

# A method for estimating the analysis error covariance with 4DVar

Hans Ngodock

Innocent Souopgui, Matthew Carrier,

Scott Smith, John Osborne and Joseph D'Addezio

# Background

- The traditional formulation of 4dvar does not provide the analysis error covariance, which is needed for analysis uncertainty.
- Bennett (2002) introduced a Monte-Carlo approach for estimating the posterior error covariance, but it requires a large number of samples (solutions of the tangent linear model).
- Recently, Moore et al. (2012) proposed a method to estimate analysis and forecast error variances using the adjoint of the 4dvar system, based on the premise of perturbing the observations and background fields.
- A new method is proposed here, that consists of an ensemble of perturbations of the optimal adjoint solution, instead of an ensemble of 4dvar analyses, e.g. Bonavita et al. (2012).
- We hypothesize that carrying out an ensemble of 4dvar solutions amounts to computing an ensemble of optimal adjoint solutions, since the forward nonlinear or linearized dynamics are not changed, and neither is the minimization process of the 4dvar system itself.
- Note that perturbing the optimal adjoint is equivalent to perturbing the innovation vector (on which the adjoint linearly depend), and the innovation vector itself depends on the observations and the background.

# Background

$$\begin{cases} \frac{\partial u}{\partial t} = L(u) + F + f, & 0 \leq t \leq T \\ u(x, 0) = I(x) + i(x), \end{cases}$$

$$y_m = H_m u(x, t_m) + \varepsilon_m, \quad 1 \leq m \leq M,$$

$$J = \int_0^T \int_{\Omega} \int_0^T \int_{\Omega} f(x, t) W_f(x, t, x', t') f(x', t') dx' dt' dx dt + \int_{\Omega} \int_{\Omega} i(x) W_i(x, x') i(x') dx' dx + \varepsilon^T W_{\varepsilon} \varepsilon,$$

$$\begin{cases} \frac{\partial \hat{u}}{\partial t} = L(\hat{u}) + F + C_f \bullet \lambda, \\ \hat{u}(x, 0) = I(x) + C_i \circ \lambda(x, 0), \\ -\frac{\partial \lambda}{\partial t} = \left[ \frac{\partial L}{\partial u}(\hat{u}) \right]^T \lambda + \sum_{m=1}^M \sum_{n=1}^M W_{\varepsilon, mn} (y_m - H_m \hat{u}) H_m^T \delta(x - x_m) \delta(t - t_m), \\ \lambda(x, T) = 0, \end{cases}$$

$$\lambda(x, t) = \int_0^T \int_{\Omega} W_f(x, t, x', t') f(x', t') dx' dt'$$

$$C_f \bullet \lambda(x, t) = \int_0^T \int_{\Omega} C_f(x, t, x', t') \lambda(x', t') dx' dt',$$

$$C_i \circ \lambda(x, 0) = \int_{\Omega} C_i(x, x') \lambda(x', 0) dx',$$

# Strong Constraints

$$\begin{cases} \frac{\partial \hat{u}}{\partial t} = L(\hat{u}) + F, \\ \hat{u}(x, 0) = I(x) + C_i \circ \lambda(x, 0), \\ -\frac{\partial \lambda}{\partial t} = \left[ \frac{\partial L}{\partial u}(\hat{u}) \right]^T \lambda + \sum_{m=1}^M \sum_{n=1}^M W_{\varepsilon, mn} (y_m - H_m \hat{u}) H_m^T \delta(x - x_m) \delta(t - t_m), \\ \lambda(x, T) = 0, \end{cases}$$

It is clear that all the corrections of the model trajectory are determined by the optimal adjoint at time 0. We can thus generate an ensemble of 4dvar analyses by

$$\lambda_n(x, 0) = \lambda(x, 0) + \mu_n(x), \quad n = 1, \dots, N$$

$N$  is the size of the ensemble. The ensemble of analyses is then computed as

$$\begin{cases} \frac{\partial \hat{u}_n}{\partial t} = L(\hat{u}_n) + F, \\ \hat{u}_n(x, 0) = I(x) + C_i \circ \lambda_n(x, 0). \end{cases}$$

