# Calculating Cardinality Ratio in Two Steps

Ronald J. Fehd, Stakana Analytics

**Abstract**

**Description:** The cardinality of a set is the number of elements in the set. The cardinality of a SAS® software data set is the number of observations of the data set, n-obs. The cardinality of a variable in a data set is the number of distinct values (levels) of the variable, n-levels. The cardinality ratio of a variable is n-levels / n-obs; the range of this value is from zero to one.

Previous algorithms combined output data sets from the frequency and contents procedures in a data step. This algorithm reduces multiple frequency procedure steps to a single call, and uses scl functions to fetch contents information in the second data step.

The output data set, a dimension table of the list of data set variable names, has variable cr-type whose values are in (few, many, unique); this variable identifies the three main types of variables in a data set, few is discrete, many is continuous, and unique is a row-identifier.

**Purpose:** The purpose of this paper is to provide a general-purpose program, ml-namex.sas, which provides enhanced information about the variables in a data set. The author uses this list-processing program in Fast Data Review, Exploratory Data Analysis (EDA) and in Test-Driven Development (TDD), a discipline of Agile and Extreme Programming.

**Audience:** programmers, data managers, database administrators

**Keywords:** SAS component language (scl) functions: attrn (nobs, nvar), close, open, varfmt, varinfmt, varlabel, varlength, varnum, vartype; frequency procedure, nlevels option

# Introduction

**Overview**

The purpose of this paper is to show an optimized algorithm for the calculations of cardinality ratio and cardinality type which reduces the number of steps from *O(n-vars)* to *O(1)* This subroutine is designed to be used by later list processing programs which provide a frequency of all the discrete variables and a summary of each of the analysis variables.

See SmryEachVar programs on page 17.

This section contains these topics.

- cardinality
- database terminology
- goal overview
- previous algorithm
- issues
- new algorithm
- output

**cardinality**

This table shows the relationship of cardinality to cardinality-ratio and cardinality-type.

| cardinality: | number of elements of a set of an array, the dimension | *n-rows(table)* |
|---|---|---|
| levels: | number of values in a column | *n-levels(column)* |
| cardinality ratio (cr): | $\frac{n-levels(column)}{n-rows(table)}$ | range: $(0:1]$ |

|  | | cardinality-type | | |
|---|---|---|---|---|
|  | few |  | many | unique |
| values: 0 | ............. | mean(cr) | ........... | 1 |
|  | distinct | | continuous | |
|  | by \| class var | | analysis var | |

**Notes:** The mean of cardinality ratio, *mean(cr)*, is used to separate category *few* from *many*; this statistic is theoretical and has been validated as empirically true from the test suite of the `sashelp` library.

**database terminology**

The goal of this program is to produce a database *dimension table*, the essential columns of which are a *primary key*, the row number, and information in three sets: 1. member name, 2. cardinality information, and 3. variable information.

This table, from Fehd [3], shows which database tables contain which columns.

| columns | | tables | | snapshots | |
|---|---|---|---|---|---|
| name | type | dimension | fact | periodic | accumulating |
| keys: | | | | | |
| primary | integer | row-number[1] | row-id | | entity-id |
| foreign[2] | integer | | • | | |
| composite[3] | | | | • | |
| boolean | integer | | | | • |
| facts | real | | • | statistic | sum |
| text | char | information | | | |

**Notes:**
[1] a *natural key* is a row-number in (1:n-rows)
[2] foreign key is the primary key of a dimension table
[3] composite key: combination of foreign keys

---

**goal overview**

When given a new data set, the first task is to place each variable into one of three categories:   is it the row identifier?      a classification variable? or an analysis variable?

| task | common name | var type c | n | cr-type |
|---|---|---|---|---|
| data review | row-id | c | | unique |
| | row-number | | n | unique |
| by\|classification vars | discrete | c | n | few |
| analysis vars | continuous | | n | many |
| problem: empty | n-levels=1 | | | |

---

**previous algorithm**

This is pseudo-code of the previous algorithm first shown in the SmryEachVar suite, 2008: Fehd [2], and further developed in 2008: Fehd [3], 2013: Fehd [6], and 2014: Fehd [8].

1. make list of var names                          contents/sql out=list-names

2. for each variable
      make list of values                          proc freq out=freq-of-variable
      save n-obs (n-levels)                                            data n-obs
                                                    proc append to list-of-n-obs

3. join list-names with list-of-n-obs                            data list-names
      calculate cardinality ratio

4. calculate mean of cardinality ratio                              proc summary

5. calculate cardinality type (cr-type)                          data list-names

---

**issues**                      The new algorithm addresses these issues.

         type :   variable type has three sets of values
                  contents: (1=n, 2=c);          scl: (C,N);          sql: (char,num)
                  choice: `$char1` in (c,n) to conform to style guide: use lowcase

    attributes :  data set attributes n-obs and n-vars are in local symbol table
                  copy to macro variables in the global symbol table with scl `attrn(...)`
                  for use in array allocation in second step

   calculation :  mean of cardinality ratio previously calculated by summary procedure
                  is calculated in an array

  side effects :  enhancements for next step:
                  macro variables with lists of *few*, and *many* are saved
                  and written to log

---

**new algorithm**              • step 1.1: make data set `list-names` with `proc freq nlevels`

                                • step 1.2: copy data set attributes *n-obs* and *n-vars*
                                     to macro variables in global symbol table
                                     for use in array allocation and calculation

                                • step 2: data structure: attribute, array
                                     loop: calculate cardinality ratio
                                     calculate mean(cardinality ratio)
                                     loop: read `list-names`,
                                        fetch var-name information with scl functions
                                        calculate cardinality type
                                        make lists from cardinality type

