

Clustering

Predictive Analytics

Acknowledgments

These slides are based on materials by:
Prof. Bizer (Universität Mannheim)

Example Applications in which Co-Occurrence Matters

- We are often interested in co-occurrence relationships

- **Marketing**

1. identify items that are bought together by sufficiently many customers
2. use this information for marketing or supermarket shelf management purposes



- **Inventory Management**

1. identify parts that are often needed together for repairs
2. use this information to equip your repair vehicles with the right parts



- **Usage Mining**

1. identify words that frequently appear together in search queries
2. use this information to offer auto-completion features to the user



1. Correlation Analysis
2. Association Analysis
 1. Frequent Itemset Generation
 2. Rule Generation
 3. Handling Continuous and Categorical Attributes
 4. Interestingness Measures

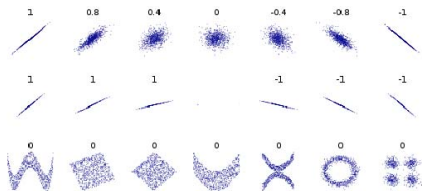
1. Correlation Analysis

- Correlation analysis measures the degree of dependency between **two variables**
 - Continuous variables: Pearson's correlation coefficient (PCC)
 - Binary variables: Phi coefficient

$$PCC(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}}$$

$$Phi(x, y) = \frac{f_{11}f_{00} - f_{01}f_{10}}{\sqrt{f_{1+}f_{+1}f_{0+}f_{+0}}}$$

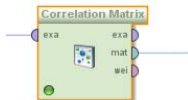
- Value range [-1,1]
 - 1 : positive correlation
 - 0 : variables independent
 - -1 : negative correlation



Correlations between Products in Shopping Baskets

	P1	P2	P3	P4	P5
Basket 1	1	1	0	1	1
Basket 2	1	0	0	1	1
Basket 3	1	0	0	0	1

1 : always bought together
 0 : sometimes bought together
 -1 : never bought together



Correlation Matrix (Correlation Matrix)										
	<input checked="" type="radio"/> Table View <input type="radio"/> Pairwise Table <input type="radio"/> Plot View <input type="radio"/> Annotations									
Attributes	ThinkPad X...	Asus EeePC	HP Laserjet...	2 GB DDR3...	8 GB DDR3...	Lenovo Tab...	Netbook-Sc...	HP CE50 T...	LT Laser M...	LT Minimaus
ThinkPad X2	1	-1	0.356	-0.816	0.612	0.583	-0.667	0.356	0.167	-0.408
Asus EeePC	-1	1	-0.356	0.816	-0.612	-0.583	0.667	-0.356	-0.167	0.408
HP Laserjet	0.356	-0.356	1	-0.218	-0.327	0.356	-0.535	1	-0.089	-0.655
2 GB DDR3	-0.816	0.816	-0.218	1	-0.500	-0.816	0.816	-0.218	0	0.200
8 GB DDR3	0.612	-0.612	-0.327	-0.500	1	0.102	-0.408	-0.327	0.102	0
Lenovo Tabl	0.583	-0.583	0.356	-0.816	0.102	1	-0.667	0.356	-0.250	0
Netbook-Sc	-0.667	0.667	-0.535	0.816	-0.408	-0.667	1	-0.535	0.167	0.408
HP CE50 To	0.356	-0.356	1	-0.218	-0.327	0.356	-0.535	1	-0.089	-0.655
LT Laser Ma	0.167	-0.167	-0.089	0	0.102	-0.250	0.167	-0.089	1	-0.408
LT Minimaus	-0.408	0.408	-0.655	0.200	0	0	0.408	-0.655	-0.408	1

Shortcoming: Measures correlation only between two items but not between multiple items, e.g. {ThinkPad, Cover} → {Minimaus}

2. Association Analysis

- Association analysis can find **multiple item co-occurrence relationships** (descriptive method)
- focuses on occurring items, not absent items
- first algorithms developed in the early 90s at IBM by Agrawal & Srikant
- initially used for **shopping basket analysis** to find how items purchased by customers are related
- later extended to more complex data structures
 - sequential patterns
 - subgraph patterns
- and other application domains
 - web usage mining, social science, life science

Association Analysis

Given a set of transactions, **find rules** that will predict the occurrence of an item based on the occurrences of other items in the transaction.

