

## Домашнее задание №1 Отчет

### Задание 1.1.

В качестве текстов для коллекции были выбраны отрывки из книг Д.К. Роулинг о Гарри Поттере. Отрывок для анализа ключевых слов был выбран случайным образом.

#### Выбранные ключевые слова:

Лексема	Частота по частотному словарю	Частота в коллекции	Частота в тексте
Сириус	-	0.000557	0.015337
Гарри	-	0.020475	0.018881
наложить	16.4	7.6e-05	0.000699
заклятие	-	0.00027	0.001399
Думбльдор	-	0.001841	0.002098
Уизли	-	0.001098	0.00979
палочка	21.1	0.00185	0.002797
волшебный	28.2	0.000802	0.001399
Гермиона	-	0.0031	0.001399
рыжий	43.2	0.000101	0.001399
Вольдеморт	-	0.000693	0.000699
Хогварц	-	0.000422	0.000699
смерть	284.1	0.000633	0.000699
заколдовать	-	3.4e-05	0.000699

1) Есть ли среди выбранных вами ключевых слов редкие слова?

Если рассматривать с точки зрения частотного словаря, то некоторые ключевые слова были попросту не найдены, так как многие из ключевых слов являются именами персонажей: Сириус, Гарри и тд. Самыми редкими словами и для одного текста, и для целой коллекции стали: заколдовать и наложить. Выбор на них пал из-за связи с тематикой книги.

2) Есть ли среди выбранных вами слов слова, вошедшие в топ 500 по частоте?

По частотному словарю - нет.

3) К каким частям речи относятся выбранные вами слова, слов какой части речи больше?

В основном это существительные, имена собственные. Так же есть по два прилагательных и глагола.

4) Какие слова встретились во всех или в большинстве документов? Каковы их грамматические характеристики.

Из списка практически во всех документах встретились слова Гарри, Гермiona и палочка (от самого встречающегося к менее). Гарри и Гермiona - существительные, имена собственные. Гарри - м. род, одушевлённое, неизменяется. Гермiona - ж.род, одушевлённое, изменяется.

## Задание 1.2.

Составьте матрицу term x document для всей коллекции и еще для трех высокочастотных слов (по частотному словарю для русского языка).

Была составлена матрица term x document для каждого слова и для каждого текста в коллекции. Также была составлена матрица, каждая ячейка в такой таблице рассчитывалась по формуле  $tf * idf$ , где  $tf$  = Кол-во вхождений слова в тексте/ Кол-во слов в тексте, а  $idf$  = Логарифм от Кол-ва документов в коллекции/ Кол-во документов с данным словом. Столбики в такой таблице - леммы слов, строчки - документ. Всего 58 документов в коллекции.

**Term x document для всей коллекции:**

Так как таблица для всех слов слишком большая, прилагаю только первые 5 строк.

	adairbert	all	alohomora	argus	arsenic	bagshot	bathilda	beat	beater	blot	...	новый	остреб	остребтетиравантик	это	щек	он	йокуть	йля	йл
0	0	0	0	0	0	0	0	0	0	0	...	0	0		0	0	0	0	1	0
1	0	0	0	0	0	0	0	0	0	0	...	0	0		0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	...	0	0		0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	...	1	0		0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	...	0	0		0	0	0	0	0	0

5 rows x 12177 columns

**tf \* idf для всей коллекции:**

В качестве трех высокочастотных слов для русского языка были выбраны слова: *быть, что, не.*

	X.head()																				
	adalbert	all	alohomora	argus	arsenic	bagshot	bathilda	beat	besbor	blot	...	взрыв	истраб	истраб	итеративник	кто	пидк	по	искуть	йля	йл
0	0	0	0	0	0	0	0	0	0	0	...	0	0			0	0	0	0	1	0
1	0	0	0	0	0	0	0	0	0	0	...	0	0			0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	...	0	0			0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	...	1	0			0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	...	0	0			0	0	0	0	0	0

5 rows x 12177 columns

5 rows x 2177 columns

Матрица term x document для шести слов из списка ключевых слов и трех высокочастотных слов:

