



Deep Learning School

Физтех-Школа Прикладной математики и информатики (ФПМИ) МФТИ

Text Summarization

Всем привет! Сегодня мы познакомимся с задачей суммаризации текста на примере генерации "сжатых" новостей. Рассмотрим некоторые базовые решения и познакомимся с архитектурами нейросетей для решения задачи. Датасет: gazeta.ru

Ноутбук создан на основе семинара Гусева Ильи на кафедре компьютерной лингвистики ABBYY МФТИ.

Загрузим датасет и необходимые библиотеки

```
1 !wget -q https://www.dropbox.com/s/431702z5a5i2w8j/gazeta_train.txt
2 !wget -q https://www.dropbox.com/s/k2egtsugohnb185/gazeta_val.txt
3 !wget -q https://www.dropbox.com/s/3gk15n5djs9wv6/gazeta_test.txt
```

```
1 !pip -q install razdel networkx pymorphy2 nltk rouge==0.3.1 summa
```

```
[██████████] 61KB 3.1MB/s
[██████████] 61KB 5.7MB/s
[██████████] 8.2MB 7.6MB/s
Building wheel for summa (setup.py) ... done
```

Dataset

Посмотрим на то, как устроен датасет

```
1 !head -n 1 gazeta_train.txt
2 !cat gazeta_train.txt | wc -l
3 !cat gazeta_val.txt | wc -l
4 !cat gazeta_test.txt | wc -l

{"url": "https://www.gazeta.ru/financial/2011/11/30/3852658.shtml", "text": "<По итогам 2011 года чистый отток может составить примерно $80 млрд, в следующем году – около $20 млрд. При этом мы ожидаем, что со второго полугодия 2012 52400
5265
5770
```

```
1 import json
2 import random
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(json.loads(line))
10    if sort_by_date:
11        records.sort(key=lambda x: x["date"])
12    if shuffle:
13        random.shuffle(records)
14    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")

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

Extractive RNN

BPE

Для начала сделаем BPE токенизацию

```
1 !pip install youtokentome

Collecting youtokentome
  Downloading https://files.pythonhosted.org/packages/c8/1c/224cdc3d9a32ed706c8fb1f30h491he6ea5da114ff4edc174014cc24fa43/youtokentome-1.0.6-cp37m-manylinux2010_x86_64.whl (1.7MB)
   ██████████| 1.7MB 4.2MB/s
Requirement already satisfied: Click>=7.0 in /usr/local/lib/python3.7/dist-packages (from youtokentome) (7.1.2)
Installing collected packages: youtokentome
Successfully installed youtokentome-1.0.6

1 import youtokentome as ytym
2
3 def train_bpe(records, model_path, model_type="bpe", vocab_size=10000, lower=True):
4     temp_file_name = "temp.txt"
5     with open(temp_file_name, "w") as ftemp:
```

```

5     with open(temp_file_name, 'w') as temp:
6         for record in records:
7             text, summary = record['text'], record['summary']
8             if lower:
9                 summary = summary.lower()
10                text = text.lower()
11            if not text or not summary:
12                continue
13            temp.write(text + "\n")
14            temp.write(summary + "\n")
15    ytmm.BPE.train(data=temp_file_name, vocab_size=vocab_size, model=model_path)
16
17 train_bpe(train_records, "BPE_model.bin")

```

Словарь

Составим словарь для индексации токенов

```

1 bpe_processor = ytmm.BPE('BPE_model.bin')
2 vocabulary = bpe_processor.vocab()

```

Кэш oracle summary

Закэшируем oracle summary, чтобы не пересчитывать их каждый раз

