# How does the regression step in the two-step EnKF connect to Bayesian estimation?

Ian Grooms

EnKF Workshop 2022

DART/EAKF
●○○

TWO-STEP BAYES
○○○○○○○○○○○○○

EXAMPLES
○○○○○
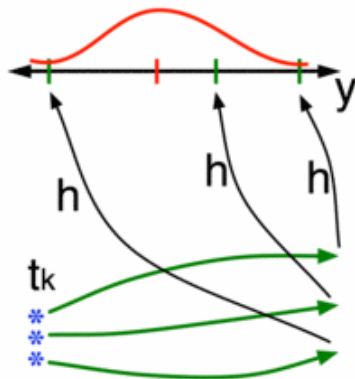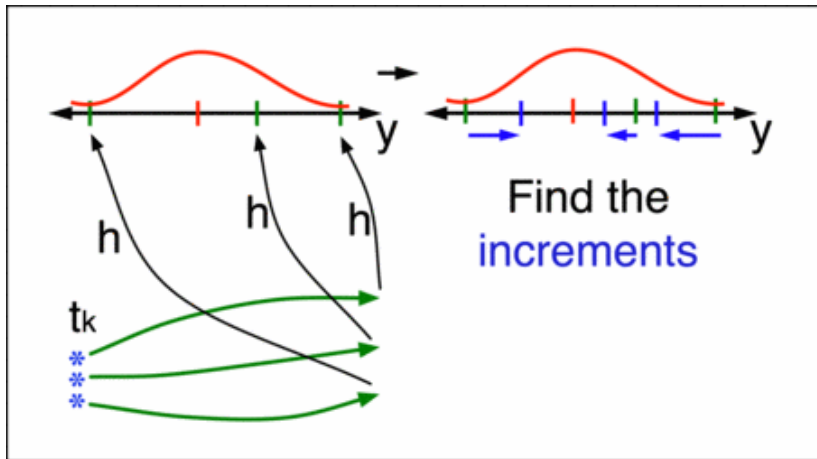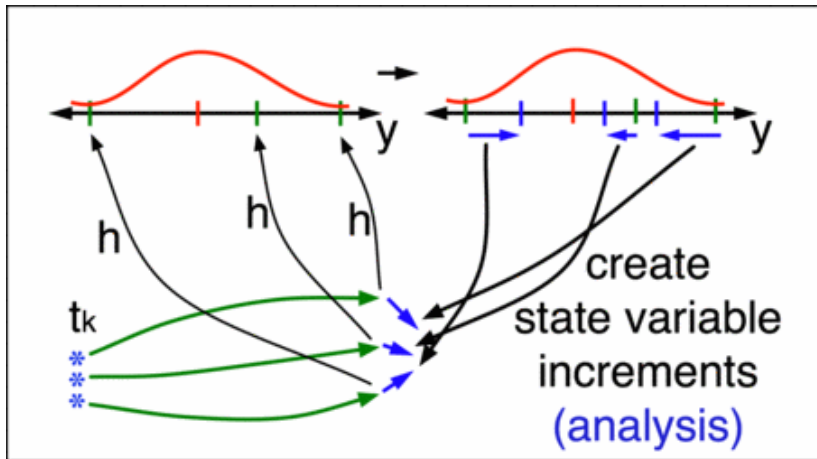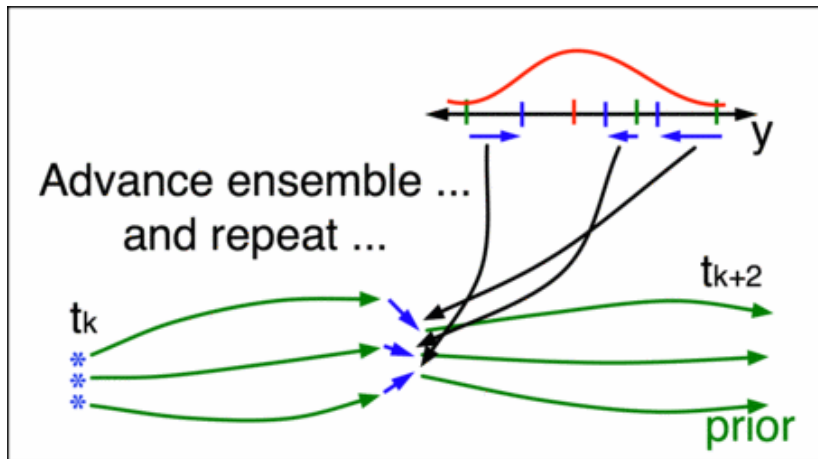
Convert each model state to an expected observation $y = h(x)$

