

Model Performance and Recommendation Systems

Sridhar Seshadri

Overview

I

Performance Evaluation of Classifiers
Recommendation Systems

Model Performance: Bias vs Variance

I

Use all the available information in dataset to build prediction model.

However, using too much information from the current data -- may not be a good idea as this may lead to poor prediction outside the training data.

Model Performance: Bias vs Variance

I

Thus, to ensure that we have not fit our model to the training data, it's a good idea to test our model performance (prediction accuracy) on a validation dataset.

Therefore, we do not use the validation dataset in fitting the model, but we do use it to choose the final model.

Model Performance: Test

I

To evaluate the performance of our final model (fitted on training data and selected on the basis of performance on the validation data), we may also keep a holdout sample (test data).

Model Performance: Test

I

Data is divided into three parts: train, validation, and test.

We must test our model accuracy, error rates, etc. on the validation set. (Rattle warns!)

Model Performance: Cross-Validate

I

Reducing the problem of overfitting (high variance) through cross-validation

Partition the dataset into, say, ten subsets.

Any nine of these ten form the training set and the last one is the hold-out set to measure model performance.

Model Performance: Cross-Validate

I

This process can be repeated ten times.

The model performance across these ten repetitions can be averaged to get expected model performance.

Metrics to Evaluate Model Performance

I

Error rate

True and False Positives and Negatives

Precision, Recall, Sensitivity, Specificity

Confusion matrix

Risk charts

Lift chart

Scoring

Classification Tree – Spam filter

Needs DAAG package

I

Objective – To classify emails as spam vs non-spam.

crl.tot total length of words in capitals

dollar number of occurrences of the \\$ symbol

bang number of occurrences of the ! symbol

money number of occurrences of the word ‘money’

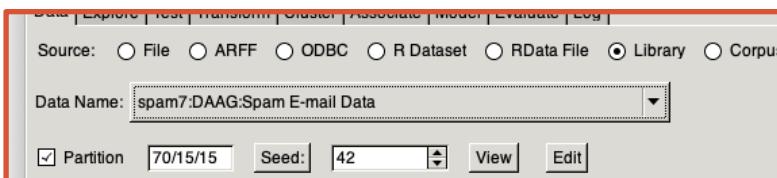
n000 number of occurrences of the string ‘000’

make number of occurrences of the word ‘make’

yesno outcome variable, a factor with levels n not spam, y spam

Source- DAAG package library R
Original source -

<http://archive.ics.uci.edu/ml/machine-learning-databases/spambase/>



	crl.tot	dollar	bang	money	n000	make	yesno
278	0.000	0.778	0.00	0.00	0.00		y
1028	0.180	0.372	0.43	0.43	0.43	0.21	y
2259	0.184	0.276	0.06	1.16	0.06		y
191	0.000	0.137	0.00	0.00	0.00		y
191	0.000	0.135	0.00	0.00	0.00		y
54	0.000	0.000	0.00	0.00	0.00		y

Model Evaluation – Spam Filter

I

We will build a Spam Filter based on Classification Tree and evaluate its performance based on the following metrics

Error rate

True and False Positives and Negatives

Classification: Scoring – to predict the class label

Precision, Recall, Sensitivity, Specificity

Confusion matrix

Risk chart

Lift chart

Output slightly changed in R 3.6.3

Classification Tree

I

Default values

R Data Miner - [Rattle (s)]

Project Tools Settings Help

Execute New Open Save Export Stop Quit

Data Explore Test Transform Cluster Associate Model Evaluate Log

Source: File ARFF ODBC R Dataset RData File Library Corpus Script

Filename: spam.csv Separator: , Decimal: . Header

Partition 70/15/15 Seed: 42 View Edit

Input Ignore Weight Calculator:
Target Data Type: Auto Categorical Numeric Survival

No. Variable	Data Type	Input	Target	Risk	Ident	Ignore	Weight	Comment
1 X	Ident	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 4,601
2 crl.tot	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 919
3 dollar	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 504
4 bang	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 964
5 money	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 143
6 n000	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 164
7 make	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 142
8 yesno	Categorical	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 2

Type Tree Forest Boost SVM Linear Neural Net Survival All

Target: yesno Algorithm: Traditional Conditional

Model Builder: rpart Include Missing

Min Split: 20 Max Depth: 30 Priors:
Min Bucket: 7 Complexity: 0.0100 Loss Matrix:

Summary of the Decision Tree model for Classification (built using 'rpart'):

n= 3220

node), split, n, loss, yval, (yprob)
* denotes terminal node

1) root 3220 1275 n (0.6040373 0.3959627)
2) dollar< 0.0555 2416 564 n (0.7665563 0.2334437)
 4) bang< 0.086 1663 166 n (0.9001804 0.0998196) *
 5) bang>=0.086 753 355 y (0.4714475 0.5285525)
 10) crl.tot< 85.5 383 117 n (0.6945170 0.3054830)
 20) bang< 0.825 304 68 n (0.7763158 0.2236842) *
 21) bang>=0.825 79 30 y (0.3797468 0.6202532) *
 11) crl.tot>=85.5 370 89 y (0.2405405 0.7594595) *
3) dollar>=0.0555 804 93 y (0.1156716 0.8843284) *

Classification tree:

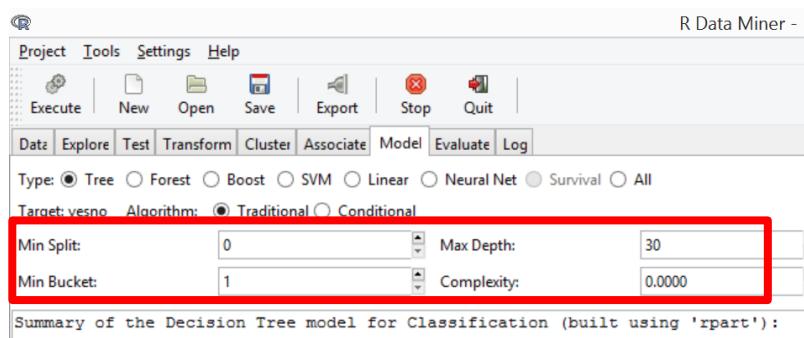
```
rpart(formula = yesno ~ ., data = crs$dataset[crs$train, c(crs$input,
  crs$target)], method = "class", model = TRUE, parms = list(split = "information"),
  control = rpart.control(usesurrogate = 0, maxsurrogate = 0))
```

Variables actually used in tree construction:
[1] bang crl.tot dollar

To see the tradeoff between the bias (under fitting) and variance (over fitting), we will first fit the model as complex as possible.

Error Rate

Error rate on training data and validation data for the most complex tree possible in Rattle



Set “min split”, “min bucket size”, and “complexity parameter” to their minimum values and “depth” to its maximum value in Rattle.

Note that complexity value lower means more complex model

Error Rate

Type: Error Matrix Risk Cost Curve Hand Lift ROC Precision Sensitivity Prv Ob Score

Model: Tree Boost Forest SVM Linear Neural Net Survival KMeans HClust

Data: Training Validation Testing Full Enter CSV File Document... R Dataset

Risk Variable:

Report: Class Probability Include: Identifiers

Error matrix for the Decision Tree model on spam.csv [**train**] (counts):

Predicted		
Actual	n	y
n	1942	3
y	100	1175
		Error
	0.2	
	7.8	

Error matrix for the Decision Tree model on spam.csv [**train**] (proportions):

Predicted		
Actual	n	y
n	60.3	0.1
y	3.1	36.5
		Error
	0.2	
	7.8	

Overall error: 3.2%, Averaged class error: 4%

Rattle timestamp: 2020-03-17 15:06:46 ashis

Type: Error Matrix Risk Cost Curve Hand Lift ROC Precision Sensitivity Prv Ob Score

Model: Tree Boost Forest SVM Linear Neural Net Survival KMeans HClust

Data: Training Validation Testing Full Enter CSV File Document... R Dataset

Risk Variable:

Report: Class Probability Include: Identifiers

Error matrix for the Decision Tree model on spam.csv [validate] (counts):

Predicted		
Actual	n	y
n	385	43
y	51	211
		Error
	10.0	
	19.5	

Error matrix for the Decision Tree model on spam.csv [validate] (proportions):

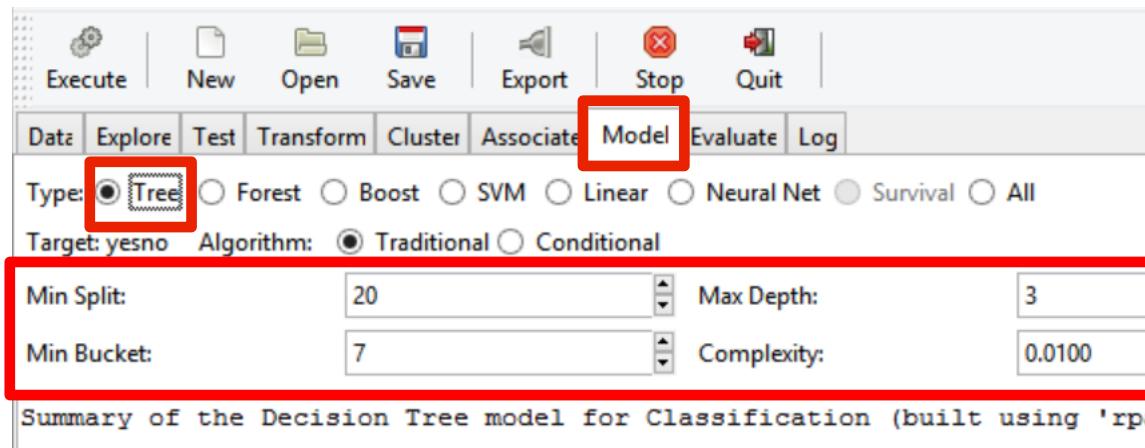
Predicted		
Actual	n	y
n	55.8	6.2
y	7.4	30.6
		Error
	10.0	
	19.5	

Overall error: 13.6%, Averaged class error: 14.75%

Rattle timestamp: 2020-03-17 15:05:57 ashis

Error Rate (Default Values except Max Depth = 3) I

Error rate on training data and validation data for less complex tree (default values in Rattle)



Error Rate (Default Values)

Type: Error Matrix Risk Cost Curve Hand Lift ROC Precision Sensitivity Prv Ob Score

Model: Tree Boost Forest SVM Linear Neural Net Survival KMeans HClust

Data: Training Validation Testing Full Enter CSV File Document... R Dataset

Risk Variable:

Report: Class Probability Include: Identifiers

```
Error matrix for the Decision Tree model on spam.csv [**train**] (counts):
Predicted
Actual   n   y   Error
n 1763 182   9.4
y  283 992  22.2

Error matrix for the Decision Tree model on spam.csv [**train**] (proportions):
Predicted
Actual   n   y   Error
n 54.8  5.7   9.4
y  8.8 30.8  22.2

Overall error: 14.4%, Averaged class error: 15.8%
Rattle timestamp: 2020-03-17 15:09:40 ashis
=====
```

Type: Error Matrix Risk Cost Curve Hand Lift ROC Precision Sensitivity Prv Ob Score

Model: Tree Boost Forest SVM Linear Neural Net Survival KMeans HClust

Data: Training Validation Testing Full Enter CSV File Document... R Dataset

Risk Variable:

Report: Class Probability Include: Identifiers

```
Error matrix for the Decision Tree model on spam.csv [validate] (counts):
Predicted
Actual   n   y   Error
n 389   39   9.1
y  59 203  22.5

Error matrix for the Decision Tree model on spam.csv [validate] (proportions):
Predicted
Actual   n   y   Error
n 56.4  5.7   9.1
y  8.6 29.4  22.5

Overall error: 14.2%, Averaged class error: 15.8%
Rattle timestamp: 2020-03-17 15:08:57 ashis
=====
```

Cutoff for Classification

I

Assign to the class with the highest probability of belonging to that class

Extension to multiple classes by comparing one class against all others

Default probability cutoff is 0.5

To demonstrate used 99:01:00 split. Later will change back

Probability

I

R Data Miner - [Rattle (spam.csv)]

Project Tools Settings Help Rattle Version 5.2.0 togawa

Execute New Open Save Export Stop Quit

Data Explore Test Transform Cluster Associate Model Evaluate Log

Type: Error Matrix Risk Cost Curve Hand Lift ROC Precision Sensitivity Prv O Score

Model: Tree Boost Forest SVM Linear Neural Net Survival KMeans HClust

Data: Training Validation Testing Full Enter CSV File Document... R Dataset

Risk Variable:

Report: Class Probability Include: Identifiers

Score

A model can be deployed on a dataset to obtain scores or classifications for each observation in the dataset.

