Replacing Markov chain Monte Carlo with generative flow neural networks



Kimmy Cushman¹

with, Yin Lin², Ciaran Hughes³, Joshua Isaacson³, George Fleming¹, and James Simone³

Yale University
 2 University of Chicago
 3 Fermi National Accelerator Laboratory

A.I. for Nuclear Physics Jefferson Lab 4 March 2020







- Lattice QCD
- MCMC and neural networks
- Neural networks for O(3) spin model
- How to believe a neural network
- Future work



Lattice QCD

Yale Kimmy Cushman | Replacing Markov chain Monte Carlo with generative flows neural networks | 4 March 2020 | 45

Lattice QCD... what is it?

deep LGT

- Gauge field simulation from first principles
- Discretized space-time
- Finite volume with periodic conditions
- Volume few times size of proton



http://www.physics.adelaide.edu.au/cssm/lattice/

Motivation for QCD



Mass comes from QCD dynamics



Spectroscopy, scattering, PDFs,

CKM matrix elements BSM spectroscopy,

Phase transitions

Ab-initio Determination of Light Hadron Masses

S. Dürr¹, Z. Fodor^{1,2,3}, J. Frison⁴, C. Hoelbling^{2,3,4}, R. Hoffmann², S. D. Katz^{2,3}, S. Krieg², T. Kurth², L. Lellouch⁴, T. Lippert^{2,5}, K.K. Szabo², G. Vulvert⁴

¹NIC, DESY Zeuthen, D-15738 Zeuthen and FZ Jülich, D-52425 Jülich, Germany.
 ²Bergische Universität Wuppertal, Gaussstr. 20, D-42119 Wuppertal, Germany.
 ³Institute for Theoretical Physics, Eötvös University, H-1117 Budapest, Hungary.
 ⁴Centre de Physique Théorique; Case 907, Campus de Luminy, F-13288 Marseille Cedex 9, France.
 ⁵Jülich Supercomputing Centre, FZ Jülich, D-52425 Jülich, Germany.

Budapest-Marseille-Wuppertal Collaboration





$$\langle \mathcal{O} | \pi(t) \ \pi^{\dagger}(0) | \mathcal{O} \rangle = C_{\pi}(t) = \langle \mathcal{O} \rangle$$
$$\langle \mathcal{O} \rangle = \frac{\int_{D} \mathcal{O}[A] \ e^{iS[A]}}{\int_{D} e^{iS[A]}}$$



$$\langle 0 | \pi(t) \pi^{\dagger}(0) | 0 \rangle = C_{\pi}(t) = \langle \mathcal{O} \rangle$$





$$\langle 0 | \pi(t) \pi^{\dagger}(0) | 0 \rangle = C_{\pi}(t) = \langle \mathcal{O} \rangle$$





 $\langle 0 | \pi(t) \pi^{\dagger}(0) | 0 \rangle = C_{\pi}(t) = \langle \mathcal{O} \rangle$







Configuration 1 \mathcal{O}_1



Configuration 2 \mathcal{O}_2



Configuration 3 \mathcal{O}_3



Configuration 4 \mathcal{O}_4

$$\langle \mathcal{O} \rangle = \frac{1}{N} \sum_{i}^{N} \mathcal{O}_{i}$$

Yale



MCMC and neural networks



- **Propose possible configuration**, *x*
- **Propose steps in a chain** $x \rightarrow x'$ **accept with**

probability
$$p = \min \left[1, \frac{e^{-S(x')}}{e^{-S(x)}}\right]$$





Drawback: thermalization and autocorrelation



Ising model



Number independent \ll generated configurations

Number configurations





Drawback: critical slowing down near critical points



Ising model ordered/disorder phases

QCD + dark matter Confined phase/quark gluon plasma

BSM Higgs Electroweak symmetry breaking transition

Deep neural networks



Fully connected network with 4 hidden layers



Deep neural networks



Fully connected network with 4 hidden layers







https://www.thispersondoesnotexist.com





https://www.thispersondoesnotexist.com





https://www.thispersondoesnotexist.com

Auto-encoders



deep



Auto-encoders



GANs





Auto-encoders GANs

Require A LOT of training data Mode collapse







Generative flow networks for O(3) spin model



Goal: Train neural network to produce configurations by learning *bijection*



See also Albergo, Kanwar, Shanahan 1904.12072 Dinh, Sohl-Dickstein, Bengio 1605.08803 Muller, McWilliams, Rousselle, Gross, Novak, 1808.03856



Why a bijection? For generation $U(x) \rightarrow q(x) (\approx p(x))$

Why $q(x) \rightarrow U(x)$?

