## 2.Exercise

### (i)

The sentence is wrong because it can happen that
$D[X_1] \subseteq D[X_2]$ but $q(X_1) \geq q(X_2)$
Let $D$ be the following database
1. abc
2. ac
3. bc
4. cd
5. acd
6. ac
Let 1,3,4,5 be the transactions marked with positive sign
Let $X_1 = ab, X_2 = ac$
So $|D[X_1]| = 1, |D^+[X_1]| = 1, |D[X_2]| = 4, |D^+[X_2]| = 2, |D^+| = 4, |D| = 6$
$q(X_1) = \sqrt{1}.|\frac{|1|}{|1|} - \frac{|4|}{|6|}| = 1.(1 - \frac{4}{6}) \approx 0.33333$
$q(X_2) = \sqrt{3}.|\frac{|2|}{|4|} - \frac{|4|}{|6|}| \approx \sqrt{3}.(0.166666) \approx 0.28867$
$q(X_2) < q(X_1)$ even $D[X_1] \subset D[X_2]$
So the statement isn't true

### (ii)

The function separates the transactions in the database into two types (interesting, not interesting).
The function calculates the quality by finding the difference between the fraction $\frac{D^+[X]}{D[X]}$ and the fraction $\frac{D^+}{D}$, and multiplying the result with a number representing how many transactions contain it.
So we can say that it tries to find the proportion of the interesting transactions containing $X$ over the total number of occurrences of $X$ in $D$.
The output can get bigger if $D[X]$ is big which means that many transactions contain $X$, or if $D^+[X]$ is close to $D[X]$ because the fraction $\frac{D^+[X]}{D[X]}$ would closer to 1 and so the difference would be bigger, which means that many of $X$ occurrences are interesting.