

ON ESTIMATING SPLINE REGRESSIONS

Lawrence C. Marsh, University of Notre Dame

Abstract

This paper presents a SAS spline polynomial regression program which pieces together polynomial regressions of different orders. These polynomial regression segments are connected at join points or so-called spline knots where the number and location of these knots are estimated by the program.

I. Introduction

Suits et al (1978) and Smith (1979) have presented clear and quite useful approaches to estimating spline functions with known knot locations. Their work was at least in part based upon the development of this method by such authors as Fuller (1969), Poirier (1973, 1975, 1976), and Buse and Lim (1977).

In his dissertation and subsequent journal article Robison (1964) discusses estimating the points of intersection of two polynomial regressions. He outlines maximum likelihood methods for estimating the regression coefficients and the location of points of intersection of the two regression equations.

Hudson (1966) provides further insights into this problem. He examines four types of join points for joining two polynomial regressions depending primarily upon whether each join point is at an abscissa data point or between two such points, and whether the regressions have equal or unequal slopes at the join points. Hudson suggests that in some cases a constrained least squares regression search routine could be used to estimate the location of the join points.

Gallant and Fuller (1973) make an important contribution in dealing with some of the issues raised by Hudson. Using the continuity and differentiability conditions they reparameterize the spline regression model to form a nonlinear regression model. This nonlinear model can then be estimated using Hartley's Modified Gauss-Newton Method of minimization to obtain least squares estimates of the regression parameters. They also point out some asymptotic properties and hypothesis testing implications of this approach.

The estimation technique to be presented herein will not be as sophisticated as the Gallant and Fuller approach but for Hudson's Type Two models will provide estimates of the number and location of spline knots using a simple extension of the Suits and Smith methods.

II. Spline Theory and SAS Models

The purpose of this paper is to demonstrate a simple method of estimating the number and location of knots in spline regressions. Discontinuities are not restricted to the highest order nonzero derivatives, but may occur for one or more derivatives at any knot location. Linear adjustments are shown to generate discontinuities only for the first derivatives. A quadratic adjustment generates discontinuities

only for the second derivatives while cubic adjustments only alter the third derivatives. Any combination of adjustments is permitted at any of the potential knot locations. However, this freedom can be restricted if one so desires.

The set of potential knot locations may consist of the values of the explanatory variable, X , or could be a predefined set of values such as one through ninety-nine for years of experience, one through thirty for years of education, or 1950 through 1980 to represent time in years. A limit of about a thousand observed or prespecified values is the upper bound on the number of potential knot locations for cubic spline models in this program. Regressions with higher order splines would be limited to fewer potential knot locations while those of lesser order could have more.

Statistical properties are not examined in this paper, but generally are the same as those commonly associated with stepwise regression procedures. The theoretical issue of what parameters are proper to estimate in a model and which ones should be specified *a priori* is not dealt with in this paper but remains a question of continuing debate in applied statistics.

There are many alternative ways of formulating spline regression models. Smith (1979) offers an approach to specifying spline models that is convenient for our purposes. In particular her "+" functions readily lend themselves to our stepwise method for estimating the number and location of spline knots.

Smith initially provides a general model for k knots in an n degree polynomial regression:

$$(1) \quad y = \sum_{j=0}^n B_{0j} X^j + \sum_{i=1}^k \sum_{j=0}^n B_{ij} (X - t_i)_+^j + e$$

This model with no continuity restrictions could be called an unrestricted dummy variable model. Each segment has its own unrestricted constant term and slope values. For example, a second degree polynomial with two knots provides for three quadratic sections as follows:

$$(2) \quad y = B_{00} + B_{01} X + B_{02} X^2 + B_{10} D1 + B_{11} (X - t_1) D1 + B_{12} (X - t_1)^2 D1 + B_{20} D2 + B_{21} (X - t_2) D2 + B_{22} (X - t_2)^2 D2 + e$$

The dummy variables $D1$ and $D2$ are turned on as X passes knots t_1 and t_2 respectively (i.e. $X < t_1$ implies $D1=0$, $X \geq t_1$ implies $D1=1$, $X < t_2$ implies $D2=0$, $X \geq t_2$ implies $D2=1$).

