

# ADB Lecture 5

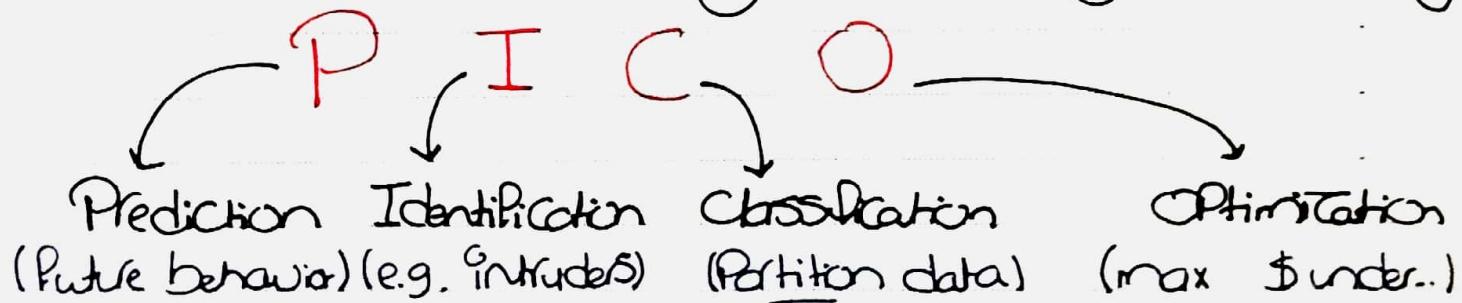
## • Data Mining

- ⇒ Discover Patterns or Rules from Large amounts of data. new info
- ⇒ Discover Interesting Structure in data
- ⇒ Automatically analyze & extract knowledge from data

## • Knowledge Discovery



## • Goals of data mining & Knowledge Discovery



# Types of Discovered Knowledge

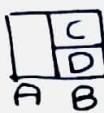
## Association Rules

- Correlate items with values or other items  
(e.g. buy milk  $\rightarrow$  buy eggs)

$$A \rightarrow B$$

## Classification Hierarchies

- goal is to create a hierarchy of classes



ABCD

m  
m

## Clustering

- Partition data into sets grouping similar items

## Time Series

- Detect similarities between items over time

## Sequential Patterns

- Predict based on sequence of events (e.g. Symptoms)

# 1. Association Rules

$\Rightarrow$  Can help businesses much

- They have big datasets of transactions (transaction id + date + list of items)
- Use it to make associations such as

98% customers who buy X, also buy Y

$\rightarrow$  helpful for marketing strategies, store layout, add-on sales, ...

## Notation

- $\rightarrow I = \{i_1, i_2, \dots, i_m\}$  set of all possible items
- $\rightarrow D = \{t_1, t_2, \dots, t_n\}$  dataset of n transactions
- each transaction has a subset of the items in I ( $t_j \subseteq I$ ) & has an ID (TID)

- An association rule Conveys Correlation between two item sets  $X, Y$  (i.e.  $X \neq Y$ )
    - \* Denoted by  $X \rightarrow Y$  where  $X, Y$  are disjoint.  $X \cap Y = \emptyset$
- e.g. {milk, bread}    e.g. {eggs, yoghurt}

→ The Strength of the association is measured by

• Confidence

$$C = \frac{\# \text{Tr. with } X \& Y}{\# \text{Tr. with } X}$$

• Support

$$S = \frac{\# \text{Tr. with } X \& Y}{n}$$

- $n$  is #transactions
- also applies for one set (e.g.  $SUP(X)$ )

Given

- Set of transactions  $D$
- minimum Support ( $minSup$ )
- minimum Confidence ( $minConf$ )

Wanted

- All association rules with support  $> minSup$  and confidence  $> minConf$

- This can be broken down into two subproblems
  1. Find "large" itemsets  
i.e., those with support exceeding minSup
  2. Find Association Rules
    - $X \cup Y$  must have high support (beyond minSup)
    - $X \rightarrow Y$  must have high confidence (beyond minConf)

