

Paper CC-04

Managing Datasets at Library Level via Dynamically Constructed, DICTIONARY Tables-Driven SAS Code

Jingxian Zhang; Quintiles; Overland Park, KS, USA

Abstract

The SAS System read-only DICTIONARY tables contain valuable information about metadata such as SAS libraries, table names, column names, etc. Integrating the information into SAS programming helps create dynamic and efficient scripts to manage datasets. In this paper, I will first discuss how the DICTIONARY tables are accessed and what information is available to SAS users. Then I will show how to dynamically construct DICTIONARY tables-driven code via three approaches - SQL select into macro-variable method, call execute method, and generate-and-include an external file method. The three macros developed by the author using the above approaches manage all datasets at the library level - capitalize all character data, dump all data into an excel file with each datasets being its own tab in excel, and query all character data with certain length. Using the basic techniques present in these macros, SAS programmers can develop their own dynamic scripts to accomplish other tasks.

Introduction

SAS programmers in Data Management of pharmaceutical industry program and format CDMS (for example, Phase Forward Inform, Medidata Rave, Oracle Clinical) database data and external vendor data such as laboratory test results and EKG data to STDM-like datasets per CDISC or sponsors' specification. A clinical trial or study typically has 20 to 30 SAS datasets. During data validation and review processes, SAS programmers often face the challenging questions from the data reviewers: Can you provide a list of all character variables that have data length ≥ 190 ? Can you change all character data to upper case? Can you present the data in excel format instead of datasets for the ease of data review? If you had only one or two small datasets, you might accomplish these tasks by updating each individual program used to generate these datasets. But if you have 30 datasets with hundreds of variables, are you going to use the above, tedious approach?

The purpose of this paper is to find a better solution that can accomplish the above tasks. I will first discuss how SAS DICTIONARY tables and their corresponding SASHELP views are accessed and what information is available to SAS users. Then I will show how to dynamically construct SAS DICTIONARY tables-driven code via three approaches – SQL select into macro-variable method, call execute method, and generate-and-include an external file method. The code generated manages datasets at the library level.

Abbreviations: CDISC, Clinical Data Interchange Standards Consortium, a non-profit group that defines clinical data standards for the pharmaceutical industry; CDMS, Clinical Data Management System, a tool used in clinical research to manage the data of a clinical trial. The clinical trial data gathered at the investigator site in the case report form are stored in the CDMS. EKG, Electrocardiogram (e-lek-tro-KAR-de-o-gram), also called ECG, a simple, painless test that records the heart's electrical activity; STDM, Study Data Tabulation Model, defines a standardized structure for data tabulations that are to be sent to the FDA as part of a regulatory submission.

SAS DICTIONARY Library

The SAS System read-only DICTIONARY tables and corresponding SASHELP views provide valuable information about SAS libraries, datasets, columns and attributes, catalogs, indexes, macros, system options, titles, views, and more. This information has been used to verify dataset structures and attributes (Teng and Wang, 2006) and data inconsistencies in multiple datasets (Murphy, 2011).

DICTIONARY is a special case of a SAS library. DICTIONARY tables are special read-only PROC SQL tables or views. DICTIONARY tables are accessed only by PROC SQL. Based on the DICTIONARY tables, SAS provides PROC SQL views that can be used in other SAS procedures and in the DATA step. These views are stored in the SASHELP library and are commonly called SASHELP views.

To see what tables the DICTIONARY has, one can use the following code, which, if ran at SAS 9.2, would return a list of 29 tables (Table 1). Table names are meaningful so that you know what they are about.

```
PROC SQL;
  SELECT UNIQUE memname FROM DICTIONARY.dictionaries;
QUIT;
```

Table 1 DICTIONARY Tables and Their Corresponding SASHELP Views (SAS 9.2)

<i>Dictionary Table</i>	<i>SASHELP Views</i>	<i>Dictionary Table</i>	<i>SASHELP Views</i>
CATALOGS	VCATALG	INFOMAPS	VINFOMP
CHECK_CONSTRAINTS	VCHKCON	LIBNAMES	VLIBNAM
COLUMNS	VCOLUMN	MACROS	VMACRO
CONSTRAINT_COLUMN_USAGE	VCNCOLU	MEMBERS	VMEMBER
CONSTRAINT_TABLE_USAGE	VCNTABU	OPTIONS	VOPTION
DATAITEMS	VDATAIT	PROMPTS	VPROMPT
DESTINATIONS	VDEST	PROMPTSXML	VPRMXML
DICTIONARIES	VDCTNRY	REFERENTIAL_CONSTRAINTS	VREFCON
ENGINES	VENGINE	REMEMBER	VREMEMB
EXTFILES	VEXTFL	STYLES	VSTYLE
FILTERS	VFILTER	TABLES	VTABLE
FORMATS	VFORMAT	TABLE_CONSTRAINTS	VTABCON
FUNCTIONS	VFUNC	TITLES	VTITLE
GOPTIONS	VGOPT	VIEWS	VVIEW
INDEXES	VINDEX		

