

Sampling

Introduction to Data Science

Team 15:

Thijs van der Knaap (s2752077)

Hatim Alsayahani (s3183696)

Sebastian Wehkamp (s2589907)

Dimitris Laskaratos (s3463702)

October 2, 2017

1 Introduction: confidence and support

1.1 Importance of lift

Given rule $X \rightarrow Y$

- **Confidence** is the number of cases in which the rule is correct relative to the number of cases in which it is applicable.
- **Support** is the percentage of transactions that contain a given itemset
- **Lift** is the ratio of the observed support to that expected if X and Y were independent.

This means that confidence does not take into account the amount of times $X \rightarrow Y$ occurs by random. With the lift ratio you can clearly see this since when the lift is one it implies that X and Y are independent of each other and no association rule can be drawn from this. If the lift is higher it means that there is at least some relationship between X and Y.

1.2 Association Rules

In order to conduct an association analysis, the first step is to identify items that are more likely to occur together.

In this case we have the items: name, age, gender, hobbies, favourite colour, income and country and more.

The item 'name' is not suitable to be included in the analysis since it is seldom associated with the rest of the items of the set.

Taking the above into consideration we can find association rules:

- between age and hobbies,
- between age and favourite colour,
- between hobbies and country,
- between income and country.

Etc.....

1.3 Wine and Cheese

Let wine=X and cheese=Y. The lift of rule

$$wine \Rightarrow cheese$$

is thus

$$Lift(X \rightarrow Y) = \frac{Support(X \cup Y)}{Support(X) \times Support(Y)}.$$

Since

$$Lift(X \rightarrow Y) = 2$$

and

$$Support(X) = 0.1$$

we get

$$0.2 = \frac{Support(X \cup Y)}{Support(Y)}$$

or

$$\frac{1}{5} = \frac{Support(X \cup Y)}{Support(Y)}.$$

That means cheese alone is 5 times more frequent in transactions than bread and cheese together.

1.4 Itemsets

Item	support
{1}	16.7%
{2}	50.0%
{3}	66.7%
{4}	50.0%
{5}	66.7%
{6}	33.3%
{7}	16.7%
{2,3}	33.3%
{2,4}	33.3%
{2,5}	16.7%
{2,6}	0%
{3,4}	50.0%
{3,5}	50.0%
{3,6}	16.7%
{4,5}	33.3%
{4,6}	16.7%
{5,6}	33.3%
{2,3,4}	33.3%
{3,4,5}	33.3%

From the above table only rules with a minimum Support of 0.3 will be used to calculate their Confidence. So we have the itemsets: $\{\{2,3\}, \{2,4\}, \{3,4\}, \{3,5\}, \{4,5\}, \{5,6\}, \{2,3,4\}, \{3,4,5\}\}$.

By applying these itemsets to the confidence equation,

$$Confidence(X \rightarrow Y) = \frac{Support(X \cup Y)}{Support(X)}$$

where X and Y can be rules on their own, we get that the association rules that produce a minimum Support of 0.3 and a minimum Confidence of 0.7 are the following:

$\{2\} \rightarrow \{3\}$
 $\{2\} \rightarrow \{4\}$
 $\{3\} \rightarrow \{4\}$
 $\{3\} \rightarrow \{5\}$
 $\{6\} \rightarrow \{5\}$
 $\{2, 3\} \rightarrow \{4\}$
 $\{2, 4\} \rightarrow \{3\}$
 $\{4, 5\} \rightarrow \{3\}$

These rules are found by considering all possible associations in a data set and consequently applying Confidence on each of those. Only the rules with

$$Confidence \geq 0.7$$

are ultimately accepted.

2 Beethoven and Iron Maiden

2.1 Reading the data

We are reading the data by using the code in the listing below.

```

1 data <- read.csv(file.choose(new = FALSE), header=FALSE, stringsAsFactors=TRUE, sep=
  "\t")
2 names(data) <- c('user', 'timestamp', 'mbid1', 'artist', 'mbid2', 'song')
3
4 keeps <- c('user', 'artist')
5 data <- data[keeps]
```

Some interesting statistics about the dataset, note that these statistics were also made using the user-profile file.

```

1          Country
2 United States :204385
3 United Kingdom: 98766
4 Canada       : 70477
5 Peru         : 57442
6 Germany      : 46712
7 Turkey       : 45335
8 (Other)      :313714
9
10         Registered
11 Mar 30, 2005: 75877
12 Mar 16, 2005: 70447
13 Nov 10, 2005: 66609
14 Feb 24, 2006: 57439
15 May 14, 2006: 51475
16 Dec 21, 2005: 45308
```

```

17 (Other)      :469676
18
19 Gender
20 :147480
21 f:289362
22 m:399989
23
24 Band
25 Kanye West   : 27115
26 Radiohead    : 6938
27 Nine Inch Nails: 6529
28 Muse         : 5459
29 ????         : 5394
30 T.I.         : 5394
31 (Other)      :780002

