

Robust initialization of variational inference through global optimization and Laplace approximations

Wyatt Bridgman¹, Reese Jones¹, Mohammad Khalil¹

¹*Sandia National Laboratories, Livermore, CA*



Sandia National Laboratories



U.S. DEPARTMENT OF
ENERGY



Conference on
Mathematics of Data Science

Sandia National Laboratories is a multimission laboratory managed and operated by National Technology and Engineering Solutions of Sandia, LLC., a wholly owned subsidiary of Honeywell International, Inc., for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA-0003525. Sand No. ???

September 21, 2022

Context: novel probabilistic strategies for transfer learning

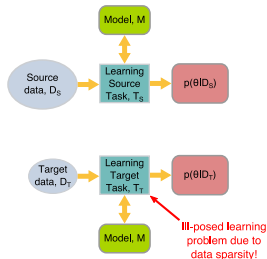
Challenge: Many Sandia mission domains defined by a lack of reliable data, preventing use of many machine learning techniques for predictive modeling:

- ▶ Expensive computer simulations.
- ▶ Prohibitive data acquisition cost.
- ▶ Limited access to classified/sensitive data.

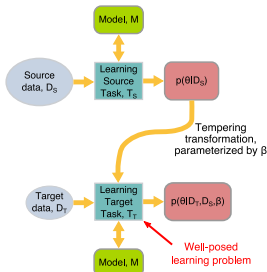
Goal: *Enhance the trust in ML within noisy and sparse data settings.*

Requirement: **High-fidelity, closed-form** approximations of parameter PDFs (as opposed to just samples) for approximation of multimodal target likelihood $p(\theta | \mathcal{D}_t)$.

Traditional TL



Probabilistic TL



Approximation of posteriors in Bayesian inference

Bayesian inference: probabilistic model $p(\mathbf{x}, \mathbf{z})$ with \mathbf{x}, \mathbf{z} observed, latent variables. Some structure specified on the joint distribution e.g.

$p(\mathbf{x}, \mathbf{z}) = p(\mathbf{x} | \mathbf{z})p(\mathbf{z})$. Here we take \mathbf{x} to be a fixed set of data \mathcal{D} . Baye's rule gives:

$$p(\mathbf{z} | \mathcal{D}) = \frac{p(\mathcal{D} | \mathbf{z})p(\mathbf{z})}{p(\mathcal{D})} = \frac{p(\mathcal{D} | \mathbf{z})p(\mathbf{z})}{\underbrace{\int p(\mathcal{D} | \mathbf{z})p(\mathbf{z}) d\mathbf{z}}_{\text{difficult to compute}}} \quad (1)$$

Sampling strategies for $p(\mathbf{z} | \mathcal{D})$ include:

- ▶ **MCMC:** Asymptotically exact but computationally expensive. Difficulties with multimodal distributions.
- ▶ **Dropout:** Adds UQ to neural networks via random perturbations of the weights. Less costly than Variational Inference.

Variational Inference (VI)

Sampling strategies often suffer from scalability issues, we also need **closed form** approximations for the TL framework.

Variational Inference

Approximate posterior $p(\mathbf{z} \mid \mathcal{D})$ by $q_{\theta}(\mathbf{z}) \in \mathcal{F}_{\theta}$ in some family \mathcal{F}_{θ} of PDFs by minimizing error measure such as KL-divergence:

$$q_{\theta}(\mathbf{z}) = \min_{q_{\theta} \in \mathcal{F}_{\theta}} D_{\text{KL}}(q_{\theta}(\mathbf{z}) \parallel p(\mathbf{z} \mid \mathcal{D})) \quad (2)$$

where θ are the variational parameters to be optimized.

- ▶ **Pros:** Obtain closed form approximation $q_{\theta}(\mathbf{z})$ whose fidelity is determined by choice of family, e.g., whether a single Gaussian or mixture. Can be scalable to large NN models depending on choice of \mathcal{F}_{θ} .
- ▶ **Cons:** Can underestimate variance, suffers from optimization pitfalls due to nonconvexity of objectives.

