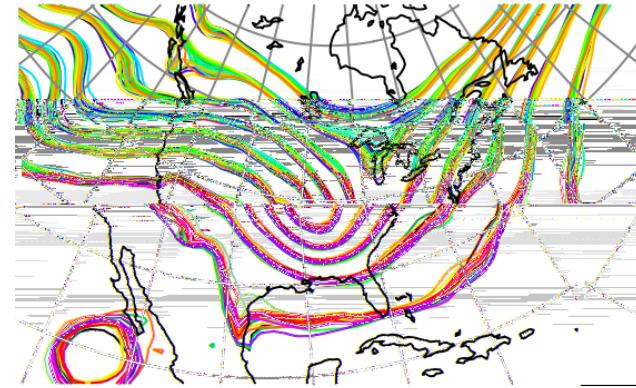


D  
A  
R  
T

ata  
ssimilation  
esearch  
estbed



## DART Tutorial Section 8: Dealing with Sampling Error



©UCAR 2014

The National Center for Atmospheric Research is sponsored by the National Science Foundation. Any opinions, findings and conclusions or recommendations expressed in this publication are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.

NCAR | National Center for  
UCAR Atmospheric Research



# Updating Additional Prior State Variables

Two primary error sources:

1. Sampling error due to noise.

Can occur even if there is a linear relation between variables.

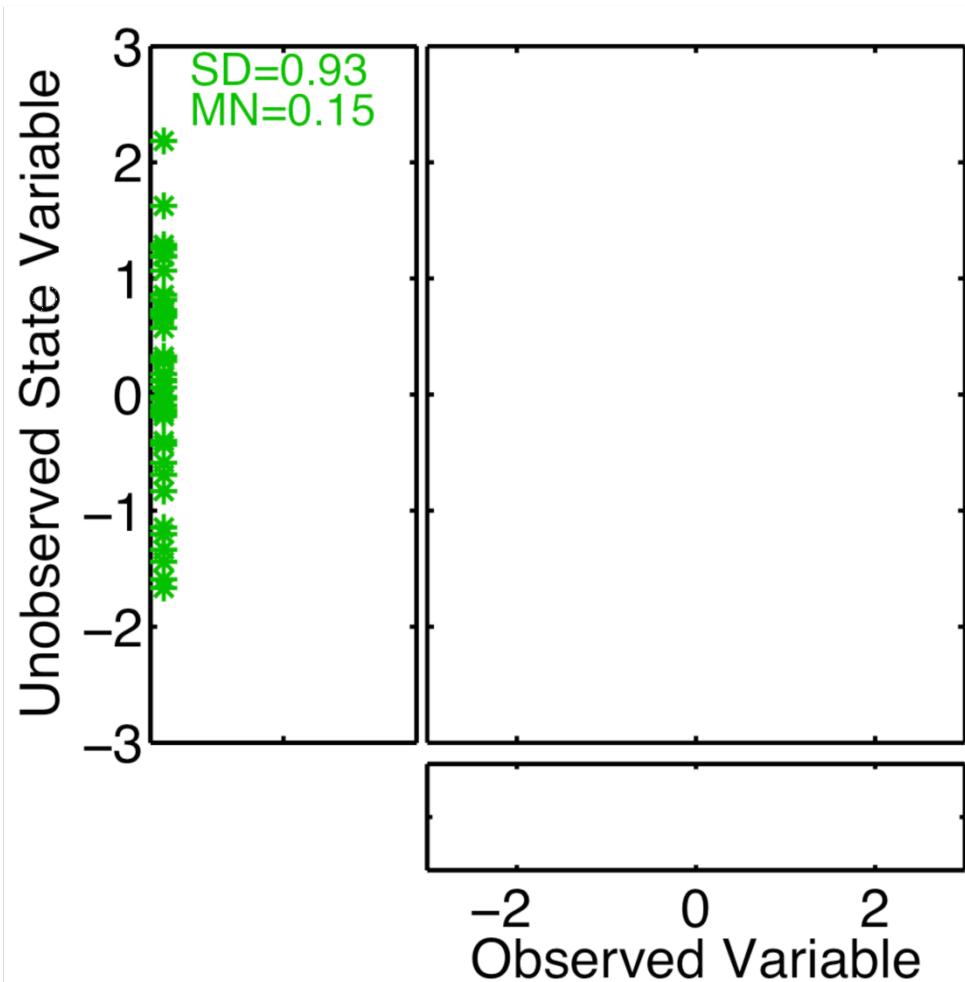
Sample regression coefficient imprecise with finite ensembles.

2. Linear approximation is invalid.

If there is substantial nonlinearity in ‘true’ relation between variables over range of prior ensemble. (see section 10).

May need to address both issues for good performance.