# Appropriate algo. used for this.
- \* Association rules have nothing to do with functional dependencies
  - Functional dependency  $X \rightarrow Y$  for a relation R is a constraint that two tuples with same value of X, must have same value for Y
    - Implies  $XZ \rightarrow Y$  (also holds for each value of  $XZ$ )
    - also  $X \rightarrow Y, Y \rightarrow Z$  implies  $X \rightarrow Z$
  - Neither of these (concatenation or association) hold for association rules
    - e.g. If  $X \rightarrow Y$  has minSup,  $XA \rightarrow Y$  may not necessarily do as well (e.g. there's only one transaction  $\{AX\}$ , if it also has  $\{Y\}$  then supp is  $\frac{1}{n}$ )
    - If  $X \rightarrow A, A \rightarrow Z$  have minConf.,  $X \rightarrow Z$  may not (e.g. only have  $\{XA\}, \{AZ\}$  then Conf for  $XZ$  is zero)

# • A Priori Algorithm

→ Basic idea:

If Itemset AB doesn't surpass minSup then so will not ABC, ABCD, ... (A Priori Principle)

- A K-itemset means an Itemset of length K
- A large Itemset means an Itemset with supp exceeding minSup

- The algorithm uses two data structures

$L_K$

has all Itemsets of length K exceeding minimum support (i.e. large)

$L_K^{main}$

$C_K$

has all Itemsets of length K that may be large

- We assume each row in dataset is  $\langle tid, i_1, i_2, \dots \rangle$  with items sorted alphabetically.  
→ is okay if relation is normalized  
 $\langle tid, item \rangle$

↑  
transaction ↑  
items in it

## • Example

Tid items

10 A,C,D

20 B,C,E

30 A,B,C,E

40 B,E

,  $minSup_{\text{cont}} = 2$  (2/4)

| $C_1 =$ | ItemSet | SUPCart | • Will ignore this step later (directly $L_1$ ) |
|---------|---------|---------|---|
|         | A       | 2       |   |
|         | B       | 3       |   |
|         | C       | 3       |   |
|         | D       | 1       |   |
|         | E       | 3       |   |

|| no. of Transactions with B

- $C_1$  has all Possible 1-items
- Violation  $K_{minSup} = 2$

| $L_1 =$ | A | 2 | • $L_K$ is $C_K$ but keeps only large item sets<br>→ hence, D discarded |
|---------|---|---|---|
|         | B | 3 |   |
|         | C | 3 |   |
|         | E | 3 |   |

| $C_2 =$ | A, B | 1 | • To get Candidate 2-item sets, do $L_1 \bowtie L_1$<br>(Keeping only 2-item sets uniquely) |
|---------|------|---|---|
|         | A, C | 2 |   |
|         | A, E | 1 |   |
|         | B, C | 2 |   |
|         | B, E | 3 |   |
|         | C, E | 2 |   |

| $L_2 =$ | A, C | 2 | • Discarded Small Items<br>$\{A, B\}, \{A, E\}$ |
|---------|------|---|---|
|         | B, C | 2 |   |
|         | B, E | 3 |   |
|         | C, E | 2 |   |

- Now you think we should do  $L_2 \bowtie L_2$  to get  $C_3$
- Not exactly. Need to carefully define  $\bowtie$  first ( $K > 2$ )

$$\cdot L_{K-1} \bowtie L_{K-1} \Rightarrow C_K$$

length<sub>i</sub> = K-1

→ Do the join by unioning only two item sets ①  
that differ in only one item (Same K-2 items)

e.g. A, B, C → A, B, C, D  
A, B, D

& justification: If they differ in two items or more  
we don't get a K-item set ( $A, B, C \rightarrow A, B, C, D, E$ )  
 $A, D, E$  length ≠ 4

& if they don't differ at all, we get a K-1 item set

- hence, this will give us the item sets of length K exactly (what we need)

→ It turns out to be sufficient to ②

- union two K-1 item sets only if all elements are the same & they differ in the last element only.

