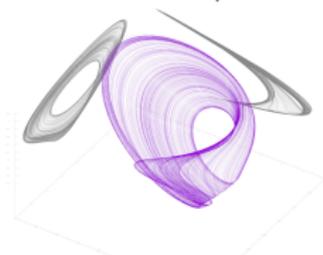


Deep learning for Physics

Lionel MATHELIN
LISN-CNRS

Dataflot / Decipher



M2 MVA 2020-2021

université
PARIS-SACLAY

1

Deep learning for Physics

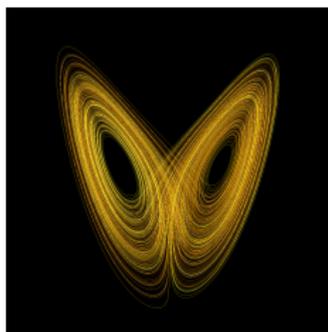
- Considerations on the size of the training set
- How to sample? Which observations?
- What to learn?
- How to learn? Structure / functional form to use
- Observability
- A continuous viewpoint: Neural ODE
- Field reconstruction

See how (deep) learning can leverage Physics/applications bottlenecks (online computational time and memory, modeling capability (super regression), etc.).

But also how Physics can help learning (bounds, first principles, speed-up learning, etc.).

A minimally complex system model

Illustration on Lorenz '63 attractor as a proxy good enough for its features.



Simplified mathematical model for atmospheric convection:

$$\begin{aligned}\dot{x} &= \sigma(y - x), && \text{convection rate,} \\ \dot{y} &= x(\rho - z) - y, && \text{horizontal temperature variations,} \\ \dot{z} &= xy - \beta z && \text{vertical temperature variations.}\end{aligned}$$

Becomes chaotic for $\rho > \sigma(\sigma + \beta + 3) / (\sigma - \beta - 1)$.

Typically: $\sigma = 10$ (Prandtl number), $\beta = 8/3$, $\rho = 28$ (Rayleigh number). Most “esthetic” values among the ones leading to chaos ...

System of nonlinear deterministic ODEs.

One wants to learn to predict from observations:

$$\dot{\mathbf{u}} = \mathbf{F}(\mathbf{u}) \quad \text{or} \quad \mathbf{u}_{t+\Delta t} = \mathbf{F}_{\Delta t}(\mathbf{u}_t).$$

Quality of the Information

Get to know your “opponent”!

How much information is necessary to describe the system/dynamics? How big should the training set be? How to design it?

Quality of the Information

Get to know your “opponent”!

How much information is necessary to describe the system/dynamics? How big should the training set be? How to design it?

Different definitions of dimension

The geometry of chaotic attractors can be complex and difficult to describe.

What is the dimension of the attractor of Lorenz'63? In other words, what is the intrinsic dimension of its data manifold?

It is therefore useful to have quantitative characterizations of such geometrical objects.

- Topological dimension \rightarrow integer,
- Metric dimension \rightarrow fractional.

Minimum number of independent dynamic variables to define the state of the system (nb degrees of liberty).

Many definitions of dimensions

- **Box Counting Dimension**

If the attractor exists in a d -dimensional phase space (d is necessarily an integer). Let $N(\epsilon)$ be the number of d -dimensional cubes of edge length ϵ from the grid that are needed to cover the attractor. Then

$$N \propto 1/\epsilon^d, \quad \Rightarrow \quad d_0 = \lim_{\epsilon \rightarrow 0_+} \frac{\log N(\epsilon)}{\log 1/\epsilon}.$$

Many definitions of dimensions

- **Box Counting Dimension**

If the attractor exists in a d -dimensional phase space (d is necessarily an integer). Let $N(\epsilon)$ be the number of d -dimensional cubes of edge length ϵ from the grid that are needed to cover the attractor. Then

$$N \propto 1/\epsilon^d, \quad \Rightarrow \quad d_0 = \lim_{\epsilon \rightarrow 0^+} \frac{\log N(\epsilon)}{\log 1/\epsilon}.$$

- **Rényi dimension**

Also called the ‘generalized dimension’, takes into account the frequency with which cubes are visited via weighting them according to their **natural measure** μ , defined as the fraction of the time that the long orbit on the attractor spends in any given ϵ -region C_j of state space. Accounts for time spent in cubes.

$$d_q = \lim_{\epsilon \rightarrow 0^+} \frac{1}{1-q} \frac{\log (\sum_j (\mu(C_j))^q)}{\log(1/\epsilon)}.$$

Special cases: d_0 : ‘box counting’ or ‘capacity dimension’, d_1 is ‘information dimension’, d_2 has been called the ‘correlation dimension’.

$$d_2 \leq d_1 \leq d_0.$$

Dimension d_2

d_2 has a nice numerical algorithm for its computation (Grassberger & Procaccia) as the slope of the following function when $\epsilon \rightarrow 0$: [Also a convenient algo (faster, less noisy with few points)]:

$$d_2 = \lim_{\epsilon \rightarrow 0_+} \frac{\log(1/N^2 \sum_{i \neq j} H(\epsilon - \|x_i - x_j\|))}{\log \epsilon} \quad \text{with } H \text{ the Heaviside function.}$$

Numerator is ratio of samples such that their distances is $< \epsilon$.

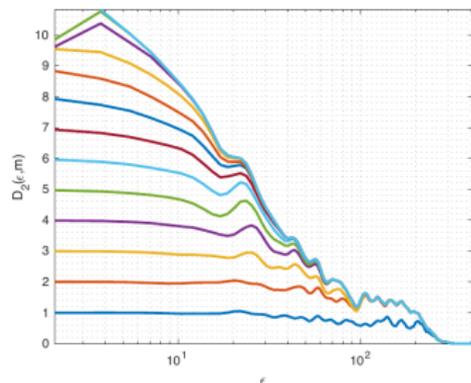
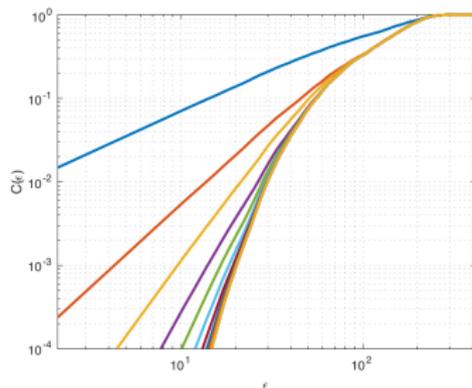


Figure: Illustration of the estimation of the d_2 dimension by varying the threshold ϵ .

Kaplan-Yorke dimension

It is defined in terms of the **Lyapunov exponents**.

Illustrative example in 1-D

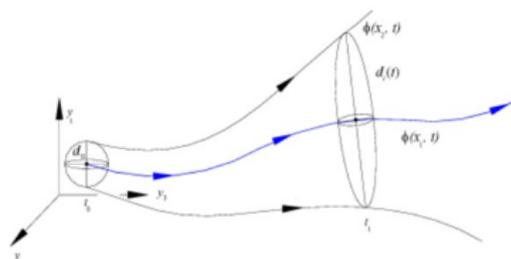
$$u_n = f(u_{n-1}),$$

$$\epsilon_n = f(u_{n-1} + \epsilon_{n-1}) - f(u_{n-1}), \quad \left| \frac{\epsilon_n}{\epsilon_0} \right| = \prod_i \left| \frac{\epsilon_i}{\epsilon_{i-1}} \right| = \prod_i |f'(u_{i-1})|,$$

$$\epsilon_n = \exp(\lambda n) \epsilon_0 \quad \Rightarrow \quad \lambda = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_i^n \ln |f'(u_{i-1})|.$$

This extends to dynamical systems via singular values of the mean Jacobian of the flow $\phi^t : \mathbf{u}_0 \rightarrow \phi^t(\mathbf{u}_0) = \mathbf{u}_t$.

