

Efficient PINNs and Natural Gradients in JAX

ML & AD in JAX for SC, Strasbourg

Nilo Schwencke

OCKHAM team (LIP – ENS Lyon • INRIA Lyon • CNRS)
Formerly (PhD): TAU team (LISN – Université Paris-Saclay • INRIA Saclay • CNRS)

June 11, 2026



Inria



- ① Parametric models and Natural Gradient in a nutshell
- ② Parametric models in JAX
- ③ Essential tools for PINNs implementation
- ④ Effective PINNs training implementation
- ⑤ Natural Gradient Implementation
- ⑥ Empirical motivation

Parametric models and Natural Gradient in a nutshell

Parametric model

$$u : \begin{cases} \mathbb{R}^P & \rightarrow \mathcal{H} \\ \theta & \mapsto u_\theta \end{cases}; \mathcal{H} \text{ Hilbert space}$$

- $\mathcal{M} := \text{Im } u = \{u_\theta : \theta \in \mathbb{R}^P\}$
- $T_\theta \mathcal{M} := \text{Im } du_\theta = \text{span}(\partial_p u_\theta)$

Quadratic regression

$$\mathcal{L}(u) := \frac{1}{2} \|u - f\|_{L^2(\Omega)}^2,$$
$$d\mathcal{L}|_u(h) = \underbrace{\langle u - f, h \rangle}_{\nabla \mathcal{L}|_u}.$$

Induces the gradient flow:

$$\begin{cases} u_0 \in L^2(\Omega) \\ \dot{u}_t = -\nabla \mathcal{L}|_{u_t} = f - u_t \end{cases}$$

Solution: $u_t = f - e^{-t}(f - u_0)$.

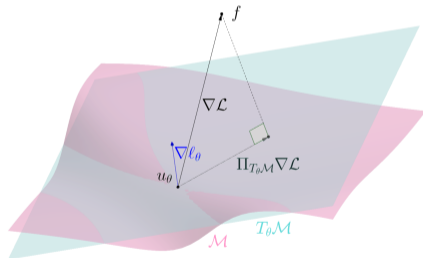
Natural gradient in functional space

But u_θ lives on the manifold \mathcal{M} .

\Rightarrow Project $\nabla \mathcal{L}$ onto $T_\theta \mathcal{M}$.

Then (Amari and Douglas 1998):

$$\theta_{t+1} \leftarrow \theta_t - \eta du_{\theta_t}^\dagger \left(\Pi_{T_{\theta_t} \mathcal{M}}^\perp \nabla \mathcal{L}|_{u_{\theta_t}} \right),$$



$$\theta_{t+1} \leftarrow \theta_t - \eta G_\theta^\dagger \nabla_\theta \ell(\theta); \quad G_{\theta_{pq}} = \langle \partial_p u_\theta, \partial_q u_\theta \rangle_{\mathcal{H}};$$
$$\ell(\theta) := \mathcal{L}(u_\theta).$$

Natural Gradient of PINNs

Key remark

The only difference between the losses:

$$\mathcal{L}_{D,B}(u) = \int_{\Omega} \|D[u] - f\|^2 + \int_{\partial\Omega} \|B[u] - g\|^2$$

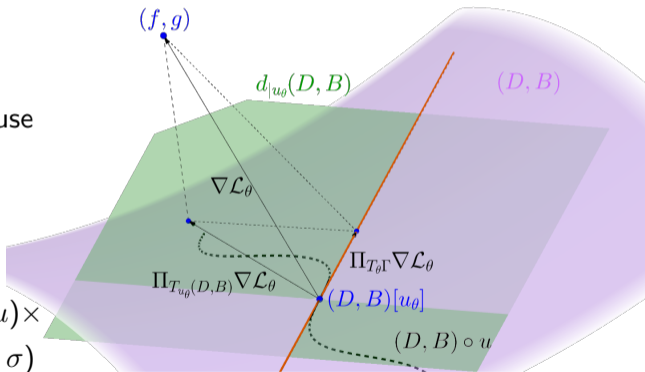
and $\mathcal{L}(u) = \frac{1}{2} \|u - f\|_{L^2(\Omega)}^2$ is the use of the operators D and B .

