

DEPARTMENT OF STATISTICS

University of Wisconsin

1210 West Dayton St.

Madison, WI 53706

TECHNICAL REPORT NO. 942

January 1995

**GRKPACK: Fitting Smoothing Spline ANOVA  
Models for Exponential Families**

by

**Yuedong Wang**

# GRKPACK: Fitting Smoothing Spline ANOVA Models for Exponential Families

Yuedong Wang\*

*Department of Biostatistics, University of Michigan,  
Ann Arbor, Michigan 48109, U.S.A.*

January 17, 1995

## Abstract

*Wahba et al (1994c) introduced Smoothing Spline ANalysis of VAriance (SS ANOVA) method for data from exponential families. Based on RKPACK, which fits SS ANOVA models to Gaussian data, we introduce GRKPACK: a collection of subroutines for binary, binomial, Poisson and Gamma data. We also show how to calculate Bayesian confidence intervals for SS ANOVA estimates.*

*Key Words:* generalized cross validation; Newton-Raphson iteration; RKPACK; smoothing parameter; smoothing spline ANOVA; unbiased risk estimate.

## 1 Introduction

Generalized linear models (GLM's) for analysis of data from exponential families have been extensively studied and widely used since 1970's (Nelder and Wedderburn, 1972; McCullagh and Nelder, 1989). As the popularity of these methods has increased, so has the need for more sophisticated model building and diagnostic checking techniques. In the context of nonparametric estimate of the GLM regression surface, O'Sullivan et al (1986) and Gu (1990) used penalized likelihood method with smoothing splines and thin plate splines. Hastie and Tibshirani (1990) used additive models. Wahba et al (1994c) introduced the SS ANOVA models using the penalized likelihood and Smoothing Spline ANalysis of Variance methods. See also Wahba et al (1994a, 1994b), Wang (1994) and Wang et al (1995) for details of SS ANOVA models. In this paper, we describe a package for estimations of the SS ANOVA models with binary, binomial, Poisson and Gamma data. We call this package as GRKPACK, which stands for generalized RKPACK.

---

\*Supported by the National Institute of Health under Grants R01 EY09946, P60 DK20572 and P30 HD18258

First, we describe the computational part of the SS ANOVA model. Suppose data have the form  $(y_i; \mathbf{t}_i)$ ,  $i = 1, 2, \dots, n$ , where  $y_i$  are independent observations and  $\mathbf{t}_i = (t_{1i}, \dots, t_{di})$ . The distribution function of  $y_i$  is from an exponential family with density function

$$g(y_i; f_i, \phi) = \exp((y_i h(f_i) - b(f_i))/a(\phi) + c(y_i, \phi)), \quad (1)$$

where  $f_i = f(\mathbf{t}_i)$  is the parameter of interest and  $h(f_i)$  is a monotone transformation of  $f_i$  known as the *canonical parameter*.  $\phi$  is an unknown scale parameter. Let  $\mathbf{t} = (t_1, \dots, t_d)$ , and let  $t_j \in \mathcal{T}^{(j)}$ , where  $\mathcal{T}^{(j)}$  is a measurable space. Let  $\mathcal{T} = \mathcal{T}^{(1)} \otimes \dots \otimes \mathcal{T}^{(d)}$ , then  $\mathbf{t} \in \mathcal{T}$ . Denote the log likelihood given  $y_i$  and  $\mathbf{t}_i$  as

$$l_i(f_i) = \log g(y_i; f_i, \phi) = (y_i h(f_i) - b(f_i))/a(\phi) + c(y_i, \phi). \quad (2)$$

The purpose is to investigate the global relationship between  $f$  and  $\mathbf{t}$ .

Let  $d\mu_j$  be a probability measure on  $\mathcal{T}^{(j)}$  and let  $\mathcal{H}^{(j)}$  be a reproducing kernel Hilbert space (RKHS) (Aronszajn, 1950) of functions on  $\mathcal{T}^{(j)}$  with  $\int_{\mathcal{T}^{(j)}} f_j(t_j) d\mu_j = 0$  for  $f_j(t_j) \in \mathcal{H}^{(j)}$ . Let  $\{1^{(j)}\}$  be the one dimensional space of constant functions on  $\mathcal{T}^{(j)}$ . Consider the RKHS

$$\begin{aligned} \mathcal{G} &= \prod_{j=1}^d (1^{(j)} \oplus \mathcal{H}^{(j)}) \\ &= \{1\} \oplus \sum_j \mathcal{H}^{(j)} \oplus \sum_{j < k} (\mathcal{H}^{(j)} \otimes \mathcal{H}^{(k)}) \oplus \dots, \end{aligned} \quad (3)$$

where  $\{1\}$  denotes the constant functions on  $\mathcal{T}$ . An element  $f_j$  in  $\mathcal{H}^{(j)}$  is called a *main effect*, an element  $f_{jk}$  in  $\mathcal{H}^{(j)} \otimes \mathcal{H}^{(k)}$  is called a *two factor interaction*, and so on.

Similar to the usual ANOVA, model space is a subspace  $\mathcal{M}$  of  $\mathcal{G}$ . By deleting some higher-order interactions, we get less flexible, but more “estimable” models. When a model is chosen, we can regroup and write the model space as

$$\mathcal{M} = \mathcal{H}^0 \oplus \sum_{\beta=1}^q \mathcal{H}^\beta, \quad (4)$$

where  $\mathcal{H}^0$  is a finite dimensional space containing functions which are not going to be penalized, usually lower order polynomials. See Wahba et al (1994c) for details. A SS ANOVA estimate is the solution to the following variational problem:

$$\min_{f \in \mathcal{M}} \left\{ - \sum_{i=1}^n l_i(f_i) + \frac{n}{2} \sum_{\beta=1}^q \lambda_\beta \|P_\beta f\|^2 \right\}. \quad (5)$$

