

Quantum Doubly Stochastic Transformers

Jannis Born Filip Skogh Kahn Rhrissorrakrai Filippo Utro Nico
Wagner Aleksandros Sobczyk
presenter: Shen Yuan



中國人民大學
RENMIN UNIVERSITY OF CHINA

高瓴人工智能学院
Gaoling School of Artificial Intelligence

Outline

Introduction

Methods

- Doubly-Stochastic Operators

- Expressivity of Doubly-Stochastic Operators

Experiments

Conclusion

Outline

Introduction

Methods

- Doubly-Stochastic Operators

- Expressivity of Doubly-Stochastic Operators

Experiments

Conclusion

Introduction

The dot product self-attention:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \mathbf{A}\mathbf{V} = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\tau}\right)\mathbf{V} \quad (1)$$

where $\mathbf{Q} := \mathbf{X}\mathbf{W}_Q$, $\mathbf{K} := \mathbf{X}\mathbf{W}_K$ and $\mathbf{V} := \mathbf{X}\mathbf{W}_V$. $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{T \times d}$

At the core of the Transformer, the softmax normalizes the attention matrix to be right stochastic. Sinkformer discovered that Transformer attention naturally converge to doubly stochastic matrices (DSMs) over training, i.e. their rows and columns sum to 1.

Introduction

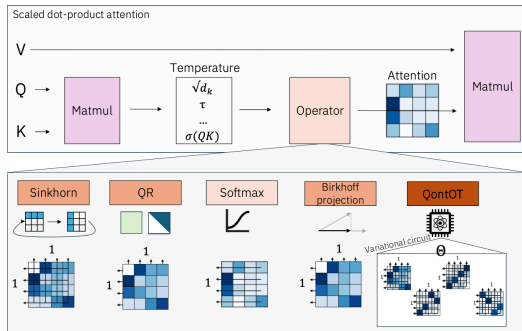


Figure 1. Doubly Stochastic Transformers. Standard scaled dot-product attention applies a Softmax activation on the query-key matrix (*top*). We study different techniques to make attention doubly stochastic attention by replacing the softmax operation (*bottom*). Our proposed Quantum Doubly Stochastic Transformer (QDSFormer) leverages QontOT, a variational quantum circuit with high expressivity.

Outline

Introduction

Methods

Doubly-Stochastic Operators

Expressivity of Doubly-Stochastic Operators

Experiments

Conclusion

Doubly Stochastic Matrices (DSMs)

We denote the n -dimensional vector of ones by $\mathbf{1}_n$ and the $n \times n$ identity matrix as \mathbf{I}_n . The *Birkhoff polytope* $\Omega_n := \mathcal{N}(\mathbf{1}_n, \mathbf{1}_n)$ defines the convex set of $n \times n$ doubly stochastic matrices (DSMs). A DSM $\mathbf{P} \in \Omega_n$ is a nonnegative matrix with row/column sum of 1

$$\mathbf{P}\mathbf{1}_n = \mathbf{1}_n, \quad \mathbf{P}^\top \mathbf{1}_n = \mathbf{1}_n, \quad \mathbf{P}_{i,j} \geq 0. \quad (2)$$

A right stochastic matrix \mathbf{R} has row sums of 1

$$\mathbf{R}\mathbf{1}_n = \mathbf{1}_n, \quad \mathbf{R}_{i,j} \geq 0. \quad (3)$$

and a left stochastic matrix \mathbf{L} has column sums of 1

$$\mathbf{L}^\top \mathbf{1}_n = \mathbf{1}_n, \quad \mathbf{L}_{i,j} \geq 0. \quad (4)$$

Doubly Stochastic Matrices (DSMs)

Birkhoff-von Neumann theorem states that the $n!$ vertices (i.e., extreme points) of the *Birkhoff polytope* Ω_n are permutation matrices, so their entries belong to $\{0, 1\}$.

Notably, every DSM $\mathbf{P} \in \Omega_n$ can be decomposed as a convex combination of permutation matrices, that is $\mathbf{P} = \sum_{i=1}^N \lambda_i \mathbf{\Pi}_i$ for some probability vector $\boldsymbol{\lambda} \in \Delta_N$, $n \times n$ permutation matrices $\{\mathbf{\Pi}_i\}$, and the number of extreme points $N \leq n^2$.

