

Multi-Dimensional Association Classification by Association

Cse352

ARTIFICIAL INTELLIGENCE

Professor Anita Wasilewska
Computer Science Department
Stony Brook University

Mining Multi-Dimensional Association

- Single-dimensional rules:

$\text{buys}(X, \text{“milk”}) \Rightarrow \text{buys}(X, \text{“bread”})$

- Multi-dimensional rules: ≥ 2 dimensions or predicates
Inter-dimension assoc. rules (*no repeated predicates*)

$\text{age}(X, \text{“19-25”}) \wedge \text{occupation}(X, \text{“student”}) \Rightarrow \text{buys}(X, \text{“coke”})$

Hybrid-dimension assoc. rules (*repeated predicates*)

$\text{age}(X, \text{“19-25”}) \wedge \text{buys}(X, \text{“popcorn”}) \Rightarrow \text{buys}(X, \text{“coke”})$

Example: Relational Data

Goal:

create multidimensional association rules

Student	Grade	Income	Buys
CS	High	Low	Milk
CS	High	High	Bread
Math	Low	Low	Bread
CS	Medium	High	Milk
Math	Low	Low	Bread

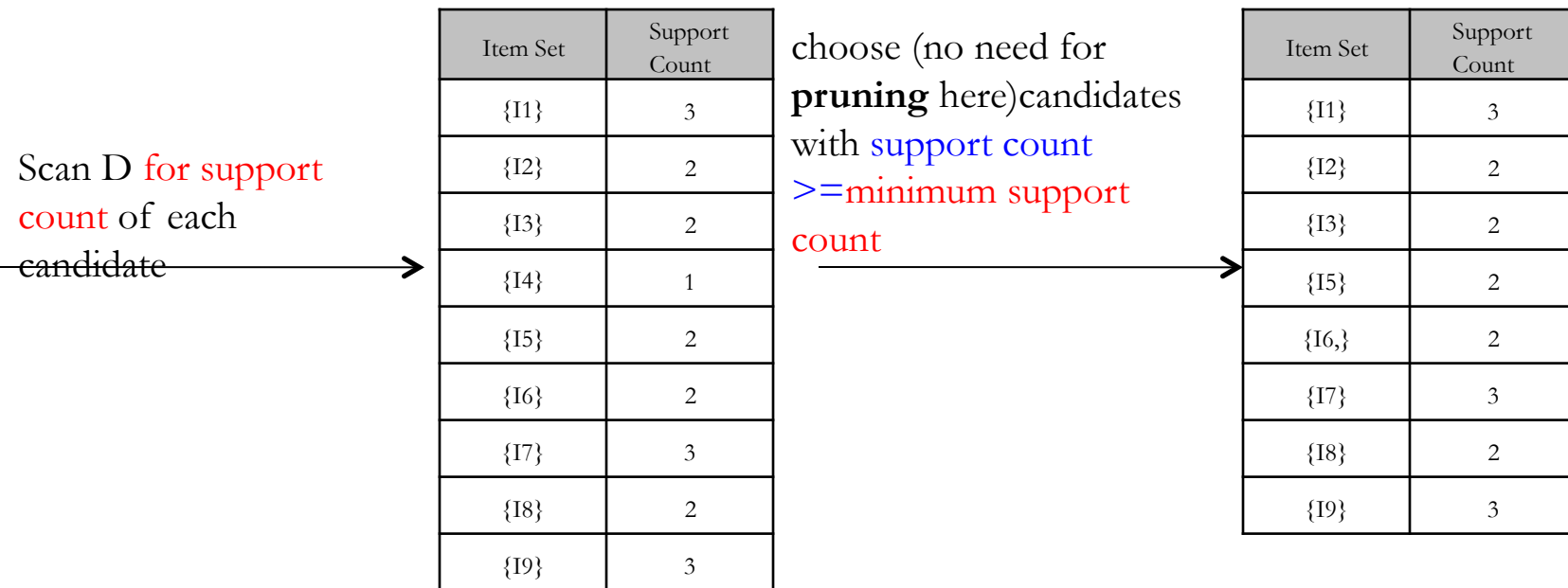
STEP 1: Data Conversion to Transaction and its count

