

Computational Advances in Data-Consistent Inversion: Measure-Theoretic Methods for Improving Predictions

Michael Pilosov

Advisor: Troy Butler

University of Colorado Denver

October 29, 2020



Department of Mathematical
& Statistical Sciences

UNIVERSITY OF COLORADO DENVER

The one where we describe why any of this matters.

Broad Goals of Uncertainty Quantification:

- Make inferences and predictions
- Quantify and reduce uncertainty (aleatoric, epistemic)
- Be *accurate* and *precise*
- Design “efficient” experiments
- Collect and use data “intelligently”



The one where we define the letters we use and what they mean.

- State variable: u (e.g. heat, energy, pressure, deflection)
- Parameters: λ (e.g. source term, diffusion, boundary data)
- Deterministic model: $\mathcal{M}(u, \lambda) = 0,$

$$\mathcal{M} : \lambda \rightarrow u(\lambda)$$

- Quantity of Interest map (**QoI**) - at least pcw differentiable

- » Functional of the solution

$$q : u(\lambda) \rightarrow \mathbb{R}$$

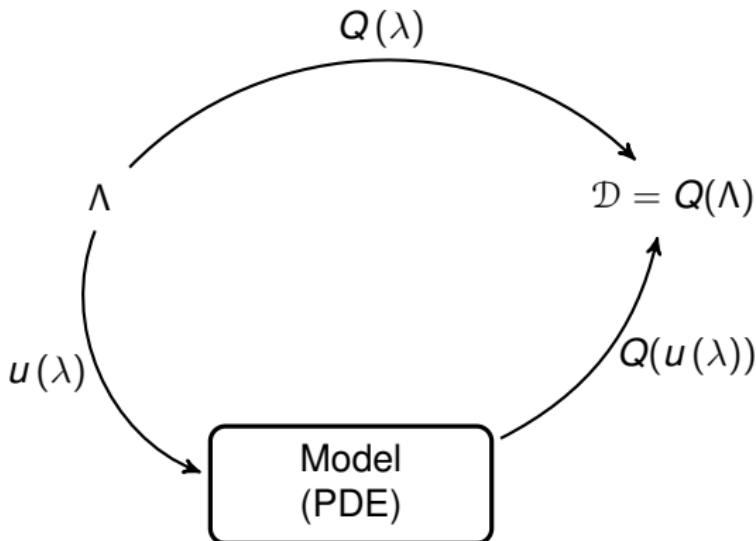
- » Can be vector valued

$$Q = \begin{bmatrix} q_1 \\ q_2 \\ \vdots \\ q_d \end{bmatrix}$$

- » $Q(\lambda) := Q(u(\lambda))$



The one where we illustrate how a QoI map relates inputs to outputs.



Defining the Quantity of Interest Map

Definition (Stochastic Forward Problem (SFP))

Given a probability measure \mathbb{P}_Λ on $(\Lambda, \mathcal{B}_\Lambda)$, and QoI map Q , the *stochastic forward problem* is to determine a measure, $\mathbb{P}_{\mathcal{D}}$, on $(\mathcal{D}, \mathcal{B}_{\mathcal{D}})$ that satisfies

$$\mathbb{P}_{\mathcal{D}}(E) = \mathbb{P}_\Lambda(Q^{-1}(E)), \quad \forall E \in \mathcal{B}_{\mathcal{D}}. \quad (1.1)$$



Definition (Stochastic Inverse Problem (SIP))

Given a probability measure, $\mathbb{P}_{\mathcal{D}}$, on $(\mathcal{D}, \mathcal{B}_{\mathcal{D}})$ the *stochastic inverse problem* is to determine a probability measure, \mathbb{P}_{Λ} , on $(\Lambda, \mathcal{B}_{\Lambda})$ satisfying

$$\mathbb{P}_{\Lambda}(Q^{-1}(E)) = \mathbb{P}_{\mathcal{D}}(E), \quad \forall E \in \mathcal{B}_{\mathcal{D}}. \quad (1.2)$$

The above is known as the *consistency condition*.

Definition (Observed Distribution)

When the measure $\mathbb{P}_{\mathcal{D}}$ in (1.2) quantifies the characterization of uncertainty in the QoI data, it is referred to as the *observed measure*, \mathbb{P}_{ob} .

If a dominating measure $\mu_{\mathcal{D}}$ exists on $(\mathcal{D}, \mathcal{B}_{\mathcal{D}})$, the *observed density* π_{ob} is given by the Radon-Nikodym derivative of \mathbb{P}_{ob} with respect to the measure $\mu_{\mathcal{D}}$.



Definition (Consistent Solution)

Any probability measure \mathbb{P}_Λ satisfying (1.2) is referred to as a *consistent solution* to the inverse problem, and (1.2).

