

# Chapter 5: Mining Frequent Patterns, Association and Correlations

**Dong-Kyu Chae**

**PI of the Data Intelligence Lab @HYU  
Department of Computer Science & Data Science  
Hanyang University**

# MaxMiner: Mining Max-patterns

## □ Review: max pattern

□ An itemset  $X$  is a **max-pattern** if  $X$  is frequent and there exists no **frequent** super-pattern  $Y \supset X$

□ i.e., no such a  $Y$

- $Y$  is a super-pattern of  $X$
- The support of  $Y$  is greater than **minSup (=frequent)**
  - The support of  $Y$  can be smaller than the support of  $X$

cf) closed pattern : The super-pattern  $Y$  should be the same with the support of  $X$  인 경우  
 $Y$  는 closed pattern.  $X$  는 Not closed pattern

## □ MaxMiner is based on the Apriori algorithm

- It tries to find only max-patterns, not all the frequent patterns
- Naïve approach: too much costly!

# MaxMiner: Mining Max-patterns

- ❑ 1<sup>st</sup> scan: find frequent items and **sort** them (**ascending order**)
  - ❑ Assume the order is A, B, C, D, E (E is most frequently occurring)
- ❑ 2<sup>nd</sup> scan: find support for 2-itemsets **with potential max-patterns**

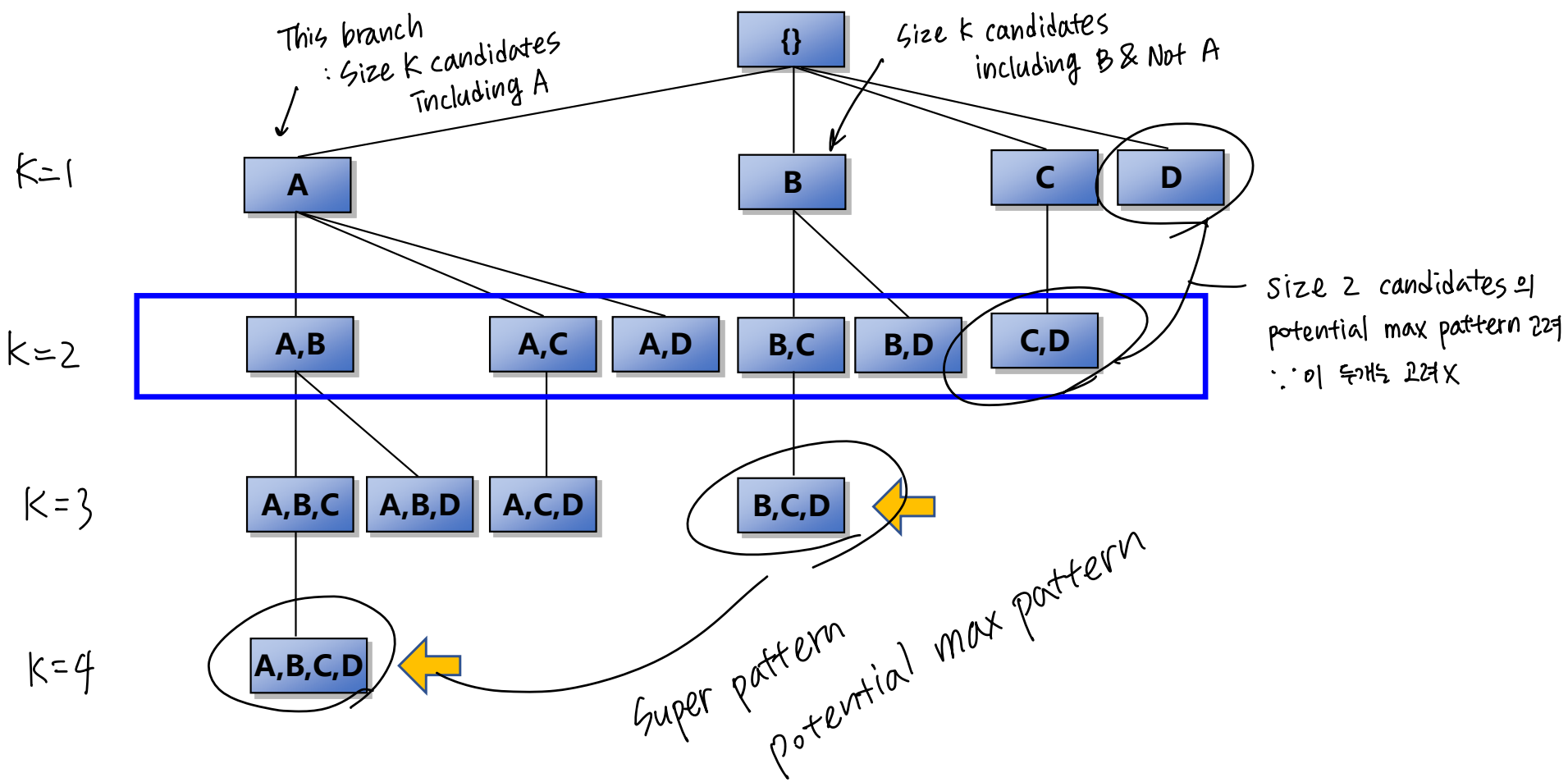
- ❑ AB, AC, AD, AE, **ABCDE** ←
  - ❑ BC, BD, BE, **BCDE** ← <sub>(~A)</sub>
  - ❑ CD, CE, **CDE** ← <sub>(~A, ~B)</sub>
  - ❑ DE
- Potential max-patterns  
for each category of  
candidates