This unrestricted dummy variable model could be estimated as is or a more refined model could be developed by imposing some continuity restrictions on the model. For example, the polynomial line segments can be made to touch by

eliminating the B_{10} and B_{20} terms. This provides a quadratic spline model which is joined at the knots but may have sharp corners at the join points due to the first derivatives being unequal. This can be smoothed out by setting B_{11} and B_{21} both equal to zero. This reduces the sharpness of the turning points at the knots by forcing the first derivatives to be equal at each knot:

$$(3) \quad y = B_{00} + B_{01} X + B_{02} X^2 + B_{12} (X - t_1)^2 D_1 + B_{22} (X - t_2)^2 D_2 + e$$

This formulation makes the values of the functions equal at the knots as well as the values of the first derivatives. The second derivatives are left unequal. Of course, if we made the second derivatives equal at the knots as well, we would end up with just one big quadratic equation covering the entire range of data.

By leaving only the highest-order, nonzero derivatives unequal, we have what Smith calls the smoothest possible spline, which she expresses in general terms as:

$$(4) \quad y = \sum_{j=0}^n B_{0j} X^j + \sum_{i=1}^k B_{in} (X - t_i)^n + e$$

This is a spline model that is as restrictive as it can be without losing its spline character. It offers some flexibility but is close to the single polynomial equation model.

If the location of the spline knots were known in advance, then we could use some selection procedure such as those available in PROC STEPWISE in SAS to estimate the degree of the polynomial within each segment and the continuity restrictions that appear to be most significant statistically. For example, if salary is considered to be related to years of education, then one might assume a spline knot at twelve years for the high school diploma, sixteen years for the B.A. degree, and eighteen for an M.B.A. The degree of the polynomial and continuity restrictions could then be determined in reference to these three knots.

But what if the number and location of knots are not known in advance? The traditional approach to such a problem is to present it as a maximum likelihood estimation problem. This is often done because of the desirable consistency and asymptotic normality properties that are often forthcoming for maximum likelihood estimators.

However, an alternative approach using stepwise regression methods can be used profitably in many cases. Suppose that instead of viewing salary as a function of years of education, we wish to view it as a function of years of experience. Experience may not offer us well defined knot locations the way education did. Instead we may want to estimate the number and location of knots for years of experience.

Say that our data set has five thousand cases but less than one hundred different values for years of experience. We could search over

all of the integers from one to one hundred for knot locations, but it may be more desirable to search over the set of actual values of experience in the data set. We can find these values by first sorting the data set by experience and then using FIRST.X or LAST.X with an OUTPUT statement to obtain the subset of unique values that experience takes on. Any time the variable of interest has repeated values in the data set, this approach will avoid the redundancy of considering the same potential knot location more than once. In any event we boil the set of candidates for possible spline knot locations down to hopefully less than one hundred but not more than a few hundred values.

Next we create a "+" function type dummy variable for each possible knot location. The X variable, experience, might take on values one through seventy-three with seventy-three corresponding knot locations $t_1 = 1$ through $t_{73} = 73$. A corresponding set of dummy variables (D_1 through D_{73}) can then be set up such that $D_i = 0$ if $X < t_i$ and $D_i = 1$ if $X \geq t_i$.

Now assume that a cubic spline model is desired. This means that three sets of spline variables will be needed. These are the linear spline variables: $(X - t_i) * D_i$, the quadratic spline variables: $(X - t_i)^2 * D_i$, and the cubic spline variables: $(X - t_i)^3 * D_i$. Altogether for a cubic spline model with 73 possible knot locations, the unrestricted dummy variable model would be:

$$(5) \quad y = B_{00} + B_{01} X + B_{02} X^2 + B_{03} X^3 + \sum_{i=1}^{73} \sum_{j=0}^3 B_{ij} (X - t_i)^j D_i + e$$

This certainly covers a lot of ground. There are 296 coefficients that might potentially be estimated here. Fortunately, PROC STEPWISE can be used to select out the ones that are statistically significant.