Lyapunov Exponent



$$\bullet \quad |\delta \mathbf{Z}(t)| = e^{\lambda t} |\delta \mathbf{Z}(0)|$$

Kaplan-Yorke dimension

For Lorenz'63: $\lambda = (0.90, 0.00, -14.57)$

The 'Lyapunov dimension' d_L or d_{KY} is given by

$$d_{KY} = k + \frac{\lambda_1 + \lambda_2 + \dots + \lambda_k}{|\lambda_{k+1}|},$$

where k is the maximum value of i such that $\lambda_1 + \dots + \lambda_i > 0$, $\lambda_1 \geq \lambda_2 \geq \dots$

The Kaplan-Yorke conjecture states that $d_1 = d_{KY}$ for 'typical' systems, that is, the information dimension is equal to the Lyapunov dimension.

This relationship is remarkable in that it relates dynamics (Lyapunov exponents) to attractor geometry and natural measure (d_1).

Hausdorff dimension

Let $N(r)$ be the minimal number of open balls of radius r necessary to cover a compact metric space X . If, when $r \rightarrow 0$, $N(r)$ grows like $\frac{1}{r^d}$, the space X is said of dimension d .

More precisely, d is the number such that $N(r)r^s \rightarrow 0$ if $s > d$, and $N(r)r^s \rightarrow \infty$ if $s < d$:

$$H^s(X) = \lim_{r \rightarrow 0} \inf_{\text{diam}(A_i) < r} \left\{ \sum_{i=1}^{\infty} \text{diam}(A_i)^s \mid X \subseteq \bigcup_{i=1}^{\infty} A_i \right\},$$
$$d_H(X) = \inf \{s \mid H^s(X) = 0\} = \sup \{s \mid H^s(X) = \infty\}.$$

For Lorenz'63: $d_2 = 2.05 \pm 0.01$, $d_H = 2.06 \pm 0.01$.

Considerations on the size of the training set

Ergodic theory

Founding idea is that the long-term statistics of a system are equivalently described in terms of invariant measure (independent of time), μ .

The Poincaré recurrence theorem guarantees that trajectories exiting a region $S \in \Omega$ of non-vanishing measure will return to it infinitely many times.

Let S be an ensemble including a point \mathbf{u}_0 of the attractor. For an **ergodic** system, the **mean recurrence time** for a given S is proportional to the measure of S :

$$\langle \tau_{S(\mathbf{u}_0)} \rangle_{\mu_S} \propto \frac{1}{\mu(S)} \sim \epsilon^{-d_A}, \quad \text{with} \quad S = B_{\mathbf{u}_0}^{d_{\text{ambient}}}(\epsilon).$$

- the mean recurrence time then becomes large when ϵ is small (# **accuracy**) and/or d_A is large: **“Curse of dimensionality”**.
- If one wants a uniform learning accuracy of order ϵ (**sampling of the natural measure**), the number of samples to be considered is then $N \sim \epsilon^{-d_A}$.

One could use a less severe bound, e.g., $N \gtrsim (1/\epsilon)^{d/2}$ suggested by Eckmann & Ruelle (1992).

Black-box approach

Structure-based consideration

If the structure was known, need to identify $\{\sigma, \beta, \rho\}$ only $\rightarrow N \geq 3$ (noiseless)!

Physics/expertise is in the structure, not the data!

White-box approach

A grey-box approach

Size of the training set is partially conditioned by the minimum complexity of

- the physical system at hand ($\sim \epsilon^{-d_A}$),
- the structure of the retained model ($\sim \#$ parameters of the network).

Choosing a model structure considerably restricts (the dimension of) the solution space and reduces the required size of the training set.

How to sample? Which observations?

A numerical experiment

One training set, three different sampling strategies:

- 1 long trajectory on the attractor, $N = 27,000$ points.
- 9 short trajectories on the attractor, $N = 9 \times 3000$.
- 3*3 short trajectories ($N = 3 \times 3 \times 3000$) initialized around the unstable fixed points: $(0, 0, 0)$,
 $(-\sqrt{\beta(\rho - 1)}, -\sqrt{\beta(\rho - 1)}, \rho - 1)$ and $(\sqrt{\beta(\rho - 1)}, \sqrt{\beta(\rho - 1)}, \rho - 1)$.
Perturbations along the three orthogonal directions of the Jacobian.

A numerical experiment

One training set, three different sampling strategies:

- 1 long trajectory on the attractor, $N = 27,000$ points.
- 9 short trajectories on the attractor, $N = 9 \times 3000$.
- 3×3 short trajectories ($N = 3 \times 3 \times 3000$) initialized around the unstable fixed points: $(0, 0, 0)$,

$$\left(-\sqrt{\beta(\rho-1)}, -\sqrt{\beta(\rho-1)}, \rho-1\right) \text{ and } \left(\sqrt{\beta(\rho-1)}, \sqrt{\beta(\rho-1)}, \rho-1\right).$$

Perturbations along the three orthogonal directions of the Jacobian.

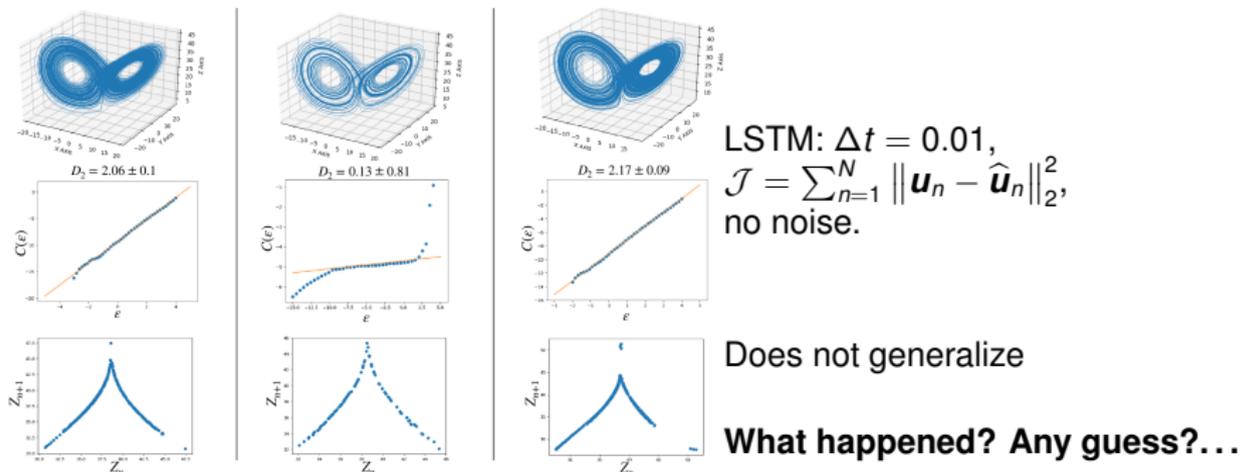


Figure: Performance of different sampling strategies.

A numerical experiment

Several viewpoints:

- The predictability of the system is better close to the equilibrium points (low entropy, e.g. Kolmogorov-Sinai entropy).
- The loss function is well suited to learn close to these equilibrium points (non-chaotic dynamic). No need to further regularize the cost landscape [Scheme smooth (equilibrium points) vs rough (attractor)].
- Curriculum learning: first learn $\dot{x} = \alpha x$, then the rest.

What to learn?