Tid	Items
10	A,B,C,D,E
20	B,C,D,E,
30	A,C,D,F,E

- ❑ Reduce a lot of candidates in later stages
  - ❑ If **BCDE** is a max-pattern, no need to check **BCD, BDE, CDE** in later scan

# MaxMiner: Mining Max-patterns

- Construct **set-enumeration tree** over four items
- Assume order is A B C D
- It helps the process of **finding max pattern candidates**



# Mining Closed Patterns: CLOSET

## □ Review: closed pattern

- An itemset  $X$  is **closed** if  $X$  is *frequent* and there exists *no super-pattern*  $Y \supset X$ , **with the same support** as  $X$
- i.e., no such a  $Y$ 
  - $Y$  is a super-pattern of  $X$
  - The support of  $Y$  is should **be the same** as that of  $X$

## □ CLOSET finds the closed patterns while running FP-growth

- The algorithm of FP-growth naturally supports finding closed patterns

# Mining Closed Patterns: CLOSET

- ❑ Naïve approach: again, too much costly!
- ❑ Find only the closed itemsets recursively during the mining process of FP-growth
  - ❑ The transactions having <sup>all</sup>**m** also has **fca** => **fcam** is a frequent **closed** pattern!

**Conditional pattern bases**

<i>item</i>	<i>cond. pattern base</i>	<i>freq</i>
<i>f</i>	{}	4
<i>c</i>	<i>f:3</i>	4
<i>a</i>	<i>fc:3</i>	3
<i>b</i>	<i>fca:1, f:1, c:1</i>	3
<i>m</i>	<i>fca:2, fcab:1</i>	3
<i>p</i>	<i>fcam:2, cb:1</i>	3

*m*-conditional pattern base:  
*fca:2, fcab:1*

{ }

*f:3*

*c:3*

*a:3*

*m*-conditional FP-tree

*min\_sup* = 3

All frequent patterns related to *m*

*m,*  
*fm, cm, am,*  
*fcm, fam, cam,*  
*fcam*



# CHARM: Mining by Exploring Vertical Data Format

## Vertical format: $t(A) = \{T_{11}, T_{25}, \dots\}$

- Represent each itemset as a **list of transaction IDs** containing the corresponding itemset

a	T1, ...
b	
c	T1, T2, T3, ...

