

# 4-valued Semantics for Data Integration under the OWA: a Deductive Database Approach

Dominique Laurent

ETIS – CNRS-University Cergy Pontoise  
FRANCE

# Motivation

- Integrating data from data sources providing **explicitly true** information and **false** information
- This leads to potential **contradictions**
- Our goal is **not to eliminate contradictions**, but rather to conduct reasoning based on them, considering that:
  - Data sources provide reliable information
  - The result is a database, with 'facts' and 'rules'
  - The Open World Assumption is made

# Running Example

- Consider the atoms  $a$ ,  $b$ ,  $c$ ,  $d$  and  $e$
- Sources  $S_1$  and  $S_2$  are such that:
  - In  $S_1$  :  $a$  and  $b$  are true, and  $c$  is false
  - In  $S_2$  :  $a$  and  $d$  are true, and  $b$  and  $c$  are false
- In the integrated database
  - $a$  is **true** and  $c$  is **false** (the sources agree)
  - $b$  is **contradictory** (the sources disagree)
  - $d$  is **true** (the sources do not disagree)
  - $e$  is **unknown** (no information in the sources)

# Running Example (cont'd)

- Consider the rules

$$(1) c \leftarrow a, \neg b$$

$$(2) e \leftarrow d$$

$$(3) c \leftarrow e$$

$$(4) \neg e \leftarrow a, d$$

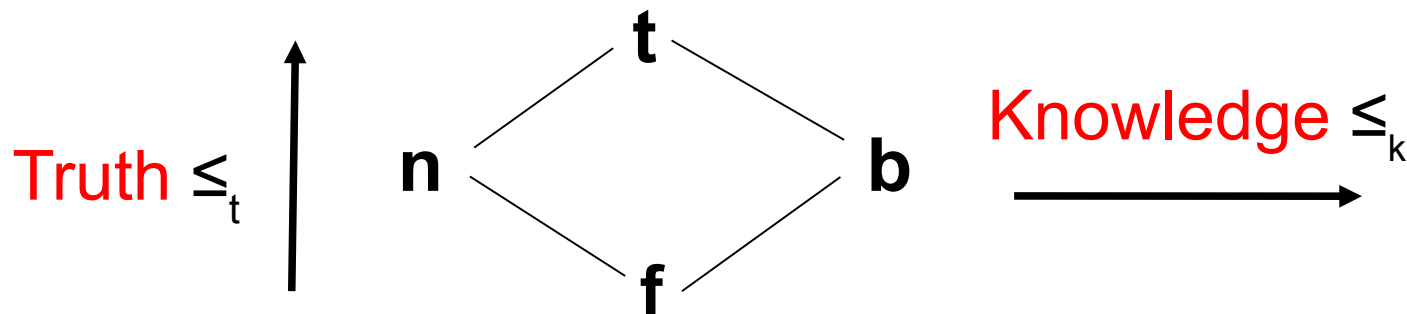
- **Rules may generate inconsistencies**
- Rules are applied based on OWA: if the body **explicitly** holds, so does the head
  - Rule (2) applies because  $d$  is known to be true
  - Then (3) applies making  $c$  true...
  - What about (1)?

# Background

- Belnap's four valued logic
  - True **t**
  - False **f**
  - **Inconsistent** **b**
  - **Unknown** **n**
- **Designated values: t and b**
  - A formula **holds** if its truth value is **t** or **b**

# Background

- Belnap's four valued logic
  - True **t**
  - False **f**
  - Inconsistent **b**
  - Unknown **n**
- Designated values: **t** and **b**
  - A formula holds if its truth value is **t** or **b**
- Two orderings



