

Лабораторная работа №7. Рекуррентные нейронные сети для анализа текста

Данные: Набор данных для предсказания оценок для отзывов, собранных с сайта imdb.com, который состоит из 50,000 отзывов в виде текстовых файлов. Отзывы разделены на положительные (25,000) и отрицательные (25,000). Данные предварительно токенизированы по принципу “мешка слов”, индексы слов можно взять из словаря (imdb.vocab). Обучающая выборка включает в себя 12,500 положительных и 12,500 отрицательных отзывов, контрольная выборка также содержит 12,500 положительных и 12,500 отрицательных отзывов, а также. Данные можно скачать по ссылке <https://ai.stanford.edu/~amaas/data/sentiment/> (<https://ai.stanford.edu/~amaas/data/sentiment/>).

```
In [0]: 1 !pip install git+https://github.com/d2l-ai/d2l-en # installing
        2 !pip install -U --pre mxnet-cu101mkl # updating mxnet to at le
        3
```

```
Uninstalling pyzmq-17.0.0:
  Successfully uninstalled pyzmq-17.0.0
Successfully installed d2l-0.11.4 pyzmq-19.0.0

Collecting mxnet-cu101mkl
  Downloading https://files.pythonhosted.org/packages/3d/4b/e51dc49ca5fe6564028e7c91b10a3f79c00d710dd691b408c77597df5883/mxnet_cu101mkl-1.6.0-py2.py3-none-manylinux1_x86_64.whl
    (https://files.pythonhosted.org/packages/3d/4b/e51dc49ca5fe6564028e7c91b10a3f79c00d710dd691b408c77597df5883/mxnet_cu101mkl-1.6.0-py2.py3-none-manylinux1_x86_64.whl) (711.0MB)
    |████████████████████████████████████████████████████████████████████████████████| 711.0MB 24kB/s
Collecting graphviz<0.9.0,>=0.8.1
  Downloading https://files.pythonhosted.org/packages/53/39/4ab213673844e0c004bed8a0781a0721a3f6bb23eb8854ee75c236428892/graphviz-0.8.4-py2.py3-none-any.whl
    (https://files.pythonhosted.org/packages/53/39/4ab213673844e0c004bed8a0781a0721a3f6bb23eb8854ee75c236428892/graphviz-0.8.4-py2.py3-none-any.whl)
Requirement already satisfied, skipping upgrade: numpy<2.0.0,>=1.16
```

```
In [0]: 1 import shutil
        2 from google.colab import drive
        3 import d2l
        4 from mxnet import gluon, np, npx, init
        5 import os
        6 from mxnet.gluon import nn, rnn
        7 from mxnet.contrib import text
        8 npx.set_np()
```

Задание 1. Загрузите данные. Преобразуйте текстовые файлы во внутренние структуры данных, которые используют индексы вместо слов.

```
In [0]: 1 !wget https://ai.stanford.edu/~amaas/data/sentiment/aclImdb_v1.
2 drive.mount('/content/drive')
3 shutil.unpack_archive("aclImdb_v1.tar.gz", "/content/aclImdb_v1
```

```
--2020-04-03 17:29:45-- https://ai.stanford.edu/~amaas/data/senti
ment/aclImdb_v1.tar.gz
(https://ai.stanford.edu/~amaas/data/sentiment/aclImdb_v1.tar.gz)
Resolving ai.stanford.edu (ai.stanford.edu)... 171.64.68.10
Connecting to ai.stanford.edu (ai.stanford.edu)|171.64.68.10|:443.
.. connected.
HTTP request sent, awaiting response... 200 OK
Length: 84125825 (80M) [application/x-gzip]
Saving to: 'aclImdb_v1.tar.gz.2'
```

```
aclImdb_v1.tar.gz.2 100%[=====>] 80.23M 85.2MB/s
in 0.9s
```

```
2020-04-03 17:29:46 (85.2 MB/s) - 'aclImdb_v1.tar.gz.2' saved [841
25825/84125825]
```

