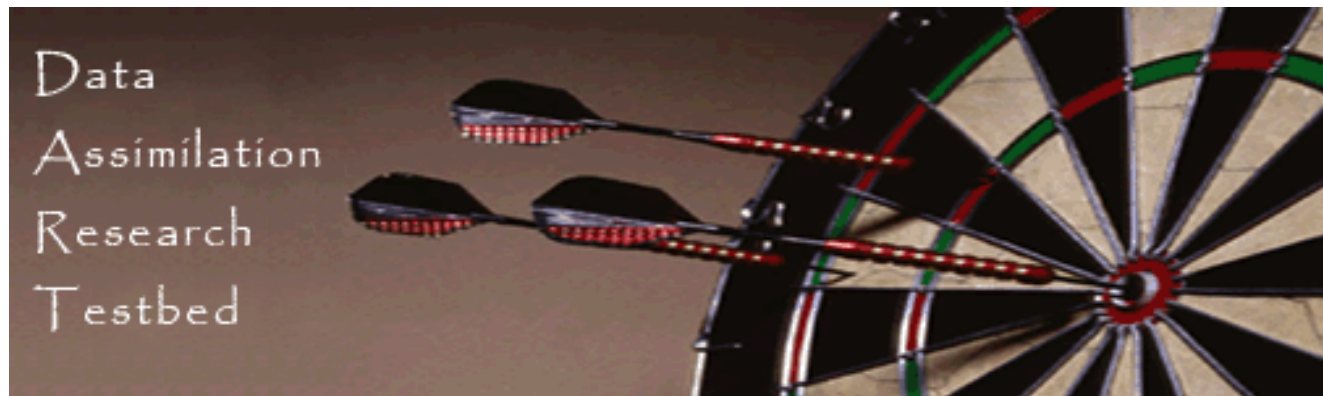


Ensemble Data Assimilation and Uncertainty Quantification

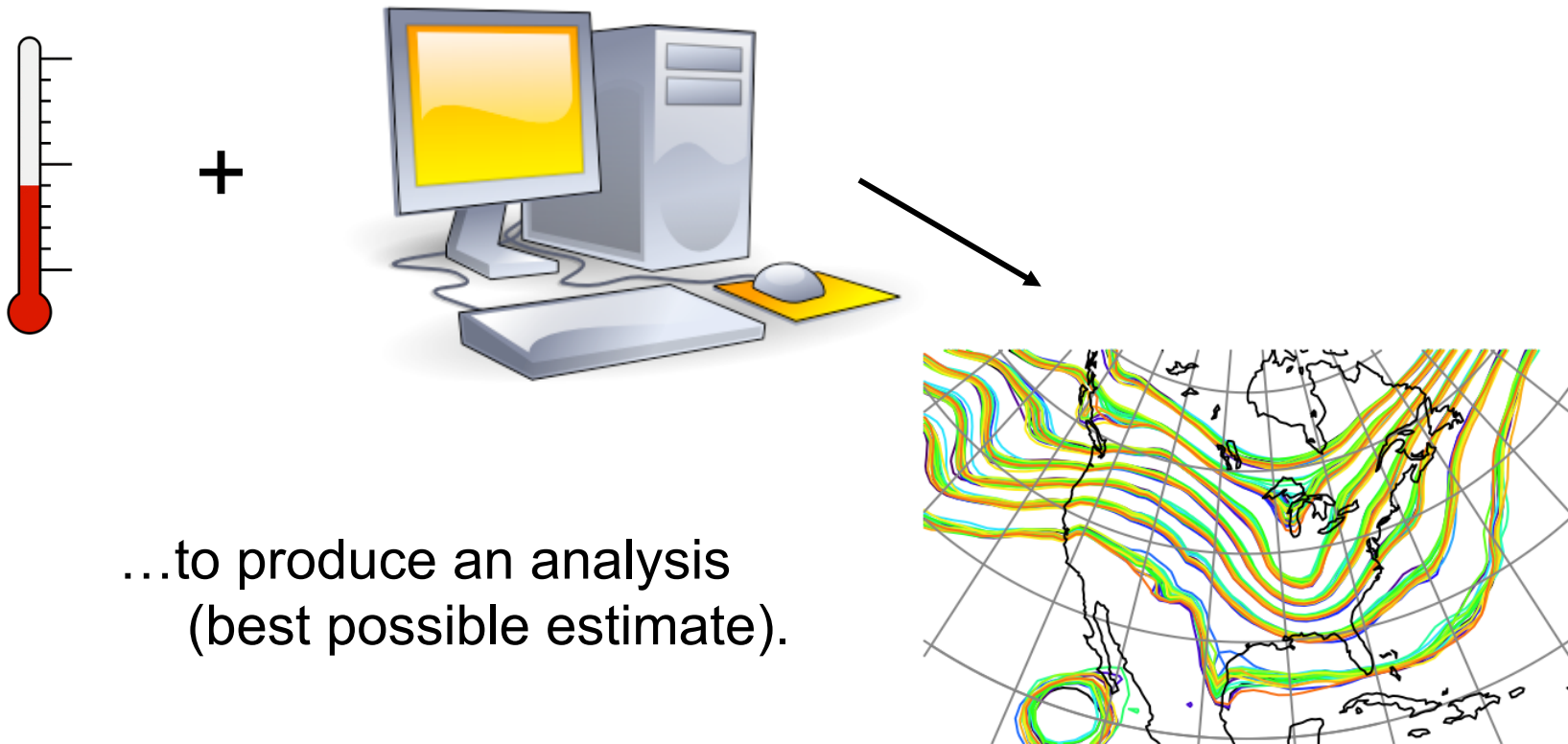


Jeff Anderson

National Center for Atmospheric Research

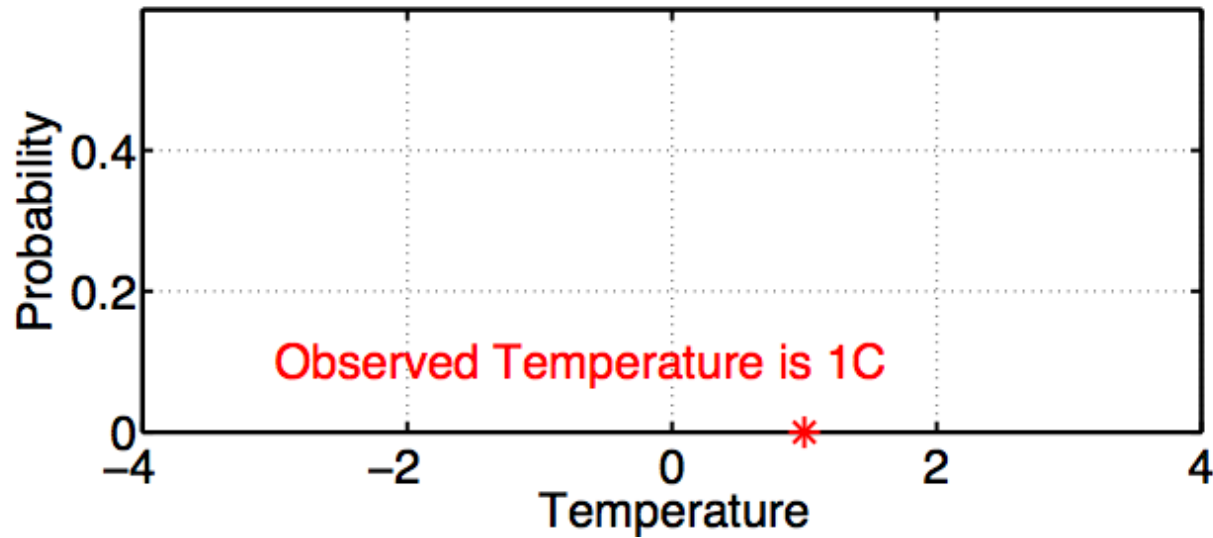
What is Data Assimilation?

Observations combined with a Model forecast...



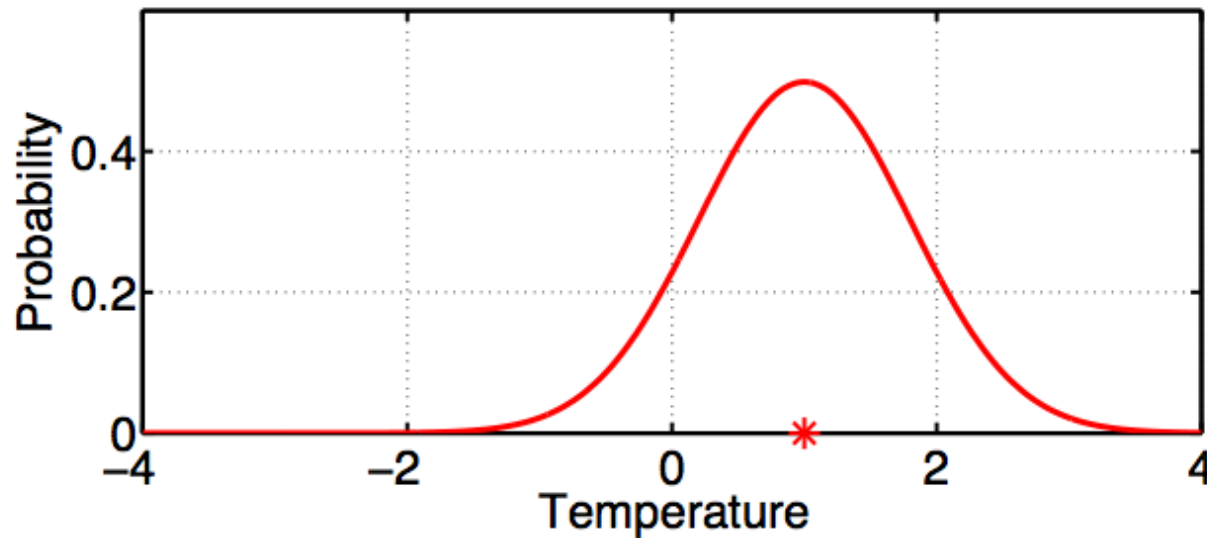
Example: Estimating the Temperature Outside

An observation has a value (*),



Example: Estimating the Temperature Outside

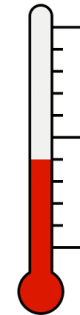
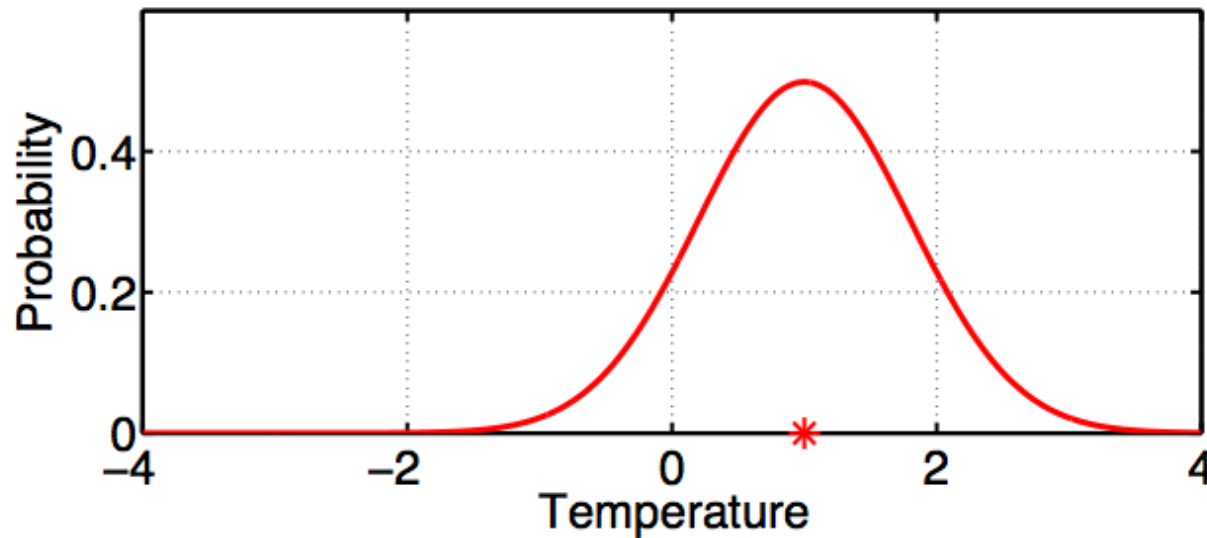
An observation has a value (*),



and an error distribution (red curve) that is associated with the instrument.

Example: Estimating the Temperature Outside

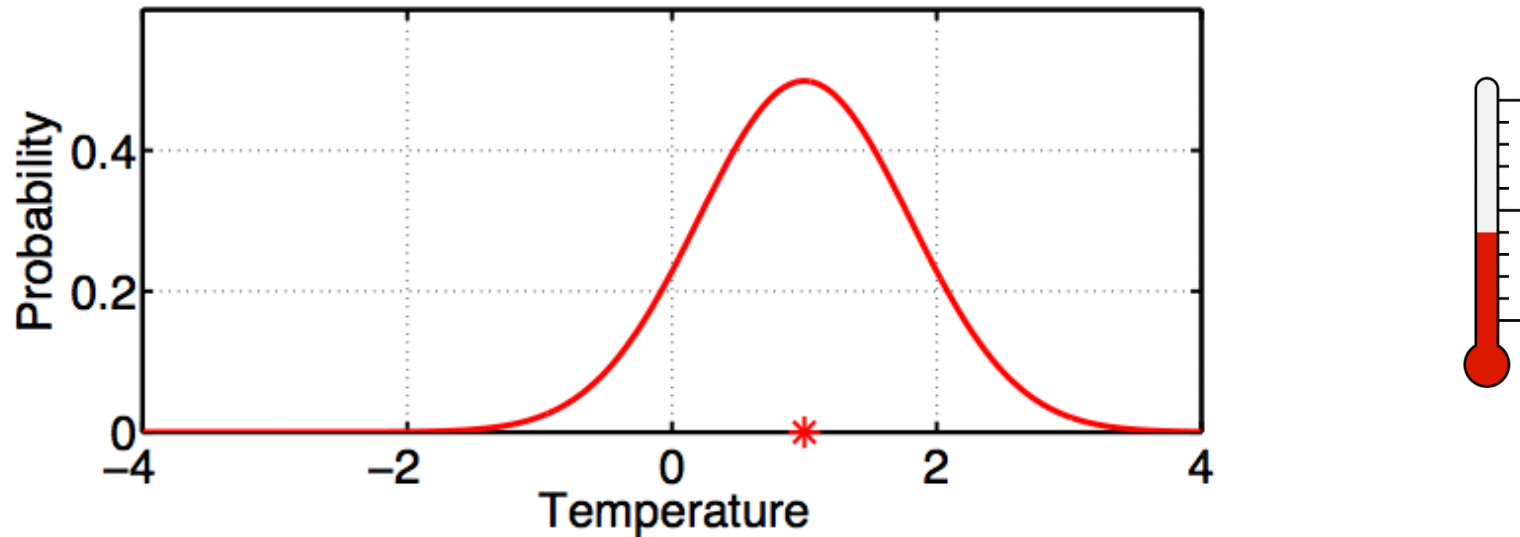
Thermometer outside measures 1C.



Instrument builder says thermometer is unbiased with $\pm 0.8\text{C}$ gaussian error.

Example: Estimating the Temperature Outside

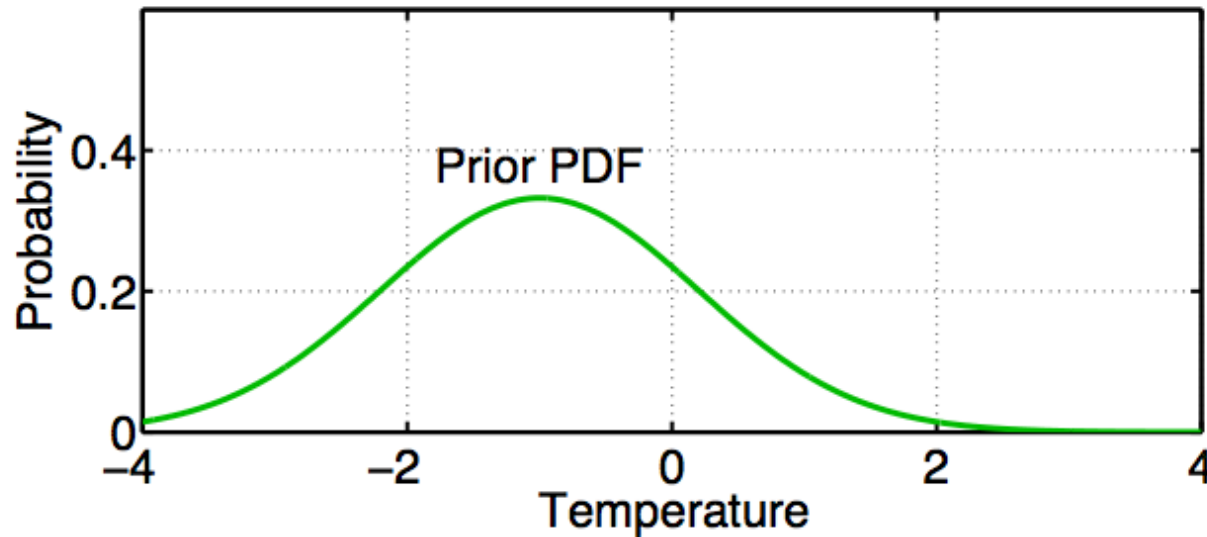
Thermometer outside measures 1C.



The red plot is $P(T | T_o)$, probability of temperature given that T_o was observed.

Example: Estimating the Temperature Outside

We also have a prior estimate of temperature.



The green curve is $P(T | C)$; probability of temperature given all available prior information C .

Example: Estimating the Temperature Outside

Prior information C can include:

1. Observations of things besides T ;
2. Model forecast made using observations at earlier times;
3. *A priori* physical constraints ($T > -273.15\text{C}$);
4. Climatological constraints ($-30\text{C} < T < 40\text{C}$).

Combining the Prior Estimate and Observation

Bayes

Theorem:

$$P(T | T_o, C) = \frac{P(T_o | T, C)P(T | C)}{\textit{Normalization}}$$

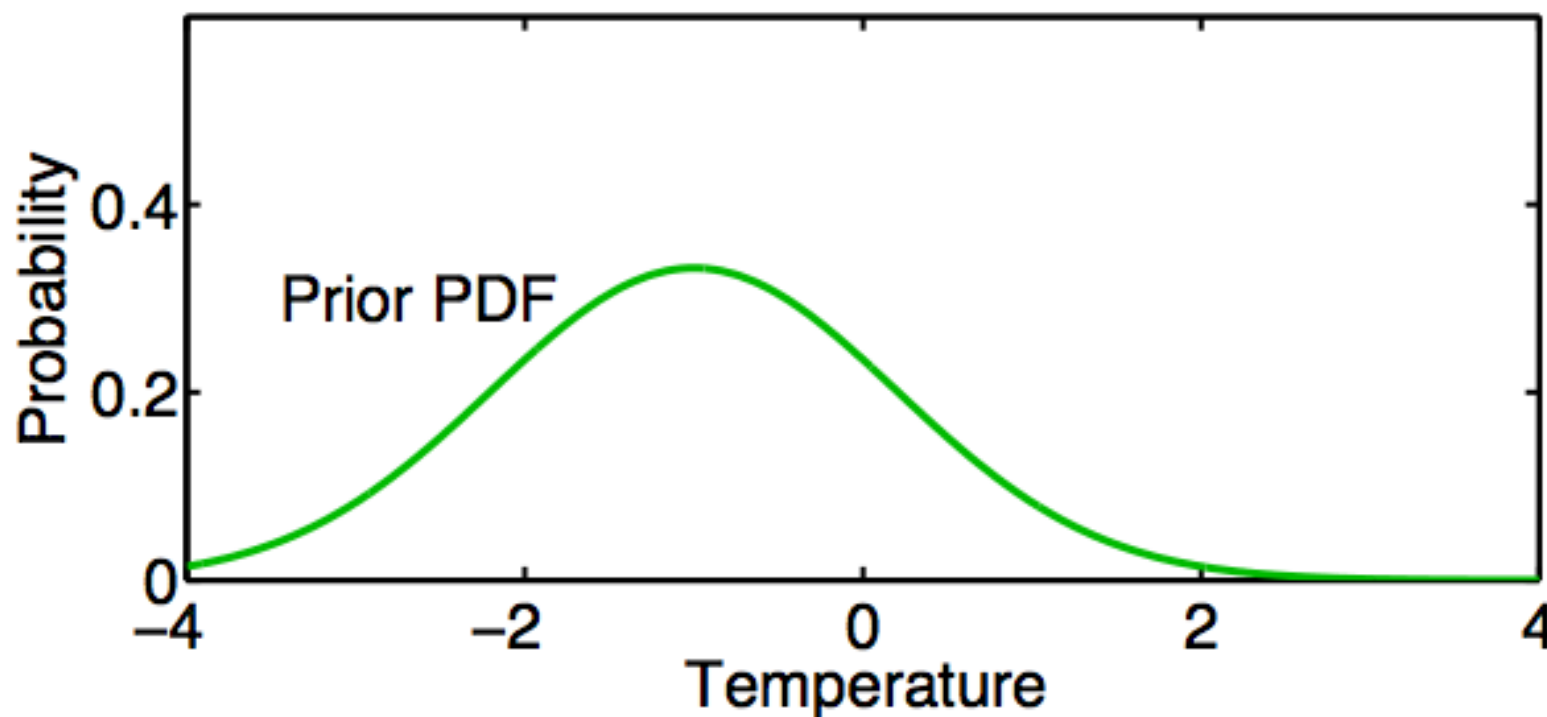
Prior

Posterior: Probability of T given observations and Prior. Also called update or analysis.

Likelihood: Probability that T_o is observed if T is true value and given prior information C.

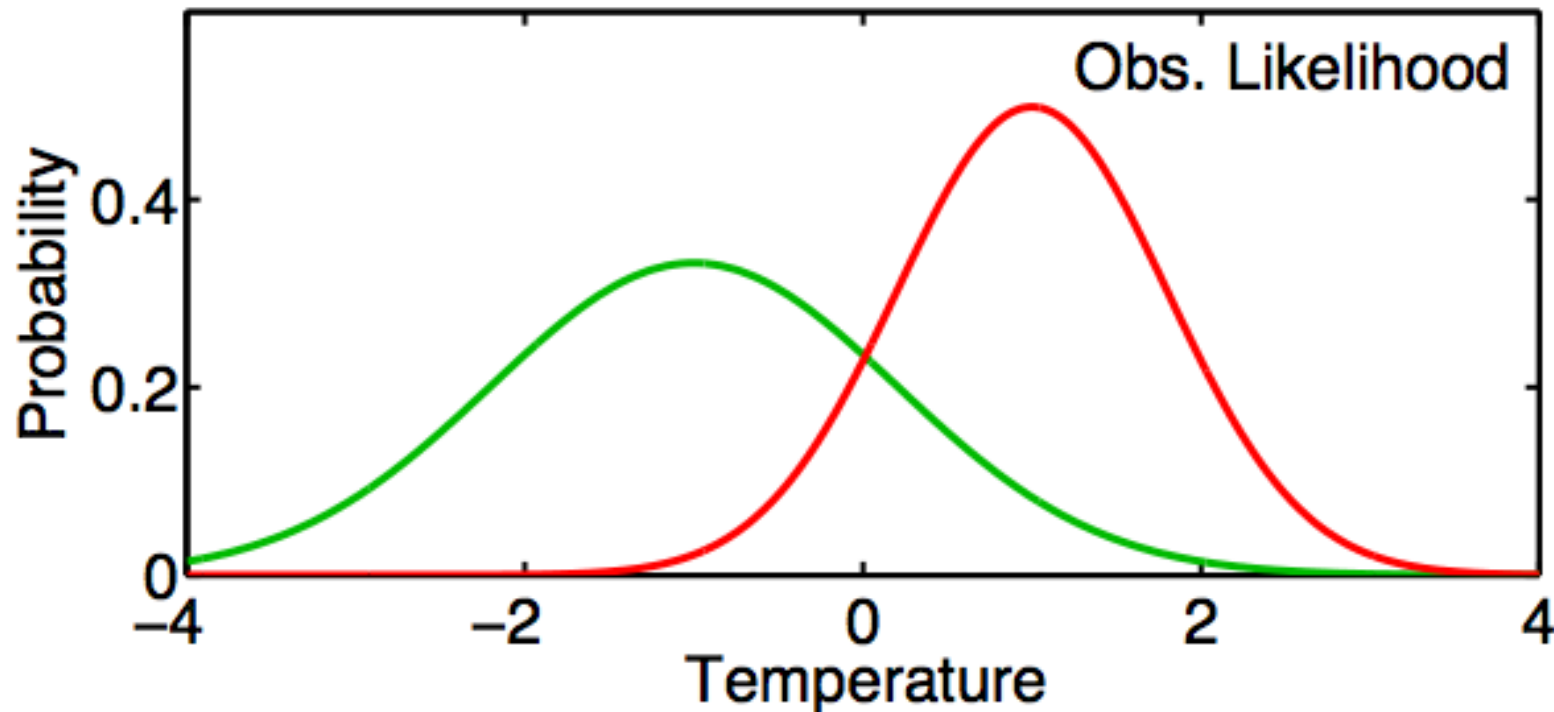
Combining the Prior Estimate and Observation

$$P(T | T_o, C) = \frac{P(T_o | T, C)P(T | C)}{\textit{normalization}}$$



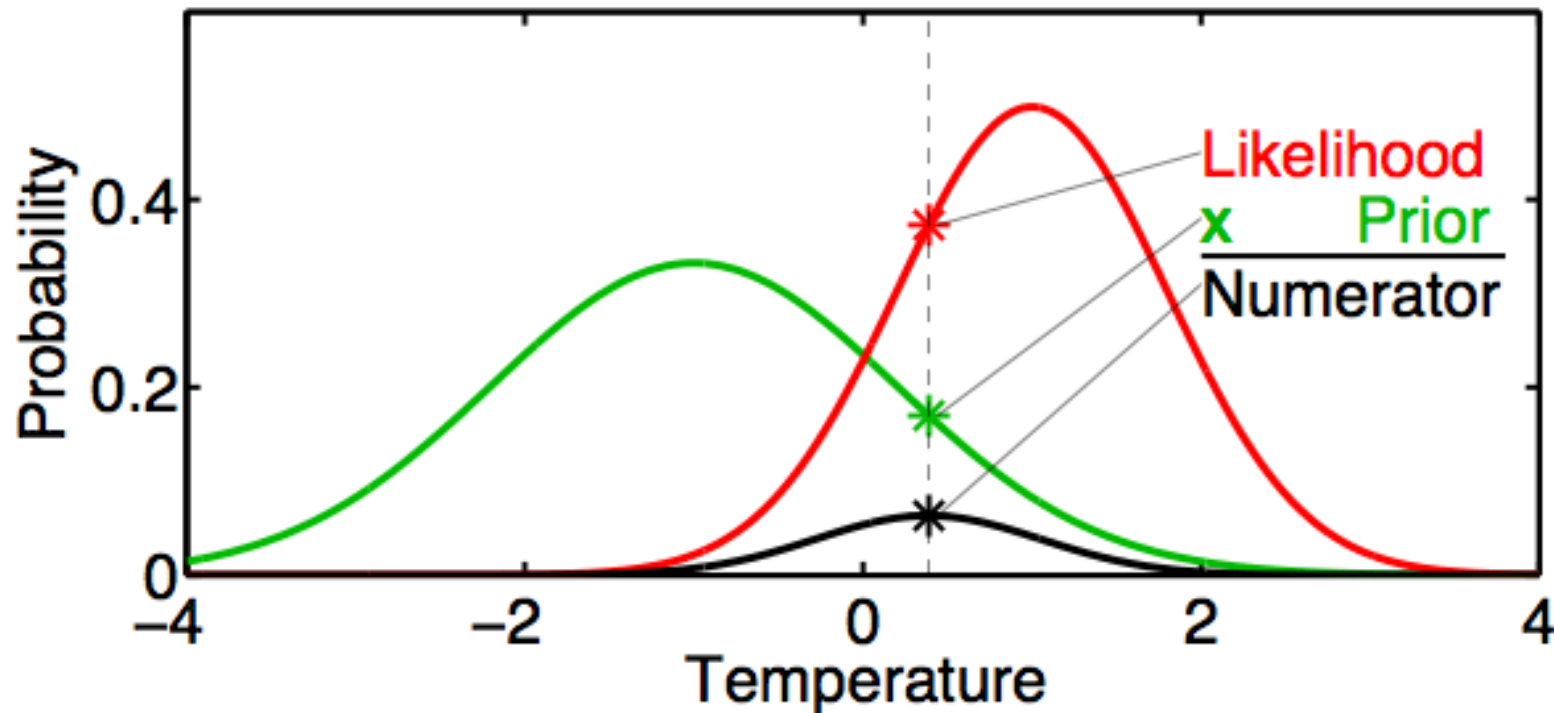
Combining the Prior Estimate and Observation