	гарри	заколдовать	волшебный	гермиона	палочка	думблдор	не	быть	что
0	0.000139	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.000288	0.000422
1	0.001041	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.000426	0.000393
2	0.000722	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.000304	0.000323
3	0.000567	0.000000	0.000495	0.000000	0.000280	0.000000	0.0	0.000281	0.000445
4	0.000898	0.000000	0.001867	0.000000	0.002848	0.000000	0.0	0.000297	0.000356
5	0.000011	0.000000	0.000520	0.000000	0.000737	0.001850	0.0	0.000113	0.000339
6	0.000032	0.000000	0.000767	0.000000	0.001448	0.009083	0.0	0.000164	0.000259
7	0.000703	0.000000	0.001360	0.000000	0.001387	0.026101	0.0	0.000135	0.000292
8	0.000860	0.000000	0.001965	0.000188	0.002784	0.043651	0.0	0.000142	0.000264
9	0.001129	0.000000	0.001042	0.000799	0.000787	0.000000	0.0	0.000136	0.000287
10	0.000980	0.000000	0.000610	0.000117	0.000173	0.000000	0.0	0.000120	0.000246
11	0.001114	0.000000	0.001273	0.000122	0.001442	0.000000	0.0	0.000184	0.000289
12	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.000000	0.000000
13	0.000053	0.000000	0.002954	0.000000	0.003349	0.000000	0.0	0.000096	0.000157
14	0.000533	0.000000	0.001966	0.000251	0.000743	0.024459	0.0	0.000135	0.000291

<b>15</b>	0.00123 5	0.00000 0	0.00038 0	0.00000 0	0.00021 5	0.00135 0	0.0	0.00013 3	0.00033 0
<b>16</b>	0.00109 1	0.00000 0	0.00228 6	0.00126 7	0.00129 6	0.00090 3	0.0	0.00011 0	0.00019 9
<b>17</b>	0.00111 5	0.00000 0	0.00034 3	0.00157 8	0.00077 7	0.00060 9	0.0	0.00014 9	0.00026 2
<b>18</b>	0.00087 2	0.00000 0	0.00061 5	0.00047 2	0.00209 1	0.00036 4	0.0	0.00014 4	0.00026 7
<b>19</b>	0.00080 3	0.00000 0	0.00093 5	0.00007 2	0.00148 4	0.00698 1	0.0	0.00014 3	0.00037 1
<b>20</b>	0.00101 3	0.00000 0	0.00252 1	0.00026 4	0.00272 8	0.00040 7	0.0	0.00014 7	0.00024 6
<b>21</b>	0.00084 5	0.00000 0	0.00083 6	0.00304 4	0.00059 2	0.00742 4	0.0	0.00016 8	0.00043 2
<b>22</b>	0.00068 2	0.00129 9	0.00081 7	0.00078 4	0.00092 7	0.00363 2	0.0	0.00009 3	0.00022 8
<b>23</b>	0.00100 6	0.00000 0	0.00061 8	0.00000 0	0.00035 0	0.00000 0	0.0	0.00017 9	0.00019 2
<b>24</b>	0.00082 8	0.00000 0	0.00060 2	0.00000 0	0.00068 2	0.00000 0	0.0	0.00027 3	0.00029 8
<b>25</b>	0.00067 2	0.00000 0	0.00000 0	0.00095 1	0.00175 6	0.00000 0	0.0	0.00023 1	0.00025 6
<b>26</b>	0.00038 5	0.00000 0	0.00000 0	0.00141 7	0.00000 0	0.00000 0	0.0	0.00025 4	0.00028 0
<b>27</b>	0.00058 8	0.00213 2	0.00000 0	0.00077 2	0.00000 0	0.00000 0	0.0	0.00020 8	0.00027 7
<b>28</b>	0.00070 4	0.00000 0	0.00000 0	0.00024 0	0.00070 9	0.00000 0	0.0	0.00020 7	0.00020 7
<b>29</b>	0.00045 3	0.00000 0	0.00063 9	0.00049 0	0.00072 4	0.00000 0	0.0	0.00025 1	0.00026 4
<b>30</b>	0.00068 8	0.00000 0	0.00061 1	0.00187 4	0.00069 2	0.00000 0	0.0	0.00024 0	0.00031 6
<b>31</b>	0.00069 9	0.00000 0	0.00000 0	0.00071 4	0.00000 0	0.00000 0	0.0	0.00025 7	0.00038 5
<b>32</b>	0.00064 8	0.00000 0	0.00000 0	0.00023 8	0.00000 0	0.00000 0	0.0	0.00021 8	0.00021 8
<b>33</b>	0.00061 4	0.00000 0	0.00000 0	0.00259 0	0.00173 9	0.00000 0	0.0	0.00022 8	0.00029 2
<b>34</b>	0.00078 2	0.00000 0	0.00000 0	0.00167 7	0.00035 4	0.00000 0	0.0	0.00023 2	0.00038 7
<b>35</b>	0.00073 4	0.00000 0	0.00000 0	0.00186 2	0.00137 5	0.00000 0	0.0	0.00026 3	0.00023 8
<b>36</b>	0.00069 4	0.00000 0	0.00000 0	0.00307 4	0.00349 4	0.00000 0	0.0	0.00021 7	0.00024 2