To see what views SASHELP views have, one can use the following syntax:

```
PROC SQL;
  SELECT DISTINCT memname
  FROM SASHELP.vsvview WHERE libname='SASHELP'AND memname like 'V%';
QUIT;
```

To see how each DICTONARY table or SASHELP view is defined, submit a DESCRIBE TABLE statement. For example, the following code would show the definition of DICTONARY.columns. The results are written to the SAS log, which has the variable names and their attribute information (Table 2).

```
PROC SQL;
  DESCRIBE TABLE DICTONARY.columns;
QUIT;
```

Table 2 Definition of DICTONARY.COLUMNS (SAS 9.2)

```
create table DICTONARY.COLUMNS(
  libname char(8) label='Library Name',
  memname char(32) label='Member Name',
  memtype char(8) label='Member Type',
  name char(32) label='Column Name',
  type char(4) label='Column Type',
  length num label='Column Length',
  npos num label='Column Position',
  varnum num label='Column Number in Table',
  label char(256) label='Column Label',
  format char(49) label='Column Format',
  informat char(49) label='Column Informat',
  idxusage char(9) label='Column Index Type',
  sortedby num label='Order in Key Sequence',
  xtype char(12) label='Extended Type',
  notnull char(3) label='Not NULL?',
  precision num label='Precision',
  scale num label='Scale',
  transcode char(3) label='Transcoded?');
```

To examine the contents of each view, one can use the following syntax:

```
PROC CONTENTS DATA = SASHELP.view_name;
RUN;
```

For more information about DICTONARY information, readers are referred to Eberhardt and Brill (2006), Thornton (2011) and CodeCrafters, Inc. (2010).

SQL Select Into Macro-Variable Approach

It is pretty easy to check data length via length function. Basic syntax of checking if a variable has data length ≥ 100 could be like: < Data step approach: DATA test; SET dataset_name; IF LENGTH(variable) ≥ 100 ; RUN;> or <SQL approach: PROC SQL; SELECT variable FROM dataset_name WHERE LENGTH (variable) ≥ 100 ; QUIT; >. These two approaches work if there are only a few variables to worry. However, it would be tedious or impossible to implement

if one needs to deal with hundreds of variables. Using the `DICTIONARY.columns` and dynamic programming techniques, I have developed the following macro `%query_length` to accomplish such tasks.

```
%MACRO query_length(mylib=, maxlen=);
  PROC SQL;
    /**** Section I ****/
  CREATE TABLE columns as
  SELECT memname, name, length, type
  FROM DICTIONARY.columns
  WHERE libname = "%upcase(&mylib)" and
         memtype="DATA" and
         type="char";

  %LET cnt = &sqlobs;

  /**** Section II ****/
  SELECT "SELECT UNIQUE " || STRIP(memname) || " as ds_name,"
  || STRIP(name) || " as var_name, subjid, length(strip("
  || STRIP(name) || ")) as max_len, "
  || STRIP(name) || " as datavalue FROM &mylib.."
  || STRIP(memname) || "
  WHERE length(strip(" || STRIP(name) || ")) >= &maxlen.; "
  INTO :select1-:select&cnt
  FROM columns;

  /**** Section III ****/
  CREATE TABLE tempds
  (DSNAME          char(10),
  VARNAME          char(10),
  SUBJID           char(20),
  DATALENGTH     num,
  DATAVALUE      char(200));

  %DO i=1 %TO &cnt;
    INSERT INTO tempds
    &&select&i;
  %END;
  QUIT;
%MEND query_length;
```

Note: For this specific macro, the variable `subjid`(subject ID) must exist in all datasets at mylib library because I want to know which `subjid` has the queried results.

The macro `%query_length` queries all character variables that have data length \geq a specified number. The `DICTIONARY.columns` table used in section I shows dataset information at the variable level (Table 2). Section I in the macro is to get the `memname` (dataset name), name

(variable name), length (variable defined length) and type (variable data type) on all character variables at a given library. Using the returned information from section I, section II dynamically constructs SQL query code which will be used to populate a dataset tempds in section III. The macro call gives output for the SELECT statement for each variable that has data length \geq &maxlen. Table 3 below shows what has been returned on the test data I used.

