

# ADAPTIVE TUNING, 4D-VAR AND REPRESENTERS IN RKHS

Grace Wahba <sup>1</sup>  
 Department of Statistics  
 University of Wisconsin  
 Madison WI USA

Summary: We (abstractly) generalize the ‘toy’ weak 4D-Var model in Gong, Wahba, Johnson & Tribbia (1998) to include adaptive tuning of a variety of parameters throughout the 4D-Var variational problem, and note issues of sensitivity and identifiability. We discuss ‘models’ for model errors which include systematic, short memory and long memory errors. Finally we remark on the role of the theory of representers in reproducing kernel Hilbert spaces in the weak 4D-Var setting.

## 1 INTRODUCTION

We first consider the general setup in the experiment in Gong et al. (1998), which is a toy weak 4D-Var model (actually one time and one space variable) with five unknown smoothing, weighting and distributed parameters, which were simultaneously adaptively tuned using generalized cross validation (GCV) calculated via the randomized trace technique. In that setup ‘model error’ was generated as the difference between a ‘nature’ model and the ‘computer’ model, but white noise model errors were assumed in the weak 4D-Var variational problem. In this paper we then (i) review the use of model errors as dual variables, (ii) review the GCV and generalized maximum likelihood (GML) tuning methods, and pinpoint sensitivity issues as tunable parameters are sprinkled liberally throughout the weak 4D-Var problem, noting that they can be studied in the influence matrix (or influence operator in the nonlinear case). Then (iii) we describe some simple models for correlated model errors and the simultaneous consideration of systematic (bias), short memory and long memory correlation. We end with (iv) a summary of some representer theory in reproducing kernel Hilbert space (RKHS) relevant to the weak 4D-Var setting.

Let  $t = 1, \dots, T$  denote discrete time and let  $\Psi_t, t = 1, \dots, T$  be a sequence of state vectors representing (some part of) nature that evolves according to

$$\Psi_{t+1} = M_t \Psi_t + N_t + \xi_t, \quad t = 1, \dots, T-1, \quad (1)$$

where  $M_t$  is the model evolution operator,  $N_t$  is a forcing function, and the  $\xi_t$  represent model errors, which we will discuss in more detail later. Here  $\Psi_*$  is the forecast for  $t = 1$  assumed to satisfy

$$\Psi_1 = \Psi_* + \epsilon_* \quad (2)$$

with  $\epsilon_* \sim \mathcal{N}(0, \sigma_f^2 Q_*)$ .  $M_t = M_t(\theta_M), N_t = N_t(\theta_M)$  are assumed to contain some tunable distributed parameters  $\theta_M$ . The observations are

$$y_t = K_t \Psi_t + \epsilon_t, \quad t \in \Lambda, \quad (3)$$

with  $\epsilon_t \sim \mathcal{N}(0, \sigma_o^2 S_t)$ .  $K_t$  is a map from state vector space to observation space at time  $t \in \Lambda$ , and  $K_t = K_t(\theta_K)$  may also contain some tunable parameters, for example calibration coefficients or bias corrections. Here  $\Lambda$  is the set of observation times, which are assumed to be a subset of the model update times  $t = 1, \dots, T$ .

In Gong et al. (1998) a toy weak 4D-Var problem was formulated as: Find  $\hat{\Psi} = (\hat{\Psi}'_1, \dots, \hat{\Psi}'_T)'$ , to minimize

$$\frac{1}{\sigma_o^2} \sum_{t \in \Lambda} \|y_t - K_t \Psi_t\|_{S_t}^2 + \frac{1}{\sigma_m^2} \sum_{t=1}^{T-1} \|\Psi_{t+1} - M_t(\theta_M) \Psi_t - N_t(\theta_M)\|_{Q_t}^2 + \frac{1}{\sigma_f^2} \|\Psi_* - \Psi_1\|_{Q_*}^2 + \frac{1}{b} \|\Psi_T\|_J^2. \quad (4)$$

<sup>1</sup>Research supported in part by NASA Grant NAG5-3769 and NSF Grant DMS9121003

