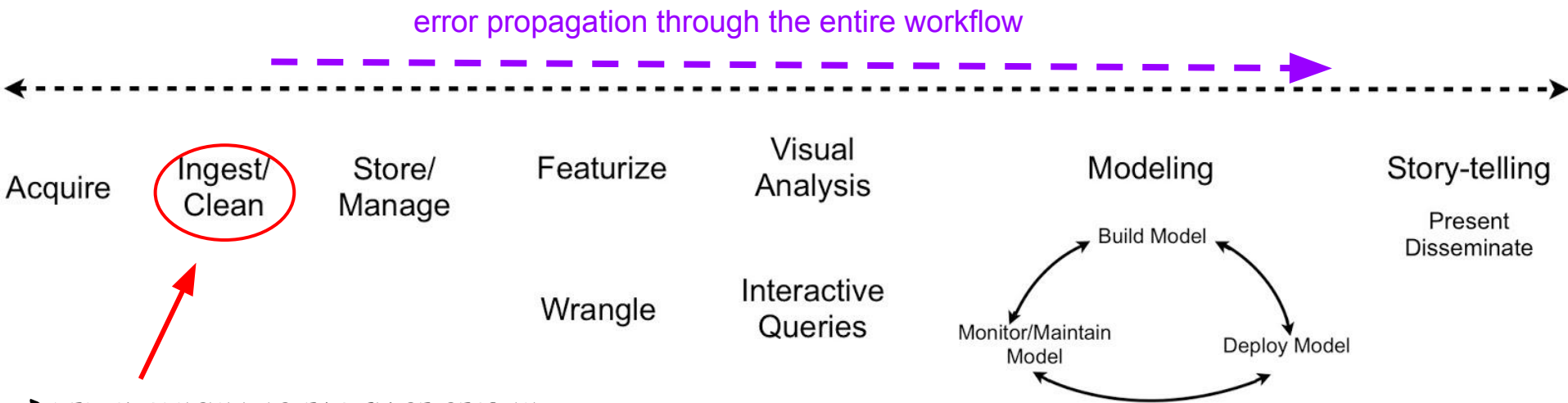

Effective Exploiting of Statistical Relational Learning for Data Cleaning

Larysa Visengeriyeva

Technische Universität Berlin, Germany

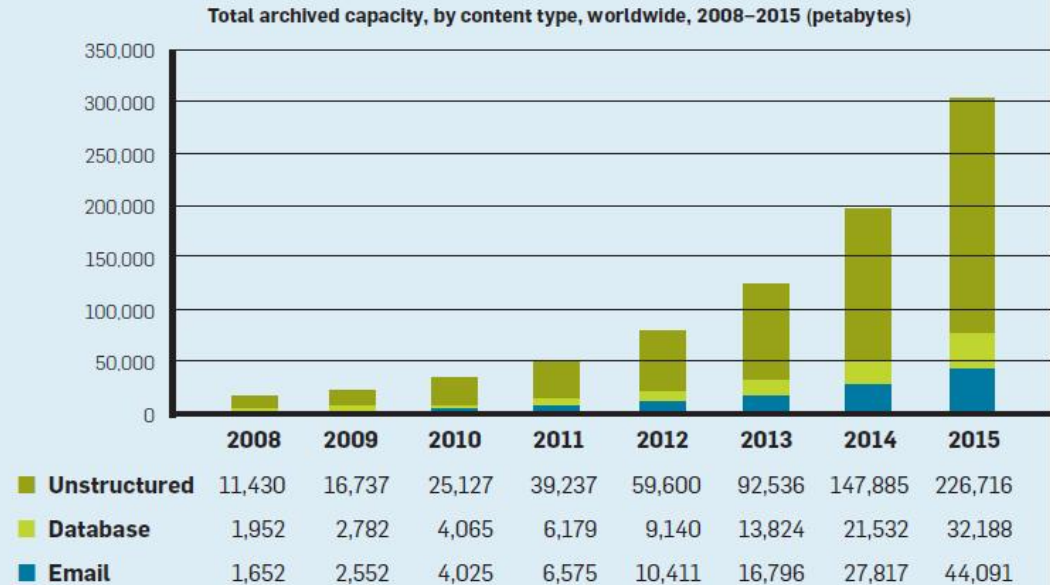
Data Science Workflow



DATA CLEANSING IS THE FIRST STEP IN ANY DATA SCIENCE WORKFLOW, AND OFTEN THE MOST IMPORTANT.

Problem Statement

Increasing need for a new approaches in **data cleaning** (for **non-relational** use-cases, e.g Web data)



Effective Data Cleaning

GOAL

1. Cleaning relational data

2. Cleaning semistructured data

SYSTEM REQUIREMENTS:

1.
 - Joint modeling for data cleaning rules
2.
 - Defining “soft” and “hard” data cleaning rules
 - Flexible rules modeling (non-constraints rules)

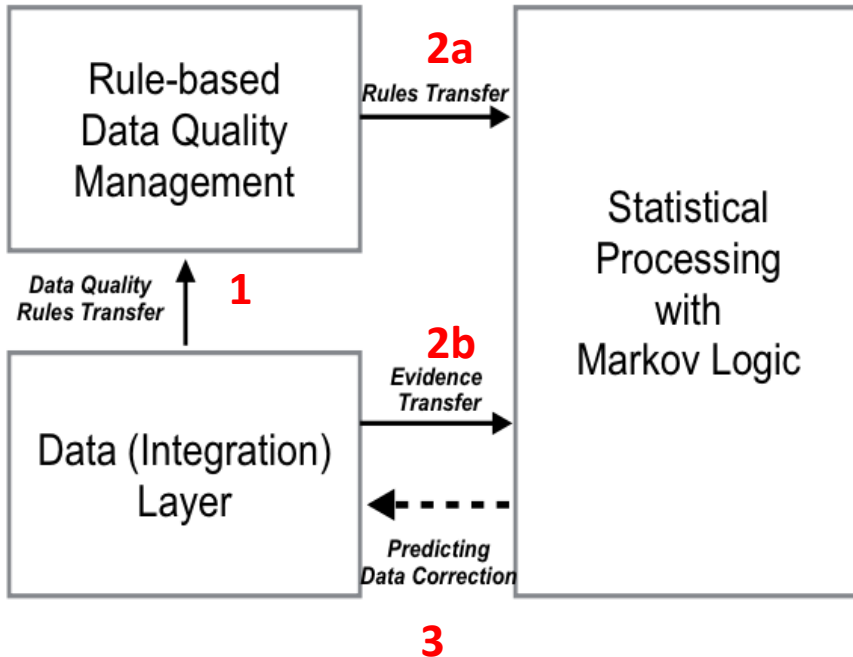
PROBLEM SPACE

Our Idea

... adapt **probabilistic graphical models** to a data cleaning.

