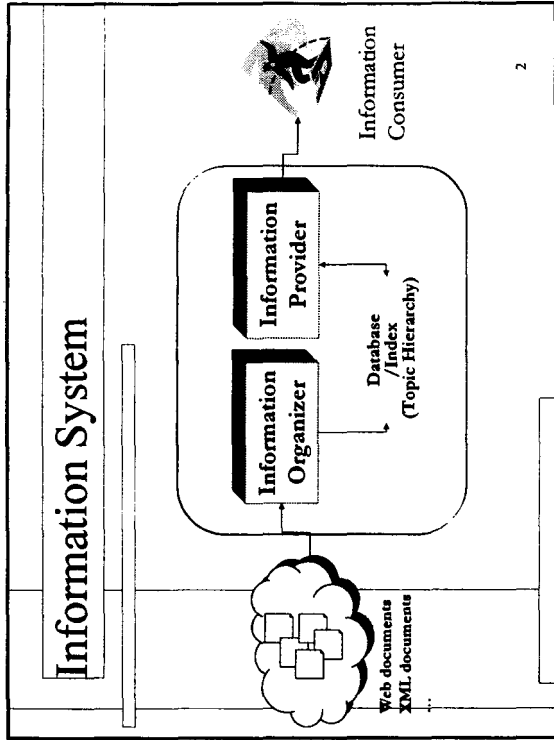


Building Topic Hierarchy of e-Documents using Text Mining Technology

Han-joon Kim
Department of Electrical and Computer Engineering
University of Seoul
khj@uos.ac.kr

