



# Deep Learning School

Физтех-Школа Прикладной математики и информатики (ФПМИ) МФТИ

---

## ▼ Задание 3

### Классификация текстов

В этом задании вам предстоит попробовать несколько методов, используемых в задаче классификации, а также понять насколько хорошо модель понимает смысл слов и какие слова в примере влияют на результат.

```
1 import pandas as pd
2 import numpy as np
3 import torch
4
5 from torchtext import datasets
6
7 from torchtext.legacy import datasets
8 from torchtext.legacy.data import Field, LabelField, BucketIterator
9
10 from torchtext.vocab import Vectors, GloVe
11
12 import torch.nn as nn
```

```
13 import torch.nn.functional as F
14 import torch.optim as optim
15 import random
16 from tqdm.autonotebook import tqdm
```

В этом задании мы будем использовать библиотеку torchtext. Она довольно проста в использовании и поможет нам сконцентрироваться на задаче, а не на написании DataLoader-a.

```
1 TEXT = Field(sequential=True, lower=True, include_lengths=True) # Поле текста
2 LABEL = LabelField(dtype=torch.float) # Поле метки
```

```
1 SEED = 0xDEAD
2
3 torch.manual_seed(SEED)
4 torch.backends.cudnn.deterministic = True
```

Датасет на котором мы будем проводить эксперименты это комментарии к фильмам из сайта IMDB.

```
1 train, test = datasets.IMDB.splits(TEXT, LABEL) # загрузим датасет
2 train, valid = train.split(random_state=random.seed(SEED)) # разобьем на части
```

```
downloading aclImdb_v1.tar.gz
aclImdb_v1.tar.gz: 100%|██████████| 84.1M/84.1M [00:03<00:00, 24.5MB/s]
```

```
1 TEXT.build_vocab(train)
2 LABEL.build_vocab(train)
```

```
1 device = "cuda" if torch.cuda.is_available() else "cpu"
2
3 train_iter, valid_iter, test_iter = BucketIterator.splits(
4     (train, valid, test),
5     batch_size = 64,
6     sort_within_batch = True,
7     device = device)
```

## ▼ RNN

Для начала попробуем использовать рекуррентные нейронные сети. На семинаре вы познакомились с GRU, вы можете также попробовать LSTM. Можно использовать для классификации как hidden\_state, так и output последнего токена.

```
1 class RNNBaseline(nn.Module):
2     def __init__(self, vocab_size, embedding_dim, hidden_dim, output_dim, n_layers,
3                   bidirectional, dropout, pad_idx):
```

```

4
5     super().__init__()
6
7     self.bidirectional = bidirectional
8     self.dropout = dropout
9
10    self.embedding = nn.Embedding(vocab_size, embedding_dim, padding_idx = pad_idx)
11
12    self.rnn = nn.LSTM(input_size=embedding_dim, hidden_size=hidden_dim,
13                        num_layers=n_layers, dropout=dropout, bidirectional=bidirectional)
14
15    # self.rnn = nn.GRU(input_size=embedding_dim, hidden_size=hidden_dim,
16    #                    num_layers=n_layers, dropout=dropout, bidirectional=bidirectional)
17    if self.bidirectional:
18        self.fc = nn.Linear(2 * hidden_dim, output_dim) # YOUR CODE GOES HERE
19    else:
20        self.fc = nn.Linear(hidden_dim, output_dim)
21
22    def forward(self, text, text_lengths):
23
24        #text = [sent len, batch size]
25
26        embedded = self.embedding(text)
27
28        #embedded = [sent len, batch size, emb dim]
29
30        #pack sequence
31        packed_embedded = nn.utils.rnn.pack_padded_sequence(embedded, text_lengths.cpu().numpy(),
32                                                             batch_first=True)
33
34        # cell arg for LSTM, remove for GRU
35        packed_output, (hidden, cell) = self.rnn(packed_embedded)
36
37        # for gru
38        # packed_output, hidden= self.rnn(packed_embedded)
39        #unpack sequence
40        output, output_lengths = nn.utils.rnn.pad_packed_sequence(packed_output)
41
42        #output = [sent len, batch size, hid dim * num directions]
43        #output over padding tokens are zero tensors
44
45        #hidden = [num layers * num directions, batch size, hid dim]
46        #cell = [num layers * num directions, batch size, hid dim]
47
48        #concat the final forward (hidden[-2,:,:]) and backward (hidden[-1,:,:]) hidden
49        #and apply dropout
50        if self.bidirectional:
51            hidden = torch.cat([hidden[-2,:,:], hidden[-1,:,:]], dim=1) # YOUR CODE GOES HERE
52        else:
53            hidden = hidden[-1,:,:]
54
55        #hidden = [batch size, hid dim * num directions] or [batch_size, hid dim * num
56        hidden = nn.Dropout(p=self.dropout)(hidden)
57
58        return self.fc(hidden)