# Weak Constraints

$$\hat{u}(x, t) = u_F(x, t) + \sum_{m=1}^M \beta_m r_m(x, t)$$

$$\begin{cases} \frac{\partial r_m}{\partial t} = Lr_m + C_f \bullet \alpha_m(x, t), \\ r_m(x, 0) = C_i \circ \alpha_m(x, 0), \\ -\frac{\partial \alpha_m}{\partial t} = L^T \alpha_m + H_m^T \delta(x - x_m) \delta(t - t_m), \\ \alpha_m(x, T) = 0, \end{cases}$$

$$\lambda_n(x, 0) = \lambda(x, 0) + \mu_n(x), \quad n = 1, \dots, N$$

$N$  is the size of the ensemble. The ensemble of analyses is then computed as

$$\begin{cases} \frac{\partial \hat{u}_n}{\partial t} = L(\hat{u}_n) + F, \\ \hat{u}_n(x, 0) = I(x) + C_i \circ \lambda_n(x, 0). \end{cases}$$

# Weak Constraints

$$\hat{u}(x, t) = u_F(x, t) + \sum_{m=1}^M \beta_m r_m(x, t)$$

$$\begin{cases} \frac{\partial r_m}{\partial t} = Lr_m + C_f \bullet \alpha_m(x, t), \\ r_m(x, 0) = C_i \circ \alpha_m(x, 0), \\ -\frac{\partial \alpha_m}{\partial t} = L^T \alpha_m + H_m^T \delta(x - x_m) \delta(t - t_m), \\ \alpha_m(x, T) = 0, \end{cases}$$

Where L denotes the linearized operator  $\left[ \frac{\partial L}{\partial u}(u) \right]$ , and the first guess is the solution of

$$\begin{cases} \frac{\partial u_F}{\partial t} = Lu_F + F, \\ u_F(x, 0) = I. \end{cases}$$

It may be shown that the representer coefficients are the solution of the linear system

$$(R + C_\varepsilon) \hat{\beta} = y - Hu_F$$

The analysis increment can be written as

$$\xi = LC_i L^T H^T (HLCL^T H^T + C_\varepsilon)^{-1} d = LCL^T H^T \beta$$

# Weak Constraints

Assuming that the prior error covariance for the first guess  $\mathbf{B}_{u_F}(x, t, x', t')$  is available  
The analysis error covariance is given by

$$\mathbf{B}_{\hat{u}}(x, t, x', t') = \mathbf{B}_{u_F}(x, t, x', t') - \mathbf{r}^T(x, t) \mathbf{P}^{-1} \mathbf{r}(x', t')$$

$$\mathbf{r}(x, t) = (r_1(x, t), r_2(x, t), \dots, r_M(x, t))^T$$

$$\mathbf{P} = \mathbf{H} \mathbf{L} \mathbf{C} \mathbf{L}^T \mathbf{H}^T + \mathbf{C}_\varepsilon$$

Obtaining the prior error covariance is difficult and tedious. It involves either propagating  $\mathbf{C}$  with the linearized and adjoint model as  $\mathbf{L} \mathbf{C} \mathbf{L}^T$  or Monte-Carlo simulations to generate an ensemble of TLM solution with perturbations of initial, boundary and model errors that sample the respective prescribed covariances

The Hessian of the cost function also provides the analysis error covariance, but it is difficult and expensive to compute, especially for high dimensions

An ensemble of data assimilation experiments (EDA) provides an ensemble of analyses from which one can compute the covariance. This is also expensive with 4dvar

# Weak Constraints

We rewrite the analysis equation in the form

$$\hat{u}_k = (u_F + \delta u_k) + L_k C L_k^T H^T (H L_k C L_k^T H^T + C_\varepsilon)^{-1} (y + \delta y_k - H(u_b + \delta u_k))$$

We make the following approximation

$$\begin{aligned} \hat{u}_k &= (u_F + \delta u_k) + L_k C L_k^T H^T (H L_k C L_k^T H^T + C_\varepsilon)^{-1} (d + \delta d_k) \\ &= (u_F + \delta u_k) + L_k C L_k^T H^T (\beta_k + \delta \beta_k) \\ &\approx u_F + L C L^T H^T (\beta + \delta \beta_k) \end{aligned}$$

We can approximate an ensemble of analyses by perturbing just the optimal representer coefficients and computing an analysis increment for each perturbation.