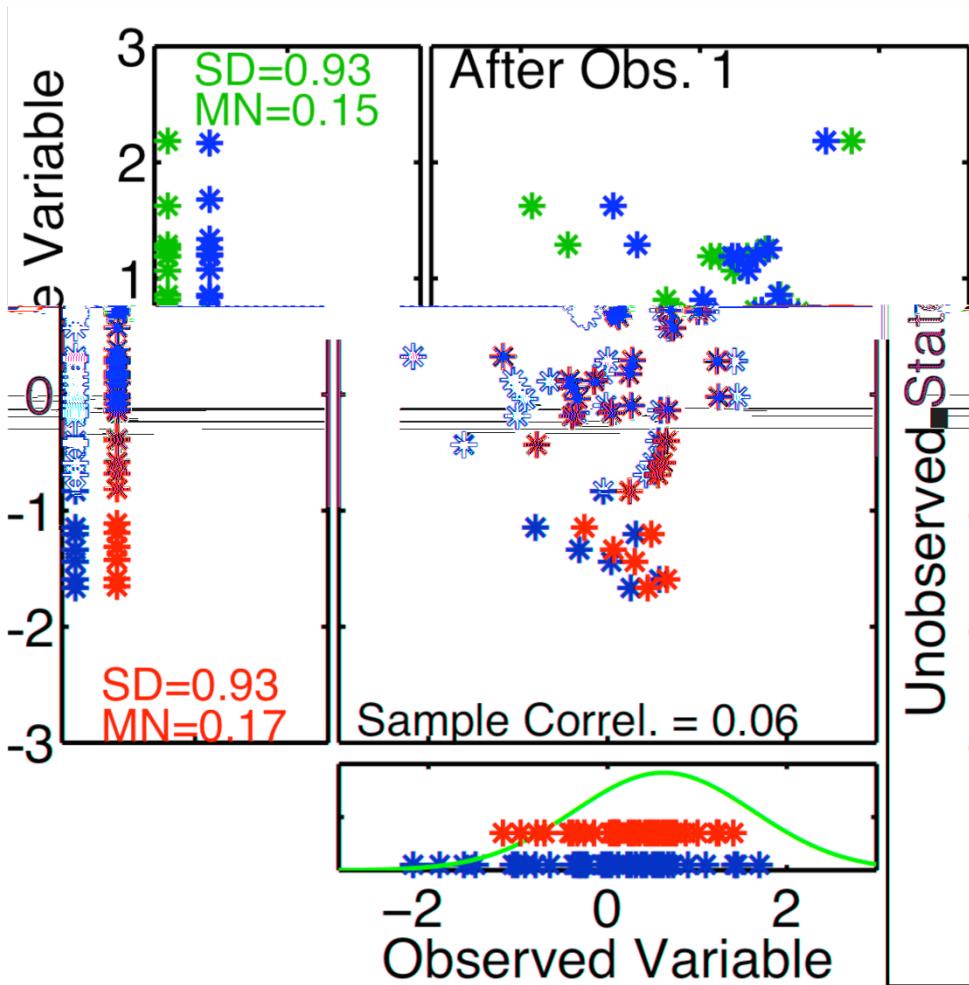
# Regression Sampling Error & Filter Divergence



Suppose unobserved state variable is known to be unrelated to set of observed variables.

Unobserved variable  
should remain unchanged.

# Regression Sampling Error & Filter Divergence

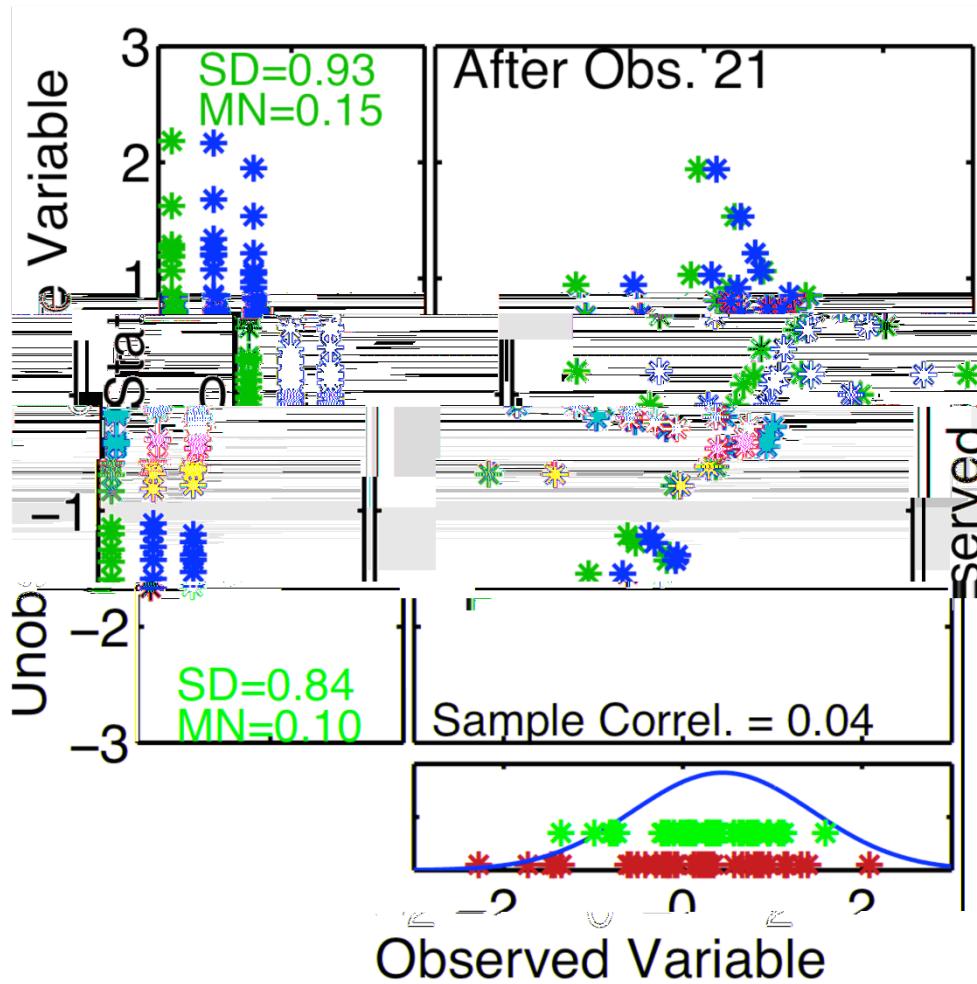


Suppose unobserved state variable is known to be unrelated to set of observed variables.

Finite samples from joint distribution will have non-zero correlation. Expected  $|\text{correl}| = 0.19$  for 20 samples.

After one observation, unobserved variable mean, standard deviation change.

