

Generative Learning of Continuous Data by Tensor Networks

Alex Meiburg^{1,2,3}, Jing Chen^{1*}, Jacob Miller^{1†}, Raphaëlle Tihon⁴, Guillaume Rabusseau^{4,5},
Alejandro Perdomo-Ortiz⁶

¹ Zapata AI, Boston, USA

² Perimeter Institute for Theoretical Physics, Waterloo, Canada

³ Institute for Quantum Computing, University of Waterloo, Canada

⁴ Mila and DIRO, Université de Montréal, Montréal, Canada

⁵ CIFAR AI Chair, Canada

⁶ Zapata AI, Toronto, Canada

* jing.chen@zapata.ai, † jacob.miller@zapata.ai

Abstract

Beyond their origin in modeling many-body quantum systems, tensor networks have emerged as a promising class of models for solving machine learning problems, notably in unsupervised generative learning. While possessing many desirable features arising from their quantum-inspired nature, tensor network generative models have previously been largely restricted to binary or categorical data, limiting their utility in real-world modeling problems. We overcome this by introducing a new family of tensor network generative models for continuous data, which are capable of learning from distributions containing continuous random variables. We develop our method in the setting of matrix product states, first deriving a universal expressivity theorem proving the ability of this model family to approximate any reasonably smooth probability density function with arbitrary precision. We then benchmark the performance of this model on several synthetic and real-world datasets, finding that the model learns and generalizes well on distributions of continuous and discrete variables. We develop methods for modeling different data domains, and introduce a trainable compression layer which is found to increase model performance given limited memory or computational resources. Overall, our methods give important theoretical and empirical evidence of the efficacy of quantum-inspired methods for the rapidly growing field of generative learning.

Copyright attribution to authors.

This work is a submission to SciPost Physics.

License information to appear upon publication.

Publication information to appear upon publication.

Received Date

Accepted Date

Published Date

1

2 Contents

3	1 Introduction	2
4	2 Background	4
5	2.1 Tensor Networks	4
6	2.2 Discrete-valued Born Machines	5
7	2.3 Related Work	6
8	3 Continuous-valued Born Machine Model	7

9	3.1 Model	7
10	3.2 Sampling	10
11	3.3 Training	10
12	4 Feature Functions	11
13	4.1 Priors from Feature Maps	12
14	5 Universal Approximation with Continuous-valued MPS	13
15	6 Compression Layer	15
16	7 Numerical Results	16
17	7.1 Rotated Hypercube	17
18	7.2 Two Moons	18
19	7.3 Iris Dataset	19
20	7.4 XY Model	19
21	7.5 Compression Test	19
22	8 Conclusions	20
23	9 Acknowledgments	21
24	References	21
25	A Generality of Isometric Feature Map Condition	25
26	B Proof of Marginal Distribution Characterization	26
27	C Proof of Universal Approximation Results	27
28	C.1 Functional Analysis Preliminaries	27
29	C.2 Proof of Theorem 2	28
30	C.3 Proof of Theorem 3	31
31	D Detailed Methods	33
32	D.1 Rotated Cube	33
33	D.2 Two Moons	34
34	D.3 Iris	34
35	D.4 XY Model	36
36	D.5 Compressable Data	36
37	E Dynamic Basis Training	36
38	<hr/>	
39		

40 1 Introduction

41 Although originally developed for the needs of quantum many-body physics [1–4], tensor net-
 42 works (TNs) have rapidly expanded to a host of other areas, where their ability to model correla-
 43 tions and reveal hidden structures within spaces of exponentially large dimension have made them
 44 an invaluable tool in such domains as quantum computing [5–7], applied mathematics [8–10], and
 45 machine learning (ML) [11–14]. In this last setting, TN models are taken as parameterized models

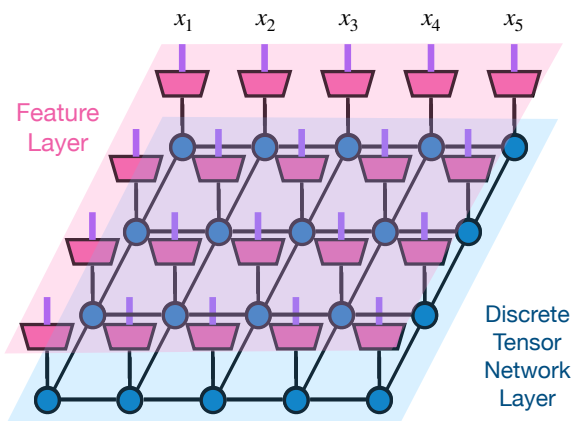


Figure 1: Continuous-valued tensor network. The feature layer (magenta) is a tensor product of feature map operators ζ defined on each site, with the thicker purple edges denoting indices associated to continuous values. The feature layer is connected to the site indices of a discrete-valued tensor network (blue). The specific network above defines a function $\Phi(\mathbf{x})$, where $\mathbf{x} = (x_1, x_2, \dots, x_{20})$.

46 for approximating functions which solve real-world tasks, where TN optimization methods such as
 47 the density matrix renormalization group (DMRG) [15] can be repurposed to formulate quantum-
 48 inspired approaches to learning the structure of naturally occurring data. This approach has opened
 49 up a string of theoretical and empirical successes, from theoretical results in previously intractable
 50 problems in learning theory [16, 17], to practical high-performance compression methods for large
 51 ML models [18, 19], to empirical successes across such tasks as image classification [11, 20], miss-
 52 ing data imputation [21, 22], and unsupervised probabilistic modeling [12, 23–26].

