Density Estimation Framework for Model Error Quantification

Khachik Sargsyan, Xun Huan, Habib Najm

Sandia National Laboratories, Livermore, CA

SIAM UQ Lausanne, Switzerland April 5-8, 2016



Sandia National Laboratories

Sandia National Laboratories is a multi-program laboratory operated by Sandia Corporation, a wholly owned subsidiary of Lockheed Martin Corporation, for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-AC04-94AL85000.

> Further acknowledgements: R. Ghanem(USC), J. Bender(LLNL), Y. Marzouk(MIT), C. Feng(MIT), M. Eldred (SNL), C. Safta (SNL), B. Debusschere (SNL).

K. Sargsyan (ksargsy@sandia.gov)

SIAM UQ 16

Density Estimation Framework for Model Error Quantification

Khachik Sargsyan, Xun Huan, Habib Najm

Sandia National Laboratories, Livermore, CA

SIAM UQ Lausanne, Switzerland April 5-8, 2016



Sandia National Laboratories

DOE Office of Advanced Scientific Computing Research (ASCR), Scientific Discovery through Advanced Computing (SciDAC) DOE Office of Basic Energy Sciences, Division of Chemical Sciences, Geosciences, & Biosciences DOD, DARPA Enabling Quantification of Uncertainty in Physical Systems (EQUiPS) program

> Further acknowledgements: R. Ghanem(USC), J. Bender(LLNL), Y. Marzouk(MIT), C. Feng(MIT), M. Eldred (SNL), C. Safta (SNL), B. Debusschere (SNL).

K. Sargsyan (ksargsy@sandia.gov)

SIAM UQ 16

Outline

- Motivation: need to quantify and propagate model error
- Issues with current state-of-the-art
- Method
- Toy cases
- Examples
 - Chemistry
 - Atmospheric transport
 - LES computations

Main target

Model error = deviation from 'truth', or from a higher-fidelity model

- Represent and estimate the error associated with
 - Simplifying assumptions, parameterizations
 - Mathematical formulation, theoretical framework
 - Numerical discretization
- ...will be useful for
 - Model validation
 - Model comparison
 - Scientific discovery and model improvement
 - Reliable computational predictions
- Inverse modeling context
 - Given experimental or higher-fidelity model data, estimate the model error



Given noisy data – Gaussian noise

•
$$y = g_{\text{true}}(x) + \epsilon$$



- Employ Bayesian inference to fit an exponential model $y_m = f(x, \lambda)$
- Discrepancy between data and prediction presumed exclusively due to *i.i.d.* Gaussian data noise $-y = f(x, \lambda) + \epsilon_d$

Plotted:

- Posterior density on the parameters
- Preditive mean and standard deviation

K. Sargsyan (ksargsy@sandia.gov)

SIAM UQ 16



- Employ Bayesian inference to fit an exponential model $y_m = f(x, \lambda)$
- Discrepancy between data and prediction presumed exclusively due to *i.i.d.* Gaussian data noise $-y = f(x, \lambda) + \epsilon_d$
- True model g(x) dashed-red differs from fit model $f(x, \lambda)$
- Actual discrepancy includes both data and model errors



- Increasing number of data points decreases posterior and predictive uncertainty
- We are increasingly sure about predictions based on the wrong model



- Increasing number of data points decreases posterior and predictive uncertainty
- We are increasingly sure about predictions based on the wrong model



- If the model has structural uncertainty, more data leads to biased and overconfident results
- We want to quantify model-vs-truth discrepancy in a rigorous and systematic way
 - Cannot ignore model error

Data-Model-Truth



$$y_i = \underbrace{f(x_i; \lambda) + \delta(x_i)}_{\text{truth } g(x_i)} + \epsilon_i^d$$

Explicit statistical modeling of model discrepancy/error $\delta(x)$

Model Error:
$$\delta(x) \sim \operatorname{GP}(\mu(x), C(x, x'))$$
Data Error: $\epsilon_i^{\mathrm{d}} \sim \operatorname{N}(0, \sigma^2)$

Estimate model parameters λ along with those of $\delta(x)$, ϵ_i^d

Additive model discrepancy: issues for physical models

$$y_i = \underbrace{f(x_i; \lambda) + \delta(x_i)}_{\text{truth}} + \epsilon_i^d$$

