

# 拉普拉斯修正的朴素贝叶斯分类器

杨启航

2022 年 4 月 13 日

## 目录

<b>1</b>	<b>实验目的</b>	<b>1</b>
<b>2</b>	<b>贝叶斯分类器</b>	<b>2</b>
2.1	最小化错误率贝叶斯最优分类器 . . . . .	2
2.2	朴素贝叶斯分类器 . . . . .	2
2.3	先验概率和条件概率的估计 . . . . .	2
2.4	拉普拉斯修正 . . . . .	3
<b>3</b>	<b>Python 实现</b>	<b>3</b>
3.1	NaiveBayes 类 . . . . .	3
3.2	测试代码 . . . . .	6
3.3	测试结果 . . . . .	8
<b>4</b>	<b>实验总结</b>	<b>12</b>
4.1	拉普拉斯修正的朴素贝叶斯分类器思想 . . . . .	12
4.2	程序总结 . . . . .	12
4.2.1	程序特性 . . . . .	12
4.2.2	待改进 . . . . .	12

## 1 实验目的

了解朴素贝叶斯分类器和拉普拉斯修正的理论，使用 Python 实现，并在西瓜数据集上进行测试。

## 2 贝叶斯分类器

### 2.1 最小化错误率贝叶斯最优分类器

设假设  $h$  错将分类为  $y_i$  的样本  $x$  分类为任意  $y_j \neq y_i$  的风险  $\lambda_{ji}$  均相等，则可以在后验概率上运用贪心思想，得出最小化分类错误率的贝叶斯最优分类器为：

$$h^*(\mathbf{x}) = \arg \max_{c \in \mathcal{Y}} P(c|\mathbf{x}) \quad (2.1)$$

### 2.2 朴素贝叶斯分类器

根据贝叶斯定理，可以由数据集估计先验概率  $P(c)$  和条件概率  $P(\mathbf{x}|c)$  以计算后验概率，则可以将式 (2.2) 改写为：

$$h^*(\mathbf{x}) = \arg \max_{c \in \mathcal{Y}} \frac{P(c)P(\mathbf{x}|c)}{P(\mathbf{x})} \quad (2.2)$$

上式中先验概率  $P(\mathbf{x}|c)$  往往因为组合爆炸问题，在有限的数据集上难以进行有效统计。

朴素贝叶斯分类器 (naive Bayes classifier) 采用了“属性条件独立性假设”，有  $P(\mathbf{x}|c) = \prod_{i=1}^n P(x_i|c)$ 。单个属性的条件概率  $P(x_i|c)$  在估计时不容易失真，而式 (2.2) 可以写为：

$$h^*(\mathbf{x}) = \arg \max_{c \in \mathcal{Y}} \frac{P(c)}{P(\mathbf{x})} \prod_{i=1}^n P(x_i|c) \quad (2.3)$$

### 2.3 先验概率和条件概率的估计

主要利用大数定理和极大似然估计，并假设连续属性服从高斯分布。

令  $D_c$  表示训练集  $D$  中第  $c$  类样本组成的结合，则有先验概率估计：

$$P(c) = \frac{|D_c|}{|D|} \quad (2.4)$$

对于离散属性而言，令  $D_{i,x_i}$  表示  $D_c$  在第  $i$  个属性上取值为  $x_i$  的样本集合，则有条件概率估计：

$$P(x_i|c) = \frac{|D_{i,x_i}|}{|D_c|} \quad (2.5)$$

对于连续属性而言，假定  $P(x_i|c) \sim \mathcal{N}(\mu_{c,i}, \sigma_{c,i}^2)$ ，其中  $\mu_{c,i}$  和  $\sigma_{c,i}^2$  分别是第  $c$  类样本在第  $i$  个属性上取值的均值和方差，则有条件概率估计：

$$P(x_i|c) = \frac{1}{\sqrt{2\pi}\sigma_{c,i}} e^{-\frac{(x_i - \mu_{c,i})^2}{2\sigma_{c,i}^2}} \quad (2.6)$$

