Theory-driven quantum machine learning for HEP



INFN - Sezione di Genova March 8th, 2023









S

Pre-introduction: What is ML?





The name of the game is "optimisation"



Pre-introduction: What is ML?

e.g. FC NN with 2 layers

$$f_{FC}(\mathbf{x}; \theta) = \sigma_2 \left(\mathscr{W}_2 \cdot \sigma_1 \left(\mathscr{W}_1 \cdot \mathbf{x} + \mathscr{B}_1 \right) + \mathscr{B}_2 \right)$$

$$\forall \theta \in \mathscr{W}_i \text{ or } \mathscr{B}_i$$

σ_i := Activation Function

Pre-introduction: What is ML?

Pre-introduction: What does ML learn?

Pre-introduction: What does ML learn?

 \sim

Pre-introduction: What does ML learn?

Each "ring" corresponds to an output of a filer based on the polynomial jet distribution $y^m \phi^n$.

Sales pitch of the talk!

- We more or less know how to get a well-performing Neural Network to classify jets, LHC events, and even cats and dogs...
- What we don't know is what this network learns.
- Can we use Quantum Mechanics to have more insight into the learning process?
 - What has a model learned?
 - What is learning?
 - How do we develop "insightful" algorithms?
 - How to perform this on a Quantum device?

Can an ML problem be formulated as a quantum manybody system?

Classification as a quantum many-body problem

Hamiltonian learning for anomaly detection

Conclusion

Hello world of HEP-ML: Top tagging

Classification as a quantum many-body problem

Tensor Networks: Origins

$$\begin{split} |\Psi\rangle &= \sum_{\phi_1, \dots, \phi_n = 0} \mathscr{W}_{\phi_1 \dots \phi_n} |\phi_1\rangle \otimes |\phi_2\rangle \otimes \dots \otimes |\phi_n\rangle \\ \\ \forall |\phi_i\rangle &\in \mathscr{H}^{\otimes 2^N} \quad \rightarrow \quad |\phi_i\rangle \in \big\{|\uparrow\rangle, |\downarrow\rangle\big\} \end{split}$$
The computational cost of a rank-N tensor is $\mathcal{O}(d^N)$!!! Computational cost is $\mathcal{O}(d^{N-1}\chi^2)$!!!

Types of Tensor Networks (some of them)

Matrix Product States for Classification

How How

How How

Hov

Sub-Outline
w to embed the data?
w to form a network?
w to train the network?

$$\mathscr{L} = \frac{1}{N} \sum_{x \in \mathbf{x}^{N}} q^{\text{truth}} \log \left(p(x^{(i)}; \theta) \right)$$

$$\mathscr{L} = \frac{1}{N} \sum_{x \in \mathbf{x}^{N}} q^{\text{truth}} \log \left(p(x^{(i)}; \theta) \right)$$

$$\mathscr{L} = \frac{1}{N} \sum_{x \in \mathbf{x}^{N}} q^{\text{truth}} \log \left(p(x^{(i)}; \theta) \right)$$

$$\mathscr{L} = \frac{1}{N} \sum_{x \in \mathbf{x}^{N}} q^{\text{truth}} \log \left(p(x^{(i)}; \theta) \right)$$

$$\mathscr{L} = \frac{1}{N} \sum_{x \in \mathbf{x}^{N}} q^{\text{truth}} \log \left(p(x^{(i)}; \theta) \right)$$

$$\mathscr{L} = \frac{1}{N} \sum_{x \in \mathbf{x}^{N}} q^{\text{truth}} \log \left(p(x^{(i)}; \theta) \right)$$

$$\mathscr{L} = \frac{1}{N} \sum_{x \in \mathbf{x}^{N}} q^{\text{truth}} \log \left(p(x^{(i)}; \theta) \right)$$

$$\mathscr{L} = \frac{1}{N} \sum_{x \in \mathbf{x}^{N}} q^{\text{truth}} \log \left(p(x^{(i)}; \theta) \right)$$

$$\mathscr{L} = \frac{1}{N} \sum_{x \in \mathbf{x}^{N}} q^{\text{truth}} \log \left(p(x^{(i)}; \theta) \right)$$

$$\mathscr{L} = \frac{1}{N} \sum_{x \in \mathbf{x}^{N}} q^{\text{truth}} \log \left(p(x^{(i)}; \theta) \right)$$

$$\mathscr{L} = \frac{1}{N} \sum_{x \in \mathbf{x}^{N}} q^{\text{truth}} \log \left(p(x^{(i)}; \theta) \right)$$

$$\mathscr{L} = \frac{1}{N} \sum_{x \in \mathbf{x}^{N}} q^{\text{truth}} \log \left(p(x^{(i)}; \theta) \right)$$

$$\mathscr{L} = \frac{1}{N} \sum_{x \in \mathbf{x}^{N}} q^{\text{truth}} \log \left(p(x^{(i)}; \theta) \right)$$

$$\mathscr{L} = \frac{1}{N} \sum_{x \in \mathbf{x}^{N}} q^{\text{truth}} \log \left(p(x^{(i)}; \theta) \right)$$

$$\mathscr{L} = \frac{1}{N} \sum_{x \in \mathbf{x}^{N}} q^{\text{truth}} \log \left(p(x^{(i)}; \theta) \right)$$

Traditionally NNs are trained with SGD, but MPS is trained with Density Matrix Renormalisation Group Algorithm

JYA, Spannowsky; JHEP '21, arXiv: 2106.08334

Why finding a quantitative measure is important?