$$P(T | T_o, C) = \frac{P(T_o | T, C)P(T | C)}{\textit{normalization}}$$



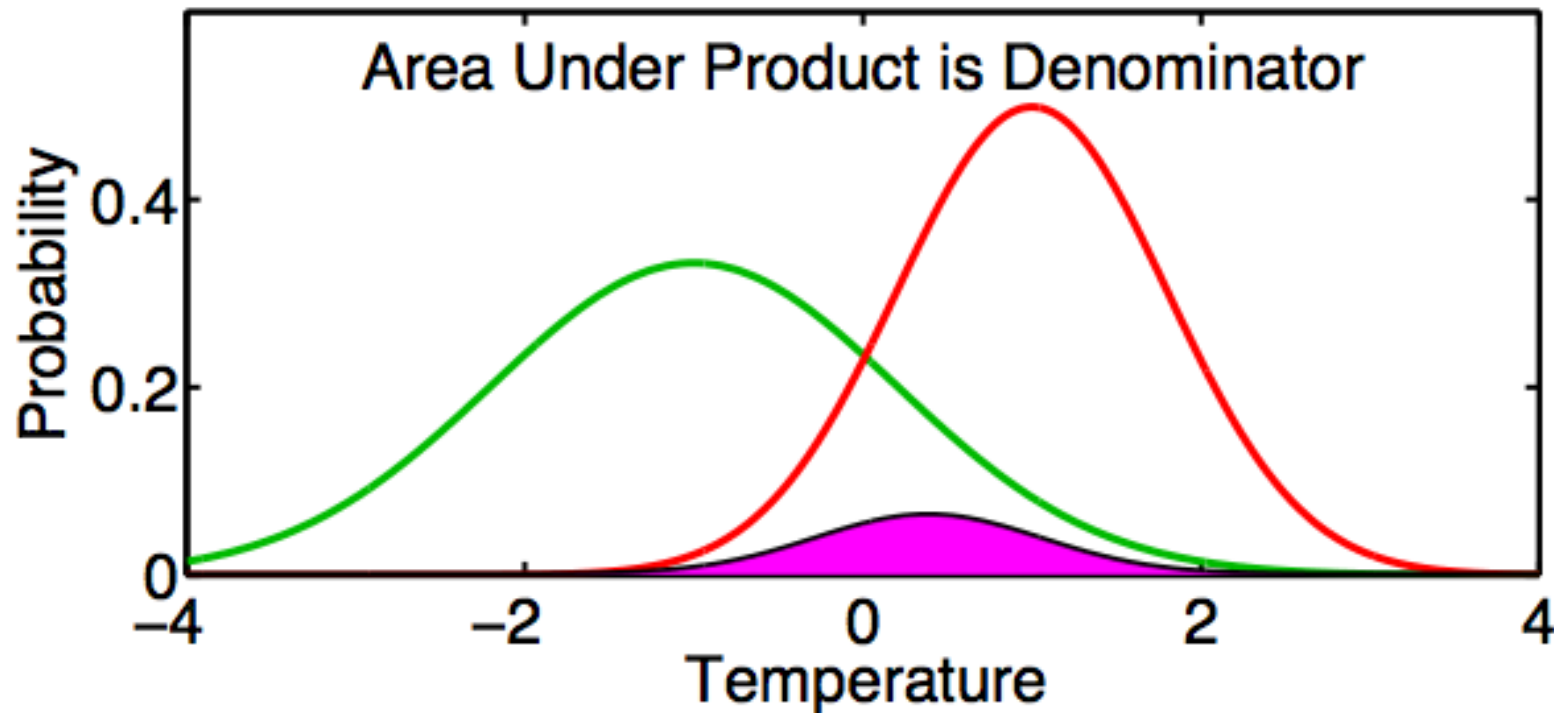
Combining the Prior Estimate and Observation

$$P(T | T_o, C) = \frac{P(T_o | T, C) P(T | C)}{\text{normalization}}$$



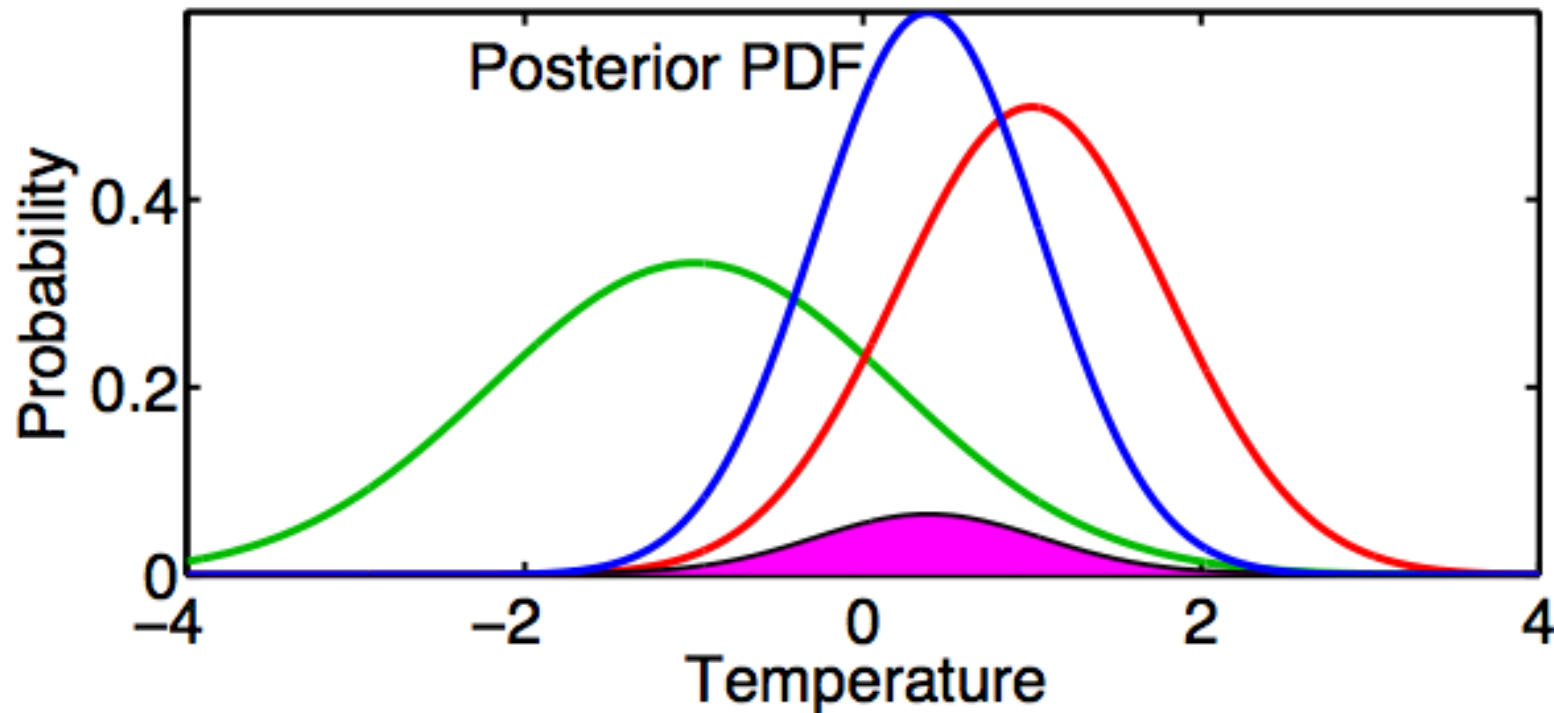
Combining the Prior Estimate and Observation

$$P(T | T_o, C) = \frac{P(T_o | T, C)P(T | C)}{\textit{normalization}}$$



Combining the Prior Estimate and Observation

$$P(T|T_o, C) = \frac{P(T_o|T, C)P(T|C)}{\textit{normalization}}$$



Consistent Color Scheme Throughout Tutorial

Green = Prior

Red = Observation

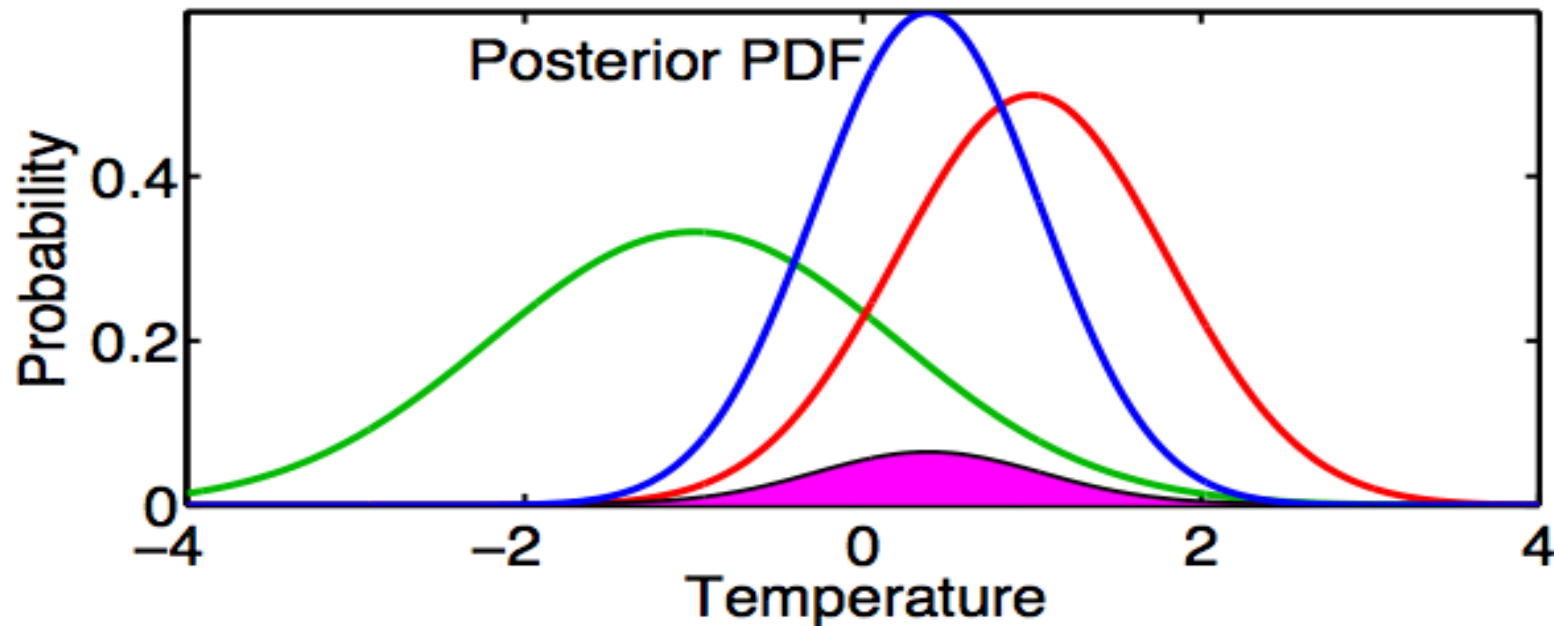
Blue = Posterior

Black = Truth

Combining the Prior Estimate and Observation

$$P(T|T_o, C) = \frac{P(T_o | T, C)P(T|C)}{\textit{normalization}}$$

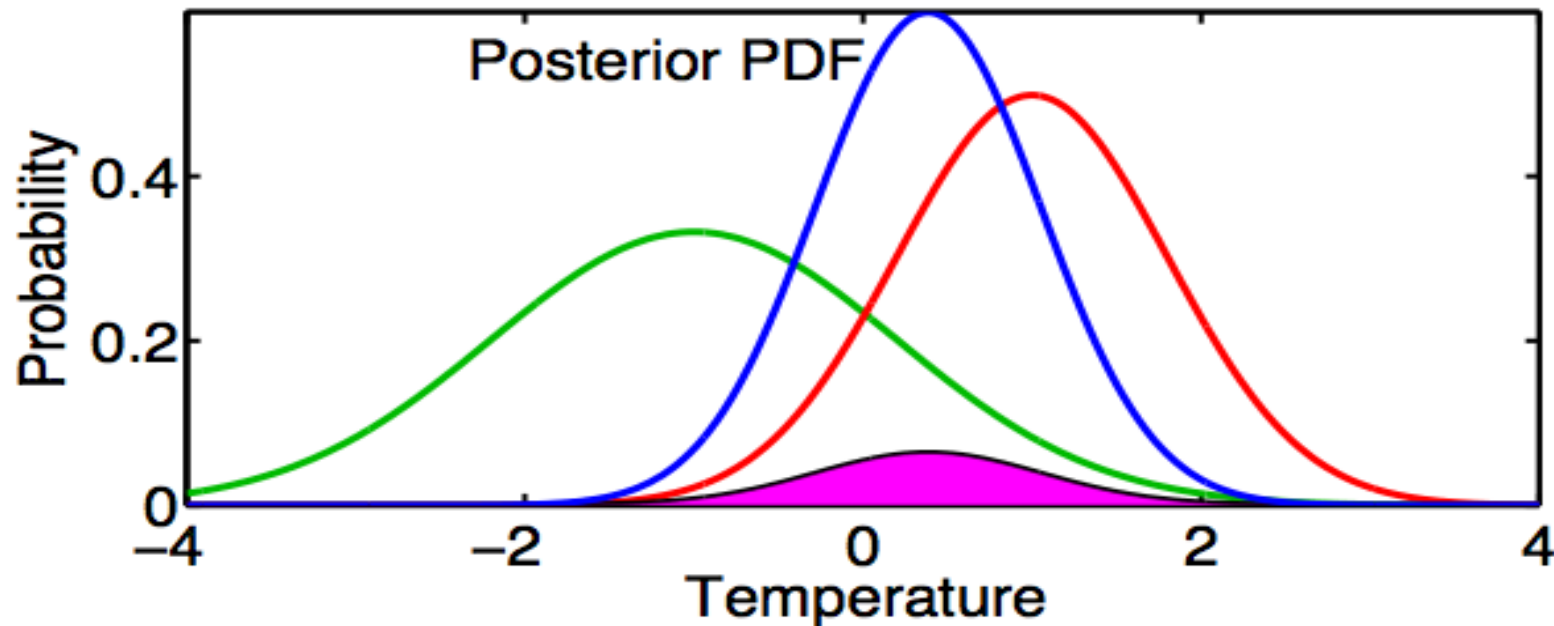
Generally no analytic solution for Posterior.



Combining the Prior Estimate and Observation

$$P(T|T_o, C) = \frac{P(T_o | T, C)P(T|C)}{\textit{normalization}}$$

Gaussian Prior and Likelihood -> Gaussian Posterior



Combining the Prior Estimate and Observation

For Gaussian prior and likelihood...

Prior $P(T | C) = \text{Normal}(T_p, \sigma_p)$

Likelihood $P(T_o | T, C) = \text{Normal}(T_o, \sigma_o)$

Then, Posterior $P(T | T_o, C) = \text{Normal}(T_u, \sigma_u)$

$$\sigma_u = \sqrt{(\sigma_p^{-2} + \sigma_o^{-2})^{-1}}$$

With

$$T_u = \sigma_u^2 [\sigma_p^{-2} T_p + \sigma_o^{-2} T_o]$$

The One-Dimensional Kalman Filter

1. Suppose we have a linear forecast model L
 - A. If temperature at time $t_1 = T_1$, then temperature at $t_2 = t_1 + \Delta t$ is $T_2 = L(T_1)$
 - B. Example: $T_2 = T_1 + \Delta t T_1$

The One-Dimensional Kalman Filter

1. Suppose we have a linear forecast model L .
 - A. If temperature at time $t_1 = T_1$, then temperature at $t_2 = t_1 + \Delta t$ is $T_2 = L(T_1)$.
 - B. Example: $T_2 = T_1 + \Delta t T_1$.
2. If posterior estimate at time t_1 is $Normal(T_{u,1}, \sigma_{u,1})$ then prior at t_2 is $Normal(T_{p,2}, \sigma_{p,2})$.

$$T_{p,2} = T_{u,1} + \Delta t T_{u,1}$$

$$\sigma_{p,2} = (\Delta t + 1) \sigma_{u,1}$$

The One-Dimensional Kalman Filter

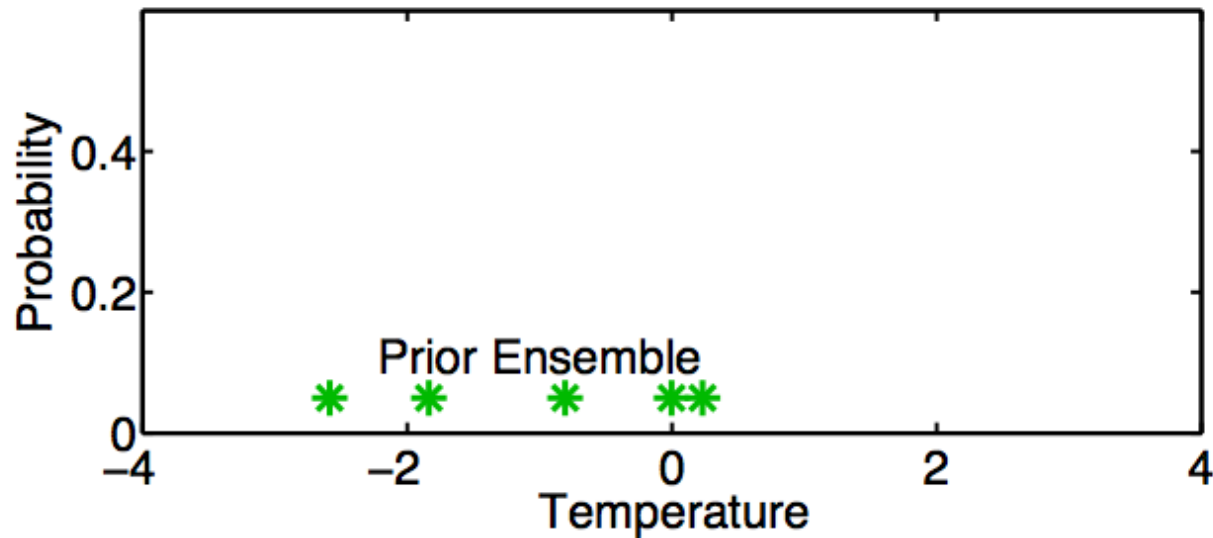
1. Suppose we have a linear forecast model L .
 - A. If temperature at time $t_1 = T_1$, then temperature at $t_2 = t_1 + \Delta t$ is $T_2 = L(T_1)$.
 - B. Example: $T_2 = T_1 + \Delta t T_1$.
2. If posterior estimate at time t_1 is $Normal(T_{u,1}, \sigma_{u,1})$ then prior at t_2 is $Normal(T_{p,2}, \sigma_{p,2})$.
3. Given an observation at t_2 with distribution $Normal(t_o, \sigma_o)$ the likelihood is also $Normal(t_o, \sigma_o)$.

The One-Dimensional Kalman Filter

1. Suppose we have a linear forecast model L .
 - A. If temperature at time $t_1 = T_1$, then temperature at $t_2 = t_1 + \Delta t$ is $T_2 = L(T_1)$.
 - B. Example: $T_2 = T_1 + \Delta t T_1$.
2. If posterior estimate at time t_1 is $Normal(T_{u,1}, \sigma_{u,1})$ then prior at t_2 is $Normal(T_{p,2}, \sigma_{p,2})$.
3. Given an observation at t_2 with distribution $Normal(t_o, \sigma_o)$ the likelihood is also $Normal(t_o, \sigma_o)$.
4. The posterior at t_2 is $Normal(T_{u,2}, \sigma_{u,2})$ where $T_{u,2}$ and $\sigma_{u,2}$ come from page 18.

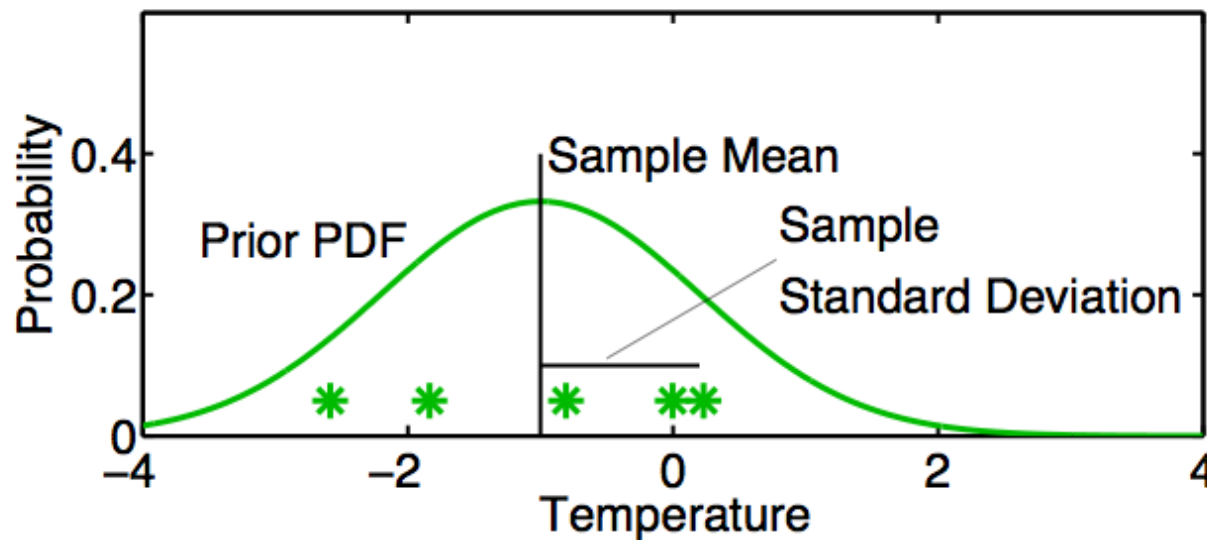
A One-Dimensional Ensemble Kalman Filter

Represent a prior pdf by a sample (ensemble) of N values:



A One-Dimensional Ensemble Kalman Filter

Represent a prior pdf by a sample (ensemble) of N values:



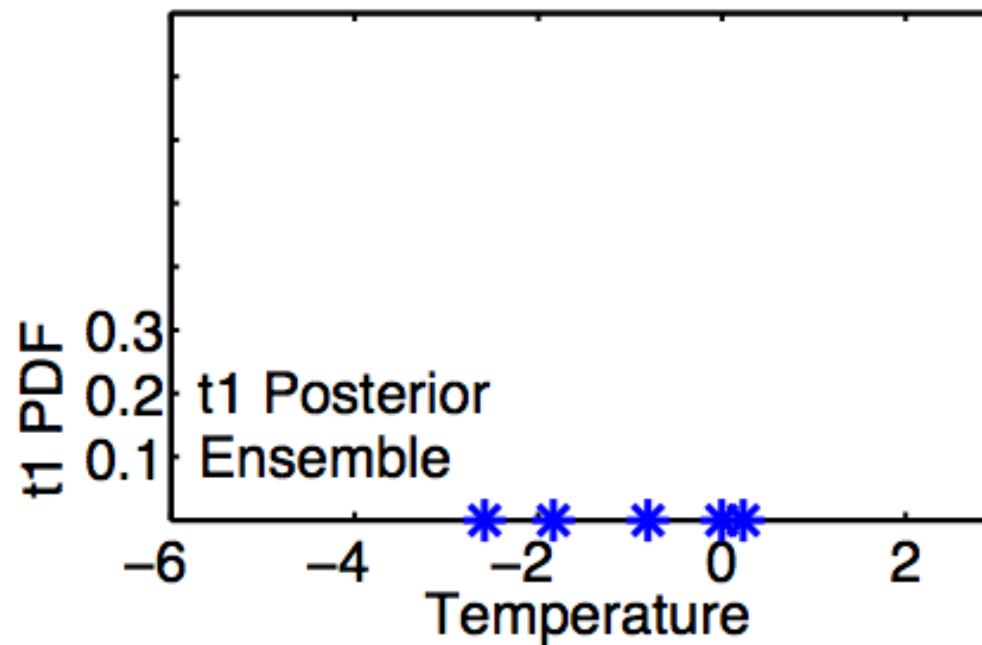
Use sample mean $\bar{T} = \sum_{n=1}^N T_n / N$

and sample standard deviation $\sigma_T = \sqrt{\sum_{n=1}^N (T_n - \bar{T})^2 / (N - 1)}$

to determine a corresponding continuous distribution $Normal(\bar{T}, \sigma_T)$

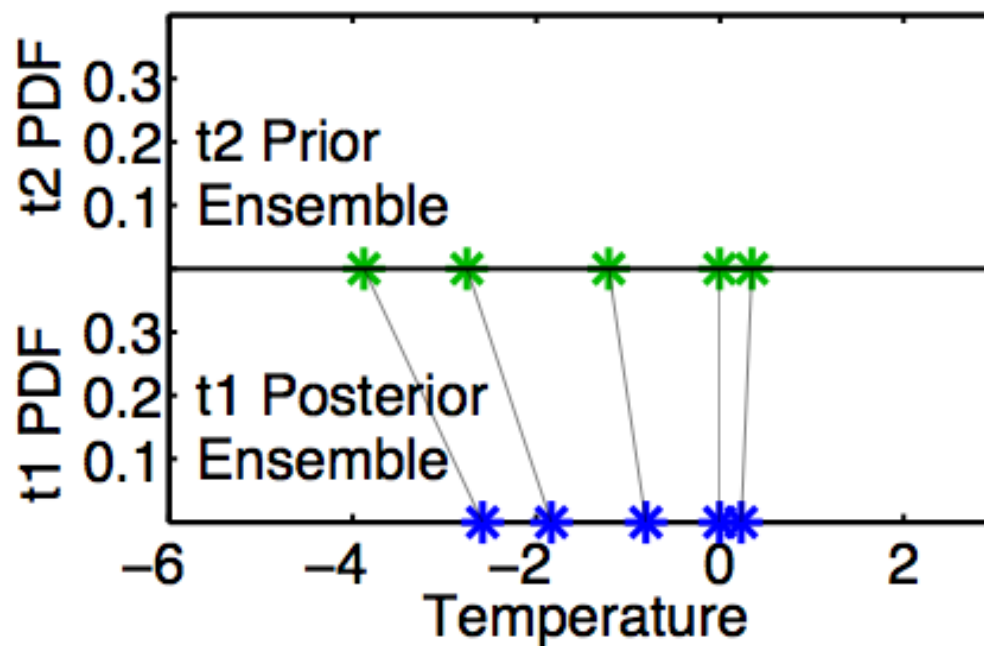
A One-Dimensional Ensemble Kalman Filter: Model Advance

If posterior ensemble at time t_1 is $T_{1,n}$, $n = 1, \dots, N$



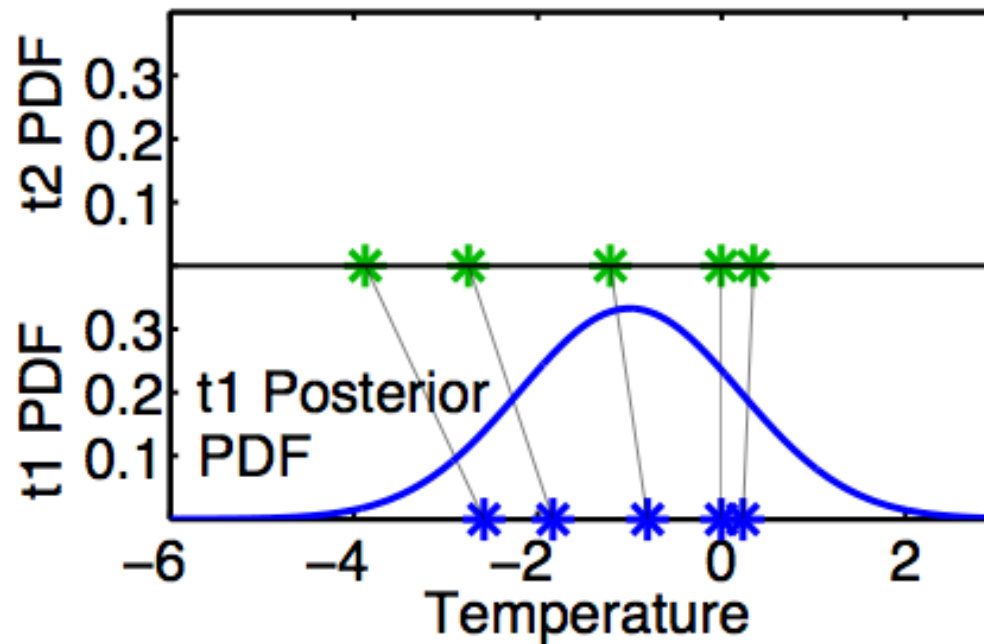
A One-Dimensional Ensemble Kalman Filter: Model Advance

If posterior ensemble at time t_1 is $T_{1,n}$, $n = 1, \dots, N$,
advance each member to time t_2 with model, $T_{2,n} = L(T_{1,n})$ $n = 1, \dots, N$.



