

SAS® GLOBAL EORUN $\cap \cap \cap$

MARCH 29 - APRIL 1 WASHINGTON, DC

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 trademarks of their respective companies.



TAP TO GO ΒΑϹΚ ΤΟ KIOSK MENU





Abstract

Introduction SAS Macro Results Discussion

Please use the headings above to navigate through the different sections of the poster

In the DATA step, merging data sets with common variables that are not included as BY variables can yield undesirable results. Specifically, the value of a common variable can be overwritten with an incorrect value. To prevent this from happening, you must ensure that the variable is read from only one "master" data set, by either dropping or renaming the variable in the other data sets. When working with data sets with just a few variables, you can quickly check which variables appear in more than one data set. However, as the number of data sets and variables increases, the chance of missing a common variable also increases. The SAS[®] macro CHECK VAR EXIST was written to identify variables that exist in more than one data set more efficiently and accurately. The macro prints all common variables, which data sets they appear in, and other pertinent information. You can then use the list to drop or rename variables where they are not relevant, thereby reducing the chance of unintentionally overwriting a large number of variables.

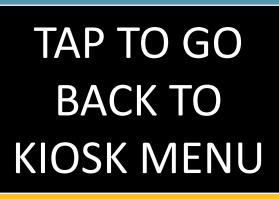
DON'T OVERWRITE ME! A SAS® MACRO TO IDENTIFY VARIABLES THAT EXIST IN MORE THAN ONE DATA SET

Andrea Barbo

Yale New Haven Health Services Corporation/Center for Outcomes Research and Evaluation (CORE)

Abstract:







Please use the headings above to navigate through the different sections of the poster

Q SAS programmers are commonly taught that when you merge datasets in the DATA step, variables in the dataset listed later on the MERGE statement replace the values of variables that also exist in a previously listed dataset.

- BY variables.
- dataset.
- involved.
- variables.

DON'T OVERWRITE ME! A SAS® MACRO TO IDENTIFY VARIABLES THAT EXIST IN MORE THAN ONE DATA SET

Andrea Barbo

Yale New Haven Health Services Corporation/Center for Outcomes Research and Evaluation (CORE)

Introduction:

This may be true for one-to-one merging, but not for one-to-many merging, because of how the Program Data Vector works. • As such, you need to be careful when combining multiple datasets that have variables in common, and not all of them are included as

The best way to avoid seeing unexpected results is to drop or rename common variables so that they only show up in one

• Figuring out the common variables can be done easily if you're working with just a couple of datasets with few variables. However, it gets more cumbersome the more datasets and variables are

The SAS[®] macro CHECK_VAR_EXIST, which will be described in the next slides, provides an automated way of identifying common

	Measure_Name	Measure_ID	Compa 🔺	•	Measure_Name	Measure_ID	1
1	Central Line Associated Bloodstream Infection (ICU + select Wards): Lower	HAI_1_CILOWER	No Diff		e of complications for hip/knee acement patients	COMP_HIP_KNEE	
	Confidence Limit				th rate for heart attack patients	MORT_30_AMI	
2	Central Line Associated Bloodstream Infection (ICU + select Wards): Upper Confidence Limit	HAI_1_CIUPPER	No Diff		th rate for CABG surgery patients	MORT_30_CABG	
		HAI_I_CIUFFER	NO DIN		th rate for COPD patients	MORT_30_COPD	
3	Central Line Associated Bloodstream Infection: Number of Device Days	HAL 1 DOPC	No Diff		th rate for heart failure patients	MORT_30_HF	
		INALI_DOIC			th rate for pneumonia patients	MORT_30_PN	
	Central Line Associated Bloodstream	HAI_1_ELIGCASES	No Diff		th rate for stroke patients	MORT_30_STK	
4	Infection (ICU + select Wards): Predicted Cases				toperative Acute Kidney Injury juiring Dialysis Rate	PSI_10_POST_KID	
1		>			toperative Respiratory Failure Rate	PSI_11_POST_RE	
					ous blood clots after surgery	PSI_12_POSTOP_	
		<u> </u>	- 11	BIO	od stream infection after surgery	PSI_13_POST_SE	
1	Acute Care Hospitals	Government - Hospital I Authority	12		wound that splits open after surgery on abdomen or pelvis	PSI_14_POSTOP_	
2	Acute Care Hospitals	Government - Hospital I Authority	13	Acc	cidental cuts and tears from medical	PSI_15_ACC_LAC	
3	Acute Care Hospitals	Government - Hospital I Authority	14	-	ssure sores	PSI_3_ULCER	
4	Acute Care Hospitals	Voluntary non-profit - Pr	15 Dea		aths among Patients with Serious	SI_4_SURG_COMP	