1



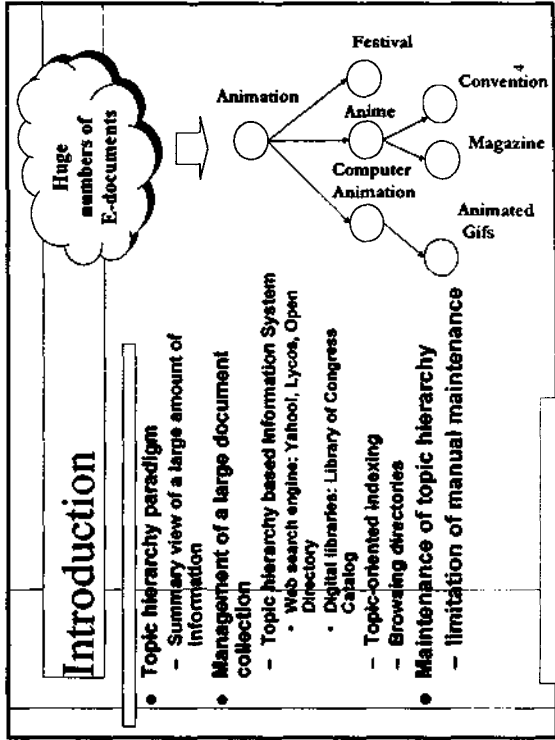
Information Organization

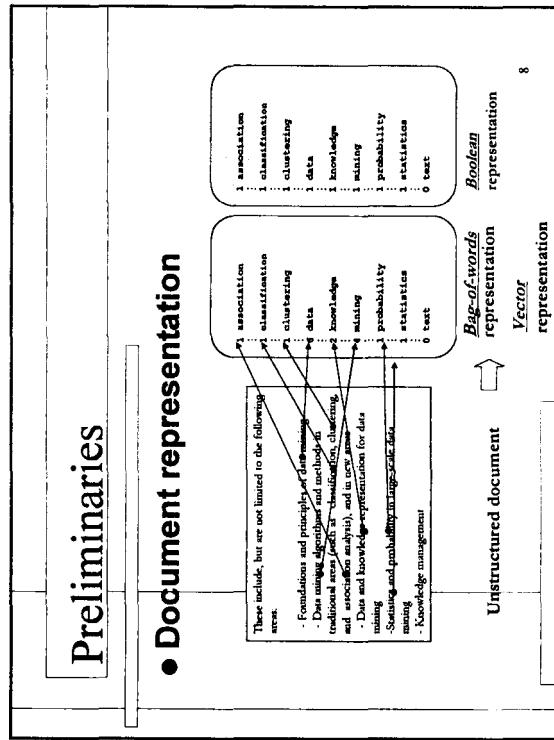
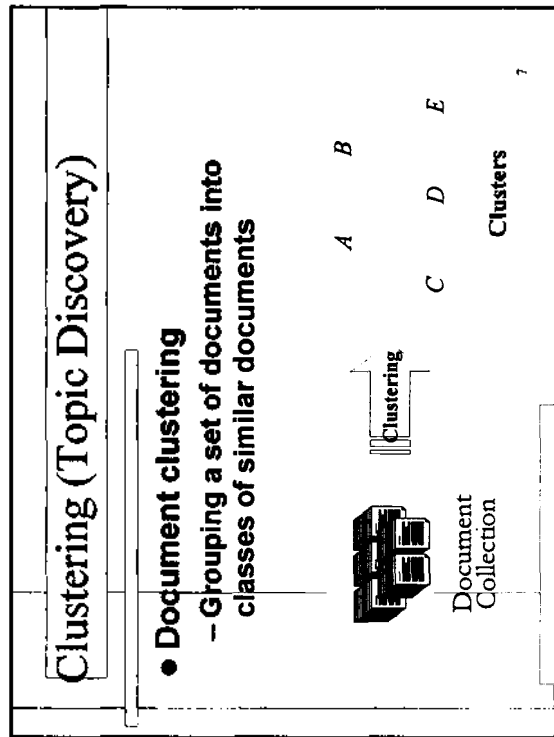
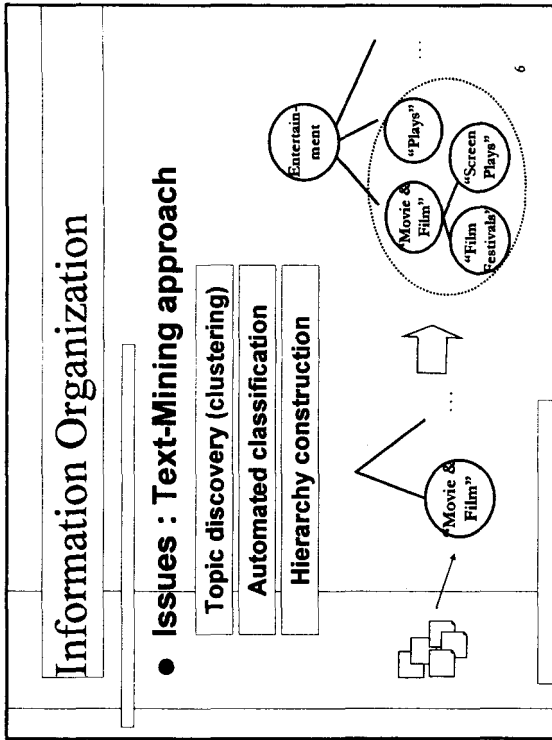
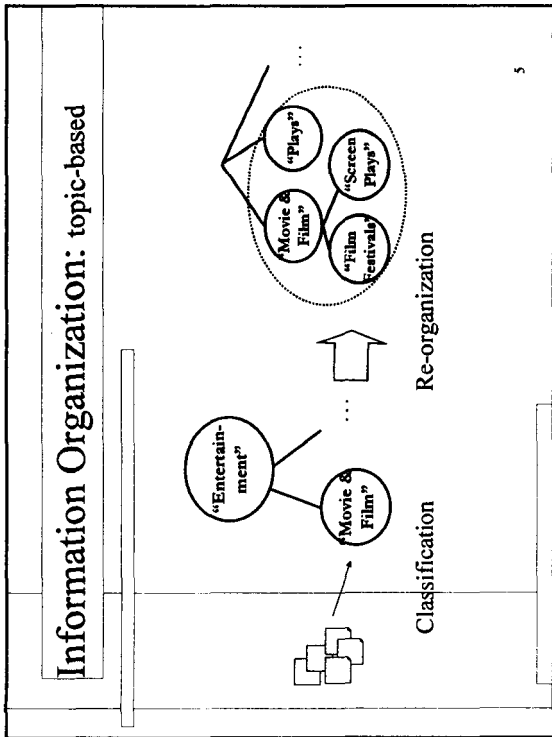
- **Indexing**
 - Key word-based
 - Polysemy/Synonymy problem
 - Topic-based
 - Topic hierarchy: Yahoo-style