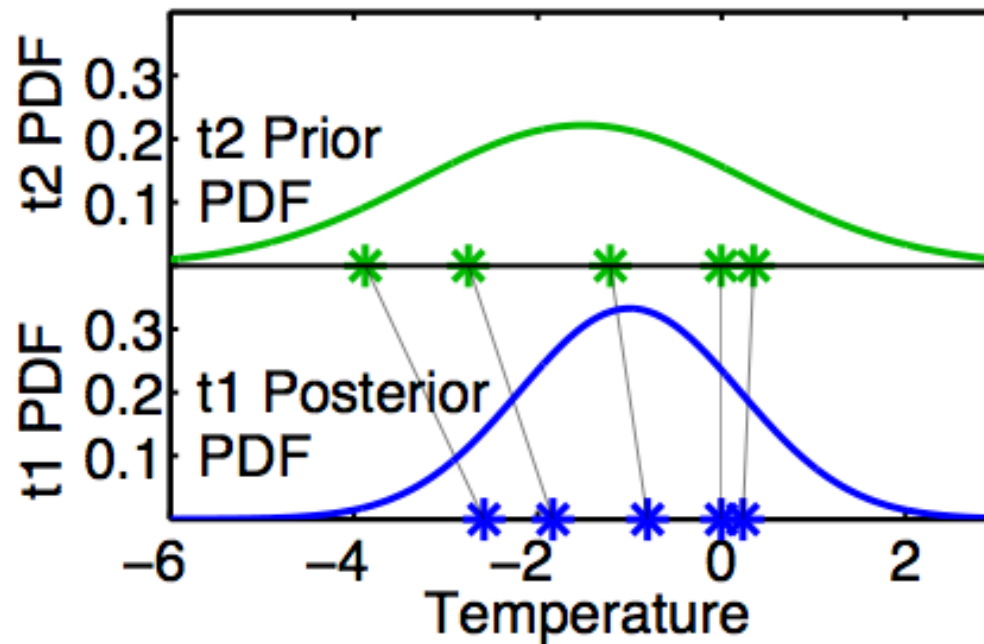
A One-Dimensional Ensemble Kalman Filter: Model Advance

Same as advancing continuous pdf at time t_1 ...

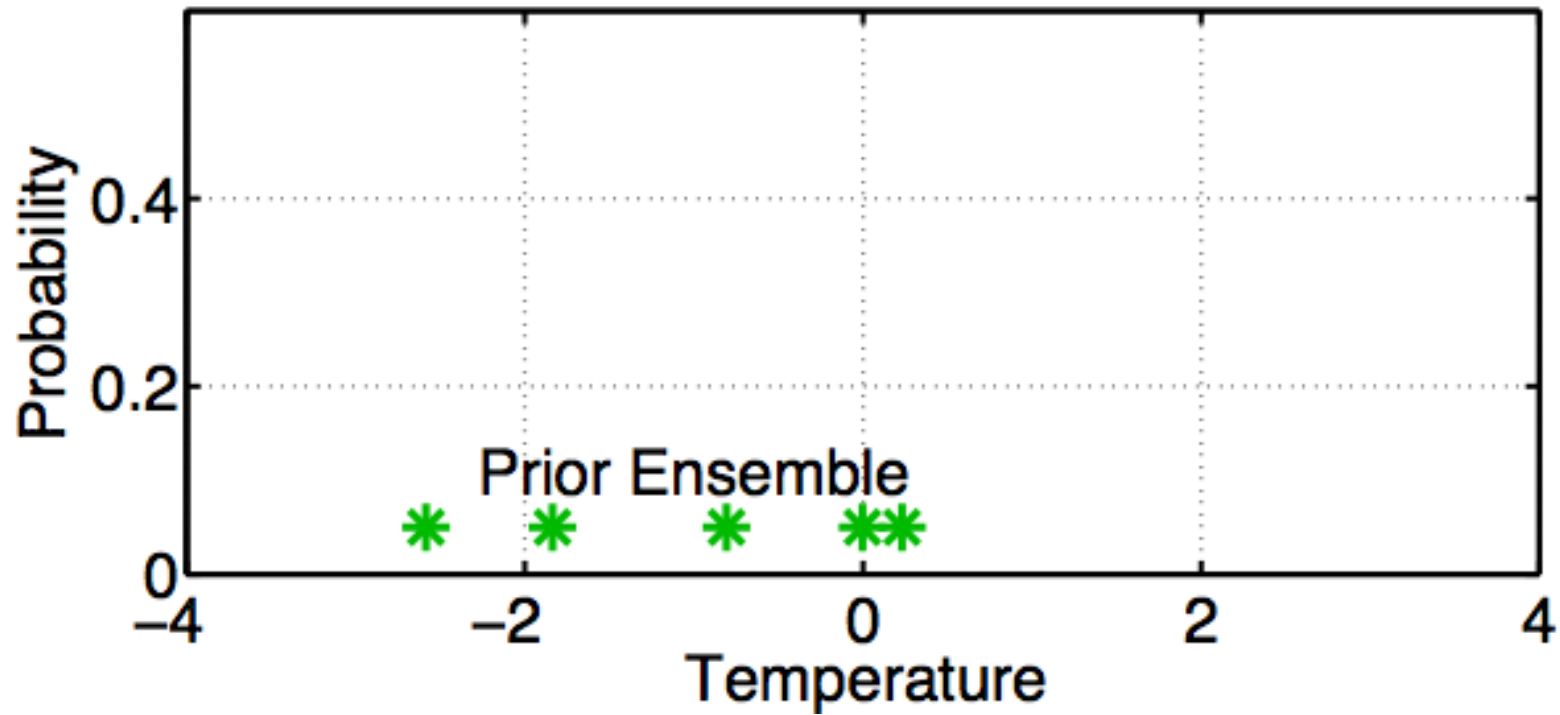


A One-Dimensional Ensemble Kalman Filter: Model Advance

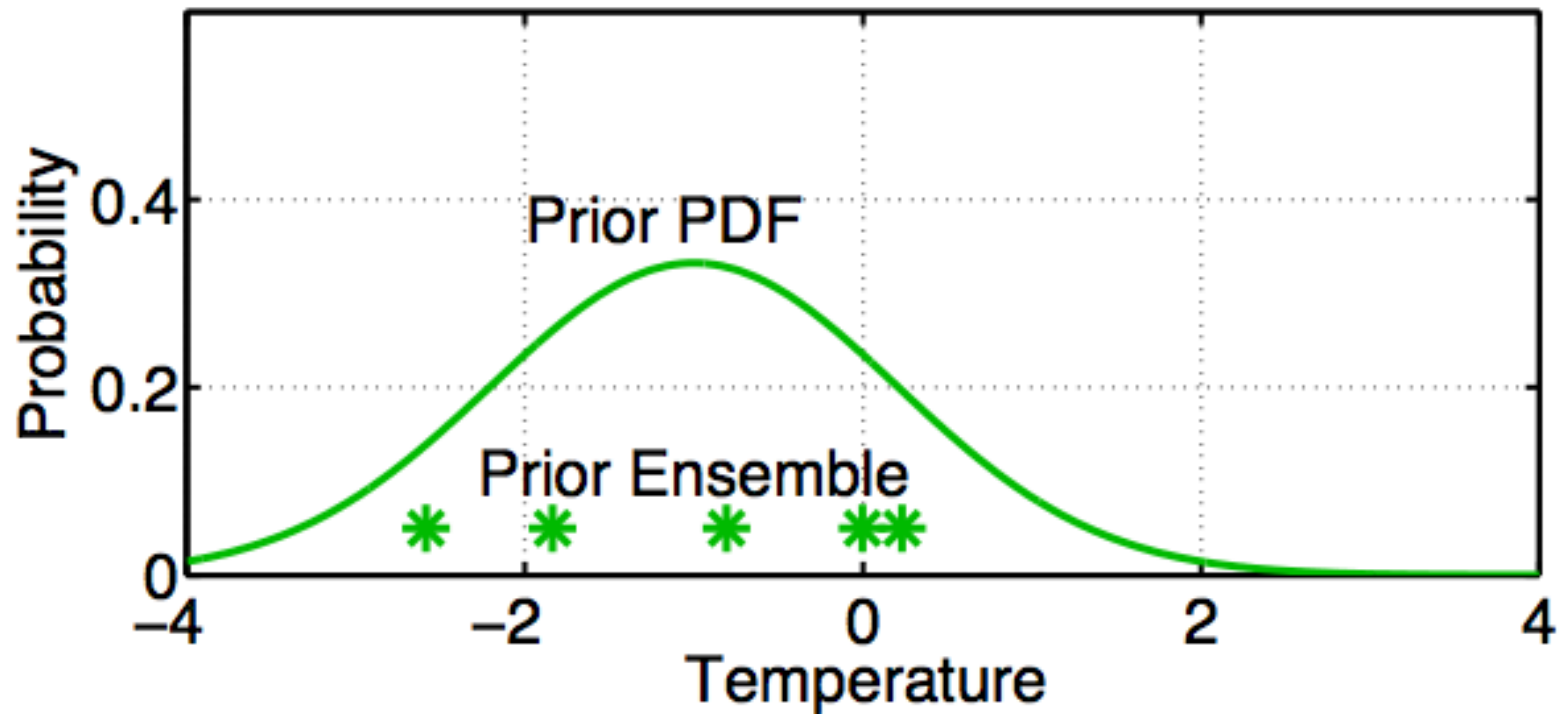
Same as advancing continuous pdf at time t_1
to time t_2 with model L.



A One-Dimensional Ensemble Kalman Filter: Assimilating an Observation

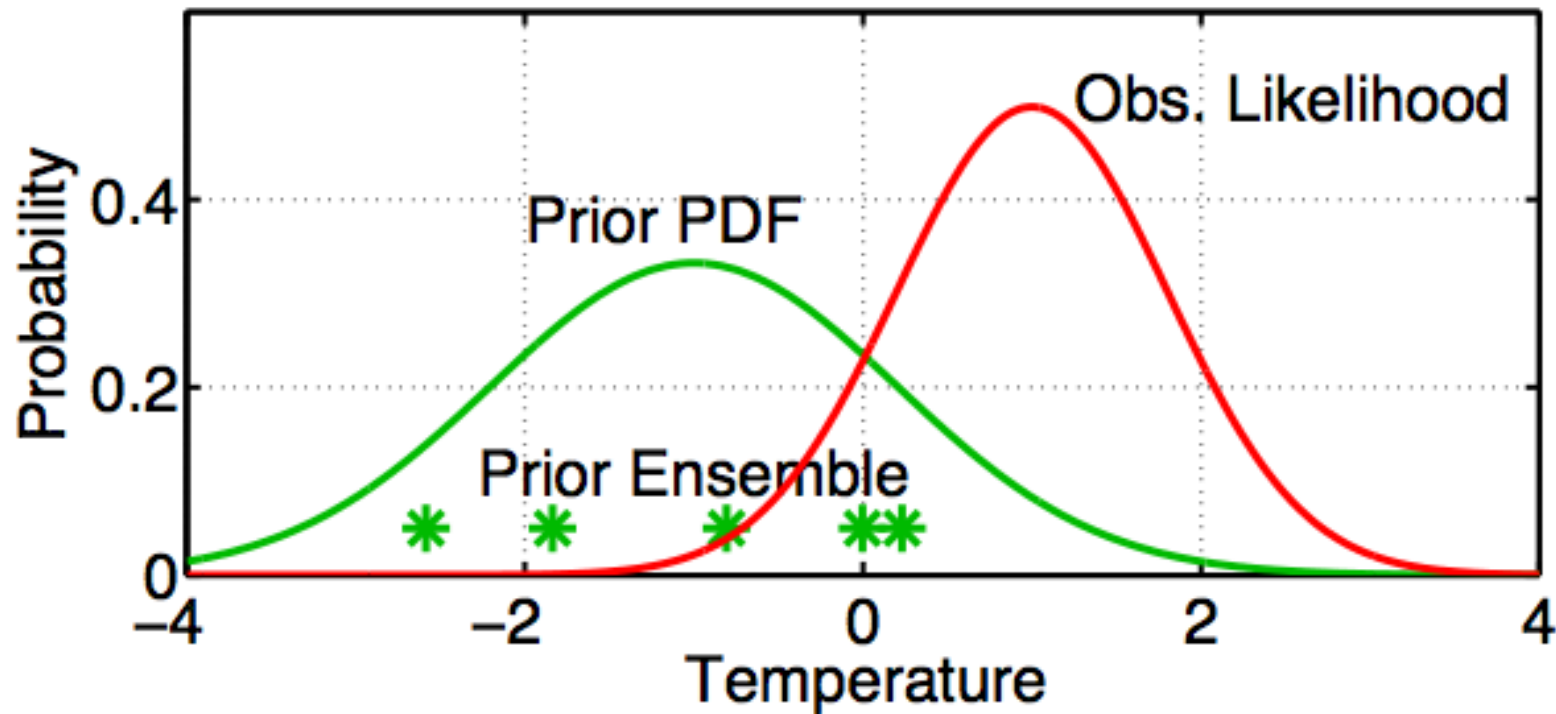


A One-Dimensional Ensemble Kalman Filter: Assimilating an Observation



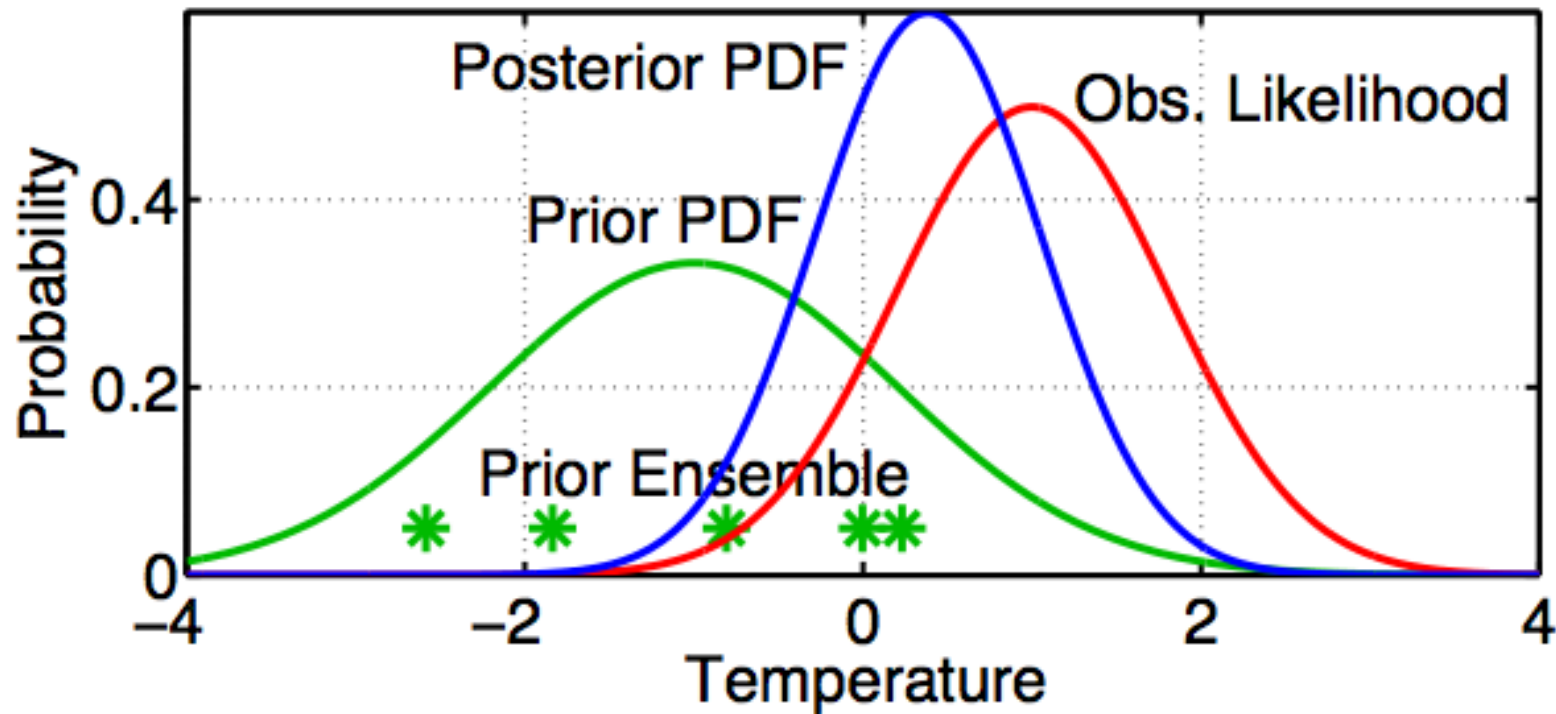
Fit a Gaussian to the sample.

A One-Dimensional Ensemble Kalman Filter: Assimilating an Observation



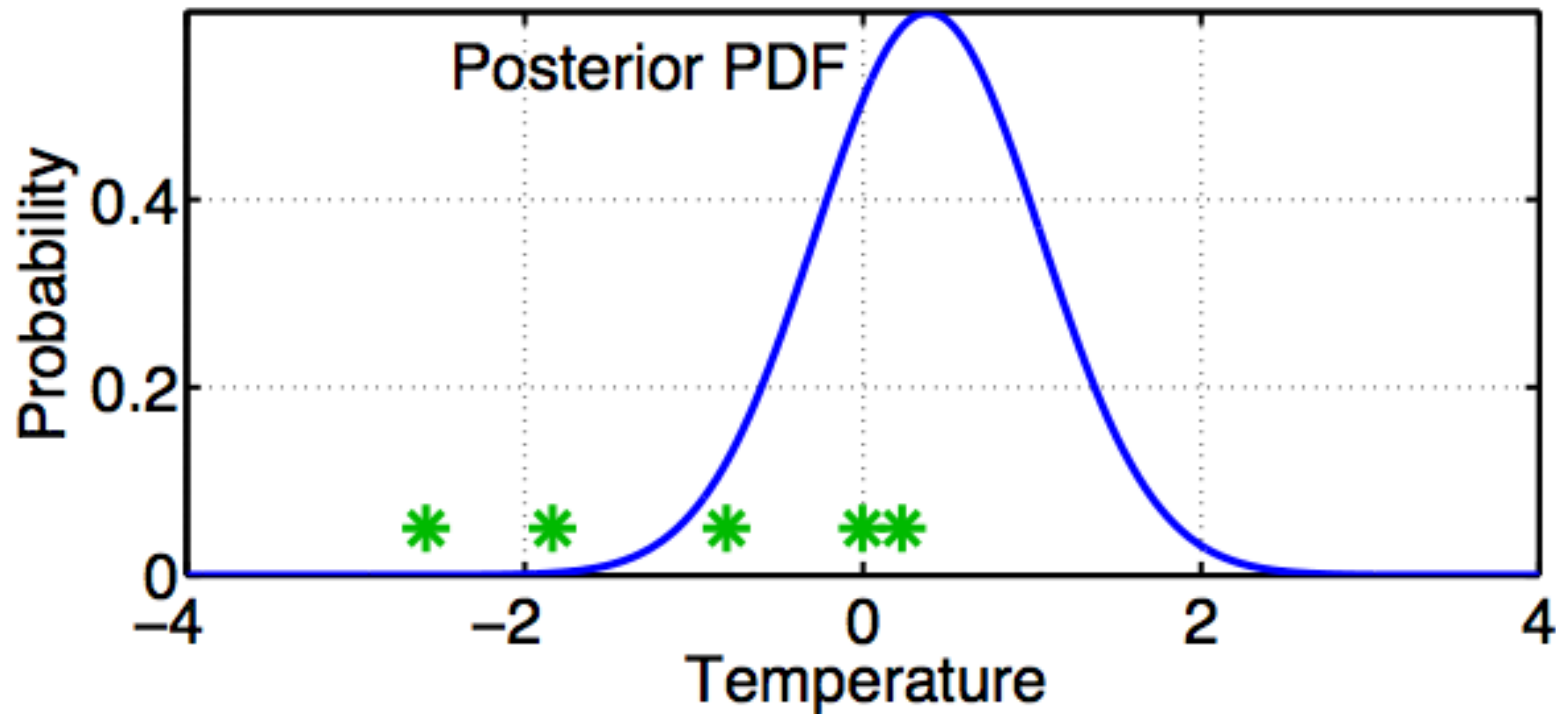
Get the observation likelihood.

A One-Dimensional Ensemble Kalman Filter: Assimilating an Observation



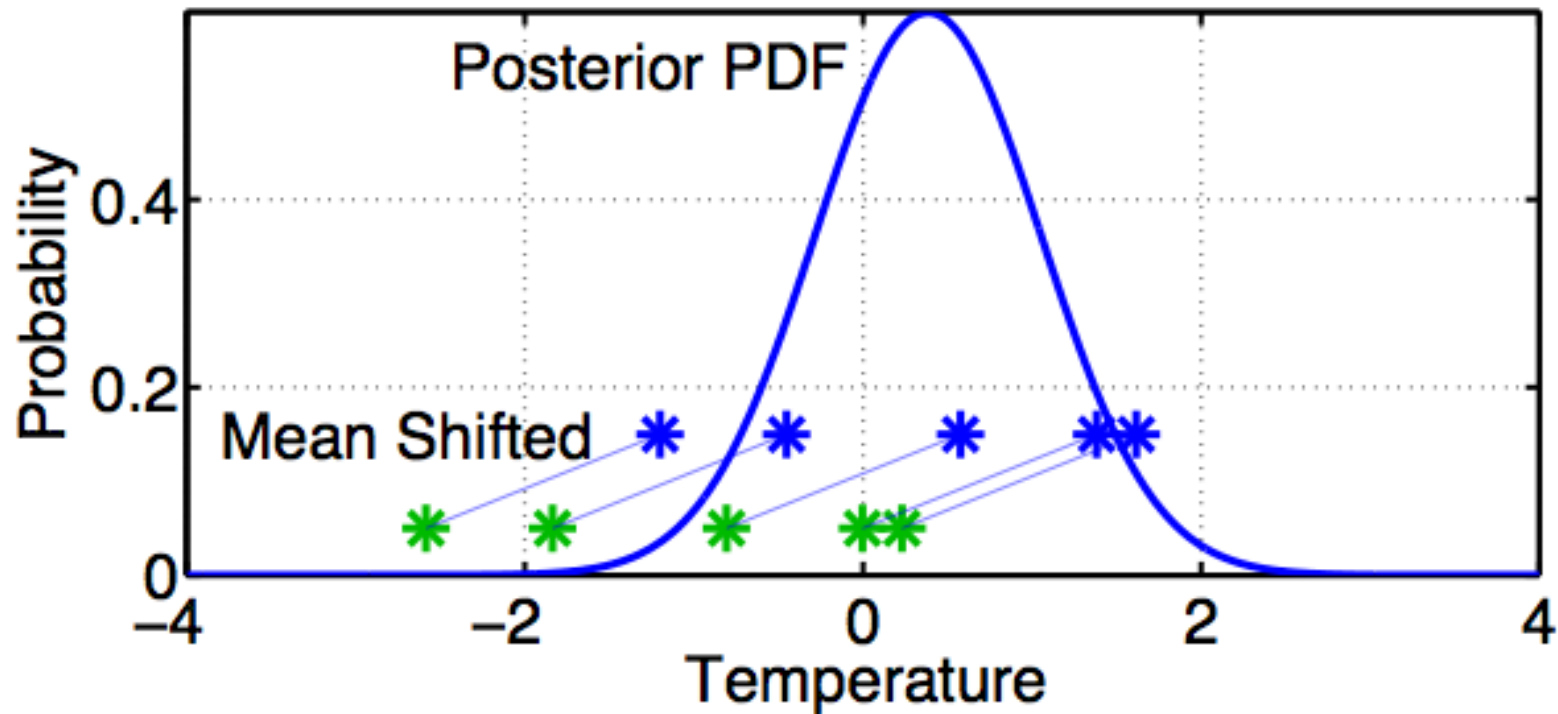
Compute the continuous posterior PDF.

