

```
In [1]: 1 import torch
2 import torch.autograd as autograd
3 import torch.nn as nn
4 import torch.optim as optim
5 import numpy as np
6 torch.manual_seed(1)
7 from sklearn.metrics import roc_auc_score
8 from sklearn.metrics import f1_score
9 import copy
10 import sys
11 from utils import preprocessing #using the same preprocessing method from ht
```

```
In [2]: 1 # Authors: Haocheng Zhang and Kehang (Fred) Chang
2 # portion of codes came from authors in https://github.com/tiantiantu/KSI
```

```
In [3]: 1 # !pip install numpy --upgrade
2 print(np.__version__)
```

1.19.5

```
In [4]: 1 # modify the default parameters of np.load
2 np_load_old = np.load
3 np.load = lambda *a,**k: np_load_old(*a, allow_pickle=True, **k)
```

```
In [5]: 1 # choose CPU if GPU is not available
2 device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
3 print(device)
```

cuda:0

```
In [6]: 1 # For consistency, import the data like other modals.
2 label_to_ix=np.load('label_to_ix.npy').item()
3 ix_to_label=np.load('ix_to_label.npy')
4 training_data=np.load('training_data.npy')
5 test_data=np.load('test_data.npy')
6 val_data=np.load('val_data.npy')
7 word_to_ix=np.load('word_to_ix.npy').item()
8 ix_to_word=np.load('ix_to_word.npy')
9 newwikivec=np.load('newwikivec.npy')
10 wikivoc=np.load('wikivoc.npy').item()
```

```
In [7]: 1 #init global vars
2 wikisize=newwikivec.shape[0]
3 rvocsize=newwikivec.shape[1]
4 wikivec=autograd.Variable(torch.FloatTensor(newwikivec))
```

```
In [8]: 1 # Use the same hyper params
        2 batchsize=32
        3 Embeddingsize=100
        4 topk=10
        5 padding_idx=0
        6 lr=0.01 #updated Lr
        7 epochs=1000
        8 dropout=0.2
        9 hidden_dim=200
        10 min_good_models=5
```

```
In [9]: 1 # Use the same preprocessing methods to get training, test and val dataset
        2 batchtraining_data=preprocessing(training_data, label_to_ix, word_to_ix, wik
        3 batchtest_data=preprocessing(test_data, label_to_ix, word_to_ix, wikivoc, ba
        4 batchval_data=preprocessing(val_data, label_to_ix, word_to_ix, wikivoc, batc
```

/home/hzhan147/utils.py:18: VisibleDeprecationWarning: Creating an ndarray from ragged nested sequences (which is a list-or-tuple of lists-or-tuples-or ndarrays with different lengths or shapes) is deprecated. If you meant to do this, you must specify 'dtype=object' when creating the ndarray

```
new_data=np.array(new_data)
```