How do we perturb?

$$d = y - H u_F \quad E(dd^T) = C_\varepsilon + H B_{u_F} H^T \quad \beta = (H L C L^T H^T + C_\varepsilon)^{-1} d \quad B_{u_F} = L C L^T$$

$$E(\beta \beta^T) = (C_\varepsilon + H L C L^T H^T)^{-1}$$

The optimal representer coefficients are correlated, with variance smaller than that of both the observations and the prior errors



# Application

Lorenz-2005 type II (Lorenz, 2005)

$$\frac{du_n}{dt} = [u, u]_{K,n} - u_n + F, \quad n = 1, \dots, N$$

$$[u, u]_{K,n} = \sum_{j=-J}^J \sum_{i=-J}^J \frac{1}{K^2} \left( -u_{n-2K-i} u_{n-K-j} + u_{n-K+j-i} u_{n+K+j} \right)$$

$N=240, K=8, J=K/2, F=15, dt=0.025$

The model has a doubling time of 16 time steps

Spin-up 2000 and 2060 time steps for the background and true solutions

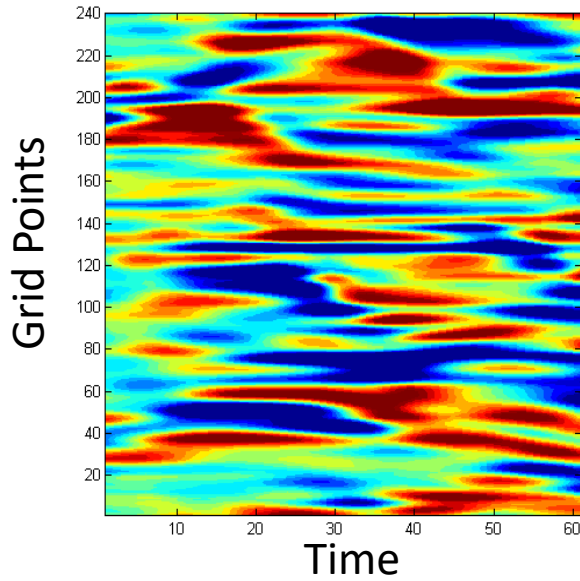
Observations sampled from true solution every 4 and 3 space and time steps respectively

Assimilation: 99 cycles of 60 time steps each,  $F=14.99$

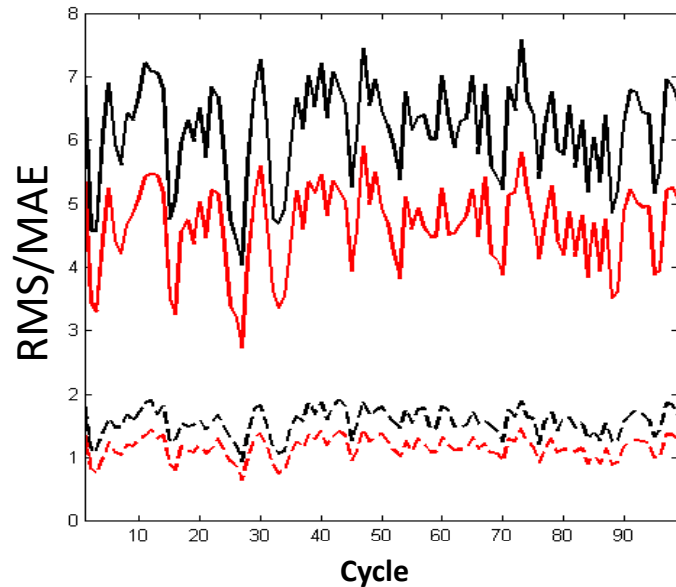
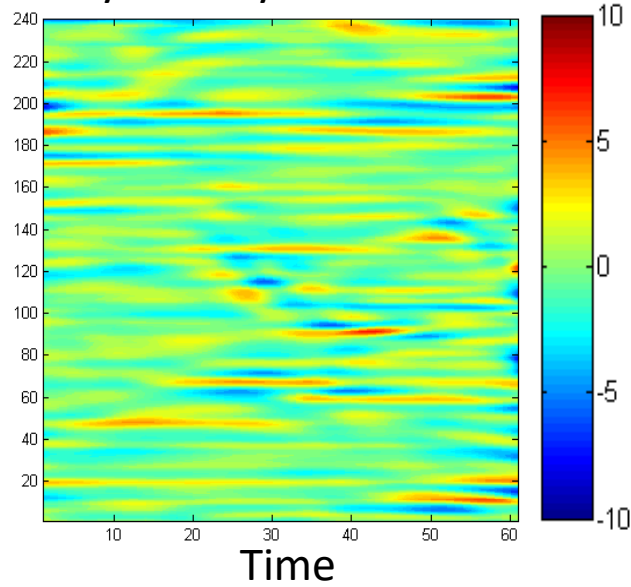
Covariances: Gaussian in space (scale 10 grid points), Markov in time (scale 20 steps)