Drive already mounted at /content/drive; to attempt to forcibly re-mount, call drive.mount("/content/drive", force_remount=True).

```
In [0]: 1 def read_imdb(data_dir, is_train):
2     data, labels = [], []
3     for label in ('pos', 'neg'):
4         folder_name = os.path.join(data_dir, 'train' if is_train
5                                     else label)
6         for file in os.listdir(folder_name):
7             with open(os.path.join(folder_name, file), 'rb') as
8                 f:
9                 review = f.read().decode('utf-8').replace('\n',
10                 data.append(review)
11                 labels.append(1 if label == 'pos' else 0)
12     return data, labels
13 train_data = read_imdb('/content/aclImdb_v1/aclImdb', is_train=
14 print('trainings:', len(train_data[0]))
```

```
trainings: 25000
```

```
In [0]: 1 train_tokens = d2l.tokenize(train_data[0], token='word')
2 vocab = d2l.Vocab(train_tokens, min_freq=5, reserved_tokens=['<
3
4 d2l.set_figsize((3.5, 2.5))
5 d2l.plt.hist([len(line) for line in train_tokens], bins=range(0
```

<Figure size 252x180 with 1 Axes>

```
In [0]: 1 num_steps = 500 # sequence length
2 train_features = np.array([d2l.truncate_pad(
3     vocab[line], num_steps, vocab['<pad>']) for line in train_t
4 train_features.shape
```

Out[13]: (25000, 500)

```
In [0]: 1 train_iter = d2l.load_array((train_features, train_data[1]), 64
2 'batches:', len(train_iter))
```

Out[17]: ('batches:', 391)

```
In [0]: 1 def load_data_imdb(batch_size, num_steps=500):
2     data_dir = '/content/aclImdb_v1/aclImdb'
3     train_data = read_imdb(data_dir, True)
4     test_data = read_imdb(data_dir, False)
5     train_tokens = d2l.tokenize(train_data[0], token='word')
6     test_tokens = d2l.tokenize(test_data[0], token='word')
7     vocab = d2l.Vocab(train_tokens, min_freq=5)
8     train_features = np.array([d2l.truncate_pad(
9         vocab[line], num_steps, vocab.unk) for line in train_to
10    test_features = np.array([d2l.truncate_pad(
11        vocab[line], num_steps, vocab.unk) for line in test_tok
12    train_iter = d2l.load_array((train_features, train_data[1])
13    test_iter = d2l.load_array((test_features, test_data[1]), b
14                                is_train=False)
15    return train_iter, test_iter, vocab
```

Задание 2. Реализуйте и обучите двунаправленную рекуррентную сеть (LSTM или GRU). Какого качества классификации удалось достичь?

```
In [0]: 1 batch_size = 64
2 train_iter, test_iter, vocab = load_data_imdb(batch_size)
```

```
In [0]: 1 class BiRNN(nn.Block):
2         def __init__(self, vocab_size, embed_size, num_hiddens,
3                     num_layers, **kwargs):
4             super(BiRNN, self).__init__(**kwargs)
5             self.embedding = nn.Embedding(vocab_size, embed_size)
6             self.encoder = rnn.LSTM(num_hiddens, num_layers=num_layers,
7                                   bidirectional=True, input_size=embed_size)
8             self.decoder = nn.Dense(2)
9
10        def forward(self, inputs):
11            embeddings = self.embedding(inputs.T)
12            outputs = self.encoder(embeddings)
13            encoding = np.concatenate((outputs[0], outputs[-1]), axis=1)
14            outs = self.decoder(encoding)
15            return outs
```

```
In [0]: 1 embed_size, num_hiddens, num_layers, ctx = 100, 100, 2, d2l.try_gpu(0)
2 net = BiRNN(len(vocab), embed_size, num_hiddens, num_layers)
3 net.initialize(init.Xavier(), ctx=ctx)
```

```
In [0]: 1 lr, num_epochs = 0.01, 5
2 trainer = gluon.Trainer(net.collect_params(), 'adam', {'learning_rate': lr})
3 loss = gluon.loss.SoftmaxCrossEntropyLoss()
4 d2l.train_ch13(net, train_iter, test_iter, loss, trainer, num_epochs)
```

loss 0.095, train acc 0.970, test acc 0.823
580.3 examples/sec on [gpu(0)]

<Figure size 252x180 with 1 Axes>

```
In [0]: 1 predict_sentiment(net, vocab, 'this movie is so great')
```

Out[54]: 'positive'

```
In [0]: 1 predict_sentiment(net, vocab, 'this movie is so bad')
```

Out[55]: 'negative'

Задание 3. Используйте индексы слов и их различное внутреннее представление (word2vec, glove). Как влияет данное преобразование на качество классификации?

```
In [0]: 1 glove_embedding = text.embedding.create('glove', pretrained_fil
```

