

Word2Vec CBOW实现

彭博 519030910366

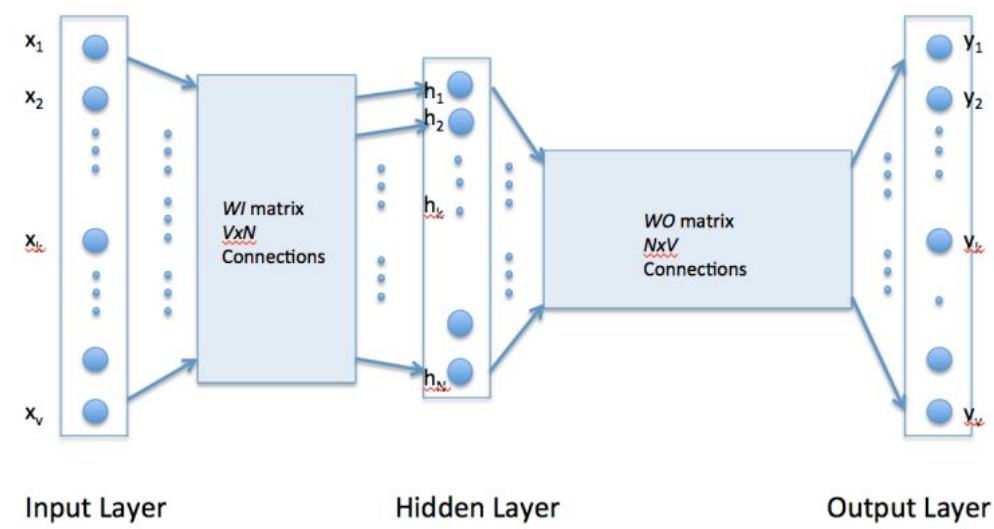
只需补全给定代码中的 `train_one_step` 函数，下面简要介绍Word2Vec和CBOW原理

CBOW

为了更接近自然语言的本质，希望在训练过程中蕴含上下文信息，但是由于引入上下文在过去的模型中容易导致词的向量表示维度过大，为后续运算带来了很大的麻烦，因此引入CBOW，使得输入向量可以满足：

- 1. 携带上下文信息
- 2. 稠密

在介绍CBOW之前，首先简要介绍Word2Vec模型，如下图所示



输入层X的大小V即训练所用到的词表大小，隐藏层神经元数目为N，（训练时设定），输出层Y大小同样为V，模型主要流程为：输入层 $X \times WI$ 矩阵，其中WI矩阵的大小为 $V \times N$ ，得到隐藏层H，隐藏层 $H \times WO$ 得到输出层U，输出U后需要进行softmax层得到最终的预测。

训练网络过程中，计算出softmax后的向量和ground truth的one hot vector的loss，然后BP反向更新权重。

CBOW结构和Word2Vec基本一致，不同之处仅仅是此时的输入层X并不是来源于一个单词而是多个单词的结合。其中的多个单词主要由一个window采样得到，（可以想象成一个窗在整个文段中移动然后每次移动后把窗中的内容放入一个词袋，作为一个target vector的一个上下文组