---

**output**                     This is the data structure of the dimension table, list-names.

```
create table WORK.LIST_NAMES
  (label='memname=class,obs=19,vars=5')
  (memname     char(32),
   varnum      num     label='var num',
   cr_type     char(10) label='card. ratio type',
   card_ratio  num      label='card. ratio',
   n_levels    num      label='n-levels nobs=19',
   name        char(32) label='name',
   type        char(1),
   length      num,
   format      char(49),
   informat    char(49),
   label       char(256)
```

This is an example report for the data set `sashelp.class`.

| memname | var num | cr_type | card_ ratio | n_ levels | name | type | length etc. |
|---------|---------|---------|-------------|-----------|------|------|-------------|
| class | 1 | .unique | 1 | 19 | Name | c | 8 |
| class | 2 | few | 0.10526 | 2 | Sex | c | 1 |
| class | 3 | few | 0.31579 | 6 | Age | n | 8 |
| class | 4 | many | 0.89474 | 17 | Height | n | 8 |
| class | 5 | many | 0.78947 | 15 | Weight | n | 8 |

## Research for Data Structure: Contents, SQL, SCL, N-levels

**Proc Contents**

**Overview**                    This section contains these topics.

- contents program
- describe contents output
- contents listing
- contents output data set

**contents program**

```
1   %let libname = sashelp;
2   %let memname = class;
3   PROC contents data =  &libname..&memname
4                 out  =   contents;
5   PROC sql;     describe table &syslast;
6                 quit;
7   PROC print    data =        &syslast noobs;
8                 var    varnum name type length;
```

**describe contents output**

```
create table WORK.CONTENTS
  (LIBNAME  char(8)   label='Library Name',
   MEMNAME  char(32)  label='Library Member Name',
   NAME     char(32)  label='Variable Name',
   TYPE     num       label='Variable Type',         <---<<<
   LENGTH   num       label='Variable Length',
   VARNUM   num       label='Variable Number',
   LABEL    char(256) label='Variable Label',
   FORMAT   char(32)  label='Variable Format',        <---<<<
```

**Notes:**  contents.type is numeric; length of format is $char32.

**contents listing**

```
Alphabetic List of Variables and Attributes
#     Variable    Type    Len
-     --------    ----    ---
3     Age         Num     8
4     Height      Num     8
1     Name        Char    8
2     Sex         Char    1
5     Weight      Num     8
```

**Notes:**  #: variable number is the row number;
the contents listing is alphabetically ordered by `Variable` (name)

**contents output data set**

```
VARNUM    NAME    TYPE*   LENGTH
------    ----    ----    ------
   3      Age      1       8
   4      Height   1       8
   1      Name     2       8
...
```

**!** →   * contents.type in (1=Num,2=Char).

5

**Proc SQL**

**Overview**                     This section contains these topics.

- sql program
- describe dictionary.columns
- describe sashelp.class
- sql report

**sql program**

```
1   %let libname = sashelp;
2   %let memname = class;
3   PROC sql; describe table dictionary.columns;
4           describe table &libname..&memname;
5           select  memname, varnum, name, type, length
6           from    dictionary.columns
7           where      libname eq "%upcase(&libname)"
8                   and memname eq "%upcase(&memname)";
9           quit;
```

**describe dictionary.columns**

```
create table DICTIONARY.COLUMNS
  (libname char(8)   label='Library Name',
   memname char(32)  label='Member Name',
   name    char(32)  label='Column Name',    unique: row-id,      primary key
   type    char(4)   label='Column Type',    <---<<<
   length  num       label='Column Length',
   varnum  num       label='Column Number',  unique: row-number, primary key
   label   char(256) label='Column Label',
   format  char(49)  label='Column Format',  <---<<<
```

**Notes:**  type is $char4; length of format is $char49.;
both `name` and `varnum` are unique, and therefore each is a primary key.

**describe sashelp.class**

```
create table SASHELP.CLASS( label='Student Data' )
  (Name   char(8),
   Sex    char(1),
   Age    num,
   Height num,
   Weight num
```

**sql report**

```
          Column
Member    Number    Column  Column  Column
Name      in Table  Name    Type    Length
------    --------  ------  ------  ------
CLASS            1  Name    char         8
CLASS            2  Sex     char         1
CLASS            3  Age     num          8
CLASS            4  Height  num          8
CLASS            5  Weight  num          8
```

**Notes:**  sql output is ordered by `Column Number (varnum)`;    type is in (char,num);
`varnum` is a *natural key*, its values are in *(1:n-vars)*.

**SAS Component Language (scl)**

**Overview**

This section contains these topics.

- scl program
- scl output data set

**scl program**

```
1   %let libname = sashelp;
2   %let memname = class;
3   DATA list_names_scl;
4       attrib varnum length =   8
5              name   length = $32
6              type   length = $ 1;
7       drop   _:; *** _temporary vars;
8
9   *see the varnum function for this example;
10  _dsid   = open ("&libname..&memname");
11  _n_vars = attrn(_dsid,'nvars');
12
13  do _i = 1 to _n_vars;
14     name   = varname(_dsid,_i);
15     varnum = varnum (_dsid,name);
16     type   = vartype(_dsid,_i);
17     output;
18     end;
19  _rc = close (_dsid);
20  stop;
21  run;
22  PROC print data = &syslast;
```

**Notes:** The scl functions `open` and `close` allow access to the other scl functions `varname`, `varnum`, and `vartype`;
`varname` and `vartype` depend on the loop index variable, `_i`,
while `varnum`, which is equal to the loop index variable, depends on the variable `name`.