```

## Поиграйтесь с гиперпараметрами

```
1 vocab_size = len(TEXT.vocab)
2 emb_dim = 100
3 hidden_dim = 256
4 output_dim = 1
5 n_layers = 2
6 bidirectional = False
7 dropout = 0
8 PAD_IDX = TEXT.vocab.stoi[TEXT.pad_token]
9 patience=3
```

```
1 model = RNNBaseline(
2     vocab_size=vocab_size,
3     embedding_dim=emb_dim,
4     hidden_dim=hidden_dim,
5     output_dim=output_dim,
6     n_layers=n_layers,
7     bidirectional=bidirectional,
8     dropout=dropout,
9     pad_idx=PAD_IDX
10 )
```

```
1 model = model.to(device)
```

```
1 opt = torch.optim.Adam(model.parameters())
2 loss_func = nn.BCEWithLogitsLoss(reduction='sum', )
3
4 max_epochs = 20
```

Обучите сетку! Используйте любые вам удобные инструменты, Catalyst, PyTorch Lightning или свои велосипеды.

```
1 %time
2 import numpy as np
3
4 min_loss = np.inf
5
6 cur_patience = 0
7
8 for epoch in range(1, max_epochs + 1):
9     train_loss = 0.0
10    model.train()
11    pbar = tqdm(enumerate(train_iter), total=len(train_iter), leave=False)
12    pbar.set_description(f"Epoch {epoch}")
13    for it, batch in pbar:
14        #YOUR CODE GOES HERE
15        opt.zero_grad()
16        output = model(batch.text[0].to(device), batch.text[1].to(device))
```

```

17     loss = loss_func(output.squeeze(1), batch.label)
18     loss.backward()
19     train_loss += loss.item()
20     opt.step()
21
22     train_loss /= len(train_iter)
23     val_loss = 0.0
24     model.eval()
25     pbar = tqdm(enumerate(valid_iter), total=len(valid_iter), leave=False)
26     pbar.set_description(f"Epoch {epoch}")
27     for it, batch in pbar:
28         # YOUR CODE GOES HERE
29         output = model(batch.text[0], batch.text[1])
30         val_loss = loss_func(output.squeeze(1), batch.label).item()
31
32
33     val_loss /= len(valid_iter)
34     if val_loss < min_loss:
35         min_loss = val_loss
36         best_model = model.state_dict()
37     else:
38         cur_patience += 1
39         if cur_patience == patience:
40             cur_patience = 0
41             break
42
43     print('Epoch: {}, Training Loss: {}, Validation Loss: {}'.format(epoch, train_loss,
44 model.load_state_dict(best_model)

CPU times: user 2 µs, sys: 0 ns, total: 2 µs
Wall time: 5.25 µs
Epoch: 1, Training Loss: 43.59233463593643, Validation Loss: 0.07081056853472176
Epoch: 2, Training Loss: 40.25798383420401, Validation Loss: 0.06330045603089414
Epoch: 3, Training Loss: 30.986384597137896, Validation Loss: 0.057714442075309104
Epoch: 4, Training Loss: 20.135748187990956, Validation Loss: 0.016463057469513456
Epoch: 5, Training Loss: 12.470059507084589, Validation Loss: 0.007022834935430753
Epoch: 6, Training Loss: 7.142795572968295, Validation Loss: 0.0020590225013635928
Epoch: 7, Training Loss: 7.073025200706329, Validation Loss: 0.010574527716232558
Epoch: 8, Training Loss: 6.931653290116874, Validation Loss: 0.0014015540985737817
Epoch: 9, Training Loss: 1.879023780053767, Validation Loss: 0.004321309974638082
Epoch: 10, Training Loss: 0.8456091671692629, Validation Loss: 0.00018960947833829007
<All keys matched successfully>

```

Посчитайте f1-score вашего классификатора на тестовом датасете.