Please use the headings above to navigate through the different sections of the poster

SAS[®] Macro CHECK_VAR_EXIST:

- Identifies variables that exist in more than one dataset.
- □ Ideal to use before merging 2+ datasets as a check to prevent incorrect variables from overwriting correct ones with the same name.
- Input parameters: **DTA** is a list of datasets to check (preceded by a libref if stored as a permanent dataset), LINK VAR is a list of variables that should be excluded from the checking (usually the ones used as BY variables in the MERGE statement). • Output: list of variables that appear in more than one dataset,
 - with additional info like length & type, in the Results Window.

DON'T OVERWRITE ME! A SAS® MACRO TO IDENTIFY VARIABLES THAT EXIST IN MORE THAN ONE DATA SET

Andrea Barbo

Yale New Haven Health Services Corporation/Center for Outcomes Research and Evaluation (CORE)

quit;

%mend check_var_exist;

```
%macro check_var_exist(dta=,link_var=);
 data _null_;
   /*remove excess blank characters from list of datasets*/
   _var="&dta";
   dta_list=tranwrd(compbl(strip(_var)),". ",".");
   call symputx("dta_list",dta_list);
   /*count how many datasets to check for overlapping variables*/
   cnt_dta=count(strip(dta_list), " ")+1;
   call symputx("cnt_dta", cnt_dta);
   /*list of variables to exclude from checking*/
   list_var=lowcase("'"||tranwrd(compbl(strip("&link_var"))," ","',")||"'");
   call symputx("list_var",list_var);
 run;
 %put &dta_list &cnt_dta &list_var;
  /*output variables that exist in more than 1 dataset*/
 proc sql;
   select *
   from (select distinct upcase(name) as name label="Column Name",type,length,libname,memname
          from sashelp.vcolumn
      %if %sysfunc(find(%scan(%sysfunc(lowcase(&dta_list)),1,' '),.))>0 %then %do;
          where ( (lowcase(libname)="%scan(%scan(%sysfunc(lowcase(&dta_list)),1,' '),1,'.')" and
lowcase(memname)="%scan(%scan(%sysfunc(lowcase(&dta_list)),1,' '),2,'.')")
      %end;
      %else %do;
          where ( (lowcase(libname)="work" and lowcase(memname)="%scan(%sysfunc(lowcase(&dta_list)),1,' ')")
      %end;
          %do i=2 %to &cnt_dta;
           %if %sysfunc(find(%scan(%sysfunc(lowcase(&dta_list)),&i,' '),.))>0 %then %do;
            or (lowcase(libname)="%scan(%scan(%sysfunc(lowcase(&dta_list)),&i,' '),1,'.')" and
lowcase(memname)="%scan(%scan(%sysfunc(lowcase(&dta_list)),&i,' '),2,'.')")
           %end;
           %else %do;
           or (lowcase(libname)="work" and lowcase(memname)="%scan(%sysfunc(lowcase(&dta_list)),&i,' ')")
           %end;
          %end;
           and lowcase(name) not in (&list_var)
   group by name
   having count(*)>1
   order by name, libname, memname
```





Please use the headings above to navigate through the different sections of the poster

- Services (CMS).
- this from the check.
- %check_var_exist(dta =

DON'T OVERWRITE ME! A SAS® MACRO TO IDENTIFY VARIABLES THAT EXIST IN MORE THAN ONE DATA SET

Andrea Barbo

Yale New Haven Health Services Corporation/Center for Outcomes Research and Evaluation (CORE)