The main goal is to translate the constraint-based **data quality rules** into a **predictive model** in order to infer errors in uncertain data and their sources by using **joint inference**.

Method



1. DQ Rules are data agnostic and therefore derived from data. (We assume that constraints are already defined for the data).

2a. Create predictive model: DQ Rules are expressed in first-order logic formulas and then as Markov logic formulas (which is a template for Markov Networks)

2b. Existing data are evidence for the predictive model created in (2a).

3. Running inference on Markov Networks predicts the data correction.

Rules Taxonomie

"HARD" RULES

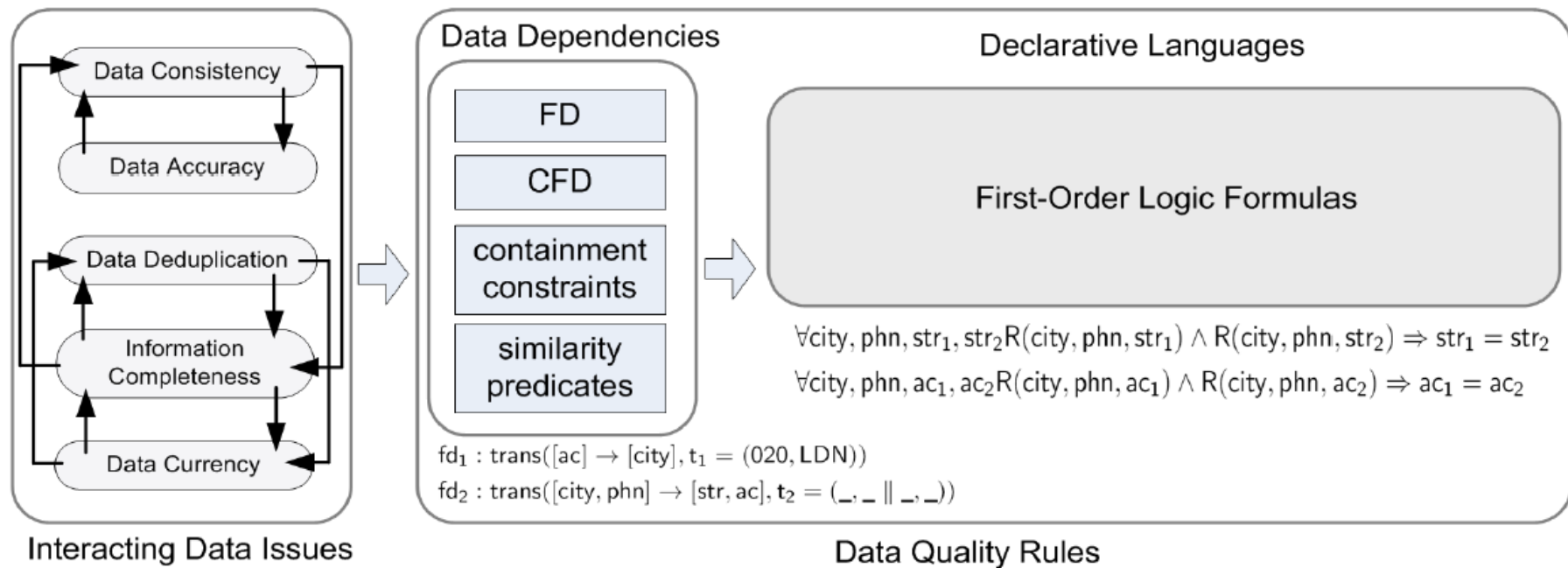
- ❖ Primary Key constraints
- ❖ Constant CFDs

"SOFT" RULES

(FOR SOME RULES MAY NOT EXIST
SUFFICIENT KNOWLEDGE)

- ★ non-Primary Key constraints
- ★ CDFs, MDs, CINDs
- ★ "Soft" FDs
- ★ Similarity predicates

How to compile data quality rules into Markov logic?



Why Markov Logic?

RECAP TWO REQUIREMENTS:

- 1 ○ **Joint modeling** for data cleaning rules
- 2 ○ Defining “**soft**” and “**hard**” data cleaning rules
- Modeling **conflicting rules**
- **Flexible rules** modeling (non-constraints rules)



MARKOV LOGIC:

- **Joint inference over hidden predicates** (by using Markov Networks)
- Defining “**soft**” and “**hard**” first-order formulae
- Modeling **contradicting formulas**
- **Declarative formalism**

Whole Stack in one Example:

INPUT:

Dirty data D_0

id	name	zip	city
1	Dirk	10000	Hamburg
2	Dirk	10000	Berlin
3	Dirk	20000	Berlin

Data cleaning based on conditional functional dependencies (\mathbf{cfd}_1 \mathbf{cfd}_2) and matching rule (\mathbf{md}_3):

\mathbf{cfd}_1 : $[\text{zip}=10000] \rightarrow [\text{city}=\text{Berlin}]$

\mathbf{cfd}_2 : $[\text{city}=\text{Hamburg}] \rightarrow [\text{zip}=20000]$

\mathbf{md}_3 : $R[\text{name}, \text{zip}] = R[\text{name}, \text{zip}] \rightarrow R[\text{city}] \not\Leftarrow R[\text{city}]$

SCENARIO I: CFD, + MD₃

Dirty data D₀

id	name	zip	city
1	Dirk	10000	Hamburg
2	Dirk	10000	Berlin
3	Dirk	20000	Berlin

Dirty data D₁

cf_d₁: [zip=10000] → [city=Berlin]

id	name	zip	city
1	Dirk	10000	Berlin
2	Dirk	10000	Berlin
3	Dirk	20000	Berlin

Cleaned data D₂

md₃: R[name, zip]=R[name, zip] → R[city] ⇐ R[city]

id	name	zip	city
1+2	Dirk	10000	Berlin
3	Dirk	20000	Berlin

SCENARIO II: $CFD_2 + MD_3$

Dirty data D_0

id	name	zip	city
1	Dirk	10000	Hamburg
2	Dirk	10000	Berlin
3	Dirk	20000	Berlin

Dirty data D_1

cfd_2 : $[city=Hamburg] \rightarrow [zip=20000]$

id	name	zip	city
1	Dirk	20000	Hamburg
2	Dirk	10000	Berlin
3	Dirk	20000	Berlin

md_3 : $R[name, zip]=R[name, zip] \rightarrow R[city] \rightleftharpoons R[city]$

