

# Project-Join-Repair

An Approach to Consistent Query Answering  
Under Functional Dependencies

Jef Wijsen

Université de Mons-Hainaut, Belgium

# Preliminary I: Repairing

- Let  $\Sigma$  be a set of fd's and  $I$  an inconsistent relation. A **repair** of  $I$  is a maximal (under set inclusion) consistent subset of  $I$ .

$I$	<u>Name</u>	Birth	Sex	ZIP	City	
	An	1964	F	7000	Mons	Name $\rightarrow$ Birth
	Ed	1962	M	7000	Bergen	Name $\rightarrow$ Sex
						Name $\rightarrow$ ZIP
						ZIP $\rightarrow$ City

# Preliminary I: Repairing

- Let  $\Sigma$  be a set of fd's and  $I$  an inconsistent relation.  
A **repair** of  $I$  is a maximal (under set inclusion) consistent subset of  $I$ .

$I$	<table style="border-collapse: collapse; width: 100%;"> <tr> <th style="text-align: left; padding: 5px;"><u>Name</u></th> <th style="text-align: left; padding: 5px;">Birth</th> <th style="text-align: left; padding: 5px;">Sex</th> <th style="text-align: left; padding: 5px;">ZIP</th> <th style="text-align: left; padding: 5px;">City</th> </tr> <tr> <td style="padding: 5px;">An</td> <td style="padding: 5px;">1964</td> <td style="padding: 5px;">F</td> <td style="padding: 5px;">7000</td> <td style="padding: 5px;">Mons</td> </tr> <tr> <td style="padding: 5px;">Ed</td> <td style="padding: 5px;">1962</td> <td style="padding: 5px;">M</td> <td style="padding: 5px;">7000</td> <td style="padding: 5px;">Bergen</td> </tr> </table>	<u>Name</u>	Birth	Sex	ZIP	City	An	1964	F	7000	Mons	Ed	1962	M	7000	Bergen	<p><b>Name</b> <math>\rightarrow</math> <b>Birth</b></p> <p><b>Name</b> <math>\rightarrow</math> <b>Sex</b></p> <p><b>Name</b> <math>\rightarrow</math> <b>ZIP</b></p> <p><b>ZIP</b> <math>\rightarrow</math> <b>City</b></p>
<u>Name</u>	Birth	Sex	ZIP	City													
An	1964	F	7000	Mons													
Ed	1962	M	7000	Bergen													

{	$R_1$	<table style="border-collapse: collapse; width: 100%;"> <tr> <th style="text-align: left; padding: 5px;"><u>Name</u></th> <th style="text-align: left; padding: 5px;">Birth</th> <th style="text-align: left; padding: 5px;">Sex</th> <th style="text-align: left; padding: 5px;">ZIP</th> <th style="text-align: left; padding: 5px;">City</th> </tr> <tr> <td style="padding: 5px;">An</td> <td style="padding: 5px;">1964</td> <td style="padding: 5px;">F</td> <td style="padding: 5px;">7000</td> <td style="padding: 5px;">Mons</td> </tr> </table>	<u>Name</u>	Birth	Sex	ZIP	City	An	1964	F	7000	Mons	$\rightsquigarrow$ Ed is left out...
	<u>Name</u>	Birth	Sex	ZIP	City								
An	1964	F	7000	Mons									
$R_2$	<table style="border-collapse: collapse; width: 100%;"> <tr> <th style="text-align: left; padding: 5px;"><u>Name</u></th> <th style="text-align: left; padding: 5px;">Birth</th> <th style="text-align: left; padding: 5px;">Sex</th> <th style="text-align: left; padding: 5px;">ZIP</th> <th style="text-align: left; padding: 5px;">City</th> </tr> <tr> <td style="padding: 5px;">Ed</td> <td style="padding: 5px;">1962</td> <td style="padding: 5px;">M</td> <td style="padding: 5px;">7000</td> <td style="padding: 5px;">Bergen</td> </tr> </table>	<u>Name</u>	Birth	Sex	ZIP	City	Ed	1962	M	7000	Bergen	$\rightsquigarrow$ An is left out...	
<u>Name</u>	Birth	Sex	ZIP	City									
Ed	1962	M	7000	Bergen									

# Prelim. II: Consistent Query Answering

- Let  $\{R_1, R_2, \dots, R_n\}$  be the set of all repairs of  $I$  w.r.t.  $\Sigma$ . The **consistent answer** to a query  $q$  is defined by:

$$q_{\Sigma}(I) := \bigcap_{i=1}^n q(R_i)$$

- For example, "SELECT Name FROM I" yields {}.