A One-Dimensional Ensemble Kalman Filter: Assimilating an Observation



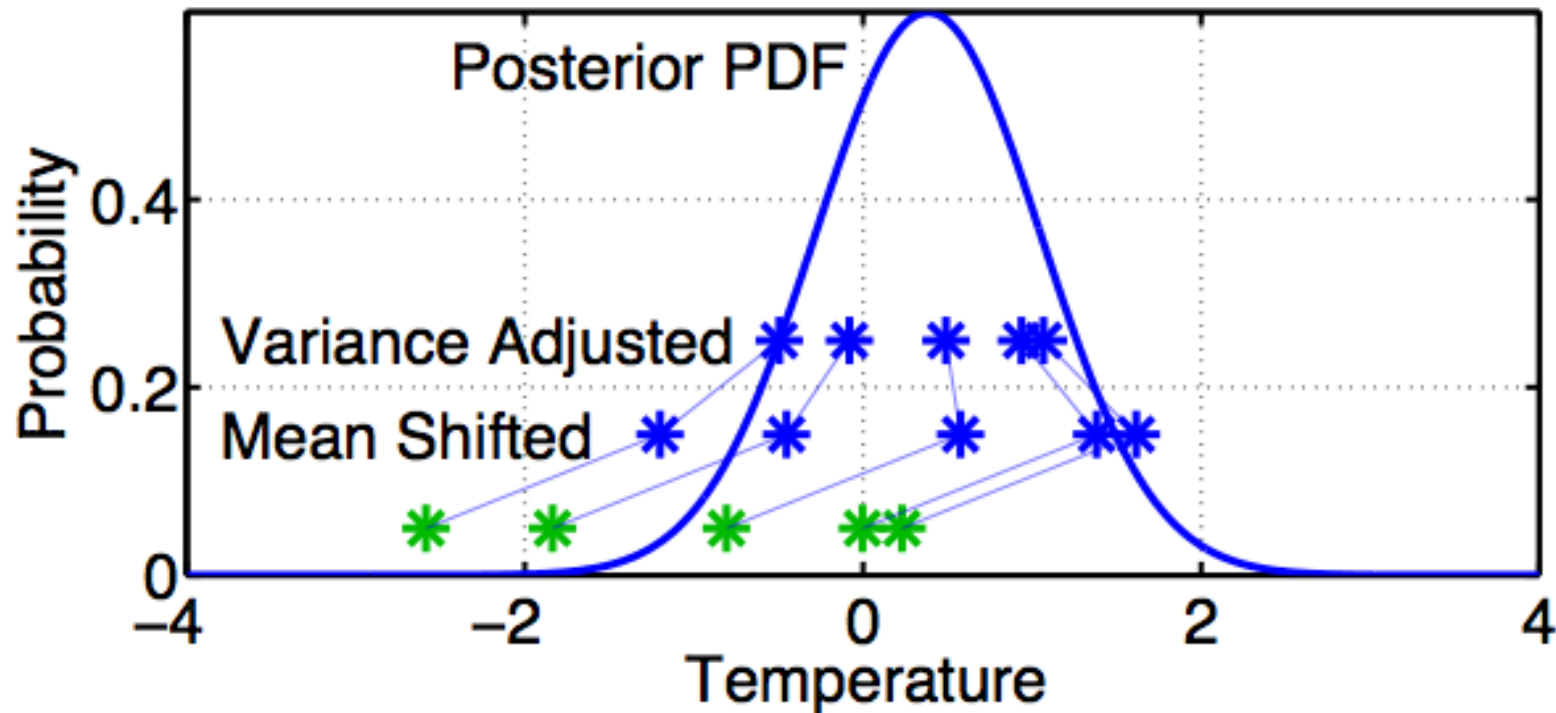
Use a deterministic algorithm to 'adjust' the ensemble.

A One-Dimensional Ensemble Kalman Filter: Assimilating an Observation



First, 'shift' the ensemble to have the exact mean of the posterior.

A One-Dimensional Ensemble Kalman Filter: Assimilating an Observation



First, 'shift' the ensemble to have the exact mean of the posterior.
Second, linearly contract to have the exact variance of the posterior.
Sample statistics are identical to Kalman filter.

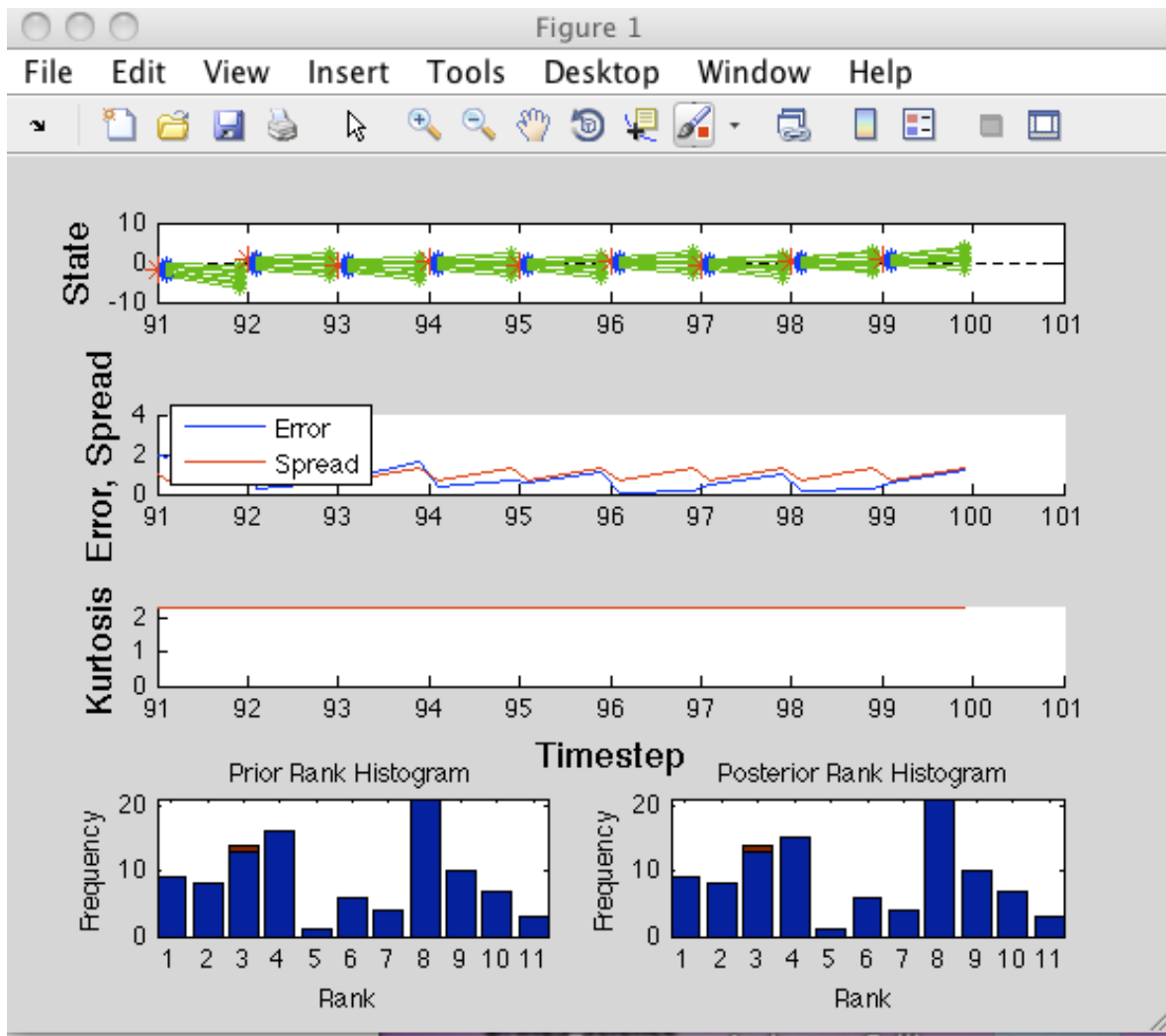
UQ from an Ensemble Kalman Filter

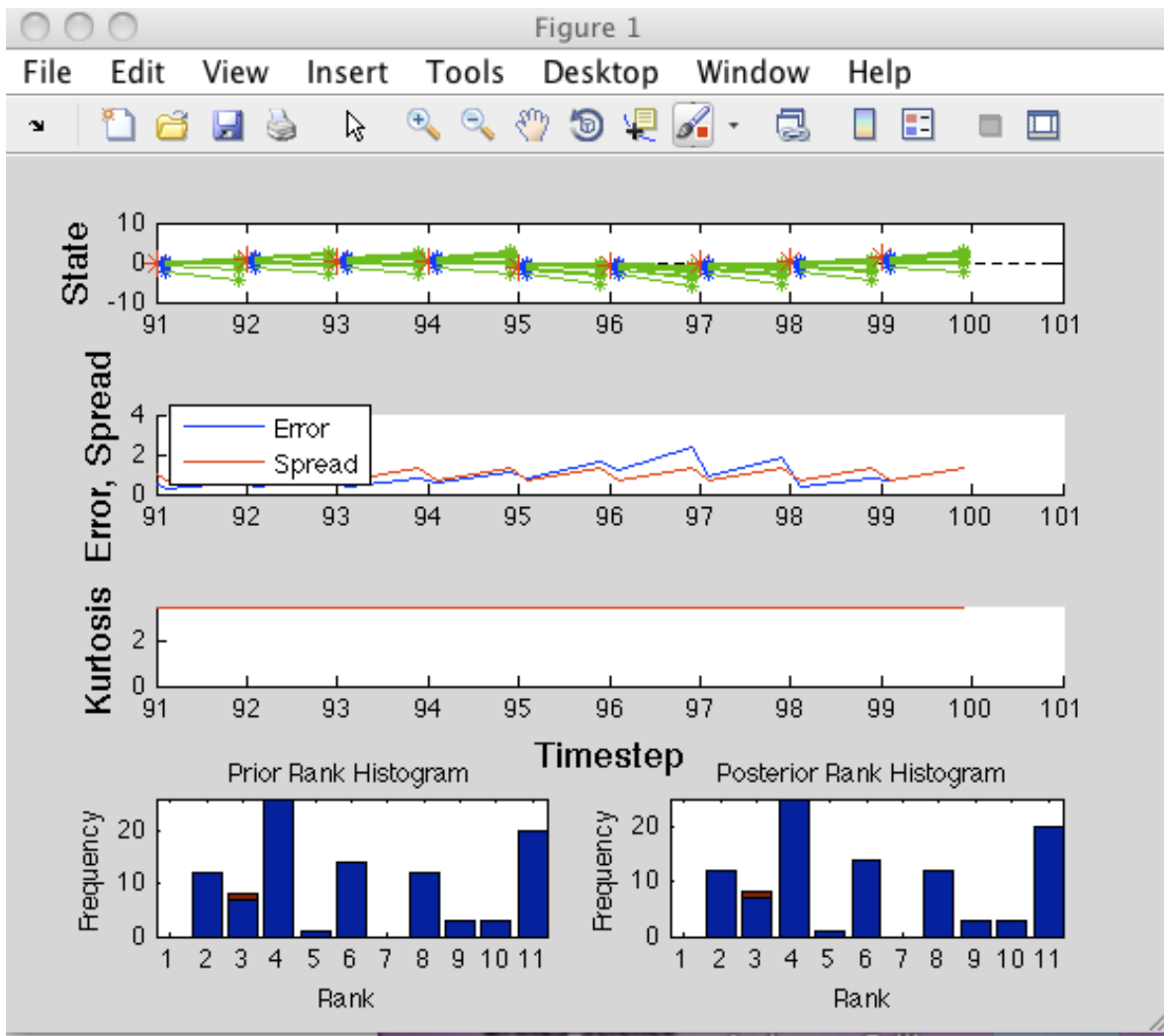
- (Ensemble) KF optimal for linear model, gaussian likelihood.
- In KF, only mean and variance have meaning.
- **Variance defines the uncertainty.**

- Ensemble allows computation of many other statistics.
- What do they mean? Not entirely clear.

- Example: Kurtosis. Completely constrained by initial ensemble.
It is problem specific whether this is even defined!

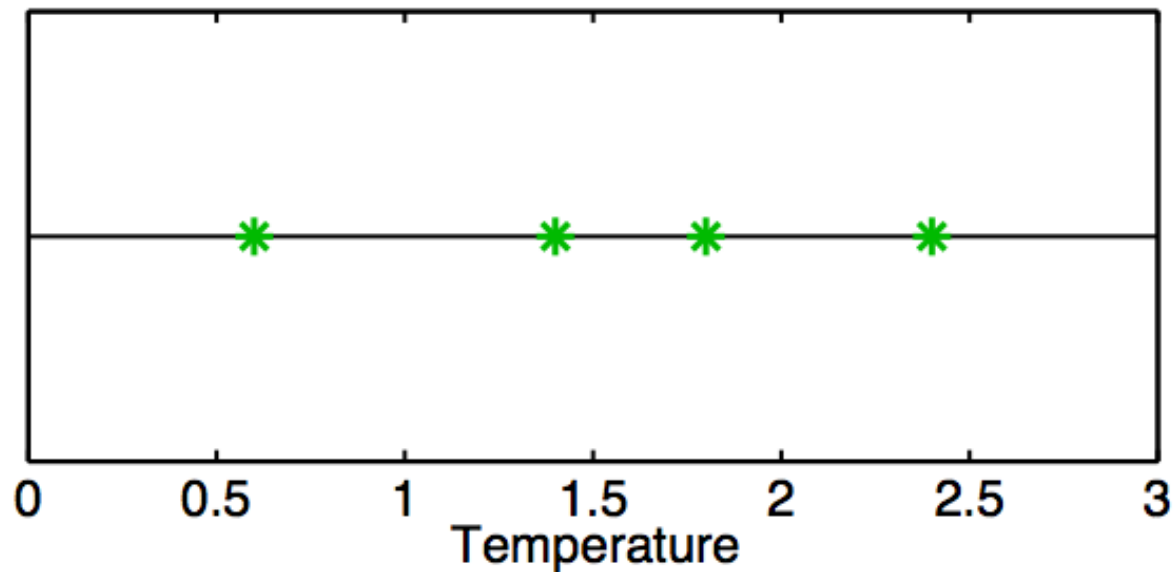
(See example set 1 in Appendix)





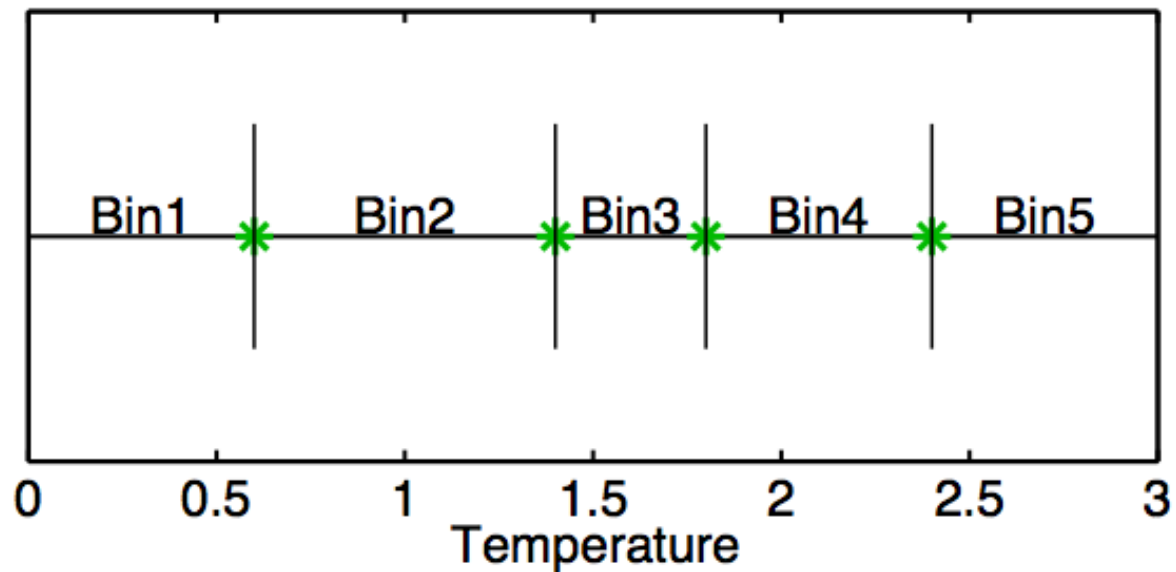
The Rank Histogram: Evaluating Ensemble Performance

Draw 5 values from a real-valued distribution.
Call the first 4 'ensemble members'.



The Rank Histogram: Evaluating Ensemble Performance

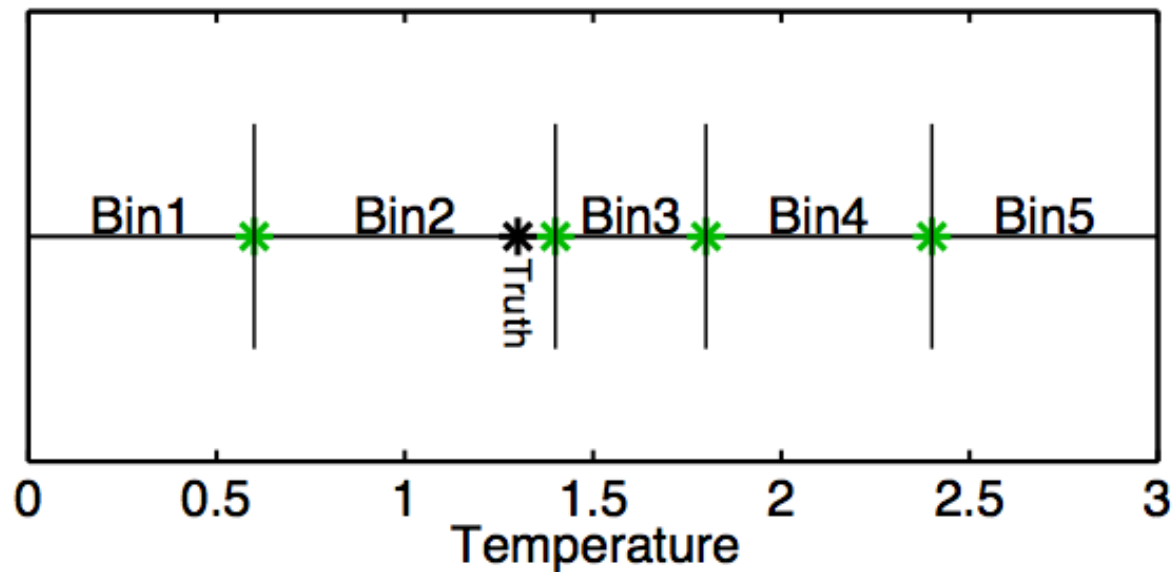
These partition the real line into 5 bins.



The Rank Histogram: Evaluating Ensemble Performance

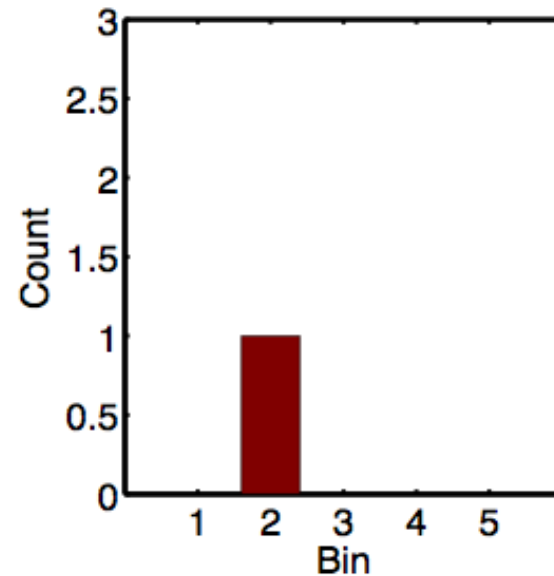
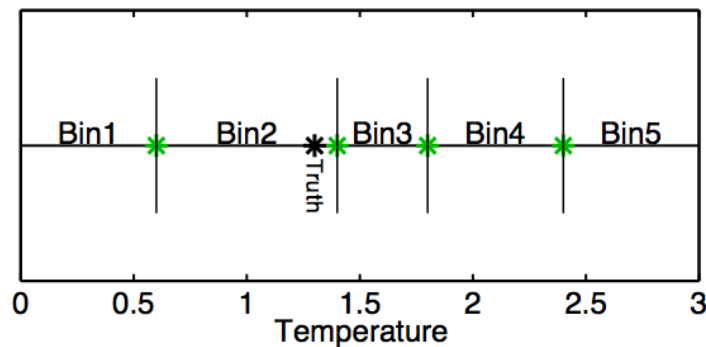
Call the 5th draw the 'truth'.

1/5 chance that this is in any given bin.



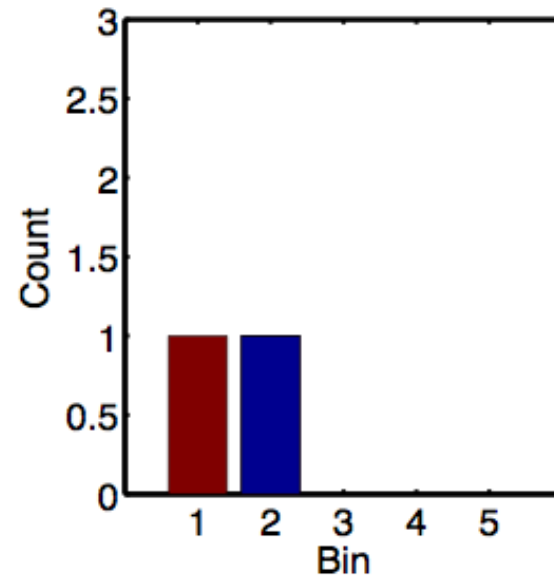
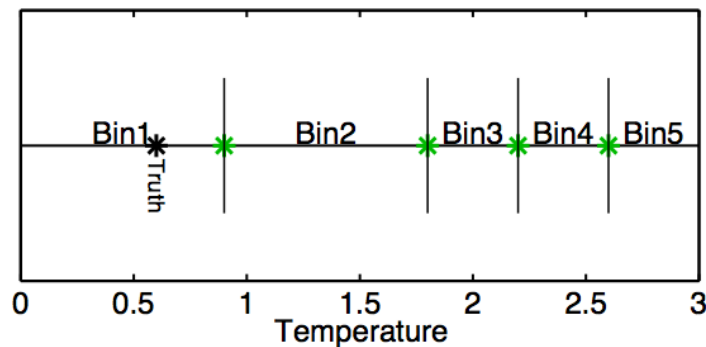
The Rank Histogram: Evaluating Ensemble Performance

Rank histogram shows the frequency of the truth in each bin over many assimilations.



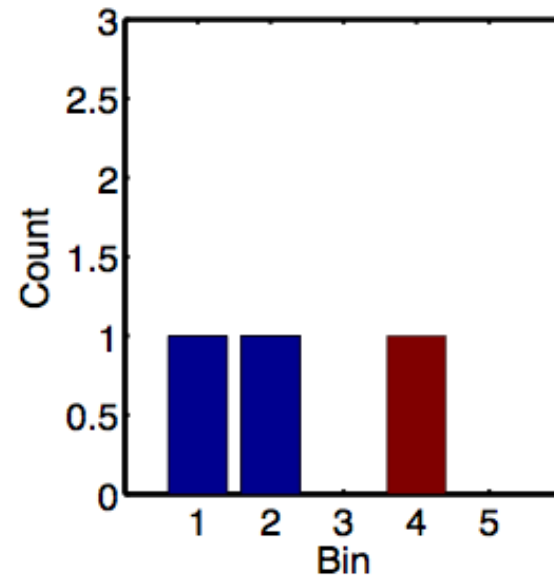
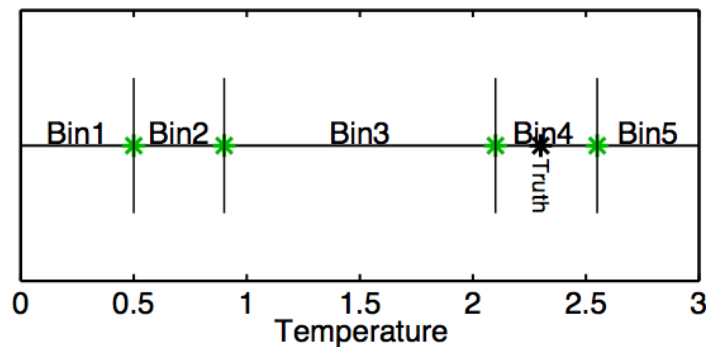
The Rank Histogram: Evaluating Ensemble Performance

Rank histogram shows the frequency of the truth in each bin over many assimilations.



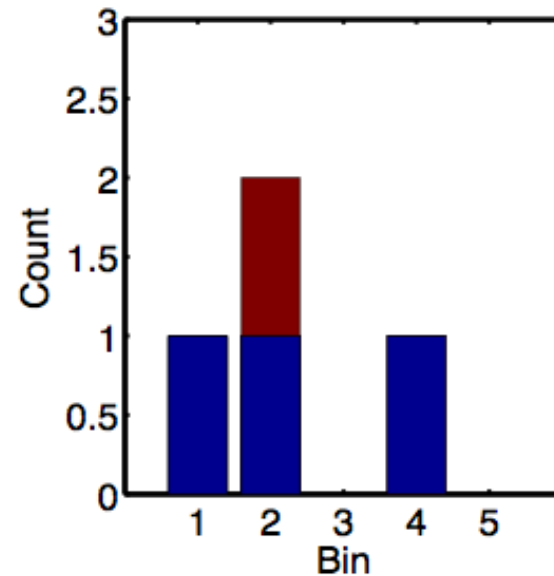
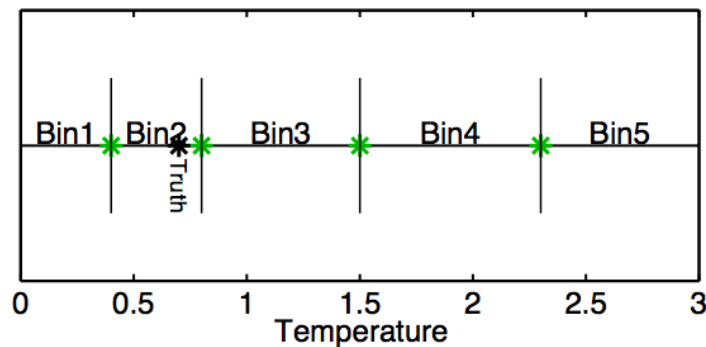
The Rank Histogram: Evaluating Ensemble Performance

Rank histogram shows the frequency of the truth in each bin over many assimilations.



The Rank Histogram: Evaluating Ensemble Performance

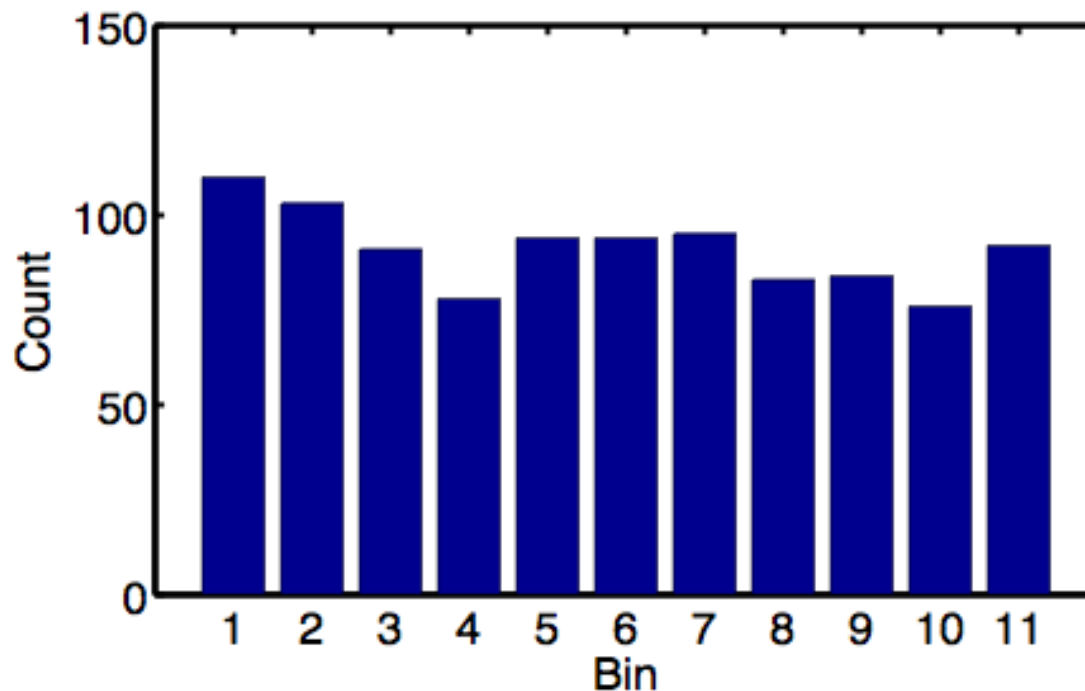
Rank histogram shows the frequency of the truth in each bin over many assimilations.