# Background

- $(\{\mathbf{t}, \mathbf{b}, \mathbf{n}, \mathbf{f}\}, \leq_t, \leq_k)$  is a distributive bi-lattice
- Basic Connectors
  - **Disjunction**: lub wrt  $\leq_t$        $\mathbf{b} \vee \mathbf{n}$     gives     $\mathbf{t}$
  - **Conjunction**: glb wrt  $\leq_t$        $\mathbf{b} \wedge \mathbf{n}$     gives     $\mathbf{f}$
  - **Negation** wrt  $\leq_t$  :               $\neg \mathbf{t} = \mathbf{f}, \neg \mathbf{f} = \mathbf{t}, \neg \mathbf{b} = \mathbf{b}, \neg \mathbf{n} = \mathbf{n}$
  - **Weak-negation** wrt  $\leq_k$  :         $-\mathbf{t} = \mathbf{b}, -\mathbf{f} = \mathbf{n}, -\mathbf{b} = \mathbf{t}, -\mathbf{n} = \mathbf{f}$
- Additional connector  $(\neg \_ \neg \_)$ 
  - **Complementation**  $\sim$                $\sim \mathbf{t} = \mathbf{f}, \sim \mathbf{f} = \mathbf{t}, \sim \mathbf{b} = \mathbf{u}, \sim \mathbf{u} = \mathbf{b}$

# Implication

- Defining  $P \supset Q$  by  $\neg P \vee Q$  **does not work** because  $(P \wedge (P \supset Q)) \supset Q$  is **not always valid**
  - It is **f** when  $P$  is **b** and  $Q$  is **f**

$P \supset Q$	t	b	n	f
t	t	b	n	f
b	t	b	t	b
n	t	t	n	n
f	t	t	t	t

Should **not** be designated

Should be designated



# Implication

- Defining  $P \supset Q$  by  $\neg P \vee Q$  does not work because  $(P \wedge (P \supset Q)) \supset Q$  is not always valid
- This is not the case for  
–  $P \Rightarrow Q$  defined by  $\sim P \vee Q$

$P \Rightarrow Q$	t	b	n	f
t	t	b	n	f
b	t	t	n	n
n	t	b	t	b
f	t	t	t	t

# Implication

- Defining  $P \supset Q$  by  $\neg P \vee Q$  does not work because  $(P \wedge (P \supset Q)) \supset Q$  is not always valid
- This is not the case for
  - $P \Rightarrow Q$  defined by  $\sim P \vee Q$
  - $P \rightarrow Q$  defined below

$P \Rightarrow Q$	t	b	n	f
t	t	b	n	f
b	t	t	n	n
n	t	b	t	b
f	t	t	t	t

$P \rightarrow Q$	t	b	n	f
t	t	b	n	f
b	t	b	n	f
n	t	t	t	t
f	t	t	t	t

$P \Rightarrow Q$  or  $P \rightarrow Q$  ???

- $v(P \Rightarrow Q)$  and  $v(P \rightarrow Q)$  are not designated iff  $v(P)$  is designated and  $v(Q)$  is not

$P \Rightarrow Q$  or  $P \rightarrow Q$  ???

- $v(P \Rightarrow Q)$  and  $v(P \rightarrow Q)$  are not designated iff  $v(P)$  is designated and  $v(Q)$  is not
- $v(P) \leq_t v(Q)$  iff  $v(P \Rightarrow Q) = \mathbf{t}$

$P \Rightarrow Q$  or  $P \rightarrow Q$  ???

- $v(P \Rightarrow Q)$  and  $v(P \rightarrow Q)$  are not designated iff  $v(P)$  is designated and  $v(Q)$  is not
- $v(P) \leq_t v(Q)$  iff  $v(P \Rightarrow Q) = \mathbf{t}$
- $v(P \rightarrow Q) = v(Q)$  iff  $v(P)$  is designated
- None does satisfy the above two properties

# $P \Rightarrow Q$ or $P \rightarrow Q$ ???

- $v(P \Rightarrow Q)$  and  $v(P \rightarrow Q)$  are not designated iff  $v(P)$  is designated and  $v(Q)$  is not
- $v(P) \leq_t v(Q)$  iff  $v(P \Rightarrow Q) = \mathbf{t}$
- $v(P \rightarrow Q) = v(Q)$  iff  $v(P)$  is designated
- None does satisfy the above two properties,  
**neither do they satisfy contraposition**
  - $P \Rightarrow Q$  is *not* equivalent to  $\neg Q \Rightarrow \neg P$
  - $P \rightarrow Q$  is *not* equivalent to  $\neg Q \rightarrow \neg P$

