Tensor Networks for Machine Learning and Applications



E.M. Stoudenmire

Mar 2021 - IPAM Tensor Methods



SIMONS

FOUNDATION

Overview



Tensor networks are a natural way to parameterize interesting and powerful machine learning models

Perspective on this area from physics point of view

Today:

- Overview
- Tensor Networks Architectures & Applications
- Algorithms & Future Directions



Tensors in Machine Learning

Where can tensors appear in machine learning applications?

Multi-Dimensional Data



Image Data

Medical Data



Neural Network Weight Layers



Possible to interpret as a very high-order tensor (not just a matrix) Linear Weights of High-Dimensional Models

$$f(\mathbf{x}) = W \cdot \Phi(\mathbf{x})$$



For certain cases of kernel learning and Gaussian processes, weights are naturally a *high-order tensor*

Why Tensor Networks?

Tensor network = factorization of huge tensor into contracted product of smaller tensors



Tensor network = factorization of huge tensor into contracted product of smaller tensors



Benefits:

- exponential reduction in memory needed
- exponential speedup of computations (addition, product)
- theoretical insight and interpretation
- estimation of missing or corrupted entries
- many optimization algorithms & strategies

Notation – Tensor Diagrams

N-index tensor = shape with N lines



Joining lines means contraction





Best understood tensor network is the *matrix product state (MPS)*^{1,2} or tensor train ³



[1] Östlund, Rommer, PRL 75, 3537 (1995)
[2] Vidal, PRL 91, 147902 (2003)
[3] Oseledets, SIAM J. Sci. Comp. 33, 2295 (2011)

Adjustable parameter of matrix product state (MPS) is bond dimension χ



If modest χ yields good approximation, obtain massive compression:

$$d^N \longrightarrow N d \chi^2$$

Can efficiently sum MPS in compressed form:

multiply by other networks:



and perfectly sample:



In quantum physics, have rich theory of which tensor networks are suited for particular "data"



(Here "data" = samples/measurements of a quantum wavefunction)

Tensor networks a general tool for linear algebra in exponentially high-dimensional spaces



For example, entanglement entropy really just a measure of multilinear tensor rank

Architectures & Applications

Tensor networks beyond finite MPS Most straightforward application of tensor networks to machine learning is using MPS

 $f(\mathbf{x}) = \mathbf{0} - \mathbf{0} - \mathbf{0} - \mathbf{0} - \mathbf{0}$

 map data to tensor product features

2. evaluate MPS model

 optimize MPS tensors for machine learning objective (supervised, unsupervised) Since 2016, tensor network machine learning now successfully "ported" to other tensor net architectures

Infinite MPS



Locally purified states



PEPS

MERA

Infinite MPS

Miller, Rabusseau, Terilla, "Tensor Networks for Probabilistic Sequence Modeling", arxiv:2003.01039



- used to generate model languages with various grammars
- very few parameters and parallel optimization
- superior results to LSTM in many cases, equal in most others
- can *generalize* from training on shorter sequences to correct results on longer sequences (so really learning the grammar)

Locally purified states

Anomaly Detection with Tensor Networks

arxiv:2006.02516



Table 3: Mean AUROC scores (in %) and standard errors on ODDS datasets.

Dataset	OC-SVM	IF	GOAD	DAGMM	TNAD
Wine	60.0	46.0 ± 8.4	48.2 ± 24.7	51.7 ± 19.3	97.3 ± 4.5
Glass	62.0	57.2 ± 1.6	53.5 ± 13.6	52.5 ± 12.9	81.8 ± 7.3
Thyroid	98.8	99.0 ± 0.1	95.8 ± 1.3	88.8 ± 6.8	99.0 ± 0.1
Satellite	79.9	77.2 ± 0.9	60.6 ± 5.3	72.1 ± 4.7	81.3 ± 0.5
Forest	97.7	71.7 ± 2.6	64.6 ± 4.7	60.9 ± 8.9	$\textbf{98.8} \pm \textbf{0.6}$