Here  $\|\mathbf{v}\|_C^2 = \mathbf{v}'\mathbf{C}\mathbf{v}$  for  $\mathbf{C}$  a non-negative definite matrix. This formulation corresponds to  $\xi_t \sim \mathcal{N}(0, \sigma_m^2 \mathbf{Q}_t)$ , independent from time to time, and a prior belief that  $\|\Psi_T\|_J^2$  is 'small' where  $\mathbf{J}$  is a quadratic penalty representing a toy version of e. g. a penalty for lack of balance. Letting  $\gamma = \sigma_o^2/\sigma_f^2$ ,  $\alpha = \sigma_o^2/\sigma_m^2$  and  $\eta = \sigma_o^2/b$ , then the minimizer of (4) is the same as the minimizer of

$$\begin{aligned} & \sum_{t \in \Lambda} \|\mathbf{y}_t - \mathbf{K}_t \Psi_t\|_{\mathbf{S}_t}^2 + \alpha \sum_{t=1}^{T-1} \|\Psi_{t+1} - \mathbf{M}_t(\theta_M) \Psi_t - \mathbf{N}_t(\theta_M)\|_{\mathbf{Q}_t}^2 + \gamma \|\Psi_* - \Psi_1\|_{\mathbf{Q}_*}^2 + \eta \|\Psi_T\|_J^2 \\ & = J_o + J_m + J_b + J_c, \end{aligned} \quad (5)$$

say. The  $J$  terms have their usual meaning as observation, model, forecast and constraint except that the coefficients in front have been scaled relative to observations. The major tuning parameters  $\{\gamma, \alpha, \eta\}$  and a two coefficient distributed parameter  $\theta_M$  in the model were simultaneously tuned by the GCV method, using the randomized trace technique, which allows the computation of the cross-validation function to be carried out by rerunning the model with perturbed data. 'Nature' was simulated using the barotropic vorticity equation on a latitude circle, solved with a high order method on a fine time and space grid, and noisy observations were generated from nature using a random number generator.  $\mathbf{S}_t$  was taken as  $\mathbf{I}$  both in generating the data and in analyzing it via (5), and  $\epsilon_*$  was generated as a zero mean random Gaussian vector with covariance a multiple of  $\mathbf{Q}_*$ . 'The model' was based on a cruder integration of the barotropic vorticity equation, so that 'model error' may be thought of as the difference between 'nature' and 'the model'. In this case, as in nature, model error is not readily describable in terms of means and covariances, nevertheless, they may be a convenient, although crude way of dealing with model error that is not well understood. In the experiment  $\mathbf{Q}_t$  was taken as the identity matrix. The five tuning parameters were selected for adaptive tuning via the GCV method because it was believed, as a result of some preliminary experiments, (plus guesswork) that the solution and the 'predicted' observations computed from the solution, were sensitive to them. It turned out that the predictive mean square error, based on  $\mathbf{y}_{true} - \mathbf{y}_{fitted}$ , where  $\mathbf{y}_{true}$  is what would have been observed if there were no errors anywhere, (known only in a simulation, of course), was sensitive to all five parameters, but, on comparable scales,  $\alpha$ , the parameter relating observational to model error, had a much broader, flat minimum. One possible explanation is that the 'white noise' assumption for model error was not a very good representation for model error. Although the sensitivity to  $\alpha$  was not great, the weak constraint estimate gave better results than the strong constraint ( $\alpha \rightarrow \infty$ ). Dealing with model error is an open scientific issue, according to Courtier (1997) and others. Before going on to some speculative discussion of approaches to model error in the weak 4D-Var problem, we note that it was clear in the experiments in Gong et al. (1998) that the five parameters being tuned interacted with one another. For example the optimal value of one of the physical parameters in  $\theta_M$  was systematically larger than the 'true' or 'nature' value but it depended on the choice of  $(\gamma, \alpha, \eta)$ .

With regard to  $\theta_M$ , in the present study,  $\theta_M$  contributed essentially two degrees of freedom to the fit. In practice,  $\theta_M$  may be widely distributed. If it contributes many degrees of freedom, then in general it will be appropriate to include a (tunable) penalty term, say  $\delta \|\theta_M\|_D^2$  to the variational problem, see Wahba (1990a), O'Sullivan (1991), Navon (1998), Evensen, Dee & Schroter (1998).

