Learning Physics from Data with Multiscale Interaction Models

ENS

Stéphane Mallat,

Collège de France École Normale Supérieure

Learning Physics and Probabilities

• Unsupervised learning: estimate p(x) for $x \in \mathbb{R}^d$ from $\{x_i\}_i$ Physical systems at equilibrium: learn the energy U(x)

 $p(x) = \mathcal{Z}^{-1} e^{-U(x)}$

Curse of dimensionality unless strong priors



Learn from few samples

Applications:

- Understand the physics from the energy U(x).
- Data generation by sampling p(x)

Learning Physics and Probabilities

• Unsupervised learning: estimate p(x) for $x \in \mathbb{R}^d$ from $\{x_i\}_i$ Physical systems at equilibrium: learn the energy U(x)

$$p(x) = \mathcal{Z}^{-1} e^{-U(x)}$$

Three Problems:

- Find $\Phi(x) = \{\Phi_k(x)\}_k$ which can linearly approximate U(x): $U(x) \approx \theta^T \Phi(x)$ (Ansatz) $\Leftrightarrow p_{\theta}(x) = \mathcal{Z}^1_{\theta} e^{-\theta^T \Phi(x)}$ closely approximates p(x).
- Optimise θ to minimise the approximation error.
- Sample $p_{\theta}(x)$ to generate new data x.

Difficult if the energy U(x) is non convex: usually the case. A solution: scale separation and renormalisation group.



Interaction of d variables x(u): pixels, particles, agents...



Regroupement of d interactions in $O(\log d)$ multiscale terms. Fast multipole algorithms.

Multiscale Interactions Claims

- 1. Long range spatial dependencies are captured by short range dependencies over wavelet coefficients : reduction of dimension of $\Phi(x)$
- 2. Can specify $\Phi(x)$ from interaction energies across scales, with the *renormalisation group*, and which are "nearly" convex.
- 3. Interaction energies across scales are non-linear, but can be partly linearised with a modulus or a ReLU by eliminating phases to define $\Phi(x)$ (similar to deep networks).



- I: Sampling, max likelihood and score matching estimations
- II: Renormalisation group and energy interactions across scales
- III: Scattering-Gaussian models of scale interactions: turbulences

I. Sampling by Langevin (or MCMC)

• Computing a sample x of p_{θ} with a Langevin dynamics

 $x_{t+1} - x_t = \epsilon \nabla_x \log p_\theta(x_t) + \sqrt{2\epsilon} z_t$ where $z_t \sim \mathcal{N}(0, Id)$

Converges but very slow if $-\log p_{\theta}(x)$ is non-convex in x.



Exploration / Exploitation

Exponential convergence if the energy is strictly convex,

The rate is the inverse of the condition number of $-\nabla_x^2 \log p_\theta(x)$. (Hessian of energy)

Pb: we want $p_{\theta} \approx p$, but $-\nabla_x^2 \log p$ is usually not positive, or badly conditioned.

• Estimated by minimising Kullback-Liebler divergences $KL(p||p_{\theta}) = \int p(x) \log \frac{p(x)}{p_{\theta}(x)} dx = \mathbb{E}_p[\log p(x) - \log p_{\theta}(x)]$

• Optimise θ by gradient descent (likelihood ascent)

 $\theta_{t+1} - \theta_t = \epsilon \nabla_{\theta} \mathbb{E}_p(\log p_{\theta})$

• For $-\log p_{\theta}(x) = \log \mathcal{Z}_{\theta} + \theta^T \Phi(x)$ the optimisation is convex

with $\nabla_{\theta} \mathbb{E}_p(\log p_{\theta}) = \mathbb{E}_p(\Phi(x)) - \mathbb{E}_{p_{\theta}}(\Phi(x))$: moment matching

which requires to compute samples of p_{θ} : expensive.





$-\log p_{\theta}(x) = \log \mathcal{Z}_{\theta} + \theta^T \Phi(x)$

• The score matching eliminates the influence of \mathcal{Z}_{θ} Minimisation of the Fisher Information divergence: $FI(p||p_{\theta}) = \mathbb{E}_p[||\nabla_x \log p_{\theta}(x) - \nabla_x \log p(x)||^2]$

Quadratic function of θ :

$$FI(p||p_{\theta}) = \mathbb{E}_p\left[\frac{1}{2}\|\theta^T \nabla_x \Phi(x)\|^2 - \theta^T \Delta_x \Phi(x)\right]$$

Calculated from data, no need to sample p_{θ} : fast algorithm.





$-\log p_{\theta}(x) = \log \mathcal{Z}_{\theta} + \theta^T \Phi(x)$

• The score matching eliminates the influence of \mathcal{Z}_{θ} Minimisation of the Fisher Information divergence:

 $FI(p||p_{\theta}) = \mathbb{E}_p[||\nabla_x \log p_{\theta}(x) - \nabla_x \log p(x)||^2]$

Theorem: If $-\log p_{\theta}(x)$ and $-\log p(x)$ are strictly convex then score matching estimation is identical to maximum likelihood.



Eliminating constants does not allow to estimate local minima levels

II. Renormalisation Group

 \mathcal{X}

 x_{j-1}

 x_{j}

 x_{j+1}

 \mathbf{X}_{I}

Kadanoff, Wilson 1970

- Given p(x) = Z⁻¹e^{-U(x)} compute p(x_j) at all scales 2^j: regular evolution. x_{j-1} → x_j: average, subsample, normalise.
 Ansatz : p(x_j) = Z⁻¹_je^{-θ^T_jΦ(x_j)} how to find it ?
 - $\theta_j = f(\theta_{j-1})$: RG equation on couplings At phase transitions: $\theta_j = \theta_{j-1}$.
- Calculated by decomposing x_{j-1} into (x_j, \overline{x}_j) :

$$\Rightarrow p(x_j) = \int p(x_{j-1}) d\bar{x}_j = \int p(x_j) p(\bar{x}_j | x_j) d\bar{x}_j$$

Wilson: approximation with Gaussian integration in \bar{x}_j

Inverse Renormalization Group

- T. Marchand, M. Ozawa, G. Biroli
- **Pb**: estimate physical energies from data. Coarse to fine:

 $p(x_{i-1}) = p(x_i) p(\bar{x}_i/x_i)$

with normalised variances $\bar{x}_i \to \bar{x}_i / \sigma_i$

$$\Rightarrow p(x) = p(x_J) \prod_{j=1}^J p(\overline{x}_j/x_j)$$

Claims 1,2: $p(\bar{x}_i | x_i)$ has short range dependencies $f_{p(x_{i-1}/x_i)}$ in \bar{x}_i , and its energy is "nearly" convex in \bar{x}_i . $p(x_J)$ may be non-convex but low-dimensional. It is usually Gaussian. x_{j+1}

 x_{j-1}

 \mathcal{X}_{j}

 \bar{x}_J

Wavelet Orthogonal Bases

Decompose each x_{j-1} into (x_j, \bar{x}_j) : in what basis ?

translated, and dilated local wavelets

 $\left\{\psi_{j,n}^{m}\right\}_{m,j,n}$ can define an orthonormal basis of $\mathbf{L}^{2}(\mathbb{R}^{d})$

 $\overline{x}_j(n,m) = \langle x, \psi_{j,n}^m \rangle = x * \psi_j^m(2^j n)$

Fast wavelet transform: $p(\bar{x}_j|x_j)$ has short dependencies



Wavelet Orthogonal Bases

Decompose each x_{j-1} into (x_j, \bar{x}_j) : in what basis ?

translated, and dilated local wavelets

 $\left\{\psi_{j,n}^{m}\right\}_{m,j,n}$ can define an orthonormal basis of $\mathbf{L}^{2}(\mathbb{R}^{d})$

$$\overline{x}_j(n,m) = \langle x, \psi_{j,n}^m \rangle = x * \psi_j^m(2^j n)$$

Fast wavelet transform: $p(\bar{x}_j|x_j)$ has short dependencies



Wavelet Conditional RG Models

T. Marchand, M. Ozawa, G. Biroli $p(x) = p(x_J) \prod_{j=1}^{J} p(\overline{x}_j/x_j)$

Define parameterised models of each $p(\overline{x}_j/x_j)$

$$p_{\bar{\theta}_j}(\bar{x}_j/x_j) = \bar{\mathcal{Z}}_j^{-1} e^{-\bar{\theta}_j^T \Phi(x_{j-1})}$$

with same Φ but different $\overline{\theta}_j$ at all scales

Claims 1,2:

1. Short range dependencies: $\overline{\theta}_j$ has few non-zero parameters

2. Convexity of $\overline{\theta}_j^T \Phi(x_{j-1})$ in \overline{x}_j and well conditioned Hessian.

 \Rightarrow fast sampling of $p_{\overline{\theta}_j}$ and fast optimisation of $\overline{\theta}_j$.

Multiscale Sampling & Energy

 \mathcal{X}

 x_{j-1}

 x_{j}

 $x_J \wedge \overline{x}_J$

T. Marchand, M. Ozawa, G. Biroli

$$p_{\theta}(x) = p_{\theta_J}(x_J) \prod_{j=1}^J p_{\bar{\theta}_j}(\overline{x}_j/x_j)$$

Multiscale energy estimation \overline{x}_{i-1} $p_{\theta}(x) = \mathcal{Z}_{\theta}^{-1} e^{-U_{\theta}(x)}$ $U_{\theta}(x) = \theta_J^T \Phi(x_J) + \sum \bar{\theta}_j^T \Phi(x_{j-1})$ \overline{x}_{j} i=1sample $p_{\bar{\theta}_i}(\bar{x}_j/x_j)$ $x_{j+1} \mathbf{x}_{j+1}$ sample $p_{\bar{\theta}_{j+1}}(\bar{x}_{j+1}/x_{j+1})$ sample $p_{\bar{\theta}_{I}}(\bar{x}_{J}/x_{J})$ sample $p_{\theta_I}(x_J)$

Multiscale Gaussian Models

$$p_{\theta}(x) = \mathcal{Z}_{\theta}^{-1} e^{-\frac{1}{2}x^{T} \theta x} = \mathcal{Z}_{\theta}^{-1} e^{-\theta^{T} \Phi(x)} \text{ with } \Phi(x) = xx^{T}$$

Turubulent Flow

At convergence $\nabla_x^2 \log p_\theta(x) = \theta = K \ge 0$ where K^{-1} is the process covariance.

Stationarity $\Rightarrow K^{-1}$ is Toeplitz diagonalised in Fourier

Power spectrum = covariance eigenvalues, has a power-lay decay for multiscale fields. \Rightarrow long range dependencies.

Gaussian eigenvalues $K / \omega |^{\eta}$ model $\begin{array}{c} K^{-1} \\ \searrow |\omega|^{-\gamma} \end{array}$ condition number $\sim L^{\eta}$

Wavelet Transform in Fourier

Orthogonal wavelets decomposes in different frequency bands



Wavelet coefficients

Frequency (Fourier) domain $\begin{array}{c} & & & \\ &$

Well Conditioned Hessian

• For fields having a power spectrum with power law decay

$$\nabla^2 \log p(\bar{x}_j | x_j) = K_j = \bar{\theta}_j$$



Wavelet representation of singular operators

Theorem: For Gaussian models of a spectrum with power-law decay each K_j is a band matrix with a condition number ~ 1 .

⇒ fast sampling of $p_{\bar{\theta}_j}(\bar{x}_j/x_j)$ with Langevin (or MCMC). and short range dependencies: low-dimensional model.

Potential Energies in Physics

 $p(x) = \mathcal{Z}^{-1} e^{-U(x)}$ with $U(x) = \frac{1}{2}x^T K x + V(x)$

where $K = \beta \Delta$ is the kinetic energy

V(x) is the non-convex potential, creating scale interactions

 $\nabla^2 \log p(x) = K + \nabla^2 V(x)$ is usually not convex.

The non-convexity appears at lower frequencies



Scalar Potential Models

 $\varphi^4 \mod : U(x) = \frac{\beta}{2} x^T \Delta x + \sum_n V(x_n)$



• Polynomial parametric model:

$$V(\varphi) = \sum_{k} \theta_{k} P_{k}(\varphi) \implies U(x) = \theta^{T} \Phi(x)$$

with $\Phi(x) = \left(xx^{T}, \sum_{n} P_{k}(x_{n})\right)_{k}$

• Stationarity $\Rightarrow U(x) = \tilde{\theta} * \tilde{\Phi}(x)$

Convolutional scalar potential Ansatz



• Mixing time for sampling by MCMC Metropolis:



Convexity and score matching J. Bruna, F. Guth, E. Lempereur The interaction energies $\bar{\theta}_j \star \Phi(\bar{x}_{j-1})$ of all $p_{\bar{\theta}_j}(x_j|x_j)$ are "nearly" convex with well conditioned Hessians:

Histograms of Hessian eigenvalues





 $2 10^{0}$

 10^{1}

 10^{2}



Wavelet Conditional RG





Not sufficient to model





III- Modeling Scale Dependencies

 x_0

Geometric structures:

٠





- Sparse wavelet coefficients
- Strongly dependant across scales and angles.

 $p(\overline{x}_j/x_j) = p(\overline{x}_j/\overline{x}_\ell, \, \ell > j)$

Modeling Scale Dependencies

Must model the energy of $p(\bar{x}_j/x_j) = p(\bar{x}_j/\bar{x}_{j'}, j'>j)$ Linear models $\theta^T \Phi$ are matching moments: $\mathbb{E}_{p_{\theta}}(\Phi(x)) = \mathbb{E}_p(\Phi(x))$ What moments should we choose ?

• Gaussian model: $\Phi(x) = x x^T \Leftrightarrow \Phi(x) = W x (W x)^T$ in a wavelet basis $W x = (\bar{x}_j)_j$.

For x stationary, the $\overline{x}_j = x \star \psi_j^m$ are not correlated across scales:

$$\mathbb{E}\Big(\overline{x}_j(u)\ \overline{x}_{j'}(u-\tau)\Big)\approx 0 \quad \text{if} \quad j\neq j'$$

because their phase oscillate at different rate since ψ_j^m and $\psi_{j'}^{m'}$ are supported in different frequency bands.

 \Rightarrow Gaussian models do not capture scale dependencies.

Scattering-Gaussian Models

- Eliminates phases with modulus: $(x, |Wx|) = (x, |\bar{x}_j|)_j$
- Scattering-Gaussian model: Etienne Lempereur

 $p_{\theta}(x) = \mathcal{Z}_{\theta}^{-1} e^{-(x,|Wx|) \theta (x,|Wx|)^{T}}$

 $p_{\theta} = \mathcal{Z}_{\theta}^{-1} e^{-\theta^{T} \Phi(x)} \text{ with } \Phi(x) = (x, |Wx|) (x, |Wx|)^{T}$ $\Rightarrow \Phi(x) = \begin{pmatrix} x(u) x(u-\tau) & x(u) |\bar{x}_{j}(u-\tau)| \\ x(u) |\bar{x}_{j}(u-\tau)| & |\bar{x}_{j}(u)| |\bar{x}_{j'}(u-\tau)| \end{pmatrix}_{j,j',\tau}$ Large matrix because modulus have long range correlations.

Scattering-Gaussian Models

- Eliminates phases with modulus: $(x, |Wx|) = (x, |\bar{x}_j|)_j$
- Scattering-Gaussian model: Etienne Lempereur

 $p_{\theta}(x) = \mathcal{Z}_{\theta}^{-1} e^{-(x,|Wx|) \theta (x,|Wx|)^{T}}$

• Wavelet: W(x, |Wx|) eliminates long range correlations $p_{\theta} = \mathcal{Z}_{\theta}^{-1} e^{-\tilde{\theta}^T \tilde{\Phi}(x)}$ with $\tilde{\Phi}(x) = (Wx, W|Wx|) (Wx, W|Wx|)^T$ and $\tilde{\theta} = W\theta W^T$ of reduced dimension

• $\tilde{\theta}$ and $\tilde{\Phi}(x)$ are reduced to $O(\log^3 d)$ diagonal coefficients: $\widetilde{\Phi}(x) = \begin{pmatrix} |\bar{x}_j(u)|^2 & \bar{x}_j(u) |\bar{x}_{j'}(u)| \\ |\bar{x}_j| \star \psi_k(u) |\bar{x}_{j'}| \star \psi_k(u) \end{pmatrix}_{j,j',k}$

Generation from Scattering Covar.

E. Allys, F. Boulanger, A. Brochard, J. Bruna, S. Chen, B. Ménard, R. Morel, G. Rochette, S. Zhang Original images of dimension $d = 5 \, 10^4$: only 1 sample



Generated with models having 500 parameters from 1 sample Reproduces moments of order 3 (bispectrum) and 4 (trispectrum)

Scattering-Gaussian Model



We can replace covariances by random projections and ReLU $$\rho$$ A low-dimensional deep networks '

Deep Scattering Network



 $-\log p_{\theta}(x) = \theta^T \Phi(x) + \log \mathcal{Z}_{\theta}$: learn only last linear filter



Conclusion

- Build models from interaction energies across scales : wavelet conditional renormalisation group.
- When are interaction energies convex (or nearly)?
- *Scattering-Gaussian models:* interaction energies captured by wavelet modulus correlations across scales: turbulence, active matter...
 - Wavelet Conditional Renormalisation Group : Phys. Rev. X, T. Marchand, M. Ozawa, G. Biroli, S.M.
 - Conditionally Strongly Log Concave Generative Models, ICML 2023 J. Bruna, E. Lempereur, F. Guth. S. M.
- Scattering Spectra Models for Physics, arXiv:2306.17210 S. Cheng, R. Morel, E. Allys, B. Menard, S. M.