53 Generative modeling, where a parameterized model is trained to draw from an unknown prob-
 54 ability distribution based on a dataset of previous samples, represents a particularly promising area
 55 for the use of TN models in ML [12, 23–26]. Beyond the significant intrinsic value of generative
 56 modeling for everyday applications (as evidenced by the recent explosion of popular interest in
 57 generative AI), TN models using the Born machine (BM) formalism [12, 27] present several dis-
 58 tinctive benefits within this domain that remain elusive with other classical methods. Some of
 59 these benefits arise from the use of tools from entanglement theory, with examples including pow-
 60 erful architecture design methods [28, 29] and rigorous expressivity relationships [30, 31], while
 61 other benefits, such as perfect sampling [32] (and variable-length generalizations [24]), follow
 62 from the distinctive mathematical composition of TN models.

63 TN generative models are not without their limitations however, the most significant of which
 64 is their near-exclusive application to distributions of *discrete* random variables. This restriction
 65 can be best understood within the BM formalism, which is often thought by physicists as de-
 66 scribing many-body wavefunctions. In this context, a TN model can be viewed as a “synthetic”
 67 many-body wavefunction, with the number of possible values of each random variable setting the
 68 dimension of the associated local spin. Because the BM formalism is primarily used in many-body
 69 quantum physics, where a TN describes a discrete “orbital” or site space, it seems natural from
 70 that standpoint to use them exclusively for discrete variables. Continuous random variables would
 71 necessitate infinite-dimensional local spins, which have received less attention in the many-body
 72 TN community. This restriction to discrete random variables severely limits the applicability of
 73 TN models in real-world generative modeling, where the majority of problems involve data with
 74 continuous features.

75 In this paper, we present a framework for employing TNs in generative modeling problems
 76 involving continuous variables. We make use of vector-valued *feature maps* as a means of map-

ping from the infinite-dimensional space associated with each continuous variable to a finite-dimensional feature space associated with each core of a TN. Using the matrix product state (MPS) ansatz for concreteness, we show how restricting these feature maps to be isometries permits an extension of the standard MPS canonical form to the continuous-valued setting, which in turn allows the use of DMRG update and sweep schemes and perfect sampling algorithms within this new setting. We find that the choice of feature maps and the dimension of the associated feature spaces have a large impact on the behavior of the associated generative model, and develop analytical methods to identify suitable choices of these parameters for different input datasets. Despite the simplicity of this continuous-valued MPS model, which contains the same variational parameters as a standard MPS, we nonetheless prove a *universal approximation theorem* demonstrating that it can approximate any sufficiently smooth probability density function to arbitrary precision, given sufficiently large bond dimensions and feature dimensions. On top of this basic model, we develop a novel *compression layer* that permits the feature map itself to be learned from data, which we show gives significant improvements in the performance of the model for a given number of variational parameters. These methods are empirically evaluated on various synthetic and real datasets containing combinations of discrete and continuous variables, where they are found to reliably capture the features of the dataset in each case.

2 Background

Before discussing the continuous case, we first give a brief overview of TNs and BM models in the setting of probability distributions over discrete variables. For a more detailed introduction to TNs and BMs, we refer the interested reader to [4, 12].

2.1 Tensor Networks

Tensor networks (TNs) are a general mathematical formalism for representing large multidimensional arrays as the contraction of smaller tensor *cores*. The collection of the model's tensor cores comprise the parameters of the model, whose elements can be varied to achieve high performance in optimization or learning tasks. Because much of the historical development of TNs took place in the setting of condensed matter physics, the multidimensional arrays in question are often thought of by physicists as describing many-body wavefunctions, with the indices of these arrays corresponding to individual spins (e.g. bosons, fermions, or qubits). In machine learning settings though, the tensors in question will describe multivariate functions to be learned from data, in which case the indices will correspond to individual variables, such as those of a multivariate probability distribution. Other use cases of TNs for ML can be found in supervised learning [11, 33], tensor regression [34, 35], and combinatorial optimization [36, 37].

Typical TN models, including all those considered here, use N separate tensor cores $\{A^{(i)}\}_{i=1}^N$ to encode an N th order tensor $\psi \in \mathbb{K}^{d_1 \times d_2 \times \dots \times d_N}$, where \mathbb{K} refers to either the real or complex numbers. Each core of the TN contains one *site index* of dimension d_i , with the other *bond indices* of $A^{(i)}$ being associated to edges of a graph describing the *network* of tensor contractions connecting the cores of the TN. Different graphs define different TN models, and the graphical structure associated to a TN constrains the correlations achievable between different regions of the model via *area laws* [3, 38]. For a given TN structure, the dimensions of the hidden indices are hyperparameters known as the *bond dimensions* of the model, which determine a trade-off between the expressivity of the model (i.e. the range of tensors which can be represented), and the computational cost of its operation.