37	0.00086 9	0.00000 0	0.00000 0	0.00363 2	0.00000 0	0.00000 0	0.0	0.00024 8	0.00040 4
38	0.00084 2	0.00000 0	0.00000 0	0.00258 2	0.00104 1	0.00000 0	0.0	0.00016 4	0.00022 8
39	0.00056 8	0.00196 6	0.00000 0	0.00213 6	0.00035 1	0.00000 0	0.0	0.00021 7	0.00046 0
40	0.00099 9	0.00000 0	0.00184 3	0.00047 1	0.00000 0	0.00000 0	0.0	0.00027 9	0.00024 1
41	0.00084 0	0.00000 0	0.00000 0	0.00000 0	0.00000 0	0.00000 0	0.0	0.00033 1	0.00026 9
42	0.00097 1	0.00000 0	0.00000 0	0.00234 9	0.00000 0	0.00000 0	0.0	0.00026 6	0.00044 3
43	0.00080 2	0.00000 0	0.00000 0	0.00452 3	0.00070 3	0.00000 0	0.0	0.00025 7	0.00032 1
44	0.00058 8	0.00000 0	0.00000 0	0.00235 1	0.00034 7	0.00000 0	0.0	0.00027 9	0.00032 9
45	0.00052 9	0.00000 0	0.00000 0	0.00340 5	0.00000 0	0.00000 0	0.0	0.00027 5	0.00030 1
46	0.00071 9	0.00000 0	0.00000 0	0.00283 3	0.00000 0	0.00000 0	0.0	0.00033 1	0.00039 4
47	0.00068 1	0.00000 0	0.00000 0	0.00313 2	0.00033 1	0.00000 0	0.0	0.00031 3	0.00033 7
48	0.00064 4	0.00000 0	0.00000 0	0.00205 0	0.00000 0	0.00000 0	0.0	0.00025 8	0.00033 1
49	0.00083 1	0.00000 0	0.00000 0	0.00231 5	0.00000 0	0.00000 0	0.0	0.00023 7	0.00051 2
50	0.00074 3	0.00000 0	0.00000 0	0.00377 2	0.00034 8	0.00000 0	0.0	0.00016 5	0.00041 9
51	0.00072 8	0.00000 0	0.00000 0	0.00484 9	0.00068 2	0.00000 0	0.0	0.00019 9	0.00026 1
52	0.00106 5	0.00202 9	0.00000 0	0.00612 2	0.00000 0	0.00000 0	0.0	0.00023 8	0.00019 8
53	0.00065 6	0.00000 0	0.00000 0	0.00000 0	0.00000 0	0.00000 0	0.0	0.00018 2	0.00031 2
54	0.00069 6	0.00000 0	0.00000 0	0.00022 9	0.00000 0	0.00000 0	0.0	0.00024 6	0.00044 3
55	0.00071 1	0.00000 0	0.00000 0	0.00242 1	0.00000 0	0.00000 0	0.0	0.00026 1	0.00032 6
56	0.00044 6	0.00000 0	0.00000 0	0.00109 4	0.00000 0	0.00000 0	0.0	0.00017 7	0.00019 2
57	0.00000 0	0.00000 0	0.00096 6	0.00000 0	0.00000 0	0.00000 0	0.0	0.00020 0	0.00014 0

