Modeling and Querying Possible Repairs in Duplicate Detection

Based on

"Modeling and Quering Possible Repairs in Duplicate Detection" PVLDB(2)1: 598-609 (2009) by

George Beskales, Mohamed A. Soliman

Ihab F. Ilyas, Shai Ben-David

Data Cleaning

Real-world data is dirty

□ Examples:

- Syntactical Errors: e.g., Micrsoft
- Heterogeneous Formats: e.g., Phone number formats
- Missing Values (Incomplete data)
- Violation of Integrity Constraints: e.g., FDs/INDs
- Duplicate Records

Data Cleaning is the process of improving data quality by removing errors, inconsistencies and anomalies of data

02/19/2010

Duplicate Elimination Process



Probabilistic Data Cleaning



Probabilistic Data Cleaning



One-shot Cleaning and Probabilistic Cleaning compared

Motivation

Probabilistic Data Cleaning:

- 1. Avoid deterministic resolution of conflicts during data cleaning
- 2. Enrich query results by considering all possible cleaning instances
- 3. Allows specifying query-time cleaning requirements

Outline

- Generating and Storing the Possible Clean Instances
- Querying the Clean Instances
- Experimental Evaluation



Uncertainty in Duplicate Detection

- Uncertain Duplicates: Determining what records should be clustered is uncertain due to noisy similarity measurements
- A possible repair of a relation is a clustering (partitioning) of the unclean relation

	F	Person			Possible Repairs			
ID	Name	ZIP	Income		X ₁	X ₂	X ₃	
P1	Green	51519	30k		{P1}	{P1,P2}	{P1,P2,P5}	
P2	Green	51518	32k		{ P2 }	{ P3,P4 }	{P3,P4}	
P3	Peter	30528	40k	Uncontain	{ P3,P4 }	{ P5 }	{ P6 }	
P4	Peter	30528	40k	Clustering	{ P5 }	{ P6 }		
P5	Gree	51519	55k		{ P6 }			
P6	Chuck	51519	30 k					

Challenges

1. The space of all possible clusterings (repairs) is exponentially large

2. How to efficiently and reasonably generate, store and query the possible repairs?

The Space of Possible Repairs



Cleaning Algorithm



 The model should store possible repairs in lossless way

Cleaning Algorithm

- Allow efficient answering of important queries (e.g. queries frequently encountered in applications)
- Provide materializations of the results of costly operations (e.g. clustering procedures) that are required by most queries
- Small space complexity to allow efficient construction, storage and retrieval of the possible repairs, in addition to efficient query processing

Algorithm – Dependent Model

A – Algorithm; R – starting unclean relation; P – set of possible parameters for A

au - continuous random variable for A from interval [τ^I, τ^u] f_τ – probability density function of τ (given by user or learned)

Applying A to R using parameter $t \in [\tau^{I}, \tau^{u}]$ generates possible clustering (i.e. repair) denoted as A(R, t)

 \mathcal{X} – set of all possible repairs, defined as {A(R,t) : t \in [T^I, T^U]}

02/19/2010

Algorithm – Dependent Model

Probability of a specific repair $X \in \mathcal{X}$:

$$\Pr(X) = \int_{\tau^l}^{\tau^u} f_{\tau}(t) \cdot h(t, X) dt$$

Where h(t, X) = 1 if A(R,t) = X, and 0 otherwise.

