# 自然语言处理（4）——手推word2vec（下）

dreamer　build ai dream　2019-08-13

上次写到了softmax的实现，然后看一下它的梯度求解：

然后根据上图再继续推导关于词向量的梯度.在以softmax为假设函数的word2vec中：

Assume you are given a predicted word vector v c corresponding to the center word c for
skipgram, and word prediction is made with the softmax function found in word2vec models

$$\hat{y}_o = p(o \mid c) = \frac{\exp(u_o^\top v_c)}{\sum_{w=1}^{W} \exp(u_w^\top v_c)}$$

关于中心词与output 词向量的梯度推导为：



注意保证矩阵的维数正确即可 :这部分代码如下。predicted 表示中心词，taget表示需要预测的那个词的编号，outputvec表示输出的那个d*V大小的矩阵 。

```
1  def softmaxCostAndGradient(predicted, target, outputVectors, dataset):
2      """
3      Softmax cost function for word2vec models
4      Implement the cost and gradients for one predicted word vector
5      and one target word vector as a building block for word2vec
6      models, assuming the softmax prediction function and cross
7      entropy loss.
8      Arguments:
```

```
 9        predicted -- numpy ndarray, predicted word vector (\hat{v} in
10                    the written component)
11        target -- integer, the index of the target word
12        outputVectors -- "output" vectors (as rows) for all tokens
13        dataset -- needed for negative sampling, unused here.
14
15        Return:
16        cost -- cross entropy cost for the softmax word prediction
17        gradPred -- the gradient with respect to the predicted word
18              vector
19        grad -- the gradient with respect to all the other word
20              vectors
21
22        We will not provide starter code for this function, but feel
23        free to reference the code you previously wrote for this
24        assignment!
25        """
26
27        Vc=predicted.reshape(-1,1) # d*1
28        z=np.dot( outputVectors,Vc)   #V*d . d*1=>v*1
29        yhat=softmax(z).reshape(-1,1) # make sure the shape is  a column    v*1
30        cost=-np.log(yhat[target])
31        yhat[target]-=1
32        gradPred=np.dot(outputVectors.T, yhat) # d*v .v*1=> d*1
33        grad=np.dot(yhat,Vc.T)    # v*d
34 ### END YOUR CODE
35        return cost, gradPred, grad
```

```
1
2 def skipgram(currentWord, C, contextWords, tokens, inputVectors, outputVector
3 dataset, word2vecCostAndGradient=softmaxCostAndGradient):
4
5     """ Skip-gram model in word2vec
6     Implement the skip-gram model in this function
7 Arguments:
```

```
8   currentWord -- a string of the current center word
9   C -- integer, context size
10  contextWords -- list of no more than 2*C strings, the context words
11  tokens -- a dictionary that maps words to their indices in
12  the word vector list
13  inputVectors -- "input" word vectors (as rows) for all tokens
14  outputVectors -- "output" word vectors (as rows) for all tokens
15  word2vecCostAndGradient -- the cost and gradient function for
16  a prediction vector given the target
17  word vectors, could be one of the two
18                              cost functions you implemented above.
19  Return:
20  cost -- the cost function value for the skip-gram model
21  grad -- the gradient with respect to the word vectors
22      """
23      cost = 0.0
24      gradIn = np.zeros(inputVectors.shape)  # w
25      gradOut = np.zeros(outputVectors.shape) #w'
26      center= tokens[currentWord]
27      predicted=inputVectors[center]
28      for word in contextWords:
29          target=tokens[word]
30          cost_sub, gradPred, grad=word2vecCostAndGradient(predicted, target, 
31          cost+=cost_sub
32          gradIn[ center] +=np.squeeze(gradPred )
33          gradOut+=grad
34      return cost, gradIn, gradOut
35   print( skipgram("c", 3, ["a", "b", "e", "d", "b", "c"],
36          dummy_tokens, dummy_vectors[:5,:], dummy_vectors[5:,:], dataset)  )
37
```