The first part in (5) is the negative log likelihood. It measures the goodness of fit. In the second part,  $P_\beta$  is the orthogonal projector in  $\mathcal{M}$  onto  $\mathcal{H}^\beta$  and  $\|P_\beta f\|^2$  is a quadratic roughness penalty.  $\lambda_\beta$ 's are a set of smoothing parameters. They control the trade-off

between the goodness of fit and the roughness of the estimate. Writing  $\lambda_\beta = \lambda/\theta_\beta$ , (5) becomes

$$\min_{f \in \mathcal{M}} \left\{ -\sum_{i=1}^n l_i(f_i) + \frac{n}{2} \lambda \|P_* f\|_\Theta^2 \right\}, \quad (6)$$

where  $P_* = \sum_{\beta=1}^q P_\beta$  is the orthogonal projection in  $\mathcal{M}$  onto  $\mathcal{H}_* = \sum_{\beta=1}^q \mathcal{H}^\beta$  and

$$\|f\|_\Theta^2 = \|P_0 f\|^2 + \sum_{\beta=1}^q \theta_\beta^{-1} \|P_\beta f\|^2, \quad (7)$$

is a modified norm indexed by  $\Theta = (\theta_1, \dots, \theta_q)$ . We denote  $R_\beta$  as the reproducing kernel (RK) (Aronszajn, 1950) for  $\mathcal{H}^\beta$  under the original norm. The RK for  $\sum_{\beta=1}^q \mathcal{H}^\beta$  under  $\|\cdot\|_\Theta$  is

$$R_\Theta = \sum_{\beta=1}^q \theta_\beta R_\beta. \quad (8)$$

The solution to (6) has the form (Wahba 1990, O'Sullivan et al. 1986)

$$f_{\lambda, \Theta}(\mathbf{t}) = \sum_{v=1}^M d_v \phi_v(\mathbf{t}) + \sum_{i=1}^n c_i \left( \sum_{\beta=1}^q \theta_\beta R_\beta(\mathbf{t}_i, \mathbf{t}) \right) = \boldsymbol{\phi}(\mathbf{t})^T \mathbf{d} + \boldsymbol{\xi}(\mathbf{t})^T \mathbf{c}, \quad (9)$$

where  $\{\phi_v\}_{v=1}^M$  is a set of basis functions of  $\mathcal{H}^0$ ,  $M = \dim(\mathcal{H}^0)$ ,  $\boldsymbol{\phi}^T(\mathbf{t}) = (\phi_1(\mathbf{t}), \dots, \phi_M(\mathbf{t}))$ ,  $\boldsymbol{\xi}^T(\mathbf{t}) = (R_\Theta(\mathbf{t}_1, \mathbf{t}), \dots, R_\Theta(\mathbf{t}_n, \mathbf{t}))$ .  $\mathbf{c}_{n \times 1}$  and  $\mathbf{d}_{M \times 1}$  are vectors of coefficients to be estimated. Substituting (9) into (6), we can estimate  $\mathbf{c}$  and  $\mathbf{d}$  by minimizing

$$I(\mathbf{c}, \mathbf{d}) = -\sum_{i=1}^n l_i(\boldsymbol{\phi}^T(\mathbf{t}_i) \mathbf{d} + \boldsymbol{\xi}^T(\mathbf{t}_i) \mathbf{c}) + \frac{n}{2} \lambda \mathbf{c}^T Q_\Theta \mathbf{c}, \quad (10)$$

where  $Q_\Theta = \sum_{\beta=1}^q \theta_\beta Q_\beta$  and  $Q_\beta$ 's are  $n \times n$  matrices with  $Q_\beta(i, k) = R_\beta(\mathbf{t}_i, \mathbf{t}_k)$ . Since  $l_i$ 's are not quadratic, (10) can not be solved directly. But if all  $l_i(f_i)$ 's are strictly concave, we can use Newton-Raphson procedure to compute  $\mathbf{c}$  and  $\mathbf{d}$  for fixed  $\lambda$  and  $\Theta$ . Let  $u_i = -dl_i/df_i$ ,  $w_i = -d^2 l_i/df_i^2$ . Let

$$\mathbf{u}^T = (u_1, \dots, u_n), \quad (11)$$

$$W = \text{diag}(w_1, \dots, w_n), \quad (12)$$

$$S = (\boldsymbol{\phi}(\mathbf{t}_1), \dots, \boldsymbol{\phi}(\mathbf{t}_n))^T. \quad (13)$$

Then

$$\partial I / \partial \mathbf{c} = Q_\Theta \mathbf{u} + n \lambda Q_\Theta \mathbf{c}, \quad (14)$$

$$\partial I / \partial \mathbf{d} = S^T \mathbf{u}, \quad (15)$$

$$\partial^2 I / \partial \mathbf{c} \partial \mathbf{c}^T = Q_\Theta W Q_\Theta + n \lambda Q_\Theta, \quad (16)$$

$$\partial^2 I / \partial \mathbf{c} \partial \mathbf{d}^T = Q_\Theta W S, \quad (17)$$

$$\partial^2 I / \partial \mathbf{d} \partial \mathbf{d}^T = S^T W S. \quad (18)$$

The Newton-Raphson iteration satisfies the linear system

$$\begin{pmatrix} Q_{\Theta}W_{-}Q_{\Theta} + n\lambda Q_{\Theta} & Q_{\Theta}W_{-}S \\ S^TW_{-}Q_{\Theta} & S^TW_{-}S \end{pmatrix} \begin{pmatrix} \mathbf{c} - \mathbf{c}_{-} \\ \mathbf{d} - \mathbf{d}_{-} \end{pmatrix} = \begin{pmatrix} -Q_{\Theta}\mathbf{u}_{-} - n\lambda Q_{\Theta}\mathbf{c}_{-} \\ -S^T\mathbf{u}_{-} \end{pmatrix}, \quad (19)$$

where the subscript minus indicates quantities evaluated at the previous Newton-Raphson iteration. Denote

$$\mathbf{f} = S\mathbf{d} + Q_{\Theta}\mathbf{c} \quad (20)$$

as the vector of estimates of  $f$  at design points. Let

$$\tilde{Q}_{\Theta} = W_{-}^{1/2}Q_{\Theta}W_{-}^{1/2}, \quad (21)$$

$$\tilde{\mathbf{c}} = W_{-}^{-1/2}\mathbf{c}, \quad (22)$$

$$\tilde{S} = W_{-}^{1/2}S, \quad (23)$$

$$\tilde{\mathbf{d}} = \mathbf{d}, \quad (24)$$

$$\tilde{\mathbf{y}} = W_{-}^{-1/2}(W_{-}\mathbf{f}_{-} - \mathbf{u}_{-}). \quad (25)$$

