



MEDNARODNA  
PODIPLOMSKA ŠOLA  
JOŽEFA STEFANA

INFORMATION AND COMMUNICATION TECHNOLOGIES  
PhD study programme

# Data Mining and Knowledge Discovery

Petra Kralj Novak

January 13, 2020

[http://kt.ijs.si/petra\\_kralj/dmtm2.html](http://kt.ijs.si/petra_kralj/dmtm2.html)

# In previous episodes ...

- 23-Oct-19
  - **Data**, data types
  - Interactive **visualization** (Orange)
  - **Classification** with decision trees (root, leaves, rules, entropy, info gain, TDIDT, ID3)
- 6-Nov-19
  - Classification: train – test (evaluate) - apply
  - **Decision tree** example (on blackboard)
  - Decision tree language bias (Orange workflow)
  - Homework:
    - InfoGain questions
    - Orange workflow
    - Reading “Classification and regression by randomForest” by Liaw & Wiener, 2002
- 25-Nov-19
  - **Evaluation:**
    - Methods: train-test, leave-one-out, randomized sampling,...
    - Metrics: accuracy, confusion matrix, precision, recall, F1,...
  - Homework: XOR, questions, precision and recall

# ... continued ...

- 2-Dec-19
  - Evaluation: **ROC**
  - **Naïve Bayes** classifier
  - Probability estimation: relative frequency, Laplace estimate
  - **Numeric prediction** (linear regression, regression tree, model tree, KNN) and **evaluation** (MSE, MAE, RMSE)
- 16-Dec-19
  - Clustering
  - K-means
  - Silhouette coefficient
  - Agglomerative clustering, dendrogram
  - DB-scan
  - Similarity, distance

# Data mining techniques

## Predictive induction

## Descriptive induction

### Classification

Decision trees

Classification rules

Naive Bayes classifier

KNN

SVM

ANN

...

### Numeric prediction

Linear regression

Regression / model trees

KNN

SVM

ANN

...

### Association rules

Apriori

FP-growth

...

### Clustering

Hierarchical

K-means

Dbscan

...

Or harried dads rewarding themselves with impulse buys





# Association rules

# Association rules – Market basket analysis

- What do customers buy together?
  - Which items imply the purchase of other items?
- \* Terminology from market basket analysis (transactions, items, itemsets, ...)
- Determine associations between groups of items bought by customers.
  - No predefined target variable(s).
  - Find interesting, useful patterns and relationships.
  - Data mining, business intelligence.



# Confidence and support

- The dataset consists of  $n$  transactions
- We have an association rule  $A \rightarrow B$

The **support** of an itemset  $A$  is defined as the fraction of the transactions in the database  $T = \{T_1 \dots T_n\}$  that contain  $A$  as a subset.

$$\text{supp}(A) = \frac{|A|}{n}$$
$$\text{supp}(A \rightarrow B) = \frac{|A \wedge B|}{n}$$

The **confidence** of the rule  $A \rightarrow B$  is the conditional probability of  $A$  and  $B$  occurring in a transaction, given that the transaction contains  $A$ .

$$\text{conf}(A \rightarrow B) = \frac{|A \wedge B|}{|A|} = P(B|A)$$



# Exercise

tid	Set of items	Binary representation
1	{ <i>Bread, Butter, Milk</i> }	110010
2	{ <i>Eggs, Milk, Yogurt</i> }	000111
3	{ <i>Bread, Cheese, Eggs, Milk</i> }	101110
4	{ <i>Eggs, Milk, Yogurt</i> }	000111
5	{ <i>Cheese, Milk, Yogurt</i> }	001011

$$\text{supp}(A) = \frac{|A|}{n}$$

$$\text{supp}(A \rightarrow B) = \frac{|A \wedge B|}{n}$$

$$\text{conf}(A \rightarrow B) = \frac{|A \wedge B|}{|A|} = P(B|A)$$

Supp ({Bread}) =

Supp ({Milk, Yogurt}) =

Conf ({Milk}  $\rightarrow$  {Yogurt}) =

Conf ({Yogurt}  $\rightarrow$  {Milk}) =

# Association rules

- Rules  $A \rightarrow B$ , where A and B are conjunctions of items
- Task: Find all association rules that satisfy the minimum support and minimum confidence constraints

- **Support:**

$$\text{supp}(A \rightarrow B) = \frac{|A \wedge B|}{n}$$

- **Confidence:**

$$\text{conf}(A \rightarrow B) = \frac{|A \wedge B|}{|A|} = P(B|A)$$

# Hard problem

- In practice
  - Millions of transactions
  - Many (thousands) of items
- Too many possible combinations
  - 1000 items for sale  $\rightarrow 2^{1000} - 1$  candidate market baskets
- Solution
  - *Apriori algorithm*



# Frequent itemsets: intuition

- We have  $n$  transactions containing (at least) {gloves, scarf and hat}
- What can we say about the number of transaction containing {gloves and scarf}?

