

Cleaning up Very Large Databases and Keeping Them Clean

Priscilla Broberg

U.S. Caden (a division of Manpower)
currently working at Agilent Technologies

Executive Summary

This presentation shows a real-world example of how a very large Customer database was cleansed and de-duplicated to shrink it down to a manageable size. The techniques used to do this are shown, as well as the processes that were implemented to maintain the new level of data cleanliness. The tricks and techniques are applicable to customer files or databases of any size in any business. Actual before and after data examples are shown.

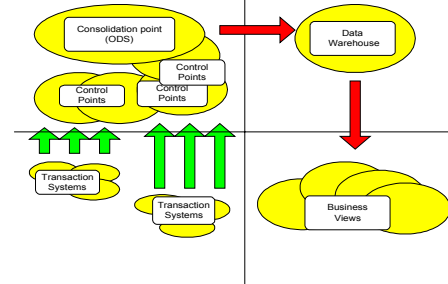
Topics covered include:

-
- Typical customer data flows, from data entry to reporting
- Proper placement of data cleansing and merging in the data flow
- Techniques to maximize effectiveness of merge/purge (de-duplication)
- Ideas for maintaining a higher level of data cleanliness, and minimizing data duplication.

Cleaning up Very Large Databases and Keeping Them Clean

The story of how a customer database got very large and very messy, then got small and clean again.

Data Warehousing Architecture



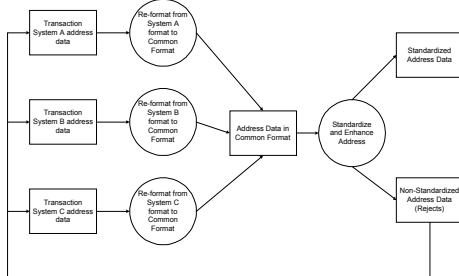
The Consolidation Point Provides Clean, De-Dupped Data to the Warehouse

- Cleanses data
- Standardizes data
- Enhances data (e.g. zip+4)
- Eliminates duplicates (merge/purge)
- Communicates back to transaction systems
 - rejected transactions
 - successfully loaded transactions

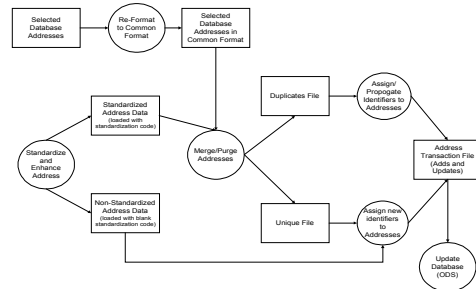
WHY DO WE NEED TO MERGE/PURGE CUSTOMER DATA?

- Data from separate transaction systems is entered and identified differently
- Need for company-wide view of customers ("Master list")
- Need to consolidate customer information Worldwide
 - avoid double counting
 - save on database storage
 - able to identify one customer with one unique identifier (cross-referenced to source systems)

Standardization Process Flow



MERGE/PURGE PROCESS FLOW



How Did We Get Into This Mess?

- ODS Database designed in late 1980's to cleanse and load a single type of customer data - Order Processing Customers. Data only went to one application for reporting. ALL records were required to be loaded, regardless of data quality!
- Later, additional sources of data, as well as receiving applications were added for Direct Marketing. These were allowed to be rejected, if they did not meet data quality standards.
- Merge/Purge rules changed.
- Moved from Mainframe to Unix platform, and changed cleansing and merge/purge tools.
- ODS Database had no delete capability. All data was added or updated, then remained there forever!
- Only incoming transactions were cleansed and merge/purged against the database.
- Once data was loaded, it was never re-cleansed or re-merge/purged.

Other Contributing Factors

- Many records were coded with the wrong country code. Only those with US country codes (US, and affiliates such as Puerto Rico, Guam, etc) went through standardization!
- We have no edit for verifying the country code against the address. We just accepted what was input.
- Once a record is loaded as non-standard, it NEVER participates in the merge/purge.
- Non-standardized records contributed to a lot of data duplication.
- Split from Hewlett-Packard caused us to inherit a database full of HP customers as well as Agilent customers. There was no attribute of the customer data to tell them apart.

Preparing for The BIG CLEAN-UP

Step #1: Pre Clean-up

- Removed data associated with Direct Marketing
 - Identified by Source Number
 - Had to make sure data was not also associated with active sources.
- Documented current Merge/Purge rules and reviewed with users
- Using a Marketing Reporting Tool (the Data Warehouse recipient of our customer data), we were able to identify which customers belonged to Agilent by reporting customer numbers on orders with Agilent product lines.
- Identified customers who had been active in the past two years, and deleted all others.
- Number of site (address records) after clean-up went from approximately 11 million rows to 1.1 million rows.
- This became the starting point for our re-standardization and re-merge/purge.

Preparing for the BIG CLEAN-UP

Step #2: Analyze remaining data

- Determine how much data is US, how much Canada, and how much non-US.
- Country code not reliable. However, we used ACE to discover this, and locate the incorrectly coded records!
- Perform test merge/purge runs on non-US/non-Canada data, using line1, line2, etc. method.
- Adjust merge/purge parameters based on results of test runs.