## 2 DUAL VARIABLES, NONLINEAR FORWARD OPERATORS AND MODELS, CORRELATED MODEL ERRORS

Following Bennett (1997), Courtier (1997), notice that, assuming that  $\Psi_*$  and the  $\mathbf{N}_t$  known, that  $\Psi_1, \dots, \Psi_T$  are determined by  $\epsilon_*$  and  $\xi = (\xi'_1, \dots, \xi'_{T-1})'$  and vice versa. Thus one may change

variables from  $\Psi_1, \dots, \Psi_T$  in (5) to  $(\epsilon_*, \xi)$  to get

$$\sum_{t \in \Lambda} \|\mathbf{y}_t - \mathbf{K}_t \Psi_t(\epsilon_*, \xi)\|_{S_t^{-1}}^2 + \alpha \sum_{t=1}^{T-1} \|\xi_t\|_{Q_t^{-1}}^2 + \gamma \|\epsilon_*\|_{Q_*^{-1}}^2 + \eta \|\Psi_T(\epsilon_*, \xi)\|_J^2, \quad (6)$$

and solve the variational problem for  $\epsilon_*$  and  $\xi$  instead of for  $\Psi$ . To take a closer look at this problem, let  $\tilde{K}_t = \mathbf{K}_t \Psi_t$ ,  $\mathbf{x} = (\epsilon_*', \xi)'$ , and redesign  $J_c$  so that it is quadratic in  $\mathbf{x}$ , letting  $J_c(\mathbf{x}) = \|\mathbf{x}\|_J^2$ ,  $J = J(\theta_c)$ . Concatenate the  $\mathbf{y}_t, \tilde{K}_t, S_t$  and  $Q_t$  in an obvious way, and allow more tuning parameters to get  $\tilde{K}(\theta_K)$ ,  $S = S(\theta_o)$  and  $Q = Q(\theta_\xi)$ , and let  $Q_* = Q_*(\theta_b)$ . Furthermore, we specifically do not want to restrict  $Q$  to be block diagonal, so that we can allow for model errors correlated from time to time. Let  $\lambda^{-1}\Sigma = \lambda^{-1}\Sigma(\theta_\Sigma)$  be a quadratic form standing in collectively for  $\alpha Q, \gamma Q_*$  and  $\eta J$ . The result is the variational problem

$$(\mathbf{y} - \tilde{K}\mathbf{x})'S^{-1}(\mathbf{y} - \tilde{K}\mathbf{x}) + \lambda\mathbf{x}'\Sigma^{-1}\mathbf{x}, \quad (7)$$

with tuning parameters  $\theta = (\theta_K, \theta_M, \theta_o, \theta_\xi, \theta_b, \theta_c)$ . We note that (7) is not changed if the  $\mathbf{K}_t$  and  $M_t$  are nonlinear, in that case  $\tilde{K}$  is a nonlinear map from  $\mathbf{x}$  to  $\mathbf{y}$ , but under the assumption about  $J_c$ , the second term is quadratic.

If  $\tilde{K} = K$ , meaning  $\tilde{K}$  is linear, then the minimizer  $\mathbf{x}_\lambda$  of (7) is<sup>2</sup>

$$\mathbf{x}_\lambda = (K'S^{-1}K + \lambda\Sigma^{-1})^{-1}K'S^{-1}\mathbf{y} = \Sigma K'(\Sigma K' + \lambda S)^{-1}\mathbf{y} = \sum c_i \eta_i, \quad (8)$$

where  $\eta_i$  is the  $i$ th column of  $\Sigma K'$ , and, letting  $\mathbf{c} = (c_1, \dots, c_n)'$ ,  $(\Sigma K' + \lambda S)\mathbf{c} = \mathbf{y}$ . This is a trivial example of representer theory (implemented, for example in PSAS, Cohn, daSilva, Guo, Sienkiewicz & Lamich (1998)), where a system of size the dimension of  $\mathbf{y}$  is to be solved even if the dimension of  $\mathbf{x}_\lambda$  is much bigger than the dimension of  $\mathbf{y}$ .