# Prelim. II: Consistent Query Answering

- Let  $\{R_1, R_2, \dots, R_n\}$  be the set of all repairs of  $I$  w.r.t.  $\Sigma$ . The **consistent answer** to a query  $q$  is defined by:

$$q_{\Sigma}(I) := \bigcap_{i=1}^n q(R_i)$$

- For example, "SELECT Name FROM I" yields  $\{\}$ .
- This semantics is old (at least 1984):

$q_{\Sigma}(I)$  is the **certain answer** to  $q$  on the **incomplete relation**  $\{R_1, R_2, \dots, R_n\}$ .

# Prelim. II: Consistent Query Answering

- Let  $\{R_1, R_2, \dots, R_n\}$  be the set of all repairs of  $I$  w.r.t.  $\Sigma$ . The **consistent answer** to a query  $q$  is defined by:

$$q_{\Sigma}(I) := \bigcap_{i=1}^n q(R_i)$$

- For example, "SELECT Name FROM I" yields  $\{\}$ .
- This semantics is old (at least 1984):

$q_{\Sigma}(I)$  is the **certain answer** to  $q$  on the **incomplete relation**  $\{R_1, R_2, \dots, R_n\}$ .

- Consistent query answering** is the complexity of the set:

$$\text{CQA}(\Sigma, q) := \{I \mid I \text{ is a relation and } q_{\Sigma}(I) = \{\}\}$$

# Project-Join-Repair (Example)

Before repairing, apply the project-join dependency

$$\sigma = \bowtie [\{\text{Name, Birth, Sex, ZIP}\}, \{\text{ZIP, City}\}] :$$

<i>I</i>	<u>Name</u>	Birth	Sex	ZIP	City
	An	1964	F	7000	Mons
	Ed	1962	M	7000	Bergen

<u>Name</u>	Birth	Sex	ZIP	ZIP	City
An	1964	F	7000	7000	Mons
Ed	1962	M	7000	7000	Bergen

$\sigma(I)$	<u>Name</u>	Birth	Sex	ZIP	City
	An	1964	F	7000	Mons
	Ed	1962	M	7000	Bergen
	An	1964	F	7000	Bergen
	Ed	1962	M	7000	Mons

# Project-Join-Repair (Ex. Contd.)

$\sigma(I)$	<u>Name</u>	Birth	Sex	ZIP	City
	An	1964	F	7000	Mons
	Ed	1962	M	7000	Bergen
	An	1964	F	7000	Bergen
	Ed	1962	M	7000	Mons

$J_1$	<u>Name</u>	Birth	Sex	ZIP	City
	An	1964	F	7000	Mons
	Ed	1962	M	7000	Mons

has two repairs:

$J_2$	<u>Name</u>	Birth	Sex	ZIP	City
	An	1964	F	7000	Bergen
	Ed	1962	M	7000	Bergen

The effect: repairing by value modification!



# Project-Join-Repair

- $\Sigma$  = set of fd's                       $\sigma$  = join dependency (jd)  
 $I$  = possibly inconsistent relation
- We repair  $\sigma(I)$ —rather than  $I$ . Note that  $\sigma(I) \supseteq I$ .
- We require  $\Sigma \models \sigma$ .

Therefore, if  $I$  satisfies  $\Sigma$ , then  $\sigma(I) = I$ .

# Project-Join-Repair

- $\Sigma$  = set of fd's                       $\sigma$  = join dependency (jd)  
 $I$  = possibly inconsistent relation
- We repair  $\sigma(I)$ —rather than  $I$ . Note that  $\sigma(I) \supseteq I$ .
- We require  $\Sigma \models \sigma$ .

Therefore, if  $I$  satisfies  $\Sigma$ , then  $\sigma(I) = I$ .