Shopping Transactions

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

Examples of Association Rules

$\{\text{Diaper}\} \rightarrow \{\text{Beer}\}$

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

$\{\text{Milk, Bread}\} \rightarrow \{\text{Eggs, Coke}\}$

**Implication means
co-occurrence,
not causality!**

Definition: Support and Frequent Itemset

– Itemset

- collection of one or more items
- example: {Milk, Bread, Diaper}
- k-itemset: An itemset that contains k items

– Support count (σ)

- frequency of occurrence of an itemset
- e.g. $\sigma(\{\text{Milk, Bread, Diaper}\}) = 2$

– Support (s)

- fraction of transactions that contain an itemset
- e.g. $s(\{\text{Milk, Bread, Diaper}\}) = 2/5 = 0.4$

– Frequent Itemset

- an itemset whose support is greater than or equal to a minimal support (*minsup*) threshold specified by the user

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

Definition: Association Rule

– Association Rule

- an implication expression of the form $X \rightarrow Y$, where X and Y are itemsets
- an association rule states that when X occurs, Y occurs with certain **probability**.
- Example:
 $\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}$
Condition **Consequent**

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

– Rule Evaluation Metrics

- **Support** (s)
fraction of transactions
that contain both X and Y
- **Confidence** (c)
measures how often items
in Y appear in transactions
that contain X

$$s(X \rightarrow Y) = \frac{|X \cup Y|}{|T|}$$

$$s = \frac{\sigma(\text{Milk, Diaper, Beer})}{|T|} = \frac{2}{5} = 0.4$$

$$c(X \rightarrow Y) = \frac{\sigma(X \cup Y)}{\sigma(X)}$$

$$c = \frac{\sigma(\text{Milk, Diaper, Beer})}{\sigma(\text{Milk, Diaper})} = \frac{2}{3} = 0.67$$

Main Challenges concerning Association Analysis

1. Mining associations from large amounts of data can be **computationally expensive**
 - algorithms need to apply smart pruning strategies
2. Algorithms often discover a **large number of associations**
 - many of them are uninteresting or redundant
 - the user needs to select the subset of the associations that is relevant given her task at hand

The Association Rule Mining Task

- Given a set of transactions T , the goal of association rule mining is to **find all rules** having
 1. support \geq *minsup* threshold
 2. confidence \geq *minconf* threshold
 - *minsup* and *minconf* are provided by the user.
 - Brute Force Approach:
 1. list all possible association rules
 2. compute the support and confidence for each rule
 3. remove rules that fail the *minsup* and *minconf* thresholds
- ⇒ **Computationally prohibitive** due to large number of candidates!

Mining Association Rules

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

Example rules:

$\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}$ ($s=0.4, c=0.67$)
 $\{\text{Milk, Beer}\} \rightarrow \{\text{Diaper}\}$ ($s=0.4, c=1.0$)
 $\{\text{Diaper, Beer}\} \rightarrow \{\text{Milk}\}$ ($s=0.4, c=0.67$)
 $\{\text{Beer}\} \rightarrow \{\text{Milk, Diaper}\}$ ($s=0.4, c=0.67$)
 $\{\text{Diaper}\} \rightarrow \{\text{Milk, Beer}\}$ ($s=0.4, c=0.5$)
 $\{\text{Milk}\} \rightarrow \{\text{Diaper, Beer}\}$ ($s=0.4, c=0.5$)

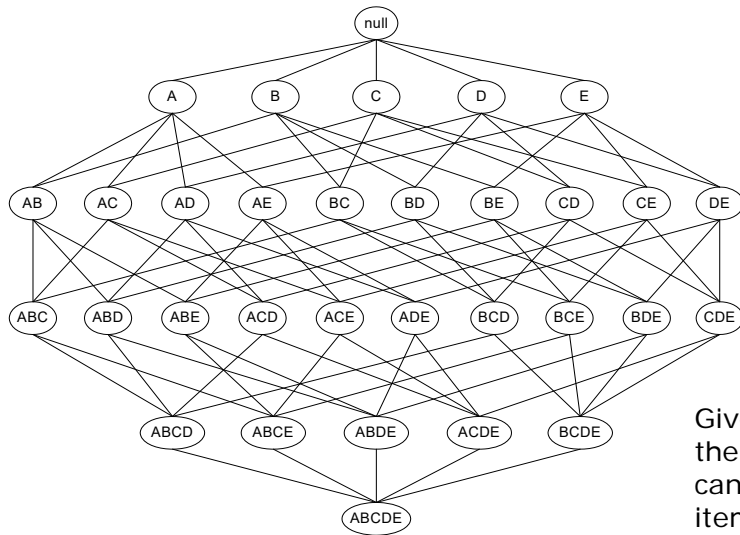
Observations:

- All the above rules are binary partitions of the same itemset:
 $\{\text{Milk, Diaper, Beer}\}$
- Rules originating from the same itemset have identical support but can have different confidence.
- Thus, we may decouple the support and confidence requirements.

Mining Association Rules

- Two-step approach:
 1. Frequent Itemset Generation
 - generate all itemsets whose support \geq minsup
 2. Rule Generation
 - generate high confidence rules from each frequent itemset, where each rule is a binary partitioning of a frequent itemset
- Frequent itemset generation is still computationally expensive

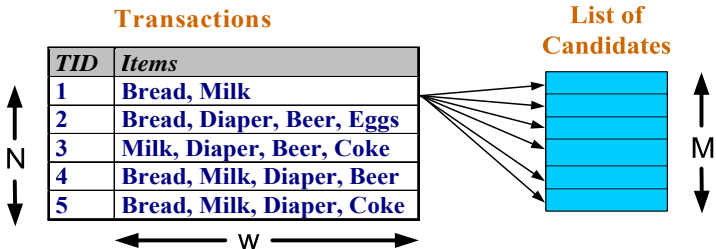
2.1 Frequent Itemset Generation



Given d items,
there are 2^d
candidate
itemsets!

Brute Force Approach

- Each itemset in the lattice is a **candidate** frequent itemset
- Count the support of each candidate by scanning the database
- Match each transaction against every candidate



- Complexity $\sim O(NMw)$ \rightarrow **Expensive since $M = 2^d$!!!**
- A smarter algorithm is required

Example: Brute Force Approach

- Example:
 - Amazon has 10 million books (i.e., Amazon Germany, as of 2011)
- That is $2^{10.000.000}$ possible itemsets
- As a number:
 - $9.04981... \times 10^{3.010.299}$
 - that is: a number with 3 million digits!
- However:
 - most itemsets will not be important at all, e.g., books on Chinese calligraphy, Inuit cooking, and data mining bought together
 - thus, smarter algorithms should be possible
 - intuition for the algorithm: All itemsets containing Inuit cooking are likely infrequent



Reducing the Number of Candidates

– Apriori Principle

If an itemset is frequent, then all of its subsets must also be frequent.

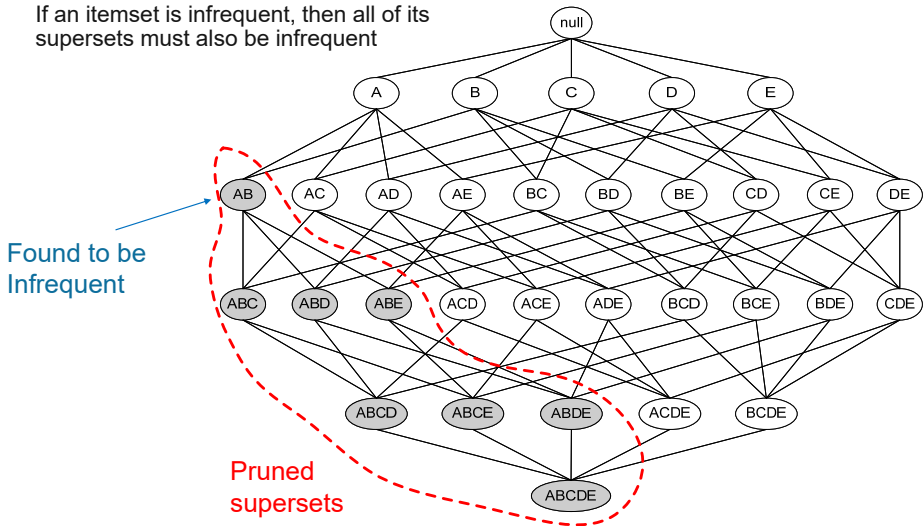
- Apriori principle holds due to the following property of the support measure:

$$\forall X, Y : (X \subseteq Y) \Rightarrow s(X) \geq s(Y)$$

- support of an itemset never exceeds the support of its subsets
- this is known as the **anti-monotone** property of support