**Ответ:**

```

1 pred_labels = []
2 true_labels = []
3 for batch in test_iter:
4     pred_labels_batch = list((model(batch.text[0],batch.text[1]) > 0.5).float().cpu()).r
5     true_labels_batch = list(batch.label.cpu())
6

```

```

7     pred_labels += pred_labels_batch
8     true_labels += true_labels_batch

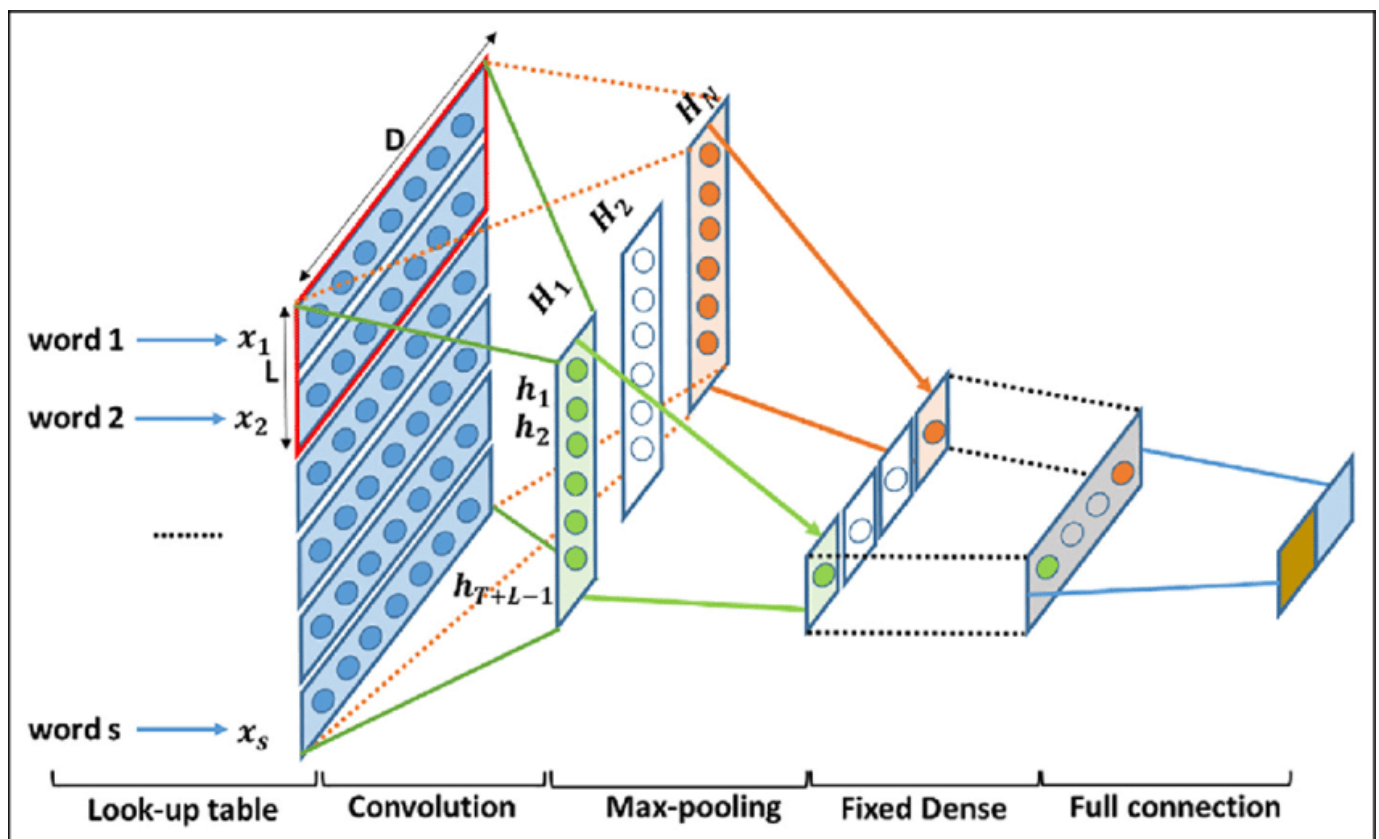
1 from sklearn.metrics import f1_score
2
3 f1_score(true_labels, pred_labels)

0.7941511251203922

```

lstm with bi - 0.8318078698271141  
lstm without bi - 0.7941511251203922`  
gru with bi 0.8255858138062064  
gru without bi 0.5913393537155914

## ▼ CNN



Для классификации текстов также часто используют сверточные нейронные сети. Идея в том, что как правило sentiment содержат словосочетания из двух-трех слов, например "очень хороший фильм" или "невероятная скука". Проходясь сверткой по этим словам мы получим какой-то большой скор и выхватим его с помощью MaxPool. Далее идет обычная полносвязная сетка. Важный момент: свертки применяются не последовательно, а параллельно. Давайте попробуем!

```

1 TEXT = Field(sequential=True, lower=True, batch_first=True) # batch_first тк мы исполь
2 LABEL = LabelField(batch_first=True, dtype=torch.float)
3
4 train, tst = datasets.IMDB.splits(TEXT, LABEL)