Compare with observation and observational error distribution

Find the increments

If

- ► The update in observation space is Gaussian and
- ► The regression is linear and estimated using ordinary least squares (OLS)

Then this two-step is equivalent to a one-step EnKF where

$$\mathbf{B}\mathbf{H}^T = \frac{1}{N-1} \sum_{n=1}^{N} (\boldsymbol{x}_i - \bar{\boldsymbol{x}})(\boldsymbol{H}(\boldsymbol{x}_i) - \overline{\boldsymbol{H}(\boldsymbol{x})})^T$$

$$\mathbf{H}\mathbf{B}\mathbf{H}^T = \frac{1}{N-1} \sum_{n=1}^{N} (\boldsymbol{H}(\boldsymbol{x}_i) - \overline{\boldsymbol{H}(\boldsymbol{x})})(\boldsymbol{H}(\boldsymbol{x}_i) - \overline{\boldsymbol{H}(\boldsymbol{x})})^T$$

**Is this just a different way of implementing an EnKF?**

**If we replace EAKF or linear regression, do we lose the connection to Bayesian estimation?**

We have two random variables $X$ and $Y$ whose joint pdf is denoted

$$[\, x,\, y\,]$$

Our goal is to sample from the conditional distribution

$$[\, x\,|\, Y = y_o\,]$$

which can also be written as a Bayesian posterior

$$[\, x\,|\, Y = y_o\,] = \frac{[\, Y = y_o\,|\, x\,]}{[\, Y = y_o\,]}[\, x\,]$$

($y_o$ is the actual value of the observation; $y$ denotes any value that the random variable $Y$ could take.)

DART/EAKF
000

TWO-STEP BAYES
○●○○○○○○○○○○

EXAMPLES
○○○○○

Next, introduce a new random variable $Z$. With this new random variable we have a new joint distribution

$$[x, y, z].$$

The old one is just the marginal

$$[x, y] = \int [x, y, z] \mathrm{d}z.$$

The posterior that we care about is just the marginal of a new posterior:

$$[x \mid Y = y_o] = \int [x, z \mid Y = y_o] \mathrm{d}z = \int \frac{[Y = y_o \mid x, z]}{[Y = y_o]} [x, z] \mathrm{d}z.$$

DART/EAKF
000

TWO-STEP BAYES
○○●○○○○○○○○○○

EXAMPLES
○○○○○

For the general two-step framework we must make the following assumption about $Z$:

$$[\,Y = y_o \,|\, x,\, z\,] = [\,Y = y_o \,|\, z\,].$$

In the standard DART two-step, we have $Y = H(X) + \epsilon$ and usually set $Z = H(X)$.

We **do not** need an observation model of the form $Y = H(X) + \epsilon$. It's just mentioned here for illustration.

Use Bayes' rule to expand the posterior inside the integral, and use our one assumption about $Z$:

$$
\begin{aligned}
[\, x \,|\, Y = y_o \,] &= \int \frac{[\, Y = y_o \,|\, x, \, z \,]}{[\, Y = y_o \,]} [\, x, \, z \,] \mathrm{d}z \\
&= \int \frac{[\, Y = y_o \,|\, z \,]}{[\, Y = y_o \,]} [\, x, \, z \,] \mathrm{d}z.
\end{aligned}
$$

Now expand the prior $[\, x, \, z \,]$ as marginal times conditional

$$
[\, x \,|\, Y = y_o \,] = \int \left( \frac{[\, Y = y_o \,|\, z \,]}{[\, Y = y_o \,]} [\, z \,] \right) [\, x \,|\, z \,] \mathrm{d}z.
$$

$$[\,x\,|\,Y = y_o\,] = \int \left( \frac{[\,Y = y_o\,|\,z\,]}{[\,Y = y_o\,]} [\,z\,] \right) [\,x\,|\,z\,]\mathrm{d}z.$$

Notice that the quantity in parenthesis is a posterior distribution

$$\frac{[\,Y = y_o\,|\,z\,]}{[\,Y = y_o\,]} [\,z\,] = [\,z\,|\,Y = y_o\,]$$

so

$$[\,x\,|\,Y = y_o\,] = \int [\,x\,|\,z\,][\,z\,|\,Y = y_o\,]\mathrm{d}z.$$

How can we draw samples $\{x_i^+\}_{i=1}^N$ from a distribution of this form?

