

Extracting Structured Data from Web Pages

Summarized by Jaehyok Chong

Algorithm: ExAlg – under revision

Input: Set of HTML pages I, Output: Template, Schema, Values

Function ECGM (Equivalence Class Generation Module)

 DiffFormat(I, T); // differentiate roles using HTML format

 while (T is not empty) {

 FindEquiv(T, L, Q); // find equivalence classes from T & put LFEQs in L

 HandInv(L, Q); // handle invalid equivalence classes in L

 DiffEq(Q, T); // differentiate roles using equivalence classes

 }

End ECGM

Function AM (Analysis Module)

 ConstTemp(L, S);

 ExVal(I, S, V);

End AM

Main procedure ExAlg(I)

 ECGM(I);

 AM(L);

End ExAlg

Example

```
<html>1<body>2
  <b>3Book4Name5</b>6Databases
  <b>7Reviews8</b>9
  <ol>10
    <li>11
      <b>12Reviewer13Name14</b>15John
      <b>16Rating17</b>18 7
      <b>19Text20</b>21 ...
    </li>22
  </ol>23
</body>24</html>25
```

```
<html>1<body>2
  <b>3Book4Name5</b>6Query Opt.
  <b>7Reviews8</b>9
  <ol>10
    <li>11
      <b>12Reviewer13Name14</b>15John
      <b>16Rating17</b>18 8
      <b>19Text20</b>21 ...
    </li>22
  </ol>23
</body>24</html>25
```

```
<html>1<body>2
  <b>3Book4Name5</b>6Data Mining
  <b>7Reviews8</b>9
  <ol>10
    <li>11
      <b>12Reviewer13Name14</b>15Jeff
      <b>16Rating17</b>18 2
      <b>19Text20</b>21 ...
    </li>22
    <li>11
      <b>12Reviewer13Name14</b>15Jane
      <b>16Rating17</b>18 6
      <b>19Text20</b>21 ...
    </li>22
  </ol>23
</body>24</html>25
```

```
<html>1<body>2
  <b>3Book4Name5</b>6Transactions
  <b>7Reviews8</b>9
  <ol>10
  </ol>23
</body>24</html>25
```

DiffFormat

- ❑ Extract tokens from HTML pages
- ❑ Differentiate roles of tokens using parse tree (structural context)
 - `()`
 - ❑ `1(1): <html><body>()`
 - ❑ `2(2): <html><body>()`
 - Name
 - ❑ `Name1: <html><body>Name`
 - ❑ `Name2: <html><body>Name`

DiffFormat

Tokens		Tokens	
<html>		Reviewer	
<body>		Name² (<html><body>Name)	
¹ (<html><body>)		² (<html><body>)	
Book		John	
Name¹ (<html><body>Name)		Jeff	
¹ (<html><body>)		Jane	
Databases		Rating	
Data		7	
Mining		2	
Query		6	
Opt.		8	
Transactions		Text	
Reviews			
			
		</body>	
² (<html><body>)		</html>	

FindEquiv – 1st iteration

- Make occurrence vectors for each tokens generated by DiffFormat()
 - Find Equivalence classes using occurrence vectors
 - Find LFEQs using size & support thresholds
 - SizeThre = 3, SupThres = 3
-