Case: Variational Bayes

- ▶ Take $q(\mathbf{z}) = \prod_{i=1}^M q(\mathbf{z}_i)$, i.e., $q(\mathbf{z})$ is from the space \mathcal{F} of PDFs that factor over partition of latent variables into $\mathbf{z} = \mathbf{z}_1, \dots, \mathbf{z}_M$.
- ▶ Can be shown that optimal solution of functional $\min_{q(\mathbf{z}) \in \mathcal{F}} D_{\text{KL}}(q(\mathbf{z}) \parallel p(\mathbf{z} \mid \mathcal{D}))$ is given by

$$q_j^*(\mathbf{z}_j) = \frac{e^{\mathbb{E}_{i \neq j} [\log p(\mathbf{z}, \mathcal{D})]}}{\int e^{\mathbb{E}_{i \neq j} [\log p(\mathbf{z}, \mathcal{D})]} d\mathbf{z}_j} \quad (3)$$

- ▶ Can sometimes determine $q_j^*(\mathbf{z}_j)$ to be known PDF whose parameters satisfy set of simultaneous nonlinear equations \rightarrow solve iteratively.
- ▶ For ML models $\mathbb{E}_{i \neq j} [\log p(\mathbf{z}, \mathcal{D})]$ hard to compute as $p(\mathbf{z}, \mathcal{D})$ parameterized by nonlinear NN \rightarrow would have to solve set of integral equations.

VI via gradient-based optimization

- The KL-divergence plus the Evidence Lower Bound (ELBO) \mathcal{L}_θ differ by a constant

$$\underbrace{\log p(\mathcal{D})}_{\text{const. w.r.t. } \theta} = \underbrace{\mathcal{L}_\theta(\mathcal{D})}_{\text{ELBO}} + D_{\text{KL}}(q_\theta(\mathbf{z}) \parallel p(\mathbf{z} \mid \mathcal{D})) \quad (4)$$

so that we can minimize $D_{\text{KL}}(\parallel)$ by minimizing

$$-\mathcal{L}_\theta(\mathcal{D}) = \underbrace{D_{\text{KL}}(q_\theta(\mathbf{z}) \parallel p(\mathbf{z}))}_{\text{KL-div from prior}} - \underbrace{\mathbb{E}_{q_\theta} [\log p(\mathcal{D} \mid \mathbf{z})]}_{\text{expected data fit}} \quad (5)$$

- If $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^{N_s}$, likelihood given by IID Gaussian over model predictions $f_{\mathbf{z}}(\mathbf{x}_i)$ with noise σ , then expected data fit like a **stochastic mean-squared error (MSE)**:

$$-\mathbb{E}_{q_\theta} [\log p(\mathcal{D} \mid \mathbf{z})] = \frac{1}{2\sigma} \mathbb{E}_{q_\theta(\mathbf{z})} \left[\sum_{s=1}^{N_s} \|\mathbf{y}_i - f_{\mathbf{z}}(\mathbf{x}_i)\|^2 \right] \quad (6)$$

Example: ELBO for linear NN

- NN given by $\text{NN}_{\mathbf{W}}(\mathbf{x}) = \mathbf{W}\mathbf{x}$, $\mathbf{W} \in \mathbb{R}^{n \times n}$ where

$$q(\mu_q, \Sigma_q) = \mathcal{N}(\mathbf{W} \mid \mu_q, \Sigma_q), \quad p(\mathbf{W}) = \mathcal{N}(\mathbf{W} \mid \mu_p, \Sigma_p)$$

then the ELBO objective function has form

$$\begin{aligned} & \underbrace{\frac{1}{\sigma^2} \text{tr}\{(\mathbf{Y} - \mu_q \mathbf{X})^T (\mathbf{Y} - \mu_q \mathbf{X})\}}_{\text{least squares in } \mu_q} + \underbrace{(\mu_p - \mu_q)^T \Sigma_p^{-1} (\mu_p - \mu_q)}_{\text{regularization } \mu_q \rightarrow \mu_p} \\ & + \underbrace{\log \det(\Sigma_q^{-1} \Sigma_p) + \text{tr}(\Sigma_p^{-1} \Sigma_q)}_{\Sigma_q \rightarrow \Sigma_p} + \underbrace{\frac{1}{\sigma^2} \text{tr}\{\mathbf{V} \mathbf{X} \mathbf{X}^T\}}_{\Sigma_q \rightarrow \mathbf{0}} \end{aligned}$$

which takes the form of **least squares** in means μ_q with **quadratic regularization**. Variance Σ_q balanced between prior Σ_p and $\mathbf{0}$.