If \mathbb{P}_Λ or $\mathbb{P}_\mathcal{D}$ absolutely continuous w.r.t μ_Λ or $\mu_\mathcal{D}$, resp, then we write

$$\pi_\Lambda := \frac{d\mathbb{P}_\Lambda}{d\mu_\Lambda} \text{ or } \pi_\mathcal{D} := \frac{d\mathbb{P}_\mathcal{D}}{d\mu_\mathcal{D}}$$

to denote the Radon-Nikodym derivatives of \mathbb{P}_Λ and $\mathbb{P}_\mathcal{D}$, resp.

In such a case, we can rewrite (1.1) and (1.2) using these pdfs:

$$\mathbb{P}_\Lambda(Q^{-1}(E)) = \int_{Q^{-1}(E)} \pi_\Lambda(\lambda) d\mu_\Lambda = \int_E \pi_\mathcal{D}(Q(\lambda)) d\mu_\mathcal{D} = \mathbb{P}_\mathcal{D}(E)$$



Definition (Initial Distribution)

When \mathbb{P}_Λ in (1.1) quantifies the characterization of uncertainty in parameter variability before observations on QoI are taken into account, it is referred to as the *initial measure* \mathbb{P}_{in} .

If a dominating measure μ_Λ exists on $(\Lambda, \mathcal{B}_\Lambda)$, the *initial distribution* π_{in} is given by the Radon-Nikodym derivative of \mathbb{P}_{in} w.r.t the measure μ_Λ .

Definition (Predicted Distribution)

The *predicted distribution* (or density) is the push-forward density of π_{in} under the map Q , and is denoted as π_{pr} .

Given as the Radon-Nikodym derivative (w.r.t $\mu_{\mathcal{D}}$) of the pushforward measure

$$\mathbb{P}_{pr}(E) = \mathbb{P}_{in}(Q^{-1}(E)), \forall E \in \mathcal{B}_{\mathcal{D}}. \quad (1.3)$$



The one where we define the solution to the SIP.

We now have all of the definitions required to summarize the density-based solution to the SIP, known as the *updated density* as:

$$\pi_{\text{up}}(\lambda) := \pi_{\text{in}}(\lambda) \frac{\pi_{\text{ob}}(Q(\lambda))}{\pi_{\text{pr}}(Q(\lambda))}. \quad (1.4)$$



Practical Considerations

- We approximate π_{pr} using density estimation on forward propagation of samples from π_{in}
- May evaluate π_{up} directly for any sample of Λ (one model solve)
- Accuracy of the computed updated density is proportional to accuracy of approximation of the predicted density
- We (currently) use Gaussian KDE
 - » Let D be the dimension of \mathcal{D}
 - » Let N be the number of samples from π_{in} propagated through Q
 - » Converges at a rate of $\mathcal{O}(N^{-4/(4+D)})$ in mean-squared error
 - » Converges at a rate of $\mathcal{O}(N^{-2/(4+D)})$ in L^1 -error



The one where we distinguish ourselves from the Bayesian Inverse Problem.

Bayesian approach: modeling epistemic uncertainties in data on a QoI obtained from a true, but unknown, parameter value, λ^\dagger .

Definition (Deterministic Forward Problem (DFP))

Given a space Λ , and QoI map Q , the *deterministic forward problem* is to determine the values, $q \in \mathcal{D}$ that satisfy

$$q = Q(\lambda), \forall \lambda \in \Lambda. \quad (1.5)$$



The one where we distinguish ourselves from the Bayesian Inverse Problem.

Definition (Deterministic Inverse Problem (DIP) Under Uncertainty)

Given a noisy datum (or data-vector) $d = q + \xi$, $q \in \mathcal{D}$, the *deterministic inverse problem* is to determine the parameter $\lambda \in \Lambda$ which minimizes

$$\|Q(\lambda) - d\| \quad (1.6)$$

where ξ is a random variable (or vector) drawn from a distribution characterizing the uncertainty in observations due to measurement errors.

In the above definition, ξ is some unobservable perturbation to the true output, arising from epistemic uncertainty (e.g. the precision of available measurement equipment).



The one where we distinguish ourselves from the Bayesian Inverse Problem.

The *posterior* is a conditional density, denoted by $\pi_{\text{post}}(\lambda | d)$, proportional to the product of the prior and data-likelihood function [3, 2, 1, 4]:

$$\pi_{\text{post}}(\lambda) := \pi_{\text{prior}}(\lambda) \frac{L_{\mathcal{D}}(q|\lambda)}{C}, \quad (1.7)$$

where we emphasize the use of π_{post} to distinguish the *posterior* from the updated density π_{up} in (1.4).

evidence term C ensures the posterior density integrates to one; given by

$$C = \int_{\Lambda} \pi_{\text{prior}}(\lambda) L_{\mathcal{D}}(q|\lambda) d\lambda.$$



The one where we provide an illustrative example.

- Suppose $\Lambda = [-1, 1] \subset \mathbb{R}$ and $Q(\lambda) = \lambda^5$ so that $\mathcal{D} = [-1, 1]$
- $\pi_{\text{in}} \sim \mathcal{U}([-1, 1])$ and $\pi_{\text{ob}} \sim N(0.25, 0.1^2)$
- $d \in \mathcal{D}$ with $d = Q(\lambda^\dagger) + \xi$ where $\xi \sim N(0, 0.1^2)$
- We then construct $\pi_{\text{post}}(\lambda | d)$ for this example assuming a uniform prior (to match the initial density) with an assumed observed value of $d = 0.25$ so that the data-likelihood function matches the observed density.



