

```
In [1]: 1 import torch
2 import torch.autograd as autograd
3 import torch.nn as nn
4 import torch.nn.functional as F
5 import torch.optim as optim
6 import numpy as np
7 torch.manual_seed(1)
8 from sklearn.metrics import roc_auc_score
9 from sklearn.metrics import f1_score
10 import copy
```

```
In [2]: 1 label_to_ix=np.load('label_to_ix.npy', allow_pickle=True).item()
2 ix_to_label=np.load('ix_to_label.npy', allow_pickle=True)
3 training_data=np.load('training_data.npy', allow_pickle=True)
4 test_data=np.load('test_data.npy', allow_pickle=True)
5 val_data=np.load('val_data.npy', allow_pickle=True)
6 word_to_ix=np.load('word_to_ix.npy', allow_pickle=True).item()
7 ix_to_word=np.load('ix_to_word.npy', allow_pickle=True)
8 newwikivec=np.load('newwikivec.npy', allow_pickle=True)
9 wikivoc=np.load('wikivoc.npy', allow_pickle=True).item()
```

```
In [3]: 1 wikisize=newwikivec.shape[0]
2 rvocsize=newwikivec.shape[1]
3 wikivec=autograd.Variable(torch.FloatTensor(newwikivec))
```

```
In [4]: 1 batchsize=32
```

```

In [5]: 1 def preprocessing(data):
2
3     new_data=[]
4     for i, note, j in data:
5         templabel=[0.0]*len(label_to_ix)
6         for jj in j:
7             if jj in wikivoc:
8                 templabel[label_to_ix[jj]]=1.0
9         templabel=np.array(templabel,dtype=float)
10        new_data.append((i, note, templabel))
11    new_data=np.array(new_data)
12
13    lenlist=[]
14    for i in new_data:
15        lenlist.append(len(i[0]))
16    sortlen=sorted(range(len(lenlist)), key=lambda k: lenlist[k])
17    new_data=new_data[sortlen]
18
19    batch_data=[]
20
21    for start_ix in range(0, len(new_data)-batchsize+1, batchsize):
22        thisblock=new_data[start_ix:start_ix+batchsize]
23        mybsize= len(thisblock)
24        numword=np.max([len(ii[0]) for ii in thisblock])
25        main_matrix = np.zeros((mybsize, numword), dtype= np.int)
26        for i in range(main_matrix.shape[0]):
27            for j in range(main_matrix.shape[1]):
28                try:
29                    if thisblock[i][0][j] in word_to_ix:
30                        main_matrix[i,j] = word_to_ix[thisblock[i][0][j]]
31
32                except IndexError:
33                    pass # because initialize with 0, so you pad with 0
34
35        xxx2=[]
36        yyy=[]
37        for ii in thisblock:
38            xxx2.append(ii[1])
39            yyy.append(ii[2])
40
41        xxx2=np.array(xxx2)
42        yyy=np.array(yyy)
43        batch_data.append((autograd.Variable(torch.from_numpy(main_matrix)),
44        return batch_data
45    batchtraining_data=preprocessing(training_data)
46    batchtest_data=preprocessing(test_data)
47    batchval_data=preprocessing(val_data)

```

/opt/conda/lib/python3.7/site-packages/ipykernel_launcher.py:11: 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

This is added back by InteractiveShellApp.init_path()

```
##### Create the  
model:
```

In [6]:

```

1 Embeddingsize=100
2 hidden_dim=200
3 class CAML(nn.Module):
4
5     def __init__(self, batch_size, vocab_size, tagset_size):
6         super(CAML, self).__init__()
7         self.hidden_dim = hidden_dim
8         self.word_embeddings = nn.Embedding(vocab_size+1, Embeddingsize, pad
9         self.embed_drop = nn.Dropout(p=0.2)
10
11
12         self.convs1 = nn.Conv1d(Embeddingsize,300,10,padding=5)
13         self.H=nn.Linear(300, tagset_size )
14         self.final = nn.Linear(300, tagset_size)
15
16         self.layer2 = nn.Linear(Embeddingsize, 1)
17         self.embedding=nn.Linear(rvocsize,Embeddingsize,bias=False)
18         self.vattention=nn.Linear(Embeddingsize,Embeddingsize)
19
20         self.sigmoid = nn.Sigmoid()
21         self.tanh = nn.Tanh()
22
23         self.dropout = nn.Dropout(0.2)
24
25     def forward(self, vec1, nvec, wiki, simlearning):
26
27
28         thisembeddings=self.word_embeddings(vec1)
29         thisembeddings = self.embed_drop(thisembeddings)
30         thisembeddings=thisembeddings.transpose(1,2)
31
32
33         thisembeddings=self.tanh(self.convs1(thisembeddings).transpose(1,2))
34
35         alpha=self.H.weight.matmul(thisembeddings.transpose(1,2))
36         alpha=F.softmax(alpha, dim=2)
37
38         m=alpha.matmul(thisembeddings)
39
40         myfinal=self.final.weight.mul(m).sum(dim=2).add(self.final.bias)
41
42         if simlearning==1:
43             nvec=nvec.view(batchsize,1,-1)
44             nvec=nvec.expand(batchsize,wiki.size()[0],-1)
45             wiki=wiki.view(1,wiki.size()[0],-1)
46             wiki=wiki.expand(nvec.size()[0],wiki.size()[1],-1)
47             new=wiki*nvec
48             new=self.embedding(new)
49             vattention=self.sigmoid(self.vattention(new))
50             new=new*vattention
51             vec3=self.layer2(new)
52             vec3=vec3.view(batchsize,-1)
53
54
55         if simlearning==1:
56             tag_scores = self.sigmoid(myfinal.detach()+vec3)

