break print('Loaded %s pretrained word vectors' % (len(self.word2vec))) def most similar(self, w, K=5): # K most similar words: self.score - np.argsort keys = list(self.word2vec.keys()) values = list(self.word2vec.values()) scores = list(map(lambda x: self.score(x,w) if x!=w else -1,keys)) arg = np.argsort(scores) return [keys[i] for i in np.flipud(arg[-K:])] def score(self, w1, w2): # cosine similarity: np.dot - np.linalg.norm vec1 = self.word2vec[w1] vec2 = self.word2vec[w2] return np.dot(vec1, vec2) / (np.linalg.norm(vec1) \*np.linalg.norm(vec2)) In [4]: | w2v = Word2vec(os.path.join(PATH TO DATA, 'crawl-300d-200k.vec'), nmax=100000) # You will be evaluated on the output of the following: for w1, w2 in zip(('cat', 'dog', 'dogs', 'paris', 'germany'), ('dog', 'pet', 'cats', 'france', 'berlin')): print(w1, w2, w2v.score(w1, w2)) for w1 in ['cat', 'dog', 'dogs', 'paris', 'germany']: print(w2v.most\_similar(w1)) Loaded 100000 pretrained word vectors cat dog 0.6716836662792491 dog pet 0.6842064029669219 dogs cats 0.7074389328052404 paris france 0.7775108541288561 germany berlin 0.7420295235998392 ['cats', 'kitty', 'kitten', 'feline', 'kitties'] ['dogs', 'puppy', 'Dog', 'doggie', 'canine'] ['dog', 'pooches', 'Dogs', 'doggies', 'canines'] ['france', 'Paris', 'london', 'berlin', 'tokyo'] ['austria', 'europe', 'german', 'berlin', 'poland'] In [5]: class BoV(): def init\_\_(self, w2v): self.w2v = w2vdef encode(self, sentences, idf=False): # takes a list of sentences, outputs a numpy array of sentence embeddings # see TP1 for help sentemb = [] for sent in sentences: if idf is False: # mean of word vectors keys = list(self.w2v.word2vec.keys()) sentemb.append(np.mean([self.w2v.word2vec[w] for w in sent if w in keys]\ , axis=0)) #assert False, 'TODO: fill in the blank' else: # idf-weighted mean of word vectors keys = list(self.w2v.word2vec.keys()) vec = np.sum([idf[w]\*self.w2v.word2vec[w] for w in sent if w in keys]\ , axis=0)/np.sum([idf[w] for w in sent if w in keys]) sentemb.append(vec) #assert False, 'TODO: fill in the blank' return np.vstack(sentemb) def most\_similar(self, s, sentences, idf=False, K=5): # get most similar sentences and \*\*print\*\* them #keys = self.encode(sentences, idf) #query = self.encode([s], idf)scores = list(map(lambda x: self.score(x,s) if x!=s else -1, sentences)) arg = np.argsort(scores) most\_sim = [sentences[i] for i in np.flipud(arg[-K:])] print('The most similar', K, 'sentences to: "',' '.join(s),'" are:\n') for i in range(len(most sim)): print('Rank ',i+1 , ': ', ' '.join(most\_sim[i])) def score(self, s1, s2, idf=False): s1 encode = self.encode([s1], idf) s2 encode = self.encode([s2], idf) return float(np.dot(s1\_encode, s2\_encode.T) / (np.linalg.norm(s1\_encode) \*np.linalg.norm(s2 encode))) def build idf(self, sentences): # build the idf dictionary: associate each word to its idf value for sent in sentences: for w in set(sent): # Using set to avoid word repetition in a sentence idf[w] = idf.get(w, 0) + 1for w in list(idf.keys()): idf[w] = max(1, np.log10(len(sentences) / (idf[w])))return idf In [687]: L = train\_sent+dev\_sent+test\_sent In [6]: | w2v = Word2vec(os.path.join(PATH TO DATA, 'crawl-300d-200k.vec'), nmax=100000) s2v = BoV(w2v)# Load sentences in "PATH\_TO\_DATA/sentences.txt" with open('/Users/adilrhoulam/Downloads/data-2/sentences.txt') as f: for i, line in enumerate(f): sent = line.rstrip().split() sentences.append(sent) # Build idf scores for each word idf = s2v.build idf(sentences) # You will be evaluated on the output of the following: s2v.most\_similar('' if not sentences else sentences[10], sentences) # BoV-mean score = s2v.score('' if not sentences else sentences[7], '' if not sentences else sentences[13] print('\n- Cosinus score using BoV-mean between "',' '.join(sentences[7])\ ,'" and "',' '.join(sentences[13]),'" is :', score,' $\n'$ ) s2v.most\_similar('' if not sentences else sentences[10], sentences, idf) # BoV-idf score = s2v.score('' if not sentences else sentences[7], '' if not sentences else sentences[13] print('\n- Cosinus score using BoV-idf between "',' '.join(sentences[7])\ ,'" and "',' '.join(sentences[13]),'" is :', score) Loaded 100000 