TWO POSSIBLE REPAIRS

id	name	zip	city
1	Dirk	20000	Hamburg
2	Dirk	10000	Berlin
3	Dirk	20000	Berlin

$D_{2.1}$

id	name	zip	city
1	Dirk	20000	Berlin
2	Dirk	10000	Berlin
3	Dirk	20000	Berlin

$D_{2.2}$

DATA CLEANING RULES

cmd₁: [zip=10000] \rightarrow [city=Berlin]

cmd₂: [city=Hamburg] \rightarrow
[zip=20000]

md₃:
 $R[\text{name}, \text{zip}] = R[\text{name}, \text{zip}] \rightarrow R[\text{city}] \Leftrightarrow R[\text{city}]$

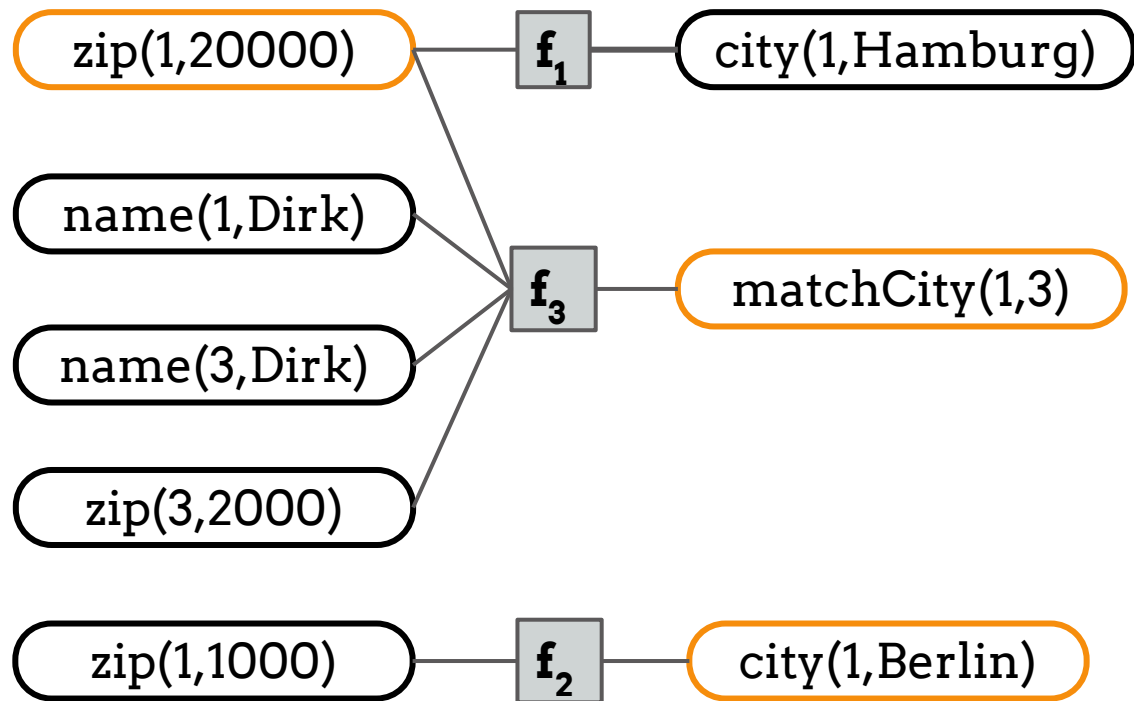
FIRST-ORDER LOGIC TRANSLATION

f₁: $\text{zip}(\text{id}, 10000) \Rightarrow \text{city}(\text{id}, \text{Berlin})$

f₂: $\text{city}(\text{id}, \text{Hamburg}) \Rightarrow \text{zip}(\text{id}, 20000)$

f₃: $\text{name}(\text{id}_1, n) \wedge \text{zip}(\text{id}_1, z) \wedge \text{name}(\text{id}_2, n) \wedge$
 $\text{zip}(\text{id}_2, z) \Rightarrow \text{matchCity}(\text{id}_1, \text{id}_2)$

MARKOV LOGIC TO FACTOR-GRAPH TRANSLATION



EVIDENCE NODES



PREDICTION NODES

Semistructured Data Model and Schema

DATA MODEL

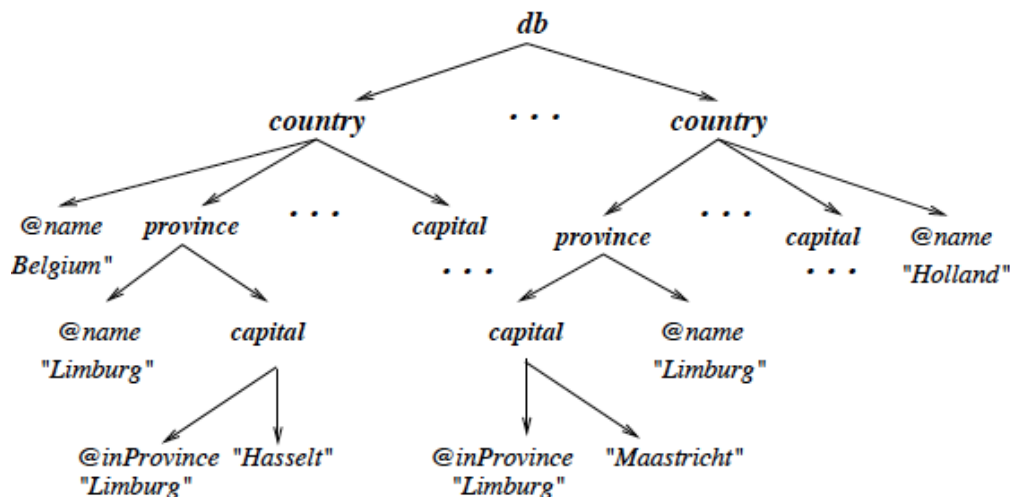
```
<value> ::= <atom> | <array> | <record>
<atom>  ::= <string> | <binary> | <double> |
           <date> | <boolean> | 'null' | ...
<array> ::= '[' ( <value> (',' <value>)* )? ']'
<record> ::= '{' ( <field> (',' <field>)* )? '}'
<field>  ::= <name> ':' <value>
```

SCHEMA

```
<schema> ::= <basic> '??' ( '|' <schema> )*
<basic>  ::= <atom> | <array> | <record>
           | 'nonnull' | 'any'
<atom>   ::= 'string' | 'double' | 'null' | ...
<array>  ::= '[' ( <schema> (',' <schema>)*
                  '...'? )? ']'
<record> ::= '{' ( <field> (',' <field>)* )? '}'
<field>  ::= ( <name> '??' | '*' ) ( ':' <schema> )?
```

Semistructured Data: FD Definition

Functional dependencies on semistructured data [21] is an expression of the form $\varphi : L \rightarrow R$, where L and R are paths (similar to the XPATH notation). Semistructured data satisfies φ iff for any tree tuples t_1, t_2 , if $t_1.L = t_2.L$, then $t_1.R = t_2.R$.