Sinkhorn's algorithm

Sinkhorn's algorithm is based on Sinkhorn's theorem, stating that for any square strictly positive matrix $\mathbf{R} \in \mathbb{R}_+^{T \times T}$, there exist (strictly) positive diagonal matrices $\mathbf{P} = \mathbf{D}_1, \mathbf{D}_2$ s.t., $\mathbf{D}_1 \mathbf{M} \mathbf{D}_2 \in \Omega_T$.

Sinkhorn's algorithm is an approximation procedure that iteratively (K times) normalizes the mass of the rows and the columns of $\mathbf{M} := \mathbf{Q} \mathbf{K}^\top$.

However, it may not converge with few iterations, especially if $\mathbf{Q} \mathbf{K}^\top$ contains large numeric values. Therefore, the Sinkformer is only an approximately doubly stochastic Transformer.

Projection on the Birkhoff polytope

Alternatively, one can project \mathbf{M} directly on Ω_T via $\mathbf{P} = \arg \min_{\mathbf{X} \in \Omega_T} \|\mathbf{X} - \mathbf{M}\|_F^2$, where the set for \mathbf{X} and the objective are convex. We chose to minimize the Frobenius norm here.

The resulting problem is a positive-definite convex quadratic program and can be rewritten as

$$\min_{s.t., \mathbf{x}_i \geq 0, \mathbf{A}\mathbf{x} = \mathbf{1}_{2n}} \frac{1}{2} \mathbf{x}^\top \mathbf{x} - \mathbf{q}^\top \mathbf{x}, \quad \mathbf{A} = \begin{pmatrix} \mathbf{1}_n^\top & \mathbf{0}^\top & \dots & \mathbf{0}^\top \\ \mathbf{0}_n^\top & \mathbf{1}^\top & \dots & \mathbf{0}^\top \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0}^\top & \mathbf{0}^\top & \dots & \mathbf{1}_n^\top \\ \mathbf{I}_n & \mathbf{I}_n & \dots & \mathbf{I}_n \end{pmatrix} \quad (5)$$

where $\mathbf{x} = \text{vec}(\mathbf{X}^\top)$, $\mathbf{q} = \text{vec}(\mathbf{M}^\top)$, $\mathbf{A} \in \mathbb{R}^{2n \times n^2}$.

QontOT

QontOT is a parameterized (variational) quantum circuit that was conceived very recently for conditional prediction of optimal transport plans.

For any unitary matrix $\mathbf{U} : \mathbf{U} \odot \bar{\mathbf{U}} \in \Omega_n$. Given the circuit parameters θ (typically in the hundreds) and $p \in \mathbb{R}$, QontOT obtains a DSM via $\mathbf{U}(p; \theta) \odot \bar{\mathbf{U}}(p; \theta)$.

Furthermore, QontOT requires the DSM dimension n to be a power of 2. While this may be prohibitive within a Transformer (because sequence length T may differ), it can be mitigated by padding. Padding to powers of two is a common technique to maximize hardware efficiency.

QR Decomposition

As highlighted above, any unitary \mathbf{U} can provide a DSM by taking $\mathbf{U} \odot \bar{\mathbf{U}}$. While there are many ways to obtain a basis, we choose a QR decomposition $\mathbf{M} = \mathbf{U}\mathbf{R}$, in which case \mathbf{R} is upper triangular.

QR has the advantage that, when implemented with Gram-Schmidt, it is differentiable if \mathbf{M} is full-rank. However note that for long-context applications the rank is rarely full because the query and key matrix have $d = \frac{d_{embed}}{n_{heads}}$ rows and typically $\mathbf{M} \in \mathbb{R}^{T \times T}$ has rank $\min\{d, T\}$, implying that \mathbf{M} only has full rank when $d \geq T$.

Soundness and Completeness

- ▶ **Soundness:** Does the operator always produce a DSM? Given that $U \odot \bar{U} \in \Omega_n$, QontOT always yields a DSM. Similarly for the QR decomposition. Instead, Sinkhorn's algorithm (SA) may fail to produce a DSM if the input matrix is not positive.
- ▶ **Completeness:** Can the operator produce all possible DSMs?

Empirical analysis

For a $n \times n$ matrix and a discretization step $d \in \mathcal{N}_+$, we sample each column from a discretized n -dimensional hypercube with d^n points, yielding d^{n^2} unique matrices. For $n = 4$ and $d = 3$ we obtain $3^{16} \approx 43M$ matrices and computed the DSM for each input, before rounding to third decimal place.

