

# 【NLP】基于机器学习的文本分类！

机器学习初学者 8月10日

以下文章来源于Datawhale，作者李露



**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：李露，西北工业大学，Datawhale优秀学习者**

据不完全统计，网民们平均每人每周收到的垃圾邮件高达10封左右。垃圾邮件浪费网络资源的同时，还消耗了我们大量的时间。大家对此深恶痛绝，于是识别垃圾邮件并对其进行过滤成为各邮件服务商的重要工作之一。

垃圾邮件识别问题本质上是一个文本分类问题，给定文档 $p$ （可能含有标题 $t$ ），将文档分类为 $n$ 个类别中的一个或多个。文本分类一般有两种处理思路：基于机器学习的方法和基于深度学习的方法。

本文主要基于机器学习的方法，介绍了特征提取+分类模型在文本分类中的应用。具体目录如下：

## Contents

- 1 数据及背景
- ▼ 2 文本表示方法
  - 2.1 One-hot
  - 2.2 Bags of Words
  - 2.3 N-gram
  - 2.4 TF-IDF
- ▼ 3 基于机器学习的文本分类
  - 3.1 导入相关的包
  - 3.2 读取数据
  - 3.3 文本分类对比
- ▼ 4 研究参数对模型的影响
  - 4.1 正则化参数对模型的影响
  - 4.2 max\_features对模型的影响
  - 4.3 ngram\_range对模型的影响
- ▼ 5 考虑其他分类模型
  - 5.1 LogisticRegression
  - 5.2 SGDClassifier
  - 5.3 SVM

## 一、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）

## 二、文本表示方法

在机器学习算法的训练过程中，假设给定 $N$ 个样本，每个样本有 $M$ 个特征，这样就组成了 $N \times M$ 的样本矩阵。在计算机视觉中可以把图片的像素看作特征，每张图片都可以视为 $height \times width \times 3$ 的特征图，然后用一个三维矩阵带入计算。

但是在自然语言领域，上述方法却不可行，因为文本的长度是不固定的。文本分类的第一步就是将不定长的文本转换到定长的空间内，即词嵌入。

## 2.1 One-hot

One-hot方法将每一个单词使用一个离散的向量表示，将每个字/词编码成一个索引，然后根据索引进行赋值。One-hot表示法的一个例子如下：

句子1：我 爱 北 京 天 安 门

句子2：我 喜 欢 上 海

首先对句子中的所有字进行索引

{'我': 1, '爱': 2, '北': 3, '京': 4, '天': 5, '安': 6, '门': 7, '喜': 8, '欢': 9, '上': 10, '海': 11}

一共11个字，因此每个字可以转换为一个11维的稀疏向量：

我：[1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]

爱：[0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0]

...

海：[0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1]

## 2.2 Bags of Words

Bags of Words，也称为Count Vectors，每个文档的字/词可以使用其出现次数来进行表示。例如对于：

句子1：我 爱 北 京 天 安 门

句子2：我 喜 欢 上 海

直接统计每个字出现的次数，并进行赋值：

句子1：我 爱 北 京 天 安 门

转换为 [1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0]

句子2：我 喜 欢 上 海

转换为 [1, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1]

可以利用sklearn的CountVectorizer来实现这一步骤。

```
1 from sklearn.feature_extraction.text import CountVectorizer
2 corpus = [
3     'This is the first document.',
```

```
4 'This document is the second document.',
5 'And this is the third one.',
6 'Is this the first document?',
7 ]
8 vectorizer = CountVectorizer()
9 vectorizer.fit_transform(corpus).toarray()
```

输出为：

```
[[0, 1, 1, 1, 0, 0, 1, 0, 1],
 [0, 2, 0, 1, 0, 1, 1, 0, 1],
 [1, 0, 0, 1, 1, 0, 1, 1, 1],
 [0, 1, 1, 1, 0, 0, 1, 0, 1]]
```

### 2.3 N-gram

N-gram与Count Vectors类似，不过加入了相邻单词组合为新的单词，并进行计数。如果N取值为2，则句子1和句子2就变为：

句子1：我爱 爱北 北京 京天 天安 安门

句子2：我喜 喜欢 欢上 上海

### 2.4 TF-IDF

TF-IDF分数由两部分组成：第一部分是词语频率(Term Frequency)，第二部分是逆文档频率(Inverse Document Frequency)

$$TF(t) = \frac{\text{该词语在当前文档出现的次数}}{\text{当前文档中词语的总数}}$$

$$IDF(t) = \log_e \left( \frac{\text{文档总数}}{\text{出现该词语的文档总数}} \right)$$

# TF-IDF

TF-IDF is a measure of originality of a word by comparing the number of times a word appears in a doc with the number of docs the word appears in.

