

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

```

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

/opt/conda/lib/python3.7/site-packages/ipykernel_launcher.py:30: 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 LSTM(nn.Module):
4
5     def __init__(self, batch_size, vocab_size, tagset_size):
6         super(LSTM, self).__init__()
7         self.hidden_dim = hidden_dim
8         self.word_embeddings = nn.Embedding(vocab_size+1, Embeddingsize, pad
9         self.lstm = nn.LSTM(Embeddingsize, hidden_dim)
10        self.hidden2tag = nn.Linear(hidden_dim, tagset_size)
11        self.hidden = self.init_hidden()
12
13
14        self.layer2 = nn.Linear(Embeddingsize, 1,bias=False)
15        self.embedding=nn.Linear(rvocsize,Embeddingsize)
16        self.vattention=nn.Linear(Embeddingsize,Embeddingsize,bias=False)
17
18        self.softmax = nn.Softmax()
19        self.sigmoid = nn.Sigmoid()
20        self.tanh = nn.Tanh()
21        self.embed_drop = nn.Dropout(p=0.2)
22
23    def init_hidden(self):
24        return (autograd.Variable(torch.zeros(1, batchsize, self.hidden_dim)
25                autograd.Variable(torch.zeros(1, batchsize, self.hidden_dim)
26
27
28    def forward(self, vec1, nvec, wiki, simlearning):
29
30        thisembeddings=self.word_embeddings(vec1).transpose(0,1)
31        thisembeddings = self.embed_drop(thisembeddings)
32
33        if simlearning==1:
34            nvec=nvec.view(batchsize,1,-1)
35            nvec=nvec.expand(batchsize,wiki.size()[0],-1)
36            wiki=wiki.view(1,wiki.size()[0],-1)
37            wiki=wiki.expand(nvec.size()[0],wiki.size()[1],-1)
38            new=wiki*nvec
39            new=self.embedding(new)
40            vattention=self.sigmoid(self.vattention(new))
41            new=new*vattention
42            vec3=self.layer2(new)
43            vec3=vec3.view(batchsize,-1)
44
45
46
47        lstm_out, self.hidden = self.lstm(
48            thisembeddings, self.hidden)
49
50        lstm_out=lstm_out.transpose(0,2).transpose(0,1)
51
52        output1=nn.MaxPool1d(lstm_out.size()[2])(lstm_out).view(batchsize,-1)
53
54        vec2 = self.hidden2tag(output1)
55        if simlearning==1:
56            tag_scores = self.sigmoid(vec2.detach()+vec3)

```

```
57         else:
58             tag_scores = self.sigmoid(vec2)
59
60
61         return tag_scores
```

Train the model:

```

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

```

```

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

```

```

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

```

```

start_training
0
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.42865152425063296
1
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.4937347971844057
2
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.5270009420629901
3
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.5647751687941437

```



```
4
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.6070973919648759
5
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.639208033930959
6
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.6591861390178926
7
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.6720871705039838
8
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.681262838674693
9
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.6952146312264086
10
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7048960649029746
11
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.709358782967236
12
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7133218440095446
13
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7188960280605515
14
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7238117730247063
15
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7336735645029677
16
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7360798024341662
17
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.739616167232288
18
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7447816545167835
19
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7469938538969276
20
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7477263734691333
21
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7520601792137857
22
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7541944230919304
```

```
23
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7524959970805342
24
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7563050451094965
25
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.754255788129518
26
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.757657513173346
27
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7557928129346213
28
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7579869778923267
29
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7589916913320011
30
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7606378455246832
31
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7589054716831615
32
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7592891570709687
33
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7614059946409589
34
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7627144764777577
35
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7627529203510008
36
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7618910381368089
37
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7627580779315861
38
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7659574118049436
39
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7631046816385884
40
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7647194594740986
41
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7645101229076846
```

```
42
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7651831962837605
43
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7651884552127561
[0.42865152425063296, 0.4937347971844057, 0.5270009420629901, 0.5647751687941
437, 0.6070973919648759, 0.639208033930959, 0.6591861390178926, 0.67208717050
39838, 0.681262838674693, 0.6952146312264086, 0.7048960649029746, 0.709358782
967236, 0.7133218440095446, 0.7188960280605515, 0.7238117730247063, 0.7336735
645029677, 0.7360798024341662, 0.739616167232288, 0.7447816545167835, 0.74699
38538969276, 0.7477263734691333, 0.7520601792137857, 0.7541944230919304, 0.75
24959970805342, 0.7563050451094965, 0.754255788129518, 0.757657513173346, 0.7
557928129346213, 0.7579869778923267, 0.7589916913320011, 0.7606378455246832,
0.7589054716831615, 0.7592891570709687, 0.7614059946409589, 0.762714476477757
7, 0.7627529203510008, 0.7618910381368089, 0.7627580779315861, 0.765957411804
9436, 0.7631046816385884, 0.7647194594740986, 0.7645101229076846, 0.765183196
2837605, 0.7651884552127561] 38
start_training
0
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7837404161967941
1
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7909443516204245
2
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7909998188906014
3
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7925701125486476
4
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7929621953952032
5
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7905842937402352
6
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7899778174580904
7
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7857877460663435
8
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7835120129886308
9
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7787964887872976
[0.7837404161967941, 0.7909443516204245, 0.7909998188906014, 0.79257011254864
76, 0.7929621953952032, 0.7905842937402352, 0.7899778174580904, 0.78578774606
63435, 0.7835120129886308, 0.7787964887872976] 4
LSTM alone:
test loss 0.03404521197080612
top- 10 0.7633712690864152
macro AUC 0.8337684970090297
micro AUC 0.9695393405722679
```

```
macro F1 0.1951925657098338
micro F1 0.6436915441878222
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
KSI+LSTM:
test loss 0.03159036487340927
top- 10 0.7901006497748037
macro AUC 0.868673097613561
micro AUC 0.976214050352376
macro F1 0.24646961919586005
micro F1 0.6546062107507058
```