General considerations on the loss

Crucial to ask the right question (objective) to hope getting a good answer (model), i.e., a model relevant for our application.

General considerations on the loss

Crucial to ask the right question (objective) to hope getting a good answer (model), i.e., a model relevant for our application.

Other problem: ease the learning.

The loss function induces the cost landscape and conditions how easy it is to get a good solution.

Non-chaotic system (near equilibrium points) \rightarrow the cost landscape \mathcal{J} as a function of the model parameters and/or of the initial condition is smooth.

However, opposite for a chaotic system and the output (loss function) does not exhibit a smooth dependence with respect to inputs \rightarrow learning via gradient is difficult, high sensitivity.

Workarounds

- **Restrict the solution space** (parameters to be learned) by penalizing irrelevant subspaces. For instance, by penalizing the deviation from the Jacobian:

$$\mathcal{J} = \sum_{n=1}^N \|\mathbf{u}_n - \hat{\mathbf{u}}_n\|_2^2 + \alpha \sum_{n=1}^N \|\nabla \mathbf{u}_n - \nabla \hat{\mathbf{u}}_n\|_2^2.$$

Here: $L^2 \rightarrow H^1$. This can be generalized to higher order losses.

\rightarrow promotes couplings (subspaces) between parameters and hence reduces the variance of the loss as a function of variations of the parameters \Rightarrow smoother landscape.

Making it practical (no need for $\nabla \mathbf{u}$):

$$\mathcal{J} = \sum_{n=1}^N \|\mathbf{u}_n - \hat{\mathbf{u}}_n\|_2^2 + \alpha \sum_{n=1}^N \|\nabla \hat{\mathbf{u}}_n\|_2^2.$$

Should fit \mathbf{u} without going crazy ($\|\nabla \hat{\mathbf{u}}_n\|$).

Workarounds

- **Restrict the solution space** (parameters to be learned) by penalizing irrelevant subspaces. For instance, by penalizing the deviation from the Jacobian:

$$\mathcal{J} = \sum_{n=1}^N \|\mathbf{u}_n - \hat{\mathbf{u}}_n\|_2^2 + \alpha \sum_{n=1}^N \|\nabla \mathbf{u}_n - \nabla \hat{\mathbf{u}}_n\|_2^2.$$

Here: $L^2 \rightarrow H^1$. This can be generalized to higher order losses.

\rightarrow promotes couplings (subspaces) between parameters and hence reduces the variance of the loss as a function of variations of the parameters \Rightarrow smoother landscape.

Making it practical (no need for $\nabla \mathbf{u}$):

$$\mathcal{J} = \sum_{n=1}^N \|\mathbf{u}_n - \hat{\mathbf{u}}_n\|_2^2 + \alpha \sum_{n=1}^N \|\nabla \hat{\mathbf{u}}_n\|_2^2.$$

Should fit \mathbf{u} without going crazy ($\|\nabla \hat{\mathbf{u}}_n\|$).

- **Penalizing “accelerations”**. The L^2 loss does not allow to penalize zero-mean accelerations.

Workarounds (2)

- Irregular sampling or more realistic loss (Ordinary Least Squares vs Total Least Squares [SCHEME]): DTW.

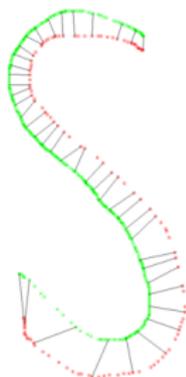


Figure: Sketch of the alignment of two irregularly sampled time-series via DTW.

Dynamic Time Warping (DTW) computes the best possible alignment between two discrete time series of respective length N and M by computing first the $N \times M$ pairwise distance matrix between these points to solve then a dynamic program (DP) using Bellman's recursion with a quadratic cost NM .

→ Robust to shifts, dilatation, etc. It is based on optimal transport.

Soft DTW makes this loss differentiable and efficient to optimize (gradient-based technique).

Physics-motivated losses

Physics-agnostic machine learned systems are fragile, for instance wrt noisy, small training datasets.

Example with stability

If one knows the system is stable \Rightarrow to be enforced during learning.

Several definitions of stability:

- Stability in machine learning: robustness wrt perturbations in data or hyperparameters.
- Stability in Physics: given \mathbf{u}_0 , the system remains within a ball of size $r \forall t > 0$ around \mathbf{u}_0 .
- Stability: $\forall \epsilon > 0, \exists \delta > 0 \setminus |\mathbf{u}_0| < \delta \Rightarrow \left| \mathbf{f}^k(\mathbf{u}_0) \right| < \epsilon, k = 1, 2, \dots$

Physics-motivated losses — Stability

Difficult to verify numerically (costly) \Rightarrow Stability in the sense of Lyapunov.

$\mathbf{u}_0 = \mathbf{0}$ is stable in the sense of Lyapunov if it exists a continuous function $\mathcal{V} : \mathbb{R}^d \rightarrow \mathbb{R}$ such that:

$$\begin{aligned}\mathcal{V}(\mathbf{0}) = 0, \quad \mathcal{V}(\mathbf{u}) > 0, \quad \forall \mathbf{u} \neq \mathbf{0}, \\ \mathcal{V}(\mathbf{f}(\mathbf{u})) - \mathcal{V}(\mathbf{u}) \leq 0, \quad \forall \mathbf{u} \neq \mathbf{0}, \quad |\mathbf{u}| < r.\end{aligned}$$

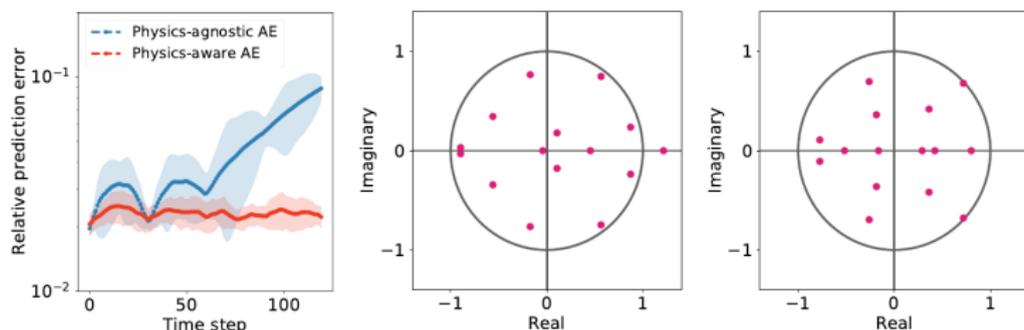
One can interpret \mathcal{V} as a generalized energy function as it is often related to the total energy of the system. Dissipative system: \mathcal{V} will decrease along trajectories.

Physics-motivated losses — Stability

If governing equations are linear, $f(\mathbf{u})$ reduces to a matrix \Rightarrow study of the eigenvalues. Stability if modulus always ≤ 1 . Penalization of overflows in the loss.

One wants to impose these rules to the learned model: form of regularization.

Example with learning the (stable) flow around a circular cylinder via an auto-encoder, Erichson *et al.* (2019).



(A) With LR $1e-2$ and WD $1e-6$. (B) Physics-agnostic model in (a). (C) Physics-aware model in (a).

Figure: Performance and stability properties of Physics-agnostic and Physics-aware learning.

Illustration with conservation laws

Expertise on Physics allows us to know there exist many conservation laws.

Allows to further reduce the size of the search space. These conservation laws are counterparts of invariances with respect to some symmetry groups.

Illustration with conservation laws

Expertise on Physics allows us to know there exist many conservation laws.

Allows to further reduce the size of the search space. These conservation laws are counterparts of invariances with respect to some symmetry groups.

