

TextCNN pytorch实现

原创

郑不凡

于 2021-11-02 16:59:12 发布

256

收藏 2

版权

文章标签:

pytorch

cn

深度学习

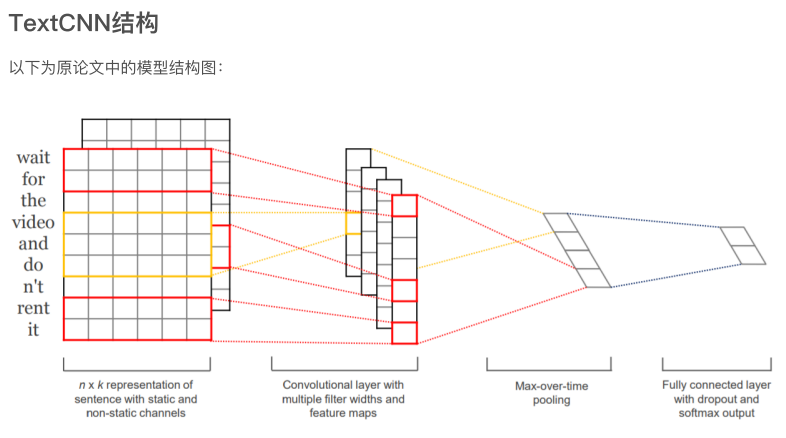


Figure 1: Model architecture with two channels for an example sentence @Owenmmmm

1. embedding
2. 卷积层
3. MaxPooling
4. Flatten
5. 全连接层

- 卷积操作
- 在卷积神经网络中仅涉及离散卷积的情形。
  - 卷积运算的作用就类似与滤波, 因此也称卷积核为filter滤波器。
  - 卷积神经网络的核心思想是捕捉局部特征 (n-gram)。CNN的优势在于能够自动地对g-gram特征进行组合和筛选, 获得不同抽象层次的语义信息。
  - 下图为用于文本分类任务的TextCNN结构描述 (这里详细解释了TextCNN架构以及词向量矩阵是如何做卷积的)



## 目录

### TextCNN结构

#### 卷积操作

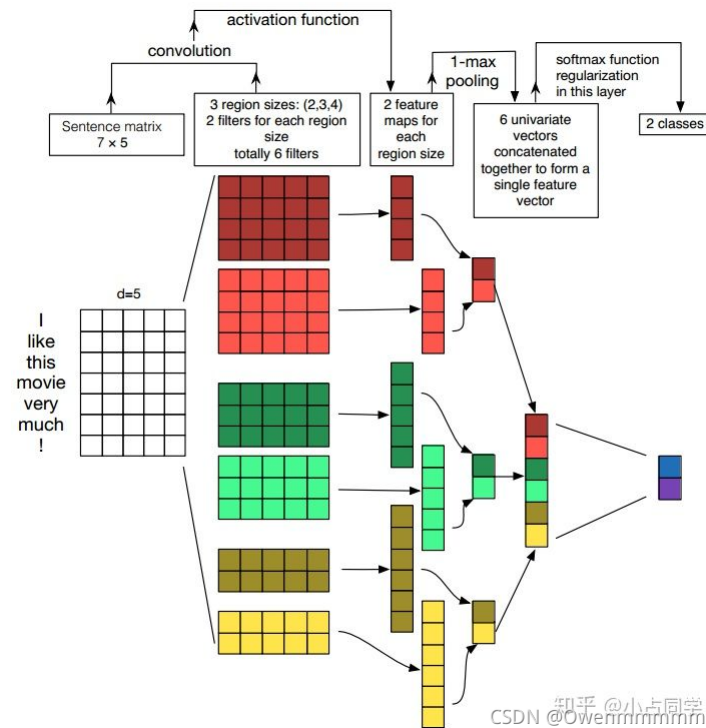
#### 简单代码实现

#### 维度变换

#### 用TextCNN进行IMDB电影评论情感分析

1. 数据预处理
2. 构建模型
3. 训练模型
4. 模型验证

#### 优化



1. 输入层： $n * k$ 的矩阵， $n$ 为句子中的单词数， $k$ 为embedding\_size。（为了使向量长度一致，对原句进行了padding操作）
2. 卷积层：在NLP中输入层是一个由词向量拼成的词矩阵，且卷积核的宽与该词矩阵的宽相同，该宽度即为词向量大小，且卷积核只会在高度方向移动。输入层的矩阵与我们的filter进行convolution，然后经过激活函数得到feature map。filter这里有三种大小（3，4，5）。
3. 池化层：max-pooling
4. softmax输出结果。

