

Uncertainty-aware surrogate models for inverse problems

Aretha Teckentrup

School of Mathematics, University of Edinburgh

Joint work with:

Tianming Bai, Kostas Zygalakis (University of Edinburgh)

Christian Jimenez Beltran, Antonio Vergari, Kostas Zygalakis (University of Edinburgh)

Uncertainty Quantification for High-Dimensional Problems



THE UNIVERSITY of EDINBURGH
School of Mathematics



MAXWELL INSTITUTE FOR
MATHEMATICAL SCIENCES

Outline

- 1 Motivation: Application in Groundwater Flow
- 2 Bayesian inverse problems
- 3 Gaussian process regression
- 4 Bayesian neural networks
- 5 Conclusions

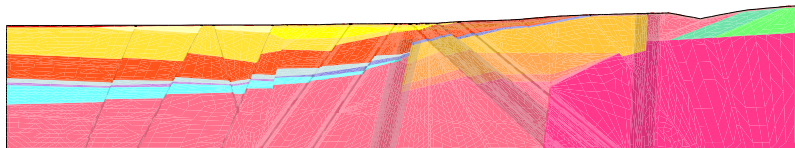
Motivation: Application in Groundwater Flow

- Modelling and simulation of groundwater flow are essential in many applications, e.g. CO₂ capture/storage
- Darcy's law for an incompressible fluid leads to the diffusion equation

$$-\nabla \cdot (k(x) \nabla p(x)) = g(x), \quad x \in D \subseteq \mathbb{R}^3,$$

with hydraulic conductivity k , source/sink terms g , and resulting pressure head p of groundwater

- Can be extended to multi-phase, time-dependent flow. . .



EDZ
CROWN SPACE
WASTE VAULTS
FAULTED GRANITE
GRANITE
DEEP SKIDDAW
N-S SKIDDAW
DEEP LATTERBARROW
N-S LATTERBARROW
FAULTED TOP M-F BVG
TOP M-F BVG
FAULTED BLEAETH BVG
BLEAETH BVG
FAULTED F-H BVG
F-H BVG
FAULTED UNDIFF BVG
UNDIFF BVG
FAULTED N-S BVG
N-S BVG
FAULTED CARB LST
CARB LST
COLLYHURST
FAULTED BROCKRAM
BROCKRAM
SHALES + EVAP
FAULTED BNHM
BOTTOM NHM
FAULTED DEEP ST BEES
DEEP ST BEES
FAULTED N-S ST BEES
N-S ST BEES
FAULTED VN-S ST BEES
VN-S ST BEES
FAULTED DEEP CALDER
DEEP CALDER
FAULTED N-S CALDER
N-S CALDER
FAULTED VN-S CALDER
VN-S CALDER
MERCIA MUDDSTONE
QUATERNARY

Motivation: Application in Groundwater Flow II

Challenge:

- To **simulate** groundwater flow using the Darcy model, we need to know the conductivity k in the entire domain D !
- Typically we only have sparse, noisy data available of k and p , which means there is **significant uncertainty** in k .

Aim:

- We will in this talk focus on combining expert prior knowledge with available data to **infer** k in a Bayesian statistical framework.
- A typical parametrisation chosen for k is **piece-wise constant**, corresponding to different rock types. In particular, we want to infer the values $\theta = \{\theta_1, \dots, \theta_{d_\theta}\}$ that k takes in D .

Bayesian inverse problems (see e.g. [Kaipio, Somersalo '04])

- We choose a **prior density** π_0 on θ , incorporating any expert knowledge such as the hydraulic conductivity being positive.
- Using Bayes' Theorem, we obtain the **posterior density** π^y on $\theta|y$ given by

$$\underbrace{\pi^y(\theta)}_{\text{posterior on } \theta} \approx \underbrace{\exp\left(-\frac{1}{2\gamma^2}\|y - \mathcal{G}(\theta)\|_2^2\right)}_{\text{likelihood on } y} \underbrace{\pi_0(\theta)}_{\text{prior on } \theta},$$

where

- ▶ y is the **observed data**, e.g. noisy point values of p given by $y = \{p(x_i; \theta) + \eta_i\}_{i=1}^{d_y} =: \mathcal{G}(\theta) + \eta$ and $\eta_i \sim N(0, \gamma^2 \mathbf{I})$
- ▶ \mathcal{G} is the parameter-to-observation map, which involves the **solution operator of the PDE**,
- ▶ the prior incorporates expert knowledge, the likelihood fits to the observed data, and the posterior is a combination of both.

