

HELSINGIN YLIOPISTO HELSINGFORS UNIVERSITET UNIVERSITY OF HELSINKI

#### 582364 Data mining, 4 cu Lecture 3: Association analysis

Spring 2010 Lecturer: Juho Rousu Teaching assistant: Taru Itäpelto





### **Frequent Pattern Discovery Application: Recommendation systems**

- Many recommendation web sites are based on collecting data on frequently co-occurring items
- If you liked the film "Aliens" you probably like "Avatar"
- People that buy book X, frequently buy book Y

IMDD	Search	All		¢						Go
Internet Movie Database		Movi	es -	TV -	News -	Videos	- C	ommu	ınity -	IMD
SCARS	10>	<u>c</u>	Oscar I	<u>Photos</u>	List of Winn	<u>iers Sh</u>	iow Rec	ар	6	
				New	s Flash » Ac	tor Co	rey Ha	aim I	s Dea	ad at 3
<u>Db</u> > <u>Avatar (2009)</u> > <b>Reco</b>	mmendatio	ons								
	Recom	menc	lations	s for						
	Ava	ntar	(200	)9) More	at IMDbPro »					
			(20)	<u>,</u>	activity of the w					
LIE F	How do t	nese re	ecomm	endations	WORK?					
	Sugge	sted b	by the	databa	se	Shov	Showtimes		lable	User
						(US	only)	@Am	nazon	Rating
AVATAR	Aliens (	1986)				8	à	₽V₽	WHS	8.5
Own the rights?	Alien: R	lesurrec	ction (1	997)				200	WHS	6.1
Own the rights?	The Lor	d of the	Rings	: The Ret	urn of the King		ň	DVD	WHS	8.8
Buy it at Amazon	(2003)									
	The Tra	nsform	ers: Th	e Movie (	1986)			₽V₽	WHS	7.2
Discuss in Boards	Star Wa	ars: Epi	sode II	I - Reven	ge of the Sith		ă	DVD	WHS	7.9
More at IMDb Pro	(2005)					_				
Add to My Movies	Termina	tor Sal	vation	(2009)				BVD BVD		6.9
Update Data	Despera	ate Livir	ng (197	7)				BVD	WHS	7.0
de biliete	The Lor	d of the	Rings	: The Two	Towers (2002)			BVD	WHS	8.7
	Star Wa	urs (197	7)					₽V₽	WHS	8.8
recommendations	Serenity	(2005)	)			63	ě	NO.	WHS	8.0
op Links	т	<b>ip</b> : if you	want to	see if a mov	ie is showing in a c	inema near y	ou, click th	ne film rol	I. (USA or	nly)
collers and uidean										

#### Frequently Bought Together



Total List Price: \$263.00 Price For All Three: \$201.66 Add all three to Cart Add all three to Wish List

Show availability and shipping details

This item: Introduction to Data Mining by Pang-Ning Tan

- Data Mining: Practical Machine Learning Tools and Techniques, Second Edition (Morgan Kaufmann Series in Data Management Systems) by Eibe Frank
- The Elements of Statistical Learning: Data Mining, Inference, and Prediction, Second Edition (Springer Series in Statistics) by Robert Tibshirani



Goal: Given a set of transactions, find

items that occur frequently together (Frequent itemsets)

- "Introduction to Data Mining" & "Elements of Statistical Learning " are frequently bought together
- rules that will predict the occurrence of an item based on the occurrences of other items in the transaction (Association rules)
  - People that bought "Introduction to Data Mining", often buy "Elements of Statistical Learning " as well

#### Frequently Bought Together



 Total List Price: \$263.90

 Price For All Three: \$201.66

 Image: Add all three to Cart

 Add all three to Cart

 Show availability and shipping details

This item: Introduction to Data Mining by Pang-Ning Tan

- Data Mining: Practical Machine Learning Tools and Techniques, Second Edition (Morgan Kaufmann Series in Data Management Systems) by Eibe Frank
- The Elements of Statistical Learning: Data Mining, Inference, and Prediction, Second Edition (Springer Series in Statistics) by Robert Tibshirani



- A collection of one or more items
  - Example: {Milk, Bread, Diaper}
- k-itemset
  - An itemset that contains k items
- Support
  - Support count (σ): Count of occurrences of an itemset
    - E.g.  $\sigma(\{Milk, Bread, Diaper\}) = 2$
  - Support: Fraction of transactions that contain an itemset
    - E.g. s({Milk, Bread, Diaper}) = 2/5

