

Data Mining

Association Rules: Advanced Concepts and Algorithms

Lecture Notes for Chapter 7

Introduction to Data Mining

Fang Zhou

Continuous and Categorical Attributes

How to apply association analysis to non-symmetric binary variables?

Gender	...	Age	Annual Income	No of hours spent online per week	No of email accounts	Privacy Concern
Female	...	26	90K	20	4	Yes
Male	...	51	135K	10	2	No
Male	...	29	80K	10	3	Yes
Female	...	45	120K	15	3	Yes
Female	...	31	95K	20	5	Yes
Male	...	25	55K	25	5	Yes
Male	...	37	100K	10	1	No
Male	...	41	65K	8	2	No
Female	...	26	85K	12	1	No
...

Example of Association Rule:

$\{\text{Gender}=\text{Male}, \text{Age} \in [21,30)\} \rightarrow \{\text{No of hours online} \geq 10\}$

Handling Categorical Attributes

- Example: Internet Usage Data

Gender	Level of Education	State	Computer at Home	Online Auction	Chat Online	Online Banking	Privacy Concerns
Female	Graduate	Illinois	Yes	Yes	Daily	Yes	Yes
Male	College	California	No	No	Never	No	No
Male	Graduate	Michigan	Yes	Yes	Monthly	Yes	Yes
Female	College	Virginia	No	Yes	Never	Yes	Yes
Female	Graduate	California	Yes	No	Never	No	Yes
Male	College	Minnesota	Yes	Yes	Weekly	Yes	Yes
Male	College	Alaska	Yes	Yes	Daily	Yes	No
Male	High School	Oregon	Yes	No	Never	No	No
Female	Graduate	Texas	No	No	Monthly	No	No
...

{Level of Education=Graduate, Online Banking=Yes}
→ {Privacy Concerns = Yes}

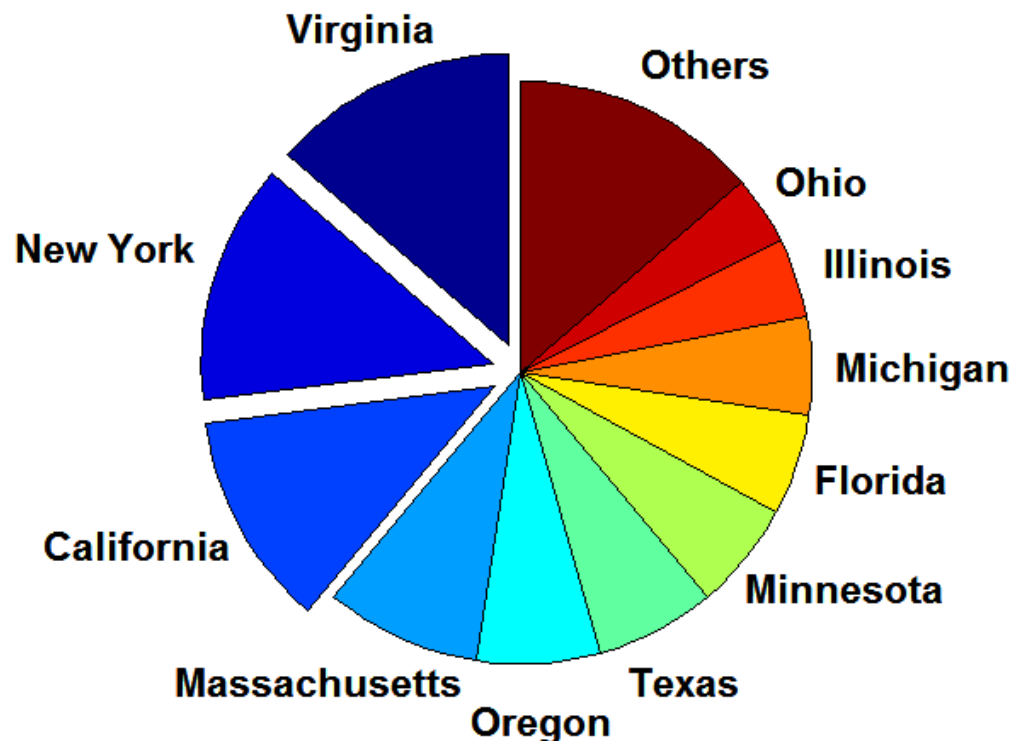
Handling Categorical Attributes

- Introduce a new “item” for each distinct attribute-value pair

Male	Female	Education = Graduate	Education = College	Education = High School	...	Privacy = Yes	Privacy = No
0	1	1	0	0	...	1	0
1	0	0	1	0	...	0	1
1	0	1	0	0	...	1	0
0	1	0	1	0	...	1	0
0	1	1	0	0	...	1	0
1	0	0	1	0	...	1	0
1	0	0	0	0	...	0	1
1	0	0	0	1	...	0	1
0	1	1	0	0	...	0	1
...