**scl output data set**

```
varnum    name       type
------    ------     ---
    1     Name       C
    2     Sex        C
    3     Age        N
    4     Height     N
    5     Weight     N
```

**Notes:** Scl.type is in upcase(C,N).

7

## Proc Freq N-levels

**Overview**                    This section contains these topics.

- proc freq program
- describe table sashelp.air
- data structure of freq nlevels

**proc freq program**

```
1   %let data = sashelp.air;
2   PROC freq   data    = &data
3               nlevels ;
4               ods        output
5               nlevels = list_names_nlevels;
6   PROC sql;  describe table &data;
7              describe table &syslast;
8              quit;
```

**describe table sashelp.air**

```
create table SASHELP.AIR
  (label='airline data (monthly: JAN49-DEC60)')
  (DATE num format=MONYY.,
   AIR  num label='international airline travel (thousands)'
```

**Notes:**    `sashelp.air` has variables with labels.

**data structure of freq nlevels**

```
create table WORK.LIST_NAMES_NLEVELS
  (TableVar      char(4)  label='Table Variable',
   TableVarLabel char(40) label='Table Variable Label',
   NLevels       num      format=BEST8.
                          label='Number of Levels'
```

**Notes:**    Variable names, `TableVar`, `NLevels`, are different from the contents and sql procedures, their variable name for the *row-identifier* is `name`.

## Research Summary

**Data Structure Issues**           The above research highlights these issues in the new algorithm.

| | |
|---|---|
| order : | `memname`, `varnum` *(row-number)*, `name`, `type`, `length`, etc. |
| freq n-levels : | variable names `TableVar` and `NLevels` to be renamed to `name`, `n_levels` |
| type : | scl lowcase(c,n) |
| length : | of format and informat: sql $char49. |

## Explanation

**Overview**

The explanation of the program contains this list of topics.
The program listing of `ml-namex.sas` is on page 13.

**autoexec**

All programs shown in this paper depend on an *autoexec* which has two *filerefs*: one named *project* for the folder containing the programs shown here, and the other named *site_inc* for the folder containing the program `ml-namex.sas`.

```
1   ** name:  autoexec.sas;
2   filename  project  '.'; *** contains programs listed here;
3   *author's useage: folder containing ml-namex.sas;
4   *filename site_inc '<...>\SAS-site\includes';
5   *for your testing:;
6   filename  site_inc '.';*folder containing ml-namex.sas;
```

**testing program**

This is the program used to test `ml-namex.sas`.

```
1   options mprint source2;
2   %let libname = sashelp;
3   %let memname = class;
4   %include site_inc(ml-namex);
5   proc print data = &syslast    noobs;
6   proc sql; describe table &syslast;
7           quit;
```

**proc freq**

This is step 1.1. The frequency output data set from a data set with variable labels contains extra variables; keep the two desired variables and rename them.

```
1   PROC freq data = &libname..&memname     nlevels;
2       ods   output nlevels=
3       list_names(keep  = tablevar       nlevels
4                   rename=(tablevar=name  nlevels= n_levels));
```

**copy n-obs, n-vars**

This is step 1.2. These macro variable assignment statements copy the data set attributes, *nobs* and *nvar* into the global symbol table for use in the next step. N-obs is used as denominator of the calculation of cardinality ratio; n-vars is used for the dimension of the array of variable names.

```
1   ****  copy     local n-obs and n-vars to global symbol table;
2   %let _dsid   = %sysfunc(open (&libname..&memname,i));
3   %let _n_obs  = %sysfunc(attrn(&_dsid,nobs));
4   %let _n_vars = %sysfunc(attrn(&_dsid,nvar));
5   %let _rc     = %sysfunc(close(&_dsid));
```

**data structure**

The data structure contains three sets of information:
1. *memname*, the relation to other members in the *libref*,
2. information about cardinality ratio,
3. information about variable.

```
1   DATA &syslast;
2       attrib memname      length = $32
3              varnum       length =   8 label = 'var num'
4              cr_type      length = $  %length(n-levels=1)  %*** label of 1st cr-type;
5                           label =    'card. ratio type'
6              card_ratio   length =   8     %*range=(0:1];
7                           format =    bestd7.5
8                           label =    'card. ratio'
9              n_levels     length =   8
10                          label =    "n-levels nobs=&_n_obs"
11             name         length = $32 label = 'name'
```

**array of cardinality ratios**

In the previous algorithm cardinality ratio was calculated in one data step and the summary procedure was used to calculate the mean; this can be accomplished in an array.

**!** → Note the use of the global macro variables, n_vars and n_obs.

```
1   array  _cr(&_n_vars);                          *** <---<<< global n_vars;
2   *...;
3   ** loop: for each row, calculate cardinality ratio;
4   do _i = 1 to dim(_cr);
5     set &syslast   (keep = n_levels)   point = _i;
6     _cr(_i) = n_levels  /  &_n_obs;              *** <---<<< global n_obs;
7     end;
8
9   ***** calculate mean for select card_ratio to cr_type;
10  _mean_cr = mean(of _cr(*));
```

**read data structure**  The scl function `open` make the data structure functions `varnum` and `vartype` available. This data-set-reading loop performs two tasks:
1. read the frequency output data set variables: name and n-levels;
2. read the variable information: `varnum`, `type`, etc.

```
1   ** read syslast(name n_levels), fetch data structure info;
2   _dsid = open("&libname..&memname");
3   do _i = 1 to dim(_cr);
4      set &syslast   point = _i;
5      varnum     =         varnum (_dsid,name);
6      type       = lowcase(vartype(_dsid,varnum));
```

**Notes:**  The lookup of `varnum` is based on the *primary-key* of the frequency-n-levels table which is the variable `name`. In fact the index variable `_i` is the same as `varnum`, but that is an assumption, which is the reason that the lookup uses the variable `name`.