## 2.4 拉普拉斯修正

在有限的数据集上，即使使用朴素贝叶斯方法，对于条件概率较小的属性，在估计的时候也很有可能出现  $P(x_i|c) = 0$  的情况，此时式 (2.3) 中的连乘为 0，信息失真，分类趋于极端化。若出现这种情况，则在估计时引入“拉普拉斯修正” (Laplacian correction)：

$$\hat{P}(c) = \frac{|D_c| + 1}{|D| + N} \quad (2.7)$$

$$\hat{P}(x_i|c) = \frac{|D_{i,x_i}| + 1}{|D_c| + N_i} \quad (2.8)$$

其中  $N$  表示训练集  $D$  出现的类别种数， $N_i$  表示第  $i$  个属性可能的取值数。

# 3 Python 实现

## 3.1 NaiveBayes 类

```

1 from numpy import float64
2 import pandas as pd
3 from scipy import stats as st
4
5
6 class NaiveBayes:
7     class U:
8         def __init__(self) -> None:
9             self.is_discrete = True
10            self.cnt = {}
11
12            is_discrete, cnt, = None, None

```

```

13
14     data = None
15     Y = None
16     tot = None
17     ycnt = None
18     Ni = None
19
20     def malloc( self ):
21         self.data = []
22         self.Y = []
23         self.tot = 0
24         self.ycnt = {}
25         self.Ni = []
26
27     def __init__(self, D) -> None:
28         self.malloc()
29         self.tot = len(D)
30         self.Y = set(D.iloc[:, -1])
31         Dy = {}
32         for y in self.Y:
33             Dy[y] = D.loc[D[D.keys()][-1] == y]
34             self.ycnt[y] = len(Dy[y])
35         for attr in D.keys()[:-1]:
36             self.Ni.append(len(set(D[attr])))
37             self.data.append(self.U())
38             if D[attr].dtype == float64:
39                 self.data[-1].is_discrete = False
40                 for y in self.Y:
41                     = Dy[y][attr].mean()
42                     s = Dy[y][attr].std()
43                     self.data[-1].cnt[y] = [ , s]
44             else:
45                 self.data[-1].is_discrete = True

```

```

46         for y in self.Y:
47             self.data[-1].cnt[y] = {}
48             for xi in Dy[y][attr]:
49                 self.data[-1].cnt[y][xi] = \
50                     self.data[-1].cnt[y].get(xi, 0) + 1
51
52     def predict(self, x):
53         laplace = False
54         res, maxp = None, 0.0
55         # is need laplacian correction
56         for y in self.Y:
57             for xi, i in zip(x, range(len(x))):
58                 if self.data[i].is_discrete:
59                     if self.data[i].cnt[y].get(xi, -1) == -1:
60                         laplace = True
61
62         for y in self.Y:
63             p = 1.0
64             for xi, i in zip(x, range(len(x))):
65                 if self.data[i].is_discrete:
66                     if laplace:
67                         p *= (self.data[i].cnt[y].get(xi, 0) + 1) / (
68                             self.ycnt[y] + self.Ni[i]
69                         )
70                     else:
71                         p *= self.data[i].cnt[y][xi] / self.ycnt[y]
72             else:
73                 p *= st.norm.pdf(
74                     x=xi,
75                     loc=self.data[i].cnt[y][0],
76                     scale=self.data[i].cnt[y][1]
77                 )
78

```

```

79         if laplace :
80             p *= (self.ycnt[y] + 1) / ( self.tot + len( self.Y))
81         else :
82             p *= self.ycnt[y] / self.tot
83
84         if p > maxp:
85             maxp, res = p, y
86     return res

```

封装 NaiveBayes 类，在对象化的时候需要传入数据集，此时会立刻完成统计并保存统计结果，之后只需要调用该对象的 Predict(x) 函数，传入待预测样本，返回预测结果。

NaiveBayes 对象以属性为单元储存统计信息，并封装在类内类 U 中，U 对象内记录了对应属性是否为离散属性，并对于离散属性记录  $|D_{c,x_i}|$ ，对于连续属性记录  $\mu_{c,i}$  和  $\sigma_{c,i}$ 。

