Association Analysis IE500 Data Mining





Outline



- What is Association Analysis?
- Frequent Itemset Generation
- Rule Generation
- Interestingness Measures
- Handling Continuous and Categorical Attributes
- Subgroup Discovery

Association Analysis



- First algorithms developed in the early 90s at IBM by Agrawal & Srikant
- Motivation
 - Availability of barcode cash registers





Association Analysis



- Initially used for Market 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
 - Life science
 - Social science
 - Web usage mining

Simple Approaches



- To find out if two items x and y are bought together, we can compute their correlation
- E.g., Pearson's correlation coefficient:

predicted value p_i actual value a_i

$$PCC = \frac{\sum_{i=1}^{n} (p_i - \bar{p}) * (a_i - \bar{a})}{\sqrt{\sum_{i=1}^{n} (p_i - \bar{p})^2} * \sqrt{\sum_{i=1}^{n} (a_i - \bar{a})^2}}$$

- Numerical coding:
 - 1: item was bought
 - 0: item was not bought
- \bar{p} average of p (i.e., how often x was bought)

Correlation Analysis in Python



• e.g., using Pandas:



University of Mannheim | IE500 Data Mining | Association Analysis | Version 06.04.2025

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.
- Examples of Association Rules
 - {Diaper} → {Beer} {Beer, Bread} → {Milk} {Milk, Bread} → {Eggs, Coke}



University of Mannheim | IE500 Data Mining | Association Analysis | Version 06.04.2025



TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Break, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke
S	hopping Transactions

Definition: Frequent Itemset



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
S	hopping Transactions

• 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. σ ({Milk, Bread, Diaper}) = 2
- Support (s)
 - Fraction of transactions that contain an itemset
 - e.g. s({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

Definition: Association Rule

- Association Rule
 - An implication of the form $X \rightarrow Y$, where X and Y are itemsets
 - Interpretation: when X occurs,Y occurs with a certain probability



TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Break, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke
S	honning Transactions

More formally, it's a *conditional probability* P(Y|X) given X, what is the probability of Y?

Definition: Association Rule



	TID	Items		
	1	Bread, Milk		
	2	Bread, Diaper, Beer, Eggs		
	3	Milk, Diaper, Beer, Coke		
	4	Break, Milk, Diaper, Beer		
	5	Bread, Milk, Diaper, Coke		
$s(X \to Y) = \frac{ X \cup Y }{ T }$ Shopping Transactions $\sigma(\{Milk, Diaper, Beer\})$				
$s({Milk, Diaper} \rightarrow {$	Beer)	$= \frac{1}{ T }$ $= \frac{2}{5} = 0.4$		
$c(X \to Y) = \frac{\sigma(X \cup Y)}{\sigma(X)}$				
$c(\{Milk, Diaper\} \rightarrow \{$	[Beer)]	$= \frac{\sigma(\{Milk, Diaper, Beer\})}{\sigma(\{Milk, Diaper\})}$ $= \frac{2}{3} = 0.67$		

Association Rule

- Example: {Milk, Diaper} \rightarrow {Beer} Condition Consequent
- **Rule Evaluation Metrics**
 - Support s: Fraction of total transactions which contain both X and Y
 - Confidence **c**: Measures how often items in Y appear in transactions that contain X

The Association Rule Mining Task



- Given a set of transactions T,
 - the goal of association rule mining is to find all rules having
 - support ≥ *minsup* threshold
 - confidence ≥ *minconf* threshold
- *minsup* and *minconf* are provided by the user
- Brute Force Approach:
 - List all possible association rules
 - Compute the support and confidence for each rule
 - Remove rules that fail the *minsup* and *minconf* thresholds

⇒ Computationally prohibitive due to large number of candidates!