Minimizing ELBO for nonlinear model

- ▶ Minimize $D_{\text{KL}}(q_{\theta}(\mathbf{z}) \parallel p(\mathbf{z} \mid \mathcal{D}))$ via gradient descent requires $\nabla_{\theta}(-\mathcal{L}_{\theta})$.
- ▶ **Score function / black-box VI:**

$$\nabla_{\theta}(-\mathcal{L}_{\theta}) = \mathbb{E}_{q_{\theta}(\mathbf{z})} [(\nabla_{\theta} q_{\theta}(\mathbf{z})) \log p(\mathcal{D} \mid \mathbf{z})]$$

- ▶ **Reparametrization gradients:** Express $\mathbf{z} = t(\epsilon, \theta)$, $\epsilon \sim p(\epsilon)$ then gradient and expectation commute

$$\nabla_{\theta} \mathbb{E}_{q_{\theta}(\mathbf{z})} [\log p(\mathcal{D} \mid \mathbf{z})] = \mathbb{E}_{p(\epsilon)} [\underbrace{\nabla_{\mathbf{z}} \log p(\mathcal{D} \mid \mathbf{z}) \nabla_{\theta} \mathbf{z}}_{\text{backprop. gradient}}]$$

Lower variance than score method but reparametrization more difficult for complex distributions like GMMs (Graves, 2016), (Figueroa 2018).

Challenges with VI for high-fidelity distributions

- ▶ VI doesn't scale well with high-fidelity posterior approximations such as Gaussian mixture models or even full covariance Gaussians.
- ▶ ELBO is **nonconvex**, optimizers can find poor local minima (Kingma, Welling 2019).
- ▶ Some approaches to address this issues include annealing (Bowman, 2016),(Sonderby et. al., 2006) and **good initialization strategies** (Rossi et. al 2019).
- ▶ Growing body of literature suggesting **Laplace approximations (LAs)** perform well in a variety of ML/UQ tasks:

$$p(\mathbf{z} \mid \mathcal{D}) \approx \frac{1}{Z_g} \exp \left[-\frac{1}{2}(\mathbf{z} - \mathbf{z}_{\text{MAP}})^T \Sigma^{-1}(\mathbf{z} - \mathbf{z}_{\text{MAP}}) \right]$$

where $\Sigma = -\mathbf{H}_{\log \phi}^{-1}(\mathbf{z}_{\text{MAP}})$, \mathbf{z}_{MAP} is maximum a posteriori estimate

Global optimization and LA

Proposed approach: *Approximate multimodal PDF with global optimization and LAs. Can be used to initialize VI or, possibly, as an alternative approximation strategy.*

Outline of proposed method:

Unnormalized posterior distribution $\tilde{p}(\mathbf{z})$

- ▶ Global optimization carried out on \tilde{p} to find modes $\mathbf{z}_1^*, \dots, \mathbf{z}_K^*$ taken as centers $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K$ of Gaussian components.
- ▶ Laplace approximation formed at each mode:

$$\mathcal{N}(\mathbf{z} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k = \mathbf{H}_{-\log \tilde{p}}^{-1}(\boldsymbol{\mu}_k))$$

- ▶ Fit the weights via constrained least squares:

$$\arg \min_{\boldsymbol{\pi}} \sum_{i=1}^N \left\{ \tilde{p}(\mathbf{z}_i) - \sum_{k=1}^K \tilde{\pi}_k \mathcal{N}(\mathbf{z}_i \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\} \quad \text{s.t. } \tilde{\pi}_k \geq 0$$

then $\int \tilde{p}(\mathbf{z}) d\mathbf{z} \approx \sum_k \tilde{\pi}_k$.

Scalability: VI vs. global opt. & LA

VI with Gaussian Mixture Models

- ▶ VI with $q_{\theta}(\mathbf{z}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{z} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$, $\mathbf{z} \in \mathbb{R}^d$ has variational parameters

$$\boldsymbol{\theta} = (\pi_1, \dots, \pi_K, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_K)$$

so that $\boldsymbol{\theta} \in \mathbb{R}^{K+Kd+K(d+d^2)/2} \longrightarrow$ grows like $\mathcal{O}(d^2)$.

- ▶ Loss function nonconvex means multiple optimization runs are needed to avoid poor local minima.

Global opt. & LA

- ▶ Carry out many local optimizations in smaller parameter space \mathbb{R}^d instead of several expensive optimization runs in larger variational parameter space $\mathbb{R}^{\mathcal{O}(d^2)}$.
- ▶ Enhance scalability with low-rank Hessian approximations.
- ▶ A variety of global optimization techniques can be used such as MLSL which purport to reduce number of local searches needed.