pretrained word vectors The most similar 5 sentences to: " 1 smiling african american boy . " are: Rank 1: an african american man smiling . Rank 2: a little african american boy and girl looking up . Rank 3: an afican american woman standing behind two small african american children . Rank 4: an african american man is sitting. Rank 5: a girl in black hat holding an african american baby . - Cosinus score using BoV-mean between " 1 man singing and 1 man playing a saxophone in a co ncert . " and " 10 people venture out to go crosscountry skiing . " is : 0.5726258859719606 The most similar 5 sentences to: " 1 smiling african american boy . " are: Rank 1: an african american man smiling . Rank 2: a little african american boy and girl looking up . Rank 3: an afican american woman standing behind two small african american children . Rank 4: an african american man is sitting. Rank 5: a girl in black hat holding an african american baby . - Cosinus score using BoV-idf between " 1 man singing and 1 man playing a saxophone in a con cert . " and " 10 people venture out to go crosscountry skiing . " is : 0.4751450875368781 2) Multilingual (English-French) word embeddings Let's consider a bilingual dictionary of size V\_a (e.g French-English). Let's define **X** and **Y** the **French** and **English** matrices. They contain the embeddings associated to the words in the bilingual dictionary. We want to find a **mapping W** that will project the source word space (e.g French) to the target word space (e.g English). Procrustes: W\* = argmin || W.X - Y || s.t W^T.W = Id has a closed form solution: W = U.V^T where U.Sig.V^T = SVD(Y.X^T) In [7]: # 1 - Download and load 50k first vectors of https://s3-us-west-1.amazonaws.com/fasttext-vectors/wiki.en.vec https://s3-us-west-1.amazonaws.com/fasttext-vectors/wiki.fr.vec # TYPE CODE HERE nmax = 50000engvec = {} with io.open('/Users/adilrhoulam/Downloads/wiki.en.vec', encoding='utf-8') as f: for i, line in enumerate(f): word, vec = line.split(' ', 1) engvec[word] = np.fromstring(vec, sep=' ') **if** i == (nmax - 1): frvec = {} with io.open('/Users/adilrhoulam/Downloads/wiki.fr.vec', encoding='utf-8') as f: for i, line in enumerate(f): word, vec = line.split(' ', 1) frvec[word] = np.fromstring(vec, sep=' ') **if** i == (nmax - 1): break In [8]: | # 2 - Get words that appear in both vocabs (= identical character strings) # Use it to create the matrix X and Y (of aligned embeddings for these words) # TYPE CODE HERE identical words = set() keys eng = set(engvec.keys()) keys fr = set(frvec.keys()) for w in keys eng: if w in keys fr: identical words.add(w) X = []Y = []for w in identical words: X.append(frvec[w]) Y.append(engvec[w]) X,Y = np.array(X), np.array(Y)In [9]: # 3 - Solve the Procrustes using the scipy package and: scipy.linalg.svd() and get the optimal # Now W\*French vector is in the same space as English vector # TYPE CODE HERE R = np.dot(Y, X.T)U, S, V = np.linalg.svd(R) #Python3 svd has become in the numpy.linalg package W = np.dot(U, V)In [10]: # 4 - After alignment with W, give examples of English nearest neighbors of some French words # You will be evaluated on that part and the code above # TYPE CODE HERE words = list(identical words) # French words nearest neighbors \* \*\*\*\*\*\* ) A = np.dot(W, X)def nearest neighbors fr(A,Y,index,words,K=5): vec = A[index]scores = list(map(lambda x: float(np.dot(x,vec)/(np.linalg.norm(x)\*np.linalg.norm(vec)))\ ,list(Y))) arg = np.argsort(scores) return [words[i] for i in np.flipud(arg[-K:])] index = words.index('place') nearest = nearest neighbors fr(A,Y,index,words) print('\nThe nearest 5 english words neighbors to the french word "place" are:\n '\ , nearest) index = words.index('costumes') nearest = nearest neighbors fr(A,Y,index,words) print(' $\n$ The nearest 5 english words neighbors to the french word "costumes" are: $\n$  ' $\n$ index = words.index('bizarre') nearest = nearest neighbors fr(A,Y,index,words) print('\nThe nearest 5 english words neighbors to the french word "bizarre" are:\n '\ index = words.index('programme') nearest = nearest neighbors fr(A,Y,index,words) print('\nThe nearest 5 english words neighbors