Mining Association Rules

• Example 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)



TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Break, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke
S	nopping Transactions

- Observations:
 - All the above rules are binary partitions of the same itemset: {Milk, Diaper, Beer}
 - Rules originating from the same itemset have identical support s
- $s(X \to Y) = \frac{|X \cup Y|}{|T|}$

- but can have different confidence
- Thus, we may decouple the support and confidence requirements

University of Mannheim | IE500 Data Mining | Association Analysis | Version 06.04.2025

Apriori Algorithm: Basic Idea



- Two-step approach:
 - 1. Frequent Itemset Generation
 - Generate all itemsets whose support ≥ 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

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 ~ O(NMw) → Expensive since M = 2^d

Brute-force Approach



amazon

- Amazon sells 12M different products (as of 2023)
 - That is 2^{12.000.000} possible itemsets
 - That's a 3.6M digit number
 - Today's supercomputers: 1,200 Petaflops,
 i.e., 1.2x 10¹⁸ floating point operations per second
 - Even if an itemset could be checked with one single floating point operation this would take ~ 10^{3,612,334} years (age of universe: 1.4x10¹⁰ years)
- However:
 - Most itemsets will not be frequent 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 University of Mannheim | IE500 Data Mining | Association Analysis | Version 06.04.2025



Anti-Monotonicity of Support

- What happens when an itemset gets larger?
 - $s({Milk}) = 0.8$
 - s({Milk,Diaper}) = 0.6
 - s({Milk,Diaper,Beer}) = 0.4
 - $s({Bread}) = 0.8$
 - s({Bread,Milk}) = 0.6
 - s({Bread,Milk,Diaper}) = 0.4
- There is a pattern here!



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

Reducing the Number of Candidates



- There is a pattern here!
 - It is called anti-monitonicity of support
- If X is a subset of Y

s(Y) is at most as large as s(X)

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

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

- Consequence for frequent itemset search (aka Apriori principle):
 - If Y is frequent, X also has to be frequent
 - i.e.: all subsets of frequent itemsets are frequent

Using the Apriori Principle for Pruning



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



University of Mannheim | IE500 Data Mining | Association Analysis | Version 06.04.2025

The Apriori Algorithm



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

Illustrating the Apriori Principle

Minimum Support Count = 3



ltems (1-	itemsets	5)			TID	Items
Item	Count				1	Bread, Milk
Bread	4				2	Bread, Diaper, Beer, Eggs
Milk	4	No need to generat	te		3	Milk, Diaper, Beer, Coke
Beer	3	candidates involvin	g		4	Break, Milk, Diaper, Beer
Diaper	4	Coke or Eggs			5	Bread, Milk, Diaper, Coke
00-		Pairs (2-items	sets)	_		
		Itemset	Count			
		{Bread, Milk}	3			
		{Bread, Beer}	2	Ne		
		{Bread, Diaper}	3	INO	need	to generate
		{Milk, Beer}	2	can	didate	e {Milk, Diaper, Beer}
		{Milk, Diaper}	3	as o	count	{Milk, Beer} = 2
		{Beer, Diaper}	3			
					Trip	l ets (3-itemsets)
				Ite	mset	Count

University of Mannheim | IE500 Data Mining | Association Analysis | Version 06.04.2025

3

{Bread, Milk, Diaper}

Illustrating the Apriori Principle

- In the example, we had six items, and examined
 - Six 1-itemsets
 - Six 2-itemsets
 - One 3-itemset
 - i.e., 13 in total
- vs. possible itemsets: 2^6 = 64

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



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

From Frequent Itemsets to Rules

- Given a frequent itemset L, find all non-empty subsets f ⊂ L such that f → L \ f satisfies the minimum confidence requirement
- Example Frequent Itemset L:
 - {Milk,Diaper,Beer}
- Example Rule:

$$- [{Milk, Diaper}] \rightarrow {Beer}$$

$$f$$

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



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

Challenge: Combinatorial Explosion



• Given a 4-itemset {A,B,C,D}, we can generate

 $\begin{array}{ll} \{A\} \rightarrow \{B,C,D\}, & \{A,B\} \rightarrow \{C,D\}, & \{B,C\} \rightarrow \{A,D\}, & \{C,D\} \rightarrow \{A,B\}, & \{A,B,C\} \rightarrow \{D\} \\ \{B\} \rightarrow \{A,C,D\}, & \{A,C\} \rightarrow \{B,D\}, & \{B,D\} \rightarrow \{A,C\}, & \{A,B,D\} \rightarrow \{C\} \\ \{C\} \rightarrow \{A,B,D\}, & \{A,D\} \rightarrow \{B,C\}, & \{A,C,D\} \rightarrow \{B\} \\ \{D\} \rightarrow \{A,B,C\}, & \{B,C\}, & \{B,C\} \rightarrow \{A\} \end{array}$

- i.e., a total of 14 rules for just one itemset!

