

Supplementary Material: SQL queries for CFD^ps

Shuai Ma, Liang Duan, Wenfei Fan, Chunming Hu, and Wenguang Chen

Below we show how the SQL queries Q_i^c and Q_i^v are generated for validating CFD^ps in Σ_{cfdp}^i , which is an extension of the SQL techniques for CFDs and eCFDs discussed in [2] and [1], respectively.

The queries Q_i^c and Q_i^v for the violations of Σ_{cfdp}^i are given as follows, which capitalize on the data table enc_L , enc_R and enc_{\neq} that encode CFD^ps in Σ_{cfdp}^i .

Q_i^c : **select** $R_i.*$ **from** $R_i, \text{enc}_L L, \text{enc}_R R, \text{enc}_{\neq} N$
where $L.\text{cid} = R.\text{cid}$ **and** $R_i.X \asymp L$ **and** $R_i.X \asymp N$ **and**
not ($R_i.Y \asymp R$ **and** $R_i.Y \asymp N$)

Q_j^v : **select distinct** X_L
from (**select** $L.\text{cid}$ **as** cid , X_L, Y_R **from** $R_i, \text{enc}_L L, \text{enc}_R R, \text{enc}_{\neq} N$
where $L.\text{cid} = R.\text{cid}$ **and** $R_i.X \asymp L$ **and**
 $R_i.X \asymp N$ **and** $R.Y = ' _ '$) **as** M
group by cid, X_L **having count** (**distinct** Y_R) > 1

Here (1) $X = \{A_1, \dots, A_{m1}\}$ and $Y = \{B_1, \dots, B_{m2}\}$ are the sets of attributes in LHS and RHS of Σ_{cfdp}^i respectively; (2) $R_i.X \asymp L$ is the conjunction of

$L.A_j$ **is null** **or** $R_i.A_j = L.A_j$ **or** ($L.A_j = ' _ '$
and ($L.A_j >$ **is null** **or** $R_i.A_j > L.A_j >$)
and ($L.A_j \geq$ **is null** **or** $R_i.A_j \geq L.A_j \geq$)
and ($L.A_j <$ **is null** **or** $R_i.A_j < L.A_j <$)
and ($L.A_j \leq$ **is null** **or** $R_i.A_j \leq L.A_j \leq$))

for each $j \in [1, m_1]$; (3) $R_i.Y \asymp R$ is defined similarly for attributes in Y ; (4) $R_i.X \asymp N$ is the conjunction of

not exists (**select** $*$ **from** N
where $L.\text{cid} = N.\text{cid}$ **and** $N.\text{pos} = \text{'LHS'}$ **and**
 $N.\text{att} = A_j$ **and** $R_i.A_j = N.\text{val}$)

for each $j \in [1, m_1]$; (5) $R_i.Y \asymp N$ is defined similarly, but with $N.\text{pos} = \text{'RHS'}$; (6) X_L is the set of following attributes

(**case when** $L.A_j$ **is not null** **then** $R_i.A_j$ **end**) **as** A_{Lj}

for each $j \in [1, m_1]$; (7) Similarly, Y_R is the set of

(**case when** $R.B_k$ **is not null** **then** $R_i.B_k$ **end**) **as** B_{Rk}

for each $k \in [1, m_2]$; (8) $R.Y = ' _ '$ is the disjunction of $R.B_k = ' _ '$ for each $k \in [1, m_2]$.

Intuitively, detecting violations of CFD^ps is a two-step process. First, query Q_i^c detects single-tuple violations, i.e.,

- S. Ma, L. Duan and C. Hu are with the SKLSDE lab, School of Computer Science and Engineering, Beihang University, China.
E-mail: {mashuai, duanl, hucm}@act.buaa.edu.cn.
- W. Fan is with the RCBd center, Beihang University, China and the School of Informatics, Edinburgh University, UK.
E-mail: wenfei@inf.ed.ac.uk.
- W. Chen is with the Department of Information Management, Peking University, China.
E-mail: chenwg@pku.edu.cn.

Manuscript received XXX, 2014; revised XXX, 2014.

the tuples t in I_i that match the LHS of a CFD^p in Σ_{cfdp}^i , but do not match its RHS. Second, query Q_i^v finds multi-tuple violations, i.e., the tuples t in I_i such that (a) there exists another tuple t' in I_i , t and t' match and agree on the LHS of a CFD^p in Σ_{cfdp}^i , but do not agree on the RHS of the CFD^p.