The one where we provide an illustrative example.

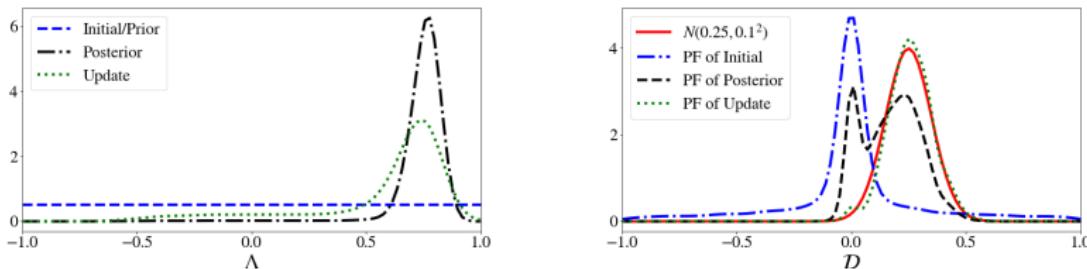


Figure: (Left) The initial/prior PDF π_{in} (blue solid curve), updated PDF π_{up} (black dashed curve), and posterior PDF π_{post} (green dashed-dotted curve) on Λ . (Right) The push-forward (PF) of the initial/prior PDF π_{pr} (blue solid curve), observed/likelihood PDF (red solid curve), PF of the updated PDF π_{up} (black dashed curve), and the PF of the posterior PDF π_{post} (green dashed-dotted curve) for the QoI.

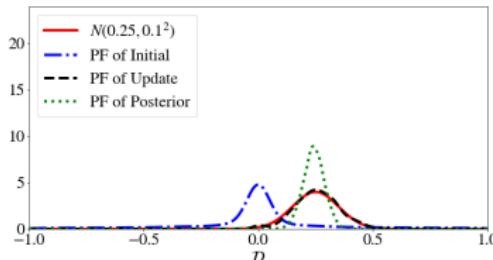
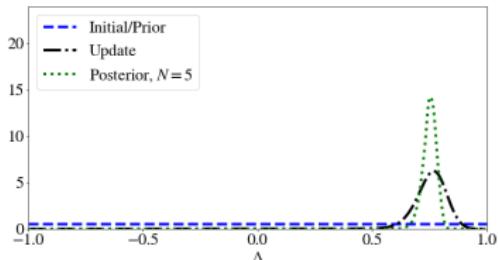
The one where we provide an illustrative example.

What happens as we collect more data?

One approach:

SIP: Use N to estimate mean of observed

DIP: likelihood function incorporates more data



SIP and DIP solutions for varying N for comparison.

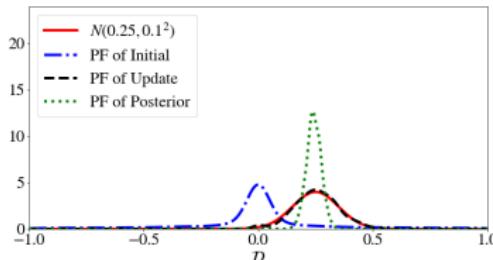
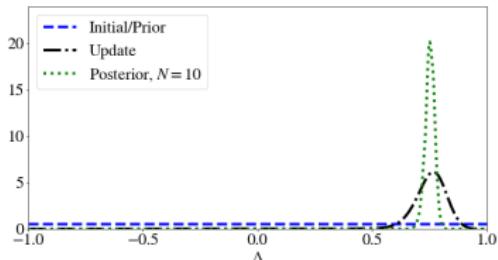
The one where we provide an illustrative example.

What happens as we collect more data?

One approach:

SIP: Use N to estimate mean of observed

DIP: likelihood function incorporates more data



SIP and DIP solutions for varying N for comparison.

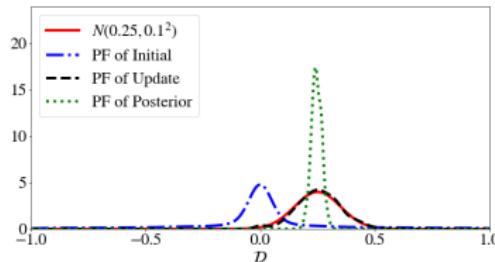
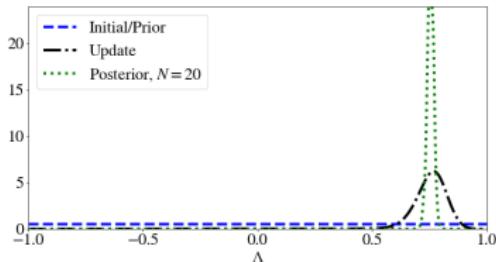
The one where we provide an illustrative example.

What happens as we collect more data?