可以看到对于每一个中心，会将其需要预测的背景词遍历一遍，将cost和梯度都进行加一遍。

```
1  def test_word2vec():
2  """ Interface to the dataset for negative sampling """
3      dataset = type('dummy', (), {})()
4  def dummySampleTokenIdx():
5  return random.randint(0, 4)
```

```
6
7    def getRandomContext(C):
8            tokens = ["a", "b", "c", "d", "e"]
9    return tokens[random.randint(0,4)], \
10              [tokens[random.randint(0,4)] for i in range(2*C)]
11       dataset.sampleTokenIdx = dummySampleTokenIdx
12       dataset.getRandomContext = getRandomContext
13       random.seed(31415)
14       np.random.seed(9265)
15       dummy_vectors = normalizeRows(np.random.randn(10,3)) # d=3 , n=10/2=5
16       dummy_tokens = dict([("a",0), ("b",1), ("c",2),("d",3),("e",4)])
17    print ("==== Gradient check for skip-gram ====")
18       gradcheck_naive(lambda vec: word2vec_sgd_wrapper(
19           skipgram, dummy_tokens, vec, dataset, 5, softmaxCostAndGradient),
20           dummy_vectors)
21       gradcheck_naive(lambda vec: word2vec_sgd_wrapper(
22           skipgram, dummy_tokens, vec, dataset, 5, negSamplingCostAndGradient),
23           dummy_vectors)
24    print ("\n==== Gradient check for CBOW      ====")
25       gradcheck_naive(lambda vec: word2vec_sgd_wrapper(
26           cbow, dummy_tokens, vec, dataset, 5, softmaxCostAndGradient),
27           dummy_vectors) #梯度检查
28       gradcheck_naive(lambda vec: word2vec_sgd_wrapper(
29           cbow, dummy_tokens, vec, dataset, 5, negSamplingCostAndGradient),
30           dummy_vectors)
31
32       print( "\n=== Results ===")
33       print( skipgram("c", 3, ["a", "b", "e", "d", "b", "c"],
34           dummy_tokens, dummy_vectors[:5,:], dummy_vectors[5:,:], dataset)  )
35    print (skipgram("c", 1, ["a", "b"],
36           dummy_tokens, dummy_vectors[:5,:], dummy_vectors[5:,:], dataset,
37           negSamplingCostAndGradient))
38    print (cbow("a", 2, ["a", "b", "c", "a"],
39           dummy_tokens, dummy_vectors[:5,:], dummy_vectors[5:,:], dataset))
40    print (cbow("a", 2, ["a", "b", "a", "c"],
41           dummy_tokens, dummy_vectors[:5,:], dummy_vectors[5:,:], dataset,
42           negSamplingCostAndGradient))
```

1

```
2   #为了梯度检查，返回函数值和其根据求导公式推出来的数值
3   def word2vec_sgd_wrapper(word2vecModel, tokens, wordVectors, dataset, C,
4   word2vecCostAndGradient=softmaxCostAndGradient):
5   batchsize = 50
6   cost = 0.0
7   grad = np.zeros(wordVectors.shape)
8   N =int( wordVectors.shape[0])
9       # print(N/2)
10  inputVectors = wordVectors[:int(N/2),:]
11  outputVectors = wordVectors[int(N/2):,:]
12  for i in range(batchsize):
13  C1 = random.randint(1,C)
14  centerword, context = dataset.getRandomContext(C1)
15
16  if word2vecModel == skipgram:
17  denom = 1
18  else:
19  denom = 1
20
21  c, gin, gout = word2vecModel(
22  centerword, C1, context, tokens, inputVectors, outputVectors,
23  dataset, word2vecCostAndGradient)
24  cost += c / batchsize / denom
25  grad[:int(N/2), :] += gin / batchsize / denom
26  grad[int(N/2):, :] += gout / batchsize / denom
27
28  return cost, grad
```
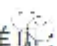
之前讲的都是正常的采样方法，这样的方法需要遍历整个语料库中的每个词然后预测每个词的上下文。

可以根据其求导公式看，梯度下降时运算量很大，要用整个d*W矩阵，当W很大时，运算量就很大。所以提出负采样方法：

目标函数：

$$J_t(\theta) = \log \sigma\left(u_o^T v_c\right) + \sum_{i=1}^{k} \mathbb{E}_{j \sim P(w)}\left[\log \sigma\left(-u_j^T v_c\right)\right]$$

