

# 使用贝叶斯方法进行推断

2021 年 4 月 11 日

## 1 独立性和条件独立性

### 1.1 独立性

对于事件  $A$  和事件  $B$ ，我们说事件  $A$  和事件  $B$  是互相独立的，如果：

$$\Pr(A \cap B) = \Pr(A)\Pr(B). \quad (1)$$

### 1.2 条件概率

当已知事件  $A$  发生时，事件  $B$  发生的概率，称为事件  $B$  关于事件  $A$  的条件概率，记做：

$$\Pr(B|A). \quad (2)$$

容易证明，对任意事件  $A$ 、事件  $B$ ，恒有：

$$\Pr(A)\Pr(B|A) = \Pr(A \cap B) \quad (3)$$

如此一来，我们得到了一种计算条件概率的方法：

$$\Pr(B|A) = \frac{\Pr(A \cap B)}{\Pr(A)} = \frac{\Pr(A|B)\Pr(B)}{\Pr(A)} \quad (4)$$

式中，我们要求  $\Pr(A) \neq 0$ .

事件的独立性可以用条件概率来表述，具体地，我们有

**定理 1.** 设  $A, B$  是事件，则

$$\Pr(A \cap B) = \Pr(A)\Pr(B) \iff \Pr(A) = \Pr(A|B). \quad (5)$$

证明. 充分性. 设给定  $\Pr(A) = \Pr(A|B)$ . 那么

$$\Pr(A)\Pr(B) = \Pr(A|B)\Pr(B) = \Pr(A \cap B) \quad (6)$$

从而充分性成立.

必要性. 设给定  $\Pr(A \cap B) = \Pr(A)\Pr(B)$ . 那么, 假设  $\Pr(B) \neq 0$ , 则

$$\frac{\Pr(A \cap B)}{\Pr(B)} = \frac{\Pr(A)\Pr(B)}{\Pr(B)} = \Pr(A) \quad (7)$$

依式4, 我们得到

$$\Pr(A|B) = \Pr(A) \quad (8)$$

从而必要性成立.  $\square$

### 1.3 条件独立性

条件独立性是说, 当某个事件发生时, 另外两个 (或多个) 事件互相独立. 具体地, 设  $A, B, C$  是事件, 则事件  $A$ , 事件  $B$  关于事件  $C$  条件独立是指:

$$\Pr(A \cap B|C) = \Pr(A|C)\Pr(B|C) \quad (9)$$

例如: 「在吸烟的人群中, 男性患肺癌的风险和女性患肺癌的风险相当」就是一个条件独立性的陈述. 对于一个受访者, 用  $A$  表示受访者是男性,  $\neg A$  表示受访者是女性, 用  $B$  表示受访者患肺癌,  $\neg B$  表示受访者不患肺癌, 用  $C$  表示吸烟, 则上述论断可写为:

$$\Pr(A \cap B|C) = \Pr(A|C)\Pr(B|C) \quad (10)$$

我们亦有条件独立性的等价表述:

**定理 2.** 设  $A, B, C$  是事件, 则

$$\Pr(A \cap B|C) = \Pr(A|C)\Pr(B|C) \iff \Pr(A|B \cap C) = \Pr(A|C) \quad (11)$$

证明. 充分性. 设给定  $\Pr(A|B \cap C) = \Pr(A|C)$ . 注意到

$$\Pr(A|B \cap C) = \frac{\Pr(A \cap B \cap C)}{\Pr(B \cap C)} = \frac{\Pr(A \cap B \cap C)}{\Pr(B|C)\Pr(C)} \quad (12)$$

所以

$$\frac{\Pr(A \cap B \cap C)}{\Pr(B|C)\Pr(C)} = \Pr(A|C) \quad (13)$$

所以

$$\frac{\Pr(A \cap B \cap C)}{\Pr(C)} = \Pr(A|C)\Pr(B|C) \quad (14)$$

对上列等式左边应用式4, 得

$$\Pr(A \cap B|C) = \Pr(A|C)\Pr(B|C) \quad (15)$$

于是充分性得证. 从上述证明过程倒推回去可得到必要性.  $\square$

## 2 朴素贝叶斯分类器

设我们有表 1 所示的数据集:

表 1: 数据集

$\mathbf{x}_1$	$\cdots$	$\mathbf{x}_m$	$\mathbf{y}$
$x_1^{(1)}$	$\cdots$	$x_m^{(1)}$	$y_1$
$\vdots$		$\vdots$	$\vdots$
$x_1^{(N)}$	$\cdots$	$x_m^{(N)}$	$y_N$

我们希望利用表 1 所提供的数据训练得到一个分类器, 它能够对于一个新的输入  $\mathbf{x}^*$ , 判断 (预测) 出该输入所对应的标签  $y^*$ .