```

```

5 trn, vld = train.split(random_state=random.seed(SEED))
6
7 TEXT.build_vocab(trn)
8 LABEL.build_vocab(trn)
9
10 device = "cuda" if torch.cuda.is_available() else "cpu"

```

```

1 train_iter, val_iter, test_iter = BucketIterator.splits(
2     (trn, vld, tst),
3     batch_sizes=(128, 256, 256),
4     sort=False,
5     sort_key= lambda x: len(x.src),
6     sort_within_batch=False,
7     device=device,
8     repeat=False,
9 )

```

Вы можете использовать Conv2d с `in_channels=1, kernel_size=(kernel_sizes[0], emb_dim)` или Conv1d с `in_channels=emb_dim, kernel_size=kernel_size[0]`. Но хорошенько подумайте над shape в обоих случаях.

```

1 class CNN(nn.Module):
2     def __init__(
3         self,
4         vocab_size,
5         emb_dim,
6         out_channels,
7         kernel_sizes,
8         dropout=0.5,
9     ):
10         super().__init__()
11
12         self.embedding = nn.Embedding(vocab_size, emb_dim)
13
14         self.conv_0 = nn.Conv1d(emb_dim, out_channels, kernel_size=kernel_sizes[0], pad
15
16         self.conv_1 = nn.Conv1d(emb_dim, out_channels, kernel_size=kernel_sizes[1], pad
17
18         self.conv_2 = nn.Conv1d(emb_dim, out_channels, kernel_size=kernel_sizes[2], pad
19
20         self.fc = nn.Linear(len(kernel_sizes) * out_channels, 1)
21
22         self.dropout = nn.Dropout(dropout)
23
24
25     def forward(self, text):
26
27         embedded = self.embedding(text)
28
29         embedded = embedded.permute(0, 2, 1) # may be reshape here
30

```

```

31         convded_0 = F.relu(self.conv_0(embedded)) # may be reshape here
32         convded_1 = F.relu(self.conv_1(embedded)) # may be reshape here
33         convded_2 = F.relu(self.conv_2(embedded)) # may be reshape here
34
35         pooled_0 = F.max_pool1d(convded_0, convded_0.shape[2]).squeeze(2)
36         pooled_1 = F.max_pool1d(convded_1, convded_1.shape[2]).squeeze(2)
37         pooled_2 = F.max_pool1d(convded_2, convded_2.shape[2]).squeeze(2)
38
39         cat = self.dropout(torch.cat((pooled_0, pooled_1, pooled_2), dim=1))
40
41         return self.fc(cat)

```

```

1 kernel_sizes = [3, 4, 5]
2 vocab_size = len(TEXT.vocab)
3 out_channels=64
4 dropout = 0.5
5 dim = 300
6
7 model = CNN(vocab_size=vocab_size, emb_dim=dim, out_channels=out_channels,
8             kernel_sizes=kernel_sizes, dropout=dropout)

```

```

1 model.to(device)

```

```

CNN(
  (embedding): Embedding(201630, 300)
  (conv_0): Conv1d(300, 64, kernel_size=(3,), stride=(1,))
  (conv_1): Conv1d(300, 64, kernel_size=(4,), stride=(1,))
  (conv_2): Conv1d(300, 64, kernel_size=(5,), stride=(1,))
  (fc): Linear(in_features=192, out_features=1, bias=True)
  (dropout): Dropout(p=0.5, inplace=False)
)

```

```

1 opt = torch.optim.Adam(model.parameters())
2 loss_func = nn.BCEWithLogitsLoss()

```

```

1 max_epochs = 30

```

Обучите!

```

1 min_loss = np.inf
2
3 cur_patience = 0
4
5 for epoch in range(1, max_epochs + 1):
6     train_loss = 0.0
7     model.train()
8     pbar = tqdm(enumerate(train_iter), total=len(train_iter), leave=False)
9     pbar.set_description(f"Epoch {epoch}")
10    for it, batch in pbar:
11        #YOUR CODE GOES HERE
12        opt.zero_grad()
13        output = model(batch.text)