Using the Apriori Principle for Pruning

If an itemset is infrequent, then all of its supersets must also be infrequent



Example: Using the Apriori Principle for Pruning

Minimum Support Count = 3

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1

No need to generate candidates involving Coke or Eggs

Items (1-itemsets)



Itemset	Count
{Bread,Milk}	3
{Bread,Beer}	2
{Bread,Diaper}	3
{Milk,Beer}	2
{Milk,Diaper}	3
{Beer,Diaper}	3

Pairs (2-itemsets)



Triplets (3-itemsets)

Item set	Count
{Bread,Milk,Diaper}	3

No need to generate candidate {Milk, Diaper, Beer} as count {Milk, Beer} = 2



The Apriori Algorithm

1. Let $k=1$
2. Generate frequent itemsets of length 1
3. Repeat until no new frequent itemsets are identified
 1. **Generate** length $(k+1)$ candidate itemsets from length k frequent itemsets
 2. **Prune** candidate itemsets that can not be frequent because they contain subsets of length k that are infrequent (Apriori Principle)
 3. **Count** the support of each candidate by scanning the DB
 4. **Eliminate** candidates that are infrequent, leaving only those that are frequent

Example: Apriori Algorithm

itemset:count

minsup=0.5

Dataset T

TID	Items
T100	1, 3, 4
T200	2, 3, 5
T300	1, 2, 3, 5
T400	2, 5

1. scan T

- Cand₁: {1}:2, {2}:3, {3}:3, {4}:1, {5}:3
- Frequ₁: {1}:2, {2}:3, {3}:3, {5}:3
- Cand₂: {1,2}, {1,3}, {1,5}, {2,3}, {2,5}, {3,5}

2. scan T

- Cand₂: {1,2}:1, {1,3}:2, {1,5}:1, {2,3}:2, {2,5}:3, {3,5}:2
- Frequ₂: {1,3}:2, {2,3}:2, {2,5}:3, {3,5}:2
- Cand₃: {2, 3, 5}

3. scan T

- C₃: {2, 3, 5}:2
- F₃: {2, 3, 5}

Frequent Itemset Generation in Rapidminer and Python

RapidMiner

The diagram shows a workflow in RapidMiner. It starts with a 'Read Excel' node (purple) with inputs 'fil' and 'out'. This is followed by an 'FP-Growth' node (green) with inputs 'exa' and 'fre'. The 'FP-Growth' node has three outputs labeled 'res'.

Parameters

FP-Growth

input format	items in dummy code...
min requirement	support
min support	0.2
min items per itemset	1
max items per itemset	0

Python

```
from orangecontrib.associate.fpgrowth import frequent_itemsets  
  
# Calculate frequent itemsets  
itemsets = dict(frequent_itemsets(dataset.values, 0.20))
```

FP-Growth

Alternative frequent itemset generation algorithm which compresses data into tree structure in memory. Details Tan/Steinbach/Kumar: Chapter 4.6

Frequent Itemsets in Rapidminer

Result History		FrequentItemSets (FP-Growth)				
Data	No. of Sets: 83	Size	Support ↓	Item 1	Item 2	Item 3
	Total Max. Size: 4					
Annotations	Min. Size: <input type="text" value="1"/>	1	0.600	Asus EeePC		
	Max. Size: <input type="text" value="4"/>	1	0.500	LT Minimaus		
	Contains item: <input type="text"/>	1	0.500	2 GB DDR3 RAM		
	<input type="button" value="Update View"/>	2	0.500	Asus EeePC	2 GB DDR3 RAM	
		1	0.400	ThinkPad X220		
		1	0.400	Netbook-Schutzhülle		
		1	0.400	Lenovo Tablet Sleeve		
		1	0.400	LT Laser Maus		
		2	0.400	Asus EeePC	LT Minimaus	
		2	0.400	Asus EeePC	Netbook-Schutzhülle	
		2	0.400	2 GB DDR3 RAM	Netbook-Schutzhülle	
		3	0.400	Asus EeePC	2 GB DDR3 RAM	Netbook-Schutzhülle
		1	0.300	HP Laserjet P2055		
		1	0.300	HP CE50 Toner		
		2	0.300	LT Minimaus	2 GB DDR3 RAM	
		2	0.300	LT Minimaus	Netbook-Schutzhülle	
		2	0.300	ThinkPad X220	Lenovo Tablet Sleeve	
		2	0.300	HP Laserjet P2055	HP CE50 Toner	
		3	0.300	Asus EeePC	LT Minimaus	2 GB DDR3 RAM