Our work utilizes the matrix product state (MPS) model, which is defined by a 1D line graph connecting adjacent sites. The bond dimensions of an MPS can in principle vary for each bond

122 connecting adjacent sites, but here will be assumed to be some constant value $\chi \geq 1$. In this case,
 123 the tensor ψ encoded by the N cores of an MPS is defined by the relation

$$\psi_{\mathbf{s}} = A^{(1)}[i_1]A^{(2)}[i_2]\cdots A^{(N)}[i_N] \quad (1)$$

124 where $\mathbf{s} = (i_1, i_2, \dots, i_N)$ is the joint value of all site indices of the tensor, and the RHS of Eq. 1
 125 describes the multiplication of a χ -dimensional row vector $A^{(1)}[i_1]$, $N-2$ different $\chi \times \chi$ matrices
 126 $A^{(2)}[i_2]\cdots A^{(N-1)}[i_{N-1}]$, and a χ -dimensional column vector $A^{(N)}[i_N]$. This MPS can therefore
 127 be completely described by 2 matrices of dimension $\chi \times d_1$ and $\chi \times d_N$, along with $N-2$ third-
 128 order tensors of shape $\{\chi \times \chi \times d_i\}_{i=2}^{N-1}$.

129 There are typically two different optimization and update strategies. One approach involves
 130 updating all tensors incrementally using gradient-based algorithms, as is commonly employed to
 131 train neural networks in machine learning settings. The other approach targets one site or two
 132 adjacent sites, optimizing them fully before moving to the next target. This method involves inter-
 133 actively sweeping and targeting tensors from left to right and then right to left, inspired by DMRG
 134 sweeps used in calculating ground states. At each step for a given target, we use gradient descent
 135 methods to update the bond tensors until convergence, thereby avoiding the frequent recalculation
 136 of environment tensor contractions.

137 Similar to DMRG schemes, we can target one or two adjacent sites for optimization. In the
 138 one-site update approach, the bond dimension is fixed and predetermined. For the two-site update,
 139 the two tensors are contracted to form a bond tensor, which is then optimized via gradient-based
 140 methods until convergence. The bond tensor can then be factorized back into two adjacent tensors,
 141 with the dimension of the newly factorized bond dynamically adjusted based on the singular value
 142 spectra occurring in the decomposition. We will refer to this approach as the DMRG two-site
 143 scheme in the following discussion. However, unlike traditional DMRG methods for ground state
 144 problems, this approach will not involve solving an eigenvalue problem.

145 2.2 Discrete-valued Born Machines

146 While TNs such as MPS allow the description of arbitrary tensors ψ , in the context of probabilistic
 147 modeling we would like our models to describe probability distributions, whose entries are non-
 148 negative and sum to 1. The Born machine (BM) model represents a natural way of doing so,
 149 which also permits the use of concepts from quantum information within the setting of classical
 150 probabilistic modeling. A BM is parameterized by a TN describing a “synthetic wavefunction”
 151 ψ over the values $\mathbf{s} = (i_1, i_2, \dots, i_N)$ of the N discrete random variables, where the elements $\psi_{\mathbf{s}}$
 152 of ψ can either be real or complex. In either case, the probability distribution defined by the TN
 153 parameterization of ψ is taken to be that given by the Born rule of quantum mechanics, namely

$$P(\mathbf{s}) = \frac{1}{Z} |\psi_{\mathbf{s}}|^2, \quad (2)$$

154 where the partition function Z is defined by

$$Z = \sum_{\mathbf{s}} |\psi_{\mathbf{s}}|^2. \quad (3)$$

155 Eqs. 2 and 3 guarantee that $P(\mathbf{s}) \geq 0$ for all \mathbf{s} , and that $\sum_{\mathbf{s}} P(\mathbf{s}) = 1$, thus ensuring a valid
 156 probability distribution. Although the naive summation in Eq. 3 is exponential in the number
 157 of variables N , for BMs defined over MPS this can be carried out in time $\mathcal{O}(Nd\chi^3)$, where
 158 $d = \max_i d_i$. Alternately, TN *canonical forms* can be used to constrain the tensor cores of the
 159 MPS to always satisfy $Z = 1$, in which case the evaluation of probabilities in Eq. 2 only has cost
 160 $\mathcal{O}(N\chi^2)$.

161 BMs are often used in the context of *density estimation*, where the goal is to learn a probability
 162 distribution P which approximates a target distribution Q using a finite data set $\mathcal{D} = \{\mathbf{s}^{(j)}\}_{j=1}^T$

163 of T samples from \mathbf{Q} . A conventional approach for optimizing the TN cores of the BM is by
 164 minimizing the Kullback-Liebler (KL) divergence $\mathbf{KL}(\mathbf{Q}, \mathbf{P}) = \sum_{\mathbf{s}} \mathbf{Q}(\mathbf{s}) \log(\mathbf{Q}(\mathbf{s})/\mathbf{P}(\mathbf{s}))$ between
 165 \mathbf{P} and \mathbf{Q} , which is equivalent to minimizing the cross-entropy loss

$$L(\mathbf{Q}, \mathbf{P}) = - \sum_{\mathbf{s}} \mathbf{Q}(\mathbf{s}) \log(\mathbf{P}(\mathbf{s})) \approx \frac{1}{T} \sum_{\mathbf{s} \in \mathcal{D}} -\log(\mathbf{P}(\mathbf{s})). \quad (4)$$

166 Although the first summation in Eq. 4 ranges over all possible values of \mathbf{s} , leading to an exponential
 167 cost with increasing number of variables N , the second summation only depends on the size of the
 168 dataset \mathcal{D} , and can therefore be used to efficiently train the model to minimize Eq. 4. In this form,
 169 we will refer to the finite sum on the right of Eq. 4 as the negative log likelihood (NLL) loss of the
 170 model on the dataset \mathcal{D} . We note that the same functional definitions as above will be used later for
 171 defining the KL divergence and NLL loss for probability density functions of continuous random
 172 variables \mathbf{x} . While the NLL loss is always non-negative for discrete-valued probabilistic models,
 173 in the continuous-valued case it is possible for this quantity to become negative for a sufficiently
 174 peaked density \mathbf{P} .