row: itemsets,  
columns: Transaction ID

T1	a c e f k
T2	c f g k i
T3	c g z

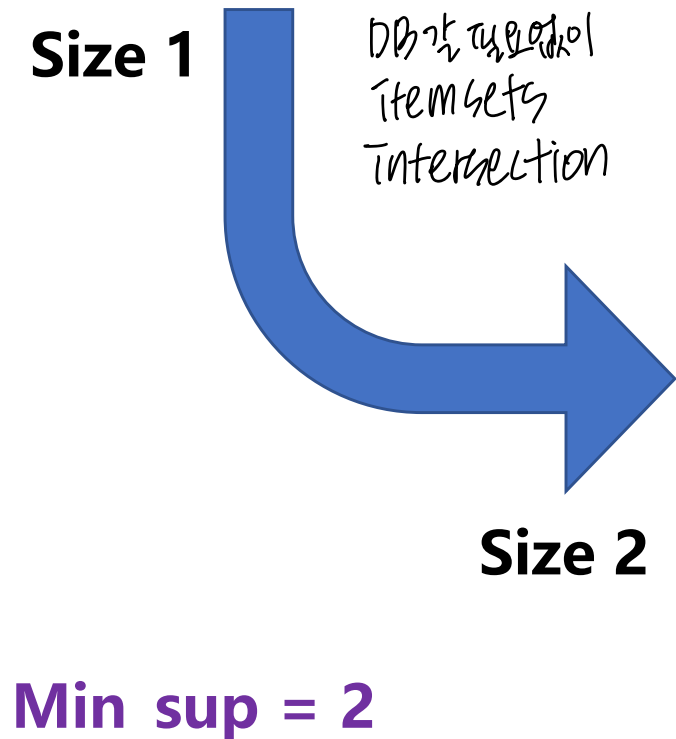
## Algorithm

- Transform a horizontally formatted data to a **vertical format** by scanning the dataset once, considering only frequent size-1 itemsets.
  - Easy: the number of items is much smaller than that of transactions
- Starting with  $k=1$ , find  $(k+1)$  size frequent itemsets from frequent  $k$ -size itemsets
  - Using the **intersection operation** and **Apriori property**
- Repeat this process until no frequent itemsets can be found in the next  $k+1$  step

# CHARM: Mining by Exploring Vertical Data Format

ITEMSET	TID SET
I1	{T101, T401, T501, T701, T801, T901}
I2	{T101, T201, T301, T401, T601, T801, T901}
I3	{T301, T501, T601, T701, T801, T901}
I4	{T201, T401}
I5	{T101, T801}

→ Transection ID 갯이 = frequency of each itemset



ITEMSET	TID SET
{I1,I2}	{T101, T401, T801, T901}
{I1,I3}	{T501, T701, T801, T901}
<del>{I1,I4}</del>	<del>{T401}</del>
{I1,I5}	{T301, T601, T801, T901}
{I2,I3}	{T201, T401}
{I2,I5}	{T101, T801}
<del>{I3,I5}</del>	<del>{T801}</del>
....	



# CHARM: Mining by Exploring Vertical Data Format

## Effectiveness:

- The **frequency** of each itemset is equal to **the length of its TID list**
- The simple **intersection operation** brings us the k+1 length candidates

ITEMSET	TID SET
{I1,I2}	{T101,T401,T801,T901}
{I1,I3}	{T501,T701,T801,T901}
<del>{I1,I4}</del>	<del>{T401}</del>
{I1,I5}	{T301,T601,T801,T901}
{I2,I3}	{T201,T401}
{I2,I5}	{T101,T801}
<del>{I3,I5}</del>	<del>{T801}</del>

....

**Size 2**

**Min\_sup = 2**

**Size 3**



ITEMSET	TID SET
{I1,I2,I3}	<del>{T801,T901}</del>
{I1,I2,I5}	<del>{T101,T901}</del>
....	



