



Extracting Schema from Semi-structured Data

Svetlozar Nestorov
Serge Abiteboul
Rajeev Motwani

Stanford University
ACM SIGMOD Conference, 1998.

“Inferring Structure in Semi-structured Data”
*Workshop on Management of Semi-structured
Data, 1997.*

Yi-Hung Wu
1999/7/14



Outline

- **Introduction**
- **Lattice Approach**
- **Program Approach**
- **Experiment**
- **Conclusion**

2/15
Information



Introduction

■ Motivation

- ◆ homepages of group members
 - ☞ name, e-mail, address, photo
 - ☞ similar information, but irregular schema
- ◆ perfect typing
 - ☞ too large size
 - ☞ query optimization, graphical query interface

■ Goal

- ◆ approximate typing
 - ☞ tradeoff between the quality of a typing and its compactness



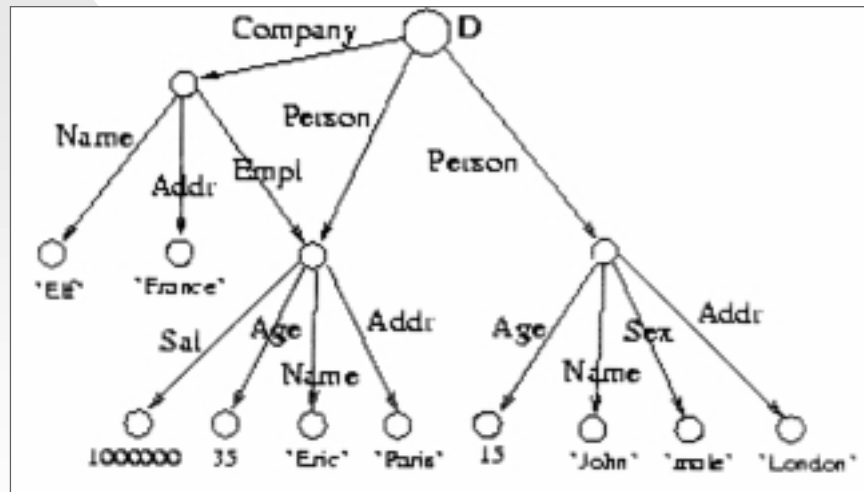
Introduction

■ Preliminary

- ◆ data model
 - ☞ Object Exchange Model (OEM)
 - ☞ labeled directed graph
- ◆ notation
 - ☞ $\text{attribute}(o) = \{\text{labels on outgoing edges}\}$
 - ☞ $\text{role}(o) = \{\text{labels on incoming edges}\}$
 - ☞ $\text{at}(S) = \{o \mid \text{attribute}(o) = S\}$
 - ☞ $\text{above}(S) = \{o \mid \text{attribute}(o) \supseteq S\}$
- ◆ definition
 - ☞ $\text{jump}(S) = \text{at}(S) / \text{above}(S)$

Introduction

■ OEM Database D



5/15
Information

Lattice Approach

■ Algorithm

- ◆ identify candidate types
- ◆ select types from candidates and organize them into a type hierarchy
- ◆ derive the typing rules
- ◆ validate the type hierarchy against the data

■ Counting lattice: $at(S)$

- ◆ significant jump: $jump(S) \geq \theta$

6/15
Information

Lattice Approach

- **Counting lattice**
 - ◆ $\text{jump}(\{\text{Addr Age Name Sal Sex}\})=1$
 - ◆ $\text{jump}(\{\text{Addr Empl Name Sub}\})=1$
 - ◆ $\text{jump}(\{\text{Addr Empl Name}\})=0.8$
 - ◆ $\text{jump}(\{\text{Addr Name Sal Sex}\})=0.6$

7/15
Information

Lattice Approach

- **Type hierarchy construction**
 - ◆ primary role: $p\text{-role}(S)$
 - ☞ the most frequent label in $\text{role}(o)$ for all candidates
 - ◆ select candidate T if there does not exist T' such that $T' \subset T$ and $p\text{-role}(T')=p\text{-role}(T)$

8/15
Information



Program Approach

■ Base relations

■ Typing rule

◆ $c(X):-A1 \& \dots \& A_n$

■ Type links

◆ $\text{link}(Y,X,l) \& c'(Y)$

◆ $\text{link}(X,Y,l) \& c'(Y)$

◆ $\text{link}(X,Y,l) \& \text{atomic}(Y,Z)$

■ Example typing rule

◆ $\text{person}(X):-\text{link}(X,Y,\text{is-manager-of}) \& \text{firm}(Y) \& \text{link}(X,Y',\text{name}) \& \text{atomic}(Y',Z)$