# Regression Sampling Error & Filter Divergence

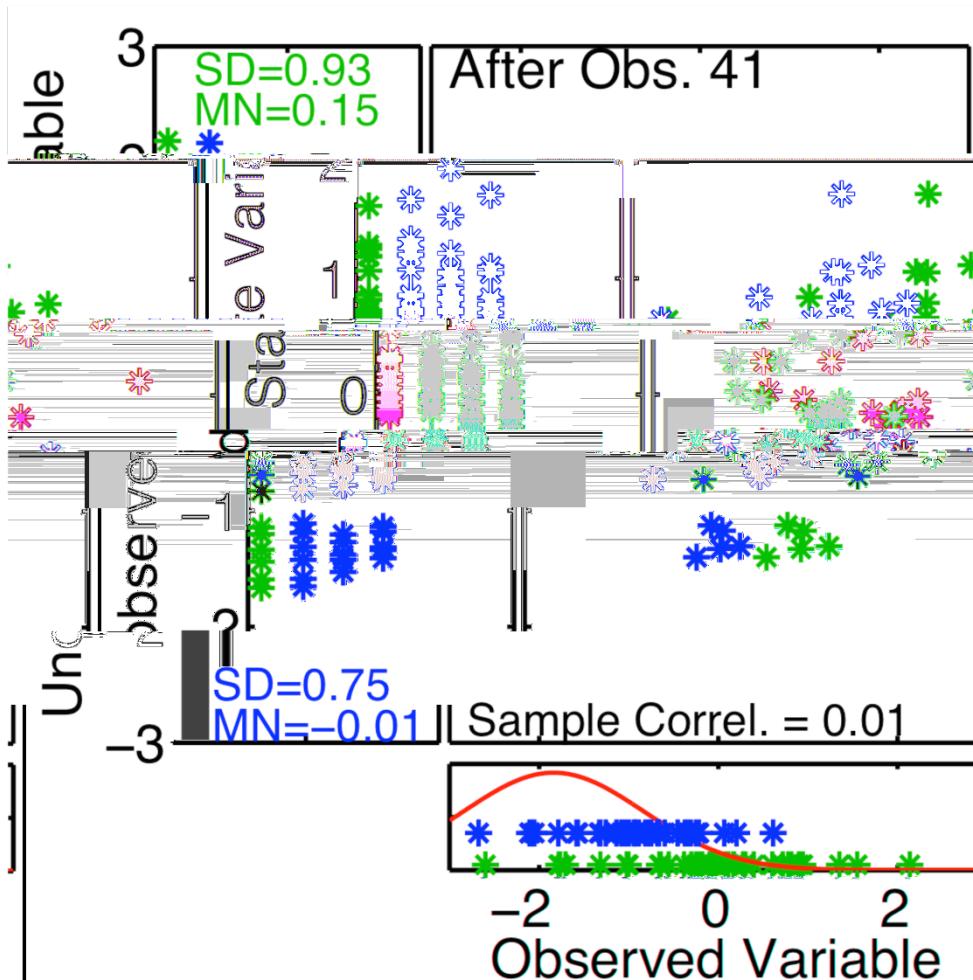


Suppose unobserved state variable is known to be unrelated to set of observed variables.

Unobserved variable should remain unchanged.

Unobserved mean follows a random walk as more observations are used.

# Regression Sampling Error & Filter Divergence



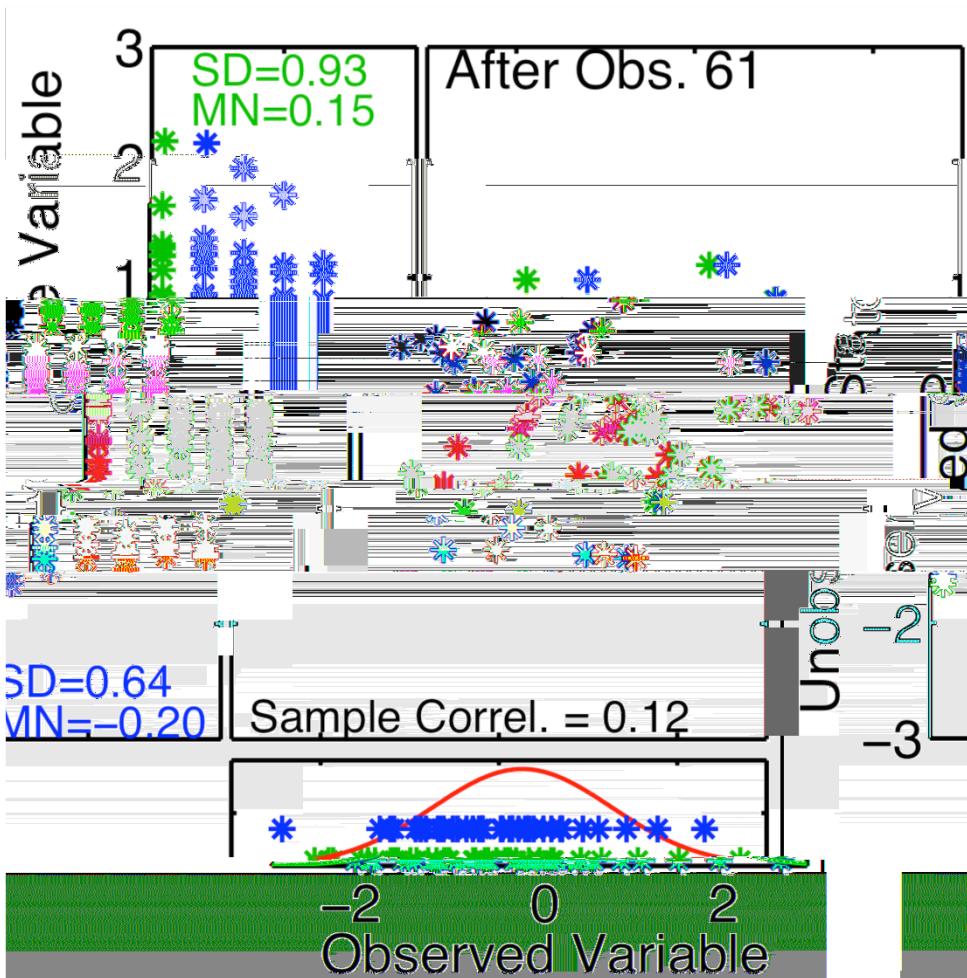
Suppose unobserved state variable is known to be unrelated to set of observed variables.