175 While BMs are trained in a similar manner to other classical probabilistic models, they possess
 176 several distinct advantages. Besides being efficient models for density estimation, BMs are also
 177 generative models whose underlying TN factorization allows efficient sampling from the exact
 178 distribution \mathbf{P} , without the need for Monte Carlo or other approximate sampling methods. The
 179 existence of such *perfect sampling* [32] methods is closely linked to the efficient TN computation
 180 of the partition function \mathbf{Z} in Eq. 3, and extends to any BM whose underlying TN has an acyclic
 181 graph structure [27]. Furthermore, the interpretation of samples from \mathbf{P} as outcomes of a projective
 182 measurement on the underlying wavefunction ψ permits the application of tools from quantum
 183 information within the setting of classical probabilistic modeling, something which has been used
 184 as a powerful theoretical tool for characterizing the expressivity of different model families [30,
 185 31], as well as answering model design questions based solely on the underlying dataset \mathcal{D} [29].

186 2.3 Related Work

187 The notion of feature functions has previously been used in tensor network models, primarily in
 188 the context of classification tasks [11, 13, 39], although also with some applications in function
 189 regression [40] and generative modeling [41]. As we discuss later, our interpretation of the fea-
 190 ture functions as isometric maps permits straightforward conditional generation and training of
 191 the continuous-valued BM model. Refs. [42, 43] look at the question of MPS approximations of
 192 continuous functions, but where increasingly fine discretizations of the function are approximated
 193 using discrete-valued MPS. Ref. [44] shows how to combine a similar style of discretization with
 194 certain analytically tractable feature functions. Ref. [45] presents a universal approximation result
 195 for functional tensor trains that we use as an important building block in the development of the
 196 universal approximation theorems of Sec. 5. Refs. [46, 47] present similar continuous generaliza-
 197 tions of tensor train (TT) models, but whose optimization is handled by very different algorithms.

198 The work of [48] studies density modeling of continuous data (phrased in terms of TTs rather
 199 than MPS), with the “squared tensor train density estimation” variation of their model having
 200 many similarities to ours. The distinct origin and focus of the TT and MPS communities lead
 201 to several important differences between the model of [48] and the one introduced here. While
 202 the model of [48] is similarly capable of perfect conditional and unconditional sampling, this
 203 requires the computation of explicit marginals that are trivial in our case owing to the use of an
 204 MPS canonical form. This use of canonical forms also allows us to optimize the model using
 205 a DMRG update and sweep approach, in contrast to updating all tensors simultaneously by the
 206 gradient-based optimization used in [48]. Our compression layer architecture is novel, as are the
 207 universal approximation results Theorems 2 and 3 proving that continuous-valued MPS permit the
 208 approximation of any (sufficiently smooth) wave function or probability density function.

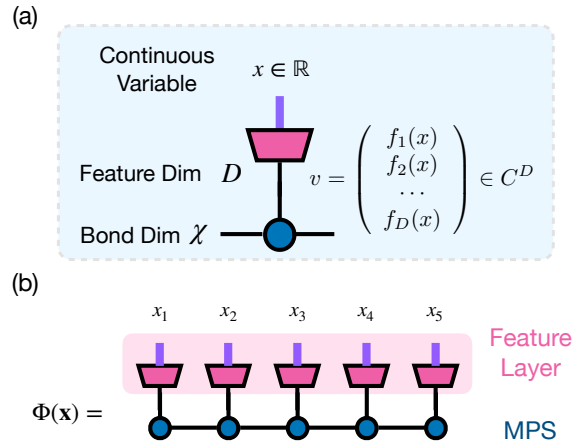


Figure 2: Continuous-valued MPS. (a) For the feature layer, the input at each site $\mathbf{x} \in \mathbb{R}$ is a continuous variable, after mapping, it outputs a discrete vector of feature dimension D , which is directly connected to the tensor network layer (blue). χ and D are hyper parameters controlling the dimensions of different bonds. (b) Graphical depiction of the continuous-valued function Φ defined in Eq. 6.

209 Our notion of “continuous-valued MPS” should not be confused with the continuous matrix
 210 product states introduced in [49]. These models utilize a continuous *spatial* dimension, and can be
 211 thought of as the limit of an infinite number of site indices, but with the site dimensions remaining
 212 constant and discrete. This type of model has applications in quantum field theory, and is not of
 213 interest in this context. The setting we consider here uses a fixed number of indices, but with each
 214 index varying over a continuous domain.