预测时，先判断是否需要拉普拉斯变换，即是否有  $\exists xi \exists c |D_{c,xi}| = 0$ ，后选用合适的估计方式以计算各个分类  $c_i$  的后验概率，取后验概率最大者返回。

### 3.2 测试代码

```

1 from NaiveBayes import NaiveBayes
2 import numpy as np
3 from sklearn import datasets
4 from time import time
5 import pandas as pd
6
7
8 df = pd.read_csv("watermelon3_0_Ch.csv")
9 df.drop("编号", axis=1, inplace=True)
10 t0 = time()
11 nb = NaiveBayes(df)
12 print("统计消耗时间: {:.3f}S".format(time() - t0))
13 precision = 0.0

```

```

14 t0 = time()
15 for i in range(len(df)):
16     x = df.iloc[i, :].tolist()
17     y = nb.predict(x[:-1])
18     print("{}: 真实分类: {}, 预测结果: {}".format(i + 1, x[:-1], y))
19     if x[-1] == y:
20         precision += 1 / (len(df) + 1)
21 print("准确率: {:.3f}".format(precision))
22 print("预测消耗时间: {:.3f}S".format(time() - t0))
23
24 iris = datasets.load_iris()
25 df = pd.DataFrame(
26     data=np.c_[iris["data"], iris["target"]],
27     columns=iris["feature_names"] + ["target"]
28 )
29 print(df)
30 t0 = time()
31 nb = NaiveBayes(df)
32 print("统计消耗时间: {:.3f}S".format(time() - t0))
33 precision = 0.0
34 t0 = time()
35 for i in range(len(df)):
36     x = df.iloc[i, :].tolist()
37     y = nb.predict(x[:-1])
38     print("{}: 真实分类: {}, 预测结果: {}".format(i + 1, x[:-1], y))
39     if x[-1] == y:
40         precision += 1 / (len(df) + 1)
41 print("准确率: {:.3f}".format(precision))
42 print("预测消耗时间: {:.3f}S".format(time() - t0))

```

测试中除了西瓜数据集，还测试了 Sk-learn 库中自带的鸢尾花数据集，以检验 NaiveBayes 类处理其它任务的能力，且鸢尾花数据集更加大，相比西瓜数据集更适合探究分类器性能。