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



#### Frequent Itemset

- An itemset whose support is greater than or equal to a *minsup* threshold
- Example:
  - We set minsup = 0.5
  - Frequent itemsets:
    - Bread (support = 0.8 > 0.5)
    - Milk (support = 0.8)
    - Diaper (support = 0.8)
    - Beer (support = 0.6)
    - {Bread, Milk} (support = 0.6)
    - {Bread, Diaper} (support = 0.6)
    - {Beer, Diaper} (support = 0.6)
    - {Milk, Diaper} (support = 0.6)

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



An expression of the form

 $X \rightarrow Y$ ,

where X and Y are itemsets

Semantics: "When X happens, Y frequently happens as well"

Examples:

 ${Milk, Diaper} \rightarrow {Beer}$ 

{"Introduction to Data Mining"} → {"Elements of Statistical Learning"}

Different from

- Iogical implication: "When X happens, Y always happens as well"
- causal relation: "X causes Y to happen"



#### **Definition: Association Rule**

The strength of an association rule X → Y

is measured by its *support* and *confidence* 

Support s(X → Y): Fraction of transactions that contain both the set X and the set Y

■  $s(X \rightarrow Y) = \sigma(X \cup Y)/|T|$ 

■ Confidence c(X → Y): how often items in Y appear in transactions that contain X

• 
$$c(X \rightarrow Y) = \sigma(X \cup Y)/\sigma(X)$$

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:

 $\{Milk, Diaper\} \Rightarrow Beer$ 

$$s = \frac{\sigma(\text{Milk}, \text{Diaper}, \text{Beer})}{|\mathsf{T}|} = \frac{2}{5} = 0.4$$
$$c = \frac{\sigma(\text{Milk}, \text{Diaper}, \text{Beer})}{\sigma(\text{Milk}, \text{Diaper})} = \frac{2}{3} = 0.67$$



Given a set of transactions T, the goal of association rule mining is to find all rules having

- support ≥ *minsup* threshold
- confidence ≥ *minconf* threshold

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



{Bread}→{Milk} (s=0.6,c=0.75) {Bread}→{Diaper} (s=0.6,c=0.75)

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

# Association Rule Mining: brute-force approach

- Brute-force approach:
  - List all possible association rules
  - Compute the support and confidence for each rule
  - Prune rules that fail the *minsup* and *minconf* thresholds
- How much time this would take?
  - For d unique items, there are 3<sup>d</sup>-2<sup>d+1</sup>+1 possible rules
    - For d = 100, this gives ca.  $5 \times 10^{47}$  possible rules to check against the database
    - Compare to age of universe ca. 5x10<sup>17</sup> seconds
- Will not work except for toy examples



#### **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 of Rules:

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

#### **Observations:**

- All the above rules are binary partitions of the same itemset: {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 ≥ minsup
      - e.g. {Milk, Diaper, Beer}
  - 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
  - There are 2<sup>d</sup> itemsets, which is still too much to enumerate



Reduce the number of candidate itemsets (M)

- Complete search: M=2<sup>d</sup>
- Use pruning techniques to reduce M
- Reduce the number of comparisons (NM)
  - Use efficient data structures to store the candidate itemsets or transactions
  - No need to match every candidate against every transaction



#### Reducing the number of candidate itemsets: Apriori principle

#### Apriori principle:

- If an itemset is frequent, then all of its subsets must also be frequent
- i.e. if {Milk,Diaper,Beer} is frequent, then
- {Milk,Diaper}, {Milk,Beer}, {Diaper,Beer},{Milk},{Diaper},{Beer} must also be frequent
- Why this is true (informally):
  - In every transaction, subsets always occur if the whole set occurs
  - Support of the itemset is given by the sum of occurrences over the transactions
- Converse does not hold:
  - Even though all subsets are frequent, an itemset may be infrequent
  - Need to check against the transactions to find the true support



#### Reducing the number of candidate itemsets: Apriori principle

#### Apriori principle:

- If an itemset is frequent, then all of its subsets must also be frequent
- Formally, the principle follows from the anti-monotone property of the support function:
  - When we put more items into the itemset, the support can only decrease

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



- Itemsets that can be constructed from a set of items have a *partial order* with respect to the subset operator
  - i.e. a set is larger than its proper subsets
- This induces a lattice where nodes correspond to itemsets and arcs correspond to the subset relation
- The lattice is called the *itemset* lattice
- For d items, the size of the lattice is 2<sup>d</sup>





#### Frequent itemsets on the itemset lattice

- The Apriori principle is illustrated on the itemset lattice
  - The subsets of a frequent itemset are frequent
  - They span a sublattice of the original lattice (the grey area)



Figure 6.3. An illustration of the *Apriori* principle. If  $\{c, d, e\}$  is frequent, then all subsets of this itemset are frequent.