Create the model:

```

In [10]: 1 class RNN(nn.Module):
2
3     def __init__(self, batch_size, vocab_size, tagset_size, padding_idx=0):
4         super(RNN, self).__init__()
5         self.hidden_dim = hidden_dim
6         self.word_embeddings = nn.Embedding(vocab_size+1, Embeddingsize, pad
7         self.rnn = nn.GRU(Embeddingsize, hidden_dim)
8         self.hidden2tag = nn.Linear(hidden_dim, tagset_size)
9         self.hidden = self.init_hidden()
10
11
12         self.layer2 = nn.Linear(Embeddingsize, 1,bias=False)
13         self.embedding=nn.Linear(rvocsize,Embeddingsize)
14         self.vattention=nn.Linear(Embeddingsize,Embeddingsize,bias=False)
15
16         self.sigmoid = nn.Sigmoid()
17         self.tanh = nn.Tanh()
18         self.embed_drop = nn.Dropout(p=dropout)
19
20     #init hidden layers and encapsulate it to a method, so that we can re-in
21     def init_hidden(self):
22         return autograd.Variable(torch.zeros(1, batchsize, self.hidden_dim).
23
24
25     def forward(self, vec1, nvec, wiki, simlearning):
26
27         thisembeddings=self.word_embeddings(vec1).transpose(0,1)
28         thisembeddings = self.embed_drop(thisembeddings)
29
30         #to match what authors' research, we use the SAME KSI algo.
31         if simlearning==1:
32             nvec=nvec.view(batchsize,1,-1)
33             nvec=nvec.expand(batchsize,wiki.size()[0],-1)
34             wiki=wiki.view(1,wiki.size()[0],-1)
35             wiki=wiki.expand(nvec.size()[0],wiki.size()[1],-1)
36             new=wiki*nvec
37             new=self.embedding(new)
38             vattention=self.sigmoid(self.vattention(new))
39             new=new*vattention
40             vec3=self.layer2(new)
41             vec3=vec3.view(batchsize,-1)
42
43         #Super simple RNN architecture: Sigmoid -> Linear -> MaxPool1d -> ta
44         rnn_out, self.hidden = self.rnn(thisembeddings, self.hidden)
45         rnn_out = self.tanh(rnn_out)
46         rnn_out=rnn_out.transpose(0,2).transpose(0,1)
47         output1=nn.MaxPool1d(rnn_out.size()[2])(rnn_out).view(batchsize,-1)
48
49         vec2 = self.hidden2tag(output1)
50         if simlearning==1:
51             tag_scores = self.sigmoid(vec2.detach()+vec3)
52         else:
53             tag_scores = self.sigmoid(vec2)
54
55
56         return tag_scores

```



```

In [11]: 1 def trainmodel(model, sim):
2         print ('start_training')
3         modelsaved=[]
4         modelperform=[]
5
6
7         bestresults=-1
8         bestiter=-1
9         for epoch in range(epochs):
10
11             model.train()
12
13             lossestrain = []
14             recall=[]
15             for mysentence in batchtraining_data:
16                 model.zero_grad()
17                 #re-init hidden layers on each train
18                 model.hidden = model.init_hidden()
19                 targets = mysentence[2].cuda()
20                 # train model
21                 tag_scores = model(mysentence[0].cuda(),mysentence[1].cuda(),wik
22                 # calc loss
23                 loss = loss_function(tag_scores, targets)
24                 # backprob
25                 loss.backward()
26                 # update params
27                 optimizer.step()
28                 # record loss for later calc
29                 lossestrain.append(loss.data.mean())
30             print (epoch)
31
32             # save model since we are tracking model improvements... If no impro
33             modelsaved.append(copy.deepcopy(model.state_dict()))
34             print ("XXXXXXXXXXXXXXXXXXXXXXXXXXXX")
35             model.eval()
36
37             recall=[]
38             for inputs in batchval_data:
39                 #re-init hidden layers on each eval
40                 model.hidden = model.init_hidden()
41                 targets = inputs[2].cuda()
42                 # eval model
43                 tag_scores = model(inputs[0].cuda(),inputs[1].cuda() ,wikivec.cu
44
45                 #calc loss
46                 loss = loss_function(tag_scores, targets)
47
48                 targets=targets.data.cpu().numpy()
49                 tag_scores= tag_scores.data.cpu().numpy()
50
51                 #calc recall based on top-K scores
52                 for idx in range(0,len(tag_scores)):
53                     temp={}
54                     for score_idx in range(0,len(tag_scores[idx])):
55                         temp[score_idx]=tag_scores[idx][score_idx]
56                     temp1=[(k, temp[k]) for k in sorted(temp, key=temp.get, reve