```

```
57         else:
58             tag_scores = self.sigmoid(myfinal)
59
60
61         return tag_scores
```

In [7]: 1 topk=10

```

In [8]: 1 def trainmodel(model, sim):
2         print ('start_training')
3         modelsaved=[]
4         modelperform=[]
5         topk=10
6
7
8         bestresults=-1
9         bestiter=-1
10        for epoch in range(1000):
11            model.train()
12
13            lossestrain = []
14            recall=[]
15            for mysentence in batchtraining_data:
16                model.zero_grad()
17
18                targets = mysentence[2].cuda()
19                tag_scores = model(mysentence[0].cuda(),mysentence[1].cuda(),wik
20                loss = loss_function(tag_scores, targets)
21                loss.backward()
22                optimizer.step()
23                lossestrain.append(loss.data.mean())
24            print (epoch)
25            modelsaved.append(copy.deepcopy(model.state_dict()))
26            print ("XXXXXXXXXXXXXXXXXXXXXXXXXXXX")
27            model.eval()
28
29            recall=[]
30            for inputs in batchval_data:
31
32                targets = inputs[2].cuda()
33                tag_scores = model(inputs[0].cuda(),inputs[1].cuda() ,wikivec.cu
34
35                loss = loss_function(tag_scores, targets)
36
37                targets=targets.data.cpu().numpy()
38                tag_scores= tag_scores.data.cpu().numpy()
39
40
41            for iii in range(0,len(tag_scores)):
42                temp={}
43                for iiii in range(0,len(tag_scores[iii])):
44                    temp[iiii]=tag_scores[iii][iiii]
45                temp1=[(k, temp[k]) for k in sorted(temp, key=temp.get, reve
46                thistop=int(np.sum(targets[iii]))
47                hit=0.0
48                for ii in temp1[0:max(thistop,topk)]:
49                    if targets[iii][ii[0]]==1.0:
50                        hit=hit+1
51                if thistop!=0:
52                    recall.append(hit/thistop)
53
54            print ('validation top-',topk, np.mean(recall))
55
56

```

```

57
58     modelperform.append(np.mean(recall))
59     if modelperform[-1]>bestresults:
60         bestresults=modelperform[-1]
61         besttiter=len(modelperform)-1
62
63     if (len(modelperform)-besttiter)>5:
64         print (modelperform,besttiter)
65     return modelsaved[besttiter]

```

```

In [9]: 1 model = CAML(batchsize, len(word_to_ix), len(label_to_ix))
        2 model.cuda()

```

```

Out[9]: CAML(
  (word_embeddings): Embedding(47961, 100, padding_idx=0)
  (embed_drop): Dropout(p=0.2, inplace=False)
  (convs1): Conv1d(100, 300, kernel_size=(10,), stride=(1,), padding=(5,))
  (H): Linear(in_features=300, out_features=344, bias=True)
  (final): Linear(in_features=300, out_features=344, bias=True)
  (layer2): Linear(in_features=100, out_features=1, bias=True)
  (embedding): Linear(in_features=12173, out_features=100, bias=False)
  (vattention): Linear(in_features=100, out_features=100, bias=True)
  (sigmoid): Sigmoid()
  (tanh): Tanh()
  (dropout): Dropout(p=0.2, inplace=False)
)

```

```

In [10]: 1 loss_function = nn.BCELoss()
        2 optimizer = optim.Adam(model.parameters())