```

1 from rouge import Rouge
2 import razdel
3 from tqdm.notebook import tqdm
4
5 import copy
6
7 def build_oracle_summary_greedy(text, gold_summary, calc_score, lower=True, max_sentences=30):
8     ...
9     # Жадное построение oracle summary
10    ...
11    gold_summary = gold_summary.lower() if lower else gold_summary
12    # Делим текст на предложения
13    sentences = [sentence.text.lower() if lower else sentence.text for sentence in razdel.sentenceize(text)][::max_sentences]
14    n_sentences = len(sentences)
15    oracle_summary_sentences = set()
16
17    score = -1.0
18    summaries = []
19    for _ in range(n_sentences):
20        for i in range(n_sentences):
21            if i in oracle_summary_sentences:
22                continue
23            current_summary_sentences = copy.copy(oracle_summary_sentences)
24            # Добавляем какое-то предложение к уже существующему summary
25            current_summary_sentences.add(i)
26            current_summary = " ".join([sentences[index] for index in sorted(list(current_summary_sentences))])
27            # Считаем метрики
28            current_score = calc_score(current_summary, gold_summary)
29            summaries.append((current_score, current_summary))
30
31    # Если получилось улучшить метрики с добавлением какого-либо предложения, то пробуем добавить ещё
32    # Иначе на этом заканчиваем
33    best_summary_score, best_summary_sentences = max(summaries)
34    if best_summary_score <= score:
35        break
36    oracle_summary_sentences = best_summary_sentences
37    score = best_summary_score
38
39    oracle_summary = " ".join([sentences[index] for index in sorted(list(oracle_summary_sentences))])
39
40 def calc_single_score(pred_summary, gold_summary, rouge):
41     return rouge.get_scores([pred_summary], [gold_summary], avg=True)['rouge-2']['f']
42
43 def add_oracle_summary_to_records(records, max_sentences=30, lower=True, nrows=1000):
44     rouge = Rouge()
45     for i, record in tqdm(enumerate(records)):
46         if i >= nrows:
47             break
48         text = record["text"]
49         summary = record["summary"]
50
51         summary = summary.lower() if lower else summary
52         sentences = [sentence.text.lower() if lower else sentence.text for sentence in razdel.sentenceize(text)][::max_sentences]
53         oracle_summary, sentences_indices = build_oracle_summary_greedy(text, summary, calc_score=lambda x, y: calc_single_score(x, y,
54                                         lower=lower, max_sentences=max_sentences))
55         record["sentences"] = sentences
56         record["oracle_sentences"] = list(sentences_indices)
57         record["oracle_summary"] = oracle_summary
58
59     return records[:nrows]
60
61 ext_train_records = add_oracle_summary_to_records(train_records, nrows=2048)
62 ext_val_records = add_oracle_summary_to_records(val_records, nrows=256)
63 ext_test_records = add_oracle_summary_to_records(test_records, nrows=256)

2048? [05:58<00:00, 2.75us]
256? [00:53<00:00, 3.70us]
256? [00:57<00:00, 3.14us]

```

Составление батчей

```

1 import torch
2 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')