Bayesian inverse problems II

Computational challenges:

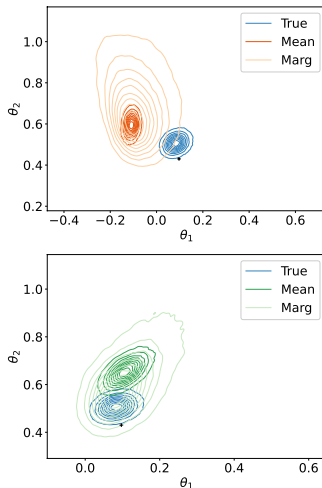
- Algorithms to compute the posterior density, such as Markov chain Monte Carlo, require repeated evaluation of the likelihood, **often $10^4 - 10^6$ evaluations** in practical applications.
- The **evaluation of the likelihood is very costly**, since \mathcal{G} involves the solution of a partial differential equation.

A solution:

- To make computations feasible, we **approximate \mathcal{G} by a surrogate model** (emulator, reduced order model, meta model...).
- This could be e.g. a coarse numerical solver, or a neural network approximation, or
- This leads to an approximate posterior distribution that is feasible to sample from, BUT we have to consider its accuracy!

Bayesian inverse problems III (see e.g. [Bai, T, Zygalakis '24])

- Simply plugging the surrogate model into the posterior typically leads to **biased and overconfident predictions**.
- We need surrogate models that are:
 - ▶ uncertainty aware, to have the right level of confidence,
 - ▶ PDE-informed, to improve accuracy and avoid spurious predictions.
- We do this with:
 - ▶ PDE-constrained Gaussian processes, or
 - ▶ PDE-constrained Bayesian neural networks.



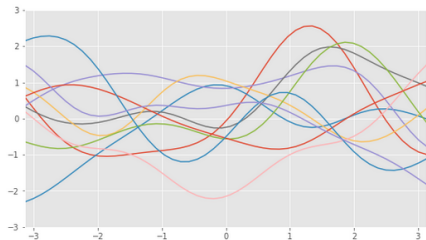
Example: 1d diffusion equation with
 $d_\theta = 2$

Gaussian process regression

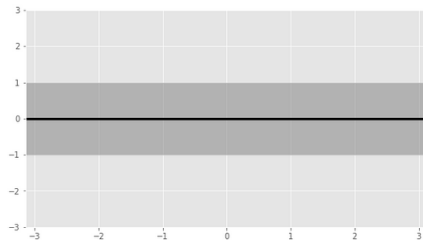
Set-up [Rasmussen, Williams '06]

- Gaussian process regression is a Bayesian methodology to emulate a function $f : \mathcal{T} \rightarrow \mathbb{R}$, e.g. $f = \Phi$ or $f = \mathcal{G}_j$, $j = 1, \dots, d_y$.
- We put a **Gaussian process prior** $\text{GP}(0, k)$ on f , where k is chosen to reflect properties of f .

For $\{\theta_i\}_{i=1}^m \subseteq \mathcal{T}$, the random variables $\{f(\theta_i)\}_{i=1}^m$ follow a joint Gaussian distribution with $\mathbb{E}[f(\theta_i)] = 0$ and $\mathbb{C}[f(\theta_i), f(\theta_j)] = k(\theta_i, \theta_j)$.



Sample paths



Mean and standard deviation

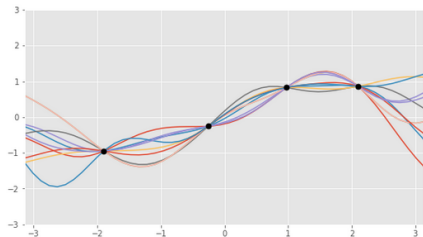
Gaussian process regression

Predictive distribution [Rasmussen, Williams '06]