设标签变量  $\mathbf{y}$  的取值范围是在  $\{c_1, c_2, \cdots, c_k\}$ , 而变量  $\mathbf{x}_1, \cdots, \mathbf{x}_m$  的取值都是离散的, 那么每当得到一组输入  $\mathbf{x}^* = (x_1^*, \cdots, x_m^*)$ , 我们能够计算出后验概率:

$$\Pr(c_j|\mathbf{x}^*). \quad (16)$$

$y_*$  则是选取自

$$y_* = \arg \max_{c \in \{c_1, \cdots, c_k\}} \Pr(c|\mathbf{x}^*) \quad (17)$$

所以说现在的重点, 就是如何去计算  $\Pr(c|\mathbf{x}^*)$ .

## 2.1 条件独立性假设

首先我们继续上一节的计算过程，展开  $\Pr(c|\mathbf{x}^*)$ ：

$$\Pr(c|\mathbf{x}^*) = \frac{\Pr(c)\Pr(\mathbf{x}^*|c)}{\Pr(\mathbf{x}^*)} \quad (18)$$

$$\propto \Pr(c)\Pr(\mathbf{x}^*|c) \quad (19)$$

$$= \Pr(c)\Pr\left(\prod_{j=1}^m x_j^*|c\right) \quad (20)$$

在贝叶斯方法中，有这样一个条件独立性假设：

$$\Pr\left(\prod_{j=1}^m x_j^*|c\right) = \prod_{j=1}^m \Pr(x_j^*|c) \quad (21)$$

于是就有

$$\Pr(c|\mathbf{x}^*) \propto \Pr(c) \prod_{j=1}^m \Pr(x_j^*|c) \quad (22)$$

于是就有

$$y^* = \arg \max_{c \in \{c_1, \dots, c_k\}} \Pr(c) \prod_{j=1}^m \Pr(x_j^*|c). \quad (23)$$

这就是贝叶斯分类器的基本工作原理。

## 2.2 拉普拉斯平滑

为了避免连乘操作中出现的某个  $\Pr(x_j^*|c)$  为 0 使得整个连乘式的结果为 0，可以引入拉普拉斯平滑操作。具体地，当在计算  $\Pr(x_j^*|c)$  时，需要先找到数据集（也就是表 1）中，那些  $\mathbf{y}$  列取值为  $c$  的行，然后再在这些行里边，找出变量  $\mathbf{x}_j$  取值为  $x_j^*$  的哪些行，并且统计这些行的行数，当统计出来的行数为 0 时，就将 0 替换为 1，这就是拉普拉斯平滑操作。

# 3 在西瓜数据集上进行推断

我们从互联网上收集来了西瓜数据集（表 2）：

表 2: 西瓜数据集

id	color	root	knock	texture	navel	touch	quality
1	青绿	蜷缩	浊响	清晰	凹陷	硬滑	是
2	乌黑	蜷缩	沉闷	清晰	凹陷	硬滑	是
3	乌黑	蜷缩	浊响	清晰	凹陷	硬滑	是
4	青绿	蜷缩	沉闷	清晰	凹陷	硬滑	是
5	浅白	蜷缩	浊响	清晰	凹陷	硬滑	是
6	青绿	稍蜷	浊响	清晰	稍凹	软粘	是
7	乌黑	稍蜷	浊响	稍糊	稍凹	软粘	是
8	乌黑	稍蜷	浊响	清晰	稍凹	硬滑	是
9	乌黑	稍蜷	沉闷	稍糊	稍凹	硬滑	否
10	青绿	硬挺	清脆	清晰	平坦	软粘	否
11	浅白	硬挺	清脆	模糊	平坦	硬滑	否
12	浅白	蜷缩	浊响	模糊	平坦	软粘	否
13	青绿	稍蜷	浊响	稍糊	凹陷	硬滑	否
14	浅白	稍蜷	沉闷	稍糊	凹陷	硬滑	否
15	乌黑	稍蜷	浊响	清晰	稍凹	软粘	否
16	浅白	蜷缩	浊响	模糊	平坦	硬滑	否
17	青绿	蜷缩	沉闷	稍糊	稍凹	硬滑	否

id 那一列是每个西瓜的唯一编号，不用做变量，quality 那列则是目标变量（也就是标签），第 2 列到第 7 列都是自变量。有了这些数据，再加上我们刚刚学过的朴素贝叶斯方法，使我们能够实现一个模型，这个模型能够根据西瓜的外表（颜色、根蒂、敲声、纹理等）来预判一个西瓜是否是好瓜。