- A 50% reduction in the number of pixels used and a 91% reduction in the number of parameters lead to the same classification quality!
- Understanding the network gives the ability to build better training algorithms.
- Scientific data is largely sparse; if we know where the information comes from, we can get rid of large amounts of data.
- Suppress the noise (and for pile-up mitigation to be confirmed)!

JYA, Spannowsky; JHEP '21, arXiv: 2106.08334

15

Types of Tensor Networks (some of them)

Yet another small intro...

Du, Hsieh, Liu, Tao; Phys. Rev. `20

Yet another small intro...

Types of Tensor Networks (some of them)

What can we gain if we adopt VQC?

Experimenting with 6-Qubits

Ansatz	D	χ	# Parameters	AUG
TTN	2	5	235	0.75
	2	10	1320	0.80
	2	20	9040	0.84
	5	10	1950	0.87
	10	20	14800	0.89
MPS	2	5	230	0.81
	2	10	860	0.81
	2	20	3320	0.81
	5	10	2150	0.89
MERA	2	5	1225	0.85
	2	10	13400	0.84
	2	20	181600	0.84
	5	10	18200	0.90
Q-TTN	-	-	9	0.89
Q-MPS	-	-	9	0.88
Q-MERA	-	-	17	0.91

Loss landscape for classical TNs becomes exponentially flat!

Ability to construct dynamic hybrid architectures!

JYA, Spannowsky; PRA '22, arXiv: 2202.10471

Near-term quantum devices are quite limited; hence hybrid quantum-classical systems are essential

Tensor Network nodes can be dynamically converted into qubits, as more will be available in the future!

Hamiltonian learning for anomaly detection

What has Hamiltonian to do with data?

JYA, Spannowsky; arXiv: 2211.03803

What has Hamiltonian to do with data?

JYA, Spannowsky; arXiv: 2211.03803

What has Hamiltonian to do with data?

JYA, Spannowsky; arXiv: 2211.03803

 $\arg\min_{\theta,\phi} \mathscr{L}_{\theta,\phi}(\sigma_D) \simeq S(\sigma_D)$

Hamiltonian as a discriminator!

JYA, Spannowsky; arXiv: 2211.03803

Trotter-Suzuki approximation
$$e^{-iT\hat{K}_{\theta}} = \prod_{i=1}^{N} e^{-i\Delta i}$$

$$\langle \hat{K} \rangle_{\theta,\phi} = \frac{1}{N_{\text{smp}}} \sum_{\sigma}^{N_{\text{smp}}} \langle \sigma_T^i | \hat{K}_{\theta} | \sigma_T^i \rangle$$

 $t\hat{K}_{\theta}$

Hamiltonian as a discriminator!

JYA, Spannowsky; arXiv: 2211.03803

Trotter-Suzuki approximation
$$e^{-iT\hat{K}_{\theta}} = \prod_{i=1}^{N} e^{-i\Delta i}$$

$t\hat{K}_{\theta}$

Conclusion

Conclusion

- The name of the game is optimisation. Techniques developed for field theory computations are easily transferable for ML applications!
- Designing quantifiable measures from the ansätze can allow us to improve training procedures and can be used for feature selection. (Thought: maybe helpful to lock on a symmetry during training?)
- Tensor Networks are the tool for the near future to understand quantum computing until the machinery is ready for more significant problems.

Correlations by SU(2) generators

Jack Y. Araz - Tensor Networks

Top Tagging through MPS

Experimenting with 4-Qubits

Experimenting with 4-Qubits

Background Rejection $(1/\epsilon_B)$ 10^{2} 10^{1}

Jack Y. Araz - Classical vs Quantum

Singular Value Decomposition

 λ_i also known as Schmidt values

Singular Value Decomposition

Singular Value Decomposition

Computational cost is $\mathcal{O}(d^{N-1}\chi^2)$!!!

Matrix Product States for Classification

×

Data Embedding

$$\Phi^{p_1 \cdots p_n}(\mathbf{x}) = \phi^{p_1}(x_1) \otimes \phi^{p_2}(x_2) \otimes \cdots \otimes \phi^{p_n}(x_n)$$
$$\phi^{p_i}(x_i) = \begin{bmatrix} \cos(x_i \ \pi/2) \\ \sin(x_i \ \pi/2) \end{bmatrix} \text{ or } \phi^{p_i}(x_i) = \begin{bmatrix} 1 \\ x_i \\ x_i^2 \end{bmatrix} \text{ or } \cdots$$

$$\bigcup_{n \in \mathbb{N}} \bigcup_{n \in \mathbb{N}} \bigoplus_{n \in \mathbb{N}} \bigoplus_$$

Density Matrix Renormalization Group Algorithm

Durham Initial

*

Density Matrix Renormalization Group Algorithm

X

Top Tagging through MPS

$$\sum_{i}^{\chi} \lambda_{\alpha} | \alpha \rangle_{A} | \alpha \rangle_{B} \rightarrow \lambda_{\alpha} := \text{Schmidt values}$$

$$\sum_{i}^{\chi} \lambda_i^2 \log \lambda_i^2$$

Top Tagging through MPS

Fisher Information & Effective Dimensions

44