### 3 TUNING METHODS

With the dimension of  $\mathbf{x}$  orders of magnitude greater than the dimension of  $\mathbf{y}$ , it is clear that the number of tunable parameters in  $\tilde{K}, S$  and  $\Sigma$  is limited by the amount of information in  $\mathbf{y}$ , and by the possibility of aliasing/identifiability. (Of course to the extent that desirable values of these parameters do not vary over long periods of time, historical information may be collected). In any case, it is neither possible, nor even desirable to have models for the model error covariance matrix  $Q$ , the forecast error matrix  $Q_*$  or the constraint functional  $J_c$  to have an overabundance of free parameters. Furthermore, the solution should be sensitive to any parameters considered for adaptive tuning.

Equation (7) may be interpreted as the variational problem associated with the statistical assumptions

$$\mathbf{y} = K\mathbf{x} + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_o^2 S), \quad \mathbf{x} \sim \mathcal{N}(0, b\Sigma), \quad \lambda = \sigma_o^2/b. \quad (9)$$

The influence matrix  $A(\lambda, \theta)$ , which maps  $S^{-1/2}\mathbf{y}$  into  $S^{-1/2}K\mathbf{x}_\lambda$ <sup>3</sup> is given by

$$A = S^{-1/2}K(K'S^{-1}K + \lambda\Sigma^{-1})^{-1}K'S^{-1/2} = B(B + I)^{-1}, \quad (10)$$

<sup>2</sup>This formula assumes that the quadratic form  $\mathbf{x}'\Sigma^{-1}\mathbf{x}$  is of full rank. If it is not, see Kimeldorf & Wahba (1971), Wahba (1990b). Having already abused notation, replace  $\mathbf{x}'\Sigma^{-1}\mathbf{x}$  by  $\mathbf{x}'P\mathbf{x}$ . For (7) to have a unique minimizer in the linear case it is necessary and sufficient that  $\mathbf{x}'P\mathbf{x} = 0$  and  $K\mathbf{x} = 0$  imply that  $\mathbf{x} = 0$ . In practice users should make sure that the null space of  $P$  is not too big!

<sup>3</sup>In the nonlinear case use this to define the influence operator.

where  $B = B(\lambda, \theta) = \frac{1}{\lambda} S^{-1/2} K \Sigma K' S^{-1/2}$ . It's clear that  $B$  must be sensitive to a component of  $\theta$  in order for estimation of that component to make sense and aliasing of parameters inside  $B$  is to be avoided<sup>4</sup>. The GML estimate for  $\lambda, \theta$  is the minimizer of

$$M(\lambda, \theta) = \frac{y' S^{-1/2} (I - A(\lambda, \theta)) S^{-1/2} y}{[\det S^{-1}]^{1/n} [\det(I - A(\lambda, \theta))]^{1/n}} \quad (11)$$

and the GML estimate for  $\sigma_o^2$  is  $\frac{1}{n} y' S^{-1/2} (I - A) S^{-1/2} y$ . The GML estimate has certain optimality properties when the stochastic model (9) is correct up to the unknown parameters. Little is known concerning its robustness to this assumption, see Wahba (1985). Note that there are other forms of maximum likelihood estimates, depending on which unknowns are included and how they enter into the formulas. The GCV estimate is the minimizer of

$$V(\lambda, \theta) = \frac{\|(I - A) S^{-1/2} y\|^2}{(\frac{1}{n} \text{trace}(I - A))^2} \equiv \frac{\|y - Kx_\lambda\|_{S^{-1}}^2}{(\frac{1}{n} \text{trace}(I - A))^2} \quad (12)$$