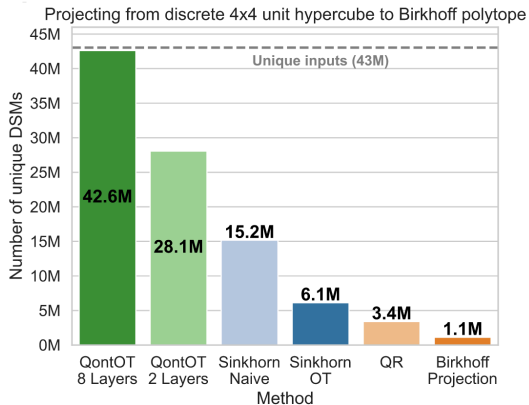


Figure 2. Number of unique DSMs obtained with each operator after exhaustively iterating over a discretized unit hypercube. With only 8 layers, QontOT produces a unique DSM for every possible input, unlike all other methods.

Empirical analysis

A powerful operator needs to possess two further characteristics:

- The information has to be preserved. Obtaining unique DSMs is useless if they destroy information from the input matrix. To assess this, we measured the Frobenius norm of the residuals between input and output matrix.
- The low entropy has to be avoided because it causes vanishing gradients and destabilizes Transformer training.

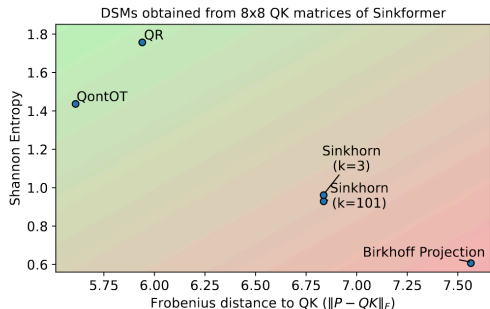


Figure 3. Shannon entropy of different flavors of doubly stochastic attention against a measure of distance preservation – the Frobenius norm of the difference between unnormalized attention QK and their obtained DSM P . QontOT uses 2 circuit layers.

Outline

Introduction

Methods

Doubly-Stochastic Operators

Expressivity of Doubly-Stochastic Operators

Experiments

Conclusion

Experimental Setup

Running the circuit on quantum hardware requires $\Omega(n^2/\epsilon^2)$ shots to obtain satisfactory sampling error. Assuming a precision of $\epsilon = 0.01$, this is in the order of $640k$ shots per sample. Since quantum hardware operates on kHz frequency, execution and online optimization on hardware is unfortunately not (yet) feasible. Therefore, we perform exact statevector simulation with *Qiskit* and implement three circuit training strategies:

- ▶ **Differentiable:** Joint optimization of circuit and Transformer parameters through backpropagation, akin to integrating the circuit as a neural network layer.
- ▶ **Mixed:** A mixed strategy where Transformer training is interleaved with 200 steps of gradient-free circuit optimization with *Nevergrad* on a per-epoch basis.
- ▶ **Static:** The circuit is used in pure inference mode with parameters obtained from a 24-qubit DSM prediction experiment on quantum hardware.

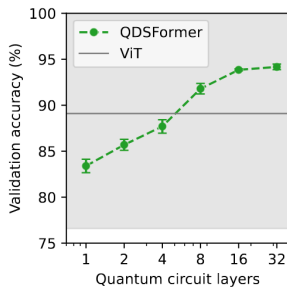
Datasets

We evaluate the different ViTs on MNIST, Fashion MNIST, seven datasets from the MedMNIST benchmark and a compositional task requiring multistep reasoning.

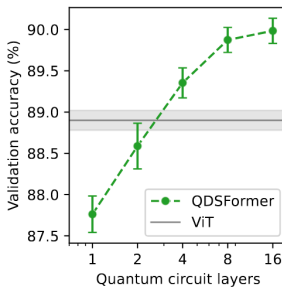
In that task, a 2×2 grid contains two MNIST digits (upper left and lower right) and two FashionMNIST items (upper right and lower left). If the digits have equal value, the label is the upper right fashion item, otherwise it is the bottom left fashion item. Performance typically ramps up quickly to $\sim 50\%$ because the model learns to attend one (and only one) of the FashionMNIST images. Upon continued training with a long saturation phase, a ViT suddenly grasps the relationship of the MNIST digits to the classification task and then climbs rapidly to a 90 – 95% accuracy. The moment of abrupt improvement is called "Eureka moment".



Experiments



(a) 1-ViT-Layer on MNIST.



(b) 2-ViT-Layer on FashionMNIST.

Figure 4. Comparison of ViT and QDSFormer while varying the circuit depth. Mean/std from 5 trainings are shown. Within (a) and (b) all models use the same number of trainable parameters.

