

Word2Vec

我看了一些Word2Vec的一些相关文章，主要分为两类。一是有关于Word2Vec的发展，主要是以下4篇文章起到奠基性的作用

1. A Neural Probabilistic Language Model.2003 (NNLM)
2. Recurrent neural network based language model.2010 (RNNLM)
3. Distributed Representations of Words and Phrases and their Compositionality.2013 (Skip Gram Model and CBOW)
4. Efficient Estimation of Word Representations in Vector Space.2013 (Skip Gram Model and CBOW)

二是关于Word2Vec应用的一些文章，一些有趣的idea像是如下的文章

1. Linguistic Regularities in Continuous Space Word Representations.2013
2. Exploiting Similarities among Languages for Machine Translation.2013
3. Distributed Representations of Sentences and Documents.2014

Development

NNLM: A Neural Probabilistic Language Model.2003

这篇文章主要解决在此之前的自然语言模型是统计语言模型和基于统计语言模型n-gram模型的维度灾难问题。

statistical language model:

$$\hat{P}(w_1^T) = \prod_{t=1}^T \hat{P}(w_t | w_1^{t-1}),$$

where w_t is the t -th word, and writing sub-sequence $w_i^j = (w_i, w_{i+1}, \dots, w_{j-1}, w_j)$.

N-gram model:

$$\hat{P}(w_t | w_1^{t-1}) \approx \hat{P}(w_t | w_{t-n+1}^{t-1}).$$

统计语言模型的基本想法就是对于一句话，在给定的前几个词的情况下，统计出现下一个词的概率。这样一句话的出现概率就是第一个词出现的概率 $P(W_1)$ 乘上在第一个词给定的情况下出现第二个词的概率 $P(W_2|W_1)$ ，依此类推，一句话的概率就是上图第一行的联合条件概率乘积。

N-gram模型就是假设一个词出现的概率只考察前后该词前后n个词，以此来降低复杂度。

这些模型的问题就是复杂度非常高，例如：

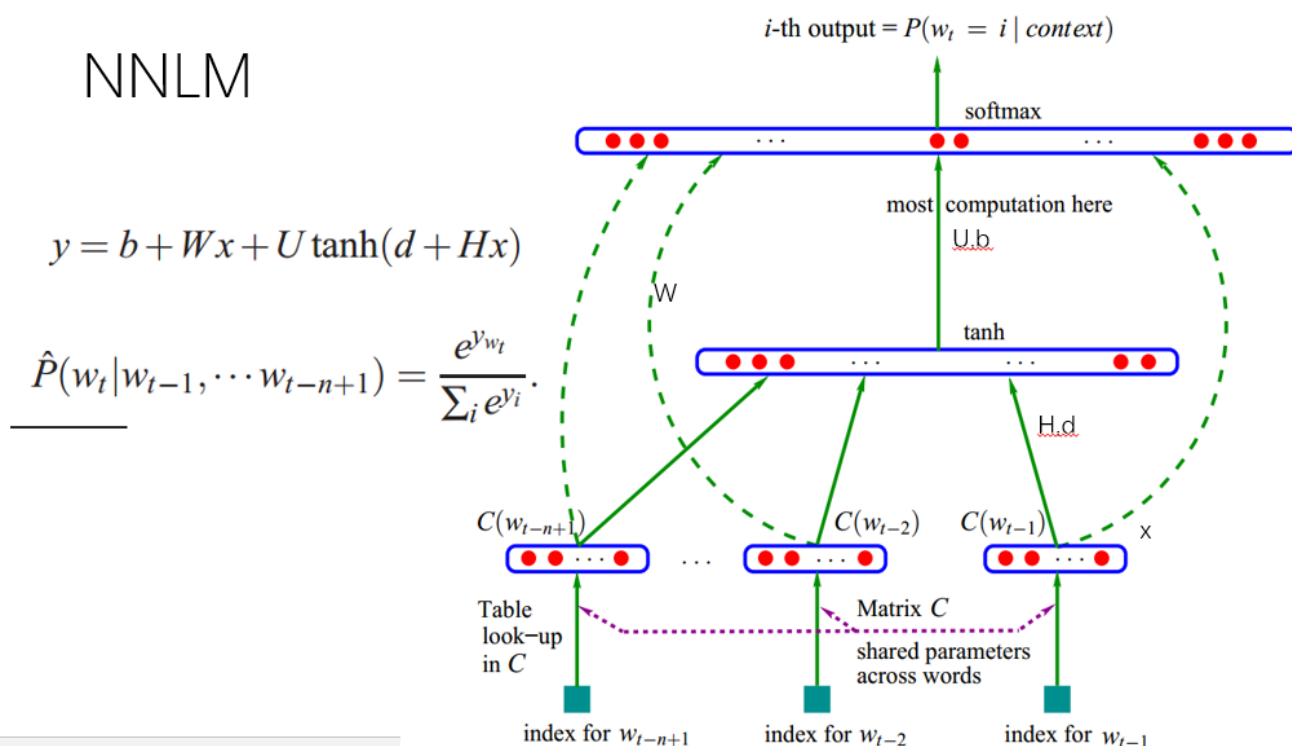
curse of dimensionality. It is particularly obvious in the case when one wants to model the joint distribution between many discrete random variables (such as words in a sentence, or discrete attributes in a data-mining task). For example, if one wants to model the joint distribution of 10 consecutive words in a natural language with a vocabulary V of size 100,000, there are potentially $100\,000^{10} - 1 = 10^{50} - 1$ free parameters. When modeling continuous variables, we obtain gen-

