# Information in Models and Data

#### **Grey Nearing**

National Center for Atmospheric Research NASA Goddard Space Flight Center University of Maryland Baltimore County

#### Grey Nearing - 26/04/2016 Schneefernerhaus

#### What is a Model?

"The totality of our so-called knowledge or beliefs, from the most casual matters of geography and history to the profoundest laws of atomic physics or even of pure mathematics and logic, is a man-made fabric which impinges on experience only along the edges." - Van Quine (1951)





## Ontology & Epistemology in Geophysical Models



#### The Description of an Experiment



Before running the experiment: $p(D_1 \land D_2) = p(D_2 | D_1) \times p(D_1)$ After running the experiment: $d(D_1 \land D_2) = d(D_2) + d(D_1) - d(D_1 \lor D_2)$ 

### The Open Problem

**The Demarcation Problem**: Science is either well-defined and impractical or un-defined and practical.

- Hume: Induction cannot be rigorously supported.
- **Popper:** Therefore, science must be deductive.
- Salmon: Falsification is not practical because few (all) models are falsified.
- Jaynes: Bayesianism is at least consistent with the axioms of deductive logic, however fundamentally inductive.

To reconcile the falsification criteria with Bayesian model evaluation: Does the model contain as much information as the observations?

#### How well are we doing right now?

Our best physically-based surface hydrology models cannot beat linear regressions that have <u>no state memory</u>, and which are <u>trained out-of-sample</u> and <u>extrapolated globally</u>.



M. Best et al. (2015) "The Plumbing of Land Surface Models" Journal of Hydrometeorology

#### Information Theory – Basic Principles





**Consequence 2**: Bounded under transformations.

$$I(D; U) \geq I(D; M(U))$$

#### Measuring Information



$$U \rightarrow \mathcal{M} \rightarrow Y$$
  
 $\stackrel{Premise: No system is isolated. A scientist perturbs a system and measures it's response$ 

**Definition:** The information content of data is defined by our ability to derive asymptotic relationships between measured perturbations and responses.

$$D = \mathcal{R}_e(U)$$

### Example 1: Uncertainty Segregation



**Definition:** The information content of data is defined by our ability to derive asymptotic relationships between measured perturbations and responses.

$$D = \mathcal{R}_e(U)$$





I(D; M) = H(D) - H(D|M)

### Example 1: Uncertainty Segregation



**Definition:** The information content of data is defined by our ability to derive asymptotic relationships between measured perturbations and responses.

$$D = \mathcal{R}_e(U)$$





G. Nearing et al. (2016) "Benchmarking NLDAS-2 to Separate Uncertainty Contributions" JHM

Model:  $d\mathbf{x} = \mu(\mathbf{x}, \mathbf{u})dt + \sigma(\mathbf{x}, \mathbf{u})dB_t$ Data Assimilation:  $p(\mathbf{x}_t | \mathbf{y}_t) \propto h(\mathbf{y}_t | \mathbf{x}_t)m(\mathbf{x}_t | \mathbf{u}_{1:t})$ 





Model:  $d\mathbf{x} = \mu(\mathbf{x}, \mathbf{u})dt + \sigma(\mathbf{x}, \mathbf{u})dB_t$ Data Assimilation:  $p(\mathbf{x}_t | \mathbf{y}_t) \propto h(\mathbf{y}_t | \mathbf{x}_t)m(\mathbf{x}_t | \mathbf{u}_{1:t})$ 







- AMSR-E Soil Moisture Retrievals
- NOAH-MP Model
- Ensemble Kalman Filter

The Ensemble Kalman Filter is only about 30% efficient in this experiment.

| Measurement   | Interpretation                              | Value<br>[nats/nats] |
|---|---|----------------------|
| Information in Noah simulations                         | Model<br>H(a)<br>Evaluation Data<br>H(a)    | 0.17                 |
| Information in LPRM (AMSR-E)<br>retrievals              | Nodel<br>H(a)<br>Retrievals<br>H(y)         | 0.24                 |
| Total information in Noah and LPRM<br>(AMSR-E) together | Model<br>H(a)<br>R(z)<br>Retrievals<br>H(y) | 0.61                 |
| Information from Data Assimilation                      | Availables<br>H(p)                          | 0.18                 |
| EnKF Efficiency   | Model<br>H(A)<br>Retrievals<br>H(y)         | 0.29                 |

### Measuring Information



| General Definition   | $I(D; M) = E\left[f\left(\frac{p_{D M}}{p_D}\right)\right]$                             |
|----------------------|---|
| Specific Definition: | $f(\xi) = -\ln(\xi)$<br>$I(D; M) = E\left[\ln(p_{D M})\right] - E\left[\ln(p_D)\right]$ |
| Linearity Property:  | I(D; M) = H(D) - H(D M)   |



Information *quality* is related to whether the probabilities move in the right direction.

#### Example 3: Information from Hypotheses



"it must be demonstrated that the model physics actually adds information to the prediction system." - van den Hurk et al. (2011; BAMS)



#### Example 3: Information from Hypotheses



"it must be demonstrated that the model physics actually adds information to the prediction system." - van den Hurk et al. (2011; BAMS)



#### Summary



The ontological model cannot be separated from the epistemological model.



Models translate information.



The model of an experiment is a logarithm.



This model of an experiment yields a deductive science.



Information is easier to work with than probabilities.