- Explicit additive statistical model for model error δ(x) Kennedy-O'Hagan (2001).
- Calibrated predictive model: $f(x; \lambda) + \delta(x)$ or $f(x; \lambda)$?
- Potential violation of physical constraints
- Disambiguation of model error $\delta(x_i)$ and data error ϵ_i^d
- Calibration of model error on measured observable does not impact the quality of model predictions on other QoIs
- Physical scientists are unlikely to augment their model with a statistical model error term on select outputs

Model Error – Challenges with current methods



• Ignoring model error $\delta(x)$ leads to incorrect predictive errors

• Conventional statistical modeling (Kennedy and O'Hagan, 2001)

- makes it difficult to disambiguate model/data errors
- may violate physical constraints
- not meaningful for prediction of other Qols
- Issue is highlighted in model-to-model calibration ($\epsilon_i = 0$)
 - no a priori knowledge of the statistical structure of the discrepancy

Model error embedding: key idea

Ideally, modelers want predictive *errorbars*: inserting randomness on the outputs has issues, so...

• Cast input parameters λ as a random variable Λ

$$y_i = f(x_i; \Lambda) + \epsilon_i^d$$

Generalize parameter forms,

Random field
$$y_i = f(x_i; \Lambda(x_i)) + \epsilon_i^d$$

More generally, explore additional parameterizations,

Extra 'physics'

$$y_i = \tilde{f}(x_i; \lambda, \Theta) + \epsilon_i^{\mathrm{d}}$$

K. Sargsyan (ksargsy@sandia.gov)

Model error embedding: key idea

Cast input parameters λ as a random variable Λ

 $y_i = f(x_i; \lambda) + \delta(x_i) + \epsilon_i \longrightarrow y_i = f(x_i; \Lambda) + \epsilon_i$

Embed model error in specific submodel phenomenology

- a modified transport or constitutive law
- a modified formulation for a material property
- Allows placement of model error term in locations where key modeling assumptions and approximations are made
 - as a correction or high-order term
 - as a possible alternate phenomenology
- Naturally preserves model structure and physical constraints

Model error embedding: Bayesian formulation

 $y_i = f(x_i; \Lambda) + \epsilon_i^{\mathrm{d}}$

- Parameter estimation of λ turns into PDF estimation of Λ
- Fixed PDF form $\pi_{\Lambda}(\cdot; \alpha)$ or Polynomial Chaos $\Lambda = \sum_{k=0}^{K} \alpha_k \Psi_k(\xi)$
- Back to parameter estimation, now for $\alpha = (\alpha_0, \dots, \alpha_K)$



Model Error – Bayesian density estimation

 $y_i = f(x_i; \Lambda) + \epsilon_i$

Parametrize embedded random variable Λ:

• PDF form $\pi_{\Lambda}(\cdot; \alpha)$

• Polynomial Chaos (PC): $\Lambda = \sum_{k} \alpha_k \Psi_k(\xi)$

• Multivariate Normal (MVN):
$$\begin{cases} \Lambda_1 = \alpha_{10} + \alpha_{11}\xi_1 \\ \Lambda_2 = \alpha_{20} + \alpha_{21}\xi_1 + \alpha_{22}\xi_2 \\ \vdots \\ \Lambda_d = \alpha_{d0} + \alpha_{d1}\xi_1 + \alpha_{d2}\xi_2 + \dots + \alpha_{dd}\xi_d \end{cases}$$

- Inverse modeling context
 - Parameter estimation of $\lambda \Rightarrow \mathsf{PDF}$ estimation of $\Lambda \Rightarrow$ parameter estimation of α
 - Bayesian formulation

 $\underbrace{p(\alpha|y)}_{x} \propto \underbrace{L_y(\alpha)}_{y} \underbrace{p(\alpha)}_{y}$ Posterior Likelihood Prior

K. Sargsyan, H. Najm, and R. Ghanem, "On the Statistical Calibration of Physical Models".

International Journal for Chemical Kinetics, 47(4): pp 246–276, 2015.



K. Sargsyan, H. Najm, and R. Ghanem, "On the Statistical Calibration of Physical Models".

International Journal for Chemical Kinetics, 47(4): pp 246–276, 2015.