上图的**free parameters**就是指前面语言模型的各个概率 P . 统计语言模型就是要统计所有文本将所有概率 P 确定下来。

文章的作者要解决这一问题，采用**distributed vectors**来表示每一个词，一句话的上下文语境用训练好的**vector**的值和网络参数来表达。如下：

1. associate with each word in the vocabulary a distributed *word feature vector* (a real-valued vector in \mathbb{R}^m),
2. express the joint *probability function* of word sequences in terms of the feature vectors of these words in the sequence, and
3. learn simultaneously the *word feature vectors* and the parameters of that *probability function*.

作者给出的网络结构如下：



如图，首先有一个全局的矩阵 C ,通过 C 将一个词 w 转换为向量 $C(w)$ 。然后用这些向量作为一个三层神经网络的输入，最终训练出参数和所有词向量的矩阵 C 。

作者首次使用神经网络来表示自然语言模型，大大降低了计算复杂度（复杂度在后面统一比较）。作者实际上已经使用了词向量但是没有发现这些向量蕴含的语义和语法内容，只是用它来作为自然语言模型的工具。

RNNLM: Recurrent neural network based language model.2010

NNLM有一个问题就是对于一个词，训练的时候它的上下文只取该词前面的n个词，如果n过大，训练复杂度增大；n过小，覆盖的上下文不够。本文作者就利用RNN的特性来解决这一问题。

RNN:

$$x(t) = w(t) + s(t-1)$$

$$s_j(t) = f\left(\sum_i x_i(t)u_{ji}\right)$$

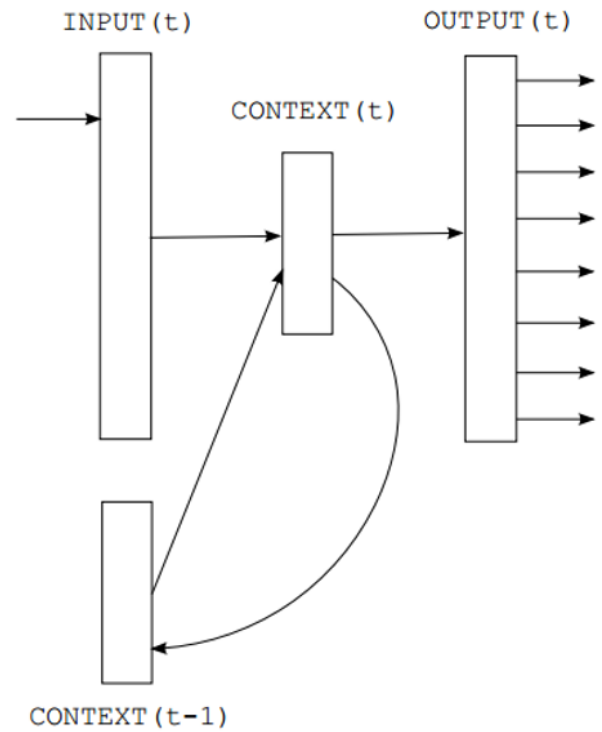
$$y_k(t) = g\left(\sum_j s_j(t)v_{kj}\right)$$

