

eXtra Large Joins

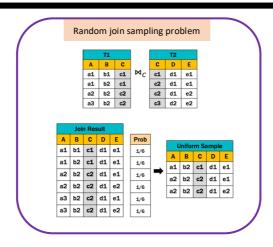
Ali Mohammadi Shanghooshabad

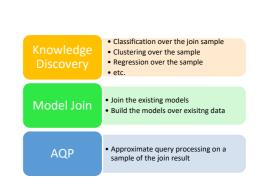
Department of Computer Science University of Warwick Coventry, UK



Skip the full join generation







PGMs as core idea

Schema and an example join query

nation (nationkey NK)
supplier (nationkey NK, suppkey SK)
customer (custkey CK, nationkey NK)
orders (orderkey OK, custkey CK)
lineitem (orderkey OK, linenumber LN)

SELECT n.nationkey as NK, s.suppkey as SK, c.custkey as CK,

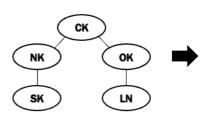
SELECT n.nationkey as NK, s.suppkey as SK, c.custkey as CK,
o.orderkey as OK, l.linenumber as LN

FROM patients a suppliers suppliers and second seco

FROM nation n, supplier s, customer c, orders o, lineitem I
WHERE n.nationkey = s.nationkey

AND s nationkey = c.nationkey

AND s.nationkey = c.nationkey AND c.custkey = o.custkey AND o.orderkey = l.orderkey;



$$P(x) = \frac{1}{Z} \prod_{c \in C} \psi_c(x_c)$$

$$Z = \sum_{x} \prod_{c \in C} \psi_c(x_c)$$

Sum-Product
Message Passing Alg.



Ancestral sampling

Uniform Sample				
СК	NK	ок	SK	LN
a1	b2	c1	d1	e1
a2	b2	c2	d1	e1
a2	b2	c2	d1	e2

Model Join

The result of a *model join* should be "similar" to the result we would have obtained if we joined the underlying tables.

Why data is absent sometimes?

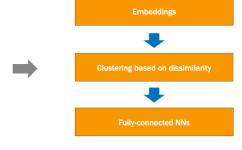
- Storage problem
- Privacy preserving
- Accurate, light and fast in-memory models (DBEst, DeepDB etc.)

Solution?

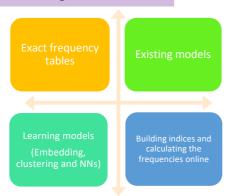
 The idea is that any edge of the PGM graph could be a model providing needed information.

Learning Challenges in tabular data

- Higher Number of distinct values
- Categorical attributes

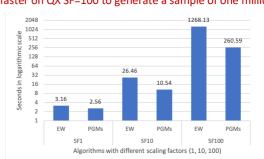


Strategies in the edges of the PGMs



Preliminary results

5X faster on QX SF=100 to generate a sample of one million



 $\,$ EW algorithm is the SOTA approach introduced in Zhao et al, SIGMOD 2018 $\,$