# Application

Cycle 99 background residuals



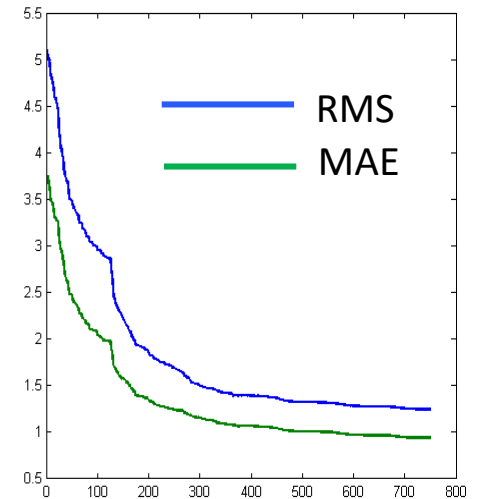
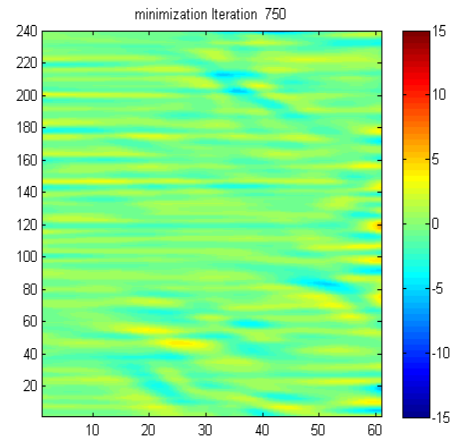
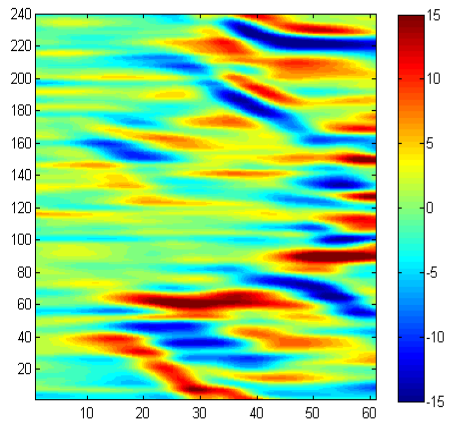
Cycle 99 analysis residuals