### 3.3 测试结果

```
1 统计消耗时间: 0.002S
2 1: 真实分类: 是, 预测结果: 是
3 2: 真实分类: 是, 预测结果: 是
4 3: 真实分类: 是, 预测结果: 是
5 4: 真实分类: 是, 预测结果: 是
6 5: 真实分类: 是, 预测结果: 是
7 6: 真实分类: 是, 预测结果: 是
8 7: 真实分类: 是, 预测结果: 否
9 8: 真实分类: 是, 预测结果: 是
10 9: 真实分类: 否, 预测结果: 否
11 10: 真实分类: 否, 预测结果: 否
12 11: 真实分类: 否, 预测结果: 否
13 12: 真实分类: 否, 预测结果: 否
14 13: 真实分类: 否, 预测结果: 是
15 14: 真实分类: 否, 预测结果: 否
16 15: 真实分类: 否, 预测结果: 是
17 16: 真实分类: 否, 预测结果: 否
18 17: 真实分类: 否, 预测结果: 否
19 准确率: 0.778
20 预测消耗时间: 0.015S
21
22 统计消耗时间: 0.035S
23 1: 真实分类: 0.0, 预测结果: 0.0
24 2: 真实分类: 0.0, 预测结果: 0.0
25 3: 真实分类: 0.0, 预测结果: 0.0
26 4: 真实分类: 0.0, 预测结果: 0.0
27 5: 真实分类: 0.0, 预测结果: 0.0
28 6: 真实分类: 0.0, 预测结果: 0.0
29 7: 真实分类: 0.0, 预测结果: 0.0
30 8: 真实分类: 0.0, 预测结果: 0.0
31 9: 真实分类: 0.0, 预测结果: 0.0
32 10: 真实分类: 0.0, 预测结果: 0.0
```



```
33 11: 真实分类: 0.0, 预测结果: 0.0
34 12: 真实分类: 0.0, 预测结果: 0.0
35 13: 真实分类: 0.0, 预测结果: 0.0
36 14: 真实分类: 0.0, 预测结果: 0.0
37 15: 真实分类: 0.0, 预测结果: 0.0
38 16: 真实分类: 0.0, 预测结果: 0.0
39 17: 真实分类: 0.0, 预测结果: 0.0
40 18: 真实分类: 0.0, 预测结果: 0.0
41 19: 真实分类: 0.0, 预测结果: 0.0
42 20: 真实分类: 0.0, 预测结果: 0.0
43 21: 真实分类: 0.0, 预测结果: 0.0
44 22: 真实分类: 0.0, 预测结果: 0.0
45 23: 真实分类: 0.0, 预测结果: 0.0
46 24: 真实分类: 0.0, 预测结果: 0.0
47 25: 真实分类: 0.0, 预测结果: 0.0
48 26: 真实分类: 0.0, 预测结果: 0.0
49 27: 真实分类: 0.0, 预测结果: 0.0
50 28: 真实分类: 0.0, 预测结果: 0.0
51 29: 真实分类: 0.0, 预测结果: 0.0
52 30: 真实分类: 0.0, 预测结果: 0.0
53 31: 真实分类: 0.0, 预测结果: 0.0
54 32: 真实分类: 0.0, 预测结果: 0.0
55 33: 真实分类: 0.0, 预测结果: 0.0
56 34: 真实分类: 0.0, 预测结果: 0.0
57 35: 真实分类: 0.0, 预测结果: 0.0
58 36: 真实分类: 0.0, 预测结果: 0.0
59 37: 真实分类: 0.0, 预测结果: 0.0
60 38: 真实分类: 0.0, 预测结果: 0.0
61 39: 真实分类: 0.0, 预测结果: 0.0
62 40: 真实分类: 0.0, 预测结果: 0.0
63 41: 真实分类: 0.0, 预测结果: 0.0
64 42: 真实分类: 0.0, 预测结果: 0.0
65 43: 真实分类: 0.0, 预测结果: 0.0
```

66	44: 真实分类: 0.0, 预测结果: 0.0
67	45: 真实分类: 0.0, 预测结果: 0.0
68	46: 真实分类: 0.0, 预测结果: 0.0
69	47: 真实分类: 0.0, 预测结果: 0.0
70	48: 真实分类: 0.0, 预测结果: 0.0
71	49: 真实分类: 0.0, 预测结果: 0.0
72	50: 真实分类: 0.0, 预测结果: 0.0
73	51: 真实分类: 1.0, 预测结果: 1.0
74	52: 真实分类: 1.0, 预测结果: 1.0
75	53: 真实分类: 1.0, 预测结果: 2.0
76	54: 真实分类: 1.0, 预测结果: 1.0
77	55: 真实分类: 1.0, 预测结果: 1.0
78	56: 真实分类: 1.0, 预测结果: 1.0
79	57: 真实分类: 1.0, 预测结果: 1.0
80	58: 真实分类: 1.0, 预测结果: 1.0
81	59: 真实分类: 1.0, 预测结果: 1.0
82	60: 真实分类: 1.0, 预测结果: 1.0
83	61: 真实分类: 1.0, 预测结果: 1.0
84	62: 真实分类: 1.0, 预测结果: 1.0
85	63: 真实分类: 1.0, 预测结果: 1.0
86	64: 真实分类: 1.0, 预测结果: 1.0
87	65: 真实分类: 1.0, 预测结果: 1.0
88	66: 真实分类: 1.0, 预测结果: 1.0
89	67: 真实分类: 1.0, 预测结果: 1.0
90	68: 真实分类: 1.0, 预测结果: 1.0
91	69: 真实分类: 1.0, 预测结果: 1.0
92	70: 真实分类: 1.0, 预测结果: 1.0
93	71: 真实分类: 1.0, 预测结果: 2.0
94	72: 真实分类: 1.0, 预测结果: 1.0
95	73: 真实分类: 1.0, 预测结果: 1.0
96	74: 真实分类: 1.0, 预测结果: 1.0
97	75: 真实分类: 1.0, 预测结果: 1.0
98	76: 真实分类: 1.0, 预测结果: 1.0