3

Introduction

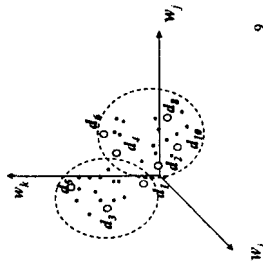
- **Topic hierarchy paradigm**
 - Summary view of a large amount of information
- **Management of a large document collection**
 - Topic hierarchy based Information System
 - Web search engine: Yahoo!, Lycos, Open Directory
 - Digital libraries: Library of Congress Catalog
 - Topic-oriented indexing
 - Browsing directories
- **Maintenance of topic hierarchy**
 - limitation of manual maintenance





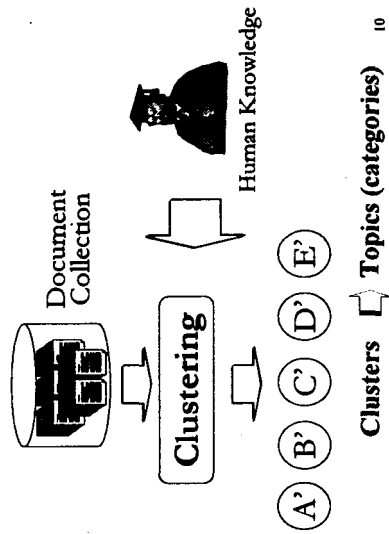
Clustering (Topic Discovery)

- Conventional clustering methods
 - Hierarchical clustering methods
 - Depends on the standard Euclidean distance metric
 - Not suitable for topic discovery
- Goal
 - Generate user-specific clusters!



9

Supervised Clustering (Topic Discovery)



10

Distance Functions

- Euclidean distance function

$$dist_w(\vec{d}_x, \vec{d}_y) = \sqrt{(\vec{d}_x - \vec{d}_y)^T \cdot (\vec{d}_x - \vec{d}_y)}$$

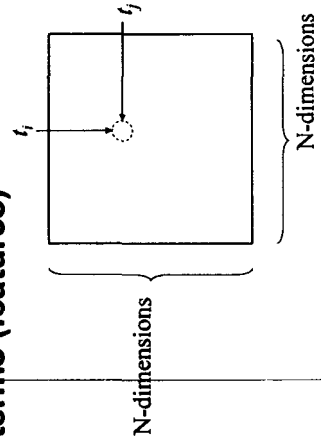
- Quadratic distance function

$$dist_w(\vec{d}_x, \vec{d}_y) = \sqrt{(\vec{d}_x - \vec{d}_y)^T \cdot \mathbf{W} \cdot (\vec{d}_x - \vec{d}_y)}$$

11

Weight Matrix

- Representing the inter-correlation of terms (features)



12

Quadratic Distance Function

Examples

- Vocabulary = {'digital', 'image', 'picture'}
- "e-book with digital images"

$$d_1 = (1, 1, 0)$$

$$d_2 = (0, 0, 1)$$

Weight Matrix

- Distance values

$$dist_{\text{euclidean}}(d_1, d_2) = \sqrt{3}$$

$$dist_{\text{manhattan}}(d_1, d_2) = \sqrt{1.2}$$

$$dist_{\text{minkowski}}(d_1, d_2) = \sqrt{3.5}$$

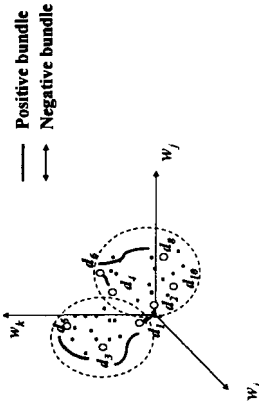
$$W_{\text{manhattan}} = \begin{pmatrix} d & i & p \\ 1.0 & 0.0 & 0.0 \\ 0.0 & 1.0 & 0.8 \\ 0.0 & 0.8 & 1.0 \end{pmatrix}$$

$$W_{\text{euclidean}} = \begin{pmatrix} d & i & p \\ 1.0 & 0.5 & 0.0 \\ 0.5 & 1.0 & 0.0 \\ 0.0 & 0.0 & 1.0 \end{pmatrix}$$