这里$t$是某个窗口，$k$是采样个数，$P(w)$是一个unigram分布，详见：build ai dream

就是故意取K个不是预测目标的词语，negative sampling实现skip-gram。这是一种采样子集简化运算的方法。具体做法是，对每个正例（中央词语及上下文中的一个词语）采样几个负例（中央词语和其他随机词语），训练binary logistic regression（也就是二分类器）。



```
1   def getNegativeSamples(target, dataset, K):  #负采样函数，取得负样本的编号
2       """ Samples K indexes which are not the target """
3
4       indices = [None] * K
5       for k in range(K):
6           newidx = dataset.sampleTokenIdx()
7           while newidx == target:
8               newidx = dataset.sampleTokenIdx()
9           indices[k] = newidx
10      return indices
11
12
13  def negSamplingCostAndGradient(predicted, target, outputVectors, dataset,
14                                 K=10):
15      """ Negative sampling cost function for word2vec models
16
17      Implement the cost and gradients for one predicted word vector
18      and one target word vector as a building block for word2vec
19      models, using the negative sampling technique. K is the sample
20      size.
21
22      Note: See test_word2vec below for dataset's initialization.
```

```
23
24      Arguments/Return Specifications: same as softmaxCostAndGradient
25      """
26
27      # Sampling of indices is done for you. Do not modify this if you
28      # wish to match the autograder and receive points!
29      indices = [target]
30      indices.extend(getNegativeSamples(target, dataset, K))
31
32      ### YOUR CODE HERE
33      grad = np.zeros(outputVectors.shape)
34      gradPred = np.zeros(predicted.shape)
35      cost = 0
36      z = sigmoid(np.dot(outputVectors[target], predicted))
37
38      cost -= np.log(z)
39      grad[target] += predicted * (z - 1.0)
40      gradPred += outputVectors[target] * (z - 1.0)
41
42       for k in  indices[1:] :
43          z = sigmoid(np.dot(outputVectors[k], predicted))
44          cost -= np.log(1.0 - z)
45          grad[k] += predicted * z
46          gradPred += outputVectors[samp] * z
47      ### END YOUR CODE
48      return cost, gradPred, grad
```

上面讲完了skip-gram 以及 两种训练方法的实现。下面讲一下cbow的实现：

INPUT    PROJECTION    OUTPUT          INPUT    PROJECTION    OUTPUT

w(t-2)
w(t-1)        SUM        w(t)
w(t+1)
w(t+2)

w(t)        w(t-2) w(t-1) w(t+1) w(t+2)

**CBOW**                          **Skip-gram**

其实训练过程和梯度更新都一样，就是predicted 和target 有所不同：

```
1   def cbow(currentWord, C, contextWords, tokens, inputVectors, outputVectors,
2           dataset, word2vecCostAndGradient=softmaxCostAndGradient):
3       """CBOW model in word2vec
4
5       Implement the continuous bag-of-words model in this function.
6
7       Arguments/Return specifications: same as the skip-gram model
8
9       Extra credit: Implementing CBOW is optional, but the gradient
10      derivations are not. If you decide not to implement CBOW, remove
11      the NotImplementedError.
12      """
13
14      cost = 0.0
15      gradIn = np.zeros(inputVectors.shape)
16      gradOut = np.zeros(outputVectors.shape)
17
18      context_id=[tokens[context] for context in contextWords ]
19      predicted_vec=inputVectors[context_id]
20      predicted=np.sum(predicted_vec,axis=0)    # 求和函数
21      target= tokens[currentWord]
22      cost_sub, gradPred, grad=word2vecCostAndGradient(predicted, target, outpu
```

```
23          cost+=cost_sub
24          gradOut+=grad
25          # 注意这儿：
26          for i in  context_id:
27              gradIn[i]+=np.squeeze(gradPred)
28      return cost, gradIn, gradOut
29
```

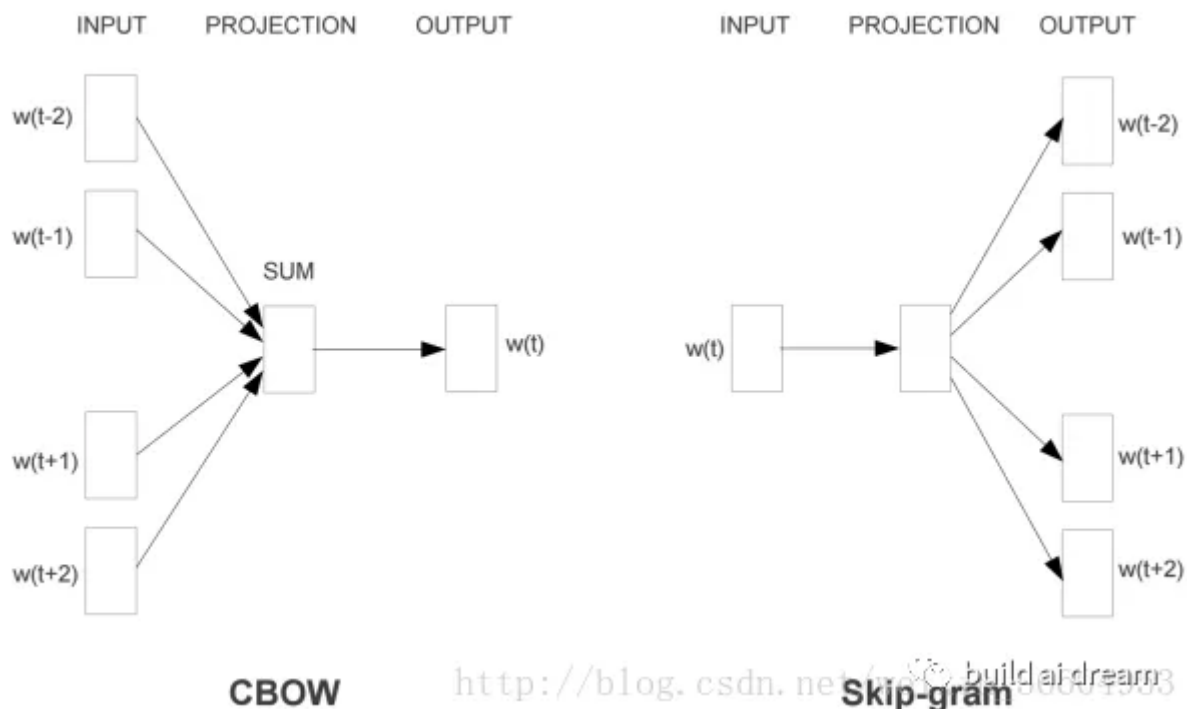　　这里上下文向量是所有上下文词向量的和，却把梯度应用到了所有的上下文词向量上去。虽然看上去挺疯狂，但就是管用。

