

# Efficient Data Mining for Path Traversal Patterns

Qiang Ye

Nov. 29, 2000

## Content

1. Introduction
2. Identifying Maximal Forward References
3. Algorithm on Full Scan(FS)
4. Algorithm on Selective Scan(SS)
5. Performance Evaluation
6. Conclusion

## Introduction

- Objective

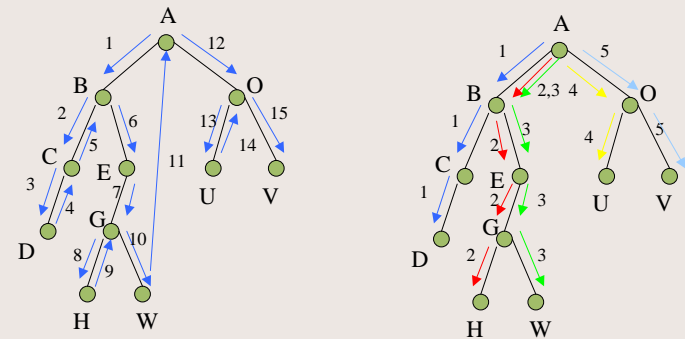
- To understand user access patterns in a distributed information-provided environment.

- Benefit

- Improving the system design: e.g., providing efficient access between highly correlated objects.

- Leading to better marketing decisions: e.g., putting advertisements in proper places.

## Maximal Forward Reference



Traversal path for a user:

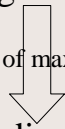
{A,B,C,D,C,B,E,G,H,G,W,A,O,U,O,V}

Maximal Forward References

{ABCD,ABEGH,ABEGW,AOU,AOV}

## Problem of finding frequent traversal patterns

Using the concept of maximal forward references



Problem of finding frequent occurring consecutive subsequences among all maximal forward references

- **Large reference sequence:** a reference sequence that appeared in a sufficient number of times.
- **Maximal reference sequence:** a large reference sequence that is not contained in any other maximal reference sequence.

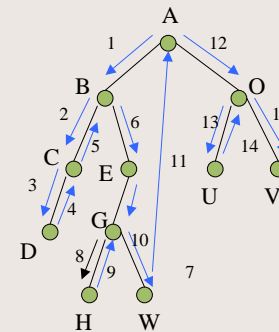
## Procedure of mining traversal patterns

- **Step 1:** Determine **maximal forward references** from the original log data.
- **Step 2:** Determine **large reference sequences** from the set of maximal forward references.
- **Step 3:** Determine **maximal reference sequences** from large reference sequences.

```

Step 1: Set  $i = 1$  and string  $Y$  to null for initialization,
where string  $Y$  is used to store the current forward
reference path. Also, set the flag  $F = 1$  to indicate a
forward traversal.
Step 2: Let  $A = s_i$  and  $B = d_j$ .
If  $A$  is equal to null then
/* this is the beginning of a new traversal */
begin
  Write out the current string  $Y$  (if not null) to the
  database  $D_F$ ;
  Set string  $Y = B$ ;
  Go to Step 5.
end
Step 3: If  $B$  is equal to some reference (say the  $j$ th refer-
ence) in string  $Y$  then
/* this is a cross-referencing back to a previous
reference */
begin
  If  $F$  is equal to 1 then write out string  $Y$  to
  database  $D_F$ ;
  Discard all the references after the  $j$ th one in
  string  $Y$ ;
   $F = 0$ ;
  Go to Step 5.
end
Step 4: Otherwise, append  $B$  to the end of string  $Y$ .
/* we are continuing a forward traversal */
If  $F$  is equal to 0, set  $F = 1$ .
Step 5: Set  $i = i + 1$ . If the sequence is not completed
scanned then go to Step 2.
    
```

## Example of Algorithm MF



Traversal path for a user:

{A,B,C,D,C,B,E,G,H,G,W,A,O,U,O,V}

move	string Y	output to $D_F$
1	AB	-
2	ABC	-
3	ABCD	-
4	ABC	ABCD
5	AB	-
6	ABE	-
7	ABEG	-
8	ABEGH	-
9	ABEG	ABEGH
10	ABEGW	-
11	A	ABEGW
12	AO	-
13	AOU	-
14	AO	AOU
15	AOV	AOV(end)

## Characteristics of FS

Using key ideas of the DHP ( Direct Hashing and Pruning) algorithm, FS has two distinct characteristics:

- Hash Technique: efficient for the generation of candidate reference sequence
- Pruning technique: progressively reduce the transaction database size.