The SAS programming requirements for this type of model are fairly well defined. The program should be able to handle an unknown number of cases with an unknown number of unique knot locations where some maximum value is specified for the degree of the polynomial. The stepwise procedure can then be used to select the statistically significant knot locations and the degree of the polynomial appropriate within each spline segment defined by these knot locations. In other words, any combination of polynomials of various degrees may be found to fit the data. No *a priori* restriction is made on the degree of the polynomial within each segment except that it not exceed the overall maximum set in advance. Since cubic splines seem fashionable the maximum degree could be set at three. However, there is nothing to prevent that maximum from being set at four or five or even nine or ten if it were desired. Of course, one eventually reaches the limits of reasonable CPU core and time usage. In practice, however, these limits are rarely reached for the types of research

problems of interest here.

III. A SAS Spline Regression Program

The set of potential knot locations may be specified in advance or may be defined as the set of unique values taken on by a particular observed variable. For example, the latter approach might be appropriate if one wished to restrict the search for knots to the actually observed values of years of education or experience for the individuals in the sample. The former approach might be better if time itself was being used as the variable of interest, and one wished to check each and every year from, say, 1930 through 1980 for knots. Of course, if year is the variable of interest and there is one and only one observation per year, then these two approaches provide the same set of potential knot locations. In that case, the number of potential knot locations would be equal to the sample size. As noted above, if the X variable has repeated values, a subset of unique values can be obtained by using the PROC SORT; BY X; statements and the IF LAST.X THEN OUTPUT; statement.

Once the subset of unique, potential knot location values has been found, then PROC TRANSPOSE can be used to create the corresponding number of dummy variables needed to identify each potential knot location for PROC STEPWISE. Alternatively, the transpose function in PROC MATRIX could be used for this purpose.

The large number of dummy variables thus created can then be used to create the corresponding large number of linear, quadratic, and cubic spline terms (and higher order terms if desired). To do this efficiently, array statements must be used since hundreds of variables are to be created. Since the length of these arrays is not known in advance, a method must be devised to create arrays of unknown length. This must be done in such a way that the string of variables thus created can be referred to without knowing the number of variables involved. This may be demonstrated using the SAS spline regression program that follows.

After using PROC TRANSPOSE to transpose the subset of potential knot location values, the single observation that results has an unknown number of newly created variables. That unknown number is, of course, the number of potential knot location values and corresponds to the number of unique values of X in the data set. If one wished to search over a prespecified set of values instead of using the observed X values, a variable Z containing those values could simply be loaded into the data set called REDUCED just prior to the first PROC TRANSPOSE.

In order to be able to refer to a string of variables in an ARRAY statement, we need to know the name of the first variable and the last variable. The first variable name is simple enough. It is the prefix with the number 1 attached, which is KNOT1 in this case. The last variable is the problem. There is an unknown number of variables here. We don't know what the name of that last variable will be, so we just throw in a variable whose name we do know. By creating a single additional variable at this

point (e.g. KNOTEND=1), the string of variables KNOT1--KNOTEND may be referred to without knowing the number of variables in the string.

In a similar manner the corresponding dummy variables and linear, quadratic and cubic variables can be referenced by array statements without knowing their length. The initial values loaded into these array data sets are meaningless except for the KNOTS data set. The other data sets are needed only to create the variable names. The proper values are assigned to these variables in the final DO loop.

Once all of the appropriate spline regression variables have thus been created, PROC STEPWISE can be used to pick out the knots that are statistically significant. The line segments that are fitted in this manner may represent a combination of linear, quadratic and cubic equations that may be linked or unlinked at the knots. Linked line segments may merely involve the equality of the functions at the knots or may involve the equality of some of the derivatives as well. Any combination is possible ranging from the completely unrestricted dummy variable model with many segments to the single polynomial function model with no knots.