• Full Likelihood: $L(\alpha) = p(y|\alpha) = p(y_1, \dots, y_N|\alpha) = \pi(y)$

- Degenerate if no data noise
- Requires multivariate KDE or high-d integration
- Gaussian approximation: $L(\alpha) \propto \exp\left(-\frac{1}{2}(y-\mu(\alpha))^T \Sigma^{-1}(\alpha)(y-\mu(\alpha))\right)$
- Non-intrusive spectral projection with Polynomial Chaos relieves the expense and provides easy access to mean $\mu(\alpha)$ and covariance $\Sigma(\alpha)$

K. Sargsyan (ksargsy@sandia.gov)

K. Sargsyan, H. Najm, and R. Ghanem, "On the Statistical Calibration of Physical Models".

International Journal for Chemical Kinetics, 47(4): pp 246–276, 2015.



- Marginalized Likelihood: $L(\alpha) = p(y|\alpha) \approx \prod_{i=1}^{N} p(y_i|\alpha) = \prod_{i=1}^{N} \pi(y_i)$
 - Requires univariate KDE
 - Neglects built-in correlations
 - Gaussian approximation:

$$L(\alpha) \propto \exp\left(-\frac{1}{2}\sum_{i=1}^{N}\sum_{i=1}^{-1}(\alpha)(y_i - \mu_i(\alpha))^2\right)$$

K. Sargsyan, H. Najm, and R. Ghanem, "On the Statistical Calibration of Physical Models".

International Journal for Chemical Kinetics, 47(4): pp 246–276, 2015.



• Approximate Bayesian Computation (ABC): $L(\alpha) = \frac{1}{\epsilon} K\left(\frac{\rho(S_{\mathcal{M}}, S_{\mathcal{D}})}{\epsilon}\right)$

- Mean of $f(x_i; \Lambda)$ is "centered" on the data
- The width of the distribution of *f*(*x_i*; Λ) is consistent with the spread of the data around the nominal model prediction

$$L(\alpha) \propto \exp\left(-\frac{1}{2\epsilon^2}\sum_{i=1}^{N}\left[\left(\mu_i(\alpha) - y_i\right)^2 + \left(\sqrt{\Sigma_{ii}(\alpha)} - \gamma|\mu_i(\alpha) - y_i|\right)^2\right]\right)$$

K. Sargsyan (ksargsy@sandia.gov)

Likelihood construction: data model

• Data
$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$$

• Model $f(x; \Lambda)$

• Model parameters $\Lambda = \lambda(\alpha; \xi) = \sum_k \alpha_k \Psi_k(\xi_1, \dots, \xi_d)$

• Data generation model; to aid likelihood $p(\mathcal{D}|\hat{\alpha})$ construction

$$y_{i} = f(x_{i}, \Lambda) + \epsilon_{i}^{d} =$$

$$= f\left(x_{i}, \sum_{k} \alpha_{k} \Psi_{k}(\xi_{1}, \dots, \xi_{d})\right) + \sigma_{\mathcal{D}} \xi_{d+i} =$$

$$\stackrel{NISP}{\approx} \sum_{k} f_{ik}(\alpha) \Psi_{k}(\xi_{1}, \dots, \xi_{d}) + \sigma_{\mathcal{D}} \xi_{d+i} =$$

$$= h_{i}(\hat{\xi}; \hat{\alpha})$$

• Infer $\hat{\alpha} = (\alpha, \sigma_{\mathcal{D}})$ • Full PC germ $\hat{\xi} = (\underbrace{\xi_1, \dots, \xi_d}_{\text{Model error}}, \underbrace{\xi_{d+1}, \dots, \xi_{d+N}}_{\text{Data noise}})$

$$y_i = \underbrace{\sum_{k} f_{ik}(\alpha) \Psi_k(\xi_1, \dots, \xi_d) + \sigma_{\mathcal{D}} \xi_{d+i}}_{h_i(\hat{\xi}; \hat{\alpha})}$$