13

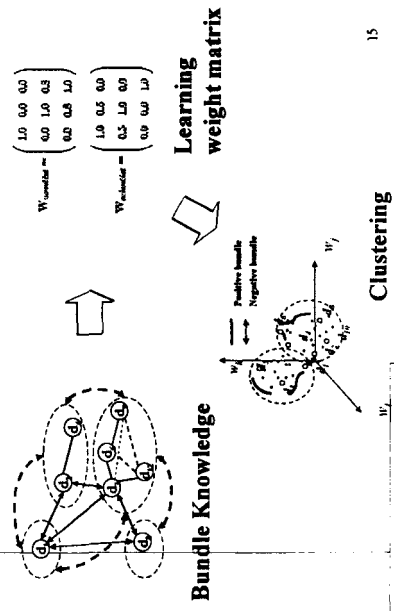
Human Knowledge

- Document bundles (or bundles)
 - Positive/Negative bundles



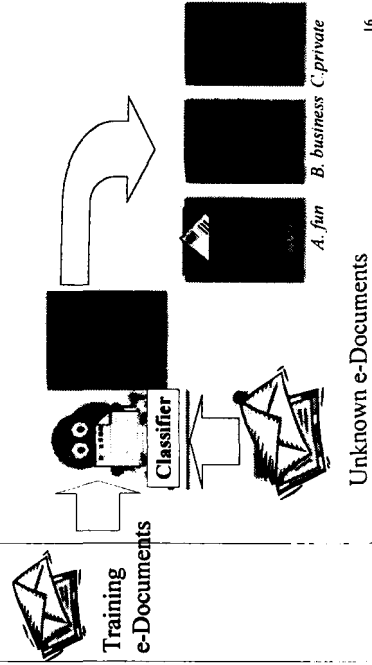
14

Supervised Clustering: overview



15

Automated Classification



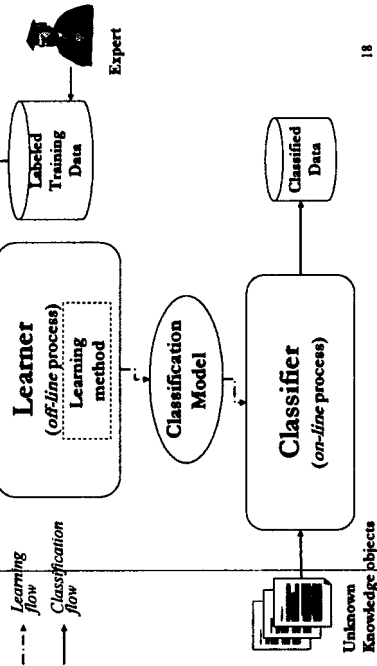
16

Automated Classification

Factors to Consider	Manual	Automated
Cost (per document)	High	Low
Speed	Slow	Fast
Consistency	Variable	High
Quality	Variable	Variable

17

Machine Learning based approach (Basic architecture)



18

Automated Classification

- **Related algorithms**
 - Decision trees
 - Neural networks
 - Support vector machines
 - Bayesian statistics
 - etc

19

Automated Classification

- **Critical Issue**
 - Difficulty in obtaining good quality training data
- **Related researches**
 - **Active learning**
 - Actively selecting the best training data out of unclassified ones
 - **EM algorithms**
 - Augment of labeled (training) documents with classifying unlabeled documents