By default the testing dataset (if any) will be scored. Otherwise the training dataset is scored. As an alternative, a CSV file can be loaded and scored. This choice of what is scored is controlled by the radio button options.

For binary models a probability score can be recorded. For regression models a value is recorded for each observation. Otherwise a class will be recorded for each observation. This can be controlled by the Class and Probability radio buttons.

Source: Rattle GUI / Togaware

Validate Score (new slide)

	A	B	C
1	YesNo	rpart	Predict
2	y	0.761628	y
3	y	0.761628	y
4	y	0.322761	n
5	y	0.881402	y
6	y	0.881402	y
7	y	0.761628	y
8	y	0.881402	y
9	y	0.881402	y
10	y	0.881402	y
11	y	0.881402	y
12	y	0.881402	y
13	y	0.881402	y
14	y	0.761628	y
15	y	0.881402	y
16	y	0.761628	y
17	y	0.881402	y
18	y	0.881402	y
19	y	0.881402	y
20	y	0.100042	n
21	y	0.881402	y

	H	I	J	K	L	M
	Predict			0	1	Error
Actual	0	18	2		0.1	
	1	5	22	0.185185		
Cutoff	0.5					

spam_validate_score_identsNew.xlsx

CUTOFF = 0.5

Validate Score (new slide)

	A	B	C
1	YesNo	rpart	Predict
2	y	0.761628	n
3	y	0.761628	n
4	y	0.322761	n
5	y	0.881402	y
6	y	0.881402	y
7	y	0.761628	n
8	y	0.881402	y
9	y	0.881402	y
10	y	0.881402	y
11	y	0.881402	y
12	y	0.881402	y
13	y	0.881402	y

H	I	J	K	L	M
Predict					
Actual	0	1	Error		
0	19	1	0.05		
1	11	16	0.407407		
Cutoff	0.8				

spam_validate_score_identsNew.xlsx

CUTOFF = 0.8

True and False Negatives

Training Set

Data Explore Test Transform Cluster Associate Model Evaluate Log

Type: Error Matrix Risk Cost Curve Hand Lift ROC Precision Sensitivity

Model: Tree Boost Forest SVM Linear Neural Net Survival KMeans HClust

Data: Training Validation Testing Full Enter CSV File servizi R Dat

Risk Variable: Report: Class Probability Include: Identifiers A

Error matrix for the Decision Tree model on spam.csv [train] (**counts**):

Predicted	n	y	Error
Actual n	1763	182	9.4
y	283	992	22.2

Error matrix for the Decision Tree model on spam.csv [**train**] (proportions):

Predicted	n	y	Error
Actual n	54.8	5.7	9.4
y	8.8	30.8	22.2

Overall error: 14.4%, Averaged class error: 15.8%

Rattle timestamp: 2020-03-17 15:17:44 ashis

True Negative

True Positive

Validation Set

Data Explore Test Transform Cluster Associate Model Evaluate Log

Type: Error Matrix Risk Cost Curve Hand Lift ROC Precision Sensitivity

Model: Tree Boost Forest SVM Linear Neural Net Survival KMeans HClust

Data: Training Validation Testing Full Enter CSV File servizi R Dat

Risk Variable: Report: Class Probability Include: Identifiers A

Error matrix for the Decision Tree model on spam.csv [validate] (**counts**):

Predicted	n	y	Error
Actual n	389	39	9.1
y	59	203	22.5

Error matrix for the Decision Tree model on spam.csv [validate] (proportions):

Predicted	n	y	Error
Actual n	56.4	5.7	9.1
y	8.6	29.4	22.5

Overall error: 14.2%, Averaged class error: 15.8%

Rattle timestamp: 2020-03-17 15:15:00 ashis

Confusion Matrix

True and False Negatives

True positive - model predicts Yes in agreement with the actual outcome (actual outcome = Yes)

True Negative - model predicts No in agreement with the actual outcome (actual outcome = No)

False Positive - model predicts Yes in **DIS**agreement with the actual outcome (actual outcome = No)

False Negative - model predicts No in **DIS**agreement with the actual outcome(actual outcome = Yes)

Specificity and Sensitivity

Confusion Matrix Validation Data

		Predicted		Error
Actual	n	y		
n	389	39	9.1	
y	59	203	22.5	

$$\text{Precision} = \text{TP} / (\text{TP} + \text{FP}) = 203 / (203 + 39) = 0.84$$

$$\text{Recall/Sensitivity/True Positive Rate} = \text{TP} / (\text{TP} + \text{FN}) = 203 / (203 + 59) = 0.77$$

$$\text{Specificity/True Negative Rate} = \text{TN} / (\text{TN} + \text{FP}) = 389 / (389 + 39) = 0.91$$

Risk

Data Explore Test Transform Cluster Associate Model Evaluate Log

Type: Error Matrix Risk Cost Curve Hand Lift ROC Precision Sensitivity Prv Ob Score

Model: Tree Boost Forest SVM Linear Neural Net Survival KMeans HClust

Data: Training Validation Testing Full Enter CSV File Document... R Dataset

Risk Variable: Report: Class Probability Include: Identifiers All

Summary Decision Tree model (built using rpart) on spam.csv [validate] by probability cutoffs.

	Recall	Caseload	Precision
0.0998196031269	1.0000000	1.0000000	0.3797101
0.3054830287206	0.8740458	0.4594203	0.7223975
0.7594594594595	0.7748092	0.3507246	0.8388430
0.884328358209	0.5343511	0.2347826	0.8641975
1.0	0.0000000	0.0000000	1.0000000

Rattle timestamp: 2020-03-17 15:19:34 ashis

The area under the Risk and Recall curves for Decision Tree model

Area under the Recall (green) curve: 91% (0.907)

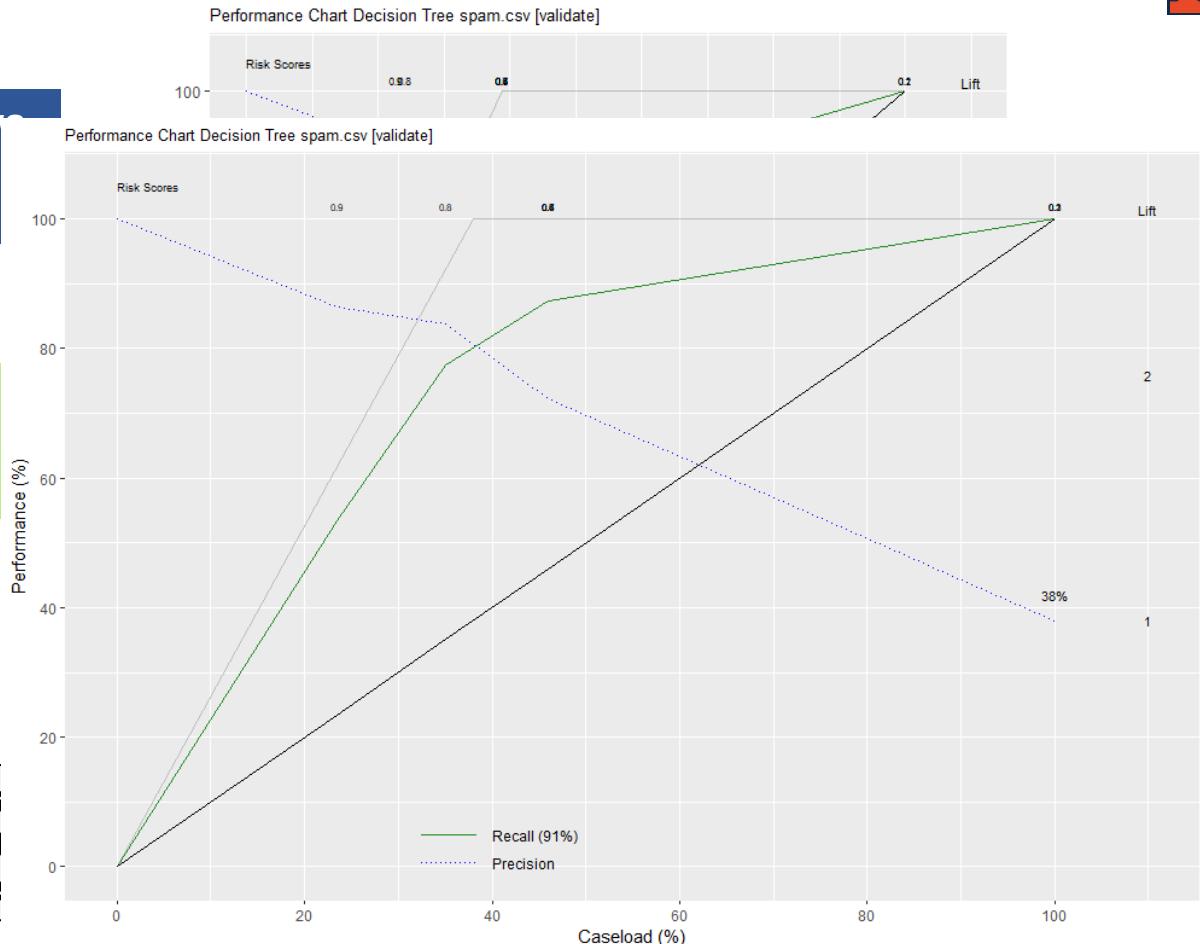
Rattle timestamp: 2020-03-17 15:19:34 ashis