Noether theorem

To every differentiable symmetry generated by local actions there corresponds a conserved current.

Each “invariance” illustrates the fact that physical laws do not change when an experiment undergoes the associated transformation, hence there is no absolute reference to carry it out.

Example: invariance wrt spatial translation \rightarrow conservation of momentum $\mathbf{p} = m \mathbf{v}$.
Invariance in time \rightarrow conservation of total energy.

Imposing constraints in the form of Partial Differential Equations

One knows that the underlying physical system typically obeys certain equations, e.g. Navier-Stokes. Introduction in the loss definition.

Example with Burgers, Raissi *et al.* (2017):

$$\begin{aligned}\frac{\partial \mathbf{u}}{\partial t} &= \alpha \mathbf{u} \cdot \frac{\partial \mathbf{u}}{\partial \mathbf{x}} + \beta \Delta \mathbf{u}, \\ \mathcal{J} &= \|\mathbf{u} - \hat{\mathbf{u}}\|_2^2 + \lambda \|\mathbf{r}\|_2^2, \\ \mathbf{r} &= \frac{\partial \mathbf{u}}{\partial t} - \alpha \mathbf{u} \cdot \frac{\partial \mathbf{u}}{\partial \mathbf{x}} - \beta \Delta \mathbf{u}.\end{aligned}$$

Spatial gradient term is easy to get from NN.

This loss mimics $\|\mathbf{u} - \hat{\mathbf{u}}\|_2^2 + \lambda \|\dot{\mathbf{u}} - \hat{\dot{\mathbf{u}}}\|_2^2$.

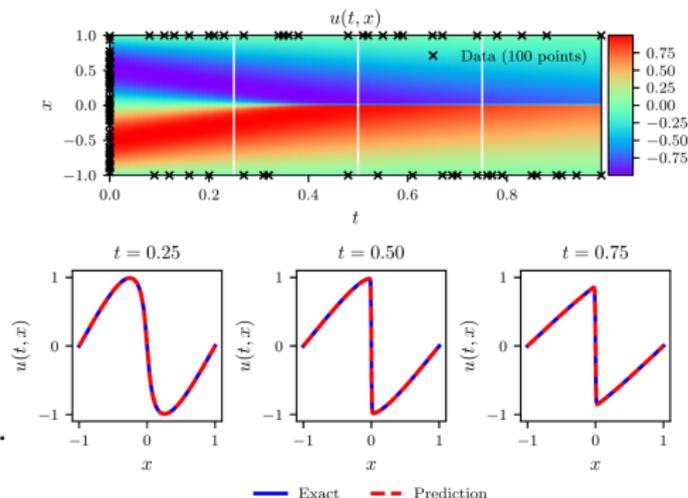


Figure: Reconstructed solution.

$$\begin{aligned}\frac{\partial \mathbf{u}(t, \mathbf{x})}{\partial t} + N(\mathbf{u}(t, \mathbf{x})) &= 0, & (t, \mathbf{x}) \in [0, T] \times \Omega, \\ \mathbf{u}(t=0, \mathbf{x}) &= \mathbf{u}_0(\mathbf{x}), \\ \mathbf{u}(t, \mathbf{x}) &= \mathbf{g}(t, \mathbf{x}), & \mathbf{x} \in \partial\Omega.\end{aligned}$$

Deep neural network approximates $\mathbf{u}(t, \mathbf{x})$ from just (t, \mathbf{x}) . No need to know the true \mathbf{u} .

$$\mathcal{J} = \left\| \frac{\partial \mathbf{u}(t, \mathbf{x})}{\partial t} + N(\mathbf{u}(t, \mathbf{x})) \right\|_{[0, T] \times \Omega}^2 + \|\mathbf{u}(0, \mathbf{x}) - \mathbf{u}_0(\mathbf{x})\|_{\Omega}^2 + \|\mathbf{u}(t, \mathbf{x}) - \mathbf{g}(t, \mathbf{x})\|_{[0, T] \times \partial\Omega}^2$$

No need for a mesh. Training on randomly sampled points in space and time.

Operators such as $\frac{\partial}{\partial t}$ or N are easy to implement on top of the network providing \mathbf{u} . Norms are approximated with random points \rightarrow Stochastic gradient descent.

Predicting with the right variable, Ayed *et al.* (2019).

Aim: predict I_{k+1} from previous snapshots $\{I_k\}_k$. One knows that I obeys

$$\frac{\partial I(t, \mathbf{x})}{\partial t} + \mathbf{w}(t, \mathbf{x}) \cdot \nabla I(t, \mathbf{x}) = D \Delta I.$$

Initial conditions not known, nor the boundary conditions, nor the diffusion parameter D . Instead, the displacement field \mathbf{w} is learned from previous snapshots with the loss:

$$\mathcal{J} = \sum_{\mathbf{x} \in \Omega} \left| I_{t+1}(\mathbf{x}) - \hat{I}_{t+1}(\mathbf{x}) \right| + \lambda_1 (\nabla \cdot \hat{\mathbf{w}}_t(\mathbf{x}))^2 + \lambda_2 \|\hat{\mathbf{w}}_t(\mathbf{x})\|_2^2 + \lambda_3 \|\nabla \hat{\mathbf{w}}_t(\mathbf{x})\|_2^2.$$

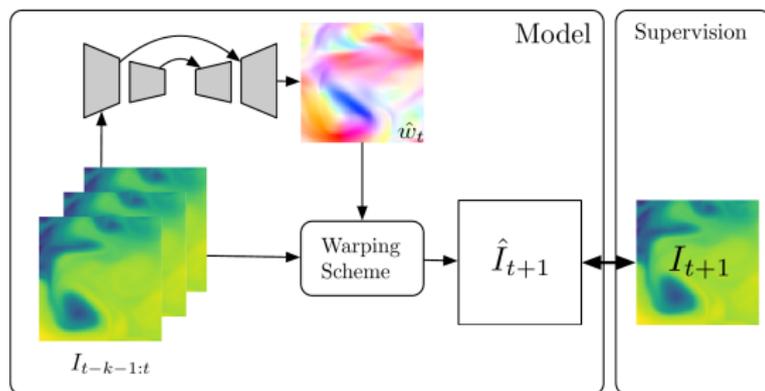


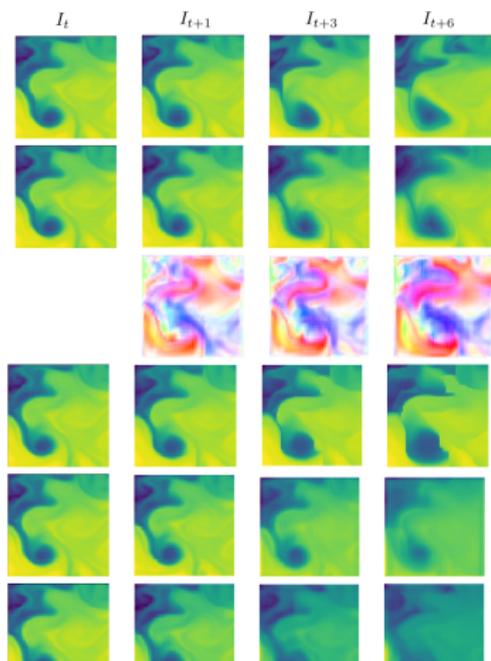
Figure: Sketch of the strategy to predict the next image \hat{I}_{t+1} .

Predicting with the right variable, Ayed *et al.* (2019).