(19) can be simplified to

$$\begin{aligned} (\tilde{Q}_{\Theta} + n\lambda I)\tilde{\mathbf{c}} + \tilde{S}\tilde{\mathbf{d}} &= \tilde{\mathbf{y}}, \\ \tilde{S}^T\tilde{\mathbf{c}} &= \mathbf{0}. \end{aligned} \quad (26)$$

Choosing appropriate smoothing parameters is crucial for effectively estimating the true function from data by fitting smoothing spline models. The generalized cross validation (GCV) method estimates smoothing parameters by minimizing the GCV score

$$V(\lambda, \Theta) = \frac{1/n\|(I - A(\lambda, \Theta))\tilde{\mathbf{y}}\|^2}{[(1/n)\text{tr}(I - A(\lambda, \Theta))]^2}, \quad (27)$$

where  $A(\lambda, \Theta)$  satisfies

$$(w_{1-}^{1/2}f_{\lambda, \Theta}(\mathbf{t}_1), \dots, w_{n-}^{1/2}f_{\lambda, \Theta}(\mathbf{t}_n))^T = A(\lambda, \Theta)\tilde{\mathbf{y}}, \quad (28)$$

and  $f_{\lambda, \Theta}(\mathbf{t}_i)$ 's are computed from the solution of (26).

The unbiased risk (UBR) method estimates smoothing parameters by minimizing the following unbiased risk estimate

$$\tilde{U}(\lambda, \Theta) = \frac{1}{n}\|(I - A(\lambda, \Theta))\tilde{\mathbf{y}}\|^2 + 2\frac{\hat{\sigma}^2}{n}\text{tr}A(\lambda, \Theta), \quad (29)$$

where  $\hat{\sigma}^2 = 1/n\sum_{i=1}^n u_{i-}^2/w_{i-}$ , an estimate of dispersion parameter. If the dispersion parameter is known to be 1, such as in the case of binary data and Poisson data, a better estimate is

$$U(\lambda, \Theta) = \frac{1}{n}\|(I - A(\lambda, \Theta))\tilde{\mathbf{y}}\|^2 + \frac{2}{n}\text{tr}A(\lambda, \Theta). \quad (30)$$

See Wang (1994) and Wang et al (1995) for detail discussions about the GCV and UBR methods.

## 2 The Algorithm

A generic code RKPACk (Gu, 1989; Gu and Wahba, 1991) is available to solve (26) and estimate  $\lambda$  and  $\Theta$  via GCV (option V) or the UBR method at the same time. When using the UBR method, we can either specify  $\sigma^2 = 1$  (option U) or estimate  $\sigma^2$  (option U<sup>~</sup>). This suggests the following algorithm:

**Algorithm.** Given the matrices  $S$ ,  $Q_\beta$ 's, the response vector  $\mathbf{y}$  and the starting vector  $\mathbf{f}_0$ :

1. Compute  $\mathbf{u}_-$  and  $W_-$ . Compute the transformations  $\tilde{S}$ ,  $\tilde{Q}_\beta = W_-^{1/2} Q_\beta W_-^{1/2}$  and  $\tilde{\mathbf{y}}$ ;
2. Call RKPACk with inputs  $\tilde{S}$ ,  $\{\tilde{Q}_\beta, \beta = 1, \dots, q\}$  and  $\tilde{\mathbf{y}}$ . That is, solve (26) and choose  $\lambda$  and  $\Theta$  by GCV (option V) or the UBR method (option U or option U<sup>~</sup>);
3. Compute the new  $\mathbf{f}$ . Stop if the algorithm converges under some criteria (for example,  $\sum_{i=1}^n w_{i-} ((f_i - f_{i-}) / (1 + |f_i|))^2 / \sum_{i=1}^n w_{i-} < p$  for some prespecified  $p > 0$  is used in our programs) or the number of iterations exceeds some prespecified number  $L$ ; otherwise go to step 1.

The starting value  $\mathbf{f}_0$  may be a constant function, a GLM fit or some other estimates. We usually let  $p = 10^{-6}$ . The algorithm usually takes 5 to 15 iterations to converge. We believe  $L = 30$  is big enough for most applications. Since changing  $\lambda$  and  $\Theta$  at each iteration means modifying the problem successively, convergence is not guaranteed. Nevertheless, the algorithm converges most of the time.

## 3 Approximate Bayesian Confidence Intervals for SS ANOVA Estimates

In an observation study, often design is not balanced. It is desirable to construct confidence intervals for the SS ANOVA estimate and decide a region in which the estimates are deemed to be reliable. Confidence intervals for components like main effects and interactions are also useful for model selection. Wahba et al (1994c) derived approximate Bayesian confidence intervals for SS ANOVA estimates and showed how to use them in practice.

Let the prior for  $f(\mathbf{t})$  be

$$F_\xi(\mathbf{t}) = \sum_{\nu=1}^M \tau_\nu \phi_\nu(\mathbf{t}) + b^{\frac{1}{2}} \sum_{\beta=1}^q \sqrt{\theta_\beta} Z_\beta(\mathbf{t}), \quad (31)$$

where  $\boldsymbol{\tau} = (\tau_1, \dots, \tau_M)^T \sim N(0, \xi I)$ ,  $Z_\beta$  are independent, zero mean Gaussian stochastic processes, independent of  $\boldsymbol{\tau}$ , with  $E Z_\beta(\mathbf{t}) Z_\beta(\mathbf{z}) = R_\beta(\mathbf{t}, \mathbf{z})$ . With the log likelihood in (2) and  $\xi \rightarrow \infty$ , Wahba et al (1994c) derived the approximate posterior mean and covariance. We list them in the following theorem.