215 3 Continuous-valued Born Machine Model

216 Despite the desirable properties of standard BMs, one obvious limitation is their restriction to
 217 modeling discrete-valued probability distributions. While this is rarely an issue in the setting of
 218 many-body physics, within classical machine learning such a restriction is extremely limiting, as
 219 most datasets used in unsupervised learning contain continuous features described by a probability
 220 density function (PDF).

221 To remedy this limitation, we show here how discrete-valued Born machines can be naturally
 222 generalized to the setting of probability distributions over any combination of discrete and contin-
 223 uous variables, as depicted in Fig. 2. This generalization is made possible by the use of *feature*
 224 *maps* which convert points in the continuous domain into finite-dimensional vectors which can
 225 be contracted with the underlying discrete-valued TN. We show in detail how this generaliza-
 226 tion preserves all of the convenient properties and standard algorithms for BMs, including perfect
 227 sampling, density evaluation at specific points in the domain, and efficient computations of the par-
 228 tition functions and marginals. For convenience of presentation, in the following we assume the
 229 use of identical feature maps for all N sites of the MPS, which are assumed to possess a common
 230 feature dimension D , but the generalization to site-dependent feature functions is straightforward.

231 3.1 Model

232 At a high level, the continuous-valued BM introduced here uses a feature map $\zeta : \mathcal{I} \rightarrow \mathbb{K}^D$
 233 to convert variables \mathbf{x} from a continuous domain \mathcal{I} to real or complex D -dimensional vectors
 234 $\mathbf{v} = \zeta(\mathbf{x}) \in \mathbb{K}^D$. Once such a map has been defined for each site of a discrete-valued MPS, it can

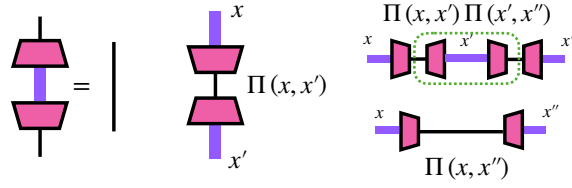


Figure 3: Graphical formulation of the isometric condition on the feature map ζ , which is equivalent to the orthonormality requirement on feature functions expressed in Eq. 6.

235 be used to convert the MPS into a function of N continuous variables.

236 A more concrete manner of representing our mapping ζ comes from picking a basis of D
 237 feature functions $\mathcal{F} = \{f_i\}_{i=1}^D$, with each $f_i : \mathcal{I} \rightarrow \mathbb{K}$ equal to the projection of ζ onto one of the
 238 D vectors $\{\mathbf{e}_i\}_{i=1}^D$ forming an orthonormal basis for \mathbb{K}^D , so that $f_i(\mathbf{x}) = \langle \mathbf{e}_i, \zeta(\mathbf{x}) \rangle$. We require ζ
 239 to be an isometry, meaning that the feature functions satisfy the relations

$$\begin{aligned} (\zeta \zeta^\dagger)_{i,j} &= \int_{-\infty}^{\infty} f_i^*(x) f_j(x) dx = \delta_{ij}, \\ (\zeta^\dagger \zeta)_{x,x'} &= \sum_{i=1}^D f_i(x) f_i^*(x') = \Pi(x, x'), \end{aligned} \quad (5)$$

240 where $\Pi(x, x')$ is a kernel function satisfying $\int_{x'} \Pi(x, x') \Pi(x', x'') dx' = \Pi(x, x'')$ (see Fig. 3).
 241 This isometry requirement is invaluable for extending the convenient properties of discrete-valued
 242 MPS and BMs to the continuous-valued setting, and can be made without loss of generality, as
 243 any feature map can be converted into an isometric form (see Appendix A for details).

244 Given a mapping ζ satisfying the above conditions, any tensor ψ containing N discrete indices
 245 $\mathbf{s} = (i_1, i_2, \dots, i_N)$ can be promoted into a continuous-valued function Φ of N continuous vari-
 246 ables $\mathbf{x} = (x_1, x_2, \dots, x_N)$ by contracting each site index with the corresponding vector $\zeta(\mathbf{x}_k)$, as
 247 described by

$$\Phi(\mathbf{x}) = \sum_{i_1, i_2, \dots, i_N} \left(\prod_{k=1}^N f_{i_k}(x_k) \right) \psi_{i_1, i_2, \dots, i_N}. \quad (6)$$

248 A graphical representation of Eq. 6 is shown in Fig. 2, where the tensor ψ is taken to be given by
 249 a discrete-valued MPS.

250 Just as with discrete-valued BMs in Eq. 2, the continuous-valued BM PDF P is given by the
 251 elementwise norm squared of the underlying function $\Phi : \Omega \rightarrow \mathbb{K}$,

$$P(\mathbf{x}) = |\Phi(\mathbf{x})|^2. \quad (7)$$

252 $P(\mathbf{x})$ is clearly non-negative everywhere in its domain of definition $\Omega = \mathcal{I}^N$ and, owing to the
 253 isometry conditions of Eq. 6, is guaranteed to satisfy the normalization condition $\int_{\mathbf{x} \in \Omega} P(\mathbf{x}) d\mathbf{x} = 1$
 254 whenever the underlying discrete-valued MPS ψ satisfies the condition $\sum_{i_1, \dots, i_N} |\psi_{i_1, \dots, i_N}|^2 = 1$.
 255 For the case of an unnormalized MPS ψ , the normalization factor required to ensure the proper
 256 normalization of Φ is precisely the squared norm of ψ , which can be efficiently computed using
 257 standard MPS methods. In contrast to the discrete-valued case however, it is possible for the PDF
 258 P to take values $P(\mathbf{x}) > 1$ at some points $\mathbf{x} \in \Omega$.