Proposition

PINNs are a quadratic regression problem with model: $(D, B) \circ u$:

$$\begin{cases} \mathbb{R}^P & \rightarrow \mathcal{H} & \rightarrow L^2(\Omega \rightarrow \mathbb{R}, \mu) \times L^2(\partial\Omega \rightarrow \mathbb{R}, \sigma) \\ \theta & \mapsto u_{\theta} & \mapsto (D[u_{\theta}], B[u_{\theta}]) \end{cases}$$

Figure: Illustration of PINNs Natural Gradient



Parametric models in JAX

Implementing a parametric model in JAX

Parametric model

$$u : \begin{cases} \mathbb{R}^P & \rightarrow \mathcal{H} \\ \theta & \mapsto u_\theta \end{cases} ; \mathcal{H} \text{ Hilbert space}$$

Functional pov: matches JAX spirit!

Linear example

```
1 def linear_model_factory(function_basis):
2     def u(theta, x):
3         return jnp.sum([jnp.dot(a, f(x)) for
4             ↪ a, f in zip(theta,
5             ↪ function_basis)])
6     return u
```

Abstract implementation

From a JAX perspective, a parametric model is a function u with signature:

```
u: [PyTree, Array] -> Array #type Pmodel
```

where the first argument contains the trainable parameters and the second argument contains the evaluation points.

Typical usage for evaluation

```
1 # Given:
2 # -theta (PyTree)
3 # -samples (Array)
4 v_u = jax.vmap(u, (None, 0))
5 u_eval = v_u(theta, samples)
```

Differentiating a parametric model in JAX

Differential of a parametric model

$$du_{\theta} : \begin{cases} \mathbb{R}^P & \rightarrow \mathcal{H} \\ h & \mapsto du_{\theta}(h) \end{cases} ; \mathcal{H} \text{ Hilbert space}$$

How to compute du_{θ} ? **with AD !**

How to do it in practice?

Use `jax.linearize`:

```
1 # Given:
2 # -u (Pmodel)
3 # -theta (PyTree)
4 # -samples (Array)
5 _, du_theta = jax.linearize(lambda p:
  ↪ jax.vmap(u, (None, 0))(p, samples), theta)
```

Given any PyTree `h` with the same structure as `theta`,

$$du_theta(h) = \left(du_{\theta}(h)(x) \right)_{x \in \text{samples}}$$

Implementation remarks

Only efficient if we differentiate in a given direction h (not for the full Jacobian; see later)

Essential tools for PINNs implementation

Continuous problem

$$\begin{cases} \mathcal{D}[u] = f & \text{in } \Omega, \\ \mathcal{B}[u] = g & \text{on } \partial\Omega. \end{cases}$$

$$\mathcal{L}(u) = \frac{1}{2} \|\mathcal{D}[u] - f\|_{L^2(\Omega)}^2 + \frac{\lambda}{2} \|\mathcal{B}[u] - g\|_{L^2(\partial\Omega)}^2$$

Discretization of the losses

Given a parametric model u :

$$\begin{aligned} \hat{\mathcal{L}}(\theta) &= \frac{1}{2} \frac{1}{N_\Omega} \sum_i (\mathcal{D}[u_\theta] - f)(x_i)^2 \\ &\quad + \frac{\lambda}{2} \frac{1}{N_{\partial\Omega}} \sum_j (\mathcal{B}[u_\theta] - g)(z_j)^2. \end{aligned}$$

Implementation ingredients

- θ : PyTree of parameters;
- $(x_i), (z_j)$: Arrays of samples;
- \mathcal{D}, \mathcal{B} built by autodiff (see later).

We start with an abstract notion of Operator (that encompasses \mathcal{D}, \mathcal{B})

Definition

```

1 Field = Array -> Array #type
2 op: Field -> Field #type Operator

```

How to insert a Pmodel u ?

```

1 def op_u_factory(u: Pmodel, op:
  ↪ Operator):
2     def op_u(theta: PyTree, x: Array):
3         return op(lambda z: u(theta,
  ↪ z))(x)
4     return op_u

```

Key fact: `op_u` is also a Pmodel!

