# Penalized Splines with an Application in Economics

**Chaojiang Wu**
*Drexel University, USA*

## INTRODUCTION

Nonparametric smoothing has become the norm of statistics research in the last few decades. The advantage of nonparametric smoothing is its data driven nature. It does not assume any predetermined functional form but allow data to speak for themselves. One particularly popular nonparametric smoothing technique is penalized splines or P-splines. This technique is popular partly because it is computationally expedient and flexible enough to be adapted in complex models. Despite of its popularity in statistics methodological research, its applications in other areas such as finance and economics are limited. One reason for such limited applications is that P-splines remain unknown to most economics researchers. This chapter introduces the general idea of penalized splines and illustrates its application in forecasting US real GNP. Some references and resources are pointed out for future researches.

## BACKGROUND

Nonparametric smoothing or nonparametric regression is a technique for fitting curves, where the form of the curve is not predetermined but estimated through data. In its simplest form, given data pairs ($x_i$, $y_i$), $i = 1,\ldots, n$, the classical simple linear regression model is

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i,$$

where $\beta_0$ is the intercept and $\beta_1$ is the slope parameter, $\varepsilon_i$ is error term, usually assumed to be i.i.d. $N(0,\sigma^2)$. In the above model, a predetermined linear functional form describes the dependence of response variable $y$ on predictor $x$; the parameters $\beta_0$ and $\beta_1$ are to be estimated. However, in many cases, it is inappropriate to fit a straight line. Alternatively, response variable y takes the following form

$$y_i = g\left(x_i\right) + \varepsilon_i,$$

where $g(\cdot)$ is an unspecified function and $\varepsilon_i$ is again i.i.d. $N(0,\sigma^2)$. No parameters are explicitly present in the model. In other words, we are trying to estimate the function rather than estimating the parameters in any predetermined form.

The most popular ways of estimating $g(\cdot)$ are local polynomial methods (local methods) and spline methods (global methods). Some excellent references on local methods can be found in Wand and Jones (1995) and Fan and Gijbels (1996). The method introduced in this chapter (P-spline) is one type of global methods. There are generally three types of spline methods: smoothing splines, regression splines and penalized splines. Penalized splines or P-splines can be viewed as a generalization of smoothing splines with a more flexible choice of bases and penalties. Alternatively, P-splines can be viewed as least square regression splines with a roughness penalty. Some discussions on smoothing spline methods can be seen in Wahba (1990) and Eubank (1988). General discussions on penalized splines can be found in Ruppert, Wand and Carroll (2003). Applications of penalized splines have gained popularity only recently. Some examples can be found in Jarrow, Ruppert and Yu (2004) on estimating the term

structure of corporate bonds, Brunauer, Lang, Wechselberger and Bienert (2010) on estimating the nonlinear effects and smooth time trend in modeling the rent in Vienna, and Çetinyürek and Lambert (2011) on estimating the survival function and hazard ratios for interval-censored data in Medicine.

## MAIN FOCUS

This section introduces P-spline estimation ideas, computational issues (bases choice, smoothing parameter selection and computing resource), and illustrate the method by one application in forecasting real United States GNP.

## Penalized Splines Estimation

The basic idea of splines estimation of an unknown smooth function *g(x)* is to express the smooth function as a linear combination of piece-wise polynomial functions, or splines. For mathematical convenience, the function is expressed as a linear combination of the spline basis. That is

$$g(x) = \beta_0 + \beta_1 x + \cdots + \beta_p x^p + \sum_{k=1}^{K} \beta_{p+k} (x - x_k)_+^p,$$

where *p* is the degree of spline (polynomials), $(x - x_k)_+^p = \max\{(x - x_k)^p, 0\}$, and $\beta = (\beta_0, \ldots, \beta_{p+k})^T$ is the spline coefficients to be estimated. Define $X = \left(1, x, \ldots, x^p, (x - x_1)_+^p, \ldots, (x - x_k)_+^p\right)$ as the spline basis (truncated power basis in this case), with break points at $x_k$, *k=1,...,K*. *X* is the "design matrix" similar to that of linear regressions. Then the unknown smooth function can be written as $g(x) = X\beta$.

The series of break points $x_1 < \ldots < x_k$ are known as knots of the splines. The choice of number (*K*) and locations ($x_k$) of knots is crucial

and usually the process is complicated and computationally intensive. One school of thought is the regression splines, which select the knots and locations carefully. Such knots selection usually requires sophisticated algorithm such as Multivariate Adaptive Regression Splines (MARS) algorithm (Friedman, 1991), which selects knots by stepwise type regressions. Smoothing splines, on the other hand, place knots at each distinct observation of *x*. Smoothing splines are theoretically appealing because of its mathematical nice properties. Specifically, the natural cubic spline with knots placed in distinct observations of *x* is the solution to the minimizer of roughness penalized least square deviance, over the class of twice differentiable functions. However, the computation of smoothing splines is intensive for such knots placement, especially when the application is complex beyond univariate smoothing. In between the two types of splines, penalized splines use sufficiently large number of knots, often equally placed in the equidistant quantiles of *x* and shrink unnecessary knots by adding penalties (Eilers & Marx, 1996). The number of knots is no longer crucial as long as the smoothing parameter is properly chosen (Ruppert, 2002).

Consider finding the smoothing estimator for *g(x)=Xβ* using spline model with *K* knots (*K* sufficiently large, say 40), one way of formulating the fit is

$$\underset{\beta}{argmin} \left\| y - X\beta \right\|^2 + \lambda \sum_{k=1}^{K} \beta_{p+k}^2,$$

where $y = (y_1, y_2, \ldots, y_n)^T$ is the vector of response variable, and *λ>0* is the smoothing parameter that controls the smoothness of the fitting curve. The first term of the objective function is the usual Ordinary Least Squares (OLS) goodness-of-fit term while the second term is the roughness penalty term which constrains the influence of unnecessary knots. The smoothing parameter *λ* balances the goodness-of-fit and roughness of the curve. Very large λ results in a linear least square

7 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: www.igi-global.com/chapter/penalized-splines-with-an-application-in-economics/107371

# Related Content

### Neighborhood Evaluation in Recommender Systems Using the Realization Based Entropy Approach

Roee Anuar, Yossi Bukchin, Oded Maimonand Lior Rokach (2014). *International Journal of Business Analytics (pp. 34-50).*
www.irma-international.org/article/neighborhood-evaluation-in-recommender-systems-using-the-realization-based-entropy-approach/119496

### Segmenting Big Data Time Series Stream Data

Dima Albergand Zohar Laslo (2014). *Encyclopedia of Business Analytics and Optimization (pp. 2126-2134).*
www.irma-international.org/chapter/segmenting-big-data-time-series-stream-data/107399

### Intelligent IoT-Enabled System in Green Supply Chain using Integrated FCM Method

Rui-Yang Chen (2015). *International Journal of Business Analytics (pp. 47-66).*
www.irma-international.org/article/intelligent-iot-enabled-system-in-green-supply-chain-using-integrated-fcm-method/126833

### Integrating Ontologies and Bayesian Networks in Big Data Analysis

Hadrian Peterand Charles Greenidge (2014). *Encyclopedia of Business Analytics and Optimization (pp. 1254-1261).*
www.irma-international.org/chapter/integrating-ontologies-and-bayesian-networks-in-big-data-analysis/107323

### Business Intelligence 2.0: The eXtensible Markup Language as Strategic Enabler

Rubén A. Mendoza (2010). *International Journal of Business Intelligence Research (pp. 63-76).*
www.irma-international.org/article/business-intelligence-extensible-markup-language/47196