Risk

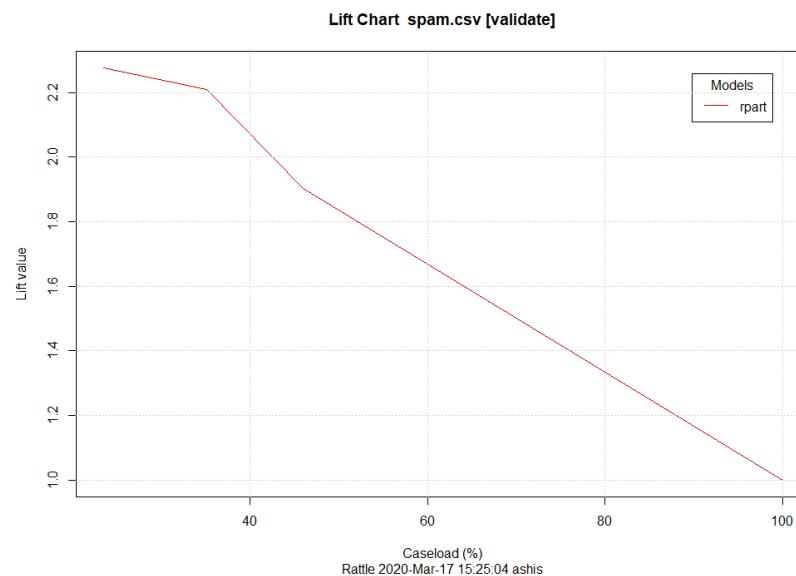
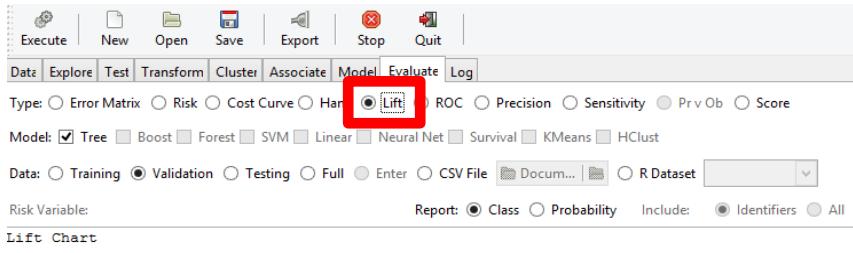
Risk Scores/Precision - This shows what percent of our predicted Spams are indeed spams.

Recall - This shows what percent spams were we able to actually predict as spams.

Lift = Risk score / % of actual order of predicted probability
the proportion of actual “yes”



Evaluation – Lift Chart



Source: Rattle GUI / Togaware

Lift = Risk score / % of actual “yes” in the overall data.

A lift value in the first 40% (sorted in descending order of predicted probability of yes) on the of data being 2 means that in the 4th decile, the proportion of actual “yes” is 2 times the proportion of “yes” in the overall data.

Predicted Class

R Data Miner - [Rattle (spam.csv)]

Project Tools Settings Help

Execute New Open Save Export Stop Quit

Data Explore Test Transform Cluster Associate Model Evaluate Log

Type: Error Matrix Risk Cost Curve Hand Lift ROC Precision Sensitivity Prv Obj Score

Model: Tree Boost Forest SVM Linear Neural Net Survival KMeans HClust

Data: Training Validation Testing Full Enter CSV File SS R Dataset

Risk Variable: Report: Class Probability Include: Identifiers All

Score

Score Files

Name: spam_validate_score_all.csv

Save in folder: SS

Cancel Save

The resulting CSV file will include just those variables having a role as Identifier (plus the Target and the Score), or else all of the variables.

The name of a CSV file into which the results will be written will be prompted for.

Predicted Class:
Spam or not

X	crl.tot	dollar	bang	money	n000	make	yesno	rpart
770	460	0.163	0.02	0.05	0.75	0.17	y	y
1255	103	0	1.417	0.71	0	0.71	y	y
3489	44	0	0	0	0	0	n	n
3781	109	0	0	0	0	0	n	n
2538	21	0	0	0	0	0	n	n

Exercise

Calculate the metrics for the **auction data** using Random Forest to obtain and comment. The table below was produced by analyzing the probability score on validate data in Microsoft Excel. The Rattle steps are as shown in next three slides.

Cutoff	Error	Sensitivity	Specificity	Precision
	0.071429	0.8	1.000	1.000
0.1	0.428571	0.8	0.444	0.444
0.2	0.357143	0.8	0.556	0.500
0.3	0.285714	0.8	0.667	0.571
0.4	0.142857	0.8	0.889	0.800
0.5	0.071429	0.8	1.000	1.000
0.6	0.071429	0.8	1.000	1.000
0.7	0.214286	0.4	1.000	1.000

Exercise - Run the RF Method on Auction Data

R Data Miner - [Rattle (auction.csv)]

Project Tools Settings Help

Rattle Version 5.2.0 togaware.com

Execute New Open Save Export Stop Quit

Date Explore Test Transform Cluster Associate Model Evaluate Log

Source: File ARFF ODBC R Dataset RData File Library Corpus Script

Filename: auction.csv Separator: , Decimal: . Header

Partition 95/5/0 Seed: 42 View Edit

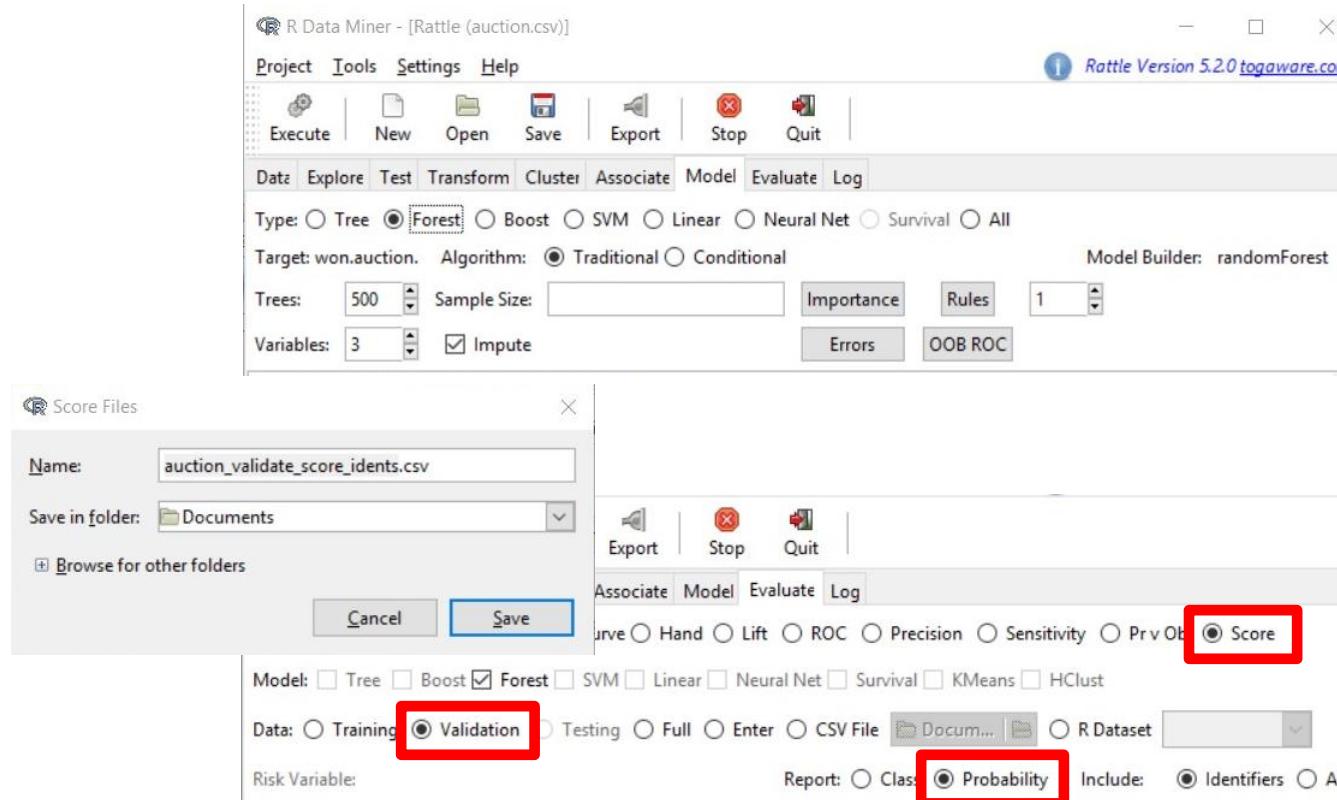
Input Ignore Weight Calculator:

Target Data Type Auto Categoric Numeric Survival

No. Variable	Data Type	Input	Target	Risk	Ident	Ignore	Weight	Comment
1 Bid	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 272
2 MSRP	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 3
3 Price	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 249
4 MSRP.Price	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 252
5 Year	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 36
6 Model.528	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 2
7 Model.526	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 2
8 Model.Baby	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 2
9 Serviced...1.0.	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 2
10 Number.of.bidders	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 20
11 won.auction.	Numeric	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 2

Source: Rattle GUI / Togaware

Exercise - Extract the Score



Source: Rattle GUI / Togaware

Output different in R 3.6.3



Exercise - The Scores

Making Recommendations

I

Data Organization

MULTIVARIATE – table of features

Traditional databases (more later)

BASKET – an **UNWEIGHTED SET** of items

Market Basket, Keyword Sets, Tag Sets,...

BAG – a **WEIGHTED SET** of items

Bag-of-Words, Bag-of-Visual-Words,...

SEQUENCE – of **SYMBOLIC** tokens

Word, Genes, Phonemes,...