Robustness of approach via global sensitivity

Variance based global sensitivity analysis:

$$f = f_0 + \sum_i f_i(X_i) + \sum_i \sum_{j>i} f_{ij}(X_i, X_j) + \dots$$

$$\mathbb{V}(f) = \mathbb{V}(f_0) + \sum_i \mathbb{V}(f_i) + \sum_i \sum_{j>i} \mathbb{V}(f_{ij}) + \dots$$

Use sensitivity analysis over ensemble of synthetic tests on GMMs to study how performance $f(d, K, \omega, c, \lambda)$ varies as a function of

Parameter	Description	Distribution
d	Dimension	$\mathcal{U}\{2, 10\}$
K	Number of mixture components	$\mathcal{U}\{2, 4\}$
ω	Exponential decay factor across weights	$\mathcal{U}[1, 2]$
c	Correlation coefficient	$\mathcal{U}[0, 0.7]$
λ	Maximum overlap between components	$\mathcal{U}[10^{-4}, 10^{-2}]$

Robustness of approach via global sensitivity

- ▶ ω controls spread of component sizes, λ controls the overlap between components measured by Dice metric.
- ▶ Accuracy measured by $D_{\text{JSD}}(\mathcal{G}(\pi, \mu, \Sigma) \parallel \mathcal{G}(\hat{\pi}, \hat{\mu}, \hat{\Sigma}))$, Jensen-Shannon divergence between true, approximate GMMs. Obtained by "symmetrizing" KL-divergence, bounded.

Global sensitivity results

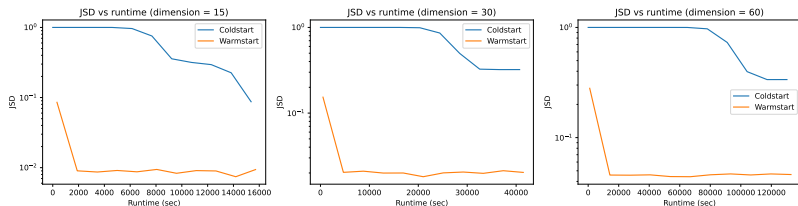
Parameter	Distribution	S	S_T
d , dim.	$\mathcal{U}\{8, 9, 10\}$	0.17 ± 10^{-3}	0.65 ± 10^{-2}
K , no. components	$\mathcal{U}\{3, 4\}$	0.13 ± 10^{-3}	0.30 ± 10^{-3}
ω , weight decay	$\mathcal{U}[1.3, 2]$	0.17 ± 10^{-2}	0.37 ± 10^{-2}
c , corr.	$\mathcal{U}[0.1, 0.7]$	0 ± 10^{-9}	0.65 ± 10^{-2}
λ , overlap	$\mathcal{U}[10^{-4}, 10^{-2}]$	0 ± 10^{-9}	0.02 ± 10^{-4}

Conclusion: *Interaction between factors which increase difficulty of global optimization have the greatest effect.*

Scalability of global optimization, LA method

Can we improve the **scalability of VI** with high-fidelity GMM surrogate posteriors using the GMM approximation scheme?

- ▶ Carry out scalability analysis in high dimensional setting on toy problems with **non-Gaussian trends**.
- ▶ **Cold start** (randomly init. VI) versus **warm start** (GMM init.)
- ▶ Generate non-Gaussian mixture models by applying nonlinear transformation $Y = I + \sigma F(Z, s, t)$ on standard normal r.v. Z where s, t control skewness and tail behavior.



Conclusion: *Using GMM approx. procedure improves scalability and achievable accuracy.*

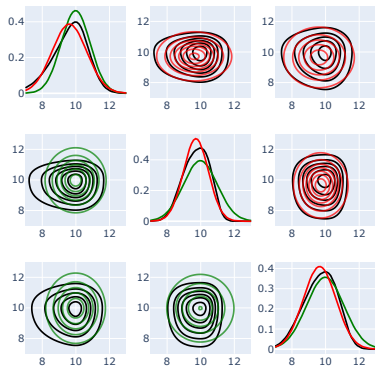
How does the Laplace approximation compare to VI?

Laplace approximation

- ▶ Captures peak and local geometry.
- ▶ Approximation away from peak worsens with increasing non-Gaussian trends.

VI-refined approximation

- ▶ refines support of modes to lie within high-probability regions of true distribution.
- ▶ Doesn't capture peaks as closely.

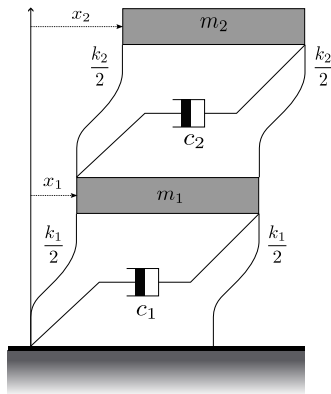


Marginals between 3 variables of 15-dim. distribution. **Black:** true, **Green:** LA, **Red:** VI.