Converted Data

Student = CS (I1)	Student =math (I2)	Grade = high (I3)	Grade =medium (I4)	Grade =low (I5)	Income =high (I6)	Income =low (I7)	Buys =milk (I8)	Buys =bread (I9)
+	-	+	-	-	-	+	+	-
+	-	+	-	-	+	-	-	+
-	+	-	-	+	-	+	-	+
+	-	-	+	-	+	-	+	-
-	+	-	-	+	-	+	-	+
3	2	2	1	2	2	3	2	3

Step 2: Apriori Algorithm

Generating 1-itemset Frequent Pattern



C1

L1

Let, the **minimum support count be 2**

Since we have 5 records \Rightarrow **minimum Support** = $2/5 = 40\%$

Let, **minimum confidence** required is **70%**

Generating 2-itemset Frequent Pattern

Generate C2 candidates from L1

Item Set
{1,12}
{1,13}
{1,14}
{1,15}
{1,16}
{1,17}
{1,18}
{1,19}
{2,13}
{2,14}
{2,15}
{2,16}
{2,17}
{2,18}
{2,19}
{3,14}
{3,15}
{3,16}
{3,17}
{3,18}
{3,19}
{4,15}
{4,16}
{4,17}
{4,18}
{4,19}
{5,16}
{5,17}
{5,18}
{5,19}
{6,17}
{6,18}
{6,19}
{7,18}
{7,19}
{8,19}

No need of pruning here-Scan D for count of each candidate

C2

Item Set	Support Count
{1,12}	0
{1,13}	2
{1,14}	1
{1,15}	0
{1,16}	2
{1,17}	1
{1,18}	2
{1,19}	1
{2,13}	0
{2,14}	0
{2,15}	2
{2,16}	0
{2,17}	2
{2,18}	0
{2,19}	2
{3,14}	0
{3,15}	0
{3,16}	1
{3,17}	1
{3,18}	1
{3,19}	1
{4,15}	0
{4,16}	1
{4,17}	0
{4,18}	1
{4,19}	0
{5,16}	0
{5,17}	2
{5,18}	0
{5,19}	2
{6,17}	0
{6,18}	1
{6,19}	0
{7,18}	1
{7,19}	2
{8,19}	0

C2

choose candidates with support count \geq minimum support count

Item Set	Support Count
{1,13}	2
{1,16}	2
{1,18}	2
{2,15}	2
{2,17}	2
{2,19}	2
{5,17}	2
{5,19}	2
{7,19}	2

L2

Generating Candidates: C_k

- **Join Step:** C_k is generated by **joining** L_{k-1} with itself
- **Prune Step:** Any $(k-1)$ -item set that is **not frequent** **cannot** be a subset of a **frequent k -item** set

Example: Joining and Pruning

1. The join step: To find C_k , a set of candidate k-itemsets is generated by joining L_{k-1} with itself.

L_k – Itemsets C_k – Candidates

For example in our case:

Considering $\{I2, I5\}$, $\{I7, I9\}$ from $L2$ to arrive at $C3$ we **Join $L2 * L2$**

and we obtain for example $\{I2, I5, I7\}$, $\{I2, I5, I9\}$ as resultant candidates in $C3$ generated from $L2$

Considering $\{I1, I3\}$, $\{I1, I6\}$ from $L2$ we generate a candidate $\{I1, I3, I6\}$ in $C3$

Example: Joining and Pruning

2. The prune step:

C_k is a superset of L_k , that is, its members **may or may not be frequent**

C_k however, **can be huge** and we **prune it** applying **Apriori Principle**
“if A is a frequent item set, then each of its subsets is a frequent item set”