Table 3 Macro QUERY_LENGTH call with &maxlen \geq 190

Dataset	Variable	Patient	Data	Length
EG	EGORRES	99994-006	ABNORMAL, BUT NOT CRITICAL. NOT DONE PREDOSE NOR BEFORE LAB DRAW. STAFF FORGOT TO HOLD AM DOSE AND PT UNCOOPERATIVE WITH INITIAL ATTEMPT TO COMPLETE EKG. ALL OTHER PROCEDURES COMPLETED THEN REA	194
FA	FACOM	99993-007	THIS IS NOT DONE DUE TO SAMPLE PROBLEM. IMPRESSION: STATUS POST VERTEBRAL BODY AUGMENTATION OF T011 AND T012. STABLE COMPLETE COLLAPSE OF THE T10 VERTEBRAL BODY. INTERVAL LOSS OF VERTEBRAL LAB ONE	197
LB	LBNAM	99993-010	THIS IS MIDWEST SAS LAB TWO FOR TEST DATA. THIS IS THE TEST LAB CENTER USED TO CREATE THE TEST DATA USED FOR THIS PRESENTATION. DUE TO REGULATIONS, THIS LAB MAY OR MAY NOT BE USED FOR PATIENTS. HOW,	199

The key code of the macro is the SELECT INTO statement that has the following syntax:

```
SELECT Column(s)
  INTO :<Macro Variable name1> -:<Macro Variable name9999>
  FROM Table-name | View-Name;
```

The SELECT statement stores returned row values in a list of user-defined macro variables. Only the required number of macro variables will be created. A number large enough to hold the number of observations returned from the SELECT statement must be specified if system macro variable SQLOBS is not used.

The syntax of string concatenation needs a little explanation. The following statement:

```
<"SELECT UNIQUE "" || STRIP(memname) || "" as ds_name,"
  || STRIP(name) || "" as var_name, subjid, length(strip("
  || STRIP(name) || ")) as max_len, "
  || STRIP(name) || "" as datavalue
  FROM &mylib.."|| STRIP(memname) || "
  WHERE length(strip("|| STRIP(name) || "))  $\geq$  &maxlen.;">
```

will be dynamically decoded to the following code after execution if memname = AE and name = STUDYID and &maxlen = 100:

```
SELECT UNIQUE 'AE' as ds_name, 'STUDYID' as var_name, subjid,
      length(strip(STUDYID)) as max_len, STUDYID as datavalue
FROM LIBREF.AE WHERE length(strip(STUDYID)) >= 100;
```

Similar SELECT statements for other variables will be dynamically generated and are kept under macro variable `selecti` (where *i* is between 1 and SQLOBS. For my test data, SQLOBS = 639). This is accomplished via dynamic variables used in the macro. Basic dynamic variables `<" || VARIABLE || ">` and `<" || VARIABLE || ">` are commonly used in dynamic sql code generation and they are different. When 'variable' is needed in the generated script, `" || VARIABLE || "` should be used. When VARIABLE is needed in the generated script, `" || VARIABLE || "` should be used. For example, if memname is AE and you want 'AE' to show up in the script, then you need to use `" || memname || "`; However, if you want AE to show up as a variable in the script, then you need to use `" || memname || "`.

Generate-and-Include an External File Approach

After I have programmed 23 datasets, I got a request to have all character data to be capitalized. Instead of modifying each individual program generating its corresponding datasets, we can use the dictionary metadata to do the capitalization. Here I have used a different approach. Using the DICTIONARY metadata information, the macro generates SAS code via put statement and then writes them to an external file. To execute the macro `%upcase_ds`, one can run it and then use the `%includes` statement to include the file (i.e., `upcase_chardata.sas`) generated from the macro call.

```
%MACRO upcase_ds(mylib=, _filepath=);
  PROC SQL;
    CREATE TABLE columns as
      SELECT memname as member, name as variable
      FROM dictionary.columns
      WHERE libname = "%upcase(&mylib)" and memtype="DATA" and type="char"
      ORDER by member, variable;
  QUIT;

  DATA _null_;
    SET work.columns end=EOF;
    BY member variable;
    FILE "&_filepath.\upcase_chardata.sas";
    dlm = byte(9);
    IF _n_ =1 then do;
      PUT 'PROC SQL;';
    END;
    PUT dlm +(-1) 'UPDATE ' "&mylib.." member ' SET ' variable '
      = UPCASE(' variable +(-1) ');';
    if EOF then do; put 'QUIT; ';
    END;
  RUN;
%MEND upcase_ds;
```

The following is part of the UPDATE statements contained in the dynamically generated file upcase_chardata.sas when the macro is called (AE is one of the datasets present in the mylib LIBREF). For my test data, I have 23 SAS datasets and 639 char variables. That means 639 update statements will be generated if the macro is called. You can see how powerful this technique is.

```
PROC SQL;
  UPDATE LIBREF.AE SET AEACN = UPCASE(AEACN);
  UPDATE LIBREF.SET AEACNOTH = UPCASE(AEACNOTH);
  UPDATE LIBREF.SET AEACNX = UPCASE(AEACNX);
  << Total 639 UPDATE statements will be generated for the test data used>>
QUIT;
```