# $P \Rightarrow Q$ or $P \rightarrow Q$ ???

- $v(P \Rightarrow Q)$  and  $v(P \rightarrow Q)$  are not designated iff  $v(P)$  is designated and  $v(Q)$  is not
- $v(P) \leq_t v(Q)$  iff  $v(P \Rightarrow Q) = \mathbf{t}$
- $v(P \rightarrow Q) = v(Q)$  iff  $v(P)$  is designated
- None does satisfy the above two properties, neither do they satisfy contraposition
  - $P \Rightarrow Q$  is *not* equivalent to  $\neg Q \Rightarrow \neg P$
  - $P \rightarrow Q$  is *not* equivalent to  $\neg Q \rightarrow \neg P$
- In fact, **the choice is not so important...**

# Our Approach

- **Valued pair** (or v-pair):  $(F, \mathbf{v})$  such that
  - $F$  is a **fact** in the underlying Herbrand Universe
  - $\mathbf{v}$  is a **truth value** in  $\{\mathbf{t}, \mathbf{b}, \mathbf{n}, \mathbf{f}\}$
- A set  $E$  of v-pairs is a **v-set**  
 $E$  is **consistent** if  $E$  contains *at most* one v-pair involving a fact  $F$
- Every consistent v-set defines a **valuation**  $v_E$  :
  - If  $(F, \mathbf{v})$  occurs in  $E$  then  $v_E(F) = \mathbf{v}$
  - Otherwise,  $v_E(F) = \mathbf{n}$



# Our Approach

- A **database**  $D$  is a pair  $(E, R)$  such that
  - $E$  is a v-set
  - $R$  is a set of rules whose bodies and/or **heads** can contain **negated** atoms
- $D$  is consistent if so is  $E$
- If  $D = (E, R)$  is consistent then for every fact  $F$ 
  - If  $F$  occurs in  $E$  then  $F$  has the truth value in the corresponding v-pair
  - Otherwise,  $F$  has the truth value **n**

# Back to the Running Example

- $D = (E, R)$  where
  - $E = \{(a, \mathbf{t}), (b, \mathbf{b}), (c, \mathbf{f}), (d, \mathbf{t})\}$
  - $R$  is the following set
    - (1)  $c \leftarrow a, \neg b$
    - (2)  $e \leftarrow d$
    - (3)  $c \leftarrow e$
    - (4)  $\neg e \leftarrow a, d$
- $v_E$  is defined by
  - $v_E(a) = \mathbf{t}$  ;  $v_E(b) = \mathbf{b}$  ;  $v_E(c) = \mathbf{f}$  ;  $v_E(d) = \mathbf{t}$  ;  $v_E(e) = \mathbf{n}$
- (1), (2) and (3) **do not hold** based on  $v_E$ 
  - $v_E(a, \neg b) = \mathbf{b}$ ,  $v_E(c) = \mathbf{f}$       thus  $v_E(\text{rule}(1)) = \mathbf{n}$  or  $\mathbf{f}$

# Back to the Running Example

- $D = (E, R)$  where

$P \Rightarrow Q$	t	b	n	<b>f</b>
t	t	b	n	f
<b>b</b>	t	t	n	<b>n</b>
n	t	b	t	b
f	t	t	t	t

$P \rightarrow Q$	t	b	n	<b>f</b>
t	t	b	n	f
<b>b</b>	t	<b>b</b>	n	<b>f</b>
n	t	<b>t</b>	t	<b>t</b>
f	t	t	t	t

= n

- (1), (2), (3), (4), (5), (6), (7), (8), (9), (10), (11), (12), (13), (14), (15), (16), (17), (18), (19), (20), (21), (22), (23), (24), (25), (26), (27), (28), (29), (30), (31), (32), (33), (34), (35), (36), (37), (38), (39), (40), (41), (42), (43), (44), (45), (46), (47), (48), (49), (50), (51), (52), (53), (54), (55), (56), (57), (58), (59), (60), (61), (62), (63), (64), (65), (66), (67), (68), (69), (70), (71), (72), (73), (74), (75), (76), (77), (78), (79), (80), (81), (82), (83), (84), (85), (86), (87), (88), (89), (90), (91), (92), (93), (94), (95), (96), (97), (98), (99), (100)
- $v_E(a, \neg b) = \mathbf{b}$ ,  $v_E(c) = \mathbf{f}$       thus  $v_E(\text{rule}(1)) = \mathbf{n}$  or  $\mathbf{f}$