```



```
14     loss = loss_func(output.squeeze(1), batch.label)
15     loss.backward()
16     train_loss += loss.item()
17     opt.step()
18
19     train_loss /= len(train_iter)
20     val_loss = 0.0
21     model.eval()
22     pbar = tqdm(enumerate(val_iter), total=len(val_iter), leave=False)
23     pbar.set_description(f"Epoch {epoch}")
24     for it, batch in pbar:
25         # YOUR CODE GOES HERE
26         output = model(batch.text)
27         val_loss += loss_func(output.squeeze(1), batch.label).item()
28
29     val_loss /= len(val_iter)
30     if val_loss < min_loss:
31         min_loss = val_loss
32         best_model = model.state_dict()
33     else:
34         cur_patience += 1
35         if cur_patience == patience:
36             cur_patience = 0
37             break
38
39     print('Epoch: {}, Training Loss: {}, Validation Loss: {}'.format(epoch, train_loss,
40 model.load_state_dict(best_model)
```

```

Epoch 1: 100%                                137/137 [00:39<00:00, 3.55it/s]
Epoch 1: 97%                                29/30 [00:07<00:00, 3.90it/s]
Epoch: 1, Training Loss: 0.5067113494350962, Validation Loss: 0.42657203574975333
Epoch 2: 100%                                137/137 [00:38<00:00, 3.56it/s]
Epoch 2: 97%                                29/30 [00:07<00:00, 3.88it/s]
Epoch: 2, Training Loss: 0.43463385213900657, Validation Loss: 0.3943417410055796
Epoch 3: 100%                                137/137 [00:39<00:00, 3.82it/s]
Epoch 3: 97%                                29/30 [00:07<00:00, 3.90it/s]
Epoch: 3, Training Loss: 0.3772826199113888, Validation Loss: 0.36494067311286926

```

Посчитайте f1-score вашего классификатора.

**Ответ:** 0.8473694279701549

```

-----
1 pred_labels = []
2 true_labels = []
3 for batch in test_iter:
4     pred_labels_batch = list((model(batch.text) > 0.5).float().cpu().numpy().reshape(-1)
5     true_labels_batch = list(batch.label.cpu().numpy())
6
7     pred_labels += pred_labels_batch
8     true_labels += true_labels_batch
9 f1_score(true_labels, pred_labels)

0.8473694279701549
Epoch: 1, Training Loss: 0.14652205117209965, Validation Loss: 0.5751664216045090

```

## ▼ Интерпретируемость

Посмотрим, куда смотрит наша модель. Достаточно запустить код ниже.

```

Epoch 3: 100%                                137/137 [00:38<00:00, 3.41it/s]

1 !pip install -q captum

|████████████████████████████████████████| 1.4 MB 5.3 MB/s

1 from captum.attr import LayerIntegratedGradients, TokenReferenceBase, visualization
2
3 PAD_IND = TEXT.vocab.stoi['pad']
4
5 token_reference = TokenReferenceBase(reference_token_idx=PAD_IND)
6 lig = LayerIntegratedGradients(model, model.embedding)

1 def forward_with_softmax(inp):
2     logits = model(inp)
3     return torch.softmax(logits, 0)[0][1]
4
5 def forward_with_sigmoid(input):
6     return torch.sigmoid(model(input))

```

```

7
8
9 # accumulate couple samples in this array for visualization purposes
10 vis_data_records_ig = []
11
12 def interpret_sentence(model, sentence, min_len = 7, label = 0):
13     model.eval()
14     text = [tok for tok in TEXT.tokenize(sentence)]
15     if len(text) < min_len:
16         text += ['pad'] * (min_len - len(text))
17     indexed = [TEXT.vocab.stoi[t] for t in text]
18
19     model.zero_grad()
20
21     input_indices = torch.tensor(indexed, device=device)
22     input_indices = input_indices.unsqueeze(0)
23
24     # input_indices dim: [sequence_length]
25     seq_length = min_len
26
27     # predict
28     pred = forward_with_sigmoid(input_indices).item()
29     pred_ind = round(pred)
30
31     # generate reference indices for each sample
32     reference_indices = token_reference.generate_reference(seq_length, device=device).u
33
34     # compute attributions and approximation delta using layer integrated gradients
35     attributions_ig, delta = lig.attribute(input_indices, reference_indices, \
36                                           n_steps=5000, return_convergence_delta=True)
37
38     print('pred: ', LABEL.vocab.itos[pred_ind], '(', '%.2f'%pred, ')', ', delta: ', abs
39
40     add_attributions_to_visualizer(attributions_ig, text, pred, pred_ind, label, delta,
41
42 def add_attributions_to_visualizer(attributions, text, pred, pred_ind, label, delta, vi
43     attributions = attributions.sum(dim=2).squeeze(0)
44     attributions = attributions / torch.norm(attributions)
45     attributions = attributions.cpu().detach().numpy()
46
47     # storing couple samples in an array for visualization purposes
48     vis_data_records.append(visualization.VisualizationDataRecord(
49         attributions,
50         pred,
51         LABEL.vocab.itos[pred_ind],
52         LABEL.vocab.itos[label],
53         LABEL.vocab.itos[1],
54         attributions.sum(),
55         text,
56         delta))