Warping scheme:

$$\hat{l}_{t+1}(\mathbf{x}) = \sum_{\mathbf{y} \in \Omega} k(\mathbf{x} - \hat{\mathbf{w}}_t, \mathbf{y}) l_t(\mathbf{y}),$$

$$k(\mathbf{x} - \hat{\mathbf{w}}_t, \mathbf{y}) = \frac{1}{4\pi \hat{D} \Delta t} \exp\left(-\frac{1}{4\hat{D} \Delta t} \|\mathbf{x} - \hat{\mathbf{w}} - \mathbf{y}\|^2\right).$$



How to learn?

Structure / functional form to use

Determinism of the reconstruction dynamics

Partial Observability

The way a system is observed is often critical and conditions the architecture and the resulting performance.

Determining d : Trajectories in the phase space cannot cross. [Sketch]

Reconstruct a flow in dimension $d_r = 2$ then $d_r = 3$, etc., until the trajectory does not cross itself anymore.

The smallest dimension this happens is the **embedding dimension** d_p .

Embeddings for dynamical systems

We seek a model for a dynamical system \Rightarrow requires two ingredients:

- a coordinate system,
- a representation for the system evolution.

The first ingredient is obtained by an **embedding**: a (diffeomorphic) function g that maps the observed data y_t to a set of coordinates \tilde{x} according to $g(\{y_t\}_t) = (\tilde{x}_0, \tilde{x}_1, \dots)$. It preserves the topology and the invariants (Lyapunov exponents, ...).

We assume the system lies on a low dimensional attractor $x(t) \in \mathbb{A} \subset \mathbb{R}^d$, a smooth submanifold of \mathbb{R}^d and has a dimension $\dim(\mathbb{A}) < d$.

Then, **Whitney embedding theorem** says that if the number of coordinates $d_c > 2 \dim(\mathbb{A})$ then the dynamics of the state can be entirely captured by the new coordinate system.

Made practical by the **Takens embedding theorem** by asserting that one can take this number of delays.

Embeddings for dynamical systems

Three common mappings:

① **Delays:** $\tilde{x}_i = y(t - i\Delta t)$ for $i = 0, 1, \dots$ and $\Delta t \in \mathbb{R}_+$. $\tilde{\mathbf{x}} = \begin{bmatrix} \mathbf{y}(t_i) \\ \mathbf{y}(t_{i-1}) \\ \dots \end{bmatrix}$.

② **Derivatives:** $\tilde{x}_i = y^{(i)}$ for $i = 0, 1, \dots$ $\tilde{\mathbf{x}} = \begin{bmatrix} \mathbf{y} \\ \dot{\mathbf{y}} \\ \ddot{\mathbf{y}} \\ \dots \end{bmatrix}$.

③ **Principal components:** an SVD is used to first transform the space followed by a use of delays or derivatives. $H = \begin{bmatrix} \mathbf{y}(t_i) \\ \mathbf{y}(t_{i-1}) \\ \dots \end{bmatrix} = U \Sigma V^*$, $\tilde{\mathbf{x}} = \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \dots \end{bmatrix}$

Embeddings for dynamical systems

In practice, how to find the embedding dimension ? False neighbor algorithm to check if topology is preserved.

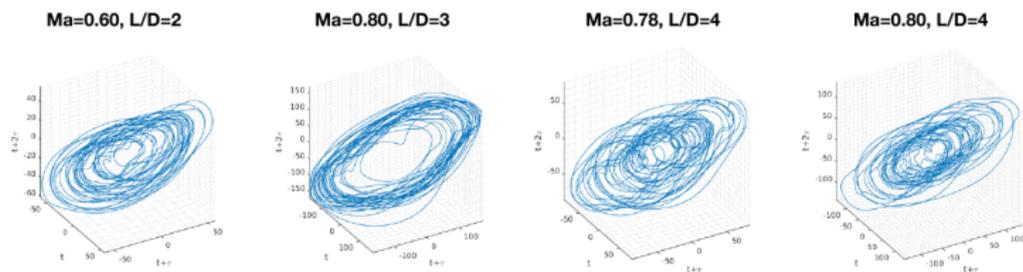


Figure: Different reconstructed phase spaces. They are topologically equivalent to the Lorenz canonical representation.

Observability

Observability

Observing the system through the measurement operator h :

- how much information can be extracted from it?
- what can be inferred on the underlying model?

For instance, how much of the state x and of its dynamics can be inferred from measurements $y(t)$?

⇒ **observability** issue (was extensively studied, in particular by the control community).

When both the model f and the observation operator h are **linear**, full observability of $x \in \mathbb{X} \subset \mathbb{R}^d$ from $y(t) \in \mathbb{R}^{n_y}$ is guaranteed whenever the Kalman rank condition is satisfied, i.e., the **observability matrix** \mathcal{O}_h is full rank $\forall x \in \mathbb{X}$:

$$\mathcal{O}_h(x) := \begin{bmatrix} h \\ hf \\ \vdots \\ hf^{d-1} \end{bmatrix}.$$

Nonlinear observability

Extends to the **nonlinear framework**, for instance considering Lie derivatives of h along the vector field f :

$$\mathcal{L}_f h(x) := \frac{\partial h}{\partial x} f(x) \quad \text{and} \quad \mathcal{L}_f^n h(x) := \mathcal{L}_f \left[\mathcal{L}_f^{n-1} h(x) \right], \quad \forall n \in \mathbb{N}, \quad \mathcal{L}_f^0 h(x) := h$$

The system governed by f is fully observable via h iff \mathcal{O}_h is full rank $\forall x \in \mathbb{X}$, with \mathcal{O}_h writing as:

$$y(t) = h(x(t)), \quad \dot{y}(t) = \frac{\partial h}{\partial x} \dot{x} \equiv \mathcal{L}_f h(x).$$

then

$$\mathcal{O}_h(x) = \begin{bmatrix} \frac{\partial \mathcal{L}_f^0 h(x)}{\partial x} \\ \frac{\partial \mathcal{L}_f^1 h(x)}{\partial x} \\ \vdots \\ \frac{\partial \mathcal{L}_f^{d-1} h(x)}{\partial x} \end{bmatrix} = \begin{bmatrix} y(t) \\ \dot{y}(t) \\ \vdots \\ y^{(d-1)}(t) \end{bmatrix}.$$

Degree of observability can be quantified based on the spectrum of the observability Gramian $\mathcal{O}_h^T \mathcal{O}_h$ via the observability index $\delta_{\mathcal{O}_h}$:

$$\delta_{\mathcal{O}_h}(x) := \frac{\lambda_{\min}(\mathcal{O}_h^T \mathcal{O}_h)}{\lambda_{\max}(\mathcal{O}_h^T \mathcal{O}_h)} \in [0, 1].$$

Choice of the observables

The above criteria are useful to quantify the relevance of the observables $y(t)$ to inform upon the state \mathbf{x} .

It can also serve to improve observability by modifying the observables so that $\delta_{\mathcal{O}_h}$ increases. \rightarrow pivotal for **sensor placement**.

While similarity transformations of coordinates do not affect the rank of the observability matrix, they can improve the observability index and provide a more robust information on the system.

Is memory necessary ?

State vector (LSTM useless) vs Partial observation (memory-capable network, LSTM, reservoir).

Alternative: embedding (delayed, differential, ODE) \rightarrow augmented input state.