**Матрица  $tf * ifd$  для шести слов из списка ключевых слов и трех высокочастотных слов:**

	гарри	заколд овать	волшеб ный	гермио на	палочк а	думбль дор	не	быть	что
0	0.000139	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.000288	0.000422
1	0.001041	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.000426	0.000393
2	0.000722	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.000304	0.000323
3	0.000567	0.000000	0.000495	0.000000	0.000280	0.000000	0.0	0.000281	0.000445
4	0.000898	0.000000	0.001867	0.000000	0.002848	0.000000	0.0	0.000297	0.000356
5	0.000001	0.000000	0.000520	0.000000	0.000737	0.001850	0.0	0.000113	0.000339
6	0.0000032	0.000000	0.000767	0.000000	0.001448	0.009083	0.0	0.000164	0.000259
7	0.000703	0.000000	0.001360	0.000000	0.001387	0.026101	0.0	0.000135	0.000292
8	0.000860	0.000000	0.001965	0.000188	0.002784	0.043651	0.0	0.000142	0.000264
9	0.001129	0.000000	0.001042	0.000799	0.000787	0.000000	0.0	0.000136	0.000287
10	0.000980	0.000000	0.000610	0.000117	0.000173	0.000000	0.0	0.000120	0.000246
11	0.001114	0.000000	0.001273	0.000122	0.001442	0.000000	0.0	0.000184	0.000289
12	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.000000	0.000000
13	0.0000053	0.000000	0.002954	0.000000	0.003349	0.000000	0.0	0.000096	0.000157
14	0.000533	0.000000	0.001966	0.000251	0.000743	0.024459	0.0	0.000135	0.000291
15	0.001235	0.000000	0.000380	0.000000	0.000215	0.001350	0.0	0.000133	0.000330
16	0.001091	0.000000	0.002286	0.001267	0.001296	0.000903	0.0	0.000110	0.000199
17	0.001115	0.000000	0.000343	0.001578	0.000777	0.000609	0.0	0.000149	0.000262
18	0.000872	0.000000	0.000615	0.000472	0.002091	0.000364	0.0	0.000144	0.000267
19	0.000803	0.000000	0.000935	0.000072	0.001484	0.006981	0.0	0.000143	0.000371
20	0.001013	0.000000	0.002521	0.000264	0.002728	0.000407	0.0	0.000147	0.000246