☞ $\text{person}(X):-\overrightarrow{\text{is-manager-of}}^{\text{firm}} \overrightarrow{\text{name}}^{\text{atomic}}$

link	FromObj	ToObj	Label
	<i>g</i>	<i>m</i>	is-manager-of
	<i>j</i>	<i>a</i>	is-manager-of
	<i>m</i>	<i>g</i>	is-managed-by
	<i>a</i>	<i>j</i>	is-managed-by
	<i>g</i>	<i>gn</i>	name
	<i>j</i>	<i>jn</i>	name
	<i>m</i>	<i>mn</i>	name
	<i>a</i>	<i>an</i>	name

atomic	Obj	Value
	<i>gn</i>	"Gates"
	<i>jn</i>	"Jobs"
	<i>mn</i>	"Microsoft"
	<i>an</i>	"Apple"



Program Approach

■ Example typing program

```

project : τ1 =  $\overrightarrow{\text{Project}}^3, \overrightarrow{\text{Project}}^4, \overrightarrow{\text{Project}}^5, \overrightarrow{\text{Project\_Member}}^3, \overrightarrow{\text{Project\_Member}}^4, \overrightarrow{\text{Name}}^0, \overrightarrow{\text{Homepage}}^0$ 
publication : τ2 =  $\overrightarrow{\text{Publication}}^3, \overrightarrow{\text{Publication}}^5, \overrightarrow{\text{Author}}^3, \overrightarrow{\text{Name}}^0, \overrightarrow{\text{Conference}}^0, \overrightarrow{\text{Postscript}}^0$ 
db-person : τ3 =  $\overrightarrow{\text{Project\_Member}}^1, \overrightarrow{\text{Group\_Member}}^5, \overrightarrow{\text{Project}}^1, \overrightarrow{\text{Birthday}}^5, \overrightarrow{\text{Degree}}^0, \overrightarrow{\text{Years\_At\_Stanford}}^0, \overrightarrow{\text{Email}}^0, \overrightarrow{\text{Home\_Page}}^0, \overrightarrow{\text{Title}}^0, \overrightarrow{\text{Name}}^0, \overrightarrow{\text{Original\_Home}}^0, \overrightarrow{\text{Personal\_Interest}}^0, \overrightarrow{\text{Research\_Interest}}^0$ 
student : τ4 =  $\overrightarrow{\text{Project\_Member}}^1, \overrightarrow{\text{Student}}^4, \overrightarrow{\text{Group\_Member}}^5, \overrightarrow{\text{Project}}^1, \overrightarrow{\text{Advisor}}^4, \overrightarrow{\text{Email}}^0, \overrightarrow{\text{Title}}^0, \overrightarrow{\text{Home\_Page}}^0, \overrightarrow{\text{Name}}^0, \overrightarrow{\text{Nickname}}^0$ 
birthday : τ5 =  $\overrightarrow{\text{Birthday}}^3, \overrightarrow{\text{Nameex}}^0, \overrightarrow{\text{Month}}^0, \overrightarrow{\text{Day}}^0, \overrightarrow{\text{Year}}^0$ 
degree : τ6 =  $\overrightarrow{\text{Degree}}^3, \overrightarrow{\text{Major}}^0, \overrightarrow{\text{School}}^0, \overrightarrow{\text{Name}}^0, \overrightarrow{\text{Year}}^0$ 

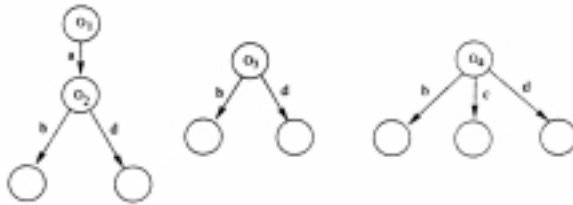
```



Program Approach

Defect

- ◆ excess
- ◆ deficit



Example

- ◆ type1: $\vec{-a^2}$
- ◆ type2: $\vec{-a^1, b^0, c^0}$
- ◆ type3: $\vec{-b^0, d^0}$
- ◆ type assignment: o_4
 - ☞ type2: excess=1, deficit=1
 - ☞ type3: excess=0, deficit=1