## 简单代码实现

```
1 import torch
2 import numpy as np
3 import torch.nn as nn
4 import torch.optim as optim
5 import torch.utils.data as Data
6 import torch.nn.functional as F
7
8 dtype = torch.FloatTensor
9 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
```

```
1 # 3 words sentences (=sequence_length is 3)
2 sentences = ["I love you", "he loves me", "she likes baseball", "I hate you", "sorr
3 labels = [1, 1, 1, 0, 0, 0] # 1 is good, 0 is not good.
4
5 embedding_size = 2
6 sequence_length = len(sentences[0])
7 num_classes = len(set(labels))
8 batch_size = 3
```

```

12 word2idx = {w:i for i,w in enumerate(vocab)}
13 vocab_size = len(vocab)

```

```

1 def make_data(sentences, labels):
2     inputs = []
3     for sen in sentences:
4         inputs.append([word2idx[n] for n in sen.split()])
5
6     targets = []
7     for out in labels:
8         targets.append(out)
9
10    return inputs, targets

```

```

1 input_batch, target_batch = make_data(sentences, labels)
2 input_batch, target_batch = torch.LongTensor(input_batch), torch.LongTensor(target_
3
4 dataset = Data.TensorDataset(input_batch, target_batch)
5 loader = Data.DataLoader(dataset, batch_size, True)

```

```

1 class TextCNN(nn.Module):
2
3     def __init__(self):
4         super(TextCNN, self).__init__()
5         self.W = nn.Embedding(vocab_size, embedding_size)
6         output_channel = 3
7         self.conv = nn.Sequential(nn.Conv2d(1, output_channel, (2, embedding_size)),
8                                   nn.ReLU(),
9                                   nn.MaxPool2d((2, 1)))
10        self.fc = nn.Linear(output_channel, num_classes)
11
12    def forward(self, X):
13        '''
14        X: [batch_size, sequence_length]
15        '''
16        batch_size = X.shape[0]
17        embedding_X = self.W(X) # [batch_size, sequence_length, embedding_size]
18        embedding_X = embedding_X.unsqueeze(1) # add channel(=1) [batch, channel(=1),
19        convded = self.conv(embedding_X) # [batch_size, output_channel, 1, 1]
20        flatten = convded.view(batch_size, -1) # [batch_size, output_channel*1*1]
21        output = self.fc(flatten)
22        return output

```

## 维度变换

1. 输入X: [batch\_size, sequence\_length]
2. embedding: 相当于把单词增加了一个维度。[batch\_size, sequence\_length, embedding\_size]; 然后我们对它做了一个unsqueeze (1) 操作, 原因是卷积操作的需要。[batch\_size, channel(=1), sequence\_length, embedding\_size]
3. convded: 我们这进行了一个二维卷积, input\_channel为1, output\_channel为3, filter\_size为 (2, embedding\_size) , 相当于bi-gram。  
[batch\_size, output\_channel, sequence\_len-1, 1]

卷积输出的高和宽的计算公式:

$$height_{out} = \frac{height_{in} - height_{kernel} + 2 \times padding}{stride} + 1$$

$$width_{out} = \frac{width_{in} - width_{kernel} + 2 \times padding}{stride} + 1$$

5. Flatten: [batch\_size, output\_channel]

```
1 | model = TextCNN().to(device)
2 | criterion = nn.CrossEntropyLoss().to(device)
3 | optimizer = optim.Adam(model.parameters(), lr=1e-3)
4 |
5 | # Training
6 | for epoch in range(5000):
7 |     for batch_x, batch_y in loader:
8 |         batch_x, batch_y = batch_x.to(device), batch_y.to(device)
9 |         pred = model(batch_x)
10 |         loss = criterion(pred, batch_y)
11 |         if (epoch + 1) % 1000 == 0:
12 |             print('Epoch:', '%04d' % (epoch + 1), 'loss =', '{:.6f}'.format(loss))
13 |
14 |         optimizer.zero_grad()
15 |         loss.backward()
16 |         optimizer.step()
```

```
1 | # Test
2 | test_text = 'i hate me'
3 | tests = [[word2idx[n] for n in test_text.split()]]
4 | test_batch = torch.LongTensor(tests).to(device)
5 | # Predict
6 | model = model.eval()
7 | predict = model(test_batch).data.max(1, keepdim=True)[1]
8 | if predict[0][0] == 0:
9 |     print(test_text, "is Bad Mean...")
10 | else:
11 |     print(test_text, "is Good Mean!!")
```