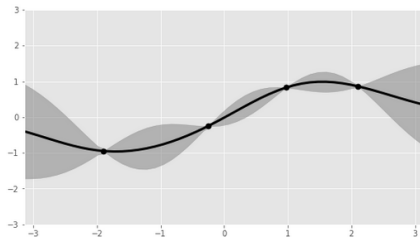
- The **Gaussian process posterior** $f_N \sim \text{GP}(m_N^f, k_N)$ on $f|d$ is obtained by conditioning the prior on function values $d = \{\theta^n, f(\theta^n)\}_{n=1}^N$:

$$m_N^f(\theta) = k(\theta, D_N)^T K(D_N, D_N)^{-1} f(D_N),$$
$$k_N(\theta, \theta') = k(\theta, \theta') - k(\theta, D_N)^T K(D_N, D_N)^{-1} k(\theta', D_N),$$

where $D_N = \{\theta^n\}_{n=1}^N$, $k(\theta, D_N) = [k(\theta, \theta^1), \dots, k(\theta, \theta^N)] \in \mathbb{R}^N$ and $K(D_N, D_N) \in \mathbb{R}^{N \times N}$ is the matrix with ij^{th} entry equal to $k(\theta^i, \theta^j)$.



Sample paths



Mean and standard deviation

Gaussian process regression

Approximations of the posterior

- Recall: $\pi^y(\theta) = \frac{1}{Z} \exp \left(- \frac{1}{2\gamma^2} \|y - \mathcal{G}(\theta)\|^2 \right) \pi_0(\theta)$
- For the remainder of the talk, assume that we approximate \mathcal{G} by Gaussian process regression. Similar results hold for $\Phi = \frac{1}{2\gamma^2} \|y - \mathcal{G}(\theta)\|^2$.
- Since the surrogate model \mathcal{G}_N is a stochastic process, a deterministic approximation of π^y is obtained:
 - by taking the **mean-based approximation**

$$\pi_{N,\text{mean}}^y(\theta) = \frac{1}{Z_{\text{mean}}^N} \exp \left(- \frac{1}{2} \|y - m_N^{\mathcal{G}}(\theta)\|_{\Gamma^{-1}}^2 \right) \pi_0(\theta),$$

- or by taking the **marginal approximation**

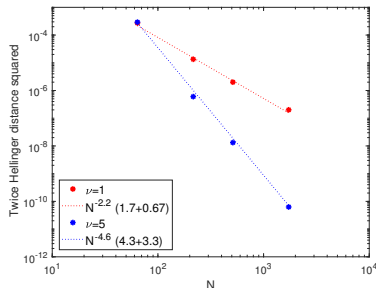
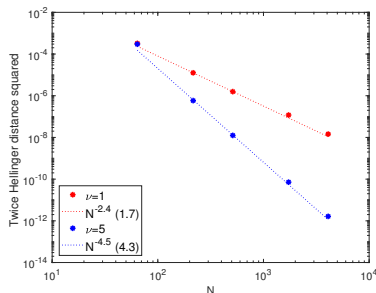
$$\pi_{N,\text{marg}}^y(\theta) = \frac{1}{\mathbb{E}(Z_N^{\text{rand}})} \mathbb{E} \left(\exp \left(- \frac{1}{2} \|y - \mathcal{G}_N(\theta)\|_{\Gamma^{-1}}^2 \right) \right) \pi_0(\theta).$$

Gaussian process regression

Convergence as $N \rightarrow \infty$ [Stuart, T '18], [T 20], [Helin, Stuart, T, Zygalakis '23]

- Both approximations converge to the true posterior π^y as $N \rightarrow \infty$.

Example: 1d diffusion equation with $d_\theta = 3$



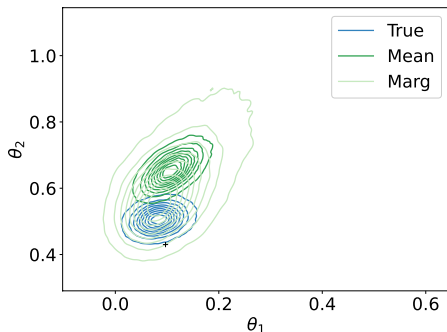
Left: Mean-based approximation. Right: Marginal approximation.

- Error in Hellinger distance depends on $\|\mathcal{G} - m_N^{\mathcal{G}}\|_{L^2(\mathcal{T}; \mathbb{R}^{d_y})}$ and $\left\| \mathbb{E} (\|\mathcal{G} - \mathcal{G}_N\|^{1+\delta})^{\frac{1}{1+\delta}} \right\|_{L^2(\mathcal{T})}$ for any $\delta > 0$, respectively.