```

2.2 Recommendations based on Beethoven

We are going to make some suggestions based on someone liking Beethoven alot. To do this we will see which user really likes Beethoven a lot using the code below. After which we will see which other bands that user likes as well. Those bands will be our recommendations.

```

1 trans <- as(data, "transactions")
2
3 rules <- apriori(trans,
4                   parameter = list(minlen=2, supp = 0.000000000001, conf =
5                                     0.0000000000000001),
6                   appearance = list(lhs = c("artist=Ludwig Van Beethoven"), default =
7                                     "rhs")
8 )
9 inspect(rules)

```

This results in the listing below

	lhs	lift	count	rhs	support	confidence
[1]	{artist=Ludwig Van Beethoven}	=>	{user=user_000020}	1.425033e-05	0.13333333	
	12.0263639	6				
[2]	{artist=Ludwig Van Beethoven}	=>	{user=user_000013}	2.375054e-06	0.02222222	
	1.7801581	1				
[3]	{artist=Ludwig Van Beethoven}	=>	{user=user_000028}	2.375054e-06	0.02222222	
	1.2611553	1				
[4]	{artist=Ludwig Van Beethoven}	=>	{user=user_000003}	1.187527e-05	0.11111111	
	5.9195945	5				
[5]	{artist=Ludwig Van Beethoven}	=>	{user=user_000006}	2.375054e-06	0.02222222	
	1.0494068	1				
[6]	{artist=Ludwig Van Beethoven}	=>	{user=user_000025}	1.187527e-05	0.11111111	
	3.2699067	5				
[7]	{artist=Ludwig Van Beethoven}	=>	{user=user_000001}	4.750109e-06	0.04444444	
	1.1252569	2				
[8]	{artist=Ludwig Van Beethoven}	=>	{user=user_000017}	9.500217e-06	0.08888889	
	2.0686516	4				

```

10 [9] {artist=Ludwig Van Beethoven} => {user=user_000031} 4.750109e-06 0.04444444
    0.9080024 2
11 [10] {artist=Ludwig Van Beethoven} => {user=user_000019} 4.750109e-06 0.04444444
    0.6798308 2
12 [11] {artist=Ludwig Van Beethoven} => {user=user_000012} 2.375054e-05 0.22222222
    2.3487577 10
13 [12] {artist=Ludwig Van Beethoven} => {user=user_000026} 1.425033e-05 0.13333333
    1.2390542 6

```

User 000020 is the most interesting user since it has the highest lift value. Now we will use the apriori algorithm again to determine which artists "User 00020" also likes.

```

1 rules <- apriori(trans,
2                   parameter = list(minlen=2, supp = 0.0001, conf = 0.0000001),
3                   appearance = list(lhs = c("user=user_000020"), default = "rhs")
4                   )

```

resulting in

	lhs	confidence	lift	rhs	count	support
1						
2	[1] {user=user_000020}	0.009640103	90.1977292	{artist=Nicky Wire}	45	0.0001068774
3	[2] {user=user_000020}	0.011139674	76.8898675	{artist=The Blood Brothers}	52	0.0001235028
4	[3] {user=user_000020}	0.009640103	42.2801856	{artist=L'Arc~En~Ciel}	45	0.0001068774
5	[4] {user=user_000020}	0.018423308	68.0439010	{artist=Iamx}	86	0.0002042547
6	[5] {user=user_000020}	0.024421594	87.8849669	{artist=Moneybrother}	114	0.0002707562
7	[6] {user=user_000020}	0.014781491	36.1839728	{artist=The Dandy Warhols}	69	0.0001638787
8	[7] {user=user_000020}	0.015209940	32.6736672	{artist=The Cooper Temple Clause}	71	0.0001686289
9	[8] {user=user_000020}	0.020565553	43.5124724	{artist=Patrick Wolf}	96	0.0002280052
10	[9] {user=user_000020}	0.026563839	43.8608566	{artist=Kasabian}	124	0.0002945067
11	[10] {user=user_000020}	0.012639246	18.7382606	{artist=Manic Street Preachers}	59	0.0001401282
12	[11] {user=user_000020}	0.042416452	33.3815895	{artist=Black Rebel Motorcycle Club}	198	0.0004702608
13	[12] {user=user_000020}	0.016495287	10.7345056	{artist=The Mars Volta}	77	0.0001828792
14	[13] {user=user_000020}	0.021636675	13.4167462	{artist=Dirty Pretty Things}	101	0.0002398805
15	[14] {user=user_000020}	0.043273350	24.1643784	{artist=Mando Diao}	202	0.0004797610
16	[15] {user=user_000020}	0.017780634	7.2966974	{artist=Queens Of The Stone Age}	83	0.0001971295
17	[16] {user=user_000020}	0.013281919	5.2312995	{artist=The Killers}	62	0.0001472534
18	[17] {user=user_000020}	0.023993145	8.6713697	{artist=Dredg}	112	0.0002660061

19	[18]	{user=user_000020} => {artist=The Strokes}	0.0002873816
		0.025921165 8.5265041 121	
20	[19]	{user=user_000020} => {artist=Babyshambles}	0.0001021273
		0.009211654 2.6366433 43	
21	[20]	{user=user_000020} => {artist=The Libertines}	0.0001805041
		0.016281063 3.5684682 76	
22	[21]	{user=user_000020} => {artist=Muse}	0.0019522947
		0.176092545 20.2963409 822	
23	[22]	{user=user_000020} => {artist=Radiohead}	0.0001021273
		0.009211654 0.8743242 43	

So we can recommend for example "Nicky Wire" and "Moneybrother" since those lift values are very high. These are sort of suprising recommendations since most of the bands listed above are alternative rockbands which is quite different from Beethoven.