259 The standard canonical form for discrete-valued MPS can be straightforwardly generalized
 260 (with the help of the isometric constraints of Eq. 6) to produce a notion of canonical form for
 261 continuous-valued MPS, as shown in Fig. 4. Just as with the usual MPS canonical form, this
 262 ensures the proper normalization of the BM distribution P throughout training, and simplifies the
 263 computation of gradients and other quantities which typically require $\mathcal{O}(\chi^3)$ time to compute.

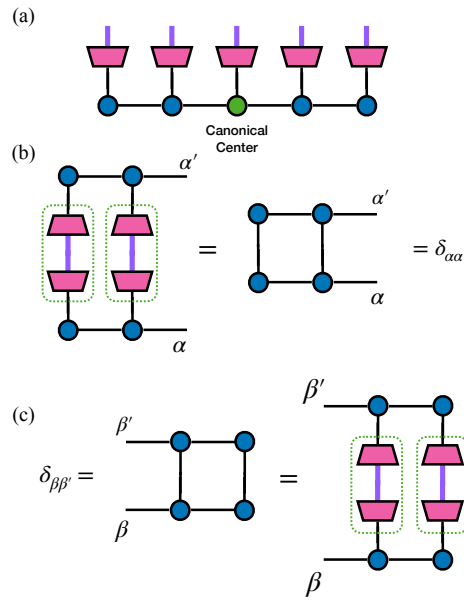


Figure 4: Continuous-valued MPS canonical form. (a) The underlying discrete-valued MPS is required to be in canonical form with an orthogonality center (green dot tensor). When the feature maps additionally satisfy the orthonormality relations of Eq. 6, then the continuous-valued MPS is said to be in continuous-valued MPS canonical form. (b-c) Graphical proof that the left (right) tensors constitute isometries from the left (right) bond spaces to the space of square-integrable functions acting on the left (right) set of continuous variables.

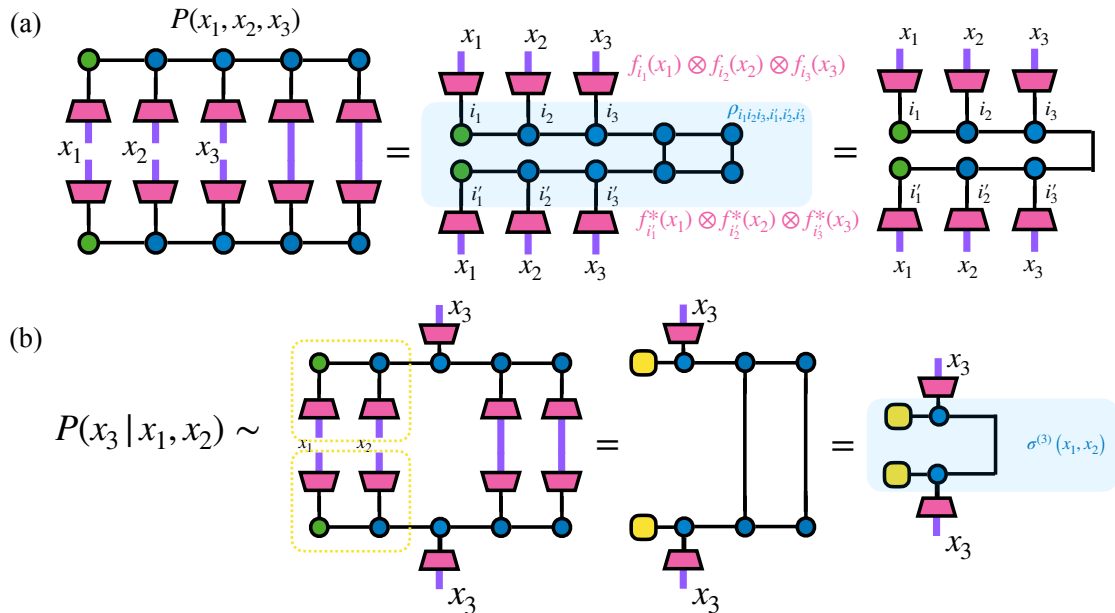


Figure 5: Tensor network diagrams depicting how calculating the probabilities of a continuous-valued MPS BM can be considerably simplified. The MPS is taken to be in canonical form, with the orthonormal center (green dot tensor) on the leftmost site. (a) The marginal distribution $P(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3)$ is given by integrating out the continuous variables $\mathbf{x}_4, \mathbf{x}_5$, which is trivial when the MPS is in continuous-valued canonical form. (b) The conditional probability $P(x_3 | x_1, x_2)$ used in the sampling process, which is facilitated by the computation of a $D \times D$ conditional density matrix $\sigma^{(3)}(x_1, x_2)$.

3.2 Sampling

Continuous-valued MPS BMs share the same perfect sampling capabilities as their discrete-valued counterparts. Sampling proceeds site by site, with the continuous random variable at each site i conditioned on those produced at previous sites $1, 2, \dots, i-1$ via contraction of a sample-dependent vector on the bond dimensions adjacent to site i .