20

Active Learning

● **Uncertainty-based sampling**

● ▲ 학습문서 비 학습문서
 — 문서 분포 함수가 판단한 경계 - - - 실제 경계

Uncertainty/Density-based Selective Sampling

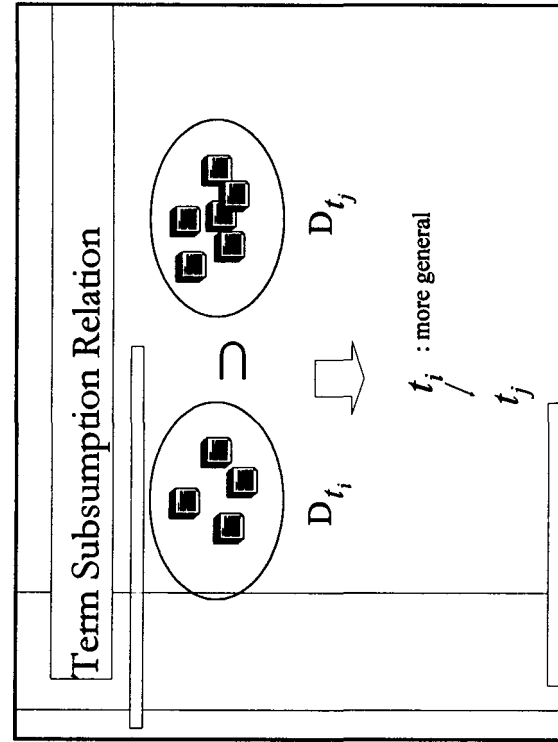
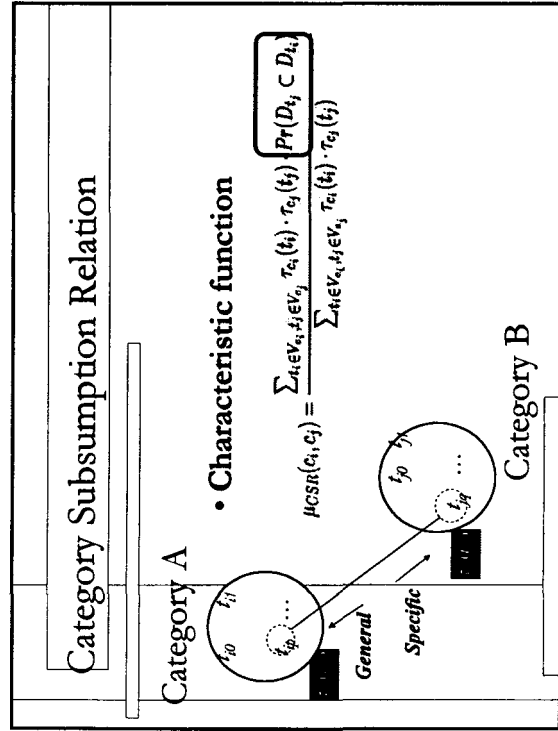
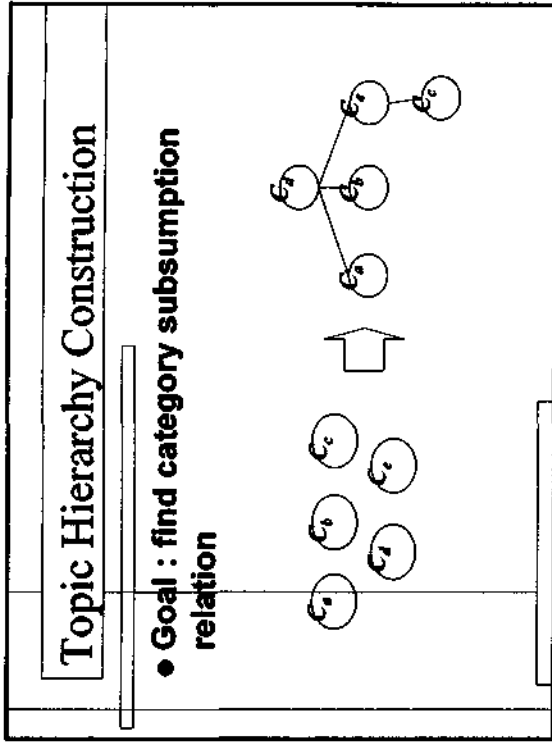
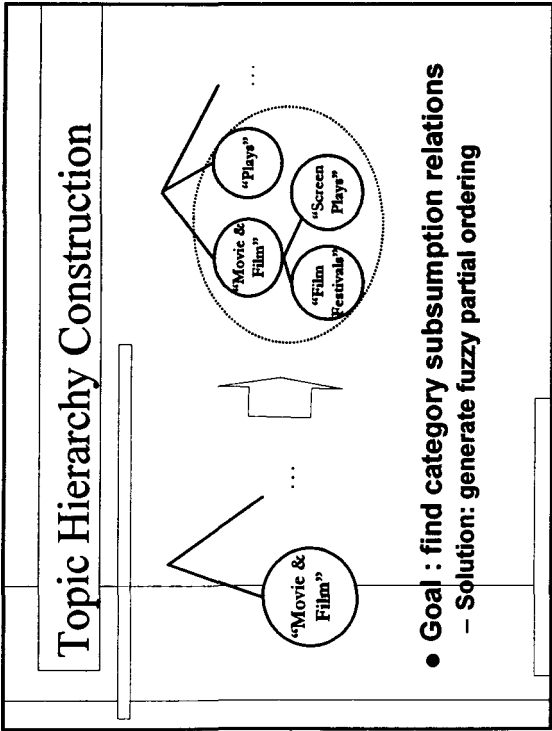
○ Unclassified data
 c_i, c_j, c_k categories
 ① ② outliers
 Uncertainty 大 Density 小
 Uncertainty 大 Density 大
 Previous decision boundary of c_i
 Current decision boundary of c_j