99	77: 真实分类: 1.0, 预测结果: 1.0
100	78: 真实分类: 1.0, 预测结果: 2.0
101	79: 真实分类: 1.0, 预测结果: 1.0
102	80: 真实分类: 1.0, 预测结果: 1.0
103	81: 真实分类: 1.0, 预测结果: 1.0
104	82: 真实分类: 1.0, 预测结果: 1.0
105	83: 真实分类: 1.0, 预测结果: 1.0
106	84: 真实分类: 1.0, 预测结果: 1.0
107	85: 真实分类: 1.0, 预测结果: 1.0
108	86: 真实分类: 1.0, 预测结果: 1.0
109	87: 真实分类: 1.0, 预测结果: 1.0
110	88: 真实分类: 1.0, 预测结果: 1.0
111	89: 真实分类: 1.0, 预测结果: 1.0
112	90: 真实分类: 1.0, 预测结果: 1.0
113	91: 真实分类: 1.0, 预测结果: 1.0
114	92: 真实分类: 1.0, 预测结果: 1.0
115	93: 真实分类: 1.0, 预测结果: 1.0
116	94: 真实分类: 1.0, 预测结果: 1.0
117	95: 真实分类: 1.0, 预测结果: 1.0
118	96: 真实分类: 1.0, 预测结果: 1.0
119	97: 真实分类: 1.0, 预测结果: 1.0
120	98: 真实分类: 1.0, 预测结果: 1.0
121	99: 真实分类: 1.0, 预测结果: 1.0
122	100: 真实分类: 1.0, 预测结果: 1.0
123	101: 真实分类: 2.0, 预测结果: 2.0
124	102: 真实分类: 2.0, 预测结果: 2.0
125	103: 真实分类: 2.0, 预测结果: 2.0
126	104: 真实分类: 2.0, 预测结果: 2.0
127	105: 真实分类: 2.0, 预测结果: 2.0
128	143: 真实分类: 2.0, 预测结果: 2.0
129	144: 真实分类: 2.0, 预测结果: 2.0
130	145: 真实分类: 2.0, 预测结果: 2.0
131	146: 真实分类: 2.0, 预测结果: 2.0

```
132 147: 真实分类: 2.0, 预测结果: 2.0
133 148: 真实分类: 2.0, 预测结果: 2.0
134 149: 真实分类: 2.0, 预测结果: 2.0
135 150: 真实分类: 2.0, 预测结果: 2.0
136 准确率: 0.954
137 预测消耗时间: 0.269S
```

## 4 实验总结

### 4.1 拉普拉斯修正的朴素贝叶斯分类器思想

贝叶斯分类器利用了贪心思想，将后验概率  $P(c|\mathbf{x})$  最大的分类  $c$  作为分类结果，而后验概率则根据数据集估计先验概率  $P(c)$  和条件概率  $P(\mathbf{x}|c)$ ，再通过贝叶斯公式计算而来。

但是样本空间的大小随属性个数指数增长，在有限的数据集上，很难对以上两个统计量进行有效估计，朴素贝叶斯和拉普拉斯修正都是为了解决这个问题而引入的方法。

### 4.2 程序总结

#### 4.2.1 程序特性

可以处理西瓜分类之外的其它分类任务，训练集以 Pandas 库 DataFrame 格式在 NaiveBayes 类实例化的时候输入，待预测样本以 Python 列表形式输入。训练集中 float64 类型会被识别为连续属性，其它类型会被识别为离散属性。

#### 4.2.2 待改进

运行效率偏低；输入输出的普适性还不够；运算样本方差的时候如果样本数量  $|D_c| = 0$  还会出现问题，不可处理极端小的数据集。