where $f(z)$ is sigmoid activation function:

$$f(z) = \frac{1}{1 + e^{-z}}$$

and $g(z)$ is softmax function:

$$g(z_m) = \frac{e^{z_m}}{\sum_k e^{z_k}}$$



如上图，RNNLM中将隐藏层视为context层，每一次迭代中，输入由新的输入和前一次的context层共同组成。这样上下文的宽度n不用很大，也可以利用到距离当前词足够远的词带来的影响。

但是本文作者同样是没有注意到词向量蕴含的语义和语法内容

Skip Gram Model and CBOW

“Distributed Representations of Words and Phrases and their Compositionality”和“Efficient Estimation of Word Representations in Vector Space”两篇文章提出了Skip Gram Model and CBOW两个模型。

CBOW模型：

CBOW Based on Hierarchical Softmax

Input Layer

Projection Layer

$$p(d_j^w | \mathbf{x}_w, \theta_{j-1}^w) = [\sigma(\mathbf{x}_w^\top \theta_{j-1}^w)]^{1-d_j^w} \cdot [1 - \sigma(\mathbf{x}_w^\top \theta_{j-1}^w)]^{d_j^w}.$$

$$\mathcal{L} = \sum_{w \in \mathcal{C}} \log \prod_{j=2}^{l^w} \{ [\sigma(\mathbf{x}_w^\top \theta_{j-1}^w)]^{1-d_j^w} \cdot [1 - \sigma(\mathbf{x}_w^\top \theta_{j-1}^w)]^{d_j^w} \}$$

$$= \sum_{w \in \mathcal{C}} \sum_{j=2}^{l^w} \{ (1 - d_j^w) \cdot \log[\sigma(\mathbf{x}_w^\top \theta_{j-1}^w)] + d_j^w \cdot \log[1 - \sigma(\mathbf{x}_w^\top \theta_{j-1}^w)] \},$$

$$\frac{\partial \mathcal{L}(w, j)}{\partial \theta_{j-1}^w} = \frac{\partial}{\partial \theta_{j-1}^w} \{ (1 - d_j^w) \cdot \log[\sigma(\mathbf{x}_w^\top \theta_{j-1}^w)] + d_j^w \cdot \log[1 - \sigma(\mathbf{x}_w^\top \theta_{j-1}^w)] \}$$

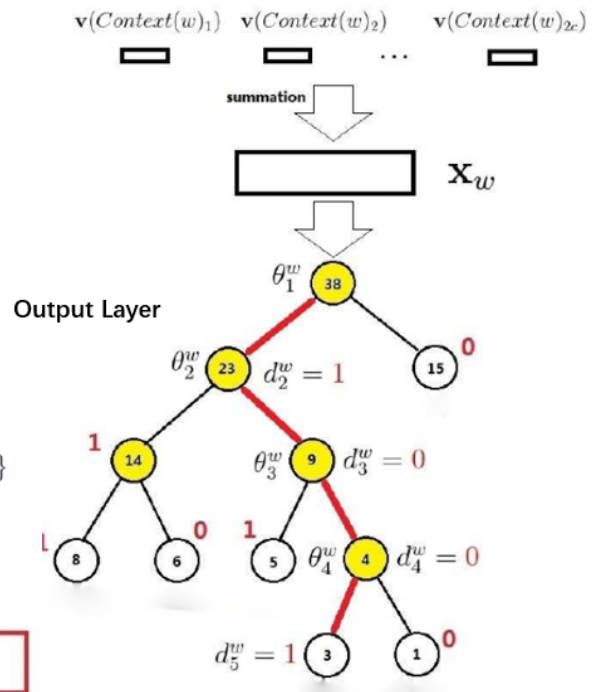
$$= (1 - d_j^w)[1 - \sigma(\mathbf{x}_w^\top \theta_{j-1}^w)] \mathbf{x}_w - d_j^w \sigma(\mathbf{x}_w^\top \theta_{j-1}^w) \mathbf{x}_w$$