## 用TextCNN进行IMDB电影评论情感分析

### 1. 数据预处理

- 设置种子SEED，保证结果可复现
- 利用torchtext构建数据集

```
1 | import torch
2 | from torchtext.legacy import data
3 | from torchtext.legacy import datasets
4 | import random
5 | import numpy as np
6 |
7 | SEED = 1234
8 |
9 | random.seed(SEED)
10 | np.random.seed(SEED)
11 | torch.manual_seed(SEED)
12 | torch.backends.cudnn.deterministic = True
13 |
14 | TEXT = data.Field(tokenize = 'spacy',
15 |                   tokenizer_language = 'en_core_web_sm',
16 |                   batch_first = True)
17 | LABEL = data.LabelField(dtype = torch.float)
18 |
19 | train_data, test_data = datasets.IMDB.splits(TEXT, LABEL)
20 |
21 | train_data, valid_data = train_data.split(random_state = random.seed(SEED))
```

- 构建vocab，加载预训练词嵌入：

```
1 | MAX_VOCAB_SIZE = 25_000
2 |
```



```

6         unk_init = torch.Tensor.normal_)
7
8     LABEL.build_vocab(train_data)

```

- 创建迭代器：

```

1 BATCH_SIZE = 64
2
3 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
4
5 train_iterator, valid_iterator, test_iterator = data.BucketIterator.splits(
6     (train_data, valid_data, test_data),
7     batch_size = BATCH_SIZE,
8     device = device)

```

## 2. 构建模型

- nn.Embedding中padding\_idx：当Embedding是随机初始化的矩阵时，会对padding\_idx所在的行进行填0。保证了padding行为的正确性。

### 维度变换

1. 输入: text [batch\_size, sequence\_len]
2. embedding和unsqueeze(1): [batch\_size,1,sequence\_len,embedding\_dim]
3. convd:[batch\_size, output\_channel,sequence\_len-filter\_size[n]+1,1]
4. pooled:[batch\_size,output\_channel,1,1]
5. Flatten: 这里的flatten在第三步和第四步进行了squeeze操作[batch\_size,output\_channel]
6. concat: [batch\_size, output\_channel\*len(filter\_sizes)]

```

1 class CNN(nn.Module):
2     def __init__(self, vocab_size, embedding_dim, n_filters, filter_sizes, output_d
3         dropout, pad_idx):
4
5         super().__init__()
6
7         self.embedding = nn.Embedding(vocab_size, embedding_dim, padding_idx = pad_
8
9         self.convs = nn.ModuleList([
10             nn.Conv2d(in_channels = 1,
11                 out_channels = n_filters,
12                 kernel_size = (fs, embedding_dim))
13             for fs in filter_sizes
14         ])
15
16         self.fc = nn.Linear(len(filter_sizes) * n_filters, output_dim)
17
18         self.dropout = nn.Dropout(dropout)
19
20     def forward(self, text):
21
22         #text = [batch size, sent len]
23
24         embedded = self.embedding(text)
25
26         #embedded = [batch size, sent len, emb dim]
27
28         embedded = embedded.unsqueeze(1)
29
30         #embedded = [batch size, 1, sent len, emb dim]
31
32         convd = [F.relu(conv(embedded)).squeeze(3) for conv in self.convs]
33
34         #convd_n = [batch size, n_filters, sent len - filter_sizes[n] + 1]
35
36         (2) for conv in convd]

```

```

38         #pooled_n = [batch_size, n_filters]
39
40         cat = self.dropout(torch.cat(pooled, dim = 1))
41
42         #cat = [batch_size, n_filters * len(filter_sizes)]
43
44         return self.fc(cat)

