

# 【NLP】基于深度学习的文本分类应用

机器学习初学者 8月13日

以下文章来源于Datawhale，作者罗美君



**Datawhale**

一个专注于AI领域的开源组织，汇集了众多领域院校和知名企业的优秀学习者，聚合了...

**作者：罗美君，算法工程师，Datawhale优秀学习者**

在基于机器学习的文本分类中，我们介绍了几种常见的文本表示方法：One-hot、Bags of Words、N-gram、TF-IDF。这些方法存在两个共同的问题：一是转换得到的向量维度很高，需要较长的训练实践；二是没有考虑到单词与单词之间的关系，只是进行了统计。

与上述表示方法不同，深度学习也可以用于文本表示，并可以将其映射到一个低维空间。fastText是Facebook2016年提出的文本分类工具，是一种高效的浅层网络。今天我们就尝试使用fastText模型进行文本分类。

## 1. 数据及背景

<https://tianchi.aliyun.com/competition/entrance/531810/information>（阿里天池-零基础入门NLP赛事）

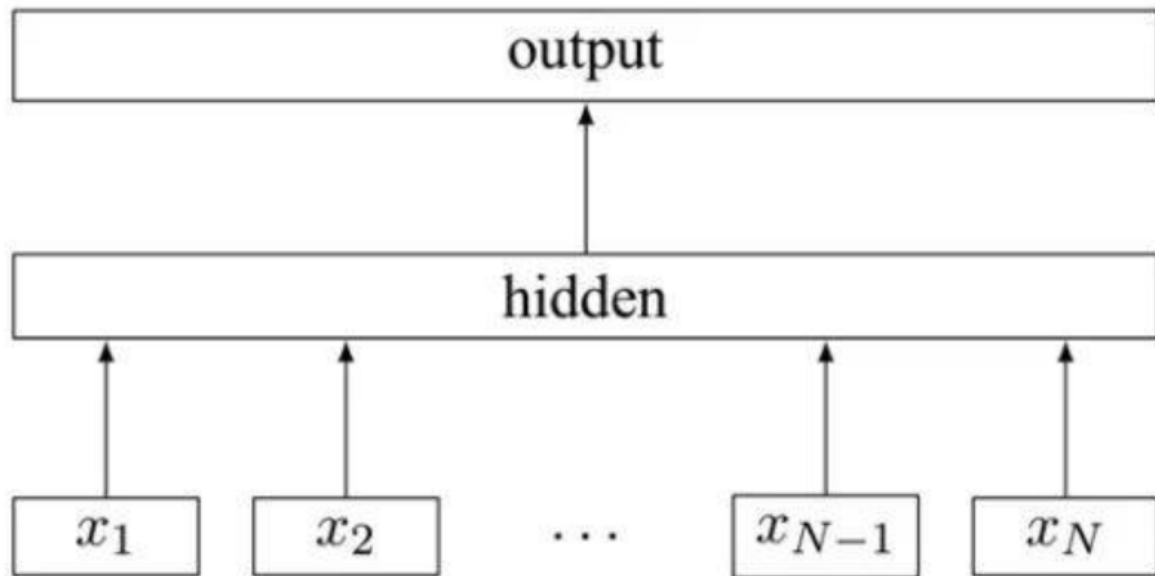
## 2. fastText模型剖析

### 2.1 概念

FastText是一种典型的深度学习词向量的表示方法，它的核心思想是将整篇文档的词及n-gram向量叠加平均得到文档向量，然后使用文档向量做softmax多分类。这中间涉及到两个技巧：字符级n-gram特征的引入以及分层Softmax分类。

### 2.2 模型框架

fastText模型架构和word2vec的CBOW模型架构非常相似。下面是fastText模型架构图：



**Figure 1:** Model architecture of `fastText` for a sentence with  $N$  ngram features  $x_1, \dots, x_N$ . The features are embedded and averaged to form the hidden variable.

注意：此架构图没有展示词向量的训练过程。可以看到，和CBOW一样，`fastText`模型也只有三层：输入层、隐含层、输出层（Hierarchical Softmax），输入都是多个经向量表示的单词，输出都是一个特定的target，隐含层都是对多个词向量的叠加平均。

不同的是，CBOW的输入是目标单词的上下文，`fastText`的输入是多个单词及其n-gram特征，这些特征用来表示单个文档；CBOW的输入单词被onehot编码过，`fastText`的输入特征是被embedding过；CBOW的输出是目标词汇，`fastText`的输出是文档对应的类标。

值得注意的是，`fastText`在输入时，将单词的字符级别的n-gram向量作为额外的特征；在输出时，`fastText`采用了分层Softmax，大大降低了模型训练时间。