# Association Rules

---

- ❑ Mining multilevel association
- ❑ Mining multidimensional association
- ❑ Mining quantitative association
- ❑ Mining interesting correlation patterns

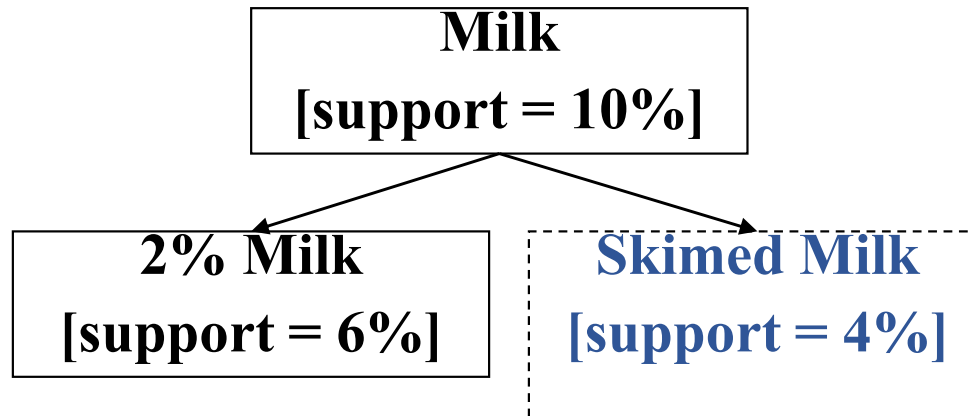
# Mining Multiple-Level Association Rules

- Items often form hierarchies (assume that it is given)
- Flexible support settings** will be needed
  - minimum
  - Items at the **lower level** are expected to have **lower support**
  - Make the probability of lower-level items and higher-level items being included in frequent patterns as similar as possible

## uniform setting

Level 1  
min\_sup = 5%

Level 2  
min\_sup = 5%



## flexible setting

Level 1  
min\_sup = 5%

Level 2  
min\_sup = 3%

# Multi-level Association: Redundancy Filtering

Some rules may be redundant due to “ancestor” relationships between items.

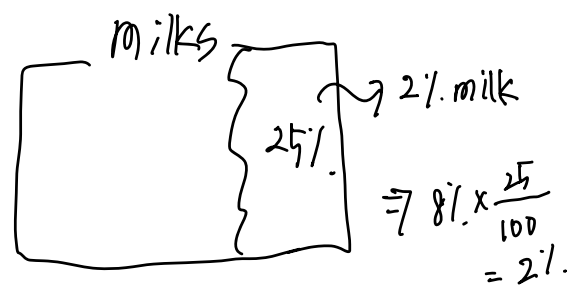
## Example

- $\text{milk} \Rightarrow \text{wheat bread}$  [support = 8%, confidence = 70%]
- ~~$2\% \text{ milk} \Rightarrow \text{wheat bread}$  [support = 2%, confidence = 72%]~~

We say the first rule is an ancestor of the second rule.

A (descendent) rule is redundant if

- Its support is close to the “expected” value, based on the ancestor rule's support value
   
 $\rightarrow$  전체 milk 카테고리에서 2% milk의 비율로 support 예측가능 할때
- Its confidence is close to that of ancestor rule





# Mining Multi-Dimensional Associations

## ❑ Single-dimensional rules:

$\text{buys}(X, \text{"milk"}) \Rightarrow \text{buys}(X, \text{"bread"})$ : **milk**  $\Rightarrow$  **bread**  
 $\nearrow$  Customer       $\nearrow$  dimension  
then, highly likely