Example Application of Frequent Itemsets

1. Take top-k frequent itemsets of size 2 containing item A
2. Rank second item according to
 - profit made by selling item
 - whether you want to reduce number of items B in stock
 - knowledge about customer preferences
3. Offer special price for combination with top-ranked second item



Wird oft zusammen gekauft



Preis für beide: EUR 138,00

Beides in den Einkaufswagen

[Verfügbarkeit und Versanddetails anzeigen](#)

- ☒ **Dieser Artikel:** Introduction to Data Mining von Pang-Ning Tan Taschenbuch EUR 85,05
- ☒ Data Mining: Concepts and Techniques (Morgan Kaufmann Series in Data Management Systems)

2.2 Rule Generation

- Given a frequent itemset L , find all non-empty subsets $f \subset L$ such that $f \rightarrow L - f$ satisfies the **minimum confidence** requirement.

Example Frequent Itemset:

{Milk, Diaper, Beer}

Example Rule:

{Milk, Diaper} \Rightarrow Beer

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

$$c = \frac{\sigma(\text{Milk, Diaper, Beer})}{\sigma(\text{Milk, Diaper})} = \frac{2}{3} = 0.67$$

Challenge: Large Number of Candidate Rules

- If $\{A,B,C,D\}$ is a frequent itemset, then the candidate rules are:

$ABC \rightarrow D,$	$ABD \rightarrow C,$	$ACD \rightarrow B,$	$BCD \rightarrow A,$
$A \rightarrow BCD,$	$B \rightarrow ACD,$	$C \rightarrow ABD,$	$D \rightarrow ABC$
$AB \rightarrow CD,$	$AC \rightarrow BD,$	$AD \rightarrow BC,$	$BC \rightarrow AD,$
$BD \rightarrow AC,$	$CD \rightarrow AB$		

- If $|L| = k$, then there are $2^k - 2$ candidate association rules (ignoring $L \rightarrow \emptyset$ and $\emptyset \rightarrow L$)

Rule Generation

- How to efficiently generate rules from frequent itemsets?
 - In general, confidence does not have an anti-monotone property
 $c(ABC \rightarrow D)$ can be larger or smaller than $c(AB \rightarrow D)$
 - But confidence of rules generated from the **same itemset** has an anti-monotone property
 - e.g., $L = \{A, B, C, D\}$:
$$c(ABC \rightarrow D) \geq c(AB \rightarrow CD) \geq c(A \rightarrow BCD)$$
 - Confidence is **anti-monotone with respect to the number of items on the right hand side** of the rule

Explanation

Confidence is anti-monotone w.r.t. number of items on the RHS of the rule

- i.e., “moving elements from left to right” cannot increase confidence

Reason:

$$c(AB \rightarrow C) := \frac{s(ABC)}{s(AB)} \quad c(A \rightarrow BC) := \frac{s(ABC)}{s(A)}$$