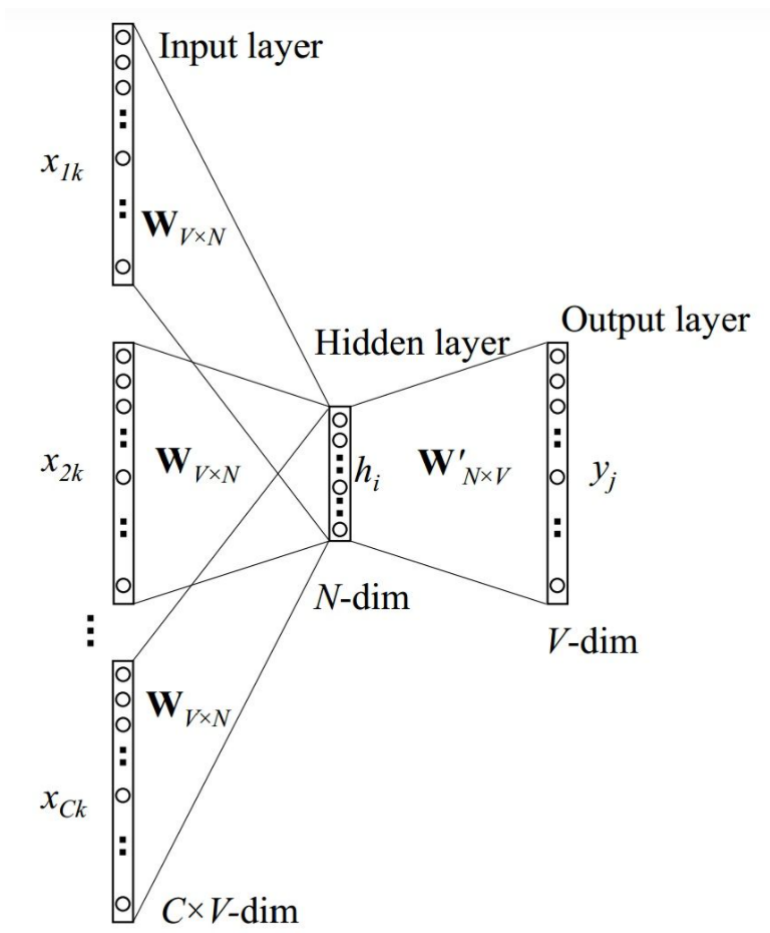


Figure 2: Continuous bag-of-word model

主要公式推导如下：

- 假设一个context_token里面有C个词(token), t_1, t_2, \dots, t_c
- 输入 $X = \sum_{i=1}^c t_i$
- 计算隐藏层 $H = \frac{W_I \cdot X}{C}$
- 输出层 $U = W_O \cdot H$
- 预测概率 $Y = Softmax(U)$
- 损失函数 $E = -\log p(w_O | w_I)$
- 权重更新公式
 - $v'_{w_j}(\text{new}) = v'_{w_j}(\text{old}) - \eta \cdot e_j \cdot \mathbf{h}$ for $j = 1, 2, \dots, V$
 - $v^{(new)}_{w_{I,c}} = v^{(old)}_{w_{I,c}} - \frac{1}{C} \cdot \eta \cdot \mathbf{E} \mathbf{H}^T$ for $c = 1, 2, \dots, C$.

代码实现

主要思路就是按照提示的四个步骤依次构建输入的one-hot向量（是context_tokens的叠加），前向传播计算隐藏层H，计算输出层U ,计算softmax，然后计算损失(loss和EH)，最后更新W1 W2 参数。

```
1 def train_one_step(self, context_tokens: List[str], target_token: str,
2   learning_rate: float) -> float:
3     """
4     Predict the probability of the target token given context tokens.
5
6     :param context_tokens: List of tokens around the target token
7     :param target_token: Target (center) token
8     :param learning_rate: Learning rate of each step
9     :return: loss of the target token
10    """
11    # == TODO: Construct one-hot vectors ==
12    V = len(self.vocab)
13    input_X = np.zeros(V) # (V,)
14    for token in context_tokens:
15        token_index = self.vocab.token_to_idx(token)
16        input_X[token_index] += 1
17    target_index = self.vocab.token_to_idx(target_token)
18    target = one_hot(V, target_index) # (V,)
19
20    # == TODO: Forward step ==
21    H = np.dot(input_X, self.W1) / len(context_tokens) # (N,)
22    U = np.dot(H, self.W2) # (V,)
23    predict_target = softmax(U) # (V,)
24
25    # == TODO: Calculate loss ==
26    loss = - np.log(np.dot(predict_target, target))
27    EH = np.dot(self.W2, np.transpose(predict_target - target)) # (N,)
28
29    # == TODO: Update parameters ==
30    self.W1 -= input_X[:, None] * EH * learning_rate / len(context_tokens)
31    self.W2 -= np.outer(H, predict_target - target) * learning_rate
32
33    return loss
```

输出结果

test1

```
1 Token number: 50
2 Vocab size: 21
3 Epoch 1, loss: 2.96. Cost 0.0 min
4 Epoch 2, loss: 1.99. Cost 0.0 min
5 Epoch 3, loss: 1.46. Cost 0.0 min
6 Epoch 4, loss: 1.16. Cost 0.0 min
7 Epoch 5, loss: 0.94. Cost 0.0 min
8 Epoch 6, loss: 0.82. Cost 0.0 min
9 Epoch 7, loss: 0.74. Cost 0.0 min
10 Epoch 8, loss: 0.70. Cost 0.0 min
11 Epoch 9, loss: 0.82. Cost 0.0 min
12 Epoch 10, loss: 1.09. Cost 0.0 min
13 [('she', 1.0), ('now', 0.9166363811901059), ('i', 0.7916710150878732), ('he',
0.7600563117313468), ('will', 0.7217754180931802)]
```

test2

```
1 Token number: 205068
2 Vocab size: 17832
3 (省略训练部分详细输出)
4 ...
5 Epoch 1, loss: 8.54. Cost 10.7 min
6 ...
7 Epoch 2, loss: 7.71. Cost 10.8 min
8 ...
9 Epoch 3, loss: 7.47. Cost 10.9 min
10 ...
11 Epoch 4, loss: 7.32. Cost 10.9 min
12 ...
13 Epoch 5, loss: 7.22. Cost 11.0 min
14 ...
15 Epoch 6, loss: 7.14. Cost 10.8 min
16 ...
17 Epoch 7, loss: 7.07. Cost 9.7 min
18 ...
19 Epoch 8, loss: 7.02. Cost 9.4 min
20 ...
21 Epoch 9, loss: 6.97. Cost 10.2 min
22 ...
23 Epoch 10, loss: 6.92. Cost 9.7 min
24 ...
25 Save model to ckpt
26 Cost 104.2 min
27 [('i', 0.9999999999999999), ('we', 0.9124169944756282), ('levy',
0.9051181652637497), ('baboon', 0.8826970358235133), ('adrift', 0.8742717194144737),
('you', 0.8730728436218511), ('jim', 0.8700385565478237), ('lounge',
0.8690574718409811), ('they', 0.8339341158975251), ('discord', 0.8330665504068829)]
```

test3

```
1 Load model from ckpt
2 spearman correlation: 0.397
3 pearson correlation: 0.551
```