*At least  $n$ .*

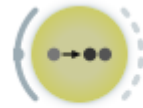
- The **anti-monotone property of support**: if we drop out an item from an itemset, support value of new itemset generated will either be the same or will increase.

$$\forall A, B : A \subseteq B \Rightarrow \text{supp}(A) \geq \text{supp}(B)$$

# Apriori



Frequent Itemsets



Association Rules

Frequent  
itemsets

- Find all itemsets within the *minSupport* constraint

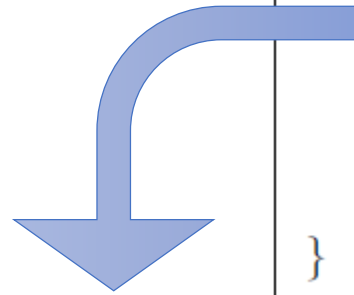
Generate rules

- For all frequent itemsets, find rules which satisfy the *minConfidence* constraint

Association  
rules

\*Frequent itemsets = large itemsets, sometimes also frequent patterns

# Apriori



```
Create  $L_1$  = set of supported itemsets of cardinality one
Set  $k$  to 2
while ( $L_{k-1} \neq \emptyset$ ) {
  Create  $C_k$  from  $L_{k-1}$ 
  Prune all the itemsets in  $C_k$  that are not
    supported, to create  $L_k$ 
  Increase  $k$  by 1
}
The set of all supported itemsets is  $L_1 \cup L_2 \cup \dots \cup L_k$ 
```

(Generates  $C_k$  from  $L_{k-1}$ )

## Join Step

Compare each member of  $L_{k-1}$ , say  $A$ , with every other member, say  $B$ , in turn. If the first  $k - 2$  items in  $A$  and  $B$  (i.e. all but the rightmost elements of the two itemsets) are identical, place set  $A \cup B$  into  $C_k$ .

## Prune Step

```
For each member  $c$  of  $C_k$  in turn {
  Examine all subsets of  $c$  with  $k - 1$  elements
  Delete  $c$  from  $C_k$  if any of the subsets is not a member of  $L_{k-1}$ 
}
```

# Apriori: Frequent itemset mining

Create  $L_1$  = set of supported itemsets of cardinality one

Set  $k$  to 2

while ( $L_{k-1} \neq \emptyset$ ) {

    Create  $C_k$  from  $L_{k-1}$

    Prune all the itemsets in  $C_k$  that are not supported, to create  $L_k$

    Increase  $k$  by 1

}

The set of all supported itemsets is  $L_1 \cup L_2 \cup \dots \cup L_k$

- The items in the sets should be ordered (alphabetically, ...)

# Apriori: constructing the next level from the previous one

- Since items in the sets are ordered (alphabetically, ...)
- Join Step:
  - Merge sets that have all the elements the same except for the rightmost one
- Prune Step:
  - Remove the set if any of its subsets are not on the previous level

(Generates  $C_k$  from  $L_{k-1}$ )

## Join Step

Compare each member of  $L_{k-1}$ , say  $A$ , with every other member, say  $B$ , in turn. If the first  $k - 2$  items in  $A$  and  $B$  (i.e. all but the rightmost elements of the two itemsets) are identical, place set  $A \cup B$  into  $C_k$ .

