


```
ERROR: boto3 1.19.30 has requirement urllib3<1.27,>=1.25.4; python_version
|██████████████████████████████████████████████████████████████████████████| 1.4MB 9.1MB/s
|██████████████████████████████████████████████████████████████████████████| 1.7MB 25.6MB/s
|██████████████████████████████████████████████████████████████████████████| 2.9MB 47.1MB/s
ERROR: allennlp 1.2.2 has requirement transformers<3.6,>=3.4, but you'll have
```

```

1 import random
2 import pandas as pd
3
4 def read_gazeta_records(file_name, shuffle=True, sort_by_date=False):
5     assert shuffle != sort_by_date
6     records = []
7     with open(file_name, "r") as r:
8         for line in r:
9             records.append(eval(line)) # Simple hack
10    records = pd.DataFrame(records)
11    if sort_by_date:
12        records = records.sort("date")
13    if shuffle:
14        records = records.sample(frac=1)
15    return records

```

```
1 train_records = read_gazeta_records("gazeta_train.txt")
2 val_records = read_gazeta_records("gazeta_val.txt")
3 test_records = read_gazeta_records("gazeta_test.txt")
```

```
1 train.records.head()
```

	url	text	title
23366	https://www.gazeta.ru/social/2017/10/11/109281...	О том, что чиновники Роскомнадзор а стали фигу...	Следователи пришли в Роскомнадзор
14768	https://www.gazeta.ru/comments/2011/06/07_e_36...	Кремлевская администрация и правительство пыта...	Сумеречное решение

```
1 train_records.shape, val_records.shape, test_records.shape
((52400, 5), (5265, 5), (5770, 5))
```

- ▼ 1 задание: TextRank (порог: 0.35 BLEU)

TextRank - unsupervised метод для составления кратких выжимок из текста. Описание метода:

1. Сплитим текст по предложениям
2. Считаем "похожесть" предложений между собой
3. Строим граф предложений с взвешенными ребрами
4. С помощью алгоритм PageRank получаем наиболее важные предложения, на основе которых делаем summary.

Функция похожести можно сделать и из нейросетевых(или около) моделей: FastText, ELMO и BERT. Выберите один метод, загрузите предобученную модель и с ее помощью для каждого предложения сделайте sentence embedding. С помощью косинусной меры определяйте похожесть предложений.

Предобученные модели можно взять по [ссылке](#).

```

1 from nltk.translate.bleu_score import corpus_bleu
2 from rouge import Rouge
3
4 def calc_scores(references, predictions, metric="all"):
5     print("Count:", len(predictions))
6     print("Ref:", references[-1])
7     print("Hyp:", predictions[-1])
8
9     if metric in ("bleu", "all"):
10         print("BLEU: ", corpus_bleu([[r] for r in references], predictions))
11     if metric in ("rouge", "all"):
12         rouge = Rouge()
13         scores = rouge.get_scores(predictions, references, avg=True)
14         print("ROUGE: ", scores)

```



```

1 !pip install catalyst

```

Requirement already satisfied: pandas>=0.22 in /usr/local/lib/python3.6/dist-packages
Requirement already satisfied: ipython in /usr/local/lib/python3.6/dist-packages
Collecting deprecation
 Downloading <https://files.pythonhosted.org/packages/02/c3/253a89ee03fc9b968>
Requirement already satisfied: PyYAML in /usr/local/lib/python3.6/dist-packages
Requirement already satisfied: numpy>=1.16.4 in /usr/local/lib/python3.6/dist-packages
Requirement already satisfied: tensorboardX in /usr/local/lib/python3.6/dist-packages
Requirement already satisfied: tqdm>=4.33.0 in /usr/local/lib/python3.6/dist-packages
Requirement already satisfied: packaging in /usr/local/lib/python3.6/dist-packages
Requirement already satisfied: torch>=1.1.0 in /usr/local/lib/python3.6/dist-packages
Collecting gitdb<5,>=4.0.1
 Downloading <https://files.pythonhosted.org/packages/48/11/d1800bca0a3bae820>
 |██| 71kB 8.1MB/s
Requirement already satisfied: pyparsing!=2.0.4,!=2.1.2,!=2.1.6,>=2.0.1 in /usr/local/lib/python3.6/dist-packages
Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3.6/dist-packages
Requirement already satisfied: python-dateutil>=2.1 in /usr/local/lib/python3.6/dist-packages
Requirement already satisfied: cycycler>=0.10 in /usr/local/lib/python3.6/dist-packages
Requirement already satisfied: joblib>=0.11 in /usr/local/lib/python3.6/dist-packages
Requirement already satisfied: scipy>=0.17.0 in /usr/local/lib/python3.6/dist-packages
Requirement already satisfied: six in /usr/local/lib/python3.6/dist-packages

```
Requirement already satisfied: six in /usr/local/lib/python3.6/dist-packages
Requirement already satisfied: retrying>=1.3.3 in /usr/local/lib/python3.6/di
Requirement already satisfied: markdown>=2.6.8 in /usr/local/lib/python3.6/di
Requirement already satisfied: tensorboard-plugin-wit>=1.6.0 in /usr/local/li
Requirement already satisfied: setuptools>=41.0.0 in /usr/local/lib/python3.6
Requirement already satisfied: google-auth<2,>=1.6.3 in /usr/local/lib/python
Requirement already satisfied: requests<3,>=2.21.0 in /usr/local/lib/python3.
Requirement already satisfied: grpcio>=1.24.3 in /usr/local/lib/python3.6/dis
Requirement already satisfied: wheel>=0.26; python_version >= "3" in /usr/loc
Requirement already satisfied: google-auth-oauthlib<0.5,>=0.4.1 in /usr/local
Requirement already satisfied: absl-py>=0.4 in /usr/local/lib/python3.6/dist-
Requirement already satisfied: protobuf>=3.6.0 in /usr/local/lib/python3.6/di
Requirement already satisfied: werkzeug>=0.11.15 in /usr/local/lib/python3.6/
Requirement already satisfied: pytz>=2017.2 in /usr/local/lib/python3.6/dist-
Requirement already satisfied: simplegeneric>0.8 in /usr/local/lib/python3.6/
Requirement already satisfied: pickleshare in /usr/local/lib/python3.6/dist-p
Requirement already satisfied: decorator in /usr/local/lib/python3.6/dist-pac
Requirement already satisfied: pexpect; sys_platform != "win32" in /usr/local
Requirement already satisfied: traitlets>=4.2 in /usr/local/lib/python3.6/dis
Requirement already satisfied: prompt-toolkit<2.0.0,>=1.0.4 in /usr/local/lib
Requirement already satisfied: pygments in /usr/local/lib/python3.6/dist-pack
Requirement already satisfied: future in /usr/local/lib/python3.6/dist-packag
Requirement already satisfied: typing-extensions in /usr/local/lib/python3.6/
Requirement already satisfied: dataclasses in /usr/local/lib/python3.6/dist-p
Collecting smmap<4,>=3.0.1
```