$$= \{ (1 - d_j^w)[1 - \sigma(\mathbf{x}_w^\top \theta_{j-1}^w)] - d_j^w \sigma(\mathbf{x}_w^\top \theta_{j-1}^w) \} \mathbf{x}_w$$

$$= [1 - d_j^w - \sigma(\mathbf{x}_w^\top \theta_{j-1}^w)] \mathbf{x}_w.$$

$$\frac{\partial \mathcal{L}(w, j)}{\partial \mathbf{x}_w} = [1 - d_j^w - \sigma(\mathbf{x}_w^\top \theta_{j-1}^w)] \theta_{j-1}^w.$$

Sample: $(Context(w), w)$



CBOW模型相比于前两个模型，去掉了隐藏层。输出不是线性的softmax，而是基于huffman树的softmax，这样首先使得输出层复杂度由N降到log(N)，其次很好的应用了huffman树的聚类效果。

Skip Gram Model与CBOW类似，只是CBOW是由当前词的上下文来推导出当前词汇，而Skip Gram Model则是反过来。

这两个模型首先是降低了计算复杂度，如下图比较：

Complexity

- N: context size
- D: distributed vector size
- H: hidden layer size
- V: output layer size

NNLM: $Q = N \times D + N \times D \times H + H \times V,$

RNNLM: $Q = H \times H + H \times V,$

CBOW Model $Q = N \times D + D \times \log_2(V).$

Skip Gram Model: $Q = C \times (D + D \times \log_2(V)).$

按照文章的说法

This makes the training extremely efficient: an optimized single-machine implementation can train on more than 100 billion words in one day

此外最重要的是，作者首次察觉到词向量蕴含了丰富的语法和语义的含义。比如

$\text{vec}(\text{"Madrid"}) - \text{vec}(\text{"Spain"}) + \text{vec}(\text{"France"})$ is closer to $\text{vec}(\text{"Paris"})$.