## Prune Step

For each member  $c$  of  $C_k$  in turn {

Examine all subsets of  $c$  with  $k - 1$  elements

Delete  $c$  from  $C_k$  if any of the subsets is not a member of  $L_{k-1}$

}



# Rules from frequent itemsets

- Generate rules with a certain confidence
- All the counts we need are in the lattice (no database scanning)
- Confidence of rules generated from the same itemset has an anti-monotone property
- No need to check all the rules, since

$$\text{Conf} ( \{A,B\} \rightarrow \{C\} ) \geq \text{Conf} ( \{A\} \rightarrow \{B,C\} )$$

$$\text{Conf}( \{A,B,C\} \rightarrow \{D\} ) \geq \text{Conf}( \{A,B\} \rightarrow \{C,D\} ) \geq \text{Conf}( \{A\} \rightarrow \{B,C,D\} )$$

\*In general, confidence does not have an anti-monotone property:  $\text{Conf}(ABC \rightarrow D)$  can be larger or smaller than  $\text{Conf}(AB \rightarrow D)$

# Exercise: Association rules

Generate frequent itemsets with support at least 2/6 and confidence at least 75%.

---

Items: **A**=apple, **B**=banana, **C**=coca-cola, **D**=doughnut

- Client 1 bought: A, B, C, D
- Client 2 bought: B, C
- Client 3 bought: B, D
- Client 4 bought: A, C
- Client 5 bought: A, B, D
- Client 6 bought: A, B, C

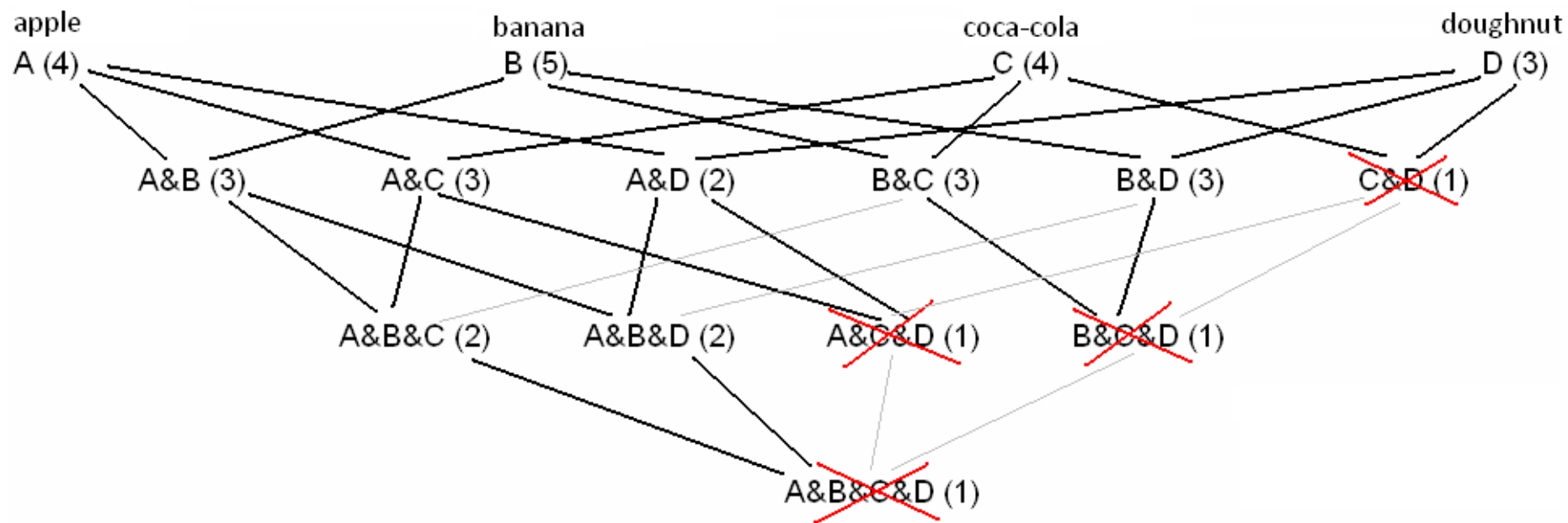
# Exercise: Frequent itemsets

To ease the counting, we transcribe into a binary representation.