- Due to anti-monotone property of support, we know

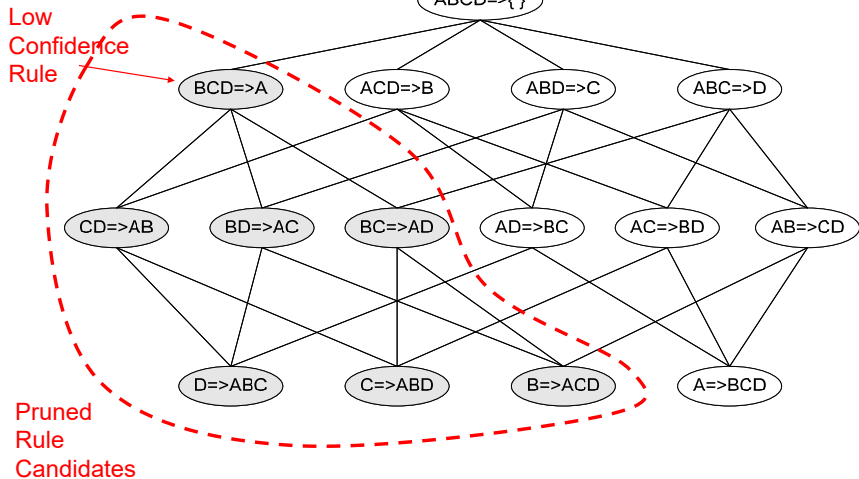
$$s(AB) \leq s(A)$$

- Hence

$$c(AB \rightarrow C) \geq c(A \rightarrow BC)$$

Candidate Rule Pruning

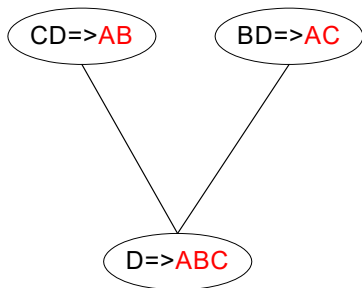
Moving elements from left to right
cannot increase confidence



Candidate Rule Generation within Apriori Algorithm

- Candidate rule is generated by merging two rules that share the same prefix in the rule consequent (right hand side of rule)

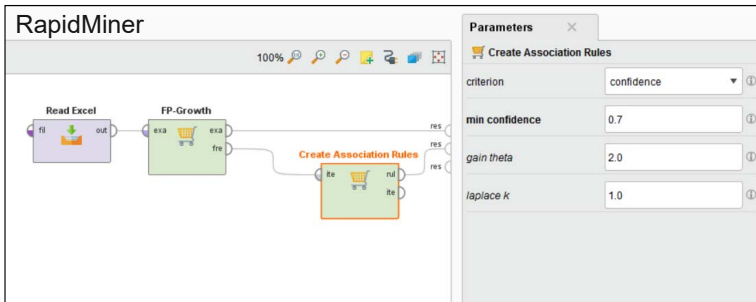
1. $\text{join}(\text{CD} \rightarrow \text{AB}, \text{BD} \rightarrow \text{AC})$
would produce the candidate
rule $\text{D} \rightarrow \text{ABC}$
2. Prune rule $\text{D} \rightarrow \text{ABC}$ if one of its
parent rules does not have
high confidence (e.g. $\text{AD} \rightarrow \text{BC}$)



- All the required information for confidence computation has already been recorded in itemset generation.
- Thus, there is no need to scan the transaction data T any more

Creating Association Rules in Rapidminer and Python

RapidMiner



The workflow diagram shows the following steps:

- Read Excel**: A node with 'fil' and 'out' ports.
- FP-Growth**: A node with 'exa' and 'fre' ports.
- Create Association Rules**: A node with 'ite' and 'rul' ports, highlighted with an orange border.

The 'Parameters' panel for the 'Create Association Rules' node is shown on the right:

Parameters	
 criterion 	confidence
 min confidence 	0.7
 gain theta 	2.0
 laplace k 	1.0

Python

```
from orangecontrib.associate.fpgrowth import association_rules  
  
# Calculate association rules from itemsets  
rules = association_rules(itemsets, 0.70)
```


Exploring Association Rules in Rapidminer

Filter by conclusion

Filter by confidence

No.	Premises	Conclusion	Support	Confiden...
58859	age = young	class = <=50K	0.072	1
58860	native-country = US, age = young	class = <=50K	0.067	1
58861	race = White, age = young	class = <=50K	0.063	1
58862	workclass = Private, age = young	class = <=50K	0.057	1
58863	hours-per-week = full-time, age = young	class = <=50K	0.044	1
58864	sex = Male, age = young	class = <=50K	0.039	1
58865	education = School, age = young	class = <=50K	0.050	1
58866	sex = Female, age = young	class = <=50K	0.032	1
58867	native-country = US, race = White, age = young	class = <=50K	0.060	1
58868	native-country = US, workclass = Private, age = young	class = <=50K	0.053	1
58869	native-country = US, hours-per-week = full-time, age ...	class = <=50K	0.041	1
58870	native-country = US, sex = Male, age = young	class = <=50K	0.037	1
58871	native-country = US, education = School, age = young	class = <=50K	0.047	1
58872	native-country = US, sex = Female, age = young	class = <=50K	0.030	1
58873	race = White, workclass = Private, age = young	class = <=50K	0.051	1
58874	race = White, hours-per-week = full-time, age = young	class = <=50K	0.039	1
58875	race = White, sex = Male, age = young	class = <=50K	0.035	1

