

基于朴素贝叶斯算法的垃圾邮件分类器

梁鑫嵘 200110619

计算机科学与技术学院

基于朴素贝叶斯算法的垃圾邮件分类器

1 理论基础

1.1 贝叶斯公式

1.2 贝叶斯公式在问题中的应用

2 实践

2.1 要点问题分析

2.1.1 如何用数据表示词语

2.1.2 正确处理“0概率”情况

2.2 程序运行流程

2.2.1 数据处理

2.2.2 训练分类器

2.2.3 测试分类器

2.2.4 运行主程序

2.3 程序运行结果

3 总结

4 附录

4.1 完整程序

4.2 数据来源

1 理论基础

1.1 贝叶斯公式

贝叶斯公式刻画了先验概率和后验概率之间的关系。

设 A_1, A_2, \dots, A_n 构成完备事件组，且对于 $\forall A_i, P(A_i) > 0$ 。 B 为样本空间的任意事件， $P(B) > 0$ ，则有：¹

$$P(A_k|B) = \frac{P(A_k)P(B|A_k)}{\sum_{i=1}^n P(A_i)P(B|A_i)}$$

1.2 贝叶斯公式在问题中的应用

对于贝叶斯公式我们可以有如下的理解²：

$$P(\text{类别}|\text{特征}) = \frac{P(\text{特征}|\text{类别})P(\text{类别})}{P(\text{特征})}$$

当我们求得 $P(\text{类别}|\text{特征})$ ，再将本样本归于概率最大的那一类，实际就完成了
一次贝叶斯分类过程。

但是上述对贝叶斯公式的理解还需要加上一个不可缺少的条件，即各变量之间相互独立，以满足“ A_1, A_2, \dots, A_n 为完备事件组”这一条件。所以在实际运用中，我们还需要假设各个特征之间是独立的。

然而实际上，能够统计到的变量之间几乎不可能是互相独立的。例如，设事件 $A =$ 电子邮件中出现“http”， $B =$ 电子邮件中包含完整链接，当 A, B 都发生的时候电子邮件很可能是一封垃圾广告邮件。但是，显然 $B \subseteq A$ ，也就是说 A, B 不互相独立。我们的程序在处理的过程中很难依据内容判断词与词之间的复杂的语义、逻辑关系，即无法判断两个词的出现事件之间是否互相独立。此外，由于人类语言的复杂性，词语之间基本都存在互相关联的情况，无法直接地抽象为互斥事件，所以利用朴素贝叶斯算法进行基于文本的垃圾邮件判断存在一定的理论误差。

在不大规模使用的情况下，朴素贝叶斯算法仍然能在简单的邮件分类作业上表现良好。

假设 A_i 为出现词汇表中的第 i 个词， B 为此邮件为垃圾邮件，

$\because A_i, A_j$ 互相独立($i \neq j$)，则：

$$\begin{aligned}
P(B|A_1, A_2, \dots, A_n) &= \frac{P(A_1|B)P(A_2|B) \cdots P(A_n|B)P(B)}{P(A_1)P(A_2) \cdots P(A_n)} \\
&= \frac{P(B)\prod_{i=1}^n P(A_i|B)}{\prod_{i=1}^n P(A_i)}
\end{aligned}$$

2 实践

2.1 要点问题分析

2.1.1 如何用数据表示词语

在程序处理的过程中，我们只需要进行数字的计算，但是我们的日常语言是使用词语等拼接起来的，为了进行这其中的转换，使得程序能够用数字的方式正确“理解”词语，我们需要一种用数字正确表示词语的方法。

常用的方法有独热码、词向量等，这里使用独热码表示。

在独热码表示法中，需要使用的数据在向量中对应位置为1，其他数据在向量中的位置为0。对应在这个分类问题中，词汇表为训练数据中所有出现的词语，向量长度为词汇表长度，每一个词语对应的独热码为 $[0, 0, \dots, 1, \dots, 0]$ 。

2.1.2 正确处理“0概率”情况

观察贝叶斯公式，分母为 $\prod_{i=1}^n P(A_i)$ ，但是如果分母为0，计算出的概率是无意义的。为了解决这样的问题，我们需要保证分母不为0，为此我们需要使用拉普拉斯平滑法处理分母数据。

拉普拉斯平滑，又称为加1平滑，是一种常用的平滑方法，经常在解决0概率问题中有所应用³。

我们令 $\alpha > 0$ ，当 $\alpha \rightarrow 0$ 时， $\prod_{i=1}^n P(A_i) \approx \prod_{i=1}^n P(A_i) + \alpha$ 。