Note: for each $\hat{\alpha}$,

the data model $h(\hat{\xi}; \hat{\alpha})$ is a multivariate random variable with cheap sampling, and easily accessible mean $\mu(\hat{\alpha})$ and covariance $\Sigma(\hat{\alpha})$

$$y_i = \underbrace{\sum_{k} f_{ik}(\alpha) \Psi_k(\xi_1, \dots, \xi_d) + \sigma_{\mathcal{D}} \xi_{d+i}}_{h_i(\hat{\xi}; \hat{\alpha})}$$

Note: for each $\hat{\alpha}$,

the data model $h(\hat{\xi}; \hat{\alpha})$ is a multivariate random variable with cheap sampling, and easily accessible mean $\mu(\hat{\alpha})$ and covariance $\Sigma(\hat{\alpha})$

- Full Likelihood: $L(\hat{\alpha}) = p(\mathcal{D}|\hat{\alpha}) = p(y_1, \dots, y_N|\hat{\alpha}) = \pi_{\mathbf{h}}(\mathbf{y})$
 - Degenerate if no data noise
 - Requires multivariate KDE
 - Gaussian approximation: $L(\hat{\alpha}) \propto \exp\left(-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu}(\hat{\alpha}))^T \Sigma^{-1}(\mathbf{y} - \boldsymbol{\mu}(\hat{\alpha}))\right)$

$$y_i = \underbrace{\sum_{k} f_{ik}(\alpha) \Psi_k(\xi_1, \dots, \xi_d) + \sigma_{\mathcal{D}} \xi_{d+i}}_{h_i(\hat{\xi}; \hat{\alpha})}$$

Note: for each $\hat{\alpha}$,

the data model $h(\hat{\xi}; \hat{\alpha})$ is a multivariate random variable with cheap sampling, and easily accessible mean $\mu(\hat{\alpha})$ and covariance $\Sigma(\hat{\alpha})$

Marginalized Likelihood:

$$L(\hat{\alpha}) = p(D|\hat{\alpha}) = \prod_{i=1}^{N} p(y_i|\hat{\alpha}) = \prod_{i=1}^{N} \pi_{h_i}(y_i)$$

- Requires univariate KDE
- Neglects built-in correlations
- Gaussian approximation: $L(\hat{\alpha}) \propto \exp\left(-\frac{1}{2}\sum_{i=1}^{N}\sum_{i=1}^{-2}(y_i \mu_i(\hat{\alpha}))^2\right)$

$$y_i = \underbrace{\sum_{k} f_{ik}(\alpha) \Psi_k(\xi_1, \dots, \xi_d) + \sigma_{\mathcal{D}} \xi_{d+i}}_{h_i(\hat{\xi}; \hat{\alpha})}$$

Note: for each $\hat{\alpha}$,

the data model $h(\hat{\xi}; \hat{\alpha})$ is a multivariate random variable with cheap sampling, and easily accessible mean $\mu(\hat{\alpha})$ and covariance $\Sigma(\hat{\alpha})$

- Approximate Bayesian Computation (ABC): $L(\hat{\alpha}) = \frac{1}{\epsilon} K\left(\frac{\rho(\mathcal{S}_{\mathcal{M}}, \mathcal{S}_{\mathcal{D}})}{\epsilon}\right)$
 - $p(y|\mathcal{D})$ is "centered" on the data
 - The width of the distribution p(y|D) is consistent with the spread of the data around the nominal model prediction

$$L(\hat{\alpha}) \propto \exp\left(-\frac{1}{2\epsilon^2}\sum_{i=1}^{N}\left[\left(\mu_i(\hat{\alpha}) - y_i\right)^2 + \left(\sqrt{\Sigma_{ii}(\hat{\alpha})} - \gamma|\mu_i(\hat{\alpha}) - y_i|\right)^2\right]\right)$$

K. Sargsyan (ksargsy@sandia.gov)

Calibrated predictions – general x

For fixed α , e.g. Maximum a Posteriori (MAP) value, Uncertain prediction: $y = f(x, \lambda(\alpha; \xi))$

Mean:

$$\mu(x,\alpha) = \mathbb{E}_{\xi}[f(x,\lambda(\alpha;\xi))]$$

Variance:

$$\sigma^2(x,\alpha) = \mathbb{V}_{\xi}[f(x,\lambda(\alpha;\xi))]$$

Average over posterior of α Posterior predictive (PP):

$$p(f|\mathcal{D}) = \int p(f|\alpha)p(\alpha|\mathcal{D})d\alpha = \mathbb{E}_{\alpha}[p(f|\alpha)]$$