• General number for a k-itemset: $2^k - 2$

- It's not 2^k since we ignore $\emptyset \rightarrow \{...\}$ and $\{...\} \rightarrow \emptyset$

Challenge: Combinatorial Explosion

- Wanted: another pruning trick like Apriori
- However
 - c({Milk,Diaper} → {Beer}) = 0.67
 - c({Milk} → {Beer}) =0.5
 - c({Diaper} → {Beer}) =0.8
 - c(ABC \rightarrow D) can be larger or smaller than c(AB \rightarrow D)

University of Mannheim | IE500 Data Mining | Association Analysis | Version 06.04.2025

• In general, confidence does not have an anti-monotone property

25



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

Challenge: Combinatorial Explosion



•	But:
•	DUL.

confidence of rules generated from the **same itemset** has an anti-monotone property

- E.g. L = {Milk,Diaper,Beer}
- {Milk,Diaper,Beer} $\rightarrow \emptyset$ c=1.0
 - {Milk,Diaper} \rightarrow {Beer} c=0.67
 - {Milk} \rightarrow {Diaper,Beer} c=0.5
 - {Diaper} \rightarrow {Milk,Beer} c=0.5
- {Milk,Beer} \rightarrow {Diaper} c=1.0
 - {Milk} \rightarrow {Diaper,Beer} c=0.5
 - {Beer} \rightarrow {Milk,Diaper} c=0.67

- e.g., L = {A,B,C,D}: c(ABC → D) ≥ c(AB → CD) ≥ c(A → BCD)

University of Mannheim | IE500 Data Mining | Association Analysis | Version 06.04.2025

Observation: moving elements from antecedent to consequence ("left to right") in the rule never increases confidence!

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

Explanation



- Confidence is anti-monotone with respect to the number of items on the right-hand side (RHS) of the rule

 i.e., "moving elements from left to right" cannot increase confidence
- Reason:

$$c(AB \to C) \coloneqq \frac{s(ABC)}{s(AB)}$$
 $c(A \to BC) \coloneqq \frac{s(ABC)}{s(A)}$

- Due to anti-monotone property of support, we know $s(AB) \le s(A)$

Hence

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

Candidate Rule Pruning



• Moving elements from left to right cannot increase confidence



University of Mannheim | IE500 Data Mining | Association Analysis | Version 06.04.2025

Rule Generation for Apriori Algorithm



- Candidate rule is generated by merging two rules that share the same prefix in the rule consequent
 - join(CD → AB, BD → AC) would produce the candidate rule D → ABC
 - Prune rule D → ABC if one of its parent rules does not have high confidence (e.g. AD → BC)



- All the required information for confidence computation has already been recorded during itemset generation
 - Thus, there is no need to see the data any more

$$c(X \to Y) = \frac{s(X \cup Y)}{s(X)}$$

University of Mannheim | IE500 Data Mining | Association Analysis | Version 06.04.2025

Association Analysis in Python



- Various packages exist
 - In the exercise, we'll use the Orange3 package
 - Frequent Itemset Generation



Creating Association Rules



Interestingness Measures



- Association rule algorithms tend to produce too many rules
 - Many of them are uninteresting or redundant
 - Redundant if {A,B,C} → {D} and {A,B} → {D}
 have same support & confidence
- 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
 - We will have a look at one: Lift
 - See Tan/Steinbach/Kumar, Chapter 6.7

Drawback of Confidence

- Association Rule:
 Tea → Coffee
- confidence(Tea \rightarrow Coffee) = $\frac{3}{4} = 0.75$

