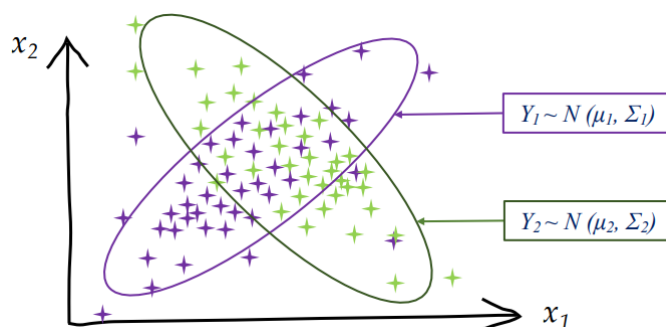


贝叶斯分类

决策树存在的问题

决策树无法解决以下情况（无法进行线性的递归的划分），但是明显两类数据存在一定的概率分布特征。



Bayes Rules

$$P(C_i|x) = \frac{p(x|C_i)P(C_i)}{p(x)}, \quad i = 1, 2$$

$P(C_i)$: **prior probability** of C_i (before observing x)

$p(C_i|x)$: **posterior probability** of C_i (after observing x)

$p(x|C_i)$: probability of x given C_i (**likelihood**)

$p(x)$: probability that x will be observed (**evidence**)

$$p(x) = P(C_1)p(x|C_1) + P(C_2)p(x|C_2)$$

$$P(\text{error}|x) = \begin{cases} P(C_1|x) & \text{if we decide } C_2 \\ P(C_2|x) & \text{if we decide } C_1 \end{cases}$$

- the **average** probability of error:

$$P(\text{error}) = \int_{-\infty}^{\infty} P(\text{error}|x)p(x)dx$$

This is minimized if for every x we ensure that $P(\text{error}|x)$ is as small as possible

- the decision rule:

Classify x into C_1 if $P(C_1|x) > P(C_2|x)$

$$P(\text{error}|x) = \min(P(C_1|x), P(C_2|x))$$

- equivalently**, classify x into C_1 if

$$\frac{p(x|C_1)p(C_1)}{p(x)} > \frac{p(x|C_2)p(C_2)}{p(x)}$$

$$p(x|C_1)p(C_1) > p(x|C_2)p(C_2)$$

Special Cases

对于 $p(x|C_1)P(C_1) > p(x|C_2)P(C_2)$

有两种特殊情况：

- $p(x|C_1) = p(x|C_2)$:
决策完全依赖于先验概率 $P(C_1)$ 和 $P(C_2)$
- $P(C_1) = P(C_2)$:
决策完全依赖于条件概率 $p(x|C_i)$

General Cases

$$\begin{aligned}P(C_i | \mathbf{x}) &= \frac{p(\mathbf{x} | C_i)P(C_i)}{p(\mathbf{x})} \\&= \frac{p(\mathbf{x} | C_i)P(C_i)}{\sum_{k=1}^K p(\mathbf{x} | C_k)P(C_k)}\end{aligned}$$

Choose C_i if $P(C_i | \mathbf{x}) = \max_k P(C_k | \mathbf{x})$

Bayes Network

Input:

- a data sample $\mathbf{x} = (x_1, x_2, \dots, x_d)$
- a fixed set of classes $C = \{C_1, \dots, C_j\}$.

Output:

- the most probable class $c \in C$:

$$\begin{aligned}c_{\text{MAP}} &= \arg \max_{c \in C} P(c | \mathbf{x}) \\&= \arg \max_{c \in C} \frac{P(\mathbf{x} | c)P(c)}{P(\mathbf{x})} \\&= \arg \max_{c \in C} P(\mathbf{x} | c)P(c) \\&= \arg \max_{c \in C} p(x_1, x_2, \dots, x_d | c) P(c)\end{aligned}$$

对于 $p(c)$ 比较容易求得，只需通过计数即可。但是对于 $p(x_1, x_2, \dots, x_d | c)$ 而言，就会非常复杂。

朴素贝叶斯

朴素贝叶斯假设 x_1, x_2, \dots, x_d 之间是独立的，即 $p(x_1, x_2, \dots, x_d | c) = p(x_1 | c)p(x_2 | c) \dots p(x_d | c)$

这样复杂的计算就得到了化简：

From:

$$c_{\text{MAP}} = \arg \max_{c \in C} p(x_1, x_2, \dots, x_d | c)P(c)$$

To:

$$c_{\text{NB}} = \arg \max_{c \in C} P(c) \prod_{i=1}^d p(x_i | c)$$

$p(x_i | c)$ 可以通过以下方式估算：

$$\hat{P}(x_i | c) \leftarrow \frac{\text{count}(x_i, c)}{\sum_{x \in |x|} \text{count}(x, c)}$$

处理Zero Counts问题

注意到，对于

$$c_{\text{NB}} = \arg \max_{c \in C} P(c) \prod_{i=1}^d p(x_i | c)$$

只要有一项 $p(x_i | c) = 0$ ，那么整体就会等于0。注意到，这是因为我们在通过计数的方式估算 $p(x_i | c)$ ，所以很容易出现数据集中在给定 c 的情况下， x_i 出现的次数为零的情况。