SERIES – of **NUMERIC** values

Time Series, Seismographic, EEG,

“Item-Set” – Basket Data

I

Retail – Market Basket of a customer

Movies – Sets of movies a user likes

Flickr – tags describing **images**

YouTube – tags describing **videos**

“Item-Set” – Basket Data

I

AdWords – tags describing **advertisements**

IMDB – tags describing **movies**

Keywords – tags describing **scientific publications**

*The definition of a basket goes beyond the
“Market Basket” – depends on application!*

Item-set Data

A Universe (Dictionary) of ALL Items:

Products | Keywords | Movies | ...

$$\boxed{V} = \left\{ v_1, v_2, \dots, v_N \right\}$$

A Collection of Item-sets:

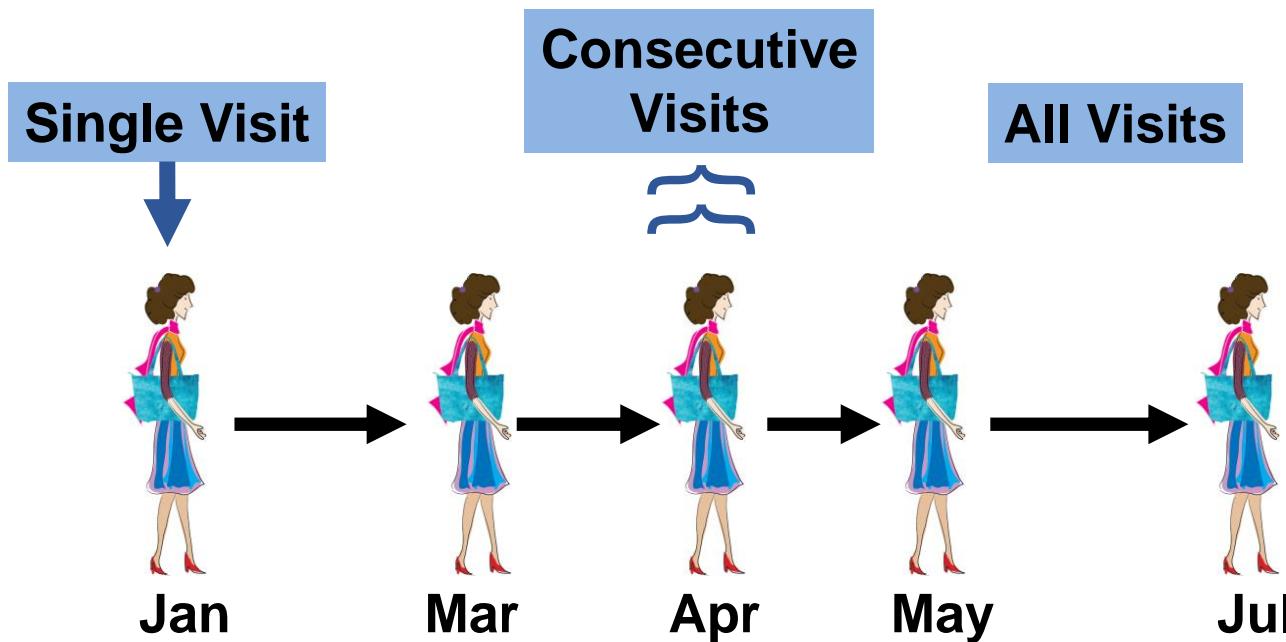
$$\boxed{X} = \left\{ \mathbf{x}_n \mid V \right\}_{n=1}^N$$

a, b, c, d
a, b, d, e
b, e, f
a, e, f, g
b, d, e, f
a, f, g
b, c, e, f

Market Basket

I

Defining a Market Basket



Frequent Item-Set Mining

I

“Frequent itemset mining has been one of the early algorithms that almost gave birth to the field of “data mining.” It was the first breakthrough of its kind in mining such itemset data and since then, there have been a number of improvements in smart data structures to store the candidate and frequent itemsets to make it faster and more scalable. It has also been applied to areas beyond retail data mining for which it was originally invented. It has been used to discover “higher order features” of type “sets of items” in various domains including computer vision.

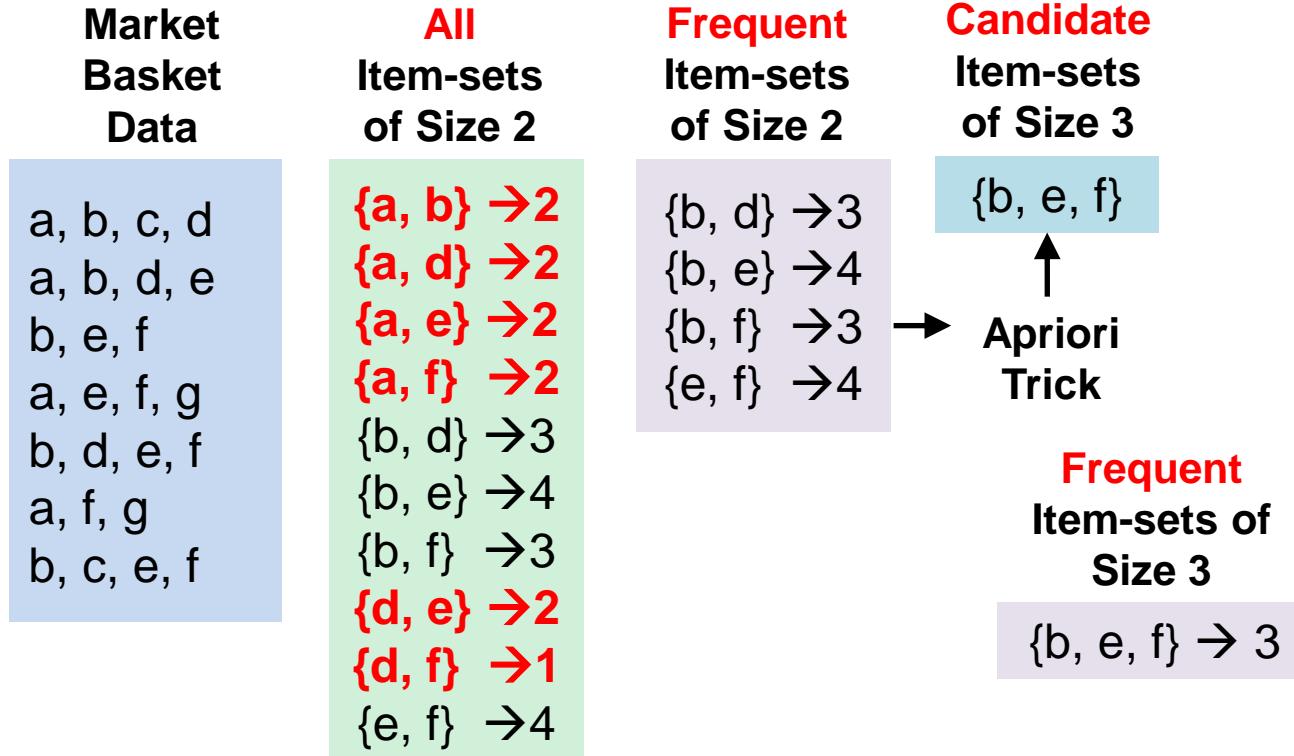
Apriori Algorithm

Support Threshold = 3

Market Basket Data	All Item-sets of Size 1	Frequent Item-sets of Size 1	Candidate Item-sets of Size 2	All Item-sets of Size 2
a, b, c, d	{a} → 4	{a} → 4	{a, b}	{a, b} → 2
a, b, d, e	{b} → 5	{b} → 5	{a, d}	{a, d} → 2
b, e, f	{c} → 2	{d} → 3	{a, e}	{a, e} → 2
a, e, f, g	{d} → 3	{e} → 5	{a, f}	{a, f} → 2
b, d, e, f	{e} → 5	{f} → 5	{b, d}	{b, d} → 3
a, f, g	{f} → 5		{b, e}	{b, e} → 4
b, c, e, f	{g} → 2		{b, f}	{b, f} → 3
		↓	{d, e}	{d, e} → 2
		Apriori Trick →	{d, f}	{d, f} → 1
			{e, f}	{e, f} → 4

Support({a}) = 4
 Support({g}) = 2

Apriori Algorithm



Frequent Itemsets → Association Rules

I

Market Basket Data	Frequent Item-sets of Size 1	Frequent Item-sets of Size 2	Association Rules
a, b, c, d	{a} → 4	{b, d} → 3	{b} → {d}
a, b, d, e	{b} → 5	{b, e} → 4	{d} → {b}
b, e, f	{d} → 3	{b, f} → 3	{b} → {e}
a, e, f, g	{e} → 5	{e, f} → 4	{e} → {b}
b, d, e, f	{f} → 5		{b} → {f}
a, f, g			{f} → {b}
b, c, e, f			{e} → {f}
			{f} → {e}

$$Confidence(\{b\} \rightarrow \{f\}) \equiv P(\{f\} | \{b\}) = \frac{Support(\{b, f\})}{Support(\{b\})} = \frac{3}{5}$$

Frequent Itemsets → Association Rules

I

Market Basket Data	Frequent Item-sets of Size 2	Frequent Item-sets of Size 3	Association Rules
a, b, c, d	{b, d} → 3		
a, b, d, e	{b, e} → 4		
b, e, f	{b, f} → 3		
a, e, f, g	{e, f} → 4		
b, d, e, f			
a, f, g			
b, c, e, f			

$$Conf(\{b, f\} \rightarrow \{e\}) = \frac{Support(\{b, e, f\})}{Support(\{b, f\})} = ?? \quad Conf(\{b\} \rightarrow \{e, f\}) = ??$$

Retail Data – A “Crazy Haystack”