<b>21</b>	0.00084 5	0.00000 0	0.00083 6	0.00304 4	0.00059 2	0.00742 4	0.0	0.00016 8	0.00043 2
<b>22</b>	0.00068 2	0.00129 9	0.00081 7	0.00078 4	0.00092 7	0.00363 2	0.0	0.00009 3	0.00022 8
<b>23</b>	0.00100 6	0.00000 0	0.00061 8	0.00000 0	0.00035 0	0.00000 0	0.0	0.00017 9	0.00019 2
<b>24</b>	0.00082 8	0.00000 0	0.00060 2	0.00000 0	0.00068 2	0.00000 0	0.0	0.00027 3	0.00029 8
<b>25</b>	0.00067 2	0.00000 0	0.00000 0	0.00095 1	0.00175 6	0.00000 0	0.0	0.00023 1	0.00025 6
<b>26</b>	0.00038 5	0.00000 0	0.00000 0	0.00141 7	0.00000 0	0.00000 0	0.0	0.00025 4	0.00028 0
<b>27</b>	0.00058 8	0.00213 2	0.00000 0	0.00077 2	0.00000 0	0.00000 0	0.0	0.00020 8	0.00027 7
<b>28</b>	0.00070 4	0.00000 0	0.00000 0	0.00024 0	0.00070 9	0.00000 0	0.0	0.00020 7	0.00020 7
<b>29</b>	0.00045 3	0.00000 0	0.00063 9	0.00049 0	0.00072 4	0.00000 0	0.0	0.00025 1	0.00026 4
<b>30</b>	0.00068 8	0.00000 0	0.00061 1	0.00187 4	0.00069 2	0.00000 0	0.0	0.00024 0	0.00031 6
<b>31</b>	0.00069 9	0.00000 0	0.00000 0	0.00071 4	0.00000 0	0.00000 0	0.0	0.00025 7	0.00038 5
<b>32</b>	0.00064 8	0.00000 0	0.00000 0	0.00023 8	0.00000 0	0.00000 0	0.0	0.00021 8	0.00021 8
<b>33</b>	0.00061 4	0.00000 0	0.00000 0	0.00259 0	0.00173 9	0.00000 0	0.0	0.00022 8	0.00029 2
<b>34</b>	0.00078 2	0.00000 0	0.00000 0	0.00167 7	0.00035 4	0.00000 0	0.0	0.00023 2	0.00038 7
<b>35</b>	0.00073 4	0.00000 0	0.00000 0	0.00186 2	0.00137 5	0.00000 0	0.0	0.00026 3	0.00023 8
<b>36</b>	0.00069 4	0.00000 0	0.00000 0	0.00307 4	0.00349 4	0.00000 0	0.0	0.00021 7	0.00024 2
<b>37</b>	0.00086 9	0.00000 0	0.00000 0	0.00363 2	0.00000 0	0.00000 0	0.0	0.00024 8	0.00040 4
<b>38</b>	0.00084 2	0.00000 0	0.00000 0	0.00258 2	0.00104 1	0.00000 0	0.0	0.00016 4	0.00022 8
<b>39</b>	0.00056 8	0.00196 6	0.00000 0	0.00213 6	0.00035 1	0.00000 0	0.0	0.00021 7	0.00046 0
<b>40</b>	0.00099 9	0.00000 0	0.00184 3	0.00047 1	0.00000 0	0.00000 0	0.0	0.00027 9	0.00024 1
<b>41</b>	0.00084 0	0.00000 0	0.00000 0	0.00000 0	0.00000 0	0.00000 0	0.0	0.00033 1	0.00026 9
<b>42</b>	0.00097 1	0.00000 0	0.00000 0	0.00234 9	0.00000 0	0.00000 0	0.0	0.00026 6	0.00044 3

43	0.00080 2	0.00000 0	0.00000 0	0.00452 3	0.00070 3	0.00000 0	0.0	0.00025 7	0.00032 1
44	0.00058 8	0.00000 0	0.00000 0	0.00235 1	0.00034 7	0.00000 0	0.0	0.00027 9	0.00032 9
45	0.00052 9	0.00000 0	0.00000 0	0.00340 5	0.00000 0	0.00000 0	0.0	0.00027 5	0.00030 1
46	0.00071 9	0.00000 0	0.00000 0	0.00283 3	0.00000 0	0.00000 0	0.0	0.00033 1	0.00039 4
47	0.00068 1	0.00000 0	0.00000 0	0.00313 2	0.00033 1	0.00000 0	0.0	0.00031 3	0.00033 7
48	0.00064 4	0.00000 0	0.00000 0	0.00205 0	0.00000 0	0.00000 0	0.0	0.00025 8	0.00033 1
49	0.00083 1	0.00000 0	0.00000 0	0.00231 5	0.00000 0	0.00000 0	0.0	0.00023 7	0.00051 2
50	0.00074 3	0.00000 0	0.00000 0	0.00377 2	0.00034 8	0.00000 0	0.0	0.00016 5	0.00041 9
51	0.00072 8	0.00000 0	0.00000 0	0.00484 9	0.00068 2	0.00000 0	0.0	0.00019 9	0.00026 1
52	0.00106 5	0.00202 9	0.00000 0	0.00612 2	0.00000 0	0.00000 0	0.0	0.00023 8	0.00019 8
53	0.00065 6	0.00000 0	0.00000 0	0.00000 0	0.00000 0	0.00000 0	0.0	0.00018 2	0.00031 2
54	0.00069 6	0.00000 0	0.00000 0	0.00022 9	0.00000 0	0.00000 0	0.0	0.00024 6	0.00044 3
55	0.00071 1	0.00000 0	0.00000 0	0.00242 1	0.00000 0	0.00000 0	0.0	0.00026 1	0.00032 6
56	0.00044 6	0.00000 0	0.00000 0	0.00109 4	0.00000 0	0.00000 0	0.0	0.00017 7	0.00019 2
57	0.00000 0	0.00000 0	0.00096 6	0.00000 0	0.00000 0	0.00000 0	0.0	0.00020 0	0.00014 0

