National Aeronautics and
Space Administration

**Jet Propulsion Laboratory**
California Institute of Technology
Pasadena, California

# Statistics in UQ and UQ in Statistics

Amy Braverman

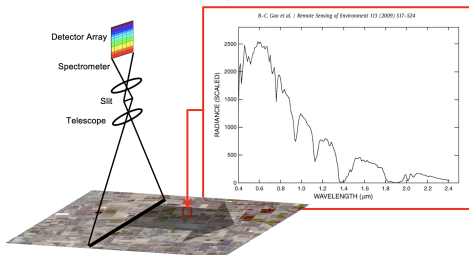Jet Propulsion Laboratory, California Institute of Technology

February 26, 2024

National Aeronautics and
Space Administration

**Jet Propulsion Laboratory**
California Institute of Technology
Pasadena, California

Outline
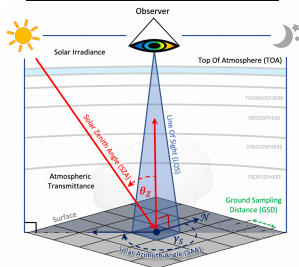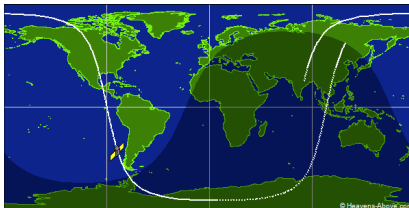
- Problem statement

- Background

- UQ formalism

- Roles of math and stat

- Modifying the UQ formalism

- Model discrepancy

- Methodology

- Summary/discussion

- ► Develop methods to quantify uncertainty in remote sensing data products delivered by the Orbiting Carbon Observatory 2 and 3 missions.

- ► The methods must be "off-line" and not interfere with operational data processing.

- ► They must be computationally efficient enough to keep up with the data stream.

- ► This problem and our solution are discussed in detail in Braverman et al., 2021. (doi: 10.1137/19M1304283).

National Aeronautics and
Space Administration

**Jet Propulsion Laboratory**
California Institute of Technology
Pasadena, California

▶ Passive remote sensing instruments measure photon counts in bins of a discretized electromagnetic spectrum.

▶ The sun provides incoming photons, which are scattered and absorbed in ways that depend on the media (atmosphere or surface) with which they interact.

▶ May also be complicated by thermal emission.

▶ The instrument discretizes the spatial field into "footprints" and aggregates photons over both footprint and spectral bin.

# Background

National Aeronautics and
Space Administration

**Jet Propulsion Laboratory**
California Institute of Technology
Pasadena, California

Remote sensing levels of data processing:

- ► Level 0: raw photon counts direct from satellite

- ► Level 1: georectified and calibrated radiances

- ► **Level 2: estimates of geophysical state**

- ► Level 3: "statistical summaries" of Level 2 on uniform space-time grid

- ► Level 4: output of models or data assimilation

**Level 2 "data" aren't "data"; they are inferences!**

*When drawing scientific conclusions or making policy decisions, it is crucial to take account of uncertainties in these inferences.*

National Aeronautics and
Space Administration

**Jet Propulsion Laboratory**
California Institute of Technology
Pasadena, California

## Remote sensing observing system:

| True state | Forward function | Noiseless radiance | Instru-ment | Obser-vation | Retrieval | State estimate |
|---|---|---|---|---|---|---|

$$\mathbf{X} \longrightarrow \boxed{\mathbf{F}_0(\cdot, \mathbf{B}_0)} \longrightarrow \mathbf{Y}_0 \longrightarrow \boxed{\mathbf{Y}_0 + \epsilon} \longrightarrow \mathbf{Y} \longrightarrow \boxed{\mathbf{R}(\cdot, \mathbf{F}_1, \mathbf{B}_1 \ldots)} \longrightarrow \hat{\mathbf{X}}$$

$F_0$ = nature's true forward function; $\mathbf{B}_0$ = other true quantities.

$F_1$ = forward model used in retrieval, $R$; $\mathbf{B}_1$ = other retrieval inputs.

$\epsilon$ = instrument measurement error.

$\ldots$ = other retrieval algorithm inputs.

National Aeronautics and
Space Administration

**Jet Propulsion Laboratory**
California Institute of Technology
Pasadena, California

UQ formalism

## VVUQ:



Adapted from Wu et al, (2018). DOI: 10.1016/j.nucengdes.2018.06.004.

