

به نام خدا



تحلیل داده های حجیم
تمرین سری دوم

استاد: دکتر غلامپور
دانشجو: سجاد هاشم بیکی

پاییز 1401

سوال 1)

الف:

$$conv(A \rightarrow B) = \frac{1 - S(B)}{1 - conf(A \rightarrow B)} = \frac{p(A) * p(B')}{p(A \cup B')}$$

$p(B')$ is the probability that B does not appear in a basket.

Conviction compares the probability that A appears without B if they were dependent with the actual frequency of the appearance of X without B. Unlike confidence, conviction factors in both $p(A)$ and $p(B)$ and always has a value of 1 when the relevant items are completely unrelated. In contrast to lift, conviction is a directed measure because it also uses the information of the absence of the consequent. Hence, conviction is monotone in confidence and lift.

ب:

Lift is not sensitive to rule direction:

$$lift(A \rightarrow B) = \frac{conf(A \rightarrow B)}{sup(B)} = \frac{sup(A \cup B)}{sup(B) sup(A)} = \frac{conf(B \rightarrow A)}{sup(A)} = lift(B \rightarrow A)$$

ج:

ضعف پارامتر confidence بدین صورت است که ممکن است یک rule مقدار confidence بالایی داشته باشد ولی اصلا interest نباشد. مثلا B را ایتم شیر در نظر بگیریم، احتمالا confidence بالایی به ازای A های مختلف داشته باشد چرا که، ایتم شیر پرفروش است (در اکثر سبدها است) و احتمالا مستقل از A خریداری میشود، که در این صورت confidence بالاست اما interest نیست.

اما در دیگر پارامترها مانند conviction, lift, interest این موضوع در نظر گرفته شده است. بطوریه احتمال B (سپورت B تقسیم بر تعداد سبدها) به رابطه آنها اضافه شده است.