# Can We Probabilistically Generate Uniformly Distributed Relation Instances Efficiently?

### Joachim Biskup    Marcel Preuß

Technische Universität Dortmund
Germany

Can We Probabilistically Generate Uniformly Distributed Relation Instances Efficiently?
└─Problem: Can We ... ?

technische universität
dortmund

# Problem: Can We ... ?

Can We Probabilistically Generate Uniformly Distributed Relation Instances Efficiently?
└─ Problem: Can We ... ?

technische universität
dortmund

# Simplified Single-FD Scenario and Probabilistic Instance Generation Task

**Inputs**

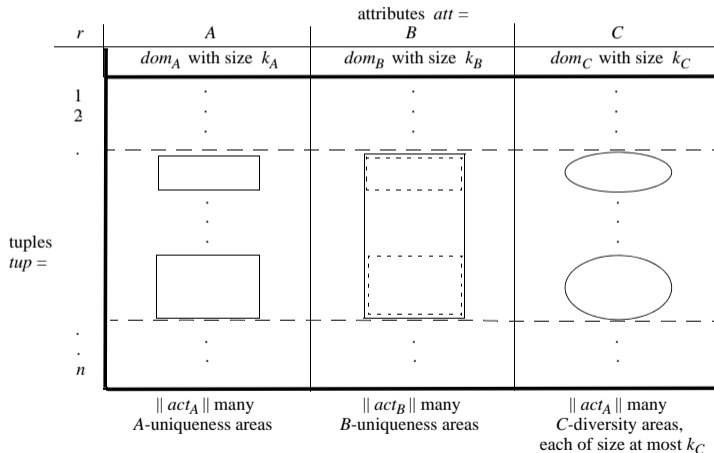| | |
|---|---|
| $(R(\{A, B, C\}), A \rightarrow B)$ | relational schema with **one functional dependency** |
| $k_{att} = \|dom_{att}\| \geq 2$ | size of domain for attribute $att \in \{A, B, C\}$ |
| $n > 0$ | required size (number of tuples) of an instance |

**Output**

$r$          relation instance as an (unordered, duplicate-free) set, to be generated with **uniform probability distribution**

# Uniqueness and Diversity Areas for Single-FD $A \to B$



| $r$ | attributes $att =$ | | |
|---|---|---|---|
| | $A$ | $B$ | $C$ |
| | $dom_A$ with size $k_A$ | $dom_B$ with size $k_B$ | $dom_C$ with size $k_C$ |

tuples $tup =$

$\| act_A \|$ many
$A$-uniqueness areas

$\| act_B \|$ many
$B$-uniqueness areas

$\| act_A \|$ many
$C$-diversity areas,
each of size at most $k_C$

# Combinatorial Analysis

## Number of Single-FD Instances: $n$-$ak_A$-$k_C$-Representations

$$\sum_{\lceil \frac{n}{k_C} \rceil \leq ak_A \leq \min(k_A, n)} \quad \sum_{\substack{1 \leq k \leq n, \\ (s_1, m_1), \ldots, (s_i, m_i), \ldots, (s_k, m_k): \\ 1 \leq s_i < s_{i+1} \leq k_C; \\ 1 \leq m_i \leq n; \\ n = \sum_{i=1}^k s_i \cdot m_i; \\ ak_A = \sum_{i=1}^k m_i}} \left[ \prod_{i=1}^k \binom{k_A - \sum_{1 \leq j < i} m_j}{m_j} \cdot k_B^{ak_A} \cdot \prod_{i=1}^k m_i \cdot \binom{k_C}{s_i} \right]$$

- $1 \leq k \leq n$            number of different sizes of $A$-uniqueness areas
- $(s_1, m_1), \ldots, (s_i, m_i), \ldots, (s_k, m_k)$    sequence of such sizes with their multiplicity
- $1 \leq s_i < s_{i+1} \leq k_C$      such sizes strictly ordered, bounded by cardinality of $dom_C$
- $1 \leq m_i \leq n$        such multiplicities bounded by overall size of instance
- $n = \sum_{i=1}^k s_i \cdot m_i$     such multiplied sizes sum up to overall size of instance
- $ak_A = \sum_{i=1}^k m_i$     such multiplicities sum up to size of selected active $A$-domain

# $n$-$ak_A$-$k_C$-Representations

| tuples/lines | uniqueness areas | sizes | multiplicities | multiplied sizes |
|---|---|---|---|---|
| $1$ | | $s_1$ | | |
| $s_1$ | | | $m_1$ | $s_1 \cdot m_1$ |
| $s_1 \cdot m_1$ | | $s_1$ | | |
| $\cdot$ | | | | |
| $\cdot$ | | | | |
| $\cdot$ | | | | |
| | | $s_k$ | | |
| | | | $m_k$ | $s_k \cdot m_k$ |
| | | $s_k$ | | |
| $n$ | | $s_k$ | | |
| | | | $\Sigma = ak_A$ | $\Sigma = n$ |