```

```

1 interpret_sentence(model, 'It was a fantastic performance !', label=1)
2 interpret_sentence(model, 'Best film ever', label=1)

```

```
3 interpret_sentence(model, 'Such a great show!', label=1)
4 interpret_sentence(model, 'It was a horrible movie', label=0)
5 interpret_sentence(model, 'I\'ve never watched something as bad', label=0)
6 interpret_sentence(model, 'It is a disgusting movie!', label=0)
7 interpret_sentence(model, 'an awfully good movie!', label=1)
8 interpret_sentence(model, 'movie is so good that I slept', label=0)
```

```
pred: pos ( 0.95 ), delta: tensor([3.5385e-05], device='cuda:0', dtype=torch.float)
pred: neg ( 0.02 ), delta: tensor([6.6953e-08], device='cuda:0', dtype=torch.float)
pred: neg ( 0.39 ), delta: tensor([2.9517e-05], device='cuda:0', dtype=torch.float)
pred: neg ( 0.00 ), delta: tensor([4.1141e-05], device='cuda:0', dtype=torch.float)
pred: neg ( 0.17 ), delta: tensor([6.1732e-05], device='cuda:0', dtype=torch.float)
pred: neg ( 0.09 ), delta: tensor([2.6798e-05], device='cuda:0', dtype=torch.float)
pred: neg ( 0.02 ), delta: tensor([0.0002], device='cuda:0', dtype=torch.float64)
pred: pos ( 0.61 ), delta: tensor([1.2505e-05], device='cuda:0', dtype=torch.float)
```



Попробуйте добавить свои примеры!

```
1 print('Visualize attributions based on Integrated Gradients')
2 visualization.visualize_text(vis_data_records_ig)
```



# Visualize attributions based on Integrated Gradients

Legend: ☐ Negative ☐ Neutral ☐ Positive

| True Label | Predicted Label | Attribution Label | Attribution Score | Word Importance                         |
|------------|-----------------|-------------------|-------------------|---|
| pos        | pos (0.95)      | pos               | 1.64              | It was a fantastic performance ! pad    |
| pos        | neg (0.02)      | pos               | 1.58              | Best film ever pad pad pad pad          |
| pos        | neg (0.39)      | pos               | 0.82              | Such a great show! pad pad pad          |
| neg        | neg (0.00)      | pos               | -0.27             | It was a horrible movie pad pad         |
| neg        | neg (0.17)      | pos               | 0.81              | I've never watched something as bad pad |
| neg        | neg (0.09)      | pos               | 0.74              | It is a disgusting movie! pad pad       |
| pos        | pos (0.64)      | pos               | 1.69              | Ужасно хороший фильм! pad pad pad pad   |
| pos        | neg (0.02)      | pos               | 0.11              | an awfully good movie! pad pad pad      |
| pos        | pos (0.95)      | pos               | 1.64              | It was a fantastic performance ! pad    |
| pos        | neg (0.02)      | pos               | 1.58              | Best film ever pad pad pad pad          |
| pos        | neg (0.39)      | pos               | 0.82              | Such a great show! pad pad pad          |
| neg        | neg (0.00)      | pos               | -0.27             | It was a horrible movie pad pad         |
| neg        | neg (0.17)      | pos               | 0.81              | I've never watched something as bad pad |
| neg        | neg (0.09)      | pos               | 0.74              | It is a disgusting movie! pad pad       |
| pos        | neg (0.02)      | pos               | 0.11              | an awfully good movie! pad pad pad      |
| neg        | neg (0.00)      | pos               | -0.61             | Incredibly bad movie pad pad pad pad    |
| neg        | pos (0.67)      | pos               | 2.35              | movie is so good that I slept           |
| neg        | pos (0.61)      | pos               | 1.77              | movie is so good that I slept           |
| pos        | pos (0.95)      | pos               | 1.64              | It was a fantastic performance ! pad    |
| pos        | neg (0.02)      | pos               | 1.58              | Best film ever pad pad pad pad          |
| pos        | neg (0.39)      | pos               | 0.82              | Such a great show! pad pad pad          |
| neg        | neg (0.00)      | pos               | -0.27             | It was a horrible movie pad pad         |
| neg        | neg (0.17)      | pos               | 0.81              | I've never watched something as bad pad |
| neg        | neg (0.09)      | pos               | 0.74              | It is a disgusting movie! pad pad       |
| pos        | neg (0.02)      | pos               | 0.11              | an awfully good movie! pad pad pad      |
| neg        | pos (0.61)      | pos               | 1.77              | movie is so good that I slept           |