```

```
57         thistop=int(np.sum(targets[idx]))
58         hit=0.0
59         for ii in temp1[0:max(thistop,topk)]:
60             if targets[idx][ii[0]]==1.0:
61                 hit=hit+1
62         if thistop!=0:
63             recall.append(hit/thistop)
64
65     print ('validation top-',topk, np.mean(recall))
66
67
68     #track model performances here based on recalls mean.
69     #if current one is better, update best recalls mean and set best idx
70     modelperform.append(np.mean(recall))
71     if modelperform[-1]>bestresults:
72         bestresults=modelperform[-1]
73         besttiter=len(modelperform)-1
74
75     #use the best idx (besttiter) to track if we have minimum models afte
76     if (len(modelperform)-besttiter)>min_good_models:
77         print (modelperform,besttiter)
78         return modelsaved[besttiter]
79     else:
80         print('Not enough min models, keep training...')
```

```

In [12]: 1 def testmodel(modelstate, sim):
2         #-----reload model static params-----#
3         model = RNN(batchsize, len(word_to_ix), len(label_to_ix))
4         model.cuda()
5         model.load_state_dict(modelstate)
6         loss_function = nn.BCELoss()
7         model.eval()
8         #-----#
9
10        recall=[]
11        lossestest = []
12
13        y_true=[]
14        y_scores=[]
15
16
17        for inputs in batchtest_data:
18            #re-init hidden layers on each test
19            model.hidden = model.init_hidden()
20            targets = inputs[2].cuda()
21
22            #test model
23            tag_scores = model(inputs[0].cuda(),inputs[1].cuda() ,wikivec.cuda())
24            #calc loss
25            loss = loss_function(tag_scores, targets)
26
27            targets=targets.data.cpu().numpy()
28            tag_scores= tag_scores.data.cpu().numpy()
29
30            #tracking Loss
31            lossestest.append(loss.data.mean())
32            y_true.append(targets)
33            y_scores.append(tag_scores)
34
35            #calc recall based on top-K scores
36            for idx in range(0,len(tag_scores)):
37                temp={}
38                for score_idx in range(0,len(tag_scores[idx])):
39                    temp[score_idx]=tag_scores[idx][score_idx]
40                temp1=[(k, temp[k]) for k in sorted(temp, key=temp.get, reverse=
41                thistop=int(np.sum(targets[idx]))
42                hit=0.0
43                for ii in temp1[0:max(thistop,topk)]:
44                    if targets[idx][ii[0]]==1.0:
45                        hit=hit+1
46                if thistop!=0:
47                    recall.append(hit/thistop)
48        y_true=np.concatenate(y_true,axis=0)
49        y_scores=np.concatenate(y_scores,axis=0)
50        y_true=y_true.T
51        y_scores=y_scores.T
52        temprue=[]
53        tempscores=[]
54
55        #prepare trues and scores for later performance calc
56        for col in range(0,len(y_true)):