PP mean:

$$\mu_{\rm PP}(x) = \mathbb{E}_{\xi} \mathbb{E}_{\alpha}[f(x, \lambda(\alpha; \xi))] = \mathbb{E}_{\alpha} \mu(x, \alpha)$$

PP variance:

$$\sigma_{\rm PP}^2(x) = \underbrace{\mathbb{E}_{\alpha}[\sigma^2(x,\alpha)]}_{} + \underbrace{\mathbb{V}_{\alpha}[\mu(x,\alpha)]}_{}$$

data noise

model error

Calibrated predictions – compare to data at x_i

For fixed α , e.g. Maximum a Posteriori (MAP) value, Uncertain prediction:

$$y_i = f(x_i, \lambda(\alpha; \xi)) + \epsilon_i^{d} \equiv h(x_i, \lambda(\alpha; \xi))$$

Mean:

$$\mu(x_i,\alpha) = \mathbb{E}_{\xi}[f(x_i,\lambda(\alpha;\xi))]$$

Variance:

$$\sigma^2(x_i,\alpha) = \mathbb{V}_{\xi}[f(x_i,\lambda(\alpha;\xi))] + \sigma_{\mathcal{D}}^2$$

Average over posterior of α Posterior predictive (PP):

$$p(h|\mathcal{D}) = \int p(h|\alpha)p(\alpha|\mathcal{D})d\alpha = \mathbb{E}_{\alpha}[p(h|\alpha)]$$

PP mean:

$$\mu_{\mathrm{PP}}(x_i) = \mathbb{E}_{\xi} \mathbb{E}_{\alpha}[f(x_i, \lambda(\alpha; \xi))] = \mathbb{E}_{\alpha} \mu(x_i, \alpha)$$

PP variance:

$$\sigma_{\rm PP}^2(x_i) = \underbrace{\mathbb{E}_{\alpha}[\sigma^2(x_i,\alpha)]}_{\text{model error}} + \underbrace{\mathbb{V}_{\alpha}[\mu(x_i,\alpha)] + \mathbb{E}_{\sigma_{\mathcal{D}}}[\sigma_{\mathcal{D}}^2]}_{\text{data noise}}$$

Calibrated predictions – compare to data at x_i

For fixed α , e.g. Maximum a Posteriori (MAP) value, Uncertain prediction:

$$y_i = h(x_i, \lambda(\alpha; \hat{\xi})) = f(x_i, \lambda(\alpha; \xi)) + \sigma_{\mathcal{D}}\xi_{d+i}$$

Mean:

$$\mu(x_i,\alpha) = \mathbb{E}_{\xi}[f(x_i,\lambda(\alpha;\xi))]$$

Variance:

$$\sigma^2(x_i,\alpha) = \mathbb{V}_{\xi}[f(x_i,\lambda(\alpha;\xi))] + \sigma_{\mathcal{D}}^2$$

Average over posterior of α

Pushed-forward posterior:

$$p(h|\mathcal{D}) = \int p(h|\alpha, \sigma_{\mathcal{D}}) p(\alpha, \sigma_{\mathcal{D}}|\mathcal{D}) d\alpha d\sigma_{\mathcal{D}} = \mathbb{E}_{\alpha, \sigma_{\mathcal{D}}}[p(h|\alpha, \sigma_{\mathcal{D}})]$$

Pushed-forward mean:

$$\mu_{\text{PFP}}(x_i) = \mathbb{E}_{\xi} \mathbb{E}_{\alpha}[f(x_i, \lambda(\alpha; \xi))] = \mathbb{E}_{\alpha} \mu(x_i, \alpha)$$

Pushed-forward variance:

$$\sigma_{\text{PFP}}^{2}(x_{i}) = \underbrace{\mathbb{E}_{\alpha}[\sigma^{2}(x_{i},\alpha)]}_{\text{model error}} + \underbrace{\mathbb{V}_{\alpha}[\mu(x_{i},\alpha)] + \mathbb{E}_{\sigma_{\mathcal{D}}}[\sigma_{\mathcal{D}}^{2}]}_{\text{data noise}}$$
K. Sargsyan (ksargsy@sandia.gov) SIAM UQ 16 April 7, 2016