# Back to the Running Example

- $D = (E, R)$  where
  - $E = \{(a, \mathbf{t}), (b, \mathbf{b}), (c, \mathbf{f}), (d, \mathbf{t})\}$
  - $R$  is the set of the following rules
    - (1)  $c \leftarrow a, \neg b$
    - (2)  $e \leftarrow d$
    - (3)  $c \leftarrow e$
    - (4)  $\neg e \leftarrow a, d$
- $v_E$  is defined by
  - $v_E(a) = \mathbf{t}$  ;  $v_E(b) = \mathbf{b}$  ;  $v_E(c) = \mathbf{f}$  ;  $v_E(d) = \mathbf{t}$  ;  $v_E(e) = \mathbf{n}$
- (1), (2) and (4) **do not hold** based on  $v_E$ 
  - $v_E(d) = \mathbf{t}$ ,  $v_E(e) = \mathbf{n}$       thus  $v_E(\text{rule}(2)) = \mathbf{n}$

# Back to the Running Example

- $D = (E, P)$  where

$P \Rightarrow Q$	t	b	<b>n</b>	f
<b>t</b>	t	b	<b>n</b>	f
b	t	t	n	n
n	t	b	t	b
f	t	t	t	t

$P \rightarrow Q$	t	b	<b>n</b>	f
<b>t</b>	t	b	<b>n</b>	f
b	t	<b>b</b>	n	<b>f</b>
n	t	<b>t</b>	t	<b>t</b>
f	t	t	t	t

- (1), (2) and (4) do not hold based on

–  $v_E(d) = \mathbf{t}$ ,  $v_E(e) = \mathbf{n}$       thus  $v_E(\text{rule}(2)) = \mathbf{n}$

# Back to the Running Example

- $D = (E, R)$  where
  - $E = \{(a, \mathbf{t}), (b, \mathbf{b}), (c, \mathbf{f}), (d, \mathbf{t})\}$
  - $R$  is the set of the following rules
    - (1)  $c \leftarrow a, \neg b$
    - (2)  $e \leftarrow d$
    - (3)  $c \leftarrow e$
    - (4)  $\neg e \leftarrow a, d$
- $v_E$  is defined by
  - $v_E(a) = \mathbf{t}$  ;  $v_E(b) = \mathbf{b}$  ;  $v_E(c) = \mathbf{f}$  ;  $v_E(d) = \mathbf{t}$  ;  $v_E(e) = \mathbf{n}$
- (1), (2) and (4) **do not hold** based on  $v_E$ 
  - $v_E(a, d) = \mathbf{t}$ ,  $v_E(\neg e) = \mathbf{n}$       thus  $v_E(\text{rule}(4)) = \mathbf{n}$

# Database Semantics

- A **model  $M$**  of  $D = (E, R)$  is a v-set such that
  - The content of  $E$  is preserved, i.e.  $E \subseteq M$
  - The rules of  $R$  are valid
- Not possible unless rules have **exceptions**
  - **An instantiated rule whose head involves a fact in  $E$  does **not** apply**

# Database Semantics

- A model  $M$  of  $D = (E, R)$  is a v-set such that
  - The content of  $E$  is preserved, i.e.  $E \subseteq M$
  - The rules of  $R$  are valid
- Not possible unless rules have **exceptions**
- For  $E = \{(a, \mathbf{t}), (b, \mathbf{b}), (c, \mathbf{f}), (d, \mathbf{t})\}$  and  $R$ 
  - (1)  $c \leftarrow a, \neg b$
  - (2)  $e \leftarrow d$
  - (3)  $c \leftarrow e$
  - (4)  $\neg e \leftarrow a, d$ $M = \{(a, \mathbf{t}), (b, \mathbf{b}), (c, \mathbf{f}), (d, \mathbf{t}), (e, \mathbf{b})\}$  is a model
  - **(1) and (3) do not apply**
  - (2) and (4) are valid