Legend: ☐ Negative ☐ Neutral ☐ Positive

True Predicted Attribution Attribution

## ▼ Эмбединги слов

Вы ведь не забыли, как мы можем применить знания о word2vec и GloVe. Давайте попробуем!

```
neg neg (0.00) pos -0.27 It was a horrible movie pad pad
1 TEXT.build_vocab(trn, vectors=GloVe())# YOUR CODE GOES HERE
2 # подсказка: один из импортов пока не использовался, быть может он нужен в строке выше
3 LABEL.build_vocab(trn)
4
5 word_embeddings = TEXT.vocab.vectors
6
7 kernel_sizes = [3, 4, 5]
8 vocab_size = len(TEXT.vocab)
9 dropout = 0.5
10 dim = 300
```

```
.vector_cache/glove.840B.300d.zip: 2.18GB [07:00, 5.18MB/s]
100%|██████████| 2196016/2196017 [05:08<00:00, 7129.39it/s]
```

```
neg neg (0.00) pos -0.27 It was a horrible movie pad pad
```

```
1 train, tst = datasets.IMDB.splits(TEXT, LABEL)
2 trn, vld = train.split(random_state=random.seed(SEED))
3
4 device = "cuda" if torch.cuda.is_available() else "cpu"
5
6
7 train_iter, val_iter, test_iter = BucketIterator.splits(
8     (trn, vld, tst),
9     batch_sizes=(128, 256, 256),
10    sort=False,
11    sort_key= lambda x: len(x.src),
12    sort_within_batch=False,
13    device=device,
14    repeat=False,
15 )

1 model = CNN(vocab_size=vocab_size, emb_dim=dim, out_channels=64,
2             kernel_sizes=kernel_sizes, dropout=dropout)
3
4 word_embeddings = TEXT.vocab.vectors
5
6 prev_shape = model.embedding.weight.shape
7
8 model.embedding.weight.data.copy_(word_embeddings) # инициализируйте эмбединги
9
10 assert prev_shape == model.embedding.weight.shape
11 model.to(device)
12
13 opt = torch.optim.Adam(model.parameters())
```

Вы знаете, что делать.

```

1 import numpy as np
2
3 min_loss = np.inf
4
5 cur_patience = 0
6
7 for epoch in range(1, max_epochs + 1):
8     train_loss = 0.0
9     model.train()
10    pbar = tqdm(enumerate(train_iter), total=len(train_iter), leave=False)
11    pbar.set_description(f"Epoch {epoch}")
12    for it, batch in pbar:
13        #YOUR CODE GOES HERE
14        opt.zero_grad()
15        output = model(batch.text)
16        loss = loss_func(output.squeeze(1), batch.label)
17        loss.backward()
18        train_loss += loss.item()
19        opt.step()
20
21    train_loss /= len(train_iter)
22    val_loss = 0.0
23    model.eval()
24    pbar = tqdm(enumerate(val_iter), total=len(valid_iter), leave=False)
25    pbar.set_description(f"Epoch {epoch}")
26    for it, batch in pbar:
27        # YOUR CODE GOES HERE
28        output = model(batch.text)
29        val_loss += loss_func(output.squeeze(1), batch.label).item()
30
31    val_loss /= len(val_iter)
32    if val_loss < min_loss:
33        min_loss = val_loss
34        best_model = model.state_dict()
35    else:
36        cur_patience += 1
37        if cur_patience == patience:
38            cur_patience = 0
39            break
40
41    print('Epoch: {}, Training Loss: {}, Validation Loss: {}'.format(epoch, train_loss,
42 model.load_state_dict(best_model)

```

```

Epoch 1: 100%                                137/137 [00:39<00:00, 3.30it/s]
Epoch 1: 25%                                30/118 [00:07<00:23, 3.80it/s]
Epoch: 1, Training Loss: 0.30204641536204485, Validation Loss: 0.2985475242137909
Epoch 2: 100%                                137/137 [00:39<00:00, 3.53it/s]
Epoch 2: 25%                                30/118 [00:07<00:22, 3.86it/s]
Epoch: 2, Training Loss: 0.1839325635211311, Validation Loss: 0.28009043584267296
Epoch 3: 100%                                137/137 [00:38<00:00, 3.58it/s]
Epoch 3: 25%                                30/118 [00:07<00:22, 3.87it/s]