Downloading <https://files.pythonhosted.org/packages/b0/9a/4d409a6234eb940e6>.

```
Requirement already satisfied: importlib-metadata; python_version < "3.8" in
Requirement already satisfied: cachetools<5.0,>=2.0.0 in /usr/local/lib/pytho
Requirement already satisfied: pyasn1-modules>=0.2.1 in /usr/local/lib/python
Requirement already satisfied: rsa<5,>=3.1.4; python_version >= "3" in /usr/l
Requirement already satisfied: urllib3!=1.25.0,!1.25.1,<1.26,>=1.21.1 in /us
Requirement already satisfied: idna<3,>=2.5 in /usr/local/lib/python3.6/dist-
Requirement already satisfied: chardet<4,>=3.0.2 in /usr/local/lib/python3.6/
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.6
Requirement already satisfied: requests-oauthlib>=0.7.0 in /usr/local/lib/pyt
Requirement already satisfied: ptyprocess>=0.5 in /usr/local/lib/python3.6/di
Requirement already satisfied: ipython-genutils in /usr/local/lib/python3.6/d
Requirement already satisfied: wcwidth in /usr/local/lib/python3.6/dist-packa
Requirement already satisfied: zipp>=0.5 in /usr/local/lib/python3.6/dist-pac
Requirement already satisfied: pyasn1<0.5.0,>=0.4.6 in /usr/local/lib/python3
Requirement already satisfied: oauthlib>=3.0.0 in /usr/local/lib/python3.6/di
```

```
1 import torch
2 import transformers
3 from tqdm import tqdm
4 from sklearn.metrics.pairwise import cosine_similarity
5 from catalyst.utils import set_global_seed
6
7 seed = 42
8 set_global_seed(seed)
```

```
1 cosine_similarity([[1,2,3,4,5]], [[2,3,4,5,6]]), cosine_similarity([[1,2,3,4,5]]
    (array([[0.99493668]]), array([[0.63636364]]), array([[ -1.]])
```

```
1 # загрузим предварительно обученные модели для токенизации и получения эмбедингов
2
3 # pretrained_model_name = "DeepPavlov/rubert-base-cased-sentence"
4 # tokenizer = transformers.AutoTokenizer.from_pretrained(pretrained_model_name)
```

```

5 # model = transformers.AutoModelWithLMHead.from_pretrained(pretrained_model_name

1 # !wget http://files.deeppavlov.ai/deeppavlov_data/bert/sentence_ru_cased_L-12_H-768_A-12.pt.tar.gz

--2020-12-06 17:08:46-- http://files.deeppavlov.ai/deeppavlov_data/bert/sentence_ru_cased_L-12_H-768_A-12.pt.tar.gz
Resolving files.deeppavlov.ai (files.deeppavlov.ai)... 93.175.29.74
Connecting to files.deeppavlov.ai (files.deeppavlov.ai)|93.175.29.74|:80... connected
HTTP request sent, awaiting response... 200 OK
Length: 661614603 (631M) [application/octet-stream]
Saving to: 'sentence_ru_cased_L-12_H-768_A-12_pt.tar.gz'

sentence_ru_cased_L 100%[=====>] 630.96M  3.02MB/s   in 3m 11s

2020-12-06 17:11:59 (3.31 MB/s) - 'sentence_ru_cased_L-12_H-768_A-12_pt.tar.gz'

1 # !tar -xvzf sentence_ru_cased_L-12_H-768_A-12_pt.tar.gz

sentence_ru_cased_L-12_H-768_A-12_pt/
sentence_ru_cased_L-12_H-768_A-12_pt/pytorch_model.bin
sentence_ru_cased_L-12_H-768_A-12_pt/bert_config.json
sentence_ru_cased_L-12_H-768_A-12_pt/vocab.txt

1 !wget http://files.deeppavlov.ai/deeppavlov_data/bert/sentence_multi_cased_L-12_H-768_A-12.pt.tar.gz

--2020-12-06 17:44:28-- http://files.deeppavlov.ai/deeppavlov_data/bert/sentence_multi_cased_L-12_H-768_A-12.pt.tar.gz
Resolving files.deeppavlov.ai (files.deeppavlov.ai)... 93.175.29.74
Connecting to files.deeppavlov.ai (files.deeppavlov.ai)|93.175.29.74|:80... connected
HTTP request sent, awaiting response... 200 OK
Length: 662398225 (632M) [application/octet-stream]
Saving to: 'sentence_multi_cased_L-12_H-768_A-12_pt.tar.gz'

sentence_multi_cased 100%[=====>] 631.71M  5.53MB/s   in 3m 17s

2020-12-06 17:47:46 (3.21 MB/s) - 'sentence_multi_cased_L-12_H-768_A-12_pt.tar.gz'

1 !tar -xvzf sentence_multi_cased_L-12_H-768_A-12_pt.tar.gz

sentence_multi_cased_L-12_H-768_A-12_pt/
sentence_multi_cased_L-12_H-768_A-12_pt/pytorch_model.bin
sentence_multi_cased_L-12_H-768_A-12_pt/bert_config.json
sentence_multi_cased_L-12_H-768_A-12_pt/vocab.txt

1 import transformers
2 import torch
3 class Seq2Vec():
4     def __init__(self, vocabName="vocab.txt", bertConfig = "bert_config.json", bertName="sentence_multi_cased_L-12_H-768_A-12.pt"):
5         # https://huggingface.co/transformers/pretrained_models.html
6         self.tokenizer = transformers.BertTokenizer.from_pretrained(tokenizer_name)
7         self.simple_vocab = {}
8         with open(vocabName, "r") as vocab:
9             all_text = vocab.read().split('\n')
10            for i in range(len(all_text)):
11                self.simple_vocab[all_text[i]] = i