---

**make cr-type**  *Cardinality-type* is assigned within the data-set-reading loop.

```
1      card_ratio = _cr(_i);
2      select;
3        when(n_levels   eq 1        ) cr_type = 'n-levels=1';
4        when(card_ratio eq 1        ) cr_type = '.unique';
5        when(card_ratio gt _mean_cr) cr_type = 'many';
6        otherwise                     cr_type = 'few';
7        end;
```

**Notes:**  Category *n-levels=1* identifies a useless variable with only one value.
Category *unique* has a dot in front of the value which places it at the top of an ordered listing.

---

**robust length of lists**  The data structure includes temporary variables for lists of variable names. The maximum length of each of these variables is *(32+1)\*n-vars*.
The maximum length of a character variable is $2^{15} - 1 = 32767$.
For the case where a data set contains more than $\frac{32767}{33} = 992$ variables, the allocation of the length of these variables fails. This problem is solved by reducing the length to 32767 in the extreme case.

```
1      %*** _temp vars for macro variables;
2      %*** 33: length(var-name)=32 +1 for delimiter;
3      %let _max_length_list = %sysfunc(min(
4                      %eval(33*&_n_vars),32767));
5          _list_few    length = $&_max_length_list
```

---

**make lists**  This paragraph creates global macro variables of lists of the variables in the *few* and *many* categories for use by later list-processing routines.

```
1      ** make list_many_c, list_many_n, for mvars;
2      if      cr_type eq 'many' then do;
3        if      type = 'c' then
4            _list_many_c = catx(' ',_list_many_c,name);
5        else if type = 'n' then
6            _list_many_n = catx(' ',_list_many_n,name);
7        end;
```

---

11

**make lengths**

This paragraph creates global macro variables of lengths of the variable in the *few* category for use by later list-processing routines.

```
1    if type eq 'c' then do;
2        _max_length = max(_max_length,length);
3        _length_few = sum(_length_few,length);
4        end;
5    else _length_few = sum(_length_few,%length(&_n_obs)); *<---<<< global n_obs;
6    *** add one for delimiter = space;
7    _length_few = sum(1,_length_few);
```

**make list-few**

This paragraph creates a global macro variable of the variable names in the *few* category in sorted order for use by later list-processing routines. This listing has these paragraphs:

1. allocation of the array, the length of the item is the the number of digits in the macro variable n-obs plus one for the delimiter and 32 for the length of the variable name

2. assignment of value: n-levels with leading zeroes, colon, and variable name

3. sorting the array into n-levels order

4. loop to copy each variable name into the list

```
1    *1   *****  _lfs: list-few-sorted value=000:name;
2        array  _lfs(&_n_vars) $%eval(%length(&_n_obs)+33);
3    *2 ...;
4      else if cr_type eq 'few'  then do;
5        **** save for sort: lfs(i)= '000:name';
6        _lfs(_i) = catx(':',put(n_levels
7                           ,z%length(&_n_obs).),name);
8    *3 ...;
9    **** loop: make list-few ordered by n-levels;
10   call sortc(of _lfs(*));
11   *4;
12   do _i = 1 to dim(_lfs);
13       ****        _lfs(_i)=000:name;
14       name = scan(_lfs(_i),-1,':');
15       _list_few = catx(' ',_list_few,name);
16       end;
```

**make macro variables**

The call symput function is used to create the set of macro variables which are used by succeeding routines.

```
1    call symputx('_list_few'    ,_list_few   );
2    call symputx('_list_many_c' ,_list_many_c);
3    call symputx('_list_many_n' ,_list_many_n);
```

**echo macro variables**

This information is written to the log at the end of the subroutine.

```
1    %put echo: &=memname &=_n_obs &=_n_vars;
2    %put info: &=_list_few;
3    %put info: &=_list_many_c;
4    %put trace: ml-namex make-list-names-cr ending;
```

result:

```
1    echo: MEMNAME=class _N_OBS=19 _N_VARS=5
2    info: _LIST_FEW=Sex Age
3    info: _LIST_MANY_N=Height Weight
4    trace: ml-namex make-list-names-cr ending
```