2.3 Handling Continuous and Categorical Attributes

- How to apply association analysis to attributes that are not asymmetric binary variables?

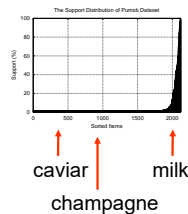
Session Id	Country	Session Length (sec)	Number of Web Pages viewed	Gender	Browser Type	Buy
1	USA	982	8	Male	Chrome	No
2	China	811	10	Female	Chrome	No
3	USA	2125	45	Female	Firefox	Yes
4	Germany	596	4	Male	IE	Yes
5	Australia	123	9	Male	Firefox	No
...

- Example Rule:

$\{\text{Number of Pages} \in [5, 10) \wedge (\text{Browser} = \text{Firefox})\} \rightarrow \{\text{Buy} = \text{No}\}$

Handling Categorical Attributes

- Transform categorical attribute into asymmetric binary variables
- Introduce a **new “item” for each distinct attribute-value pair**
 - e.g. replace “Browser Type” attribute with
 - attribute: “Browser Type = Chrome”
 - attribute: “Browser Type = Firefox”
 -
- Issues
 1. What if attribute has many possible values?
 - many of the attribute values may have very low support
 - potential solution: aggregate low-support attribute values
 2. What if distribution of attribute values is highly skewed?
 - example: 95% of the visitors have Buy = No
 - most of the items will be associated with (Buy=No) item
 - potential solution: drop the highly frequent item



Handling Continuous Attributes

- Transform continuous attribute into binary variables using discretization

- equal-width binning
- equal-frequency binning



- Issue: Size of the discretization intervals affects support & confidence

$\{\text{Refund} = \text{No}, (\text{Income} = \$51,251)\} \rightarrow \{\text{Cheat} = \text{No}\}$

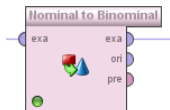
$\{\text{Refund} = \text{No}, (60K \leq \text{Income} \leq 80K)\} \rightarrow \{\text{Cheat} = \text{No}\}$

$\{\text{Refund} = \text{No}, (0K \leq \text{Income} \leq 1B)\} \rightarrow \{\text{Cheat} = \text{No}\}$

- If intervals are too small
 - itemsets may not have enough support
- If intervals too large
 - rules may not have enough confidence
 - e.g. combination of different age groups compared to a specific age group

Attribute Transformation in RapidMiner and Python

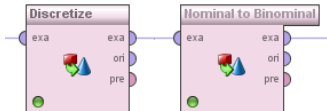
- Categorical attribute values to binary attributes



```
from sklearn.preprocessing import OneHotEncoder

# Apply one-hot encoding
encoder = OneHotEncoder(sparse=False)
onehot_data = encoder.fit_transform(dataset)
```

- Continuous attribute values to binary attributes



```
from sklearn.preprocessing import KBinsDiscretizer

# Discretize and one-hot encode dataset
discretizer = KBinsDiscretizer(n_bins=5, encode='onehot', strategy='quantile')
discretized_onehot_data = discretizer.fit_transform(dataset)
```

2.4 Interestingness Measures

- Association rule algorithms tend to produce **too many rules**
 - many of them are uninteresting or redundant
 - redundant if $\{A,B,C\} \rightarrow \{D\}$ and $\{A,B\} \rightarrow \{D\}$ have same support & confidence
- Interestingness of patterns **depends on application**
 - one man's rubbish may be another's treasure
- Interestingness measures can be used to prune or rank the derived rules.
- In the original formulation of association rules, support & confidence were the only interestingness measures used.
- Later, various other measures have been proposed
 - See Tan/Steinbach/Kumar, Chapter 6.7
 - We will have a look at one: Lift