```

```

12     self.model = transformers.BertModel.from_pretrained(bertName, config = bertC
13     self.model.to("cuda")
14     self.model.eval()
15     self.data = {}
16     self.PAD = self.simple_vocab['[PAD]']
17     self.SEP = self.simple_vocab['[SEP]']
18     self.CLS = self.simple_vocab['[CLS]']
19 def get_ids(self,tokens):
20     result = [self.CLS]
21     for k in tokens:
22         if (k in self.simple_vocab):
23             result.append(self.simple_vocab[k])
24         else: result.append(0)
25     result.append(self.SEP)
26     if (result == [self.CLS]): return [self.CLS,0,self.SEP]
27     return result
28 def get_vec(self,sequence):
29     with torch.no_grad():
30         if("|".join(sequence) in self.data):
31             return self.data["|".join(sequence)]
32         tokens = self.tokenizer.tokenize(" ".join(sequence))
33         tokens = self.get_ids(tokens)
34
35         embeddings = self.model(torch.LongTensor([tokens]).to("cuda")).last_hidde
36
37         res = embeddings.cpu().reshape(-1).detach().numpy()
38         self.data["|".join(sequence)] = res
39         return res
40
41 model = Seq2Vec()

```

```

1 def embed_str(sentence):
2     '''
3     Функция для вывода массива эмбедингов предложения, полученных с пом.
4     модели BERT. Эти эмбединги – выход скрытого слоя модели.
5     '''
6     # print(sentence)
7     if len(sentence) > 1:
8         sentence = ' '.join(sentence)
9     # if sentence in embed_sentence:
10    #     return embed_sentence[sentence]
11
12    input_ids = torch.tensor(tokenizer.encode(sentence)).unsqueeze(0)
13    outputs = model(input_ids)
14    last_hidden_states = [outputs[0].detach().numpy().squeeze()[-1, :]]
15    # embed_sentence[sentence] = last_hidden_states
16    return last_hidden_states

```

```

1 from itertools import combinations
2 import networkx as nx
3 import numpy as np
4 import pymorphy2
5 import razdel
6

```

```

7 i = 0
8
9 def unique_words_similarity(words1, words2):
10     '''
11     Функция подсчёта близости предложений на основе пересечения слов
12     '''
13     global i
14     i += 1
15     words1 = set(words1)
16     words2 = set(words2)
17     if not len(words1) or not len(words2):
18         return 0.0
19     return len(words1.intersection(words2))/(np.log10(len(words1)) + np.log10(len(words2)))
20
21
22 def your_super_words_similarity(words1, words2):
23     global i
24     i += 1
25     # print(i)
26
27     # words1 = np.squeeze(words1)
28     # words2 = np.squeeze(words2)
29     # print('your_super_words_similarity')
30     words1, words2 = model.get_vec(words1), model.get_vec(words2)
31     # print(words1)
32     # print(words1.shape)
33     # print(words2)
34     # print(words2.shape)
35     dist = cosine_similarity([words1], [words2]) #[:, 0]
36     # print('dist', dist)
37     return dist[0]
38     # words1 = words1/np.linalg.norm(words1)
39     # words2 = words2/np.linalg.norm(words2)
40
41     # return words1.dot(words2)
42
43
44 def gen_text_rank_summary(text, calc_similarity=unique_words_similarity, summary_rank=0.5):
45     '''
46     Составление summary с помощью TextRank
47     '''
48     # Разбиваем текст на предложения
49     sentences = [sentence.text for sentence in razdel.sentenize(text)]
50     n_sentences = len(sentences)
51
52     # Токенизируем предложения
53     sentences_words = [[token.text.lower() if lower else token.text for token in sentence.tokenize()] for sentence in sentences]
54
55     # При необходимости лемматизируем слова
56     if morph is not None:
57         sentences_words = [[morph.parse(word)[0].normal_form for word in words] for words in sentences_words]
58
59     # Для каждой пары предложений считаем близость
60     pairs = combinations(range(n_sentences), 2)
61     # print('len pairs ', len(pairs))