12

## Program Listings

**ml-namex.sas**                    This is the listing of the program `ml-namex.sas`.

```
1    %put trace: ml-namex make-list-names-cr beginning;
2    %put echo parameters: &=libname &=memname;
3     /*   name: ...\SAS-site\includes\ml-namex.sas
4         author: Ronald J. Fehd 2016
5        ----------------------------------------------------
6    Summary:  description: make list of variable names,
7                           cardinality ratio and card-type
8                  purpose: for list processing routines
9                  -------------------------------------------
10   Contexts: program group: procedure returns data set
11            program  type: subroutine
12            SAS      type: parameterized include
13            uses routines: n/a
14       ----------------------------------------------------
15   Specifications: input   : macro variables
16                            libname: libref
17                            memname: data set name
18                  process: proc freq n-levels
19                           copy memname.nobs, .nvars
20                                from local to global
21                           data: read var information
22                                with scl functions
23                           calculate cardinality ratio
24                                and card-ratio-type
25                  output : list_names, a dimension table
26                           macro-variables: _length_few
27                           _list_few    _list_many_c
28                           _list_unique _list_many_n
29                           _max_length_c
30                             for data structure attrib
31                             valu_c length=$&_max_length_c
32   note: output is ordered by varnum
33       ----------------------------------------------------
34   usage:   autoexec:
35   filename site_inc '<...\SAS-site\includes>';
36   - - - ml-names-x-test.sas - - -
37   %let libname = sashelp;
38   %let memname = class;
39   %include site_inc(ml-namex);
40   proc sql; describe table list_names;
41           quit;
42   proc print data = list_names;
43       -----------------------------------------------------
44   this is a Derivative Work of the SmryEachVar suite
45   www.sascommunity.org/wiki/SmryEachVar_A_Data_Review_Suite
46   this program is available on:
47   http://www.sascommunity.org/wiki/Making_Lists
48   http://www.mwsug.org/2016-proceedings.html paper TT03
49   it is designed to work with parameterized includes:
50   CxInclude, and SmryHuge; macros: CallMacro and CallText
51       ----------------------------------------------------
52   date-time: 2016-08-24 4:00:50 PM
53   word count      words:  734
54                   lines:  195
55   characters(no   spaces): 5410
56   characters(with spaces): 7893
57   **** .................... */
58   **** make dimension table: name + n-levels;
59   ODS  exclude all; * noprint;
60   PROC freq data   = &libname..&memname nlevels;
61        ods  output
62             nlevels= list_names
63           (keep   = tablevar          nlevels
64             rename =(tablevar= name     nlevels=n_levels));
65   run;
66   ODS select all;
```

```
67
68    ****  copy       n-obs, n-vars from local to global;
69    %let _dsid   = %sysfunc(open (&libname..&memname,i));
70    %let _n_obs  = %sysfunc(attrn(&_dsid,nobs));
71    %let _n_vars = %sysfunc(attrn(&_dsid,nvar));
72    %let _rc     = %sysfunc(close(&_dsid));
73    %symdel _dsid _rc;
74
75    **** enhance dimension table with variable information,
76         cardinality ratio, card-ratio-type: cr-type;
77    DATA &syslast(label=
78        "memname=&memname,obs=&_n_obs,vars=&_n_vars");
79        attrib memname      length = $32
80            %* primary key;
81                varnum       length =   8 label = 'var num'
82            %* cardinality information;
83                cr_type      length = $  %length(n-levels=1)
84                             label =    'card. ratio type'
85                card_ratio   length =   8      %*range=(0:1];
86                             format =   bestd7.5
87                             label =    'card. ratio'
88                n_levels     length =   8
89                             label =    "n-levels nobs=&_n_obs"
90            %* variable information;
91                name         length = $32 label = 'name'
92                type         length = $ 1 %*range:(c,n) from scl;
93                length       length =   8
94                format       length = $49
95              informat       length = $49
96                label        length =$256
97                _length_few  length =   4
98          %*** _temp vars for macro variables of lists;
99          %*** 33: length(var-name)=32 +1 for delimiter;
100         %let _max_length_list = %sysfunc(min(
101                       %eval(33*&_n_vars),32767));
102             _list_few    length = $&_max_length_list
103             _list_many_c length = $&_max_length_list
104             _list_many_n length = $&_max_length_list
105             _list_unique length = $&_max_length_list;
106           %put echo calculation: &=_max_length_list;
107         %symdel _max_length_list;
108         array  _cr(&_n_vars);
109         *****  _lfs: list-few-sorted value=000:name;
110         array  _lfs(&_n_vars) $%eval(%length(&_n_obs)+33);
111         drop   _:;           * _temporary variables;
112         retain memname       "&memname"
113                _length_few 0 _max_length 1
114                _testing       %eval(
115                %sysfunc(getoption(mprint )) eq MPRINT
116             and %sysfunc(getoption(source2)) eq SOURCE2);
117
118   ** loop: for each row, calculate cardinality ratio;
119   do _i = 1 to dim(_cr);
120      set &syslast    (keep = n_levels)   point = _i;
121      _cr(_i) = n_levels / &_n_obs;
122      end;
123
124   ***** calculate mean for select card_ratio to cr_type;
125   _mean_cr = mean(of _cr(*));
126
127   *****   loop: read syslast(name n_levels), fetch var info;
128   _dsid = open("&libname..&memname");** for scl functions;
129   do _i = 1 to dim(_cr);
130      set &syslast    point = _i; *** primary-key=name;
131      if _testing then putlog name= n_levels=;
132      varnum      =          varnum  (_dsid,name  );
133      type        = lowcase(vartype  (_dsid,varnum));
134      length      =          varlength(_dsid,varnum);
135      format      =          varfmt   (_dsid,varnum);
136      informat    =          varinfmt (_dsid,varnum);
137      label       =          varlabel (_dsid,varnum);
138      card_ratio = _cr(_i);
139      select;
140        when(n_levels   eq 1) cr_type = 'n-levels=1';
141        when(card_ratio eq 1) cr_type = '.unique';%*row-id;
142        when(card_ratio gt
143                  _mean_cr) cr_type = 'many'; %*analysis;
```

```
144       otherwise              cr_type = 'few';  %*by|class;
145        end;
146     ** make list_many_c, _many_n, for mvars;
147     if      cr_type eq 'many' then do;
148        if     type eq 'c' then
149             _list_many_c = catx(' ',_list_many_c,name);
150         else if type eq 'n' then
151             _list_many_n = catx(' ',_list_many_n,name);
152        end;
153     else if cr_type eq 'few'  then do;
154        **** save for sort: lfs(i)= '000:name';
155        _lfs(_i) = catx(':',put(n_levels
156                           ,z%length(&_n_obs).),name);
157        *** add one for delimiter = space;
158        _length_few = sum(1,_length_few);
159        if type eq 'c' then do;
160             _max_length=max(_max_length,length);
161             _length_few=sum(_length_few,length);
162             end;
163        else _length_few=sum(_length_few,%length(&_n_obs));
164        end;
165     else if cr_type eq '.unique'  then
166         _list_unique = catx(' ',_list_unique,name);
167     output;
168     end;
169 _rc = close(_dsid);
170
171 **** loop: make list-few ordered by n-levels;
172 call sortc(of _lfs(*));
173 do _i = 1 to dim(_lfs);
174     if _testing and _lfs(_i) ne ' ' then putlog _lfs(_i)=;
175     ****        _lfs(_i)=0000:name;
176     name = scan(_lfs(_i),-1,':');
177     _list_few = catx(' ',_list_few,name);
178     end;
179
180 call symputx('_length_few' ,_length_few );
181 call symputx('_list_few'   ,_list_few   );
182 call symputx('_list_many_c' ,_list_many_c);
183 call symputx('_list_many_n' ,_list_many_n);
184 call symputx('_list_unique' ,_list_unique);
185 call symputx('_max_length_c',_max_length );
186 stop;
187 run;
188 %put info: &=memname &=_n_obs &=_n_vars;
189 %put info: &=_length_few;
190 %put info: &=_list_few;
191 %put info: &=_list_many_c;
192 %put info: &=_list_many_n;
193 %put info: &=_list_unique;
194 %put info: &=_max_length_c;
195 %put trace: ml-namex make-list-names-cr ending;
```