相比于支持向量机、随机森林、决策树和神经网络，朴素贝叶斯分类器的训练非常简单——仅仅是将数据原样保存就可以了，计算工作主要是在推断是进行的。我们下面基于 Pandas 和 Numpy 实现一个简单的朴素贝叶斯分类器，它能够接受离散取值的自变量和离散取值的因变量。它的初始化过程是这样的：读入一个 `pd.DataFrame` 并将其保存，以第 1 列到倒数第 2 列作为自变量，最后 1 列作为因变量。

```
1 import pandas as pd
2 from typing import Union
```

```
3 from typing import Any
4 import numpy as np
5
6 watermelons = pd.read_csv('watermelon.csv')
7 watermelons.iloc[:, 0] = None
8
9 del watermelons['id']
10 del watermelons['density']
11 del watermelons['sugar']
12
13 # 朴素贝叶斯分类器
14 class NaiveBayesianClassifier:
15
16     # 目标变量默认为最后 1 列
17     # 第 1 列到倒数第 2 列是自变量
18     def __init__(self, data: pd.DataFrame):
19         self.data = data
20
21     # target 为目标变量取值
22     # 返回先验概率 Pr(target)
23     def prior_prob(self, target: Union[str, int]) -> float:
24         # 总变量个数
25         n_var = self.data.shape[1]
26
27         # 标签为 target 的记录个数
28         n_target = (self.data.iloc[:, n_var-1] == target).sum()
29
30         # 样本量
31         n_samples = self.data.shape[0]
32
33         # 先验概率
34         pr_target = n_target / n_samples
35
36         return pr_target
37
38     # target 为目标变量的取值
39     # input_variables 为输入自变量向量
40     # 输出为 Pr(target | input_variable) * Pr(input_variable)
41     def probability(
```

```
42         self,
43         target: Any,
44         input_variables: np.array
45     ) -> float:
46
47         # 先验概率
48         pr_target = self.prior_prob(target)
49
50         return pr_target
51
52     # 根据条件独立性计算条件概率 Pr(input_variables | target)
53     def cond_prob(
54         self,
55         target: Any,
56         input_variables: np.array
57     ) -> float:
58         col_idx = 0
59         n_cols = self.data.shape[1]
60         prob = 1
61         while col_idx <= n_cols-2:
62             prob = prob * self.single_cond_prob(
63                 col_idx,
64                 input_variables[col_idx],
65                 target
66             )
67             col_idx = col_idx + 1
68
69         return prob
70
71     # 计算单个变量的条件概率 Pr(input_variable | target)
72     def single_cond_prob(
73         self,
74         col_idx: int,
75         input_variable: Any,
76         target: Any
77     ) -> float:
78
79         n_cols = self.data.shape[1]
80         last_col = self.data.iloc[:, n_cols-1]
```

```
81         selector = last_col == target
82         subset = self.data.loc[selector, :].iloc[:, col_idx]
83
84         match_selector = subset == input_variable
85         n_matches = match_selector.sum() or 1
86         subset_size = subset.shape[0]
87
88         return n_matches / subset_size
89
90     # 根据输入预测所属类别
91     def predict(self, input_variables: np.array) -> Any:
92
93         n_cols = self.data.shape[1]
94         unique_classes = self.data.iloc[:, n_cols-1].unique()
95         odds = []
96         for c in unique_classes:
97             odds.append(
98                 self.cond_prob(c, input_variables)
99             )
100         odds_np = np.array(odds)
101         max_idx = odds_np.argmax()
102         return unique_classes[max_idx]
103
104     f = NaiveBayesianClassifier(watermelons)
105
106     row_idx = 0
107     n_cols = f.data.shape[1]
108     input_variables = f.data.iloc[row_idx, range(0, n_cols-1)]
109
110     print(f.predict(input_variables))
```