In theory it may not be suitable for estimating sensitive parameters inside  $S$  (i.e.  $\theta_o$ ), since it is theoretically based on assuming that the problem is being scaled so that  $S^{-1/2} \epsilon \sim \mathcal{N}(0, \sigma_o^2 I)$ , where  $S$  is assumed reasonably correct. Subsets of observations where this is not true, (for example radiance data), may be excluded from  $V$  by partial GCV, see Wahba, Johnson, Gao & Gong (1994), there  $(I - A)$  is replaced by  $E(I - A)$  where  $E$  is a possibly weighted indicator matrix for the observations to be included. However the GCV estimate is robust to various assumptions about  $K, \Sigma$  and  $x$ . We remark that both the GCV and GML estimate can be defined when  $K$  is nonlinear. In either the linear or nonlinear case the trace of  $A$  may be estimated by the randomized trace method without having  $A$  explicitly, given a 'black box' which produces  $Kx_\lambda$  given  $y$ , see Wahba, Johnson, Gao & Gong (1995). Dee & daSilva (1998) and Dee, Gaspari, Redder, Rukhovets & daSilva (1998) have used maximum likelihood methods to estimate parameters in forecast error covariances in several practical examples, and have compared some of the results with GCV estimates, obtaining generally similar results in the examples tried. It may be possible to combine the strengths of both methods by iterating back and forth, using likelihood methods for parameters in  $S$  and GCV for parameters in  $K$  and  $\Sigma$ , this is speculative at the moment. The ordinary cross-validation 'warhorse' of leaving out a subset of the observations may also be used (with care).

## 4 MODELS FOR MODEL ERROR, DISCRETE TIME CASE

We outline some classes of models for time dependent model error,  $t = 1, \dots, T$ . Let  $g$  be a generic index,  $g = (lat, long, z, type)$  where *type* indexes the analysis variables, i. e. *type* = surface temperature, *type* = vorticity, etc. See Wahba (1992) for more on generic indices. We list some stochastic models for  $\xi_t(g)$  which remove the restriction that  $E\xi_t = 0, E\xi'_s \xi_t = 0, s \neq t$ . A fairly general class of models is

$$\xi_t(g) = \mu_t(g) + \sum_{k=1}^{\infty} z_k(t) \sqrt{\lambda_k(t)} \Phi_k(g) \quad (13)$$

where  $\mu_t(g) = \sum_{\nu=1}^M d_\nu F_\nu(t, g)$  is a mean function (bias term) specified except for a modest number of coefficients  $d_\nu, \nu = 1, \dots, M$ , to be found, the  $\lambda_k$  and  $\Phi_k$  are specified up to some parameters  $\theta_\xi$ , and  $Ez_k(t) = 0, Ez_k(s)z_l(t) = r_{k,l}(s, t)$ . Then

$$\text{cov}\xi_s(g)\xi_t(h) = \sum_{k,l} r_{k,l}(s, t) \sqrt{\lambda_k(s)} \sqrt{\lambda_l(t)} \Phi_k(g) \Phi_l(h). \quad (14)$$

<sup>4</sup>e. g. In principle at least, the Hessian of  $B$  with respect to  $(\lambda, \theta)$  should be well conditioned.

The simplest generalization over model error independent from time to time is the tensor product case,  $r_{k,l}(s,t) = 0, k \neq l, r_{k,k} = r(s,t)$ , independent of  $k$ , and  $\lambda_k(s) = \lambda_k$  independent of  $s$ , which gives

$$E\xi_s(g)\xi_t(h) = r(s,t)R(g,h), \quad (15)$$

where  $R(g,h) = \sum_k \lambda_k \Phi_k(g)\Phi_k(h)$ . This gives the penalty term for model error as

$$\sum_{s,t=1}^{T-1} r^{st}(\xi_s - \mu_s)'R^{-1}(\xi_t - \mu_t), \quad (16)$$

where  $r^{st}$  is the  $(s,t)$ th entry of the inverse of the matrix with  $(s,t)$  entry  $r(s,t)$ . Here the coefficients in  $\mu_t$  become part of the variational problem. Simple examples include  $z(\cdot)$  an autoregressive scheme or moving average where the correlation structure of the  $z(\cdot)$  process can be defined as short or long memory. A model like

$$E\xi_s(g)\xi_t(h) = \sum_{\alpha} r_{\alpha}(s,t)R_{\alpha}(g,h) \quad (17)$$

would allow for different time scales in the  $r_{\alpha}$ . Luo, Wahba & Johnson (1998) consider (15) in a simple situation with  $g$  on the sphere,  $R$  an isotropic covariance on the sphere and  $z(\cdot)$  a second order difference scheme forced by white noise. They used this model to estimate the linear time trend as a function of space given historical data irregular in time and space, by using the fact that the time trend is obtained as an orthogonal projection of the fit onto the relevant subspace of the implied RKHS. It may be possible to use similar techniques to diagnose model error. Griffith & Nichols (1998) have recently examined some simplified dynamical models for model error.

