

CSCI 4502/5502

Data Mining

Fall 2020
Lecture 07 (Sep 15)

Reminders

- ◆ Homework 2

- ◆ due at 9:30am, Th, Sep 17

- ◆ SUBMIT your attempt in Canvas before deadline

- ◆ Computing and Software Career & Internship Fair

- ◆ 11am-4pm, Tu, Sep 15, virtual on Handshake

Review

♦ Chap 4 & 5: Data Warehouse, Data Cube

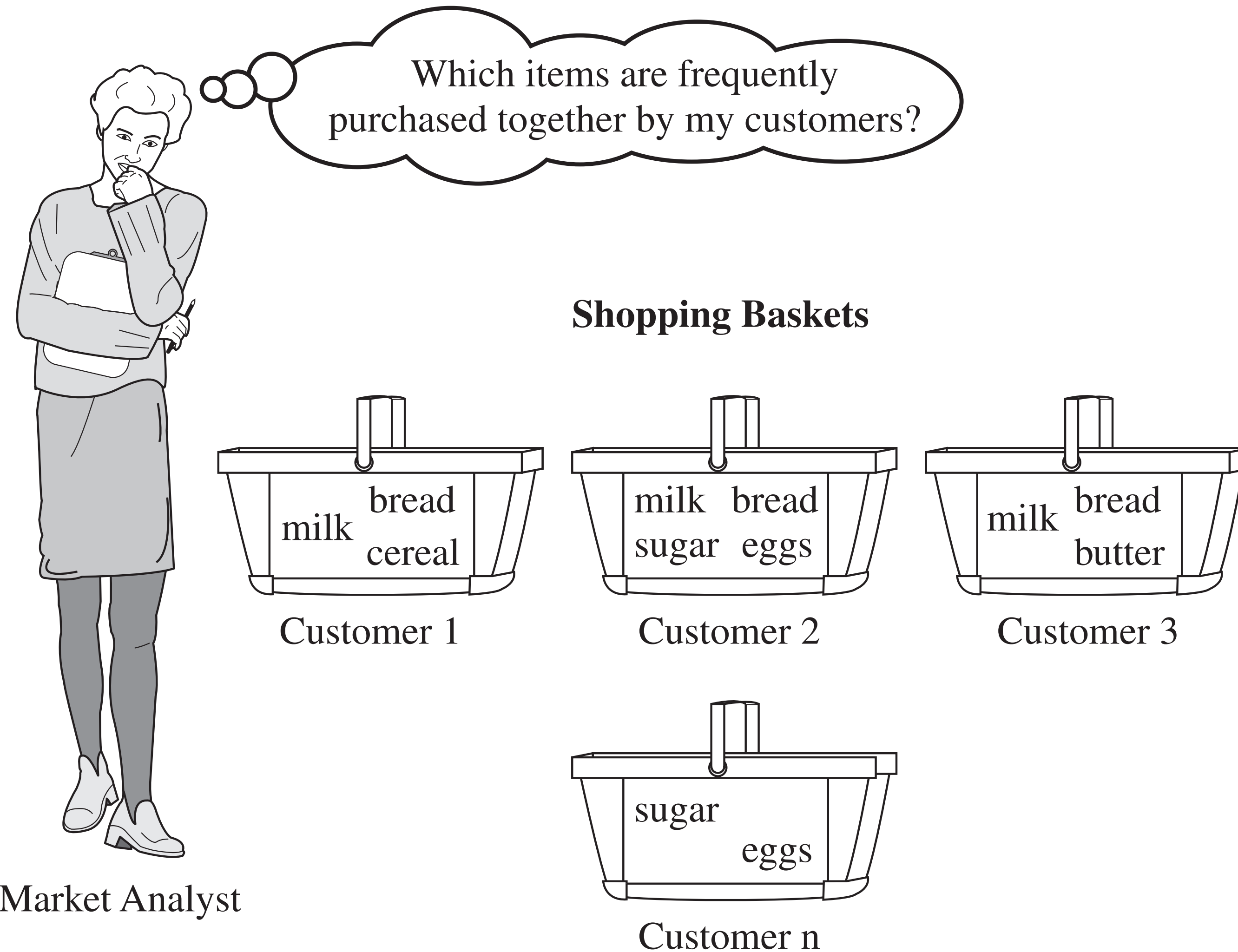
- ♦ what is data warehouse?
- ♦ OLTP vs. OLAP
- ♦ what is data cube?
- ♦ data cube operations
- ♦ data cube computation

Review: Part I

- ◆ Chapter 1: Introduction
- ◆ Chapter 2: Getting to Know your Data
- ◆ Chapter 3: Data Preprocessing
- ◆ Chapter 4: Data Warehousing & Online Analytical Processing
- ◆ Chapter 5: Data Cube Technology
- ◆ Part 2: Core DM Techniques
- ◆ Part 3: Mining Complex Data, DM Trends

Chapter 6: Mining Frequent Patterns, Associations & Correlations

Market Basket Analysis



http://www.information-drivers.com/images/beer_and_baby.gif

Frequent Pattern Analysis

- ◆ Frequent patterns in a data set

- ◆ a set of items

- ◆ subsequences

- ◆ substructures

- ◆ Examples

- ◆ Web log

- ◆ Road traffic

Basic Concepts

- ♦ Frequent itemset

- ♦ $X = \{x_1, x_2, \dots, x_k\}$

- ♦ Association rule $X \Rightarrow Y$

- ♦ support: probability that a transaction contains $X \cup Y$

- ♦ confidence: conditional probability that a transaction containing X also contains Y