- Can be vectorized (`jax.vmap`)
- Can be auto-differentiated

What if `op_u` has parameters ?

```

p_op: [Field, PyTree] -> Field #type POperator

```

- Offline

```

1 # Given alpha (PyTree)
2 op: Operator = lambda v: p_op(v, alpha)

```

- Online

```

1 def p_op_up_factory(u: Pmodel, op: Operator):
2     def p_op_u(theta: PyTree, x: Array, alpha:
  ↪ PyTree):
3         return op(lambda z: u(theta, z),
  ↪ alpha)(x)
4     return p_op_u

```

Computing differential operators in JAX

Partial derivatives by AD

Computation of $\partial_i v$, $v : \mathbb{R}^n \rightarrow \mathbb{R}^m$:

```
1 def dd(v: Field, i) -> Field:
2   def dv_i(x):
3     e_i = jnp.zeros_like(x).at[i].set(1.)
4     _, dv_ix = jax.jvp(v, (x,), (e_i,))
5     return dv_ix
6   return dv_i
```

Higher orders by induction

Computation of $\partial_{ij} v$:

```
1 # Given v a Field
2 v_ij = dd(dd(v,j), i)
```

Example: heat operator

$$\mathcal{D}[u](t, x) = \partial_t u(t, x) - \kappa \partial_{xx} u(t, x).$$

```
1 def heat_operator(v: Field, kappa) -> Field:
2   v_t = dd(v, 0)
3   v_xx = dd(dd(v, 1), 1)
4
5   def heat_v(tx):
6     return v_t(tx) - kappa * v_xx(tx)
7   return heat_v
```

Computational remarks

- Chaining `dd` may become expensive. `jax.experimental.jet` is an appealing alternative.
- We may also use more efficient techniques for some PDEs (e.g., the Laplacian, cf. practical session)

Effective PINNs training implementation

$$\text{PINNs loss: } \hat{\mathcal{L}}(\theta) = \frac{1}{2} \frac{1}{N_{\Omega}} \sum_i \left(\mathcal{D}[u_{\theta}] - f \right) (x_i)^2 + \frac{\lambda}{2} \frac{1}{N_{\partial\Omega}} \sum_j \left(\mathcal{B}[u_{\theta}] - g \right) (z_j)^2$$

is a special case of an abstract loss: $\hat{\mathcal{L}}_e(\theta) := \sum_{s=1}^e \sum_{i=1}^{N_s} \frac{w_s}{2N_s} \left(\mathcal{O}_s[u_{\theta}](x_i^s) - f_k(x_i^s) \right)^2$

Implementation

```

1 def loss_factory(u: Pmodel, ops: Iterable[Operators], fs: Iterable[Fields]):
2     v_op_us = [jax.vmap(op_u_factory(u, op), (None, 0)) for op in ops]
3     v_fs = [jax.vmap(f, (0,)) for f in fs]
4     def loss(theta: PyTree, samples: Array, weights: Iterable[float]):
5         cumulative_loss = 0.
6         for v_op_u, v_f, sample, weight in zip(v_op_us, v_fs, samples, weights):
7             op_u_eval = v_op_u(theta, sample)
8             f_eval = v_f(sample)
9             residual = op_u_eval - f_eval
10            loss_op = jnp.mean(residual ** 2)
11            cumulative_loss += weight * loss_op
12        return cumulative_loss / 2
13    return loss

```

Simply apply `jax.grad`

```
1 # Given:
2 # -loss (obtained with loss_factory)
3 # -theta (PyTree)
4 # -samples (Array)
5 # -weights (Iterable[float])
6 grad_loss = jax.grad(loss)
7 grad_loss_eval = grad_loss(theta, samples,
  ↪ weights) #PyTree
```

Last issue: update theta

Use `jax.flatten_util.ravel_pytree`
(should be imported separately!)

```
1 from jax.flatten_util import ravel_pytree
2 flat_theta, unravel_theta =
  ↪ ravel_pytree(theta)
```

- `flat_theta`: Array representing theta (arithmetic is applicable)

... and its gradient

- `unravel_theta`: Array \rightarrow PyTree, such that

```
unravel_theta(flat_theta) == theta #True
```

We can now train our PINN !

```
1 # minimal training loop
2 # Given:
3 # -etas (Iterable): learning rates
4 # -nsteps (int): number of training steps
5 for i in range(nsteps):
6     eta = etas[i]
7     grad_loss_eval = grad_loss(theta, samples,
  ↪ weights)
8     flat_theta, unravel_theta =
  ↪ ravel_pytree(theta)
9     flat_grad, _ = ravel_pytree(grad_loss_eval)
10    theta = unravel_theta(theta - eta *
  ↪ flat_grad)
```

Natural Gradient Implementation

Computing the Jacobian for Natural Gradient

What we need?