## ❑ Multi-dimensional rules: $\geq 2$ dimensions

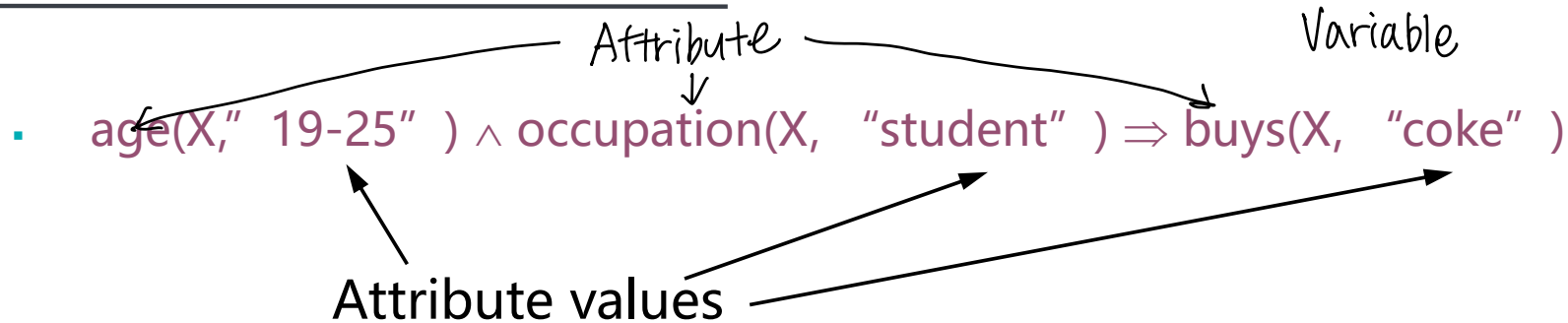
### ❑ Inter-dimension assoc. rules (*no repeated dimensions*)

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

### ❑ Hybrid-dimension assoc. rules (*repeated dimensions*)

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

# Attribute (Feature, Dimension) Types



## □ Categorical Attributes

- Finite number of possible values, no ordering among values

ex) occupation, buys

## □ Quantitative Attributes

- Numeric, implicit ordering among values
- Typically discretization is required (will be explained in the next several chapters)