Gaussian process regression

Mean-based and marginal approximations [Bai, T, Zygalakis '24]

- For small N , the difference between $\pi_{N,\text{mean}}^y$ and $\pi_{N,\text{marg}}^y$ can be significant.
- Using $\pi_{N,\text{mean}}^y$ can lead to biased predictions with high confidence.
- Only $\pi_{N,\text{marg}}^y$ uses the **uncertainty in \mathcal{G}_N** , modelling the error in the surrogate model.



Example: 1d diffusion equation with $d_\theta = 2$ and $N = 4$

Gaussian process regression

Marginal likelihood [Bai, T, Zygalakis '24]

- With $\mathcal{G}_N \sim \text{GP}(m_N^{\mathcal{G}}, k_N)$, we can analytically compute the marginal likelihood $\mathbb{E} \left(\exp \left(- \frac{1}{2\gamma^2} \|y - \mathcal{G}_N(\theta)\|^2 \right) \right)$.
- We have $\mathcal{G}_N(\theta) = m_N^{\mathcal{G}}(\theta) + \xi$, with $\xi \sim \text{N}(0, k_N(\theta, \theta))$. Hence

$$\begin{aligned} & \mathbb{E} \left(\exp \left(- \frac{1}{2} \|y - \mathcal{G}_N(\theta)\|_{\Gamma^{-1}}^2 \right) \right) \\ &= \frac{1}{\sqrt{(2\pi)^{d_y} \det(\Sigma(\theta))}} \int_{\mathbb{R}^{d_y}} \exp \left(- \frac{\|y - m_N^{\mathcal{G}}(\theta) - \xi\|_{\Gamma^{-1}}^2}{2} \right) \exp \left(- \frac{\|\xi\|_{\Sigma^{-1}(u)}^2}{2} \right) d\xi \\ &\propto \frac{1}{\sqrt{\det(\Gamma + \Sigma(\theta))}} \exp \left(- \frac{\|y - m_N^{\mathcal{G}}(\theta)\|_{(\Gamma + \Sigma(\theta))^{-1}}^2}{2} \right), \end{aligned}$$

where $\Sigma(\theta) = k_N(\theta, \theta)$.

Gaussian process regression

Variance inflation

- Compared to the mean-based likelihood

$$\frac{1}{\sqrt{\det(\Gamma)}} \exp\left(-\frac{\|y - m_N^{\mathcal{G}}(\theta)\|_{\Gamma^{-1}}^2}{2}\right),$$

the marginal likelihood

$$\frac{1}{\sqrt{\det(\Gamma + \Sigma(\theta))}} \exp\left(-\frac{\|y - m_N^{\mathcal{G}}(\theta)\|_{(\Gamma + \Sigma(\theta))^{-1}}^2}{2}\right),$$

is a form of **variance inflation**.

Gaussian process regression

Variance inflation II

- Variance inflation is an emerging tool to improve Bayesian inference in complex models, see e.g. [Conrad et al '17], [Calvetti et al '18], [Cui, Fox, Neumayer '20].
- It is closely related to the well-established inclusion of **modelling error** [Kennedy, O'Hagan '01]:

$$y = \mathcal{G}(\theta) + \eta + \tilde{\eta},$$

with $\tilde{\eta} \sim \mathcal{N}(m, C)$.

- Using Gaussian process regression, we have
 - ▶ a **parameter-dependent variance inflation** $\Sigma(\theta)$, rather than assuming that the error in the (surrogate) model is independent of θ .
 - ▶ an **explicit** model for $\Sigma(\theta)$ that is **readily tuned**.

Markov chain Monte Carlo methods

- In practice, we need to use **sampling methods** such as MCMC to sample from target density $\pi = \pi_{N,\text{mean}}^y$ or $\pi = \pi_{N,\text{marg}}^y$.

ALGORITHM 1. (Metropolis Hastings)

- Choose $\theta^{(1)}$ with $\pi(\theta^{(1)}) > 0$.
- At state $\theta^{(i)}$, sample a proposal θ' from density $q(\theta' | \theta^{(i)})$.
- Accept sample θ' with probability

$$\alpha(\theta' | \theta^{(i)}) = \min \left(1, \frac{\pi(\theta') q(\theta^{(i)} | \theta')}{\pi(\theta^{(i)}) q(\theta' | \theta^{(i)})} \right),$$

i.e. $\theta^{(i+1)} = \theta'$ with probability $\alpha(\theta' | \theta^{(i)})$; otherwise stay at $\theta^{(i+1)} = \theta^{(i)}$.

