

干货分享 | 机器学习竞赛必备基础知识—Word2Vec

一碗数据汤 3月25日

以下文章来源于kaggle竞赛宝典，作者尘沙杰少



kaggle竞赛宝典

数据竞赛加分骚操作 & 数据分析方法 & 实践机器学习 & Kaggle + 天池 + 其他

上

转自：kaggle竞赛宝典

1.简介

本文我们主要介绍词嵌入中一种非常经典的算法，Word2Vec,早期Word2Vec主要被用在文本类的问题中,但是现在做比赛的朋友应该都发现了,几乎一半的传统数据竞赛都会用到Word2Vec,所以这边我们必须得好好看看Word2Vec究竟在学习什么,这样今后也可以举一反三更好地使用这些技术。

本文我们先简单介绍词嵌入模型，然后详细介绍Word2Vec,包括Word2Vec在做什么,为什么用Word2Vec(它与老的BOW模型比有什么创新)，Word2Vec常见的两种框架以及选择,最后我们会给出基于Pytorch版本的Word2Vec的两种实现。

本文我们不会罗列太多的数学,主要讲框架以及实现,大家可以**思考将其如何的用到自己的问题中**,在词嵌入模型中,我们也会列举些许相关的应用。

2.什么是词嵌入(Word Embedding)

词嵌入(Word Embedding)是一类词的表示方法,它可以通过很多机器学习模型将原先不同的词转化为不同的实数向量，目前的Word Embedding的技术非常多，例如Google的Word2Vec,Stanford的Glove, Facebook的Fastext等等。

Word Embedding有的时候也被称作为分布式语义模型或向量空间模型等,所以从名字和其转换的方式我们就可以明白, Word Embedding技术可以将相同类型的词归到一起,例如苹果, 芒果香蕉等, 在投影之后的向量空间距离就会更近, 而书本, 房子这些则会与苹果这些词的距离相对较远。

3.什么时候使用词嵌入模型

目前为止, Word Embedding可以用到特征生成, 文件聚类, 文本分类和自然语言处理等任务, 例如:

- 计算相似的词: Word Embedding可以被用来寻找与某个词相近的词。
- 构建一群相关的词: 对不同的词进行聚类, 将相关的词聚集到一起;
- 用于文本分类的特征: 在文本分类问题中, 因为词没法直接用于机器学习模型的训练, 所以我们将词先投影到向量空间,这样之后便可以基于这些向量进行机器学习模型的训练;
- 用于文件的聚类

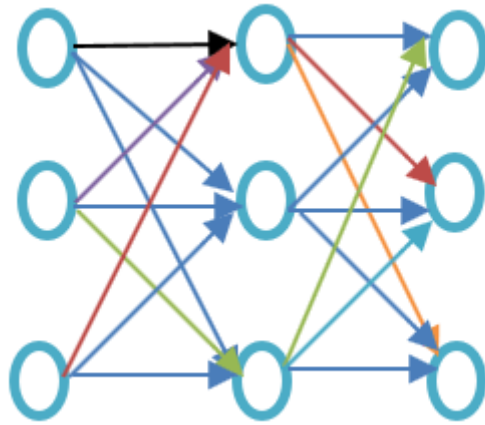
上面列举的是文本相关任务,当然目前词嵌入模型已经被扩展到方方面面。典型的, 例如:

- 在微博上面,每个人都用一个词来表示,对每个人构建Embedding,然后计算人之间的相关性,得到关系最为相近的人;
- 在推荐问题里面,依据每个用户的购买的商品记录,对每个商品进行Embedding,就可以计算商品之间的相关性,并进行推荐;
- 在此次天池的航海问题中,对相同经纬度上不同的船进行Embedding, 就可以得到每个船只的向量,就可以得到经常在某些区域工作的船只;

可以说,词嵌入为寻找物体之间相关性带来了巨大的帮助。现在基本每个数据竞赛都会见到Embedding技术。本文我们主要关注用的最多的Word2Vec模型。

4.什么是Word2Vec模型?