- ♦ minimum support, minimum confidence

Example

♦ Let $\text{min_sup} = 50\%$, $\text{min_conf} = 50\%$

♦ Frequent patterns

♦ A, B, D, E, AD

♦ Association rules

♦ $A \Rightarrow D$ (%, %)

♦ $D \Rightarrow A$ (%, %)

Tid	Items
1	A, B, D
2	A, C, D
3	A, D, E
4	B, E, F
5	B, C, D, E, F

Example

♦ Let $\text{min_sup} = 50\%$, $\text{min_conf} = 50\%$

♦ Frequent patterns

♦ A 3, B 3, D 4, E 3, AD 3

♦ Association rules

♦ $A \Rightarrow D$ (60 %, 100 %)

♦ $D \Rightarrow A$ (60 %, 75 %)

Tid	Items
1	A, B, D
2	A, C, D
3	A, D, E
4	B, E, F
5	B, C, D, E, F

Mining Association Rules

- ♦ Two-step process
 - ♦ find all **frequent itemsets** (w/ min_sup)
 - ♦ generate **strong association rules** from the frequent itemsets (min_sup, min_conf)
- ♦ A long pattern contains a combinatorial **number of subpatterns** (e.g., 100 items)

$$\binom{100}{1} + \binom{100}{2} + \dots + \binom{100}{100} = 2^{100} - 1 \approx 1.27 \times 10^{30}$$

Closed & Max Patterns

- ◆ Solution: mine closed patterns & max-patterns
- ◆ Closed pattern X
 - ◆ no super-pattern $Y \supset X$ w/ the same support
- ◆ Max-pattern X
 - ◆ no super-pattern $Y \supset X$
- ◆ Closed pattern is a lossless compression of frequent patterns
 - ◆ reducing the number of patterns and rules

Example

- ♦ $\{ \langle a_1, \dots, a_{100} \rangle, \langle a_1, \dots, a_{50} \rangle \}$, $\text{min_sup} = 0.5$
- ♦ Frequent pattern?
 - ♦ all item combinations
- ♦ Closed pattern?
 - ♦ $\langle a_1, \dots, a_{100} \rangle$: 1
 - ♦ $\langle a_1, \dots, a_{50} \rangle$: 2
- ♦ Max-pattern?
 - ♦ $\langle a_1, \dots, a_{100} \rangle$: 1

Apriori Algorithm (I)

♦ Apriori property

- ♦ subset of a freq. itemset is also frequent
- ♦ e.g., {beer, diaper, nuts}, {beer, diaper}

♦ Apriori pruning

- ♦ if X is infrequent,
- ♦ then superset of X is pruned

Apriori Algorithm (2)

◆ Procedure

- ◆ 1. scan DB to get freq. l -itemset
- ◆ 2. generate candidate $(k+1)$ -itemsets from freq. k -itemsets
- ◆ 3. test candidate $(k+1)$ -itemsets against DB
- ◆ 4. stop when no freq. or candidate itemsets can be generated

Apriori Algorithm: Example

Tid	Items
1	A, C, D
2	B, C, E
3	A, B, C, E
4	B, E

$\text{min_sup} = 0.5$

Itemset	sup

Itemset	sup

Itemset	sup

Apriori Algorithm: Example

Tid	Items
1	A, C, D
2	B, C, E
3	A, B, C, E
4	B, E

min_sup = 0.5

Itemset	sup
{A}	0.5
{B}	0.75
{C}	0.75
{D}	0.25
{E}	0.75

Itemset	sup
{B, C, E}	0.5

Itemset	sup
{A, B}	0.25
{A, C}	0.5
{A, E}	0.25
{B, C}	0.5
{B, E}	0.75
{C, E}	0.5

Important Details

- ♦ **Self-joining** of k -itemsets to generate $(k+1)$ -itemsets
 - ♦ two k -itemsets are joined if their first $(k-1)$ items are the same
- ♦ **Pruning**: remove if subset not frequent
- ♦ Example: $L3 = \{abc, abd, acd, ace, bcd\}$
 - ♦ abc and $abd \Rightarrow abcd$
 - ♦ acd and $ace \Rightarrow acde$
 - ♦ $acde$ pruned because ade is not in $L3$

Interestingness Measure

- ♦ Association rule
 - ♦ $A \Rightarrow B$ [support, confidence]
- ♦ A strong association rule
 - ♦ play basketball \Rightarrow eat cereal [40%, 66.7%]
- ♦ The rule is misleading
 - ♦ overall, 75% of students eat cereal
 - ♦ play basketball \Rightarrow not eat cereal [20%, 33.3%]

Correlation Rules

- ◆ Correlation rule

- ◆ $A \Rightarrow B$ [support, confidence, **correlation**]

- ◆ Measure of dependent/correlated events

$$lift(A, B) = \frac{P(A \cup B)}{P(A)P(B)}$$

- ◆ lift = 1?

independent

- ◆ lift < 1?

negatively dependent

- ◆ lift > 1?

positively dependent

$$lift(A, B) = \frac{P(A \cup B)}{P(A)P(B)}$$

	basketball	not basketball	sum (row)
cereal	2000	1750	3750
not cereal	1000	250	1250
sum (col)	3000	2000	5000

$$lift(B, C) = \frac{2000/5000}{(3000/5000) \times (3750/5000)} = 0.89$$

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