此时分母已经不为0，但是仍然存在一些问题，计算得到的 $P(B|A_1, A_2, \dots, A_n)$ 可能因为分母过小而过大，所以最好对其加上一个对数函数，防止其因为计算数字过大而溢出。因为使用了对数函数，要求参数 > 0 ，因此我们可以对分子 $P(B)\prod_{i=1}^n P(A_i|B)$ 也加上 α 。最终得到的公式表示如下：

$$\begin{aligned}
Q_B &= \ln(P(B|A_1, A_2, \dots, A_n)) \\
&= \ln\left(\frac{P(B)\prod_{i=1}^n P(A_i|B)}{\prod_{i=1}^n P(A_i)}\right) \\
&\approx \ln P(B) + \sum_{i=1}^n \ln\left(\frac{P(A_i|B)}{P(A_i) + \alpha_1}\right) + \alpha_1 \quad (\alpha_1, \alpha_2 > 0) \\
&\approx \ln P(B) + \sum_{i=1}^n \ln(P(A_i|B) + \alpha_1) - \sum_{i=1}^n \ln(P(A_i) + \alpha_2)
\end{aligned}$$

2.2 程序运行流程

2.2.1 数据处理

```
1 def split_text(text: str):
2     # 将特殊符号作为切分标志进行字符串切分，即非字母、非数字
3     tokens = re.split(r'\W+', text)
4     # 除了单个字母，例如大写的I，其它单词变成小写；并且排除纯数字
5     return [tok.lower() for tok in tokens if len(tok) > 2 and not
6             tok.isdigit()]
7
8 def get_words_data():
9     """
10    获取训练用文本数据，并且格式化分割到单词
11    :return: 得到的文本数据，分类列表
12    """
13    class_types = [r for r in os.listdir('email') if
14                    os.path.isdir(os.path.join('email', r))]
15
16    def read_data(filename_: str) -> str:
17        with open(filename_, 'r', encoding='gbk') as f:
18            return f.read()
19
20    words_data = []
21    for c in class_types:
22        for filename in os.listdir(os.path.join('email', c)):
23            file_data = read_data(os.path.join(f'email/{c}', filename))
24            words_data.append((c, split_text(file_data)))
25
26    return words_data, class_types
27
28 def get_words_label(words_data: list) -> list:
29     """
30     得到当前数据集下的词汇表
31     :param words_data: 读取到的词语数据
32     :return: 词汇表
33     """
34     # 使用 set 去重
```

```

34     words_label = set({})
35     for words in words_data:
36         words_label.update(words[1])
37     res = list(words_label)
38     res.sort()
39     return res
40
41
42 def get_words_vector(words, words_label: list) -> list:
43     """
44     得到一个词语或者一个词语列表对应的词向量
45     :param words: 词语或者词语列表
46     :param words_label: 词汇表
47     :return: 词向量
48     """
49     return [(1 if val == words else 0) if not isinstance(words, list)
50             else (1 if val in words else 0) for val in words_label]

```

2.2.2 训练分类器

```

1  def native_bayes_train(words_label: list, train_data: list,
2  class_types: list, alpha: float = 1e-3):
3      """
4      训练朴素贝叶斯分类器
5      :param words_label: 词汇表
6      :param train_data: 训练数据集
7      :param class_types: 分类列表
8      :param alpha: > 0 且很小的数, 用于拉普拉斯平滑
9      :return: 训练结果, 各种类的概率
10     """
11     # 初始化就加上这个 alpha, 防止出现 Nan。
12     p_result = {c: np.array([alpha for _ in range(len(words_label))])
13     for c in class_types}
14     p_b = {c: alpha for c in class_types}
15     for data in train_data:
16         # 得到词向量
17         words_vector = np.array(get_words_vector(data[1], words_label))
18         # 记录到训练数据
19         p_result[data[0]] += words_vector + alpha

```

```

18         # 记录到对应分类的概率
19         p_b[data[0]] += 1 / len(train_data)
20     for k in p_result:
21         p_result[k] = np.log(p_result[k])
22     return p_result, p_b