• **but** support(Coffee) =
$$\frac{18}{20} = 0.9$$

- Although confidence is high, rule is misleading as the fraction of coffee drinkers is higher than the confidence of the rule
 - We want confidence($X \rightarrow Y$) > support(Y)
 - otherwise rule is misleading as X reduces probability of Y

Contingency table

Coffee

3

15

18

Tea

Tea

Coffee

1

1

2

4

16

20









- We discover a high confidence rule for tea \rightarrow coffee
 - 75% of all people who drink tea also drink coffee
 - Hypothesis: people who drink tea are likely to drink coffee
 - Implicitly: more likely than all people
- Test: Compare the confidence of the two rules

- Rule: Tea → coffee
$$c(tea → coffee) = \frac{s(\{tea\} \cup \{coffee\})}{s(\{tea\})}$$
- Default rule: all → coffee $c(all → coffee) = \frac{s(\{all\} \cup \{coffee\})}{s(\{all\})} = \frac{s(\{coffee\})}{1}$

 $= s(\{coffee\})$

• We accept a rule iff its confidence is higher than the default rule $\frac{c(tea \rightarrow coffee)}{c(all \rightarrow coffee)} = \frac{c(tea \rightarrow coffee)}{s(\{coffee\})} > 1$ Lift



• The lift of an association rule $X \to Y$ is defined as: $\frac{c(X \to Y) = \frac{s(X \cup Y)}{s(X)}}{Lift} = \frac{c(X \to Y)}{s(Y)} = \frac{s(X \cup Y)}{s(X) * s(Y)}$

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

Lift (Example)

Contingency table



Items Coffee

Coffee

Coffee Coffee

Coffee

Coffee

Coffee Coffee

Coffee

Coffee

Coffee Coffee

Coffee

Coffee Coffee

Tea

Bread

Tea, Coffee

Tea, Coffee Tea, Coffee

TID

1

2

3

4

5

6

7

8 9

10 11

12

13 14

15

16 17

18 19

20

•	Association Rule:	
	Tea \rightarrow Coffee	

	Coffee	Coffee		
Теа	3	1	4	
Теа	15	1	16	
	18	2	20	

• confidence(Tea \rightarrow Coffee) = $\frac{3}{4} = 0.75$

• **but** support(Coffee) =
$$\frac{18}{20} = 0.9$$

$$Lift(Tea \rightarrow Coffee) = \frac{c(tea \rightarrow coffee)}{c(all \rightarrow coffee)} = \frac{c(tea \rightarrow coffee)}{s(\{coffee\})}$$
$$= \frac{0.75}{0.9} = 0.833 < 1$$

lift < 1, therefore is negatively correlated and removed

Interestingness Measures



- There are lots of measures proposed in the literature
- Some measures are good for certain applications, but not for others
- Details: see literature (e.g., Tan et al.)