FindEquiv – 1st iteration

Tokens	Occurrence	Tokens	Occurrence
<html>	1, 1, 1, 1 (<i>E1</i>)	Reviewer	1, 2, 1, 0 (<i>E7</i>)
<body>	1, 1, 1, 1 (<i>E1</i>)	Name ² (<html><body>Name)	1, 2, 1, 0 (<i>E7</i>)
 ¹ (<html><body>)	2, 2, 2, 2 (<i>E2</i>)	 ² (<html><body>)	3, 6, 3, 0 (<i>E8</i>)
Book	1, 1, 1, 1 (<i>E1</i>)	John	1, 0, 1, 0 (<i>E9</i>)
Name ¹ (<html><body>Name)	1, 1, 1, 1 (<i>E1</i>)	Jeff	0, 1, 0, 0 (<i>E4</i>)
 ¹ (<html><body>)	2, 2, 2, 2 (<i>E2</i>)	Jane	0, 1, 0, 0 (<i>E4</i>)
Databases	1, 0, 0, 0 (<i>E3</i>)	Rating	1, 2, 1, 0 (<i>E7</i>)
Data	0, 1, 0, 0 (<i>E4</i>)	7	1, 0, 0, 0 (<i>E3</i>)
Mining	0, 1, 0, 0 (<i>E4</i>)	2	0, 1, 0, 0 (<i>E4</i>)
Query	0, 0, 1, 0 (<i>E5</i>)	6	0, 1, 0, 0 (<i>E4</i>)
Opt.	0, 0, 1, 0 (<i>E5</i>)	8	0, 0, 1, 0 (<i>E5</i>)
Transactions	0, 0, 0, 1 (<i>E6</i>)	Text	1, 2, 1, 0 (<i>E7</i>)
Reviews	1, 1, 1, 1 (<i>E1</i>)		1, 2, 1, 0 (<i>E7</i>)
	1, 1, 1, 1 (<i>E1</i>)		1, 1, 1, 1 (<i>E1</i>)
	1, 2, 1, 0 (<i>E7</i>)	</body>	1, 1, 1, 1 (<i>E1</i>)
 ² (<html><body>)	3, 6, 3, 0 (<i>E8</i>)	</html>	1, 1, 1, 1 (<i>E1</i>)

FindEquiv – 1st iteration

- Find LFEQs using SizeThre = 3, SupThres = 3
 - E_1 : <html> <body> Book Name¹ Reviews </body> </html>
 - E_7 : Reviewer Name² Rating Text

Equivalence Classes	SizeThre	SupThres	Equivalence Classes	SizeThre	SupThres
E_1	9	4	E_6	1	1
E_2	2	4	E_7	6	3
E_3	2	1	E_8	2	3
E_4	6	1	E_9	1	2
E_5	3	1			

HandInv – 1st iteration

- Handle invalid equivalence classes
 - [Observation 4.3] A valid equivalence class is ordered and a pair of two valid equivalence classes is nested
 - Find invalid LFEQs determined by FindEquiv()
 - In this iteration, there is no invalid LFEQs.
-

Ordered Equivalence classes

An equivalence class is **ordered**, if its tokens can be ordered $\langle t_1, \dots, t_m \rangle$, such that, for every page p_i ($1 \leq i \leq n$), and every pair of tokens t_j, t_k ($1 \leq j < k \leq m$),

- I. If t_j occurs at least l times in p_i , the l^{th} occurrence of t_j in p_i occurs before the l^{th} occurrence of t_k in p_i , and
 - II. If t_j occurs at least $(l+1)$ times in p_i , the $(l+1)^{\text{st}}$ occurrence of t_j in p_i is after the l^{th} occurrence of t_k in p_i .
-

Nesting of Equivalence classes

A pair of equivalence classes, E_i and E_j is **nested** if,

- I. The span of any occurrence of E_i does not overlap with the span of any occurrence of E_j , or
 - II. The span of all occurrences of E_j is within $\text{Pos}(p)$ of some occurrence of E_i for some fixed p ; or vice versa.
-

Remove either one arbitrarily

□ ADEFBDEFC

■ ABC

■ DEF

Causes of invalid equivalence classes

- I. Correlation of instantiation of two or more type constructors
 - II. Frequently Occurring tuples
 - III. Overlap in the tokens associated with different type constructors
- Type A (I + II) is considered to be processed by HandInv and Type B (III) is just removed by HandInv
-

Almost-Ordered

An equivalence class E is **almost-ordered** if the tokens in the equivalence class can be partitioned into ordered equivalence classes that do not overlap

- $E = \{A, B, C, D\}$
 $ABABCD, ABCD, ABABABCD$
→ divided into two equivalence classes of $\{A, B\}$ and $\{C, D\}$ which are ordered and do not overlap with each other
 - Try All possible partitions
 - Binary Partition only?
 - $ABABCD, ABCDEF, ABABABCD, ABCDEF, ABCDEF$
 - $\{A, B\}$ and $\{C, D, E, F\}$?
 - $\{A, B\}$ and $\{C, D\}$ and $\{E, F\}$?
-

Nearly-Nested

Two equivalence classes, E_i and E_j are said to be **nearly-nested** if for every token $t \in E_i$ (and vice versa) there exists a position p , such that each occurrence of t is either outside the span of any occurrence of E_j , or within $\text{Pos}(p)$ of some occurrence of E_j .

