

Q2

Rahul Atre

2023-09-26

Q2 [Association Rules Mining]

Evaluate the following candidate association rules for the British Musical Dataset introduced in the course notes.

Recall: The above five rule metrics are stated as follows:

$$Support(X \rightarrow Y) = \frac{Freq(X \cap Y)}{N} \in [0, 1]$$

$$Confidence(X \rightarrow Y) = \frac{Freq(X \cap Y)}{Freq(X)} \in [0, 1]$$

$$Interest(X \rightarrow Y) = Confidence(X \rightarrow Y) - \frac{Freq(Y)}{N} \in [-1, 1]$$

$$Lift(X \rightarrow Y) = \frac{N^2 \cdot Support(X \rightarrow Y)}{Freq(X) \cdot Freq(Y)} \in (0, N^2)$$

$$Conviction(X \rightarrow Y) = \frac{1 - Freq(Y)/N}{1 - Confidence(X \rightarrow Y)} \geq 0$$

First, we can create a function that calculates all the rule metrics, given the frequency of x, y, XAndY, and N.

```
rule_metric_func<- function(freqX, freqY, freqXandY, N) {  
  support <- paste("Support: ", freqXandY/N * 100, "%")  
  confidence <- paste("Confidence: ", freqXandY/freqX * 100, "%")  
  interest <- paste("Interest: " , freqXandY/freqX - freqY/N)  
  lift <- paste("Lift: ", (N^2 * freqXandY/N)/(freqX * freqY))  
  conviction <- paste("Conviction: ", (1 - freqY/N)/(1 - freqXandY/freqX))  
  
  output <- cat(support, confidence, interest, lift, conviction, sep = "\n")  
}
```

i) $X \rightarrow Y$

```
#Data obtained from notes  
freqX <- 3888  
freqY <- 9092  
freqXandY <- 2720
```

```
N <- 15356
```

```
rule_metric_func(freqX, freqY, freqXandY, N) #Calculating Rule Metrics
```

```
## Support: 17.7129460797083 %  
## Confidence: 69.9588477366255 %  
## Interest: 0.107507206201889  
## Lift: 1.18157508341797  
## Conviction: 1.35786645352136
```

ii) $W \rightarrow Z$

```
freqW <- 2010  
freqZ <- 6855  
freqWandZ <- 132  
N <- 15356
```

```
rule_metric_func(freqW, freqZ, freqWandZ, N) #Calculating Rule Metrics
```

```
## Support: 17.7129460797083 %  
## Confidence: 135.323383084577 %  
## Interest: 0.906828516961947  
## Lift: 3.03140170772687  
## Conviction: -1.56721875928664
```

iii) $(Y \cap W) \rightarrow X$

```
freqYandW <- 1852  
freqX <- 3888  
freqYandWandX <- 1778  
N <- 15356
```

```
rule_metric_func(freqYandW, freqX, freqYandWandX, N) #Calculating Rule Metrics
```

```
## Support: 17.7129460797083 %  
## Confidence: 146.868250539957 %  
## Interest: 1.21549157026021  
## Lift: 5.80069150023554  
## Conviction: -1.59342210613123
```

- If an individual owns a classical music album (W), then they also own a hip-hop album (Z), given that $\text{Freq}(W) = 2010$, $\text{Freq}(Z) = 6855$, and $\text{Freq}(W \cap Z) = 132$.
- If an individual owns both the Beatles' Sergeant Peppers' Lonely Hearts Club Band and a classical music album, then they were born before 1976, given that $\text{Freq}(Y \cap W) = 1852$ and $\text{Freq}(Y \cap W \cap X) = 1778$.
- Out of the 3 rules that have been established in the previous question ($X \rightarrow Y$ [course notes], $W \rightarrow Z$, and $(Y \text{ AND } W) \rightarrow X$), which do you think is more useful? Which is more surprising?

Out of the 3 rules that were established, I believe that the second rule is the most useful one. Let's assume that an entertainment company is looking to expand their business into new genres of music, but are unsure

of which ones to explore. As a business, the next logical step would be to analyze data trends in the music landscape. If the company specializes in classical music, and their findings were similar to the second rule, then they can invest capital into hiring hip-hop artists to increase the overall revenue. The company could also bundle two albums (classical and hip-hop) together for a “discounted price” but silently raise the cost of one of the two albums for a higher profit margin. They could also allow for collaboration between hip-hop and classical artists to create innovative music.

The rule that I found most surprising was the third rule, stating that if an individual owns both the Beatles’ *Sergeant Peppers’ Lonely Hearts Club Band* and a classical music album, they would be born before 1976. There seems to be a very strong correlation in the data, as less than 200 data points were lost when applying the intersection rule. Ideally, data scientists should not try to create artificial meaning from the data (i.e. Assume that the Beatles’ were popular within the older generation’s culture, therefore the total population liking Beatles’ are overwhelmingly old). However, this is quite a surprising statistic as the number of data points lost after applying the intersection is very low. Ultimately, correlation does not equal causation, as a younger individual may find classical music pleasant or enjoy listening to the Beatles’.