16/29

Model Error – Predictions

$$f(x; \Lambda) = f(x; \sum_{k} \alpha_{k} \Psi_{k}(\xi)) = \sum_{k} f_{k}(x; \alpha) \Psi_{k}(\xi)$$

- Non-intrusive spectral projection (NISP) will be employed for
 - Likelihood computation
 - Posterior/pushed-forward predictions
 - Easy access to first two moments:

$$\mu(x;\alpha) = f_0(x;\alpha), \qquad \qquad \sigma^2(x;\alpha) = \sum_{k>0} f_k^2(x;\alpha) ||\Psi_k||^2$$

Predictive mean

$$\mathbb{E}[y(x) = \mathbb{E}_{\alpha}[\mu(x;\alpha)]$$

Decomposition of predictive variance

$$\mathbb{V}[y(x)] = \underbrace{\mathbb{E}_{\alpha}[\sigma^2(x;\alpha)]}_{} + \underbrace{\mathbb{V}_{\alpha}[\mu(x;\alpha)] + \sigma_d^2}_{}$$

Model error

Poserior/Data error

Attribution of error components

$$y_{i} = \underbrace{\sum_{k} f_{ik}(\alpha) \Psi_{k}(\xi_{1}, \dots, \xi_{d}) + \sigma_{\mathcal{D}} \xi_{d+i}}_{h_{i}(\hat{\xi}; \hat{\alpha})}$$

Stochastic dimensions:

- Model error ξ_1, \ldots, ξ_d
- Measurement error $\xi_{d+1}, \ldots, \xi_{d+N}$
- Posterior uncertainty (α): can be represented via its own PC expansion (using MCMC samples and Rosenblatt transformation)

Full PC expansion: $y_i = \sum f_j \Psi_j(\hat{\hat{\xi}})$ Full stochastic *germ*:

$$\hat{\xi} = (\underbrace{\xi_1, \dots, \xi_d}_{\text{Model error}}, \underbrace{\xi_{d+1}, \dots, \xi_{d+N}}_{\text{Measurement error}}, \underbrace{\xi_{d+N+1}, \dots, \xi_{d+N+N_{\alpha}}}_{\text{Posterior uncertainty}})$$

Posterior predictive variance:

$$\sigma_{\text{PP}}^2(x_i) = \mathbb{E}_{\alpha}[\sigma^2(x_i,\alpha)] + \mathbb{E}_{\sigma_{\mathcal{D}}}[\sigma_{\mathcal{D}}^2] + \mathbb{V}_{\alpha}[\mu(x_i,\alpha)]$$

Challenges and Mitigation

- Density estimation is more challenging than parameter estimation
 - Inverse problem is ill-posed or intractable
 - \Rightarrow Employ approximate or empirical likelihoods
- Potentially a high-dimensional Bayesian problem
 - Full posterior may be inaccessible...
 - ⇒ Resort to optimization algorithms in no-noise case
 - ... or hard to sample from
 - ⇒ Adaptive MCMC algorithms, Likelihood-informed subspaces
- Sparse data or expensive high-fidelity simulations
 - With low information content, calibration may struggle
 - \Rightarrow More informative priors/regularization

Predictions account for model error

Calibrating single-exponential models with data from a double exponential model $g(x) = e^{-0.5x} + e^{-2x}$

Linear-exponential $f(x, \lambda) = e^{\lambda_1 + \lambda_2 x}$

Additive Gaussian error



K. Sargsyan (ksargsy@sandia.gov)

Predictions account for model error

Calibrating single-exponential models with data from a double exponential model $g(x) = e^{-0.5x} + e^{-2x}$

Linear-exponential $f(x, \lambda) = e^{\lambda_1 + \lambda_2 x}$



Quadratic-exponential $f_2(x, \lambda) = e^{\lambda_1 + \lambda_2 x + \lambda_3 x^2}$



Test problem – Cubic data fit by a line – ABC



- MAP predictive mean centered on data
- MAP predictive standard deviation captures range of discrepancy
- Increasing number of data points has a small effect on both predictive mean and stdev.