FDs in Clausal Form



$$\forall P_{c1}, P_{t1}, P_{c2}, P_{t2}, L, R_1, R_2 \left(P_{c1}(P_{t1}(L \rightarrow \epsilon)) \wedge P_{c2}(P_{t2}(L \rightarrow \epsilon)) \right. \\ \left. \wedge P_{c1}(P_{t1}(R_1 \rightarrow \epsilon)) \wedge P_{c2}(P_{t2}(R_2 \rightarrow \epsilon)) \Rightarrow R_1 = R_2 \right)$$

The *clausal form* of φ will be useful for the later translation into the Markov logic program:

$$\forall P_{c1}, P_{t1}, P_{c2}, P_{t2}, L, R_1, R_2 \left(\neg P_{c1}(P_{t1}(L \rightarrow \epsilon)) \vee \right. \\ \neg P_{c2}(P_{t2}(L \rightarrow \epsilon)) \vee \neg P_{c1}(P_{t1}(R_1 \rightarrow \epsilon)) \vee \neg P_{c2}(P_{t2}(R_2 \rightarrow \epsilon)) \\ \left. \vee R_1 = R_2 \right)$$

Compilation into Markov Logic

FD

$\varphi: L \rightarrow R$

CLAUSAL FORM

$$\forall P_{c1}, P_{t1}, P_{c2}, P_{t2}, L, R_1, R_2 \left(\neg P_{c1}(P_{t1}(L \rightarrow \epsilon)) \vee \neg P_{c2}(P_{t2}(L \rightarrow \epsilon)) \vee \neg P_{c1}(P_{t1}(R_1 \rightarrow \epsilon)) \vee \neg P_{c2}(P_{t2}(R_2 \rightarrow \epsilon)) \vee R_1 = R_2 \right)$$

ATOM

$P_c(P_t(S \rightarrow \epsilon))$

MARKOV LOGIC PREDICATE $S(\text{tree-tuple-id}(P_c // P_t), \text{value})$

Compilation into Markov Logic

$\varphi: L \rightarrow R$

$$\forall P_{c1}, P_{t1}, P_{c2}, P_{t2}, L, R_1, R_2 \left(\neg P_{c1}(P_{t1}(L \rightarrow \epsilon)) \vee \neg P_{c2}(P_{t2}(L \rightarrow \epsilon)) \vee \neg P_{c1}(P_{t1}(R_1 \rightarrow \epsilon)) \vee \neg P_{c2}(P_{t2}(R_2 \rightarrow \epsilon)) \vee R_1 = R_2 \right)$$

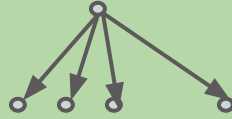


$$L(\text{id}_1(P_c // P_t), l) \wedge L(\text{id}_2(P_c // P_t), l) \wedge R(\text{id}_1(P_c // P_t), r_1) \wedge R(\text{id}_2(P_c // P_t), r_2) \Rightarrow (r_1 = r_2)$$

$$\begin{aligned} & \text{name}(\text{id}(\text{country1} // \text{province1}), n) \wedge \text{name}(\text{id}(\text{country1} // \text{province2}), n) \wedge \\ & \text{capital}(\text{id}(\text{country1} // \text{province1}), \text{cname}_1) \wedge \text{capital}(\text{id}(\text{country1} // \text{province2}), \text{cname}_2) \\ & \Rightarrow \text{sameCapital}(\text{id}(\text{country1} // \text{province1}), \text{id}(\text{country1} // \text{province2})) \end{aligned}$$

STRUCTURED DATA

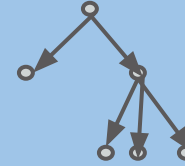
DATA FORMAT



CLEANING RULES

- Data cleaning rules (based on integrity constraints)
- Data & schema agnostic

SEMISTRUCTURED DATA



- Use integrity constraints for semistructured data
- Schema is not obviously specified
- Flexible fields/attributes
- Data cleaning rules are path/navigation based



Markov logic interface (First-order logic)

CLEANING ENGINE

inference alg 1

inference alg 2

inference alg 3

...

inference alg N

Evaluation Strategy

PARAMETERS

Noise %

Data
size

Inference
parameters

Various data
cleaning rules

METRICS

Method quality (F_1 -
score)



Method accuracy



Efficiency (runtime
and memory
consumption)



Initial Experiments

METHOD QUALITY:

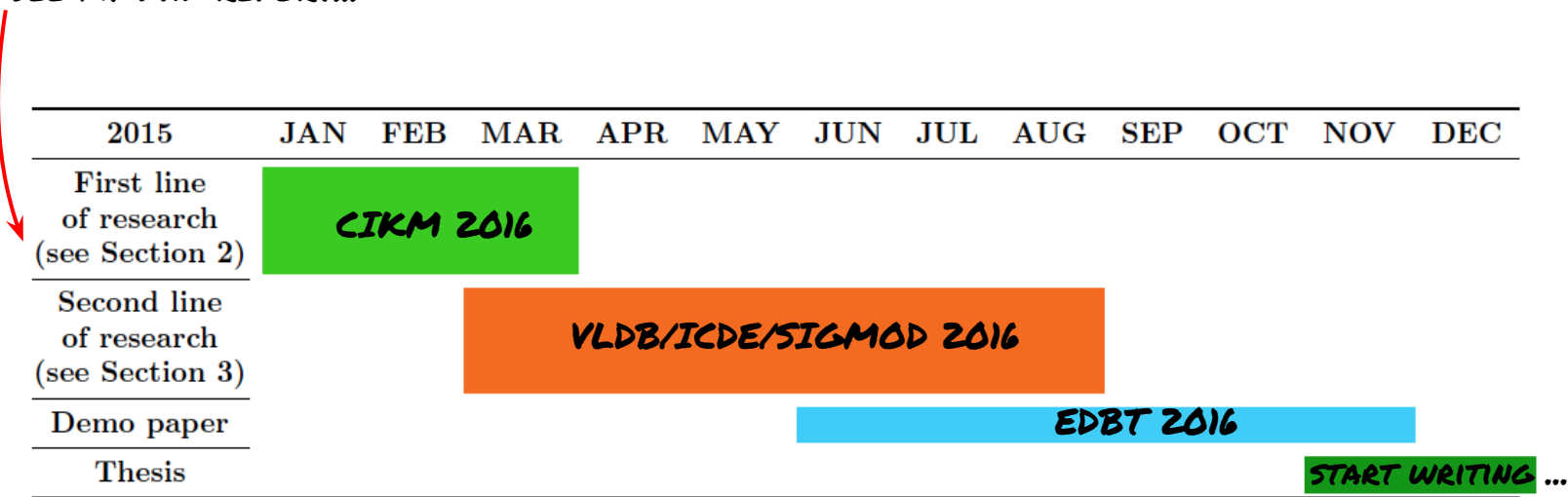
Cleaning relational data: Deduplication (Matching Dependencies) + Improving Accuracy (Conditional Functional Dependencies)