```

```
In [11]: 1 basemodel= trainmodel(model, 0)
          2 torch.save(basemodel, 'CAML_model')
```

```
start_training
0
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.6028159410819657
1
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7314154270403253
2
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7724270743774396
3
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7882633693612231
4
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7974735633097602
5
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8054670879999204
6
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8080145665457611
7
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8089998811603987
8
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8116785014218055
9
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8111941689491625
10
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8129467900513655
11
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8127608618111046
12
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8126343213339993
13
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8113614322259091
14
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8121865293281534
15
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8125544505383168
[0.6028159410819657, 0.7314154270403253, 0.7724270743774396, 0.788263369361223
1, 0.7974735633097602, 0.8054670879999204, 0.8080145665457611, 0.80899988116039
87, 0.8116785014218055, 0.8111941689491625, 0.8129467900513655, 0.8127608618111
046, 0.8126343213339993, 0.8113614322259091, 0.8121865293281534, 0.812554450538
3168] 10
```


In [12]:

```

1 model = CAML(batchsize, len(word_to_ix), len(label_to_ix))
2 model.cuda()
3 model.load_state_dict(basemodel)
4 loss_function = nn.BCELoss()
5 optimizer = optim.Adam(model.parameters())
6 KSImodel= trainmodel(model, 1)
7 torch.save(KSImodel, 'KSI_CAML_model')

```

```
start_training
```

```

0
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.815991490746274
1
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8178064030108383
2
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.820040526535236
3
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8200844425246018
4
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8206450266622128
5
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8189880367972983
6
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8187936850792691
7
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8158021972977897
8
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8129353409313335
9
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8130162588698469
[0.815991490746274, 0.8178064030108383, 0.820040526535236, 0.8200844425246018,
0.8206450266622128, 0.8189880367972983, 0.8187936850792691, 0.8158021972977897,
0.8129353409313335, 0.8130162588698469] 4

```

```

In [13]: 1 def testmodel(modelstate, sim):
2         model = CAML(batchsize, len(word_to_ix), len(label_to_ix))
3         model.cuda()
4         model.load_state_dict(modelstate)
5         loss_function = nn.BCELoss()
6         model.eval()
7         recall=[]
8         lossestest = []
9
10        y_true=[]
11        y_scores=[]
12
13
14        for inputs in batchtest_data:
15
16            targets = inputs[2].cuda()
17
18            tag_scores = model(inputs[0].cuda(),inputs[1].cuda() ,wikivec.cuda())
19
20            loss = loss_function(tag_scores, targets)
21
22            targets=targets.data.cpu().numpy()
23            tag_scores= tag_scores.data.cpu().numpy()
24
25
26            lossestest.append(loss.data.mean())
27            y_true.append(targets)
28            y_scores.append(tag_scores)
29
30            for iii in range(0,len(tag_scores)):
31                temp={}
32                for iiii in range(0,len(tag_scores[iii])):
33                    temp[iiii]=tag_scores[iii][iiii]
34                    temp1=[(k, temp[k]) for k in sorted(temp, key=temp.get, reverse=
35                    thistop=int(np.sum(targets[iii]))
36                    hit=0.0
37
38                    for ii in temp1[0:max(thistop,topk)]:
39                        if targets[iii][ii[0]]==1.0:
40                            hit=hit+1
41                    if thistop!=0:
42                        recall.append(hit/thistop)
43            y_true=np.concatenate(y_true,axis=0)
44            y_scores=np.concatenate(y_scores,axis=0)
45            y_true=y_true.T
46            y_scores=y_scores.T
47            temptrue=[]
48            tempscores=[]
49            for col in range(0,len(y_true)):
50                if np.sum(y_true[col])!=0:
51                    temptrue.append(y_true[col])
52                    tempscores.append(y_scores[col])
53            temptrue=np.array(temptrue)
54            tempscores=np.array(tempscores)
55            y_true=temptrue.T
56            y_scores=tempscores.T

```

```

57 y_pred=(y_scores>0.5).astype(np.int)
58 print ('test loss', torch.stack(losses).mean().item())
59 print ('top-',topk, np.mean(recall))
60 print ('macro AUC', roc_auc_score(y_true, y_scores,average='macro'))
61 print ('micro AUC', roc_auc_score(y_true, y_scores,average='micro'))
62 print ('macro F1', f1_score(y_true, y_pred, average='macro') )
63 print ('micro F1', f1_score(y_true, y_pred, average='micro') )

```

```

In [14]: 1 print ('CAML alone:          ')
          2 testmodel(basemodel, 0)
          3 print ('XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX')
          4 print ('KSI+CAML:          ')
          5 testmodel(KSImodel, 1)

```

```

CAML alone:
test loss 0.03200652822852135
top- 10 0.8082940181977147
macro AUC 0.8528814144184469
micro AUC 0.9780247595956819
macro F1 0.2661915812890757
micro F1 0.6589893901367232
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
KSI+CAML:
test loss 0.031134463846683502
top- 10 0.8176190394571022
macro AUC 0.8884265849082379
micro AUC 0.9805001312980839
macro F1 0.2999242538289451
micro F1 0.6662273234847237

```

```

In [ ]: 1

```