$$\text{TF-IDF} = \text{TF}(t, d) \times \text{IDF}(t)$$

Term frequency

Inverse document frequency

Number of times term  $t$  appears in a doc,  $d$

$$\log \frac{1 + n}{1 + \text{df}(d, t) + 1}$$

# of documents

Document frequency of the term  $t$

## 三、基于机器学习的文本分类

接下来我们将研究文本表示对算法精度的影响，对比同一分类算法在不同文本表示下的算法精度，通过本地构建验证集计算F1得分。

### 3.1 导入相关的包

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.feature_extraction.text import CountVectorizer
4 from sklearn.feature_extraction.text import TfidfVectorizer
5 from sklearn.linear_model import RidgeClassifier
6 import matplotlib.pyplot as plt
7 from sklearn.metrics import f1_score
```

### 3.2 读取数据

```
1 train_df = pd.read_csv('./data/train_set.csv', sep='\t', nrows=15000)
```

## 3.3 文本分类对比

### 3.3.1 Count Vectors + RidgeClassifier

```
1 vectorizer = CountVectorizer(max_features=3000)
2 train_test = vectorizer.fit_transform(train_df['text'])
3
4 clf = RidgeClassifier()
5 clf.fit(train_test[:10000], train_df['label'].values[:10000])
6
7 val_pred = clf.predict(train_test[10000:])
8 print(f1_score(train_df['label'].values[10000:], val_pred, average='macro'))
```

输出为  $F1 = 0.654418775812$ .

### 3.3.2 TF-IDF + RidgeClassifier

```
1 tfidf = TfidfVectorizer(ngram_range=(1,3), max_features=3000)
2 train_test = tfidf.fit_transform(train_df['text'])
3
4 clf = RidgeClassifier()
5 clf.fit(train_test[:10000], train_df['label'].values[:10000])
6
7 val_pred = clf.predict(train_test[10000:])
8 print(f1_score(train_df['label'].values[10000:], val_pred, average='macro'))
```

输出为  $F1 = 0.87193721737$ .

## 四、研究参数对模型的影响

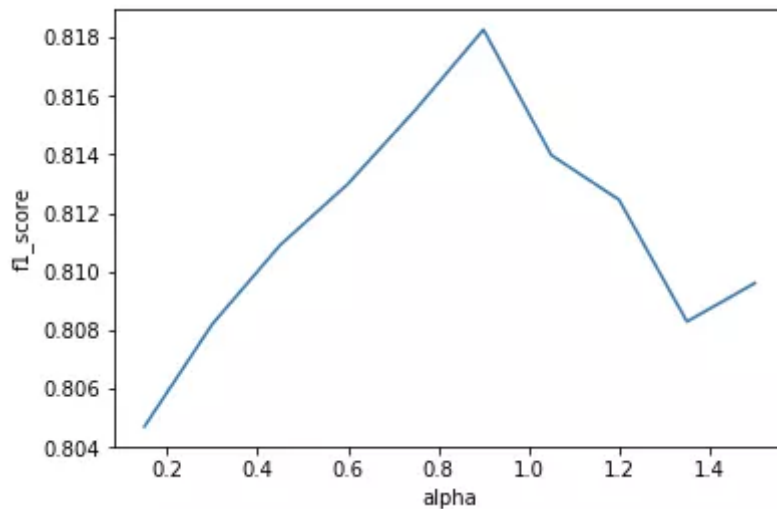
### 4.1 正则化参数对模型的影响

取大小为5000的样本，保持其他参数不变，令 $\alpha$ 从0.15增加至1.5，画出关于 $\alpha$ 和 $F1$ 的图像

```
1 sample = train_df[0:5000]
2 n = int(2*len(sample)/3)
3 tfidf = TfidfVectorizer(ngram_range=(2,3), max_features=2500)
4 train_test = tfidf.fit_transform(sample['text'])
5 train_x = train_test[:n]
6 train_y = sample['label'].values[:n]
7 test_x = train_test[n:]
8 test_y = sample['label'].values[n:]
9
10 f1 = []
11 for i in range(10):
12     clf = RidgeClassifier(alpha = 0.15*(i+1), solver = 'sag')
```

