

Index Structures of User Profiles for Efficient Web Page Filtering Services

Y. H. Wu

A. L. P. Chen

Database Laboratory
Department of Computer Science
National Tsing Hua University

4/13/2000

Outline

- Introduction**
- Related Approaches**
- Our Approaches**
- Comparisons**
- Conclusion**

Introduction

□ Motivation

- ▶ Searching problem on the WWW
 - ☞ search engine
 - ☞ meta-search engine
- ▶ The performance may get worse if the number of web pages grows rapidly

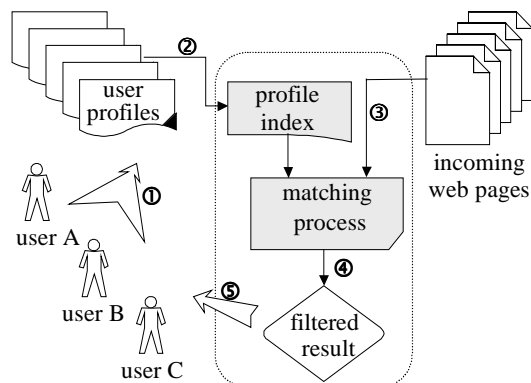
□ Goal

- ▶ Filtering approach
 - ☞ find the matched profiles for each web page

1

Introduction

□ A Web Page Filtering Service



2

Introduction

□ An Example

- ▶ Conjunction of keywords
- ▶ Boolean model
- ▶ Matches
 - ☞ $\{P_1, P_4\}$

Profile	Keyword
P ₁	a b
P ₂	a d
P ₃	a d e
P ₄	b f
P ₅	c d e f

Example page
a b c f

3

Related Approaches

□ The Counting Method

- ▶ Keyword array: inverted lists
- ▶ Profile arrays: TOTAL, COUNT
- ▶ Matching criteria: COUNT=TOTAL

Profile	TOTAL	COUNT	Keyword
P ₁	2	2	a → P ₁ P ₂ P ₃
P ₂	2	1	b → P ₁ P ₄
P ₃	3	1	c → P ₅
P ₄	2	2	d → P ₂ P ₃ P ₅
P ₅	4	2	e → P ₃ P ₅
			f → P ₄ P ₅

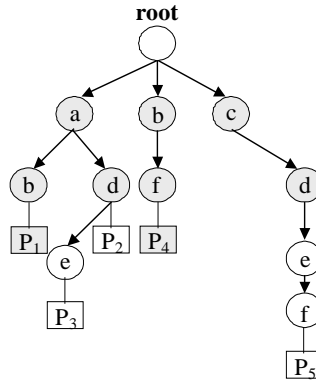
Example page
a b c f

4

Related Approaches

□ The Tree Method

- ▶ K-nodes: internal nodes
- ▶ P-nodes: leaf nodes
- ▶ External path
 - ☞ root → p-node
 - ☞ a profile
- ▶ Matches
 - ☞ root → a → b → P₁
 - ☞ root → b → f → P₄



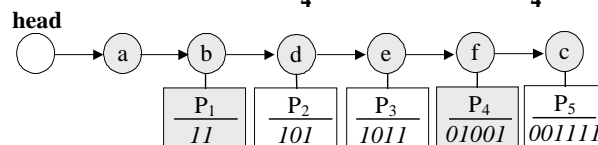
Example page
a b c f

5

Our Approaches

□ Method 1

- ▶ Index path with path signatures
- ▶ Path signature of the example page: 110011
- ▶ Matches
 - ☞ at node b: P₁ AND 11 = P₁
 - ☞ at node f: P₄ AND 11001 = P₄



Example page
a b c f

6

Our Approaches

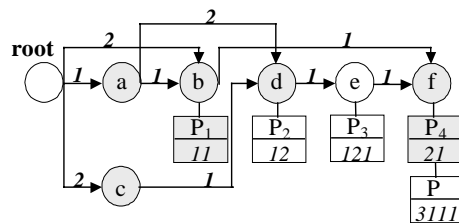
□ Method 2

▶ Index graph with path signatures

▶ Matches

☞ root → a → b → P₁: 11

☞ root → b → f → P₄: 21



Example page

a b c f

7

Our Approaches

□ Method 3

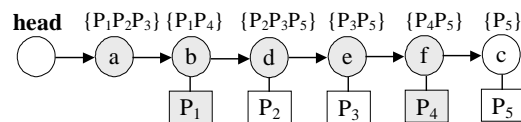
▶ Index path with profile sets

▶ Candidate set

▶ Target set \leftarrow candidate set \cap profile set

☞ at node b: $T = \{P_1P_2P_3P_4P_5\} \cap \{P_1P_4\} = \{P_1P_4\}$

☞ at node d: $T = \{P_2P_3P_4P_5\} \cap \{P_2P_3P_5\} = \{P_2P_3P_5\}$



Example page

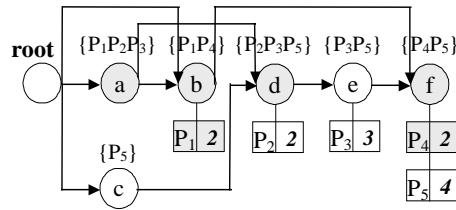
a b c f

8

Our Approaches

□ Method 4

- ▶ Index graph with profile sets
- ▶ Path length
- ▶ Matching criteria: matching of keywords \wedge profile id \subset target set \wedge equal path length



Example page

a b c f

9

Comparisons

□ Notation

Symbol	Description
P	The set of all profiles
K	The set of all distinct keywords
n	Average number of keywords in a profile
f	Average number of profiles in which a specific keyword is specified
m	Average number of keywords to represent a web page

10

Comparisons

□ Summary

Approaches Criteria	Counting Method	Tree Method	Method 1	Method 2	Method 3	Method 4
Duplication of information	profile	keyword	no	no	no	no
Sorting of keywords	no	yes	no	yes	no	yes
Storage space	$O(P +f K)$	$O(n P)$	$O(P + K)$	$O(P + K)$	$O(P + K)$	$O(P + K)$
Insertion/Deletion time	$O(nf)$	$O(nf)$	$O(n+f)$	$O(nf)$	$O(n+f)$	$O(n+f)$
Matching time	$O(mf+ P)$	$O(mf)$	$O(mf)$	$O(mf)$	$O(mf)$	$O(mf)$
Modification time	$O(nf)$	$O(nf)$	$O(n+f)$	$O(nf)$	$O(n+f)$	$O(n+f)$

11

Conclusion

□ Contribution

- ▶ Four new methods for profile indexing
- ▶ Comparisons by complexity analyses
- ▶ Efficient web page filtering service

□ Future Work

- ▶ Prototype system for real data
- ▶ Dissemination and display of the filtered results
- ▶ More predicates for specifying the user profiles

12