**Theorem 1** Let  $\mathbf{c}$  and  $\mathbf{d}$  be a solution to (10). Let  $\lambda, \Theta$  be smoothing parameters at convergence and  $W, Q_\Theta$  be matrix based on converged values. Let  $n\lambda = \hat{\sigma}^2/b$  and  $M = Q_\Theta + n\lambda W^{-1}$ . For  $g_{0,\nu}(\mathbf{t}) = \tau_\nu \phi_\nu(\mathbf{t}), g_\beta(\mathbf{t}) = b^{\frac{1}{2}} \sqrt{\theta_\beta} Z_\beta(\mathbf{t}), \nu = 1, \dots, M, \beta = 1, \dots, q$ , we have

$$\begin{aligned} E(g_{0,\nu}(\mathbf{t})|\mathbf{y}) &\approx d_\nu \phi_\nu(\mathbf{t}), \\ E(g_\beta(\mathbf{t})|\mathbf{y}) &\approx \sum_{i=1}^n c_i \theta_\beta R_\beta(\mathbf{t}, \mathbf{t}_i), \\ \frac{1}{b} \text{Cov}(g_{0,\nu}(\mathbf{z}), g_{0,\mu}(\mathbf{t})|\mathbf{y}) &\approx \phi_\nu(\mathbf{z}) \phi_\mu(\mathbf{t}) \mathbf{e}_\nu^T (S^T M^{-1} S)^{-1} \mathbf{e}_\mu, \\ \frac{1}{b} \text{Cov}(g_\beta(\mathbf{z}), g_{0,\nu}(\mathbf{t})|\mathbf{y}) &\approx -d_{\nu,\beta}(\mathbf{z}) \phi_\nu(\mathbf{t}), \\ \frac{1}{b} \text{Cov}(g_\beta(\mathbf{z}), g_\beta(\mathbf{t})|\mathbf{y}) &\approx \theta_\beta R_\beta(\mathbf{z}, \mathbf{t}) - \sum_{i=1}^n c_{i,\beta}(\mathbf{z}) \theta_\beta R_\beta(\mathbf{t}, \mathbf{t}_i), \\ \frac{1}{b} \text{Cov}(g_\gamma(\mathbf{z}), g_\beta(\mathbf{t})|\mathbf{y}) &\approx -\sum_{i=1}^n c_{i,\gamma}(\mathbf{z}) \theta_\beta R_\beta(\mathbf{t}, \mathbf{t}_i), \end{aligned}$$

where  $\mathbf{e}_\nu$  is the  $\nu$ th unit vector, and  $(d_{1,\beta}(\mathbf{z}), \dots, d_{M,\beta}(\mathbf{z})) = d_\beta(\mathbf{z})^T$  and  $(c_{1,\beta}(\mathbf{z}), \dots, c_{n,\beta}(\mathbf{z})) = c_\beta(\mathbf{z})^T$  are given by

$$d_\beta(\mathbf{z}) = (S^T M^{-1} S)^{-1} S^T M^{-1} \begin{pmatrix} \theta_\beta R_\beta(\mathbf{z}, \mathbf{t}_1) \\ \vdots \\ \theta_\beta R_\beta(\mathbf{z}, \mathbf{t}_n) \end{pmatrix}, \quad (32)$$

$$c_\beta(\mathbf{z}) = [M^{-1} - M^{-1} S (S^T M^{-1} S)^{-1} S^T M^{-1}] \begin{pmatrix} \theta_\beta R_\beta(\mathbf{z}, \mathbf{t}_1) \\ \vdots \\ \theta_\beta R_\beta(\mathbf{z}, \mathbf{t}_n) \end{pmatrix}. \quad (33)$$

## 4 A Special Case: Tensor Product of $W_2$

To illustrate how to use SS ANOVA method and how to construct Bayesian confidence intervals, consider the special case  $\mathcal{T}^{(j)} = [0, 1], j = 1, \dots, d$  and  $d \geq 2$  (we usually transform all continuous variables into  $[0, 1]$  for fitting and then transform them back). Take the component space on  $[0, 1]$  as the RKHS

$$W_2 = \{f : f \text{ and } f^{(1)} \text{ abs. cont.}, \int_0^1 (f^{(2)})^2 < \infty\} \quad (34)$$

with norm

$$f^2 = \left(\int_0^1 f\right)^2 + \left(\int_0^1 f^{(1)}\right)^2 + \int_0^1 (f^{(2)})^2. \quad (35)$$

We decompose  $W_2 = \mathcal{N} \oplus \mathcal{L} \oplus \mathcal{S}$ , where  $\mathcal{N}$  is the space of constants, with the square norm  $(\int_0^1 f)^2$ ;  $\mathcal{L}$  is the space of linear functions which integrate to zero, with the square norm

