

Derby/S A DBMS for Sample-Based Query Answering

Faculty of Computer Science Database Technology Group

Anja Klein anja.klein@sap.com

Rainer Gemulla gemulla@inf.tu-dresden.de

1

3

(5)

(7)

Philipp Rösch philipp.roesch@tu-dresden.de Wolfgang Lehner wolfgang.lehner@tu-dresden.de

2

4

6

8

Sample Definition

- declarative sample definition language
 - $\circ\,$ choice of sampling scheme is left to the system
 - based on information about expected queries and the data set itself
- parameters
 - o base data (as SELECT-statement) & sample size
 - o grouping columns, aggregate functions for optimization

Catalog Tables

- one logical sample may consist of multiple physical samples (e.g., stratification)
- additional system tables

SYSPHYSICALSAMPLE

unique ID
reference to logical
sample
reference to table
with sample data
sampling scheme
current sample size
specific parameters
(optional)
none/group-by/
aggregation
list of aggregates or
grouping attributes

SYSLOGICALSAMPLE

LogicalSampleID	unique ID	
Name	name of sample	
SchemaID	schema of sample	
SQL	CREATE statement	
ResultSetSize	size of base data	
SampleSize	desired sample size	
SizeType	tuples/fraction	
Adaptive	additional memory	
	for optimization (y/n)	
Only	fix optimization	
	attributes (y/n)	
Manager	system/user	
RefreshType	immediate/fast/	
	complete	
RefreshTime	on demand/periodic	
PeriodLength	period length	
RefreshPeriod	unit of period length	

Maintenance

- automatic sample maintenance
 - o samples are kept up-to-date
 - \circ immediate and deferred refresh supported
 - o incremental strategies are used whenever possible

	Incremental	Complete
Immediate	IMMEDIATE	
Deferred	FAST	COMPLETE

• specification (RefreshClause of CREATE-Statement):



Maintenance Behind the Curtain

- intercepting DML operations
 - o capture modifications of base data (INSERT, UPDATE, and DELETE)
 - o compute deltas of base data
 - o derive net effect on the sample
 - \circ depending on the maintenance strategy
 - apply to sample (immediate refresh)
 - write to staging table (deferred refresh)
- REFRESH command
 manually refresh sample

Approximate Queries

- SQL-like queries for approximate query processing
 - queries against base tables
 - distinction between existential (SOME) and statistical (APPROXIMATE) errors



• optional specification of accepted error (future work)

Rewriting the Query Tree

- transparent replacement of base tables by appropriate samples
- example: existing sample over lineitem and orders (lo_sample)



INTERVAL and CONFIDENCE

- INTERVAL) - (=) I - %)

- computation of error bounds with user-defined parameters
 - CONFIDENCE returns the confidence for a user-defined interval

CONFIDENCE - p %

 INTERVAL returns the interval with a user-defined confidence



- implemented as additional aggregation functions
- large-sample confidence intervals

Example



http://wwwdb.inf.tu-dresden.de/research/