National Aeronautics and
Space Administration

**Jet Propulsion Laboratory**
California Institute of Technology
Pasadena, California

UQ formalism

$\mathbf{F}(\cdot)$ → Forward UQ → $\hat{\mathbf{y}} = \mathbf{F}(\mathbf{x}^*)$ (ensemble)

Distribution of $\hat{\mathbf{y}}$ describes uncertainty in predictions made by $\mathbf{F}$

sample: $\mathbf{x}^*$ (ensemble)

Dist($\mathbf{x}$) ← Inverse UQ ← Direct observations: $\mathbf{y}_1^{(\mathrm{R})}, \ldots, \mathbf{y}_N^{(\mathrm{R})}$

$P\left(\mathbf{x} | \mathbf{y}_1^{(\mathrm{R})}, \ldots, \mathbf{y}_N^{(\mathrm{R})}\right)$

$\mathbf{y}_i^{(\mathrm{R})} = \mathbf{F}^{\mathbf{True}}(\mathbf{x}) + \boldsymbol{\epsilon}_i$

National Aeronautics and
Space Administration

**Jet Propulsion Laboratory**
California Institute of Technology
Pasadena, California

▶ Statistical methods: inference from observations about unknown probabilistic model

  ▶ estimation and hypothesis testing

  ▶ exploratory data analysis, density estimation, unsupervised learning

  ▶ regression, supervised learning, to uncover significant relationships

  ▶ uncover, test, and quantify relationships from data

  ▶ use estimated model to make statistical predictions with uncertainty.

▶ Statistical models inherently carry uncertainties with them.

# Roles of Math and Stat



$$P\left(\mathbf{x}|\mathbf{y}_1^{(R)}, \ldots, \mathbf{y}_N^{(R)}\right)$$

$$L\left(\mathbf{x};\mathbf{y}_1^{(R)}, \ldots, \mathbf{y}_N^{(R)}\right)$$

$$\mathbf{y}_i^{(R)} = \mathbf{F}^{\mathbf{True}}(\mathbf{x}) + \epsilon_i$$

National Aeronautics and
Space Administration

**Jet Propulsion Laboratory**
California Institute of Technology
Pasadena, California

- ▶ Mathematical UQ: mathematical approaches for understanding sources of uncertainty in **F** and facilitating efficient forward UQ.

    - ▶ exploit structure and properties of **F** to guide forward UQ

    - ▶ alternatives to brute-force Monte Carlo forward UQ

    - ▶ numerical and other approximations for speed and efficiency

    - ▶ optimization!

- ▶ Uncertainty expressed through probability distributions, and driven by probabilistic description of input uncertainties.

National Aeronautics and
Space Administration

**Jet Propulsion Laboratory**
California Institute of Technology
Pasadena, California

# Roles of Math and Stat



Figure content labels: $\mathbf{F}(\cdot)$; Forward UQ; $\hat{\mathbf{y}} = \mathbf{F}(\mathbf{x}^*)$ (ensemble); Mathematical guts of $\mathbf{F}$; Etc.; Sparse grids; Polynomial expansions; Reduced order models/surrogates; sample: $\mathbf{x}^*$ (ensemble); Local and Global sensitivity analysis; Distribution of $\hat{\mathbf{y}}$ describes uncertainty in predictions made by $\mathbf{F}$; Dist($\mathbf{x}$); Inverse UQ; $P\left(\mathbf{x}|\mathbf{y}_1^{(\mathrm{R})},\ldots,\mathbf{y}_N^{(\mathrm{R})}\right)$; Direct observations: $\mathbf{y}_1^{(\mathrm{R})},\ldots,\mathbf{y}_N^{(\mathrm{R})}$; $\mathbf{y}_i^{(\mathrm{R})} = \mathbf{F}^{\mathsf{True}}(\mathbf{x}) + \epsilon_i$

National Aeronautics and
Space Administration

**Jet Propulsion Laboratory**
California Institute of Technology
Pasadena, California

Roles of Math and Stat

► Inverse problems: infer the state of a system from noisy, indirect measurements.

  ► heavy use of Bayesian methods

  ► overlaps substantially with statistics, but more focussed on this class of problems

  ► emphasis on algorithms/samplers

  ► because result is a distribution, easy forward propagation

# Roles of math and stat



$F(\cdot)$

Forward UQ

$\hat{y} = F(x^*)$
(ensemble)

Design and analysis
of computer experiments

Distribution of $\hat{y}$
describes uncertainty in
predictions made by $F$

sample: $x^*$
(ensemble)

Data
Assimilation

Model
discrepancy

MCMC and
friends

Etc.