DART/EAKF
000

TWO-STEP BAYES
00000●000000

EXAMPLES
00000

Consider the following analogy. Suppose we have the following dynamics

$$X^{k+1} = M\left(X^k\right) + W^k$$

where $W^k$ is a random variable with pdf $[\,w\,]$. We know how to sample from $[\,X^{k+1}\,]$.

First draw an ensemble $\{x_i^k\}_{i=1}^N$ of samples of $X^k$.

Then apply the dynamics $M$ to each $x_i^k$ and finally add a sample from the noise $w_i^k$.

The pdf of $X^{k+1}$ is the marginal distribution

$$[x^{k+1}] = \int [x^{k+1}, x^k] \mathrm{d}x^k = \int [x^{k+1} \,|\, x^k][x^k] \mathrm{d}x^k$$

The first step in sampling from this distribution is to sample from $[x^k]$.

Then sample from $[x^{k+1} \,|\, X^k = x_i^k]$.

## GENERAL TWO-STEP SUMMARY

The two-step Bayesian update has exactly the same form.
Recall that we want to sample from

$$[\,x\,|\,Y = y_o\,] = \int [\,x\,|\,z\,][\,z\,|\,Y = y_o\,]\mathrm{d}z.$$

▶ Step 1: Generate an ensemble $\{z_i^+\}_{i=1}^N$ from the posterior
  $[\,z\,|\,Y = y_o\,]$.
▶ Step 2: Sample $x_i^+$ from pdf $[\,x\,|\,Z = z_i^+\,]$.

The difference between this case and the analogy is that in the
analogy we know what the dynamics are, so we know how to
sample from the conditional distribution.

DART/EAKF
000

TWO-STEP BAYES
○○○○○○○○○●○○○

EXAMPLES
○○○○○

This is where regression comes in. Propose a linear model of the form

$$\boldsymbol{X} = \boldsymbol{\beta}_0 + \boldsymbol{\beta}_1 Z + \boldsymbol{\eta}.$$

(Assume for the moment that $Z$ is scalar to make the exposition easier.)

If $\boldsymbol{\eta}$ is Gaussian with zero mean and covariance $\boldsymbol{\Sigma}$, then we are saying

$$\boldsymbol{X}|Z = z \sim \mathcal{N}\left(\boldsymbol{\beta}_0 + \boldsymbol{\beta}_1 z, \boldsymbol{\Sigma}\right)$$

We can estimate the regression coefficients using OLS.

## 2-STEP ENKF SUMMARY

- ▶ Generate a prior ensemble $\{x_i^-\}_{i=1}^N$.
- ▶ Generate a prior ensemble $\{z_i^-\}_{i=1}^N$.
- ▶ Generate a posterior ensemble $\{z_i^+\}_{i=1}^N$ using an EnKF.
- ▶ Estimate the regression coefficients as the OLS solution of the following system

$$\beta_0 + \beta_1 z_i^- = x_i^-, \quad i = 1, \ldots, N$$

- ▶ Generate samples $\eta_i = x_i^- - \beta_0 - \beta_1 z_i^-$
- ▶ Set $x_i^+ = \beta_0 + \beta_1 z_i^+ + \eta_i$.

You can write this in incremental form as: $\Delta x_i = \beta_1 \Delta z_i$.

**Non-Gaussian Generalizations**

The nice thing about the first step is that it is typically low-dimensional, so we can use methods to sample from $[\,z\,|\,Y = y_o\,]$ that work for non-Gaussian distributions but that might be impractical in higher dimensions. E.g.

- ▶ Particle Filters
- ▶ Gaussian Mixture Methods
- ▶ RHF: Anderson MWR 2010
- ▶ Gamma/Inverse-Gamma/Gaussian: Bishop QJRMS 2016
- ▶ Etc

**Non-Gaussian Generalizations**

Now that we know how the second step connects to the Bayesian problem, we can use advanced regression models

► General Linear Model $X = \eta + \sum_{j=1}^{J} \beta_j \phi_j(Z)$

► Generalized Linear Model $g(X) = \eta + \sum_{j=1}^{J} \beta_j \phi_j(Z)$. (E.g. Anderson Rank Regression 2019)

► Nonlinear Models (e.g. neural nets) $g(X; \beta_g) = \eta + f(Z; \beta_f)$

The assumptions about $\eta$ determine the form of the conditional distribution $[x \,|\, z]$, which also indicates the form of the objective function that must be minimized to find the unknown parameters.