Word2vec是一种得到词表示的方法, 它可以较好地捕捉句法和语义词汇的关系。如果仅从网络结构看的话, 它就是一个两层的浅层网络。



即一个输入层+一个隐藏层+一个输出层。相较于潜在语义分析模型,Word2vec是一种更好以及更加高效的方案。

5.Word2Vec在做什么？

Word2vec在向量空间中对词进行表示,或者说词以向量的形式表示,在词向量空间中:相似含义的单词一起出现,而不同的单词则位于很远的地方。这也被称为语义关系。

神经网络不理解文本,而只理解数字。词嵌入提供了一种将文本转换为数字向量的方法。

Word2vec就是在重建词的语言上下文。那什么是语言上下文?在一般的生活情景中,当我们通过说话或写作来交流,其他人会试图找出句子的目的。例如,“印度的温度是多少”,这里的上下文是用户想知道“印度的温度”即上下文。

简而言之,句子的主要目标是语境。围绕口头或书面语言的单词或句子(披露)有助于确定上下文的意义。Word2vec通过上下文学习单词的矢量表示。

6.为什么用Word2Vec?

6.1 在词嵌入之前

关于为什么用Word2Vec好,这就要看在Word2Vec之前大家都在用什么方法,这些方法有什么不足,然后我们就能明白用Word2Vec的好处了。

6.2 潜在语义分析方法

潜在语义分析方法是在词向量之前用的最多的方法,它使用BOW(Bag of Words)的概念,每个词都被以编码的向量所表示,每个词都是一种稀疏的表示,其中向量的维度就是词汇表的大小。如果某个词出现了,那么我们会对它进行计数。

```
1 from sklearn.feature_extraction.text import CountVectorizer
```

```
2 vectorizer = CountVectorizer()  
3 data_corpus = ['guru99 is the best size for online tutorials. I love to visit  
4 vocabulary = vectorizer.fit(data_corpus)  
5 X          = vectorizer.transform(data_corpus)  
6 print(X.toarray())
```

```
[[1 1 2 1 1 1 1 1 1 1 1]]
```

```
1 print(vocabulary.get_feature_names())
```

```
['best', 'for', 'guru99', 'is', 'love', 'online', 'size', 'the', 'to', 'tutorials', 'visit']
```

在潜在语义分析中，每一行（特征列）表示的是某个词，每个列表示的是词出现在在某个文本中的次数。

后来很多学者发现Countervector方法会忽略词在不同文本词库中的出现情况，按理说如果某些词在不同的文本中都经常出现，那么应该降低此类词汇的重要性会更好。所以就出现了TFIDF方法。

- TFIDF中，字词的重要性随着它在文件中出现的次数成正比增加，但同时会随着它在语料库中出现的频率成反比下降。

但是其实TFIDF也存在很多明显的问题，它的不足如下(参见wiki)：

对区别文档最有意义的词语应该是那些在文档中出现频率高，而在整个文档集合的其他文档中出现频率少的词语，所以如果特征空间坐标系取tf词频作为测度，就可以体现同类文本的特点。另外考虑到单词区别不同类别的能力，tf-idf法认为一个单词出现的文本频数越小，它区别不同类别文本的能力就越大。因此引入了逆文本频度idf的概念，以tf和idf的乘积作为特征空间坐标系的取值测度，并用它完成对权值tf的调整，调整权值的目的在于突出重要单词，抑制次要单词。但是在本质上idf是一种试图抑制噪声的加权，并且单纯地认为文本频率小的单词就越重要，文本频率大的单词就越无用，显然这并不是完全正确的。idf的简单结构并不能有效地反映单词的重要程度和特征词的分布情况，使其无法很好地完成对权值调整的功能，所以tf-idf法的精度并不是很高。

6.3 BOW方法的问题

不管是Countervector还是TFIDF,我们发现它们都是从全局词汇的分布来对文本进行表示,所以缺点也明显,

- 它忽略了单个文本句子中词的顺序,例如 'this is bad' 在BOW中的表示和 'bad is this'是一样的;
- 它忽略了词的上下文,假设我们写一个句子,"He loved books. Education is best found in books".我们会在处理这两句话的时候是不会考虑前一个句子或者后一个句子是什么意思,但是他们之间是存在某些关系的

为了克服上述的两个缺陷, Word2Vec被开发出来并用来解决上述的两个问题。

7.Word2Vec如何工作?

