

人工智能第七次作业

13.15 Answer the Question

After your yearly checkup, the doctor has bad news and good news. The bad news is that you tested positive for a serious disease and that the test is 99% accurate (i.e., the probability of testing positive when you do have the disease is 0.99, as is the probability of testing negative when you don't have the disease). The good news is that this is a rare disease, striking only 1 in 10,000 people of your age. Why is it good news that the disease is rare? What are the chances that you actually have the disease?

解：设 A 为是否患有该病，B 为是否测试呈阳性，计算确实患有这种病的概率为

$$\begin{aligned}P(A = 1|B = 1) &= P(A = 1 \wedge B = 1)/P(B = 1) \\&= (P(B = 1|A = 1)P(A = 1))/(P(B = 1|A = 1)P(A = 1) + P(B = 1|A = 0)P(A = 0)) \\&= (99\% \times \frac{1}{10000})/(99\% \times \frac{1}{10000} + (1 - 99\%) \times (1 - \frac{1}{10000})) \\&= \frac{1}{102} \approx 0.98\%\end{aligned}$$

也就是说因为这个“好消息”，实际上患病概率只有0.98%，所以它的确是个好消息。

13.18 Answer the Question

Suppose you are given a bag containing n unbiased coins. You are told that $n - 1$ of these coins are normal, with heads on one side and tails on the other, whereas one coin is a fake, with heads on both sides.

- **a** Suppose you reach into the bag, pick out a coin at random, flip it, and get a head. What is the (conditional) probability that the coin you chose is the fake coin?
- **b** Suppose you continue flipping the coin for a total of k times after picking it and see k heads. Now what is the conditional probability that you picked the fake coin?
- **c** Suppose you wanted to decide whether the chosen coin was fake by flipping it k times. The decision procedure returns *fake* if all k flips come up heads; otherwise it returns *normal*. What is the (unconditional) probability that this procedure makes an error?

解：分小题解答如下

- **小题a**：设 A 为是否拿到假币，B 为是否为正面

$$\begin{aligned}P(A = 1|B = 1) &= P(A = 1 \wedge B = 1)/P(B = 1) \\&= (P(B = 1|A = 1)P(A = 1))/(P(B = 1|A = 1)P(A = 1) + P(B = 1|A = 0)P(A = 0)) \\&= (1 \times \frac{1}{n})/(1 \times \frac{1}{n} + 0.5 \times (1 - \frac{1}{n})) \\&= \frac{2}{n + 1}\end{aligned}$$

- **小题b**：设 A 为是否拿到假币，B 为是否连续 k 次正面

$$\begin{aligned}
P(A=1|B=1) &= P(A=1 \wedge B=1)/P(B=1) \\
&= (P(B=1|A=1)P(A=1))/(P(B=1|A=1)P(A=1) + P(B=1|A=0)P(A=0)) \\
&= (1 \times \frac{1}{n})/(1 \times \frac{1}{n} + \frac{1}{2^k} \times (1 - \frac{1}{n})) \\
&= \frac{2^k}{n + 2^k - 1}
\end{aligned}$$

- **小题c:** 设 A 为是否拿到假币, B为是否连续 k 次正面

$$\begin{aligned}
P(A=1 \wedge B=0) + P(A=0 \wedge B=1) &= 0 + P(A=0 \wedge B=1) \\
&= P(B=1|A=0)P(A=0) \\
&= \frac{1}{2^k} \times \frac{1}{n} \\
&= \frac{n-1}{n2^k}
\end{aligned}$$

13.21 (Adapted from Pearl (1988).) Answer the Question

Suppose you are a witness to a nighttime hit-and-run accident involving a taxi in Athens. All taxis in Athens are blue or green. You swear, under oath, that the taxi was blue. Extensive testing shows that, under the dim lighting conditions, discrimination between blue and green is 75% reliable.

- **a.** Is it possible to calculate the most likely color for the taxi? (*Hint:* distinguish carefully between the proposition that the taxi is blue and the proposition that it appears blue.)
- **b.** What if you know 9 out of 10 Athenian taxis are green?

解: 分小题解答如下

- **小题a:** 不可能, 除非结合生活常识。我们假设当地绿色出租车的比例为 p , 计算这个猜测下肇事车辆为绿色出租车的概率如下

$$\begin{aligned}
P(\text{肇事为蓝车} | \text{目击为蓝车}) &= P(\text{肇事为蓝车} \wedge \text{目击为蓝车}) / P(\text{目击为蓝车}) \\
&= \frac{P(\text{目击为蓝车} | \text{肇事为蓝车}) \times P(\text{肇事为蓝车})}{P(\text{目击为蓝车} | \text{肇事为蓝车}) \times P(\text{肇事为蓝车}) + P(\text{目击为蓝车} | \text{肇事非蓝车}) \times P(\text{肇事非蓝车})} \\
&= \frac{3/4 \times p}{3/4 \times p + (1 - 3/4) \times (1 - p)} \\
&= \frac{3p}{2p + 1}
\end{aligned}$$

随着 p 值在 $(0, 1)$ 变动这个概率可以在 $(0, 1)$ 间变动, 所以无法判断

- **小题b:** 代入小题 a 的结论, 得到肇事车辆为蓝车的概率是 $0.3/1.2 = 25\%$, 不是蓝车的概率是 75% 。

13.22 Answer the Question

Text categorization is the task of assigning a given document to one of a fixed set of categories on the basis of the text it contains. Naive Bayes models are often used for this task. In these models, the query variable is the document category, and the “effect” variables are the presence or absence of each word in the language; the assumption is that words occur independently in documents, with frequencies determined by the document category.

- **a.** Explain precisely how such a model can be constructed, given as “training data” a set of documents that have been assigned to categories.
- **b.** Explain precisely how to categorize a new document.

- **c.** Is the conditional independence assumption reasonable? Discuss.

解：分小题解答如下

- **小题a:** 对每一个分好类的文档的类别，统计其各个单词出现的频率，可以构建出模型的先验概率 $P(Category = A)$ ：描述各文档中类别 A 所占的比例； $P(Word_i | Category = A)$ ：描述类别 A 的文档里单词 $Word_i$ 出现的频率
- **小题b:** 对每一个给定的新文档，统计其中各单词的出现频率。根据贝叶斯公式算出其归属于各个 $Category$ 的概率，将其归类到概率最大的一类中
- **小题c:** 不合理，单词有可能具有前后的关联，导致其出现概率不能简单看作其组成部分的概率的乘积。