Dist($x$)

Inverse problems

Direct
observations:
$y_1^{(R)}, \ldots, y_N^{(R)}$

$P\left(x | y_1^{(R)}, \ldots, y_N^{(R)}\right)$

$y_i^{(R)} = F^{True}(x^{(R)}) + \epsilon_i$

# Modifying the UQ formalism



$$P\left(\mathbf{x} \mid \mathbf{y}_1^{(\mathrm{R})}, \ldots, \mathbf{y}_N^{(\mathrm{R})}, \hat{\mathbf{y}}_1, \ldots, \hat{\mathbf{y}}_M\right)$$

$$\mathbf{y}_i^{(\mathrm{R})} = \mathbf{F}^{\mathbf{True}}(\mathbf{x}^{(\mathrm{R})}) + \boldsymbol{\epsilon}_i$$

Boxes and labels in diagram:
- $\mathbf{F}(\cdot)$
- Forward UQ
- $\hat{\mathbf{y}} = \mathbf{F}(\mathbf{x}^*)$ (ensemble)
- sample: $\mathbf{x}^*$ (ensemble)
- Dist($\mathbf{x}$)
- Inverse UQ
- Direct observations: $\mathbf{y}_1^{(\mathrm{R})}, \ldots, \mathbf{y}_N^{(\mathrm{R})}$

# Modifying the UQ formalism



$\mathbf{F}(\cdot)$

Forward UQ

$\hat{\mathbf{y}} = \mathbf{F}(\mathbf{X}^*)$
(ensemble)

This is fine if the objective is to perform UQ on $\mathbf{F}$. But it's not.

The objective is to perform UQ on the operational retrieval algorithm, $\mathbf{R}$.

sample: $\mathbf{X}^*$
(ensemble)

Dist($\mathbf{X}$)

Inverse UQ

Direct observations:
$\mathbf{y}_1^{(R)}, \ldots, \mathbf{y}_N^{(R)}$

$P\left(\mathbf{X} | \mathbf{y}_1^{(R)}, \ldots, \mathbf{y}_N^{(R)}, \hat{\mathbf{y}}_1, \ldots, \hat{\mathbf{y}}_M\right)$

$\mathbf{y}_i^{(R)} = \mathbf{F}^{\text{True}}\left(\mathbf{x}_i^{(R)}\right) + \boldsymbol{\epsilon}_i$

Note: this is not the operational retrieval algorithm!

# Modifying the UQ formalism



$\mathbf{R}(\cdot, \ldots)$

Forward UQ

$\hat{\mathbf{x}} = \mathbf{R}(\mathbf{Y}^*)$
(ensemble)

Operational retrieval algorithm

sample: $\mathbf{Y}^*$
(ensemble)

Problem: we have very few
direct observations of $\mathbf{x}^{(R)}$.

Dist($\mathbf{Y}$)

Inverse UQ

Direct
observations:
$\mathbf{x}_1^{(R)}, \ldots, \mathbf{x}_N^{(R)}$

$P\left(\mathbf{Y} | \mathbf{x}_1^{(R)}, \ldots, \mathbf{x}_N^{(R)}, \hat{\mathbf{x}}_1, \ldots, \hat{\mathbf{x}}_M\right)$

$\mathbf{x}_i^{(R)} = \mathbf{R}^{\text{True}}\left(\mathbf{y}_i^{(R)}\right) + \boldsymbol{\tau}_i$

$\mathbf{R}^{\text{True}} = \left[\mathbf{F}^{\text{True}}\right]^{-1}$

National Aeronautics and
Space Administration

**Jet Propulsion Laboratory**
California Institute of Technology
Pasadena, California

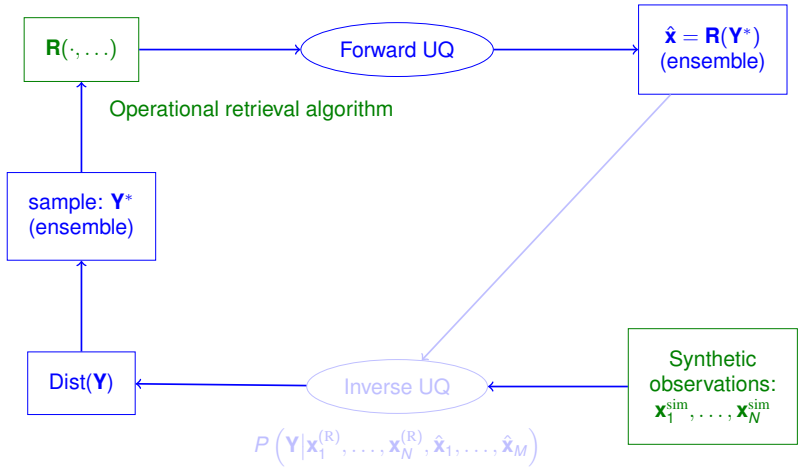**To recap,**

▶ We want to perform UQ on the operational retrieval algorithm, **R**. (Note that **F** is embedded in, and is thus part of, **R**.)