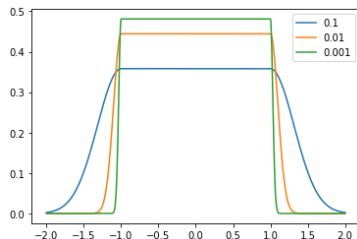
Gaussian process regression

MALA proposals [Roberts, Tweedie '96], [Bai, T, Zygalakis '24]

- In inverse problems, we often have high dimensional parameters θ , and we require a an **efficient choice of proposals** such as MALA:

$$\theta' = \theta^{(i)} + \beta \nabla \log \pi(\theta^{(i)}) + \sqrt{2\beta} \xi_i, \quad \text{where} \quad \xi_i \sim \mathcal{N}(0, \mathbf{I})$$

- For $\pi = \pi_{N,\text{mean}}^y$ and $\pi = \pi_{N,\text{marg}}^y$, the gradient of the log-likelihood exists provided $k(\cdot, \theta^n)$ is differentiable.
- For common choices of k , e.g. $k(\theta, \theta') = \sigma^2 \exp(-\frac{\|\theta - \theta'\|_2^2}{2\lambda^2})$, the **gradient can be computed explicitly**.
- Some priors, such as the uniform prior, require smoothing using Moreau–Yoshida regularisation [Pereyra '16].



Gaussian process regression

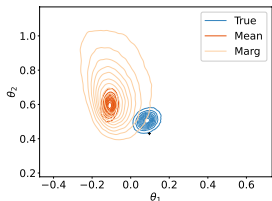
PDE constraints [Bai, T, Zygalakis '24], [Raissi, Perdikaris, Karniadakis '17]

For **linear PDEs** $\mathcal{L}_\theta p(x; \theta) = g(x)$, such as the diffusion equation with $\mathcal{L}_\theta p(x; \theta) = -\nabla \cdot (k(x; \theta) \nabla p(x; \theta))$, we can **incorporate** \mathcal{L}_θ :

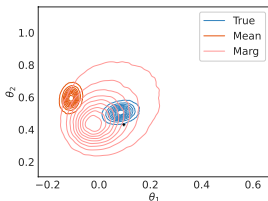
- Put a **joint Gaussian process prior** on p and g as a function of θ :

$$\begin{bmatrix} p(x_p, \theta) \\ g(x_g, \theta) \end{bmatrix} \sim \text{GP} \left(0, k_1(\theta, \theta') \begin{bmatrix} k_2(x_p, x_p) & \mathcal{L}^{\theta'} k_2(x_p, x_g) \\ \mathcal{L}^\theta k_2(x_p, x_g) & \mathcal{L}^\theta \mathcal{L}^{\theta'} k_2(x_g, x_g) \end{bmatrix} \right).$$

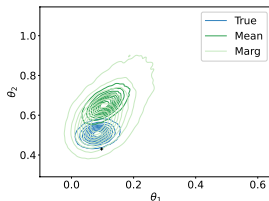
- This gives a joint prior on $\{p(x_i, \theta)\} = \mathcal{G}(\theta)$ and $\{g(\tilde{x}_j, \theta)\}$.
- We update this to a posterior by conditioning on training data in the usual way, and then use the marginal posterior on p .



Independent



Spatially correlated



PDE-constrained

Bayesian neural networks

PDE constraints [Jiminez Beltran et al, '24], [Sirignana, Spiliopoulos '18]

- The **deep Galerkin method** provides a neural network that
 - ▶ takes as inputs x in the spatial domain and parameter value θ ,
 - ▶ is trained to approximate the solution of the PDE, $f_W(x; \theta) \approx p(x; \theta)$.
- In particular, f_W is trained to minimize the following **loss**:

$$\frac{1}{K} \sum_{i=1}^K \ell(D_i; W) = \frac{1}{K} \sum_{i=1}^K (\ell_g(D_{g,i}; W) + \ell_b(D_{b,i}; W)),$$

