

Hw 7

PB21111653 李宇哲

T1 (13.15)

Question:

13.15 在一年一度的体检之后，医生告诉你一些坏消息和一些好消息。坏消息是你在一严重疾病的测试中结果呈阳性，而这个测试的准确度为 99%（即当你确实患这种病时，测试结果为阳性的概率为 0.99；而当你未患这种疾病时测试结果为阴性的概率也是 0.99）。好消息是，这是一种罕见的病，在你这个年龄段大约 10 000 人中才有 1 例。为什么“这种病很罕见”对于你而言是一个好消息？你确实患有这种病的概率是多少？

Answer:

$$P(\text{test}|\text{disease}) = 0.99$$

$$P(\neg \text{test}|\neg \text{disease}) = 0.99$$

$$P(\text{disease}) = 0.0001$$

因为

$$P(\text{disease}|\text{test}) = \frac{P(\text{test}|\text{disease})P(\text{disease})}{P(\text{test}|\text{disease})P(\text{disease}) + P(\text{test}|\neg \text{disease})P(\neg \text{disease})}$$

这种病很罕见说明 $P(\text{disease})$ 很小，即分子小，所以得病的概率 $P(\text{disease}|\text{test})$ 比较小

$$P(\text{disease}|\text{test}) = \frac{0.99 \times 0.0001}{0.99 \times 0.0001 + 0.01 \times 0.9999} = 0.009804$$

T2 (13.18)

Question:

- 13.18** 假设给你一只袋子，装有 n 个无偏差的硬币，并且告诉你其中 $n - 1$ 个硬币是正常的，一面是正面而另一面是反面。不过剩余 1 枚硬币是伪造的，它的两面都是正面。
- 假设你把手伸进口袋均匀随机地取出一枚硬币，把它抛出去，硬币落地后正面朝上。那么你取出伪币的（条件）概率是多少？
 - 假设你不停地抛这枚硬币，一共抛了 k 次，而且看到 k 次正面向上。那么你取出伪币的条件概率是多少？
 - 假设你希望通过把取出的硬币抛 k 次的方法来确定它是不是伪造的。如果抛 k 次后都是正面朝上，那么决策过程返回 *fake*（伪造），否则返回 *normal*（正常）。这个过程发生错误的（无条件）概率是多少？

Answer:

a.

$$P(\text{Fake}|\text{heads}) = \frac{2}{2+n-1} = \frac{2}{n+1}$$

b.

$$P(\text{Fake}|\text{heads}) = \frac{2^k}{2^k+n-1}$$

c.

$$P(\text{heads}|\neg \text{fake})P(\neg \text{fake}) = \frac{n-1}{2^n}$$

T3 (13.21)

Question:

- 13.21** (改编自 Pearl (1988) 的著述。) 假设你是雅典一次夜间出租车肇事逃逸的交通事故的目击者。雅典所有的出租车都是蓝色或者绿色的。而你发誓所看见的肇事出租车是蓝色的。大量测试表明, 在昏暗的灯光条件下, 区分蓝色和绿色的可靠度为 75%。
- 有可能据此计算出肇事出租车最可能是什么颜色吗? (提示: 请仔细区分命题“肇事车是蓝色的”和命题“肇事车看起来是蓝色的”。)
 - 如果你知道雅典的出租车 10 辆中有 9 辆是绿色的呢?

Answer:

a.

设 肇事车是蓝色的 为事件B, 肇事者看起来是蓝色的为事件A

$$P(A|B) = 0.75, P(\neg A|\neg B) = 0.75$$

不可能根据目前信息计算出肇事出租车最可能是什么颜色的, 需要一个先验概率, 即雅典的出租车蓝色或者绿色的比例

b.

$$P(B|A) = \frac{0.075}{0.075+0.225} = 0.25$$

$$P(\neg B|A) = \frac{0.225}{0.075+0.225} = 0.75$$

所以更有可能是绿色的

T4 (13.22)

Question:

- 13.22** 文本分类是基于文本内容将给定的一个文档分类成固定的几个类中的一类。朴素贝叶斯模型经常用于这个问题。在朴素贝叶斯模型中, 查询 (query) 变量是这个文档的类别, 而结果 (effect) 变量是语言中每个单词的存在与否; 假设文档中单词的出现是独立的, 单词的出现频率由文档类别决定。
- 给定一组已经被分类的文档, 准确解释如何构造这样的模型。
 - 准确解释如何分类一个新文档。
 - 题目中的条件独立性假设合理吗? 请讨论。

Answer:

a.