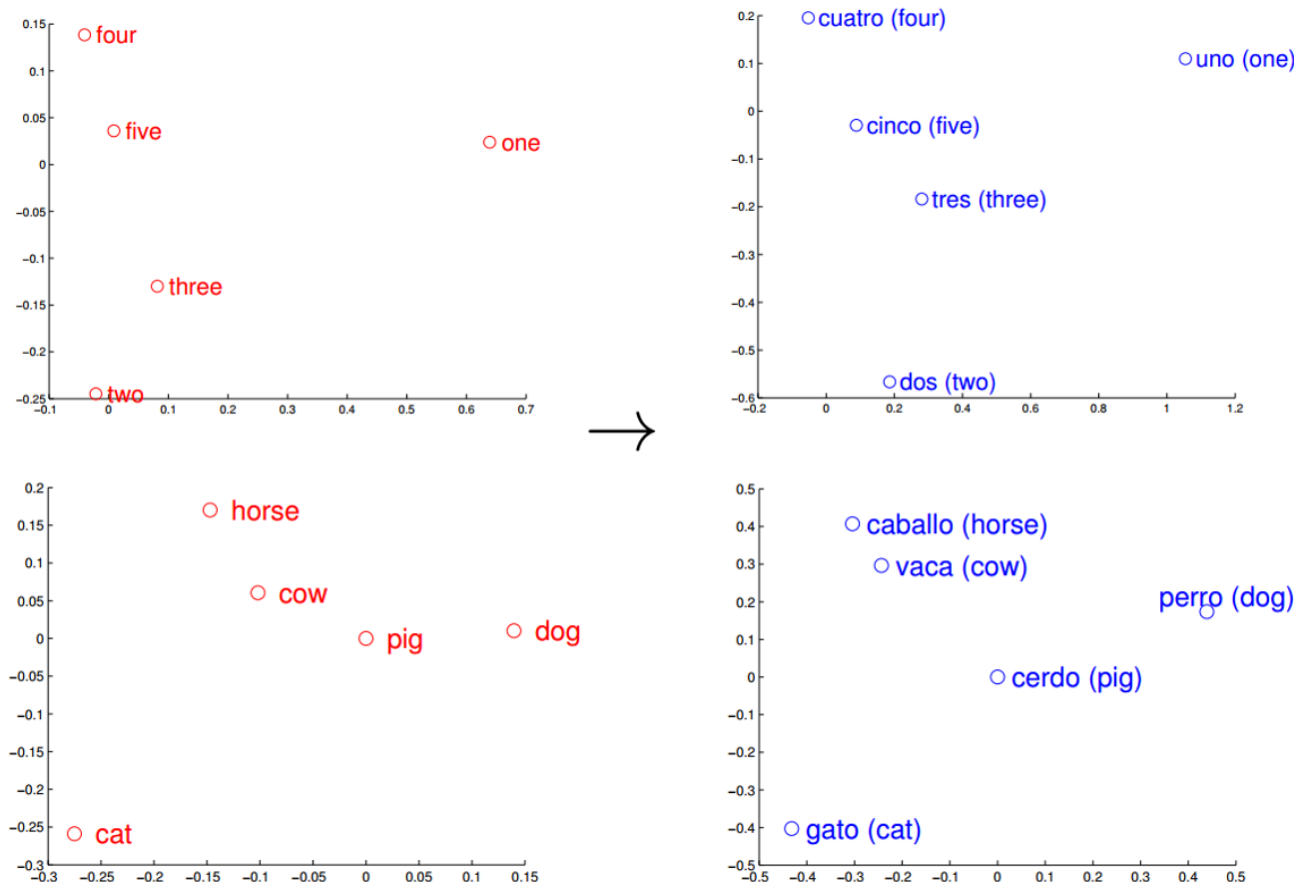
文章中用了大量篇幅和实现表现词向量的良好特性

Application

然后还有些应用词向量的文章。

Exploiting Similarities among Languages for Machine Translation

文章的想法非常简单。在一种语言中一个词的词向量为 \mathbf{X}_i ，另一种语言中其对应的词的词向量为 \mathbf{Z}_i ，作者做了一些实验发现这写词在各自的语言空间中的相对位置关系有一定的相似性，如下图：



所以作者认为，只要训练一个转移矩阵 \mathbf{W} ，使得 \mathbf{X}_i 经过变换 $\mathbf{W}\mathbf{X}_i$ 转换到与 \mathbf{Z}_i 接近的位置，再转换成对应的词，这样就完成了翻译，训练目标如下图：