- $E_1 = \{A, B, C, D\}$, $E_2 = \{E, F, G, H\}$
ABCDA(EF)BC(GH)D, A(EF)BC(GH)D, EFGH
→ divide E_1 into three equivalence classes of $\{A\}$, $\{B, C\}$, $\{D\}$ and E_2 into two equivalence classes of $\{E, F\}$ and $\{G, H\}$
 - $E_1 = \{A, B, C, D\}$, $E_2 = \{E, F, G, H\}$
A(EF)BC(GH)D, A(EF)BC(GH)D, EFGH
→ divide E_1 into three equivalence classes of $\{A\}$, $\{B, C\}$, $\{D\}$ and E_2 into two equivalence classes of $\{E, F, G, H\}$
-

DiffEq – 1st iteration

- Differentiate roles of tokens which are not in the set of LFEQs using LFEQs
- For the **tokens satisfies SupThres**, differentiate the roles
 - $\langle b \rangle^1, \langle /b \rangle^1, \langle b \rangle^2, \langle /b \rangle^2$
 - $\text{Pos}(i, E)$ represents i -th position($\text{Pos}(i)$) in $E = \langle t_1, \text{Pos}(1), t_2, \text{Pos}(2), t_3, \text{Pos}(3), \dots, \text{Pos}(n - 1), t_n \rangle$

Tokens	Differentiated roles	Dtokens
$\langle b \rangle^1$	$\text{Pos}(2, E1), \text{Pos}(4, E1)$	$\langle b \rangle^1_1, \langle b \rangle^1_2$
$\langle /b \rangle^1$	$\text{Pos}(4, E1), \text{Pos}(5, E1)$	$\langle /b \rangle^1_1, \langle /b \rangle^1_2$
$\langle b \rangle^2$	$\text{Pos}(1, E7), \text{Pos}(3, E7), \text{Pos}(4, E7)$	$\langle b \rangle^2_1, \langle b \rangle^2_2, \langle b \rangle^2_3$
$\langle /b \rangle^2$	$\text{Pos}(3, E7), \text{Pos}(4, E7), \text{Pos}(5, E7)$	$\langle /b \rangle^2_1, \langle /b \rangle^2_2, \langle /b \rangle^2_3$

An Example for DiffEq

□ ACDCBDF

- Assume ABF is in LFEQ, but C and D are not in LFEQ

□ ACDCBDF

FindEquiv – 2nd iteration

- For the newly added dtokens by DiffEq(), make occurrence vectors and equivalence classes
-

FindEquiv – 2nd iteration

Tokens	Occurrence	Tokens	Occurrence
<html>	1, 1, 1, 1 (<i>E1</i>)	Reviewer	1, 2, 1, 0 (<i>E7</i>)
<body>	1, 1, 1, 1 (<i>E1</i>)	Name ² (<html><body>Name)	1, 2, 1, 0 (<i>E7</i>)
¹₁	1, 1, 1, 1 (<i>E1</i>)	²₁	1, 2, 1, 0 (<i>E7</i>)
Book	1, 1, 1, 1 (<i>E1</i>)	²₂	1, 2, 1, 0 (<i>E7</i>)
Name ¹ (<html><body>Name)	1, 1, 1, 1 (<i>E1</i>)		
¹₁	1, 1, 1, 1 (<i>E1</i>)		
¹₂	1, 1, 1, 1 (<i>E1</i>)	Rating	1, 2, 1, 0 (<i>E7</i>)
		²₂	1, 2, 1, 0 (<i>E7</i>)
		²₃	1, 2, 1, 0 (<i>E7</i>)
		Text	1, 2, 1, 0 (<i>E7</i>)
Reviews	1, 1, 1, 1 (<i>E1</i>)	²₃	1, 2, 1, 0 (<i>E7</i>)
¹₂	1, 1, 1, 1 (<i>E1</i>)		1, 2, 1, 0 (<i>E7</i>)
	1, 1, 1, 1 (<i>E1</i>)		1, 1, 1, 1 (<i>E1</i>)
	1, 2, 1, 0 (<i>E7</i>)	</body>	1, 1, 1, 1 (<i>E1</i>)
²₁	1, 2, 1, 0 (<i>E7</i>)	</html>	1, 1, 1, 1 (<i>E1</i>)