to the french word "programme" are:\n '\ , nearest) index = words.index('populations') nearest = nearest neighbors fr(A,Y,index,words) print(' $\n$ The nearest 5 english words neighbors to the french word "populations" are: $\n$ ' ,nearest) index = words.index('automobiles') nearest = nearest\_neighbors\_fr(A,Y,index,words) print('\nThe nearest 5 english words neighbors to the french word "automobiles" are:\n '\ index = words.index('liaison') nearest = nearest\_neighbors\_fr(A,Y,index,words) print('\nThe nearest 5 english words neighbors to the french word "liaison" are:\n '\ , nearest) \*\*\*\*\*\*\*\*\*\*) def nearest neighbors eng(A,Y,index,words,K=5): vec = Y[index]scores = list(map(lambda x: float(np.dot(x,vec)/(np.linalg.norm(x)\*np.linalg.norm(vec)))\ ,list(A))) arg = np.argsort(scores) return [words[i] for i in np.flipud(arg[-K:])] index = words.index('motel') nearest = nearest neighbors eng(A,Y,index,words) print('\nThe nearest 5 french words neighbors to the english word "motel" are:\n '\ ,nearest) index = words.index('liberation') nearest = nearest neighbors eng(A,Y,index,words) print('\nThe nearest 5 french words neighbors to the english word "liberation" are:\n '\ , nearest) index = words.index('fast') nearest = nearest neighbors eng(A,Y,index,words) print('\nThe nearest 5 french words neighbors to the english word "fast" are:\n '\ , nearest) index = words.index('weapons') nearest = nearest\_neighbors\_eng(A,Y,index,words) print(' $\n$ The nearest 5 french words neighbors to the english word "weapons" are: $\n$  ' $\n$ ,nearest) index = words.index('cultures') nearest = nearest\_neighbors\_eng(A,Y,index,words) print(' $\n$ The nearest 5 french words neighbors to the english word "cultures" are: $\n$  ' $\n$ , nearest) index = words.index('food') nearest = nearest\_neighbors\_eng(A,Y,index,words) print('\nThe nearest 5 french words neighbors to the english word "food" are:\n '\ index = words.index('freedom') nearest = nearest neighbors eng(A,Y,index,words) print('\nThe nearest 5 french words neighbors to the english word "freedom" are:\n '\ The nearest 5 english words neighbors to the french word "place" are: ['place', 'places', 'spot', 'takes', 'taking'] The nearest 5 english words neighbors to the french word "costumes" are: ['costumes', 'costume', 'attire', 'dress', 'shirts'] The nearest 5 english words neighbors to the french word "bizarre" are: ['bizarre', 'strange', 'weird', 'horrible', 'grotesque'] The nearest 5 english words neighbors to the french word "programme" are: ['programme', 'programmes', 'program', 'bbc', 'itv'] The nearest 5 english words neighbors to the french word "populations" are: ['populations', 'population', 'migrations', 'migrants', 'areas'] The nearest 5 english words neighbors to the french word "automobiles" are: ['automobiles', 'automobile', 'cars', 'automotive', 'trucks'] The nearest 5 english words neighbors to the french word "liaison" are: ['liaison', 'liaisons', 'coordination', 'staff', 'aide'] \*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\* English words nearest neighbors \* The nearest 5 french words neighbors to the english word "motel" are: ['motel', 'hotel', 'hotels', 'diner', 'restaurant'] The nearest 5 french words neighbors to the english word "liberation" are: ['liberation', 'libération', 'revolutionary', 'independence', 'front'] The nearest 5 french words neighbors to the english word "fast" are: ['fast', 'slow', 'pace', 'speed', 'agile'] The nearest 5 french words neighbors to the english word "weapons" are: ['weapons', 'weapon', 'munitions', 'projectiles', 'missiles'] The nearest 5 french words neighbors to the english word "cultures" are: ['cultures', 'culture', 'traditions', 'peoples', 'civilization'] The nearest 5 french words neighbors to the english word "food" are: ['food', 'foods', 'nutrition', 'meat', 'desserts'] The nearest 5 french words neighbors to the english word "freedom" are: ['freedom', 'liberty', 'conscience', 'democracy', 'rights'] If you want to dive deeper on this subject: <a href="https://github.com/facebookresearch/MUSE">https://github.com/facebookresearch/MUSE</a> 3) Sentence classification with BoV and scikit-learn