## 5 ELEMENTS OF REPRESENTER THEORY

In this section, time is continuous,  $t \in [0, 1]$ , and some components of  $g$  (e. g. space variables) are also to be thought of as continuous. Continuous time representer theory in RKHS has recently been applied in a number of places, see Bennett (1992), Bennett (1997), Bennett, Chua & Leslie (1996), Eknes & Evensen (1997), Evensen et al. (1998), Amodei (1997), Wahba (1992). We remind the reader that for every positive definite function  $R$  on  $\mathcal{T} \otimes \mathcal{T}$ , where  $\mathcal{T}$  is an abstract index set, there exists a unique RKHS and vice versa (The Moore-Aronszajn Theorem). There also exists a well defined zero mean Gaussian stochastic process with  $R$  as its covariance, however, sample functions of the stochastic process are not, with probability 1 in the RKHS if the RKHS is infinite dimensional. See Wahba (1990b), Weinert (1982) for more on reproducing kernel Hilbert spaces.

Let

$$\xi_t(g) = \xi_t^0(g) + \xi_t^1(g), \quad (18)$$

where

$$\mathcal{L}\xi_t^0(g) = 0, \quad \mathcal{B}\xi_t^1(g) = 0, \quad (19)$$

where  $\mathcal{L}_t$  is a (linear) evolution operator, and  $\mathcal{B}$  are initial/boundary conditions which serve to make the solution of differential equation unique, so that

$$\xi_t(g) = \xi_t^0(g) + \int G(t, g; \tilde{t}, \tilde{g})u(\tilde{t}, \tilde{g})d\tilde{t}d\tilde{g} \quad (20)$$

for some  $u(\cdot)$  where  $G$  is the Green's function for  $\mathcal{L}$  and  $\mathcal{B}$ . If  $u$  is treated as though is is a zero mean Gaussian stochastic process with covariance  $R_u(s, g; t, h)$  then

$$E\xi_s^1(g)\xi_t^1(h) = \int \int G(s, g; \tilde{s}, \tilde{g})G(t, h; \tilde{t}, \tilde{h})R_u(\tilde{s}, \tilde{g}; \tilde{t}, \tilde{h})d\tilde{s}d\tilde{g}d\tilde{t}d\tilde{h} = R^1(s, g; t, h), \quad [say]. \quad (21)$$

Let

$$E\xi_s^0(g)\xi_t^0(h) = R^0(s, g; t, h), \quad (22)$$

and suppose that  $\xi_s^0$  and  $\xi_s^1$  are independent, then

$$E\xi_s(g)\xi_t(h) \equiv R(s, g; t, h) = R^0(s, g; t, h) + R^1(s, g; t, h) \equiv R_{\theta_0}^0(s, g; t, h) + R_{\theta_M}^1(s, g; t, h). \quad (23)$$

Under general circumstances  $\mathcal{L}R^0(s, g; \cdot, \cdot) = 0$ , where  $\mathcal{L}$  is applied to  $R^0$  considered as a function of  $(\cdot, \cdot)$  for each fixed  $(s, g)$ . Furthermore, the RKHS  $\mathcal{H}_R$  with reproducing kernel (RK) given by  $R = R^0 + R^1$  of (23) consists of the direct sum of the orthogonal subspaces  $\mathcal{H}_R = \mathcal{H}_{R^0} \oplus \mathcal{H}_{R^1}$ , respectively containing solutions of the homogeneous equation and solutions to the differential equation satisfying homogeneous boundary conditions. Changing notation from  $\xi_t(g)$  to  $f_t(g)$  to indicate that we are now letting  $f$  be an element of  $\mathcal{H}_R$ , we have that if  $f^1 \in \mathcal{H}_{R^1}$ , then  $\|f^1\|_{\mathcal{H}_{R^1}}^2 = \|\mathcal{L}f^1\|_{\mathcal{H}_{R_u}}^2$  where  $\|\cdot\|_{\mathcal{H}_{R_u}}^2$  is the square norm in  $\mathcal{H}_{R_u}$ . If  $u$  had instead been taken as 'white noise' then  $R_u$  would not appear in (21) and the  $\mathcal{H}_{R_u}$  norm would be replaced with the usual  $L_2$  norm. Decomposing  $f$  into  $f^0$  and  $f^1$  analogous to (18,19) gives (the obvious)  $\|f\|_{\mathcal{H}_R}^2 = \|f^0\|_{\mathcal{H}_{R^0}}^2 + \|f^1\|_{\mathcal{H}_{R^1}}^2$ .