▶ This requires a computational experiment in which we sample over **R**'s inputs to get an ensemble of outputs ($\hat{\mathbf{x}}$'s) that can be compared to direct observations, $\mathbf{x}^{(R)}$.

▶ We do not have enough instances of $\mathbf{x}^{(R)}$ to do this.

▶ Moreover, performing inverse UQ without oversimplifying (e.g., using MCMC) is computationally infeasible. (The oversimplified version *is* **R**).
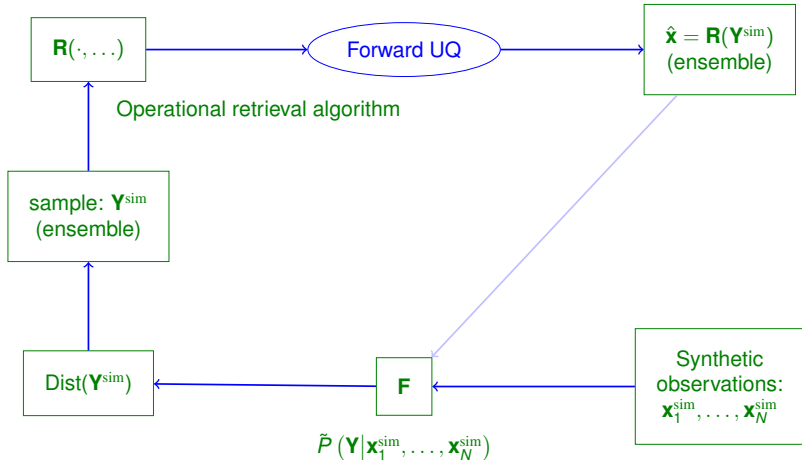
▶ So what can we do?

National Aeronautics and
Space Administration

Jet Propulsion Laboratory
California Institute of Technology
Pasadena, California

Modifying the UQ formalism

$\mathbf{R}(\cdot, \ldots)$

Forward UQ

$\hat{\mathbf{x}} = \mathbf{R}(\mathbf{Y}^*)$
(ensemble)

Operational retrieval algorithm

sample: $\mathbf{Y}^*$
(ensemble)

Dist($\mathbf{Y}$)

Inverse UQ

Synthetic
observations:
$\mathbf{x}_1^{\text{sim}}, \ldots, \mathbf{x}_N^{\text{sim}}$

$P\left(\mathbf{Y} \mid \mathbf{x}_1^{(\mathrm{R})}, \ldots, \mathbf{x}_N^{(\mathrm{R})}, \hat{\mathbf{x}}_1, \ldots, \hat{\mathbf{x}}_M\right)$

# Modifying the UQ formalism

# Modifying the UQ formalism



$\mathbf{R}(\cdot, \dots)$

Forward UQ

$\hat{\mathbf{x}} = \mathbf{R}(\mathbf{Y}^{\text{sim}})$
(ensemble)

Operational retrieval algorithm

sample: $\mathbf{Y}^{\text{sim}}$
(ensemble)

Think like statistician

$\tilde{P}\left(\hat{\mathbf{X}}^{\text{sim}} | \mathbf{X}^{\text{sim}}\right)$

Dist($\mathbf{Y}^{\text{sim}}$)

$\mathbf{F}$

Synthetic
observations:
$\mathbf{x}_1^{\text{sim}}, \dots, \mathbf{x}_N^{\text{sim}}$

$\tilde{P}\left(\mathbf{Y} | \mathbf{x}_1^{\text{sim}}, \dots, \mathbf{x}_N^{\text{sim}}\right)$

# Modifying the UQ formalism



$\mathbf{R}(\cdot, \mathbf{F}, \dots)$

Forward UQ

Synthetic estimates
$\hat{\mathbf{x}}_1^{\text{sim}}, \dots, \hat{\mathbf{x}}_N^{\text{sim}}$