where

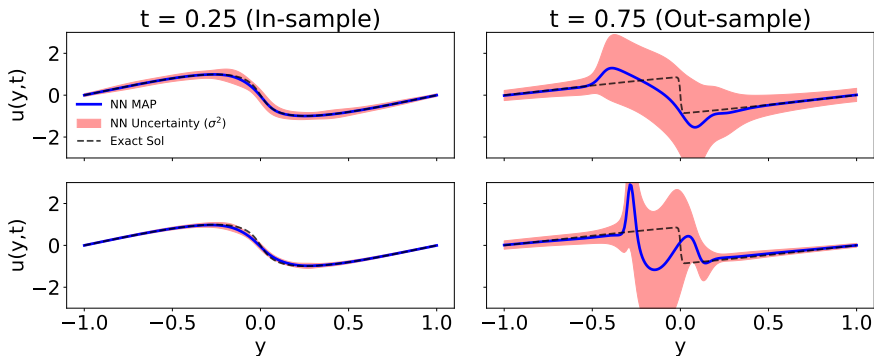
- ▶ at each iteration in training, K collocations for x_i , ∂x_i , and θ_i are sampled from densities π^p , π^b , and π^θ , respectively,
- ▶ $\{D_i\} = \{D_{g,i}\} \cup \{D_{b,i}\}$, $D_{g,i} = \{x_i, \theta_i\}$ and $D_{b,i} = \{\partial x_i, \theta_i\}$,
- ▶ $\ell_g(D_{g,i}; W) = (\mathcal{L}_{\theta_i}(x_i, f_W(x_i, \theta_i); \theta_i) - g(x_i; \theta_i))^2$ measures the error in the approximation of the differential operator, and
- ▶ $\ell_b(D_{b,i}; W) = (f_W(\partial x_i, \theta_i) - b(\partial x_i; \theta_i))^2$ measures the error in the boundary conditions.

Bayesian neural networks

PDE constraints + Bayesian [Jiminez Beltran et al, '24]

- To provide **uncertainty quantification**, we interpret the loss as a negative log-likelihood.
- By doing this, the **likelihood** $p(D|W)$ characterises how well a choice of weights W approximates the PDE solution $p(x; \theta)$, and we obtain a corresponding **posterior** $p(W|D)$ on the weights.
- For computational efficiency, we:
 - ▶ only consider weights in the **penultimate layer** of the neural network as random, and
 - ▶ use a **Laplace approximation** of the posterior $p(W|D)$,which is generally enough to deliver good uncertainty estimates.

Numerical example [Jiminez Beltran et al '24]



Burgers Equation:

$$\frac{\partial u}{\partial t} + u \frac{\partial u}{\partial y} - \theta \frac{\partial^2 u}{\partial y^2} = 0, \quad t > 0, \quad \theta > 0 \quad \text{and} \quad y \in [-1, 1],$$

with conditions $u(y, 0) = -\sin(\pi y)$, and $u(-1, t) = u(1, t) = 0$.

Conclusions

- Partial differential equations model many phenomena in science and engineering.
- Many tasks, such as parameter inference in partial differential equation models, can quickly become infeasible.
- It is common to use surrogate models to alleviate the computational burden.
- However, to avoid overconfident and biased predictions, it is crucial to use surrogate models that:
 - ▶ come with uncertainty quantification, and
 - ▶ incorporate physical constraints.

References I



T. BAI, A. L. TECKENTRUP, AND K. C. ZYGALAKIS, *Gaussian processes for Bayesian inverse problems associated with linear partial differential equations*, Statistics and Computing, 34 (2024), p. 139.



C. JIMINEZ BELTRAN, A. VERGARI, A. L. TECKENTRUP, AND K. C. ZYGALAKIS, *Galerkin meets Laplace: Fast uncertainty estimation in neural PDEs*, in ICLR 2024 Workshop on AI4DifferentialEquationsInScience, 2024.



J. KAIPIO AND E. SOMERSALO, *Statistical and computational inverse problems*, Springer, 2004.



M. RAISSI, P. PERDIKARIS, AND G. E. KARNIADAKIS, *Machine learning of linear differential equations using Gaussian processes*, Journal of Computational Physics, 348 (2017), pp. 683–693.



C. E. RASMUSSEN AND C. K. WILLIAMS, *Gaussian processes for machine learning*, (2006).



J. SIRIGNANO AND K. SPILIOPOULOS, *DGM: A deep learning algorithm for solving partial differential equations*, Journal of computational physics, 375 (2018), pp. 1339–1364.