Experiments

Table 1. L -layered ViT validation accuracy on FashionMNIST for different attention methods. QontOT uses 16 circuit layers. Mean/std computed from 5 trainings.

L	Softmax	Softmax $_{\sigma^2}$	QR	QontOT	Sinkhorn
1	<u>86.5</u> ± 0.19	75.3 ± 4.61	87.1 ± 0.26	85.6 ± 0.10	84.2 ± 3.64
2	88.9 ± 0.12	84.6 ± 2.11	<u>89.3</u> ± 0.07	90.0 ± 0.15	89.1 ± 0.73
3	<u>89.4</u> ± 0.28	86.3 ± 2.69	<u>89.4</u> ± 0.11	90.3 ± 0.13	<u>89.4</u> ± 0.77
4	<u>89.7</u> ± 0.29	87.1 ± 1.15	89.5 ± 0.07	90.3 ± 0.14	89.1 ± 1.08

Table 2. Identical to [Table 1](#) but for MNIST.

L	Softmax	Softmax $_{\sigma^2}$	QR	QontOT	Sinkhorn
1	89.1 ± 12.5	66.7 ± 22.5	96.6 ± 0.10	93.9 ± 0.11	<u>94.3</u> ± 1.97
2	98.1 ± 0.33	93.0 ± 4.57	<u>98.3</u> ± 0.13	98.4 ± 0.05	98.2 ± 0.27
3	<u>98.6</u> ± 0.11	97.7 ± 0.65	<u>98.6</u> ± 0.13	98.7 ± 0.06	<u>98.6</u> ± 0.12
4	98.8 ± 0.10	97.9 ± 0.71	<u>98.7</u> ± 0.11	98.8 ± 0.07	97.9 ± 1.57

Experiments

Table 3. Test accuracy for different MedMNIST datasets across five attention types in a 2-layer ViT. QontOT uses 16 circuit layers.

MedMNIST dataset	Softmax	Softmax _{σ^2}	QR	QontOT	Sinkhorn
OCT	64.4 ± 1.6	43.6 ± 3.0	62.5 ± 0.9	61.6 ± 0.6	55.1 ± 5.2
Pneumonia	84.2 ± 0.8	84.7 ± 2.0	84.3 ± 0.7	86.1 ± 1.0	83.0 ± 1.5
Tissue	60.0 ± 0.2	49.4 ± 1.2	59.0 ± 0.1	60.6 ± 0.1	56.9 ± 2.0
OrganA	78.8 ± 0.5	73.6 ± 1.7	78.4 ± 0.6	81.2 ± 0.3	77.0 ± 2.5
OrganC	79.8 ± 0.5	71.7 ± 7.3	79.6 ± 0.3	82.7 ± 0.5	79.7 ± 1.0
OrganS	64.4 ± 0.6	59.3 ± 0.9	62.6 ± 0.8	68.1 ± 0.6	63.5 ± 0.9
Breast	79.6 ± 2.0	78.2 ± 2.2	81.3 ± 2.9	80.0 ± 1.1	80.1 ± 0.8
Mean	73.0	65.8	72.5	74.3	70.8

Experiments

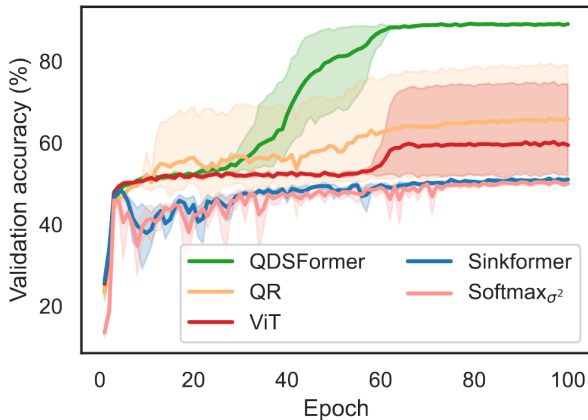


Figure 5. Validation accuracy for each training epoch, highlighting the abrupt learning referred to as Eureka Moment (EM) on the compositional dataset. Confidence bounds from 5 runs.

Outline

Introduction

Methods

- Doubly-Stochastic Operators

- Expressivity of Doubly-Stochastic Operators

Experiments

Conclusion

Conclusion

- ▶ This paper proposed the quantum doubly stochastic Transformer by connecting the centerpiece of a novel variational quantum circuit (devised for optimal transport) with the Transformer, through doubly-stochastic attention improves performance in Transformers.
- ▶ However, this method requires a large amount of data for training and is not suitable for larger models.