Handling Categorical Attributes

- Some attributes can have many possible values
 - Many of their attribute values have very low support
 - ◆ Potential solution: Aggregate the low-support attribute values



Handling Categorical Attributes

- Distribution of attribute values can be highly skewed
 - Example: 85% of survey participants own a computer at home
 - ◆ Most records have Computer at home = Yes
 - ◆ Computation becomes expensive; many frequent itemsets involving the binary item (Computer at home = Yes)
 - ◆ Potential solution:
 - discard the highly frequent items
- Computational Complexity
 - Avoid generating candidate itemsets that contain more than one item from the same attribute

Handling Continuous Attributes

Gender	...	Age	Annual Income	No of hours spent online per week	No of email accounts	Privacy Concern
Female	...	26	90K	20	4	Yes
Male	...	51	135K	10	2	No
Male	...	29	80K	10	3	Yes
Female	...	45	120K	15	3	Yes
Female	...	31	95K	20	5	Yes
Male	...	25	55K	25	5	Yes
Male	...	37	100K	10	1	No
Male	...	41	65K	8	2	No
Female	...	26	85K	12	1	No
...

Example of Association Rule:

$\{\text{Gender}=\text{Male}, \text{Age} \in [21,30)\} \rightarrow \{\text{No of hours online} \geq 10\}$

$\{\text{Age} \in [21,35), \text{Salary} \in [70\text{k}, 120\text{k})\} \rightarrow \text{Buy}$

$\{\text{Age} \in [21,30), \text{Chat Online} = \text{Yes}\}$

$\rightarrow \text{No of hours online: } \mu=14, \sigma=4$

Handling Continuous Attributes

- Different methods:
 - Discretization-based
 - Statistics-based
 - Non-discretization based

Discretization-based Methods

Gender	...	Age	Annual Income	No of hours spent online per week	No of email accounts	Privacy Concern
Female	...	26	90K	20	4	Yes
Male	...	51	135K	10	2	No
Male	...	29	80K	10	3	Yes
Female	...	45	120K	15	3	Yes
Female	...	31	95K	20	5	Yes
Male	...	25	55K	25	5	Yes
Male	...	37	100K	10	1	No
Male	...	41	65K	8	2	No
Female	...	26	85K	12	1	No
...

[illegible]

Discretization-based Methods

- Unsupervised:

- Equal-width binning $\langle 1\ 2\ 3 \rangle \langle 4\ 5\ 6 \rangle \langle 7\ 8\ 9 \rangle$
- Equal-frequency binning $\langle 1\ 2 \rangle \langle 3\ 4\ 5\ 6\ 7 \rangle \langle 8\ 9 \rangle$
- Cluster-based

- Supervised discretization

Continuous attribute, v

	1	2	3	4	5	6	7	8	9
Chat Online = Yes	0	0	20	10	20	0	0	0	0
Chat Online = No	150	100	0	0	0	100	100	150	100

