Generative AI & molecular simulation

Lei Wang (王磊) Institute of Physics, CAS https://wangleiphy.github.io





Generative AI

Lakeview with geese in the autumn panoramic artist



https://huggingface.co/spaces/stabilityai/stable-diffusion



Discriminative learning



 $y = f(\boldsymbol{x})$ or $p(y|\mathbf{x})$

Generative learning



 $p(\boldsymbol{x}, y)$



What I cannot reate, I to not understand. Why const × sort. Po TO LEARN; Bethe Ansitz Prob. Kanton Know how to solve every problem that has been solved Hall accel. Temp Non Linear Dessical Hyper





What I cannot reate, I to not understand. Why const × Sort. PO TO LEARN; Bethe Aments Prob. Kanton Know how to solve every problem that has been solved 2-D Hall accel. Temp Non Linear Dessical Hyper

"What I can not create, I do not understand"



Generative AI: a new buzz word in silicon valley

A Coming-Out Party for Generative A.I., Silicon Valley's New Craze

A celebration for Stability AI, the start-up behind the controversial Stable Diffusion image generator, represents the arrival of a new A.I. boom.

Protocol

Biz Carson October 21, 2022

New York Times

Kevin Roose

Oct. 21, 2022

Sequoia's Sonya Huang: The generative Al hype is 'absolutely justified'

She's bullish on generative AI given the "superpowers" it gives humans who work with it.





https://www.sequoiacap.com/article/generative-ai-a-creative-new-world/ by Sonya Huang, Pat Grady and GPT-3

	PRE-2020	2020	2022	2023?	2025?	2030?	
TEXT	Spam detectionTranslationBasic Q&A		Longer form Second drafts	Vertical fine tuning gets good (scientific papers, etc)	Final drafts better than the human average	Final drafts better than professional writers	
CODE	1-line auto-complete	Multi-line generation	Longer form Better accuracy	More languages More verticals	Text to product (draft)	Text to product (final better than full-time developers	
IMAGES			Art Logos Photography	Mock-ups (product design, architecture, etc.)	Final drafts (product design, architecture, etc.)	Final drafts better than professional artists, designers, photographers)	
VIDEO / 3D / GAMING			First attempts at 3D/video models	Basic / first draft videos and 3D files	Second drafts	Al Roblox Video games and movies are personalized dreams	
			Large model availability:	First attempts	Almost there	Ready for prime tir	







https://huggingface.co/spaces/stabilityai/stable-diffusion

the inner structure of an electron



Generate image



https://future.com/how-to-build-gpt-3-for-science/ How to Build a GPT-3 for Science (scientific literature and data)

Josh Nicholson

Posted August 18, 2022

Some examples of complex potential prompts are:

"Tell m "Tell m "Gener "What "Who l field?"

"Write me a scientific paper based on my data"

- "Tell me why this hypothesis is wrong"
- "Tell me why my treatment idea won't work"
- "Generate a new treatment idea"
- "What evidence is there to support social policy X?" "Who has published the most reliable research in this

Generative AI for matter engineering

latent space

Review: "Inverse molecular design using machine learning", Sanchez-Lengeling & Aspuru-Guzik, Science '18







Generative AI for statistical physics

Renormalization group Molecular simulation Lattice field theory



Li and LW, PRL '18

Li, Dong, Zhang, LW, PRX '20







Noe et al, Science '19 Wirnsberger et al, JCP '20

Albergo et al, PRD '19 Kanwar et al, PRL '20

These are principled computation: quantitatively accurate, interpretable, reliable, and generalizable even without data



Probabilistic Generative Modeling $p(\mathbf{x})$

How to express, learn, and sample from a high-dimensional probability distribution?

CHAPTER 5. MACHINE LEARNING BASICS





Figure 5.12: Sampling images uniformly at random (by randomly picking each pixel Figure 1.9: Example inputs from the MNIST dataset. The "NIST" stands for National according to a uniform distribution) gives rise to noisy images. Although there is a non-Institute of Standards and Technology, the agency that originally collected this data. zero probability to generate an image of a face or any other object frequently encountered The "M" stands for "modified," since the data has been preprocessed for easier use with in AI applications, we never actually observe this happening in practice. This suggests in AI applications, we never actually observe this happening in practice. This suggests that the images encountered in AI applications occupy a negligible proportion of the volume of image space. Of course, concentrated probability distributions are not simpler to show the probability distributions are not sincleaded " volume of image space. Of course, concentrated probability distributions are not specific to that the data lies on a reasonably small number of manifolds. We must also establish that the examples we encounter are connected to each other by other it allows machine learning researchers to study their algorithms in controlled laboratory

conditions, much as biologists often study fruit flies.

3	4	7	8	9	0	1	2	3	4	5	6	7	8	6
5	5	4	7	8	9	2	9	3	9	3	8	2	0	5
6	5	3	5	3	8	0	0	3	4	1.	5	3	0	8
/	Ĵ	8	1	1	1	З	8	9	7	6	7	4	1	6
1	8	0	6	9	4	9	9	3	7	1	9	2	2	5
4	5	6	7	8	9	0	1	2	3	4	5	6	7	0
6	7	8	9	8	1	0	5	5	/	Ŷ	0	4	1	9
8	5	0	6	5	5	3	3	3	9	8	7	4	0	6
/	7	3	2	8	8	7	8	4	6	0	2	0	3	6
P	3	R	4	9	'4	6	5	3	Z	ধ	5	9	4	/
3	4	5	6	7	ଞ	9	0	I	2	3	4	5	6	7
3	U	5	6	7	ଟ	9	6	4	2	6	4	7	5	6
Ŷ	3	S	3	8	2	0	q	8	0	5	6	0	F	0
5	4	3	4	l	5	3	0	જ	3	0	6	2	7	1
3	8	5	н	2	Û	9	7	6	7	4	1	6	8	4
7	1	9	જ	O	6	9	4	9	9	6	2	3	7	1
7	8	0	1	2	3	4	5	6	7	8	0	1	2	3
9	1	2	3	4	5	6	7	8	9	2	1	2	1	3
7	0	7	1	5	7	9	9	4	7	0	3	4	1	4
4	8	4	1	8	6	4	4	6	3	5	7	2	5	9



images

Probab

How to high-din

CHAPTER 5. MACHINE LEARNING BASICS

"... the images encountered in Al applications occupy a negligible proportion of the volume of image space."

