

(a) [15 分]

使用置信度的一个缺点是它忽略了 $Pr(B)$ 。为什么这是一个缺点？请解释一下提升度和确信度为什么不受这个缺点的影响。

$$(a). \text{coef}(A \rightarrow B) = Pr(B|A) = \frac{S_{AB}}{S_A}$$

由于 $\text{coef}(A \rightarrow B)$ 忽略了 $Pr(B)$ ，所以当 $Pr(B)$ 取一些极大值，如 1 时，B 在每个篮子中都出现。此时。

$\text{coef}(A \rightarrow B)$ 也会很高。但这是由于 B 的普遍存在导致的高 coef 而非 AB 之间的关联规则导致的高 coef 。所以置信度忽略 $Pr(B)$ 是一个缺点的原因是其没有排除一些由于物品 B 的高频率而导致高置信度的物品组合。

而提升度通过在 $\text{coef}(A \rightarrow B)$ 的基础除以 S_B ，即除以 B 的频率，消除了由于 B 的高频率所带来的高 $\text{coef}(A \rightarrow B)$ 的影响。此操作 $\text{coef}(A \rightarrow B)$ 是由于高频率 B 而导致的虚高，会受到高压缩因子 S_B 的影响而降低。故提升度不受这个缺点影响。

确信度与提升度相似， $1 - S_B$ 代表 A, B 独立时 B 不出现频率。 $1 - \text{coef}(A \rightarrow B)$ 代表实际 A 发生但 B 不发生的频率。同样通过 $1 - S_B$ 这个压缩因子避免了由于高频率 B 导致确信度虚高的问题。故确信度不受这个缺点影响。

(b) [15 分]

如果 度量($A \rightarrow B$) = 度量($B \rightarrow A$)，则称该度量是对称的。在前文介绍的几种度量中，有哪些是对称的？对于每个度量，请证明该度量是对称的，或者提供一个反例表明该度量不是对称的。

(b). 置信度：

$$\begin{aligned} &\text{取 } S(A)=0.5 \quad S(B)=1 \quad S(AB)=0.5 \\ &\text{则 } \text{Coef}(A \rightarrow B) = \Pr(B|A) = \frac{S(AB)}{S(A)} = 1. \\ &\text{Coef}(B \rightarrow A) = \Pr(A|B) = \frac{S(AB)}{S(B)} = 0.5 \\ &\text{Coef}(A \rightarrow B) \neq \text{Coef}(B \rightarrow A) \\ &\Rightarrow \text{置信度不对称。} \end{aligned}$$

确信度：

$$\begin{aligned} &\text{取 } S(A)=0.4 \quad S(B)=0.5 \quad S(AB)=0.3 \\ &\text{则 } \text{Conv}(A \rightarrow B) = \frac{1-S(B)}{1-\text{Coef}(A \rightarrow B)} = \frac{\frac{1}{2}}{1-\frac{3}{4}} = 2; \\ &\text{Conv}(B \rightarrow A) = \frac{1-S(A)}{1-\text{Coef}(B \rightarrow A)} = \frac{0.6}{1-0.6} = 1.5 \\ &\text{Conv}(A \rightarrow B) \neq \text{Conv}(B \rightarrow A) \\ &\Rightarrow \text{确信度不对称。} \end{aligned}$$

提升度：

$$\begin{aligned} &\text{lift}(A \rightarrow B) = \frac{\text{Conf}(A \rightarrow B)}{S(B)} = \frac{S(AB)}{S(A) \cdot S(B)} \\ &\text{lift}(B \rightarrow A) = \frac{\text{Conf}(B \rightarrow A)}{S(A)} = \frac{S(AB)}{S(A) \cdot S(B)} \\ &\Rightarrow \text{lift}(A \rightarrow B) = \text{lift}(B \rightarrow A) \end{aligned}$$

\Rightarrow 提升度是对称的。

(c) [20 分]

完美蕴含是指相关的条件概率为 1。例如， $P(B|A) = 1$ ，则称 $A \rightarrow B$ 为完美蕴含。如果一个度量对于所有完美蕴含的情况，都能达到其最大可实现值，则称该度量是可取的。这使得我们很容易识别出最佳规则。在上述度量中，哪些是可取的？您可以忽略 0/0 的情况，但不能忽略其他无穷大的情况（即，最大值也可以为无穷大）。同时，您可以通过一个例子来简单解释。

(c) 置信度.

$\text{conf}(A \rightarrow B)$ 的最大可实现值为 1

当 $\Pr(B|A) = 1$ 时

$$\text{conf}(A \rightarrow B) = \Pr(B|A) = 1$$

\Rightarrow 置信度是可取的。

置信度

$$\text{lift}(A \rightarrow B) = \frac{\text{conf}(A \rightarrow B)}{S(B)}$$

\therefore 当 $S(B) \rightarrow 0$ 时，其最大可实现值为 $+\infty$

反例：令 $S(B) = 0.5$

$$R | \Pr(B|A) = 1 \text{ 时}, \text{lift}(A \rightarrow B) = 2.$$

此时 $\text{lift}(A \rightarrow B)$ 未可取到其最大可实现值 $+\infty$.

\Rightarrow 置信度是不可取的。

准确度.

$$\text{conv}(A \rightarrow B) = \frac{1 - S(B)}{1 - \text{conf}(A \rightarrow B)}$$

\therefore 当 $1 - S(B) \rightarrow 0$ 而 $1 - \text{conf}(A \rightarrow B) \rightarrow 0$ 时。

$\text{conv}(A \rightarrow B)$ 的最大可实现值为 $+\infty$

当 $\Pr(B|A) = 1$ 时， $\text{conf}(A \rightarrow B) = 1 \Rightarrow 1 - \text{conf}(A \rightarrow B) =$

忽略 $1 - S(B) = 0$ 的情况

此时 $\text{conv}(A \rightarrow B) \rightarrow +\infty$. 取到最大可实现值。

\Rightarrow 准确度是可取的。