The Rank Histogram: Evaluating Ensemble Performance

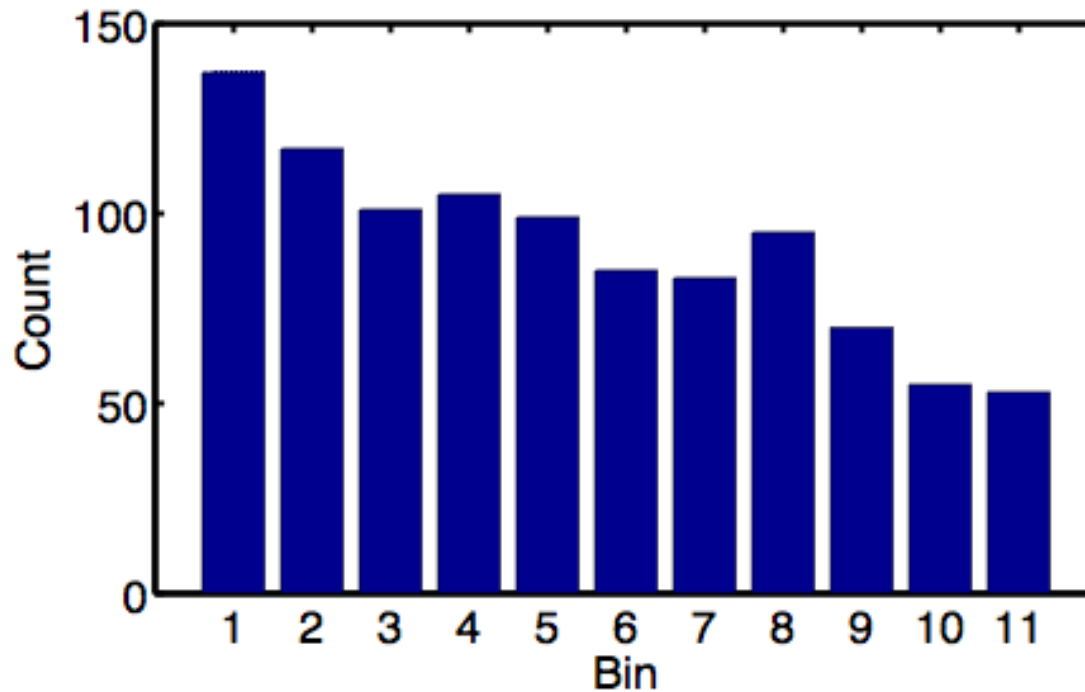
Rank histograms for good ensembles should be uniform (caveat sampling noise).

Want truth to look like random draw from ensemble.



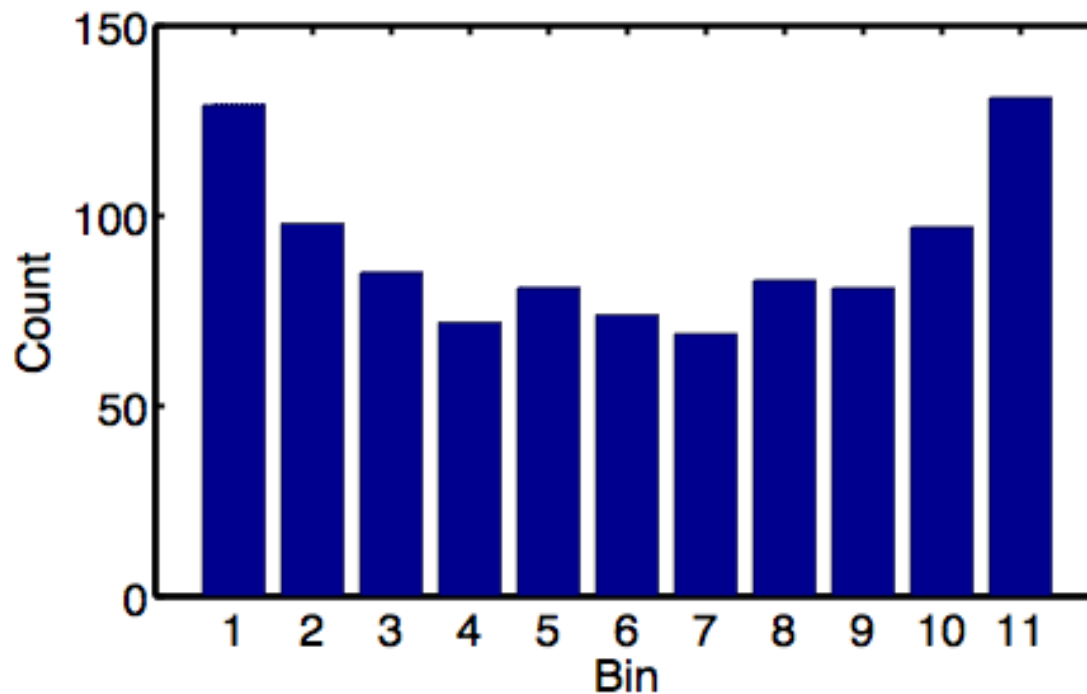
The Rank Histogram: Evaluating Ensemble Performance

A biased ensemble leads to skewed histograms.



The Rank Histogram: Evaluating Ensemble Performance

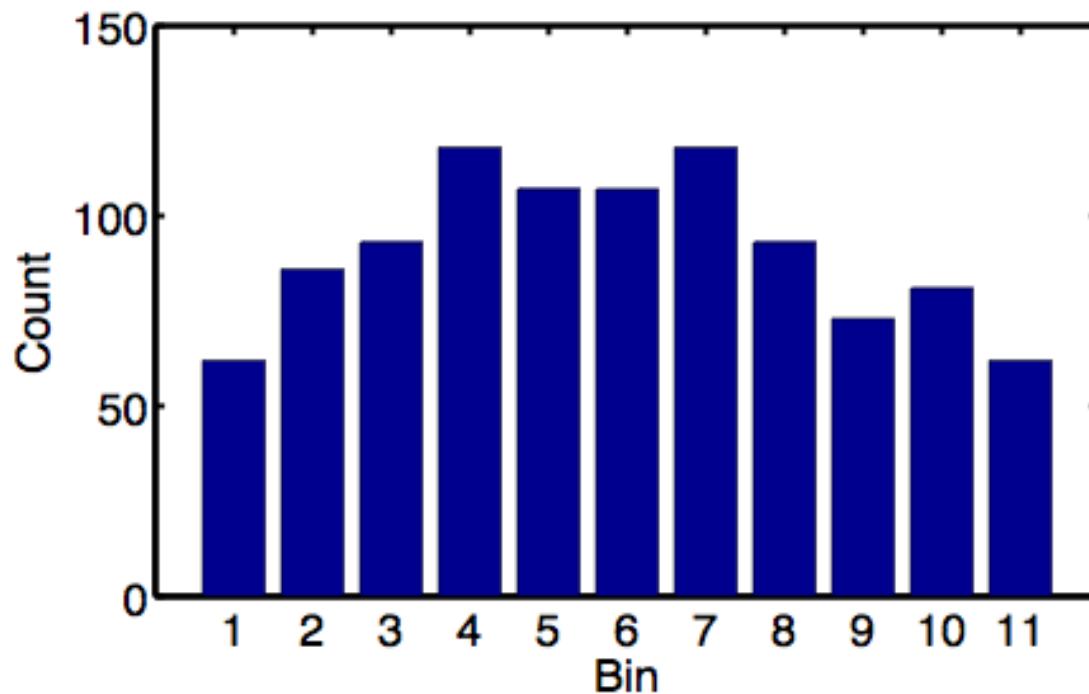
An ensemble with too little spread gives a u-shape.
This is the most common behavior for geophysics.



The Rank Histogram: Evaluating Ensemble Performance

An ensemble with too much spread is peaked in the center.

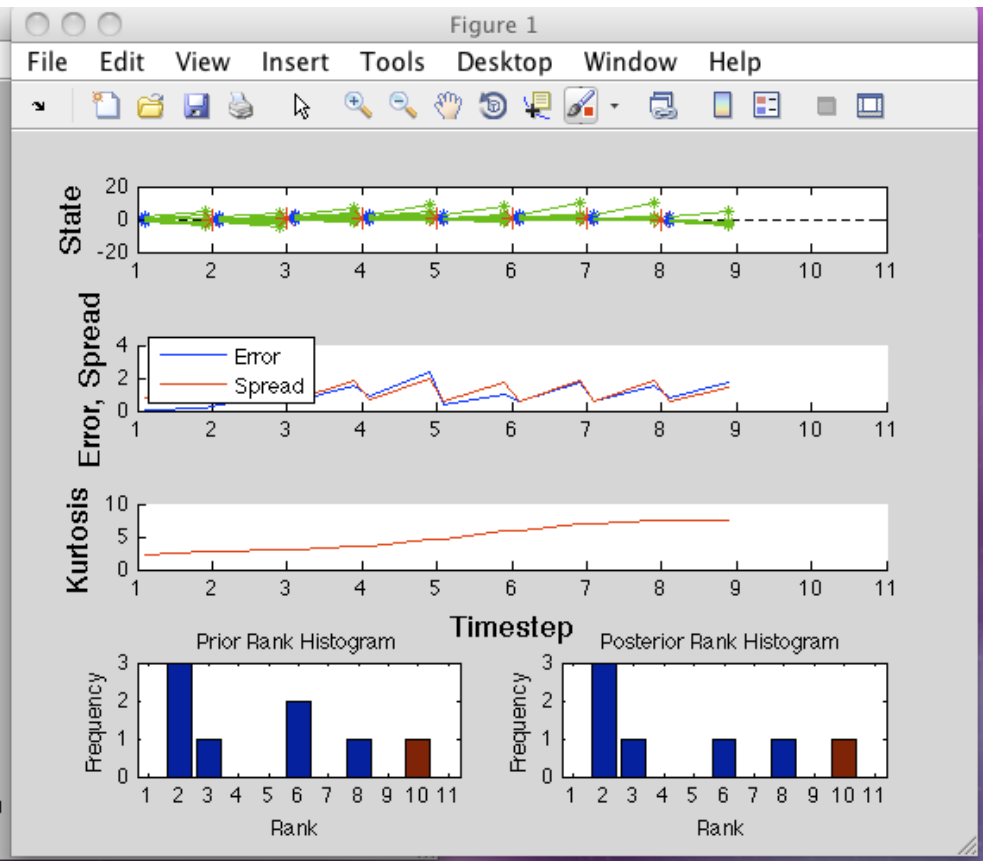
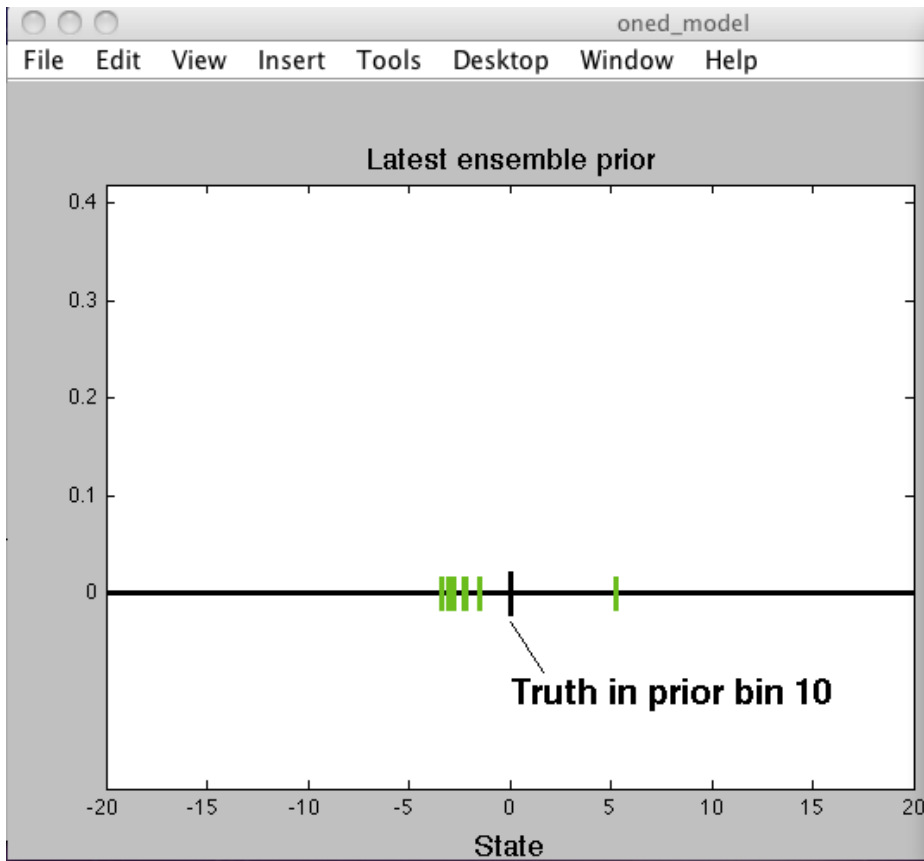
(See example set 1 in appendix)

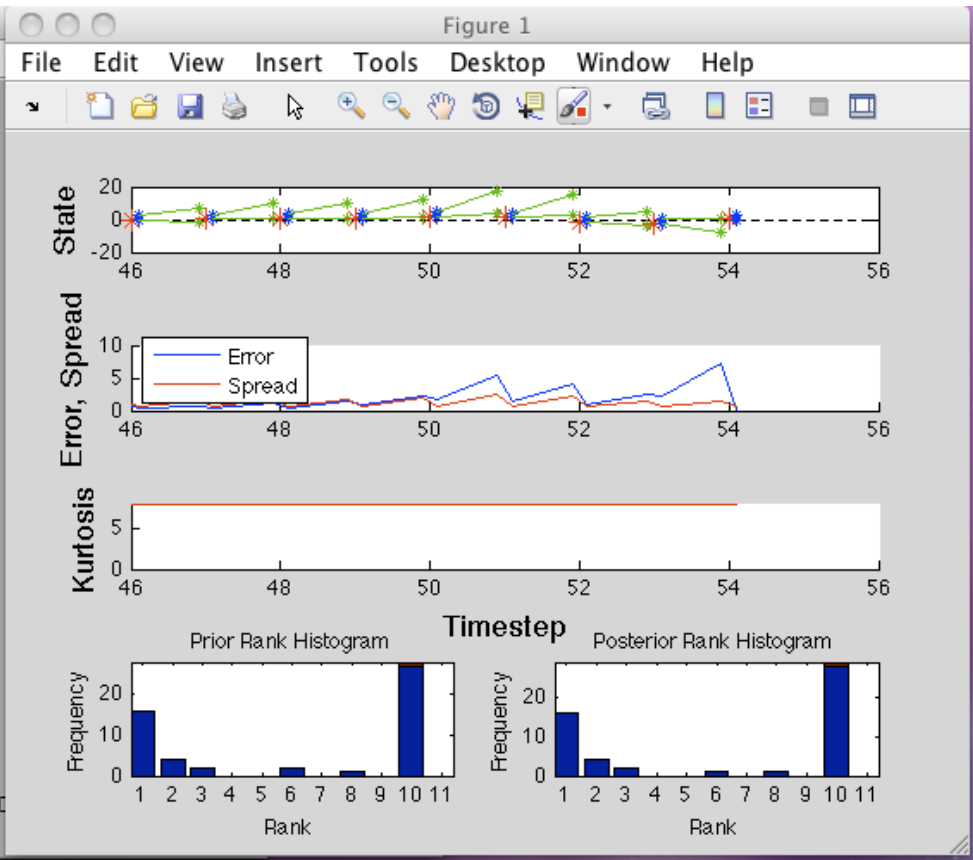
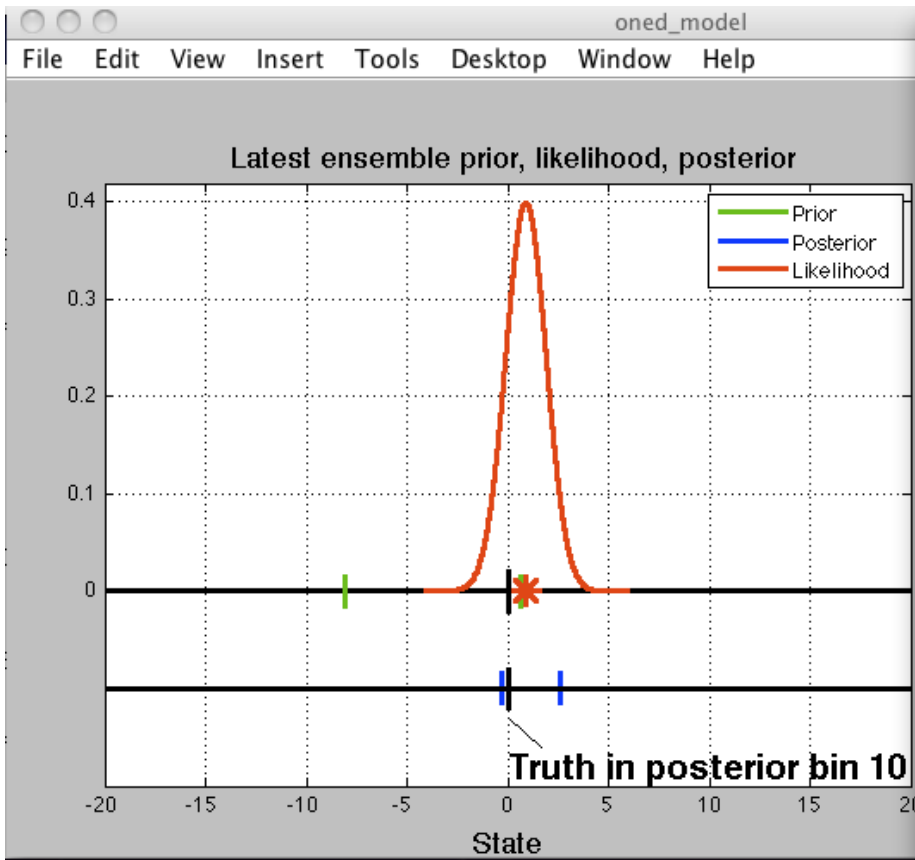


Back to 1-Dimensional (Ensemble) Kalman Filter

- All bets are off when model is nonlinear, likelihood nongaussian.
- Must assess quality of estimates for each case.
- Example: A weakly-nonlinear 1D model:
$$dx/dt = x + 0.4|x|x$$

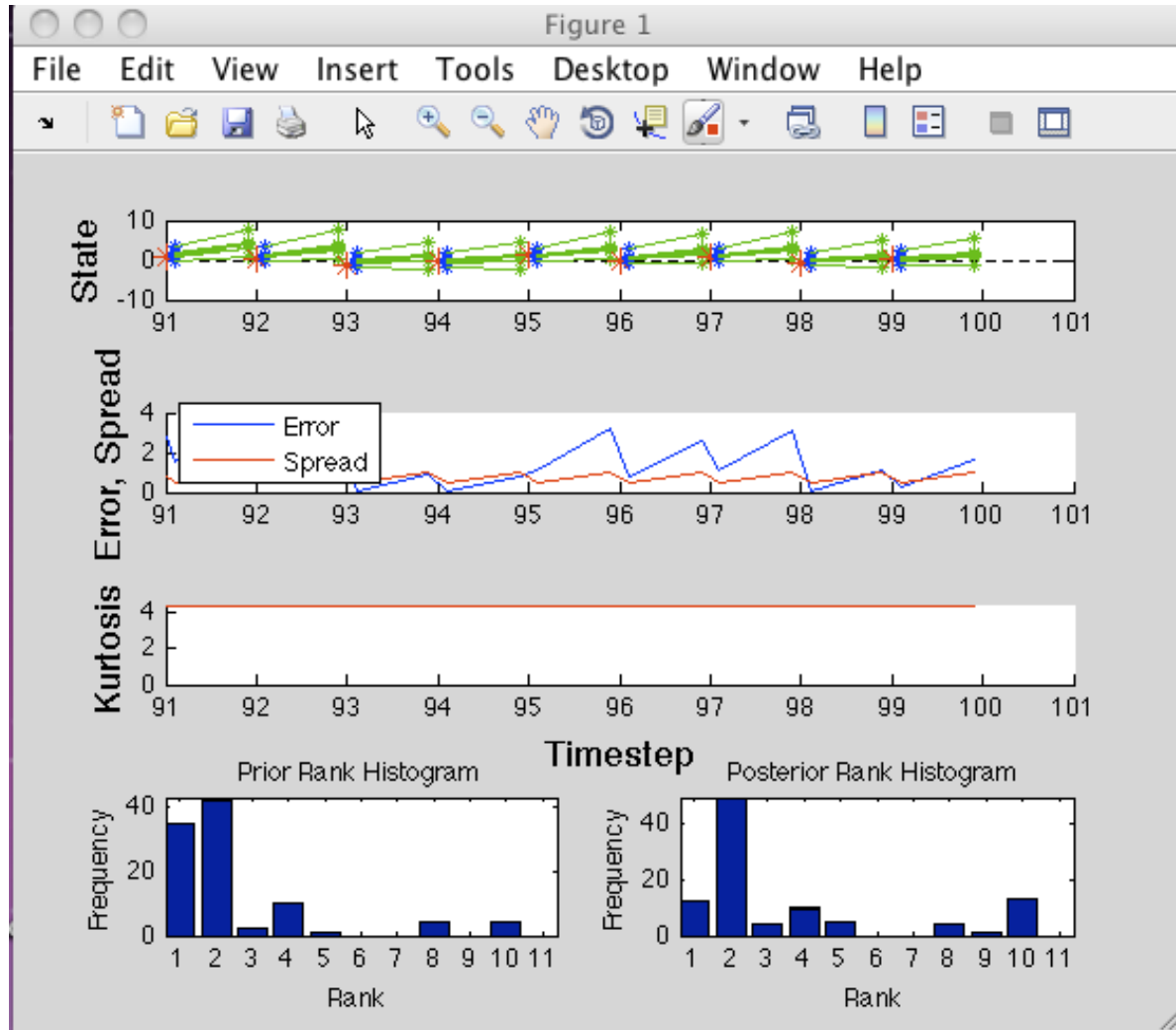
(See example set 2 in appendix)
- Ensemble statistics can become degenerate.
- Nonlinear ensemble filters may address this but beyond scope for today.

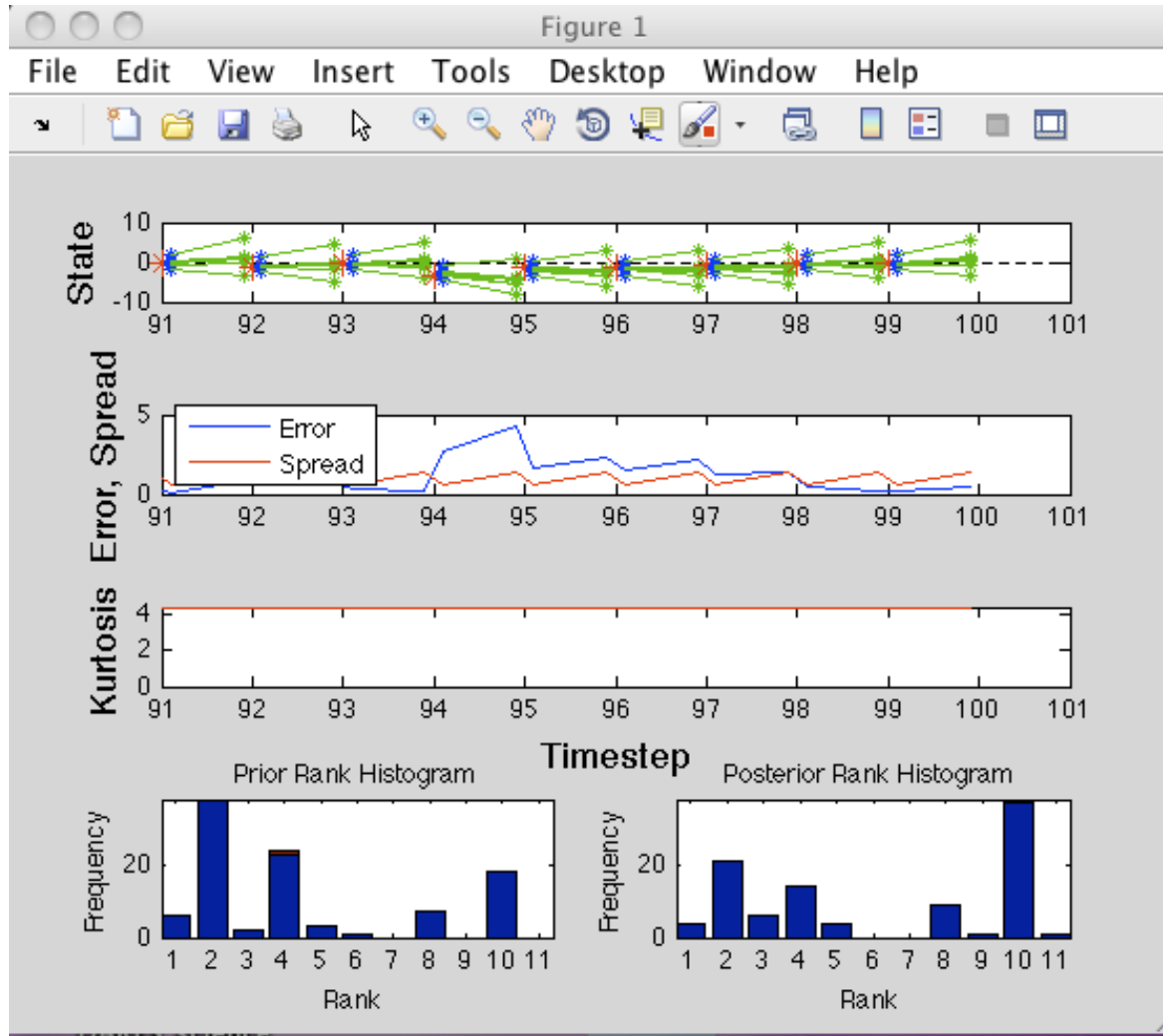




Ensemble Kalman Filter with Model Error

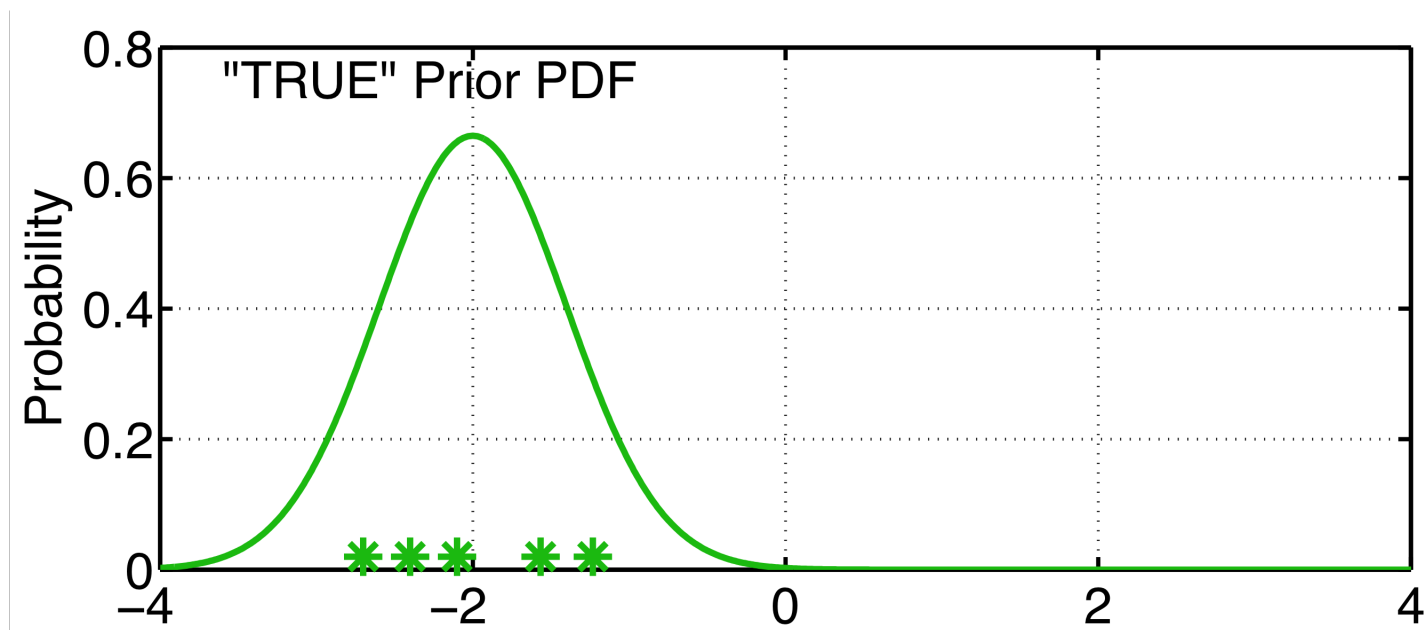
- Welcome to the real world. Models aren't perfect.
- Model (and other) error can corrupt ensemble mean and variance.
(See example set 3 part 1 in appendix)
- Adapt to this by increasing uncertainty in prior.
- Inflation is common method for ensembles.