I

Few buy a complete “logical” item-set in same basket

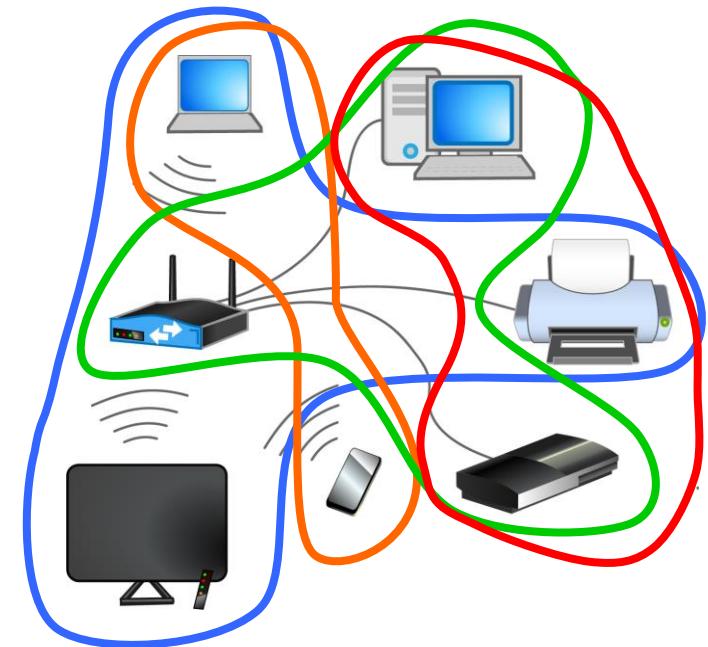
Already have other products

Buy them from another retailer

Buy them at a different time

Got them as gifts

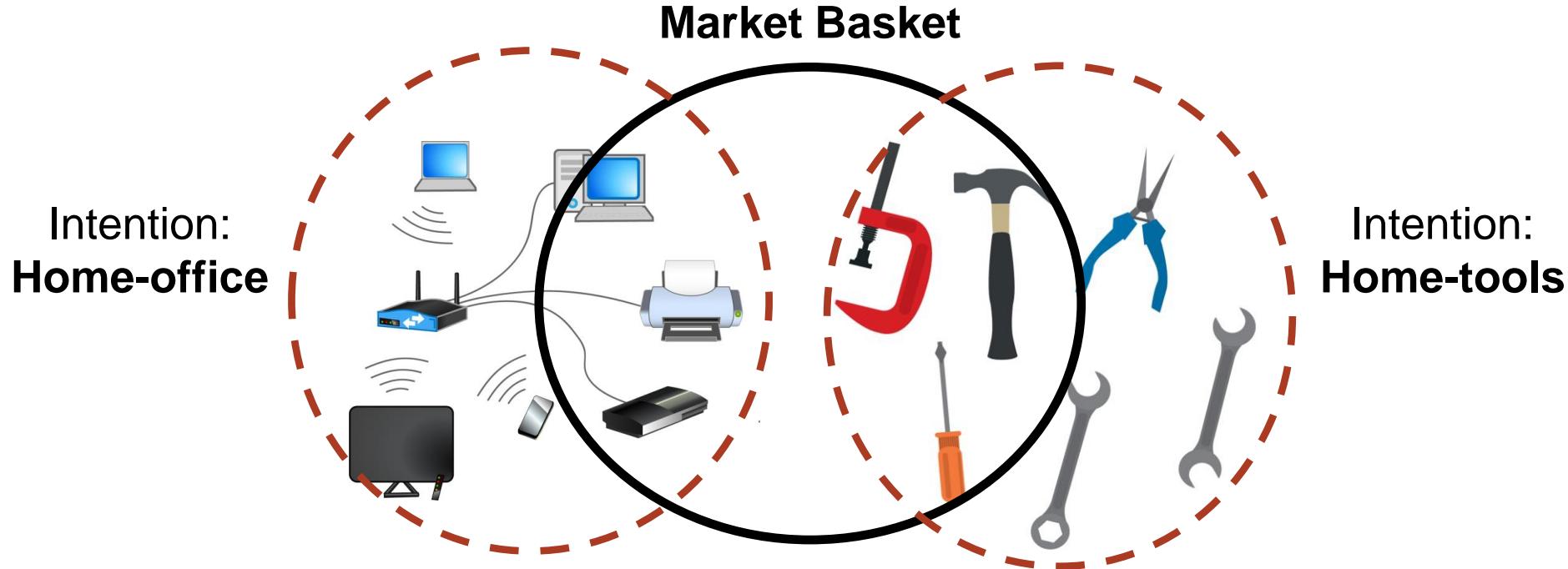
Etc



It's a projection of **latent customer** intentions

It Gets Even Crazier

I



It's a **mixture of projections** of **latent** intentions

How to Recommend Other Topics for Further Study

I

Advanced regression models

SVM

Neural networks

Deep Learning

Philosophy

Association Rule Mining

I

This algorithm can take data in two formats: Basket (each row is a transaction, where each item purchased is given as a list) and Dataframe (where each row is a transaction and each column represents an item that was bought or not)

Association Rule Mining



For Basket format, we will use the DVD_Transaction data from Rattle, and for Dataframe format we will use a subset of Groceries data R Library.

DVD Transaction

I

Predicting which movie is customer likely to order based on his movie watching history

To get the dvd_trans file into Rattle, follow the following steps:

Load rattle() from R

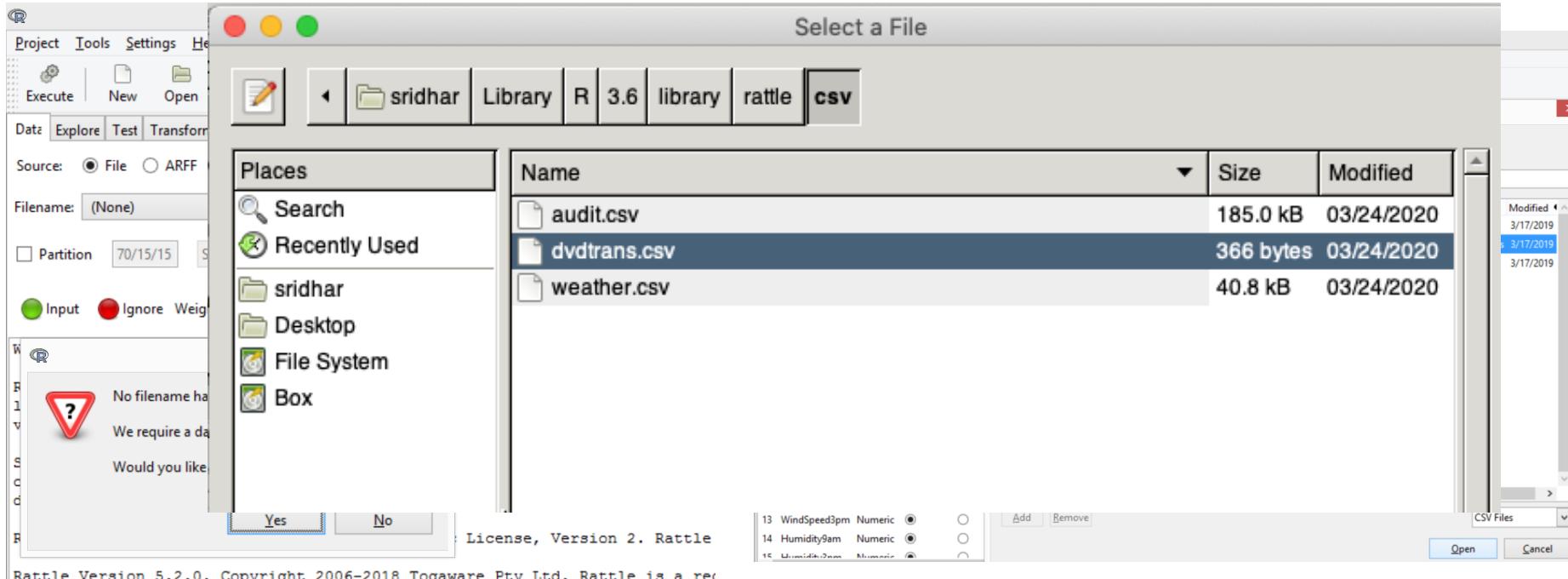
Without selecting any data file, press *Execute*

The rattle would ask permission to load default weather.csv, say Yes

Then click on the *File* button and you will see dvd_trans.csv, select it

DVD Transaction

The file location is machine dependent



DVD_Transactions

R Data Miner - [Rattle (d)]

Project Tools Settings Help

Data Explore Test Transform Cluster Associate Model Evaluate Log

Source: File ARFF ODBC R Dataset RData File Library Corpus Script

Filename: **dvdtrans.csv** Separator: , Decimal: . Header

Partition 70/15/15 Seed: 42 View Edit

Input **Ignore** Weight Calculator: Target Data Type: Auto Categorical Numeric Survival

No. Variable	Data Type	Input	Target	Risk	Ident	Ignore	Weight	Comment
1 ID	Numeric	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 10
2 Item	Categorical	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 10

R Data

Project Tools Settings Help

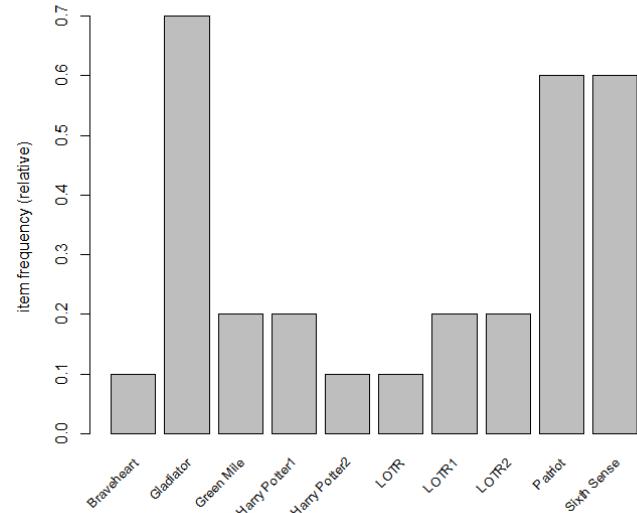
Data Explore Test Transform Cluster Associate Model Evaluate Log

Baskets Support: 0.1000 Confidence: 0.1000 Min Length: 2

Freq Plot Show Rules Sort by: Support Plot

Association Rule Analysis

Association analysis identifies relationships or affinities between observations and/or between variables. These relationships are then expressed as a collection of association rules. The approach has been particularly successful in mining very large transaction databases. It is also often referred to as basket (as in shopping basket) analysis.



Relative frequency among the items

Rules: Sorted by Support and Sorted by Lift

If sort does not work copy to excel and sort

Summary of the Transactions:

Length	Class	Mode
10 transactions	S4	

Summary of the Apriori Association Rules:

Number of Rules: 117

	lhs	rhs	support	confidence	lift	count
[1]	{Patriot}	=> {Gladiator}	0.6	1.0000000	1.4285714	6
[2]	{Gladiator}	=> {Patriot}	0.6	0.8571429	1.4285714	6
[3]	{Sixth Sense}	=> {Gladiator}	0.5	0.8333333	1.1904762	5
[4]	{Gladiator}	=> {Sixth Sense}	0.5	0.7142857	1.1904762	5
[5]	{Patriot}	=> {Sixth Sense}	0.4	0.6666667	1.1111111	4

The rule $\{\text{Gladiator, Green Mile}\} \rightarrow \{\text{LOTR}\}$ has lift =10, which means probability of watching $\{\text{LOTR}\}$ is 10 times, if we know that the person has watched $\{\text{Gladiator, Green Mile}\}$, the probability of watching $\{\text{LOTR}\}$ without any information about his watching history.

Summary of the Transactions:

Length	Class	Mode
10 transactions	S4	

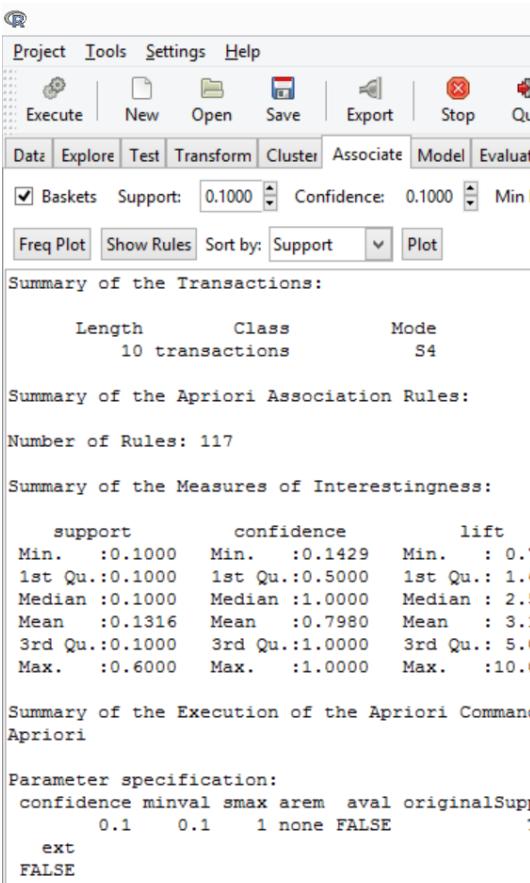
Summary of the Apriori Association Rules:

Number of Rules: 117

	lhs	rhs	support	confidence	lift	count
[1]	{Gladiator,Green Mile}	=> {LOTR}	0.1	1.0000000	10.0000000	10.0000000
[2]	{Gladiator,Green Mile,Sixth Sense}	=> {LOTR}	0.1	1.0000000	10.0000000	10.0000000
[3]	{Harry Potter2}	=> {Harry Potter1}	0.1	1.0000000	5.0000000	5.0000000
[4]	{Harry Potter1}	=> {Harry Potter2}	0.1	0.5000000	5.0000000	5.0000000
[5]	{LOTR}	=> {Green Mile}	0.1	1.0000000	5.0000000	5.0000000
[6]	{Green Mile}	=> {LOTR}	0.1	0.5000000	5.0000000	5.0000000