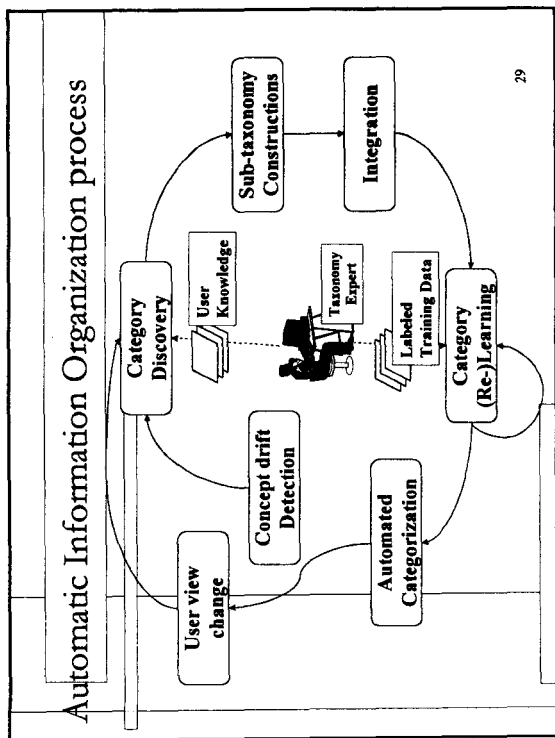
EM (Expectation Maximization) algorithm

Initial training set D^{tl} → $\hat{\theta}$ (①)
 $\hat{\theta}$ → f_{θ} (②)
 f_{θ} → D^{tc} (③)
 D^{tc} → D^{tu}

Machine Learning based approach (The proposed architecture)

Learning flow (dashed arrow)
 Classification flow (solid arrow)



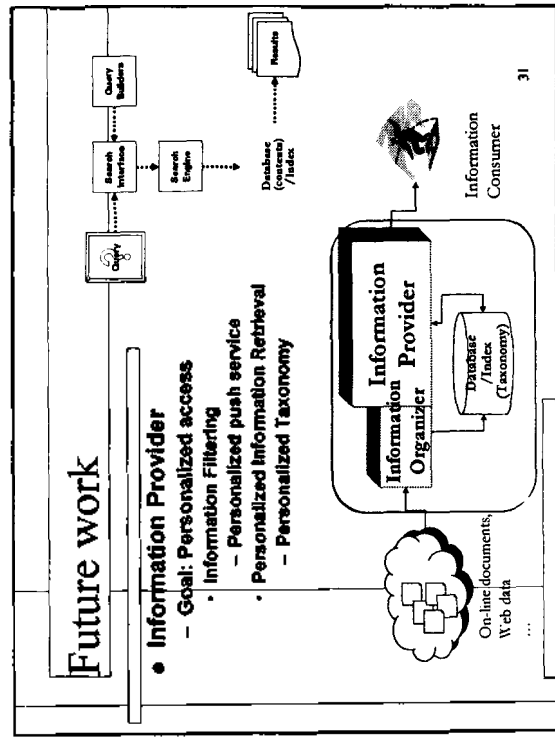


29

Summary

- Text-mining approach to e-documents organization based on topic hierarchy
 - Machine-Learning & Information Theory-based
 - 'Category (topic) discovery' problem
 - document bundle-based user-constraint document clustering
 - 'Automatic categorization' problem
 - Accelerated EM with CU-based active learning
 - 'Hierarchy Construction' problem
 - Unsupervised learning of category subsumption relation

30



31

Building Topic Hierarchy of e-Documents using Text Mining Technology

Thank you

Han-joon Kim
 Email: khj@uos.ac.kr
 Tel: 02-2210-5632

32