Results:

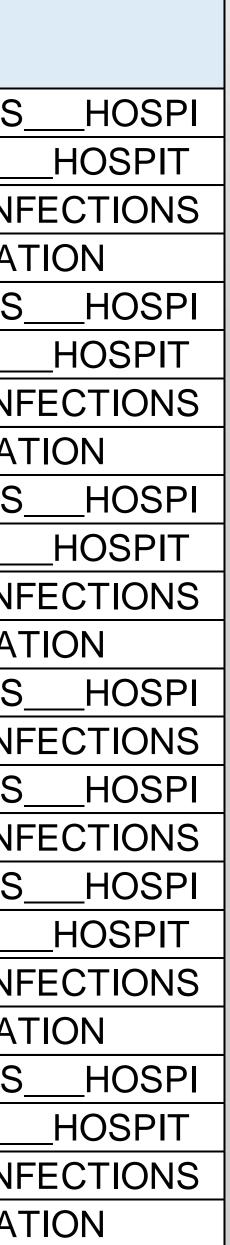
To illustrate how the macro can be used, we downloaded a few CSV files from Data.Medicare.gov and imported into SAS. **Data.Medicare.gov** is a website where consumers can freely download official healthcare-related data produced by the Centers for Medicare & Medicaid • We checked 5 datasets, 3 of which are temporary and 2 are permanent datasets, for common variables. As we're interested in merging all 5 datasets by the variable, Provider ID, we exclude

Hospital_general_information Fy_2019_ipps_fr_impact_file sasgf.Complications_and_deaths___hospi Healthcare_associated_infections sasgf.Patient_survey__hcahps____hospit , **link var** = Provider_ID)

Colun

ADDRE ADDRE ADDRE ADDRE HOSPIT HOSPIT HOSPIT HOSPIT LOCATI LOCATI LOCATI LOCATI MEASU MEASU MEASU MEASU STATE STATE STATE STATE ZIP_CO ZIP_CO ZIP_CO ZIP_COI

mn Name		Column	-	Member Name
	Туре	Length		
SS	char	51	SASGF	COMPLICATIONS_AND_DEATHS
ESS	char	50	SASGF	PATIENT_SURVEYHCAHPS
ESS	char	50	WORK	HEALTHCARE_ASSOCIATED_INF
ESS	char	50	WORK	HOSPITAL_GENERAL_INFORMA
TAL_NAME	char	71	SASGF	COMPLICATIONS_AND_DEATHS
TAL_NAME	char	71	SASGF	PATIENT_SURVEYHCAHPS
TAL_NAME	char	50	WORK	HEALTHCARE_ASSOCIATED_INF
TAL_NAME	char	50	WORK	HOSPITAL_GENERAL_INFORMA
ION	char	88	SASGF	COMPLICATIONS_AND_DEATHS
ION	char	88	SASGF	PATIENT_SURVEYHCAHPS
ION	char	86	WORK	HEALTHCARE_ASSOCIATED_INF
ION	char	89	WORK	HOSPITAL_GENERAL_INFORMA
JRE_ID	char	25	SASGF	COMPLICATIONS_AND_DEATHS
JRE_ID	char	15	WORK	HEALTHCARE_ASSOCIATED_IN
JRE_NAME	char	72	SASGF	COMPLICATIONS_AND_DEATHS
JRE_NAME	char	98	WORK	HEALTHCARE_ASSOCIATED_INF
	char	2	SASGF	COMPLICATIONS_AND_DEATHS
	char	2	SASGF	PATIENT_SURVEYHCAHPS
	char	2	WORK	HEALTHCARE_ASSOCIATED_INF
	char	2	WORK	HOSPITAL_GENERAL_INFORMA
DDE	num	8	SASGF	COMPLICATIONS_AND_DEATHS
DDE	num	8	SASGF	PATIENT_SURVEYHCAHPS
DDE	num	8	WORK	HEALTHCARE_ASSOCIATED_INF
DDE	num	8	WORK	HOSPITAL_GENERAL_INFORMA







Please use the headings above to navigate through the different sections of the poster

When variables exist in multiple datasets involved in a merge, and they're not listed as BY variables, you need to ensure they are read from a single "most correct" source, or there's a risk the incorrect value is saved.