Given a parametric model u ,

$$J_\theta = \left(\partial_{\theta_p} u_\theta(x_i) \right)_{p,i}$$

Then:

$$G_\theta = \frac{1}{N} J_\theta J_\theta^\top, \quad \nabla_{\text{NG}} \hat{\mathcal{L}} = G_\theta^\dagger \nabla_\theta \hat{\mathcal{L}}.$$

When using (abstract) operators

Compute for each operator \mathcal{O}_s :

$$J_\theta^s = \left(\partial_{\theta_p} \mathcal{O}_s[u_\theta](x_i) \right)_{p,i}; \quad G_\theta^s = \frac{1}{N_s} J_\theta^s J_\theta^{s\top}.$$

Then the full Gram reads:

$$G_\theta := \sum_s \lambda_s G_\theta^s.$$

Forward-mode implementation

```
1 def jac_linearize_factory(u):
2     def jac(theta: PyTree, samples: Array):
3         flat_theta, unravel =
4             ↪ ravel_pytree(theta)
5
6         def batch_u(p):
7             return jax.vmap(u, (None,
8                 ↪ 0))(unravel(p), samples)
9
10        _, tangent_fn = jax.linearize(batch_u,
11            ↪ flat_theta)
12
13        # P forward passes
14        can_basis = jnp.eye(flat_theta.size)
15        return jax.vmap(tangent_fn)(can_basis)
16    return jac
```

Computing the Jacobian for Natural Gradient

Efficient reverse-mode version

```
1 def jacrev_factory(u):
2     def jac(theta: PyTree, samples: Array):
3         flat_theta, unravel = ravel_pytree(theta)
4
5         def u_flat(p, x):
6             return model(unravel(p), x)
7
8         jac_x = jax.jacrev(u_flat, argnums=0)
9         J = jax.vmap(jac_x, (None, 0))(flat_theta,
10          ↪ samples)
11
12         return J.reshape((flat_theta.size, -1))
13     return jac
```

Running time comparison

```
fwd mod: 67.1 ms ± 3.32 µs per loop
rev mod: 256 µs ± 7.15 µs per loop
```

Warning

Do not apply jacrev after vmap:

```
jax.jacrev(lambda p: vmap_model(p,
↪ xs))(flat_p)
```

This materializes the whole batched Jacobian at once and may trigger:

```
RESOURCE_EXHAUSTED: Out of memory
```

Empirical motivation

1+1 D Heat equation

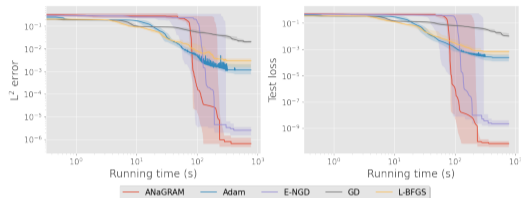


Figure: Performance comparison w.r.t running time for Heat equation in 1+1 D:

$$\begin{cases} \partial_t u - \frac{1}{4} \partial_{xx} u = 0 & \text{in } [0, 1]^2 \\ u = 0 & \text{on } [0, 1] \times \{0, 1\} \\ u = \sin(\pi x) & \text{on } \{0\} \times [0, 1] \end{cases} \quad \begin{cases} \partial_t u - 10^{-3} \partial_{xx} u = 5(u - u^3) & \text{in } \Omega = [0, 1] \times [-1, 1] \\ u = -1 & \text{on } \partial\Omega_b = [0, 1] \times \{-1, 1\} \\ u(0, x) = x^2 \cos(\pi x) & \text{on } \partial\Omega_0 = \{0\} \times [-1, 1] \end{cases}$$