<b>A</b>	<b>B</b>	<b>C</b>	<b>D</b>
<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>
	<b>1</b>	<b>1</b>	
	<b>1</b>		<b>1</b>
<b>1</b>		<b>1</b>	
<b>1</b>	<b>1</b>		<b>1</b>
<b>1</b>	<b>1</b>	<b>1</b>	

# Frequent itemsets (= large itemsets)

A	B	C	D
1	1	1	1
	1	1	
	1		1
1		1	
1	1		1
1	1	1	



# Association rules ...

Itemset (count)	Rule	Support	Confidence	Over threshold
AB (3)	$A \rightarrow B$	3/6	3/4 = 75%	✓
	$B \rightarrow A$	3/6	3/5 = 60%	
AC (3)	$A \rightarrow C$	3/6	3/4 = 75%	✓
	$C \rightarrow A$	3/6	3/4 = 75%	✓
AD (2)	$A \rightarrow D$	2/6	2/4 = 50%	
	$D \rightarrow A$	2/6	2/3 = 67%	
BC (3)	$B \rightarrow C$	3/6	3/5 = 60%	
	$C \rightarrow B$	3/6	3/4 = 75%	✓
BD (3)	$B \rightarrow D$	3/6	3/5 = 60%	
	$D \rightarrow B$	3/6	3/3 = 100%	✓

## ... association rules

ABC (2)	$AB \rightarrow C$	$2/6$	$2/3 = 67\%$	
	$AC \rightarrow B$	$2/6$	$2/3 = 67\%$	
	$BC \rightarrow A$	$2/6$	$2/3 = 67\%$	
	$A \rightarrow BC$	We do not generate these rules because transferring members of a supported <u>itemset</u> from the left-hand side of a rule to the right-hand side cannot increase the value of rule confidence.		
	$B \rightarrow AC$			
	$C \rightarrow AB$			
ABD (2)	$AB \rightarrow D$	$2/6$	$2/3 = 67\%$	
	$AD \rightarrow B$	$2/6$	$2/2 = 100\%$	✓
	$BD \rightarrow A$	$2/6$	$2/3 = 67\%$	
	$A \rightarrow BD$	$2/6$	$2/4 = 50\%$	
	$B \rightarrow AD$	We do not generate this rule.		
	$D \rightarrow AB$	$2/6$	$2/3 = 67\%$	

# Lift

- The lift of rule  $L \rightarrow R$  measures how many more times the items in L and R occur together in transactions than would be expected if the itemsets L and R were statistically independent.

$$\text{lift}(L \rightarrow R) = \frac{\text{support}(L \cup R)}{\text{support}(L) \times \text{support}(R)}$$

$$\text{lift}(L \rightarrow R) = \text{lift}(R \rightarrow L)$$

# Leverage

- The leverage of rule  $L \rightarrow R$  is the difference between the support for  $L \cup R$  (i.e. the items in  $L$  and  $R$  occurring together in the database) and the support that would be expected if  $L$  and  $R$  were independent.

$$\text{leverage}(L \rightarrow R) = \text{support}(L \cup R) - \text{support}(L) \times \text{support}(R)$$



# Literature

- Max Bramer: Principles of data mining (2007)
  - 13. Association Rule Mining II
- What is the "true story" about using data mining to identify a relation between sales of beer and diapers? <http://www.dssresources.com/newsletters/66.php>