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 class BatchIterator():
10     def __init__(self, records, vocabulary, batch_size, bpe_processor, shuffle=True, lower=True, max_sentences=30, max_sentence_length=100):
11         self.records = records
12         self.num_samples = len(records)
13         self.batch_size = batch_size
14         self.bpe_processor = bpe_processor
15         self.shuffle = shuffle
16         self.batches_count = int(math.ceil(self.num_samples / batch_size))
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.batches_count
26
27     def __iter__(self):
28         indices = np.arange(self.num_samples)
29         if self.shuffle:
30             np.random.shuffle(indices)
31
32         for start in range(0, self.num_samples, self.batch_size):
33             end = min(start + self.batch_size, self.num_samples)
34             batch_indices = indices[start:end]
35
36             batch_inputs = []
37             batch_outputs = []
38             max_sentence_length = 0
39             max_sentences = 0
40             batch_records = []
41
42             for data_ind in batch_indices:
43
44                 record = self.records(data_ind)
45                 batch_records.append(record)
46                 text = record["text"]
47                 summary = record["summary"]
48                 summary = summary.lower() if self.lower else summary
49
50                 if "sentences" not in record:
51                     sentences = [sentence.text.lower() if self.lower else sentence.text for sentence in razdel.sentenceize(text)][::self.#
52                 else:
53                     sentences = record["sentences"]
54                     max_sentences = max(len(sentences), max_sentences)
55
56                     # номера предложений, которые в нашем сэмпле
57                     if "oracle_sentences" not in record:
58                         calc_score = lambda x, y: calc_single_score(x, y, self.rouge)
59                         sentences_indices = build_oracle_summary_greedy(text, summary, calc_score=calc_score, lower=self.lower, max_sente#
60                     else:
61                         sentences_indices = record["oracle_sentences"]
62
63                     # inputs - индексы слов в предложении
64                     inputs = [bpe_processor.encode(sentence)[::self.max_sentence_length] for sentence in sentences]
65                     max_sentence_length = max(max_sentence_length, max([len(tokens) for tokens in inputs]))
66
67                     # получение метки класса предложения
68                     outputs = [int(i in sentences_indices) for i in range(len(sentences))]
69                     batch_inputs.append(inputs)
70                     batch_outputs.append(outputs)
71
72                     tensor_inputs = torch.zeros((self.batch_size, max_sentences, max_sentence_length), dtype=torch.long, device=self.device)
73                     # we add index 2 for padding
74                     # YOUR CODE
75                     tensor_outputs = torch.full((self.batch_size, max_sentences), fill_value=2.0, dtype=torch.float32, device=self.device)
76
77                     for i, inputs in enumerate(batch_inputs):
78                         for j, sentence_tokens in enumerate(inputs):
79                             tensor_inputs[i][j][:len(sentence_tokens)] = torch.tensor(sentence_tokens, dtype=torch.int64)
80
81                     for i, outputs in enumerate(batch_outputs):
82                         tensor_outputs[i][:len(outputs)] = torch.LongTensor(outputs)
83
84                     tensor_outputs = tensor_outputs.long()
85
86                     yield {
87                         'inputs': tensor_inputs,
88                         'outputs': tensor_outputs,
89                         'records': batch_records
90                     }
91
92 train_iterator = BatchIterator(ext_train_records, vocabulary, 32, bpe_processor, device=device)
93 val_iterator = BatchIterator(ext_val_records, vocabulary, 32, bpe_processor, device=device)
94 test_iterator = BatchIterator(ext_test_records, vocabulary, 32, bpe_processor, device=device)

```

Extractor - SummaRuNNer

<https://arxiv.org/pdf/1611.04230.pdf>

Homework

- В данной реализации в outputs в качестве padding используется индекс 0. Измените в функции __iter__ индекс padding, чтобы он не совпадал с классом 0 или 1, например, 2.
- В качестве criterion используйте CrossEntropyLoss вместо BCEWithLogitsLoss
- Из-за смены criterion, вы уже должны подавать на вход criterion ни одно число, а logits для каждого класса. Перед подачей logits вы можете отфильтровать предсказания для класса padding. В этом пункте вам придется изменять файл train_model.py, а именно функции train и evaluate.
- Используйте два варианта обучения: с весами в CrossEntropyLoss и без
- Также сравните inference, когда вы ранжируете logits, и когда вы выбираете предложения, у которых logits > 0, в двух вариантах обучения.
- Реализуйте дополнительно характеристику предложения novelty. Как влияет добавление novelty на качество summary?
- Постарайтесь улучшить качество модели, полученной на семинаре: $BLEU \approx 0.45$

```

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 from torch.autograd import Variable
8
9 from torch.nn.utils.rnn import pack_padded_sequence as pack
10 from torch.nn.utils.rnn import pad_packed_sequence as unpack
11
12 class SentenceEncoderRNN(nn.Module):
13     def __init__(self, input_size, embedding_dim, hidden_size, n_layers=3, dropout=0.3, bidirectional=True):
14         super().__init__()
15
16         num_directions = 2 if bidirectional else 1
17         assert hidden_size % num_directions == 0
18         hidden_size = hidden_size // num_directions
19
20         self.embedding_dim = embedding_dim
21         self.input_size = input_size
22         self.hidden_size = hidden_size
23         self.n_layers = n_layers
24         self.dropout = dropout
25         self.bidirectional = bidirectional
26
27         self.embedding_layer = nn.Embedding(input_size, embedding_dim)
28         self.rnn_layer = nn.LSTM(embedding_dim, hidden_size, n_layers, dropout=dropout, bidirectional=bidirectional, batch_first=True)
29         self.dropout_layer = nn.Dropout(dropout)
30
31     def forward(self, inputs, hidden=None):

```