③ ④

hence,

$$\begin{array}{c} A, C \\ B, C \\ B, E \\ C, E \end{array} \quad \begin{array}{c} A, C \\ B, C \\ B, E \\ C, E \end{array}$$

$L_2 \bowtie L_2$

| ItemSet         | SP |
|-----------------|----|
| $C_3 = B, C, E$ | 2  |

- to ignore duplicates (red arrow), can add condition that last element in 1st operand must be smaller than last element in 2nd operand (only black arrow lies)

- Clear as well that  $L_3 = \{B, C, E\}$
- ItemSet      Sup
- |           |   |
|-----------|---|
| $B, C, E$ | 2 |
|-----------|---|
- Nothing satisfies join condition for  $C_4$  (i.e.  $C_4 = \emptyset$ )  
and hence algorithm terminates
- Large item sets are: 4 item sets from  $L_1$ , 4 item sets from  $L_2$ , 1 itemset of  $L_3$

- In this example, we never had to do the Pruning Step

→ And what's that?

→ Suppose  $L_2 = \begin{array}{ll} \{A, B\} & 4 \\ \{A, C\} & 4 \\ \{A, E\} & 2 \\ \{B, C\} & 4 \\ \{B, D\} & 2 \\ \{B, E\} & 2 \end{array}$

then  $C_3 = \{A, B, C\}$

$\{A, B, E\}$

$\{A, C, E\}$

$\{B, C, D\}$

$\{B, C, E\}$

$\{B, D, E\}$

}

Do we have to count the support of each or can we do better

→  $C_3 = \text{Prune}(C_3)$

↑ use the A Priori Principle

- Since  $\{E, C\}$  wasn't in  $L_2 \rightarrow$  it didn't have minSup  $\rightarrow \{B, C, E\}$  and  $\{A, C, E\}$  won't have minSup

- Likewise  $\{E, D\}$  didn't have minSup  $\rightarrow \{B, D, E\}$  won't  
 $\{C, D\}$  didn't have minSup  $\rightarrow \{B, C, D\}$  won't

$C_3 = \{A, B, C, A, B, E\}$

} we only need to count minSup for these (then filter in L3)

→ In general, the Prune Step takes result from join and keeps only itemsets where all  $(K-1)$ -itemset subsets are in  $L_{K-1}$

e.g.  $A, B, C \rightarrow \begin{matrix} & & \\ \downarrow & \downarrow & \downarrow \\ A, B & B, C & A, C \end{matrix}$        $A, B, E \rightarrow \begin{matrix} & & \\ \downarrow & \downarrow & \downarrow \\ A, B & B, E & A, E \end{matrix}$  } all in  $L_2$

• Thus, our algorithm is

$L_1 = \{\text{large 1-itemsets}\}$

For ( $K=2 ; L_{K-1} = \emptyset ; K++$ )      || until  $L_{K-1}$  is empty

$C_K = \text{aPriori\_Gen}(L_{K-1})$       || generate  $K$ -itemsets  
 $L_{K-1} \rightarrow C_K$

For ( $t$  in transactions)      || For each candid.  
 For ( $c$  in SubSet( $C_K, t$ ))      ← Itemset in  $t$       Q.e.  
 $c.\text{Count}++;$

$L_K = \{c \in C_K \mid c.\text{Count} \geq \text{minSup}\}$       || include only those  $\geq \text{minSup}$

Return  $L_1 \cup L_2 \cup \dots \cup L_{K-2}$

## • A Priori Gen

1. Join Step ( $L_{k-1} \bowtie L_{k-1}$ )

\* go up 3 pages (things should make sense)

Insert into  $C_k$

Select  $P.i_1, P.i_2, \dots, P.i_{k-1}, Q.i_{k-1}$  } ABC  
From  $L_{k-1}$  as  $P$ ,  $L_{k-1}$  as  $Q$  ABD → ABCD

where  $P.i_1 = Q.i_1, \dots, P.i_{k-2} = Q.i_{k-2}, P.i_{k-1} < Q.i_{k-1}$

↳ respect alpha. order as well  
 $\leftarrow$

2. Prune Step

\* go up 2 pages

For each  $C$  in  $C_k$

For each  $k-1$  subset of  $C$

If subset not in  $L_{k-1}$

delete  $C$  from  $C_k$

Return  $C_k$

## • A Priori TID

→ just like A Priori but encodes itemsets

→ later passes, smaller # itemsets, smaller encodings & faster reads

• only need database to count supp in 1st pass ①

→ involves other modifications to make ① hold  
( $C_k$ )