One approach:

SIP: Use N to estimate mean of observed

DIP: likelihood function incorporates more data



SIP and DIP solutions for varying N for comparison.

We formally define the maximal updated density (MUD) point as

$$\lambda^{\text{MUD}} := \arg \max \pi_{\text{up}}(\lambda). \quad (2.1)$$

We motivate the use of the MUD point as an alternative to the MAP point for parameter estimation problems.



Let $\|\mathbf{x}\|_C^2 := (\mathbf{x}, \mathbf{x})_C = \mathbf{x}^\top C \mathbf{x}$.

Inverse covariances associated with non-degenerative multivariate Gaussian distributions will play the role of C .

Suppose that the initial and prior densities are both given by the same $\mathcal{N}(\lambda_0, \Sigma_{\text{init}})$ distribution.

Additionally, suppose the map Q is linear and that the data-likelihood and observed densities are both given by the same $\mathcal{N}(\mathbf{y}, \Sigma_{\text{obs}})$ distribution.

The linearity of Q implies that $Q(\lambda) = A\lambda$ for some $A \in \mathbb{R}^{d \times p}$, and that the predicted density follows a $\mathcal{N}(Q(\lambda_0), \Sigma_{\text{pred}})$ distribution where

$$\Sigma_{\text{pred}} := A\Sigma_{\text{init}}A^\top. \quad (2.2)$$



The one with the regularization equations.

$\pi_{\text{up}}(\lambda) = \pi_{\text{in}}(\lambda) \frac{\pi_{\text{ob}}(Q(\lambda))}{\pi_{\text{pr}}(Q(\lambda))}$	$\pi_{\text{post}}(\lambda d) = \frac{\pi_{\text{prior}}(\lambda) \pi_{\text{like}}(d \lambda)}{\int_{\Lambda} \pi_{\text{like}}(d \lambda) \pi_{\text{prior}}(\lambda) d\mu_{\Lambda}}$
Tikhonov	$T(\lambda) := \ Q(\lambda) - \mathbf{y}\ _{\Sigma_{\text{obs}}^{-1}}^2 + \ \lambda - \lambda_0\ _{\Sigma_{\text{init}}^{-1}}^2$
Data-Consistent	$J(\lambda) := T(\lambda) - \ Q(\lambda) - Q(\lambda_0)\ _{\Sigma_{\text{pred}}^{-1}}^2$

Table: The λ which minimizes these functionals also maximizes the updated PDF (left) and the Bayesian posterior PDF (right).

$T(\lambda)$ is the typical functional often associated with Tikhonov regularization. The $J(\lambda)$ has an additional term subtracted from $T(\lambda)$ coming from the predicted density that serves as “unregularization” in data-informed directions.

The one where an example highlights a key difference.

Consider a linear QoI map is defined by $A = [\begin{array}{cc} 1 & 1 \end{array}]$.
2-D input, 1-D output \implies rank-deficient

Parameters in the initial and observed densities are given by

$$\lambda_0 = [\begin{array}{cc} 0.25 & 0.25 \end{array}]^\top,$$
$$\Sigma_{\text{init}} = [\begin{array}{cc} 1 & -0.25 \\ -0.25 & 0.5 \end{array}],$$
$$\mathbf{y} = 1, \text{ and } \Sigma_{\text{obs}} = [\begin{array}{c} 0.25 \end{array}]$$