1) Назовите те слова, у которых мощность обратного индекса (количество документов, в которых слово встречается) равна количеству документов в коллекции

Во всех документах встретилось слова *не*.

2) Найдите тексты, удовлетворяющие запросу Word1&Word2&¬Word3

Насколько я понимаю, требуется найти тексты, в которых встречается слова 1 и 2, но не встречается слово 3. Так как в задании не уточнено, что именно подразумевается под словами 1, 2, 3, я рассмотрела комбинации слов из колонок таблиц. При составлении комбинаций не учитывался порядок. Далее привожу таблицу, где пересечение множеств не дало пустое множество.



Word1&Word2& – Word3	Document
('гарри', 'заколдовать', 'волшебный')	{27, 52, 39}
('гарри', 'заколдовать', 'палочка')	{27, 52}
('гарри', 'заколдовать', 'дumbledore')	{27, 52, 39}
('гарри', 'волшебный', 'гермиона')	{3, 4, 5, 6, 7, 13, 15, 23, 24}
('гарри', 'волшебный', 'палочка')	{40}
('гарри', 'волшебный', 'дumbledore')	{3, 4, 40, 9, 10, 11, 13, 23, 24, 29, 30}
('гарри', 'гермиона', 'палочка')	{32, 37, 40, 42, 45, 46, 48, 49, 52, 54, 55, 56, 26, 27, 31}
('гарри', 'гермиона', 'дumbledore')	{9, 10, 11, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 54, 55, 56}
('гарри', 'палочка', 'дumbledore')	{3, 4, 9, 10, 11, 13, 23, 24, 25, 28, 29, 30, 33, 34, 35, 36, 38, 39, 43, 44, 47, 50, 51}
('заколдовать', 'гермиона', 'палочка')	{27, 52}
('заколдовать', 'гермиона', 'дumbledore')	{27, 52, 39}
('заколдовать', 'палочка', 'дumbledore')	{39}
('волшебный', 'гермиона', 'палочка')	{40}
('волшебный', 'гермиона', 'дumbledore')	{40, 9, 10, 11, 29, 30}
('волшебный', 'палочка', 'дumbledore')	{3, 4, 9, 10, 11, 13, 23, 24, 29, 30}
('гермиона', 'палочка', 'дumbledore')	{33, 34, 35, 36, 38, 39, 9, 10, 11, 43, 44, 47, 50, 51, 25, 28, 29, 30}

### Задание 1.3.

Составьте список из 10 лексем/словоформ одного из документов коллекции; в список должны войти:

шесть ключевых слов,

два частотных слова из топ-100 по частотному словарю (см. Новый частотный словарь русского языка <http://dict.ruslang.ru/freq.php> или другой частотный список);

два редких слова (встретившихся во всей Вашей коллекции не больше трех раз, а в Вашем тексте - не больше одного).

Так как список ключевых слов составлялся на основе отрывка одного из документов, то наличие 6 слов из этого списка и два частотных слова могут задать множество, состоящее только из этого же документа. Поэтому именно я перебрала возможные комбинации из списка ключевых слов, для того, чтобы выбрать какой-то другой текст.

Так как  $tf * idf$  был рассчитан ранее, дополнительных вычислений на этом шаге не потребовалось.