02/19/2010

Creating U-clean Relations

e [Р	erso	onc										
bas	ID		Income		С		Р	Vehicle ^c						
ata	CP1		31k	{P	1,P2}		[1,3]		ID		Price	С	Р	
	CP2	2	40k	{P	3,P4}	[0,10]	0	CV1		5k	{V1}	[0,4]	
lea	CP3		55k	{	P5}		[0,3]	0	CV2		7k	{V2}	[0,4]	
ပူ	CP4	ł	30k	{	P6}	[0,10]	C	CV3		6k	{V1,V2}	[4,10]	
tai	CP5	;	39k	{P1,	P2,P5}	[[3,10]	C	CV4		8k	{V3}	[0,5]	
cer	CP6	5	30k	{	P1}		[0,1]	C	CV5		4k	{V4}	[0,5]	
5	CP7		32k	{	P2}		[0,1]	0	CV6		6k	{V3,V4}	[5,10]	
Г		A	, (Perso	on, τ	;: U[0,1		\mathcal{A}_2 (Vehicle, τ_2 : U[0,10])							
e	Person								Vehicle					
as	ID	Nam	e ZIP) E	Birth Date		Income		ID	Make		Model	Price	
tab	P1	Gree	n 5135	9	781310)	30k		V1	H	Ionda	Civic	5k	
Da	P 2	Gree	n 5135	8	781210		32k		V 2	(Civic		7k	
Iclean	P 3	3 Peter 3012		28	870932		40k		V3	N	Vissan	Altima	8k	
	P 4	P4 Peter 30128		28	870932		40k		V4 I	Nissan Altima		na	4k	
5	P 5	Gree	e 5135	i9	1977121	10	55 k							
	P6	Chuc	k 5135	i9	1946092	24	30k							

Hierarchical Clustering Algorithms

- Records are clustered in a form of a hierarchy: all singletons are at leaves, and one cluster is at root
- **Example**: Linkage-based Clustering Algorithm



Uncertain Hierarchical Clustering

- Hierarchical clustering algorithms can be modified to accept uncertain parameters
- The number of generated possible repairs is linear



Uncertain Hierarchical Clustering

Algorithm 1 U_Cluster $(R(A_1, \ldots, A_m), \tau^l, \tau^u)$ **Require:** $R(A_1, \ldots, A_m)$: The unclean relation **Require:** τ^l : Minimum threshold value **Require:** τ^u : Maximum threshold value 1: Define a new singleton cluster C_i for each record $r_i \in R$ (i.e., C_i contains a unique identifier of r_i) 2: $\mathcal{C} \leftarrow \{C_1, \ldots, C_{|R|}\}$ 3: for each $r_i \in R$ do Add $(r_i[A_1], ..., r_i[A_m], C_i, [\tau^l, \tau^u])$ to R^c 4: 5: end for 6: while $(|\mathcal{C}| > 1)$ and distance between the closest pair of clusters (C_i, C_i) in C is less than τ^u) do 7: $C_k \leftarrow C_i \cup C_i$ 8: Replace C_i and C_j in \mathcal{C} with C_k 9: $r_k \leftarrow \text{get_representative_record}(C_k)$ {See Section 3.3} 10: Add $(r_k[A_1], \ldots, r_k[A_m], C_k, [dist(C_i, C_j), \tau^u])$ to R^c if $(dist(C_i, C_j) < \tau^l)$ then 11: 12: Remove the c-records corresponding to C_i and C_j from R^c 13: else 14: Set the upper bounds of parameter settings of the *c*-records corresponding to C_i and C_i to $dist(C_i, C_i)$ 15: end if 16: end while 17: return R^c

Probabilities of Repairs



Representing the Possible Repairs

				ID	•••	Income	С	Р
				CP1	•••	31k	{ P1,P2 }	[1,3)
Clustering 1	Clustering 2	Clustering 3	1	CP2	•••	40k	{P3,P4}	[0,10)
{P1}	{P1,P2}	{P1,P2,P5}		CP3		55k	{P5}	[0.3)
{P2}	{P3,P4}	{P3,P4}				301/2	(2 C) (D6)	
{P3,P4}	{P5}	{P6}			•••	JUK		
{P5}	{P6}		1	CP5	•••	39k	{P1,P2,P5}	[3,10)
(\mathbf{P}_{ℓ})	(10)			CP6	•••	30k	{P1}	[0,1)
{01}				CP7	••••	32k	{P2}	[0,1)

 $0 \le \tau < 1$ $1 \le \tau < 3$ $3 \le \tau < 10$ Pr = 0.1 Pr = 0.2 Pr = 0.7 **U-clean Relation** *Person*^C

Constructing Probabilistic Repairs



Re-creating Probabilistic repairs from U-clean Relations

NN-based clustering

- Tuples represented as points in d-dimensional space
- Columns are dimensions
- Dimensions ordering is very important choice
- Clusters are created by coalescing points that are "near" to each other

• With growing t points that are further and further away are being put together into mutual clusters, also various clusters can be united

 Algorithm stops when there is only one cluster and all points belong to it or maximum value of t is reached.

