# MayBMS

**Cornell University — Computer Science**

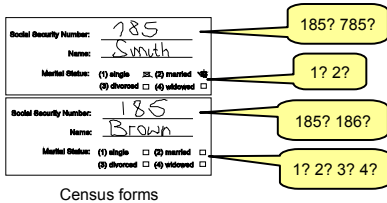# MayBMS: A System for Managing Large Uncertain and Probabilistic Databases

Lyublena Antova — Cornell University
Christoph Koch — Cornell University
Dan Olteanu — Oxford University

## 1. Motivation

Goal: manage uncertain information in different application scenarios: data integration, scientific data collections, census ...



Census forms

There are 2 * 2 * 2 * 4 = 32 possible instances of the forms information!

### Features of MayBMS

✓ scalable DBMS for supporting uncertain and probabilistic data
✓ purely relational representation of attribute-level uncertainty
✓ efficient query processing
✓ query language for probabilistic databases

## 2. U-relational Databases

$R[SSN,N,MS]$: personal information

$U_{R[SSN]}$

| TID | SSN | $V_1 \rightarrow D_1$ |
|---|---|---|
| $t_1$ | 185 | $x \rightarrow 1$ |
| $t_1$ | 785 | $x \rightarrow 2$ |
| $t_2$ | 185 | $y \rightarrow 1$ |
| $t_2$ | 186 | $y \rightarrow 2$ |

$U_{R[MS]}$

| TID | MS | $V_1 \rightarrow D_1$ |
|---|---|---|
| $t_1$ | 1 | $z \rightarrow 1$ |
| $t_1$ | 2 | $z \rightarrow 2$ |
| $t_2$ | 1 | $w \rightarrow 1$ |
| $t_2$ | 2 | $w \rightarrow 2$ |
| $t_2$ | 3 | $w \rightarrow 3$ |
| $t_2$ | 4 | $w \rightarrow 4$ |

$U_{R[N]}$

| TID | N | $V_1 \rightarrow D_1$ |
|---|---|---|
| $t_1$ | Smith | |
| $t_2$ | Brown | |

world table: prob. distribution of the variables

| W | $V \rightarrow D$ | Pr |
|---|---|---|
| | $x \rightarrow 1$ | 0.4 |
| | $x \rightarrow 2$ | 0.6 |
| | $y \rightarrow 1$ | 0.7 |
| | $y \rightarrow 2$ | 0.3 |
| | $z \rightarrow 1$ | 0.8 |
| | $z \rightarrow 2$ | 0.2 |
| | $w \rightarrow 1$ | 0.25 |
| | $w \rightarrow 2$ | 0.25 |
| | $w \rightarrow 3$ | 0.25 |
| | $w \rightarrow 4$ | 0.25 |

$S[SSN,ST]$: credit status

$U_{S[SSN,ST]}$

| TID | SSN | ST | $V_1 \rightarrow D_1$ |
|---|---|---|---|
| $s_1$ | 185 | bad | $w \rightarrow 3$ |
| $s_1$ | 185 | good | $w \rightarrow 4$ |

encode attribute alternatives and correlations with variables

Construct a possible world: pick a value for each variable

| R | SSN | Name | MS |
|---|---|---|---|
| $t_1$ | 785 | Smith | 2 |
| $t_2$ | 186 | Brown | 4 |

| S | SSN | ST |
|---|---|---|
| $s_1$ | 185 | good |

Probability of the world: 0.6*0.3*0.2*0.25=0.009

## 3. Query Language

### World-set Algebra

✓ extend relational algebra with uncertainty-specific constructs e.g.:
- **conf**: confidence computation
- **repair by key**: create the possible repairs of an instance violating a key constraint
- **assert**: remove worlds violating a constraint

✓ semantics: evaluate the query in each world

✓ properties
- generic: independent from representation details
- conservative over relational algebra: right degree of expressive power
- efficient evaluation: simple encoding of positive relational algebra + possible into positive relational algebra queries on U-relational databases

## 4. Query Evaluation

a) $possible(\pi_N(\sigma_{MS=3}(R)))$

Query on U-relational databases:

$\pi_N(\sigma_{MS=3}(U_{R[SSN]} \bowtie_{\varphi \& \psi} U_{R[MS]})))$

$merge(\pi_N R, \pi_{MS} R)$

$\varphi = l.TID = r.TID$
$\psi = (l.V_1 = r.V_1 \rightarrow l.D_1 = r.D_1) \& (l.V_2 = r.V_1 \rightarrow l.D_2 = r.D_1)$

b) $repair\text{-}key_{SSN}(R)$

new variable for each non-unique SSN value

$U_{R[SSN]}$

| TID | SSN | $V_1 \rightarrow D_1$ | $V_2 \rightarrow D_2$ |
|---|---|---|---|
| $t_1$ | 185 | $x \rightarrow 1$ | $x_1 \rightarrow 1$ |
| $t_1$ | 785 | $x \rightarrow 2$ | |
| $t_2$ | 185 | $y \rightarrow 1$ | $x_1 \rightarrow 2$ |
| $t_2$ | 186 | $y \rightarrow 2$ | |

ensure consistent variable assignment

c) $\pi_{MS}(R \bowtie_{SSN} (\sigma_{ST=bad} S))$

Query on column-stores

$\pi_{MS}(merge(\pi_{SSN} R, \pi_{MS} R) \bowtie_{SSN} \sigma_{ST=bad}(S)) \equiv$

$\pi_{MS}(merge(\pi_{\emptyset}(\pi_{SSN} R \bowtie_{SSN} \sigma_{ST=bad}(S)), \pi_{MS} R))$

push **merge** up

$T$ — $U_{T[]}$

| TID | $V_1 \rightarrow D_1$ | $V_2 \rightarrow D_2$ |
|---|---|---|
| $t_1,s_1$ | $x \rightarrow 1$ | $w \rightarrow 3$ |
| $t_2,s_1$ | $y \rightarrow 1$ | $w \rightarrow 3$ |

intermediate result:

final result:

$U_{P[MS]}$

| TID | MS | $V_1 \rightarrow D_1$ | $V_2 \rightarrow D_2$ | $V_3 \rightarrow D_3$ |
|---|---|---|---|---|
| $t_1,s_1$ | 1 | $x \rightarrow 1$ | $w \rightarrow 3$ | $z \rightarrow 1$ |
| $t_1,s_1$ | 2 | $x \rightarrow 2$ | $w \rightarrow 3$ | $z \rightarrow 2$ |
| $t_2,s1$ | 3 | $y \rightarrow 1$ | $w \rightarrow 3$ | |

## 5. Confidence Computation

$U_{R[A]}$

| A | $V_1 \rightarrow D_1$ | $V_2 \rightarrow D_2$ |
|---|---|---|
| 1 | $x \rightarrow 1$ | |
| 1 | $x \rightarrow 2$ | $y \rightarrow 1$ |
| 1 | $x \rightarrow 2$ | $z \rightarrow 1$ |
| 1 | $u \rightarrow 1$ | $v \rightarrow 1$ |
| 1 | $u \rightarrow 2$ | |

confidence of (A:1) = probability of the world-set defined by

$\{\{x \rightarrow 1\}, \{x \rightarrow 2, y \rightarrow 1\}, \{x \rightarrow 2, z \rightarrow 1\}, \{u \rightarrow 1, v \rightarrow 1\}, \{u \rightarrow 2\}\}$



## 6. Experiments

✓ extended TPC-H population generator 2.6 to generate U-relational databases
✓ parameters: scale (s), uncertainty ratio (x), correlation ratio (z), max alternatives per field (m), drop after correlation (p)
✓ each generated world has the sizes of relations and join selectivities of the original TPC-H one-world case
✓ queries translated into SQL and run on PostgresSQL

| s | z | TPC-H dbsize | #worlds | dbsize | #worlds | dbsize |
|---|---|---|---|---|---|---|
| 0.5 | 0.1 | 853 | $10^{64368.0}$ | 3843 | $10^{3.97 \times 06}$ | 5427 |
| 0.5 | 0.5 | 853 | $10^{23528.9}$ | 3856 | $10^{2.33 \times 06}$ | 6682 |
| 1 | 0.1 | 1706 | $10^{87203.0}$ | 7683 | $10^{7.94 \times 06}$ | 11264 |
| 1 | 0.5 | 1706 | $10^{51290.9}$ | 7712 | $10^{4.66 \times 06}$ | 13312 |
| | | | X=0.0 | X=0.001 | | X=0.1 |

Fig. 1: Number of worlds and size in MB of the U-relational db for different scale, uncertainty and correlation ratios.
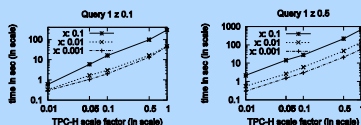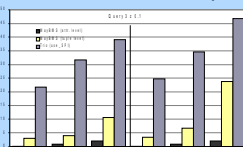


Fig. 2: Query evaluation



Fig. 3: Attribute- vs. tuple-level representation

### 7. Selected Publications

1. Fast and simple relational processing of uncertain data, ICDE'08
2. Conditioning Probabilistic Databases, Technical report '08
3. MayBMS: Managing incomplete information with probabilistic world-set decompositions, ICDE'07, Demo paper
4. 10^10^6 worlds and beyond: efficient representation and processing of incomplete information, ICDE'07
5. From complete to incomplete information and back, SIGMOD'07