### 2.3 字符级别的n-gram

`word2vec`把语料库中的每个单词当成原子的，它会为每个单词生成一个向量。这忽略了单词内部的形态特征，比如："apple" 和"apples"，"达观数据"和"达观"，这两个例子中，两个单词都有较多公共字符，即它们的内部形态类似，但是在传统的`word2vec`中，这种单词内部形态信息因为它们被转换成不同的id丢失了。

为了克服这个问题，`fastText`使用了字符级别的n-grams来表示一个单词。对于单词"apple"，假设 $n$ 的取值为3，则它的trigram有：

"<ap", "app", "ppl", "ple", "le>"

其中，<表示前缀，>表示后缀。于是，我们可以用这些trigram来表示"apple"这个单词，进一步，我们可以用这5个trigram的向量叠加来表示"apple"的词向量。

这带来两点好处：

- 对于低频词生成的词向量效果会更好。因为它们的n-gram可以和其它词共享。
- 对于训练词库之外的单词，仍然可以构建它们的词向量。我们可以叠加它们的字符级n-gram向量。

## 2.4 分层softmax

fastText的结构：

- 文本分词后排成列做输入。
- lookup table变成想要的隐层维数。
- 隐层后接huffman Tree。这个tree就是分层softmax减少计算量的精髓。

## 3. 简单实现fastText

为了简化任务：

1. 训练词向量时，我们使用正常的word2vec方法，而真实的fastText还附加了字符级别的n-gram作为特征输入；
2. 我们的输出层使用简单的softmax分类，而真实的fastText使用的是Hierarchical Softmax。

首先定义几个常量：

- VOCAB\_SIZE = 2000
- EMBEDDING\_DIM = 100
- MAX\_WORDS = 500
- CLASS\_NUM = 5
- VOCAB\_SIZE表示词汇表大小，这里简单设置为2000；

EMBEDDING\_DIM表示经过embedding层输出，每个词被分布式表示的向量的维度，这里设置为100。比如对于“达观”这个词，会被一个长度为100的类似于[ 0.97860014, 5.93589592, 0.22342691, -3.83102846, -0.23053935, ...]的实值向量来表示；

MAX\_WORDS表示一篇文档最多使用的词个数，因为文档可能长短不一（即词数不同），为了能feed到一个固定维度的神经网络，我们需要设置一个最大词数，对于词数少于这个阈值的文档，我们需要用“未知词”去填充。比如可以设置词汇表中索引为0的词为“未知词”，用0去填充少于阈值的部分；

CLASS\_NUM表示类别数，多分类问题，这里简单设置为5。

模型搭建遵循以下步骤：

1. 添加输入层（embedding层）。Embedding层的输入是一批文档，每个文档由一个词汇索引序列构成。例如：[10, 30, 80, 1000] 可能表示“我 昨天 来到 达观数据”这个短文本，其中“我”、“昨天”、“来到”、“达观数据”在词汇表中的索引分别是10、30、80、1000；Embedding层将每个单词映射成EMBEDDING\_DIM维的向量。于是：input\_shape=(BATCH\_SIZE, MAX\_WORDS), output\_shape=(BATCH\_SIZE, MAX\_WORDS, EMBEDDING\_DIM);
2. 添加隐含层（投影层）。投影层对一个文档中所有单词的向量进行叠加平均。keras提供的GlobalAveragePooling1D类可以帮我们实现这个功能。这层的input\_shape是Embedding层的output\_shape，这层的output\_shape=(BATCH\_SIZE, EMBEDDING\_DIM);
3. 添加输出层（softmax层）。真实的fastText这层是Hierarchical Softmax，因为keras原生并没有支持Hierarchical Softmax，所以这里用Softmax代替。这层指定了CLASS\_NUM，对于一篇文档，输出层会产生CLASS\_NUM个概率值，分别表示此文档属于当前类的可能性。这层的output\_shape=(BATCH\_SIZE, CLASS\_NUM)
4. 指定损失函数、优化器类型、评价指标，编译模型。损失函数我们设置为categorical\_crossentropy，它就是我们上面所说的softmax回归的损失函数；优化器我们设置为SGD，表示随机梯度下降优化器；评价指标选择accuracy，表示精度。