I set up Lorenz-96 observing *all* variables every 0.05 model time units, with three different observation models:

$$\text{Linear: } Y = X + \epsilon$$

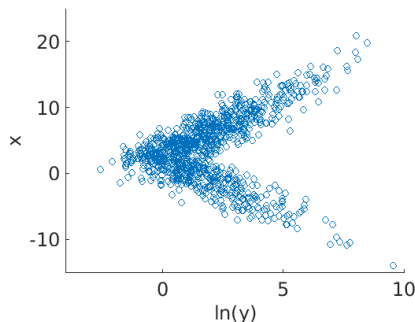$$\text{Logit-Normal: } Y = \frac{1}{1 + \exp\{0.5 \times (X - 2.5) + \epsilon\}}$$

$$\text{Log-Normal: } Y = \exp\{0.5 \times |X - 2.5| + \epsilon\}$$

where $\epsilon$ are standard normal.

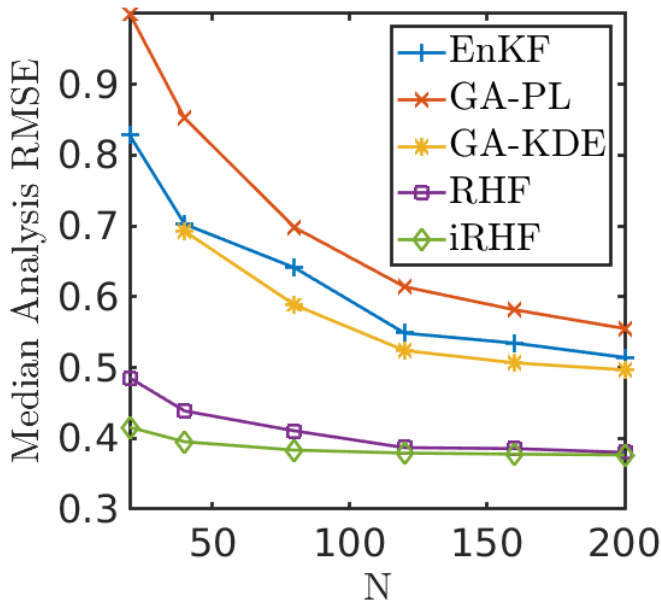I won't show results for the linear obs case.

I used multiplicative prior inflation and localization, both tuned to produce optimal results.

For the log-normal observations the likelihood is bimodal. In the standard approach ($Z = H(\mathbf{X})$) the second step requires fitting a regression to this kind of data
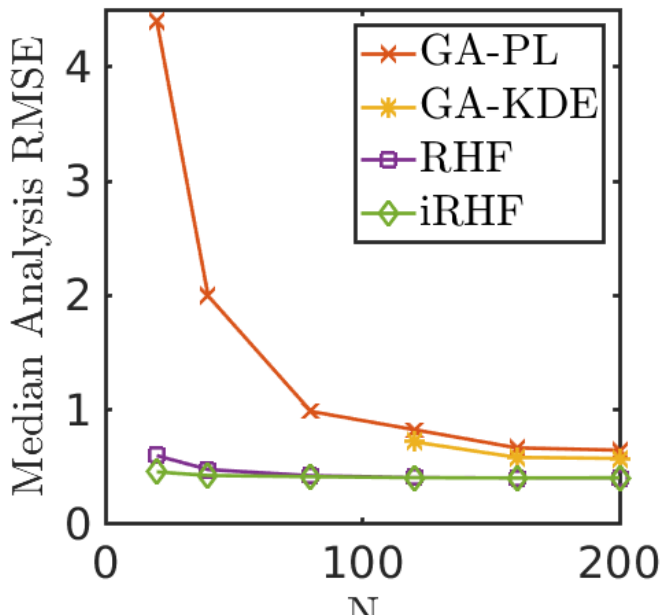


By choosing $Z = X$, I don't have to do regression through this kind of scatterplot. All the nonlinearity/bimodality is handled in the first step.

# LOGIT-NORMAL RESULTS

# LOG-NORMAL RESULTS

More details about iRHF along with a comparative discussion
of Gaussian Anamorphosis methods can be found in

Grooms, "A comparison of nonlinear extensions to the
ensemble Kalman filter" Computational Geosciences, 2022.