

# BUSINESS INTELLIGENCE

with Jackie



Jackie Jansen

## Data Skew: What is it and Why do I care?

Last month we discussed the various indexing technologies on the iSeries. This month we are going to continue on and talk about data skew. What is it, do I have it and why do I care?

If you recall, the database optimizer is the component of the database that decides how to actually execute your query. The optimizer decides whether or not the database will read your entire file (full table scan), use one of the indexes that you have built, or possibly build a temporary index. When your query is going to return only a few records from a large file, then an index is obviously the way to go. If your query is going to return a large number of records, an index may or may not be cost justified. When the database uses an index, it has to retrieve both the index pages and then the data pages from disk. If you are retrieving a large number of records, a full table scan could easily be the fastest retrieval method. By taking advantage of large blocking factors and no index I/O a full table scan may perform much faster than using an index.

With most commercial databases the optimizer assumes that your data is evenly spread out over the various keys. Let me give you a good example. In Canada we have 10 provinces and 3 territories. If we were indexing by province then most optimizers would base their decision on the assumption that each province must have 1/13 of the population. Now you and I know, that most national companies will have a few more customers in Ontario than they have in Prince Edward Island or Nunavut. If we were to query a large file by province then we would like the optimizer to be able to make a different decision on how to retrieve the data

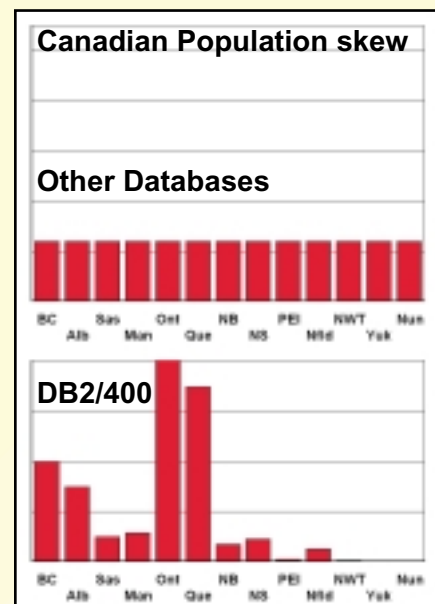
when an end user asked for customers in Ontario and Quebec as opposed to customers in Prince Edward Island. Typically databases will keep track of the number of unique key values, in our case 13, but they don't track the actual number of occurrences for each unique key value. The really good news for all of us in iSeries / AS/400 land is that, built into our system, is an optimizer that handles this for us. If this is the case then you might ask, why are you bothering to read this article. Well for a couple of reasons. First, it is always good to understand how the system works and secondly there are things that you can do to enable the optimizer to make intelligent choices. The iSeries optimizer actually calculates the data skew through indexes that have been built over the field we are interested in eg: province.

Remember our friend, the radix index from last month. If you have a radix or binary tree index over province, the optimizer will check a subset of the

nodes in the tree structure to estimate the number of records that a query will return. With EVIs or Encoded Vector Indexes, the actual count of each unique value is stored in the index. This means that with DB2/400 data skew issues are handled by the existence of indexes.

Now you see why reading this column was important. Although the optimizer will evaluate the data skew, you need to first help it by creating the appropriate indexes for the optimizer to query. This procedure holds true even when your query is based on multiple conditions. For example, if you often asked for a list of all the unmarried women in Quebec you would consider creating a radix index over province, sex and marital status. Using EVIs you would create an index over each of the 3 columns.

Constantly reevaluating the retrieval method would be expensive. The optimizer will reevaluate its retrieval method based on multiple conditions including changes to the SQL statement and changes in the size of the data file. This topic alone could fill a future column. By now you should also know "What Data Skew is and Why you care". My thanks to Amy Anderson from our iSeries Teraplex Centre for her input and insights. I willingly admit to stealing her examples and Canadianizing them. [T](#) [G](#)



Jackie Jansen is a Certified Consulting IT Specialist. She currently works in the IBM Americas Business Intelligence Solutions Centre. Jackie is a frequent speaker at AS/400 Technical Conferences and User Group meetings. Contact her at [jjansen@ca.ibm.com](mailto:jjansen@ca.ibm.com).