- Consistent query answering becomes the complexity of:

$$\text{CQAJD}(\Sigma, \sigma, q) := \{I \mid I \text{ is a relation and } q_{\Sigma}(\sigma(I)) = \{\}\}$$

Determine tractable cases of this problem!

- If  $\sigma$  is the identity jd, then  $\text{CQAJD}(\Sigma, \sigma, q) = \text{CQA}(\Sigma, q)$ .

# Tractable Cases

- Let fd's in  $\Sigma$  be of the form  $X \rightarrow A$  with  $A \notin X$ .
- Main theorem:

CQAJD( $\Sigma, \sigma, q$ ) is in **P** under the following condition:

- $\Sigma$  is acyclic (and hence has unique key  $K$ );
- no two fd's of  $\Sigma$  have the same right-hand side;
- for each  $X \rightarrow A$  of  $\Sigma$ ,  $K \subseteq X$  or  $\sigma \models X \twoheadrightarrow A$ ; and
- $q$  is a rooted, typed conjunctive query.

# Tractable Cases

- Let fd's in  $\Sigma$  be of the form  $X \rightarrow A$  with  $A \notin X$ .
- Main theorem:

CQAJD( $\Sigma, \sigma, q$ ) is in **P** under the following condition:

- $\Sigma$  is acyclic (and hence has unique key  $K$ );
- no two fd's of  $\Sigma$  have the same right-hand side;
- for each  $X \rightarrow A$  of  $\Sigma$ ,  $K \subseteq X$  or  $\sigma \models X \twoheadrightarrow A$ ; and
- $q$  is a rooted, typed conjunctive query.

- There are cases of CQA( $\Sigma, q$ ) that are **NP**-complete under these conditions!
- To conclude, repairing  $\sigma(I)$ —rather than  $I$ —can be more natural and more efficient at the same time.

# Preservation of Key Values

- Acyclic  $\Leftrightarrow$  attributes can be ordered such that if  $XA \rightarrow B$  is in a minimal cover of  $\Sigma$ , then  $A < B$ .
- Every repair preserves all key values.

$I \not\models X \rightarrow A$	<table style="border-collapse: collapse;"> <thead> <tr> <th style="border-bottom: 1px solid black;"><u>K</u></th> <th style="border-bottom: 1px solid black;">X</th> <th style="border-bottom: 1px solid black;">A</th> <th style="border-bottom: 1px solid black;">Y</th> </tr> </thead> <tbody> <tr> <td>k</td> <td>x</td> <td>a</td> <td>y</td> </tr> <tr> <td>l</td> <td>x</td> <td>b</td> <td>y'</td> </tr> </tbody> </table>	<u>K</u>	X	A	Y	k	x	a	y	l	x	b	y'	$\sigma(I) \models X \twoheadrightarrow A$	<table style="border-collapse: collapse;"> <thead> <tr> <th style="border-bottom: 1px solid black;"><u>K</u></th> <th style="border-bottom: 1px solid black;">X</th> <th style="border-bottom: 1px solid black;">A</th> <th style="border-bottom: 1px solid black;">Y</th> </tr> </thead> <tbody> <tr> <td>k</td> <td>x</td> <td>a</td> <td>y</td> </tr> <tr> <td>l</td> <td>x</td> <td>b</td> <td>y'</td> </tr> <tr> <td style="color: red;">l</td> <td style="color: red;">x</td> <td style="color: red;">a</td> <td style="color: red;">y'</td> </tr> <tr> <td colspan="4" style="text-align: center;">...</td> </tr> </tbody> </table>	<u>K</u>	X	A	Y	k	x	a	y	l	x	b	y'	l	x	a	y'	...			
<u>K</u>	X	A	Y																																
k	x	a	y																																
l	x	b	y'																																
<u>K</u>	X	A	Y																																
k	x	a	y																																
l	x	b	y'																																
l	x	a	y'																																
...																																			