Figure 5.12: Sampling images uniformly at rando according to a uniform distribution) gives rise to no zero probability to generate an image of a face or an in AI applications, we never actually observe this if that the images encountered in AI applications or volume of image space.

Of course, concentrated probability distributed that the data lies on a reasonably small num establish that the examples we encounter are e

162

deling

DEEP LEARNING

Ian Goodfellow, Yoshua Bengio, and Aaron Courville

Page 159

om a vition?



Probabilistic Generative Modeling $p(\mathbf{x})$

How to express, learn, and sample from a high-dimensional probability distribution?



https://blog.openai.com/generative-models/







Lecture Note http://wangleiphy.github.io/lectures/PILtutorial.pdf

Generative Models for Physicists

Lei Wang^{*}

Institute of Physics, Chinese Academy of Sciences Beijing 100190, China

October 28, 2018

Abstract

Generative models generate unseen samples according to a learned joint probability distribution in the highdimensional space. They find wide applications in density estimation, variational inference, representation learning and more. Deep generative models and associated techniques (such as differentiable programing and representation learning) are cutting-edge technologies physicists can learn from deep learning.

This note introduces the concept and principles of generative modeling, together with applications of modern generative models (autoregressive models, normalizing flows, variational autoencoders etc) as well as the old ones (Boltzmann machines) to physics problems. As a bonus, this note puts some emphasize on physics-inspired generative models which take insights from statistical, quantum, and fluid mechanics.

The latest version of the note is at http://wangleiphy.github.io/. Please send comments, suggestions and corrections to the email address in below.

3

1

CONTENTS

GENERATIVE MODELING 2							
1.1 Probabilistic Generative Modeling 2							
1.2 Generative Model Zoo 4							
1.2.1 Boltzmann Machines 5							
1.2.2 Autoregressive Models 8							
1.2.3 Normalizing Flow 9							
1.2.4 Variational Autoencoders 13							
1.2.5 Tensor Networks 15							
1.2.6 Generative Adversarial Networks 17							
1.2.7 Generative Moment Matching Networks 18							
1.3 Summary 20							
PHYSICS APPLICATIONS 21							
2.1 Variational Ansatz 21							
2 Renormalization Group 22							
Monte Carlo Update Proposals 22							
4 Chemical and Material Design 23							
2.5 Quantum Information Science and Beyond 24							
RESOURCES 25							
BIBLIOGRAPHY 26							

So, what is the fuss?



$$\longrightarrow p(\mathbf{x}) \ge 0$$

Normalization?

Sampling?

$$\int dx \, p(\mathbf{x}) = 1 \qquad \qquad \mathbb{E} \\ \mathbf{x} \sim p(\mathbf{x})$$

Lecture 1: sampling distributions in vanilla MD Lecture 2: free energy calculation (no rare event)



us im

Children computing the number π on the Monte Carlo beach.





Adults computing the number π at the Monte Carlo heliport.











Two sides of the same coin **Generative modeling Statistical physics**



Known: samples Unknown: generating distribution

Maximum likelihood estimation

"learn from data"

 $\mathscr{L} = -\mathbb{E}_{\mathbf{x} \sim \text{dataset}} \left| \ln p(\mathbf{x}) \right|$



Known: energy function Unknown: partition function, samples

Variational free energy

"learn from Hamiltonian"

 $F = \mathbb{E}_{\substack{\boldsymbol{x} \sim p(\boldsymbol{x})}} \left[k_B T \ln p(\boldsymbol{x}) + \boldsymbol{H}(\boldsymbol{x}) \right]$



The variational free energy principle

variational density

Difficulties in Applying the Variational Principle to Quantum Field Theories¹

Richard P. Feynman

$F[p] = \int d\mathbf{x} \, p(\mathbf{x}) \left[H(\mathbf{x}) + k_B T \ln p(\mathbf{x}) \right] \ge F$ $\downarrow \qquad \qquad \downarrow \qquad \qquad \downarrow \qquad \qquad \downarrow \qquad \qquad \downarrow$ riational density energy entropy $\mathbf{\widehat{o}}$

¹transcript of his talk in 1987

Generative models



Deep variational free energy approach

Use deep generative models as the variational density

$$F[p] = \mathbb{E}_{\substack{x \sim p(x) \\ energy}} \begin{bmatrix} H(x) + k_B T \ln p(x) \end{bmatrix}$$





Li and LW, PRL '18 Wu, LW, Zhang, PRL '19

with normalizing flow & autoregressive models



V Turning a sampling problem to an optimization problem better leverages the deep learning engine:



Deep variational free-energy in the context

E, Han, Zhang, Physics Today 2020



Objective	Model	Data	Task	
MD potential energy surface	3N-dim function	DFT energy/ force	Generalizatio	
DFT xc energy functional	3-dim functional	QMC/ CCSD/		
Variational free-energy	3N-dim functional	No	Optimizatior	



Forward KL or Reverse KL?