1+1 D Allen-Cahn equation

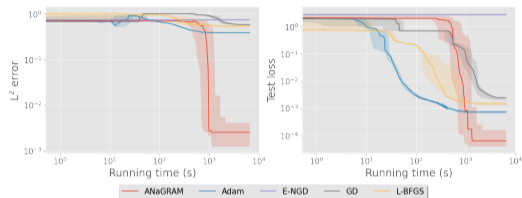


Figure: Performance comparison w.r.t running time for Allen-Cahn equation in 1+1 D:

Note: E-NGD refers to Müller and Zeinhofer2023). ANaGRAM to Schwencke and Furtlehner2025).

Experiment	Train Loss		L_2 Error	
	AMStramGRAM †	ANaGRAM	AMStramGRAM †	ANaGRAM
Heat Equation	6.29e-29 ± 6.78e-30	8.56e-11 ± 7.05e-11	2.32e-14 ± 1.14e-14	1.28e-06 ± 1.75e-06
Laplace 2D	1.46e-28 ± 1.87e-29	4.27e-13 ± 4.66e-13	2.24e-15 ± 2.52e-16	3.49e-09 ± 3.58e-09
Laplace 5D	2.04e-08 ± 1.16e-08	6.37e-08 ± 7.01e-08	2.12e-05 ± 8.15e-06	4.00e-05 ± 2.93e-05
Allen–Cahn	3.19e-11 ± 2.37e-11	2.19e-04 ± 4.16e-04	5.87e-05 ± 6.25e-06	4.32e-03 ± 5.93e-03

Experiment	Train Loss		L_2 Error	
	AMStramGRAM †	SSBroyden *	AMStramGRAM †	SSBroyden *
Burgers (1+1 D)	2.99e-12 ± 9.26e-13	2.92e-10 ± 1.45e-10	1.5e-06 ± 9.43e-7	1.59e-06 ± 1.02e-6
Non-Linear Poisson	8.51e-24 ± 2.24e-24	3.03e-16 ± 3.82e-16	6.81e-10 ± 1.41e-09	9.29e-12 ± 5.85e-12
Allen–Cahn	3.19e-11 ± 2.37e-11	6.42e-12 ± 5.52e-12	5.87e-05 ± 6.25e-06	3.94e-06 ± 1.72e-06

* refers to the order two method of Urbán et al.2025), with adaptive sampling and hard constraint enforcement on boundary conditions.

† refers to Schwencke et al.2025)

- ① Always think from a functional perspective (matches both the theory and effective implementation)
- ② Defining abstract coding concepts (`Pmodel`, `Operator`, `Field`) helps the implementation.
- ③ AD is an efficient way of defining the differential operators.
- ④ Always vectorize!
- ⑤ Use Natural Gradient!

Not covered (see practical session):

- regularization
- linesearch

Thank you for your attention ! Questions welcome.



- AMARI, S.-I. AND S. C. DOUGLAS (1998): “Why Natural Gradient?” in *Proceedings of the 1998 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP’98 (Cat. No. 98CH36181)*, IEEE, vol. 2, 1213–1216.
- MÜLLER, J. AND M. ZEINHOFER (2023): “Achieving High Accuracy with PINNs via Energy Natural Gradient Descent,” in *International Conference on Machine Learning*, PMLR, 25471–25485.
- SCHWENCKE, N. AND C. FURTLERHNER (2025): “ANaGRAM: A Natural Gradient Relative to Adapted Model for Efficient PINNs Learning,” in *The Thirteenth International Conference on Learning Representations*.
- SCHWENCKE, N., C. ROUSSELOT, A. SHILOVA, AND C. FURTLERHNER (2025): “AMStramGRAM: Adaptive Multi-Cutoff Strategy Modification for ANaGRAM,” in *arXiv Preprint*, arXiv Preprint.
- URBÁN, J. F., P. STEFANOU, AND J. A. PONS (2025): “Unveiling the Optimization Process of Physics Informed Neural Networks: How Accurate and Competitive Can PINNs Be?” *Journal of Computational Physics*, 523, 113656.