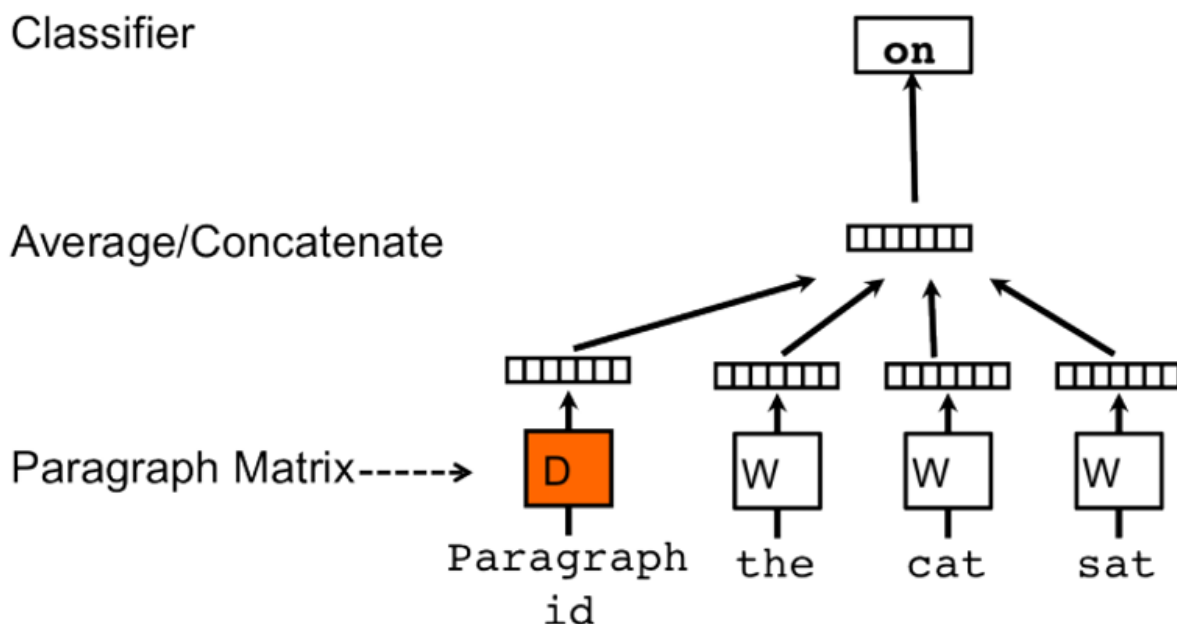
$$\min_W \sum_{i=1}^n \|Wx_i - z_i\|^2$$

作者想法很新颖，但是模型结构非常简单，训练的过程中只考虑到词，没有融入整句话的语义概念。翻译的效果比不上之前的翻译模型，只是开阔了一种新的思路。

Distributed Representations of Sentences and Documents

这一篇文章进一步扩展词向量，提出一种方法来训练句子和文章的向量。

Distributed Memory Model of Paragraph Vectors (PV-DM)



如图，作者的想法也非常简单，在原有的CBOW模型上，作者在输入层额外添加了一个paragraph vector。这个向量是由段落决定的，文章的每一段对应一个向量，在训练过程中，同一段落内部使用同一个向量。这一个向量就代表着该段落的context and semantic.

这是一个非常简单但是巧妙的想法，效果也相当不错。作者做了一个实验来来评估模型，实验是分析一段评论是 positive or negative，数据集是Stanford Sentiment Treebank Dataset。实验结果如下，错误率要低于之前的模型。

Model	Error rate (Positive/ Negative)	Error rate (Fine- grained)
Naïve Bayes (Socher et al., 2013b)	18.2 %	59.0%
SVMs (Socher et al., 2013b)	20.6%	59.3%
Bigram Naïve Bayes (Socher et al., 2013b)	16.9%	58.1%
Word Vector Averaging (Socher et al., 2013b)	19.9%	67.3%
Recursive Neural Network (Socher et al., 2013b)	17.6%	56.8%
Matrix Vector-RNN (Socher et al., 2013b)	17.1%	55.6%
Recursive Neural Tensor Network (Socher et al., 2013b)	14.6%	54.3%
Paragraph Vector	12.2%	51.3%