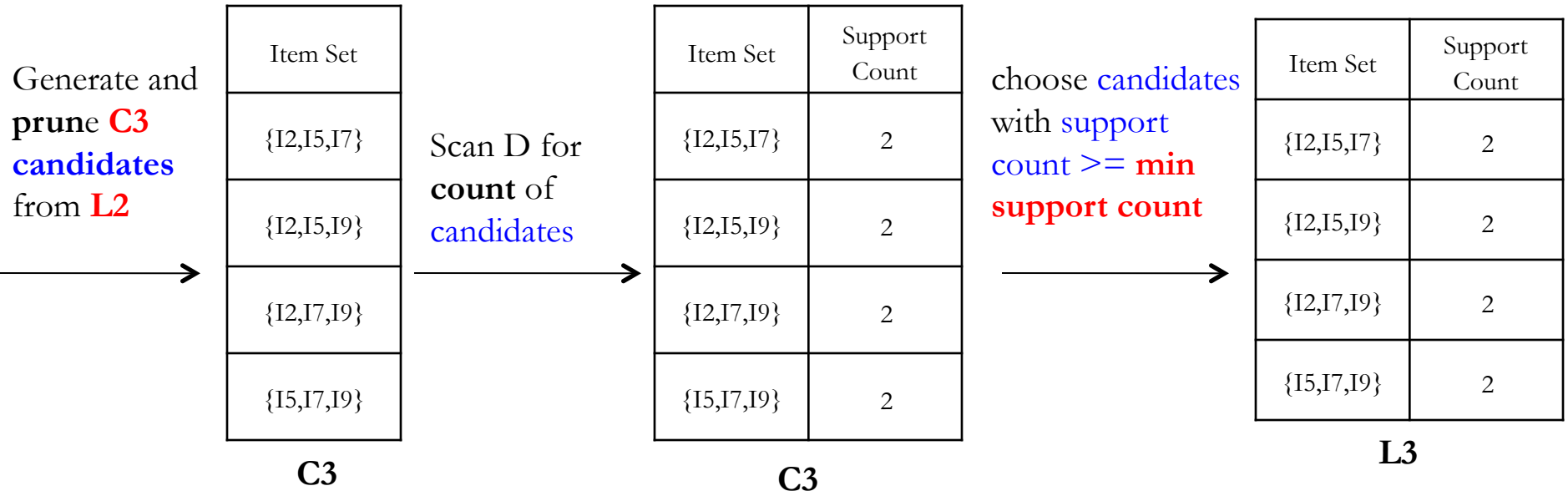
It is expressed by formulation of the

Prune Step: Any $(k-1)$ -item set that is **not frequent** cannot be a subset of a frequent k -item set

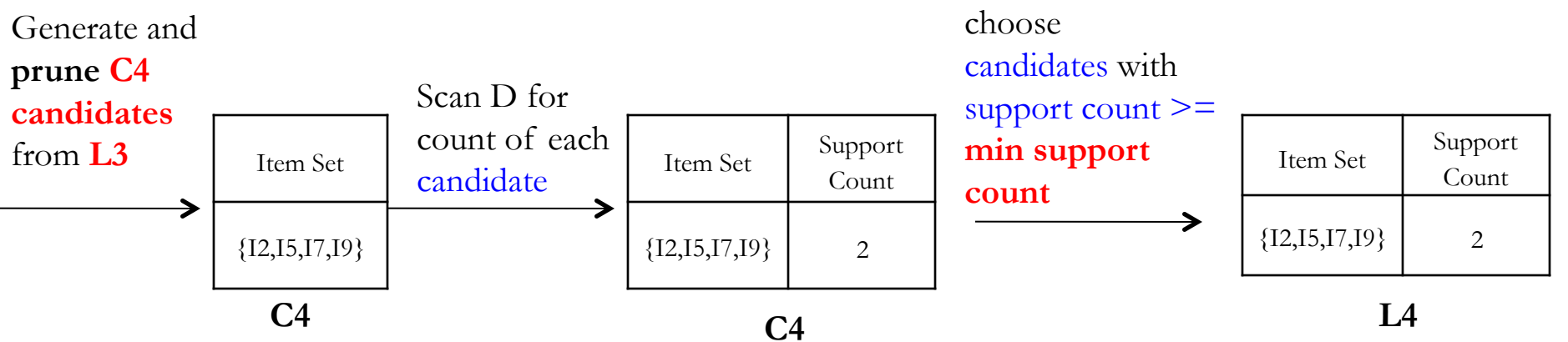
Thus, $\{I_2, I_5, I_7\}$, $\{I_2, I_5, I_9\}$ from **join step** are considered since **all their subsets are frequent**

but $\{I_1, I_3, I_6\}$ is **discarded** since its subset $\{I_3, I_6\}$ is **not frequent**, i.e. was not in **L_2**