# Sophisticated Probabilistic Generation Procedure

Can We Probabilistically Generate Uniformly Distributed Relation Instances Efficiently?
└─ Sophisticated Probabilistic Generation Procedure

technische universität
dortmund

# Probabilistic Generation Procedure: Step I. Preprocessing

1. determine and list all possible sizes $ak_A$ of an active domain $act_A$ for attribute $A$

2. for each listed $ak_A$: (solve Restricted Integer Partition Problem [Euler 1741], i.e.,) determine and list all possible $n$-$ak_A$-$k_C$-representations $S$

3. for each listed $ak_A$, for each listed $n$-$ak_A$-$k_C$-representation $S$: calculate and keep the number $Inst(ak_A, S)$ of complying instances

4. for each listed $ak_A$: calculate and keep the number $Inst(ak_A)$ of complying instances

5. determine and keep the number $Inst$ of all instances

6. annotate each listed size $ak_A$ and each listed $n$-$ak_A$-$k_C$-representation $S$ with probabilities $Inst(ak_A)/Inst$ and $Inst(ak_A, S)/Inst(ak_A)$, respectively

Can We Probabilistically Generate Uniformly Distributed Relation Instances Efficiently?
└ Sophisticated Probabilistic Generation Procedure

technische universität
dortmund

# Probabilistic Generation Procedure: Step II. Generation with Probabilities



select size of $\quad$ select $\qquad$ fill areas of array for:

active domain $\quad$ $n$-$ak_A$-$k_C$-representation $S$ $\quad$ $A$-uniqueness $\qquad$ partial $\qquad$ $C$-diversity

$ak_A$ $\qquad$ $s_1, \ldots, s_{akA}$ $\qquad$ $a_1, \ldots, a_{akA}$ $\quad$ $B$-uniqueness

pairwise different

$\dfrac{Inst(ak_A)}{Inst}$ $\qquad$ $\dfrac{Inst(ak_A,S)}{Inst(ak_A)}$ $\qquad\qquad\qquad$ $\dfrac{1}{Inst(ak_A,S)}$

$\dfrac{1}{Inst}$

Can We Probabilistically Generate Uniformly Distributed Relation Instances Efficiently?
└ Sophisticated Probabilistic Generation Procedure

technische universität
dortmund

# Probabilistic Generation Procedure: Outline of Time Complexity

- ▶ **Part I** (performed only once)
  - ▶ essentially solving the *Restricted Integer Partition Problem*: exponential

- ▶ **Part II** (optimizable if few expected collisions)
  - ▶ applying tools based on a standard *pseudo-random generator*
    - ▶ *give-next-element* operation: constant
    - ▶ *shuffle* operation: linear
    - ▶ adaptation of pseudo-randomness: linear
  - ▶ *selection of different values* by the following alternatives:
    - ▶ repeat the give-next-element operation on a collision
    - ▶ shuffle an array representation of the pertinent set of options and extract
  - ▶ in total: quadratic (plus the time for adaptations)

Can We Probabilistically Generate Uniformly Distributed Relation Instances Efficiently?
└ Failure of Naive Generation Procedures

technische universität
dortmund

# Failure of Naive Generation Procedures

Can We Probabilistically Generate Uniformly Distributed Relation Instances Efficiently?
└ Failure of Naive Generation Procedures

technische universität
dortmund

## "First Choose A-Values, then B-/C-Values" with $n = 2$ and $k_{att} = 2$



| Step 1 | Step 2a | Step 2b | probability |

$r[1,A] = r[2,A];\ act_A = 1;\ s_1 = 2$

| $r[1,A] \in \{1,2\}$ | $r[2,A] \in \{1,2\}$ | $r[1,B] = r[2,B] \in \{1,2\}$ | $\{r[1,C]\,,\,r[2,C]\} \in \{\ \{1,2\}\ \}$ | |

1/2   1 ──── 1 ──────── $\{1,2\}$   1/8

1/2   2 ──── 1 ──────── $\{1,2\}$   1/8

$r[1,A] \neq r[2,A];\ act_A = 2;\ s_1 = 1;\ s_2 = 1$

$r[1,B] \in \{1,2\}$   $r[2,B] \in \{1,2\}$   $\{r[1,C] \in \{1,2\}$   $r[2,C]\} \in \{1,2\}$

*same set* $\{1,2\}=\{2,1\}$,
*thus same instances*

**1/64**
*instance occurring twice*

# Conclusions and Open Answer

## Summary of Current Achievements

▶ systematic counting method for relation instances of required size
for one functional dependency and given domain sizes

▶ identification of the failure of naive approaches to probabilistic instance generation:
"local uniform selections" do not necessarily lead to "global uniformity"

▶ design and verification of a sophisticated probabilistic generation procedure
with uniformly distributed outputs:
  1. combinatorially adapted probabilities, to select a template structure
  2. uniform probabilities, to select actual values

▶ open answer to our question, due to challenging interactions of requirements:
  ▶ by the database schema
  ▶ for the probability distribution of the outputs

▶ even more challenging for more general scenarios