```

- 设置超参数，实例化模型

```

1 INPUT_DIM = len(TEXT.vocab)
2 EMBEDDING_DIM = 100
3 N_FILTERS = 100
4 FILTER_SIZES = [3,4,5]
5 OUTPUT_DIM = 1
6 DROPOUT = 0.5
7 PAD_IDX = TEXT.vocab.stoi[TEXT.pad_token]
8
9 model = CNN(INPUT_DIM, EMBEDDING_DIM, N_FILTERS, FILTER_SIZES, OUTPUT_DIM, DROPOUT,

```

- 查看参数

```

1 def count_parameters(model):
2     return sum(p.numel() for p in model.parameters() if p.requires_grad)
3
4 print(f'The model has {count_parameters(model):,} trainable parameters')

```

- 加载预训练词嵌入

```

1 pretrained_embeddings = TEXT.vocab.vectors
2
3 model.embedding.weight.data.copy_(pretrained_embeddings)

```

- 将unk 和 pad 初始权重归零。

```

1 UNK_IDX = TEXT.vocab.stoi[TEXT.unk_token]
2
3 model.embedding.weight.data[UNK_IDX] = torch.zeros(EMBEDDING_DIM)
4 model.embedding.weight.data[PAD_IDX] = torch.zeros(EMBEDDING_DIM)

```

### 3. 训练模型

- 初始化优化器、损失函数

```

1 import torch.optim as optim
2
3 optimizer = optim.Adam(model.parameters())
4
5 criterion = nn.BCEWithLogitsLoss()
6
7 model = model.to(device)
8 criterion = criterion.to(device)

```

- 计算精度的函数

```

1 def binary_accuracy(preds, y):
2     """
3     Returns accuracy per batch, i.e. if you get 8/10 right, this returns 0.8, NOT 8
4     """
5
6

```



```
9     acc = correct.sum() / len(correct)
10     return acc
```

- train

```
1 def train(model, iterator, optimizer, criterion):
2
3     epoch_loss = 0
4     epoch_acc = 0
5
6     model.train()
7
8     for batch in iterator:
9
10         optimizer.zero_grad()
11
12         predictions = model(batch.text).squeeze(1)
13
14         loss = criterion(predictions, batch.label)
15
16         acc = binary_accuracy(predictions, batch.label)
17
18         loss.backward()
19
20         optimizer.step()
21
22         epoch_loss += loss.item()
23         epoch_acc += acc.item()
24
25     return epoch_loss / len(iterator), epoch_acc / len(iterator)
```

- eval (注意: 同样, 由于使用的是dropout, 我们必须记住使用 `model.eval()` 来确保在评估时能够关闭 dropout。)

```
1 def evaluate(model, iterator, criterion):
2
3     epoch_loss = 0
4     epoch_acc = 0
5
6     model.eval()
7
8     with torch.no_grad():
9
10         for batch in iterator:
11
12             predictions = model(batch.text).squeeze(1)
13
14             loss = criterion(predictions, batch.label)
15
16             acc = binary_accuracy(predictions, batch.label)
17
18             epoch_loss += loss.item()
19             epoch_acc += acc.item()
20
21     return epoch_loss / len(iterator), epoch_acc / len(iterator)
```

```
1 import time
2
3 def epoch_time(start_time, end_time):
4     elapsed_time = end_time - start_time
5     elapsed_mins = int(elapsed_time / 60)
6     elapsed_secs = int(elapsed_time - (elapsed_mins * 60))
7     return elapsed_mins, elapsed_secs
```

- training



举报

```

1 | N_EPOCHS = 5
2 |
3 | best_valid_loss = float('inf')
4 |
5 | for epoch in range(N_EPOCHS):
6 |
7 |     start_time = time.time()
8 |
9 |     train_loss, train_acc = train(model, train_iterator, optimizer, criterion)
10 |    valid_loss, valid_acc = evaluate(model, valid_iterator, criterion)
11 |
12 |    end_time = time.time()
13 |
14 |    epoch_mins, epoch_secs = epoch_time(start_time, end_time)
15 |
16 |    if valid_loss < best_valid_loss:
17 |        best_valid_loss = valid_loss
18 |        torch.save(model.state_dict(), 'tut4-model.pt')
19 |
20 |    print(f'Epoch: {epoch+1:02} | Epoch Time: {epoch_mins}m {epoch_secs}s')
21 |    print(f'\tTrain Loss: {train_loss:.3f} | Train Acc: {train_acc*100:.2f}%')
22 |    print(f'\tVal. Loss: {valid_loss:.3f} | Val. Acc: {valid_acc*100:.2f}%')