A continuous viewpoint: Neural ODE

Neural ODE (ODENet), Chen *et al.* (2018)

One can interpret Resnets like Ordinary Differential Equations discretized with an Euler scheme:

$$\mathbf{u}_{t+\Delta t} = \mathbf{u}_t + \mathbf{f}(\mathbf{u}_t, t; \theta).$$

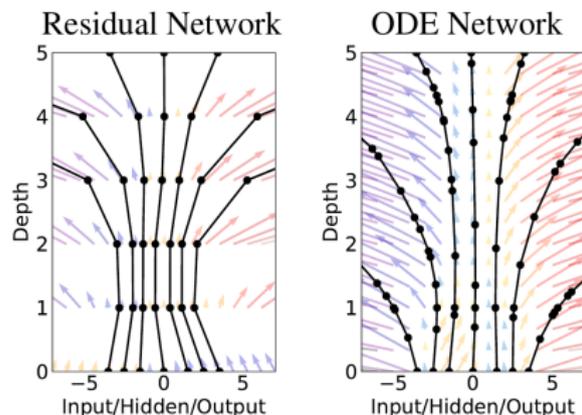
Neural ODE (ODENet), Chen *et al.* (2018)

One can interpret Resnets like Ordinary Differential Equations discretized with an Euler scheme:

$$\mathbf{u}_{t+\Delta t} = \mathbf{u}_t + \mathbf{f}(\mathbf{u}_t, t; \theta).$$

In the limit of more layers and smaller Δt :

$$\frac{d\mathbf{u}(t)}{dt} = \mathbf{f}(\mathbf{u}(t), t; \theta).$$



One can use an ODE solver to compute $\mathbf{u}(t)$ at any real continuous time $t > 0$.

More flexible since the dynamics might also depend on t , for instance using a neural network with two inputs: t and $\mathbf{u}(t)$.

Figure: In the limit of an infinite number of layers, a ResNet (left) can be interpreted as an ODE problem (right).

Neural ODE (ODENet), Chen *et al.* (2018)

Training via the adjoint Loss is defined at a final time $t = T$: $\mathcal{J}(\mathbf{u}(T))$

How to train / compute $\nabla_{\theta} \mathcal{J}$?

$$\begin{aligned}\mathcal{L} &:= J(\mathbf{u}_T) + \ll \lambda, \dot{\mathbf{u}}(t) + \mathbf{f}(\mathbf{u}; \theta) \gg, \\ &= J(\mathbf{u}_T) + \int_t \lambda(t) (\dot{\mathbf{u}}(t) + \mathbf{f}(\mathbf{u}; \theta)) dt, \\ \left\langle \frac{\partial \mathcal{L}}{\partial \mathbf{u}} \delta \mathbf{u} \right\rangle = 0 \forall \delta \mathbf{u} &= \nabla_{\mathbf{u}} J \delta \mathbf{u} + \int_t \delta \mathbf{u} \left(-\dot{\lambda} + \lambda \nabla_{\mathbf{u}} \mathbf{f} \right) dt + [\mathbf{u} \lambda]_0^T, \\ &= -\dot{\lambda} + \lambda \nabla_{\mathbf{u}} \mathbf{f}, \quad \text{with } \lambda_T = -\nabla_{\mathbf{u}} J|_T, \\ \left\langle \frac{\partial \mathcal{L}}{\partial \lambda} \delta \lambda \right\rangle = 0 \forall \delta \lambda &= \dot{\mathbf{u}}(t) + \mathbf{f}(\mathbf{u}; \theta), \\ \left\langle \frac{\partial \mathcal{L}}{\partial \theta} \delta \theta \right\rangle = 0 \forall \delta \theta &= \frac{\partial J}{\partial \theta} \delta \theta + \int_t \lambda \frac{\partial \mathbf{f}}{\partial \theta} \delta \theta dt, \quad \Rightarrow \frac{\partial J}{\partial \theta} = - \int_t \lambda \frac{\partial \mathbf{f}}{\partial \theta} dt.\end{aligned}$$

Both \mathbf{u} and λ are obtained via an ODE solver from 0 to T and from T to 0 respectively.

All this can be solved with a single call to an ODE solver considering the augmented state: $\mathbf{z} := \left[\mathbf{u}, \lambda, \frac{\partial J}{\partial \theta} \right]$.

Neural ODE (ODENet) — Example

Model and complement an **irregularly sampled** time series $(\mathbf{x}_0, \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$

- ▶ Difficult to do with recurrent networks.
- ▶ The dynamics learned / to be learned do not act directly on \mathbf{x}_t (observations) but on the (latent) internal state \mathbf{z}_t (state vector) [in order to be able to deal with partial or noisy information]

Neural ODE (ODENet) — Example

Model and complement an **irregularly sampled** time series ($\mathbf{x}_0, \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$)

- ▶ Difficult to do with recurrent networks.
- ▶ The dynamics learned / to be learned do not act directly on \mathbf{x}_t (observations) but on the (latent) internal state \mathbf{z}_t (state vector) [in order to be able to deal with partial or noisy information]

- 1 Transform the \mathbf{x}_t sequence into a guessed initial (probabilistic) internal state $q(\mathbf{z}_0 | \mathbf{x}_0, \dots, \mathbf{x}_N)$ using an RNN (going reverse in time, from \mathbf{x}_N to \mathbf{x}_0 ; the RNN has to be learned as well)
- 2 Apply the dynamics \mathbf{f}_θ previously learned to \mathbf{z}_0 using an ODE solver and produce guessed states $\mathbf{z}_{t>0}$
- 3 From $\mathbf{z}_{t>N}$, deduce missing observations $\mathbf{x}_{t>N}$.

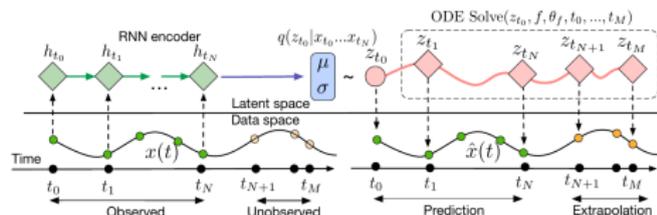
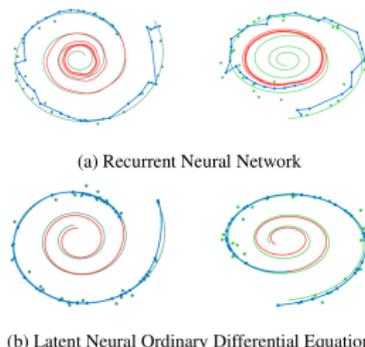


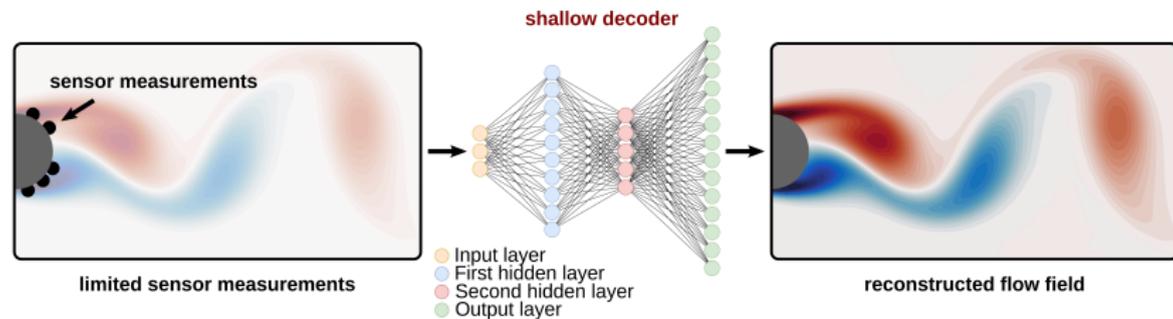
Figure: Sketch of the prediction strategy for future observations $\mathbf{x}_{t>N}$.