$(\int_0^1 f^{(1)})^2$ ; and  $\mathcal{S}$  is the space of functions with square integrable 2nd derivative and satisfy  $\int_0^1 f^{(v)} = 0$ ,  $v = 0, 1$ , with the square norm  $\int_0^1 (f^{(2)})^2$ . The RK for subspace  $\mathcal{S}$  is

$$R(t, z) = k_2(t)k_2(z) - k_4(t - z), \quad (36)$$

where  $k_\nu(\cdot) = B_\nu(\cdot)/\nu!$  and  $B_\nu(\cdot)$  is the  $\nu$ th Bernoulli polynomial. Let  $\mathcal{G}$  be the tensor product of component spaces

$$\begin{aligned} \mathcal{G} &= \otimes_{j=1}^d (\mathcal{N}^j \oplus \mathcal{L}^j \oplus \mathcal{S}^j) \\ &= \{1\} \oplus \{(\oplus_{j=1}^d \mathcal{L}^j) \oplus (\oplus_{j=1}^d \mathcal{S}^j)\} \\ &\quad \oplus \{(\oplus_{j < k} (\mathcal{L}^j \otimes \mathcal{L}^k)) \oplus (\oplus_{j \neq k} (\mathcal{L}^j \otimes \mathcal{S}^k)) \oplus (\oplus_{j < k} (\mathcal{S}^j \otimes \mathcal{S}^k))\} \oplus \dots, \end{aligned}$$

where with some abuse of notation, we are omitting factors of the form  $\otimes \{\mathcal{N}^j\}$ ;  $\{1\}$  is the space of constant functions on  $[0, 1]^d$ ;  $\mathcal{L}^j = \mathcal{N}^1 \otimes \dots \otimes \mathcal{L}^j \otimes \dots \otimes \mathcal{N}^d$  is the space of functions that is a constant on  $t_k$ ,  $k \neq j$  and a linear function on  $t_j$ . Others have similar interpretation. The terms in the three brackets are the spaces of the constant, the main effects and the 2-interactions respectively. For simplicity of notation (calculations can be easily extended to having more than one interactions), suppose we choose a model space contains the constant, all main effects and the interaction between  $t_1$  and  $t_2$ :

$$\mathcal{M} = \{1\} \oplus \{(\oplus_{j=1}^d (\mathcal{L}^j \oplus \mathcal{S}^j))\} \oplus \{(\mathcal{L}^1 \otimes \mathcal{L}^2) \oplus (\mathcal{S}^1 \otimes \mathcal{L}^2) \oplus (\mathcal{L}^1 \otimes \mathcal{S}^2) \oplus (\mathcal{S}^1 \otimes \mathcal{S}^2)\} \quad (37)$$

Therefore, we have  $M = d + 2$ ,  $q = d + 3$ . We usually take

$$\phi_1(\mathbf{t}) = 1, \quad (38)$$

$$\phi_\nu(\mathbf{t}) = t_{\nu-1} - 0.5, \quad \nu = 2, \dots, d + 1, \quad (39)$$

$$\phi_M(\mathbf{t}) = (t_1 - 0.5) \times (t_2 - 0.5) \quad (40)$$

as basis functions for  $\mathcal{H}^0$  and

$$R_\beta(\mathbf{t}, \mathbf{z}) = R(t_\beta, z_\beta), \quad \beta = 1, \dots, d, \quad (41)$$

$$R_{d+1}(\mathbf{t}, \mathbf{z}) = R(t_1, z_1) \times (t_2 - 0.5) \times (z_2 - 0.5), \quad (42)$$

$$R_{d+2}(\mathbf{t}, \mathbf{z}) = (t_1 - 0.5) \times (z_1 - 0.5) \times R(t_2, z_2), \quad (43)$$

$$R_{d+3}(\mathbf{t}, \mathbf{z}) = R(t_1, z_1) \times R(t_2, z_2) \quad (44)$$

as the RK's for  $\mathcal{H}^1, \dots, \mathcal{H}^{d+3}$  respectively.

Write

$$f(\mathbf{t}) = C + f_1(t_1) + \dots + f_d(t_d) + f_{1,2}(t_1, t_2), \quad (45)$$

where  $f_j$ 's are the main effects and  $f_{1,2}$  is the interaction between  $t_1$  and  $t_2$ . Comparing to Theorem 1, we have

$$C = g_{01}, \quad (46)$$

$$f_j(t_j) = g_{0,j+1}(\mathbf{t}) + g_j(\mathbf{t}), \quad j = 1, \dots, d, \quad (47)$$

$$f_{1,2}(t_1, t_2) = g_{0,M}(\mathbf{t}) + g_{d+1}(\mathbf{t}) + g_{d+2}(\mathbf{t}) + g_{d+3}(\mathbf{t}). \quad (48)$$



Suppose we want to calculate posterior means and standard deviations of the main effects, the interaction and the overall function on grid points  $\mathcal{Z}_1 \times \cdots \times \mathcal{Z}_d$ , where  $\mathcal{Z}_j$ 's are sets of points in  $[0, 1]$ . For simplicity of notation, suppose  $\mathcal{Z}_j = \mathcal{Z} = \{z_1, \dots, z_K\}$ . Calculations for different  $\mathcal{Z}_j$ 's are the same. Let

$$\mathbf{f}_j^T = (f_j(z_1), \dots, f_j(z_K)), \quad j = 1, \dots, d, \quad (49)$$

$$\mathbf{f}_{1,2}^T = (f_{1,2}(z_1, z_1), \dots, f_{1,2}(z_K, z_1), \dots, f_{1,2}(z_1, z_K), \dots, f_{1,2}(z_K, z_K)), \quad (50)$$

$$A = (S^T M^{-1} S)^{-1}, \quad (51)$$

$$\boldsymbol{\phi}_\nu^T = (\phi_\nu(z_1), \dots, \phi_\nu(z_K)), \quad \nu = 2, \dots, d+1, \quad (52)$$

$$\boldsymbol{\phi}_{1,2}^T = (\phi_M(z_1, z_1), \dots, \phi_M(z_K, z_1), \dots, \phi_M(z_1, z_K), \dots, \phi_M(z_K, z_K)), \quad (53)$$

$$\mathbf{d}_j^T = (d_{j+1,j}(z_1), \dots, d_{j+1,j}(z_K)), \quad j = 1, \dots, d, \quad (54)$$