#### Frequent itemsets on the itemset lattice

#### Conversely

- The supersets of an infrequent itemset are infrequent
- They also span a sublattice of the original lattice (the crossed out nodes)
- If we know that {a,b} is infrequent, we never need to check any of the supersets
  - This fact is used in supportbased pruning





Y)

### Frequent items set generation in Apriori Algorithm

- Input: set of items I, set of transactions T, number of transactions N, minimum support *minsup*
- Output: frequent k-itemsets F<sub>k</sub>, k=1,...
- Method:
  - k=1
  - Compute support for each 1-itemset (item) by scanning the transactions
  - $F_1$  = items that have support above *minsup*
  - Repeat until no new frequent itemsets are identified
    - 1.  $C_{k+1}$  = candidate *k*+1 -itemsets generated from length k frequent itemsets  $F_k$
    - 2. Compute the support of each candidate in  $C_{k+1}$  by scanning the transactions T
    - 3.  $F_{k+1}$  = Candidates in  $C_{k+1}$  that have support above *minsup*.



### **Illustrating Apriori Principle**

Item	Count	Items (1-itemsets)					
Bread	4						
Coke	2				_		
Milk	4		Itemset	Count	Pairs (2-item	nsets)	
Beer	3		{Bread,Milk}	3			
Diaper	4		{Bread,Beer}	2	(No need to generate		
Eggs 1			{Bread,Diaper}	3	candidates involving Coke		
{		{Milk,Beer}	2				
{			{Milk,Diaper}	3			
			{Beer,Diaper}	3			
Minimum Support = 3					- 	. (2.1)	. 、
					Iriple	ts (3-items	ets)
TE average		a al a va a	<b>_</b>	Itomoot		Count	I
If every subset is considered,			-	nemsei	Count		
${}^{6}C_{1} + {}^{6}C_{2} + {}^{6}C_{3} = 41$			L	{Bread,Milk,Diaper} 3			
With support-based pruning,							
6 + 6 + 1 = 13				••••			
					••		



#### **Characteristics of Apriori algorithm**

- Level-wise search algorithm:
  - traverses the itemset lattice level-by-level (1-itemsets, 2itemsets ...)
- Generate-and-test strategy:
  - new candidate itemsets are generated from smaller frequent itemsets
  - support is tested to weed out infrequent itemsets





#### **Requirements for candidate generation step**

- It should not generate too many unnecessary candidates, i.e. itemsets where at least one of the subsets is infrequent
- It should ensure the completeness of the candidate set,
   i.e. no frequent itemset is left out
- 3. It should not generate the same candidate set more than once
  - e.g. {a,b,c,d} can be generated by merging {a,b,c} with {d}
     or {b,d} with {a,c}, {a,b} with {c,d}

Candidate generation strategies: F<sub>k-1</sub> x F<sub>1</sub> method

#### F<sub>k-1</sub> x F<sub>1</sub> method: Combine frequent k-1 – itemsets with frequent 1itemsets



**Figure 6.7.** Generating and pruning candidate k-itemsets by merging a frequent (k - 1)-itemset with a frequent item. Note that some of the candidates are unnecessary because their subsets are infrequent.

# Candidate generation strategies: F<sub>k-1</sub> x F<sub>1</sub> method

Satisfaction of our requirements (#1-3):

- while many k-itemsets are left ungenerated, can still generate unnecessary candidates #
  - e.g. merging {Beer, Diapers} with {Milk} is unnecessary, since {Beer, Milk} is infrequent
- 2. method is complete: each frequent itemset consists of a frequent k-1 –itemset and a frequent 1-itemset  $\checkmark$

## Candidate generation strategies: F<sub>k-1</sub> x F<sub>1</sub> method

- 3. can generate the same set twice (\*)
  - e.g. {Bread, Diapers, Milk} can be generated by merging {Bread,Diapers} with {Milk} or {Bread,Milk} with
     {Diapers} or {Diapers, Milk} with {Bread}
  - This can be circumvented by keeping all frequent itemsets in their lexicographical order (
    - e.g. {Bread,Diapers} can be merged with {Milk} as 'Milk' comes after 'Bread' and 'Diapers' in lexicographical order
    - {Diapers, Milk} is not merged with {Bread}, {Bread, Milk} is not merged with {Diapers} as that would violate the lexicographical ordering