- RMS background
- MAE background
- - RMS Analysis
- - MAE Analysis

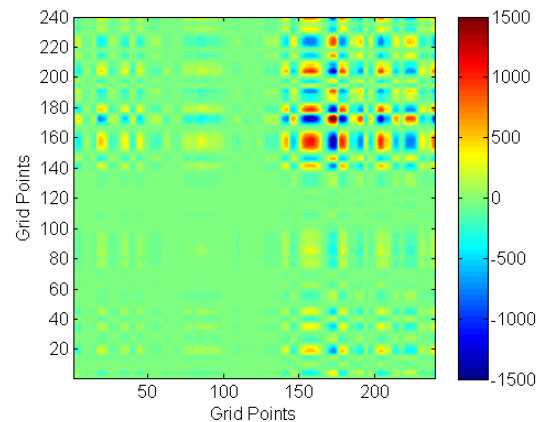
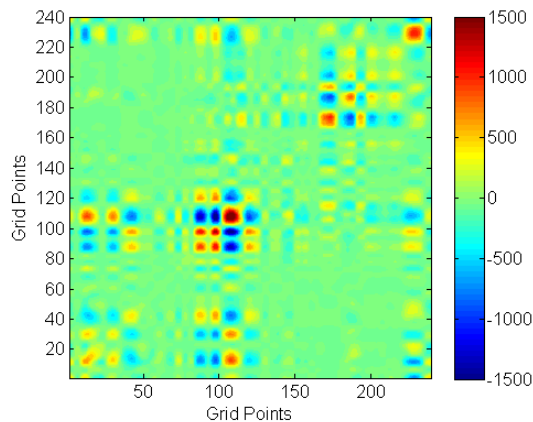
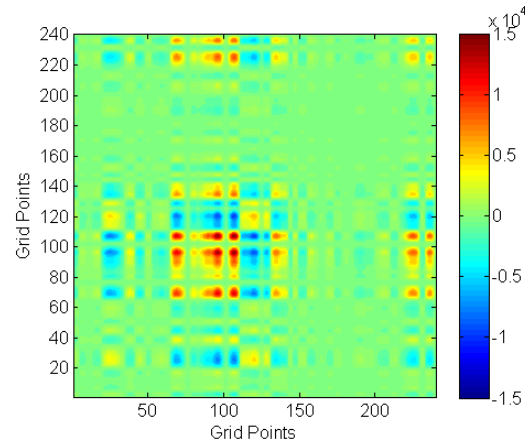
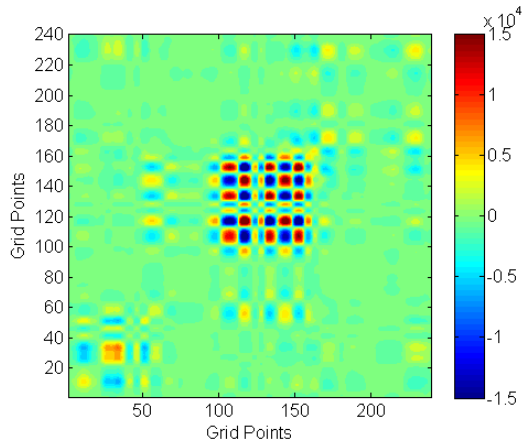
# Application

## Convergence of the minimization at cycle 80



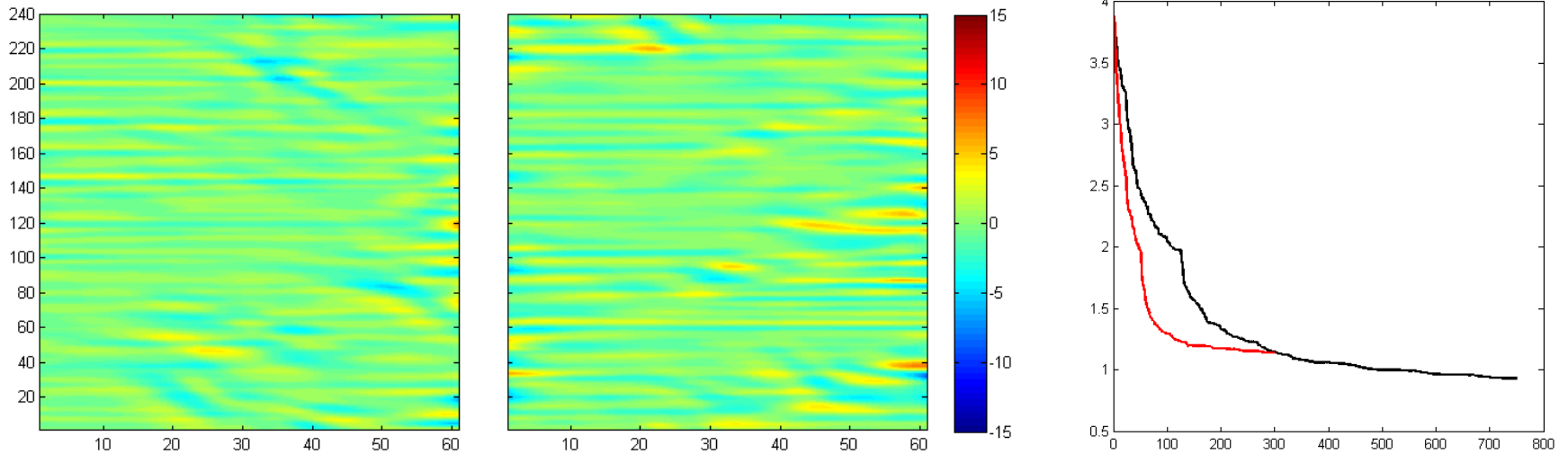
Residuals at the beginning and end of the minimization iterations, and RMS and MAE of the assimilated solution at every iteration of cycle 80.