Plots

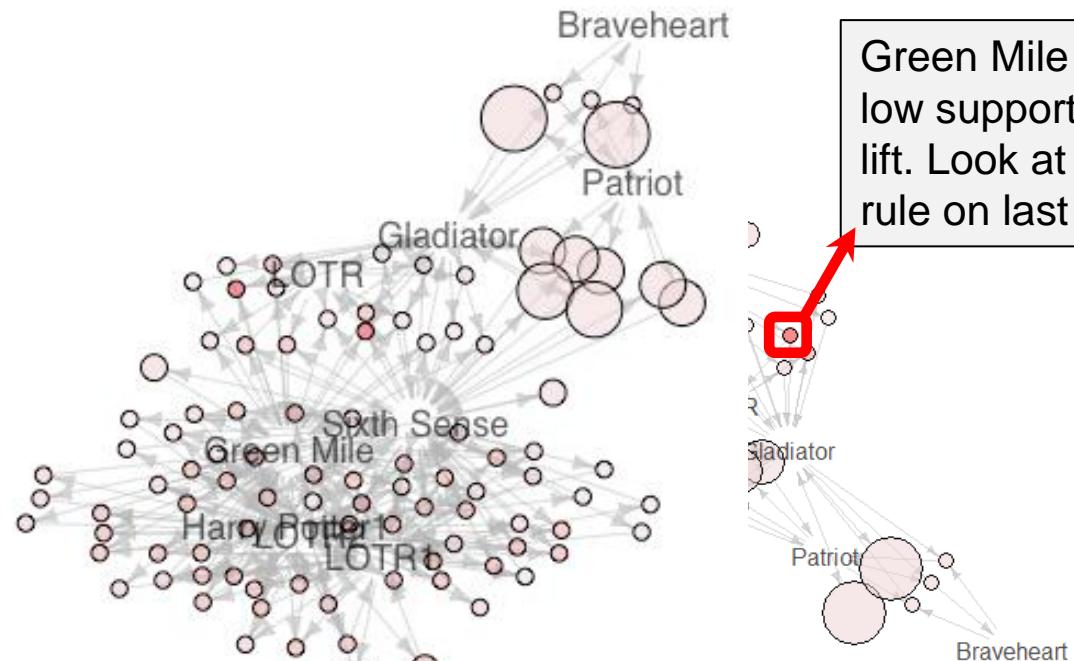
I



Graph for 100 rules

size: support (0.1 - 0.6)
color: lift (0.714 - 10) rules

size: support (0.1 - 0.6)
color: lift (0.714 - 10)



Exercise

The library file is bad. We use a subset.

The subset file is in data under lecture 8.
Remember to convert to categorical

Grocery data as exercise.

Exercise – Use the groceries subset data in R Library
([Groceries:arules:Groceries Data Set](#)) to mine association rules and comment upon them.

Recommender System

Predict ratings and create personalized recommendations for products like books, songs or movies

Other Movies You Might Enjoy

Amelie  Add ★★★☆☆ Not Interested	Y Tu Mama Tambien  Add ★★★☆☆ Not Interested
Guys and Dolls  Add ★★★☆☆ Not Interested	Mostly Martha  Add ★★★★☆ Not Interested
Only Human  Add ★★★☆☆ Not Interested	Russian Dolls  Add ★★★☆☆ Not Interested

Eiken has been added to your Queue at position 2.
This movie is available now.
[Move To Top Of My Queue](#)

[Continue Browsing](#) [Visit your Queue](#)

[Close](#)



<https://d4datascience.wordpress.com/2016/07/22/recommender-systems-101/>

<https://pixabay.com/photos/shelf-toys-shop-store-childhood-830421/>

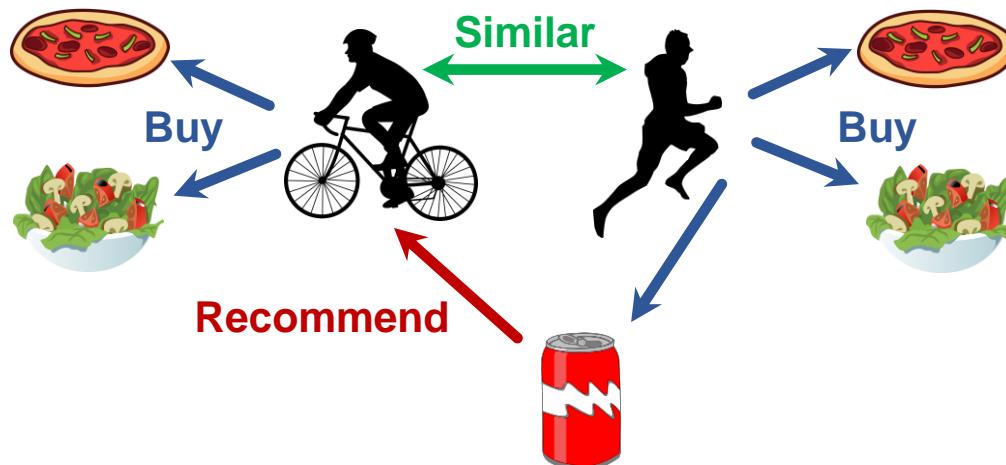
Types of Recommender Systems

Collaborative Filtering

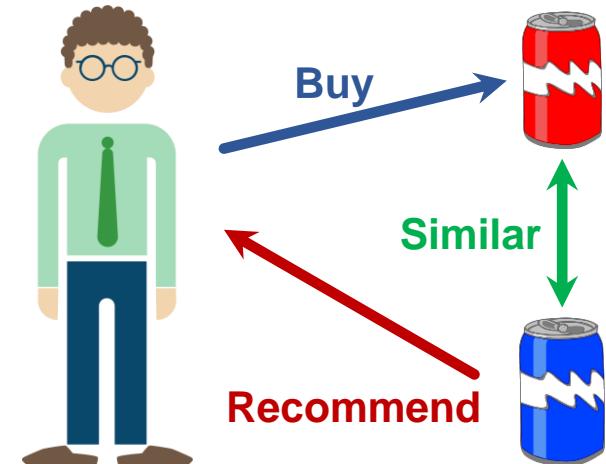
Content-Based Filtering

Hybrid Recommendation Systems – combine above two

Collaborative Filtering



Content-Based Filtering



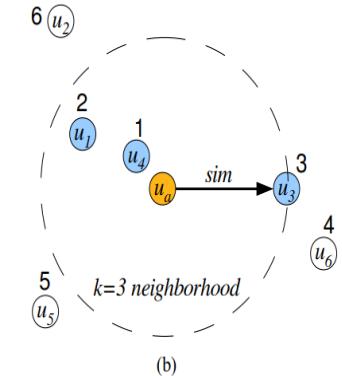
Type of Collaborative Filtering

User-based Collaborative Filtering

We want to predict the ratings for Items i_1 , i_2 , i_5 , and i_7 for the User u_a

	i_1	i_2	i_3	i_4	i_5	i_6	i_7	i_8
u_1	?	4.0	4.0	2.0	1.0	2.0	?	?
u_2	3.0	?	?	?	5.0	1.0	?	?
u_3	3.0	?	?	3.0	2.0	2.0	?	3.0
u_4	4.0	?	?	2.0	1.0	1.0	2.0	4.0
u_5	1.0	1.0	?	?	?	?	?	1.0
u_6	?	1.0	?	?	1.0	1.0	?	1.0
u_a	?	?	4.0	3.0	?	1.0	?	5.0
	3.5	4.0			1.3		2.0	

(a)



(b)

Based on the rating provided by u_a on items i_3 , i_4 , i_6 , and i_8 , we find other users in the neighborhood of U_a created on the basis of ratings on i_3 , i_4 , i_6 , and i_8 . For $k=3$, The users u_1 , u_4 , and u_3 are in the neighborhood. Now, the predicted ratings for u_a for each item is the average rating by users in the neighborhood (the avg. could also be weighted by the others users distance from u_a)

Type of Collaborative Filtering

Item-based Collaborative Filtering

We want to predict the ratings for Items i_2 , i_3 , i_4 , i_6 and i_7 for the User u_a

We create an item-to-item similarity matrix using any similarity measure (Pearson correlation and Cosine similarity). Choose a “k” and k (=3) largest entries are stored per row (these entries are marked using bold face)

In the given matrix – the three items most similar to i_3 are i_2 , i_5 , and i_8 . i_2 can't help in prediction as u_a has not rated it. Thus, we use the ratings provided by the user u_a on items i_5 and i_8 . We take weighted average of u_a 's ratings on i_5 and i_8 to predict his rating on i_3 . Where, weights are the distance between i_3 and i_5 (0.4) and i_3 and i_8 (0.5).

$$U_a \text{'s rating for } i_3 = (0.4*4 + 0.5*5)/(0.4+0.5)$$

S	i_1	i_2	i_3	i_4	i_5	i_6	i_7	i_8	\hat{r}_a
i_1	-	0.1	0	0.3	0.2	0.4	0	0.1	-
i_2	0.1	-	0.8	0.9	0	0.2	0.1	0	0.0
i_3	0	0.8	-	0	0.4	0.1	0.3	0.5	4.6
i_4	0.3	0.9	0	-	0	0.1	0	0.2	3.2
i_5	0.2	0	0.4	0	-	0.1	0.2	0.1	-
i_6	0.4	0.2	0.1	0.3	0.1	-	0	0.1	2.0
i_7	0	0.1	0.3	0	0.2	0	-	0	4.0
i_8	0.1	0	0.5	0.2	0.1	0.1	0	-	-

u_a	2	?	?	?	4	?	?	5	
-------	---	---	---	---	---	---	---	---	--

Overall Idea in the R Script

I

```
# 100 jokes evaluated by 5000 users
```

```
# Split the data into a train set of 90% and 10% of test(holdout) set
```

```
# Evaluation Scheme - We need to define an evaluation scheme to make predictions. This evaluation  
scheme creates train, test data internally for latter usage – we label this as eval
```

```
# Good rating is the cutoff value of rating above which we would predict the user to like the joke
```

```
# Once eval created:
```

```
# For the test set 15 jokes will be given to the recommender algorithm and the other jokes will # be held out  
for computing the error.
```

```
# The 15 jokes are used to find users from the train set who are similar to the users (items) from the test set  
for whom we are evaluating our predictions.
```

```
# Afterwards you may input your evaluation for fifteen jokes and predicted values for the rest of the jokes.
```

R Script

```
install.packages("pacman") # Install the package for managing other packages
library(pacman) # load the "package manage" package in the current session
p_load("recommenderlab") # install load the reqired package(s)

data(Jester5k) # load the jokes data from the recommenderlab library

Jester5k # see the datatype
# 5000 x 100 rating matrix of class 'realRatingMatrix' with 362106 ratings.

set.seed(111) # for reproducability of the results

# We can get all the ratings from this'realRatingMatrix' Jester5k through the command "getRatings"

# Inspecting the max and min in the data
# Check max rating by any user for any joke
max(getRatings(Jester5k)) # 9.9
# Check max rating by any user for any joke
min(getRatings(Jester5k)) # -9.95
```

[Eigentaste: A Constant Time Collaborative Filtering Algorithm](#). Ken Goldberg, Theresa Roeder, Dhruv Gupta, and Chris Perkins. Information Retrieval, 4(2), 133-151. July 2001.