用训练数据feed模型时，你需要：

1. 将文档分好词，构建词汇表。词汇表中每个词用一个整数（索引）来代替，并预留“未知词”索引，假设为0；
2. 对类标进行onehot化。假设我们文本数据总共有3个类别，对应的类标分别是1、2、3，那么这三个类标对应的onehot向量分别是[1, 0, 0]、[0, 1, 0]、[0, 0, 1]；
3. 对一批文本，将每个文本转化为词索引序列，每个类标转化为onehot向量。就像之前的例子，“我 昨天 来到 达观数据”可能被转化为[10, 30, 80, 1000]；它属于类别1，它的类标就是[1, 0, 0]。由于我们设置了MAX\_WORDS=500，这个短文本向量后面就需要补496个0，即[10, 30, 80, 1000, 0, 0, 0, ..., 0]。因此，batch\_xs的维度为(BATCH\_SIZE, MAX\_WORDS)，batch\_ys的维度为(BATCH\_SIZE, CLASS\_NUM)。

代码如下：

```
1 # coding: utf-8
2 from __future__ import unicode_literals
3
4 from keras.models import Sequential
5 from keras.layers import Embedding
6 from keras.layers import GlobalAveragePooling1D
```

```
7 from keras.layers import Dense
8
9 VOCAB_SIZE = 2000
10 EMBEDDING_DIM = 100
11 MAX_WORDS = 500
12 CLASS_NUM = 5
13
14
15 def build_fastText():
16     model = Sequential()
17     # 将词汇数VOCAB_SIZE映射为EMBEDDING_DIM维
18     model.add(Embedding(VOCAB_SIZE, EMBEDDING_DIM, input_length=MAX_WORDS))
19     # 平均文档中所有词的embedding
20     model.add(GlobalAveragePooling1D())
21     # softmax分类
22     model.add(Dense(CLASS_NUM, activation='softmax'))
23     # 定义损失函数、优化器、分类度量指标
24     model.compile(loss='categorical_crossentropy', optimizer='SGD', metrics=|
25     return model
26
27 if __name__ == '__main__':
28     model = build_fastText()
29     print(model.summary())
```

## 4. 使用fastText文本分类

### 4.1 加载库

```
1 import time
2 import numpy as np
3 import fasttext
4 import pandas as pd
5
6 from sklearn.metrics import f1_score
7 from sklearn.utils import shuffle
8 from sklearn.model_selection import StratifiedKFold
```

## 4.2 fastText分类

主要超参数：

- lr: 学习率
- dim: 词向量的维度
- epoch: 每轮的个数
- wordNgrams: 词的n-gram, 一般设置为2或3
- loss: 损失函数 ns(negative sampling, 负采样)、hs(hierarchical softmax, 分层softmax)、softmax、ova(One-VS-ALL)

```

1 def fasttext_model(nrows, train_num, lr=1.0, wordNgrams=2, minCount=1, epoch=
2     start_time = time.time()
3
4     # 转换为FastText需要的格式
5     train_df = pd.read_csv('../input/train_set.csv', sep='\t', nrows=nrows)
6
7     # shuffle
8     train_df = shuffle(train_df, random_state=666)
9
10    train_df['label_ft'] = '__label__' + train_df['label'].astype('str')
11    train_df[['text', 'label_ft']].iloc[:train_num].to_csv('../input/fastText
12
13    model = fasttext.train_supervised('../input/fastText_train.csv', lr=lr, w
14                                     minCount=minCount, epoch=epoch, loss=lc
15
16    train_pred = [model.predict(x)[0][0].split('__')[-1] for x in train_df.it
17    print('Train f1_score:', f1_score(train_df['label'].values[:train_num].as
18    val_pred = [model.predict(x)[0][0].split('__')[-1] for x in train_df.ilo
19    print('Val f1_score:', f1_score(train_df['label'].values[train_num:].ast
20    train_time = time.time()
21    print('Train time: {:.2f}s'.format(train_time - start_time))
22
23    # 预测并保存
24    test_df = pd.read_csv('../input/test_a.csv')
25
26    test_pred = [model.predict(x)[0][0].split('__')[-1] for x in test_df['tex
27    test_pred = pd.DataFrame(test_pred, columns=['label'])
28    test_pred.to_csv('../input/test_fastText_ridgeclassifier.csv', index=False
29    print('Test predict saved.')
30    end_time = time.time()

```

```
31     print('Predict time:{:.2f}s'.format(end_time - train_time))
32
33
34 if __name__ == '__main__':
35     nrows = 200000
36     train_num = int(nrows * 0.7)
37     lr=0.01
38     wordNgrams=2
39     minCount=1
40     epoch=25
41     loss='hs'
42
43     fasttext_model(nrows, train_num)
```