2.3 Make someone like Eminem

For this assignment we concatenated for every user the music that it likes into one colum.

```
1 merged <- dcast(setDT(data), user ~ rowid(user), value.var = c('artist'), fill = '')
```

We planned on running the apriori on this merged dataset, but due to limited stack size this was not possible. Each time we ran it the size surpassed 4.2 Gb, which is max available RAM at the lab computers. If it would have been possible to run this command we would search for a chain of artists that would persuade a Beethoven fan to like Eminem music. From Beethoven we would search for a relation to an artist that either directly or indirectly (by repeating this step) has a high chance of also liking Eminem.

2.4 Merge with user profiles

We merged both files using a python dataframe and export it again to a csv file named "merged.csv". This is done using the code from the listing below.

```
1 import pandas as pd
2
3 lastFM = pd.read_csv('lastFM.tsv', sep='\t')
4 lastFM.columns = ['User', 'Timestamp', 'MBIDBand', 'Band', 'MBIDSong', 'Song']
5 lastFM = lastFM[['User', 'Timestamp', 'Band', 'Song']]
6
7 profiles = pd.read_csv('userid-profile.tsv', sep='\t')
8 profiles.columns = ['User', 'Gender', 'Age', 'Country', 'Registered']
9
10 merged = pd.merge(lastFM, profiles, how='right', on='User')
11
12 merged.to_csv('merged.csv')
```

The generated CSV file could be used to for example check which bands are popular in a specific country, gender, or age. Below is example code which shows what is popular in the "United States".

```

1 library(arules)
2 library(data.table)
3
4 data <- read.csv(file.choose(new = FALSE), header=TRUE, stringsAsFactors=TRUE)
5
6 keeps <- c('User', 'Timestamp', 'Band', 'Song', 'Gender', 'Country', 'Registered')
7 data <- data[keeps]
8
9 rules <- apriori(data,
10                  parameter = list(minlen=2, supp = 0.0001, conf = 0.005),
11                  appearance = list(lhs = c("Country=United States"), default = "rhs"
12                                     )
13 )

```

Resulting the sample below. In this result you can see that for example the band "Sasha & John Digweed" is popular in the "United States".

	lhs	lift	count	rhs	support	confidence
1						
2	[1] {Country=United States}	4.0867325	1068	=> {Band=Sasha & John Digweed}	0.001276243	0.005225432
3	[2] {Country=United States}	4.0943856	1159	=> {Band=Matthew Good}	0.001384987	0.005670671
4	[3] {Country=United States}	3.9816314	1130	=> {Band=Gomez}	0.001350332	0.005528781
5	[4] {Country=United States}	3.7245283	1148	=> {Band=Spoon}	0.001371842	0.005616851
6	[5] {Country=United States}	4.0825264	1377	=> {Band=They Might Be Giants}	0.001645494	0.006737285
7	[6] {Country=United States}	3.0921403	1231	=> {Band=Broken Social Scene}	0.001471026	0.006022947
8	[7] {Country=United States}	3.5565638	1435	=> {Band=Elbow}	0.001714803	0.007021063
9	[8] {Country=United States}	3.9391833	1599	=> {Band=Band Of Horses}	0.001910780	0.007823470
10	[9] {Country=United States}	3.5316115	1732	=> {Band=Cut Copy}	0.002069713	0.008474203
11	[10] {Country=United States}	3.0609962	1641	=> {Band=The Verve}	0.001960969	0.008028965
12	[11] {Country=United States}	2.3050702	1238	=> {Band=Pixies}	0.001479391	0.006057196
13	[12] {Country=United States}	3.5089650	1972	=> {Band=Boards Of Canada}	0.002356509	0.009648458
14	[13] {Country=United States}	3.1612705	1843	=> {Band=Blur}	0.002202356	0.009017296
15	[14] {Country=United States}	4.0704318	2379	=> {Band=Sasha}	0.002842868	0.011639797

2.5 Give Iron Maiden CD to your friends

You could create a profile using Data Mining of a person who is very likely to like it. You could use data mining to determine artists who users like that also like Iron Maiden and some artists which

they are sure not to like. You can ask your friends to grade these bands and see who is a perfect match for Iron Maiden. You can also combine this with factual data like country, gender, and age to do some extra matching.