$$\begin{aligned} \mathbf{d}_{1,2}^T &= (d_{M,d+1}(z_1, z_1), \dots, d_{M,d+1}(z_K, z_1), \dots, d_{M,d+1}(z_1, z_K), \dots, d_{M,d+1}(z_K, z_K)) \\ &+ (d_{M,d+2}(z_1, z_1), \dots, d_{M,d+2}(z_K, z_1), \dots, d_{M,d+2}(z_1, z_K), \dots, d_{M,d+2}(z_K, z_K)) \\ &+ (d_{M,d+3}(z_1, z_1), \dots, d_{M,d+3}(z_K, z_1), \dots, d_{M,d+3}(z_1, z_K), \dots, d_{M,d+3}(z_K, z_K)), \end{aligned} \quad (55)$$

$$C_\beta = (c_\beta(z_1), \dots, c_\beta(z_K))_{n \times K}, \quad \beta = 1, \dots, d, \quad (56)$$

$$\begin{aligned} C_{1,2} &= (c_{d+1}(z_1, z_1), \dots, c_{d+1}(z_K, z_1), \dots, c_{d+1}(z_1, z_K), \dots, c_{d+1}(z_K, z_K))_{n \times K^2} \\ &+ (c_{d+2}(z_1, z_1), \dots, c_{d+2}(z_K, z_1), \dots, c_{d+2}(z_1, z_K), \dots, c_{d+2}(z_K, z_K))_{n \times K^2} \\ &+ (c_{d+3}(z_1, z_1), \dots, c_{d+3}(z_K, z_1), \dots, c_{d+3}(z_1, z_K), \dots, c_{d+3}(z_K, z_K))_{n \times K^2} \end{aligned} \quad (57)$$

$$\Lambda_\beta = (\theta_\beta R_\beta(z_k, t_{\beta i}))_{K \times n}, \quad \beta = 1, \dots, d, \quad (58)$$

$$\Lambda_{1,2} = ((\theta_{d+1} R_{d+1} + \theta_{d+2} R_{d+2} + \theta_{d+3} R_{d+3})((z_{k_1}, z_{k_2}), (t_{1i}, t_{2i})))_{K^2 \times n}, \quad (59)$$

$$\Sigma_\beta = (\theta_\beta R_\beta(z_i, z_k))_{K \times K}, \quad \beta = 1, \dots, d, \quad (60)$$

$$\Sigma_{1,2} = ((\theta_{d+1} R_{d+1} + \theta_{d+2} R_{d+2} + \theta_{d+3} R_{d+3})((z_{i1}, z_{i2}), (z_{k1}, z_{k2})))_{K^2 \times K^2}. \quad (61)$$

From Theorem 1, the posterior means and standard deviations for the constant, the main effects at points  $\mathcal{Z}$  and interaction at points  $\mathcal{Z} \times \mathcal{Z}$  are

$$E(C|\mathbf{y}) \approx d_1, \quad (62)$$

$$\text{Cov}(C|\mathbf{y}) \approx bA(1, 1), \quad (63)$$

$$E(\mathbf{f}_j|\mathbf{y}) \approx d_{j+1}\boldsymbol{\phi}_{j+1} + \Lambda_j \mathbf{c}, \quad j = 1, \dots, d, \quad (64)$$

$$\begin{aligned} \text{Cov}(\mathbf{f}_j|\mathbf{y}) &\approx b[A(j+1, j+1)\boldsymbol{\phi}_{j+1}\boldsymbol{\phi}_{j+1}^T - \mathbf{d}_j\boldsymbol{\phi}_{j+1}^T - \boldsymbol{\phi}_{j+1}\mathbf{d}_j^T + \Sigma_j - \Lambda_j C_j], \\ & \quad j = 1, \dots, d, \end{aligned} \quad (65)$$

$$E(\mathbf{f}_{1,2}|\mathbf{y}) \approx d_M\boldsymbol{\phi}_{1,2} + \Lambda_{1,2}\mathbf{c}, \quad (66)$$

$$\text{Cov}(\mathbf{f}_{1,2}|\mathbf{y}) \approx b[A(M, M)\boldsymbol{\phi}_{1,2}\boldsymbol{\phi}_{1,2}^T - \mathbf{d}_{1,2}\boldsymbol{\phi}_{1,2}^T - \boldsymbol{\phi}_{1,2}\mathbf{d}_{1,2}^T + \Sigma_{1,2} - \Lambda_{1,2}C_{1,2}]. \quad (67)$$

For any point  $\mathbf{z} \in \mathcal{Z}^d$ , the posterior mean of the overall function

$$E(f(\mathbf{z})|\mathbf{y}) = E(C|\mathbf{y}) + \sum_{j=1}^d E(f_j(\mathbf{z}|\mathbf{y})) + E(f_{1,2}|\mathbf{y}). \quad (68)$$

Table 1: Drivers in GRKPACK

data	single smoothing parameter drivers	multiple smoothing parameter drivers
binary	dbedr	dbedr
binomial	dbedr	dbedr
Poisson	dpsdr	dpedr
Gamma	dgsdr	dgedr

The posterior variance of the overall function can be calculated using the formula

$$\begin{aligned}
 \text{Var}(C + \sum_{j=1}^d f_j(\mathbf{z}) + f_{1,2}(\mathbf{z})|\mathbf{y}) = & \\
 & \text{Var}(C|\mathbf{y}) + \sum_{j=1}^d \text{Var}(f_j(\mathbf{z})|\mathbf{y}) + \text{Var}(f_{1,2}(\mathbf{z})|\mathbf{y}) \\
 & + 2 \sum_{j=1}^d \text{Cov}(C, f_j(\mathbf{z})|\mathbf{y}) + 2\text{Cov}(C, f_{1,2}(\mathbf{z})|\mathbf{y}) + 2 \sum_{j=1}^d \text{Cov}(f_j(\mathbf{z}), f_{1,2}(\mathbf{z})|\mathbf{y}). \quad (69)
 \end{aligned}$$

Each term in (69) can be read off easily from the calculations for the component posterior variances.