# Database Semantics

$P \Rightarrow Q$	t	<b>b</b>	n	f
<b>t</b>	t	<b>b</b>	n	f
b	t	t	n	n
n	t	b	t	b
f	t	t	t	t

$P \rightarrow Q$	t	<b>b</b>	n	f
<b>t</b>	t	<b>b</b>	n	f
b	t	<b>b</b>	n	<b>f</b>
n	t	<b>t</b>	t	<b>t</b>
f	t	t	t	t

(2)  $e \leftarrow d$

(4)  $\neg e \leftarrow a, d$

$M = \{(a, \mathbf{t}), (b, \mathbf{b}), (c, \mathbf{f}), (d, \mathbf{t}), (\mathbf{e}, \mathbf{b})\}$  is a model

– (1) and (3) do not apply

– (2) and (4) are valid

# Computing the Database Semantics

- Given  $D = (E, R)$ , define  $inst^E(R)$  as the set of all **instantiated rules**  $r$  such that:
  - $r$  is an instantiation of a rule in  $R$
  - **The head of  $r$  does not occur in  $E$**
- Given a v-set  $S$  define:  $\Gamma(S) = E \cup$ 
  - $\{(h, \mathbf{t}) \mid (\exists r \text{ in } inst^E(R))(h = head(r), v_S(body(r)) = \mathbf{t})\}$
  - $\{(h, \mathbf{b}) \mid (\exists r \text{ in } inst^E(R))(h = head(r), v_S(body(r)) = \mathbf{b})\}$
  - $\{(h, \mathbf{f}) \mid (\exists r \text{ in } inst^E(R))(h = \neg head(r), v_S(body(r)) = \mathbf{t})\}$
  - $\{(h, \mathbf{b}) \mid (\exists r \text{ in } inst^E(R))(h = \neg head(r), v_S(body(r)) = \mathbf{b})\}$

# Example

- $D = (E, R)$  where  $E = \emptyset$  and  $R$ 
  - (1)  $c \leftarrow a, \neg b$
  - (2)  $e \leftarrow d$
  - (3)  $c \leftarrow e$
  - (4)  $\neg e \leftarrow a, d$
- $S = \{(a, \mathbf{t}), (b, \mathbf{b}), (d, \mathbf{t})\}$
- $\Gamma(S) = \{(c, \mathbf{b}), (e, \mathbf{t}), (e, \mathbf{f})\}$

because  $v_s(\text{body}(1)) = \mathbf{b}$

because  $v_s(\text{body}(4)) = \mathbf{t}$

because  $v_s(\text{body}(2)) = \mathbf{t}$

# Computing the Database Semantics

- Given  $D = (E, R)$  and  $S$ , define  $\Sigma(S)$  as
$$\Sigma(S) = \{(F, \mathbf{v}_k) \mid F \text{ occurs in } \Gamma(S) \text{ and}$$
$$\mathbf{v}_k = \text{lub}_k\{\mathbf{v} \mid (F, \mathbf{v}) \text{ in } \Gamma(S)\}\}$$
- The following sequence is monotonic wrt  $\leq_k$ 
  - $\Sigma_0 = E$
  - For  $p > 0$ ,  $\Sigma_{p+1} = \Sigma(\Sigma_p)$
- The **semantics** of  $D = (E, R)$  is the limit of the sequence, denoted  $\Sigma^*$

# Example

- $D = (E, R)$  ;  $E = \{(a, \mathbf{t}), (b, \mathbf{b}), (d, \mathbf{t})\}$  and  $R$ 
  - (1)  $c \leftarrow a, \neg b$
  - (2)  $e \leftarrow d$
  - (3)  $c \leftarrow e$
  - (4)  $\neg e \leftarrow a, d$
- $inst^E(R) = R$  and  $\Sigma_0 = \{(a, \mathbf{t}), (b, \mathbf{b}), (d, \mathbf{t})\}$