#	Measure	Formula
1	ϕ -coefficient	$\frac{P(A,B)-P(A)P(B)}{(P(A)P(B)(1-P(A))(1-P(B)))}$
2	Goodman-Kruskal's (λ)	$\frac{\sqrt{P(A)P(B)(1-P(A))(1-P(B))}}{\sum_{j}\max_{k}P(A_{j},B_{k})+\sum_{k}\max_{j}P(A_{j},B_{k})-\max_{j}P(A_{j})-\max_{k}P(B_{k})}$
3	Odds ratio (α)	$\frac{P(A,B)P(\overline{A},\overline{B})}{P(\overline{A},\overline{B})P(\overline{A},\overline{B})}$
4	Yule's Q	$\frac{P(A,B)P(\overline{AB})-P(A,\overline{B})P(\overline{A},B)}{P(A,B)P(\overline{AB})+P(A,\overline{B})P(\overline{A},B)} = \frac{\alpha-1}{\alpha+1}$
5	Yule's Y	$\frac{\sqrt{P(A,B)P(\overline{AB})} - \sqrt{P(A,\overline{B})P(\overline{A},\overline{B})}}{\sqrt{P(A,B)P(\overline{AB})} + \sqrt{P(A,\overline{B})P(\overline{A},\overline{B})}} = \frac{\sqrt{\alpha} - 1}{\sqrt{\alpha} + 1}$
6	Kappa (κ)	$\frac{P(A,B)+P(\overline{A},\overline{B})-P(A)P(B)-P(\overline{A})P(\overline{B})}{1-P(A)P(B)-P(\overline{A})P(\overline{B})}$
7	Mutual Information (M)	$\frac{\sum_{i}\sum_{j}P(A_{i},B_{j})\log \frac{1}{P(A_{i})P(B_{j})}}{\min(-\sum_{i}P(A_{i})\log P(A_{i}),-\sum_{j}P(B_{j})\log P(B_{j}))}$
8	J-Measure (J)	$\max\left(P(A,B)\log(\frac{P(B A)}{P(B)}) + P(A\overline{B})\log(\frac{P(\overline{B} A)}{P(\overline{B})}),\right.$
		$\left(P(A,B) \log(rac{P(A B)}{P(A)}) + P(\overline{A}B) \log(rac{P(\overline{A} B)}{P(A)}) \right)$
9	Gini index (G)	$\max\left(P(A)[P(B A)^{2}+P(\overline{B} A)^{2}]+P(\overline{A})[P(B \overline{A})^{2}+P(\overline{B} \overline{A})^{2}]\right)$
		$(-P(B)^2 - P(\overline{B})^2,$
		$P(B)[P(A B)^{2} + P(\overline{A} B)^{2}] + P(\overline{B})[P(A \overline{B})^{2} + P(\overline{A} \overline{B})^{2}]$
		$-P(A)^2 - P(\overline{A})^2$
10	Support (s)	P(A,B)
11	Confidence (c)	$\max(P(B A), P(A B))$
12	Laplace (L)	$\max\left(\frac{NP(A,B)+1}{NP(A)+2},\frac{NP(A,B)+1}{NP(B)+2}\right)$
13	Conviction (V)	$\max\left(rac{P(A)P(\overline{B})}{P(A\overline{B})}, rac{P(B)P(\overline{A})}{P(B\overline{A})} ight)$
14	Interest (I)	$\frac{P(A,B)}{P(A)P(B)}$
15	cosine (IS)	$\frac{\hat{P}(A,B)}{\sqrt{P(A)P(B)}}$
16	$\operatorname{Piatetsky-Shapiro's}\left(PS ight)$	P(A,B) - P(A)P(B)
17	Certainty factor (F)	$\max\left(rac{P(B A)-P(B)}{1-P(B)},rac{P(A B)-P(A)}{1-P(A)} ight)$
18	Added Value (AV)	$\max(P(B A) - P(B), P(A B) - P(A))$
19	Collective strength (S)	$\frac{P(A,B)+P(\overline{AB})}{P(A)P(B)+P(\overline{A})P(\overline{B})} \times \frac{1-P(A)P(B)-P(\overline{A})P(\overline{B})}{1-P(A,B)-P(\overline{AB})}$
20	Jaccard (ζ)	$\frac{P(A,B)}{P(A)+P(B)-P(A,B)}$
21	Klosgen (K)	$\sqrt{P(A,B)} \max(P(B A) - P(B), P(A B) - P(A))$

University of Mannheim | IE500 Data Mining | Association Analysis | version ub.04.2025

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:

{Number of Pages \in [5,10) \land (Browser=Firefox)} \rightarrow {Buy = 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
 - What if attribute has many possible values?
 - Many of the attribute values may have very low support
 - Potential solution: aggregate low-support attribute values
 - 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
 - {Refund=No, (Income=\$51,251)} → {Cheat=No}
 - {Refund=No, (60K<= Income <=80K)} → {Cheat=No}
 - {Refund=No, (0K<= Income <=1B)} → {Cheat=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

Subgroup Discovery



- Association Rule Mining:
 - Find all patterns in the data
- Classification:
 - Identify the best patterns that can predict a target variable
 - Those need not to be all
- Subgroup Discovery:
 - Find **all patterns** that can explain a target variable

Subgroup Discovery vs. Classification



- Example: learn to classify animals
 - Two possible models
 - has Trunk
 → Elephant (acc. 98%)