Typical importance sampling

$$\begin{split} \langle \mathcal{O} \rangle &= \frac{1}{Z} \int dx \ \mathcal{O}(x) \ \mathrm{e}^{-S(x)} \\ &= \frac{1}{Z} \int dx \ \mathcal{O}(x) \ p(x) \\ &\approx \frac{1}{N} \sum_{i} \ \mathcal{O}_{i} \big|_{p} \end{split}$$



Why a bijection? For generation $U(x) \rightarrow q(x) (\approx p(x))$

Why $q(x) \rightarrow U(x)$?

Typical importance sampling

$$\langle \mathcal{O} \rangle = \frac{1}{Z} \int dx \ \mathcal{O}(x) \ e^{-S(x)}$$

= $\frac{1}{Z} \int dx \ \mathcal{O}(x) \ p(x)$

$$\approx \frac{1}{N} \sum_{i} \mathcal{O}_{i} |_{p}$$

Generative flow importance sampling

$$\begin{split} \left| \mathcal{O} \right\rangle &= \frac{1}{Z} \int dx \ \mathcal{O}(x) \ p(x) \\ &= \frac{1}{Z} \int dx \ \mathcal{O}(x) \ \frac{p(x)}{q(x)} \ q(x) \\ &\approx \frac{\sum_{i} \mathcal{O}_{i} \frac{p_{i}}{q_{i}} |_{q}}{\sum_{i} \frac{p_{i}}{q_{i}} |_{q}} \end{split}$$



Why a bijection? For generation $U(x) \rightarrow q(x) (\approx p(x))$

Why $q(x) \rightarrow U(x)$?

Typical importance sampling

$$\langle \mathcal{O} \rangle = \frac{1}{Z} \int dx \ \mathcal{O}(x) \ e^{-S(x)}$$

= $\frac{1}{Z} \int dx \ \mathcal{O}(x) \ p(x)$

Yale

 $\approx \frac{1}{N} \sum_{i} \mathcal{O}_{i} |_{p}$ Sampled according to p(x)Hard to sample from (MCMC)

Generative flow importance sampling

$$\langle \mathcal{O} \rangle = \frac{1}{Z} \int dx \ \mathcal{O}(x) \ p(x)$$

= $\frac{1}{Z} \int dx \ \mathcal{O}(x) \ \frac{p(x)}{q(x)} \ q(x)$
 $\approx \frac{\sum_{i} \mathcal{O}_{i} \frac{p_{i}}{q_{i}}|_{q}}{\sum_{i} \frac{p_{i}}{q_{i}}|_{q}}$
Sampled according to $q(x)$
Easy to sample from ($U(x)$ is easy)



Composition of coupling layers $H = h_L \circ \cdots \circ h_2 \circ h_1$ H(U(x)) = q(x) **Learned Jacobian of variable transformation** $q(x_i) = U(x_i) \left| \frac{\partial H(x; \theta)}{\partial x_i} \right|^{-1}$

Train by minimizing divergence between q(x) and p(x)

p(x)

Neural network architecture





Neural network architecture



Choose coupling layers function $C(x_A, m(x_B))$ *easily invertible*

Ex: $C(x_A; s, t) = x_A \odot e^s + t$ \vec{s} and \vec{t} learned parameters

Neural network architecture

deep LCT





deep LCT





QCD vs O(3) model





Lattice QCD (pure gauge)

$$S = \frac{\beta}{3} \sum_{n} \sum_{\mu < \nu} \operatorname{Re} \operatorname{tr} \left[1 - U_{\mu\nu} \right]$$

 $8 \times 4 N^4$ degrees of freedom

SU(3) local gauge symmetry



O(3) spin model

$$S = -\beta \sum_{\langle i,j \rangle} \vec{s}_i \cdot \vec{s}_j$$

 $2N^2$ degrees of freedom

O(3) symmetry

O(3) spin model training

- 1) Choose system size
- 2) Choose hyper parameters
- **3)** Train loss function
- 4) Validation

dee

Choose system



Choose hyper parameters



- 1) Number of samples 5000
- 2) Number of networks 5
- 3) Number of layers 6



Disney Research Neural Important Sampling

4) Coupling transform - piecewise cumulative

distribution function

5) Learning rate, ...