Dealing with systematic error: Variance Inflation

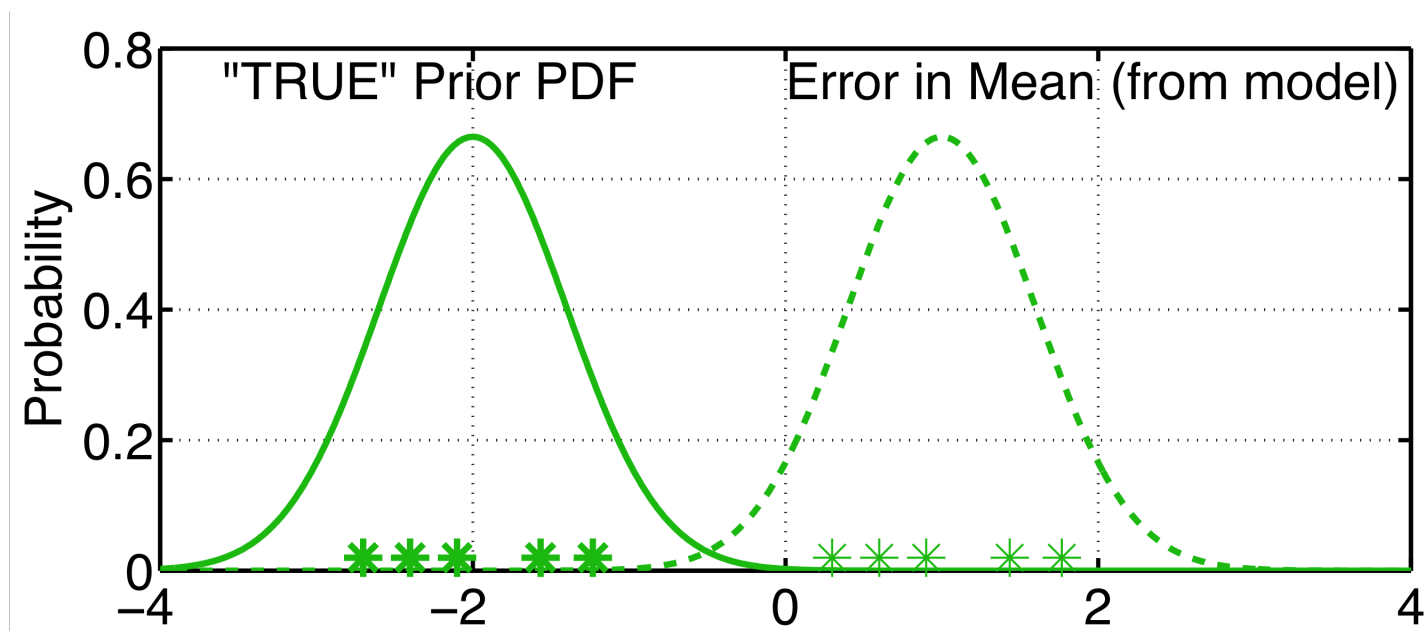
Observations + physical system => 'true' distribution.



Dealing with systematic error: Variance Inflation

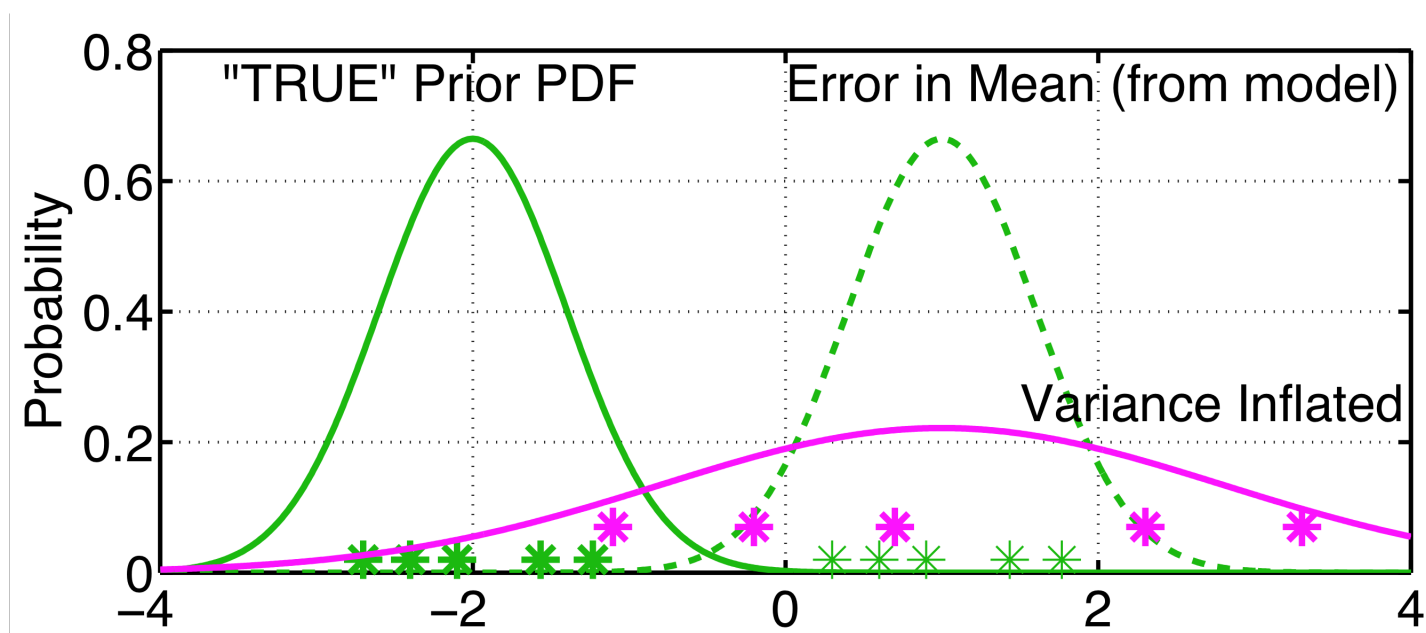
Observations + physical system => 'true' distribution.
Model bias (and other errors) can shift actual prior.

Prior ensemble is too certain (needs more spread).



Dealing with systematic error: Variance Inflation

Naïve solution: increase the spread in the prior.
Give more weight to the observation, less to prior.



Ensemble Kalman Filter with Model Error

- Errors in model (and other things) result in **too much certainty** in prior.
- Can use inflation to correct for this.
- Generally tuned on dependent data set.
- Adaptive methods exist but must be calibrated.
(See example set 3 part 2 in appendix)

- **Uncertainty errors depend on forecast length.**

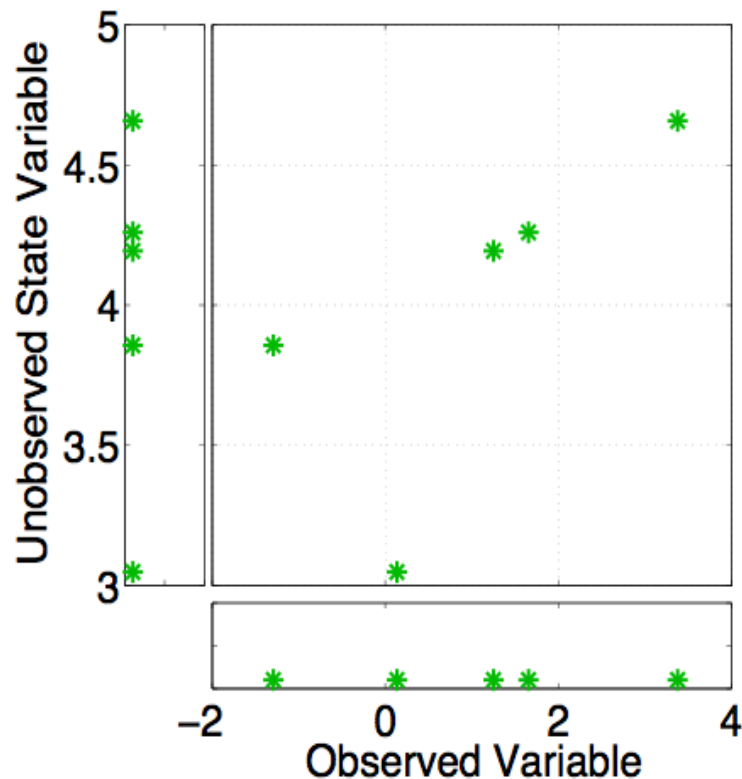
Multivariate Ensemble Kalman Filter

So far, we have an observation likelihood for single variable.

Suppose the model prior has additional variables.

Use linear regression to update additional variables.

Ensemble filters: Updating additional prior state variables

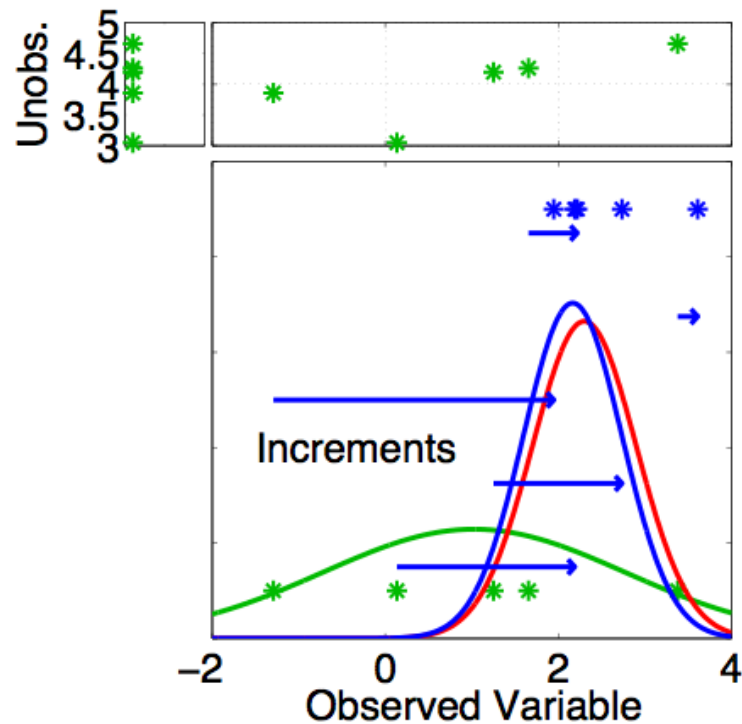


Assume that all we know is prior joint distribution.

One variable is observed.

What should happen to the unobserved variable?

Ensemble filters: Updating additional prior state variables

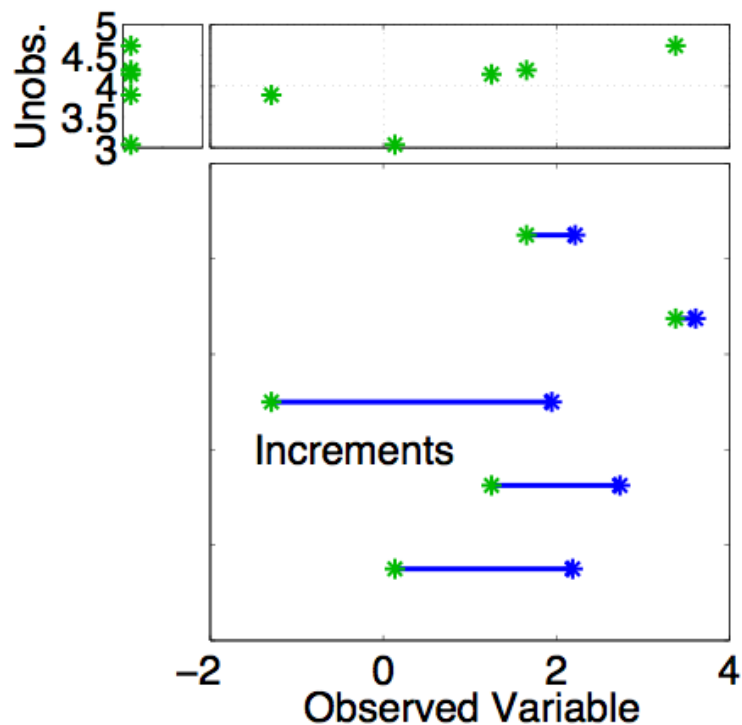


Assume that all we know is prior joint distribution.

One variable is observed.

Compute increments for prior ensemble members of observed variable.

Ensemble filters: Updating additional prior state variables

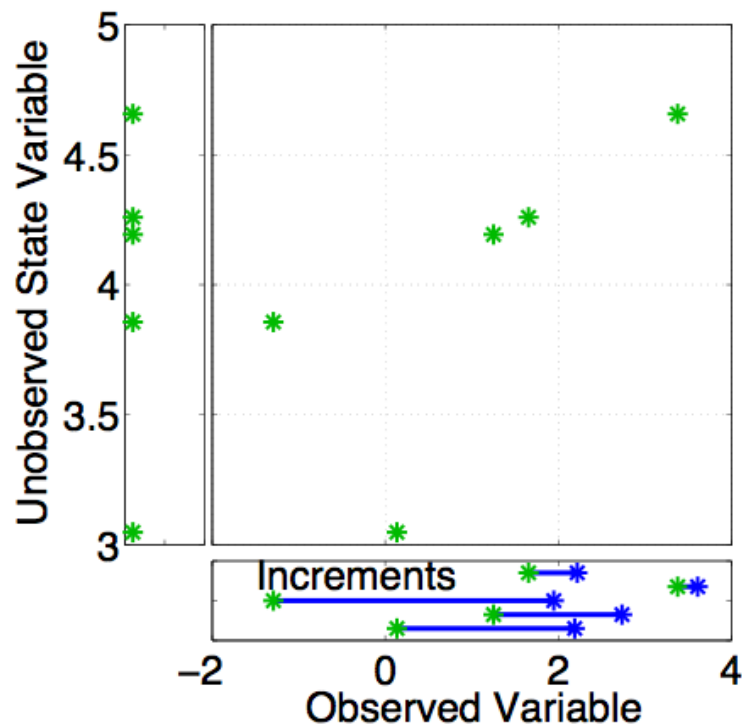


Assume that all we know is prior joint distribution.

One variable is observed.

Using only increments guarantees that if observation had no impact on observed variable, unobserved variable is unchanged (highly desirable).

Ensemble filters: Updating additional prior state variables



Assume that all we know is prior joint distribution.

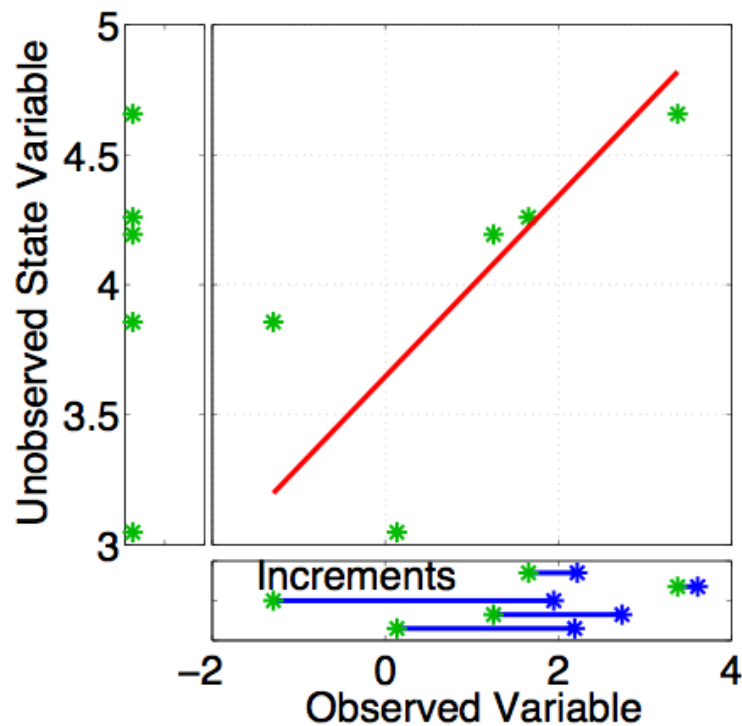
How should the unobserved variable be impacted?

First choice: least squares.

Equivalent to linear regression.

Same as assuming binormal prior.

Ensemble filters: Updating additional prior state variables



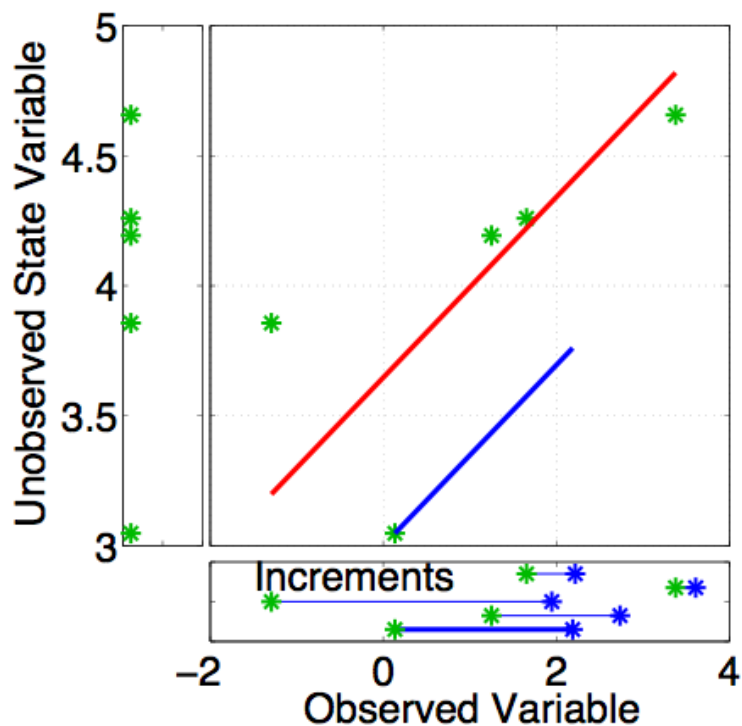
Have joint prior distribution of two variables.

How should the unobserved variable be impacted?

First choice: least squares.

Begin by finding least squares fit.

Ensemble filters: Updating additional prior state variables

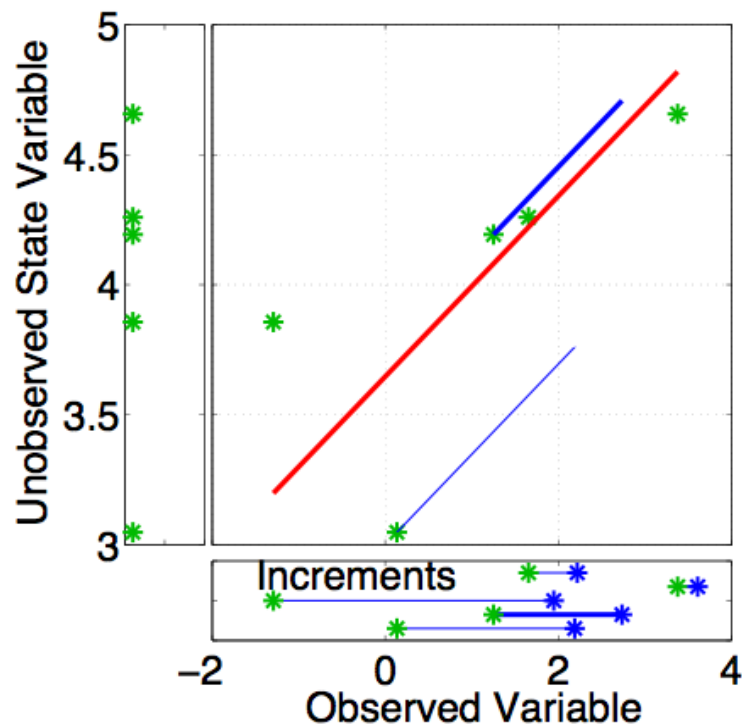


Have joint prior distribution of two variables.

Next, regress the observed variable increments onto increments for the unobserved variable.

Equivalent to first finding image of increment in joint space.

Ensemble filters: Updating additional prior state variables

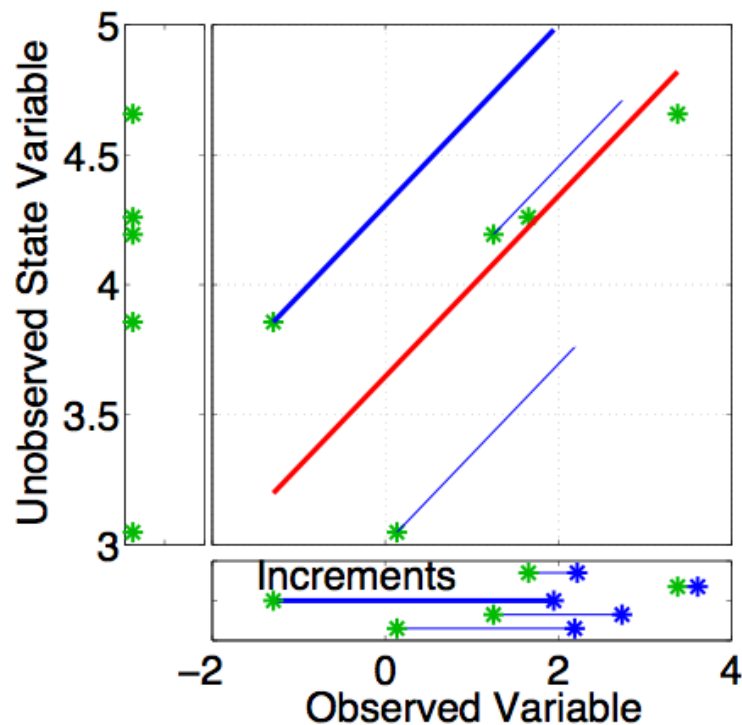


Have joint prior distribution of two variables.

Next, regress the observed variable increments onto increments for the unobserved variable.

Equivalent to first finding image of increment in joint space.

Ensemble filters: Updating additional prior state variables

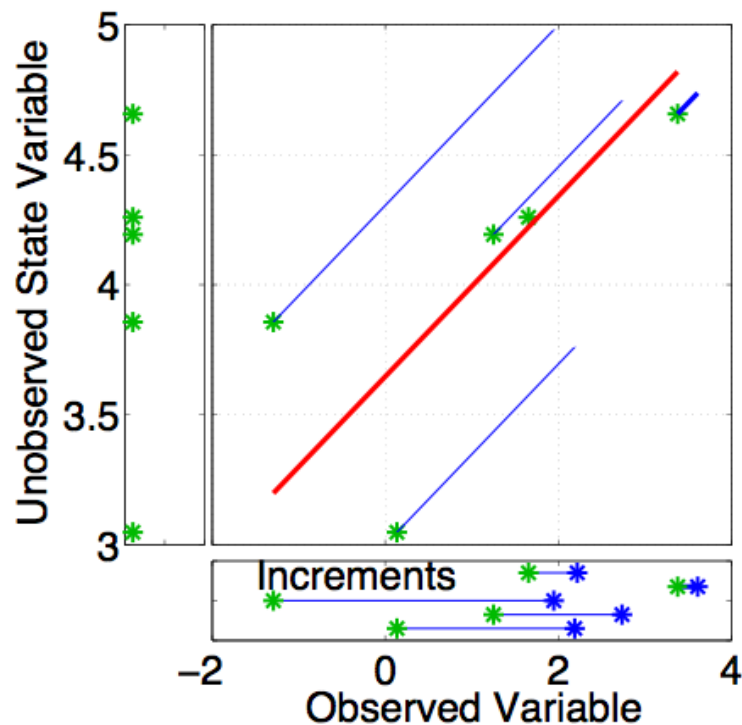


Have joint prior distribution of two variables.

Next, regress the observed variable increments onto increments for the unobserved variable.

Equivalent to first finding image of increment in joint space.

Ensemble filters: Updating additional prior state variables

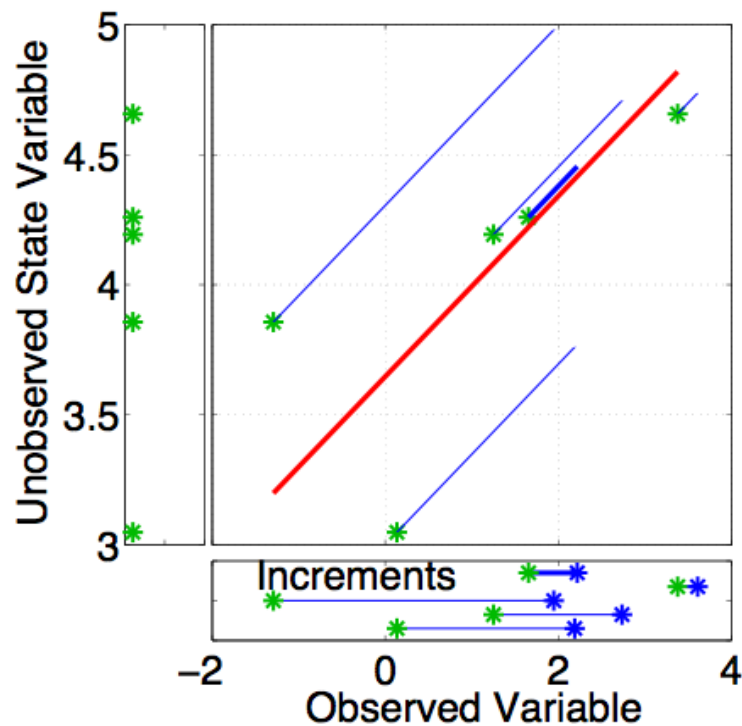


Have joint prior distribution of two variables.

Next, regress the observed variable increments onto increments for the unobserved variable.