Unobserved variable should remain unchanged.

Unobserved S.D. systematically decreases.

Expected change in  $|SD|$  is negative for any non-zero sample correlation.

# Regression Sampling Error & Filter Divergence



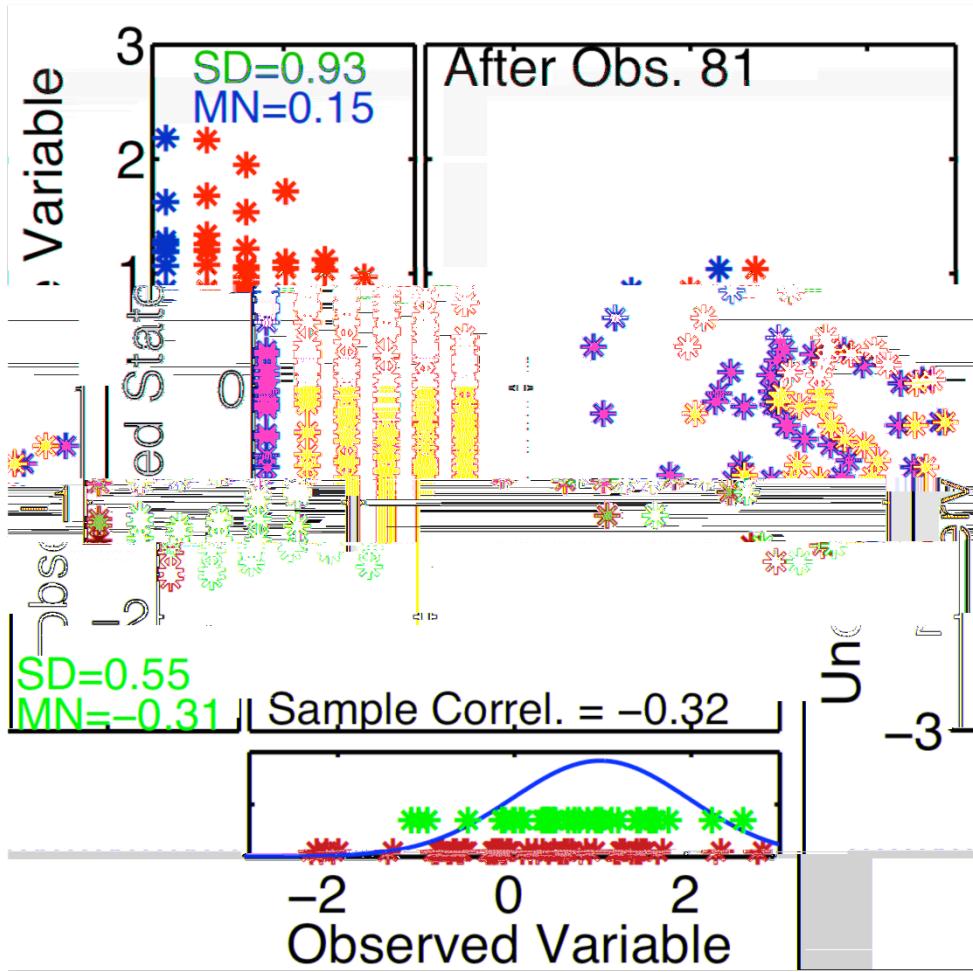
Suppose unobserved state variable is known to be unrelated to set of observed variables.

Unobserved variable should remain unchanged.

Unobserved S.D. systematically decreases.

Expected change in  $|SD|$  is negative for any non-zero sample correlation.

# Regression Sampling Error & Filter Divergence



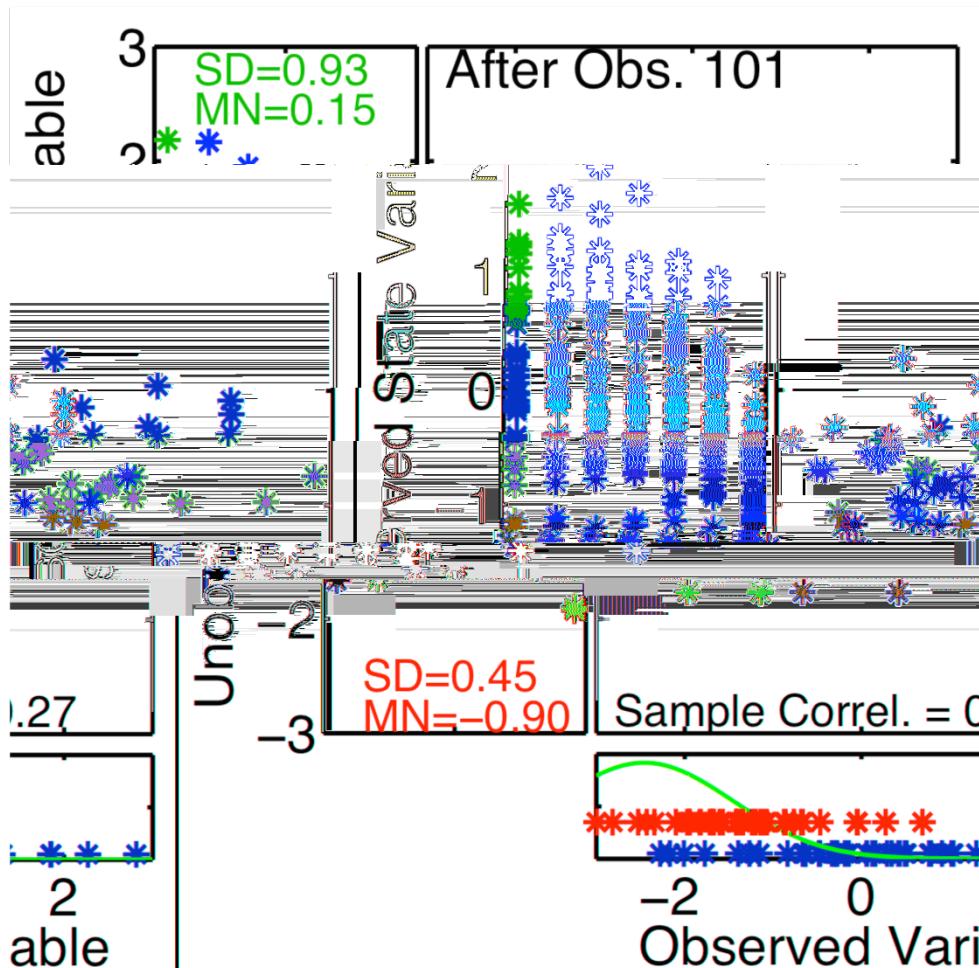
Suppose unobserved state variable is known to be unrelated to set of observed variables.