```
13 clf.fit(train_x, train_y)
14 val_pred = clf.predict(test_x)
15 f1.append(f1_score(test_y, val_pred, average='macro'))
16
17 plt.plot([0.15*(i+1) for i in range(10)], f1)
18 plt.xlabel('alpha')
19 plt.ylabel('f1_score')
20 plt.show()
```

结果如下：



可以看出 $\alpha$ 不宜取的过大，也不宜过小。 $\alpha$ 越小模型的拟合能力越强，泛化能力越弱， $\alpha$ 越大模型的拟合能力越差，泛化能力越强。

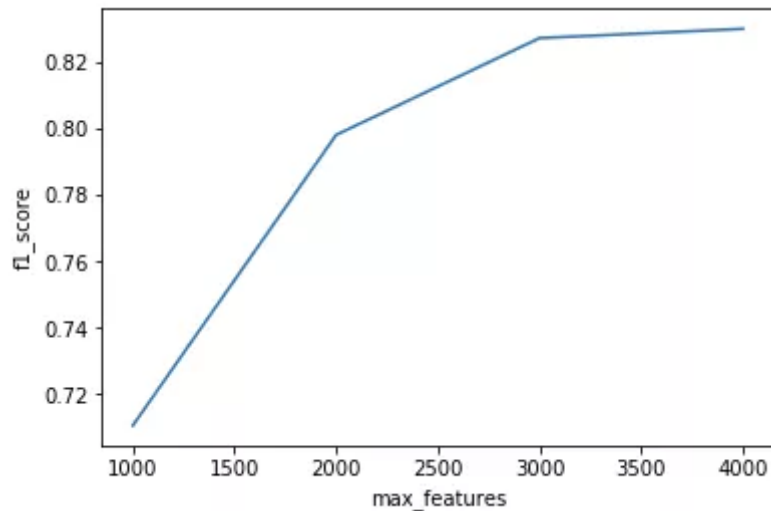
## 4.2 max\_features对模型的影响

分别取max\_features的值为1000、2000、3000、4000，研究max\_features对模型精度的影响

```
1 f1 = []
2 features = [1000,2000,3000,4000]
3 for i in range(4):
4     tfidf = TfidfVectorizer(ngram_range=(2,3), max_features=features[i])
5     train_test = tfidf.fit_transform(sample['text'])
6     train_x = train_test[:n]
7     train_y = sample['label'].values[:n]
8     test_x = train_test[n:]
9     test_y = sample['label'].values[n:]
10    clf = RidgeClassifier(alpha = 0.1*(i+1), solver = 'sag')
11    clf.fit(train_x, train_y)
12    val_pred = clf.predict(test_x)
```

```
13     f1.append(f1_score(test_y, val_pred, average='macro'))
14
15     plt.plot(features, f1)
16     plt.xlabel('max_features')
17     plt.ylabel('f1_score')
18     plt.show()
```

结果如下：



可以看出max\_features越大模型的精度越高，但是当max\_features超过某个数之后，再增加max\_features的值对模型精度的影响就不是很显著了。

### 4.3 ngram\_range对模型的影响

n-gram提取词语字符数的下边界和上边界，考虑到中文的用词习惯，ngram\_range可以在(1,4)之间选取

```
1  f1 = []
2  tfidf = TfidfVectorizer(ngram_range=(1,1), max_features=2000)
3  train_test = tfidf.fit_transform(sample['text'])
4  train_x = train_test[:n]
5  train_y = sample['label'].values[:n]
6  test_x = train_test[n:]
7  test_y = sample['label'].values[n:]
8  clf = RidgeClassifier(alpha = 0.1*(i+1), solver = 'sag')
9  clf.fit(train_x, train_y)
10 val_pred = clf.predict(test_x)
11 f1.append(f1_score(test_y, val_pred, average='macro'))
12
13 tfidf = TfidfVectorizer(ngram_range=(2,2), max_features=2000)
```

