

The SAS MAKESPL Macro

Ellen Hertzmark, Ruifeng Li, and Donna Spiegelman

Abstract

The %MAKESPL macro is a SAS macro that makes restricted cubic spline variables to be used in procedures. It is incorporated in several of the macros that test for non-linearity, but can also be used on its own to create spline variables for covariates (allowing better control for the covariates usually using up fewer degrees of freedom).

Keywords: SAS, macro, restricted cubic splines, nonlinear relation

Contents

1. Description
2. Invocation and Details
3. Examples
4. Computational Methods
5. References
6. Credits
7. See Also

1: Description

%MAKESPL is a SAS macro that makes restricted cubic spline variables for a given number of knot points or a list of knots. This macro is useful when you wish to control for confounding by a continuous variable with a possibly non-linear relationship with the outcome.

2: Invocation and Details

To call %MAKESPL, your program must know where to look for it. The most efficient way is to include the following statement (or its equivalent) at the top of your program.

```
options mautosource sasautos='/usr/local/channing/sasautos';
```

%MAKESPL has very few parameters.

DATA The name of data set on which the regression is to be run
 REQUIRED

SPLVBL The name of the variable to be "splined"
 REQUIRED

NK The number of knots, if you want %MAKESPL to place them
 automatically. The numbers of knots for which automatic
 knot placement can be done are
 3, 4, 5, 6, 7, 8, 9, 10, 17, 21, 25, and 50.
 See COMPUTATIONAL METHODS for details about knot placement.
 default is 4

KNOT1 The list of knots if you want to determine them.
 If you give a value for KNOT1, this takes precedence over
 any value NK may have.

REFVAL A value of SPLVBL to use as the reference level for procedures
 that need reference levels (LOGISTIC, log-binomial regression,
 PHREG).

The macro prints the values of all the spline variables
when SPLVBL = REFVAL.

To do this, it adds one observation to the dataset.

default is the minimum value of SPLVBL

NOTE: Sometimes the best way to control for covariates that are
originally continuous is to use 3 or 4-knot splines,
rather than sets of indicators for categorical levels
of the covariates.

For many purposes (such as plotting and estimation),
you will need to know the values of the spline variables
at a specified value for the covariate. You can find
this out by using that specified value as REFVAL.

If you wish to make spline variables for several covariates,
you must do them one at a time, using the OUTDAT of
one run as the DATA of the next (See Example 1).

If you set REFVAL=(value SPLVBL will have in ADJDAT),
the macro will print out the values.

NOTE: If the value of REFVAL is outside the range of your data,
the macro will give you the following warning message
in the .log file.

WARNING: Your REFVAL, <value of REFVAL> , is

lower than all values of the exposure in of the data.

The macro will continue, but the graph may look strange.

OR

WARNING: Your REFVAL, <value of REFVAL> , is

higher than all values of the exposure in of the data.
The macro will continue, but the graph may look strange.

'The graph may look strange' means that the point for which SPLVBL = REFVAL will be included in the plotting points. For example, if the range of SPLVBL in the data is 16-40, and REFVAL=0, and the predicted value for SPLVBL=0 is a lot higher than that for SPLVBL=16, your graph will have a large drop from 0 to 16, and all the real information may be squeezed into a small part of the graph.

- OUTDAT The name of the dataset %MAKESPL will make.
This dataset includes all the variables in the original dataset, plus the newly made spline variables.
It has the extra observation made for SPLVBL = REFVAL.
If MAKEPTS = T, it includes the plotting points (See below).
default=_splstuf
- MAKEPTS Whether you want to make 501 plotting points for a graph.
If MAKEPTS = T, the plotting points in the OUTDAT dataset have non-missing values for SPLVBL and COVAR (See below).
Values the outcome variable in your analysis are missing.
Because the outcome variable is missing, the plotting points will not be used in the estimation routine, but those procedures that predict values of the outcome based on the model will produce predicted values for the graph.
The macro makes a variable ({{\tt _ine_}}) in OUTDAT to distinguish the plotting points.
default is F
- COVAR The list of covariates in the model other than SPLVBL.
This is only required if MAKEPTS=T and there are covariates in your model and the vertical axis of your graph will be an absolute quantity (predicted probability, value of a continuous variable), rather than a relative quantity (relative risk, odds ratio).
- ADJDAT The name of a data set with one (1) observation containing the values of each covariate to be used for plotting, IF there are covariates in your model and MAKEPTS=T.
This is REQUIRED for PROC MIXED and PROC GENMOD based procedures, as well as for the PLOTPROB=T option in LGTPHCURV8.

