
50 4: docs:
51 [[0.2505 0.3433 0.0929]
52  [0.6963 0.4147 0.0988]
53  [0.0532 0.2419 0.8083]]

```

3 Выводы

Были подсчитаны вероятности отнесения документов и слов к темам алгоритмом ЕМ.

В зависимости от начальных матриц, алгоритм «уверенно относит» первый и последний документы к разным темам (1 и 2 или наоборот). С отнесением второго документа алгоритм «не может» определиться.

Похожая ситуация со словами – второе и третье слова достаточно уверенно отнесены к темам, а с отнесением первого слова алгоритм вновь не оперделился.

С распередленеем по трем темам похожая ситуация, но решения лагоритма чуть менее «уверенные».