В итоге были выбран 20 текст и следующие слова:

Слово (Лексема)	Count(wi)	Fr(Coll) Countcoll(wi)/ N(Coll)	tf
быть	37	0.012036	0.008224
гарри	103	0.020707	0.022894
думбльдор	21	0.001899	0.004668
кошачий	1	0.000017	0.000222
наложить	2	0.000077	0.000445
не	112	0.023765	0.024894
палочка	14	0.001899	0.003112
троица	1	0.000017	0.000222
уизли	3	0.001124	0.000667
хогварц	6	0.000409	0.001334

## Расчет $tf.idf$

Слово (Лексема)	Count(doc)	Count(wi)	DocLength	N	idf	tf	tf.idf
быть	57	37	4499	58	0.017392	0.008224	0.000143
гарри	56	103	4499	58	0.035091	0.022894	0.000803
думбльдор	13	21	4499	58	1.495494	0.004668	0.006981
кошачий	2	1	4499	58	3.367296	0.000222	0.000748
наложить	7	2	4499	58	2.114533	0.000445	0.000940
не	58	112	4499	58	0.000000	0.024894	0.000000
палочка	36	14	4499	58	0.476924	0.003112	0.001484

троица	2	1	4499	58	3.367296	0.000222	0.000748
уизли	29	3	4499	58	0.693147	0.000667	0.000462
хогварц	13	6	4499	58	1.495494	0.001334	0.001994

1) Соответствуют ли те слова, которые попали вверх списка, упорядоченного по убыванию tf.idf, Вашей интуиции?

Да, в топ списка попали все слова из списка ключевых слов

2) Все ли ключевые слова попали в верхнюю часть списка (в первые шесть слов), ранжированного по tf.idf?

Все, кроме одного (уизли).

3) Какие слова попали вниз ранжированного списка? Каковы их характеристики с точки зрения грамматических характеристик, семантики;

В конце списка оказались слова, выбранные из топа частотных слов для русского языка: быть и не.

Быть - глагол, несовершенный вид, непереходный, изолированное спряжение; не - частица; неизменяемое.

4) Как, по-вашему, должен быть устроен список «стоп»-слов, данные о которых нет смысла включать в таблицу?

В этот список должны попасть слова, для которых tf.idf либо 0, либо очень близкий к 0. Нам не нужны слова, которые встречаются практически во всей коллекции текстов, например, предлоги, междометия, цифры, частицы.

#### Усложненный вариант:

Определите тематический вес слов в одном из документов коллекции; отсортируйте слова из текста по tf.idf (выделить топ 20 и 20 с минимальным tf.idf).

Лексема	tf.idf largest	Лексема	tf.idf smallest
вернона	0.017175	люба	0.000061
петунья	0.015291	необходимый	0.000061
дудлить	0.015238	бы	0.000058
дядя	0.014883	забыть	0.000057
фигга	0.010351	остаться	0.000055

тётя	0.008807	кроме	0.000052
дементор	0.008596	тем	0.000048
думбльдор	0.006981	который	0.000043
миссис	0.004970	свой	0.000043
мундгнус	0.004538	один	0.000039
министерство	0.004092	момент	0.000037
дурслей	0.003852	только	0.000035
колдовство	0.003695	наш	0.000035
сын	0.003607	через	0.000031
сова	0.003554	чтобы	0.000027
кухня	0.003475	другой	0.000024
конверт	0.003313	без	0.000023
слушание	0.003292	для	0.000012
авоська	0.002994	не	0.000000
уладить	0.002708	себя	0.000000

Дополнительные вопросы:

- 1) какие слова из списка тематически значимых слов, составленного вручную, вошли в список топ 20 слов по tf.idf, а какие не вошли;

Вошло только слово *думбльдор*, остальные (*гарри, наложить, уизли, палочка, хогварц, быть, не, троица, кошачий*) не вошли.

- 2) задайте пороговое значение по tf.idf для ключевых слов;

Наверное, 0.006

- 3) какие слова, на ваш взгляд, имеют высокий tf.df (выше порогового значения), но не являются ключевыми;

Кажется, что слова дядя и тетя не совсем подходят для того, чтобы быть ключевыми, так как могут много где встречаться. Однако, зная сами рассказы, я бы предложила в качестве ключевых использовать словосочетания дядя Вернон и тетя Петунья, которые так же встретились в топе списка.