Word2Vec通过学习与其相邻的上下文来进行预测，举例来说，我们要得到句子"He Loves Football"中Loves的词向量,这边我们假设：

loves = V_{in} . $P(V_{out} / V_{in})$ is calculated

where,

V_{in} is the input word.

P is the probability of likelihood.

V_{out} is the output word.

那么词"loves"在语料中移动每个单词。词与词之间的语法和语义便可以得到编码，这非常有助于帮助我们寻找相近与相似的词汇。

8.Word2Vec框架

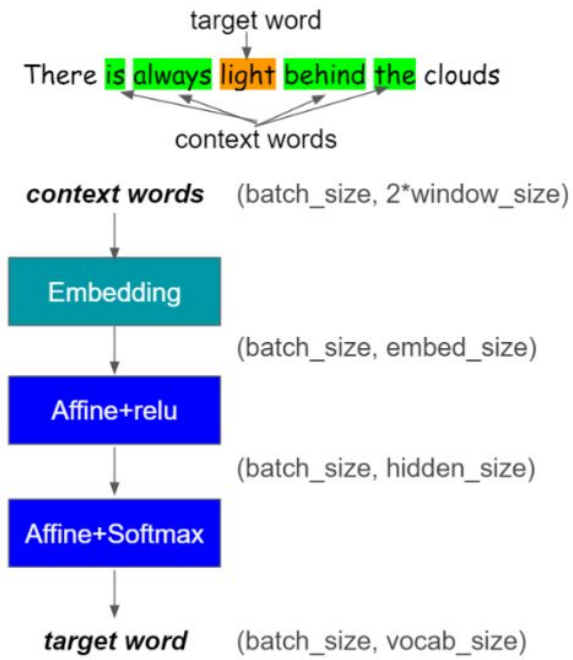
Word2Vec解决的问题就是上面说的,此处我们介绍Word2Vec的两种框架,

1. Continuous Bag of words (CBOW)
2. Skip gram

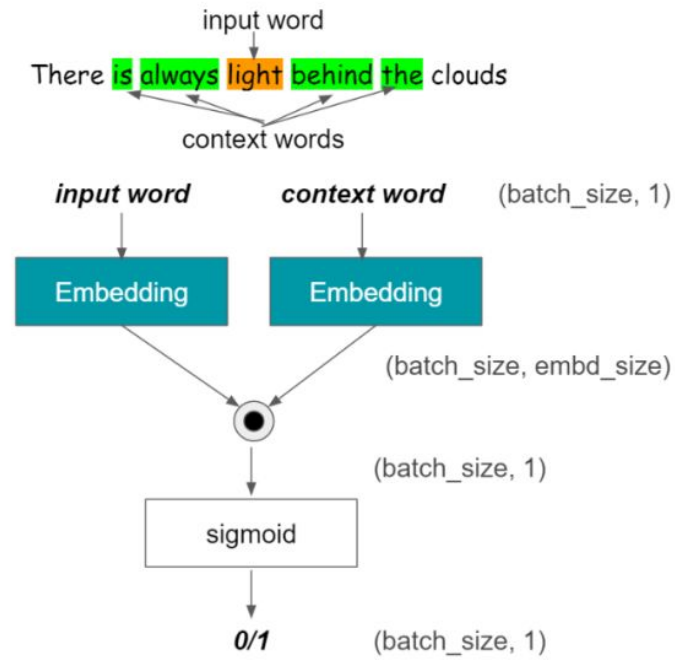
我们都知道学习词的表示是无监督的，但是如果没有targets/labels,那么我们将很难去训练该模型，Skip-gram和CBOW将无监督的表示到有监督的形式,这样便可以用于模型训练了。其中，

- CBOW,当前的词使用其周围的上下文(某个window size)的词进行预测，例如，如果 $w_{i-1}, w_{i-2}, w_{i+1}, w_{i+2}$ 是给定的词或者是上下文，那么我们的模型需要预测出是 w_i 。
- Skip-Gram则与CBOW相反,CBOW意味着它从单词中预测给定的序列或上下文。如果 w_i 是 给 定 的， 那 么 我 用 它 去 预 测 它 的 上 下 文 $w_{i-1}, w_{i-2}, w_{i+1}, w_{i+2}$ 。