## Traditional Generation Method

Database $D$		$C_1$		$L_1$	
TID	Items	Itemset	Sup.	Itemset	Sup.
100	A C D	{A}	2	{A}	2
200	B C E	{B}	3	{B}	3
300	A B C E	{C}	3	{C}	3
400	B E	{D}	1	{E}	3
		{E}	3		

$C_2$		$C_2$		$L_2$	
Itemset	Scan $D$	Itemset	Sup.	Itemset	Sup.
{A B}	→	{A B}	1	{A C}	2
{A C}	→	{A C}	2	{B C}	2
{A E}	→	{A E}	1	{B E}	3
{B C}	→	{B C}	2	{C E}	2
{B E}	→	{B E}	3		
{C E}	→	{C E}	2		

$C_3$		$C_3$		$L_3$	
Itemset	Scan $D$	Itemset	Sup.	Itemset	Sup.
{B C E}	→	{B C E}	2	{B C E}	2

## Hash Technique

$C_1$	count	
{A}	2	100 {A C}, {A D}, {C D}
{B}	3	200 {B C}, {B E}, {C E}
{C}	3	300 {A B}, {A C}, {A E}, {B C}, {B E}, {C E}
{D}	1	400 {B E}
{E}	3	

## Hash Technique

Hash Function:  $h(\{xy\}) = ((\text{order of } x) * 10 + (\text{order of } y)) \bmod 7$

Hash address	0	1	2	3	4	5	6
Hash table $H_2$		{A E}	{B C}, {B C}		{B E}, {B E}, {B E}	{A B}	{A C}, {C D}, {A C}
	3	1	2	0	3	1	3

The number of items hashed to bucket 2

Generating $C_2$	# in a bucket with the itemset	$C_2$
$L_1 \times L_1$	{A B} 1	{A C}
	{A C} 3	{B C}
	{A E} 1	{B E}
	{B C} 2	{C E}
	{B E} 3	
	{C E} 3	

## Pruning Technique

- Trimming the size of each transaction:

$\{A B C D\} \Rightarrow \{A B C\}$

- Reducing the number of the transactions in the database:

TID	Items		TID	Items
100	ACD	$\Rightarrow$	200	BCE
200	BCE			

## Pruning Technique

- Any subset of a large itemsets must be a large itemsets by itself.
- A transaction can be used to determine the set of large  $(k+1)$  itemsets only if it consists of  $(K+1)$  large  $k$ -itemsets.
- An item in transaction  $t$  can be trimmed if it does not appear in at least  $k$  of the candidate  $k$ -itemsets.

## Pruning Technique

$C_2$	count		$L_2$
{A C}	2	$\longrightarrow$	{A C}
{B C}	2		{B C}
{B E}	3		{B E}
{C E}	2		{C E}

TID	Items	
100	ACD	{AC} $\longrightarrow$ Discard
200	BCE	{BC} {BE} {CE} $\longrightarrow$ Keep {BCE}

For 100:

A	B	C
$a[0]=1 < 2$	$a[1]=1 < 2$	$a[2]=0 < 2$

## Algorithm on Selective Scan

- FS generates a **small number** of candidate two-references by using a hashing technique.
- If memory is enough, we can generate  $C_3$  from  $C_2 * C_2$  instead of  $L_2 * L_2$ , hence saving one round of database scan.
- Using this concept, we can determine all  $L_k$ s by as few as **two** scans.

## Algorithm on Selective Scan

- Using  $L_2 * L_2$  to generate  $C_3$ ,  $C_3$  would be :  $\{\{BCE\}\}$

$C_2$			$C_2$			$L_2$	
Itemset		Scan	Itemset	Sup.		Itemset	Sup.
{A B}		$D$	{A B}	1		{A C}	2
{A C}		$\rightarrow$	{A C}	2		{B C}	2
{A E}			{A E}	1		{B E}	3
{B C}			{B C}	2		{C E}	2
{B E}			{B E}	3			
{C E}			{C E}	2			

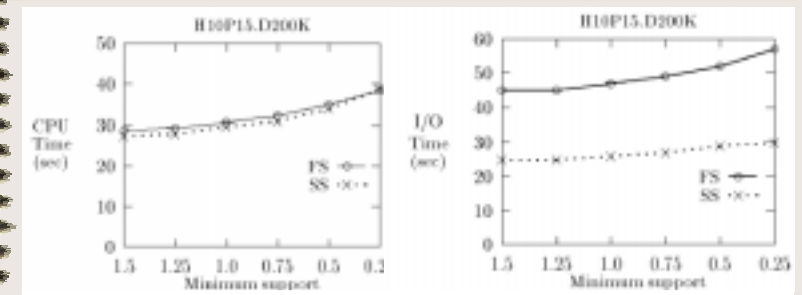
  

$C_3$			$C_3$			$L_3$	
Itemset		Scan	Itemset	Sup.		Itemset	Sup.
{B C E}		$D$	{B C E}	2		{B C E}	2
		$\rightarrow$					

- Using  $C_2 * C_2$  to generate  $C_3$ ,  $C_3$  would be :

$\{\{ABC\},\{ABE\},\{ACE\},\{BCE\}\}$

## Performance Evaluation



-SS in general outperforms FS.

-When I/O cost is taken into account, their difference becomes prominent.

## Conclusion

- Data Mining is application-dependent.
- How to choose hash functions?
- Can we and need we use hash tree to count the number of candidate reference sequence?
- Are the algorithm assumptions reasonable?