```
14 train_test = tfidf.fit_transform(sample['text'])
15 train_x = train_test[:n]
16 train_y = sample['label'].values[:n]
17 test_x = train_test[n:]
18 test_y = sample['label'].values[n:]
19 clf = RidgeClassifier(alpha = 0.1*(i+1), solver = 'sag')
20 clf.fit(train_x, train_y)
21 val_pred = clf.predict(test_x)
22 f1.append(f1_score(test_y, val_pred, average='macro'))
23
24 tfidf = TfidfVectorizer(ngram_range=(3,3), max_features=2000)
25 train_test = tfidf.fit_transform(sample['text'])
26 train_x = train_test[:n]
27 train_y = sample['label'].values[:n]
28 test_x = train_test[n:]
29 test_y = sample['label'].values[n:]
30 clf = RidgeClassifier(alpha = 0.1*(i+1), solver = 'sag')
31 clf.fit(train_x, train_y)
32 val_pred = clf.predict(test_x)
33 f1.append(f1_score(test_y, val_pred, average='macro'))
34
35 tfidf = TfidfVectorizer(ngram_range=(1,3), max_features=2000)
36 train_test = tfidf.fit_transform(sample['text'])
37 train_x = train_test[:n]
38 train_y = sample['label'].values[:n]
39 test_x = train_test[n:]
40 test_y = sample['label'].values[n:]
41 clf = RidgeClassifier(alpha = 0.1*(i+1), solver = 'sag')
42 clf.fit(train_x, train_y)
43 val_pred = clf.predict(test_x)
44 f1.append(f1_score(test_y, val_pred, average='macro'))
```

输出如下

```
[0.82713628786864823,
0.798778183959721,
0.65177193684467682,
0.84535543277498504]
```

ngram\_range取(1,3)的效果较好。

## 五、考虑其他分类模型



特征提取使用TF-IDF，与第三节中**TF-IDF + RidgeClassifier**的特征提取保持一致，再来看下其他几种分类算法的效果。

## 5.1 LogisticRegression

LogisticRegression的目标函数为：

$$\min ||\omega||_{l_i} + C \sum \log(\exp(-y_i(X_i^T \omega + c)) + 1)$$

```

1 from sklearn import linear_model
2
3 tfidf = TfidfVectorizer(ngram_range=(1,3), max_features=5000)
4 train_test = tfidf.fit_transform(train_df['text']) # 词向量 15000*max_features
5
6 reg = linear_model.LogisticRegression(penalty='l2', C=1.0, solver='liblinear')
7 reg.fit(train_test[:10000], train_df['label'].values[:10000])
8
9 val_pred = reg.predict(train_test[10000:])
10 print('预测结果中各类新闻数目')
11 print(pd.Series(val_pred).value_counts())
12 print('\n F1 score为')
13 print(f1_score(train_df['label'].values[10000:], val_pred, average='macro'))

```

输出为0.846470490043。

## 5.2 SGDClassifier

SGDClassifier使用mini-batch来做梯度下降，在处理大数据的情况下收敛更快

```

1 tfidf = TfidfVectorizer(ngram_range=(1,3), max_features=5000)
2 train_test = tfidf.fit_transform(train_df['text']) # 词向量 15000*max_features
3
4 reg = linear_model.SGDClassifier(loss="log", penalty='l2', alpha=0.0001, l1_ratio=0.1)
5 reg.fit(train_test[:10000], train_df['label'].values[:10000])
6
7 val_pred = reg.predict(train_test[10000:])
8 print('预测结果中各类新闻数目')
9 print(pd.Series(val_pred).value_counts())
10 print('\n F1 score为')
11 print(f1_score(train_df['label'].values[10000:], val_pred, average='macro'))

```

输出为0.847267047346

### 5.3 SVM

```
1 from sklearn import svm
2 tfidf = TfidfVectorizer(ngram_range=(1,3), max_features=5000)
3 train_test = tfidf.fit_transform(train_df['text']) # 词向量 15000*max_features
4
5 reg = svm.SVC(C=1.0, kernel='linear', degree=3, gamma='auto', decision_function)
6 reg.fit(train_test[:10000], train_df['label'].values[:10000])
7
8 val_pred = reg.predict(train_test[10000:])
9 print('预测结果中各类新闻数目')
10 print(pd.Series(val_pred).value_counts())
11 print('\n F1 score为')
12 print(f1_score(train_df['label'].values[10000:], val_pred, average='macro'))
```

输出为0.884240695943.

对比几种机器学习算法可以看出，在相同的TF-IDF特征提取方法基础上，用SVM得到的分类效果最好。



AI Start

## 机器学习初学者

算法讲解

学习路线

论文解读

学术技巧



长按二维码  
关注公众号

往期精彩回顾