二者的区别大家可以参考下图：



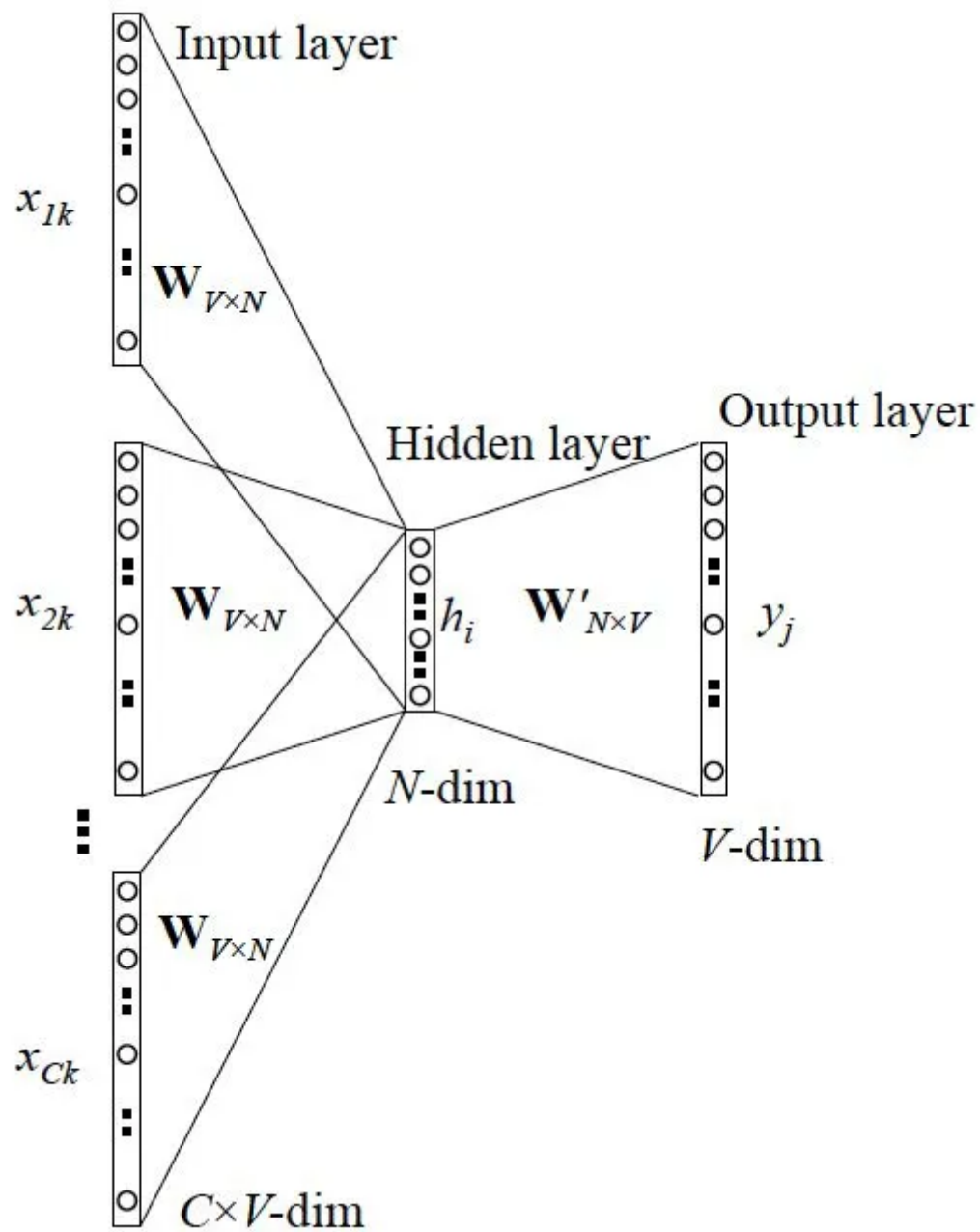
CBOW



Skip-gram

8.1 CBOW

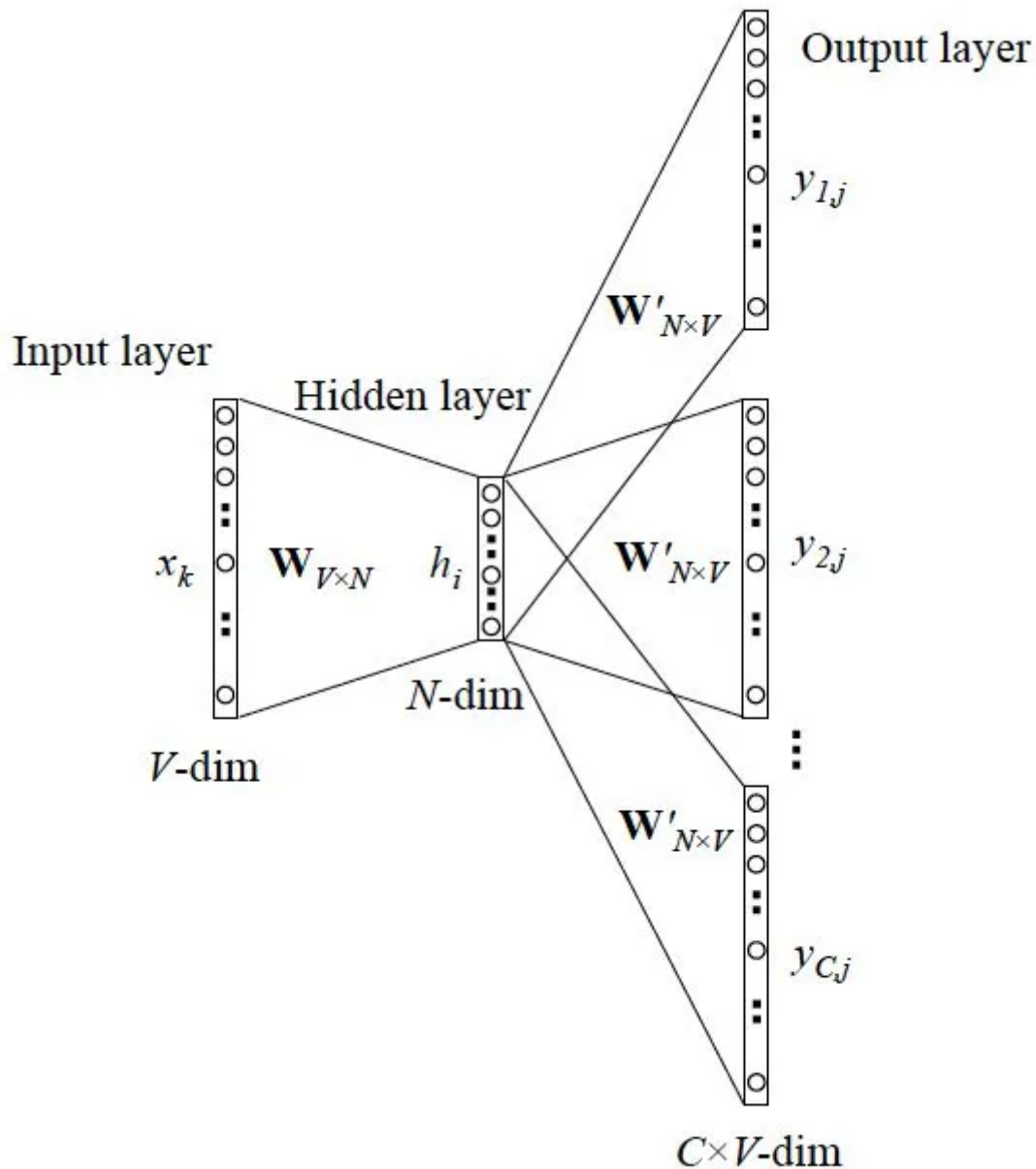
下面这个就是CBOW的图示，



我们假设 V 是词库的大小, N 是隐藏层的大小, 我们将输入定义为 $\{x_{i-1}, x_i, x_{i+1}, x_{i+2}\}$, 那么我们就得到一个权重矩阵, 它的大小是 $V * N$ 的. 即, 我们通过周围的词汇来对我们当前的词汇进行预测。