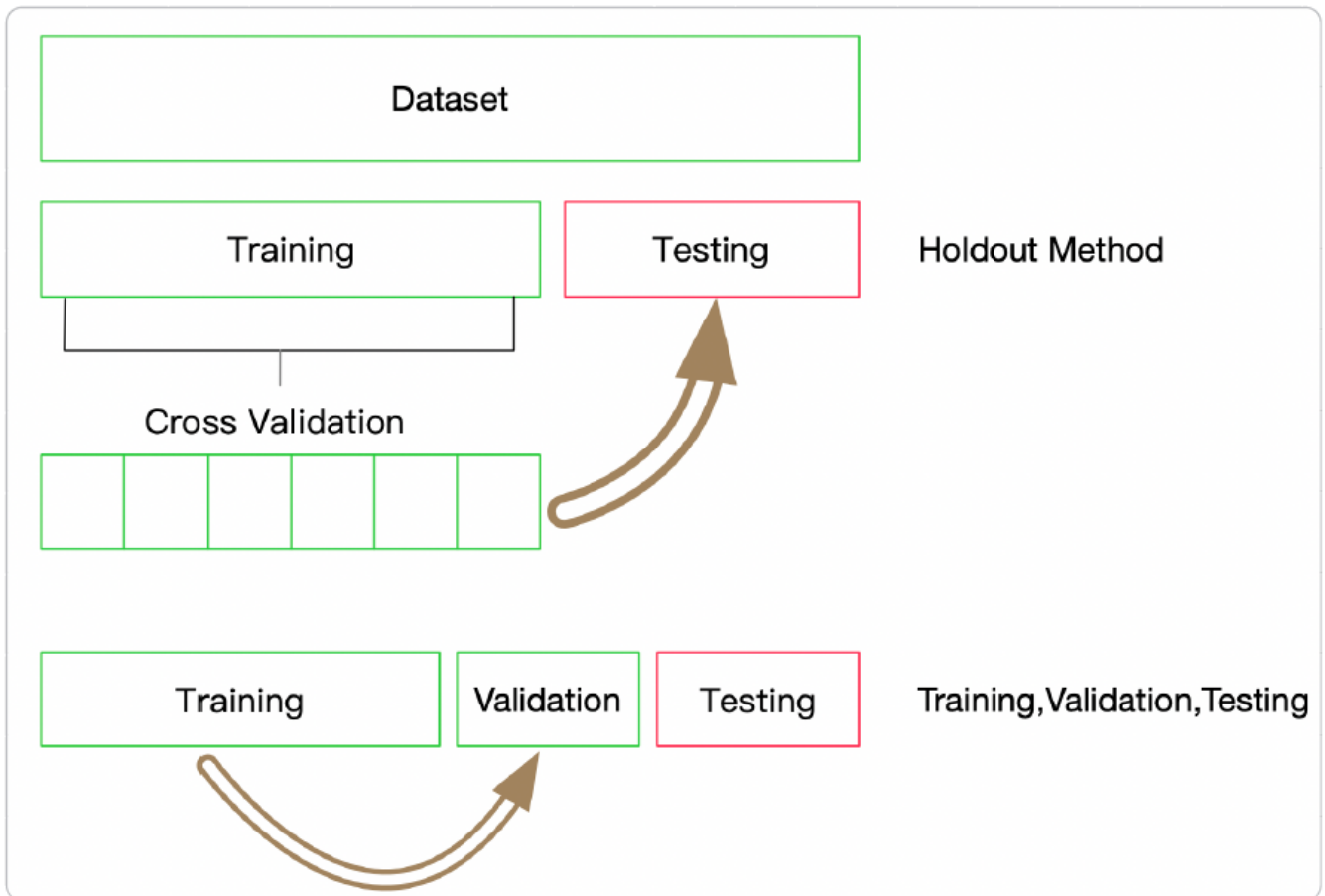
结果:

```
1 Train f1_score: 0.998663548149514
2 Val f1_score: 0.911468448971427
3 Train time: 257.32s
4 Test predict saved.
5 Predict time:13.40s
```

### 4.3 K折交叉验证

在使用FastText中，有一些模型的参数需要选择，这些参数会在一定程度上影响模型的精度，那么如何选择这些参数呢？有两种方式：

- 通过阅读文档，要弄清楚这些参数的含义，哪些参数会增加模型的复杂度；
- 通过在验证集上进行验证模型精度，找到模型是否过拟合或欠拟合。



这里我们采用第二种方法，用K折交叉验证的思想进行参数调节。注意：每折的划分必须保证标签的分布与整个数据集的分布一致。

```

1 models = []
2 scores = []
3 pred_list = []
4
5 # K折交叉验证
6 skf = StratifiedKFold(n_splits=n_splits, shuffle=True, random_state=666)
7 for train_index, test_index in skf.split(train_df['text'], train_df['label_ft']):
8
9     train_df[['text', 'label_ft']].iloc[train_index].to_csv('../input/fastText_train.csv')
10
11     model = fasttext.train_supervised('../input/fastText_train.csv', lr=lr, w=word_embeddings,
12                                     minCount=minCount, epoch=epoch, loss=loss)
13     models.append(model)
14
15     val_pred = [model.predict(x)[0][0].split('__')[-1] for x in train_df.iloc[test_index]]
16     score = f1_score(train_df['label'].values[test_index].astype(str), val_pred)
17     print('score', score)
18     scores.append(score)