3: Examples

Using a data set from HPFS, we are planning to examine the relationship of a number of risk factors to BMI, cross-sectionally in 1986. We plan to use splines for all the continuous variables.

BMI86 is the individual's BMI in 1986

age86 is the individual's age (in years) in 1986
tfat86n is the individual's daily intake of total fat
in grams per day in 1986
alco86n is the individual's daily intake of alcohol
in grams per day in 1986
smk86 is the individual's smoking status in 1986
(0=non-smoker, 1=smoker)

Example 1. Making splines for age86 and tfat86n.

To make "spline" variables for age86 and tfat86n, and to print out 3 observations from the output data set each time we ran the macro, we used the following code.

```
title2 'make 4 knot splines for age86';
%makespl(data=all, splvbl=age86, nk=4, refval=55, outdat=a1);

title2 'first 3 observations of a1';
proc print data=a1 (obs=3);
  var age86 age861 age862 tfat86n alco86n bmi86;

title2 'make 5 knot splines for tfat86n';
%makespl(data=a1, splvbl=tfat86n, nk=5, refval=70, outdat=a2);

title2 'first 3 observations of a2';
proc print data=a2 (obs=3);
  var age86 age861 age862 tfat86n tfat86n1 tfat86n2 tfat86n3
  alco86n bmi86;
run;
```

By using the output data set (OUTDAT) from the first run of %MAKESPL as the input data set (DATA) for the second run of %MAKESPL, the output data set of the second run (a2) has the spline variables for age86 and for tfat86n.

Since we are not planning to do the analysis before running we left MAKEPTS=F (the default) both times.

Since NK for age86 was 4, the "spline" variables for age86 are named age861 and age862. Since NK for tfat86n was 5, the "spline" variables for tfat86n are named tfat86n1, tfat86n2, and tfat86n3. In general the names will be variable name1 ...variable name(nk-2), where NK is the number of knots, either given to the macro or counted from a list given in KNOT1.

The results of the above bit of program are shown below.

First the macro lists the knot locations.

Because we gave the value 55 for REFVAL in the first call to %MAKESPL and 70 for REFVAL in the second call to %MAKESPL, we get printouts of the values of the "spline" variables at the reference values. If we had not given a value for REFVAL, the macro would have used the

Let t_j be the j th knot point.

Let $kd = (t_{nk} - t_1)^{2/3}$, where kd is a normalizing parameter to get the spline variables back into the original units.

For a level of the exposure x , x_j , the value of the j th spline variable (j runs from 1 to $NK-2$) is given by