SAS SPLINE REGRESSION PROGRAM

```
DATA XY; INPUT X Y @@; N+1; CARDS;
*** data cards ***
PROC MEANS NOPRINT; VAR N;
OUTPUT OUT=NUMBER MAX=NCOUNT;

PROC SORT DATA=XY; BY X;
DATA REDUCED; SET; BY X;
IF LAST.X THEN OUTPUT; KEEP X;

PROC TRANSPOSE DATA=REDUCED PREFIX=KNOT;
DATA KNOTS; SET; KNOTEND=1;
PROC TRANSPOSE DATA=REDUCED PREFIX=D;
DATA DS; SET; DEND=1;
PROC TRANSPOSE DATA=REDUCED PREFIX=L;
DATA ARRAYL; SET; LEND=1;
PROC TRANSPOSE DATA=REDUCED PREFIX=Q;
DATA ARRAYQ; SET; QEND=1;
PROC TRANSPOSE DATA=REDUCED PREFIX=C;
DATA ARRAYC; SET; CEND=1;

DATA MATCH;
MERGE KNOTS DS ARRAYL ARRAYQ ARRAYC NUMBER;
DO I=1 TO NCOUNT; OUTPUT; END;
DROP _NAME_ I NCOUNT;

DATA; MERGE XY MATCH;
ARRAY KNOT KNOT1--KNOTEND;
ARRAY D D1--DEND;
ARRAY L L1--LEND;
ARRAY Q Q1--QEND;
ARRAY C C1--CEND;
DO OVER KNOT;
IF X LT KNOT THEN D=0;
IF X GE KNOT THEN D=1;
L=D*(X-KNOT);
Q=D*(X-KNOT)**2;
C=D*(X-KNOT)**3;
END; X2=X**2; X3=X**3;

PROC STEPWISE; MODEL Y=X X2 X3 L1--LEND Q1--QEND
C1--CEND / INCLUDE=1 SLE=1 SLS=1;
```

IV. An Interest Rate Application

Estimating a spline model for interest rates on commercial bonds provides a convenient example of determining the number and location of spline knots. The New York City open market rates for four-to-six month commercial paper with Aa rating or equivalent will be used for the period 1890 through 1981.

PROC STEPWISE selected terms for the spline regression model that attained a level of significance of at least .01 or less. The Y variable is the interest rate and the X variable is the year. Variables beginning with the letter C represent cubic adjustments, those with Q represent quadratic adjustments, and L stands for linear adjustments. The number following these letters indicates the knot location and corresponds to the "i" subscript in equations (1) and (5). Table I displays the spline regression model results.

Table I.

Interest rate as a function of time in years, X.

| Variable Name | Estimated Regression Coefficient | Student t Statistic | Prob Value |
|---------------|----------------------------------|---------------------|------------|
| INTERCEPT | 144.47767186 | 9.0990 | .0001 |
| X | -.07291068 | -8.8004 | .0001 |
| C36 | -.00068578 | -5.7497 | .0001 |
| Q50 | .08271411 | 5.7691 | .0001 |
| L63 | -.64603088 | -2.9252 | .0044 |
| L89 | 2.74175672 | 6.6394 | .0001 |

N=92 $R^2=.8574$ $\bar{R}^2=.8491$ F=103.4 F Prob .0001

In general terms these results provide the following functional form:

$$(6) y = B_{00} + B_{01} X + B_{36} (X - t_{36})^3 D36 +$$

$$B_{50} (X - t_{50})^2 D50 + B_{63} (X - t_{63}) D63 +$$

$$B_{89} (X - t_{89}) D89 + e$$

where B_{00} is the intercept term, B_{01} is the coefficient of X, B_{36} is the C36 coefficient, B_{50} is the Q50 coefficient, B_{63} is the L63 coefficient, and B_{89} is the L89 coefficient.

Substituting in for the estimated coefficient values and knot locations results in the following fitted values for Y:

$$(7) y = 144 - .073 X - .0007 (X - 1925)^3 D36 + .083 (X - 1939)^2 D50 - .646 (X - 1952) D63 + 2.74 (X - 1978) D89$$

Note that observation 36 is 1925, observation 50 is 1939, observation 63 is 1952, and observation 89 is 1978. Five separate equations represent the five time period segments found by collecting terms on X in equation (7).

(8) 1890 - 1924:

$$y = 144 - .073 X$$

(9) 1925 - 1938:

$$y = 4,892,038 - 7623.8 X + 3.96 X^2 - .0007 X^3$$

(10) 1939 - 1951:

$$y = 5,203,020 - 7944.6 X + 4.04 X^2 - .0007 X^3$$

(11) 1952 - 1977:

$$y = 5,204,281 - 7945.2 X + 4.04 X^2 - .0007 X^3$$

(12) 1978 - 1981:

$$y = 5,198,858 - 7942.5 X + 4.04 X^2 - .0007 X^3$$

The initial linear relationship becomes cubic in 1925. The intercept term and the coefficient for X adjust dramatically to compensate for the introduction of the cubic and quadratic terms. These relationships can be seen in Figures I and II.