 ${\mathcal X}$

Fig. 3.6, Goodfellow, Bengio, Courville, <u>http://www.deeplearningbook.org/</u>

Probability Density



 ${\mathcal X}$

$$D_{\alpha}(p \mid\mid q) = \frac{\int_{x} \alpha p(x) + (1 - \alpha)q(x) - p(x)^{\alpha}q(x)^{1 - \alpha}dx}{\alpha(1 - \alpha)}$$





$$D_{-1}(p || q) = \frac{1}{2} \int_{x} \frac{(q(x) - p(x))^{2}}{p(x)} dx$$
$$\lim_{\alpha \to 0} D_{\alpha}(p || q) = \text{KL}(q || p)$$
$$D_{\frac{1}{2}}(p || q) = 2 \int_{x} \left(\sqrt{p(x)} - \sqrt{q(x)}\right)$$
$$\lim_{\alpha \to 1} D_{\alpha}(p || q) = \text{KL}(p || q)$$
$$D_{2}(p || q) = \frac{1}{2} \int_{x} \frac{(p(x) - q(x))^{2}}{q(x)} dx$$







GAUSSIAN-BERNOULLI RBMS WITHOUT TEARS

Renjie Liao^{*1}, Simon Kornblith², Mengye Ren³, David J. Fleet^{2,4,5}, Geoffrey Hinton^{2,4,5}

2210.10318

6 3 3 3 6 8

584419

3779376

15350ZZ

4251242

050709



 $\mathscr{L} = -\mathbb{E}_{\mathbf{x} \sim \text{dataset}} \left[\ln p(\mathbf{x}) \right] \quad p(\mathbf{x}) = e^{-E(\mathbf{x})}/Z$



GAUSSIAN-BERNOULLI RBMS WITHOUT TEARS

Renjie Liao^{*1}, Simon Kornblith², Mengye Ren³, David J. Fleet^{2,4,5}, Geoffrey Hinton^{2,4,5}

2210.10318

6 3 3 3 6 8

584419

5350ZA

050709

3779876

4251242

Boltzmann machines



 $\mathscr{L} = -\mathbb{E}_{\mathbf{x} \sim \text{dataset}} \left[\ln p(\mathbf{x}) \right] \quad p(\mathbf{x}) = e^{-E(\mathbf{x})} / Z$



GAUSSIAN-BERNOULLI RBMS WITHOUT TEARS

Renjie Liao^{*1}, Simon Kornblith², Mengye Ren³, David J. Fleet^{2,4,5}, Geoffrey Hinton^{2,4,5}

2210.10318

5

6633368

4584419

3779376

15350ZA

4251242

3050709

Boltzmann machines



 $\mathscr{L} = -\mathbb{E}_{\mathbf{x} \sim \text{dataset}} \left[\ln p(\mathbf{x}) \right] \quad p(\mathbf{x}) = e^{-E(\mathbf{x})} / Z$



GAUSSIAN-BERNOULLI RBMS WITHOUT TEARS

Renjie Liao^{*1}, Simon Kornblith², Mengye Ren³, David J. Fleet^{2,4,5}, Geoffrey Hinton^{2,4,5}

2210.10318

5

6633368

4584419

3779376

15350ZA

4251242

3050709

Boltzmann machines

Autoregressive models

Language: GPT 2005.14165



Image: PixelCNN 1601.06759



 $p(\mathbf{x}) = p(x_1)p(x_2 | x_1)p(x_3 | x_1, x_2) \cdots$

Speech: WaveNet 1609.03499



Molecular graph: 1810.11347



Autoregressive models

Language: GPT 2005.14165



Image: PixelCNN 1601.06759



 $p(\mathbf{x}) = p(x_1)p(x_2 | x_1)p(x_3 | x_1, x_2) \cdots$

Speech: WaveNet 1609.03499



Molecular graph: 1810.11347



Variational autoregressive networks



Sherrington-Kirkpatrick spin glass



Variational autoregressive network

$$p(\mathbf{x}) = \prod_{i} p(x_i | \mathbf{x}_{< i})$$

github.com/wdphy16/stat-mech-van Wu, LW, Zhang, PRL '19

Conventional approaches

Naive mean-field factorized probability $p(\mathbf{x})$

Bethe approximation pairwise interaction

$$p(\mathbf{x}) = \prod_{i} p(x_i) \prod_{(i,j) \in E} \frac{p(x_i, y_i)}{p(x_i)p(x_i)}$$



Implementation: autoregressive masks



 $p(x_2 | x_1) = \text{Bern}(\hat{x}_2)$ $p(x_3 | x_1, x_2) = \text{Bern}(\hat{x}_3)$ $p(x_1) = \operatorname{Bern}(\hat{x}_1)$

Other examples: PixelCNN, van den Oord et al, 1601.06759 Casual transformer, 1706.03762 Other ways to implement autoregressive models: recurrent networks

Masked Autoencoder Germain et al, 1502.03509



Normalizing flows



https://deepmind.com/blog/high-fidelity-speech-synthesis-wavenet/





https://blog.openai.com/glow/


Normalizing flows



https://deepmind.com/blog/high-fidelity-speech-synthesis-wavenet/





https://blog.openai.com/glow/



Normalizing flow in a nutshell

Base density

"neural net" with 1 neuron





Normalizing flow for physics: an intuition





coupled

oscillators



High-dimensional, composable, learnable, nonlinear transformations



Neural network renormalization group

Collective variables



Physical variables

Li, LW, PRL '18 lio12589/NeuralRG



Neural network renormalization group

Collective variables



Probability Transformation

$$\ln p(\mathbf{x}) = \ln \mathcal{N}(z) - \ln \left| \det \left(\frac{\partial x}{\partial z} \right) \right|$$



Physical variables

Li, LW, PRL '18 lio12589/NeuralRG



Continuous normalizing flows $\ln p(\mathbf{x}) = \ln \mathcal{N}(z) - \ln \left| \det \left(\frac{\partial x}{\partial z} \right) \right|$

Consider infinitesimal change-of-variables Chen et al 1806.07366 $\ln p(\mathbf{x}) - \ln \mathcal{N}(z) = -\ln \left| \det \left(1 + \varepsilon \frac{\partial v}{\partial z} \right) \right|$ $x = z + \varepsilon v$