# Example

- $D = (E, R)$  ;  $E = \{(a, \mathbf{t}), (b, \mathbf{b}), (d, \mathbf{t})\}$  and  $R$ 
  - (1)  $c \leftarrow a, \neg b$
  - (2)  $e \leftarrow d$
  - (3)  $c \leftarrow e$
  - (4)  $\neg e \leftarrow a, d$
- $inst^E(R) = R$  and  $\Sigma_0 = \{(a, \mathbf{t}), (b, \mathbf{b}), (d, \mathbf{t})\}$
- **Computation of  $\Sigma_1$** 
  - $\Gamma(\Sigma_0) = \Sigma_0 \cup \{(c, \mathbf{b}), (e, \mathbf{t}), (e, \mathbf{f})\}$
  - $\Sigma_1 = \{(a, \mathbf{t}), (b, \mathbf{b}), (d, \mathbf{t}), (c, \mathbf{b}), (e, \mathbf{b})\}$

# Example

- $D = (E, R)$  ;  $E = \{(a, \mathbf{t}), (b, \mathbf{b}), (d, \mathbf{t})\}$  and  $R$ 
  - (1)  $c \leftarrow a, \neg b$
  - (2)  $e \leftarrow d$
  - (3)  $c \leftarrow e$
  - (4)  $\neg e \leftarrow a, d$
- $inst^E(R) = R$  and  $\Sigma_0 = \{(a, \mathbf{t}), (b, \mathbf{b}), (d, \mathbf{t})\}$
- Computation of  $\Sigma_1$ 
  - $\Gamma(\Sigma_0) = \Sigma_0 \cup \{(c, \mathbf{b}), (e, \mathbf{t}), (e, \mathbf{f})\}$
  - $\Sigma_1 = \{(a, \mathbf{t}), (b, \mathbf{b}), (d, \mathbf{t}), (c, \mathbf{b}), (e, \mathbf{b})\}$
- $\Sigma_2 = \Sigma(\Sigma_1) = \Sigma_1$
- Thus  $\Sigma_1$  is the semantics of  $D$ , i.e.  $\Sigma^* = \Sigma_1$

# Back to the Running Example

- $D = (E, R)$  ;  $E = \{(a, \mathbf{t}), (b, \mathbf{b}), (c, \mathbf{f}), (d, \mathbf{t})\}$  and  $R$ 
  - (1)  $c \leftarrow a, \neg b$
  - (2)  $e \leftarrow d$
  - (3)  $c \leftarrow e$
  - (4)  $\neg e \leftarrow a, d$
- $inst^E(R) = \{(2), (4)\}$  and  $\Sigma_0 = E$
- Computation of  $\Sigma_1$ 
  - $\Gamma(\Sigma_0) = \Sigma_0 \cup \{(e, \mathbf{t}), (e, \mathbf{f})\}$
  - $\Sigma_1 = \{(a, \mathbf{t}), (b, \mathbf{b}), (c, \mathbf{f}), (d, \mathbf{t}), (e, \mathbf{b})\}$
- $\Sigma_2 = \Sigma(\Sigma_1) = \Sigma_1$
- Thus  $\Sigma^* = \Sigma_1$



# Properties

- For every  $D = (E, R)$ ,  $\Sigma^*$  is a model of  $D$
- This model is minimal wrt set-inclusion

# Properties

- For every  $D = (E, R)$ ,  $\Sigma^*$  is a model of  $D$
- This model is minimal wrt set-inclusion
- However
  - There are **other minimal models** wrt set-inclusion
  - This model is **not minimal wrt truth ordering**
  - This model is **not minimal wrt knowledge ordering**
- **Conjecture:**
  - All minimal models wrt set-inclusion contain the **same** valid facts

# Future Work

- **Ongoing investigations:** characterizing the computed model, updates
- **Further research:**
  - Extend the relational algebra to this model
  - Study traditional constraints such as FDs, or TGDs
- **Other future Issues**
  - Implementation
  - Apply in real world applications