Generating 3-itemset Frequent Pattern



Generating 4-itemset Frequent Pattern



Generating Multidimensional Association Rules

Let **minimum confidence** required be **70%**

- For example, let's consider 4-item frequent set
- $I = \{I_2, I_5, I_7, I_9\}$
- Its nonempty subsets needed to create rules
- (we write $\{2\}$ instead of $\{I_2\}$.. etc) are:
- $\{2\}, \{5\}, \{7\}, \{9\},$
- $\{2,5\}, \{2,7\}, \{2,9\}, \{5,7\}, \{5,9\}, \{7,9\},$
- $\{2,5,7\}, \{2,5,9\}, \{2,7,9\}, \{5,7,9\}$

We create for example some association rules as follows

$$R1: 2 \wedge 5 \wedge 7 \rightarrow 9 \quad R2: 2 \wedge 5 \wedge 9 \rightarrow 7 \quad R3: 5 \wedge 7 \rightarrow 2 \wedge 9$$

Multidimensional Association Rules

- R1 : $2 \wedge 5 \wedge 7 \rightarrow 9$

$\text{student}(x, \text{math}) \wedge \text{grade}(X, \text{low}) \wedge \text{income}(x, \text{low})$

$\Rightarrow \text{buys}(X, \text{bread})$

- R2 : $2 \wedge 5 \wedge 9 \rightarrow 7$

$\text{student}(x, \text{math}) \wedge \text{grade}(X, \text{low}) \wedge \text{buys}(X, \text{bread})$

$\Rightarrow \text{income}(x, \text{low})$

- R3 : $5 \wedge 7 \rightarrow 2 \wedge 9$

$\text{grade}(X, \text{low}) \wedge \text{income}(x, \text{low}) \Rightarrow \text{student}(x, \text{math}) \wedge \text{buys}(X, \text{bread})$

Example: Classification Data

Classification by Association

Student	Grade	Income	Buys
CS	High	Low	Milk
CS	High	High	Bread
Math	Low	Low	Bread
CS	Medium	High	Milk
Math	Low	Low	Bread

Converted Data

Student = CS (I1)	Student =math (I2)	Grade = high (I3)	Grade =medium (I4)	Grade =low (I5)	Income =high (I6)	Income =low (I7)	Buys =milk (I8)	Buys =bread (I9)
+	-	+	-	-	-	+	+	-
+	-	+	-	-	+	-	-	+
-	+	-	-	+	-	+	-	+
+	-	-	+	-	+	-	+	-
-	+	-	-	+	-	+	-	+
3	2	2	1	2	2	3	2	3

Generating **Classification Rules** by **Association**

When mining **association rules** for use in **classification** we are **only interested** in **association rules** of the form

$$i_1 \& i_2 \& \dots \& i_k \rightarrow i_c$$

where i_c is an **item associated** with a **class label c**

- The process of finding such rules is called
- **Classification by Association**

Classification by Association

- When generating **classification by association rules**
- we are **only interested** in **association rules** of the form
- $(p_1 \wedge p_2 \wedge \dots \wedge p_l) \rightarrow \text{class} = C$
- where the rule antecedent is a **conjunction of items**
- p_1, p_2, \dots, p_l **associated** with a **class label C**
- In our **example class is** either **I8** or **I9**
- as we want to **predict** whether a **student with given characteristics** **buys Milk** or **buys Bread**

Generating **Classification Rules** by Association

Let **minimum confidence** required be **70%**

We run **Apriori Algorithm** as before and

- **For example**, let's consider **4-item frequent set**
- **$I = \{I_2, I_5, I_7, I_9\}$** where **$I_9$** represents **buys-Bread**
- Its **nonempty subsets** needed to create **association rules**
- (we write **$\{2\}$** instead of **$\{I_2\}$** .. etc) are:
- **$\{2\}, \{5\}, \{7\}, \{9\},$**
- **$\{2,5\}, \{2,7\}, \{2,9\}, \{5,7\}, \{5,9\}, \{7,9\},$**
- **$\{2,5,7\}, \{2,5,9\}, \{2,7,9\}, \{5,7,9\}$**
- To create **classification rules** we consider **only subsets** that contain the **class item 9**