Novel anomaly detection framework

Results better than neural networks for tabular data

Quantum process tomography with ... tensor networks arxiv:2006.02424



Scalable learning of noisy quantum "channels" or processes from experiments

PEPS (Projected Entangled Pair States)

Cheng, Wang, Zhang, "Supervised Learning with PEPS" arxiv:2009.09932



PEPS = 2D analogue of MPS

• cost to train is high



Framework for learning



MERA for Audio Classification

Reyes, Stoudenmire, "A Multi-Scale Tensor Network Architecture for Classification and Regression" arxiv:2001.08286





Factorize weights as MERA network with trainable MPS weights on top





Fixed MERA layers encode a discrete wavelet transform

Key capability is fine-graining weights through layers:



Algorithms & Theory

including expressivity, data efficiency, model adaptation, ...



Why can tensor networks succeed?

Martyn, Vidal, et al. arxiv:2007.06082

Why can tensor networks succeed?

Two recent papers give different perspective:

- mutual information (MI) may be best quantity to consider
- strategies to estimate MI



- study synthetic Gaussian data
- realistic images have MI like that of near-neighbor correlated Gaussian



- MI of text is beyond area law scaling
- some images are area law, but others may be beyond area law

Why do tensor networks succeed at all?

If images high entanglement, how do tensor networks succeed?



Martyn, Vidal, et al. arxiv:2007.06082

Rapidly increasing entanglement as more images summed

Possible answers:

- Convy et al says: quantity studied by Martyn et al. may be overly sensitive to choice of feature map
- entropy might decrease as more images summed
- supervised learning may require much less resources

Important that we can even pose and answer these questions!

Algorithms could be what sets tensor network models apart

Available algorithms (partial list):

- ALS + gradient descent
- ALS + optimal local update (Stokes, Terilla)
- modified ALS (= 2 site DMRG)
- Riemannian optimization
- TT-cross algorithm
- density matrix algorithm (polynomial-scaling TT-SVD)

An interesting story of two algorithms for a challenging dataset

The dataset: even-parity bit strings



An interesting story of two algorithms for a challenging dataset

Algorithm #1: optimal <u>local</u> update (generative modeling)



Algorithm #1: optimal <u>local</u> update (generative modeling)

- can use geometric reasoning to develop
- excellent results up to length N=20 bitstrings
- linear scaling in data length and training set size
- can fail for poor choice of initial state





Algorithm #2: direct compression / density matrix algorithm

3. diagonalize densitymatrices to obtain MPStensors one at a time

Algorithm #2: direct compression / density matrix algorithm

- cannot get stuck & gives deterministic results
- scales quadratically in training set size
- can develop a theory of generalization:



Outlook & Future Directions

Many opportunities to fix downsides of tensor network optimization algorithms

Main downsides:

- cost to train (memory + time) can be high
- best architecture choices still being explored

Opportunities to improve:

+ use sparse (e.g. block-sparse) tensors within models

+ use infinite tensor networks more

+ use symmetries more (evidence conv. layers help, infinite TNs for translation symmetry, ...)

+ devise more data-efficient approaches (theory can help)

Push interpretability / understanding / theory frontier

Opportunities here:

+ can characterize which data is <u>learnable</u> by a given type and size of tensor network

+ theory of <u>generalization</u> (Bradley et al.) could be improved:

- tighter bounds,
- apply to other data sets
- estimate using summary statistics of real data
- + progress in theory will feed back into better training <u>algorithms</u>

One final thought:

Tensor networks likely just the "right" way to do linear algebra in very high-dimensional spaces (versus being quantum mechanics oriented)

We haven't even figured out the biggest pieces of tensor networks yet, such as analog of QR or SVD factorizations for matrix-like tensor networks

Some major progress only happened recently, such as canonical forms of PEPS (arxiv:1902.05100) or progress in randomized algorithms (arxiv:2003.05101)