Unobserved variable should remain unchanged.

Unobserved S.D. systematically decreases.

Expected change in  $|SD|$  is negative for any non-zero sample correlation.

# Regression Sampling Error & Filter Divergence



Suppose unobserved state variable is known to be unrelated to set of observed variables.

Estimates of unobserved become too confident.

Give progressively less weight to meaningful obs.

Eventually, meaningful obs are essentially ignored.

# Filter Divergence

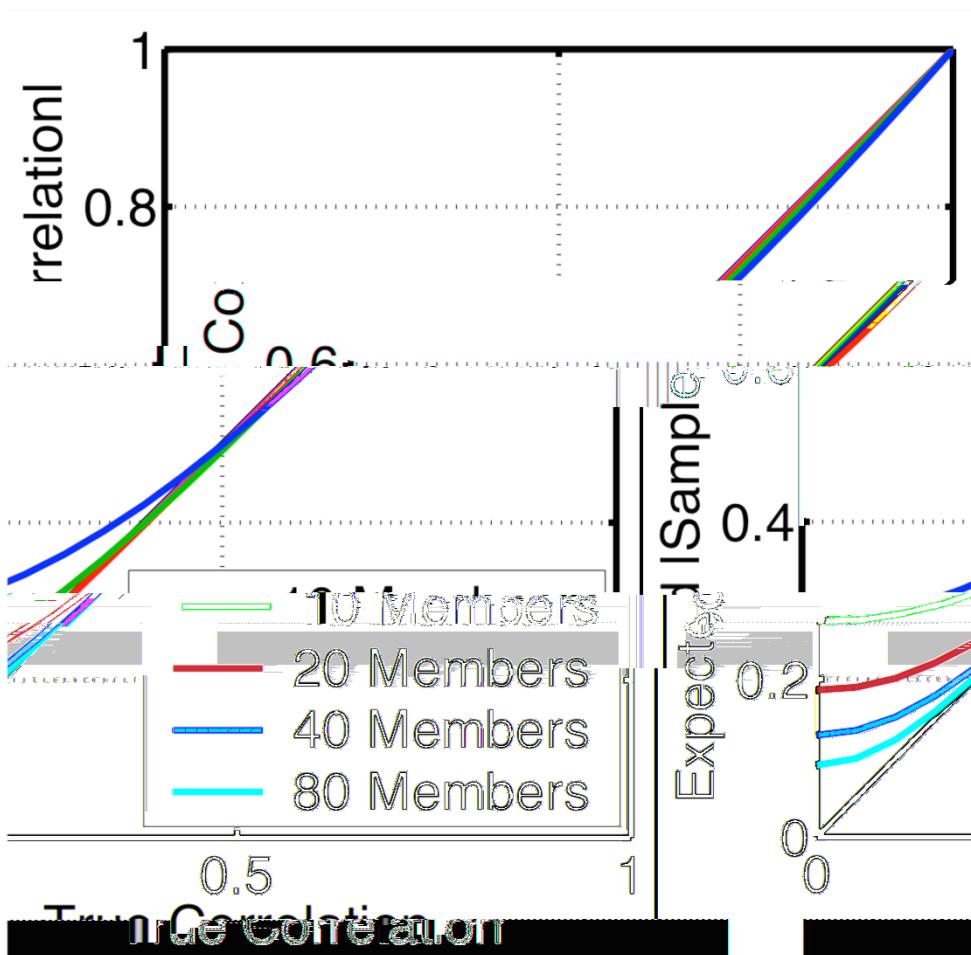
Ignoring meaningful observations due to overconfidence is a type of FILTER DIVERGENCE.

This was seen in initial Lorenz 96 (40-variable) experiment.

The spread became small => the filter thought it had a good estimate.

The error stayed large because good observations were ignored.