```

32     embedded = self.dropout_layer(self.embedding_layer(inputs))
33     outputs, _ = self.rnn_layer(embedded, hidden)
34     sentences_embeddings = torch.mean(outputs, 1)
35     # [batch_size, hidden_size]
36     return sentences_embeddings
37
38 class SentenceTaggerRNN(nn.Module):
39     def __init__(self,
40         vocabulary_size,
41         use_content=True,
42         use_saliency=True,
43         use_novelty=True,
44         token_embedding_dim=128,
45         sentence_encoder_hidden_size=256,
46         hidden_size=256,
47         bidirectional=True,
48         sentence_encoder_n_layers=2,
49         sentence_encoder_dropout=0.3,
50         sentence_encoder_bidirectional=True,
51         n_layers=2,
52         dropout=0.3):
53
54     super().__init__()
55
56     num_directions = 2 if bidirectional else 1
57     assert hidden_size % num_directions == 0
58     hidden_size = hidden_size // num_directions
59
60     self.hidden_size = hidden_size
61     self.n_layers = n_layers
62     self.dropout = dropout
63     self.bidirectional = bidirectional
64
65     self.sentence_encoder = SentenceEncoderRNN(vocabulary_size,
66                                                 token_embedding_dim,
67                                                 sentence_encoder_hidden_size,
68                                                 sentence_encoder_n_layers,
69                                                 sentence_encoder_dropout,
70                                                 sentence_encoder_bidirectional)
71
72     self.rnn_layer = nn.LSTM(sentence_encoder_hidden_size, hidden_size, n_layers, dropout=dropout,
73                             bidirectional=bidirectional, batch_first=True)
74
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.salience_linear_layer = nn.Linear(hidden_size * 2, hidden_size * 2)
79     self.novelty_linear_layer = nn.Linear(hidden_size * 2, hidden_size * 2)
80     self.tanh_layer = nn.Tanh()
81
82     self.use_content = use_content
83     self.use_saliency = use_saliency
84     self.use_novelty = use_novelty
85
86     def forward(self, inputs, hidden=None):
87         # parameters of the probability
88         content = 0
89         salience = 0
90         novelty = 0
91
92         # [batch_size, seq num, seq_len]
93         batch_size = inputs.size(0)
94         sentences_count = inputs.size(1)
95         tokens_count = inputs.size(2)
96         inputs = inputs.reshape(-1, tokens_count)
97         # [batch_size * seq num, seq_len]
98
99         embedded_sentences = self.sentence_encoder(inputs)
100        embedded_sentences = self.dropout_layer(embedded_sentences.reshape(batch_size, sentences_count, -1))
101        # [batch_size * seq num, seq_len, hidden_size] -> [batch_size, seq num, hidden_size]
102
103        outputs, _ = self.rnn_layer(embedded_sentences, hidden)
104        # [batch_size, seq num, hidden_size]
105
106        # W * h^T
107        if self.use_content:
108            content = self.content_linear_layer(outputs).squeeze(2) # 1-representation
109            # [batch_size, seq num]
110
111            # h^T * W * d
112            if self.use_saliency:
113                salience = torch.bmm(outputs, self.salience_linear_layer(document_embedding).unsqueeze(2)).squeeze(2) # 2-representation
114                # [batch_size, seq num, hidden_size] * [batch_size, hidden_size, 1] = [batch_size, seq num, ]
115
116            if self.use_novelty:
117                # at every step add novelty to prediction of the sentence
118                predictions = content + salience
119
120                # 0) initialize summary_representation and novelty by zeros
121                # YOUR CODE
122                summary_representation = torch.zeros(batch_size, self.hidden_size * 2).to(device) #[batch_size, hidden_size]
123                novelty = torch.zeros(batch_size, sentences_count)
124
125                for sentence_num in range(sentences_count):
126
127                    # 1) take sentence_num_state from outputs(representation of the sentence with number sentence_num)
128                    # 2) calculate novelty for current sentence
129                    # 3) add novelty to predictions
130                    # 4) calculate probability for current sentence
131                    # 5) add sentence_num_state with the weight which is equal to probability to summary_representation
132
133                    # YOUR CODE
134                    sentence_num_state = outputs[:, sentence_num, :]
135
136                    novelty = torch.bmm(outputs, self.novelty_linear_layer(torch.tanh(summary_representation)).unsqueeze(2)).squeeze(2) #2
137                    # [batch_size, seq num, hidden_size] * [batch_size, hidden_size, 1] = [batch_size, seq num, ]
138
139                    predictions[:, sentence_num] += novelty[:, sentence_num] #3
140
141                    probability = torch.sigmoid(predictions[:, sentence_num].unsqueeze(1)) #4
142
143                    summary_representation += probability * sentence_num_state #5
144
145        return content + salience + novelty

```

▼ Model

$$P(y_j = 1 \mid \mathbf{h}_j, \mathbf{s}_j, \mathbf{d}) = \sigma(W_c \mathbf{h}_j + \mathbf{h}_j^T W_s \mathbf{d})$$