```

```
57         if np.sum(y_true[col])!=0:
58             temptrue.append(y_true[col])
59             tempscores.append(y_scores[col])
60     temptrue=np.array(temptrue)
61     tempscores=np.array(tempscores)
62     y_true=temptrue.T
63     y_scores=tempscores.T
64
65     #extract predictions
66     y_pred=(y_scores>0.5).astype(np.int)
67
68     #print all the metrics
69     print ('test loss', torch.stack(lossesestest).mean().item())
70     print ('top-',topk, np.mean(recall))
71     print ('macro AUC', roc_auc_score(y_true, y_scores,average='macro'))
72     print ('micro AUC', roc_auc_score(y_true, y_scores,average='micro'))
73     print ('macro F1', f1_score(y_true, y_pred, average='macro') )
74     print ('micro F1', f1_score(y_true, y_pred, average='micro') )
```



```
In [13]: 1 # START all the training here
2 model = RNN(batchsize, len(word_to_ix), len(label_to_ix), padding_idx)
3 model.cuda()
4
5 #use BCE Loss as Loss function
6 loss_function = nn.BCELoss()
7 #use Adam optimizer with lr
8 optimizer = optim.Adam(model.parameters(), lr=lr)
9 #train model with mode 0 (base RNN)
10 basemodel= trainmodel(model, 0)
11 #save base RNN model as file named 'RNN_model'
12 torch.save(basemodel, 'RNN_model')
```

```
start_training
0
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.40101608663973637
Not enough min models, keep training...
1
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.38618976115425224
Not enough min models, keep training...
2
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.38970947758109376
Not enough min models, keep training...
3
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.4504764719365307
Not enough min models, keep training...
4
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.4564919365005066
Not enough min models, keep training...
5
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.4660167262287978
Not enough min models, keep training...
6
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.46967451319739933
Not enough min models, keep training...
7
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.455825257555359
Not enough min models, keep training...
8
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.4620362717945726
Not enough min models, keep training...
9
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.46059115258629557
Not enough min models, keep training...
10
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.47923415371216876
```

```
Not enough min models, keep training...
11
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.47959522091145806
Not enough min models, keep training...
12
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.4516464749271538
Not enough min models, keep training...
13
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.5100073226221283
Not enough min models, keep training...
14
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.5121169906786013
Not enough min models, keep training...
15
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.5237021301347901
Not enough min models, keep training...
16
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.5205453972325008
Not enough min models, keep training...
17
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.5219308745393592
Not enough min models, keep training...
18
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.5286124162093939
Not enough min models, keep training...
19
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.5241472545949427
Not enough min models, keep training...
20
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.5295816630662274
Not enough min models, keep training...
21
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.5341732020650285
Not enough min models, keep training...
22
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.5366185058931356
Not enough min models, keep training...
23
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.5119711288506761
Not enough min models, keep training...
24
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.5162232759963225
Not enough min models, keep training...
```

```
25
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.5134678804632867
Not enough min models, keep training...
26
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.523015572021539
Not enough min models, keep training...
27
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.5020314007431665
[0.40101608663973637, 0.38618976115425224, 0.38970947758109376, 0.4504764719365
307, 0.4564919365005066, 0.4660167262287978, 0.46967451319739933, 0.45582525755
5359, 0.4620362717945726, 0.46059115258629557, 0.47923415371216876, 0.479595220
91145806, 0.4516464749271538, 0.5100073226221283, 0.5121169906786013, 0.5237021
301347901, 0.5205453972325008, 0.5219308745393592, 0.5286124162093939, 0.524147
2545949427, 0.5295816630662274, 0.5341732020650285, 0.5366185058931356, 0.51197
11288506761, 0.5162232759963225, 0.5134678804632867, 0.523015572021539, 0.50203
14007431665] 22
```

```
In [14]: 1 #START all the KSI training here
2 model = RNN(batchsize, len(word_to_ix), len(label_to_ix), padding_idx)
3 model.cuda()
4 model.load_state_dict(basemodel)
5
6 #use BCE loss as loss function
7 loss_function = nn.BCELoss()
8 #use Adam optimizer with lr
9 optimizer = optim.Adam(model.parameters(), lr=lr)
10 #train model with mode 1 (KSI RNN)
11 KSImodel= trainmodel(model, 1)
12 #save KSI RNN model as file named 'KSI_RNN_model'
13 torch.save(KSImodel, 'KSI_RNN_model')
```

```
start_training
```

```
0
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.6862041917202336
Not enough min models, keep training...
1
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7036455748105339
Not enough min models, keep training...
2
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.6926186362212684
Not enough min models, keep training...
3
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.6931495612398194
Not enough min models, keep training...
4
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.6868922438566547
Not enough min models, keep training...
5
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.6939204668987508
Not enough min models, keep training...
6
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.683987716591683
[0.6862041917202336, 0.7036455748105339, 0.6926186362212684, 0.6931495612398194, 0.6868922438566547, 0.6939204668987508, 0.683987716591683] 1
```

```
In [15]: 1 #print separator between two models' performances for better readability
2 print ('RNN alone: ')
3 testmodel(basemodel, 0)
4 print ('XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX')
5 print ('KSI+RNN: ')
6 testmodel(KSImodel, 1)
```

RNN alone:

test loss 0.05026585981249809

top- 10 0.5301316443299815

macro AUC 0.7298834993966542

micro AUC 0.9351245553636647

macro F1 0.03837643252372856

micro F1 0.28125

XX

KSI+RNN:

test loss 0.03798191621899605

top- 10 0.6998947962814728

macro AUC 0.8281500377577606

micro AUC 0.9660602134574177

macro F1 0.10300033466534458

micro F1 0.44059065934065933

In []:

1