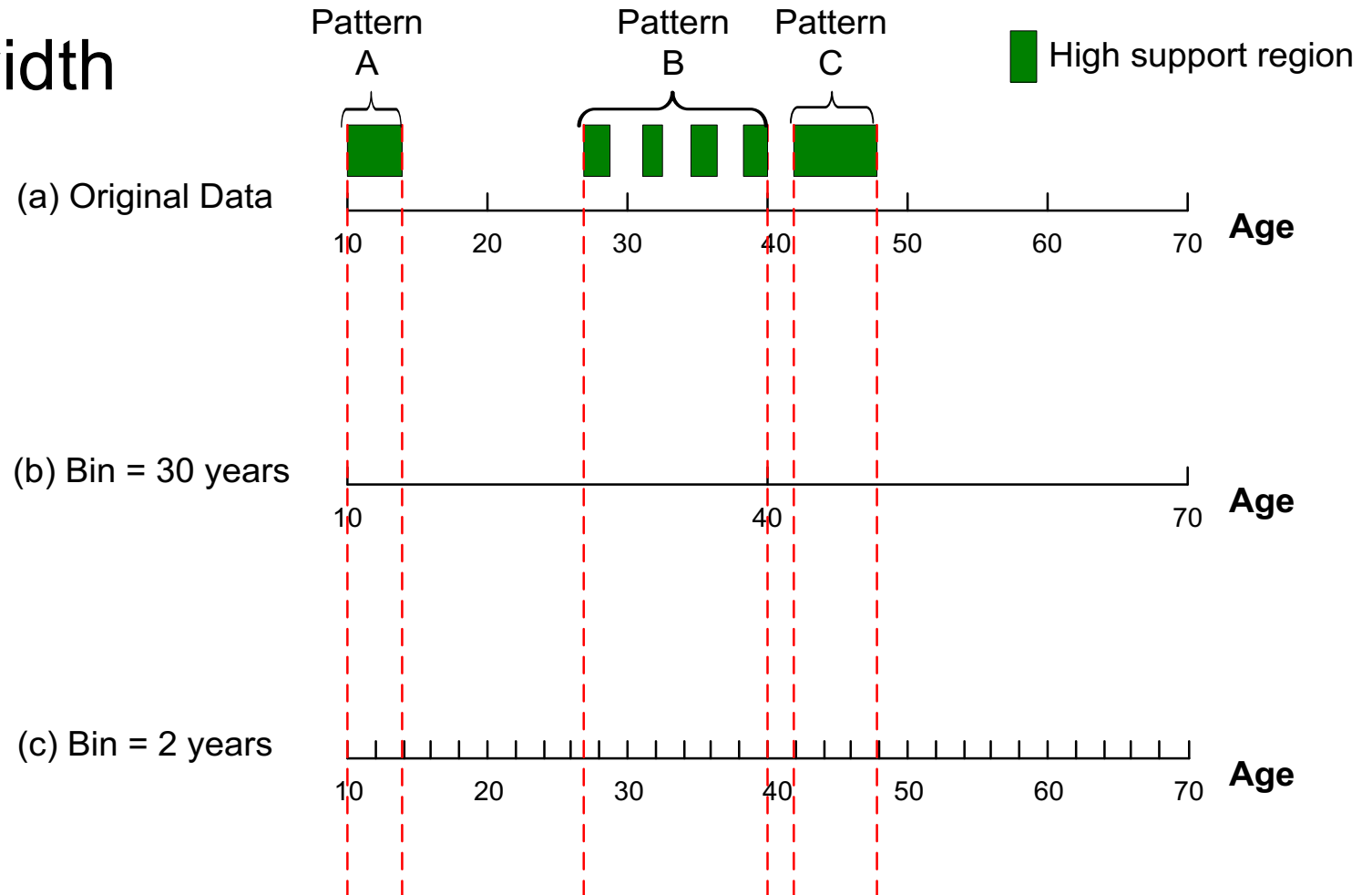
{
bin₁

{
bin₂

{
bin₃

Discretization Issues

● Interval width



Pattern A: Age $\in [10, 15)$ \longrightarrow Chat Online = Never
Pattern B: Age $\in [26, 41)$ \longrightarrow Chat Online = Never
Pattern C: Age $\in [42, 48)$ \longrightarrow Online Banking = Yes

Discretization Issues

- Interval too wide (e.g., Bin size= 30)
 - May merge several disparate patterns
 - ◆ Patterns A and B are merged together
 - May lose some of the interesting patterns
 - ◆ Pattern C **may not have enough confidence**
- Interval too narrow (e.g., Bin size = 2)
 - Pattern A is broken up into two smaller patterns
 - ◆ Can recover the pattern by merging adjacent subpatterns
 - Pattern B is broken up into smaller patterns
 - ◆ Cannot recover the pattern by merging adjacent subpatterns
 - Some windows **may not meet support threshold**

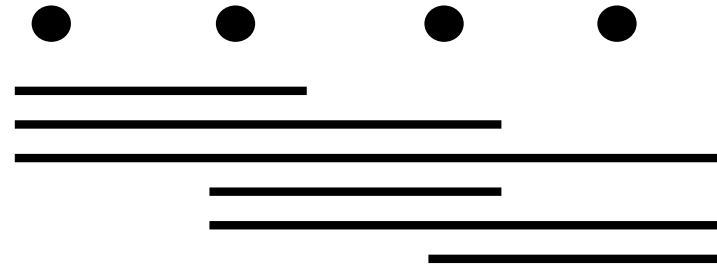
Discretization Issues

- Size of the discretized intervals affect support & confidence
 - If intervals too small
 - ◆ may not have enough support
 - If intervals too large
 - ◆ may not have enough confidence
- Potential solution: use all possible intervals