Let  $L_1, \dots, L_n$  be  $n$  bounded linear functionals on  $\mathcal{H}_R$ . Basic representer theory (see Kimeldorf & Wahba (1971), Wahba (1990b)) in RKHS tells us that the solution to the problem: find  $f \in \mathcal{H}_R$  to minimize

$$\sum_{i=1}^n (y_i - L_i f)^2 + \lambda \|f\|_{\mathcal{H}_R}^2 \quad (24)$$

is in the span of the  $n$  representers  $\eta_i$  of  $L_i$  in  $\mathcal{H}_R$ , where

$$\eta_i(s, g) = L_{i(t,h)} R(s, g; t, h), \quad (25)$$

where  $L_{i(t,h)}$  means  $L_i$  applied to what follows considered as a function of  $(t, h)$ . We may replace  $\|f\|_{\mathcal{H}_R}^2$  in (24) by e. g.  $\|f^0\|_{\mathcal{H}_{R^0}}^2 + w\|f^1\|_{\mathcal{H}_{R^1}}^2$ , the new problem is in theory solved with the aid of the RK  $R^0(s, g; t, h) + w^{-1}R^1(s, g; t, h)$ ;  $w \rightarrow \infty$  corresponds to the 'perfect model' assumption. Simple prototypes appear in Kimeldorf & Wahba (1971) and Wahba (1990b) where a seminorm penalty is also allowed.

Tuning for model error may have potential to provide diagnostic information concerning model error. Tony Weaver (personal communication) has remarked on the necessity of considering extrapolating (forecasting) the correlated part of any fit to model error. Most desirable, of course is to eliminate model error to the extent possible.

Recently Lin, Wahba, Xiang, Gao, Klein & Klein (1998), in a different (and much simpler) context, but with a relatively large, irregularly spaced data set, solved the variational problem under consideration in the span of a selected subset of the representers  $\eta_i, i = 1, \dots, n$ , with excellent results.

## 6 ACKNOWLEDGMENTS

We thank Andrew Bennett for generously providing advance copies of much of his recent work, and Phillipe Courtier for providing talk notes. We thank Tony Hollingsworth, Francois Bouttier, Heiki Jarvinen, Mike Fisher, Tony Weaver, Olivier Talagrand, Nancy Nichols and Gene Golub for helpful conversations and remarks.