Operational retrieval algorithm

Think like statistician

$\tilde{P}\left(\hat{\mathbf{X}}^{\text{sim}}, \hat{\mathbf{Y}}^{\text{sim}} | \mathbf{X}^{\text{sim}}\right)$

Synthetic
observations
$\mathbf{y}_1^{\text{sim}}, \dots, \mathbf{y}_N^{\text{sim}}$

$\mathbf{F}$

Synthetic
observations:
$\mathbf{x}_1^{\text{sim}}, \dots, \mathbf{x}_N^{\text{sim}}$

$\tilde{P}\left(\mathbf{Y} | \mathbf{x}_1^{\text{sim}}, \dots, \mathbf{x}_N^{\text{sim}}\right)$

$\boldsymbol{\epsilon}^{\text{sim}} \sim \text{MVN}\left(\boldsymbol{\mu}_{\boldsymbol{\epsilon}}, \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}}, \right)$    Measurement error

National Aeronautics and
Space Administration

**Jet Propulsion Laboratory**
California Institute of Technology
Pasadena, California

What does "think like a statistician" mean?

▶ Statisticians invent new estimators and quantify their operating characteristics.

▶ Here, we quantify the operating characteristics of the system

$$\mathbf{X}^{\mathrm{sim}} \to \mathbf{F} \to \mathbf{Y}^{\mathrm{sim}} \to \mathbf{R} \to \hat{\mathbf{X}}^{\mathrm{sim}}$$

with and empirical estimate of $P\left(\hat{\mathbf{X}}^{\mathrm{sim}}, \mathbf{X}^{\mathrm{sim}}\right)$.

▶ Proposition: uncertainty is quantified by any useful reduction of $\tilde{P}\left(\hat{\mathbf{X}}^{\mathrm{sim}}, \mathbf{X}^{\mathrm{sim}}\right)$, e.g., $\tilde{P}\left(\mathbf{X}^{\mathrm{sim}} | \hat{\mathbf{X}}^{\mathrm{sim}},\right)$.

Fine, but

1. what about the real system?

2. inverse crime: $\mathbf{F}$ is used twice (once to create $\mathbf{y}^{\text{sim}}$ and once in $\mathbf{R}$).

National Aeronautics and
Space Administration

**Jet Propulsion Laboratory**
California Institute of Technology
Pasadena, California

Fine, but

1. what about the real system?

2. inverse crime: $\mathbf{F}$ is used twice (once to create $\mathbf{y}^{\mathrm{sim}}$ and once in $\mathbf{R}$).

1. We use the learned relationship $\tilde{P}\left(\mathbf{X}^{\mathrm{sim}}|\hat{\mathbf{X}}^{\mathrm{sim}},\right)$ to quantify uncertainty is an actual instance of $\hat{\mathbf{X}}$:

$$\tilde{P}\left(\mathbf{X}^{\mathrm{True}}|\hat{\mathbf{X}}^{\mathrm{Actual}}\right).$$

National Aeronautics and
Space Administration

**Jet Propulsion Laboratory**
California Institute of Technology
Pasadena, California

Fine, but

1. what about the real system?

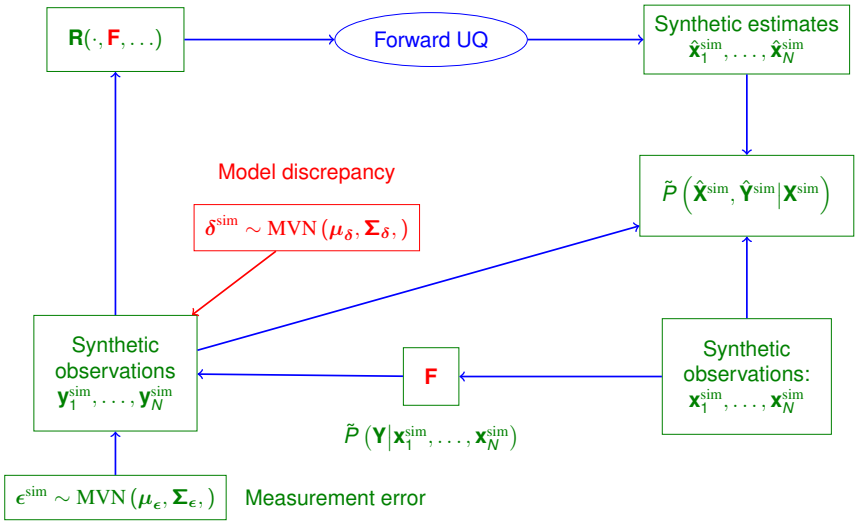2. inverse crime: $\mathbf{F}$ is used twice (once to create $\mathbf{y}^{\text{sim}}$ and once in $\mathbf{R}$).