Downloading /root/.mxnet/embeddings/glove/glove.6B.zip from <https://apache-mxnet.s3-accelerate.dualstack.amazonaws.com/gluon/embeddings/glove/glove.6B.zip...> (<https://apache-mxnet.s3-accelerate.dualstack.amazonaws.com/gluon/embeddings/glove/glove.6B.zip...>)

```
In [0]: 1 embeds = glove_embedding.get_vecs_by_tokens(vocab.idx_to_token)
2 embeds.shape
```

Out[42]: (49339, 100)

```
In [0]: 1 net.embedding.weight.set_data(embeds)
2 net.embedding.collect_params().setattr('grad_req', 'null')
```

```
In [0]: 1 lr, num_epochs = 0.01, 5
2 trainer = gluon.Trainer(net.collect_params(), 'adam', {'learning_rate': lr})
3 loss = gluon.loss.SoftmaxCrossEntropyLoss()
4 d2l.train_ch13(net, train_iter, test_iter, loss, trainer, num_e
```

loss 0.317, train acc 0.866, test acc 0.844
604.6 examples/sec on [gpu(0)]

<Figure size 252x180 with 1 Axes>

```
In [0]: 1 def predict_sentiment(net, vocab, sentence):
2     sentence = np.array(vocab[sentence.split()], ctx=d2l.try_gpu())
3     label = np.argmax(net(sentence.reshape(1, -1)), axis=1)
4     return 'positive' if label == 1 else 'negative'
```

```
In [0]: 1 predict_sentiment(net, vocab, 'this movie is so great')
```

Out[50]: 'positive'

```
In [0]: 1 predict_sentiment(net, vocab, 'this movie is so bad')
```

Out[51]: 'negative'

Задание 4. Поэкспериментируйте со структурой сети (добавьте больше рекуррентных, полносвязных или сверточных слоев). Как это повлияло на качество классификации?

```
In [0]: 1 lr, num_epochs = 0.01, 5
2 trainer = gluon.Trainer(net.collect_params(), 'adam', {'learning_rate': lr})
3 loss = gluon.loss.SoftmaxCrossEntropyLoss()
4 d2l.train_ch13(net, train_iter, test_iter, loss, trainer, num_epochs)
```

<Figure size 252x180 with 1 Axes>

```
In [2]: 1 !git clone https://github.com/bfelbo/DeepMoji.git
```

```
In [17]: 1 !pip install tensorflow==1.14
```

Страница 6 из 10

```
|████████████████████████████████████████| 3.2MB 49.4MB/s  
Requirement already satisfied: gast>=0.2.0 in /usr/local/lib/pytho  
n3.6/dist-packages (from tensorflow==1.14) (0.3.3)  
Requirement already satisfied: absl-py>=0.7.0 in /usr/local/lib/py  
thon3.6/dist-packages (from tensorflow==1.14) (0.9.0)  
Requirement already satisfied: wrapt>=1.11.1 in /usr/local/lib/pyt  
hon3.6/dist-packages (from tensorflow==1.14) (1.12.1)  
Collecting tensorflow-estimator<1.15.0rc0,>=1.14.0rc0
```

```

Requirement already satisfied: protobuf>=3.6.1 in /usr/local/lib/python3.6/dist-packages (from tensorflow==1.14) (3.10.0)
Requirement already satisfied: google-pasta>=0.1.6 in /usr/local/lib/python3.6/dist-packages (from tensorflow==1.14) (0.2.0)
Requirement already satisfied: grpcio>=1.8.6 in /usr/local/lib/python3.6/dist-packages (from tensorflow==1.14) (1.27.2)
Requirement already satisfied: numpy<2.0,>=1.14.5 in /usr/local/lib/python3.6/dist-packages (from tensorflow==1.14) (1.18.2)
Requirement already satisfied: keras-applications>=1.0.6 in /usr/local/lib/python3.6/dist-packages (from tensorflow==1.14) (1.0.8)
Requirement already satisfied: markdown>=2.6.8 in /usr/local/lib/python3.6/dist-packages (from tensorboard<1.15.0,>=1.14.0->tensorflow==1.14) (3.2.1)
Requirement already satisfied: werkzeug>=0.11.15 in /usr/local/lib/python3.6/dist-packages (from tensorboard<1.15.0,>=1.14.0->tensorflow==1.14) (1.0.1)
Requirement already satisfied: setuptools>=41.0.0 in /usr/local/lib/python3.6/dist-packages (from tensorboard<1.15.0,>=1.14.0->tensorflow==1.14) (46.0.0)
Requirement already satisfied: h5py in /usr/local/lib/python3.6/dist-packages (from keras-applications>=1.0.6->tensorflow==1.14) (2.10.0)