Training against loss



Loss functions

KL
$$L(x, \theta) = \sum \frac{p_i}{q_i} \log\left(\frac{p_i}{q_i}\right)$$

Exp $L(x, \theta) = \sum \frac{p_i}{q_i} \log\left(\frac{p_i}{q_i}\right)^2$
Chi² $L(x, \theta) = \sum \frac{(p_i - q_i)^2}{q_i^2}$
 0^{-1}
 $\beta = 2.0$
 $\beta = 1.5$
 $\beta = 1.0$
 $\beta = 0.5$
 $\beta = 0.5$
Epoch



How to believe a neural network: Validation

Check p and q















$$C(1,\beta=1) = 0.53731 \quad \langle \mathcal{O} \rangle = \sum_{i} \mathcal{O}_{i} \frac{p_{i}}{q_{i}}|_{q} \quad (1,\beta=1) = 0.53731 \quad \langle \mathcal{O} \rangle = \sum_{i} \mathcal{O}_{i} \frac{p_{i}}{q_{i}}|_{q}$$



























2) Random (uniform)
$$\langle \mathcal{O} \rangle = \frac{\sum \mathcal{O}_i p_i}{\sum p_i}$$

3) Generative flows training \langle

$$\langle \mathcal{O} \rangle = \frac{\sum \mathcal{O}_i \left. \frac{p_i}{q_i} \right|_q}{\sum \frac{p_i}{q_i}}$$

Yale







How to believe a neural network: Tests



Similar actions



Long-distance physics should not depend on short distance details of action

$$S = -\sum \beta \ s_x \cdot s_{x\pm 1}$$

$$S = -\sum \beta_1 \ s_x \cdot s_{x\pm 1} + \beta_2 \ s_x \cdot s_{x\pm 2}$$

$$S = -\sum \beta_1 \ s_x \cdot s_{x\pm 1} + \beta_2 \ s_x \cdot s_{x\pm 2} + \beta_3 \ (s_x \cdot s_{x\pm 1})^2$$

$$S = -\sum \beta_{\alpha} S_{\alpha}$$



Similar actions



1) Given 1-term action, predict β from ensemble





Similar actions



M. Hasenbusch et al., Phys Lett B 338 (1994) 308-312

1) Given 1-term action, predict β from ensemble





Similar actions



- 1) Given 1-term action, predict β from ensemble
- 2) Given 2-term action, predict β_1 , β_2 from ensemble





Similar actions



- 1) Given 1-term action, predict β from ensemble
- 2) Given 2-term action, predict β_1 , β_2 from ensemble
- **3)** Compute corresponding "truncated" actions: $\{\beta_1, \beta_2\} \rightarrow \{\beta_1, 0\}$







3) Compute corresponding "truncated" actions: $\{\beta_1, \beta_2\} \rightarrow \{\beta_1, 0\}$





Similar actions



- **3)** Compute corresponding "truncated" actions: $\{\beta_1, \beta_2\} \rightarrow \{\beta_1, 0\}$
- 4) Test neural networks





Future work

Future work

- Complete tests of invariant physics
- Train larger 2D lattices
- Attempt gauge theory Schwinger model
- Lattice QCD (small lattices)
- Cost-benefit analysis





Thank you!





Thank you!





Monte Carlo renormalization group Canonical demon algorithm

$$S = -\sum \beta_{\alpha} S_{\alpha}$$

$$S_D = +\sum \beta_{\alpha} d_{\alpha} \quad d \in [0, d_{\max}]$$



M. Hasenbusch et al., Phys Lett B 338 (1994) 308-312

Micro-canonical MCMC with $S_{\text{tot}} = S + S_D$: $\Delta S_{\text{tot}} = 0$

$$Z = \left(\prod_{\alpha} \int_{0}^{d_{\max}} \mathrm{d}d_{\alpha} \right) \int \mathrm{d}\phi \, \mathrm{e}^{-(S+S_{D})} \quad \Rightarrow \quad \langle d_{\alpha} \rangle = \frac{1}{\beta_{\alpha}} \left(1 - \frac{\beta_{\alpha}d_{\max}}{\mathrm{e}^{\beta_{\alpha}d_{\max}} - 1} \right)$$





Yale















$$C(t) = \langle 0 | \bar{\pi}(t) \pi(0) | 0 \rangle$$

=
$$\sum_{m}^{\infty} \langle 0 | \pi(0) | E_m \rangle e^{-E_m t} \langle E_m | \pi(0) | 0 \rangle$$









http://watersoundimage.yolasite.com/what-is-a-w-s-image.php

$$C(t) = \langle 0 | \bar{\pi}(t) \pi(0) | 0 \rangle$$

=
$$\sum_{m}^{\infty} \langle 0 | \pi(0) | E_m \rangle e^{-E_m t} \langle E_m | \pi(0) | 0 \rangle$$

$$\Rightarrow C(t) = \sum_{m=1}^{\infty} a_m \,\mathrm{e}^{-E_m t}$$



