

In [1]: 1 `%matplotlib inline`

In [2]: 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`

```

In [3]: 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()
10
11
12
13 wikisize=newwikivec.shape[0]
14 rvocsize=newwikivec.shape[1]
15 wikivec=autograd.Variable(torch.FloatTensor(newwikivec))
16
17 batchsize=32
18
19
20
21 def preprocessing(data):
22
23     new_data=[]
24     for i, note, j in data:
25         templabel=[0.0]*len(label_to_ix)
26         for jj in j:
27             if jj in wikivoc:
28                 templabel[label_to_ix[jj]]=1.0
29         templabel=np.array(templabel,dtype=float)
30         new_data.append((i, note, templabel))
31     new_data=np.array(new_data)
32
33     lenlist=[]
34     for i in new_data:
35         lenlist.append(len(i[0]))
36     sortlen=sorted(range(len(lenlist)), key=lambda k: lenlist[k])
37     new_data=new_data[sortlen]
38
39     batch_data=[]
40
41     for start_ix in range(0, len(new_data)-batchsize+1, batchsize):
42         thisblock=new_data[start_ix:start_ix+batchsize]
43         mybsize= len(thisblock)
44         numword=np.max([len(ii[0]) for ii in thisblock])
45         main_matrix = np.zeros((mybsize, numword), dtype= np.int)
46         for i in range(main_matrix.shape[0]):
47             for j in range(main_matrix.shape[1]):
48                 try:
49                     if thisblock[i][0][j] in word_to_ix:
50                         main_matrix[i,j] = word_to_ix[thisblock[i][0][j]]
51
52                 except IndexError:
53                     pass # because initialize with 0, so you pad with 0
54
55     xxx2=[]
56     yyy=[]

```

```
57         for ii in thisblock:
58             xxx2.append(ii[1])
59             yyy.append(ii[2])
60
61         xxx2=np.array(xxx2)
62         yyy=np.array(yyy)
63         batch_data.append((autograd.Variable(torch.from_numpy(main_matrix))),
64         return batch_data
65     batchtraining_data=preprocessing(training_data)
66     batchtest_data=preprocessing(test_data)
67     batchval_data=preprocessing(val_data)
68
69
70
```

/opt/conda/envs/python3.7/site-packages/ipykernel_launcher.py:31: 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

Create the model:

```

In [4]: 1 Embeddingsize=100
        2 hidden_dim=200
        3 class CNN(nn.Module):
        4
        5     def __init__(self, batch_size, vocab_size, tagset_size):
        6         super(CNN, 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         self.hidden2tag = nn.Linear(300, tagset_size)
        12
        13
        14         self.convs1 = nn.Conv1d(Embeddingsize,100,3)
        15         self.convs2 = nn.Conv1d(Embeddingsize,100,4)
        16         self.convs3 = nn.Conv1d(Embeddingsize,100,5)
        17
        18
        19         self.layer2 = nn.Linear(Embeddingsize, 1,bias=False)
        20         self.embedding=nn.Linear(rvocsize,Embeddingsize)
        21         self.vattention=nn.Linear(Embeddingsize,Embeddingsize)
        22
        23         self.sigmoid = nn.Sigmoid()
        24         self.tanh = nn.Tanh()
        25         self.dropout = nn.Dropout(0.2)
        26
        27     def forward(self, vec1, nvec, wiki, simlearning):
        28
        29         thisembeddings=self.word_embeddings(vec1)
        30         thisembeddings = self.embed_drop(thisembeddings)
        31         thisembeddings=thisembeddings.transpose(1,2)
        32
        33         output1=self.tanh(self.convs1(thisembeddings))
        34         output1=nn.MaxPool1d(output1.size()[2])(output1)
        35
        36         output2=self.tanh(self.convs2(thisembeddings))
        37         output2=nn.MaxPool1d(output2.size()[2])(output2)
        38
        39         output3=self.tanh(self.convs3(thisembeddings))
        40         output3=nn.MaxPool1d(output3.size()[2])(output3)
        41
        42         output4 = torch.cat([output1,output2,output3], 1).squeeze(2)
        43
        44         if simlearning==1:
        45             nvec=nvec.view(batchsize,1,-1)
        46             nvec=nvec.expand(batchsize,wiki.size()[0],-1)
        47             wiki=wiki.view(1,wiki.size()[0],-1)
        48             wiki=wiki.expand(nvec.size()[0],wiki.size()[1],-1)
        49             new=wiki*nvec
        50             new=self.embedding(new)
        51             vattention=self.sigmoid(self.vattention(new))
        52             new=new*vattention
        53             vec3=self.layer2(new)
        54             vec3=vec3.view(batchsize,-1)
        55
        56