```

1 !#!gdown https://drive.google.com/uc?id=1Mi5_iczAlcyF7zGDPY6niyeD82P0_PBH -o train_model.py
2
3 from google.colab import drive
4 drive.mount('/content/drive')
5
6 !cp /content/drive/MyDrive/hw_text_summarization/train_model.py .
7

```

```

8 import train_model
9 import imp
10 imp.reload(train_model)
11 from train_model import train_with_logs
12
13 Mounted at /content/drive

1 all_counts = torch.zeros(3)
2 for batch in iter(train_iterator):
3     all_counts += batch['outputs'].view(-1).unique(return_counts=True)[1].cpu()
4
5 class_0_count = all_counts[0]
6 class_1_count = all_counts[1]
7
8 print("class_0_count:", class_0_count)
9 print("class_1_count:", class_1_count)
10
11 weights = torch.tensor([class_1_count / (class_0_count + class_1_count), class_0_count / (class_0_count + class_1_count)])
12 weights = weights.to(device)
13 weights

    class_0_count= tensor(53811.)
    class_1_count= tensor(4620.)
    tensor([0.0791, 0.9209], device='cuda:0')

1 N_EPOCHS = 20
2 CLIP = 1
3
4 def train_all(use_class_weights, N_EPOCHS, CLIP, lr=1e-3):
5     optimizer = optim.Adam(model.parameters(), lr)
6     if use_class_weights:
7         # weights depend on the number of objects of class 0 and 1
8         criterion = nn.CrossEntropyLoss(weight=weights, ignore_index=2)
9     else:
10        criterion = nn.CrossEntropyLoss(ignore_index=2)
11     train_with_logs(model, train_iterator, val_iterator, optimizer, criterion, N_EPOCHS, CLIP)

```

▼ Inference

```

1 from train_model import punct_detokenize, postprocess
2
3 def top_k_inference_summarunner(model, iterator, top_k=3):
4     print("top_k=", top_k)
5
6     references = []
7     predictions = []
8
9     model.eval()
10    for batch in test_iterator:
11
12        preds = model(batch['inputs'])
13        sum_in = torch.argsort(preds, dim=1)[:, -top_k:]
14
15        for i in range(len(batch['outputs'])):
16
17            summary = batch['records'][i]['summary'].lower()
18            pred_summary = '.join([batch['records'][i][sentences][ind] for ind in sum_in.sort(dim=1)[0][i]])
19
20            summary, pred_summary = postprocess(summary, pred_summary)
21
22            references.append(summary)
23            predictions.append(pred_summary)
24
25    calc_scores(references, predictions)