```

```

1 pred_labels = []
2 true_labels = []
3 for batch in test_iter:
4     pred_labels_batch = list((model(batch.text) > 0.5).float().cpu().numpy().reshape(-1)
5     true_labels_batch = list(batch.label.cpu().numpy())
6
7     pred_labels += pred_labels_batch
8     true_labels += true_labels_batch
9 f1_score(true_labels, pred_labels)

0.8653351304987135

```

Посчитайте f1-score вашего классификатора.

**Ответ:** 0.8653351304987135

Проверим насколько все хорошо!

```

1 PAD_IND = TEXT.vocab.stoi['pad']
2
3 token_reference = TokenReferenceBase(reference_token_idx=PAD_IND)
4 lig = LayerIntegratedGradients(model, model.embedding)
5 vis_data_records_ig = []
6
7 interpret_sentence(model, 'It was a fantastic performance !', label=1)
8 interpret_sentence(model, 'Best film ever', label=1)
9 interpret_sentence(model, 'Such a great show!', label=1)
10 interpret_sentence(model, 'It was a horrible movie', label=0)
11 interpret_sentence(model, 'I\'ve never watched something as bad', label=0)
12 interpret_sentence(model, 'It is a disgusting movie!', label=0)
13 interpret_sentence(model, 'an awfully good movie!', label=1)
14 interpret_sentence(model, 'movie is so good that I slept', label=0)

pred: pos ( 0.92 ), delta: tensor([0.0002], device='cuda:0', dtype=torch.float64)
pred: neg ( 0.16 ), delta: tensor([8.6935e-06], device='cuda:0', dtype=torch.float64)
pred: neg ( 0.37 ), delta: tensor([2.6935e-06], device='cuda:0', dtype=torch.float64)
pred: neg ( 0.00 ), delta: tensor([1.0840e-05], device='cuda:0', dtype=torch.float64)
pred: neg ( 0.50 ), delta: tensor([4.0545e-05], device='cuda:0', dtype=torch.float64)
pred: neg ( 0.00 ), delta: tensor([1.2201e-06], device='cuda:0', dtype=torch.float64)
pred: neg ( 0.01 ), delta: tensor([4.0200e-05], device='cuda:0', dtype=torch.float64)
pred: pos ( 0.56 ), delta: tensor([0.0002], device='cuda:0', dtype=torch.float64)

```



```
1 print('Visualize attributions based on Integrated Gradients')
2 visualization.visualize_text(vis_data_records_ig)
```

Visualize attributions based on Integrated Gradients

Legend: ☐ Negative ☐ Neutral ☐ Positive

| True Label | Predicted Label | Attribution Label | Attribution Score | Word Importance                         |
|------------|-----------------|-------------------|-------------------|---|
| pos        | pos (0.92)      | pos               | 1.85              | It was a fantastic performance ! pad    |
| pos        | neg (0.16)      | pos               | 1.25              | Best film ever pad pad pad pad          |
| pos        | neg (0.37)      | pos               | 1.54              | Such a great show! pad pad pad          |
| neg        | neg (0.00)      | pos               | -0.46             | It was a horrible movie pad pad         |
| neg        | neg (0.50)      | pos               | 1.67              | I've never watched something as bad pad |
| neg        | neg (0.00)      | pos               | -0.91             | It is a disgusting movie! pad pad       |
| pos        | neg (0.01)      | pos               | -0.19             | an awfully good movie! pad pad pad      |
| neg        | pos (0.56)      | pos               | 1.80              | movie is so good that I slept           |

Legend: ☐ Negative ☐ Neutral ☐ Positive

| True Label | Predicted Label | Attribution Label | Attribution Score | Word Importance                         |
|------------|-----------------|-------------------|-------------------|---|
| pos        | pos (0.92)      | pos               | 1.85              | It was a fantastic performance ! pad    |
| pos        | neg (0.16)      | pos               | 1.25              | Best film ever pad pad pad pad          |
| pos        | neg (0.37)      | pos               | 1.54              | Such a great show! pad pad pad          |
| neg        | neg (0.00)      | pos               | -0.46             | It was a horrible movie pad pad         |
| neg        | neg (0.50)      | pos               | 1.67              | I've never watched something as bad pad |

---

✓ 0s completed at 2:39 PM