## References

- Amodei, L. (1997), Reproducing kernels of vector-valued function spaces, *in* A. LeMehaute, C. Rabut & L. Schumaker, eds, 'Surface Fitting and Multiresolution Methods', Vanderbilt University Press, Nashville TN, pp. 17–26.
- Bennett, A. (1992), *Inverse methods in physical oceanography*, Cambridge University Press, 346pp.
- Bennett, A. (1997), 'Inverse methods and data assimilation', College of Oceanic And Atmospheric Sciences, Corvallis OR. Summer School Lecture Notes.
- Bennett, A., Chua, B. & Leslie, L. (1996), 'Generalized inversion of a global weather prediction model', *Meteorology and Atmospheric Physics* **60**, 165–178.
- Cohn, S., daSilva, A., Guo, J., Sienkiewicz, M. & Lamich, D. (1998), 'Assessing the effects of data selection with the DAO physical-space statistical analysis system', *Monthly Weather Review* **126**, 2913–2926.
- Courtier, P. (1997), 'Dual formulation of four-dimensional variational assimilation', *Q. J. R. Meteorol. Soc.* **123**, 2449–2461.
- Dee, D. & daSilva, A. (1998), 'Maximum-likelihood estimation of forecast and observation error covariance parameters. part i: Methodology', *Monthly Weather Review* **in press**, xx–xx.
- Dee, D., Gaspari, G., Redder, C., Rukhovets, L. & daSilva, A. (1998), 'Maximum-likelihood estimation of forecast and observation error covariance parameters. part ii: Applications', *Monthly Weather Review* **in press**, xx–xx.
- Eknes, M. & Evensen, G. (1997), 'Parameter estimation solving a weak constraint variational formulation for an Ekman model', *J. Geophysical Research* **102(C6)**, 479–491.
- Evensen, G., Dee, D. & Schroter, J. (1998), Parameter estimation in dynamical models, *in* E. Chassignet & J. Verron, eds, 'Ocean Modeling and Parameterizations', NATO ASI, **in press**, xx–xx.
- Gong, J., Wahba, G., Johnson, D. & Tribbia, J. (1998), 'Adaptive tuning of numerical weather prediction models: simultaneous estimation of weighting, smoothing and physical parameters', *Monthly Weather Review* **125**, 210–231.
- Griffith, A. & Nichols, N. (1998), 'Adjoint methods for treating model error in data assimilation', *Proc. ICFD Conference on Numerical Methods in Fluid Dynamics* **in press**, xx–xx.
- Kimeldorf, G. & Wahba, G. (1971), 'Some results on Tchebycheffian spline functions', *J. Math. Anal. Applic.* **33**, 82–95.
- Lin, X., Wahba, G., Xiang, D., Gao, F., Klein, R. & Klein, B. (1998), Smoothing spline ANOVA models for large data sets with Bernoulli observations and the randomized GACV, Technical Report 998, Department of Statistics, University of Wisconsin, Madison WI.
- Luo, Z., Wahba, G. & Johnson, D. (1998), 'Spatial-temporal analysis of temperature using smoothing spline ANOVA', *J. Climate* **11**, 18–28.
- Navon, I. (1998), 'Practical and theoretical aspects of adjoint parameter estimation and identifiability in meteorology and oceanography', *Dynamics of Atmospheres and Oceans* **27**, 55–79.

O'Sullivan, F. (1991), 'Sensitivity analysis for regularized estimation in some system identification problems', *SIAM J. Sci. Stat. Comput.* **12**, 1266–1283.

Wahba, G. (1985), 'A comparison of GCV and GML for choosing the smoothing parameter in the generalized spline smoothing problem', *Ann. Statist.* **13**, 1378–1402.

Wahba, G. (1990a), Regularization and cross validation methods for nonlinear, implicit, ill-posed inverse problems, in A. Vogel, C. Ofoegbu, R. Gorenflo & B. Ursin, eds, 'Geophysical Data Inversion Methods and Applications', Vieweg, Weisbaden-Braunschweig, pp. 3–13.

Wahba, G. (1990b), *Spline Models for Observational Data*, SIAM. CBMS-NSF Regional Conference Series in Applied Mathematics, v. 59.

Wahba, G. (1992), Multivariate function and operator estimation, based on smoothing splines and reproducing kernels, in M. Casdagli & S. Eubank, eds, 'Nonlinear Modeling and Forecasting, SFI Studies in the Sciences of Complexity, Proc. Vol XII', Addison-Wesley, pp. 95–112.

Wahba, G., Johnson, D., Gao, F. & Gong, J. (1994), Adaptive tuning of numerical weather prediction models: Part I, randomized GCV and related methods in three and four dimensional data assimilation, Technical Report 920, Department of Statistics, University of Wisconsin, Madison, WI.

Wahba, G., Johnson, D., Gao, F. & Gong, J. (1995), 'Adaptive tuning of numerical weather prediction models: randomized GCV in three and four dimensional data assimilation', *Mon. Wea. Rev.* **123**, 3358–3369.

Weinert, H., ed. (1982), *Reproducing kernel Hilbert spaces: Application in signal processing*, Hutchinson Ross, Stroudsburg, PA.