```

2.2.3 测试分类器

```

1 def native_bayes_test(words_label: list, p_result: dict, test_data:
2     list, p_b: dict, alpha: float = 1e-3) -> float:
3     """
4     测试朴素贝叶斯分类器
5     :param words_label: 词汇表
6     :param p_result: 训练结果
7     :param test_data: 测试数据集
8     :param p_b: 各种类的概率
9     :param alpha: > 0 且很小的数, 用于拉普拉斯平滑
10    :return: 测试正确的概率
11    """
12    accurate = 0
13    for data in test_data:
14        # 获得词向量
15        words_vector = np.array(get_words_vector(data[1], words_label))
16        # 统计相对概率
17        Qs = {key: np.log(p_b[key]) + sum(p_result[key] * words_vector)
18              - np.log(sum(words_vector + alpha)) for key in
19                p_result}
20        # 选取相对概率最高的那个作为分类结果
21        classification = reduce((lambda x, y: x if x[1] > y[1] else y),
22                                ((k, Qs[k]) for k in Qs))[0]
23        print(f"result {classification == data[0]}, classification =
24              {classification}, label = {data[0]}; Qs: {Qs}")
25        # 记录正确率
26        accurate += (1 if classification == data[0] else 0) /
27        len(test_data)
28    return accurate

```

2.2.4 运行主程序

```

1 def main():
2     # 读取训练、测试数据
3     words_data, class_types = get_words_data()
4     # 打乱训练数据
5     random.shuffle(words_data)
6     # 生成词汇表
7     words_label = get_words_label(words_data)
8     # 训练朴素贝叶斯分类器，使用随机40份邮件做训练数据集
9     p_result, p_b = native_bayes_train(words_label, words_data[:40],
class_types)
10    # 测试朴素贝叶斯分类器，使用剩下的邮件做测试数据集
11    accurate = native_bayes_test(words_label, p_result,
words_data[40:], p_b)
12    print(f"DONE! accurate = {(accurate * 100):.2f}%")

```

完整程序见附录。

2.3 程序运行结果

```

1 result True, classification = ham, label = ham; Qs: {'ham':
-59.347845111308025, 'spam': -85.7843006553319}
2 result True, classification = spam, label = spam; Qs: {'ham':
-36.34975776983966, 'spam': 5.109566958100539}
3 result True, classification = ham, label = ham; Qs: {'ham':
-92.15951999754893, 'spam': -184.30347898520958}
4 result True, classification = spam, label = spam; Qs: {'ham':
-57.07097136557958, 'spam': -53.66419688826986}
5 result True, classification = spam, label = spam; Qs: {'ham':
-90.63426772205429, 'spam': 30.369671882623937}
6 result True, classification = spam, label = spam; Qs: {'ham':
-90.63426772205429, 'spam': 30.369671882623937}
7 result True, classification = spam, label = spam; Qs: {'ham':
-28.69504741287626, 'spam': -16.289914553986566}
8 result True, classification = spam, label = spam; Qs: {'ham':
-111.28152808527354, 'spam': 32.34635005125606}
9 result True, classification = spam, label = spam; Qs: {'ham':
-105.0581659869714, 'spam': 30.966905653462643}
10 result True, classification = spam, label = spam; Qs: {'ham':
-90.63426772205429, 'spam': 30.369671882623937}
11 DONE! accurate = 100.00%

```

程序判断正确率可以达到90%以上。

3 总结

本论文基于概率论与数理统计课程中的贝叶斯分类器部分的理论知识，完成了一个基于朴素贝叶斯算法的垃圾邮件分类器，展示了概率论这一学科的丰富的运用场景。

4 附录

4.1 完整程序

```
1 import os
2 import re
3 import numpy as np
4 from functools import reduce
5 import random
6
7
8 def split_text(text: str):
9     # 将特殊符号作为切分标志进行字符串切分，即非字母、非数字
10    tokens = re.split(r'\W+', text)
11    # 除了单个字母，例如大写的I，其它单词变成小写；并且排除纯数字
12    return [tok.lower() for tok in tokens if len(tok) > 2 and not
13            tok.isdigit()]
14
15 def get_words_data():
16     """
17     获取训练用文本数据，并且格式化分割到单词
18     :return: 得到的文本数据，分类列表
19     """
20    class_types = [r for r in os.listdir('email') if
21                    os.path.isdir(os.path.join('email', r))]
22
23    def read_data(filename_: str) -> str:
24        with open(filename_, 'r', encoding='gbk') as f:
25            return f.read()
26
27    words_data = []
28    for c in class_types:
```



```

28         for filename in os.listdir(os.path.join('email', c)):
29             file_data = read_data(os.path.join(f'email/{c}',
filename))
30             words_data.append((c, split_text(file_data)))
31         return words_data, class_types
32
33
34 def get_words_label(words_data: list) -> list:
35     """
36     得到当前数据集下的词汇表
37     :param words_data: 读取到的词语数据
38     :return: 词汇表
39     """
40     # 使用 set 去重
41     words_label = set({})
42     for words in words_data:
43         words_label.update(words[1])
44     res = list(words_label)
45     res.sort()
46     return res
47
48
49 def get_words_vector(words, words_label: list) -> list:
50     """
51     得到一个词语或者一个词语列表对应的词向量
52     :param words: 词语或者词语列表
53     :param words_label: 词汇表
54     :return: 词向量
55     """
56     return [(1 if val == words else 0) if not isinstance(words, list)
57             else (1 if val in words else 0) for val in words_label]
58
59
60 def native_bayes_train(words_label: list, train_data: list,
class_types: list, alpha: float = 1e-3):
61     """
62     训练朴素贝叶斯分类器
63     :param words_label: 词汇表