```

```

62 # if calc_similarity==unique_words_similarity:
63 scores = [(i, j, calc_similarity(sentences_words[i], sentences_words[j])) for i, j in combinations(range(len(sentences_words)), 2)]
64 # else:
65 #     '''
66 #     скор для варианта с расчетом косинусного расстояния м-ду эмбедингами
67 #     переведем sentences_words в эмбединги
68 #     '''
69 #     scores = [(i, j, calc_similarity(embed_str(sentences_words[i]), embed_str(sentences_words[j]))) for i, j in combinations(range(len(sentences_words)), 2)]
70
71 print('i = ', i)
72 # Строим граф с рёбрами, равными близости между предложениями
73 g = nx.Graph()
74 g.add_weighted_edges_from(scores)
75
76 # Считаем PageRank
77 pr = nx.pagerank(g)
78 result = [(i, pr[i], s) for i, s in enumerate(sentences) if i in pr]
79 result.sort(key=lambda x: x[1], reverse=True)
80
81 # Выбираем топ предложений
82 n_summary_sentences = max(int(n_sentences * summary_part), 1)
83 result = result[:n_summary_sentences]
84
85 # Восстанавливаем оригинальный их порядок
86 result.sort(key=lambda x: x[0])
87
88 # Восстанавливаем текст выжимки
89 predicted_summary = " ".join([sentence for i, proba, sentence in result])
90 predicted_summary = predicted_summary.lower() if lower else predicted_summary
91 return predicted_summary
92
93 def calc_text_rank_score(records, calc_similarity=unique_words_similarity, summary_part=0.1, lower=False):
94     references = []
95     predictions = []
96
97     for text, summary in records[['text', 'summary']].values[:nrows]:
98         summary = summary if not lower else summary.lower()
99         references.append(summary)
100
101         predicted_summary = gen_text_rank_summary(text, calc_similarity, summary_part)
102         text = text if not lower else text.lower()
103         predictions.append(predicted_summary)
104
105     return calc_scores(references, predictions)

```

```
1 # test_records
```

```
1 # calc_text_rank_score(test_records, calc_similarity=unique_words_similarity)
```

```
1 calc_text_rank_score(test_records, calc_similarity=your_super_words_similarity)
```

```
1 = 585934
```

```
i = 584124
```

```
i = 584790
```

```
i = 586065
```



```
l = 586065
i = 587193
i = 587518
i = 587924
i = 588452
i = 589313
i = 589908
i = 590404
i = 591034
i = 591737
i = 592202
i = 592553
i = 593049
i = 593325
i = 593920
i = 594866
i = 595427
i = 596168
i = 596696
i = 597399
i = 597750
i = 598026
i = 598491
i = 599157
i = 599653
i = 601423
i = 601633
i = 602374
i = 603004

i = 603599
i = 604502
i = 604967
i = 605292
i = 605958
i = 607086
i = 607362
i = 607858
i = 608599
i = 609005
i = 609440
i = 609650
i = 610028
i = 610889
i = 611592
i = 612088
i = 613414
i = 614080
i = 615971
i = 616322
i = 616757
i = 617352
i = 617505
i = 618066
i = 618732
i = 619197
i = 620523
i = 621189
```

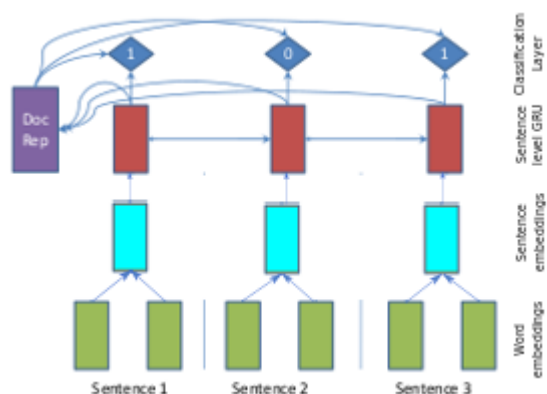
▼ 2 Задание: Extractive RNN (порог: 0.35 BLEU)

Второй метод, который вам предлагается улучшить – поиск предложений для summary с помощью RNN. В рассмотренной методе мы использовали LSTM для генерации sentence embedding. Попробуйте использовать другие архитектуры: CNN, Transformer; или добавьте предобученные модели, как и в первом задании.

P.S. Тут предполагается, что придется изменять много кода в ячейках (например, поменять токенизацию).

▼ Модель

Картинка для привлечения внимания:



Статья с оригинальным методом: <https://arxiv.org/pdf/1611.04230.pdf>

Список вдохновения:

- <https://towardsdatascience.com/understanding-how-convolutional-neural-network-cnn-perform-text-classification-with-word-d2ee64b9dd0b> Пример того, как можно применять CNN в текстовых задачах
- <https://arxiv.org/pdf/1808.08745.pdf> Очень крутой метод генерации summary без Transformers
- <https://towardsdatascience.com/super-easy-way-to-get-sentence-embedding-using-fasttext-in-python-a70f34ac5b7c> – простой метод генерации sentence embedding
- <https://towardsdatascience.com/fse-2b1ffa791cf9> – Необычный метод генерации sentence embedding
- <https://github.com/UKPLab/sentence-transformers> – BERT предобученный для sentence embedding

P.S. Выше написанные ссылки нужны только для разогрева вашей фантазии, можно воспользоваться ими, а можно придумать свой.

Комментарий к заданию: Если посмотреть на архитектуру почти SummaRuNNer, то в ней есть два главных элемента: первая часть, которая читает предложения и возвращает векторы на каждое предложение, и вторая, которая выбирает предложения для суммаризации. Вторую часть мы не трогаем, а первую меняем. На что меняем – как вы решите. Главное: она должна иметь хорошее качество и встроиться в текущую модель.

```

1 import copy
2 import random
3
4 def build_oracle_summary_greedy(text, gold_summary, calc_score, lower=True, max_
5     '''
6     Жадное построение oracle summary
7     '''
8     gold_summary = gold_summary.lower() if lower else gold_summary
9     # Делим текст на предложения
10    sentences = [sentence.text.lower() if lower else sentence.text for sentence
11    n_sentences = len(sentences)
12    oracle_summary_sentences = set()
13    score = -1.0
14    summaries = []
15    for _ in range(min(n_sentences, 2)):
16        for i in range(n_sentences):
17            if i in oracle_summary_sentences:
18                continue
19            current_summary_sentences = copy.copy(oracle_summary_sentences)
20            # Добавляем какое-то предложения к уже существующему summary
21            current_summary_sentences.add(i)
22            current_summary = " ".join([sentences[index] for index in sorted(list
23            # Считаем метрики
24            current_score = calc_score(current_summary, gold_summary)
25            summaries.append((current_score, current_summary_sentences))
26            # Если получилось улучшить метрики с добавлением какого-либо предложения, то пробуем
27            # Иначе на этом заканчиваем
28            best_summary_score, best_summary_sentences = max(summaries)
29            if best_summary_score <= score:
30                break
31            oracle_summary_sentences = best_summary_sentences
32            score = best_summary_score
33    oracle_summary = " ".join([sentences[index] for index in sorted(list(oracle_
34    return oracle_summary, oracle_summary_sentences
35
36 def calc_single_score(pred_summary, gold_summary, rouge):
37     return rouge.get_scores([pred_summary], [gold_summary], avg=True)['rouge-2']

```

```

1 from tqdm import tqdm_notebook as tqdm
2
3 def calc_oracle_score(records, n_rows=1000, lower=True):
4     references = []
5     predictions = []
6     rouge = Rouge()
7
8     for text, summary in tqdm(records[['text', 'summary']].values[:n_rows]):
9         summary = summary if not lower else summary.lower()
10        references.append(summary)
11        predicted_summary, _ = build_oracle_summary_greedy(text, summary, calc_s
12        predictions.append(predicted_summary)
13
14    calc_scores(references, predictions)
15

```

```
16 calc_oracle_score(test_records)
```

```
/usr/local/lib/python3.6/dist-packages/ipykernel_launcher.py:8: TqdmDeprecati
```

```
This function will be removed in tqdm==5.0.0
```

```
Please use `tqdm.notebook.tqdm` instead of `tqdm.tqdm_notebook`
```

```
100%
```

```
1000/1000 [02:26<00:00, 6.80it/s]
```

```
Count: 1000
```

```
Ref: челябинскую городскую больницу №6 из-за мизерных зарплат продолжают покидать врачи. 1
```

```
Нур: однако руководство больницы никак не реагирует на происходящее. тем временем в челябин
```

```
BLEU: 0.4340029354374597
```

```
ROUGE: {'rouge-1': {'f': 0.3539652034259989, 'p': 0.42831485382882784, 'r':
```



Если надо, поменяйте код загрузки токенизатора

```
1 import os
2
3 import youtokentome as yttm
4
5 def train_bpe(records, model_path, model_type="bpe", vocab_size=30000, lower=True):
6     temp_file_name = "temp.txt"
7     with open(temp_file_name, "w") as temp:
8         for text, summary in records[['text', 'summary']].values:
9             if lower:
10                 summary = summary.lower()
11                 text = text.lower()
12                 if not text or not summary:
13                     continue
14                 temp.write(text + "\n")
15                 temp.write(summary + "\n")
16     yttm.BPE.train(data=temp_file_name, vocab_size=vocab_size, model=model_path)
17
18 train_bpe(train_records, "BPE_model.bin")
```

```
1 bpe_processor = yttm.BPE('BPE_model.bin')
```

```
2 bpe_processor.encode(["октябрь богат на изменения"], output_type=yttm.OutputType.SUB
```

```
[['_октябрь', '__богат', '__на', '__изменения']])
```



Если надо, поменяйте код словаря

```
1 from collections import Counter
```

```
2 from typing import List, Tuple
```

```

3 import os
4
5 class Vocabulary:
6     def __init__(self, bpe_processor):
7         self.index2word = bpe_processor.vocab()
8         self.word2index = {w: i for i, w in enumerate(self.index2word)}
9         self.word2count = Counter()
10
11     def get_pad(self):
12         return self.word2index["<PAD>"]
13
14     def get_sos(self):
15         return self.word2index["<SOS>"]
16
17     def get_eos(self):
18         return self.word2index["<EOS>"]
19
20     def get_unk(self):
21         return self.word2index["<UNK>"]
22
23     def has_word(self, word) -> bool:
24         return word in self.word2index
25
26     def get_index(self, word):
27         if word in self.word2index:
28             return self.word2index[word]
29         return self.get_unk()
30
31     def get_word(self, index):
32         return self.index2word[index]
33
34     def size(self):
35         return len(self.index2word)
36
37     def is_empty(self):
38         empty_size = 4
39         return self.size() <= empty_size
40
41     def reset(self):
42         self.word2count = Counter()
43         self.index2word = ["<pad>", "<sos>", "<eos>", "<unk>"]
44         self.word2index = {word: index for index, word in enumerate(self.index2w

1 vocabulary = Vocabulary(bpe_processor)
2 vocabulary.size()

30000

1 from rouge import Rouge
2 import razdel
3
4 def add_oracle_summary_to_records(records, max_sentences=30, lower=True, nrow=1
5     rouge = Rouge()
6     sentences_ = []

```

```

7     oracle_sentences_ = []
8     oracle_summary_ = []
9     if nrows is not None:
10         records = records.iloc[:nrows].copy()
11     else:
12         records = records.copy()
13
14     for text, summary in tqdm(records[['text', 'summary']].values):
15         summary = summary.lower() if lower else summary
16         sentences = [sentence.text.lower() if lower else sentence.text for sentence in sentences_]
17         oracle_summary, sentences_indicies = build_oracle_summary_greedy(text, sentences, lower=lower)
18
19         sentences_ += [sentences]
20         oracle_sentences_ += [list(sentences_indicies)]
21         oracle_summary_ += [oracle_summary]
22     records['sentences'] = sentences_
23     records['oracle_sentences'] = oracle_sentences_
24     records['oracle_summary'] = oracle_summary_
25     return records
26
27 ext_train_records = add_oracle_summary_to_records(train_records, nrows=30000)
28 ext_val_records = add_oracle_summary_to_records(val_records, nrows=None)
29 ext_test_records = add_oracle_summary_to_records(test_records, nrows=None)

```

/usr/local/lib/python3.6/dist-packages/ipykernel_launcher.py:14: TqdmDeprecat

This function will be removed in tqdm==5.0.0

Please use `tqdm.notebook.tqdm` instead of `tqdm.tqdm_notebook`

100% 30000/30000 [29:04<00:00, 17.19it/s]

100% 5265/5265 [06:01<00:00, 14.58it/s]

100% 5770/5770 [05:59<00:00, 16.03it/s]

Используй pickle для сохранения записей, чтобы потом не пересоздавать их потом.

Если решаешь задание в колабе, можешь подключить свой гугл диск и сохранить данные в нём.

```

1
2 from google.colab import drive
3 drive.mount('/content/drive')

```

Mounted at /content/drive

```

1 import pickle
2
3 with open("train_records.bin", 'wb') as file:
4     pickle.dump(ext_train_records, file)
5

```

```

6 with open("val_records.bin", 'wb') as file:
7     pickle.dump(ext_val_records, file)
8
9 with open("test_records.bin", 'wb') as file:
10     pickle.dump(ext_test_records, file)

1 # pickle_train = open("/content/drive/MyDrive/train_records.bin", "rb")
2 # pickle_test = open("/content/drive/MyDrive/test_records.bin", "rb")
3 # pickle_valid = open("/content/drive/MyDrive/val_records.bin", "rb")
4
5 # ext_test_records = pickle.load(pickle_test)
6 # ext_val_records = pickle.load(pickle_valid)
7 # ext_train_records = pickle.load(pickle_train)

```

▼ (!)

Если надо, поменяйте код генератора датасета и батчевалки

```

1 import random
2 import math
3 import razdel
4 import torch
5 import numpy as np
6 from rouge import Rouge
7
8
9 from torch.utils import data
10
11
12 class ExtDataset(data.Dataset):
13     def __init__(self, records, vocabulary, bpe_processor, lower=True, max_sentences=10):
14         self.records = records
15         self.num_samples = records.shape[0]
16         self.bpe_processor = bpe_processor
17         self.lower = lower
18         self.rouge = Rouge()
19         self.vocabulary = vocabulary
20         self.max_sentences = max_sentences
21         self.max_sentence_length = max_sentence_length
22         self.device = device
23
24     def __len__(self):
25         return self.records.shape[0]
26
27     def __getitem__(self, idx):
28         cur_record = self.records.iloc[idx]
29         inputs = list(map(lambda x: x[:self.max_sentence_length], self.bpe_processor(cur_record['inputs'])))
30         outputs = [int(i in cur_record['oracle_sentences']) for i in range(len(inputs))]
31         return {'inputs': inputs, 'outputs': outputs}