In [692]: # 1 - Load train/dev/test of Stanford Sentiment TreeBank (SST) # (https://nlp.stanford.edu/~socherr/EMNLP2013\_RNTN.pdf) # TYPE CODE HERE train sent = [] train labels = [] path = '/Users/adilrhoulam/Downloads/nlp project-3/data/SST/stsa.fine.train' with io.open(path, encoding='utf-8') as f: for line in f: label, sent = line.split(' ', 1) train sent.append(sent.rstrip().split()) train labels.append(int(label)) dev sent = [] dev labels = []path = '/Users/adilrhoulam/Downloads/nlp project-3/data/SST/stsa.fine.dev' with io.open(path, encoding='utf-8') as f: for line in f: label, sent = line.split(' ', 1) dev sent.append(sent.rstrip().split()) dev labels.append(int(label)) test sent = [] path = '/Users/adilrhoulam/Downloads/nlp project-3/data/SST/stsa.fine.test.X' with io.open(path, encoding='utf-8') as f: for line in f: test\_sent.append(line.rstrip().split()) In [699]: # 2 - Encode sentences with the BoV model above # TYPE CODE HERE PATH TO DATA = "/Users/adilrhoulam/Downloads/data-2" w2v = Word2vec(os.path.join(PATH TO DATA, 'crawl-300d-200k.vec'), nmax=100000) train\_sent\_vec\_mean = s2v.encode(train\_sent) dev sent vec mean = s2v.encode(dev sent)test sent vec mean = s2v.encode(test sent) idf = s2v.build\_idf(train\_sent+dev\_sent+test\_sent) train sent vec idf = s2v.encode(train sent,idf) dev sent vec idf = s2v.encode(dev sent,idf) test\_sent\_vec\_idf = s2v.encode(test\_sent,idf) Loaded 100000 pretrained word vectors In [940]: | # 3 - Learn Logistic Regression on top of sentence embeddings using scikit-learn (consider tuning the L2 regularization on the dev set) # TYPE CODE HERE from sklearn.linear\_model import LogisticRegression C = [0.001, 0.01, 0.05, 0.1, 0.5, 1, 3, 10, 20, 30, 40, 50, 60, 70, 80, 90, 100, 200, 500, 1000]scores = [] for c in C: logReg = LogisticRegression(penalty='12',C = c) logReg.fit(train sent vec mean, train labels) score = logReg.score(dev\_sent\_vec\_mean,dev\_labels) scores.append(score) c1 = C[np.argmax(scores)] print('The best score on dev set when using mean BoV is: ',round(np.max(scores)\*100,2),'% and is given by given by a regularization parameter C = ',c1) scores = [] for c in C: logReg = LogisticRegression(penalty='12',C = c) logReg.fit(train sent vec idf,train labels) score = logReg.score(dev sent vec idf, dev labels) scores.append(score) c2 = C[np.argmax(scores)] print('The best score on dev set when using idf BoV is: ',round(np.max(scores)\*100,2),'% \ and is given by given by a regularization parameter C = ',c2) logReg = LogisticRegression(penalty='12',C = c1) logReg.fit(train\_sent\_vec\_mean,train\_labels) score = logReg.score(train sent vec mean, train labels) print('The score on train data using the tuned parameter on mean BoV is :',\ round(score\*100,2)) logReg = LogisticRegression(penalty='12',C = c2) logReg.fit(train sent vec idf,train labels) score = logReg.score(train\_sent\_vec\_idf,train\_labels) print('The score on train data using the tuned parameter on idf BoV is :',\ round(score\*100,2)) The best score on dev set when using mean BoV is: 43.69 % and is given by given by a regul arization parameter C = 0.5The best score on dev set when using idf BoV is : 42.51 % and is given by given by a regula rization parameter C = 0.5The score on train data using the tuned parameter on mean BoV is : 47.94 The score on train data using the tuned parameter on idf BoV is: 48.68 In [777]: # 4 - Produce 2210 predictions for the test set (in the same order). One line = one prediction (=0,1,2,3,4). # Attach the output file "logreg bov y test sst.txt" to your deliverable. You will be evaluated on the results of the test set. # TYPE CODE HERE path = '/Users/adilrhoulam/Downloads/nlp project-3/logreg bov y test sst.txt' thefile = open(path, 'w') logReg = LogisticRegression(penalty='12',C = 20) logReg.fit(train\_sent\_vec\_mean,train\_labels) predictions = logReg.predict(test sent vec mean) for i in range(len(test sent vec mean)): thefile.write("%s\n" % predictions[i]) In [942]: # BONUS! # 5 - Try to improve performance with another classifier # Attach the output file "XXX bov y test sst.txt" to your deliverable (where XXX = the nam e of the classifier) # TYPE