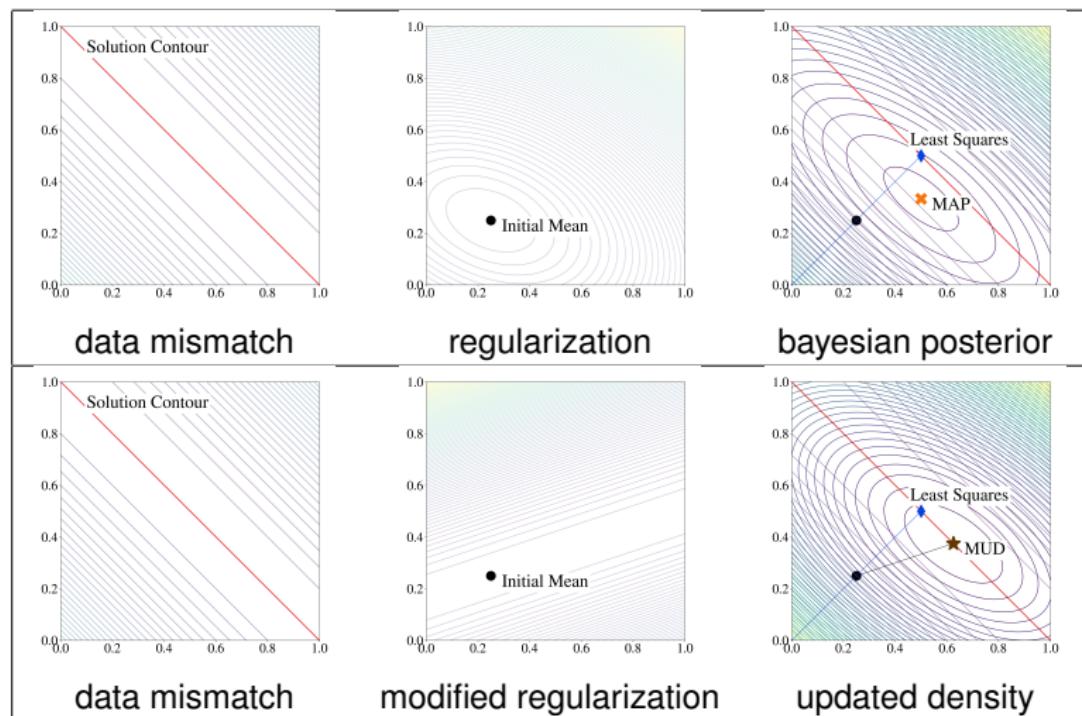


Figure: Gaussian data mismatch for a 2-to-1 linear map (left plots). Gaussian initial/prior induce different regularization terms (middle plots), which leads to different optimization functions (right plots) and parameter estimates.

The posterior covariance is formally given by

$$\Sigma_{\text{post}} := (A^\top \Sigma_{\text{obs}}^{-1} A + \Sigma_{\text{init}}^{-1})^{-1}. \quad (2.3)$$

Applying Woodbury identity and (2.2), we rewrite the posterior covariance:

$$\Sigma_{\text{post}} = \Sigma_{\text{init}} - \Sigma_{\text{init}} A^\top [\Sigma_{\text{pred}} + \Sigma_{\text{obs}}]^{-1} A \Sigma_{\text{init}} \quad (2.4)$$

Can now interpret Σ_{post} as a rank d correction (or update) of Σ_{init} .
 $\Sigma_{\text{pred}} + \Sigma_{\text{obs}}$ is invertible because it is the sum of two s.p.d matrices.
We rewrite the closed form expression for the MAP point given in [5] as

$$\lambda^{\text{MAP}} = \lambda_0 + \Sigma_{\text{post}} A^\top \Sigma_{\text{obs}}^{-1} (\mathbf{y} - b - A\lambda_0). \quad (2.5)$$



We define

$$R := \Sigma_{\text{init}}^{-1} - A^\top \Sigma_{\text{pred}}^{-1} A. \quad (2.6)$$

Using this R , rewrite $J(\lambda)$ as

$$J(\lambda) := \|\mathbf{y} - Q(\lambda)\|_{\Sigma_{\text{obs}}^{-1}}^2 + \|\lambda - \lambda_0\|_R^2. \quad (2.7)$$

In this form, we identify R as the *effective regularization* in $J(\lambda)$ due to the formulation in the data-consistent framework.

$$\Sigma_{\text{up}} := \left(A^\top \Sigma_{\text{obs}}^{-1} A + R \right)^{-1}. \quad (2.8)$$

Since R is not invertible, Woodbury's identity cannot be applied (yet).



We derive using several identities

$$\Sigma_{\text{up}} = \Sigma_{\text{init}} - \Sigma_{\text{init}} A^\top \Sigma_{\text{pred}}^{-1} [\Sigma_{\text{pred}} - \Sigma_{\text{obs}}] \Sigma_{\text{pred}}^{-1} A \Sigma_{\text{init}}. \quad (2.9)$$

Substitute Σ_{up} for Σ_{post} in (2.5) to write the point that minimizes J as:

$$\lambda^{\text{MUD}} = \lambda_0 + \Sigma_{\text{up}} A^\top \Sigma_{\text{obs}}^{-1} (\mathbf{y} - b - A\lambda_0). \quad (2.10)$$

Substituting (2.9) into (2.10) and simplifying, we have

$$\lambda^{\text{MUD}} = \lambda_0 + \Sigma_{\text{init}} A^\top \Sigma_{\text{pred}}^{-1} (\mathbf{y} - b - A\lambda_0). \quad (2.11)$$



Predictability assumption:
smallest predicted > largest observed (eigenvalues of covariances)

Theorem

Suppose $Q(\lambda) = A\lambda + b$ for some full rank $A \in \mathbb{R}^{d \times p}$ with $d \leq p$ and $b \in \mathbb{R}^d$. If $\pi_{\text{in}} \sim N(\lambda_0, \Sigma_{\text{init}})$, $\pi_{\text{ob}} \sim N(\mathbf{y}, \Sigma_{\text{obs}})$, and the predictability assumption holds, then

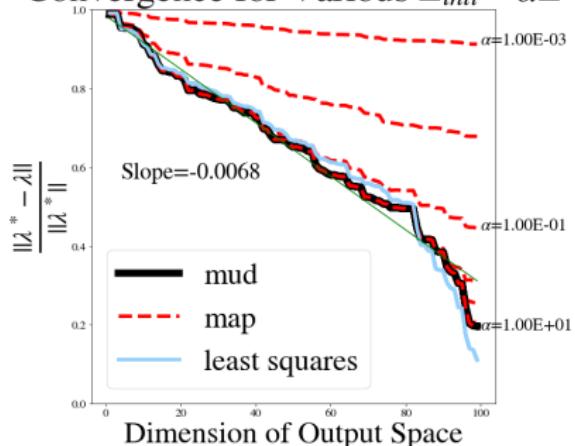
- (a) There exists a unique parameter, denoted by λ^{MUD} , that maximizes π_{up} .
- (b) $Q(\lambda^{\text{MUD}}) = \mathbf{y}$.
- (c) If $d = p$, λ^{MUD} is given by A^{-1} . If $d < p$, λ^{MUD} is given by (2.11) and the covariance associated with this point is given by (2.9).