1. We use the learned relationship $\tilde{P}\left(\mathbf{X}^{\text{sim}}|\hat{\mathbf{X}}^{\text{sim}},\right)$ to quantify uncertainty in an actual instance of $\hat{\mathbf{X}}$:

$$\tilde{P}\left(\mathbf{X}^{\text{True}}|\hat{\mathbf{X}}^{\text{Actual}}\right).$$

2. Introduce model discrepancy.

National Aeronautics and
Space Administration

**Jet Propulsion Laboratory**
California Institute of Technology
Pasadena, California

Model discrepancy

National Aeronautics and
Space Administration

**Jet Propulsion Laboratory**
California Institute of Technology
Pasadena, California

# Model discrepancy

- Model discrepancy is $\delta = \mathbf{F}^{\text{True}}\left(\mathbf{X}^{\text{True}}\right) - \mathbf{F}\left(\mathbf{X}^{\text{True}}\right)$.

- We would like to simulate from the distribution of $\delta \sim \text{MVN}\left(\boldsymbol{\mu}_{\delta}, \boldsymbol{\Sigma}_{\delta}\right)$.

- Assume this distribution is Gaussian with mean $\boldsymbol{\mu}_{\delta} \approx \text{E}(\delta^{\text{sim}})$ and covariance matrix $\boldsymbol{\Sigma}_{\delta} \approx \text{cov}(\delta^{\text{sim}})$.

- We have noisy samples, $\mathbf{Y}_i = \mathbf{F}^{\text{True}}\left(\mathbf{X}_i^{\text{True}}\right) + \boldsymbol{\epsilon}_i^{\text{True}}, \ i = 1, \dots, N$.

- We don't have $\mathbf{F}\left(\mathbf{X}^{\text{True}}\right)$, but we do have $\left[\mathbf{F}\left(\mathbf{X}^{\text{sim}}\right) - \mathbf{F}(\hat{\mathbf{X}}^{\text{sim}})\right]$, which motivates the approximation,

$$\delta^{\text{sim}} \approx \mathbf{F}^{\text{True}}\left(\mathbf{X}^{\text{True}}\right) - \mathbf{F}\left(\hat{\mathbf{X}}^{\text{Actual}}\right) - \left[\mathbf{F}(\mathbf{X}^{\text{sim}}) - \mathbf{F}(\hat{\mathbf{X}}^{\text{sim}})\right].$$

National Aeronautics and
Space Administration

**Jet Propulsion Laboratory**
California Institute of Technology
Pasadena, California

Model discrepancy

- Let $\mathbf{Y}^{\text{Actual}} \equiv \mathbf{F}^{\text{True}}\left(\mathbf{X}^{\text{True}}\right) + \boldsymbol{\epsilon}^{\text{True}}$, and $\hat{\mathbf{Y}}^{\text{Actual}} \equiv \mathbf{F}\left(\hat{\mathbf{X}}^{\text{Actual}}\right)$, and similarly for simulated. Then,

$$\boldsymbol{\delta}^{\text{sim}} \approx \left(\mathbf{Y}^{\text{Actual}} - \boldsymbol{\epsilon}^{\text{True}} - \hat{\mathbf{Y}}^{\text{Actual}}\right) - \left(\mathbf{Y}^{\text{sim}} - \hat{\mathbf{Y}}^{\text{sim}}\right),$$

$$\boldsymbol{\delta}^{\text{sim}} + \boldsymbol{\epsilon} \approx \left(\mathbf{Y}^{\text{Actual}} - \hat{\mathbf{Y}}^{\text{Actual}}\right) - \left(\mathbf{Y}^{\text{sim}} - \hat{\mathbf{Y}}^{\text{sim}}\right).$$

- Expected value:

$$\mathrm{E}(\boldsymbol{\delta}^{\text{sim}} + \boldsymbol{\epsilon}^{\text{Actual}}) \approx \mathrm{E}\left(\mathbf{Y}^{\text{Actual}} - \hat{\mathbf{Y}}^{\text{Actual}}\right) - \mathrm{E}\left(\mathbf{Y}^{\text{sim}} - \hat{\mathbf{Y}}^{\text{sim}}\right),$$

$$\tilde{\boldsymbol{\mu}}_{\boldsymbol{\delta}} + \mathbf{0} \approx \frac{1}{N}\sum_{n=1}^{N}\left(\mathbf{Y}_n^{\text{Actual}} - \hat{\mathbf{Y}}_n^{\text{Actual}}\right) - \frac{1}{M}\sum_{m=1}^{M}\left(\mathbf{Y}_m^{\text{sim}} - \hat{\mathbf{Y}}_m^{\text{sim}}\right),$$

where $n = 1, \ldots, N$ indexes actual retrievals, and $m = 1, \ldots, M$ indexes trials of the simulation.