8.2 Skip-Gram

下面这个就是 Skip-Gram 的图示,



Skip-Gram模型和CBOW是相对的,我们通过当前的词汇去预测周围的词汇。

8.3 如何选择CBOW还是Skip-Gram

- CBOW训练的时候要比Skip-Gram要快很多;
- CBOW相较于Skip-Gram对于常见的词可以提供更好的表示;
- Skip-Gram从需要少量的训练数据集中也可以表示稀有的单词或者短语;

8.4 词向量理解

Word2Vec是希望把词映射到词向量空间中,那么这中间就需要一个这样的映射,这里怎么合理的理解? 我们在训练时,先将原词映射到 V 维的空间, V 是词库中不同词的个数。比如现在输入一个 x 的 one-hot encoder: $[1, 0, 0, \dots, 0]$, 对应一个简单的词, 则在输入层到隐含层的权重更新时, 只有对应 1 这个位置的权重被激活, 这些权重的个数, 跟隐含层节点数是一致的, 从而这些权重组成一个向量 v_x 来表示 x , 而因为每个词语的 one-hot encoder 里面 1 的位置是不同的, 所以, 这个向量 v_x 就可以用来唯一表示 x 。

9.代码展示

```
1 import torch
2 from torch.autograd import Variable
3 import torch.nn as nn
4 import torch.nn.functional as F
5 import torch.optim as optim
6
7 torch.manual_seed(1)
```

构建映射,将词映射到整数上

```
1 CONTEXT_SIZE = 2 # 2 words to the left, 2 to the right
2 text = """We are about to study the idea of a computational process.
3 Computational processes are abstract beings that inhabit computers.
4 As they evolve, processes manipulate other abstract things called data.
5 The evolution of a process is directed by a pattern of rules
6 called a program. People create programs to direct processes. In effect,
7 we conjure the spirits of the computer with our spells.""".split()
8
9 split_ind = (int)(len(text) * 0.8)
10 vocab = set(text)
11 vocab_size = len(vocab)
12 print('vocab_size:', vocab_size)
13
14 w2i = {w: i for i, w in enumerate(vocab)}
15 i2w = {i: w for i, w in enumerate(vocab)}
```

构建CBOW以及SkipGram所需的数据形式

```
1 # context window size is two
2 def create_cbow_dataset(text):
3     data = []
4     for i in range(2, len(text) - 2):
5         context = [text[i - 2], text[i - 1],
6                   text[i + 1], text[i + 2]]
```

```

7         target = text[i]
8         data.append((context, target))
9     return data
10
11 def create_skipgram_dataset(text):
12     import random
13     data = []
14     for i in range(2, len(text) - 2):
15         data.append((text[i], text[i-2], 1))
16         data.append((text[i], text[i-1], 1))
17         data.append((text[i], text[i+1], 1))
18         data.append((text[i], text[i+2], 1))
19         # negative sampling
20         for _ in range(4):
21             if random.random() < 0.5 or i >= len(text) - 3:
22                 rand_id = random.randint(0, i-1)
23             else:
24                 rand_id = random.randint(i+3, len(text)-1)
25             data.append((text[i], text[rand_id], 0))
26     return data
27
28 cbow_train = create_cbow_dataset(text)
29 skipgram_train = create_skipgram_dataset(text)
30 print('cbow sample', cbow_train[0])
31 print('skipgram sample', skipgram_train[0])

```

构建模型框架

```

1 class CBOW(nn.Module):
2     def __init__(self, vocab_size, embd_size, context_size, hidden_size):
3         super(CBOW, self).__init__()
4         self.embeddings = nn.Embedding(vocab_size, embd_size)
5         self.linear1 = nn.Linear(2*context_size*embd_size, hidden_size)
6         self.linear2 = nn.Linear(hidden_size, vocab_size)
7
8     def forward(self, inputs):
9         embedded = self.embeddings(inputs).view((1, -1))
10        hid = F.relu(self.linear1(embedded))
11        out = self.linear2(hid)
12        log_probs = F.log_softmax(out)