%noise	Data repairing with CFDs									Data repair with CFD and MD
	ADDR	CITY	COND	COUNT	HOSPNAME	MEASURE	PHONE	STATE	ZIP	
2	0.9952	0.9944	0.9976	0.9941	1.0	0.9954	1.0	0.9954	1.0	0.9976
4	0.9944	0.9988	1.0	0.9978	0.9964	0.9899	0.9963	0.9957	0.9929	0.9957
6	0.9964	0.9967	0.9969	0.996	0.9963	0.997	0.9947	0.9943	0.9971	0.9992
8	0.9978	0.9931	0.9942	0.9955	0.9979	0.9976	0.997	0.9948	0.9977	0.9948
10	0.9949	0.9942	0.9955	0.9995	0.9972	0.9944	0.9951	0.9931	0.9977	0.9953
20	0.9963	0.996	0.9943	0.9957	0.995	0.9944	0.9955	0.9965	0.9959	0.9968
30	0.9963	0.9966	0.9962	0.9971	0.9954	0.9961	0.9971	0.9944	0.9962	0.9972

Table 3: Detailed HOSP evaluation: F_1 measure of each attribute and noise per cent. Repair and matching performed separately from each other.

Outlook

SEE MY PHD REPORT...







Example (how data quality rules influence each other)

Data cleaning based on functional dependencies

↓ **fd1: [A] -> [B]**
fd2: [B] -> [C]

dirty data D_0

	A	B	C
t1	a1	b1	c1
t2	a2	b2	c3
t3	a1	null	c1

repaired data D_1

	A	B	C
t1	a1	b1	c1
t2	a2	b2	c3
t3	a1	b1	c1

deduplicated data D_2

	A	B	C
t1'=t1+t3	a1	b1	c1
t2	a2	b2	c3

Null value imputation t3[B]

tuples deduplication t1 and t3

Going top-down

1. What is data quality?
2. What are constraint-based data quality rules?
3. How to compile data quality rules into the probabilistic graphical models?
4. What is Markov logic?
5. How we predict data correction with Markov logic?

What is Data Quality?

5 Dimensions of Data Quality:

1. Consistency
2. Currency
3. Accuracy
4. Deduplication
5. Information Completeness

Conditional Functional Dependency

Definition 2. A **conditional functional dependency** (CFD) defined on the relational schema $R(U)$ is a pair $R(X \rightarrow Y, T_p)$, where

1. $X \rightarrow Y$ is a standard FD;
2. T_p is a set of constraints holding on the particular subset of tuples. \square

EXAMPLE:

$\text{cfd}_1 : \text{TRANS}([\text{ZIPCODE}] \rightarrow [\text{CITY}], \text{t}_1 = (90001 \parallel \text{Los Angeles}))$

Matching Dependency

Definition 3. A **Matching Dependency** (MD) for schemas R_1 and R_2 is syntactically defined as:

$$R_1[X_1] \approx R_2[X_2] \rightarrow R_1[Y_1] \rightleftharpoons R_2[Y_2]$$

where X_1 and X_2 are pairwise compatible sets of attributes in R_1 and R_2 , respectively; similarly for Y_1 and Y_2 ; \approx indicates similar attributes and \rightleftharpoons is called the *matching operator*. \square

EXAMPLE:

$\text{md}_1 : \text{TRANS}[\text{LASTNAME}, \text{CITY}, \text{STREET}] = \text{CUST}[\text{LASTNAME}, \text{CITY}, \text{STREET}]$

$\wedge \text{TRANS}[\text{FIRSTNAME}] \approx \text{CUST}[\text{FIRSTNAME}] \rightarrow \text{TRANS}[\text{FirstName}] \rightleftharpoons \text{CUST}[\text{FirstName}]$

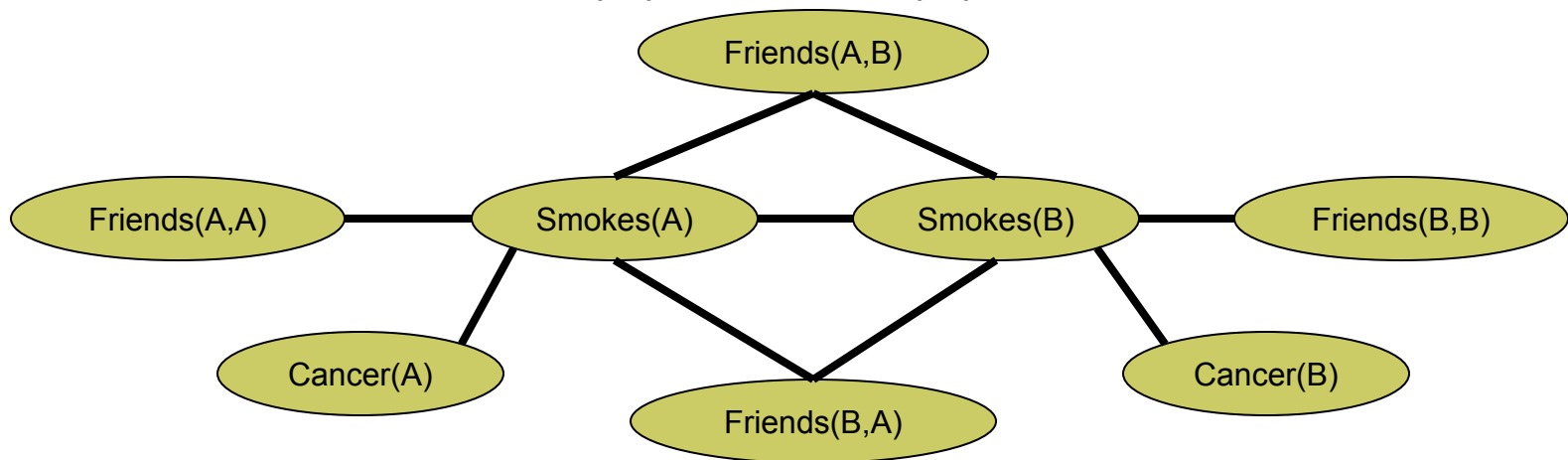
Example: Friends & Smokers



1.5 $\forall x \text{ Smokes}(x) \Rightarrow \text{Cancer}(x)$

1.1 $\forall x, y \text{ Friends}(x, y) \Rightarrow (\text{Smokes}(x) \Leftrightarrow \text{Smokes}(y))$