Example 1: Using the coding of Fig. 4, two SQL queries for checking CFD^ps φ_2 , φ_3 and φ_4 of Fig. 2 are given as follows:

Q_1^c : **select** $R_1.*$ **from** $\text{item } R_1, \text{enc}_L L, \text{enc}_R R, \text{enc}_{\neq} N$
where $L.\text{cid} = R.\text{cid}$ **and**
 $(L.\text{sale}$ **is null** **or** $R_1.\text{sale} = L.\text{sale}$ **or** $L.\text{sale} = ' _ '$) **and**
not exists (**select** $*$ **from** N
where $N.\text{cid} = L.\text{cid}$ **and** $N.\text{pos} = \text{'LHS'}$ **and**
 $N.\text{att} = \text{'sale'}$ **and** $R_1.\text{sale} = N.\text{val}$) **and**
 $(L.\text{price}$ **is null** **or** $R_1.\text{price} = L.\text{price}$ **or** ($L.\text{price} = ' _ '$ **and**
 $(L.\text{price} >$ **is null** **or** $R_1.\text{price} > L.\text{price} >$) **and**
 $(L.\text{price} \geq$ **is null** **or** $R_1.\text{price} \geq L.\text{price} \geq$))) **and**
not exists (**select** $*$ **from** N
where $N.\text{cid} = L.\text{cid}$ **and** $N.\text{pos} = \text{'LHS'}$ **and**
 $N.\text{att} = \text{'price'}$ **and** $R_1.\text{price} = N.\text{val}$) **and**
not (($R.\text{shipping}$ **is null** **or** $R_1.\text{shipping} = R.\text{shipping}$ **or**
 $R.\text{shipping} = ' _ '$) **and**
not exists (**select** $*$ **from** N
where $N.\text{cid} = R.\text{cid}$ **and** $N.\text{pos} = \text{'RHS'}$ **and**
 $N.\text{att} = \text{'shipping'}$ **and** $R_1.\text{shipping} = N.\text{val}$) **and**
 $(R.\text{price}$ **is null** **or** $R_1.\text{price} = R.\text{price}$ **or** ($R.\text{price} = ' _ '$ **and**
 $(R.\text{price} \geq$ **is null** **or** $R_1.\text{price} \geq R.\text{price} \geq$) **and**
 $(R.\text{price} <$ **is null** **or** $R_1.\text{price} < R.\text{price} <$))) **and**
not exists (**select** $*$ **from** N
where $N.\text{cid} = R.\text{cid}$ **and** $N.\text{pos} = \text{'RHS'}$ **and**
 $N.\text{att} = \text{'price'}$ **and** $R_1.\text{price} = N.\text{val}$))

Q_1^v : **select distinct** $\text{sale}_L, \text{price}_L$ **from** (
select $L.\text{cid}$ **as** cid ,
case when $L.\text{sale}$ **is not null** **then** $R_1.\text{sale}$ **end**) **as** sale_L ,
case when $L.\text{price}$ **is not null** **then** $R_1.\text{price}$ **end**) **as** price_L ,
case when $R.\text{shipping}$ **is not null** **then** $R_1.\text{shipping}$ **end**) **as** shipping_R ,
case when $R.\text{price}$ **is not null** **then** $R_1.\text{price}$ **end**) **as** price_R
from $\text{item } R_1, \text{enc}_L L, \text{enc}_R R, \text{enc}_{\neq} N$
where $L.\text{cid} = R.\text{cid}$ **and**
 $(L.\text{sale}$ **is null** **or** $R_1.\text{sale} = L.\text{sale}$ **or** $L.\text{sale} = ' _ '$) **and**
not exists (**select** $*$ **from** N
where $N.\text{cid} = L.\text{cid}$ **and** $N.\text{pos} = \text{'LHS'}$ **and**
 $N.\text{att} = \text{'sale'}$ **and** $R_1.\text{sale} = N.\text{val}$) **and**
 $(L.\text{price}$ **is null** **or** $R_1.\text{price} = L.\text{price}$ **or** ($L.\text{price} = ' _ '$ **and**
 $(L.\text{price} >$ **is null** **or** $R_1.\text{price} > L.\text{price} >$) **and**
 $(L.\text{price} \geq$ **is null** **or** $R_1.\text{price} \geq L.\text{price} \geq$))) **and**
not exists (**select** $*$ **from** N
where $N.\text{cid} = L.\text{cid}$ **and** $N.\text{pos} = \text{'LHS'}$ **and**
 $N.\text{att} = \text{'price'}$ **and** $R_1.\text{price} = N.\text{val}$) **and**
 $(R.\text{shipping} = ' _ '$ **or** $R.\text{price} = ' _ '$)) **as** M
group by $\text{cid}, \text{sale}_L, \text{price}_L$
having count (**distinct** $\text{shipping}_R, \text{price}_R$) > 1

□

ACKNOWLEDGMENTS

This work is supported in part by 973 program (No. 2014CB340300) and NSFC (No. 61322207). Fan is supported in part by 973 Program (No. 2012CB316200), NSFC (No. 61133002), Guangdong Innovative Research Team Program (2011D005), Shenzhen Peacock Program (1105100030834361) of China, and EPSRC (EP/J015377/1) of UK.

REFERENCES

- [1] L. Bravo, W. Fan, F. Geerts, and S. Ma. Increasing the expressivity of conditional functional dependencies without extra complexity. In *ICDE*, 2008.
- [2] W. Fan, F. Geerts, X. Jia, and A. Kementsietsidis. Conditional functional dependencies for capturing data inconsistencies. *TODS*, 33(2), 2008.