#### **Задание 1.4.**

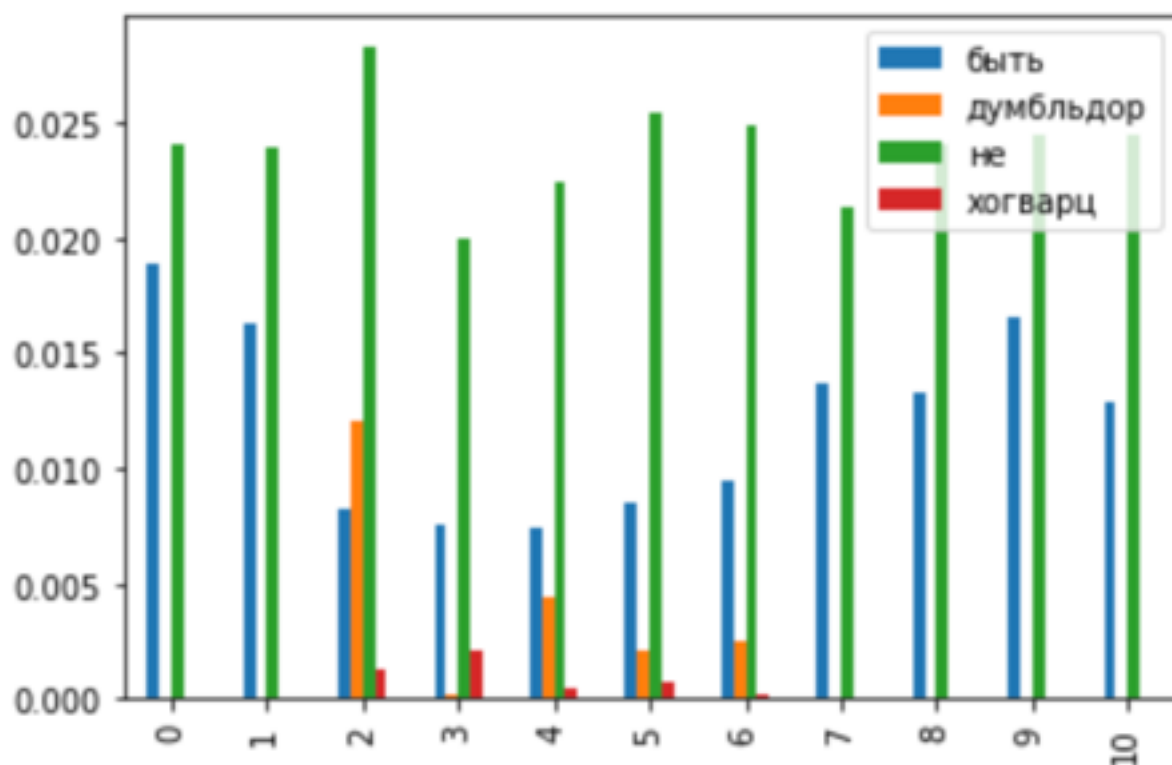
4 слова: два высокочастотных в языке - не, быть

два – с самым высоким tf.idf. - думбльдор и хогварц

ipm

	быть	думбльдор	не	хогварц
0	18823.082042	0.000000	24146.782013	0.000000
1	16332.020367	0.000000	23921.606302	0.000000
2	8168.693009	12063.069909	28305.471125	1234.802432
3	7553.303375	191.222870	19982.789942	2103.451573
4	7461.258848	4400.229577	22479.433710	478.285824
5	8452.844525	2089.467186	25453.509355	664.830468
6	9491.268033	2467.729689	24867.122248	94.912680
7	13736.263736	0.000000	21314.892005	0.000000
8	13230.013230	0.000000	24097.524098	0.000000
9	16539.923954	0.000000	24524.714829	0.000000
10	12889.773481	0.000000	24547.436262	0.000000

Диаграмма относительной частоты слова по подкорпусам



1) Отличаются ли диаграммы для самых частотных в языке слов и для слов с высоким  $tf.idf$  в Вашем списке, если отличаются, то чем?

Как видно из диаграммы, частотные слова для русского языка присутствуют во всех частях коллекции. Практически для всех частей их частотность превышает частотности слов из нашего списка (за исключением 2 текста, в котором частотность слова думбудьдор превышает частотность слова быть).