## • Problems with APriori (why APrioriTID or Hybrid)

- APriori goes over entire dataSet every Pass  $\leftarrow \text{calc. support (Lk)}$
- #Transactions is so big (millions added per hour)
- Scanning db is expensive

- A Possible Improvement is to only check transactions that can possibly contain the K-itemSet

- \* Now we've solved the 1st Part of the SubProb.
  - generating all large item sets
  - hence, time to discover rules (associations)

For every large itemset P

For every non-empty subset of P Q

$\text{if } \text{sup}(P) / \text{sup}(P-Q) >$

Form rule  $P-Q \rightarrow Q$

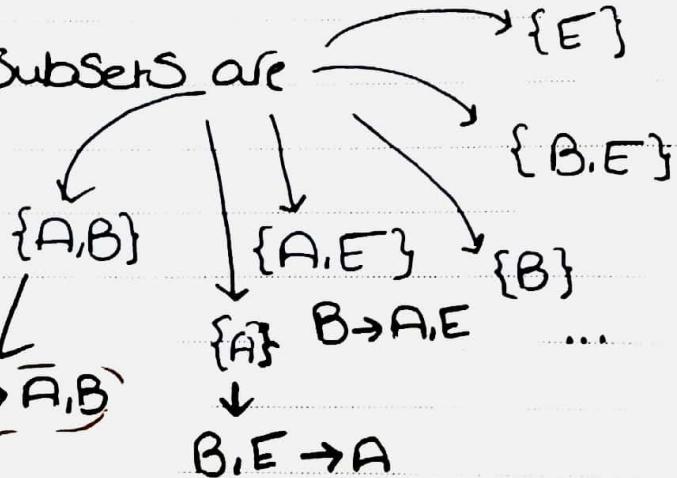
### • Example

→ Let item set P be {A, B, E}

. Then Possible Subsets are

Can be generated by  
DFS

accept if  $\frac{\text{sup}(A,B,E)}{\text{sup}(E)} < (\bar{E} \rightarrow \bar{A}, \bar{B})$   
> mincf  $\frac{\text{sup}(E)}{\text{sup}(A)}$



- In general, for  $K$ -ItemSet, we have  $2^K - 2$  Possible Subsets. (ignoring  $\emptyset$  & the Itemset Itself)
  - Means we'll have to do that much minConf Checks for each ItemSet
  - by the A Priori Principle, have all needed supports (no need for database Scanning)

» To Optimize:

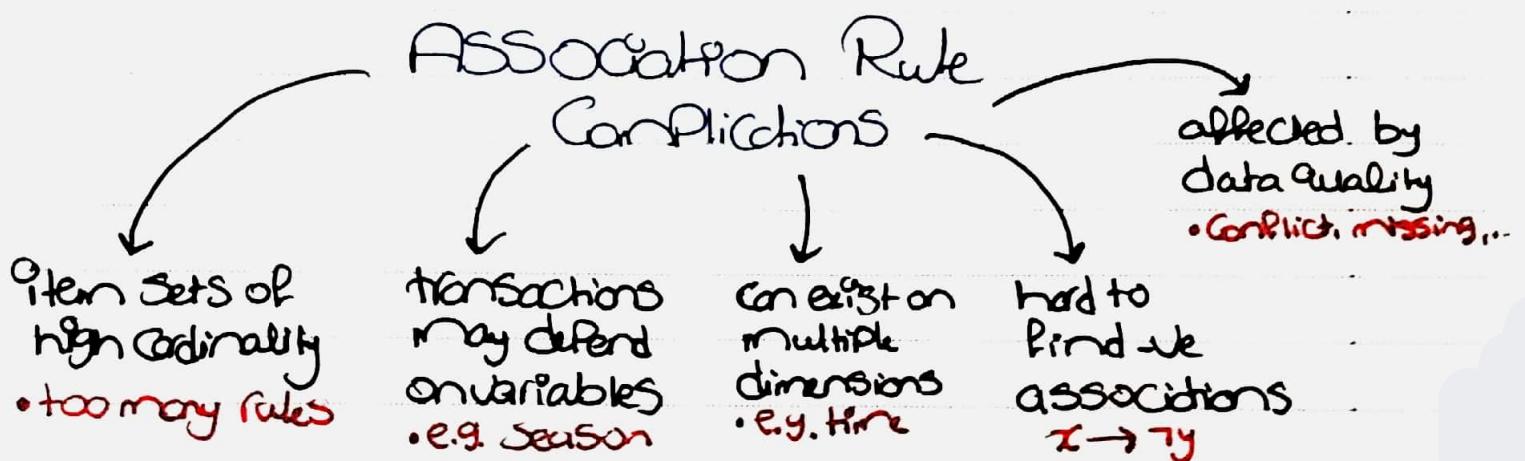
Observe that

If  $(P-q_{sub}) \rightarrow q_{sub}$  doesn't hold then so will not  $(P-q) \rightarrow q$

So we can get the rules smartly by

- Starting with 1-Item Consequents

Correspond to  $\{A\}, \{B\}, \{E\}$  resp. { }  $B, E \rightarrow A$  ← If this fails, any rule with Consequent A will also fail (i.e. ignore any  $q$  with A in it)  
 $A, E \rightarrow B$   
 $A, B \rightarrow E$   
 • No need to check  $B \rightarrow A, E$  or  $E \rightarrow A, B$



## 2. Classification

- Uses a training set & pre-classified data (Supervised learning)
- Produces a decision tree or set of rules to classify future items

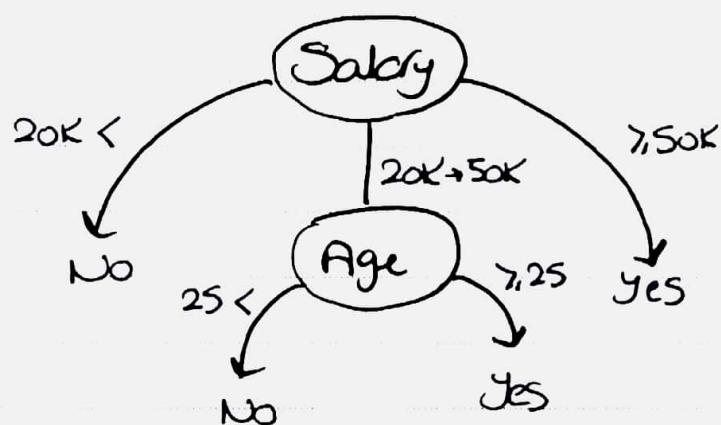
• NN Lec. 7

Sample Training data

| RID | Married | Salary | ACCT-balance | age | loan |
|-----|---------|--------|--------------|-----|------|
| 1   | no      | >50K   | <5K          | >25 | yes  |

, can be actual nos.

Class to Predict



} Example Decision Tree

- involves 4 rules you can extract.

## 3. Clustering

→ Goal is to place records in groups where the records in a group are highly similar (unsupervised learning)

• K-means

→ Start with initial cluster centers

Repeat:

|| until no change

•

Assign each example to closest center  
Recalculate centers as mean of cluster

## 4. Sequential Pattern Analysis

- Assume transactions ordered by time of purchase
- e.g.  $\{A, B, C\}$ ,  $\{B, E\}$ ,  $\{A, D, F\}$  ← **Sequence of Item Sets**
- Support of sequence is % of sequences of which **it's a Subsequence**
- Want sequences that exceed minSup  
e.g.  $S_1, S_2, S_3$  exceeding minSup means if customers buy  $S_1$ , they are likely to buy  $S_2$  then  $S_3$

## 5. Time Series Analysis

- Time Series = Sequence of events
  - e.g. Closing Price of Stock occurs every week day ( $Sun \rightarrow Tues$ )
    - Sequence of Values for a Stock → time series
    - Can use to Predict Future Trend for Stock.
  - Also helpful for **temporal data management**  
(e.g. Find the best quarter for the stock)