Drawback of Confidence

Contingency table

	Coffee	<u>Coffee</u>	
Tea	15	5	20
<u>Tea</u>	75	5	80
	90	10	100

Association Rule: Tea \rightarrow Coffee

- $\text{confidence}(\text{Tea} \rightarrow \text{Coffee}) = 0.75$
- **but** $\text{support}(\text{Coffee}) = 0.9$
- although confidence is high, **rule is misleading** as the fraction of coffee drinkers is higher than the confidence of the rule
- we want $\text{confidence}(X \rightarrow Y) > \text{support}(Y)$
- otherwise rule is misleading as X reduces probability of Y

Lift

- The lift of an association rule $X \rightarrow Y$ is defined as:

$$Lift = \frac{c(X \rightarrow Y)}{s(Y)}$$

- Confidence normalized by support of consequent
- Interpretation
 - if $lift > 1$, then X and Y are positively correlated
 - if $lift = 1$, then X and Y are independent
 - if $lift < 1$, then X and Y are negatively correlated

Example: Lift

Contingency table

	Coffee	Coffee	
Tea	15	5	20
Tea	75	5	80
	90	10	100

$$Lift = \frac{c(X \rightarrow Y)}{s(Y)}$$

Association Rule: Tea \rightarrow Coffee

- confidence(Tea \rightarrow Coffee) = 0.75
- but support(Coffee) = 0.9

$$Lift(\text{Tea} \rightarrow \text{Coffee}) = 0.75/0.9 = 0.8333$$

- lift < 1, therefore is **negatively correlated**

Exploring Association Rules in RapidMiner

Top Window: AssociationRules (Create Association Rules)

Show rules matching: all of these conclusions:

class = <=50K
education = HS-grad
class = >50K
education = Bachelors
education = Some-college
occupation = Other-service
occupation = Prof-specialty
occupation = Exec-managerial
occupation = Adm-clerical
education = Masters

No.	Premises	Conclusion	Support	Confiden...	Lift
47	occupation = Machine-op-inspct	class = <=50K	0.085	0.922	1.150
42	occupation = Adm-clerical	class = <=50K	0.080	0.854	1.064
34	occupation = Prof-specialty	class = <=50K	0.069	0.521	0.850
38	occupation = Sales	class = <=50K	0.068	0.798	0.995
52	education = 5th-6th	class = <=50K	0.066	0.946	1.179
17	class = >50K	occupation = Prof-specialty	0.064	0.321	2.417
30	occupation = Prof-specialty	class = >50K	0.064	0.479	2.417
13	class = >50K	education = Bachelors	0.058	0.295	1.758
25	education = Bachelors	class = >50K	0.058	0.348	1.758
35	occupation = Exec-managerial	class = <=50K	0.053	0.554	0.691
3	education = HS-grad	occupation = Other-service	0.051	0.211	1.428
24	occupation = Other-service	education = HS-grad	0.051	0.346	1.428
49	occupation = Handlers-cleaners	class = <=50K	0.049	0.936	1.167

Bottom Window: AssociationRules (Create Association Rules)

Show rules matching: all of these conclusions:

class = <=50K
education = HS-grad
class = >50K
education = Bachelors
education = Some-college
occupation = Other-service
occupation = Prof-specialty
occupation = Exec-managerial
occupation = Adm-clerical
education = Masters

No.	Premises	Conclusion	Support	Confidence	Lift
25	education = Bachelors	class = >50K	0.058	0.348	1.758
29	occupation = Exec-managerial	class = >50K	0.043	0.446	2.249
30	occupation = Prof-specialty	class = >50K	0.064	0.479	2.417
31	education = Masters	class = >50K	0.030	0.484	2.441

Conclusion

- The algorithm does the counting for you and finds patterns in the data
- You need to do the interpretation based on your knowledge about the application domain.
 - Which patterns are meaningful?
 - Which patterns are surprising?

Literature for this Slideset

Pang-Ning Tan, Michael Steinbach, Anuj Karpatne, Vipin Kumar: **Introduction to Data Mining**. 2nd Edition. Pearson.

**Chapter 4: Association Analysis:
Basic Concepts and Algorithms**

**Chapter 7: Association Analysis:
Advanced Concepts**