$\Downarrow$

$R_1$	<table style="border-collapse: collapse;"> <thead> <tr> <th style="border-bottom: 1px solid black;"><u>K</u></th> <th style="border-bottom: 1px solid black;">X</th> <th style="border-bottom: 1px solid black;">A</th> <th style="border-bottom: 1px solid black;">Y</th> </tr> </thead> <tbody> <tr> <td>k</td> <td>x</td> <td>a</td> <td>y</td> </tr> </tbody> </table>	<u>K</u>	X	A	Y	k	x	a	y	$J_1$	<table style="border-collapse: collapse;"> <thead> <tr> <th style="border-bottom: 1px solid black;"><u>K</u></th> <th style="border-bottom: 1px solid black;">X</th> <th style="border-bottom: 1px solid black;">A</th> <th style="border-bottom: 1px solid black;">Y</th> </tr> </thead> <tbody> <tr> <td>k</td> <td>x</td> <td>a</td> <td>y</td> </tr> <tr> <td style="color: red;">l</td> <td style="color: red;">x</td> <td style="color: red;">a</td> <td style="color: red;">y'</td> </tr> </tbody> </table>	<u>K</u>	X	A	Y	k	x	a	y	l	x	a	y'
<u>K</u>	X	A	Y																				
k	x	a	y																				
<u>K</u>	X	A	Y																				
k	x	a	y																				
l	x	a	y'																				

# Rooted, Typed Conjunctive Queries

Key positions are underlined in:

$$Answer(\vec{z}) \leftarrow R(\underline{\vec{x}}_1, \vec{y}_1), \dots, R(\underline{\vec{x}}_n, \vec{y}_n)$$

**Typed** No variable occurs at distinct positions in  $R$ -atoms.

**Rooted** If  $R(\underline{\vec{x}}_i, \vec{y}_i)$  and  $R(\underline{\vec{x}}_j, \vec{y}_j)$  with  $i \neq j$  share a common variable  $w$ , then  $w$  occurs in  $\vec{z}$  (i.e.  $w$  is free)  
or all variables of  $\underline{\vec{x}}_i, \underline{\vec{x}}_j$  occur in  $\vec{z}$ .

# Rooted, Typed Conjunctive Queries

Key positions are underlined in:

$$\text{Answer}(\vec{z}) \leftarrow R(\underline{\vec{x}}_1, \vec{y}_1), \dots, R(\underline{\vec{x}}_n, \vec{y}_n)$$

**Typed** No variable occurs at distinct positions in  $R$ -atoms.

**Rooted** If  $R(\underline{\vec{x}}_i, \vec{y}_i)$  and  $R(\underline{\vec{x}}_j, \vec{y}_j)$  with  $i \neq j$  share a common variable  $w$ , then  $w$  occurs in  $\vec{z}$  (i.e.  $w$  is free)  
or all variables of  $\vec{x}_i, \vec{x}_j$  occur in  $\vec{z}$ .

- Get names of male persons living in An's city.

$$\langle x_n \rangle \leftarrow \langle \underline{x}_n, x_b, \mathbf{M}, x_z, x_c \rangle, \langle \underline{\mathbf{An}}, y_b, y_s, y_z, x_c \rangle$$

- Man–Woman pairs with same birth and city.

$$\langle x_n, y_n \rangle \leftarrow \langle \underline{x}_n, x_b, \mathbf{M}, x_z, x_c \rangle, \langle \underline{y}_n, x_b, \mathbf{F}, y_z, x_c \rangle$$

# Computing the Consistent Answers

$\sigma(I)$	<u>Name</u>	Birth	Sex	ZIP	City
	An	1964	F	7000	Mons
	Ed	1962	M	7000	Bergen
	An	1964	F	7000	Bergen
	Ed	1962	M	7000	Mons



$N$	<u>Name</u>	Birth	Sex	ZIP	City
	An	1964	F	7000	$w$
	Ed	1962	M	7000	$w$

Under the conditions of the main theorem:

For any rooted, typed conjunctive query  $q$ ,  
 $q_{\Sigma}(\sigma(I)) = \text{the ground tuples in } q(N)$ .



# Conclusion

- Database repairing and consistent query answering under fd's.
- Applying a join dependency prior to repairing by tuple deletion:
  - is semantically meaningful;  
the effect is a sort of repairing by value modification.
  - may give you tractability.