Discretization: all possible intervals

Number of intervals = k

Total number of Adjacent intervals = $k(k-1)/2$



● Execution time

- If the range is partitioned into k intervals, there are $O(k^2)$ new items
- If an interval $[a,b)$ is frequent, then all intervals that subsume $[a,b)$ must also be frequent
 - ◆ E.g.: if $\{\text{Age} \in [21,25), \text{Chat Online}=\text{Yes}\}$ is frequent, then $\{\text{Age} \in [10,50), \text{Chat Online}=\text{Yes}\}$ is also frequent
- Improve efficiency:
 - ◆ Use maximum support to avoid intervals that are too wide

Discretization Issues

- Redundant rules

R1: {Age \in [18,20), Gender=Male} \rightarrow {Chat Online=Yes}

R2: {Age \in [18,23), Gender=Male} \rightarrow {Chat Online=Yes}

- If both rules have the same support and confidence, prune the more specific rule (R1)

Statistics-based Methods

- Example:

{Income > 100K, Online Banking=Yes} → Age: $\mu=34$

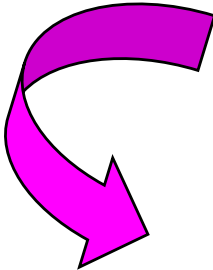
- Rule consequent consists of a continuous variable, characterized by their statistics

- mean, median, standard deviation, etc.

- Approach:

- Withhold the target attribute from the rest of the data
- Extract frequent itemsets from the rest of the attributes
 - ◆ Binarized the continuous attributes (except for the target attribute)
- For each frequent itemset, compute the corresponding descriptive statistics of the target attribute
 - ◆ Frequent itemset becomes a rule by introducing the target variable as rule consequent
- Apply statistical test to determine interestingness of the rule

Statistics-based Methods



Gender	...	Age	Annual Income	No of hours spent online per week	No of email accounts	Privacy Concern
Female	...	26	90K	20	4	Yes
Male	...	51	135K	10	2	No
Male	...	29	80K	10	3	Yes
Female	...	45	120K	15	3	Yes
Female	...	31	95K	20	5	Yes
Male	...	25	55K	25	5	Yes
Male	...	37	100K	10	1	No
Male	...	41	65K	8	2	No
Female	...	26	85K	12	1	No
...

Frequent Itemsets:

{Male, Income > 100K}
{Income < 30K, No hours ∈ [10,15]}
{Income > 100K, Online Banking = Yes}

Association Rules:

{Male, Income > 100K} → Age: $\mu = 30$
{Income < 40K, No hours ∈ [10,15]} → Age: $\mu = 24$
{Income > 100K, Online Banking = Yes}
 → Age: $\mu = 34$

Statistics-based Methods

- How to determine whether an association rule interesting?

- Compare the statistics for segment of population **covered** by the rule vs segment of population **not covered** by the rule:

$$A \Rightarrow B: \mu \quad \text{versus} \quad \bar{A} \Rightarrow B: \mu'$$

- Statistical hypothesis testing:

- ◆ Null hypothesis: $H_0: \mu' = \mu + \Delta$
- ◆ Alternative hypothesis: $H_1: \mu' > \mu + \Delta$
- ◆ Z has zero mean and variance 1 under null hypothesis

$$Z = \frac{\mu' - \mu - \Delta}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

Statistics-based Methods

- Example:

r: Browser=Mozilla \wedge Buy=Yes \rightarrow Age: $\mu=23$

- Rule is interesting if difference between μ and μ' is more than 5 years (i.e., $\Delta = 5$)
- For r, suppose $n_1 = 50$, $s_1 = 3.5$
- For r' (complement): $n_2 = 250$, $s_2 = 6.5$

$$Z = \frac{\mu' - \mu - \Delta}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} = \frac{30 - 23 - 5}{\sqrt{\frac{3.5^2}{50} + \frac{6.5^2}{250}}} = 3.11$$

- For 1-sided test at 95% confidence level, critical Z-value for rejecting null hypothesis is 1.64.
- Since Z is greater than 1.64, r is an interesting rule

Non-discretization methods

Document-term matrix:

TID	W1	W2	W3	W4	W5
D1	2	2	0	0	1
D2	0	0	1	2	2
D3	2	3	0	0	0
D4	0	0	1	0	1
D5	1	1	1	0	2

Example:

W1 and W2 tends to appear together in the same document

Non-discretization methods

- Data contains only **continuous attributes** of the same “type”
 - e.g., frequency of words in a document

TID	W1	W2	W3	W4	W5
D1	2	2	0	0	1
D2	0	0	1	2	2
D3	2	3	0	0	0
D4	0	0	1	0	1
D5	1	1	1	0	2

- Potential solution:
 - Convert into 0/1 matrix and then apply existing algorithms
 - ◆ lose word frequency information
 - Discretization does not apply as users want association among words not ranges of words

Min-Apriori

- How to determine the support of a word?
 - If we simply sum up its frequency, support count will be greater than total number of documents!
 - ◆ Normalize the word vectors – e.g., using L_1 norms
 - ◆ Each word has a support equals to 1.0

TID	W1	W2	W3	W4	W5
D1	2	2	0	0	1
D2	0	0	1	2	2
D3	2	3	0	0	0
D4	0	0	1	0	1
D5	1	1	1	0	2

Normalize



TID	W1	W2	W3	W4	W5
D1	0.40	0.33	0.00	0.00	0.17
D2	0.00	0.00	0.33	1.00	0.33
D3	0.40	0.50	0.00	0.00	0.00
D4	0.00	0.00	0.33	0.00	0.17
D5	0.20	0.17	0.33	0.00	0.33

Min-Apriori

- New definition of support:

$$\text{sup}(C) = \sum_{i \in T} \min_{j \in C} D(i, j)$$

TID	W1	W2	W3	W4	W5
D1	0.40	0.33	0.00	0.00	0.17
D2	0.00	0.00	0.33	1.00	0.33
D3	0.40	0.50	0.00	0.00	0.00
D4	0.00	0.00	0.33	0.00	0.17
D5	0.20	0.17	0.33	0.00	0.33

Example:

Sup(W1,W2,W3)

= 0 + 0 + 0 + 0 + 0.17

= 0.17

- Support increases monotonically as the normalized frequency of a word increases
- Support increases monotonically as the number of documents that contain the word increases

Anti-monotone property of Support

TID	W1	W2	W3	W4	W5
D1	0.40	0.33	0.00	0.00	0.17
D2	0.00	0.00	0.33	1.00	0.33
D3	0.40	0.50	0.00	0.00	0.00
D4	0.00	0.00	0.33	0.00	0.17
D5	0.20	0.17	0.33	0.00	0.33

Example:

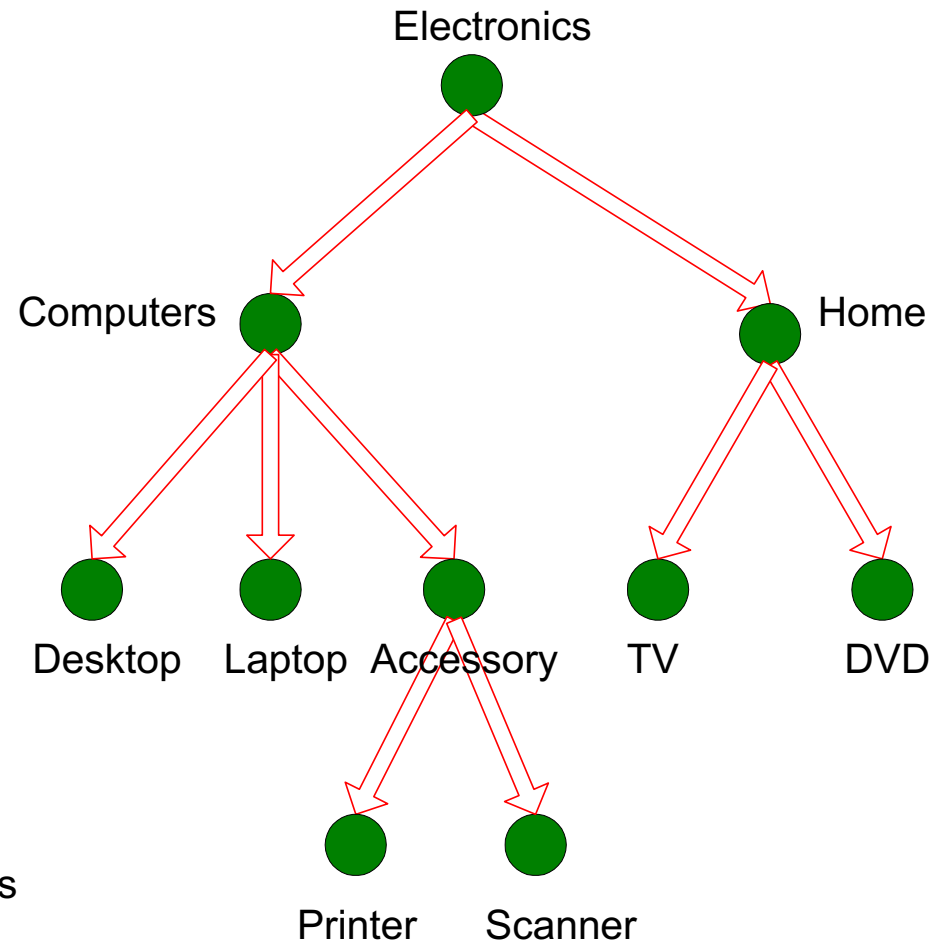
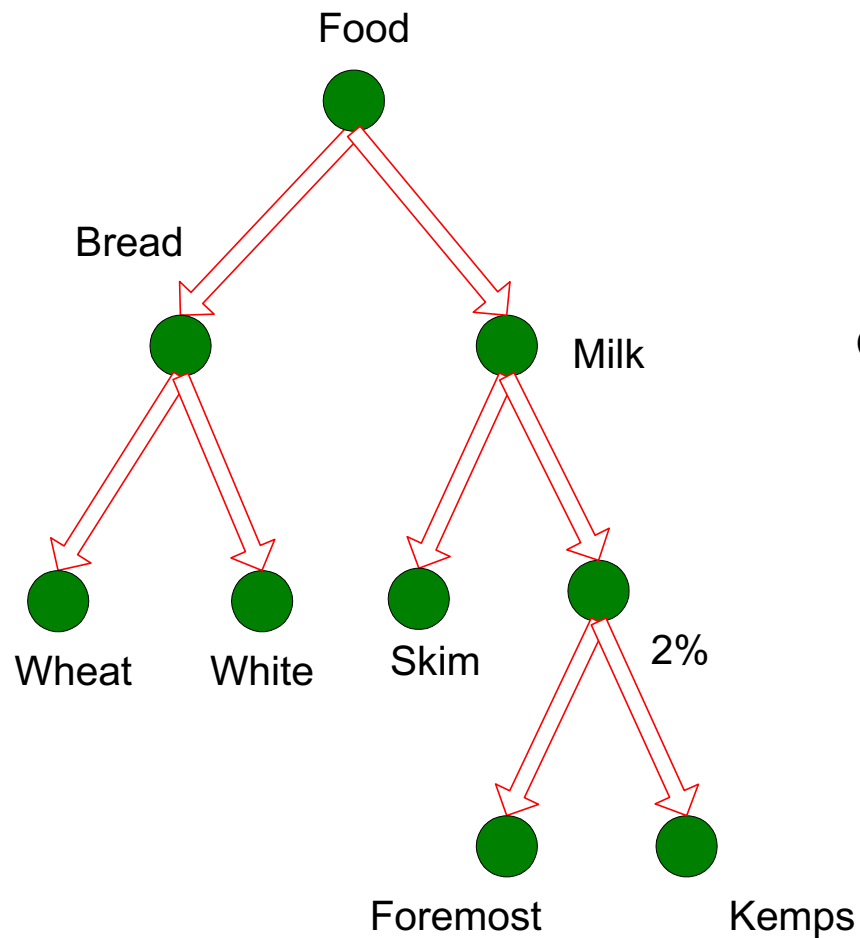
$$\text{Sup}(W1) = 0.4 + 0 + 0.4 + 0 + 0.2 = 1$$

$$\text{Sup}(W1, W2) = 0.33 + 0 + 0.4 + 0 + 0.17 = 0.9$$

$$\text{Sup}(W1, W2, W3) = 0 + 0 + 0 + 0 + 0.17 = 0.17$$

Support has an anti-monotone property.

Concept Hierarchies



Multi-level Association Rules

- Why should we incorporate concept hierarchy?
 - Rules at lower levels may not have enough support to appear in any frequent itemsets
 - Rules at lower levels of the hierarchy are overly specific
 - ◆ e.g., skim milk → white bread, 2% milk → wheat bread, skim milk → wheat bread, etc.are indicative of association between milk and bread
 - Rules at higher level of hierarchy may be too generic
 - ◆ e.g., food->electronics