Equivalent to first finding image of increment in joint space.

Ensemble filters: Updating additional prior state variables

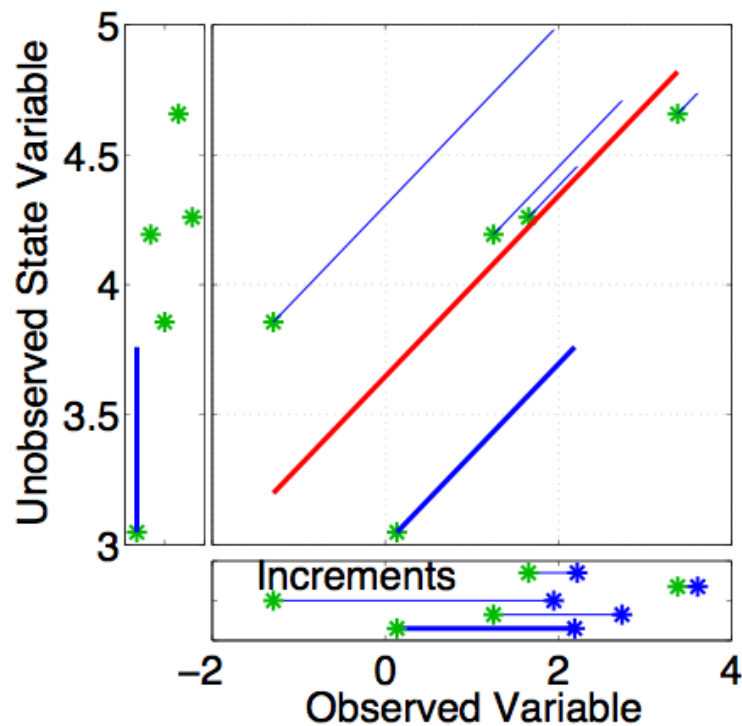


Have joint prior distribution of two variables.

Next, regress the observed variable increments onto increments for the unobserved variable.

Equivalent to first finding image of increment in joint space.

Ensemble filters: Updating additional prior state variables

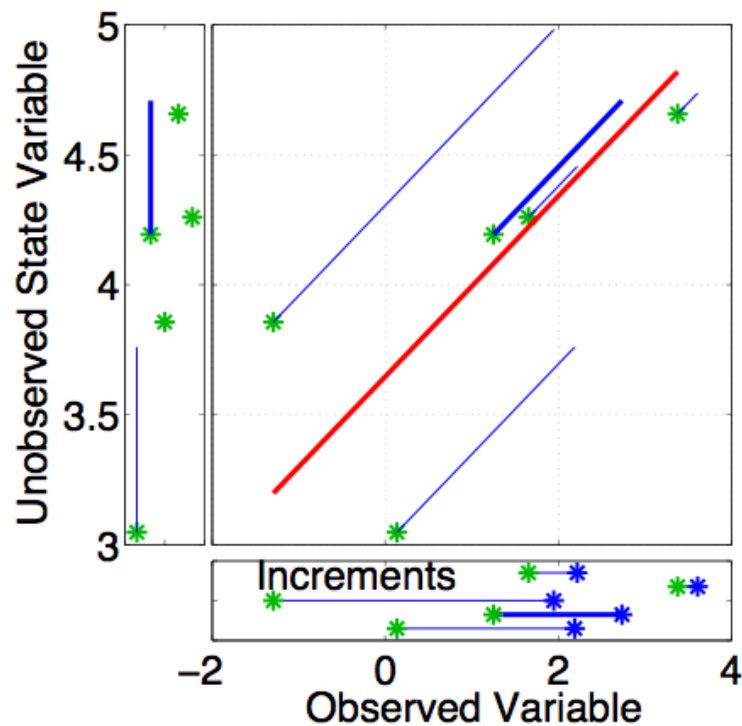


Have joint prior distribution of two variables.

Regression: Equivalent to first finding image of increment in joint space.

Then projecting from joint space onto unobserved priors.

Ensemble filters: Updating additional prior state variables

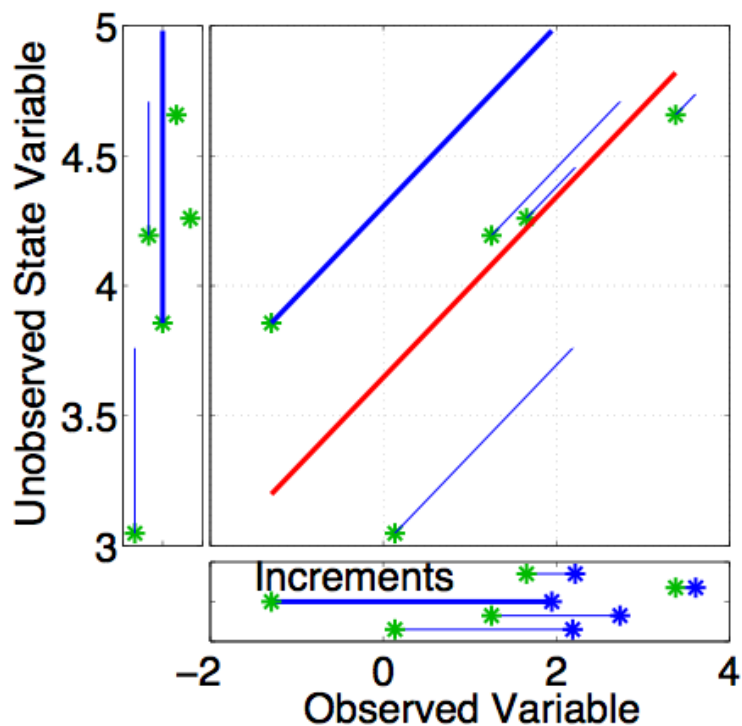


Have joint prior distribution of two variables.

Regression: Equivalent to first finding image of increment in joint space.

Then projecting from joint space onto unobserved priors.

Ensemble filters: Updating additional prior state variables

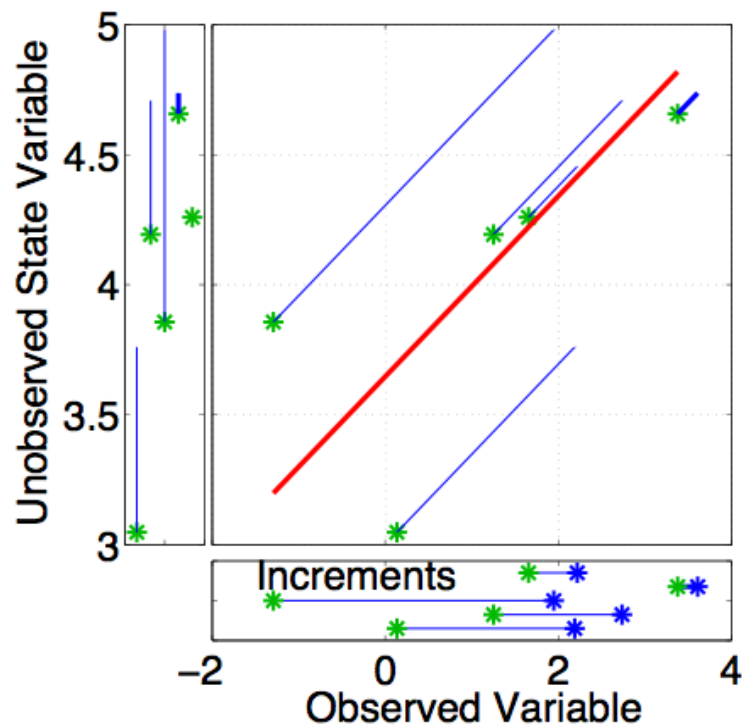


Have joint prior distribution of two variables.

Regression: Equivalent to first finding image of increment in joint space.

Then projecting from joint space onto unobserved priors.

Ensemble filters: Updating additional prior state variables

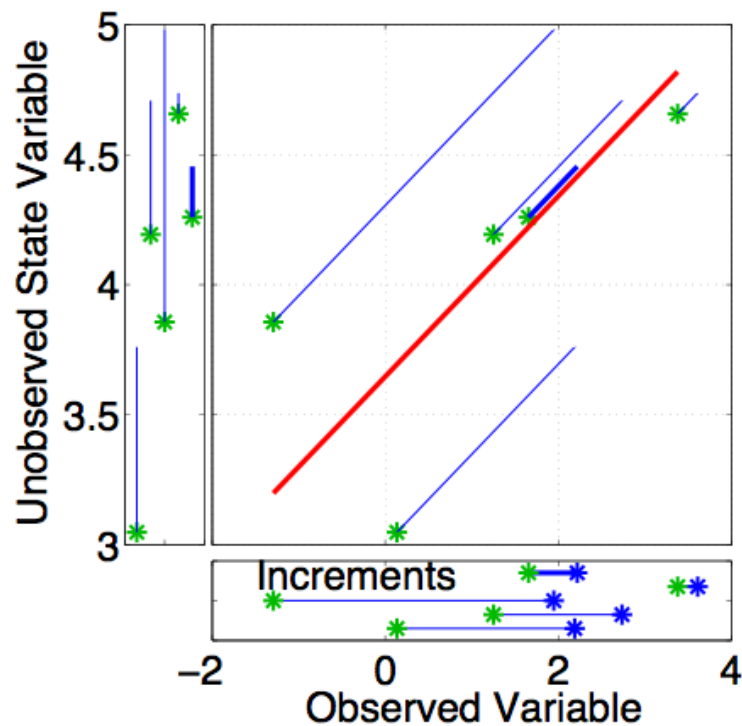


Have joint prior distribution of two variables.

Regression: Equivalent to first finding image of increment in joint space.

Then projecting from joint space onto unobserved priors.

Ensemble filters: Updating additional prior state variables



Have joint prior distribution of two variables.

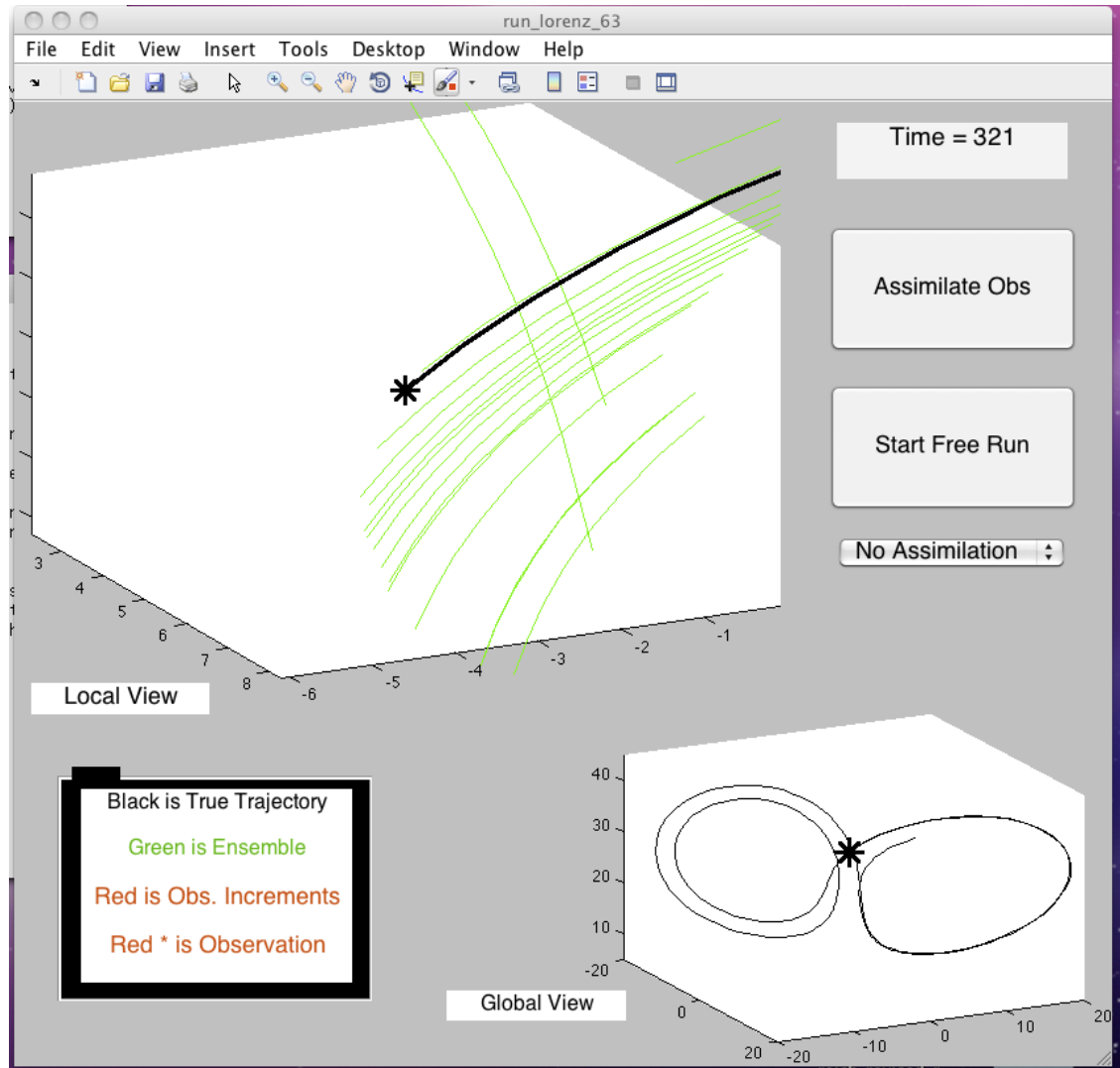
Regression: Equivalent to first finding image of increment in joint space.

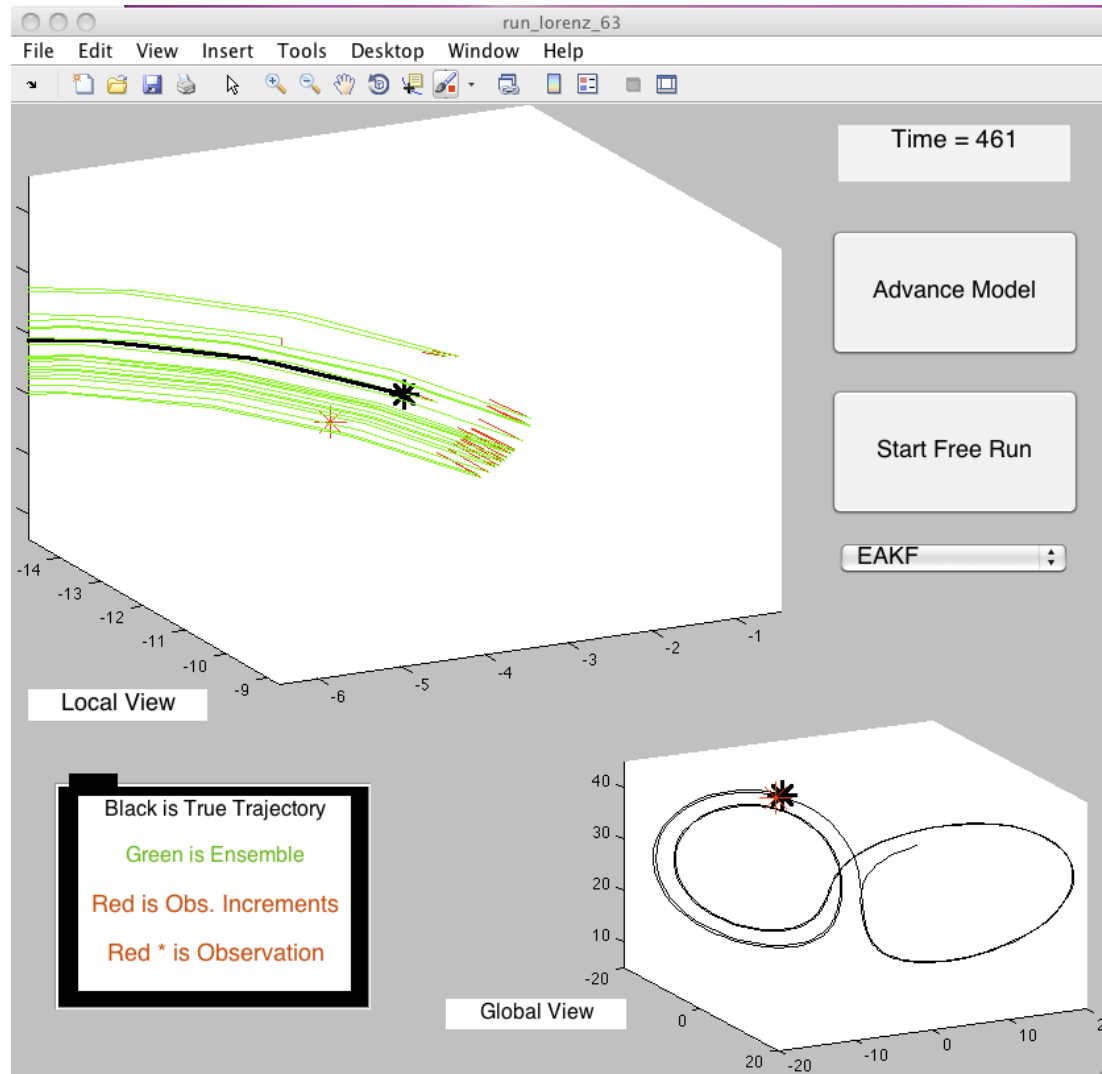
Then projecting from joint space onto unobserved priors.

The Lorenz63 3-Variable 'Chaotic' Model

- Two lobed attractor.
- State orbits lobes, switches lobes sporadically.
- Nonlinear if observations are sparse.
- Qualitative uncertainty information available from ensemble.
- Nongaussian distributions are sampled.
- **Quantitative uncertainty requires careful calibration.**

(See Example 4 in the appendix)





Ensemble Kalman Filters with Large Models

- Kalman filter is too costly (time and storage) for large models.
- Ensemble filters can be used for very large applications.
- Sampling error when computing covariances may become large.
- Spurious correlations between unrelated observations and state variables contaminate results.
- Leads to further uncertainty underestimates.
- Can be compensated by a priori limits on what state variables are impacted by an observation. Called 'localization'.
- **Requires even more calibration.**

The Lorenz96 40-Variable ‘Chaotic’ Model

- Something like weather systems around a latitude circle.
- Has spurious correlations for small ensembles.
- Basic implementation of ensemble filter leads to bad uncertainty estimates.
- Both inflation and localization can improve performance.
- **Quantitative uncertainty requires careful calibration.**

(See Example 5 in the appendix)

Uncertainty and Ensemble Kalman Filters: Conclusions

- Ensemble KF variance is exact quantification of uncertainty when things are linear, gaussian, perfect, and ensemble is large enough.
- Too little uncertainty nearly universal for geophysical applications.
- Uncertainty errors are a function of forecast lead time.
- Adaptive algorithms like inflation improve uncertainty estimates.
- Quantitative uncertainty estimates require careful calibration.
- Out of sample estimates (climate change) are tenuous.

Appendix: Matlab examples

This is an outline of the live matlab demos from NCAR's DART system that were used to support this tutorial. The matlab code is available by checking out the DART system from NCAR at:

<http://www2.image.ucar.edu/forms/dart-software-download>

The example scripts can be found in the DART directory:

DART_LAB/matlab

Some of the features of the GUI's for the exercises do not work due to bugs in Matlab releases 2010a and 2010b but should work in earlier and later releases.

Many more features of these matlab scripts are described in the tutorial files in the DART directory:

DART_LAB/presentation

These presentations are a complement and extension of things presented in this tutorial.

Appendix: Matlab examples

Example Set 1:

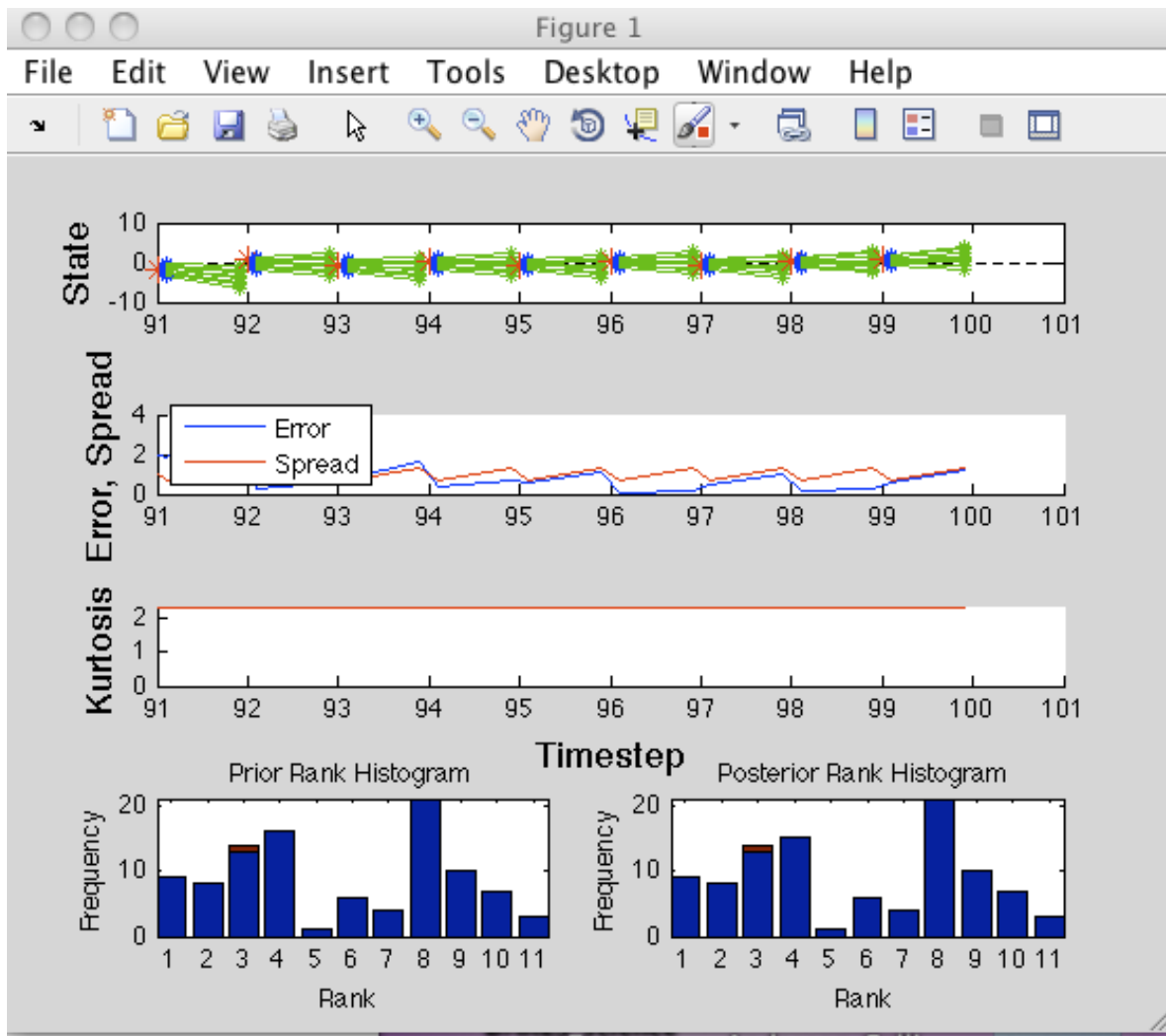
Use matlab script oned_model

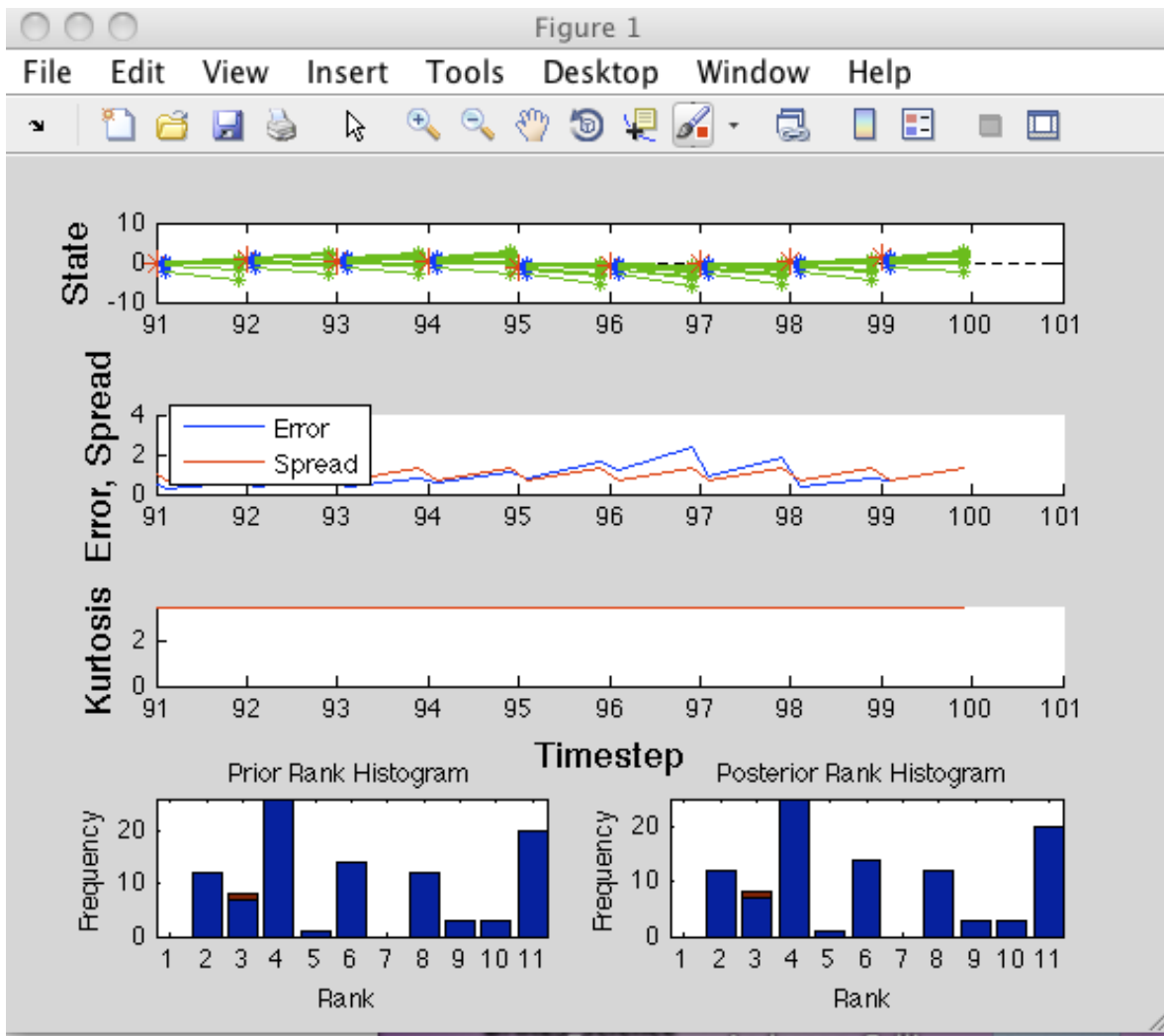
Change the “Ens. Size” to 10

Select “Start Free Run”

Observe that the spread curve in the second diagnostic panel quickly converges to have the same prior and posterior values. The rms error varies with time due to observation error, but its time mean should be the same as the spread.

Restart this exercise several times to see that the kurtosis curve in the third panel also converges immediately, but to values that change with every randomly selected initial ensemble. Also note the behavior of the rank histograms for this case. They may not be uniform even though the ensemble filter is the optimal solution. In fact, they can be made arbitrarily nonuniform by the selection of the initial conditions.