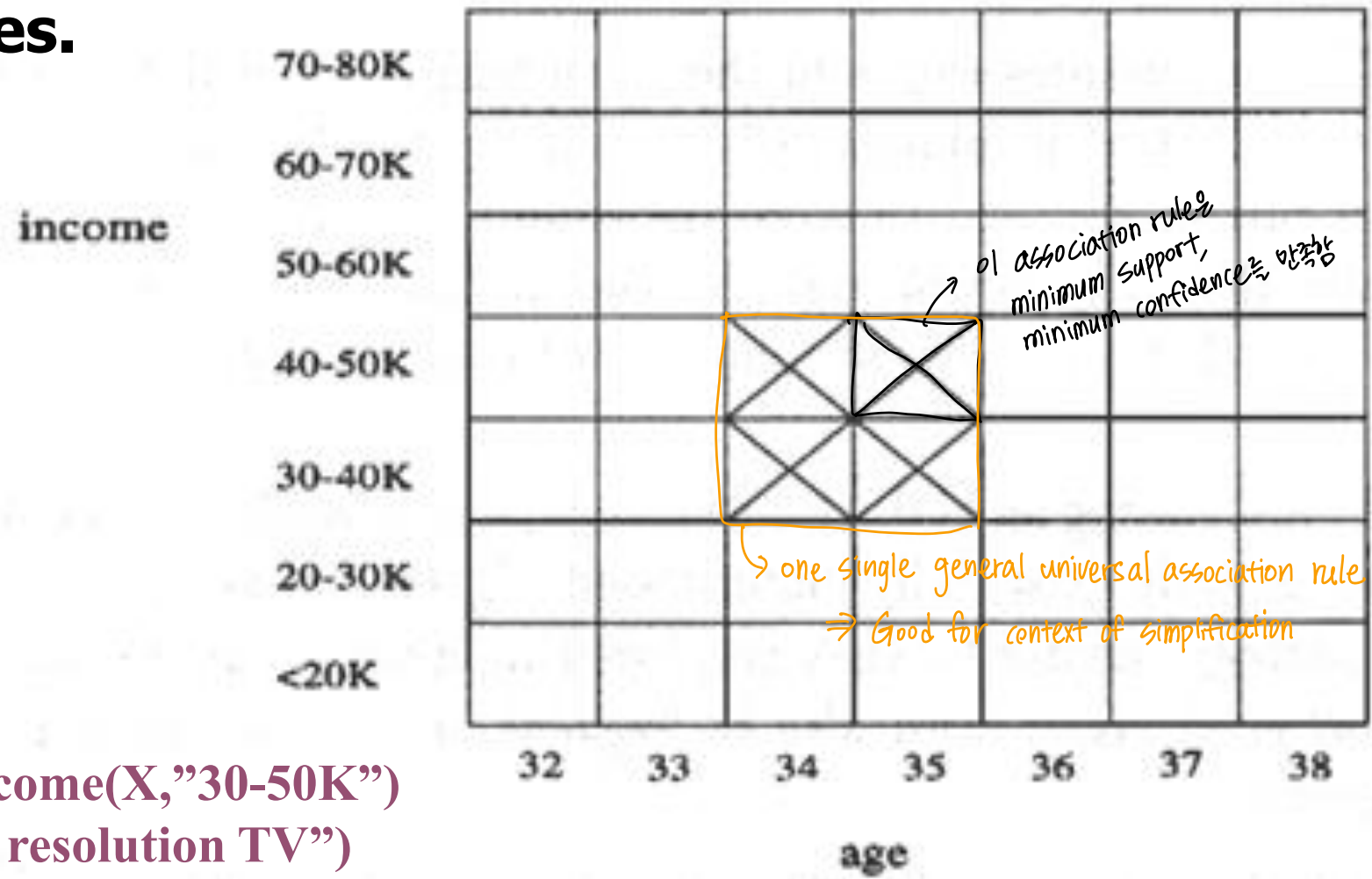
ex) age

# Quantitative Association Rules

- Numeric attributes are discretized
- Skeleton of our association rules:  $A_{\text{quan1}} \wedge A_{\text{quan2}} \Rightarrow A_{\text{cat}}$ 

quantitative

catagorical
- If needed, we **cluster adjacent association** rules to form **general rules**.
- Example



age(X,"34-35")  $\wedge$  income(X,"30-50K")  
 $\Rightarrow$  buys(X,"high resolution TV")

# From Association Mining to Correlation Analysis

- ❑ *play basketball*  $\Rightarrow$  *eat cereal* [40%, 66.7%] is misleading! Not meaningful!
  - ❑ Because, the overall % of students eating cereal is 75% ( > 66.7% )  
: general trend
- ❑ *play basketball*  $\Rightarrow$  *not eat cereal* [20%, 33.3%] is more meaningful, common sense를 학생들이 시의미를 많이 먹었다고 말하고 있는 상황에서 이 %는 높음. 의미있음 although it has lower support and confidence (limitation of the support & confidence framework)
- ❑ Measure of **interestingness** of correlated events: **lift**

Contingency table

	Basketball	Not basketball	Sum (row)
Cereal	2000	1750	3750
Not cereal	1000	250	1250
Sum(col.)	3000	2000	5000

$$lift = \frac{P(A \text{ and } B)}{P(A)P(B)}$$

*lift = 1*: Event A, B are completely independent  
No correlation at all

*lift < 1*: Two events have negative correlation  
*> 1*: " positive "

$$lift(B, C) = \frac{2000 / 5000}{3000 / 5000 * 3750 / 5000} = 0.89 \quad \rightarrow \text{negative correlation compare to the general trend}$$

$$lift(B, \neg C) = \frac{1000 / 5000}{3000 / 5000 * 1250 / 5000} = 1.33$$



# Summary

---

- ❑ **Frequent pattern & association rule mining—an important task in data mining**
- ❑ **Scalable frequent pattern mining methods**
  - ❑ **Apriori** (the basic candidate generation & test approach)
  - ❑ **FP-growth** (without the candidate generation & test)
  - ❑ Minor improvements (**DIC, Partition, Sampling, MaxMiner, CHARM**, etc... )
- **Mining a variety of rules and interesting patterns**
  - Hierarchy and multi-attribute settings, Lift, etc...

# Thank You



Data  
Intelligence  
Lab