Generating **Classification Rules** by Association

Consider 3- itemset Frequent Sets **{2,5,9}**, **{2,7,9}**, **{5,7,9}**

We create **classification** by association rules as follows

R1 : 5 ^ 7 → 9 [40%,100%]

◦ Confidence = $sc\{I5,I7,I9\} / sc\{I5,I7\} = 2/2 = 100\%$

◦ **R2** is **selected**

◦ **R3 : 2 ^ 7 → 9** [40%,100%]

◦ Confidence = $sc\{I2,I7,I9\} / sc\{I2,I7\} = 2/2 = 100\%$

◦ **R3** is **selected**

◦ **R4 : 2 ^ 5 → 9** [40%,100%]

◦ Confidence = $sc\{I2,I7,I9\} / sc\{I2,I7\} = 2/2 = 100\%$

◦ **R4** is **selected**

Generating Classification by Association Rules

Consider 2- itemset Frequent Sets $\{2,9\}$, $\{5,9\}$, $\{7,9\}$,
and $\{1,8\}$ from **L2**

We create **classification by association rules** as follows

R5 : 5 → 9 [40%,100%]

- **Confidence** = $sc\{I5,I9\} / sc\{I9\} = 2/2 = 100\%$
- **R5** is **Selected**

R6 : 2 → 9 [40%,100%]

- **Confidence** = $sc\{I2,I9\} / sc\{I9\} = 2/2 = 100\%$
- **R6** is **Selected**

R7 : 7 → 9 [40%,100%]

- **Confidence** = $sc\{I7,I9\} / sc\{I9\} = 2/2 = 100\%$
- **R7** is **Selected**

R8 : 1 → 8 [40%, 66%]

- **Confidence** = $sc\{I1,I8\} / sc\{I1\} = 2/3 = 66.66\%$
- **R8** is **Rejected**

List of Selected **Classification by Association Rules**

- $2 \wedge 5 \wedge 7 \rightarrow 9$ [40%,100%]
- $2 \wedge 5 \rightarrow 9$ [40%,100%]
- $2 \wedge 7 \rightarrow 9$ [40%,100%]
- $5 \wedge 7 \rightarrow 9$ [40%,100%]
- $5 \rightarrow 9$ [40%,100%]
- $7 \rightarrow 9$ [40%,100%]
- $2 \rightarrow 9$ [40%,100%]

- We reduce the **confidence** to **66%** to include **I8**
- $1 \rightarrow 8$ [40%,66%]

Test Data

Student	Grade	Income	Buys
Math	Low	Low	Bread
CS	Low	Low	Milk
Math	Low	Low	Milk
Math	Low	Low	Bread
CS	Medium	High	Bread

- **First Tuple**

is correctly classified by the rule

$I2 \ \& \ I5 \ \& \ I7 \ \rightarrow \ I9$

Student=math & grade=low & income=low \rightarrow buys=bread **[Success]**

- **Second Tuple:**

There is no rule for class I8: buys=bredI8 **[Error]**

- **Third Tuple:**

There is no rule for class I8: buys=bredI8 **[Error]**

Test Data

Student	Grade	Income	Buys
Math	Low	Low	Bread
CS	Low	Low	Milk
Math	Low	Low	Milk
Math	High	Low	Bread
CS	Medium	High	Bread

- **Fourth Tuple**

is correctly classify by the rule $I_2 \wedge I_7 \rightarrow I_9$ [Success]

• **Student=Math & Income=Low \rightarrow Buys=Bread**

- **Fifth Tuple**

is correctly classify by the rule $I_1 \rightarrow I_9$ [Success]

Student=CS \rightarrow Buys=Bread

Hence we have **80% predictive accuracy**

And **20% Error rate**