 $\varepsilon \to 0$

 $\frac{dx}{dt}$ = v

$$\frac{d\ln p(\boldsymbol{x},t)}{dt} = -\nabla \cdot \boldsymbol{v}$$

Continuous n

$\ln p(\mathbf{x}) = \ln \mathcal{N}$

t =

Consider infinitesimal change-of-variables Chen et al 1806.07366

 $\ln p(x)$ $x = z + \varepsilon v$

 $\frac{dx}{dt} = v$

 $\varepsilon \to 0$

$$f(z) - \ln \left| \det \left(\frac{\partial x}{\partial z} \right) \right|$$

$$\begin{aligned} \mathbf{x} &- \ln \mathcal{N}(\mathbf{z}) = -\ln \left| \det \left(1 + \varepsilon \frac{\partial \mathbf{v}}{\partial \mathbf{z}} \right) \right. \\ &\frac{1}{1 + \varepsilon} \int \mathbf{v} = 0 \\ \frac{d \ln p(\mathbf{x}, t)}{dt} = -\nabla \cdot \mathbf{v} \end{aligned}$$

Fluid physics behind flows





Simple density

Zhang, E, LW 1809.10188 wangleiphy/MongeAmpereFlow

$$\frac{d}{dt} = \frac{\partial}{\partial t} + \mathbf{v} \cdot \nabla$$
 "material
dt $\frac{\partial}{\partial t} + \mathbf{v} \cdot \nabla$ derivative"

Complex density



Neural Ordinary Differential Equations

Residual network



$$\boldsymbol{x}_{t+1} = \boldsymbol{x}_t + \boldsymbol{v}(\boldsymbol{x}_t)$$

Chen et al, 1806.07366

ODE integration



 $d\mathbf{x}/dt = \mathbf{v}(\mathbf{x})$

Harbor el al 1705.03341 Lu et al 1710.10121, E Commun. Math. Stat 17'...



Neural Ordinary Differential Equations

Residual network





Chen et al, 1806.07366

Harbor el al 1705.03341 Lu et al 1710.10121, E Commun. Math. Stat 17'...







Continuous normalizing flows implemented with NeuralODE

Chen et al, 1806.07366, Grathwohl et al 1810.01367

Target



Density





Samples





Continuous normalizing flow have no structural constraints on the transformation Jacobian

Continuous normalizing flows implemented with NeuralODE

Chen et al, 1806.07366, Grathwohl et al 1810.01367

Target



Density





Samples





Continuous normalizing flow have no structural constraints on the transformation Jacobian

The two use cases

(a) Density estimation

 $\int 0$

"learn from data" $\mathscr{L} = -\mathbb{E}_{\mathbf{x} \sim \text{dataset}} \left[\ln p(\mathbf{x}) \right]$

(b) Variational free energy

 \mathbf{f}^T

"learn from Hamiltonian" $F = \mathop{\mathbb{E}}_{x \sim p(x)} \left[k_B T \ln p(x) + H(x) \right]$



Example: Classical Coulomb gas in a harmonic trap



https://github.com/fermiflow/FermiFlow/blob/github/classical_coulomb_gas.ipynb



Fermi Flow

Xie, Zhang, LW, 2105.08644, JML '22

github.com/fermiflow

Continuous flow of electron density in a quantum dot

Fermi Flow

Xie, Zhang, LW, 2105.08644, JML '22

github.com/fermiflow

Continuous flow of electron density in a quantum dot

Training: Monte Carlo Gradient Estimators

 $\nabla_{\boldsymbol{\theta}} \mathbb{E}_{\boldsymbol{x} \sim p_{\boldsymbol{\theta}}} \left| f(\boldsymbol{x}) \right|$

Score function estimator (REINFORCE)

$$\nabla_{\theta} \mathbb{E}_{\boldsymbol{x} \sim p_{\theta}} \left[f(\boldsymbol{x}) \right] = \mathbb{E}_{\boldsymbol{x} \sim p_{\theta}} \left[f(\boldsymbol{x}) \nabla_{\theta} \ln p_{\theta}(\boldsymbol{x}) \right]$$

Pathwise estimat

tor (Reparametrization trick)
$$\boldsymbol{x} = g_{\theta}(\boldsymbol{z})$$

 $\nabla_{\theta} \mathbb{E}_{\boldsymbol{x} \sim p_{\theta}} \left[f(\boldsymbol{x}) \right] = \mathbb{E}_{\boldsymbol{z} \sim \mathcal{N}(\boldsymbol{z})} \left[\nabla_{\theta} f(g_{\theta}(\boldsymbol{z})) \right]$

Review: 1906.10652

Reinforcement learning Variational inference Variational Monte Carlo Variational quantum algorithms

 $\bullet \bullet \bullet$



10.1 Guidance in Choosing Gradient Estimators

With so many competing approaches, we offer our rules of thumb in choosing an estimator, which follow the intuition we developed throughout the paper:

- If our estimation problem involves continuous functions and measures that are continuous in the domain, then using the pathwise estimator is a good default. It is relatively easy to implement and a default implementation, one without other variance reduction, will typically have variance that is low enough so as not to interfere with the optimisation.
- If the cost function is not differentiable or a black-box function then the score-function or the measure-valued gradients are available. If the number of parameters is low, then the measurevalued gradient will typically have lower variance and would be preferred. But if we have a high-dimensional parameter set, then the score function estimator should be used.
- If we have no control over the number of times we can evaluate a black-box cost function, effectively only allowing a single evaluation of it, then the score function is the only estimator of the three we reviewed that is applicable.
- The score function estimator should, by default, always be implemented with at least a basic variance reduction. The simplest option is to use a baseline control variate estimated with a running average of the cost value.
- When using the score-function estimator, some attention should be paid to the dynamic range of the cost function and its variance, and to find ways to keep its value bounded within a reasonable range, e.g., transforming the cost so that it is zero mean, or using a baseline.
- For all estimators, track the variance of the gradients if possible and address high variance by using a larger number of samples from the measure, decreasing the learning rate, or clipping the gradient values. It may also be useful to restrict the range of some parameters to avoid extreme values, e.g., by clipping them to a desired interval.
- The measure-valued gradient should be used with some coupling method for variance reduction. Coupling strategies that exploit relationships between the positive and negative components of the density decomposition, and which have shared sampling paths, are known for the commonly-used distributions.
- If we have several unbiased gradient estimators, a convex combination of them might have lower variance than any of the individual estimators.
- If the measure is discrete on its domain then the score-function or measure-valued gradient are available. The choice will again depend on the dimensionality of the parameter space.
- In all cases, we strongly recommend having a broad set of tests to verify the unbiasedness of the gradient estimator when implemented.

Mohamed et al, 1906.10652 $\nabla_{\boldsymbol{\theta}} \mathbb{E}_{\boldsymbol{x} \sim p_{\boldsymbol{\theta}}} \left[f(\boldsymbol{x}) \right]$

When to use which?

More discussions Roeder et al, 1703.09194 Vaitl et al 2206.09016, 2207.08219





$$\eta = \nabla_{\theta} \int \mathcal{N}(x|\mu)$$

Score function Score function + variance reduction Value of the cost



https://github.com/deepmind/mc_gradients Mohamed et al, 1906.10652

 $(\mu, \sigma^2) f(x; k) dx; \quad \theta \in \{\mu, \sigma\}$







Case study: Normalizing flow for atomic solids

Variational free energy with a really deep (and a bit awkward) permutation equivariant flow



System	N	LFEP	LBAR	MBAR
LJ	$\begin{array}{c} 256 \\ 500 \end{array}$	3.10800(28)	3.10797(1)	3.10798(9)
LJ		3.12300(41)	3.12264(2)	3.12262(10

$$\ln Z = \ln \mathbb{E}_{x \sim p(x)} \left[e^{-\beta E(x) - \ln p(x)} \right]$$

free energy perturbation (Zwanzig 1954) $\ln Z_B - \ln Z_A = \ln \mathbb{E}_A \left[e^{-\beta (E_B - E_A)} \right]$

Wirnsberger et al, 2111.08696 <u>https://github.com/deepmind/flows_for_atomic_solids</u>





Normalizing flow for atomic solids

F. Hardware details and computational cost

For our flow experiments, we used 16 A100 GPUs to train each model on the bigger systems (512-particle mW and 500-particle LJ). It took approximately 3 weeks of training to reach convergence of the free-energy estimates. Obtaining 2M samples for evaluation took approximately 12 hours on 8 V100 GPUs for each of these models.

For each baseline MBAR estimate, we performed 100 separate simulations for LJ and 200 for mW, corresponding to the number of stages employed. These simulations were performed with LAMMPS [8] and each of them ran on multiple CPU cores communicating via MPI. We used 4 cores for the 64-particle and 216-particle mW experiments and 8 cores for all other systems. The MD simulations completed after approximately 11 and 14 hours for LJ (256 and 500 particles), and 7, 20 and 48 hours for mW (64, 216 and 512 particles). To evaluate the energy matrix for a single MBAR

Will the score-function gradient estimator do better? Heavy lifting is mostly due to preserving permutation. But, does it really matter?





A tale of three equations

Langevin equation (SDE)

$q(\mathbf{x}_{t+dt} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_t + \mathbf{f}dt, 2Tdt\mathbf{I})$

$$\frac{\partial p(\boldsymbol{x},t)}{\partial t} + \nabla \cdot$$

"Particle method" (ODE)

 $d\mathbf{x}$ dt

(Another way to reverse the diffusion is via the reverse-time SDE Anderson 1982)

or
$$x_{t+dt} - x_t = fdt + \sqrt{2Tdt}\mathcal{N}(0,1)$$

Fokker-Planck equation (PDE)

$$\left[p(\boldsymbol{x},t)\boldsymbol{f}\right] - T\nabla^2 p(\boldsymbol{x},t) = 0$$

 $= \boldsymbol{f} - T \nabla \ln p(\boldsymbol{x}, t)$

Maoutsa et al, 2006.00702 Song et al, 2011.13456



$$\mathcal{P}(\vec{x},t) = \int d^{3}\vec{x}' \left(\frac{1}{4\pi D\epsilon}\right)^{3/2} \exp\left[-\frac{\left(\vec{x}-\vec{x}'-\epsilon\vec{v}(\vec{x}')\right)^{2}}{4D\epsilon}\right] \mathcal{P}(\vec{x}',t-\epsilon), \quad (9.18)$$
simplified by the change of variables,

$$\vec{y} = \vec{x}' + \epsilon\vec{v}(\vec{x}') - \vec{x} \implies$$

$$d^{3}\vec{y} = d^{3}\vec{x}' \left(1 + \epsilon\nabla \cdot \vec{v}(\vec{x}')\right) = d^{3}\vec{x}' \left(1 + \epsilon\nabla \cdot \vec{v}(\vec{x}) + \mathcal{O}(\epsilon^{2})\right). \quad (9.19)$$
ping only terms at order of ϵ , we obtain

$$t, t) = \left[1 - \epsilon\nabla \cdot \vec{v}(\vec{x})\right] \int d^{3}\vec{y} \left(\frac{1}{4-\epsilon}\right)^{3/2} e^{-\frac{y^{2}}{4D\epsilon}} \mathcal{P}(\vec{x}+\vec{y}-\epsilon\vec{v}(\vec{x}),t-\epsilon)$$