- has Trunk AND weight>3000kg AND color=grey AND height>2m
 → Elephant (acc 99%)
- Which one do you prefer?
 - Occam's Razor:
 - if you have two theories that explain a phenomenon equally well, choose the simpler one (has Trunk \rightarrow Elephant)
- What is our goal?
 - Classify animals at high accuracy
 - Learn as much about elephants (more general: the data) as possible

Subgroup Discovery – Algorithms



- Early algorithms (e.g., EXPLORA, MIDOS, 1999s)
 - Learn unpruned decision tree
 - Extract rules
 - Compute measures for rules, rate and rank
- Newer algorithms
 - Based on association rule mining (APRIORI-SD and others, 2000s)
 - Based on evolutionary algorithms (2000s)



- One of the most common metrics in Subgroup Discovery is WRAcc (Weighted Relative Accuracy), using probability of subgroup (S) and target (T)
 - WRAcc = P(ST) P(S)*P(T)

	Elephant	¬Elephant	
has Trunk AND weight>3000kg AND color=grey AND height>2m	1894	0	
¬(has Trunk AND weight>3000kg AND color=grey AND height>2m)	32	54874	



University of Mannheim | IE500 Data Mining | Association Analysis | Version 06.04.2025



- One of the most common metrics in Subgroup Discovery is WRAcc (Weighted Relative Accuracy), using probability of subgroup (S) and target (T)
 - WRAcc = P(ST) P(S)*P(T) = 0.033 0.033*0.034 = 0.032

	Elephant	¬Elephant	
has Trunk AND weight>3000kg AND color=grey AND height>2m	0.033	0.0	0.033
¬(has Trunk AND weight>3000kg AND color=grey AND height>2m)	0.0006	0.966	0.967
(0.034	0.966	



WRAcc = P(ST) - P(S)*P(T)

- Observations:
 - The higher P(ST), the more examples are covered
 - i.e., higher WRAcc means high coverage (like support)
 - The lower P(S) P(ST), the more accurate the subgroup
 - i.e., the higher P(ST)-P(S), the more accurate the subgroup
 - P(T) is a constant factor anyways, given a dataset
 - i.e., higher WRacc means higher accuracy
- Bottom line: WRacc represents both coverage and accuracy



WRAcc = P(ST) - P(S)*P(T)

- Observations:
 - If P(S) and P(T) are independent, P(ST) = P(S)*P(T), i.e., WRAcc = 0.0
 - Subgroup and target do not interact, this is not interesting
 - Best case:
 - P(ST) = P(S), i.e., no non-target examples covered by subgroup
 - P(ST) = P(T), i.e., no target examples not covered by subgroup
 - i.e., optimimum is P(T) P²(T)
 - Our elephant rule: P(ST) P(S)*P(T) = 0.033 0.033*0.034 = 0.032
 - Maximum WRacc: P(T) P(T)² = 0.034 0.034 ² = 0.032844
 - i.e., our rule is pretty good!

What's Next?



- Prof. Gemulla
 - HWS: Large-Scale Data Management, Machine Learning
 - FSS: Deep Learning
- Prof. Bizer
 - HWS: Web Data Integration, Large Language Models and Agents
 - FSS: Web Mining
- Prof. Stuckenschmidt
 - HWS: Decision Support
- Prof. Ponzetto
 - HWS: Information Retrieval and Web Search
 - FSS: Advanced Methods in Text Analytics
- Prof. Keuper
 - HWS: Higher Level Computer Vision, Image Processing
 - FSS:Generative Computer Vision Models
- Prof. Rehse
 - FSS: Advanced Process Mining

University of Mannheim | IE500 Data Mining | Association Analysis | Version 06.04.2025

Questions?





Literature for this Slideset



- Pang-Ning Tan, Michael Steinbach, Karpatne, Vipin Kumar: Introduction to Data Mining.
 2nd Edition. Pearson.
- Chapter 4: Association Analysis: Basic Concepts and Algorithms
- Chapter 7: Association Analysis: Advanced Concepts