National Aeronautics and
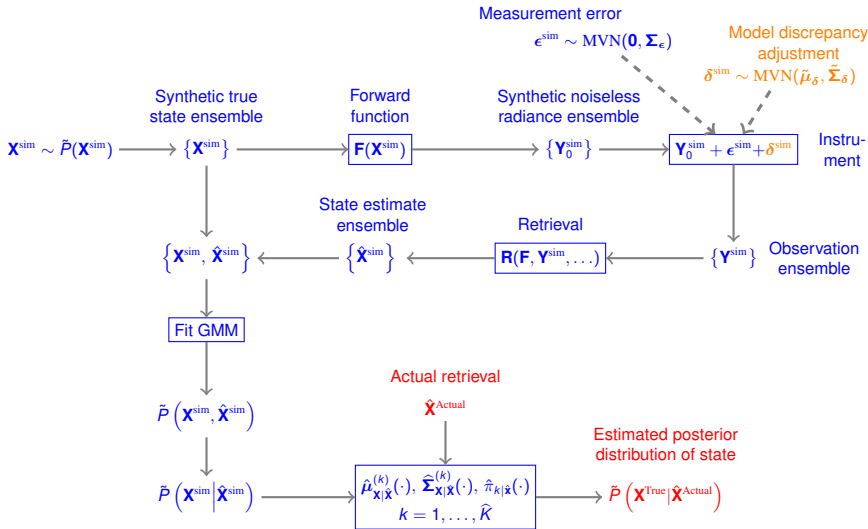Space Administration

Jet Propulsion Laboratory
California Institute of Technology
Pasadena, California

Model discrepancy

► Covariance:

$$\mathrm{cov}(\boldsymbol{\delta}^{\mathrm{sim}} + \boldsymbol{\epsilon}^{\mathrm{Actual}}) \approx \mathrm{cov}\left(\mathbf{Y}^{\mathrm{Actual}} - \hat{\mathbf{Y}}^{\mathrm{Actual}}\right) + \mathrm{cov}\left(\mathbf{Y}^{\mathrm{sim}} - \hat{\mathbf{Y}}^{\mathrm{sim}}\right)$$

$$- 2\,\mathrm{cov}\left(\mathbf{Y} - \hat{\mathbf{Y}}, \mathbf{Y}^{\mathrm{sim}} - \hat{\mathbf{Y}}^{\mathrm{sim}}\right)$$

$$\tilde{\boldsymbol{\Sigma}}_{\boldsymbol{\delta}} = \mathrm{cov}\left(\boldsymbol{\delta}^{\mathrm{sim}}\right) \leq \mathrm{cov}\left(\mathbf{Y}^{\mathrm{Actual}} - \hat{\mathbf{Y}}^{\mathrm{Actual}}\right) + \mathrm{cov}\left(\mathbf{Y}^{\mathrm{sim}} - \hat{\mathbf{Y}}^{\mathrm{sim}}\right) - \mathrm{cov}(\boldsymbol{\epsilon}^{\mathrm{Actual}}),$$

$$\approx \widehat{\mathrm{cov}}\left(\mathbf{Y}^{\mathrm{Actual}} - \hat{\mathbf{Y}}^{\mathrm{Actual}}\right) + \widehat{\mathrm{cov}}\left(\mathbf{Y}^{\mathrm{sim}} - \hat{\mathbf{Y}}^{\mathrm{sim}}\right) - \mathrm{cov}(\boldsymbol{\epsilon}^{\mathrm{Actual}}),$$

assuming $\mathrm{cov}\left(\mathbf{Y}^{\mathrm{Actual}} - \hat{\mathbf{Y}}^{\mathrm{Actual}}, \mathbf{Y}^{\mathrm{sim}} - \hat{\mathbf{Y}}^{\mathrm{sim}}\right) \geq 0$, and $\boldsymbol{\epsilon}$ and $\boldsymbol{\delta}^{\mathrm{sim}}$ are independent.