```

```
Successfully uninstalled tensorflow-2.2.0rc2
```

Successfully installed tensorboard-1.14.0 tensorflow-1.14.0 tensorflow-estimator-1.14.0

```
In [0]: 1 import os
        2 os.chdir('/content/DeepMoji')
```

```
In [3]: 1 from __future__ import print_function
        2 import deepmoji
        3 import examples.example_helper
        4 import numpy as np
        5 from keras.preprocessing import sequence
        6 from keras.datasets import imdb
        7 import deepmoji.model_def
        8
        9 nb_tokens = 20000
       10 maxlen = 80
       11 batch_size = 32
       12
       13 print('Loading data...')
       14 print(len(train_data[0]), 'train sequences')
       15 print(len(test_data[0]), 'test sequences')
       16 print('Pad sequences (samples x time)')
       17 print('X_train shape:', train_data.shape)
       18 print('X_test shape:', test_data.shape)
       19
       20 print('Build model...')
       21 model = deepmoji.model_def.deepmoji_architecture(nb_classes=2,
       22 model.summary()
       23
       24 model.compile(loss='binary_crossentropy',
       25               optimizer='adam',
       26               metrics=['accuracy'])
       27
       28 print('Train...')
       29 model.fit(train_data[0], train_data[1], batch_size=batch_size,
       30         validation_data=(test_data[0], test_data[0]))
       31 score, acc = model.evaluate(test_data[0], test_data[1], batch_s
       32 print('Test score:', score)
       33 print('Test accuracy:', acc)
```

```
Loading data...
25000 train sequences
25000 test sequences
Pad sequences (samples x time)
X_train shape: (25000, 80)
X_test shape: (25000, 80)
```


Build model...
Model: "DeepMoji"

| Layer (type) connected to | Output Shape | Param # | C |
|--|------------------|---------|---|
| ===== | | | |
| input_2 (InputLayer) | (None, 80) | 0 | |
| ===== | | | |
| embedding (Embedding) input_2[0][0] | (None, 80, 256) | 5120000 | i |
| ===== | | | |
| activation_2 (Activation) embedding[0][0] | (None, 80, 256) | 0 | e |
| ===== | | | |
| bi_lstm_0 (Bidirectional) activation_2[0][0] | (None, 80, 1024) | 3149824 | a |
| ===== | | | |
| bi_lstm_1 (Bidirectional) bi_lstm_0[0][0] | (None, 80, 1024) | 6295552 | b |
| ===== | | | |
| concatenate_2 (Concatenate) bi_lstm_1[0][0] | (None, 80, 2304) | 0 | b |
| ===== | | | |
| i_lstm_0[0][0] | | | b |
| ===== | | | |
| activation_2[0][0] | | | a |
| ===== | | | |
| attlayer (AttentionWeightedAverage) concatenate_2[0][0] | (None, 2304) | 2304 | c |
| ===== | | | |
| softmax (Dense) attlayer[0][0] | (None, 1) | 2305 | a |
| ===== | | | |
| Total params: 14,569,985 | | | |
| Trainable params: 14,569,985 | | | |
| Non-trainable params: 0 | | | |

Train...
Train on 25000 samples, validate on 25000 samples
Epoch 1/5
25000/25000 [=====] - 6155s 246ms/step -
loss: 0.4263 - acc: 0.8038 - val_loss: 0.3644 - val_acc: 0.8446

```
Epoch 2/5
25000/25000 [=====] - 5949s 238ms/step -
loss: 0.2494 - acc: 0.9028 - val_loss: 0.3690 - val_acc: 0.8428
Epoch 3/5
25000/25000 [=====] - 5676s 227ms/step -
loss: 0.1577 - acc: 0.9453 - val_loss: 0.4109 - val_acc: 0.8388
Epoch 4/5
25000/25000 [=====] - 5633s 225ms/step -
loss: 0.0913 - acc: 0.9715 - val_loss: 0.5532 - val_acc: 0.8300
Epoch 5/5
25000/25000 [=====] - 5886s 235ms/step -
loss: 0.0647 - acc: 0.9823 - val_loss: 0.7103 - val_acc: 0.8306
25000/25000 [=====] - 1033s 41ms/step
Test score: 0.7102789056491852
Test accuracy: 0.83064
```