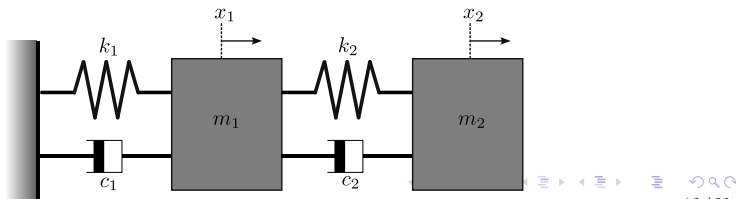
How do these approx. differences translate into predictions?

Structural dynamics problem

- ▶ Two-story shear frame structure.
- ▶ Subject 2nd floor to initial displacement.
- ▶ Inverse problem of determining damping coefficients c_1 , c_2 while observing only the first floor's motion.
- ▶ Can obtain multimodal posterior over c_1, c_2 .



Equivalent to mass-spring-damper system:



Structural dynamics problem

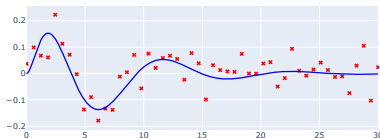
- Equations of motion:

$$\frac{d}{dt} \begin{bmatrix} \mathbf{x} \\ \mathbf{v} \end{bmatrix} = \begin{bmatrix} \mathbf{0} & \mathbf{I} \\ -\mathbf{M}^{-1}\mathbf{K} & -\mathbf{M}^{-1}\mathbf{C} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{v} \end{bmatrix}$$

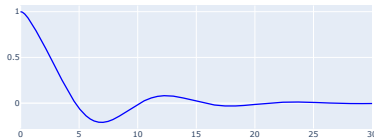
- Log likelihood from matrix exponential

$$\frac{1}{\sigma^2} \sum_{i=1}^{N_D} (y_i - \mathbf{H} \exp(\mathbf{A}(c_1, c_2)t_i) \bar{\mathbf{x}}_0)^2$$

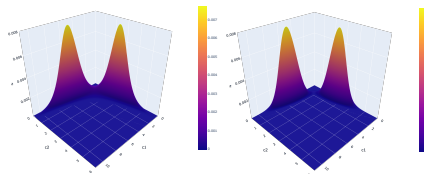
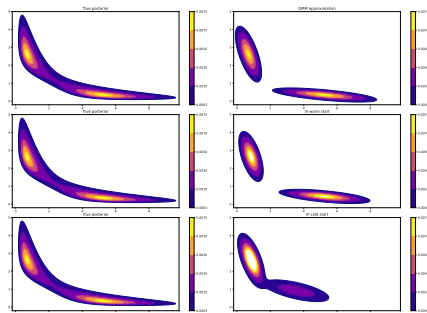
- Noisy observations of first floor



- Second floor displacement



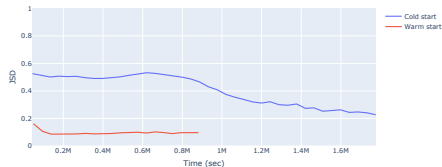
Posterior and GMM approximation



True multimodal posterior (left) and two-component GMM approximation (right).

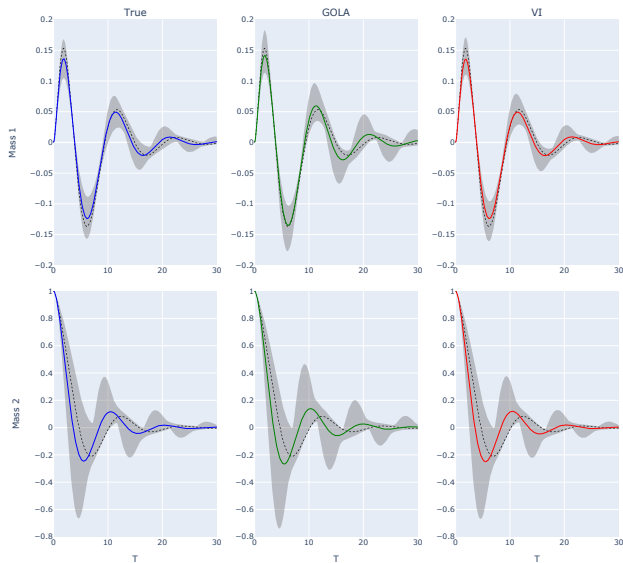
Contour plots of:

- ▶ **(Top):** GMM approximation
- ▶ **(Middle):** GMM approx. refined with VI
- ▶ **(Bottom):** Example of randomly initialized VI solution. Gets stuck in local min.



JS-divergence vs wall-clock time for warm-start, cold-start.

Posterior pushforward



Posterior pushforward of True (left), Global opt., LA (middle), VI-refined (right)

The End

Thank you!