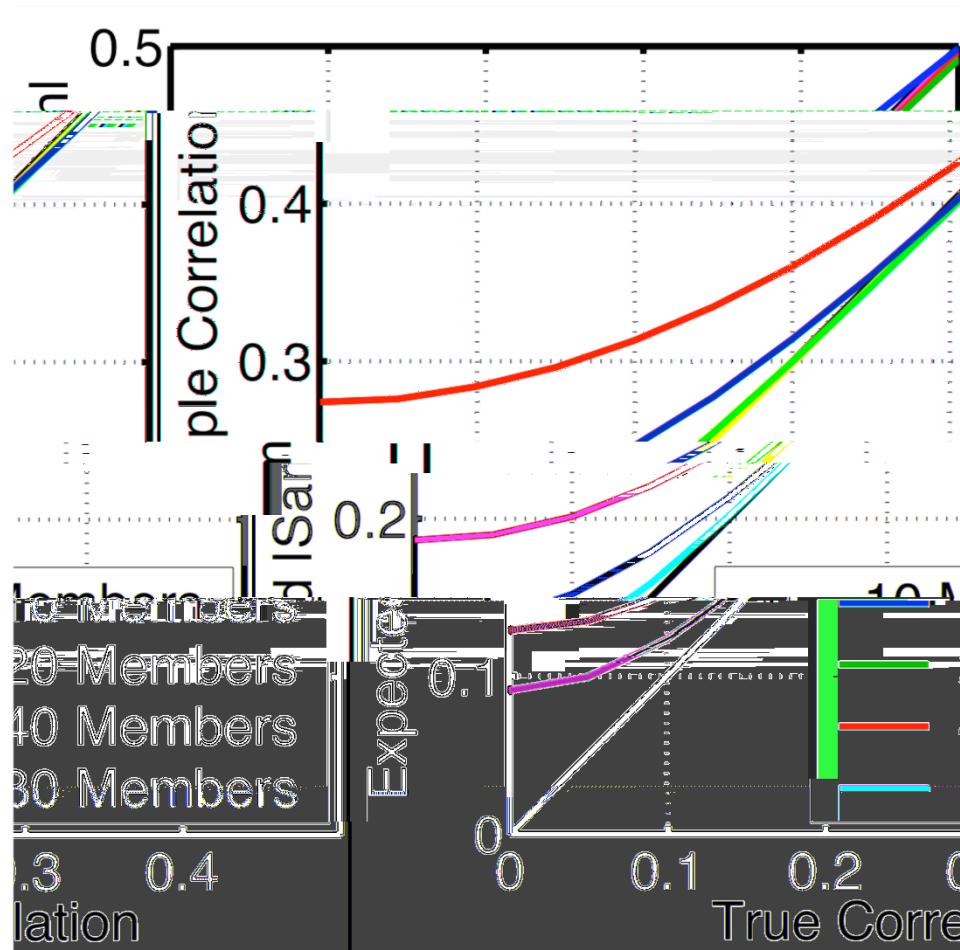
# Regression Sampling Error & Filter Divergence



Plot shows expected absolute value of sample correlation versus true correlation.

Error decreases with sample size and for larger |real correlations|.

# Regression Sampling Error & Filter Divergence



Plot shows expected absolute value of sample correlation versus true correlation.

For small true correlations, errors are still undesirably large even for 80 member ensembles.

# Dealing with Regression Sampling Error

1. Ignore it: if number of unrelated observations is small and there is some way of maintaining variance in priors.  
We did this in the 3 and 9 variable models.
2. Use larger ensembles to limit sampling error (test in lorenz\_96).  
This can get expensive for big problems.  
Try modifying *ens\_size* in *&filter\_nml* (try 40, 80, 160).

Note: For ensemble sizes greater than 80, set  
*&filter\_nml: perturb\_from\_single\_instance = .true.*

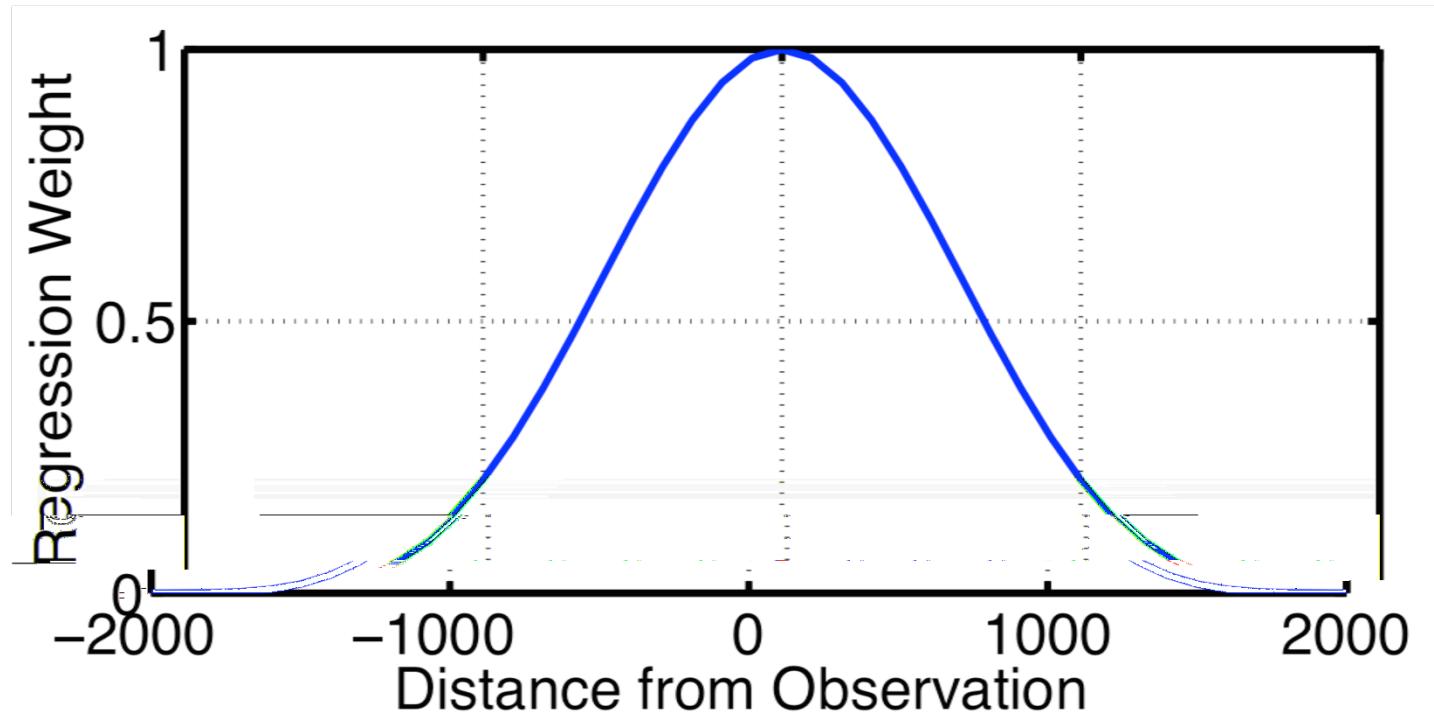
This tells DART to generate initial ensemble members using a random number generator, rather than reading them from an input file (which only contains 80 members in this directory).