Clean-up Steps

Data-Cleansing

- Country Code clean-up must be done first. Since this is part of the match-key, re-calculate match-key.
- Re-standardize US and Canada. After re-standardization, re-calculate match key again. (Postal code is also a component of match-key)
- Update database with new country codes and match-keys, as well as newly-standardized addresses.
- Our match-key algorithm: First letter of Business Name, followed by first four numbers of address, followed by first 3 bytes of postal code, followed by 3-byte country code. '@' used as filler where no data exists.
- Example: IBM 123 Main Street, Anytown, Anystate, 99999 would be coded as: !123@999000 ('000' is our country code of US).

Example -Before Standardization

3009 NW 75 AVENUE	MIAMI FLORIDA 33122	ATTN: LILIANA P. VELAZQUEZ 351	000000333351
7200 NW 7 STREET 2ND FLOOR	MIAMI FLORIDA 33126	351	000000333351
C/O FLEMING 105255	7051 NW 37 STREET	MIAMI, FL. 33156-6559	301 1315000001
14413 IMPORT DR.	LARDO, TX 78041 USA	201	07804144201
C/O MIAMI PANALPINA INT	3505 N.W. 107TH AVE.	MIAMI, FL 33178	355 C313733355
2100 BLUE LAGOON DRIVE SUITE 1050	MIAMI FLORIDA 33126	355	00000105255
10777 WESTHEIMER STE 625	HOUSTON TX 77042	351	00000107351
1900 CONCOURSE DRIVE	SAN JOSE, CA 95131	223	00000109223
420-B FARMERICHSON DRIVE	EL PASO, TX 79907	201	07910709201
2429 TERMINAL BLVD.	MOUNTAIN VIEW, CA 94043	357	00000262357
400 REIMANN AVENUE	SANDWICH IL 60548-0900	412	C0554000412
1310 MEMOREX DRIVE	SANTA CLARA, CA 95050	583	00000133583
SUITE 118	990 RICHARD AVENUE	SANTA CLARA CA 95050	583 0050505083
2712 EAST MERALOMA AVE	690 MAIN STREET	STRATFORD, CT 06615-0129	489 05500000489
ACCOUNTS PAYABLE	ANARHEIM, CA 92806	549	100000273549
1101 CYPRESS CREEK ROAD	2712 E MERALOMA AVENUE	ANARHEIM CA 92806	549 1273273549
6780-R SIERRA COURT	CEDAR PARK TX 78613	405	07861278405
ATTN: TIMOTHY BOB SMITH	DUBLIN, CA 94568	427	00000945427
ACCOUNTS PAYABLE	DUBLIN CA 94568	427	00000945427
	13105 E. ALGONQUIN ROAD	SCHAUMBURG IL 60196	428 06105130428
	1301 E ALGONQUIN ROAD	SCHAUMBURG IL 60196	428 060105130428

Clean-up Steps Eliminating Duplicate Data

- If required to maintain record of eliminated data, use the dup groups to create elimination transactions. An elimination transaction is basically like a "change of address" transaction. All that is needed is the old address identifier and the new (surviving) address identifier.
- If this is NOT required, delete any addresses not in your "mail" file, and you are done!

Clean-up Steps Eliminating Duplicate Data (cont)

- Steps to performing eliminations:
 - 1) Create elimination transactions from dup groups
 - 2) Apply eliminations to all tables in which the address identifier is used. For example, our database uses this identifier in X tables. Change old identifier to new identifier, based on transaction.
 - 3) Once all tables have been updated, create a row in an "elimination table" to keep track of this change (old ID --> new ID)
 - 4) Finally, delete old (eliminated) address record

Sample Elimination Transactions (Created from Sample Dups File)

Old ID	New ID
017588008	017508378
018679713	018660122
014707877	014268700
017626543	017037130
017784498	017037130
017818072	017037130
017907210	017037130
007083284	003488990
017023274	010279043
017095742	010279043
017819768	010279043
018408512	018180551
007083034	002223280
015338690	002223280
017140834	017038804

Lessons Learned

- It is important to understand your current data flow and processing. If you haven't documented it thoroughly, start now!
- Make sure your users understand the data-cleansing and merge/purge rules. They own these!
- Know what data you have control over, and what data you do not. For example, we can clean-up data, but we cannot force the source systems to send us clean data.
- For best results, re-standardize all addresses in database whenever you get a new zip+4 update file from Firstlogic.
- Re-merge/purge entire database at least 4 times a year.
- Use Firstlogic tools to analyze your data, as well as to cleanse it in production.
- Don't assume you cannot merge/purge non-US data. It can be done quite effectively using the user-definable fields (Merg_Purg1, Merg_Purg2, etc).
- Read the Firstlogic Software Update Bulletins and Customer Care Bulletins that come with your upgrades. There may be new features you can take advantage of!

Improvements/Benefits

- Reduced address rows in database from 11 million to < 2 million
 - Benefits:
 - Less disk space usage
 - Easier database administration
 - Faster processing times, as data merge/purges against fewer rows
 - Improved data quality, as duplicates are eliminated
 - Better decision making, as user confidence in data improved
 - Improved processing times on downstream systems, as less data is passed to them

Cost Savings

- Support went from 3 full-time programmers rotating on-call duty (24/7), to Call-center, with 1 on-call "deep support" programmer.
- Call-center support much less expensive
- Support programmers became available to work on new projects.
- Went from one full-time DBA to one part-time DBA.
- Lowered disk space costs
- Lowered processing (machine time) costs
- Estimated total annual savings: \$500,000