and

$$\begin{aligned} p(\vec{x},t) &= \int d^{3}\vec{x}' \left(\frac{1}{4\pi D\epsilon}\right)^{3/2} \exp\left[-\frac{\left(\vec{x}-\vec{x}'-\epsilon\vec{v}(\vec{x}')\right)^{2}}{4D\epsilon}\right] \mathcal{P}(\vec{x}',t-\epsilon), \end{aligned} \tag{9.18} \\ \text{mplified by the change of variables,} \\ \vec{y} &= \vec{x}' + \epsilon\vec{v}(\vec{x}') - \vec{x} \implies \\ d^{3}\vec{y} &= d^{3}\vec{x}' \left(1 + \epsilon\nabla \cdot \vec{v}(\vec{x}')\right) = d^{3}\vec{x}' \left(1 + \epsilon\nabla \cdot \vec{v}(\vec{x}) + \mathcal{O}(\epsilon^{2})\right). \end{aligned} \tag{9.19} \\ \text{ng only terms at order of } \epsilon, \text{ we obtain} \\ t) &= \left[1 - \epsilon\nabla \cdot \vec{v}(\vec{x})\right] \int d^{3}\vec{y} \left(\frac{1}{4-\epsilon}\right)^{3/2} e^{-\frac{y^{2}}{4D\epsilon}} \mathcal{P}(\vec{x}+\vec{y}-\epsilon\vec{v}(\vec{x}),t-\epsilon) \end{aligned}$$

Keep

$$\mathcal{P}(\vec{x},t) = \left[1 - \epsilon \nabla \cdot \vec{v}(\vec{x})\right] \int d^{3}\vec{y} \left(\frac{1}{4\pi D\epsilon}\right)^{3/2} e^{-\frac{y^{2}}{4D\epsilon}} \mathcal{P}(\vec{x} + \vec{y} - \epsilon \vec{v}(\vec{x}), t - \epsilon)$$

$$= \left[1 - \epsilon \nabla \cdot \vec{v}(\vec{x})\right] \int d^{3}\vec{y} \left(\frac{1}{4\pi D\epsilon}\right)^{3/2} e^{-\frac{y^{2}}{4D\epsilon}}$$

$$\times \left[\mathcal{P}(\vec{x},t) + (\vec{y} - \epsilon \vec{v}(\vec{x})) \cdot \nabla \mathcal{P} + \frac{y_{i}y_{j} - 2\epsilon y_{i}v_{j} + \epsilon^{2}v_{i}v_{j}}{2} \nabla_{i}\nabla_{j}\mathcal{P} - \epsilon \frac{\partial \mathcal{P}}{\partial t} + \mathcal{O}(\epsilon^{2})\right]$$

$$= \left[1 - \epsilon \nabla \cdot \vec{v}(\vec{x})\right] \left[\mathcal{P} - \epsilon \vec{v} \cdot \nabla + \epsilon D \nabla^{2}\mathcal{P} - \epsilon \frac{\partial \mathcal{P}}{\partial t} + \mathcal{O}(\epsilon^{2})\right].$$
(9.20)
Equating terms at order of ϵ leads to the Fokker–Planck equation,

$$\frac{\partial \mathcal{P}}{\partial t} + \nabla \cdot \vec{J} = 0, \quad \text{with} \quad \vec{J} = \vec{v} \,\mathcal{P} - D \nabla \mathcal{P}.$$
(9.21)

$$\frac{\partial \mathcal{P}}{\partial t} + \nabla \cdot \vec{J} = 0$$
, with $\vec{J} = \vec{v} \,\mathcal{P} - D\nabla \mathcal{P}$.

from Langevin to Fokker-Planck





The conditional trick (originated from denoising score matching Vincent 2011)



$$\begin{aligned} \mathbf{x}_{t} | \mathbf{x}_{0} \rangle \| \mathbf{s}_{\boldsymbol{\theta}}(\mathbf{x}_{t}, t) - \nabla_{\mathbf{x}_{t}} \log q_{t}(\mathbf{x}_{t} | \mathbf{x}_{0}) \|_{2}^{2} \\ \mathbf{x}_{t} & \mathbf{neural} & \mathbf{score of diffused} \\ \mathbf{x}_{t} & \mathbf{network} & \mathbf{data sample} \end{aligned}$$

https://cvpr2022-tutorial-diffusion-models.github.io/



Lessons from diffusion models

Continuous normalizing flow still has great potential.

Going beyond maximum likelihoo

Break the loss into small pieces,)sample them \mathbf{X}_{0}

The conditional trick (originated from denoising score matching Vincent 2011)

$$\min_{\boldsymbol{\theta}} \mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{\mathbf{x}_t \sim q_t(\mathbf{x}_t)} || \mathbf{s}_{\boldsymbol{\theta}}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t) ||_2^2$$

$$\text{diffusion diffused neural score of data } \mathbf{x}_t$$

$$\max_{t \in T} t = t$$

$$\nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t)$$



(marginal)

https://cvpr2022-tutorial-diffusion-models.github.io/





Claim:

$$\mathbb{E}_{x \sim q(x)} |s_{\theta}(x) - \nabla_{x} \ln q(x)|^{2} = \mathbb{E}_{x_{0} \sim q_{0}(x_{0})} \mathbb{E}_{x \sim q(x|x_{0})} |s_{\theta}(x) - \nabla_{x} \ln q(x|x_{0})|^{2} + \text{const.}$$

$$q(x) = \int q(x|x_{0})q_{0}(x_{0})dx_{0}$$
Independent of θ

$$\mathbb{E}_{x_{0} \sim q_{0}(x_{0})} \mathbb{E}_{x \sim q(x|x_{0})} |s|^{2} = \int dx_{0} \int dx q_{0}(x_{0}) q(x|x_{0}) |s|^{2} = \int dx q(x) |s|^{2} = \mathbb{E}_{x \sim q(x)} |s^{2}|$$