Time and Space Complexity

- Hierarchical clustering arranges records in N-ary tree, records being the leaves
- Maximum possible number of nodes of a tree that has n leaves (assuming each non-leaf node has at least 2 children) is 2n-1
- Therefore number of possible clusters is bounded by 2n-1, so it is linear w.r.t. to number of starting tuples.
- In general above algorithms have asymptotic complexity exactly as the original ones on which they build, since only constant amount of work is added to each iteration

Outline

- Generating and Storing the Possible Clean Instances
- Querying the Clean Instances
- Experimental Evaluation



Queries over U-Clean Relations

We adopt the *possible worlds semantics* to define queries on U-clean relations



Example: Selection Query

Person ^C								
ID	•••	Income	С	Р				
CP1	•••	31k	{P1,P2}	[1,3)				
CP2	•••	40k	{P3,P4}	[0,10)				
CP3	•••	55k	{P5}	[0,3)				
CP4	•••	30k	{P6}	[0,10)				
CP5	•••	39k	{P1,P2,P5}	[3,10)				
CP6	•••	30k	{P1}	[0,1)				
CP7	•••	32k	{P2}	[0,1)				



Example: Projection Query



Example: Join Query

			Person		Vehicle ^c					
ID	•••	Income	С	Р	ID	Price	С	Р		
CP1	•••	31k	{ P1,P2 }	τ ₁ :[1,3)	CV5	4k	{V4}	$\tau_2:[3,5)$		
CP2	•••	40 k	{ P3,P4 }	τ ₁ :[0,10)	CV6	6k	{V3,V4}	$\tau_2:[5,10)$		
•••	•••	•••	•••	•••	•••	••••	••••	••••		
	SELECT Income, Price FROM Person ^c , Vehicle ^c WHERE Income/10 >= Price									
		Incom	e Price	С		Р				
	40k 4k {P3,			{P3,P4} ^ {V4	-} τ ₁ :[0	$, 10) ^{7} \tau_{2}$:[3,5)			
		•••	•••	•••		•••				

Aggregation Queries

Person ^c									
ID	•••	Income	С	Р					
CP1	•••	31k	{P1,P2}	[1,3)					
CP2		40k	{P 3 ,P 4 }	[0,10)					
CP3		55k	{P5}	[0,3)					
CP4		30k	{P6}	[0,10)					
CP5		39k	{P1,P2,P5}	[3,10)					
CP6		30k	{P1}	[0,1)					
CP7		32k	{P2}	[0,1)					

D

SELECT Sum(Income) FROM Person^c



Other Meta-Queries

- 1. Obtaining the most probable clean instance
- 2. Obtaining the α -certain clusters
- Obtaining a clean instance corresponding to a specific parameters of the clustering algorithms
- Obtaining the probability of clustering a set of records together

Outline

- Generating and Storing the Possible Clean Instances
- Querying the Clean Instances
- Experimental Evaluation

- □ A prototype as an extension of PostgreSQL
- Synthetic data generator provided by Febrl (freely extensible biomedical record linkage)
- □ Two hierarchical algorithms:
 - Single-Linkage (S.L.)
 - Nearest-neighbor based clustering algorithm (N.N.) [Chaudhuri et al., ICDE'05]

- Machine used: SunFire X4100, Dual Core 2.2GHz, 8GB RAM
- □ 10% of records were duplicates
- Each query executed 5 times and average time taken







- Exucution time overhead for S.L. Is 30%, and for NN less than 5%, and it is negligibly correlated with percentage of duplicates
- □ Space overhead is 8.35%
- Extracting clean instance for 100,000 records requeries only 1.5sec, which means this approach is more efficient than restarting deduplication algorithm whenever a new parameter setting is requested

Conclusion

- We allow representing and querying multiple possible clean instances
- We modified hierarchical clustering algorithms to allow generating multiple repairs
- We compactly store the possible repairs by keeping the lineage information of clusters in special attributes
- New (probabilistic) query types can be issued against the population of possible repairs