# Application



A comparison of the analytical and estimated covariances for cycle 70 (a, b) and cycle 80 (c, d). The analytical covariance used 1000 samples for the prior, and the estimated covariance used 300 members

# Application



Cycle 70 residuals from the prescribed covariance at the end of the minimization (750 iterations), and the from the estimated covariance (300 iterations).

# Application

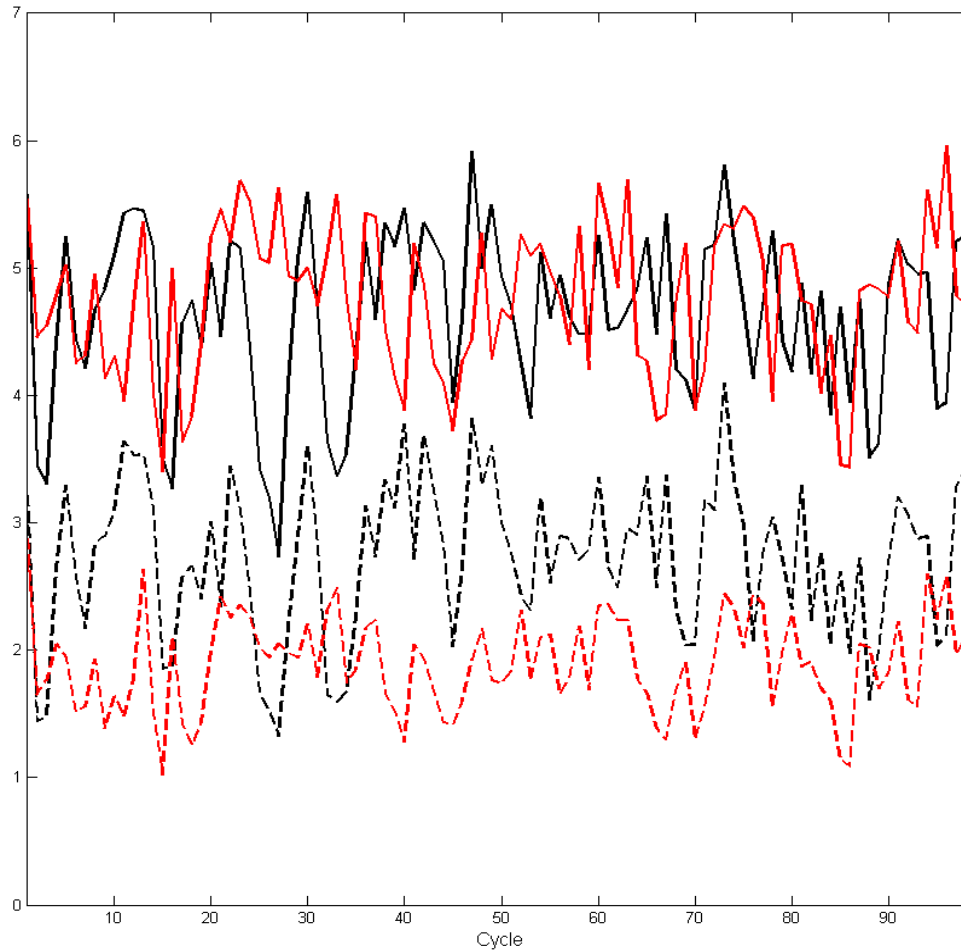


Figure 5: Mean absolute error of the background (solid lines) and analysis (dashed) using the prescribed covariances (black), and the estimated covariances (red). The analyses are computed after 100 minimization iterations.

# Conclusion

A method for estimating the analysis error covariance in a 4dvar data assimilation is proposed. It consists of perturbing the optimal representer coefficients and generating an ensemble of pseudo analyses

The method was applied to the Lorenz-2005 model, in a twin-data assimilation experiment.

The mean of the pseudo analyses ensemble, as well as the analysis that used the estimated covariance, has residuals comparable to the original data assimilation experiment using the prescribed covariance

The experiment using the estimated covariance converged much faster than the one using prescribed covariances.

There is still work to be done in the proper perturbation of the representer coefficients