Two constants: **Anna (A)** and **Bob (B)**



Markov Logic Networks



- MLN (weights+FO formulas) is **template** for ground Markov nets
- Probability of a world x :

$$P(x) = \frac{1}{Z} \exp \left(\sum_i w_i n_i(x) \right)$$

WEIGHT OF FORMULA i

NO. OF TRUE GROUNDINGS OF FORMULA i IN x

- **Typed** variables and constants greatly reduce size of ground Markov net
- Functions, existential quantifiers, etc.
- Infinite and continuous domains

FD:[C_CUSTKEY -> C_NAME]

O_ORDERKEY	...	C_CUSTKEY	C_NAME	C_ADDR	C_NATION_KEY	C_PHONE
O_ID_1		C_ID_1	Customer#1	A1	N_ID_15	25-989-741-2988
O_ID_100		C_ID_1	NULL	A1	N_ID_15	25-989-741-2988
O_ID_101		C_ID_1	Customer#1xy	A1	N_ID_15	25-989-741-2988
O_ID_128		C_ID_28	Customer#28	A28	N_ID_8	18-774-241-1462
O_ID_229		C_ID_28	Customer#28	A2XXX	N_ID_8	18-774-241-1462
O_ID_3		C_ID_3	Customer#3	A3	N_ID_1	11-719-748-3364
O_ID_103		C_ID_3	Customer#3	A3	N_ID_1	11-719-748-3364
O_ID_255		C_ID_55	Customer#55	A55	N_ID_3	18-774-241-1462
O_ID_149		C_ID_49	Customer#49	A49	N_ID_3	18-774-241-1462
O_ID_14		C_ID_14	Customer#14	A14	N_ID_5	18-774-241-1462

customer 1
(C_NAME contains
noisy data and
missing values)

customer 28
(C_ADDR contains
noisy data)

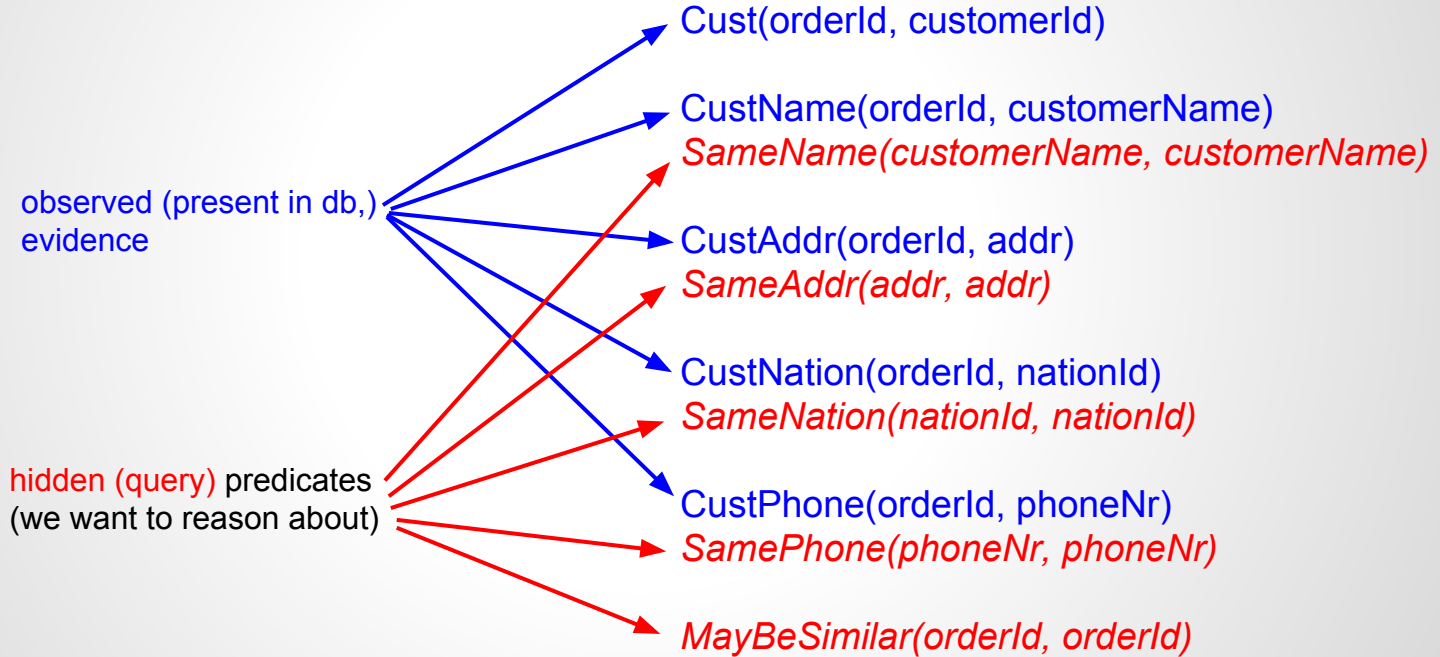
customer 3
(no errors)

MODELING INTEGRITY CONSTRAINTS: PREDICATES DECLARATION

Database attributes -> Predicates

Type -> Predicate type

Data -> domain data



MODELING INTEGRITY CONSTRAINTS: FIRST-ORDER LOGIC FORMULAS

1.0 $\text{Cust}(\text{id1}, x) \wedge \text{Cust}(\text{id2}, x) \wedge \text{CustName}(\text{id1}, n1) \wedge \text{CustName}(\text{id2}, n2) \Rightarrow (\text{MaybeSimilar}(\text{id1}, \text{id2}) \wedge \text{SameName}(n1, n2))$

0.9 $\text{Cust}(\text{id1}, x) \wedge \text{Cust}(\text{id2}, x) \wedge \text{MaybeSimilar}(\text{id1}, \text{id2}) \wedge (\text{CustName}(\text{id1}, n) \vee \text{CustName}(\text{id2}, n)) \Rightarrow (\text{MayHaveName}(\text{id1}, n) \wedge \text{MayHaveName}(\text{id2}, n))$