- done.
- log after.

DON'T OVERWRITE ME! A SAS® MACRO TO IDENTIFY VARIABLES THAT EXIST IN MORE THAN ONE DATA SET

Andrea Barbo

Yale New Haven Health Services Corporation/Center for Outcomes Research and Evaluation (CORE)

Discussion:

□ The SAS macro CHECK_VAR_EXIST was written to aid programmers in identifying more efficiently which variables could be wrongly overwritten even before the merging is

• The output of the macro is used to determine where to include a DROP or KEEP statement. It can also be used to determine the maximum length for each common variable, which could be handy when concatenating datasets using the SET statement, to prevent the truncation of the variable. Another use is to determine if any of the common variables have different types (character vs numeric).

• A simpler but less efficient way to check for common variables is by using **OPTIONS MSGLEVEL=I**. Setting MSGLEVEL to I will make the log display additional notes pertaining to the merge processing. However, this requires you to run the DATA step merging first and then check the

```
124 data hospital_results;
125 merge Hospital_general_information Healthcare_associated_infections;
126 by Provider_ID;
127 run;
                  5334 observations read from the data set WORK.HOSPITAL_GENERAL_INFORMATION.
171288 observations read from the data set WORK.HEALTHCARE_ASSOCIATED_INFECTIONS.
      The data set WORK.HOSPITAL_RESULTS has 171864 observations and 36 variables
     : DATA statement used (Total process time):
                            0.42 seconds
       real time
       cou time
                            0.46 seconds
     options msglevel=i;
     data hospital results;
     merge Hospital_general_information Healthcare_associated_infections;
    by Provider_ID;
132
133
    run;
INFO: The variable Hospital_Name on data set WORK.HOSPITAL_GENERAL_INFORMATION will be
      overwritten by data set WORK.HEALTHCARE_ASSOCIATED_I\bar{\text{NFECTIONS}}
INFO: The variable Address on data set WORK.HOSPITAL_GENERAL_INFORMATION will be overwritten by
      data set WORK.HEALTHCARE_ASSOCIATED_INFECTIONS
INFO: The variable City on data set WORK.HOSPITAL_GENERAL_INFORMATION will be overwritten by data
      set WORK.HEALTHCARE_ASSOCIATED_INFECTIONS.
INFO: The variable State on data set WORK.HOSPITAL_GENERAL_INFORMATION will be overwritten by
      data set WORK.HEALTHCARE_ASSOCIATED_INFECTIONS
INFO: The variable ZIP_Code on data set WORK.HOSPITAL_GENERAL_INFORMATION will be overwritten by
      data set WORK.HEALTHCARE_ASSOCIATED_INFECTIONS
INFO: The variable County_Name on data set WORK.HOSPITAL_GENERAL_INFORMATION will be overwritten
      by data set WORK.HEALTHCARE_ASSOCIATED_INFECTIONS.
INFO: The variable Phone_Number on data set WORK.HOSPITAL_GENERAL_INFORMATION will be overwritten
      by data set WORK.HEALTHCARE_ASSOCIATED_INFECTIONS.
INFO: The variable Location on data set WORK.HOSPITAL_GENERAL_INFORMATION will be overwritten by
      data set WORK.HEALTHCARE_ASSOCIATED_INFECTIONS
NOTE: There were 5334 observations read from the data set WORK.HOSPITAL_GENERAL_INFORMATION
NOTE: There were 171288 observations read from the data set WORK.HEALTHCARE_ASSOCIATED_INFECTIONS.
NOTE: The data set WORK.HOSPITAL_RESULTS has 171864 observations and 36 variables.
NOTE: DATA statement used (Total process time):
                            0.60 seconds
       real time
                            0.61 seconds
       cpu time
```





SAS® GLOBAL FORUM 2020

USERS PROGRAM

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 trademarks of their respective companies.