The first, second, and third derivatives can be derived from equation (6) as follows:

$$(13) \frac{dy}{dX} = B_{01} + 3 B_{36} (X - t_{36})^2 D36 +$$

$$2 B_{50} (X - t_{50}) D50 + B_{63} D63 + B_{89} D89$$

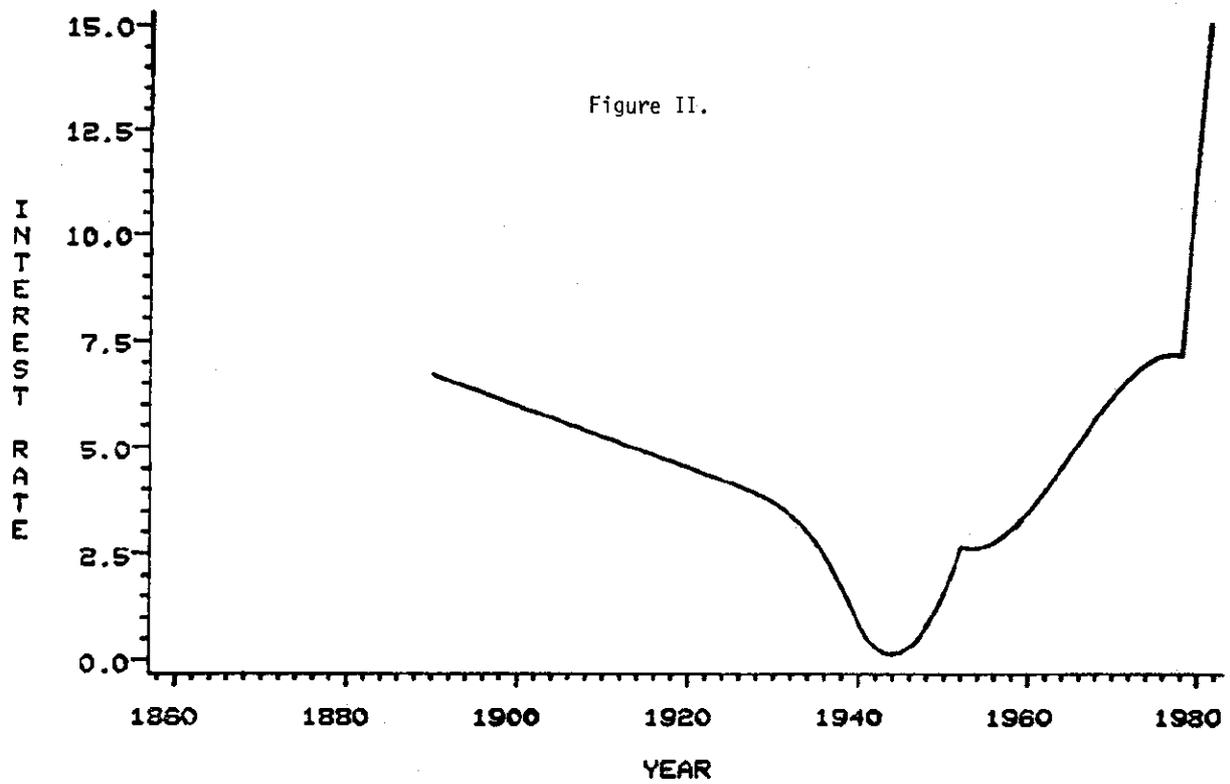
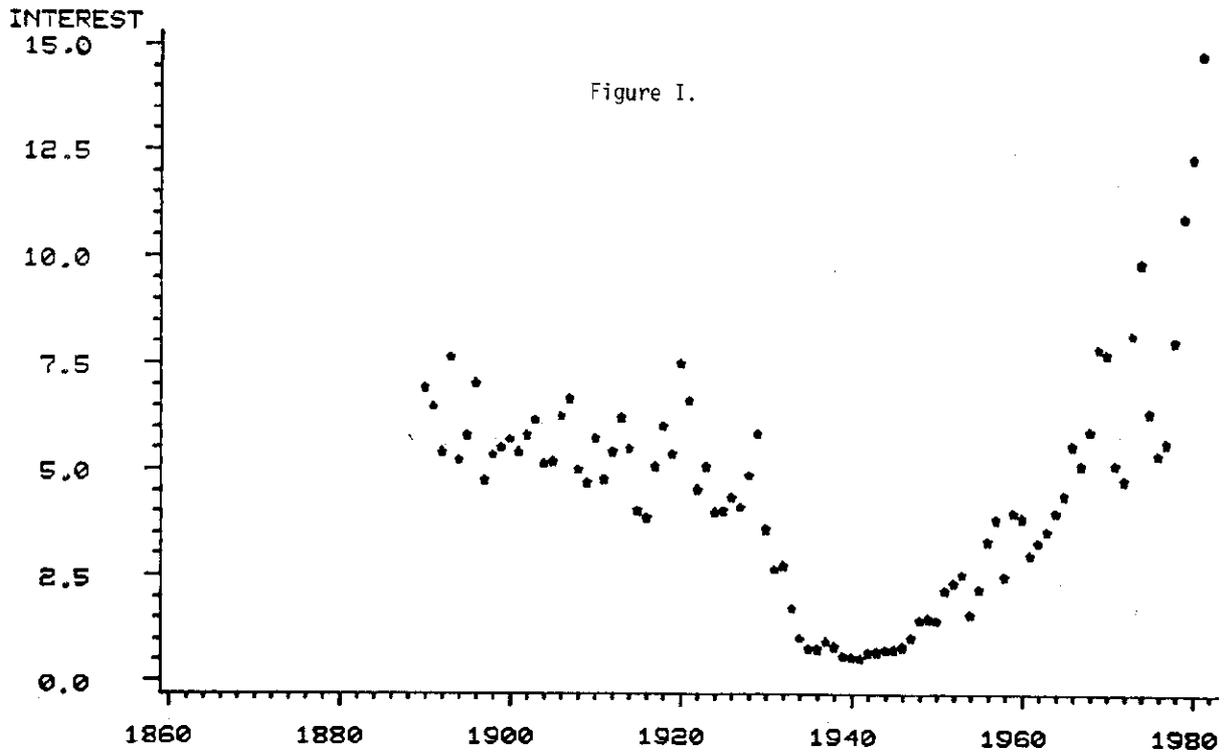
$$(14) \frac{d^2y}{dX^2} = 6 B_{36} (X - t_{36}) D36 + 2 B_{50} D50$$