CODE HERE #Neural Network from sklearn.svm import SVC clf = SVC(C=200, cache size=250, class weight=None, coef0=0.0,  $\setminus$ decision function shape='sigmoid', degree=4, gamma='auto', kernel='rbf',\ max iter=-1, probability=False, random state=None, shrinking=True,\ tol=0.1, verbose=**False**) clf.fit(train sent vec mean, train labels) clf.score(dev sent vec mean, dev labels) Out [942]: 0.43869209809264303 In [974]: path = '/Users/adilrhoulam/Downloads/nlp project-3/SVM bov y test sst.txt' thefile = open(path, 'w') predictions = clf.predict(test sent vec mean) for i in range(len(test\_sent\_vec\_mean)): thefile.write("%s\n" % predictions[i]) 4) Sentence classification with LSTMs in Keras 4.1 - Preprocessing In [860]: import keras In [861]: # 1 - Load train/dev/test sets of SST PATH\_TO\_DATA = '/Users/adilrhoulam/Downloads/nlp\_project-3/data/' # TYPE CODE HERE train sent = [] train\_labels = [] path = PATH\_TO\_DATA+'SST/stsa.fine.train' with io.open(path, encoding='utf-8') as f: for line in f: label, sent = line.split(' ', 1) train\_sent.append(sent) train labels.append(int(label)) dev sent = [] dev labels = []path = PATH TO DATA+'SST/stsa.fine.dev' with io.open(path, encoding='utf-8') as f: for line in f: label, sent = line.split(' ', 1) dev sent.append(sent) dev labels.append(int(label)) test sent = [] path = PATH TO DATA+'SST/stsa.fine.test.X' with io.open(path, encoding='utf-8') as f: for line in f: test\_sent.append(line) In [862]: # 2 - Transform text to integers using keras.preprocessing.text.one hot function # https://keras.io/preprocessing/text/ # TYPE CODE HERE from keras.preprocessing import text train\_sent = list(map(lambda x: text.one\_hot(x,len(x)),train\_sent)) dev\_sent = list(map(lambda x: text.one\_hot(x,len(x)),dev\_sent)) test sent = list(map(lambda x: text.one hot(x,len(x)),test sent)) Padding input data Models in Keras (and elsewhere) take batches of sentences of the same length as input. It is because Deep Learning framework have been designed to handle well Tensors, which are particularly suited for fast computation on the GPU. Since sentences have different sizes, we "pad" them. That is, we add dummy "padding" tokens so that they all have the same length. The input to a Keras model thus has this size: (batchsize, maxseglen) where maxseglen is the maximum length of a sentence in the batch. In [864]: | # 3 - Pad your sequences using keras.preprocessing.sequence.pad\_sequences # https://keras.io/preprocessing/sequence/ # TYPE CODE HERE from keras.preprocessing import sequence vec len = np.vectorize(len) maxseqlen = max(max(vec len(train sent)), max(vec len(dev sent)), max(vec len(test sent))) x train = sequence.pad sequences(train sent, maxlen=maxseqlen) x\_dev = sequence.pad\_sequences(dev\_sent,maxlen=maxseqlen) x\_test = sequence.pad\_sequences(test\_sent,maxlen=maxseqlen) 4.2 - Design and train your model In [954]: # 4 - Design your encoder + classifier using keras.layers # In Keras, Torch and other deep learning framework, we create a "container" which is the Sequential() module. Then we add components to this contained : the lookuptable, the LSTM, the classifier et All of these components are contained in the Sequential() and are trained together. # ADAPT CODE BELOW from keras.models import Sequential from keras.layers import Embedding, LSTM, Dense, Activation max vec = np.vectorize(max) embed\_dim = 32 # word embedding dimension nhid = 64 # number of hidden units in the LSTM vocab size = max(max(max\_vec(train\_sent)), max(max\_vec(dev\_sent)), max(max\_vec(test\_sent))) + 1 # size of the vocabulary  $n_{classes} = 5$ model = Sequential() model.add(Embedding(vocab size, embed dim)) model.add(LSTM(nhid, dropout=0.1, recurrent dropout=0.1)) model.add(Dense(n\_classes, activation='sigmoid')) In [955]: # 5 - Define your loss/optimizer/metrics # MODIFY CODE BELOW loss classif = 'categorical crossentropy' # find the right loss for multi-class classific ation optimizer = 'adam' # find the right optimizer metrics classif = ['accuracy'] # Observe how easy (but blackboxed) this is in Keras model.compile(loss=loss\_classif, optimizer=optimizer, metrics=metrics classif) print(model.summary()) Layer (type) Output Shape Param # \_\_\_\_\_\_ embedding 42 (Embedding) (None, None, 32) 8992 1stm 34 (LSTM) (None, 64) 24832 dense 44 (Dense) (None, 5) \_\_\_\_\_\_ Total params: 34,149 Trainable params: 34,149 Non-trainable params: 0 None In [956]: # 6 - Train your model and find the best hyperparameters for your dev set # you will be evaluated on the quality of your predictions on the test set # ADAPT CODE BELOW from keras.utils import np utils bs = 64n = 20y train = np utils.to categorical(train labels) y dev = np utils.to categorical(dev labels) history = model.fit(x\_train, y\_train, batch\_size=bs, epochs=n\_epochs, validation\_data=(x\_dev, y\_dev)) Train on 8544 samples, validate on 1101 samples loss: 1.5713 - val acc: 0.2534 Epoch 2/20 1\_loss: 1.5718 - val\_acc: 0.2652 Epoch 3/20 l loss: 1.5689 - val acc: 0.2507 Epoch 4/20 l loss: 1.5709 - val acc: 0.2534 l loss: 1.5750 - val acc: 0.2389 Epoch 6/20 l loss: 1.5857 - val acc: 0.2352 Epoch 7/20 8544/8544 [= l loss: 1.5929 - val acc: 0.2398 Epoch 8/20 1 loss: 1.5964 - val acc: 0.2452 l loss: 1.6253 - val acc: 0.2416 Epoch 10/20 l loss: 1.6208 - val acc: 0.2434 Epoch 11/20 l loss: 1.6082 - val acc: 0.2452 Epoch 12/20 l loss: 1.6187 - val acc: 0.2480 Epoch 13/20 l loss: 1.6319 - val acc: 0.2425 Epoch 14/20 l loss: 1.6393 - val acc: 0.2407 Epoch 15/20 l loss: 1.6345 - val acc: 0.2371 Epoch 16/20 l\_loss: 1.6501 - val acc: 0.2507 Epoch 17/20 l loss: 1.6485 - val acc: 0.2498 Epoch 18/20 l loss: 1.6583 - val acc: 0.2589 Epoch 19/20 l loss: 1.6720 - val acc: 0.2407 Epoch 20/20 l loss: 1.6653 - val acc: 0.2470 In [961]: # 7 - Generate your predictions on the test set using model.predict(x test) https://keras.io/models/model/ Log your predictions in a file (one line = one integer: 0,1,2,3,4) Attach the output file "logreg 1stm y test sst.txt" to your deliverable. # TYPE CODE HERE path = '/Users/adilrhoulam/Downloads/nlp project-3/logreg lstm y test sst.txt' thefile = open(path, 'w') predictions = model.predict classes(x test) for i in range(len(test\_sent\_vec\_mean)): thefile.write("%s\n" % predictions[i]) 4.3 -- innovate! In [969]: # 8 - Open question: find a model that is better on your dev set # (e.g: use a 1D ConvNet, use a better classifier, pretrain your lookup tables ..) # you will get point if the results on the test set are better: be careful of not overfitt ing your dev set too much.. Attach the output file "XXX XXX y test sst.txt" to your deliverable. # TYPE CODE HERE from keras.models import Sequential from keras.layers import Embedding, LSTM, Dense, Activation, Conv1D, Dropout, GlobalMaxPooling max vec = np.vectorize(max) embed dim = 32 # word embedding dimension = 64 # number of hidden units in the LSTM vocab size = max(max(max vec(train sent)), max(max vec(dev sent)), max(max vec(test sent))) + 1 # size of the vocabulary n classes = 5model = Sequential() model.add(Embedding(vocab\_size,embed\_dim)) model.add(Dropout(0.2)) model.add(Conv1D(64,3,padding='valid',activation='relu',strides=1)) model.add(GlobalMaxPooling1D()) model.add(Dense(64)) model.add(Dropout(0.3)) model.add(Activation('relu')) model.add(Dense(n\_classes)) model.add(Activation('softmax')) print(model.summary()) model.compile(loss='categorical\_crossentropy', optimizer='adam', metrics=['accuracy']) Param # Layer (type) Output Shape \_\_\_\_\_\_ embedding\_46 (Embedding) 8992 (None, None, 32) dropout\_17 (Dropout) (None, None, 32) conv1d 9 (Conv1D) (None, None, 64) 6208 global\_max\_pooling1d\_8 (Glob (None, 64) 4160 dense\_51 (Dense) (None, 64) dropout\_18 (Dropout) (None, 64) activation\_17 (Activation) (None, 64) dense\_52 (Dense) 325 (None, 5) activation\_18 (Activation) (None, 5) \_\_\_\_\_\_ Total params: 19,685 Trainable params: 19,685 Non-trainable params: 0 None In [970]: from keras.utils import np\_utils bs = 64n = pochs = 15y\_train = np\_utils.to\_categorical(train\_labels) y\_dev = np\_utils.to\_categorical(dev\_labels) history = model.fit(x\_train, y\_train, batch\_size=bs, epochs=n\_epochs, validation\_data=(x\_dev, y\_dev)) Train on 8544 samples, validate on 1101 samples 1\_loss: 1.5730 - val\_acc: 0.2534 Epoch 2/15 1\_loss: 1.5736 - val\_acc: 0.2534