**Demonstration Programs, Fast Library Review**

**Overview**                    This section contains listings of these programs.

- report memnames information

- routine ml-namex, call, append

- program make-list-memnames with cardinality ratio

- make-list-memnames output

**report memnames**              This report lists the information — n-obs, n-vars, data set label — of all data
**information**                  sets in a *libref*.

```
1   %let libname = sashelp;
2   %put echo parameters: &=libname;
3   PROC sql; describe table dictionary.tables;
4           create   table list_memnames as
5           select   memname, nobs, nvar, memlabel
6           from     dictionary.tables
7           where    libname eq "%upcase(&libname)"
8             and    memtype eq 'DATA';
9           describe table &syslast;
10          select * from  &syslast;
11          title3  "&libname members=&sqlobs";
12          quit;
13  %put info: &=sqlobs;
```

```
                            memname  nobs   nvar  memlabel
                            -------
        output listing      AACOMP   1544    4
                            AARFM      61     4
                            ADSMSG    426     6
                            AFMSG    1090     6
                            AIR       144     2   airline data (monthly: JAN49-DEC60)
```

**routine**                      This program calls the subroutine `ml-namex.sas` and appends the output
**ml-namex-call-append**         to the data set `list_memnames`.

```
1   *ml-namex-call-append.sas;
2   %include site_inc(ml-namex);
3   PROC append data = &syslast
4               base = list_memnames;
5   run;
```

**make-list-memnames
with cardinality ratio**

```
1   *make-list-memnames-with-cr;
2   options mprint source2;
3   %let libname = sashelp;
4   PROC sql; select catt('%let memname=',memname
5                ,';%include project(ml-namex-call-append);')
6           into :list_memnames    separated by ' '
7           from   dictionary.tables
8           where  libname eq "%upcase(&libname)"
9             and  memtype eq 'DATA'
10            and  nobs and nvar;
11          quit;
12  %put info: &=sqlobs;
13  &list_memnames
14  run;
15  PROC print data = &syslast;
16          title3 &syslast;
17          title4 "&libname contains &sqlobs members";
18          by      memname;
19          id      memname;
```

**make-list-memnames
output**

This is an example of the output of make-list-memnames-with-cr.sas.

```
         var          card-   n-
memname  num  cr-type         ratio  levels  name    type  length
-------  ---  ----------    -------  ------  ------  ----  ------
AACOMP    1   few           0.00583       9  locale   c        5
          2   few           0.11917     184  key      c       60
          3   n-levels=1    0.00065       1  lineno   n        4
          4   many          0.87111    1345  text     c     1200

AARFM     1   n-levels=1    0.01639       1  locale   c        5
          2   .unique             1      61  key      c       60
          3   n-levels=1    0.01639       1  lineno   n        4
          4   many          0.96721      59  text     c     1200
```

**!** → This data set can be sorted by `lowcase(name)` to identify variables that
have different types.

---

**Programs, Summary-Each-Var, version 2016.08**

---

**Overview**

This section has listings of programs for the SmryEachVar suite, version
2016.08.