```

- test

```

1 | model.load_state_dict(torch.load('tut4-model.pt'))
2 |
3 | test_loss, test_acc = evaluate(model, test_iterator, criterion)
4 |
5 | print(f'Test Loss: {test_loss:.3f} | Test Acc: {test_acc*100:.2f}%')

```

#### 4. 模型验证

```

1 | import spacy
2 | nlp = spacy.load('en_core_web_sm')
3 |
4 | def predict_sentiment(model, sentence, min_len = 5):
5 |     model.eval()
6 |     tokenized = [tok.text for tok in nlp.tokenizer(sentence)]
7 |     if len(tokenized) < min_len:
8 |         tokenized += ['<pad>'] * (min_len - len(tokenized))
9 |     indexed = [TEXT.vocab.stoi[t] for t in tokenized]
10 |    tensor = torch.LongTensor(indexed).to(device)
11 |    tensor = tensor.unsqueeze(0)
12 |    prediction = torch.sigmoid(model(tensor))
13 |    return prediction.item()

```

- 输入需要评测的句子

```

1 | predict_sentiment(model, "This film is terrible")
2 | predict_sentiment(model, "This film is great")

```

#### 优化

- 在TextCNN的实践中，有很多地方可以优化：

Filter尺寸：这个参数决定了抽取n-gram特征的长度，这个参数主要跟数据有关，平均长度在50以内的话，用10以下就可以了，否则可以长一些。在调参时可以先用一个尺寸grid search，找到一个最优尺寸，然后尝试最优尺寸和附近尺寸的组合

Filter个数：这个参数会影响最终特征的维度，维度太大的话训练速度就会变慢。这里在100-600之间调参即可





正则化：指对CNN参数的正则化，可以使用dropout或L2，但能起的作用很小，可以试下小的dropout率(<0.5)，L2限制大一点

Pooling方法：根据情况选择mean、max、k-max pooling，大部分时候max表现就很好，因为分类任务对细粒度语义的要求不高，只抓住最大特征就好了

Embedding表：中文可以选择char或word级别的输入，也可以两种都用，会提升些效果。如果训练数据充足（10w+），也可以从头训练  
蒸馏BERT的logits，利用领域内无监督数据

加深全连接：原文文只使用了一层全连接，而加到3、4层左右效果会更好

TextCNN是很适合中短文本场景的强baseline，但不太适合长文本，因为 卷积核 尺寸通常不会设很大，无法捕获长距离特征。同时max-pooling也存在局限，会丢掉一些有用特征。另外再仔细想的话，TextCNN和传统的n-gram词袋模型本质是一样的，它的好效果很大部分来自于词向量的引入，解决了词袋模型的稀疏性问题。

参考文献  
TextCNN 的 PyTorch 实现  
深入TextCNN（一）详述CNN及TextCNN原理  
深度学习文本分类模型综述+代码+技巧

Pytorch 实现文本分类 06-11  
文本分类的标准代码，Pytorch实现 数据集Dataset - IMDB - SST - Trec ### 模型 - FastText - BasicCNN (KimCNN,MultiLay...  
文本分类系列(1):textcnn及其pytorch实现 u014514939的博客 8853  
textcnn 原理：核心点在于使用卷积来捕捉局部相关性，具体到文本分类任务中可以利用CNN来提取句子中类似 n-gram 的...