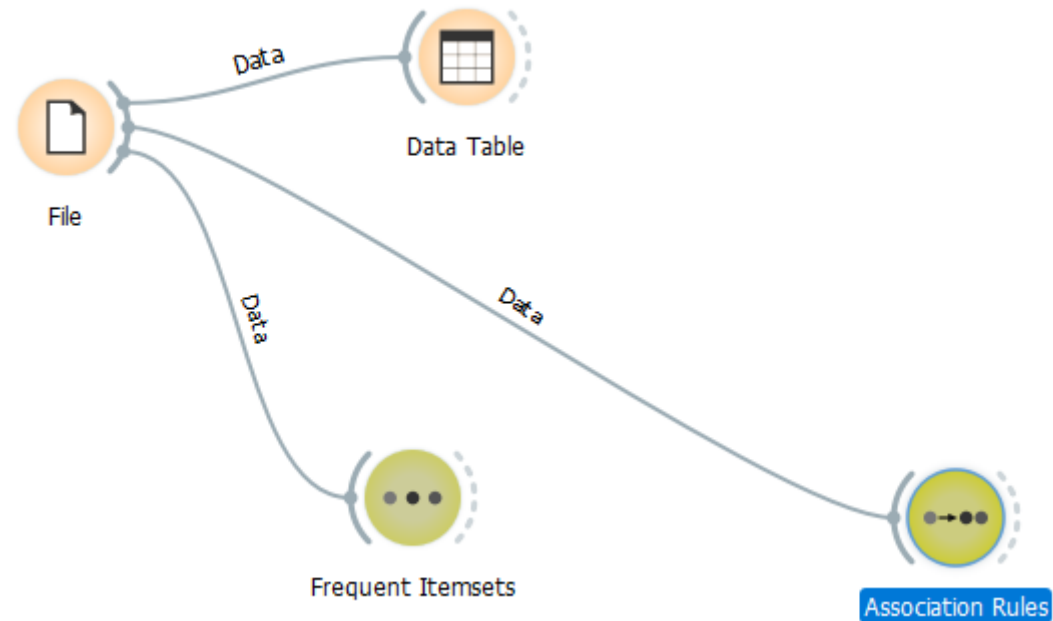
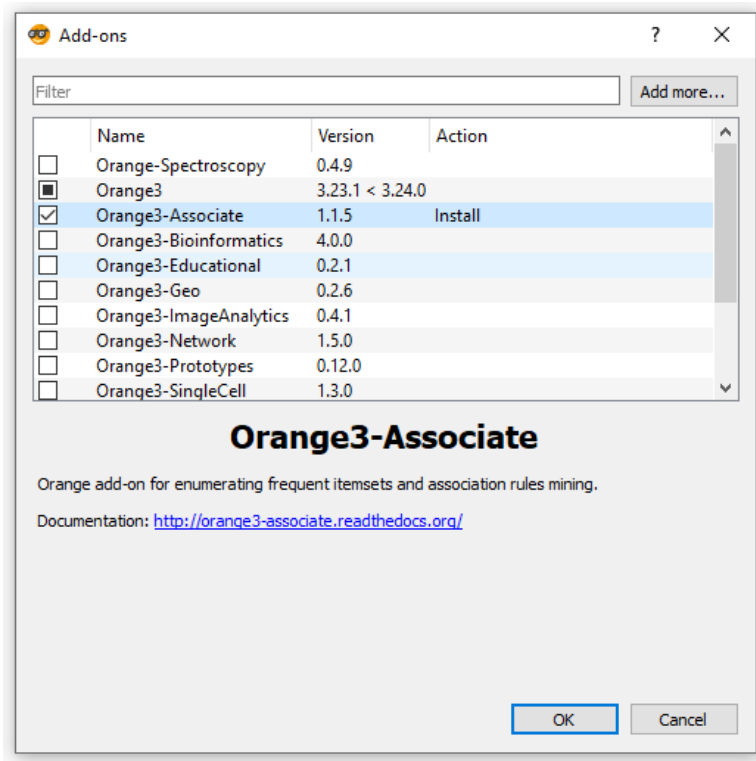
# Homework

1. Transformation of an attribute-value dataset to a transaction dataset.
2. What would be the association rules for a dataset with two items A and B, each of them with support 80% and appearing in the same transactions as rarely as possible?
  - a. minSupport = 50%, min conf = 70%
  - b. minSupport = 20%, min conf = 70%
3. What if we had 4 items: A,  $\neg$ A, B,  $\neg$  B
4. Compare decision trees and association rules regarding handling an attribute like "PersonID". What about attributes that have many values (eg. Month of year)

A	B
Green	White
Green	White
Green	Blue
Green	Blue
Green	Blue
Green	Blue
Green	Blue
Green	Blue
White	Blue
White	Blue

# Association rules: Orange workflow

## 1. Install Add-on Orange3-Associate



\* Start with a small minSupport and we increase it gradually (to avoid running out of memory)

# Association rules quality measures in Orange

Supp	Conf	Covr	Strg	Lift	Levr	Antecedent	Consequent
0.050	0.178	0.283	0.618	1.017	0.001	Fresh Vegetables	Fresh Fruit
0.050	0.287	0.175	1.619	1.017	0.001	Fresh Fruit	Fresh Vegetables

- **support, confidence, lift, leverage**
- **coverage:** how often antecedent items are found in the data set (support of antecedent/data)
- **strength:** (support of consequent/support of antecedent)

# Lab exercise

## Datasets

- <https://biolab.si/core/foodmart.basket>
- [https://github.com/digizeph/data\\_mining/blob/master/data/FoodMart.csv](https://github.com/digizeph/data_mining/blob/master/data/FoodMart.csv)
- <http://file.biolab.si/datasets/voting.tab>

1. Compare the two datasets (files)
2. Generate frequent itemsets and association rules for both datasets. What is the difference?
3. Frequent itemsets and association rules for „Voting.tab“