R Script – Continued

I

```
##### Inspecting a user's rating pattern #####
```

```
# Check the ratings by any given user
```

```
rowCounts(Jester5k[1,]) # u2841 - 81 --> User u2841 has rated 81 jokes
```

```
as(Jester5k[1,], "list") # to see the jokes rated by user u2841
```

```
# u2841's rating for joke 1 is "7.91"
```

```
rowMeans(Jester5k[1,]) # Average of all rating for this user is 3.855185
```

```
##### Inspecting Data Summary #####
```

```
# If we use the ratings as provided by the user, the model may have row bias, which is an individual's
```

```
# tendency to rate every joke very high or very low
```

```
# To remove the row bias, we can normalize the ratings - There are many options for normalization
```

```
# such as row centering (difference from the row mean) or Z-score normalization
```

```
# Below we show the distribution of the non-normalized and normalized ratings
```

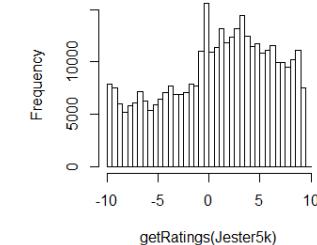
```
hist(getRatings(Jester5k), breaks=50) # top figure
```

```
hist(getRatings(normalize(Jester5k)), breaks=50) # Normalized 2nd from top figure
```

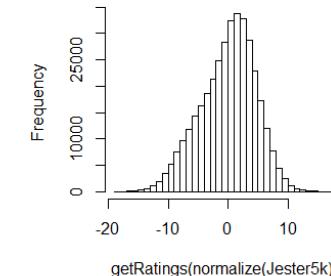
```
hist(rowCounts(Jester5k), breaks = 50) # distribution of # of jokes rated by a user (bottom left)
```

```
hist(colMeans(Jester5k), breaks = 50) # distribution of average ratings per joke(bottom right)
```

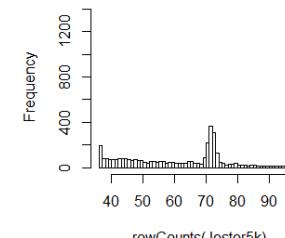
Histogram of getRatings(Jester5k)



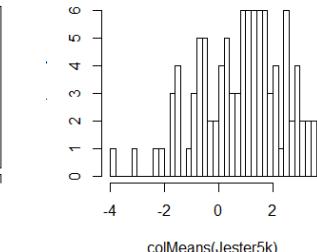
istogram of getRatings(normalize(Jest



Histogram of rowCounts(Jester5k)



Histogram of colMeans(Jester5k)



R Script – Continued

```
# Evaluation Scheme - We need to define an evaluation scheme to make predictions. This evaluation scheme creates train, test data internally for
# later usage
eval <- evaluationScheme(Jester5k[1:2000], method="split", train=0.9, given=15, goodRating=5) # we will use only 2000 users for this exercise out of which 1800 will be in
train set
```

```
# Split the data into a train set of 90% and 10% of test(holdout) set
# For the test set 15 jokes will be given to the recommender algorithm and the other jokes will
# be held out for computing the error. The 15 jokes are used to find users from the train set who are similar
# to the users from the test set for whom we are evaluating our predictions.
# Based on the similar users we find ratings for the remaining jokes for the test users.
# Good rating is the cutoff value of rating above which we would predict the user to like the joke
```

```
eval # This is a scheme that we will use to train the model and evaluate the model's performance
##### Building the Recommender System #####
# A recommender is created using the creator function Recommender()
# Here "known" are the 15 jokes rating we had given in the test to find similar users from the train set
```

```
r_ubcf <- Recommender(getData(eval, "train"), "UBCF") # Recommendation Model generated by user-based collaboration filtering method
pred_ubcf <- predict(r_ubcf, getData(eval, "known"), type="ratingMatrix") # predictions based on UBCF | Predict returns predicted ratings on test data
as(pred_ubcf, "matrix")[1:10,1:10] # Check the predicted ratings of 10 users for first 10 jokes.
```

```
r_ibcf <- Recommender(getData(eval, "train"), "IBCF") # Recommendation Model generated by item-based collaboration filtering method
pred_ibcf <- predict(r_ibcf, getData(eval, "known"), type="ratingMatrix") # predictions based on IBCF | Predict returns predicted ratings on test data
as(pred_ibcf, "matrix")[1:10,1:10] # Check the predicted ratings of 10 users for first 10 jokes.
```

Assessing prediction accuracy of the two models

```
error <- rbind( UBCF = calcPredictionAccuracy(pred_ubcf, getData(eval, "unknown")), # "unkown" are the held out ratings in the test set
                 IBCF = calcPredictionAccuracy(pred_ibcf, getData(eval, "unknown")))
error
# unknown is the jokes not used in making predictions for test data
```

> error

	RMSE	MSE	MAE
UBCF	4.676325	21.86801	3.70405
IBCF	5.334554	28.45746	4.19891

Summary

I

Model performance

Classification error

Frequent Item-set

Recommendation Systems

Course Wrap Up

I

Groceries Data

I

Groceries dataset from Rattle

We use only a subset of the data for illustrative purposes. Data - groceries_subset.csv

transaction id	frankfurter	sausage	ham	finished products	chicken	turkey
1	0	0	0	0	0	0
2	0	0	0	0	0	0
3	0	0	0	0	0	0
4	0	0	0	0	0	0
5	0	0	0	0	0	0
6	0	0	0	0	0	0
7	0	0	0	0	0	0
8	0	0	0	0	0	0
9	0	0	0	0	0	0
10	0	0	0	0	0	0
11	0	0	0	0	0	0
12	0	0	0	0	0	0
13	0	0	0	0	0	0
14	1	0	0	0	0	0
15	0	0	0	0	1	0

Groceries_subset

I

R Data Miner - [Rattle (grocer...)]

Project Tools Settings Help

Execute New Open Save Export Stop Quit

Data Explore Test Transform Cluster Associate Model Evaluate Log

Source: File ARFF ODBC R Dataset RData File Library Corpus Script

Filename: groceries_subset.... Separator: , Decimal: . Header

Partition 70/15/15 Seed: 42 View Edit

Input Ignore Weight Calculator: Target Data Type: Auto Categorical Numeric

No. Variable	Data Type	Input	Target	Risk	Ident	Ignore	Weight	Comment
1 transaction.id	Numeric	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 200
2 frankfurter	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 2
3 sausage	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 2
4 ham	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 2
5 finished.products	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 2
6 chicken	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 2
7 turkey	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 2

Project Tools Settings Help

Execute the current tab to have it take effect (shortcut is F2).

Data Explore Test Transform Cluster Associate Model

Type: Rescale Impute Recode Cleanse

Binning: Quantiles KMeans Equal Width As Categorical Indicator Variable Join Categoricals

No. Variable Data Type and Number Missing

1 transaction.id	Numeric [1 to 200; unique=200]
2 frankfurter	Numeric [0 to 1; unique=2; mean=0.05; std=0.22]
3 sausage	Numeric [0 to 1; unique=2; mean=0.05; std=0.22]
4 ham	Numeric [0 to 1; unique=2; mean=0.05; std=0.22]
5 finished.products	Numeric [0 to 1; unique=2; mean=0.05; std=0.22]
6 chicken	Numeric [0 to 1; unique=2; mean=0.05; std=0.22]
7 turkey	Numeric [0 to 1; unique=2; mean=0.05; std=0.22]

Target Data Type: Auto Categorical Numeric Survival

No. Variable	Data Type	Input	Target	Risk	Ident	Ignore	Weight	Comment
1 transaction.id	Numeric	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 200
2 frankfurter	Numeric	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	Unique: 2
3 sausage	Numeric	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	Unique: 2
4 ham	Numeric	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	Unique: 2
5 finished.products	Numeric	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	Unique: 2
6 chicken	Numeric	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	Unique: 2
7 turkey	Numeric	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	Unique: 2
8 TFC_frankfurter	Categorical	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 2
9 TFC_sausage	Categorical	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 2
10 TFC_ham	Categorical	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 2
11 TFC_finished.products	Categorical	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 2
12 TFC_chicken	Categorical	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 2
13 TFC_turkey	Categorical	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 2

Source: Rattle GUI / Togaware

Association Rule Mining

R Data Miner - [Rattle (groceries_subset.RData)]

Project Tools Settings Help

Execute New Open Save Export Stop Quit

Data Explore Test Transform Cluster Associate Model Evaluate Log

Baskets Support: 0.1000 Confidence: 0.5000 Min Length: 2

Freq Plot Show Rules Sort by: Support Plot

Summary of the Transactions:

Length	Class	Mode
200 transactions		S4

Summary of the Apriori Association Rules:

Number of Rules: 186

Summary of the Measures of Interestingness:

support	confidence	lift	count
Min. :0.805	Min. :0.9048	Min. :0.9942	Min. :161.0
1st Qu.:0.835	1st Qu.:0.9323	1st Qu.:0.9987	1st Qu.:167.0
Median :0.875	Median :0.9830	Median :0.9994	Median :175.0
Mean :0.876	Mean :0.9650	Mean :1.0028	Mean :175.2
3rd Qu.:0.900	3rd Qu.:0.9891	3rd Qu.:1.0068	3rd Qu.:180.0
Max. :0.985	Max. :0.9949	Max. :1.0252	Max. :197.0

Summary of the Execution of the Apriori Command:

Apriori

Parameter specification:

confidence	minval	smax	arem	aval	originalSupport	maxtime	support	minlen	maxlen	target	ext
0.5	0.1	1	none	FALSE		TRUE	5	0.1	2	10	rules FALSE