$$\mathbb{E}_{x_{0} \sim q_{0}(x_{0})} \mathbb{E}_{x \sim q(x|x_{0})} [s \cdot \nabla \ln q(x|x_{0})] = \int dx_{0} \int dx q_{0}(x_{0}) q(x|x_{0}) \frac{s \cdot \nabla q(x|x_{0})}{q(x|x_{0})}$$

$$= \int dx_{0} \int dx q_{0}(x_{0}) s \cdot \nabla q(x|x_{0})$$

$$\mathbb{E}_{x_{0} \sim q_{0}(x_{0})} \mathbb{E}_{x \sim q(x|x_{0})} |s|^{2} = \int dx_{0} \int dx q_{0}(x_{0}) q(x|x_{0}) |s|^{2} = \int dx q(x) |s|^{2} = \mathbb{E}_{x \sim q(x)} |s^{2}|$$

$$\mathbb{E}_{x_{0} \sim q_{0}(x_{0})} \mathbb{E}_{x \sim q(x|x_{0})} [s \cdot \nabla \ln q(x|x_{0})] = \int dx_{0} \int dx q_{0}(x_{0}) q(x|x_{0}) \frac{s \cdot \nabla q(x|x_{0})}{q(x|x_{0})}$$

$$= \int dx_{0} \int dx q_{0}(x_{0}) s \cdot \nabla q(x|x_{0})$$

$$= \int dx s \cdot \nabla q(x) = \mathbb{E}_{x \sim q(x)} [s \cdot \nabla \ln q(x)]$$

要让 $s_{\theta}(x_t, t)$ 等于 $\nabla_{x_t} \log p(x_t | x_0)$ 的加权平均【即 $\nabla_{x_t} \log p(x_t)$ 】只需要最小化 $\|s_{\theta}(x_t, t) - \nabla_{x_t} \log p(x_t | x_0)\|^2$ 的加权平均 https://spaces.ac.cn/archives/9209





Flow matching



 $\mathscr{L}_{\mathrm{FM}} = \mathbb{E}_{t \sim \mathscr{U}(0,1)} \mathbb{E}_{\mathbf{x}}$

$$\mathbf{x} \sim p(\mathbf{x},t) \left\| \mathbf{v}_{\theta}(\mathbf{x},t) - \mathbf{u}(\mathbf{x},t) \right\|^{2}$$

Albergo et al, 2209.15571, Lipman et al, 2210.02747



$$\frac{\partial p(\boldsymbol{x} \,|\, \boldsymbol{x}_1, t)}{\partial t} + \nabla$$

$$p(\boldsymbol{x}, t) = \int p(\boldsymbol{x} | \boldsymbol{x}_1, t) q(\boldsymbol{x}_1) d\boldsymbol{x}_1$$

$$\mathscr{L}_{\text{CFM}} = \mathbb{E}_{t \sim \mathscr{U}(0,1)} \mathbb{E}_{\boldsymbol{x}_1 \sim q(\boldsymbol{x}_1)} \mathbb{E}_{\boldsymbol{x} \sim p(\boldsymbol{x}|\boldsymbol{x}_1,t)} \left| \boldsymbol{v}_{\theta}(\boldsymbol{x},t) - \boldsymbol{u}(\boldsymbol{x}|\boldsymbol{x}_1,t) \right|^2$$

 $\nabla_{\theta} \mathcal{L}$

$$\boldsymbol{x} = (1 - t)\boldsymbol{x}_0$$
$$\boldsymbol{x}_0 \sim \mathcal{N}(0, I)$$

an example:

Conditional flow matching

$$\nabla \cdot \left[p(\boldsymbol{x} \,|\, \boldsymbol{x}_1, t) \boldsymbol{u}(\boldsymbol{x} \,|\, \boldsymbol{x}_1, t) \right] = 0$$

$$p(\boldsymbol{x}, t)\boldsymbol{u}(\boldsymbol{x}, t) = \int p(\boldsymbol{x} \,|\, \boldsymbol{x}_1, t) \boldsymbol{u}(\boldsymbol{x} \,|\, \boldsymbol{x}_1, t) \, q(\boldsymbol{x}_1) d\boldsymbol{x}_1$$

$$\mathscr{L}_{\mathrm{FM}} = \nabla_{\theta} \mathscr{L}_{\mathrm{CFM}}$$

 $p(\boldsymbol{x} | \boldsymbol{x}_1, t) = \mathcal{N}\left(t\boldsymbol{x}_1, (1-t)^2\right)$ $+tx_1$ $u(x | x_1, t) = dx/dt = x_1 - x_0$

 $\nabla_{\theta} \mathscr{L}_{\text{FM}} = \nabla_{\theta} \mathscr{L}_{\text{CFM}}$ Claim:

where
$$\mathscr{L}_{\text{FM}} = \mathbb{E}_{t \sim \mathscr{U}(0,1)} \mathbb{E}_{\boldsymbol{x} \sim p(\boldsymbol{x},t)} \left| \boldsymbol{v}_{\theta}(\boldsymbol{x},t) - \boldsymbol{u}(\boldsymbol{x},t) \right|^{2}$$

 $\mathscr{L}_{\text{CFM}} = \mathbb{E}_{t \sim \mathscr{U}(0,1)} \mathbb{E}_{\boldsymbol{x}_{1} \sim q(\boldsymbol{x}_{1})} \mathbb{E}_{\boldsymbol{x} \sim p(\boldsymbol{x}|\boldsymbol{x}_{1},t)} \left| \boldsymbol{v}_{\theta}(\boldsymbol{x},t) - \boldsymbol{u}(\boldsymbol{x}|\boldsymbol{x}_{1},t) \right|^{2}$
 $p(\boldsymbol{x},t) = \int p(\boldsymbol{x}|\boldsymbol{x}_{1},t) q(\boldsymbol{x}_{1}) d\boldsymbol{x}_{1} \quad p(\boldsymbol{x},t) \boldsymbol{u}(\boldsymbol{x},t) = \int p(\boldsymbol{x}|\boldsymbol{x}_{1},t) \boldsymbol{u}(\boldsymbol{x}|\boldsymbol{x}_{1},t) q(\boldsymbol{x}_{1}) d\boldsymbol{x}_{1}$