Candidate generation strategies: F<sub>k-1</sub> x F<sub>k-1</sub> method

#### F<sub>k-1</sub> x F<sub>k-1</sub> method: Combine a frequent k-1 –itemset with another frequent k-1 -itemset



**Figure 6.8.** Generating and pruning candidate k-itemsets by merging pairs of frequent (k-1)-itemsets.



## Candidate generation strategies: F<sub>k-1</sub> x F<sub>k-1</sub> method

- Items are stored in lexicographical order in the itemset
   When considering merging, only pairs that share first k-2 items are considered
  - e.g. {Bread, Diapers} is merged with {Bread,Milk}
  - if the pairs share fewer than k-2 items, the resulting itemset would be larger than k, so we do not need to generate it yet
- The resulting k-itemset has k subsets of size k-1, which will be checked against support threshold
  - The merging ensures that at least two of the subsets are frequent
  - An additional check is made that the remaining k-2 subsets are frequent as well

Data mining, Spring 2010 (Slides adapted from Tan, Steinbach Kumar)



Satisfaction of our requirements (#1-3):

- 1. avoids the generation of many unnecessary candidates that are generated by the  $F_{k-1} \times F_1$  method
  - e.g. will not generate {Beer, Diapers, Milk} as {Beer, Milk} is infrequent
- method is complete: every frequent k-itemset can be formed of two frequent k-1 –itemsets differing in their last item.
- 3. each candidate itemset is generated only once



- Given the candidate itemsets  $C_k$  and the set of transactions T, we need to compute the support counts  $\sigma(X)$  for each itemset X in  $C_k$
- Brute-force algorithm would compare each transaction against each itemset → large amount of comparisons
- An alternative approach
  - divide the candidate itemsets  $C_k$  into buckets by using a hash function
  - for each transaction t:
    - hash the itemsets contained in t into buckets using the same hash function
    - compare the corresponding buckets of candidates and the transaction
    - increment the support counts of each matching candidate itemset
  - A hash tree is used to implement the hash function

Generate Hash Tree

Suppose you have 15 candidate itemsets of length 3:

 $\{1 \ 4 \ 5\}, \{1 \ 2 \ 4\}, \{4 \ 5 \ 7\}, \{1 \ 2 \ 5\}, \{4 \ 5 \ 8\}, \{1 \ 5 \ 9\}, \{1 \ 3 \ 6\}, \{2 \ 3 \ 4\}, \{5 \ 6 \ 7\}, \{3 \ 4 \ 5\}, \{3 \ 5 \ 6\}, \{3 \ 5 \ 7\}, \{6 \ 8 \ 9\}, \{3 \ 6 \ 7\}, \{3 \ 6 \ 8\}$ 

You need:

• Hash function e.g.  $h(p) = p \mod 3$ 

• Max leaf size: max number of itemsets stored in a leaf node (if number of candidate itemsets exceeds max leaf size, split the node)



#### Matching the transaction against candidates





#### Matching the transaction against candidates





#### Matching the transaction against candidates





- From a frequent itemset, we still need to generate the association rules
- From each k-itemset, one can produce 2<sup>k</sup>-2 association rules
  - e.g. {Beer,Bread,Diapers} generates {Beer} → {Bread,Diapers}, {Bread}
     → {Beer, Diapers}, Diapers → {Beer,Bread}, {Beer,Bread} -> {Diapers},
     {Beer,Diapers} → {Bread}, {Bread, Diapers} → {Beer}
- All of the association rules generated from the same frequent itemset have the same support
- The confidence of the rules will be different, however
- We want to find efficiently the rules that have high confidence



- The confidence of any association rule can be computed from the support counts of the frequent itemsets:
  - c(Beer -> Bread, Diapers) =  $\sigma$ (Beer, Bread, Diapers)/ $\sigma$ (Beer)
  - We don't need to scan the transactions to find the high-confidence rules
- The confidence does not have a similar anti-monotone property as support has:
  - e.g. if c({Beer,Milk} → {Bread,Diapers}) exceed the confidence threshold *minconf*, it does not follow that {Beer} → {Bread} satisfies the confidence threshold



- However, between rules generated from the same frequent itemset we have the following property: if X → Y-X does not satisfy the confidence threshold then no rule X' → Y-X', where X' is a subset of X, satisfies the confidence threshold
  - e.g. if {Beer,Bread} → {Milk} does not satisfy the confidence threshold,
     {Beer} → {Bread,Milk} and {Bread} → {Beer,Milk} also do not
- This property follows from the anti-monotone property of the support count  $\sigma(X') \ge \sigma(X)$ , thus

 $c(X \rightarrow Y-X) = \sigma(X \cup Y-X)/\sigma(X) \ge \sigma(X' \cup Y-X')/\sigma(X') = c(X' \rightarrow Y-X')$ 



- Apriori uses level-wise search for rule generation
- It starts from empty right-hand side and all items in the lefthand side
- To generate a rule in the next level it merges the left-hand sides of two confident rules on the previous level
- When a non-confident rule is found, an entire subgraph (grey area) is pruned



Figure 6.15. Pruning of association rules using the confidence measure.