Test problem – Cubic data fit by a quadratic – ABC



- MAP predictive mean centered on data
- MAP predictive standard deviation captures range of discrepancy
- Increasing number of data points has a small effect on both predictive mean and stdev.

Test problem – Cubic data fit by a cubic – ABC



- MAP predictive mean centered on data
- MAP predictive standard deviation captures range of discrepancy
- Increasing number of data points has a small effect on both predictive mean and stdev.

More data leads to 'leftover' model error

Calibrating a quadratic $f(x) = \lambda_0 + \lambda_1 x + \lambda_2 x^2$ w.r.t. 'truth' $g(x) = 6 + x^2 - 0.5(x+1)^{3.5}$ measured with noise $\sigma = 0.1$.



Summary of features:

- Well-defined model-to-model calibration
- Model-driven discrepancy correlations
- Respects physical constraints
- Disambiguates model and data errors
- Calibrated predictions of multiple Qols



Chemistry problem – ABC

- Homogeneous ignition, methane-air mixture
- Single-step global reaction model calibrated against a detailed chemical kinetic model
- Data: ignition time; range of initial T & equivalence ratio
- Single-step model:

$$CH_4 + 2O_2 \rightarrow CO_2 + 2H_2C$$

$$\Re = [CH_4][O_2]k_f$$

$$k_f = A \exp(-E/R^o T)$$

•
$$(\ln A, E) = \sum_k \alpha_k \Psi_k(\xi)$$



Quality of Uncertain Calibrated Model Predictions

Calibrated uncertain fit model is consistent with the detailed-model data.

Over the range of (T^0, Φ) :

- MAP predictive mean ignition-time is centered on the data
- MAP predictive stdv is consistent with the scatter of the data



K. Sargsyan, H.N. Najm, and R. Ghanem "On the Statistical Calibration of Physical Models" Int. J. Chem. Kin., 47(4): 246-276, 2015

TransCom3 Experiment of CO2 Flux Inversion

[Gurney et al., Tellus B, 2003]

- Observations d at *N* = 77 sites around the world
- Inverse problem: find fluxes s at M = 22 locations
- Linearized 'response' model R, such that $d\approx Rs$

 $\mathbf{d} = \mathbf{R}\mathbf{s} + \boldsymbol{\epsilon}_{\mathbf{d}}$

- Model **R** is never perfect thus contaminating the inversion
- The inferred values of s compensate for model deficiencies
- *ϵ*_d is meant to capture data errors, but is 'entangled' with model errors

Consider 14 different response models R



Infer fluxes s, given measurements d to satisfy $d\approx Rs$

- Conventional additive Gaussian error (least-squares): $\mathbf{d} = \mathbf{Rs} + \xi$
- Embed probabilistic model for fluxes s:

 $\mathbf{d} = \mathbf{R}(\mu_{\mathbf{s}} + \mathbf{C}_{\mathbf{s}}\xi)$

Consider 14 different response models R



Infer fluxes s, given measurements d to satisfy $d\approx Rs$

- Conventional additive Gaussian error (least-squares): $\mathbf{d} = \mathbf{Rs} + \xi$
- Embed probabilistic model for fluxes s:

 $\mathbf{d} = \mathbf{R}(\mu_{\mathbf{s}} + \mathbf{C}_{\mathbf{s}}\xi)$















Preliminary results – embed model err in C_R

Calibrate with TKE data, predict both TKE and Pressure Pushed forward posterior



Summary and Future

- Represent, quantify and propagate physical model errors
- Parameter estimation \Rightarrow density estimation
- Bayesian machinery to find parameters of the PDFs
- Approximate/empirical likelihoods impose constraints of interest
- Differentiates from data noise; allows model-to-model calibration
- Implemented in UQTk
- Applied in chemistry, climate modeling, LES computations
- K. Sargsyan, H. Najm, and R. Ghanem. "On the Statistical Calibration of Physical Models". *International Journal for Chemical Kinetics*, 47(4): 246-276, 2015.
- Optimal design for maximum information
- Bayesian problem still hard
- Hierarchical Bayesian formulation
- More intrusive embedding; problem specific

Additional Material

Likelihood construction - variants

Full Likelihood

$$L(\alpha) = p(D|\alpha) = p(y_{\text{data},1}, \dots, y_{\text{data},N}|\alpha)$$