The one where we show how rank and dimension impact our solutions.

Example: scaling random diagonal initial covariances

Convergence for Various $\Sigma_{init} = \alpha \Sigma$



Convergence for Various $\Sigma_{init} = \alpha \Sigma$

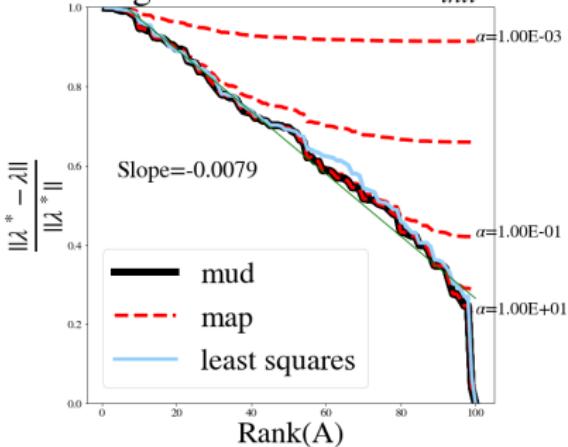


Figure: Relative errors between λ^\dagger and (i) the least squares solution obtained through numpy's linalg.pinv module, (ii) the closed-form solution for the MUD point given in Eq (2.11), and (iii) the MAP point. (Left): Error for increasing dimensions of D for A taken to be a Gaussian Random Map. (Right): Error for increasing row-rank of A , generated with Gaussian vectors and a SVD.

The one where we leverage this framework for general streams of data.

Suppose $\exists d$ measurement devices generating repeated noisy data.

For each $1 \leq j \leq d$, denote by $\mathcal{M}_j(\lambda^\dagger)$ the j th measurement device.
 N_j is number of noisy data obtained for $\mathcal{M}_j(\lambda^\dagger)$.

$d_{j,i}$ is the i th noisy datum for the j th measurement, where $1 \leq i \leq N_j$.

Assume an unbiased additive error model for the measurement noise,
with independent identically distributed (i.i.d.) Gaussian errors so that

$$d_{j,i} = M_j(\lambda^*) + \xi_i, \quad \xi_i \sim N(0, \sigma_j^2), \quad 1 \leq i \leq N_j. \quad (2.12)$$

We now construct a d -dimensional vector-valued map from data
obtained on the d measurement devices.



The one with the Weighted Mean Error (WME) map.

The weighted mean error (WME) map, denoted by $Q_{\text{WME}}(\lambda)$ has j th component, denoted by $Q_{\text{WME},j}(\lambda)$, given by

$$Q_{\text{WME},j}(\lambda) := \frac{1}{\sqrt{N_j}} \sum_{i=1}^{N_j} \frac{M_j(\lambda) - d_{j,i}}{\sigma_j}. \quad (2.13)$$

$Q_{\text{WME},j}(\lambda^\dagger)$ is the sample avg of N_j random draws from an i.i.d. $N(0, N_j)$. By assumption, the observed data are generated according to the fixed true physical parameter vector given by λ^\dagger in (2.12).

Subsequently, each component of $Q_{\text{WME}}(\lambda^\dagger)$ is a random draw from an $N(0, 1)$ distribution.

Therefore, with this choice of data-defined QoI map, we specify π_{ob} as a $N(\mathbf{0}_{d \times 1}, \mathbf{I}_{d \times d})$ distribution.

The one where measurements impact the predictability assumption.

The j th diagonal component of the predicted covariance matrix is given by the predicted variance associated with using the scalar-valued

$$Q_{\text{WME},j}.$$

Then, the associated predicted variance is given by

$$\frac{N_j}{\sigma_j^2} M_j \Sigma_{\text{init}} M_j^\top \quad (2.14)$$

Since Σ_{init} is assumed to be non-degenerative and M_j is a non-trivial row vector, this predicted variance grows linearly with N_j .

In other words, the j th diagonal component of the predicted covariance has the form $\beta_j N_j$ for some $\beta_j > 0$.



Let $N_{\min,j}$ denote the minimum N_j for $1 \leq j \leq N$ necessary to make the j th diagonal components sufficiently large so that the smallest eigenvalue of the predicted covariance is larger than 1.

The following result is now an immediate consequence of Theorem 2.1:

Corollary

If $\pi_{\text{in}} \sim N(\lambda_0, \Sigma_{\text{init}})$ and data are obtained for d linearly independent measurements on Λ with an additive noise model with i.i.d. Gaussian noise for each measurement, then **there exists a minimum number of data points obtained for each of the measurements such that there exists a unique λ^{MUD} and $Q_{\text{WME}}(\lambda^{\text{MUD}}) = 0$.**

The one where we show that our approach works even when some assumptions are violated.