For any site k , the conditional PDF of the random variable \mathbf{x}_k satisfies

$$P(\mathbf{x}_k | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{k-1}) = \frac{P(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{k-1}, \mathbf{x}_k)}{P(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{k-1})}, \quad (8)$$

where the marginal PDFs are defined for any k as

$$P(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k) = \sum_{\substack{i_1, \dots, i_k \\ i'_1, \dots, i'_k}} \left(\prod_{\ell=1}^k f_{i_\ell}^*(\mathbf{x}_\ell) \right) \rho_{i_1, \dots, i_k, i'_1, \dots, i'_k} \left(\prod_{\ell=1}^k f_{i'_\ell}(\mathbf{x}_\ell) \right). \quad (9)$$

In the above, $\rho_{i_1, \dots, i_k, i'_1, \dots, i'_k}$ represents the discrete reduced density matrix resulting from integrating over all remaining variables $\mathbf{x}_{k+1}, \dots, \mathbf{x}_N$. Although the summation involved in Eq. 9, as well as the integrations needed to compute the reduced density matrix, are prohibitively expensive to implement directly, Fig. 5 shows how the tensor network representation of P can be used to remedy this situation.

When the underlying MPS is in canonical form, tracing out the rightmost variables $\mathbf{x}_{k+1}, \dots, \mathbf{x}_N$ can be performed efficiently, and computing value of the conditional probability distribution $P(\mathbf{x}_k | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{k-1})$ can be accomplished with complexity $\mathcal{O}(D^2 \chi^2)$. This process is facilitated by a $D \times D$ conditional density matrix $\sigma^{(k)}(\mathbf{x}_1, \dots, \mathbf{x}_{k-1})$ associated to site k , shown in Fig. 5b. The conditional distribution in question is then given by

$$P(\mathbf{x}_k | \mathbf{x}_1, \dots, \mathbf{x}_{k-1}) = Z_k^{-1} \sum_{i_k, i'_k=1}^D f_{i_k}^*(\mathbf{x}) \sigma_{i_k i'_k}^{(k)}(\mathbf{x}_1, \dots, \mathbf{x}_{k-1}) f_{i'_k}(\mathbf{x}), \quad (10)$$

where the normalization constant Z_k is chosen such that $\int_{\mathbf{x}_k \in \mathcal{I}} P(\mathbf{x}_k | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{k-1}) d\mathbf{x}_k = 1$. This in turn can be computed via numerical or closed-form integration over \mathbf{x}_k , to obtain a cumulative distribution function $F(\mathbf{x}_k) = \int_{\mathbf{x}' \leq \mathbf{x}_k} P(\mathbf{x}') d\mathbf{x}'$. This permits a random sample to be produced using inverse transform sampling, by sampling a uniformly random $\mathbf{z} \sim [0, 1]$ and then applying the inverse of the cumulative distribution F to yield the random sample $\mathbf{x}_k = F^{-1}(\mathbf{z})$. Continuing this process for $k = 1, 2, \dots, N$ yields an exact sample from the BM PDF $P(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$, with $\mathcal{O}(N(\chi^2 D + \chi D^2))$ complexity.

3.3 Training

In the simplest formulation of a continuous-valued MPS, the feature functions $\mathcal{F} = \{f_i\}_{i=1}^D$ are chosen in advance and unchanged throughout training. Only the core tensors of the discrete-valued MPS representation of ψ are taken as tunable parameters, and are trained to minimize the model's NLL on a dataset of unlabeled samples.

Given a dataset with continuous data, each datum can be mapped to a tensor product of vectors associated with the corresponding features at each site. For a dataset \mathcal{D} with N continuous features,

295 the j 'th sample $\mathbf{x}^{(j)} = (x_1^{(j)}, x_2^{(j)}, \dots, x_N^{(j)})$ is mapped into

$$\zeta(x_1^{(j)}) \otimes \zeta(x_2^{(j)}) \otimes \dots \otimes \zeta(x_N^{(j)}) = \bigotimes_{k=1}^N \begin{pmatrix} f_1(x_k^{(j)}) \\ \vdots \\ f_D(x_k^{(j)}) \end{pmatrix}, \quad (11)$$

296 where $\zeta(x_k^{(j)})$ is the vector representation of the k 'th feature of the j 'th sample of \mathcal{D} .

297 Computing the NLL requires a summation over all data samples in \mathcal{D} , in each of which the
 298 site indices of ψ are contracted with the feature vectors given in Eq. 11. The MPS can then be
 299 trained to learn this dataset by any conventional means, such as gradient descent on the NLL of the
 300 distribution, or an adapted version of DMRG [12]. In this latter method, the cores for a pair of sites
 301 $(i, i+1)$ are trained by first contracting the tensor ψ with the feature vectors at sites $1, 2, \dots, i-1$
 302 and $i+2, i+3, \dots, n$, then optimizing the remaining bond tensor to minimize the NLL according
 303 to the procedure described in [12].