We need to calculate  $(S^T M^{-1} S)^{-1}$ ,  $c_\beta(\mathbf{t})$  and  $d_\beta(\mathbf{t})$ . Gu and Wahba (1993) discussed how to calculate these quantities when  $W = I$ . Let  $\tilde{Q}_\Theta = W^{1/2} Q_\Theta W^{1/2}$ ,  $\tilde{S} = W^{1/2} S$ , and  $\tilde{M} = \tilde{Q}_\Theta + n\lambda I$ . We can calculate  $(\tilde{S}^T \tilde{M}^{-1} \tilde{S})^{-1}$ ,  $\tilde{d}_\beta(\mathbf{t})$  and  $\tilde{c}_\beta(\mathbf{t})$  the same way as Gu and Wahba (1993) with their  $R_\beta(\mathbf{t}, \mathbf{t}_i)$  replaced by  $\tilde{R}_\beta(\mathbf{t}, \mathbf{t}_i) = \sqrt{w_i} R_\beta(\mathbf{t}, \mathbf{t}_i)$ . We then have  $(S^T M^{-1} S)^{-1} = (\tilde{S}^T \tilde{M}^{-1} \tilde{S})^{-1}$ ,  $d_\beta(\mathbf{t}) = \tilde{d}_\beta(\mathbf{t})$  and  $c_\beta(\mathbf{t}) = W^{1/2} \tilde{c}_\beta(\mathbf{t})$ . Two utility routines `dcrdr`, `dsms` in RKPACK can be used to calculate these quantities.

The  $(1-\alpha)100\%$  Bayesian confidence interval for a function  $g$  ( $g$  could be the constant  $C$ , the main effects  $f_j$ 's, the interaction  $f_{1,2}$  and the overall function  $f$ ) at a given point  $\mathbf{z}$  is

$$(E(g(\mathbf{z})|\mathbf{y}) - z_{\alpha/2} \sqrt{\text{Var}(g(\mathbf{z})|\mathbf{y})}, E(g(\mathbf{z})|\mathbf{y}) + z_{\alpha/2} \sqrt{\text{Var}(g(\mathbf{z})|\mathbf{y})}), \quad (70)$$

where  $z_{\alpha/2}$  is the upper  $\alpha/2$  percentile of the standard Normal distribution.

## 5 GRKPACK

GRKPACK is a collection of subroutines using the SS ANOVA algorithm in section 2 for binary, binomial, Poisson and Gamma data. Users can modify these routines for other distributions in exponential families. We are developing routines for ordinal data.

Subroutines in GRKPACK are listed in Table 1. Correspondence of notations between this report and GRKPACK are listed in Table 2.

Table 2: GRKPACK notation correspondence

report	code	type	explanation
<b>options for smoothing parameters</b>			
option V	vmu='v'	character	GCV
option U	vmu='u'	character	UBR with dispersion parameter $\sigma^2 = 1$
option U~	vmu='u~'	character	UBR and estimate the dispersion parameter
<b>model parameters</b>			
$M$	nnull	integer	dimension of the null space $\mathcal{H}^0$
$q$	nq	integer	number of RKHS's to be penalized
<b>design, RK's and their transformations</b>			
$S$	s	matrix	the design matrix of $\mathcal{H}^0$
$Q_\beta$	q( $\cdot, \cdot, \beta$ )	matrix	RK's evaluated at design points
$\tilde{S}$	swk	matrix	$\tilde{S} = W^{1/2}S$
$\tilde{Q}_\beta$	qwk( $\cdot, \cdot, \beta$ )	matrix	$\tilde{Q}_\beta = W^{1/2}Q_\beta W^{1/2}$
<b>precision parameters</b>			
$p$	prec	scale	stop criteria for the algorithm
$L$	maxiter	integer	maximum number of iterations for the algorithm
<b>data and their transformations</b>			
$n$	nobs	integer	number of observations
$\mathbf{y}$	y	vector	observation vector
$\tilde{\mathbf{y}}$	ywk	vector	$\tilde{\mathbf{y}} = W^{-1/2}(W\mathbf{f} - \mathbf{u})$
<b>coefficients, estimates and derivatives</b>			
$\mathbf{c}$	c	vector	$\mathbf{c}$ vector in (9)
$\mathbf{d}$	d	vector	$\mathbf{d}$ vector in (9)
$\mathbf{f}$	eta	vector	estimate at design points
$\mathbf{u}$	u	vector	first derivative vector at design points
$W$	w	vector	second derivatives, $W = \text{diag}(w(1), \dots, w(n))$
<b>summary information</b>			
$\lambda$	nlaht	scale	main smoothing parameter, nlaht = $\log_{10}(n\lambda)$
$\Theta$	theta	vector	subsidiary smoothing parameters in $\log_{10}$ scale
$V U \tilde{U}$	score	scale	minimum GCV and UBR score
$\hat{\sigma}^2$	varht	scale	estimate of dispersion parameter

Table 3: Notation correspondence for calculations of posterior means and variances

report	code	type	explanation
<b>data and prediction points</b>			
$d$	ncov	integer	number of continuous covariates
$j$	cov(j)	integer	cov(j) is the column in data corresponding to $t_j$
$t_{ji}$	x(i,cov(j))	scale	jth covariate at the ith design point
$K$	nplot	integer	number of points in $\mathcal{Z}$
$\mathcal{Z}$	z	vector	points for calculation of predictions
<b>posterior means and variances</b>			
$b$	b	scale	$b = \hat{\sigma}^2/n\lambda$
$E(C \mathbf{y})$	d(1)	scale	posterior mean of the constant
$\text{Var}(C \mathbf{y})$	b × sms(1,1)	scale	posterior variance of the constant
$E(\mathbf{f}_j \mathbf{y})$	main(·,j)	vector	posterior mean vector of the $t_j$ main effect
$\text{Cov}(\mathbf{f}_j \mathbf{y})$	mainsd(·,j)	vector	posterior variances of the $t_j$ main effect
$E(\mathbf{f}_{1,2} \mathbf{y})$	inter	matrix	posterior mean of the interaction between $t_1$ and $t_2$
$\text{Cov}(\mathbf{f}_{1,2} \mathbf{y})$	intersd	matrix	posterior variances of the interaction
<b>matrices for calculations</b>			
$(S^T M^{-1} S)^{-1}$	sms	matrix	see (51)
$\Lambda_\beta^T$	r(·,·,β+ncov)	matrix	see (58)
$W^{1/2} \Lambda_\beta^T$	r(·,·,β)	matrix	transform of (58)
$\Lambda_{1,2}^T$	rr(·,·,2)	matrix	see (59)
$W^{1/2} \Lambda_{1,2}^T$	rr(·,·,1)	matrix	transform of (59)
$C_\beta$	cr(·,·,β)	matrix	see (56)
$C_{1,2}$	crr(·,·)	matrix	see (57)
$\mathbf{d}_j$	dr(j+1,·,j)	vector	see (54)
$\mathbf{d}_{1,2}$	drr(M,·)	vector	see (55)