$$(15) \frac{d^3y}{dX^3} = 6 B_{36} D36$$

The functions themselves as well as their first and second derivatives are equal at the 1925 knot. The third derivatives are not equal at this first join point.

A quadratic adjustment takes place in 1939. At this second join point the first and third derivatives are equal but not the second derivatives. This allows for an interesting minimizing loop in the predicted interest rates around 1940. In 1952 a discontinuity takes place in the first derivative due to a linear adjustment. However, the functions and their corresponding second and third derivatives are equal at the

INTEREST RATE ON COMMERCIAL BONDS



1952 knot.

Another final linear adjustment takes place in 1978 which substantially straightens out the fitted relationships. The general effect on the derivatives of this linear adjustment is by necessity the same as the 1952 linear adjustment. In other words, linear adjustments result in discontinuities only for the first derivatives. Similarly, quadratic adjustments cause discontinuities only for the second derivatives while cubic adjustments only generate discontinuities for the third derivatives.

V. Conclusions

This paper provided a method for estimating the number and location of spline knots (join points) for spline regression models. A SAS program for carrying out this estimation has been included. This program makes special use of PROC TRANSPOSE and PROC STEPWISE as well as the DO OVER capability associated with ARRAY statements. There is no particular limit on the number of cases that this program can handle, but since SAS has generally limited the number of variables that can be accessed at any one time to four thousand, a cubic spline model should not try to search over much more than about a thousand possible knot locations. In any case, a search over such a large number of potential knot locations would require a very large REGION parameter size for CPU space. Hopefully, most problems will require a search of only a few hundred knot locations or less.

The interest rate example has shown the power of this technique for fitting spline regressions. An analysis of the derivatives shows how tightly or loosely the polynomial line segments are connected at each join point. The interest rate for prime commercial paper provided some dramatic changes that demonstrated the need for the flexibility of the polynomial spline fitting regression technique. An intuitive interpretation of these interest rate changes is left for the reader and his/her financial advisor. Bear in mind, however, that these are nominal interest rates and real interest rates might tell an entirely different story.