```

```
57         vec2 = self.hidden2tag(output4)
58         if simlearning==1:
59             tag_scores = self.sigmoid(vec2.detach()+vec3)
60         else:
61             tag_scores = self.sigmoid(vec2)
62
63
64         return tag_scores
```

Train the model:

```

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

```

```

57
58
59
60     modelperform.append(np.mean(recall))
61     if modelperform[-1]>bestresults:
62         bestresults=modelperform[-1]
63         besttiter=len(modelperform)-1
64
65     if (len(modelperform)-besttiter)>5:
66         print (modelperform,besttiter)
67         return modelsaved[besttiter]
68
69 model = CNN(batchsize, len(word_to_ix), len(label_to_ix))
70 model.cuda()
71
72 loss_function = nn.BCELoss()
73 optimizer = optim.Adam(model.parameters())
74
75 basemodel= trainmodel(model, 0)
76 torch.save(basemodel, 'CNN_model')
77
78 model = CNN(batchsize, len(word_to_ix), len(label_to_ix))
79 model.cuda()
80 model.load_state_dict(basemodel)
81 loss_function = nn.BCELoss()
82 optimizer = optim.Adam(model.parameters())
83 KSImodel= trainmodel(model, 1)
84 torch.save(KSImodel, 'KSI_CNN_model')
85
86 def testmodel(modelstate, sim):
87     model = CNN(batchsize, len(word_to_ix), len(label_to_ix))
88     model.cuda()
89     model.load_state_dict(modelstate)
90     loss_function = nn.BCELoss()
91     model.eval()
92     recall=[]
93     lossestest = []
94
95     y_true=[]
96     y_scores=[]
97
98
99     for inputs in batchtest_data:
100
101         targets = inputs[2].cuda()
102
103         tag_scores = model(inputs[0].cuda(),inputs[1].cuda() ,wikivec.cuda())
104
105         loss = loss_function(tag_scores, targets)
106
107         targets=targets.data.cpu().numpy()
108         tag_scores= tag_scores.data.cpu().numpy()
109
110
111         lossestest.append(loss.data.mean())
112         y_true.append(targets)
113         y_scores.append(tag_scores)

```

```

114
115     for iii in range(0,len(tag_scores)):
116         temp={}
117         for iiii in range(0,len(tag_scores[iii])):
118             temp[iiii]=tag_scores[iii][iiii]
119         temp1=[(k, temp[k]) for k in sorted(temp, key=temp.get, reverse=
120         thistop=int(np.sum(targets[iii]))
121         hit=0.0
122
123         for ii in temp1[0:max(thistop,topk)]:
124             if targets[iii][ii[0]]==1.0:
125                 hit=hit+1
126             if thistop!=0:
127                 recall.append(hit/thistop)
128     y_true=np.concatenate(y_true,axis=0)
129     y_scores=np.concatenate(y_scores,axis=0)
130     y_true=y_true.T
131     y_scores=y_scores.T
132     temptrue=[]
133     tempscores=[]
134     for col in range(0,len(y_true)):
135         if np.sum(y_true[col])!=0:
136             temptrue.append(y_true[col])
137             tempscores.append(y_scores[col])
138     temptrue=np.array(temptrue)
139     tempscores=np.array(tempscores)
140     y_true=temptrue.T
141     y_scores=tempscores.T
142     y_pred=(y_scores>0.5).astype(np.int)
143     print ('test loss', torch.stack(lossesest).mean().item())
144     print ('top-',topk, np.mean(recall))
145     print ('macro AUC', roc_auc_score(y_true, y_scores,average='macro'))
146     print ('micro AUC', roc_auc_score(y_true, y_scores,average='micro'))
147     print ('macro F1', f1_score(y_true, y_pred, average='macro') )
148     print ('micro F1', f1_score(y_true, y_pred, average='micro') )
149
150 print ('CNN alone:          ')
151 testmodel(basemodel, 0)
152 print ('XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX')
153 print ('KSI+CNN:              ')
154 testmodel(KSImodel, 1)