MODELING EQUALITY

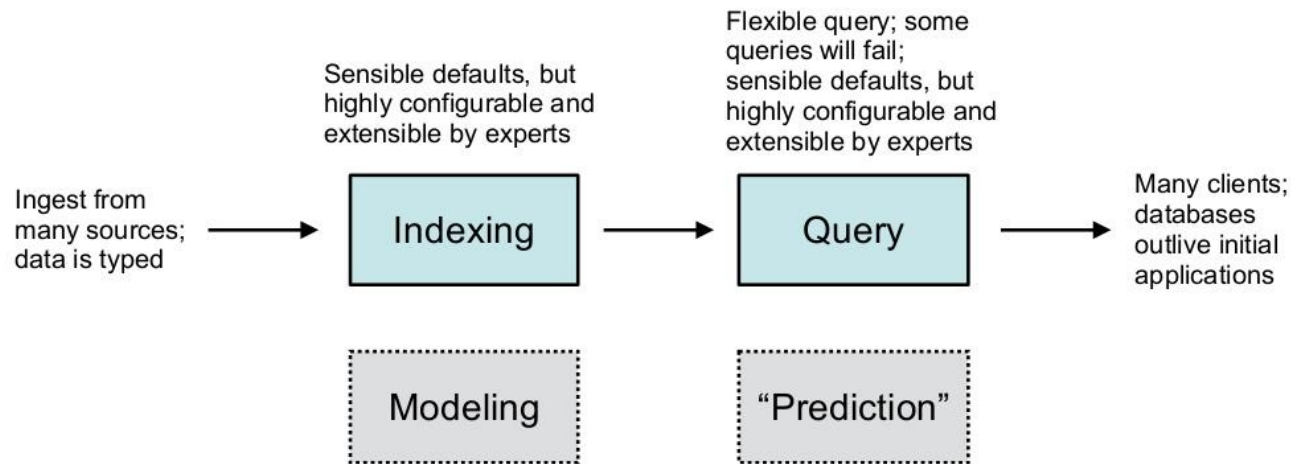


OUR TRICK



How we predict data correction with Markov logic?

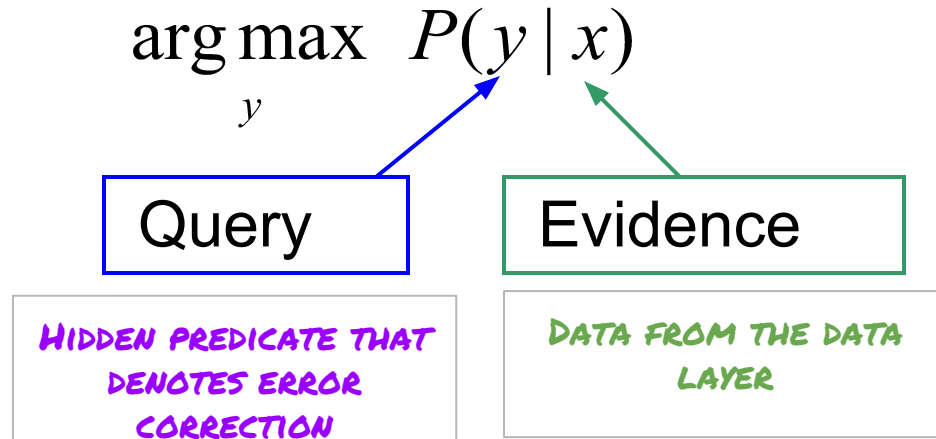
The Database Analogy



MAP/MPE Inference



- **Problem:** Find most likely state of world given evidence



MAP/MPE Inference



- **Problem:** Find most likely state of world given evidence

$$P(X = x) = \frac{1}{Z} \exp \left(\sum_i w_i n_i(x) \right) \quad \Rightarrow \quad \arg \max_y \sum_i w_i n_i(x, y)$$

Dirty Data:
**CUSTOMERS NAME
 IS MISSING OR
 INCORRECT**

O_ORDERKEY	...	C_CUSTKEY	C_NAME	C_ADDR	C_NATION_KEY	C_PHONE
O_ID_1		C_ID_1	Customer#1	A1	N_ID_15	25-989-741-2988
O_ID_100		C_ID_1	NULL	A1	N_ID_15	25-989-741-2988
O_ID_101		C_ID_1	Customer#1xy	A1	N_ID_15	25-989-741-2988

DATA (EVIDENCE) FROM DATABASE

```

Customer(O_ID_1, C_ID_1)
CustName(O_ID_1, Customer#1)
CustAddr(O_ID_1, A1)
CustNation(O_ID_1, N_ID_15)
CustPhone(O_ID_1, PHONE_25-989-741-2988)
---
Customer(O_ID_100, C_ID_1)
/* name is missing */
CustAddr(O_ID_100, A1)
CustNation(O_ID_100, N_ID_15)
CustPhone(O_ID_100, PHONE_25-989-741-2988)
---
Customer(O_ID_101, C_ID_1)
CustName(O_ID_101, Customer#1xy)
CustAddr(O_ID_101, A1)
CustNation(O_ID_101, N_ID_15)
CustPhone(O_ID_101, PHONE_25-989-741-2988)

```

INFERENCE (HIDDEN PREDICATES) RESULTS

```

MaybeSimilar(O_ID_101,O_ID_100) 0.817574
MaybeSimilar(O_ID_100,O_ID_1) 0.817574

```

```

SameName(Customer#1xy, Customer#1) 0.622459

```

```

MayHaveName(O_ID_100, Customer'000000001) 0.588274
MayHaveName(O_ID_100, Customer'000000001xy) 0.588274

```

**JOINT
 DEDUPLICATION AND
 NULL VALUE
 COMPUTATION**

Markov Logic Advantages for Data Cleaning

1. Defining “soft” and “hard” data cleaning rules
2. Joint inference of the hidden predicates
(where are errors?)
3. Ability to model contradicting rules

Conflicting rules

CFD: $R([A] \rightarrow [B], T_0)$

where T_0 consists of the two pattern tuples $(_ || b)$ and $(_ || c)$ with $b \neq c$.

Meaning: The first pattern requires $t[B]=b$, at the same time the second pattern says that $t[B]$ must be c .

Can we model such conflicting rules with Markov Logic?

Modeling *extended* conditional dependencies (I)

CFDs extended allowing disjunction, negation and pattern tuples:

eCFD: $R([city] \rightarrow [AC], (NY \mid \mid \{212, 718, 646\}))$

“for tuples t such that $t[city]=NY$, their area code should come from a set $\{212, 718, 646\}$ ”

Can we model such extended constraints with Markov Logic?

Modeling *extended* conditional dependencies (II)

CFDs extended allowing predicate-patterns:

CFD: $R([\text{book}] \rightarrow [\text{price}], (_ \mid \mid > 0))$

“for every book, its price should be greater than zero”

Can we model such extended constraints with Markov Logic?

Modeling *extended* conditional dependencies (III)

CFDs with cardinality constraints and synonym rules:

$\text{CFD}^c: R([\text{country}, \text{zip}] \rightarrow [\text{street}], (\text{UK}, _ || _), c) \ c \subseteq \mathbb{N}$

“For all tuples t such that $t[\text{country}, \text{zip}] \asymp (\text{UK}, _)$, all tuples that are synonymous to t in country and zip can only have at most c distinct street values.”