Field reconstruction

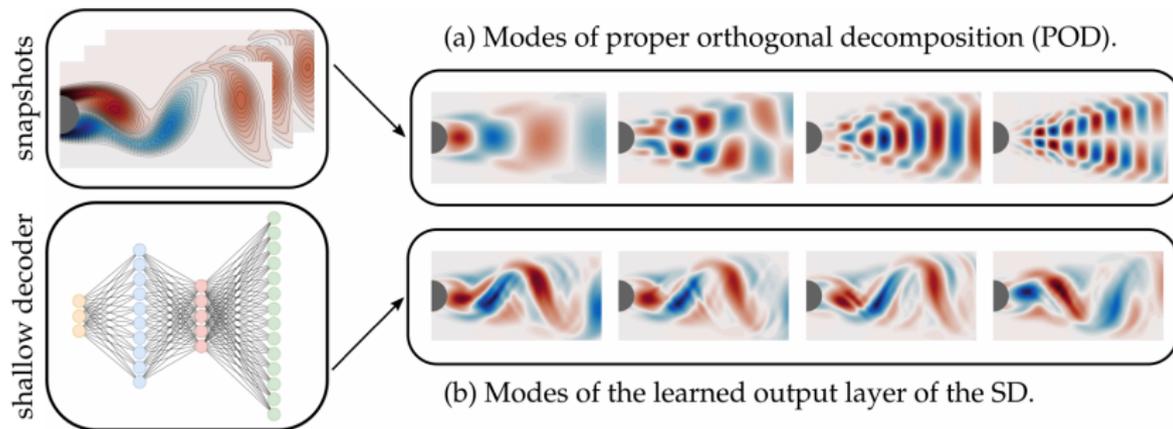
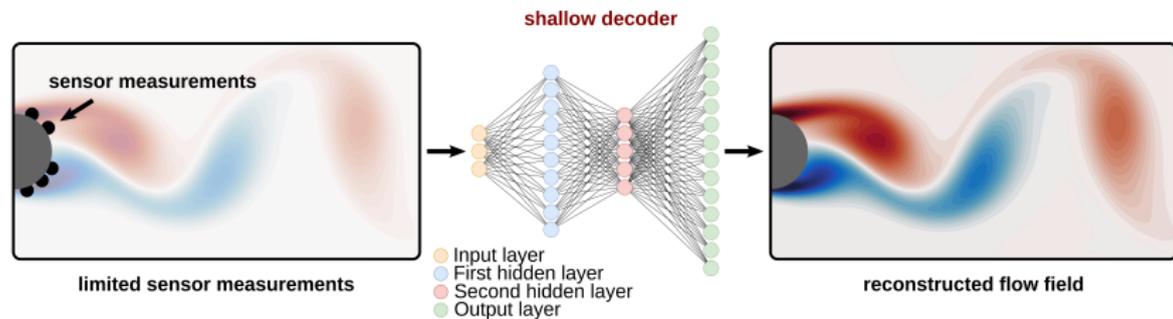
A shallow decoder for flow reconstruction, Erichson *et al.* (2020)

From sensors $\mathbf{s} \in \mathbb{R}^5 \rightarrow$ to field estimation $\hat{\mathbf{x}} \in \mathbb{R}^{78406}$

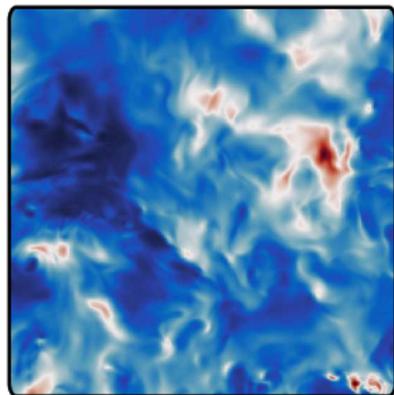


A shallow decoder for flow reconstruction, Erichson *et al.* (2020)

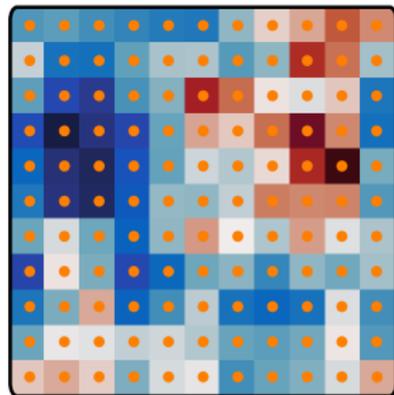
From sensors $\mathbf{s} \in \mathbb{R}^5 \rightarrow$ to field estimation $\hat{\mathbf{x}} \in \mathbb{R}^{78406}$



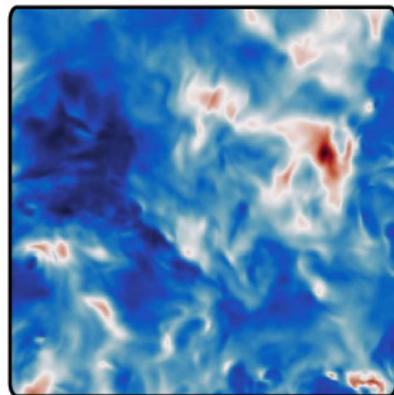
Super-resolution — Rise...



Snapshot

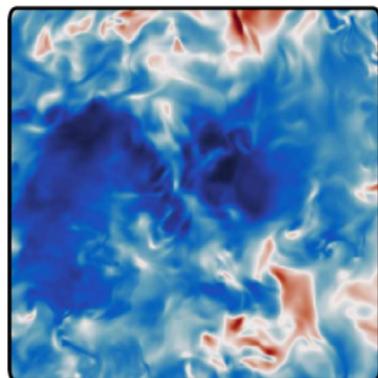


Low resolution

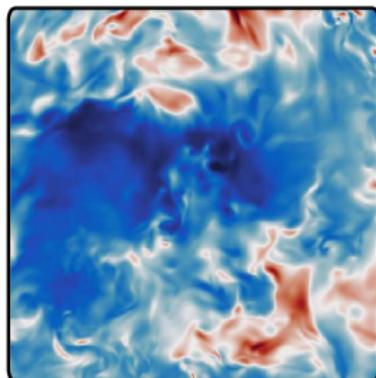


Shallow decoder

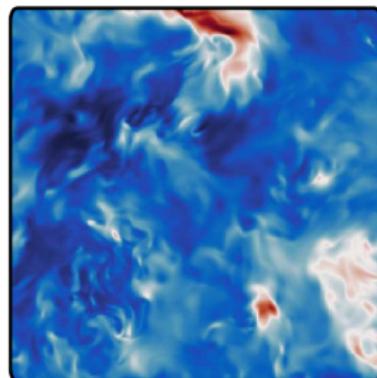
... and fall — On the requirement of satisfying sufficient statistics