Epoch 3/15

Epoch 4/15

Epoch 5/15

Epoch 7/15

Epoch 8/15

Epoch 10/15

Epoch 11/15

Epoch 12/15

l loss: 1.5712 - val acc: 0.2607

1 loss: 1.5726 - val\_acc: 0.2625

l\_loss: 1.5723 - val\_acc: 0.2598

l loss: 1.5741 - val acc: 0.2543

l loss: 1.5868 - val acc: 0.2598

1\_loss: 1.5848 - val\_acc: 0.2589

1\_loss: 1.6009 - val\_acc: 0.2470

l\_loss: 1.6043 - val\_acc: 0.2489

l loss: 1.6236 - val acc: 0.2461

l\_loss: 1.6395 - val\_acc: 0.2480

**Deep Learning for NLP - Project** 

• 4 / 20 points will be allocated to the clarity of your code

DO NOT INCLUDE THE DATASETS IN THE DELIVERABLE..

In [2]: PATH TO DATA = "/Users/adilrhoulam/Downloads/data-2"

def \_\_init\_\_(self, fname, nmax=100000):
 self.load wordvec(fname, nmax)

def load wordvec(self, fname, nmax):

self.word2vec = {}

next(f)

All cells should be runnable (modulo trivial compatibility bugs that we'd fix)

1) Monolingual (English) word embeddings

self.word2id = dict.fromkeys(self.word2vec.keys())
self.id2word = {v: k for k, v in self.word2id.items()}
self.embeddings = np.array(self.word2vec.values())

with io.open(fname, encoding='utf-8') as f:

word, vec = line.split(' ', 1)

self.word2vec[word] = np.fromstring(vec, sep=' ')

for i, line in enumerate(f):

**if** i == (nmax - 1):

· Do not create any additional cell

· Efficient code will have a bonus

· the predictions of the SST test set

Fill in the blanks

**DELIVERABLE**:

In [1]: import io

import os

In [3]: class Word2vec():

import scipy

import numpy as np

· this notebook

RULES: