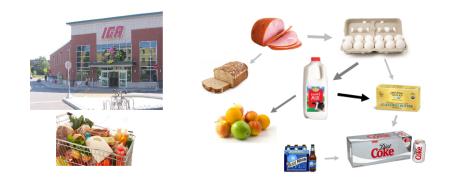
RSI

1.13. Application: Association Rules, i.e. Market Basket Analysis

Q

9

Market Basket Analyses are a common application of association rules. One goal of a market basket analysis is to understand the association between items purchases. The relationship between items purchased at a grocery store will be considered in this handout.



An association rule highlights the fact that some items are more (or less) indicative of the purchase of others. For example, purchasing cereal increases the likelihood of purchasing milk. These types of analyses may also reveals that liquor and milk are rarely purchased together.

Rule	Item	#	# Milk +	Confidence	Lift
{kitchen utensil} -> {Milk}	kitchen utensil	4	3	0.750	2.935
{honey} -> {Milk}	honey	15	11	0.733	2.870
{cereals} -> {Milk}	cereals	56	36	0.643	2.516
{rice} -> {Milk}	rice	75	46	0.613	2.400
{rubbing alcohol} -> {Milk}	rubbing alcohol	10	6	0.600	2.348
{cocoa drinks} -> {Milk}	cocoa drinks	22	13	0.591	2.313
{pudding powder} -> {Milk}	pudding powder	23	13	0.565	2.212
{jam} -> {Milk}	jam	53	29	0.547	2.142
{baking powder} -> {Milk}	baking powder	174	91	0.523	2.047
{cooking chocolate} -> {Milk}	cooking chocolate	25	13	0.520	2.035
{preservation products} -> {Milk}	preservation products	2	1	0.500	1.957
{baby cosmetics} -> {Milk}	baby cosmetics	6	3	0.500	1.957
{butter} -> {Milk}	butter	545	271	0.497	1.946
{candy} -> {Milk}	candy	294	81	0.276	1.078
{hair spray} -> {Milk}	hair spray	11	3	0.273	1.067
{seasonal products} -> {Milk}	seasonal products	140	37	0.264	1.034
{specialty chocolate} -> {Milk}	specialty chocolate	299	79	0.264	1.034
{photo/film} -> {Milk}	photo/film	91	23	0.253	0.989
{frozen fruits} -> {Milk}	frozen fruits	12	3	0.250	0.978
{canned beer} -> {Milk}	canned beer	764	87	0.114	0.446
{liquor} -> {Milk}	liquor	109	6	0.055	0.215
{baby food} -> {Milk}	baby food	1	0	0.000	0.000

Association Rules are used to uncover associations or relationships that exist between items. Often these rules are constructed to identify relationships between items purchased, i.e. Market Basket Analysis.

Procedural Steps

- 1. Determine how often items are purchased
- 2. Determine how often items are purchased in conjunction with other items
- 3. Identify which purchased items are indicative of others being purchased

Data Technologies

- 1. Filtering/Subsets
- 2. Creating Tables
- 3. Applications of Summaries

Consider the following subset of data from a collection of transactions from a grocery store.

Transaction ID	Items Purchased
1	{Bread, Milk}
2	{Eggs, Ham}
3	{Bread, Fruit, Milk}
4	{Beer, Bread, Butter, Fruit, Soda}
5	{Bread, Fruit, Milk, Soda}

Association rules are developed under the following guiding principles.

Association Rule Principles

1. Items should be purchased somewhat (**Support**)

2. Reliability, i.e. the degree to which one set of items predicts the purchase of another set of items (**Confidence**)

Consider the following association rule – the purchase of Bread indicates the purchase of Milk.

Rule 1

 $\{\text{Bread}\} \rightarrow \{\text{Milk}\}$

Compute the support and confidence for this rule.

 $Support(Bread AND Milk) = rac{\# Bread AND Milk}{\# Transactions} =$ Support(Bread AND Milk)

 $Confidence \ of \ Rule \ \#1 = \ rac{Support(Bread \ AND \ Milk)}{Support(Bread)} =$

Consider a second association rule for the purchase of Milk.

Rule 2

 $\{\mathrm{Fruit}\} \ o \{\mathrm{Milk}\}$

Compute the support and confidence for this rule.

Support(Fruit AND Milk) =

 $Confidence \ of \ Rule =$

Question

1. Why might Rule #1 be considered "better" than Rule #2 when interest lies in the purchase of

Milk? Milk? Consider a third association rule for the purchase of Milk. Rule 3 $\{Bread, Fruit\} \rightarrow \{Milk\}$ Compute the support and confidence for this rule.

 $Support(Bread, \ Fruit, \ AND \ Milk) =$

$Confidence\ of\ Rule =$

Lift is another measure often considered when evaluating rules of association.

$Lift(\{Bread\} \rightarrow \{Milk\}) =$	Confidence(Bread AND Milk)	P(Milk Bread)
$Lifi(\{Dreau\} \rightarrow \{Wink\}) =$	Support(Milk)	= $P(Milk)$

For our example, realize that the support for Milk is fairly large. i.e, Milk was purchased in 60% of the transactions. This provides a baseline value for confidence. That is, rules that exceed this value indicate gains when considering the association provided by the rule. When the lift of a rule is near 1, then the rule provides little information to understanding the purchase of the item.

• Lift>1 implies positive association between items

• Liftpprox 1 implies no association between items

• Lift < 1 implies negative association between items

Rule	Support	Confidence	Lift
$\{\text{Bread}\} \ \rightarrow \{\text{Milk}\}$	$\frac{3}{5}$	$rac{rac{3}{5}}{rac{4}{5}}=rac{3}{4}$	$rac{rac{3}{4}}{rac{3}{5}} = 1.25$
$\{ \mathrm{Fruit} \} \ o \{ \mathrm{Milk} \}$	$\frac{2}{5}$	$rac{rac{2}{5}}{rac{3}{5}}=rac{2}{3}$	$\frac{rac{2}{3}}{rac{3}{5}} = 1.11$
$\{Bread,\ Fruit\}\ o \{ ext{Milk}\}$	$\frac{2}{5}$	$rac{rac{2}{5}}{rac{3}{5}}=rac{2}{3}$	$rac{rac{2}{3}}{rac{3}{5}} = 1.11$

Some Comments

- Association rules with no support have zero confidence. E.g. Beer is never purchased with Milk, so the rule $\{Beer\} \rightarrow \{Milk\}$ should not be considered.
- The confidence of a rule should not be considered independent of it's support. For example, the rule $\{Eggs\} \rightarrow \{Ham\}$ has Confidence = 1. That is, 100% of the time eggs were purchased, so was Ham. However, this rule has very low support as Eggs and Ham were only purchased once.
- Association rules are not invariant. For example, the confidence for the rule $\left(\mathbb{N}^{(1)}\right)$
- ${Bread} \rightarrow {Milk}$ is different than the confidence of the rule ${Milk} \rightarrow {Bread}$.

1.13.1. Common Data Structure

	List			В	inary	Repr	esent	tatior	n (Ma	trix)	
Transaction ID	Items Purchased		ID	Beer	Bread	Butter	Eggs	Fruit	Ham	Milk	Soda
1	{Bread, Milk}	\rightarrow	1	0	1	0	0	0	0	1	0
2	{Eggs, Ham}		2	0	0	0	1	0	1	0	0
3	{Bread, Fruit, Milk}		3	0	1	0	0	1	0	1	0
4	{Beer, Bread, Butter, Fruit, Soda}		4	1	1	1	0	1	0	0	1
5	{Bread, Fruit, Milk, Soda}		-		-	1	0	-	0	0	1
L	1		5	0	1	0	0	1	0	1	1

Next, consider the complete grocery dataset. This dataset contains 9835 transactions and 169 unique items. This dataset can be downloaded from the Workshop website.

Data Source						
Address	http://www.statsclass.org					
Description	Groceries Dataset Michael Hahsler, Kurt Hornik, and Thomas Reutterer (2006) Implications of probabilistic data modeling for mining association rules. In M. Spiliopoulou, R. Kruse, C. Borgelt, A. Nuernberger, and W. Gaul, editors, From Data and Information Analysis to Knowledge Engineering, Studies in Classification, Data Analysis, and Knowledge Organization, pages 598–605. Springer- Verlag.					

Open the Groceries dataset in Excel. The binary representation of this market basket dataset is provided in this Excel file. A snippet is shown here.

	A	в	С	D	E	F	G	н	1	J	K	L	M	N	0	P	Q	R
1	frankfurter	sausag	liver loaf	ham	meat	finished products	organic sausag	chicker	turkey	pork	beef	hamburger meat	fish	citrus fruit	tropical fruit	pip fruit	grapes	berries
2	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	(
3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	(
4	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	(
5	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	(
6	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	(
7	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	(
8	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	(
9	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	(
10	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	(
11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	(
12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	(
13	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	0	0	(
14	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	(
15	1	Π	Π	Π	Π	n	Π	Π	Π	Π	Π	Π	Π	Π	n	Π	Π	ſ

Spreadsheets consist of rows and columns. Datasets also consist of rows and columns as well, but also contain information that is not data, e.g. variable names. Excel does not differentiate the header row from actual data unless you convert the collection of rows and columns into a **Table**.

Putting Data into the Table structure in Excel



Give your table a name for easy referencing



Short-cuts for cursor

n	movement in Excel							
	Short-cuts in Excel							
	Ctrl Home	Upper Left Corner						
	Ctrl End	Lower Right Corner						
	Ctrl ←	Move to Left edge						
	Ctrl →	Move to Right edge						
	Ctrl ↑	Move to top						
	Ctrl ↓	Move to bottom						

The following snippet shows the Groceries dataset specified as a table.

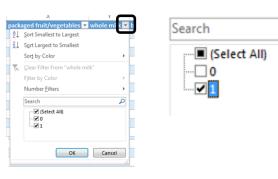
	Δ	В		C	D	F	F	G	н	Т	1	к
1		-	e 🔻	liver loaf 💌	-	-	finished products 💌	-		turkey 🔻	pork 🔻	
2		0	0	0	0	0	0	0	0	0	0	(
3		0	0	0	0	0	0	0	0	0	0	(
4		0	0	0	0	0	0	0	0	0	0	(
5		0	0	0	0	0	0	0	0	0	0	(
6		0	0	0	0	0	0	0	0	0	0	(
7		0	0	0	0	0	0	0	0	0	0	
8		0	0	0	0	0	0	0	0	0	0	
9		0	0	0	0	0	0	0	0	0	0	
10		0	0	0	0	0	0	0	0	0	0	
11		0	0	0	0	0	0	0	0	0	0	(
12		0	0	0	0	0	0	0	0	0	0	(
13		0	0	0	0	0	0	0	0	0	0	(
14		0	0	0	0	0	0	0	0	0	0	
15		1	0	0	0	0	0	0	0	0	0	(

The drop-down arrows provided for each variable (or field) are called Filters. Filters in Excel allow you to subset rows.

Filter on Whole Milk Select V

Select Whole Milk = 1 to identify transactions that purchased whole milk

Q



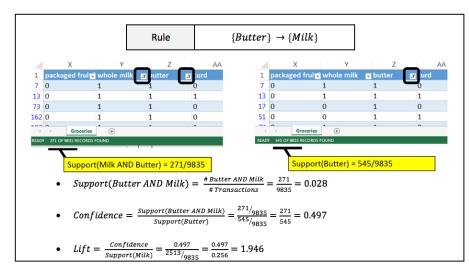
After a Filter is applied, certain rows are hidden from view. Excel indicates this fact with changing the row label color to blue.

		W	Х	Y
	1	other vegetables 💌	packaged fruit/vegetables 💌	whole milk 耳 k
Transactions for	4	0	0	1
Transactions for	6	1	0	1
which	7	0	0	1
whole milk = 1	11	0	0	1
whole mik = 1	13	0	0	1
	24	0	0	1
	34	1	0	1

The status bar in Excel, the bar across the bottom of the Excel file, provides simple summaries for columns of the table. For example, if the Whole Milk column (column Y) is highlighted, the following summaries are shown.

Groceries (+)	READY	2513 OF 9835 RECORDS FOUND	 OUNT: 2514 SUM: 2513
		Groceries (+)	E 4

Applying filters to columns Whole Milk and Butter allows one to easy compute the support and confidence for the rule $\{Butter\} \rightarrow \{Milk\}$.



The **=COUNT()** function in Excel can used to count the number of nonblank rows in a column. Excel functions also work with tables and variable names. The following will provide a count of the number of transactions in the Groceries dataset, i.e. 9835. The use of the table and variable names is preferred as this avoids the need to highlight an exact range of cells in Excel.

Counting the	e number of rows in table
Using Range	=COUNT(Y:Y)
Using Table	=COUNT(Groceries[whole milk])

The **=**COUNTIF() function in Excel provides a count of only the cells that satisfy some condition. The following can be used to compute the support for butter.

=COUNTIF(Groceries[butter] , 1)

If more than one condition is needed, the =COUNTIFS() function can be used. COUNTIFS is necessary to compute Support(Butter AND Whole Milk).

=COUNTIFS(Groceries[butter] , 1 , Groceries[whole milk] , 1)

A brief description of the **COUNTIFS** function in Excel is provide here.

_	1 st Condition		2 nd Condition	
=COUNTIFS(C	Groceries[butter],1	, Groceries[whole milk]	,1)
	Which variable?	Condition?	Which variable?	Condition?

Move to far right of the Groceries table in Excel. You can use Ctrl \rightarrow to move quickly to the far right edge. Enter the following function in Excel to compute the counts necessary for measuring support for the rule $\{Butter\} \rightarrow \{Milk\}.$

	FL	FM	FN	FO	FP
1	shop	i bags 🛛 🖵			
2	0	0			Automating Counting in Excel
3	0	0		Number of Transactions	=COUNT(Groceries[whole milk])
4	0	0		# (Butter)	=COUNTIF(Groceries[butter],1)
5	0	0			
6	0	0		# (Butter AND Whole Milk)	=COUNTIFS(Groceries[butter],1,Groceries[whole milk],1)
7	0	0			
8	0	0			

Use the value computed above to compute the Confidence and Lift for this rule.

	FL	FM	FN	FO	FP
1	shop	i bags 🖵			
2	0	0			Automating Counting in Excel
3	0	0		Number of Transactions	=COUNT(Groceries[whole milk])
4	0	0		# (Butter)	=COUNTIF(Groceries[butter],1)
5	0	0			
6	0	0		# (Butter AND Whole Milk)	=COUNTIFS(Groceries[butter],1,Groceries[whole milk],1)
7	0	0			
8	0	0		Confidence	= FP6 / FP4
9	0	0			
10	0	0		# (Milk)	=COUNTIF(Groceries[whole milk],1)
11	0	0		Lift	= FP8 / (FP10 / FP3)
10	0	0			

Verify that these formulas are correct by comparing them to the output provided below.

FL FM FN FO FP

1	si∓ I	v 60		
2	0 0		Automating Counting in E	xcel
3	0	0	Number of Transactions	<mark>98</mark> 35
4	0	0	# (Butter)	545
5	0	0		
6	0	0	# (Butter AND Whole Milk)	271
7	0	0		
8	0	0	Confidence	0.497
9	0	0		
10	0	0	# (Milk)	2513
11	0	0	Lift	1.946

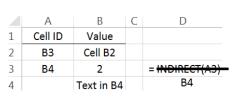
1.13.2. Evaluating Several Rules

The procedure provided above lack efficiencies and does not scale well when several rules need to be evaluated. For example, to evaluate the rule $\{Yogurt\} \rightarrow \{Milk\}$, the formulas for support will need to be changed. The =INDIRCT() function in Excel will help increase the efficiency in computing the support, confidence, and lift for several rules.

INDIRECT() Function

			01				
	А	В	Step) 1: Obta	ain value	tro	om another c
L	Cell ID	Value		А	В	С	D
t	B3	Cell B2	1	Cell ID	Value		
			2	B3	Cell B2		
;							
	B4	2	3	(B4)	2		= INDIRECT(A3)

Step 2: Use value in specified cell in formula



The **=INDIRECT()** function can be used in the following manner to automatically update the variable names when computing the support for several rules.

=COUNTIF(INDIRECT (" Groceries[" & B2 & "] "), 1)

The following setup is used to evaluate six different association rules for Milk.

	A	В	C
1	Rule	Item	#
2	{butter} -> {Milk}	butter	=COUNTIF(INDIRECT("Groceries["&B2&"]"),1)
3	{yogurt} -> {Milk}	yogurt	
4	{whipped/sour cream} -> {Milk}	whipped/sour cream	
5	{cereals} -> {Milk}	cereals	
6	{canned beer} -> {Milk}	canned beer	
7	{make up remover} -> {Milk}	make up remover	

This formula can be copied down in Excel to evaluate the support for the remaining rules. The confidence and lift are computed for these rules as well.

	А	В	С	D	E	F	G
1	Rule	Item	#	# Milk +	Confidence		Lift
2	{butter} -> {Milk}	butter	545	271	0.497		1.946
3	{yogurt} -> {Milk}	yogurt	1372	551	0.402		1.572
4	{whipped/sour cream} -> {Milk}	whipped/sour cream	705	317	0.450		1.760
5	{cereals} -> {Milk}	cereals	56	36	0.643		2.516
6	{canned beer} -> {Milk}	canned beer	764	87	0.114		0.446
7	{make up remover} -> {Milk}	make up remover	8	2	0.250		0.978

Questions

1. The Lift for $\{Cereal\} \rightarrow \{Milk\}$ is about 2.5 which is fairly high. Thus, given that the transaction includes cereal, there is 2.5 fold increase in the likelihood of milk being

purchased.

a. Compute Support(Cereal AND Milk).

- b. This value is fairly low. Why does a low support value negate the usefulness of a rule? 2. The Lift value for the rule $\{Canned \setminus Beer\} \rightarrow \{Milk\}$ is lowest on this list. What can be said about the purchase of Canned Beer AND Milk?
- 3. Which of these rules is least useful in the prediction of Milk? Explain how you made this determination.

Task

Use Excel to obtain the Confidence and Lift for all association rules for Whole Milk where only single items are considered on the left.

- Copy all variable names and paste them into a single column. This can be done using Paste Special specify Values and Transpose when pasting.
- The =CONCATENATE() function can be used to create the Rule column, i.e.
- =CONCATENATE("{",B2,"} -> {Milk}").

Specify Values and Transpose under Paste Special

frankfurter		
sausage	Paste Special	
liver loaf		
ham	Paste	
meat	○ AII	All using So
finished products	© <u>F</u> ormulas	All except bo
organic sausage	Values	Column wid
chicken	Formats	Formulas an
	Comments	O Values and r
turkey	Validation	All merging
pork	Operation	
beef	None	Multiply
hamburger meat	© Add	Divide
fish	Subtract	Opinde
citrus fruit	O Sabriaci	
tropical fruit	Skip <u>b</u> lanks	Transpose
pip fruit		
grapes	Paste Link	ОК
herries		

Output for Rules

Rule	Item	#	# Milk +	Confidence	Lift
{kitchen utensil} -> {Milk}	kitchen utensil	4	3	0.750	2,935
			-		
{honey} -> {Milk}	honey	15	11	0.733	2.870
{cereals} -> {Milk}	cereals	56	36	0.643	2.516
{rice} -> {Milk}	rice	75	46	0.613	2.400
{rubbing alcohol} -> {Milk}	rubbing alcohol	10	6	0.600	2.348
{cocoa drinks} -> {Milk}	cocoa drinks	22	13	0.591	2.313
{pudding powder} -> {Milk}	pudding powder	23	13	0.565	2.212
{jam} -> {Milk}	jam	53	29	0.547	2.142
{baking powder} -> {Milk}	baking powder	174	91	0.523	2.047
{cooking chocolate} -> {Milk}	cooking chocolate	25	13	0.520	2.035
{preservation products} -> {Milk}	preservation products	2	1	0.500	1.957
{baby cosmetics} -> {Milk}	baby cosmetics	6	3	0.500	1.957
{butter} -> {Milk}	butter	545	271	0.497	1.946
{candy} -> {Milk}	candy	294	81	0.276	1.078
{hair spray} -> {Milk}	hair spray	11	3	0.273	1.067
{seasonal products} -> {Milk}	seasonal products	140	37	0.264	1.034
{specialty chocolate} -> {Milk}	specialty chocolate	299	79	0.264	1.034
{photo/film} -> {Milk}	photo/film	91	23	0.253	0.989
{frozen fruits} -> {Milk}	frozen fruits	12	3	0.250	0.978
{canned beer} -> {Milk}	canned beer	764	87	0.114	0.446
{liquor} -> {Milk}	liquor	109	6	0.055	0.215
{baby food} -> {Milk}	baby food	1	0	0.000	0.000