GRKPACK also contains some application and simulation programs we used in Wahba et al (1994a) and Wang et al (1994). Correspondence of notations between the report and calculations of Bayesian confidence intervals in these programs are listed in Table 3. Notice that  $mainsd(\cdot, j)$  contains the diagonal elements of  $\text{Cov}(\mathbf{f}_j|\mathbf{y})$ ;  $inter$  is a  $K \times K$  matrix with  $inter(i, j) = E(\mathbf{f}_{1,2}|\mathbf{y})((i-1) \times K + j)$ ;  $intersd$  is a  $K \times K$  matrix with  $intersd(i, j) = \text{Cov}(\mathbf{f}_{1,2}|\mathbf{y})((i-1) \times K + j, (i-1) \times K + j)$ .

All subroutines in RKPACk and GRKPACK are self-documented. RKPACk is available from [netlib@research.att.com](mailto:netlib@research.att.com) and [statlib@lib.stat.cmu.edu](mailto:statlib@lib.stat.cmu.edu). GRKPACK is available from the author at [yuedong@umich.edu](mailto:yuedong@umich.edu) or anonymous ftp to <ftp.stat.wisc.edu> and go to directory /pub/wahba/software. We are going to submit them to [netlib@research.att.com](mailto:netlib@research.att.com) and [statlib@lib.stat.cmu.edu](mailto:statlib@lib.stat.cmu.edu) later. Any suggestions are appreciated and should be for-

warded to the author.

## 6 Acknowledgements

The author is indebted to Grace Wahba, Douglas Bates and Chong Gu for their advice and encouragement. The author thanks Dong Xiang and Zhen Luo for testing out the routines and suggestions.

## References

- [1] Aronszajn, N. (1950). “Theory of reproducing kernels”. *Trans. Amer. Math. Soc.*, 68, 337–404.
- [2] Gu, C. (1989). “RKPACK and its applications: Fitting smoothing spline models”. *Proceedings of the Statistical Computing Section, ASA*, 42–51.
- [3] Gu, C. (1990). “Adaptive spline smoothing in non-Gaussian regression models”. *J. Amer. Stat. Asso.*, 85, 801–807.
- [4] Gu, C. and Wahba, G. (1991). “Minimizing GCV/GML scores with multiple smoothing parameters via the Newton method”. *SIAM J. Sci. Stat. Comput.*, 12, 383–398.
- [5] Gu, C. and Wahba, G. (1993). “Smoothing spline ANOVA with component-wise Bayesian confidence intervals”. *Journal of Computational and Graphical Statistics*, 2, 97–117.
- [6] Hastie, T. and Tibshirani, R. (1990). *Generalized Additive Models*. Chapman and Hall.
- [7] McCullagh, P. and Nelder, J. A. (1989). *Generalized Linear Models*. Chapman and Hall.
- [8] Nelder, J. A. and Wedderburn, R. W. M. (1972). “Generalized linear interactive models”. *J. Royal Stat. Soc., Ser. B*, 135, 370–384.
- [9] O’Sullivan, F., Yandell, B. and Raynor, W. (1986). “Automatic smoothing of regression functions in generalized linear models”. *J. Amer. Stat. Assoc.*, 81, 96–103.
- [10] Wahba, G. (1990). *Spline Models for Observational Data*. SIAM, Vol. 59.
- [11] Wahba, G., Gu, C., Wang, Y. and Chappell, R. (1994a). “Soft classification, a.k.a. risk estimation, via penalized log likelihood and smoothing spline analysis of variance”. *The Mathematics of Generalization, Santa Fe Institute Studies in the Science of Complexity*, Vol. XX, 329–360. D. Wolpert eds, Addison Wesley.
- [12] Wahba, G., Wang, Y., Gu, C., Klein, R. and Klein, B. (1994b). “Structured machine learning for ‘soft’ classification with smoothing spline ANOVA and stacked turning, testing and evaluation”. *Advances in Neural Information Processing*, 6, 415–422. Cowan, J.D., Tesauro, G. and Alspector, J. eds. Morgan Kaufmann Publishers.

- [13] Wahba, G., Wang, Y., Gu, C., Klein, R. and Klein, B. (1994c). “Smoothing spline ANOVA for exponential families, with application to Wisconsin Epidemiological study of diabetic retinopathy”. Technical Report 940, University of Wisconsin-Madison, Dept. of Statistics.
- [14] Wang, Y., Wahba, G.,(1994). ”Bootstrap Confidence Intervals for Smoothing Splines and Their Comparison to Bayesian Confidence Intervals”. To appear in *J. of Stat. Comp. and Simu.*
- [15] Wang, Y. (1994). *Smoothing Spline Analysis of Variance of Data from Exponential Families*. Ph.D Thesis, University of Wisconsin-Madison, Dept. of Statistics. Technical Report 928, University of Wisconsin-Madison, Dept. of Statistics.
- [16] Wang, Y., Wahba, G., Chappell, R. and Gu, C. (1995). “Simulation studies of smoothing parameter estimates and Bayesian confidence intervals in Bernoulli SS ANOVA models”.