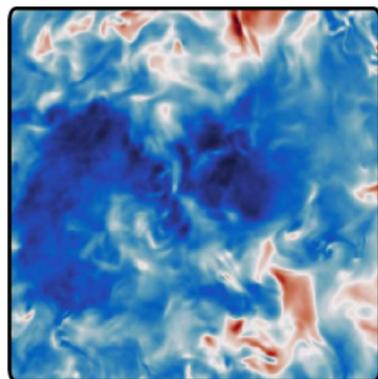
(a) Test snapshot $t = 1$



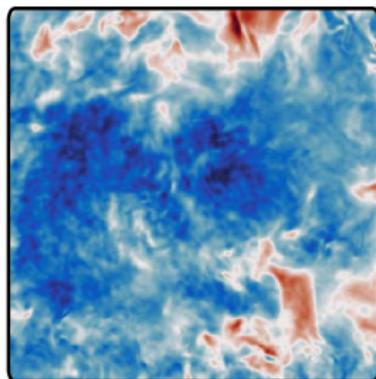
(c) Test snapshot $t = 20$



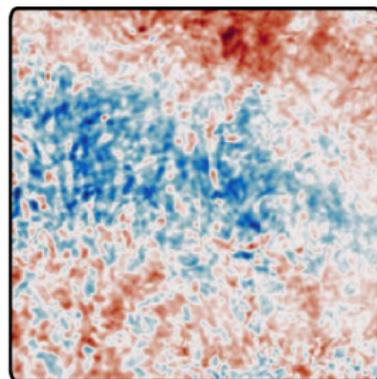
(e) Test snapshot $t = 50$



(b) Reconstruction



(d) Reconstruction



(f) Reconstruction

High-dimensional boundary value problem

$$\begin{aligned} \frac{\partial u}{\partial t} - \nabla (a \nabla u) + b \nabla u + cu + d &= 0, && \text{in } \Omega \times [0, T], \\ u(\mathbf{x}, t) &= g(\mathbf{x}, t), && \text{on } \partial\Omega \times [0, T], \\ u(\mathbf{x}, 0) &= h(\mathbf{x}), && \text{on } \Omega. \end{aligned}$$

Weak solution of PDEs — Zang *et al.* (2020)

High-dimensional boundary value problem

$$\begin{aligned} \frac{\partial u}{\partial t} - \nabla (a \nabla u) + b \nabla u + cu + d &= 0, & \text{in } \Omega \times [0, T], \\ u(\mathbf{x}, t) &= g(\mathbf{x}, t), & \text{on } \partial\Omega \times [0, T], \\ u(\mathbf{x}, 0) &= h(\mathbf{x}), & \text{on } \Omega. \end{aligned}$$

Weak formulation

For time-independent problems, *for instance*

$$\langle A[u], \varphi \rangle := \int_{\Omega} a \nabla u \nabla \varphi + b \nabla u \varphi + cu \varphi + d \varphi \, d\mathbf{x} = 0, \quad \forall \varphi \in H_0^1.$$

Approximation basis for the solution is given by a neural network \rightarrow

$$u(\mathbf{x}, t) \approx u_{\theta}(\mathbf{x}, t) \in H^1(\Omega; \mathbb{R}).$$

Test functions $\{\varphi\} \in H_0^1(\Omega; \mathbb{R})$

The choice of the projection $\langle \cdot \rangle$ is *yours*.

Induced operator norm minimization

Let $A[u] : H_0^1(\Omega) \rightarrow \mathbb{R}$ such that $A[u](\varphi) := \langle A[u], \varphi \rangle$.

Its operator norm is

$$\|A[u]\|_{\text{op}} := \max \left\{ \langle A[u], \varphi \rangle / \|\varphi\|_2 \mid \varphi \in H_0^1, \varphi \neq 0 \right\}.$$

u is a weak solution of the PDE iff $\|A[u]\|_{\text{op}} = 0$.

Induced operator norm minimization

Let $A[u] : H_0^1(\Omega) \rightarrow \mathbb{R}$ such that $A[u](\varphi) := \langle A[u], \varphi \rangle$.

Its operator norm is

$$\|A[u]\|_{\text{op}} := \max \left\{ \langle A[u], \varphi \rangle / \|\varphi\|_2 \mid \varphi \in H_0^1, \varphi \neq 0 \right\}.$$

u is a weak solution of the PDE iff $\|A[u]\|_{\text{op}} = 0$.

Hence, a formulation for an approximation of the weak solution u_θ could be

$$\min_{u_\theta \in H^1} \|A[u_\theta]\|_{\text{op}}^2 \iff \min_{u_\theta \in H^1} \max_{\varphi \in H_0^1} |\langle A[u_\theta], \varphi \rangle|^2 / \|\varphi\|_2^2.$$

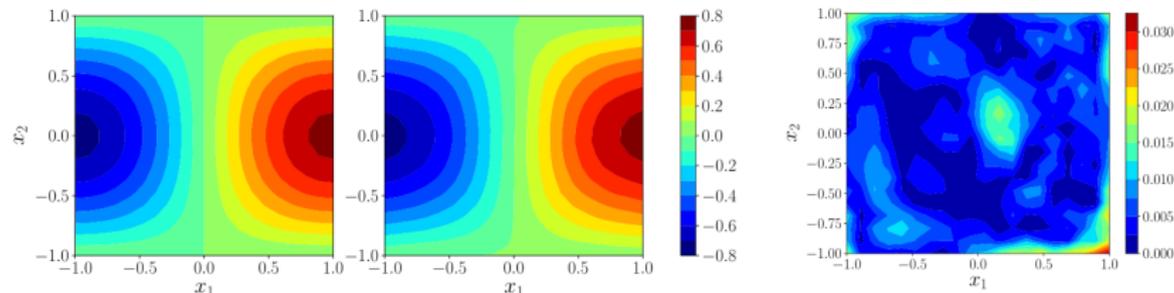
- A deep neural network parameterized by $\theta \rightarrow u_\theta$,
- An adversarial network parameterized by η tending to maximize $\langle A[u_\theta], \varphi_\eta \rangle \rightarrow \varphi_\eta$,

Example with a parabolic equation

Nonlinear reaction-diffusion equation

$$\begin{aligned} \frac{\partial u}{\partial t} - \Delta u - u^2 &= 0, & \text{in } \Omega \times [0, T], \\ u(\mathbf{x}, t) &= g(\mathbf{x}, t), & \text{on } \partial\Omega \times [0, T], \\ u(\mathbf{x}, 0) &= h(\mathbf{x}), & \text{on } \Omega. \end{aligned}$$

Crank-Nicolson time scheme. 6 hidden layers, 40 neurons per layer. Trained on collocation points in Ω .



(a) True $u^*(x, T)$ (left) vs estimated $u_\theta(x, T)$ (right)

(b) $|u_\theta(x, T) - u^*(x, T)|$

- AYED IBRAHIM, DE BÉZENAC EMMANUEL, PAJOT ARTHUR, BRAJARD JULIEN & GALLINARI PATRICK, Learning Dynamical Systems from Partial Observations, *arXiv e-prints*, p. arXiv:1902.11136, 2019.
- CHEN RICKY T. Q., RUBANOVA YULIA, BETTENCOURT JESSE & DUVENAUD DAVID, Neural Ordinary Differential Equations, *arXiv e-prints*, p. arXiv:1806.07366, 2018.
- ERICHSON N.B., MUEHLEBACH M. & MAHONEY M.W., 2019 Physics-informed autoencoders for Lyapunov-stable fluid flow prediction. Preprint: arXiv:1905.10866.
- ERICHSON N. BENJAMIN, MATHÉLIN LIONEL, YAO ZHEWEI, BRUNTON STEVEN, MAHONEY MICHAEL & KUTZ J. NATHAN, Shallow neural networks for fluid flow reconstruction with limited sensors, *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, **476** (2238), p. 20200097, 2020.
- RAISSI M., PERDIKARIS P. & KARNIADAKIS G.EM, 2017 Physics Informed Deep Learning (Part I): Data-driven Solutions of Nonlinear Partial Differential Equations. Preprint: arXiv:1711.10561.
- SIRIGNANO JUSTIN & SPILIOPOULOS KONSTANTINOS, DGM: A deep learning algorithm for solving partial differential equations, *Journal of Computational Physics*, **375**, p. 1339–1364, 2018.
- ZANG YAOHUA, BAO GANG, YE XIAOJING & ZHOU HAOMIN, Weak adversarial networks for high-dimensional partial differential equations, *Journal of Computational Physics*, **411**, p. 109409, 2020.