```

```

64     :param train_data: 训练数据集
65     :param class_types: 分类列表
66     :param alpha: > 0 且很小的数, 用于拉普拉斯平滑
67     :return: 训练结果, 各种类的概率
68     """
69     # 初始化就加上这个 alpha, 防止出现 Nan。
70     p_result = {c: np.array([alpha for _ in range(len(words_label))])}
71     for c in class_types:
72         p_b = {c: alpha for c in class_types}
73         for data in train_data:
74             # 得到词向量
75             words_vector = np.array(get_words_vector(data[1],
76 words_label))
77             # 记录到训练数据
78             p_result[data[0]] += words_vector + alpha
79             # 记录到对应分类的概率
80             p_b[data[0]] += 1 / len(train_data)
81         for k in p_result:
82             p_result[k] = np.log(p_result[k])
83         return p_result, p_b
84
85
86 def native_bayes_test(words_label: list, p_result: dict, test_data:
87 list, p_b: dict, alpha: float = 1e-3) -> float:
88     """
89     测试朴素贝叶斯分类器
90     :param words_label: 词汇表
91     :param p_result: 训练结果
92     :param test_data: 测试数据集
93     :param p_b: 各种类的概率
94     :param alpha: > 0 且很小的数, 用于拉普拉斯平滑
95     :return: 测试正确的概率
96     """
97     accurate = 0
98     for data in test_data:
99         # 获得词向量
100         words_vector = np.array(get_words_vector(data[1],
101 words_label))

```

```

98         # 统计相对概率
99         Qs = {key: np.log(p_b[key]) + sum(p_result[key] *
words_vector) - np.log(sum(words_vector + alpha)) for key in
100             p_result}
101         # 选取相对概率最高的那个作为分类结果
102         classification = reduce((lambda x, y: x if x[1] > y[1] else
y), ((k, Qs[k]) for k in Qs))[0]
103         print(f"result {classification == data[0]}, classification =
{classification}, label = {data[0]}; Qs: {Qs}")
104         # 记录正确率
105         accurate += (1 if classification == data[0] else 0) /
len(test_data)
106         return accurate
107
108
109 def main():
110     # 读取训练、测试数据
111     words_data, class_types = get_words_data()
112     # 打乱训练数据
113     random.shuffle(words_data)
114     # 生成词汇表
115     words_label = get_words_label(words_data)
116     # 训练朴素贝叶斯分类器，使用随机40份邮件做训练数据集
117     p_result, p_b = native_bayes_train(words_label, words_data[:40],
class_types)
118     # 测试朴素贝叶斯分类器，使用剩下的邮件做测试数据集
119     accurate = native_bayes_test(words_label, p_result,
words_data[40:], p_b)
120     print(f"DONE! accurate = {(accurate * 100):.2f}%")
121
122
123 if __name__ == '__main__':
124     main()
125

```

运行目录结构：

```
1 | bayes.py
2 |─ email/
3 |   |─ham
4 |   |   *.txt
5 |   |─spam
6 |   |   *.txt
```

4.2 数据来源

https://github.com/Asia-Lee/Naive_Bayes

1. 王勇, 田波平. 概率论与数理统计[M]. 第二版. 北京: 科学出版社, 2005. [↗](#)
2. 忆臻.[EB/OL]. <https://zhuanlan.zhihu.com/p/26262151>. 2017年04月10日-2021年11月23日. [↗](#)
3. 郑万通 .[EB/OL]. <https://blog.csdn.net/zhengwantong/article/details/72403808>. 2017年05月17日-2021年11月23日. [↗](#)