Marginalized Likelihood

$$L(\alpha) = p(D|\alpha) = \prod_{i=1}^{N} p(y_{\text{data},i}|\alpha)$$

- Approximate Bayesian Computation (ABC)
 - seek to satisfy the constraints:
 - p(y|D) is "centered" on the data
 - The width of the distribution p(y|D) is "consistent" with the spread of the data around the nominal model prediction

Full Likelihood

$$L(\alpha) = p(D|\alpha) = \pi_f(y_{\text{data},1}, \dots, y_{\text{data},N}|\alpha)$$

where:

 $\pi_f(\cdot, \alpha)$: *N*-variate density of the random variable (f_1, \ldots, f_N) with $f_i = f(x_i, \lambda(\alpha))$

Problem: $\pi_f(\cdot)$ is degenerate in general when N > M

Consider a case with M = 1, $\lambda \sim N(\mu, \sigma^2)$, and $f = \lambda x$ Let N = 2, hence $(f_1, f_2) = (\lambda x_1, \lambda x_2)$ for any λ sample With $f_1/x_1 = f_2/x_2 = \lambda$, (f_1, f_2) are dependent and $\pi_f(\cdot|\mu, \sigma)$ is non-zero only along the line $f_2 = (x_2/x_1)f_1$

hence $\pi_f(y_{data,1}, y_{data,2} | \mu, \sigma)$ is non-zero only along the line $y_{data,2}/x_2 = y_{data,1}/x_1$

Marginalized Likelihood

$$L(\alpha) = p(D|\alpha) = \prod_{i=1}^{N} \pi_{f_i}(y_{\text{data},i}|\alpha)$$

where $\pi_{f_i}(\cdot, \alpha)$ is the univariate density of the RV $f_i = f(x_i, \lambda(\alpha))$

Problem: the likelihood has multiple singularities corresponding to α values leading to vanishing marginal variances at each x_i

Gaussian example: Let $f_i \sim N(\mu_i(\alpha), \sigma_i(\alpha)^2)$, then

$$L(\alpha) = \frac{1}{(2\pi)^{N/2}} \prod_{i=1}^{N} \frac{1}{\sigma_i(\alpha)} \exp\left(\frac{(\mu_i(\alpha) - y_{\text{data},i})^2}{2\sigma_i(\alpha)^2}\right)$$

Multiple singularities, $\sigma_i(\alpha) = 0, i = 1, ..., N$

Posterior maximization always finds one of these singularities, fitting one point perfectly, while misfitting the rest $(\Rightarrow \text{ priors})$

Approximate Bayesian Computation (ABC)

Employ a kernel density as a pseudo-likelihood to enforce select constraints:

• Uncertain prediction p(y|D) is centered on the data

• With $\mu_i(\alpha) = \mathbf{E}_{\xi}[f(x_i, \lambda(\xi; \alpha))]$: minimize $\|\mu_i(\alpha) - y_{\text{data}, i}\|_2^2$

- The width of the distribution p(y|D) is consistent with the spread of the data around the nominal model prediction
 - With $\sigma_i(\alpha)^2 = V_{\xi}[f(x_i, \lambda(\xi, \alpha))]$: minimize $\|(\sigma_i(\alpha) - \gamma | \mu_i(\alpha) - y_{\text{data},i}|)\|_2^2$
 - γ is a factor that specifies the desired match between σ_i and the discrepancy |μ_i(α) y_{data,i}|, on average

ABC Likelihood

With $\rho(S)$ being a metric of the statistic S, use the kernel function as an ABC likelihood:

$$L_{\text{ABC}}(\alpha) = \frac{1}{\epsilon} K\left(\frac{\rho(\mathcal{S})}{\epsilon}\right)$$

where ϵ controls the severity of the consistency control

Propose the Gaussian kernel density:

$$L_{\epsilon}(\alpha) = \frac{1}{\epsilon\sqrt{2\pi}} \prod_{i=1}^{N} \exp\left(-\frac{(\mu_i(\alpha) - y_{\mathrm{d},i})^2 + (\sigma_i(\alpha) - \gamma|\mu_i(\alpha) - y_{\mathrm{d},i}|)^2}{2\epsilon^2}\right)$$