```

```

start_training
0
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.4530655205021126
1
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.529073778538287
2
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.5752237865553222
3
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.6170548555043338
4

```



```
XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.6467904998800238
5
XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.6673111386677142
6
XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.6812441074304616
7
XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.6933596956867344
8
XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7052186240553053
9
XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7117489006787834
10
XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.722211327818717
11
XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7265657968833666
12
XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7334972301355923
13
XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7345698864388871
14
XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7347385604593721
15
XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7408095989262027
16
XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7420631058486372
17
XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7462089367848072
18
XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7501318022145766
19
XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7528999172460481
20
XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7534285951564181
21
XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7498198864222271
22
XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7532788638686233
23
```

```
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7540027411409296
24
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7573006911292467
25
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7541092517584836
26
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7552991371761425
27
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7589011476791745
28
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7577853193291497
29
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7583405388902464
30
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7560631161889517
31
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7577583138492285
32
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7598690505225836
33
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7560855487936969
34
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.756507376786185
35
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7564485341974435
36
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7546294812010665
37
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7564334654936617
[0.4530655205021126, 0.529073778538287, 0.5752237865553222, 0.6170548555043338,
0.6467904998800238, 0.6673111386677142, 0.6812441074304616, 0.6933596956867344,
0.7052186240553053, 0.7117489006787834, 0.722211327818717, 0.7265657968833666,
0.7334972301355923, 0.7345698864388871, 0.7347385604593721, 0.7408095989262027,
0.7420631058486372, 0.7462089367848072, 0.7501318022145766, 0.7528999172460481,
0.7534285951564181, 0.7498198864222271, 0.7532788638686233, 0.7540027411409296,
0.7573006911292467, 0.7541092517584836, 0.7552991371761425, 0.7589011476791745,
0.7577853193291497, 0.7583405388902464, 0.7560631161889517, 0.7577583138492285,
0.7598690505225836, 0.7560855487936969, 0.756507376786185, 0.7564485341974435,
0.7546294812010665, 0.7564334654936617] 32
start_training
0
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
```

```
validation top- 10 0.7669781847454861
1
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7734568574207574
2
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7768123484218621
3
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7797745185621022
4
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.780834126290448
5
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7800837381359998
6
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7795087269621294
7
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7794720793975843
8
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7781662562211387
9
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7783857329042567
[0.7669781847454861, 0.7734568574207574, 0.7768123484218621, 0.779774518562102
2, 0.780834126290448, 0.7800837381359998, 0.7795087269621294, 0.779472079397584
3, 0.7781662562211387, 0.7783857329042567] 4
CNN alone:
test loss 0.03949494659900665
top- 10 0.753674492427492
macro AUC 0.8306164538048904
micro AUC 0.9666265882181901
macro F1 0.21344930909581109
micro F1 0.6257646033743762
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
KSI+CNN:
test loss 0.03750767931342125
top- 10 0.7747399300742249
macro AUC 0.8609588009174831
micro AUC 0.9723801813466743
macro F1 0.24036851193298106
micro F1 0.6399044546248596
```

In []:

1