Appendix: Matlab examples

Example Set 2:

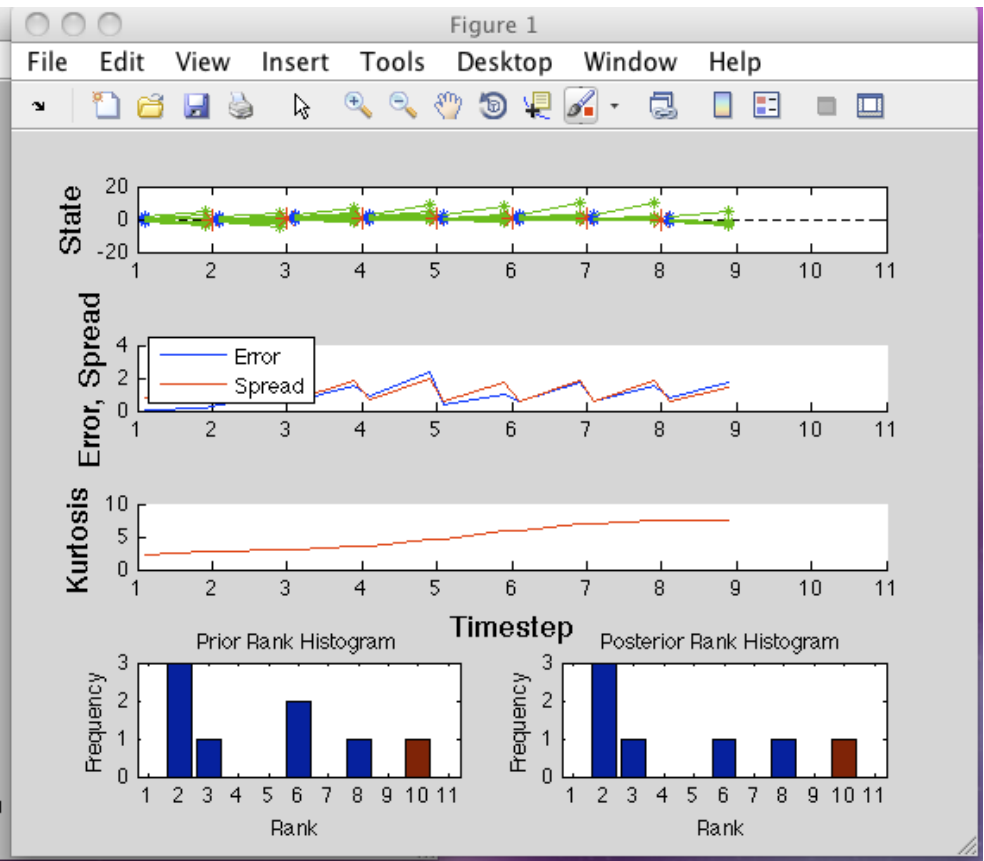
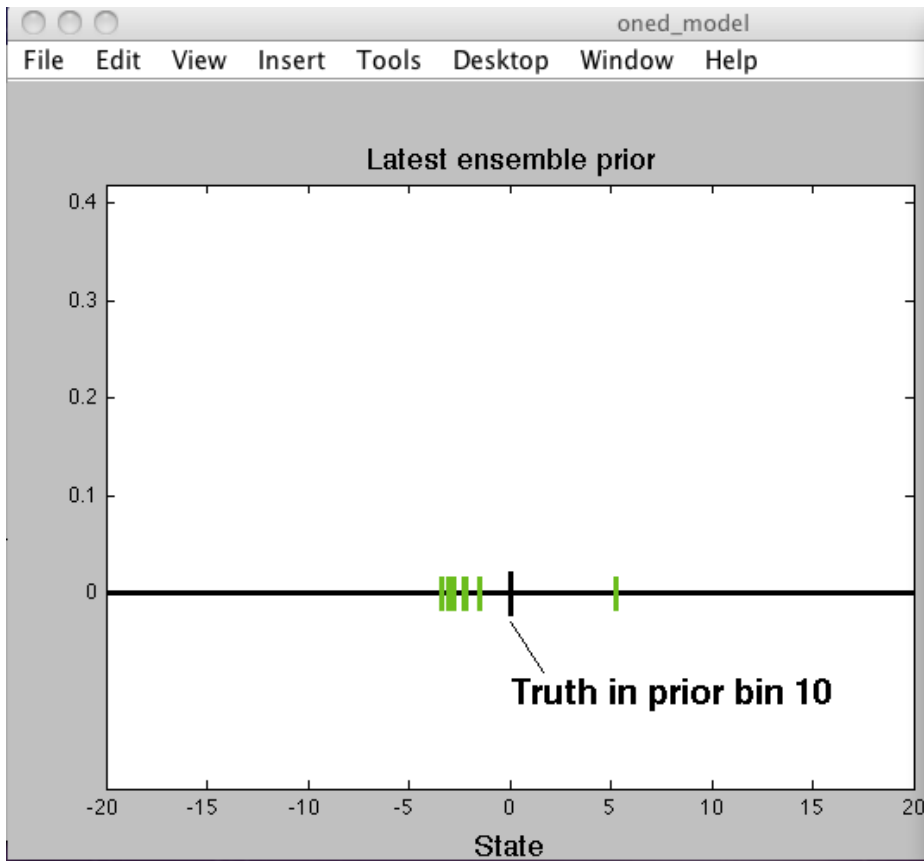
Use matlab script oned_model

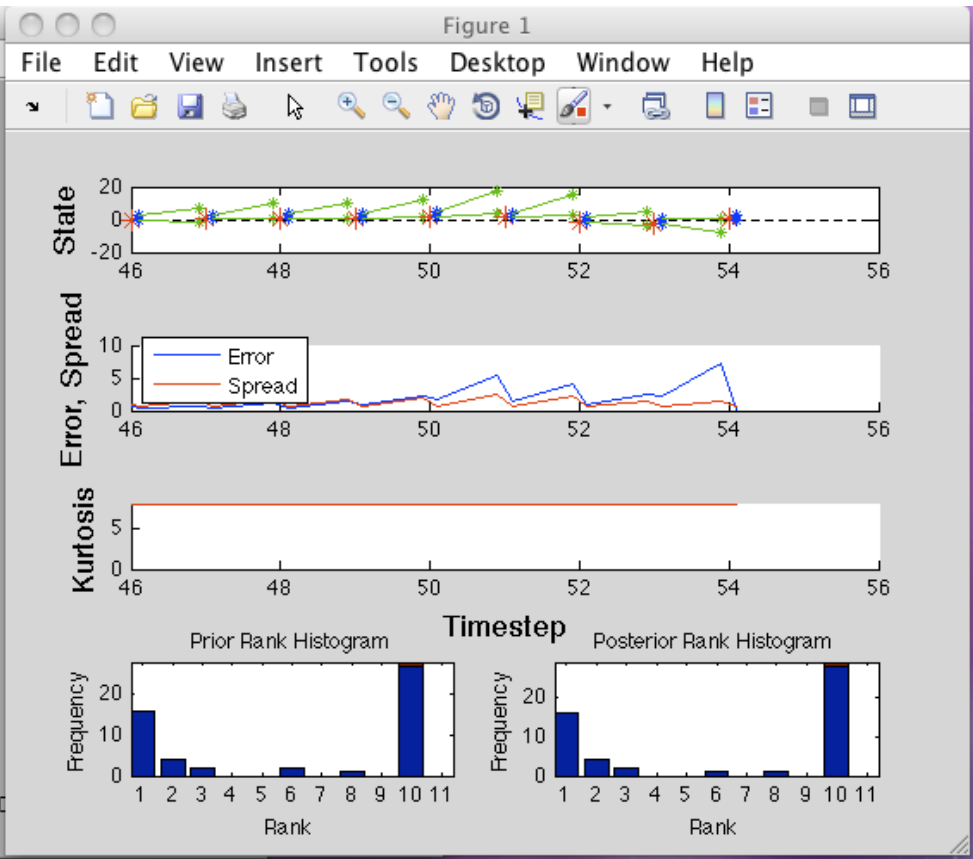
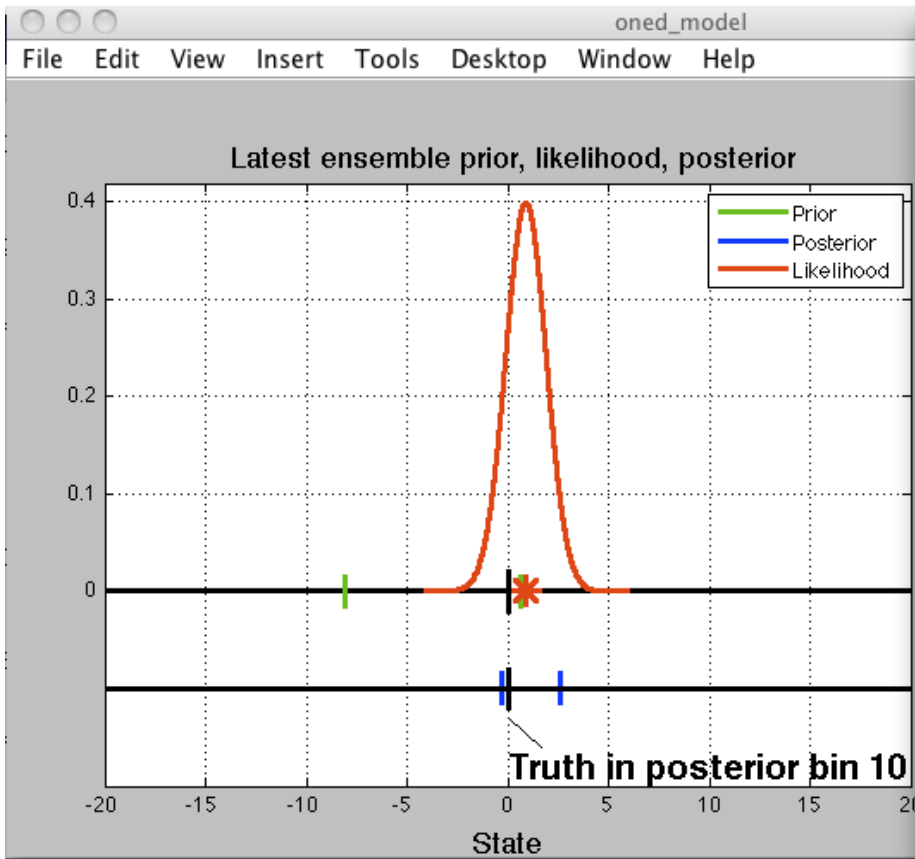
Change the “Ens. Size” to 10

Change “Nonlin. A” to 0.4

Select “Start Free Run”

This examines what happens when the model is nonlinear. Observe the evolution of the rms error and spread in the second panel and the kurtosis in the third panel. Also observe the ensemble prior distribution (green) in the first panel and in the separate plot on the control panel (where you changed the ensemble size). In many cases, the ensemble will become degenerate with all but one of the ensemble members collapsing while the other stays separate. Observe the impact on the rank histograms as this evolves.





Appendix: Matlab examples

Example Set 3:

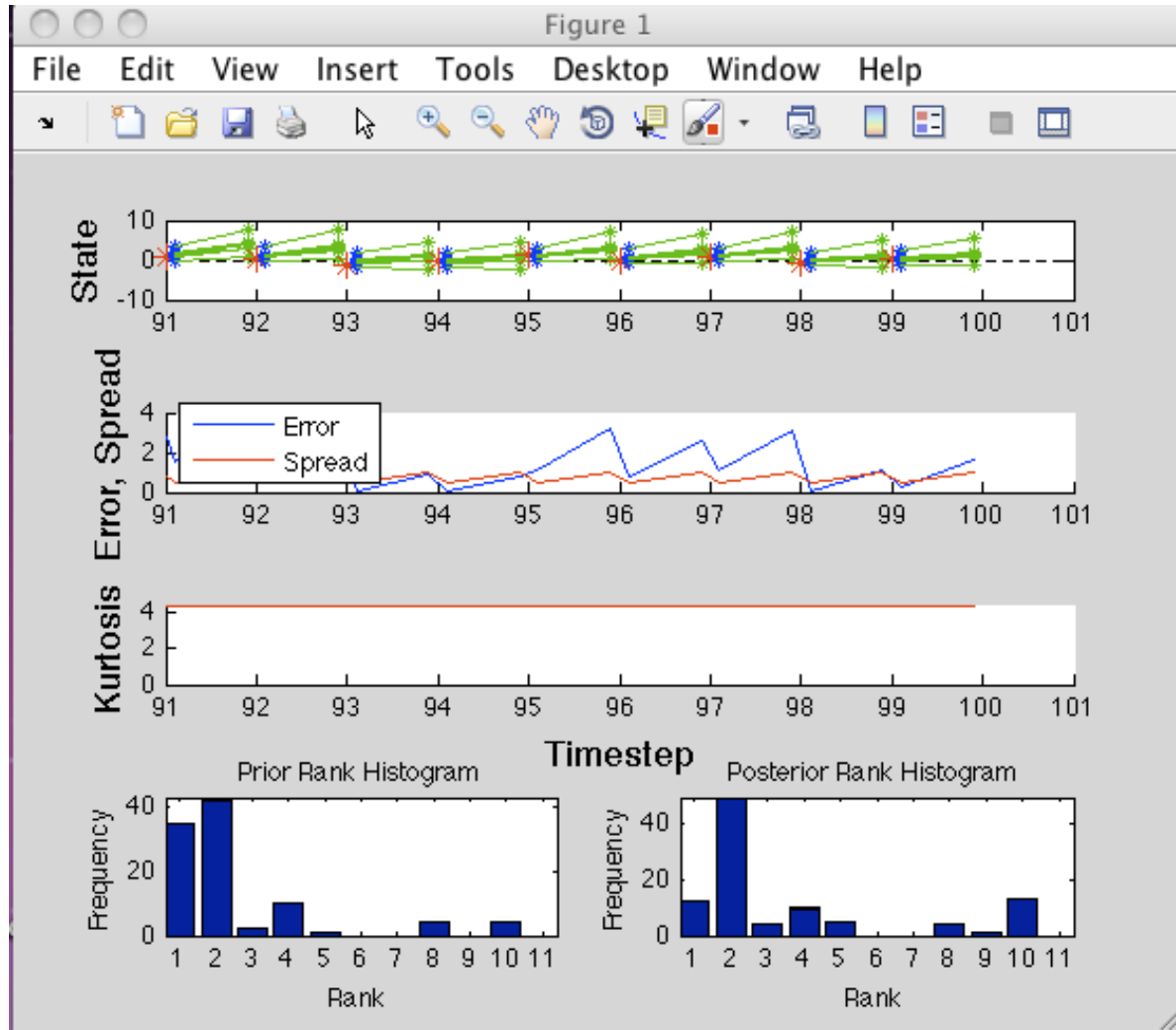
Use matlab script oned_model

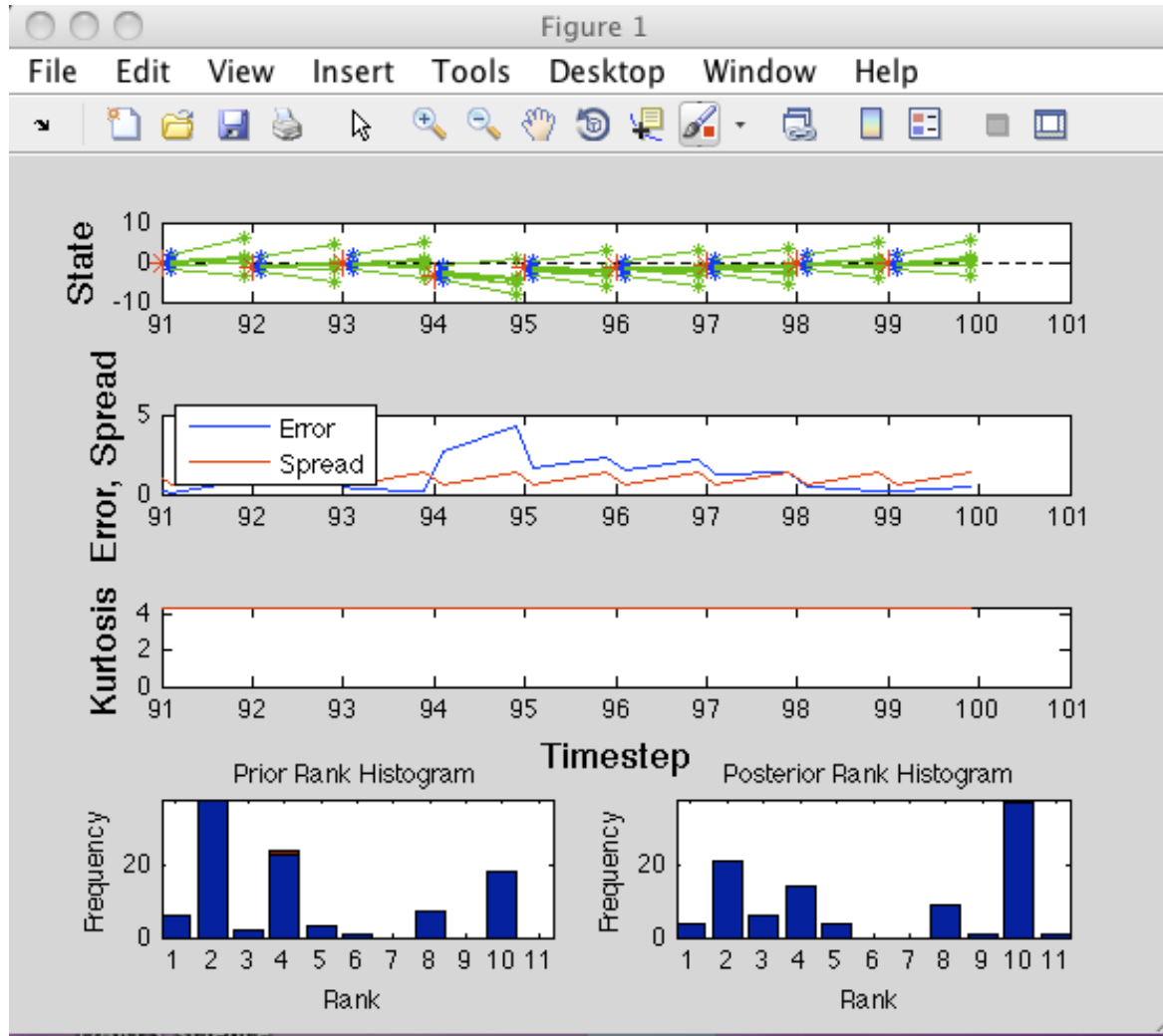
Change the “Ens. Size” to 10

Change “Model Bias” to 1

Select “Start Free Run”

This examines what happens when the model used to assimilate observations is biased compared to the model that generated the observations (an imperfect model study). The model for the assimilation adds 1 to what the perfect model would do at each time step. Observe the evolution of the rms error and spread in the second panel and the rank histograms. Now, try adding some uncertainty to the prior ensemble using inflation by setting “Inflation” to 1.5. This should impact all aspects of the assimilation.





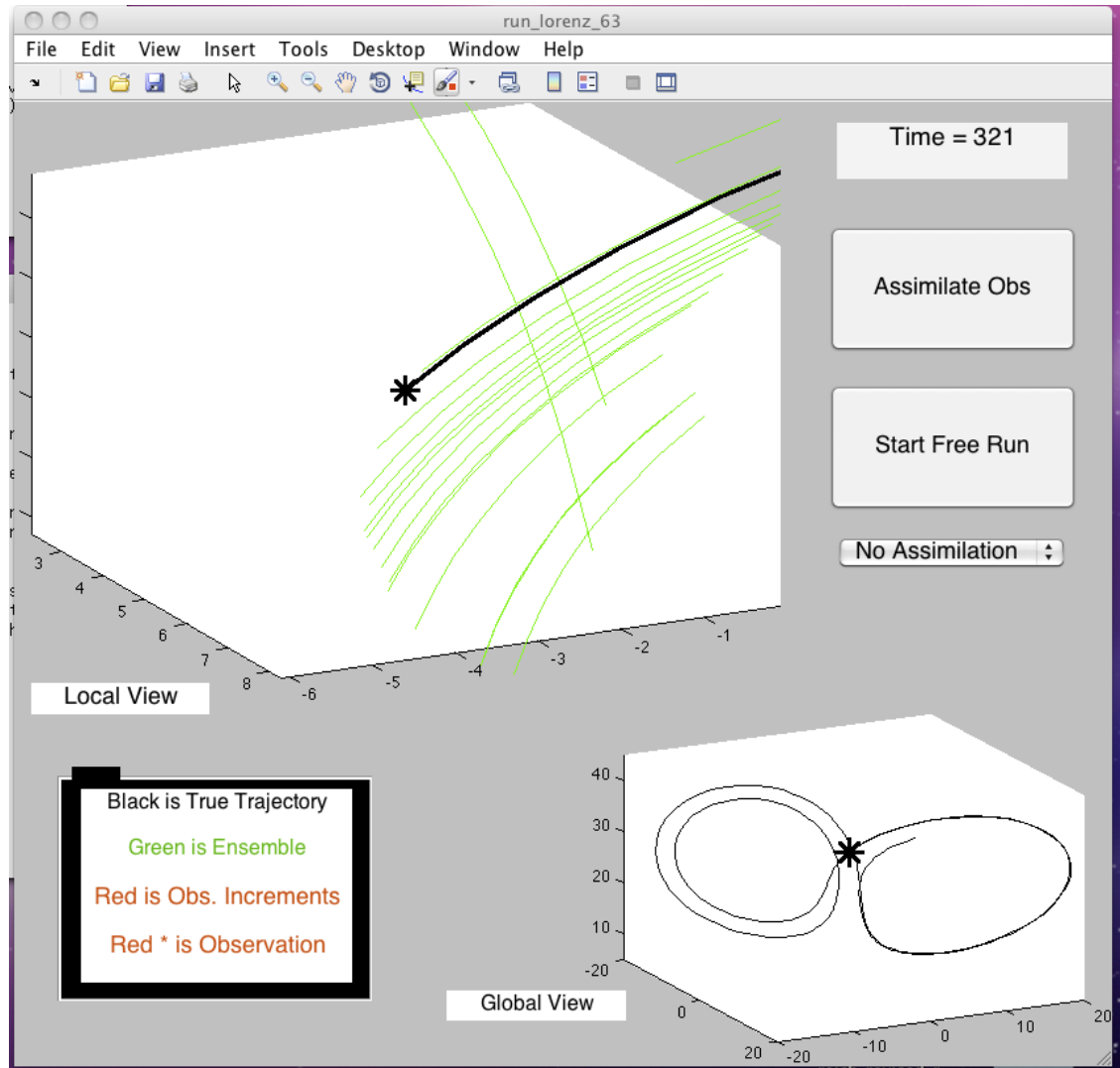
Appendix: Matlab examples

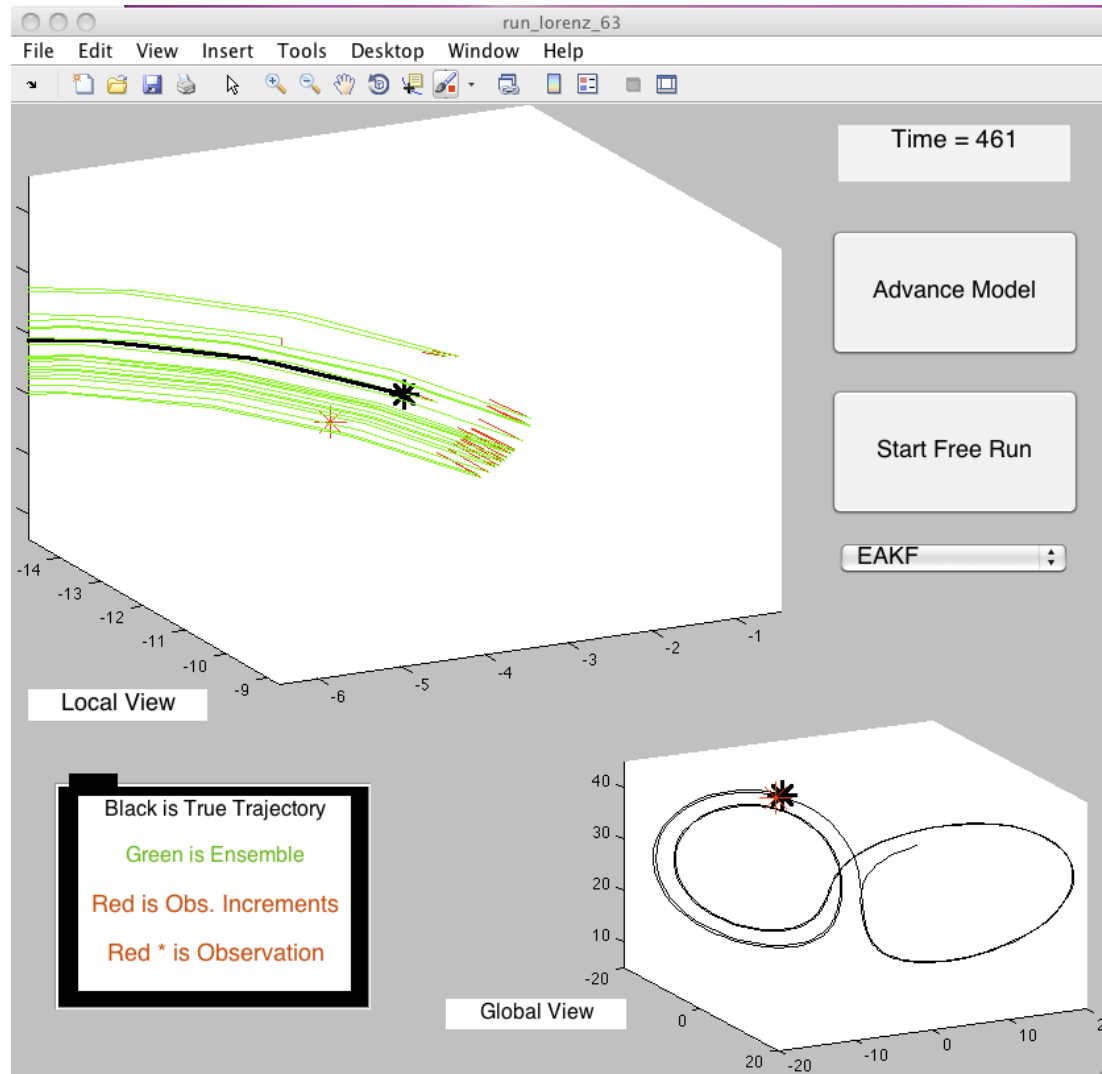
Example Set 4:

Use matlab script run_lorenz_63

Select “Start Free Run”

This run without assimilation shows how uncertainty grows in an ensemble forecast in the 3 variable Lorenz model. After a few trips around the attractor select “Stop Free Run”. Then turn on an ensemble filter assimilation by changing “No Assimilation “ to “EAKF”. Continue the run by selecting “Start Free Run” and observe how the ensemble evolves. You can stop the advance again and use the matlab rotation feature to study the structure of the ensemble.





Appendix: Matlab examples

Example Set 5:

Use matlab script run_lorenz_96

Select “Start Free Run”

This run without assimilation shows how uncertainty grows in an ensemble forecast in the 40 variable model. After about time=40, select “Stop Free Run” and start an assimilation by changing “No Assimilation” to “EAKF” and selecting “Start Free Run”. Exploring the use of a localization of 0.2 and/or an inflation of 1.1 helps to show how estimates of uncertainty can be improved by modifying the base ensemble Kalman filter.