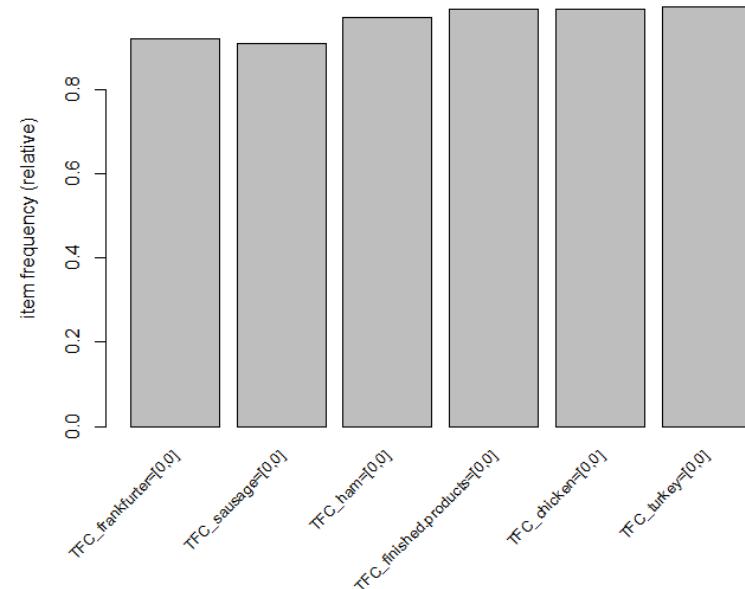
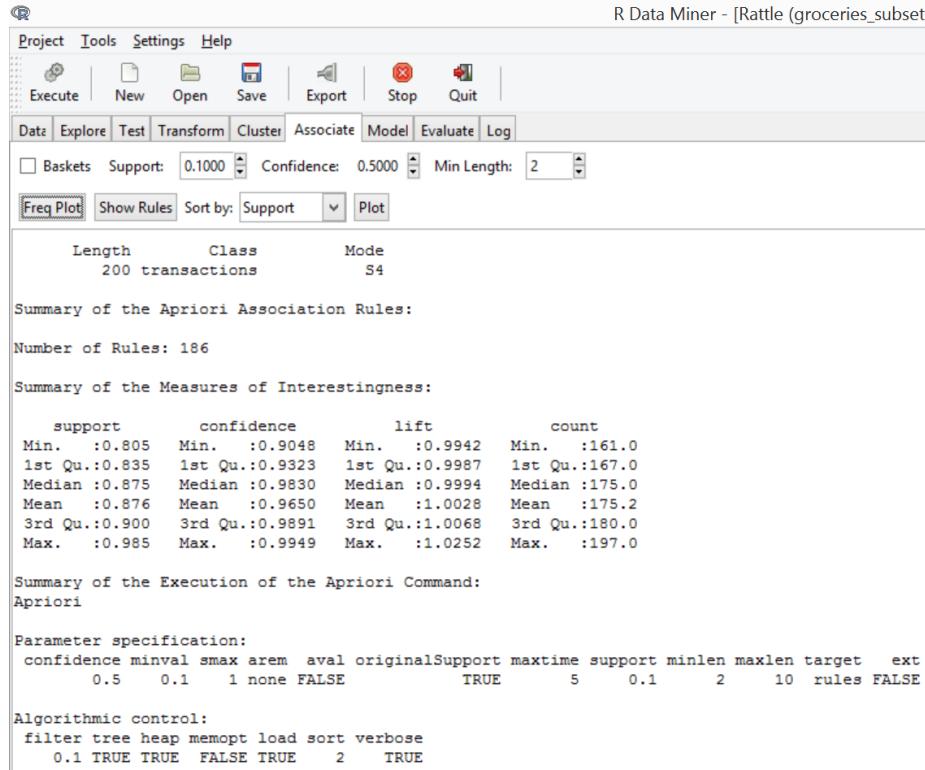
Algorithmic control:

filter tree heap memopt load sort verbose

<-----

The Association Rules model has been built. Time taken: 0.00 secs

Frequency Plots



Association Rules Sorted on Support

I

R Data Miner - [Rattle (groceries_subset.csv)]

Project Tools Settings Help

Execute New Open Save Export Stop Quit

Data Explore Test Transform Cluster Associate Model Evaluate Log

Baskets Support: 0.1000 Confidence: 0.5000 Min Length: 2

Freq Plot Show Rules Sort by: Support Plot

All Rules

lhs	rhs	support	confidence	lift	count
{TFC_finished.products=[0,0]}	=> {TFC_turkey=[0,0]}	0.985	0.9949495	0.9999492	197
{TFC_turkey=[0,0]}	=> {TFC_finished.products=[0,0]}	0.985	0.9899497	0.9999492	197
{TFC_chicken=[0,0]}	=> {TFC_turkey=[0,0]}	0.985	0.9949495	0.9999492	197
{TFC_turkey=[0,0]}	=> {TFC_chicken=[0,0]}	0.985	0.9899497	0.9999492	197
{TFC_finished.products=[0,0]}	=> {TFC_chicken=[0,0]}	0.980	0.9898990	0.9998980	196
{TFC_chicken=[0,0]}	=> {TFC_finished.products=[0,0]}	0.980	0.9898990	0.9998980	196

Sorted on support

R Data Miner - [Rattle (groceries_subset.csv)]

Project Tools Settings Help

Execute New Open Save Export Stop Quit

Data Explore Test Transform Cluster Associate Model Evaluate Log

Baskets Support: 0.1000 Confidence: 0.5000 Min Length: 2

Freq Plot Show Rules Sort by: Lift Plot

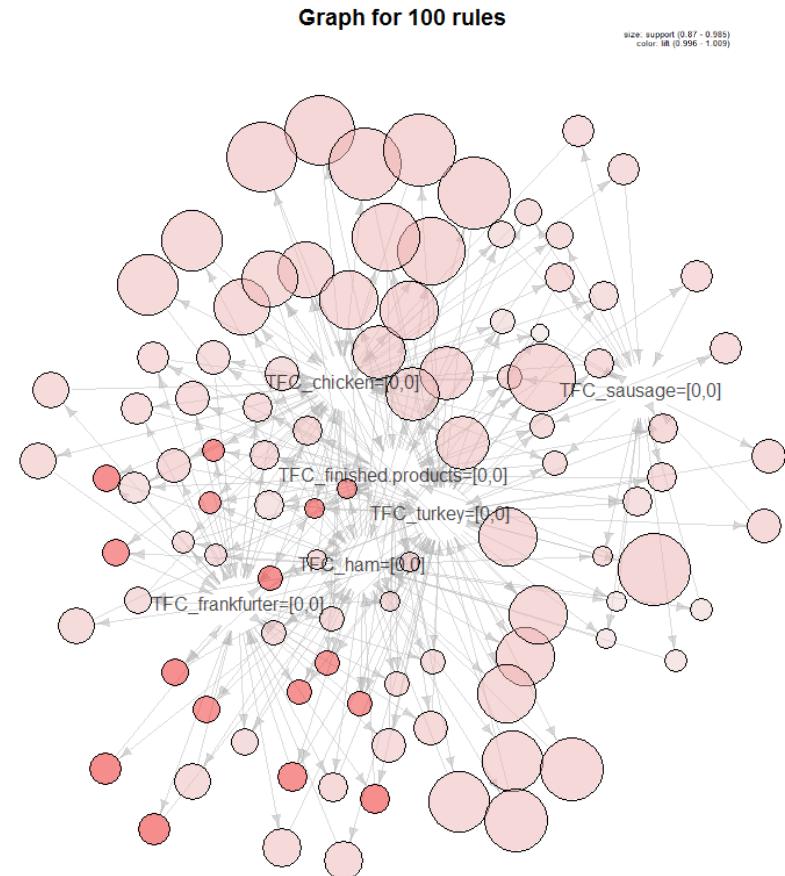
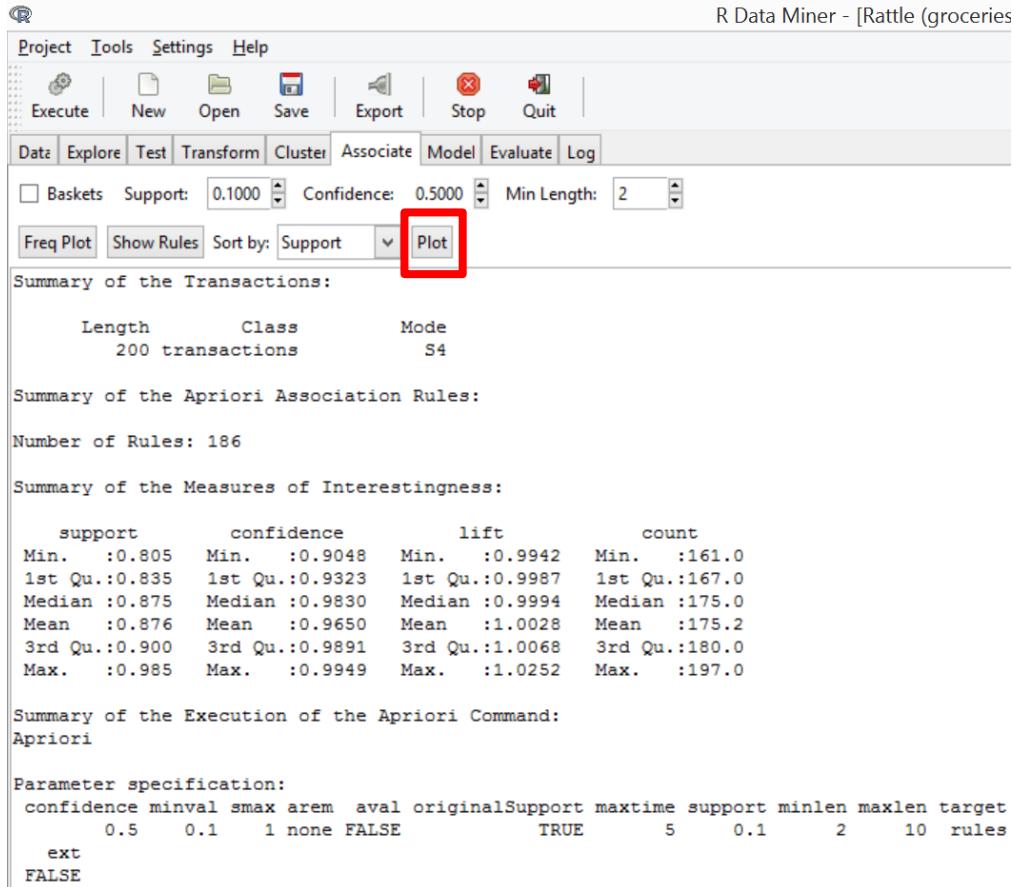
All Rules

lhs	rhs	support	confidence	lift	count
{TFC_sausage=[0,0], TFC_ham=[0,0]}	=> {TFC_frankfurter=[0,0]}	0.830	0.9431818	1.0251976	166
{TFC_sausage=[0,0], TFC_ham=[0,0], TFC_turkey=[0,0]}	=> {TFC_frankfurter=[0,0]}	0.825	0.9428571	1.0248447	165
{TFC_sausage=[0,0], TFC_ham=[0,0], TFC_finished.products=[0,0]}	=> {TFC_frankfurter=[0,0]}	0.820	0.9425287	1.0244878	164

Sorted on lift

Plot – Display Rules Graphically

I



Source: Rattle GUI / Togaware

References

- Hahsler, M. (2015). *recommenderlab: A framework for developing and testing recommendation algorithms*. Retrieved May 22, 2019, from <https://cran.r-project.org/web/packages/recommenderlab/vignettes/recommenderlab.pdf>
- Goldberg, K., Roeder, T., Gupta, D., & Perkins, C. (2001). [Eigentaste: A Constant Time Collaborative Filtering Algorithm.](#) *Information Retrieval*, 4(2).
- Rattle GUI / Togaware (<https://rattle.togaware.com/>)
- Shailesh, K. (2019) Unsupervised machine learning. In B. Pochiraju and S. Seshadri (Eds.), *Essentials of Business Analytics: An Introduction to the Methodology and its Applications* (p. 502). Springer.