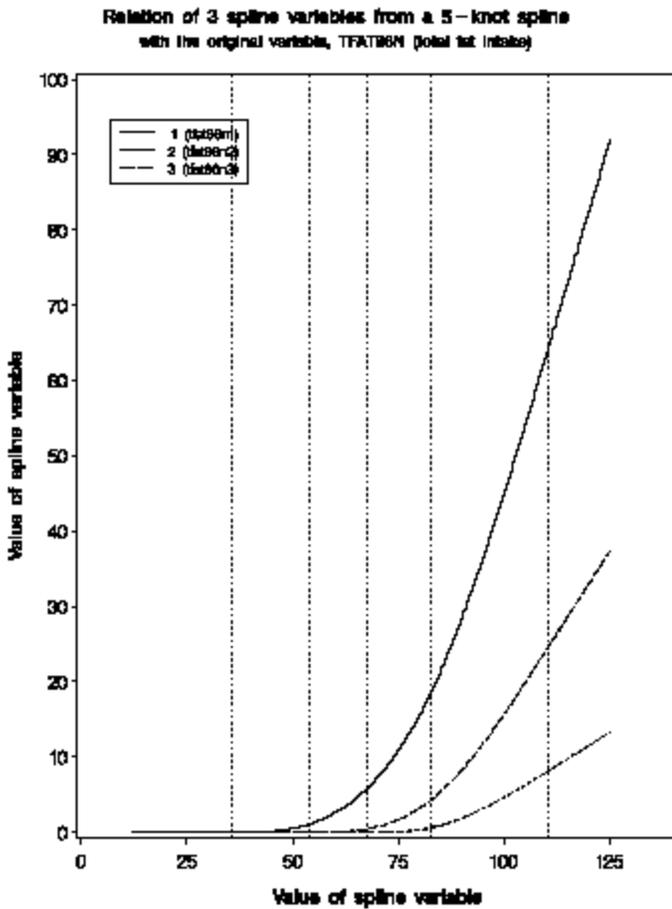
$$\begin{aligned} x_j = & \max((x - t_j)/kd, 0)^3 \\ & + (t_{nk-1} - t_j) * \max((x - t_{nk})/kd, 0)^3 \\ & - (t_{nk} - t_j) * \max((x - t_{nk-1})/kd, 0)^3 / (t_{nk} - t_{nk-1}) \end{aligned}$$

For $x < t_j$ the value of the j th spline variable is 0 (as are the 'higher' spline variables) (because all the 'max' values are 0, since $x < t_j < t_{nk-1} < t_{nk}$). As x gets larger, it has more and more nonzero spline variables.

As an example, we show the values of the 3 "spline" variables for the 5 knot splines of `tfat86n`. The graph below was made using the 'plotting points' using the following code.

```
%makespl( data=allid, splvbl=tfat86n, nk=5, makepts=T, outdat=a99);  
data inest; set a99;  
if _ine_ eq 1;  
run;
```

This code uses the automatic `_ine_` variable to select the 'plotting points.'



The reference lines are at the knot locations.

In fact, `tfat86n1` becomes nonzero at the first knot, `tfat86n2` becomes nonzero at the second knot, and `tfat86n3` becomes nonzero at the third knot. Because of the large range of the vertical axis, this is not visible on the graph.

5: References

ovindarajulu U.S., Spiegelman D., Thurston S.W., Ganguli B., Eisen E.A.: "Comparing smoothing techniques in Cox models for exposure-response relationships". *Statistics in Medicine*, 2007; 26:3735-3752.

Govindarajulu U.Malloy, E.J., Ganguli, B., Spiegelman, D, Eisen, E.A.: A comparison of fitted non-linear exposure-response relationships in Cox models using smoothing methods through simulations. In preparation, 2007.

Smith, Patricia L.: Splines as a useful and convenient statistical tool. *The American Statistician* 33(2): 57-, 1979.

Harrell, Frank E, Jr., Lee, Kerry L., Pollock, Barbara G.: Regression models in clinical studies: determining relationships between predictors and response. *JNCI* 80: 1198-1202, 1988.

Durrleman, Sylvain, and Simon, Richard: Flexible regression models with cubic splines. *Statistics in Medicine* 8: 551-561, 1989.

Devlin TF, Weeks BJ (1986): Spline functions for logistic regression modeling. Proc Eleventh Annual SAS Users Group International. Cary NC: SAS Institute, Inc., pp. 646-51.

Stone CJ, Koo CY (1985): Additive splines in statistics. *Proc Stat Comp Sect Am Statist Assoc*, pp. 45-8.

Stone CJ (1986): Comment, pp. 312-314, to paper by Hastie T. and Tibshirani R. (1986): Generalized additive models. *Statist Sciences* 1:297-318.

6: Credits

Written by Ellen Hertzmark and Donna Spiegelman for the Channing Laboratory. Based on a macro by Frank Harrell. Questions can be directed to Ellen Hertzmark, stleh@channing.harvard.edu, (617) 432-4957.

7: See Also

Channing SAS macro %LGTPHCURV9