至此：W2V的两种方法我们都实现了：一般来说一个词语有两个向量，可以将其拼接或者相加。相加用得更多

下面为常见的面试题目：
　q :word2vec有哪两种结构？有什么区别

　ans:
　　word2vec包含两种结构，一种是skip-gram结构，一种是cbow结构，skip-gram结构是利用中间词预测邻近词，cbow模型是利用上下文词预测中间词，网络结构如下图：



　　区别:CBOW - continous bag-of-words， 是根据相邻单词来计算当前单词的词向量，因而不考虑window 内的单词的顺序；skip gram 则相反，是根据当前单词来计算

其邻近单词的词向量，因而考虑window 内的单词语序（距离当前词比较近的单词被赋予的权重就大）；CBOW 训练速度比 skip gram 快，但skip gram 对于低频词(多个背景词对其训练)的效果比较好.

q：　word2vec中有没有非线性层
ans：

　word2vec两种模型有三层，输入层，映射层，输出层，隐藏层是线性结构，这也是为什么相对于语言模型训练较快的原因。

q: CBOW模型输入是什么，输出是什么，初始化出入是怎么做的?
ans：

　CBOW模型输入是上下文的词向量，这个词向量是随机初始化的产生的，将这个上下文词向量加起来得到映射层，输出层是层次softmax(hierarchical Softmax)，最终输出中心词的最大概率。

q: 层次softmax（hierarchical Softmax）在word2vec中是怎么应用的，为什么要用，用了有哪些好处?

ans：

　面试官问此问题的目的：主要是考察应聘者对word2vec优化的理解及对word2vec更深层次的理解。在回答这个问题的时候应聘者应先阐述普通softmax会产生什么样的问题，然后再阐述使用层次softmax后能够达到什么样的效果。
　　相应答案（仅供参考）
　假设我们采用普通的softmax，那么在输出层需要得到一个词典大小的概率分布，假设公式我们采用softmax（Wx+b）（这里一定要在纸或者白板上写出），那么W将会有一个维度是10000维，这是一个很大的权重矩阵，那么我们整个优化目标是10000个，这个工作量的巨大的，会严重影响计算效率。如果采用层次softmax，那么我们的输出层将是一个霍夫曼二叉树，这个二叉树的最大深度是log(D)-1,每个节点处进行一个二分类，二分类的权重矩阵W2有一个维度是2，相对于词典大小10000来说就是很小的数。采用层次softmax好处是将一个多分类问题转换成一个一个的二分类问题，降低的计算复杂度，提高了运算效率。

q:在word2vec优化中采用层次softmax的哈夫曼树如何构建，叶子节点代表什么?

ans:

面试官问此问题的目的：主要是考察应聘者对word2vec优化底层原理的理解。

相应答案（仅供参考）

在word2vec中构建哈夫曼树中，我们以词表中的全部词作为叶子节点，词频作为节点的权，构建Huffman树，叶子节点为词表中的全部词，权重是词频。

（层次化索引在后续会补充上来）

refer: 微信公众号：AI壹号堂

https://blog.csdn.net/weixin_40547993/article/details/88779922