```

```
1 # Это батчевалка
```

```
2 def collate_fn(records):
```

```

2 def collate_fn(records):
3     max_length = max(len(sentence) for record in records for sentence in record['inputs'])
4     max_sentences = max(len(record['outputs']) for record in records)
5
6     new_inputs = torch.zeros((len(records), max_sentences, max_length))
7     new_outputs = torch.zeros((len(records), max_sentences))
8     for i, record in enumerate(records):
9         for j, sentence in enumerate(record['inputs']):
10             new_inputs[i, j, :len(sentence)] += np.array(sentence)
11             new_outputs[i, :len(record['outputs'])] += np.array(record['outputs'])
12     return {'features': new_inputs.type(torch.LongTensor), 'targets': new_outputs.type(torch.LongTensor)}

```

```

1 import numpy as np
2
3 import torch
4 import torch.nn as nn
5 import torch.nn.functional as F
6 import torch.optim as optim
7
8 from torch.nn.utils.rnn import pack_padded_sequence as pack
9 from torch.nn.utils.rnn import pad_packed_sequence as unpack
10
11
12 class YourSentenceEncoder(nn.Module):
13     # Место для вашего Sentence Encoder-а. Разрешается использовать любые методы, которые вы знаете.
14     def __init__(self, input_size, embedding_dim, hidden_size, n_layers=3, dropout=0.5, bidirectional=False):
15         super().__init__()
16
17         num_directions = 2 if bidirectional else 1
18         assert hidden_size % num_directions == 0
19         hidden_size = hidden_size // num_directions
20
21         self.embedding_dim = embedding_dim
22         self.input_size = input_size
23         self.hidden_size = hidden_size
24         self.n_layers = n_layers
25         self.dropout = dropout
26         self.bidirectional = bidirectional
27
28         self.embedding_layer = nn.Embedding(input_size, embedding_dim)
29         self.rnn_layer = nn.LSTM(embedding_dim, hidden_size, n_layers, dropout=dropout, batch_first=True)
30         self.dropout_layer = nn.Dropout(dropout)
31
32     def forward(self, inputs, hidden=None):
33         embedded = self.embedding_layer(inputs)
34         outputs, _ = self.rnn_layer(embedded, hidden)
35         sentences_embeddings = torch.mean(outputs, 1)
36
37         return sentences_embeddings
38
39
40 class SentenceTaggerRNN(nn.Module):
41     def __init__(self,
42                 vocabulary_size,
43                 token_embedding_dim=256,

```



```

44         sentence_encoder_hidden_size=256,
45         hidden_size=256,
46         bidirectional=True,
47         sentence_encoder_n_layers=2,
48         sentence_encoder_dropout=0.3,
49         sentence_encoder_bidirectional=True,
50         n_layers=1,
51         dropout=0.3):
52     super(SentenceTaggerRNN, self).__init__()
53
54     num_directions = 2 if bidirectional else 1
55     assert hidden_size % num_directions == 0
56     hidden_size = hidden_size // num_directions
57
58     self.hidden_size = hidden_size
59     self.n_layers = n_layers
60     self.dropout = dropout
61     self.bidirectional = bidirectional
62
63     # Your sentence encoder model
64     self.sentence_encoder = YourSentenceEncoder(vocabulary_size, token_embed
65                                                  sentence_encoder_hidden_size,
66                                                  sentence_encoder_dropout, ser
67
68     self.rnn_layer = nn.LSTM(
69         sentence_encoder_hidden_size,
70         hidden_size,
71         n_layers,
72         dropout=dropout,
73         bidirectional=bidirectional,
74         batch_first=True)
75     self.dropout_layer = nn.Dropout(dropout)
76     self.content_linear_layer = nn.Linear(hidden_size * 2, 1)
77     self.document_linear_layer = nn.Linear(hidden_size * 2, hidden_size * 2)
78     self.saliency_linear_layer = nn.Linear(hidden_size * 2, hidden_size * 2)
79     self.tanh_layer = nn.Tanh()
80
81     def forward(self, inputs, hidden=None):
82         batch_size = inputs.size(0)
83         sentences_count = inputs.size(1)
84         tokens_count = inputs.size(2)
85         inputs = inputs.reshape(-1, tokens_count)
86         embedded_sentences = self.sentence_encoder(inputs)
87         embedded_sentences = embedded_sentences.reshape(batch_size, sentences_co
88         outputs, _ = self.rnn_layer(embedded_sentences, hidden)
89         outputs = self.dropout_layer(outputs)
90         document_embedding = self.tanh_layer(self.document_linear_layer(torch.m
91         content = self.content_linear_layer(outputs).squeeze(2)
92         saliency = torch.bmm(outputs, self.saliency_linear_layer(document_embedd
93         return content + saliency
94
95 model = SentenceTaggerRNN(vocabulary.size())

```

/usr/local/lib/python3.6/dist-packages/torch/nn/modules/rnn.py:61: UserWarnin

dropout option adds dropout after all but last recurrent layer, so non-zero d:

▼ Обучение

```

1 device = torch.device('cuda')
2
3 loaders = {
4     'train': data.DataLoader(
5         ExtDataset(
6             ext_train_records,
7             vocabulary,
8             bpe_processor=bpe_processor
9         ),
10        batch_size=64,
11        collate_fn=collate_fn
12    ),
13    'valid': data.DataLoader(
14        ExtDataset(
15            ext_val_records,
16            vocabulary,
17            bpe_processor=bpe_processor
18        ),
19        batch_size=64,
20        collate_fn=collate_fn
21    ),
22    'test': data.DataLoader(
23        ExtDataset(
24            ext_test_records,
25            vocabulary,
26            bpe_processor=bpe_processor
27        ),
28        batch_size=64,
29        collate_fn=collate_fn
30    ),
31 }
32
33 lr = 1e-4
34 num_epochs = 5
35
36 optimizer = torch.optim.Adam(model.parameters(), lr=lr)
37 criterion = nn.BCEWithLogitsLoss()
38 # Maybe adding scheduler?
39
40
41 from tqdm.notebook import trange, tqdm
42
43 def train():
44     model.to(device)
45     pbar_loader = tqdm(len(loaders["train"]) + len(loaders["valid"]), desc=f"Total steps")
46     for e in trange(num_epochs, desc="Epoch"):
47         train_loss = 0
48         valid_loss = 0

```

```
9     valid_loss = 0
10     train_it = 0
11     valid_it = 0
12
13     model.train()
14     for batch in loaders["train"]:
15         features = batch["features"].to(device)
16         targets = batch["targets"].to(device)
17
18         logits = model(features)
19
20         loss = criterion(logits, targets)
21         train_loss += loss.item()
22         train_it += 1
23
24         optimizer.zero_grad()
25         loss.backward()
26         optimizer.step()
27         # Maybe adding scheduler?
28
29         pbar_loader.update()
30         pbar_loader.set_description(
31             f"Train Loss: {train_loss / train_it:.3}, Valid Loss: {0}"
32         )
33
34     model.eval()
35     with torch.no_grad():
36         for batch in loaders["valid"]:
37             features = batch["features"].to(device)
38             targets = batch["targets"].to(device)
39
40             logits = model(features)
41
42             loss = criterion(logits, targets)
43             valid_loss += loss.item()
44             valid_it += 1
45
46             pbar_loader.update()
47             pbar_loader.set_description(
48                 f"Train Loss: {train_loss / train_it:.3}, "
49                 f" Valid Loss: {valid_loss / valid_it:.3}"
50             )
51     print(
52         f"Epoch {e}; Train Loss: {train_loss / train_it:.3}, "
53         f" Valid Loss: {valid_loss / valid_it:.3}"
54     )
55     pbar_loader.reset()

```

1 train()

Train Loss: 0.178, Valid Loss: 0.179: 0%

0/552 [00:00<05:04, 1.81it/s]

Epoch: 100%

5/5 [36:18<00:00, 435.65s/it]

Epoch 0; Train Loss: 0.204, Valid Loss: 0.183

Epoch 1; Train Loss: 0.189, Valid Loss: 0.181

```

1 device = torch.device("cuda")
2
3 top_k = 3
4
5 def postprocess(ref, hyp, is_multiple_ref=False, detokenize_after=False, tokenize_after=False):
6     if is_multiple_ref:
7         reference_sents = ref.split(" s_s ")
8         decoded_sents = hyp.split("s_s")
9         hyp = [w.replace("<", "&lt;").replace(">", "&gt;").strip() for w in decoded_sents]
10        ref = [w.replace("<", "&lt;").replace(">", "&gt;").strip() for w in reference_sents]
11        hyp = " ".join(hyp)
12        ref = " ".join(ref)
13    ref = ref.strip()
14    hyp = hyp.strip()
15    if detokenize_after:
16        hyp = punct_detokenize(hyp)
17        ref = punct_detokenize(ref)
18    if tokenize_after:
19        hyp = hyp.replace("@@UNKNOWN@", "<unk>")
20        hyp = " ".join([token.text for token in razdel.tokenize(hyp)])
21        ref = " ".join([token.text for token in razdel.tokenize(ref)])
22    return ref, hyp
23
24 references = []
25 predictions = []
26
27 model.eval()
28 for num, batch in enumerate(loaders["test"]):
29
30     logits = model(batch["features"].to(device))
31     in_summary = torch.argsort(logits, dim=1)[ :, -top_k:]
32     for i in range(len(batch['targets'])):
33
34         summary = ext_test_records.iloc[i]['summary']
35         summary = summary.lower()
36         predicted_summary = ' '.join([ext_test_records.iloc[i]['sentences'][idx] for idx in in_summary])
37         summary, predicted_summary = postprocess(summary, predicted_summary)
38
39         references.append(summary)
40         predictions.append(predicted_summary)
41
42 calc_scores(references, predictions)

```

Count: 5770

Ref: следователи возбудили дело против жительницы сургута , подозреваемой в убийстве 49-летнего

Hyp: жительницу сургута подозревают в убийстве 49-летнего сожителя . по версии ведомства ,

BLEU: 0.4482006046197879

ROUGE: {'rouge-1': {'f': 0.30057638825676386, 'p': 0.28523669319955836, 'r':