```

```

13         return log_probs
14
15 class SkipGram(nn.Module):
16     def __init__(self, vocab_size, embd_size):
17         super(SkipGram, self).__init__()
18         self.embeddings = nn.Embedding(vocab_size, embd_size)
19
20     def forward(self, focus, context):
21         embed_focus = self.embeddings(focus).view((1, -1))
22         embed_ctx = self.embeddings(context).view((1, -1))
23         score = torch.mm(embed_focus, torch.t(embed_ctx))
24         log_probs = F.logsigmoid(score)
25
26         return log_probs

```

进行模型训练

```

1  embd_size = 100
2  learning_rate = 0.001
3  n_epoch = 30
4
5  def train_cbow():
6      hidden_size = 64
7      losses = []
8      loss_fn = nn.NLLLoss()
9      model = CBOW(vocab_size, embd_size, CONTEXT_SIZE, hidden_size)
10     print(model)
11     optimizer = optim.SGD(model.parameters(), lr=learning_rate)
12
13     for epoch in range(n_epoch):
14         total_loss = .0
15         for context, target in cbow_train:
16             ctx_idxs = [w2i[w] for w in context]
17             ctx_var = Variable(torch.LongTensor(ctx_idxs))
18
19             model.zero_grad()
20             log_probs = model(ctx_var)
21
22             loss = loss_fn(log_probs, Variable(torch.LongTensor([w2i[target]])))
23

```

```
24         loss.backward()
25         optimizer.step()
26
27         total_loss += loss.data
28         losses.append(total_loss)
29     return model, losses
30
31 def train_skipgram():
32     losses = []
33     loss_fn = nn.MSELoss()
34     model = SkipGram(vocab_size, embd_size)
35     print(model)
36     optimizer = optim.SGD(model.parameters(), lr=learning_rate)
37
38     for epoch in range(n_epoch):
39         total_loss = .0
40         for in_w, out_w, target in skipgram_train:
41             in_w_var = Variable(torch.LongTensor([w2i[in_w]]))
42             out_w_var = Variable(torch.LongTensor([w2i[out_w]]))
43
44             model.zero_grad()
45             log_probs = model(in_w_var, out_w_var)
46             loss = loss_fn(log_probs[0], Variable(torch.Tensor([target])))
47
48             loss.backward()
49             optimizer.step()
50
51             total_loss += loss.data
52             losses.append(total_loss)
53     return model, losses
54
55 cbow_model, cbow_losses = train_cbow()
56 sg_model, sg_losses = train_skipgram()
```

模型测试

```
1 def test_cbow(test_data, model):
2     print('====Test CBOW====')
3     correct_ct = 0
```

```
4     for ctx, target in test_data:
5         ctx_idx = [w2i[w] for w in ctx]
6         ctx_var = Variable(torch.LongTensor(ctx_idx))
7
8         model.zero_grad()
9         log_probs = model(ctx_var)
10        _, predicted = torch.max(log_probs.data, 1)
11        #     print(predicted, int(predicted[0]))
12        predicted_word = i2w[int(predicted[0])]
13        print('predicted:', predicted_word)
14        print('label      :', target)
15        if predicted_word == target:
16            correct_ct += 1
17
18        print('Accuracy: {:.1f}% ({:d}/{:d})'.format(correct_ct/len(test_data)*100, correct_ct, len(test_data)))
19
20 def test_skipgram(test_data, model):
21     print('====Test SkipGram====')
22     correct_ct = 0
23     for in_w, out_w, target in test_data:
24         in_w_var = Variable(torch.LongTensor([w2i[in_w]]))
25         out_w_var = Variable(torch.LongTensor([w2i[out_w]]))
26
27         model.zero_grad()
28         log_probs = model(in_w_var, out_w_var)
29         _, predicted = torch.max(log_probs.data, 1)
30         predicted = predicted[0]
31         if predicted == target:
32             correct_ct += 1
33
34     print('Accuracy: {:.1f}% ({:d}/{:d})'.format(correct_ct/len(test_data)*100, correct_ct, len(test_data)))
35
36 test_cbow(cbow_train, cbow_model)
37 print('-----')
38 test_skipgram(skipgram_train, sg_model)
```