# Machine Learning for Inverse Problems in Climate Science

Rebecca Willett, University of Chicago



Video credit: Samuli Siltanen <u>https://www.youtube.com/watch?v=q7Rt\_OY\_7tU</u>

# Can machine learning help reconstruct images? Train deep neural network to reconstruct CT images from sinogram measurements



This approach can require *many* training samples.

# It also ignores everything we know about the data collection process.

Zhu, Liu, Rosen, Rosen, 2017; Arridge, Maass, Öktem, Schönlieb, 2019; Ongie, Jalal, Metzler, Baraniuk, Dimakis, Willett, 2020; Akçakaya, Yaman, Chung, Ye, 2022; Sahel, Bryan, Cleary, Farhi, Eldar, 2022; Kamilov, Bouman, Buzzard, Wohlberg, 2022 There are *many* settings in which we have both training data and physical models.







Also fluid dynamics, turbulence, particle accelerators, scattering, automatic control...



Physics-based models can inform neural network architectures and training, improving the reliability, efficiency and interpretability of ML systems.

#### Example: linear inverse problems in imaging

 $y = Hx + \varepsilon$ Recover *x* from *y* Observe: Goal:



Image reconstruction by supervised learning

1. Collect training data pairs  $(x_i, y_i)$  using a known forward model:

$$y_i = Hx_i + \varepsilon_i$$

2. Train a reconstruction network  $f_{\theta}$  by minimizing over a loss; e.g.

$$\min_{\theta} \sum_{i} \|x_i - f_{\theta}(y_i)\|_2^2$$

3. Reconstruct new measurements y by  $\hat{x} = f_{\theta}(y)$ 





Arridge, Maass, Öktem, Schönlieb, 2019 ; Ongie, Jalal, Metzler, Baraniuk, Dimakis, Willett, 2020; Monga, Li, Eldar, 2021

# Classical approach to solving inverse problems



Data fit term measures how well image x fits observation y, taking physical model H into account

Regularization function measures to what extent an image *x* has expected geometry (e.g. smoothness or sharp edges)

# Learning to reconstruct



#### Instead of using choosing R(x) a priori based on smoothness or geometric models, can we learn a regularizer using training data?

### **Optimization framework**

$$y \longrightarrow \underset{x}{\text{minimize } ||Hx - y||^2 + R(x) \longrightarrow \hat{x}}$$

for 
$$k = 1, 2, ...$$
  
 $z^{(k)} = x^{(k)} - \eta H^{\mathsf{T}}(Hx^{(k)} - y)$   
 $x^{(k+1)} = \operatorname{regularize}(z^{(k)}, R)$ 

data consistency step regularization step (e.g. proximal operator)



repeat until convergence

# Deep Unrolling

$$y \longrightarrow \underset{x}{\text{minimize } ||Hx - y||^2 + R(x)} \longrightarrow \hat{x}$$

for 
$$k = 1, 2, ...$$
  
 $z^{(k)} = x^{(k)} - \eta H^{\mathsf{T}}(Hx^{(k)} - y)$   
 $x^{(k+1)} = \mathsf{CNN}(z^{(k)})$ 

data consistency step regularization step



"Unroll" K iterations, train end-to-end in a supervised manner

# Physics-guided neural network architecture



# Physics-guided neural network architecture

One big neural network



 $\eta H^{+}$ 

 $x^{(1)}$ 

Some weights to be learned from training data

Physical models, inverse problem methods and optimization theory lead to novel architectures and highly effective learning methods

 $\hat{\chi}$ 

### Example: MRI reconstruction



#### **Original Image**

Machine learning method higher accuracy 16 s to compute

Classical method lower accuracy 350 s to compute

# Data assimilation

European Centre for Medium-Range Weather Forecasts:

"To make a forecast we need to know the current state of the atmosphere and the Earth's surface (land and oceans). The weather forecasts produced at ECMWF use data assimilation to estimate initial conditions for the forecast model from meteorological observations."

Here dim $(x_t) = 10^8$  and dim $(y_t) = 10^6$ , and data is collected in 6-hour windows.

Additional applications in tracking, molecular chemistry, robotics, phylogenetics, economics, geosciences, and much more.



https://www.ecmwf.int/en/research/data-assimilation

# State space modeling







# Machine learning for data assimilation

- Making accurate forecasts requires having a good model of underlying dynamics
- These models may have unknown parameters or only be approximate and sometimes we have no model at all!



#### Data assimilation: estimating dynamics from indirect data



# Our approach

• General strategy: choose  $\theta := (\alpha, \beta)$  (e.g. neural network weights) that maximizes likelihood of observations

$$\hat{\theta} = \arg \max_{\theta} \mathscr{L}(\theta)$$
 where  $\mathscr{L}(\theta) := \log p(y_{1:T} | \theta)$ 

- Challenge: log-likelihood generally does not have closed-form expression and must be numerically approximated
- Insight 1: Stochastic filtering tools like the Ensemble Kalman Filter (EnKF) yield approximate likelihood
- Insight 2: Using auto-differentiation to calculate likelihood gradients improves accuracy of learning optimal parameters

Our method: gradient ascent on approximate likelihood calculated using EnKF

# Stochastic filtering & prediction

- **Goal**: estimate the current state of an evolving dynamical system observed via indirect measurements
- **Example**: tracking
  - State is position and velocity at time *t*
  - But we only observe noisy location at each time t
  - **Filtering**: At each time, given past observations, estimate true location and velocity
  - **Prediction**: At each time, given past observations, predict future location and velocity



Image from Li, Wang, Wang, & Li 2010

## Ensemble Kalman Filter — Forecast step

Let  $x_{t-1}^n$  be the value of the  $n^{\text{th}}$  particle at time t - 1. First, for each particle, we predict where it will be at the next time given dynamics  $(F_{\alpha}, Q_{\beta})$ :

$$\hat{x}_{t}^{n} = F_{\alpha}(x_{t-1}^{n}) + Q_{\beta}^{1/2}\xi_{t}^{n}, \qquad \xi_{t}^{n} \sim \mathcal{N}(0,I)$$



# Ensemble Kalman Filter — Analysis step

Next, we observe  $y_t$  and use this observation plus knowledge of H to improve our estimate of each  $x_t^n$ :

$$\begin{aligned} x_t^n &= \hat{x}_t^n + \hat{K}_t(y_t + R^{1/2}\gamma_t^n - H\hat{x}_t^n), \\ \gamma_t^n &\sim \mathcal{N}(0, I) \end{aligned}$$

where

$$\hat{K}_t = \hat{C}_t H^\top (H\hat{C}_t H^\top + R)^{-1}$$

is called the Kalman gain



The **Ensemble Kalman Filter** sequentially estimates filtering distributions  $p_{\theta}(x_{1:t} | y_{1:t}), 1 \le t \le T.$ 

The distributions are represented using a collection of N particles.

Using ensemble mean and covariance, we can approximate the likelihood  $\mathscr{L}(\theta)$ 

Expectation Maximization Brajard, Carassi, Bocquet, & Bertino, 2020

- for k = 1, 2, ...
  - $x_{0:T}^{1:N} = \text{EnKF}(\theta^k, y_{1:T})$

• 
$$\theta^{k+1} = \theta^k + \eta \nabla_{\theta} \left[ \sum_{n,t} \log \mathcal{N} \left( y_t^n; F_{\theta^k}(x_{t-1}^n), Q_{\theta^k} \right) \right]$$

Brajard, Carassi,



# Key insight

In contrast to the EM method, we

- treat the particles as functions of  $\theta$  and
- compute likelihood gradients that reflect this dependence
- by leveraging **automatic differentiation**

Automatic differentiation is different from numerical differentiation

- autodiff uses compositions of elementary functions whose derivatives are known
- autodiff to compute gradients incurs negligible extra computational cost compared to evaluation of likelihoods
- finite difference approximations cause discretization errors

Expectation Maximization Brajard, Carassi, Bocquet, & Bertino, 2020

• for k = 1, 2, ...

• 
$$x_{0:T}^{1:N} = \text{EnKF}(\theta^k, y_{1:T})$$

• 
$$\theta^{k+1} = \theta^k + \eta \nabla_{\theta} \left[ \sum_{n,t} \log \mathcal{N} \left( y_t^n; F_{\theta^k}(x_{t-1}^n), Q_{\theta^k} \right) \right]$$

Brajard, Carassi,

# Our Autodiff-EnFK approach

• for k = 1, 2, ...

• 
$$x_{0:T}^{1:N}(\theta^k) = \text{EnKF}(\theta^k, y_{1:T})$$

• 
$$\theta^{k+1} = \theta^k + \eta \nabla_{\theta} \left[ \sum_{n,t} \log \mathcal{N}\left(y_t^n; F_{\theta^k}(x_{t-1}^n(\theta^k)), Q_{\theta^k}\right) \right]$$

#### **Compute accurate gradients using autodifferentiation**





#### Lorenz-96 system

Let  $F^*$  be flow map of vector field:

$$\frac{dx}{ds} = f^*(x), \qquad F^* : x(s) \mapsto x(s + \Delta_s)$$

with 
$$f^{*(i)}(x) = -x^{(i-1)}(x^{(i-2)} - x^{(i+1)}) - x^{(i)} + 8, i = 1,...,40$$



- Common test model for filtering algorithms and low-frequency climate models
- Prototypical turbulent dynamical system; our setting: strong chaotic turbulence
- Dynamics exhibit strong energy-conserving non-linearities
- Can be defined for any desired state space dimension  $d_x$

# Example on turbulent dynamics with partial obs.

We only observe 66% of the state and use a neural network to estimate the underlying dynamics



# Example — estimating simulation parameters

- Climate simulator takes parameters x and outputs simulation y = H(x)
- Given observations  $y_{\text{ODS}}$ , what are the corresponding parameters x?
- Similar to previous inverse problem settings, but now
  - *H* is nonlinear
  - we don't have an explicit form for *H*, can only access through simulations
  - we want to quantify uncertainty about *x*



# Past approaches

**Classical method:** For some predefined moment function *m*,

$$\hat{x} = \arg\min_{x} ||m(y) - m(H(x))||_{\Sigma[m(y)].}^2$$

- need expert knowledge to choose m,
- requires repeated (slow) runs of H for each new observation y,
- gradient-free optimization methods like ensemble Kalman inversion highly dependent on prior  $p_x$ .

**Supervised regression:** Learn a neural network  $f_{\theta}$  so that

$$\hat{x} = f_{\theta}(y);$$

difficult to get accurate uncertainty estimates.

Schneider, Tapio, Lan, Shiwei, Stuart, Andrew, et al., 2017

# Standard Emulator Approach

low-dimensional parameter x

 $\hat{H}_{ heta}$ 

high-dimensional dynamics y

$$\hat{x} = \arg\min_{x} \|m(y) - m(\hat{H}_{\theta}(x))\|_{\Sigma[m(y)]}$$

# Embed & Emulate (ours) low-dimensional parameter $x \quad \hat{g}_{\theta}$ high-dimensional dynamics y $f_{\theta}$ embedded parameters and dynamics $\approx f_{\theta} \circ H$ $\hat{x} = \arg \min_{x} ||f_{\theta}(y) - \hat{g}_{\theta}(x)||_{2}$

# Leveraging ideas from computer vision

- CLIP = Contrastive Language-Image Pre-Training
- Learns
   embedding so
   that
  - if an image and text go together, their embeddings are similar
  - if an image and text are unrelated, their embeddings are far apart



# Embed & Emulate key ideas



Inter-domain contrastive learning scheme: Diagonals are dot products between representations of "positive" pairs  $(x_i, y_i)$ .

- We design an "emulator" that fits well in the context of parameter estimation problem.
- We use CLIP-wise loss to align the metric space of the "emulator" and the embedding network.
- We use contrastive loss to capture intra-domain structural information to learn meaningful embeddings.

# Contrastive losses

Contrastive losses like  $\ell_{YY}$  identify positive and negative pairs based on simulation parameters x in the training data, making pairs robust to chaotic effects

 $f_{\theta}(\mathbf{y}_1)$ 

 $f_{\theta}(y_2)$ 

 $f_{\theta}(y_3)$ 

 $f_{\theta}(y_4)$ 

 $f_{\theta}(\tilde{y}_1)$ 

 $\mathcal{L}_{YY}(\theta)$ 

encoder  $f_{ heta}$ trajectory

V



# Example estimating Lorenz-96 parameters



Computation time for 500 training samples + 200 testing samples (including time to generate training data, reported in minutes).

# Embed & Emulate architecture



# Regression head informs EnKI prior



Regression head helps guide contrastive learning

Using regression head to set prior used by EnKI further reduces errors Physics-based models can inform neural network architectures and training, improving the reliability and interpretability of ML systems.

# Al & Science @ UChicago

#### We envision **AI as an integral component of the scientific method**, guiding the construction of hypotheses, designing sequences of experiments, and analyzing data to develop new hypotheses, while advancing core AI principles.









# Key challenges of AI in science

- Effective ML training for small or sparse datasets
- Incorporation of physical models into AI structure
- Creation of **surrogates**
- Dimensional reduction, data synthesis and compression, and reduced order models
- Control and AI-enabled experimental design
- **Operation** of experimental facilities
- Robustness, inference, and calibration
- Real-time decision-making
- Predictive maintenance and event prediction
- Training data variability and **noise**
- Interpretable models and algorithms
- Coupling simulations and experiment

2020 DOE Report "Opportunities and Challenges from Artificial Intelligence and Machine Learning for the Advancement of Science, Technology, and the Office of Science Missions"



New Schmidt Futures fellowship at UChicago to foster next generation of AIdriven scientists





# **E** SCHMIDT **FUTURES**

https://aiscience.uchicago.edu/

# Thank you!