这一篇训练sentence vector还可以应用到前一篇翻译的工作中，可以弥补其没有融入整句话的语义概念的不足之处。

附：CBOW核心代码

这是目前工作中用到的一段CBOW核心代码，参数还没有调好。不过实现过程中，对于论文中写的比较概括模糊的地方有了比较清晰的认识。

```

1  # -*- coding: utf-8 -*-
2  """
3  Created on Tue Jan 10 17:07:40 2017
4
5  @author: duocai
6  """
7
8  import numpy as np
9  import json
10 import math
11 import datetime
12 import random
13 import threading
14
15 #控制内积范围，超过范围的再加sigmoid很小，可以丢弃
16 MAX_EXP = 8
17
18 # 向量维度
19 layer_size = 100
20 # 单元location or word..个数
21 unit_num = 182968
22 # 单元向量表
23 vec_table = np.random.uniform(-1,1,layer_size*unit_num).reshape(unit_num,layer_size)
24 # 参数表，参数个数即内部节点个数为叶节点个数（unit_num）-1，初始化为正态分布取值
25 para_table = np.random.normal(0.0,1,layer_size*(unit_num-1)).reshape(unit_num-
1,layer_size)
26 # 所有用户，相当于所有文章
27 users = json.load(open('../user_new.json'))
28 # 所有地点，相当于所有单词
29 locats = json.load(open('../location_tree.v3.json'))
30 if len(locats.keys()) > unit_num:
31     print('地点超出范围')
32 # window size
33 window = 5
34 # 学习率
35 start_alpha = 0.03
36 # 开始时间
37 start_time = datetime.datetime.now()
38
39 def sigmoid(x):
40     return 1/(1 + math.exp(-x))
41
42 def show_message():
43     print('vector of location 0:',vec_table[0])
44
45 # 训练index指定的用户
46 def train(index):
47     global layer_size
48     global window
49     global MAX_EXP
50     global start_alpha
51     global start_time
52     sent = users[index]['locations']#文章，用户访问地点序列

```



```

53 sent_len = len(sent)
54 neul = np.zeros(layer_size)#隐藏层
55 neule = np.zeros(layer_size)# 记录隐藏层累积变化量
56 for pos in range(sent_len):
57     loc_id = sent[0] # 地点id string
58     if (loc_id not in locats):#不存在该地点则忽略
59         continue
60     # 计算context向量和
61     num = random.randint(0,window)#随机起点，并不是严格从0开始
62     start = num
63     while start < 2 * window + 1 :
64         cur_pos = pos - window + start
65         #左右几个词不包含当前词
66         # or 现在位置超过范围，则跳过
67         if start == window or cur_pos < 0 or cur_pos >= sent_len:
68             start+=1
69             continue
70         # 计算context向量和
71         neul += vec_table[int(sent[cur_pos])]
72         start+=1
73     #利用霍夫曼树计算
74     loc = locats[loc_id] # 当前地点
75     points = loc['points']
76     codes = loc['codes']
77     codes_len = len(codes)
78     for i in range(codes_len):
79         # 计算内积
80         dot = np.dot(neul,para_table[int(points[i])])
81         ##内积不在范围内直接丢弃
82         if dot > MAX_EXP or dot < -MAX_EXP:
83             continue
84         ##simoid
85         dot = sigmoid(dot)
86         # 偏导数的公用部分*学习率alpha
87         g = (1 - int(codes[i]) - dot)*start_alpha
88
89         # 反向更新参数
90         # 先更新隐藏层
91         neule += g*para_table[int(points[i])]
92         # 后更新参数
93         para_table[int(points[i])] += g*neul
94
95     ## 将更新传递到词向量
96     start = num
97     while start < 2 * window + 1 :
98         cur_pos = pos - window + start
99         #左右几个词不包含当前词
100        # or 现在位置超过范围，则跳过
101        if start == window or cur_pos < 0 or cur_pos >= sent_len:
102            start+=1
103            continue
104        # 更新词向量
105        vec_table[int(sent[cur_pos])] += neule

```

```
106         start+=1
107         # 输出现在的时间
108         if index%200 == 0:
109             print('user: ', index, ',time: ', (datetime.datetime.now() - start_time).seconds,
110                   's')
111             show_message()
112     ## 利用多线程训练
113     threads = []
114     for i in range(len(users)):
115         threads.append(threading.Thread(target=train,args=(i,)))
116     ## run
117     for t in threads:
118         t.setDaemon(True)
119         t.start()
120     t.join()#等待子线程结束
```