National Aeronautics and
Space Administration

**Jet Propulsion Laboratory**
California Institute of Technology
Pasadena, California



Measurement error
$\epsilon^{\mathrm{sim}} \sim \mathrm{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_\epsilon)$

Model discrepancy
adjustment
$\delta^{\mathrm{sim}} \sim \mathrm{MVN}(\tilde{\boldsymbol{\mu}}_\delta, \tilde{\boldsymbol{\Sigma}}_\delta)$

Synthetic true
state ensemble

Forward
function

Synthetic noiseless
radiance ensemble

$\mathbf{X}^{\mathrm{sim}} \sim \tilde{P}(\mathbf{X}^{\mathrm{sim}})$ $\longrightarrow$ $\{\mathbf{X}^{\mathrm{sim}}\}$ $\longrightarrow$ $\boxed{\mathbf{F}(\mathbf{X}^{\mathrm{sim}})}$ $\longrightarrow$ $\{\mathbf{Y}_0^{\mathrm{sim}}\}$ $\longrightarrow$ $\boxed{\mathbf{Y}_0^{\mathrm{sim}} + \epsilon^{\mathrm{sim}} + \delta^{\mathrm{sim}}}$ Instru-
ment

State estimate
ensemble

Retrieval

$\{\mathbf{X}^{\mathrm{sim}}, \hat{\mathbf{X}}^{\mathrm{sim}}\}$ $\longleftarrow$ $\{\hat{\mathbf{X}}^{\mathrm{sim}}\}$ $\longleftarrow$ $\boxed{\mathbf{R}(\mathbf{F}, \mathbf{Y}^{\mathrm{sim}}, \ldots)}$ $\longleftarrow$ $\{\mathbf{Y}^{\mathrm{sim}}\}$ Observation
ensemble

$\boxed{\text{Fit GMM}}$

$\tilde{P}\left(\mathbf{X}^{\mathrm{sim}}, \hat{\mathbf{X}}^{\mathrm{sim}}\right)$

Actual retrieval

$\hat{\mathbf{X}}^{\mathrm{Actual}}$

Estimated posterior
distribution of state

$\tilde{P}\left(\mathbf{X}^{\mathrm{sim}} \big| \hat{\mathbf{X}}^{\mathrm{sim}}\right)$ $\longrightarrow$ $\boxed{\begin{array}{c} \hat{\mu}_{\mathbf{X}|\hat{\mathbf{X}}}^{(k)}(\cdot), \ \hat{\boldsymbol{\Sigma}}_{\mathbf{X}|\hat{\mathbf{X}}}^{(k)}(\cdot), \ \hat{\pi}_{k|\hat{\mathbf{x}}}(\cdot) \\ k = 1, \ldots, \hat{K} \end{array}}$ $\longrightarrow$ $\tilde{P}\left(\mathbf{X}^{\mathrm{True}} \big| \hat{\mathbf{X}}^{\mathrm{Actual}}\right)$

National Aeronautics and
Space Administration

**Jet Propulsion Laboratory**
California Institute of Technology
Pasadena, California

Summary/discussion

▶ UQ community traditionally focusses on uncertainties of deterministic models' output.

▶ Stat community traditionally focusses on building statistical models, which carry uncertainty with them, but do not explicitly encode mechanistic knowledge.

▶ Remote sensing is a good example of a problem that combines elements of both.

▶ Other examples?

National Aeronautics and
Space Administration

**Jet Propulsion Laboratory**
California Institute of Technology
Pasadena, California

If the UQ community (as presently constituted) is to expand towards more statistics and statisticians, what parts of that audience should we target?

- ▶ design and analysis of computer experiments, and experimental design in general (ASA's UQ Interest Group)

- ▶ spatial/spatio-temporal statistics uses GP's and other models with spatial location/time as inputs; leverage this in more general UQ settings (e.g., emulators)

- ▶ machine learning is often inherently statistical but uncertainty not emphasized- could we do more?

- ▶ inverse problems community that intersects with UQ is not well-represented in mainstream statistics- another opportunity?

National Aeronautics and
Space Administration

**Jet Propulsion Laboratory**
California Institute of Technology
Pasadena, California

► Mainstream statisticians contemplate a wider range of applications in which computational models are not necessarily the focus, and exist along side data collected for other purposes.

► Expanding UQ territory will be accomplished by young researchers willing to think "outside the box". It will require a cultural shift.

National Aeronautics and
Space Administration

**Jet Propulsion Laboratory**
California Institute of Technology
Pasadena, California