```



```

19
20 print('mean score: ', np.mean(scores))
21 train_time = time.time()
22 print('Train time: {:.2f}s'.format(train_time - start_time))

```

## 所有代码

```

1 def fasttext_kfold_model(nrows, train_num, n_splits, lr=1.0, wordNgrams=2, minCount=10, epoch=10, loss=1.0):
2     start_time = time.time()
3
4     # 转换为FastText需要的格式
5     train_df = pd.read_csv('../input/train_set.csv', sep='\t', nrows=nrows)
6
7     # shuffle
8     train_df = shuffle(train_df, random_state=666)
9
10    train_df['label_ft'] = '__label__' + train_df['label'].astype('str')
11
12    models = []
13    train_scores = []
14    val_scores = []
15
16    # K折交叉验证
17    skf = StratifiedKFold(n_splits=n_splits, shuffle=True, random_state=666)
18    for train_index, test_index in skf.split(train_df['text'], train_df['label']):
19        train_df[['text', 'label_ft']].iloc[train_index].to_csv('../input/fastText_train.csv')
20
21        model = fasttext.train_supervised('../input/fastText_train.csv', lr=lr, minCount=minCount, epoch=epoch, loss=loss)
22
23        models.append(model)
24
25        train_pred = [model.predict(x)[0][0].split('__')[-1] for x in train_df['text'].iloc[train_index]]
26        train_score = f1_score(train_df['label'].values[train_index].astype('str'), train_pred)
27        # print('Train length: ', len(train_pred))
28        print('Train score: ', train_score)
29        train_scores.append(train_score)
30
31        val_pred = [model.predict(x)[0][0].split('__')[-1] for x in train_df['text'].iloc[test_index]]

```

```

32     val_score = f1_score(train_df['label'].values[test_index].astype(str),
33     # print('Val length: ', len(val_pred))
34     print('Val score', val_score)
35     val_scores.append(val_score)
36
37     print('mean train score: ', np.mean(train_scores))
38     print('mean val score: ', np.mean(val_scores))
39     train_time = time.time()
40     print('Train time: {:.2f}s'.format(train_time - start_time))
41
42     return models
43
44 def fasttext_kfold_predict(models, n_splits):
45
46     pred_list = []
47
48     start_time = time.time()
49     # 预测并保存
50     test_df = pd.read_csv('../input/test_a.csv')
51
52     # 消耗时间较长
53     for model in models:
54         test_pred = [model.predict(x)[0][0].split('__')[-1] for x in test_df['text']]
55         pred_list.append(test_pred)
56
57     test_pred_label = pd.DataFrame(pred_list).T.apply(lambda row: np.argmax(row), axis=1)
58     test_pred_label.columns='label'
59
60     test_pred_label.to_csv('../input/test_fastText_ridgeclassifier.csv', index=False)
61     print('Test predict saved.')
62     end_time = time.time()
63     print('Predict time:{:.2f}s'.format(end_time - start_time))
64
65
66 if __name__ == '__main__':
67     nrows = 200000
68     train_num = int(nrows * 0.7)
69     n_splits = 3
70     lr=0.1
71     wordNgrams=2

```

```
72 minCount=1
73 epoch=25
74 loss='hs'
75 dim=200
76
77 """
78 Train score: 0.9635013320936988
79 Val score 0.9086640111428032
80 Train score: 0.9623510782430645
81 Val score 0.9094998879044359
82 Train score: 0.9628121318772955
83 Val score 0.9096191534698315
84 mean train score: 0.9628881807380196
85 mean val score: 0.9092610175056901
86 Train time: 740.60s
87 """
88
89 models = fasttext_kfold_model(nrows, train_num, n_splits, lr=lr, wordNgrams=wordNgrams)
90 fasttext_kfold_predict(models, n_splits=n_splits)
```

K折交叉验证能增加模型的稳定性，尤其时间有限，验证的结果仅达0.909，有时间的朋友可以调整超参数，获得更高的准确率。



**AI Start**

## 机器学习初学者

算法讲解

论文解读

学习路线

学术技巧



长按二维码  
关注公众号