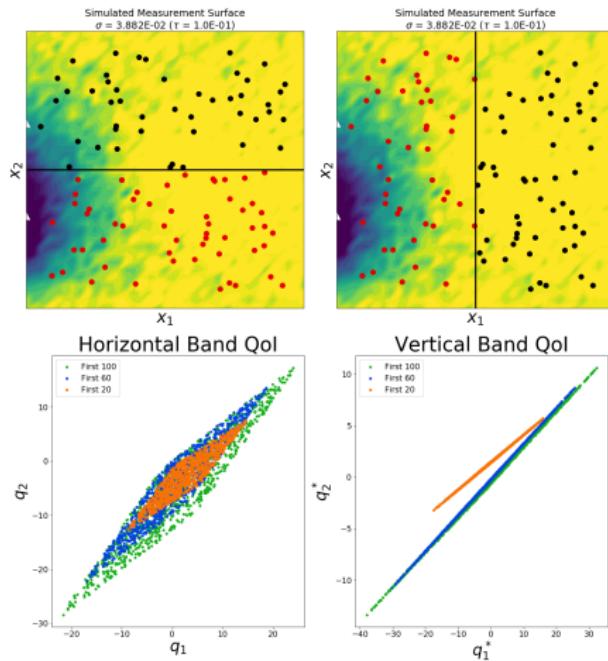


Figure: $N = 1000$ parameter evaluations for both methods of partitioning Ω .

The one with the small problems in many batches.

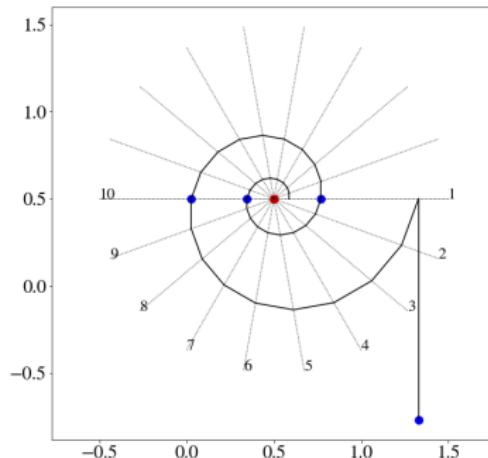
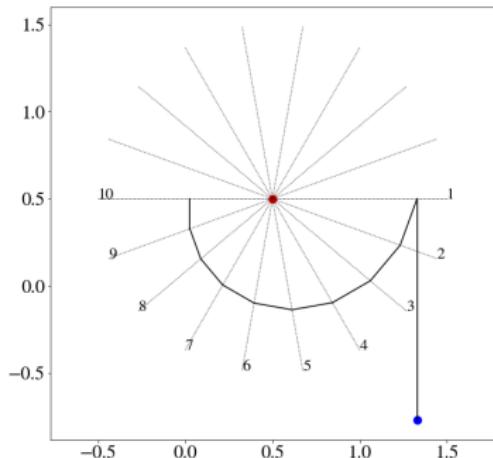


Figure: Dotted lines show the solution contours for each row of the operator A . (Left): First epoch for iterating through 10 QoI. (Right): Three more epochs allows our estimate to get much closer to the true value.

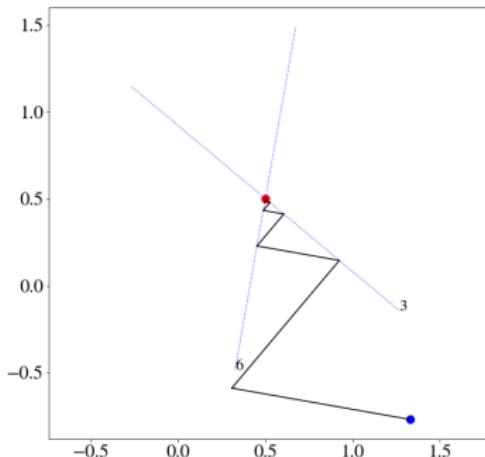
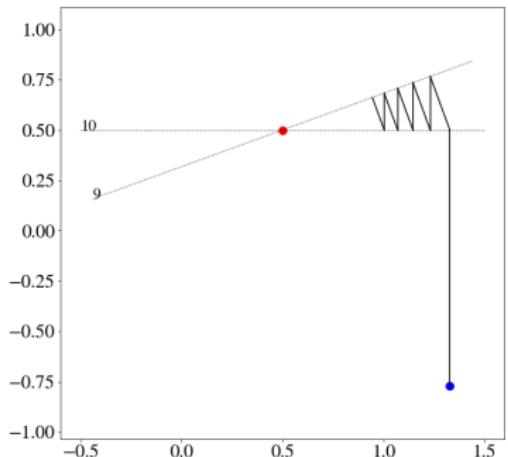


Figure: Iterating through five epochs of two QoI, each formed by picking two of the ten available rows of A at random. The random directions chosen on the left exhibit more redundancy than those on the right, so the same amount of iteration results in less accuracy.