In case readers wonder how this task can be done using SQL select into macro-variable method, here is the code:

```
%MACRO upcase_chardata(mylib);
PROC SQL;
  CREATE TABLE columns as
  SELECT memname, name, length, type
  FROM DICTIONARY.COLUMNS
  WHERE libname = "%upcase(&mylib)" and memtype="DATA" and type="char";
  %let cnt = &sqllobs;

  SELECT "update &mylib.||strip(memname)||" set "||name||" = upcase("||strip(name)||");"
  INTO :update1-:update&cnt
  FROM columns;

  %DO i=1 %TO &cnt;
    &&update&i;
  %END;

QUIT;
%MEND upcase_chardata;
```

Call Execute Approach

When exporting multiple SAS datasets to Excel files, the traditional method is to write multiple steps as below.

```
PROC EXPORT DATA = LIBREF.ae DBMS =excel2000
  OUTPUT = "_outpath\ae.xls" REPLACE; SHEET = "AE";
RUN;
PROC EXPORT DATA = LIBREF.cm DBMS =excel2000
  OUTPUT = "_outpath\cm.xls" REPLACE; SHEET = "CM";
RUN;
```

Consider the scenario that one wants to dump all the datasets in the library to one excel file with each tab corresponding to a dataset. Instead of writing this block of code many times for exporting each datasets, I have developed a sas_to_excel macro (mylib= , _outpath = , _project =) by using SASHELP view VTABLE and CALL EXECUTE routine, where mylib is libref, _outpath specifies output folder and _project is the excel output file name.

```
%MACRO sas_to_excel(mylib= , _outpath = , _project =);  
  %Macro printds(libname,dsname) ;  
    PROC EXPORT DATA = &libname.&dsname DBMS = excel2000  
      OUTFILE = "&_outpath.\&_project.xls" REPLACE; SHEET = "&dsname";  
    RUN;  
  %MEND printds;  
  
  DATA _null_ ;  
  SET sashelp.vtable ;  
  WHERE libname = "%upcase(&mylib)";  
  CALL EXECUTE('%printds('||strip(libname)||','||strip(memname)||')' ) ;  
  RUN;  
%MEND sas_to_excel;
```

Note that CALL EXECUTE is used within a data step and has the following syntax: CALL EXECUTE (argument). The macro uses the %printds as the argument for CALL EXECUTE and gets LIBREF and MEMNAME from SASHELP.VTABLE. Any DICTIONARY tables that have the LIBREF and MEMNAME can be used here to replace VTABLE. For other uses of CALL EXECUTE, readers are referred to Ruelle and Moses (2006) and Michel(2005).

Conclusion

The author used three case study macros to demonstrate how SAS code is dynamically built based on the DICTIONARY metadata. The techniques presented maximize procedural efficiency and lighten the workload of routine processing. Based on their expertise, SAS programmers may select a proper method and develop their own dynamic scripts to accomplish programming tasks.

References

CodeCrafters, Inc., 2010. Summary of SAS DICTIONARY Tables and Views.
<http://www.codecraftersinc.com/pdf/DICTIONARYTablesRefCard.pdf>

Eberhardt, P. and Brill, I. (2006) How Do I Look it Up If I Cannot Spell It: An Introduction to SAS® Dictionary Tables. SUGI 31 Proceedings 2006, San Francisco, California

Michel, D. 2005. CALL EXECUTE: A Powerful Data Management Tool. SUGI 30 Proceedings, 2005, Philadelphia, Pennsylvania

Murphy, W.C. 2011. Who Do You Have? Where Are They? SAS Global Forum 2011, Las Vegas, Nevada.

Ruelle, A. and Moses, K. (2006). CALL EXECUTE: A Primer. PharmaSUG 2006, Bonita Springs, Florida.

Teng, C and Wang, W. 2006. Simple Ways to Use PROC SQL and SAS DICTIONARY TABLES to Verify Data Structure of the Electronic Submission Data Sets. PharmaSUG 2006, Bonita Springs, Florida.

Thornton, P. 2011. SAS® DICTIONARY: Step by Step. SAS Global Forum 2011, Las Vegas, Nevada.

Acknowledgements

I would like to thank David Corliss, Coders Corner Section Chair, for accepting my abstract and paper and Jill Jenia for reviewing this paper and providing her helpful remarks.

Contact information

Your comments and questions are valued and encouraged. Contact the author at:

Jingxian Zhang
Quintiles, Inc.
6700 W. 115th St.
Overland Park, KS 66223
Phone: 913-708-6674
Fax: 913-871-9569
Email: jing.zhang@quintiles.com

SAS and all other SAS Institute Inc. product or service names are registered trademarks or trademarks of SAS Institute Inc. in the USA and other countries. ® indicates USA registration. Other brand and product names are registered trademarks or trademarks of their respective companies.