# Dealing with Regression Sampling Error

1. Ignore it: if number of unrelated observations is small and there is some way of maintaining variance in priors.  
We did this in the 3 and 9 variable models.
2. Use larger ensembles to limit sampling error (test in lorenz\_96).  
This can get expensive for big problems.  
Try modifying *ens\_size* in `&filter_nml` (try 40, 80, 160).
3. Use additional a priori information about relation between observations and state variables.  
Don't let an observation impact state if they are known to be unrelated.
4. Try to determine the amount of sampling error and correct for it.  
There are many ways to do this; some simple, some complex.

# Dealing with Regression Sampling Error

3. Use additional a priori information about relation between observations and state variables.



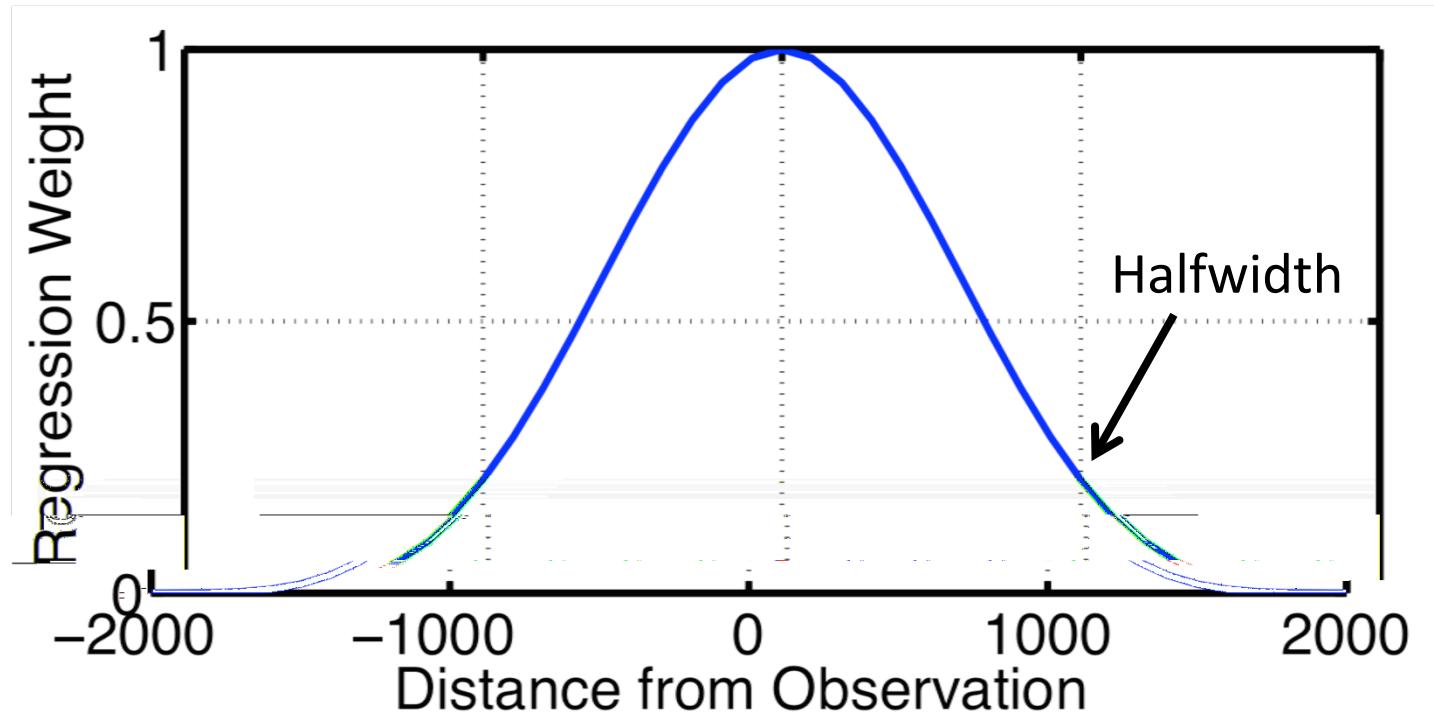
For atmospheric assimilation problems:

Weight regression as function of horizontal *distance* from observation.

Gaspari-Cohn: 5th order compactly supported polynomial.

# Dealing with Regression Sampling Error

3. Use additional a priori information about relation between observations and state variables.



Can use other functions to weight regression.

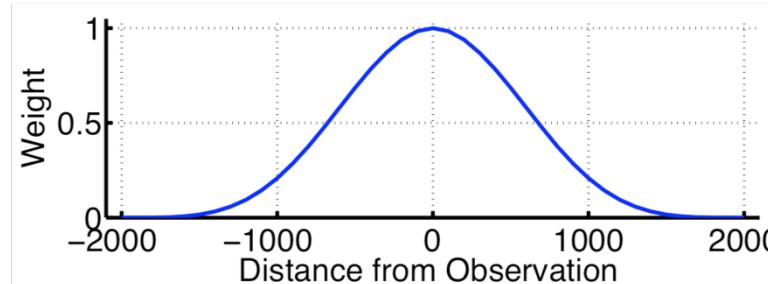
Unclear what distance means for some obs./state variable pairs.

Referred to as **LOCALIZATION**.

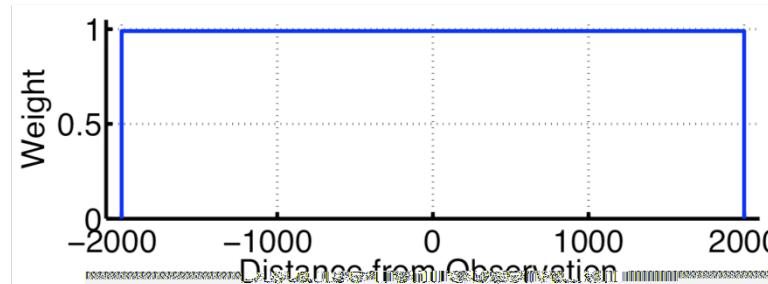
# DART provides several localization options

1. Different shapes for the localization function are available.  
Controlled by *select\_localization* in *&cov\_cutoff\_nml*.