1 def threshold_inference_summarunner(model, iterator, threshold=0):
2     print("threshold", threshold)
3
4     references = []
5     predictions = []
6
7     model.eval()
8     for batch in test_iterator:
9
10        preds = model(batch['inputs'])
11
12        for i in range(len(batch['outputs'])):
13            sentences = (preds[i] > threshold).nonzero(as_tuple=True)[0].unique()
14
15            summary = batch['records'][i]['summary'].lower()
16            pred_summary = '.join([batch['records'][i][sentences][ind] for ind in sentences if ind < len(batch['records'][i]['sentences'])])
17
18            if len(pred_summary) == 0:
19                pred_summary = '.join(batch['records'][i]['sentences'])
20
21            summary, pred_summary = postprocess(summary, pred_summary)
22
23            references.append(summary)
24            predictions.append(pred_summary)
25
26    calc_scores(references, predictions)

1 vocab_size = len(vocabulary)

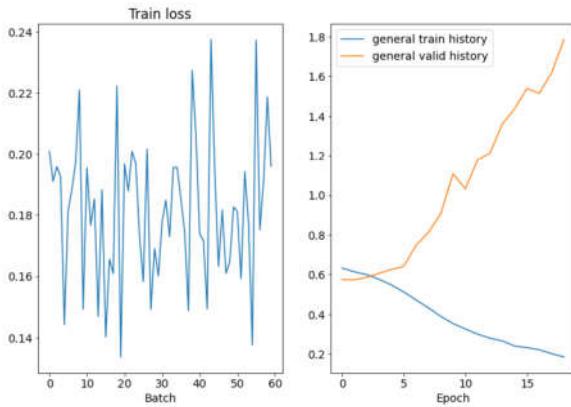
```

▼ Обучим модель без весов и оценим ее качество

```

1 model = SentenceTaggerRNN(vocab_size, use_novelty=False).to(device)
2
3 train_all(False, N_EPOCHS, CLIP)
4
5 model.load_state_dict(torch.load('best-val-model.pt'))
6
7 top_k_inference_summarunner(model, test_iterator, 1)
8 top_k_inference_summarunner(model, test_iterator, 2)
9 top_k_inference_summarunner(model, test_iterator, 3)
10 top_k_inference_summarunner(model, test_iterator, 5)
11 threshold_inference_summarunner(model, test_iterator, 0.2)
12 threshold_inference_summarunner(model, test_iterator, 0.3)
13 threshold_inference_summarunner(model, test_iterator, 0.25)
14 threshold_inference_summarunner(model, test_iterator, 0.4)
15 threshold_inference_summarunner(model, test_iterator, 0.5)

```

Epoch: 28 | Time: 0m 8s
 Train Loss: 0.180 | Train PPL: 1.198
 Val. Loss: 1.905 | Val. PPL: 6.720

top_k= 1
 Count: 256

Ref: неандертальцы хоронили сабреты , соблюдая определенные ритуалы , считают ученые . новая находка позволяет предположить , что неандертальцы создавали подобия могил , а более ранние данные указывают на то , что , возможно , при

Нур: неандертальцы хоронили мертвого с особыми ритуалами , считает команда исследователей под руководством специалистов из кембрийского университета .

BLEU: 0.1974612847722521

ROUGE: {'rouge-1': {'f': 0.2816651701712057, 'p': 0.44265913216514763, 'r': 0.2191164861780385}, 'rouge-2': {'f': 0.14088066542198771, 'p': 0.22634890714244005, 'r': 0.10919055716826959}, 'rouge-l': {'f': 0.2109489572734253, 'p': 0.3

top_k= 2
 Count: 256

Ref: убийство иранского генерала сулеймани , барбарская пила от вильфанди , поражение сборной россии в финале молодежного чемпионата мира по хоккею – что интересовало телеграм-каналы на прошлой неделе .

Нур: начнем с главного события новогодних дней – убийства американцами иранского генерала сулеймани на « стекломой » называет его культовой фигурой шиитского мира . « с его деятельностью связывают практически всю « силовую » политику

BLEU: 0.431322726295193

ROUGE: {'rouge-1': {'f': 0.3087075074591505, 'p': 0.34382112511256513, 'r': 0.2977321270053813}, 'rouge-2': {'f': 0.14450677226244584, 'p': 0.1593672053278618, 'r': 0.1431857683675643}, 'rouge-l': {'f': 0.2610573441648967, 'p': 0.3

top_k= 3
 Count: 256

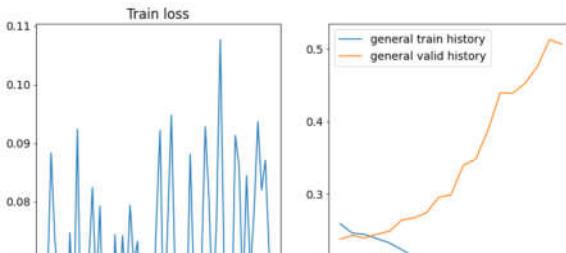
Ref: президент россии владимир путин одобрил идею создания парламентского измерения « нормандского формата ». этот план он обсудил на встрече с лидером украинской партии « оппозиционная платформа – за жизнь » виктором медведчуком .

Нур: президент россии владимир путин поддержал идею создать парламентское измерение « нормандского формата » с участием представителей россии , украины , германии и франции . « конечно , если парламентарии на украине , в россии , в

▼ Обучим модель с novelty и без весов и оценим ее качество

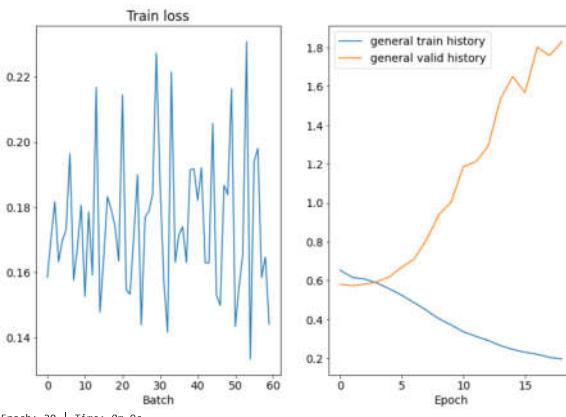
```
Ref: политизацию письма гапки и магаз макла  вопреки стечению моментных гипотетических пишущих  снова пополнила  на этот раз использовать ими же ими же ими же  выпустивший облизанный поклонительский фильм о своей работе