Can we model such extended constraints with Markov Logic?

\asymp is an operator, defining that either $[\text{country}, \text{zip}] = (\text{UK}, _)$ or $[\text{country}, \text{zip}] = (_, _)$

Denial Constraints

$\neg(G(g,f,n,r,c,a,s), G(g',f',n',r',c',a',s'), (r=r'), (c='NYC'), (c'\neq'NYC'), (s'>s))$

states: *Every time there are two employees with the same rank, one in NYC and one in a different city, there is a violation if the salary of the second is greater than the salary of the first.*

Can we model denial constraints with Markov Logic?

Denial constraints are the first-order formulae:

$\forall \mathbf{x} \neg(R_1(\mathbf{x}_1) \wedge \dots \wedge R_n(\mathbf{x}_n) \wedge P_1 \wedge \dots \wedge P_m)$ where R_i is a relational atom and P_j is a predicate.

Quality Rules with Operators

(Build-in) predicates containing the following operators: $=, <, >, \approx, \leq, \geq, \neq$

CAN BE EXPRESSED AS DENIAL CONSTRAINTS

Quality improvement for schemaless data

ONGOING WORK..

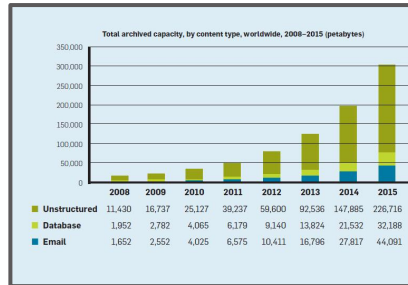
Cleaning semi-structured data

data analyzed have
fundamental quality
problems

data is in
semi-structured
form

need for a new
approaches in
data cleaning
for **non-**
relational use-
cases

“garbage in, garbage out.”



Schema vs. Schemaless

Database Schema	Schemaless data
Describes a logical structure and semantic of data	Denotes data where the logical structure and semantic are (yet) unknown.

	Title	Year	Length	Filmtype
r_0	Amelie	2001	NULL	color
r_1	Babe	1995	89	color
r_2	Cocktail	1988	NULL	color
r_4	NULL	1996	134	color
r_5	Frenzy	NULL	116	color
r_6	Gaslight	1944	114	bw
r_7	Hamlet	1990	242	NULL
r_8	NULL	1958	100	color
r_9	John Q	2002	116	color

	Title	Year	Filmtype			
r_0	Amelie	2001	color			
	Title	Year	Length	Filmtype		
r_1	Babe	1995	89	color		
	Title	Year	Filmtype	Noise-4		
r_2	Cocktail	1988	color	noiseval		
	Length					
r_3	75					
	Year	Length	Filmtype	Noise-2		
r_4	1996	134	color	noiseval		
	Title	Filmtype	Length	Noise-4	Noise-5	
r_5	Frenzy	116	color	noiseval	noiseval	
	Title	Year	Length	Filmtype	Noise-1	
r_6	Gaslight	1944	114	bw	noiseval	
	Title	Year	Length			
r_7	Hamlet	1990	242			
	Year	Filmtype	Length	Noise-5		
r_8	1958	color	100	noiseval		
	Title	Year	Length	Filmtype	Noise-2	Noise-3
r_9	John Q	2002	116	color	noiseval	noiseval

What is Markov Logic?

Markov Logic is a knowledge representation system that combines:

- Weighted First-Order Logic
- Markov Networks (undirected prob. graph. models)

What are constraint-based data quality rules?

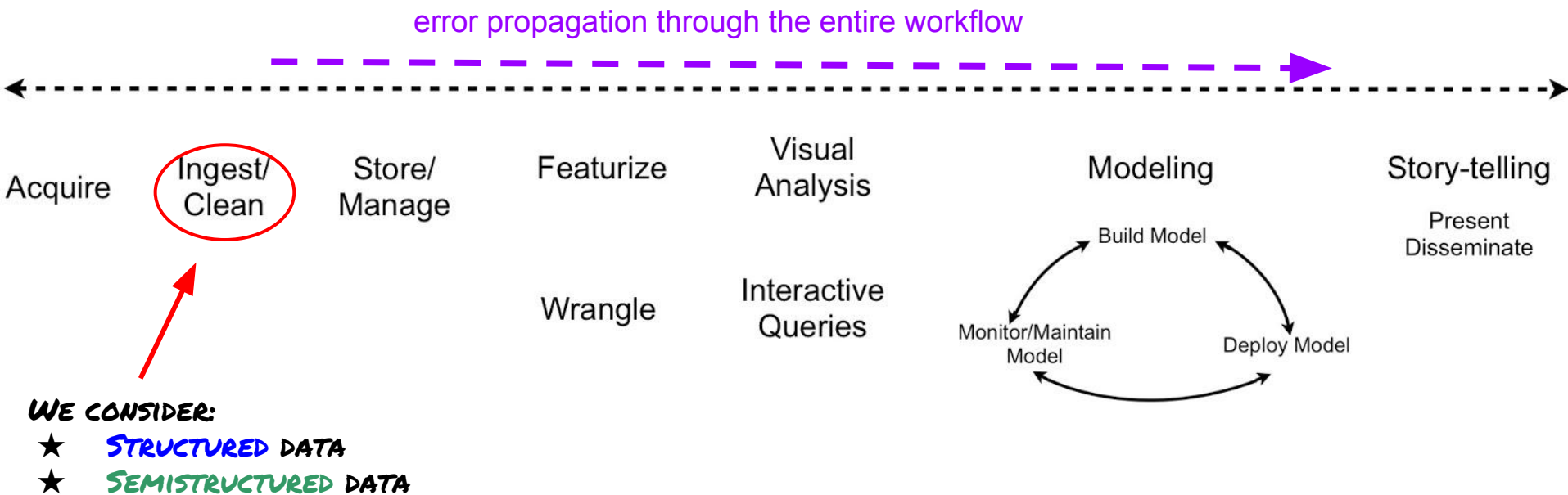
We use data dependencies to formulate data quality rules.

ADVANTAGE:

Data dependencies specify **semantic** of data in a **declarative way**

Example: Functional Dependencies: $[A \rightarrow B_1, B_2, \dots]$

Two Challenges



Two Requirements on Data Cleaning Systems:

1

- **Joint modeling** for data cleaning rules

2

- Defining “**soft**” and “**hard**” data cleaning rules
- Modeling **conflicting** rules
- **Flexible rules** modeling (non-constraints rules)