Program Approach

Minimal perfect typing

◆ step1

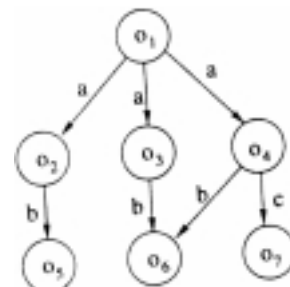
- ☞ type1: $\vec{-a^2, a^3, a^4}$
- ☞ type2: $\vec{-a^1, b^0}$
- ☞ type3: $\vec{-a^1, b^0}$
- ☞ type4: $\vec{-a^1, b^0, c^0}$

◆ step2

- ☞ equivalent class: type2=type3

◆ step3

- ☞ type1: $\vec{-a^2, a^3}$
- ☞ type2: $\vec{-a^1, b^0}$
- ☞ type3: $\vec{-a^1, b^0, c^0}$





Program Approach

■ Clustering

- ◆ merge similar types

- ☞ type1: \vec{a}^0, \vec{b}^3

- ☞ type2: \vec{a}^0, \vec{b}^4

- ☞ type3: \vec{a}^0, \vec{b}^1

- ☞ type4: \vec{a}^0, \vec{b}^2

- ◆ the order of coalescing has a significant effect on the quality of the results

- ◆ distance function

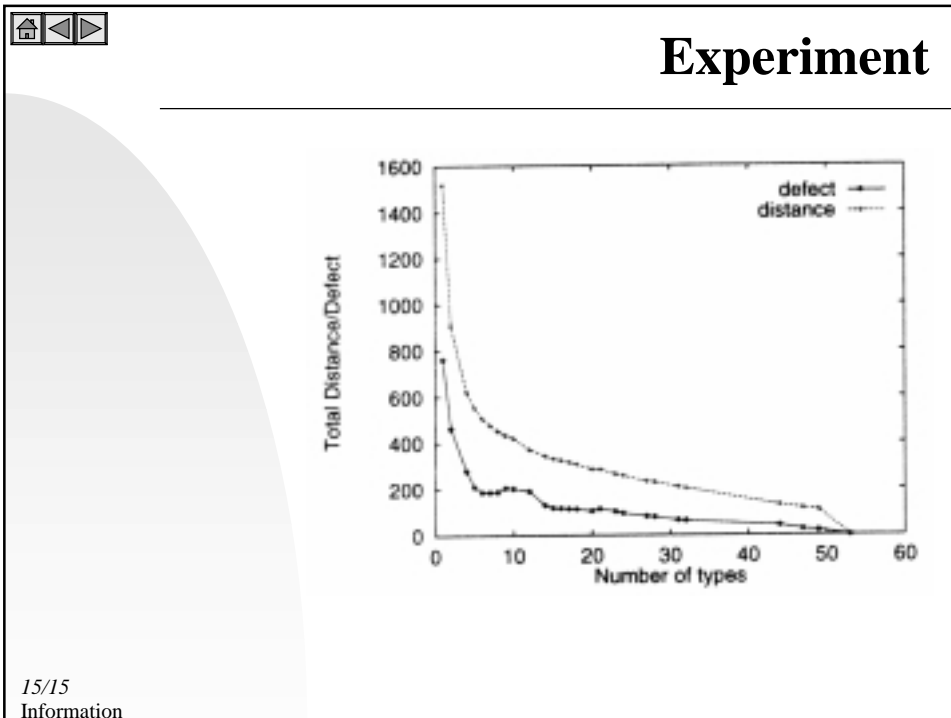
- ☞ symmetric difference

■ Recasting



Experiment

DB No	Synthetic Data						Typing		
	Bipartite ?	Overlap ?	Perturb ?	Intended Types	Objects	Links	Perfect Types	Optimal Types	Defect
1	Y	N	N	10	1500	2909	30	10	225
2	Y	N	Y	10	1500	2958	52	10	307
3	Y	Y	N	6	950	2409	19	6	239
4	Y	Y	Y	6	950	2442	35	6	283
5	N	N	N	5	400	726	317	5	181
6	N	N	Y	5	400	749	341	5	310
7	N	Y	N	5	400	775	375	5	291
8	N	Y	Y	5	400	795	381	5	333



- Navigation icons: Home, Previous, Next
- ## Conclusion
- **Contribution**
 - ◆ an algorithm for approximate typing of semi-structured data
 - ◆ each object may have more than one roles
- 16/15
Information