In [1]:
```python
%matplotlib inline
```

In [2]:
```python
import torch
import torch.autograd as autograd
import torch.nn as nn
import torch.nn.functional as F
import torch.optim as optim
import numpy as np
torch.manual_seed(1)
from sklearn.metrics import roc_auc_score
from sklearn.metrics import f1_score
import copy
```

In [3]:

```python
label_to_ix=np.load('label_to_ix.npy', allow_pickle=True).item()
ix_to_label=np.load('ix_to_label.npy', allow_pickle=True)
training_data=np.load('training_data.npy', allow_pickle=True)
test_data=np.load('test_data.npy', allow_pickle=True)
val_data=np.load('val_data.npy', allow_pickle=True)
word_to_ix=np.load('word_to_ix.npy', allow_pickle=True).item()
ix_to_word=np.load('ix_to_word.npy', allow_pickle=True)
newwikivec=np.load('newwikivec.npy', allow_pickle=True)
wikivoc=np.load('wikivoc.npy', allow_pickle=True).item()

wikisize=newwikivec.shape[0]
rvocsize=newwikivec.shape[1]
wikivec=autograd.Variable(torch.FloatTensor(newwikivec))

batchsize=32



def preprocessing(data):

    new_data=[]
    for i, note, j in data:
        templabel=[0.0]*len(label_to_ix)
        for jj in j:
            if jj in wikivoc:
                templabel[label_to_ix[jj]]=1.0
        templabel=np.array(templabel,dtype=float)
        new_data.append((i, note, templabel))
    new_data=np.array(new_data)

    lenlist=[]
    for i in new_data:
        lenlist.append(len(i[0]))
    sortlen=sorted(range(len(lenlist)), key=lambda k: lenlist[k])
    new_data=new_data[sortlen]

    batch_data=[]

    for start_ix in range(0, len(new_data)-batchsize+1, batchsize):
        thisblock=new_data[start_ix:start_ix+batchsize]
        mybsize= len(thisblock)
        numword=np.max([len(ii[0]) for ii in thisblock])
        main_matrix = np.zeros((mybsize, numword), dtype= np.int)
        for i in range(main_matrix.shape[0]):
            for j in range(main_matrix.shape[1]):
                try:
                    if thisblock[i][0][j] in word_to_ix:
                        main_matrix[i,j] = word_to_ix[thisblock[i][0][j]]

                except IndexError:
                    pass        # because initialze with 0, so you pad with 0

        xxx2=[]
        yyy=[]
        for ii in thisblock:
            xxx2.append(ii[1])
```

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

/opt/conda/lib/python3.7/site-packages/ipykernel_launcher.py:29: VisibleDepreca
tionWarning: 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 de
precated. If you meant to do this, you must specify 'dtype=object' when creatin
g the ndarray

Create the model:

In [4]:

```python
Embeddingsize=100
hidden_dim=200
class LSTMattn(nn.Module):

    def __init__(self, batch_size, vocab_size, tagset_size):
        super(LSTMattn, self).__init__()
        self.hidden_dim = hidden_dim
        self.word_embeddings = nn.Embedding(vocab_size+1, Embeddingsize, pad
        self.lstm = nn.LSTM(Embeddingsize, hidden_dim)
        self.hidden = self.init_hidden()

        self.H=nn.Linear(hidden_dim, tagset_size )
        self.final = nn.Linear(hidden_dim, tagset_size)

        self.layer2 = nn.Linear(Embeddingsize, 1,bias=False)
        self.embedding=nn.Linear(rvocsize,Embeddingsize)
        self.vattention=nn.Linear(Embeddingsize,Embeddingsize,bias=False)

        self.softmax = nn.Softmax()
        self.sigmoid = nn.Sigmoid()
        self.embed_drop = nn.Dropout(p=0.2)

    def init_hidden(self):
        return (autograd.Variable(torch.zeros(1, batchsize, self.hidden_dim)
                autograd.Variable(torch.zeros(1, batchsize, self.hidden_dim)


    def forward(self, vec1, nvec, wiki, simlearning):


        thisembeddings=self.word_embeddings(vec1).transpose(0,1)
        thisembeddings = self.embed_drop(thisembeddings)


        if simlearning==1:
            nvec=nvec.view(batchsize,1,-1)
            nvec=nvec.expand(batchsize,wiki.size()[0],-1)
            wiki=wiki.view(1,wiki.size()[0],-1)
            wiki=wiki.expand(nvec.size()[0],wiki.size()[1],-1)
            new=wiki*nvec
            new=self.embedding(new)
            vattention=self.sigmoid(self.vattention(new))
            new=new*vattention
            vec3=self.layer2(new)
            vec3=vec3.view(batchsize,-1)


        lstm_out, self.hidden = self.lstm(
            thisembeddings, self.hidden)



        lstm_out=lstm_out.transpose(0,1)

        alpha=self.H.weight.matmul(lstm_out.transpose(1,2))
        alpha=F.softmax(alpha, dim=2)
```

```
57
58          m=alpha.matmul(lstm_out)
59
60          myfinal=self.final.weight.mul(m).sum(dim=2).add(self.final.bias)
61
62
63          if simlearning==1:
64              tag_scores = self.sigmoid(myfinal.detach()+vec3)
65          else:
66              tag_scores = self.sigmoid(myfinal)
67
68
69          return tag_scores
```

Train the model:

In [5]:

```python
topk=10

def trainmodel(model, sim):
    print ('start_training')
    modelsaved=[]
    modelperform=[]
    topk=10


    bestresults=-1
    bestiter=-1
    for epoch in range(1000):

        model.train()

        lossestrain = []
        recall=[]
        for mysentence in batchtraining_data:
            model.zero_grad()
            model.hidden = model.init_hidden()
            targets = mysentence[2].cuda()
            tag_scores = model(mysentence[0].cuda(),mysentence[1].cuda(),wik
            loss = loss_function(tag_scores, targets)
            loss.backward()
            optimizer.step()
            lossestrain.append(loss.data.mean())
        print (epoch)
        modelsaved.append(copy.deepcopy(model.state_dict()))
        print ("XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX")
        model.eval()

        recall=[]
        for inputs in batchval_data:
            model.hidden = model.init_hidden()
            targets = inputs[2].cuda()
            tag_scores = model(inputs[0].cuda(),inputs[1].cuda() ,wikivec.cu

            loss = loss_function(tag_scores, targets)

            targets=targets.data.cpu().numpy()
            tag_scores= tag_scores.data.cpu().numpy()


            for iii in range(0,len(tag_scores)):
                temp={}
                for iiii in range(0,len(tag_scores[iii])):
                    temp[iiii]=tag_scores[iii][iiii]
                temp1=[(k, temp[k]) for k in sorted(temp, key=temp.get, reve
                thistop=int(np.sum(targets[iii]))
                hit=0.0
                for ii in temp1[0:max(thistop,topk)]:
                    if targets[iii][ii[0]]==1.0:
                        hit=hit+1
                if thistop!=0:
                    recall.append(hit/thistop)
```

```python
57            print ('validation top-',topk, np.mean(recall))
58
59
60
61          modelperform.append(np.mean(recall))
62          if modelperform[-1]>bestresults:
63              bestresults=modelperform[-1]
64              bestiter=len(modelperform)-1
65
66          if (len(modelperform)-bestiter)>5:
67              print (modelperform,bestiter)
68              return modelsaved[bestiter]
69
70  model = LSTMattn(batchsize, len(word_to_ix), len(label_to_ix))
71  model.cuda()
72  loss_function = nn.BCELoss()
73  optimizer = optim.Adam(model.parameters())
74
75  basemodel= trainmodel(model, 0)
76  torch.save(basemodel, 'LSTMattn_model')
77
78  model = LSTMattn(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_LSTMattn_model')
85
86  def testmodel(modelstate, sim):
87      model = LSTMattn(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         model.hidden = model.init_hidden()
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)
```

```python
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(lossestest).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 ('LSTMattn alone:              ')
151 testmodel(basemodel, 0)
152 print ('XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX')
153 print ('KSI+LSTMattn:            ')
154 testmodel(KSImodel, 1)
```

```
start_training
0
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.37988820337654555
1
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.5104482603359948
2
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.6266978811026137
3
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.6595165288696112
4
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
```

```
                    validation top- 10 0.6859925783877213
            5
            XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
                    validation top- 10 0.7074581327207593
            6
            XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
                    validation top- 10 0.7279008211626297
            7
            XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
                    validation top- 10 0.7452480736993783
            8
            XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
                    validation top- 10 0.759016435522099
            9
            XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
                    validation top- 10 0.7637212570536341
            10
            XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
                    validation top- 10 0.7705378830020023
            11
            XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
                    validation top- 10 0.7759446183764552
            12
            XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
                    validation top- 10 0.7818910104256545
            13
            XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
                    validation top- 10 0.7873337243834775
            14
            XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
                    validation top- 10 0.7897661533039431
            15
            XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
                    validation top- 10 0.7896431993888611
            16
            XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
                    validation top- 10 0.7931866940027331
            17
            XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
                    validation top- 10 0.7941931676870783
            18
            XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
                    validation top- 10 0.7928390644189163
            19
            XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
                    validation top- 10 0.7963826490295235
            20
            XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
                    validation top- 10 0.7961022934012223
            21
            XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
                    validation top- 10 0.7951516205182724
            22
            XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
                    validation top- 10 0.7973918889129703
            23
            XXXXXXXXXXXXXXXXXXXXXXXXXXXXX
```

```
validation top- 10 0.7948993508422044
24
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7938015279237891
25
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7935765049264573
26
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7928097294806639
27
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.7928827429214348
[0.37988820337654555, 0.5104482603359948, 0.6266978811026137, 0.6595165288696
112, 0.6859925783877213, 0.7074581327207593, 0.7279008211626297, 0.7452480736
993783, 0.759016435522099, 0.7637212570536341, 0.7705378830020023, 0.77594461
83764552, 0.7818910104256545, 0.7873337243834775, 0.7897661533039431, 0.78964
31993888611, 0.7931866940027331, 0.7941931676870783, 0.7928390644189163, 0.79
63826490295235, 0.7961022934012223, 0.7951516205182724, 0.7973918889129703,
0.7948993508422044, 0.7938015279237891, 0.7935765049264573, 0.792809729480663
9, 0.7928827429214348] 22
start_training
0
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8044314961620841
1
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8067662100662395
2
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8087511717170007
3
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8093848491817769
4
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8091355231047173
5
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8092732886573508
6
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8074876173077044
7
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8054115864744192
8
XXXXXXXXXXXXXXXXXXXXXXXXXXXX
validation top- 10 0.8041572170163643
[0.8044314961620841, 0.8067662100662395, 0.8087511717170007, 0.80938484918177
69, 0.8091355231047173, 0.8092732886573508, 0.8074876173077044, 0.80541158647
44192, 0.8041572170163643] 3
LSTMattn alone:
test loss 0.03530406951904297
top- 10 0.7900533901928264
macro AUC 0.8469555433080055
micro AUC 0.974202738572247
```

```
macro F1 0.253806770132122
micro F1 0.6438496873279481
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
KSI+LSTMattn:
test loss 0.03251442313194275
top- 10 0.804811802206705
macro AUC 0.8848782941414216
micro AUC 0.9778077616901655
macro F1 0.3016705252258659
micro F1 0.6616050763935742
```

Type *Markdown* and LaTeX: $\alpha^2$