参与评论 您还未登录，请先 登录 后发表或查看评论

TextCNN的PyTorch实现\_数学家是我理想的博客\_textcnn... 3-28  
target\_batch=torch.LongTensor(input\_batch),torch.LongTensor(target\_batch)dataset=Data.TensorDataset(input\_batch,targ...  
TextCNN\_pytorch实现\_\_Wooden\_的博客 2-17  
TextCNN\_pytorch实现 importnumpyasnmp importtorch fromtorch.functionalimportsplit importtorch.nnasnn importtorch.optim...  
Pytorch实现TextCNN模型对IMDB数据集进行分类 红雨瓢泼的博客 5943  
文章目录1 TextCNN简介1.1 论文简介1.2 模型简介2 数据处理 1 TextCNN简介 1.1 论文简介 本文实现的TextCNN模型来源...  
pytorch学习之textCNN实现 yellow\_red\_people的博客 1万+  
最近在学pytorch，所以尝试使用pytorch实现textCNN，ps(git 上有其他人textCNN的实现)。pytorch比tensorflow好的一个地...  
TextCNN的PyTorch实现 mathor的博客 1398  
本文主要介绍一篇将CNN应用到NLP领域的一篇论文 Convolutional Neural Networks for Sentence Classification，然后给...  
pytorch实现textCNN 热门推荐 无所知的博客 2万+  
pytorch实现textCNN1. 原理2. 数据预处理2.1 转换为csv格式2.2 观察数据分布2.3 由文本得到训练用的mini-batch数据3. 模...  
pytorch 实现 textCNN 杂文集 3040  
textCNN 模型 textCNN模型主要使用了一维卷积层和时序最大池化层。假设输入的文本序列由nn个词组成，每个词用dd维...  
使用pytorch实现简单的TextCNN（步骤详解） qq\_37822951的博客 6668  
本文中将根据原理图所示，使用pytorch搭建简单的TextCNN网络结构。尽可能详细的描述每个步骤。TextCNN的详细过程...  
Pytorch-IMDB电影评论情感分析 每天起床第一句要给自己打打气 5047  
原文链接: http://chenhao.space/post/a5b86241.html Pytorch-情感分析 第一步：导入IMDB电影数据集，只有训练集和测试...  
【pytorch模型实现4】TextCNN lyj223061的博客 75  
TextCNN模型实现 NLP模型代码github仓库: https://github.com/lyj157175/Models import torch import torch.nn as nn impor...  
gbz3300255的博客 305  
函数。是log\_softmax和 nll\_loss...

NLP学习之使用pytorch搭建TextCNN模型进行中文文本分类	持久决心的博客	2万+
最近花周末两天时间利用pytorch实现了TextCNN进行了中文文本分类，在此进行记录。数据获取 中文数据是从https://github...		
使用pytorch进行IMDB情感分析	最新发布	weixin_43991446的博客
使用pytorch进行IMDB情感分析 建议：将代码整合到main()函数中。1. 配置 1.1 设置cuda和随机种子 # 设置cuda device = ...		
chinese_text_cnn: TextCNN Pytorch实现中文文本分类情感分析		02-03
TextCNN Pytorch实现中文文本分类 论文 参考 依赖项 python3.5 pytorch == 1.0.0 torchtext == 0.3.1 jieba == 0.39 词向量...		
pytorch搭建TextCNN与使用案例		呆萌的代Ma
TextCNN算法流程 整体流程是将词拼接在一起，一句话构成一个特征图 根据卷积核得到多个特征向量 每个特征向量全局池...		
pytorch构建LSTM分类器用于IMDB情感分类		Lyric1的博客
本文基于pytorch构建LSTM情感分类分类器，在IMDB数据集上进行测试，涉及文本预处理、数据集加载、模型训练、保存...		
情感分析-火炬: IMDb数据集上的火炬情感分析		02-04
使用PyTorch进行情感分析 存储库将引导您完成构建完整的情感分析模型的过程，该模型将能够预测给定评论的极性（无论...		
pytorch实现的TextCNN (Dataset, DataLoader的使用)		SinGain的博客
主要是Dataset, DataLoader的使用 (1) 数据处理，生成Batch和向量化词表 import torch import numpy as np from tqdm i...		
pytorch实现IMDB数据集情感分类 (全连接层的网络、LSTM)		Delusional的博客
任务描述 使用Pytorch相关API，设计两种网络结构，一种网络结构中只有全连接层，一种使用文本处理中最为常用的LSTM...		
pytorch使用套路 (以RNN分析imdb文本情感为例)		JJY的博客
这里写目录标题1.导入包2.预定义参数3.加载数据4.构建网络5.定义训练函数6.定义测试函数7.开始训练 在使用过keras和pyt...		

“相关推荐”对你有帮助么？

非常没帮助

没帮助

一般

有帮助

非常有帮助

©2022 CSDN 皮肤主题：数字20 设计师：CSDN官方博客 返回首页

举报

郑不凡

关注

1

0

2