1 model = SentenceTaggerRNN(vocab_size, use_novelty=True).to(device)
2
3 train_all(False, N_EPOCHS, CLIP)
4
5 model.load_state_dict(torch.load('best-val-model.pt'))
6
7 top_k_inference_summarunner(model, test_iterator, 1)
8 top_k_inference_summarunner(model, test_iterator, 2)
9 top_k_inference_summarunner(model, test_iterator, 3)
10 top_k_inference_summarunner(model, test_iterator, 5)
11 threshold_inference_summarunner(model, test_iterator, 0.2)
12 threshold_inference_summarunner(model, test_iterator, 0.3)
13 threshold_inference_summarunner(model, test_iterator, 0.35)
14 threshold_inference_summarunner(model, test_iterator, 0.4)
15 threshold_inference_summarunner(model, test_iterator, 0.5)
```



- Обучим модель с novelty и весами и оценим ее качество

```
1 model = SentenceTaggerRNN(vocab_size, use_novelty=True).to(device)
2
3 train_all(True, N_EPOCHS, CLIP)
4
5 model.load_state_dict(torch.load('best-val-model.pt'))
6
7 top_k_inference_summarunner(model, test_iterator, 1)
8 top_k_inference_summarunner(model, test_iterator, 2)
9 top_k_inference_summarunner(model, test_iterator, 3)
10 top_k_inference_summarunner(model, test_iterator, 5)
11 threshold_inference_summarunner(model, test_iterator, 0.2)
12 threshold_inference_summarunner(model, test_iterator, 0.3)
13 threshold_inference_summarunner(model, test_iterator, 0.35)
14 threshold_inference_summarunner(model, test_iterator, 0.4)
15 threshold_inference_summarunner(model, test_iterator, 0.5)
```



Вывод:

- использование весов улучшает результаты, когда используем threshold-based inference_summarunner
 - использование novelty не улучшает результат
 - максимальное значение BLEU, которого удалось достичь - это 0.47