FindEquiv – 2nd iteration

- Find LFEQs using SizeThre = 3, SupThres = 3
 - $E1$: `<html> <body> 11 Book Name1 11 12 Reviews 12 </body> </html>`
 - $E7$: ` 21 Reviewer Name2 21 22 Rating 22 23 Text 23 `

Equivalence Classes	SizeThre	SupThres	Equivalence Classes	SizeThre	SupThres
$E1$	13	4			
			$E7$	12	3

HandInv – 2nd iteration

- In this iteration, there is no invalid LFEQs.
-

DiffEq – 2nd iteration

- This time, there is no tokens to be differentiated by the contexts
 - Returns 2 LFEQs
 - $E1$: `<html> <body> 11 Book Name1 11 12 Reviews 12 </body> </html>`
 - $E7$: ` 21 Reviewer Name2 21 22 Rating 22 23 Text 23 `
-

ConstTemp

- Using the set of LFEQs found by ECGM (Equivalence Class Generation Module), construct the template for the input Web Pages
- From the root equivalence class (with occurrence vector $\langle 1, 1, \dots \rangle$), for each non-empty positions of each equivalence classes, check the *PosStrings* for the patterns like below table

	Pattern	Template
1	$E_j E_j E_j E_j E_j \dots$	$\{T_{E_j}\}_e$
2	$E_j S E_j S E_j S E_j S \dots E_j$	$\{T_{E_j}\}_S$
3	E_j or E_k	$T_{E_j} \mid T_{E_k}$
4	e or E_j	$(T_{E_j})?$
5	String of dtokens and empty equivalence classes	T_B
6	Unknown	T_B

ConstTemp

- Find root equivalence class: E_1
- Partition E_1 with non-empty positions
 - E_1 : (`<html><body>Book Name`)_{c₁} **PosString(E1, 6)**
 (`Reviews`)_{c₂} **PosString(E1, 10)**
 (`</body></html>`)_{c₃}
 - E_7 : (`Reviewer Name`)_{c₁} **PosString(E7, 1)**
 (`Rating`)_{c₂} **PosString(E7, 2)**
 (`Text`)_{c₃} **PosString(E7, 3)** (``)_{c₄}
- Check the PosStings to find the patterns

PosString	Pattern	Template
PosString(E_1 , 6)	String of dtokens	T_B
PosString(E_1 , 10)	$E_7 E_7 \dots$	$\{T_{E_7}\}_e$
PosString(E_7 , 5)	String of dtokens	T_B
PosString(E_7 , 8)	String of dtokens	T_B
PosString(E_7 , 11)	String of dtokens	T_B

ConstTemp

- Returns the template

- $\langle (\langle \text{html} \rangle \langle \text{body} \rangle \langle \text{b} \rangle \text{Book Name} \langle \text{/b} \rangle)_{c_1} \text{*}(\text{A})$
 $(\langle \text{b} \rangle \text{Reviews} \langle \text{/b} \rangle \langle \text{ol} \rangle)_{c_2} \{ (\langle \text{li} \rangle \langle \text{b} \rangle \text{Reviewer Name} \langle \text{/b} \rangle)_{c_4}$
 $\text{*}(\text{B}) (\langle \text{b} \rangle \text{Rating} \langle \text{/b} \rangle)_{c_5} \text{*}(\text{C}) (\langle \text{b} \rangle \text{Text} \langle \text{/b} \rangle)_{c_6} \text{*}(\text{D})$
 $(\langle \text{/li} \rangle)_{c_7} \}_e (\langle \text{/ol} \rangle \langle \text{/body} \rangle \langle \text{/html} \rangle)_{c_3} \rangle$
-

ExVal

- Using the template constructed by ConstTemp, extract values of each web pages

Page	A	B	C	D	
1	Databases	John	7	...	
2	Data Mining	Jeff	2	...	
		Jane	6	...	
3	Query Opt.	John	8	...	
4	Transactions	<i>e</i>	<i>e</i>	<i>e</i>	