- cx-include-list-names

- proc-freq

- proc-smry

- demo-SmryEachVar-2016

- SmryEachVar output

---

**cx-include-list-names.sas**

This program is a custom version of the routine CallXinc, Fehd [5].

```
1   %put trace: cx-include-list-names beginning;
2   %put echo parameters &=cx_data &=cx_include;
3   DATA _null_;
4       attrib statement length = $
5       %sysfunc(max(%eval(
6        %length(*let name=)+32+1+%length(*let type=C)+1)
7       ,%length(*include &cx_include;)  ));
8   do until(endofile);
9      set &cx_data %*may contain where(...);
10        end = endofile;
11      statement =  catt('%let name=',name,';'
12                    ,'%let type=',type,';');
13      call execute(catt('%nrstr(',statement,')'));
14      statement = "%include &cx_include;";
15      call execute(catt('%nrstr(',statement,')'));
16      end;
17  stop;
18  run;
19  %put trace: cx-include-list-names ending;
```

**proc-freq.sas**

This program standardizes the output data set from the frequency procedure so that it can contain both character and numeric variables, for `cr-type eq 'few'`.

```
1   %put trace: proc-freq beginning;
2   %put echo parameters:&=libname &=memname &=name &=type;
3   %put echo parameter :&=_max_length_c from ml-namex;
4   PROC freq data   = &libname..&memname;
5           tables   &name  / list missing noprint
6           out     =  out_freq
7          (rename =(&name = var_&type) );
8   run;
9   %put echo parameter: &=_max_length_c;
10  DATA &syslast;
11      attrib memname length = $32
12              name length = $32
13            var_c   length = $&_max_length_c
14            var_n   length = 8;
15      if 0 then set &syslast(keep = count percent);
16      retain memname "&memname"
17              name "&name"
18            var_c   '.'
19            var_n   .;
20  do until(endofile);
21     set &syslast end = endofile;
22     output;
23     end;
24  stop;
25  run;
26  PROC append data = &syslast
27           base = list_freq;
28  run;
29  %put trace: proc-freq ending;
```

**proc-smry.sas**                    This program provides a simple set of summary statistics, one row for each
                                     continuous variable, for `cr-type eq 'many'` and `type eq 'n'`.

```
1    %put trace: proc-smry beginning;
2    %put echo parameters:&=libname &=memname &=name &=type;
3    PROC summary data   = &libname..&memname;
4                var      &name;
5               output
6                 out = out_summary
7             ( drop =_type_  _freq_)
8       mean    (&name) = mean    %*average;
9       std     (&name) = std_dev
10      min     (&name) = min     %*p000;
11      median  (&name) = median  %*p050;
12      max     (&name) = max     %*p100;
13        ;
14   run;
15   DATA &syslast;
16       attrib memname length = $32
17               name length = $32;
18      retain memname "&memname"
19              name "&name";
20   do until(endofile);
21     set &syslast end = endofile;
22     output;
23     end;
24   stop;
25   run;
26   PROC append data = &syslast
27             base = list_smry;
28   run;
29   %put trace: proc-smry ending;
```

**demo-SmryEachVar-2016.sas**

This program shows the basic methods of the SmryEachVar suite:

1. call subroutine to make list of variable names

2. call proc frequency for each of the by or class variables

3. call proc summary for each of the analysis variables

```
1    *name: SmryEachVar, version 2016.08;
2    *options mprint source2;
3    %let libname = sashelp;
4    %let memname = bweight;
5    %let memname = class;
6
7    ** 1 make list of variable names;
8    %include site_inc(ml-namex);
9    PROC print  data =   &syslast noobs;
10       title3 "%cmpres(&syslast unique: &_list_unique)";
11
12   ** 2 freq of few, either char or num;
13   %let cx_data    = list_names(where=(cr_type eq 'few'));
14   %let cx_include = project(proc-freq);
15   %include project(cx-include-list-names);
16   PROC print  data =   &syslast noobs;
17       title3 "%cmpres(&syslast few: &_list_few)";
18             by memname notsorted name;
19             id memname           name;
20
21   ** 3 summary of many, only numeric;
22   %let cx_data    = list_names(where=(cr_type eq 'many'
23                                   and type eq 'n'));
24   %let cx_include = project(proc-smry);
25   %include project(cx-include-list-names);
26   PROC print  data =   &syslast noobs;
27       title3 "%cmpres(&syslast many.n: &_list_many_n)";
28   run;
```

## SmryEachVar output

```
Calculating Cardinality Ratio and -type in 2 steps
demo-SmryEachVar-2016
list_names unique: Name
                               card_
memname     varnum    cr_type   ratio    n_levels    name     type    length    format    informat    label
-------     ------    -------   -------   --------    ------   ----    ------
 class         1      .unique      1            19    Name      c         8
 class         2      few       0.10526         2    Sex       c         1
 class         3      few       0.31579         6    Age       n         8
 class         4      many      0.89474        17    Height    n         8
 class         5      many      0.78947        15    Weight    n         8


list_freq few: Sex Age
memname     name     var_c    var_n    COUNT    PERCENT
-------     ----     -----    -----    -----    -------
 class      Sex      F          .        9      47.3684
                     M          .       10      52.6316
 class      Age      .         11        2      10.5263
                     .         12        5      26.3158
                     .         13        3      15.7895
                     .         14        4      21.0526
                     .         15        4      21.0526
                     .         16        1       5.2632


list_smry many.n: Height Weight
memname     name      mean      std      min    median    max
-------     ------    ------    ------    ----   ------    ---
 class      Height    62.337    5.1271    51.3    62.8      72
 class      Weight    100.026   22.7739   50.5    99.5     150
```

# Summary

**Commentary**                     The log of this subroutine's use of the frequency procedure has been re-
duced by an order of magnitude, from *O(n-vars)* to *O(1)* While this ap-
pears impressive theoretically, in fact the frequency procedure still does
the counting of n-levels for each variable and the overall time to run the sin-
gle call of frequency-nlevels is approximately equal to the sum of the time
for frequencies of each variable.

**Conclusion**                     By analyzing the data structure of both the input and output data sets from
the procedures used in previous algorithms we have found an optimum
data structure for the desired output and moved calculations from several
data steps and procedures into a single data step with arrays.