Proof:

$$\mathbb{E}_{\boldsymbol{x}_{1}\sim q(\boldsymbol{x}_{1})}\mathbb{E}_{\boldsymbol{x}\sim p(\boldsymbol{x}|\boldsymbol{x}_{1},t)}\left|\boldsymbol{v}_{\theta}\right|^{2} = \int d\boldsymbol{x}_{1} \int d\boldsymbol{x}q(\boldsymbol{x}_{1})p(\boldsymbol{x}|\boldsymbol{x}_{1},t)\left|\boldsymbol{v}_{\theta}\right|^{2} = \int d\boldsymbol{x}p(\boldsymbol{x},t)\left|\boldsymbol{v}_{\theta}\right|^{2} = \mathbb{E}_{\boldsymbol{x}\sim p(\boldsymbol{x},t)}$$

$$\mathbb{E}_{\boldsymbol{x}_1 \sim q(\boldsymbol{x}_1)} \mathbb{E}_{\boldsymbol{x} \sim p(\boldsymbol{x}|\boldsymbol{x}_1,t)} \left[\boldsymbol{v}_{\theta} \cdot \boldsymbol{u}(\boldsymbol{x} \mid \boldsymbol{x}_1,t) \right] = \int d\boldsymbol{x}_1 \int d\boldsymbol{x}_1 d\boldsymbol{x}_$$

$$= \int d\mathbf{x} p(\mathbf{x}, t) \mathbf{v}_{\theta} \cdot \mathbf{u}(\mathbf{x}, t) = \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}, t)} \left[\mathbf{v}_{\theta} \cdot \mathbf{u}(\mathbf{x}, t) \right]$$

 $\int d\mathbf{x} q(\mathbf{x}_1) p(\mathbf{x} \,|\, \mathbf{x}_1, t) \left[\mathbf{v}_{\theta} \cdot \mathbf{u}(\mathbf{x} \,|\, \mathbf{x}_1, t) \right]$



Flow matching is all you need!

This framework contains various diffusion models as special cases

(Lipman et al, 2210.02747)

Fun to try: flow matching for computing free energy difference

- Greater freedom (and optimal transport theory) in the interpolation path
- 400x speedup compared to continuous normalizing flow (Albergo et al, 2209.15571)
- Surpasses diffusion model on Imagenet in likelihood and sample quality

Fun to try: Train Riemannian flows with this Part II Optimal transport and Riemannian geometry



🙆 Springer

Likelihood free simulator

Prone to mode collapse

More tricky to train than others

Performance have been surpassed by diffusion models

I found it is less useful to scientific applications

GAN

Close connection to variational calculus we have just learned

 $p(x) = \frac{e^{-E(x)}}{Z}$

Variational free energy

0423 ۲ 2200

Approximate sampling and estimation of partition functions using neural networks

> **George T. Cantwell** Santa Fe Institute, 1399 Hyde Park Road, Santa Fe, NM, 87501 gcant@umich.edu

VAE

 $p(z \mid x) = \frac{p(x \mid z)p(z)}{p(x)}$

Variational Bayes

We consider the closely related problems of sampling from a distribution known up to a normalizing constant, and estimating said normalizing constant. We show how variational autoencoders (VAEs) can be applied to this task. In their standard applications, VAEs are trained to fit data drawn from an unknown and intractable distribution. We invert the logic and train the VAE to fit a simple and tractable distribution, on the assumption of a complex and intractable latent distribution, specified up to normalization. This procedure constructs approximations without the use of training data or Markov chain Monte Carlo sampling. We illustrate our method on three examples: the Ising model, graph clustering, and ranking.



Deep variational free energy for dense hydrogen Xie, Li, Wang, Zhang, LW, 2209.06095

 $F = \mathbb{E}_{\substack{S \sim p(S)}} k_B T \ln p($

Classical protons coupled to ground state electrons



$$\psi(S) + \mathbb{E}_{R \sim |\psi_S(R)|^2} \left[\frac{H\psi_S(R)}{\psi_S(R)} \right]$$



Deep variational free energy for dense hydrogen Xie, Li, Wang, Zhang, LW, 2209.06095

 $F = \mathbb{E}_{\substack{S \sim p(S)}} k_B T \ln p($

Classical protons coupled to ground state electrons



$$\psi(S) + \mathbb{E}_{R \sim |\psi_S(R)|^2} \left[\frac{H\psi_S(R)}{\psi_S(R)} \right]$$



The dense hydrogen problem



Generative model for proton probability density distribution Deep neural network (Ferminet) for electron wavefunction

Xie, Li, Wang, Zhang, LW, 2209.06095






Dense hydrogen in the sky and in the lab

Jupiter interior



Inertial confinement fusion



Equation-of-state is the input for hydrodynamics simulations





Generative models on a manifold

Issues

Symmetry in generative models

Coarse graining to slow/collective variables

Accelerated sampling in/with generative models

ML accelerated sampling

- 1. Cheap surrogate function for MC weight Neal 96' Jun. S Liu of A recommender engine for MC updates when the surrogate is a generative model: Huang, LW, 1610.02746, Liu, Qi, Meng, Fu, 1610.03137
- 2. Reinforcement learning the transition kernel: Song et al, State Transition 1706.07561, Levy et al 1711.09268, Cusumano-Towner et al 1801.03612, Bojesen, 1808.09095
- 3. Performs MC in the variationally learned disentangled representation: Wavelet MC, Ismail 03'