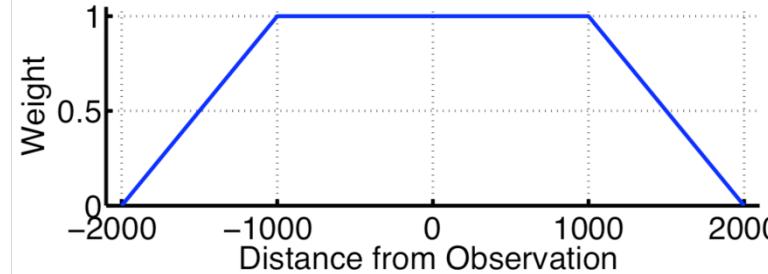
1=> Gaspari-Cohn



2=> Boxcar



3=> Ramped Boxcar



2. Halfwidth of localization function set by *cutoff* in *&assim\_tools\_nml*

# Experimenting with Lorenz 96

The lorenz\_96 domain is mapped to a [0, 1] periodic range.

Try a variety of half widths for a Gaspari Cohn localization by

```
&assim_tools_nml  
  filter_kind      = 1  
  cutoff           = 1000000.0  
  
...  
&filter_nml  
  ens_size = 20  
  perturb_from_single_instance = .false.  
...
```

Change

This just makes sure you start from  
same conditions each time.

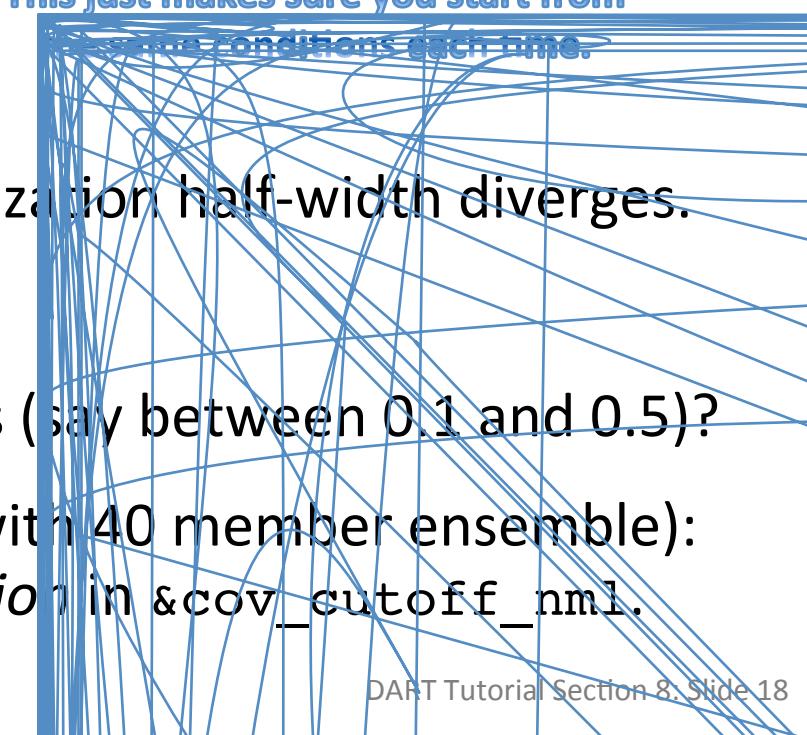
We already know that a very large localization half-width diverges.

What happens for a very small value?

What happens with intermediate values (say between 0.1 and 0.5)?

Can also try changing the shape, (best with 40 member ensemble):

Try option 2 or 3 for *select\_localization* in &cov\_cutoff\_nml.



# Dealing with Regression Sampling Error

4. Try to determine the amount of sampling error and correct for it.

Many ways to do this. DART implements one naive way:

1. Take set of increments from a given observation,
2. Suppose this observation and a state variable are not correlated,
3. Compute the expected decrease in spread given not correlated,
4. Add this amount of spread back into the state variable.

The expected decrease in spread is computed by off-line Monte Carlo.  
Results of off-line simulation are tabulated and applied.

(This can be a very useful technique when you're analytically clueless).

Try this algorithm: set

```
&assim_tools_nml: spread_restoration = .true.
```

How does it work with 20 ensemble members, no localization?

# Dealing with Regression Sampling Error

4. Try to determine the amount of sampling error and correct for it.

Many ways to do this. DART also implements a sampling error correction algorithm that can reduce but not eliminate problems. This algorithm ALMOST ALWAYS IMPROVES large model results.

Try this algorithm: set

```
&assim_tools_nml: sampling_error_correction = .true.
```

How does it work with 20 ensemble members, no localization?

# DART Tutorial Index to Sections

1. Filtering For a One Variable System
2. The DART Directory Tree
3. DART Runtime Control and Documentation
4. How should observations of a state variable impact an unobserved state variable?  
Multivariate assimilation.
5. Comprehensive Filtering Theory: Non-Identity Observations and the Joint Phase Space
6. Other Updates for An Observed Variable
7. Some Additional Low-Order Models
8. Dealing with Sampling Error
9. More on Dealing with Error; Inflation
10. Regression and Nonlinear Effects
11. Creating DART Executables
12. Adaptive Inflation
13. Hierarchical Group Filters and Localization
14. Quality Control
15. DART Experiments: Control and Design
16. Diagnostic Output
17. Creating Observation Sequences
18. Lost in Phase Space: The Challenge of Not Knowing the Truth
19. DART-Compliant Models and Making Models Compliant
20. Model Parameter Estimation
21. Observation Types and Observing System Design
22. Parallel Algorithm Implementation
23. Location module design (not available)
24. Fixed lag smoother (not available)
25. A simple 1D advection model: Tracer Data Assimilation