**Suggested Reading**

predecessors :   Fehd [1], Macro FreqAll; Fehd [2], SmryEachVar;
Fehd [3], Database Vocabulary: dimension table;
Fehd [9], Read Column into Row;
Fehd [6], Data Review with N-Levels or Cardinality Ratio

routines :   that use this subroutine for setup
Fehd [5], List Processing Routine CallXinc;
Fehd [4], Using SmryEachVar
Fehd [7], Macro CallMacro;
Fehd [8], Using Cardinality Ratio

scl :   Fraeman [10], using scl functions to get data structure information;
Yindra [11], review of data step and scl functions for use with `sysfunc`

programs :   sas.community.org:        Cardinality Ratio definition,        Making Lists,

Summarize Each Var,                         Routine CxInclude (CallXinc),

Routine CallXinc,

**Author Information**             Ronald J. Fehd                          Ron.Fehd.macro.maven@gmail.com

About the author:   sco.wiki                http://www.sascommunity.org/wiki/Ronald_J._Fehd
LinkedIn                               www.linkedin.com/Ronald.Fehd
affiliation                     Stakana Analytics, Senior Maverick
also known as              `macro maven` on SAS-L, Theoretical Programmer

Programs:                      **http://www.sascommunity.org/wiki/**

## References

[1] Ronald J. Fehd. Journeymen's tools: Data review macro FreqAll — using Proc SQL list processing with Dictionary.Columns to eliminate macro do loops. In *SAS Global Forum Annual Conference Proceedings*, 2007. URL `http://www2.sas.com/proceedings/forum2007/028-2007.pdf`. Coder's Corner, 10 pp.; attributes, dictionary.columns, metadata, proc append, proc freq, proc sql, program header; bibliography.

[2] Ronald J. Fehd. SmryEachVar: A data-review routine for all data sets in a libref. In *SAS Global Forum Annual Conference Proceedings*, 2008. URL `http://www2.sas.com/proceedings/forum2008/003-2008.pdf`. Applications Development, 24 pp.; call execute, data review, data structure, dynamic programming, list processing, parameterized includes, utilities (writattr, writvalu) to repair missing elements in data structure; best contributed paper in ApDev.

[3] Ronald J. Fehd. Database vocabulary: Is your data set a dimension (lookup) table, a fact table or a report? In *Western Users of SAS Software Annual Conference Proceedings*, 2008. URL `http://wuss.org/proceedings08/08WUSS%2520Proceedings/papers/dmw/dmw04.pdf`. Databases and Warehouses, 8 pp.; cardinality ratio, composite key, database design, foreign key, grain, nlevels, normal forms, primary key, relational database, snapshots: accumulating or periodic.

[4] Ronald J. Fehd. Using the SmryEachVar data review suite. In *Western Users of SAS Software Annual Conference Proceedings*, 2009. URL `http://www.lexjansen.com/wuss/2009/app/APP-Fehd1.pdf`. Applications Development, 20 pp.; call execute, data review, data structure, dynamic programming, list processing, parameterized includes, examples.

[5] Ronald J. Fehd. List processing routine CallXinc: Calling parameterized include programs using a data set as list of parameters. In *Western Users of SAS Software Annual Conference Proceedings*, 2009. URL `www.lexjansen.com/wuss/2009/app/APP-Fehd2.pdf`. Applications Development, 20 pp.; call execute, data review, data structure, dynamic programming, list processing, parameterized includes, examples.

[6] Ronald J. Fehd. Data review information: N-levels or cardinality ratio. In *SAS Global Forum Annual Conference Proceedings*, 2013. URL `http://support.sas.com/resources/papers/proceedings13/299-2013.pdf`. Statistics and Data Analysis, 6 pp.; using proc freq nlevels and nobs to calculate cardinality ratio — range in (0:1) — of a variable to determine its type in (continuous, discrete, unique, worthless).

[7] Ronald J. Fehd. List processing macro call-macro. In *MidWest SAS Users Group Annual Conference Proceedings*, 2014. URL `www.mwsug.org/proceedings/2014/BB/MWSUG-2014-BB04.pdf`. Coders Corner, 19 pp.; using %sysfunc with SCL functions to read a list, a control data set, and for each observation, call a macro with variable names and values as named parameters.

[8] Ronald J. Fehd. Using cardinality ratio for fast data review. In *Western Users of SAS Software Annual Conference Proceedings*, 2014. URL `http://www.lexjansen.com/wuss/2014/98_Final_Paper_PDF.pdf`. Coders Corner, 14 pp.; successor of SmryEachVar, macros to calculate cardinality ratio and process discrete and continuous variables with procs freq, mode, and summary.

[9] Ronald J. Fehd. Reading a column into a row to count n-levels, calculate cardinality ratio and create frequency and summary output in one step. In *MidWest SAS Users Group Annual Conference Proceedings*, 2015. URL `http://www.lexjansen.com/mwsug/2015/BB/MWSUG-2015-RF-04.pdf`. Rapid Fire, 10 pp.

[10] Kathy Hardis Fraeman. Get into the groove with %sysfunc: Generalizing SAS(R) macros with conditionally executed code. In *SAS Users Group International Annual Conference Proceedings*, 2008. URL `http://www.lexjansen.com/nesug/nesug08/po/po06.pdf`. Posters, 8 pp.; using scl functions to get data set information, 4 examples.

[11] Chris Yindra. %Sysfunc — the brave new macro world. In *SAS Users Group International Annual Conference Proceedings*, 1998. URL `http://www2.sas.com/proceedings/sugi23/Advtutor/p44.pdf`. Advanced Tutorials, 7 pp.; review of data step and scl functions available to %sysfunc, 10 examples.