所以我们可以使用**Laplace Smoothing**:

$$\hat{P}(x_i | c) \leftarrow \frac{\text{count}(x_i, c) + 1}{\sum_{x \in |x|} (\text{count}(x, c) + 1)}$$

本质上就是在之前的式子上下同时加1做了一个近似，但是这保证了 $p(x_i|c) \neq 0$ 恒成立。

整体流程

Naive.Bayes.Learn(examples)

```
begin
  for each class c do
     $\hat{p}(c) \leftarrow \text{estimate } p(c)$ 
    for each attribute value  $x_i$  of each attribute  $x$  do
       $\hat{p}(x_i|c) \leftarrow \text{estimate } p(x_i|c)$ 
    end
  end
end
```

Classify.New.Instance(x)

```
begin
   $c_{NB} = \arg \max_{c \in C} P(c) \prod_{i=1}^d p(x_i|c)$ 
end
```

Example

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

$$\begin{aligned}
 P(\text{PlayTennis} = y) &= 9/14 & P(\text{PlayTennis} = n) &= 5/14 \\
 P(\text{Outlook} = \text{sunny}|y) &= 2/9 & P(\text{Outlook} = \text{sunny}|n) &= 3/5 \\
 P(\text{Outlook} = \text{overcast}|y) &= 4/9 & P(\text{Outlook} = \text{overcast}|n) &= 0/5 \\
 P(\text{Outlook} = \text{rain}|y) &= 3/9 & P(\text{Outlook} = \text{rain}|n) &= 2/5 \\
 P(\text{Temp} = \text{hot}|y) &= 2/9 & P(\text{Temp} = \text{hot}|\text{PlayTennis} = n) &= 2/5 \\
 P(\text{Temp} = \text{mild}|y) &= 4/9 & P(\text{Temp} = \text{mild}|n) &= 2/5 \\
 P(\text{Temp} = \text{cool}|y) &= 3/9 & P(\text{Temp} = \text{cool}|n) &= 1/5 \\
 P(\text{Humidity} = \text{high}|y) &= 3/9 & P(\text{Humidity} = \text{normal}|n) &= 1/5 \\
 P(\text{Humidity} = \text{normal}|y) &= 6/9 & P(\text{Humidity} = \text{high}|n) &= 4/5 \\
 P(\text{Wind} = \text{strong}|y) &= 3/9 & P(\text{Wind} = \text{strong}|n) &= 3/5 \\
 P(\text{Wind} = \text{weak}|y) &= 6/9 & P(\text{Wind} = \text{weak}|n) &= 2/5
 \end{aligned}$$

New instance : $\langle \text{sunny}, \text{cool}, \text{high}, \text{strong} \rangle$

$$\begin{aligned}
 P(y)P(\text{sunny}|y)P(\text{cool}|y)P(\text{high}|y)P(\text{strong}|y) &= .005 \\
 P(n)P(\text{sunny}|n)P(\text{cool}|n)P(\text{high}|n)P(\text{strong}|n) &= .021 \\
 \rightarrow v_{NB} &= n
 \end{aligned}$$

处理连续值

把连续值拆分成长度相等的区间即可。

Pros and Cons

- 优点
 - 训练和测试都非常快。训练的计算量很小。
 - 当独立的假设成立的时候，NB的表现会更好；在多类型预测的表现上也不错；
 - 易于维护，对于删除和添加数据集的数据的情况比较好处理（重新计数即可）。

- 缺点
 - 朴素贝叶斯的核心在于我们假设属性 x_i 之间是相互独立的，但往往事实并非如此，显然两个属性之间可能存在一定联系，这会带来一定误差；
 - 朴素贝叶斯分类器的复杂度是固定的而且比较低，所以可能会导致欠拟合（为了解决这种方法可以重新引入贝叶斯网络）；

应用场景

- 实时预测（因为NB的运行效率很高）；
- 多类型预测；
- 文字分类，比如垃圾邮件的判定。

Example of Spam Filtering

Input:

- a training set of m hand-labeled documents $(d_1, c_1), \dots, (d_m, c_m)$ where $d_i = \{w_1, w_2, \dots, w_{|d|}\}$, and $c_i \in C = \{c_1, \dots, c_{|C|}\}$

Training:

- from training corpus, extract $Vocab$
- calculate $P(c_i)$ terms
 - for each c_j in C do
 - $docs_j \leftarrow$ all docs with class $= c_j$
 - $P(c_j) \leftarrow \frac{|docs_j|}{|total \# documents|}$
- calculate $P(w_k | c_j)$ terms
 - $Text_j \leftarrow$ single doc containing all $docs_j$
 - for each word w_k in $Vocab$ do
 - $n_k \leftarrow$ # of occurrences of w_k in $Text_j$
 - $P(w_k | c_j) \leftarrow \frac{n_k + \alpha}{n + \alpha |Vocab|}$

Test:

- for $d = \{w_1, w_2, \dots, w_{|d|}\}$
- For each $c \in C$, calculate $score(c) = p(c)p(w_1|c)p(w_2|c)\dots p(w_{|d|}|c)$
- Output c with the maximum score.