这样的模型需要

- 先验概率 $P(Category)$
- 条件概率 $P(Word_i|Category)$

对于每个领域 (category), $P(Category = c)$ 表示属于领域c的所有文档的概率

$P(Word_i = true|Category = c)$ 为包含单词i的类别c文档的概率

b.

判断任意一个给定的单词是否出现在某个文档的分类中, 选择概率最高的那一个类别

c.

不合理

某一个特定的单词, 比如 computer science 出现的概率 可能比 computer 和 science 这两个单词概率乘积高

T5 (14.12)

Question:

14.12 两个来自世界上不同地方的宇航员同时用他们自己的望远镜观测了太空中某个小区域内恒星的数目 N 。他们的测量结果分别为 M_1 和 M_2 。通常，测量中会有不超过 1 颗恒星的误差，发生错误的概率 e 很小。每台望远镜可能出现（出现的概率 f 更小一些）对焦不准确的情况（分别记作 F_1 和 F_2 ），在这种情况下科学家会少数三颗甚至更多的恒星（或者说，当 N 小于 3 时，连一颗恒星都观测不到）。考虑图 14.22 所示的三种贝叶斯网络结构。

- 这三种网络结构哪些是对上述信息的正确（但不一定高效）表示？
- 哪一种网络结构是最好的？请解释。
- 当 $N \in \{1, 2, 3\}$, $M_1 \in \{0, 1, 2, 3, 4\}$ 时，请写出 $P(M_1 | N)$ 的条件概率表。概率分布表里的每个条目都应该表达为参数 e 和/或 f 的一个函数。
- 假设 $M_1 = 1$, $M_2 = 3$ 。如果我们假设 N 取值上没有先验概率约束，可能的恒星数目是多少？
- 在这些观测结果下，最可能的恒星数目是多少？解释如何计算这个数目，或者，如果不可能计算，请解释还需要什么附加信息以及它将如何影响结果。

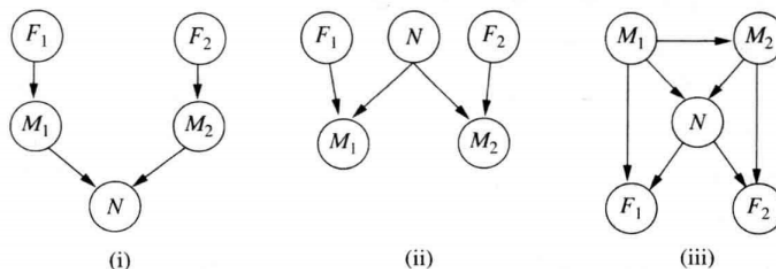


图 14.22 望远镜问题的三种可能网络

Answer:

a.

(ii) 和 (iii) 是对上述信息的正确表示

b.

(ii) 是最好的网络结构

c.

	N = 1	N = 2	N = 3
$M_1 = 0$	$f + e(1-f)$	f	f
$M_1 = 1$	$(1-2e)(1-f)$	$e(1-f)$	0.0
$M_1 = 2$	$e(1-f)$	$(1-2e)(1-f)$	$e(1-f)$
$M_1 = 3$	0.0	$e(1-f)$	$(1-2e)(1-f)$
$M_1 = 4$	0.0	0.0	$e(1-f)$

d.

$N = 2, 4$, 或者 $N \geq 6$

e.

最可能的 恒星数目 $N = 2$

不可能计算，需要知道先验分布 $P(N)$

T6 (14.13)

Question:

14.13 考虑图 14.22(ii)的网络, 假设两个望远镜完全相同。 $N \in \{1, 2, 3\}$, $M_1, M_2 \in \{0, 1, 2, 3, 4\}$, CPT 表和习题 14. 12 所描述的一样。使用枚举算法 (图 14. 9) 计算概率分布 $\mathbf{P}(N | M_1=2, M_2=2)$ 。

Answer:

$$P(N|M_1 = 2, M_2 = 2) = \alpha \sum_{f_1, f_2} P(f_1, f_2, N, M_1 = 2, M_2 = 2) = \alpha \sum_{f_1, f_2} P(f_1)P(f_2)P(N)P(M_1 = 2|f_1, N)P(M_2 = 2|f_2, N)$$

唯一可能的情况是 $f_1 = f_2 = false$

$$P(N|M_1 = 2, M_2 = 2) = \alpha(1 - f)(1 - f) < p_1, p_2, p_3 > < e, (1 - 2e), e > < e, (1 - 2e), e > = \alpha' < p_1 e^2, p_2 (1 - 2e)^2, p_3 e^2 >$$