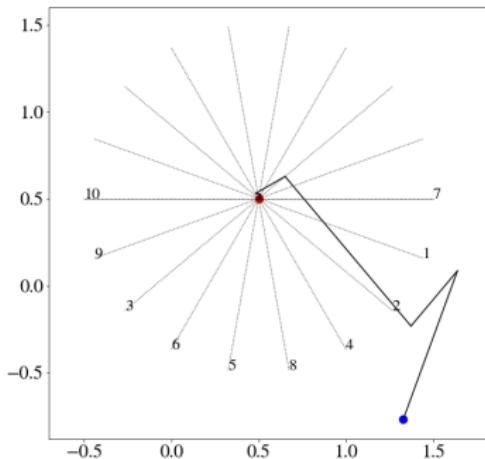
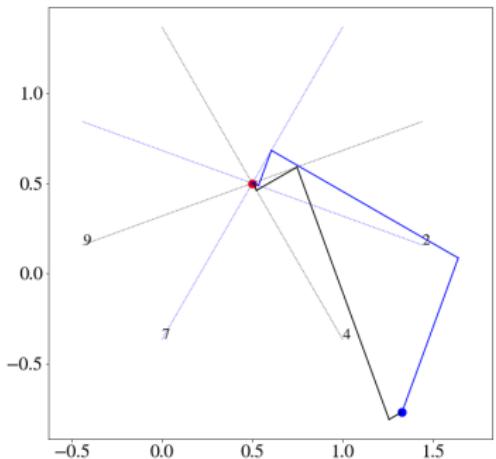


Figure: (Left): Subsets of available QoI components can be chosen to exhibit minimal redundancy and lead to expedited convergence. (Right): Random components of the QoI map used for each iterative step. This leads to an overall similar level of precision in this example, without the need to use gradients.

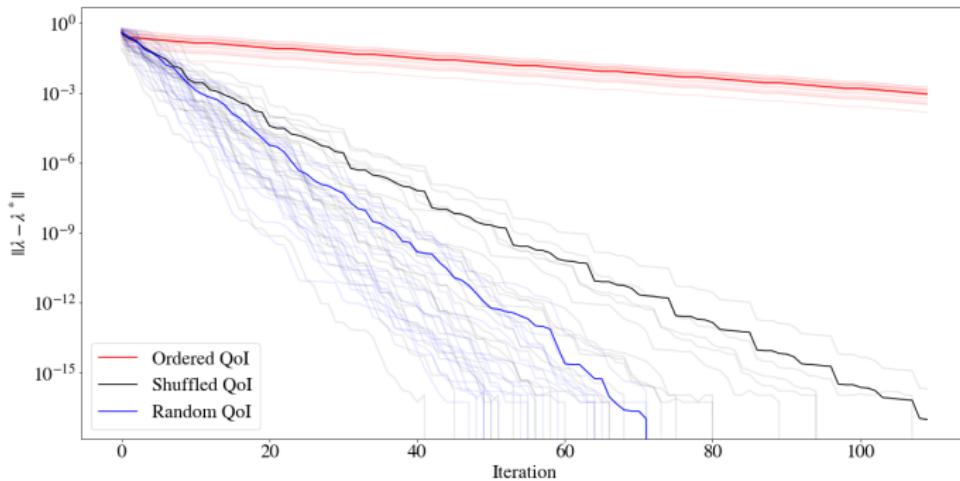


Figure: 20 initial means are chosen and iterated on for three approaches for ordering QoI. Individual experiments are transparent and the mean error is shown as a solid line for each approach.

How do I know I can trust you?

You don't. But I enabled you to check for yourself.

- Public repository hosted on Github.com
(github.com/mathematicalmichael/thesis)
- Github Actions implements Continuous Integration / Deployment
- Each change is validated for reproducibility
- makefile for convenience (make <filename>)
 - » dissertation + presentation (L^AT_EX, themes, style files)
 - » every example, convergence result (Python)
 - » every image in every figure
- PyPi published implementation of main methods: pip install mud
- Unit tests aid in ensuring integrity of functions
- Docker guarantees software runtime (ran on x86 and arm)
docker pull mathematicalmichael/python:thesis(latex:thesis)



 M. Allmaras, W. Bangareth, J.M. Linhart, J. Polanco, F. Wang, K. Wang, J. Webster, and S. Zedler.

Estimating parameters in physical models through Bayesian inversion: A complete example.

2013.

 J.O. Berger.

Statistical Decision Theory and Bayesian Analysis.
Springer-Verlag, 1985.

 S. Myers R. Walpole, R. Myers and K. Ye.

Probability & Statistics for Engineers & Scientists.
Pearson Education, 2007.

 Ralph C. Smith.

Uncertainty Quantification: Theory, Implementation, and Applications.

Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 2013.

 Albert Tarantola.

*Inverse Problem Theory and Methods for Model Parameter
Estimation.*
siam, 2005.