VI. Bibliography

- Bellman, Richard and Robert Roth, "Curve Fitting by Segmented Straight Lines", Journal of the American Statistical Association, vol. 64, no. 327, September 1969, pages 1079-1084.
- Bookstein, Fred L., "On a Form of Piecewise Linear Regression", The American Statistician, vol. 29, no. 3, August 1975, pages 116-117.
- Brunelle, Rocco L. and David W. Johnson, "The Use of a Linear Spline Model in the Analysis of a Repeated Measure Experiment Through SAS", SUGI, vol. 5, 1980, pages 236-240.
- Buse, A. and L. Lim, "Cubic Splines as a Special Case of Restricted Least Squares", Journal of the American Statistical Association, vol. 72, no. 357, March 1977, pages 64-68.
- Capizzi, Thomas and Robert D. Small, "Using Spline Functions and the Bootstrap to Fit Differential Equation Models to Data", SUGI, vol. 7, 1982, pages 580-585.
- Ertel, J.E. and E.B. Fowlkes, "Methods for Fitting Linear Spline and Piecewise Multiple Linear Regression", Proceedings of Computer Science and Statistics, vol. 8, 1975, pages 222-227.
- Fuller, Wayne A., "Grafted Polynomials as Approximating Functions", Australian Journal of Agricultural Economics, June 1969, pages 35-46.
- Gallant, A.R. and W.A. Fuller, "Fitting Segmented Polynomial Regression Models Whose Join Points Have to be Estimated", Journal of the American Statistical Association, vol. 68, 1973, pages 144-147.
- Hudson, Derek J., "Fitted Segmented Curves Whose Join Points Have to be Estimated", Journal of the American Statistical Association, vol. 61, December 1966, pages 1097-1129.
- Irvine, John Michael, Changes in Regime in Regression Analysis, Ph.D. Dissertation, Yale University, May 1982.
- _____, "Testing for Changes in Regime in Regression Models", working paper, ASA/Census Fellowship Program, Washington, D.C., 1982.
- Mehta, Hina and Thomas Capizzi, "Evaluating Linear Model Representations of Cubic Splines Using Proc Reg", SUGI, vol. 7, 1982, pages 562-567.
- Poirier, Dale J., "Piecewise Regression Using Cubic Splines", Journal of the American Statistical Association, vol. 68, September 1973, pages 515-524.
- _____, "On the Use of Bilinear Splines in Economics", Journal of Econometrics, vol. 3, February 1975, pages 23-24.
- _____, The Econometrics of Structural Change with Special Emphasis on Spline Functions, Amsterdam: North Holland Publishing Company, 1976.
- Quandt, Richard E., "The Estimation of the Parameters of a Linear Regression System Obeying Two Separate Regimes", Journal of the American Statistical Association, vol. 53, 1958, pages 873-880.
- _____, "Tests of the Hypothesis that a Linear Regression System Obey Two Separate Regimes", Journal of the American Statistical Association, vol. 55, 1960, pages 324-330.
- Robb, A. Leslie, "Accounting for Seasonality with Spline Functions", Review of Economics and Statistics, vol. 62, no. 2, May 1980, pages 321-323.
- Robison, D.E., "Estimates for the Points of Intersection of Two Polynomial Regressions", Journal of the American Statistical Association, vol. 59, 1964, pages 214-224.
- Smith, Patricia L., "Splines as a Useful and Convenient Statistical Tool", The American Statistician, vol. 33, no. 2, May 1979, pages 57-62.
- Suits, Daniel B., Andrew Mason and Louis Chan, "Spline Functions Fitted by Standard Regression Methods", Review of Economics and Statistics, vol. 60, no. 1, February 1978, pages 132-139.
- Tishler, Asher and Israel Zang, "A Maximum Likelihood Method for Piecewise Regression Models with a Continuous Dependent Variable", Applied Statistics, vol. 30, no. 2, 1981, pages 116-124.
- _____, "A New Maximum Likelihood Algorithm for Piecewise Regression", Journal of the American Statistical Association, vol. 76, no. 376, December 1981, pages 980-987.
- _____, "A Switching Regression Method Using Inequality Conditions", Journal of Econometrics, vol. 11, 1979, pages 259-274.