304 4 Feature Functions

305 In a setting with discrete data, the possible values of the dataset's categorical features determine
 306 the sizes of the site indices of the TN, so that a feature taking d possible values is always associated
 307 with a site dimension of d . In the continuous-valued setting however, the feature functions and
 308 feature dimension D represent new hyperparameters with a significant impact on the inductive
 309 bias and expressiveness of the model. The following are all feature maps we assess numerically
 310 in Sec. 7, which are natural choices for different types of continuous domains. We describe the
 311 component functions of each map, along with their behavior under isometrization (i.e. imposing
 312 the isometry conditions of Eq. 6).

313 **Fourier** The complex exponentials $e^{i2\pi kx}$ for $k = 0, 1, \dots$ restricted to the compact interval
 314 $[0, 1]$, which already satisfy Eq. 6.

315 **Legendre** Polynomials of degree $k = 0, 1, \dots$ restricted to the compact interval $[-1, 1]$. Isometriza-
 316 tion leads these to be proportional to the Legendre polynomials.

317 **Laguerre** Polynomials of degree $k = 0, 1, \dots$ multiplied by the exponential $e^{-x/2}$, and defined
 318 on the half interval $\{x \in \mathbb{R} | x \geq 0\}$. Isometrization leads these to be proportional to the
 319 Laguerre polynomials multiplied by $e^{-x/2}$.

320 **Hermite** Polynomials of degree $k = 0, 1, \dots$ multiplied by the Gaussian $e^{-x^2/2}$, and defined on
 321 all of \mathbb{R} . Isometrization leads these to be proportional to the Hermite polynomials multiplied
 322 by $e^{-x^2/2}$.

323 Beyond these particular cases, the framework we use permits many other possible feature
 324 maps, including the discretization of continuous variables into categorical ones by binning. Con-
 325 sider the D feature functions f_k defined as

$$f_k(x) = \begin{cases} 1, & \lambda_{k-1} \leq x \leq \lambda_k \\ 0, & \text{otherwise.} \end{cases} \quad (12)$$

326 These indicator functions serve as ‘‘one-hot’’ encodings of the categorical variable associated to
 327 the placement of x into one of D separate bins with bin edges $\lambda_0 < \lambda_1 < \dots < \lambda_D$, and satisfy
 328 Eq. 6 up to normalization.

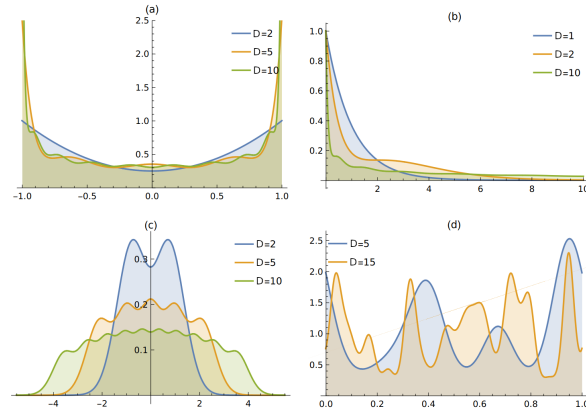


Figure 6: (a-c): The expected univariate distribution P_{init} when a D -dimensional feature map is used. (a) Legendre polynomials, which converge to an arcsin distribution in the limit of $D \rightarrow \infty$. (b) Laguerre functions, for which $\mathbb{E}[\mathbf{x}]$ increases with increasing D . (c) Hermite functions, which progressively broaden with increasing D . (d) The expected distribution for Fourier functions is uniform for all D , but we instead illustrate the univariate distributions associated with specific random MPS at two values of D , for bond dimension $\chi = 5$.

329 One important and obvious criterion when choosing a feature map is the domain of data be-
 330 ing studied. The Lagrange and Fourier feature maps can be used (with appropriate shifting and
 331 scaling) to describe data in any connected compact interval $[a, b]$, with the latter also permitting
 332 features on a periodic domain (for example, angular data). The Laguerre feature map is suitable
 333 for data taking nonnegative real values without an obvious upper limit, while the Hermite feature
 334 map is suitable for data which can range over the whole real line.

335 4.1 Priors from Feature Maps

336 Beyond simply constraining the domain of input data, the choice of feature map for a given
 337 continuous-valued BM sets the inductive bias of the model in a manner which can be precisely
 338 quantified, in the form of univariate marginal distributions at initialization. A common initializa-
 339 tion method for MPS BMs is to choose the elements of each MPS core to be independent identi-
 340 cally distributed (IID) random variables, and in this case the following Theorem characterizes the
 341 single-site marginal distributions over each continuous random variable.

342 **Theorem 1.** Consider a continuous-valued MPS with feature dimension D and an isometric fea-
 343 ture map $\zeta : \mathcal{I} \rightarrow \mathbb{K}^D$ at site i characterized by feature functions $\mathcal{F} = \{f_1, f_2, \dots, f_D\}$. Given
 344 an initialization of all MPS core elements by IID random variables of zero mean and fixed vari-
 345 ance, the expected single-site marginal distribution $P_{\text{init}}(\mathbf{x}_i)$ of the randomly initialized MPS BM
 346 is given by

$$P_{\text{init}}(\mathbf{x}_i) = \frac{1}{D} \|\zeta(\mathbf{x}_i)\|^2 = \frac{1}{D} \sum_{k=1}^D |f_k(\mathbf{x}_i)|^2. \quad (13)$$

347 The proof of Theorem 1 is given in Appendix B, and is based on a simple characterization
 348 of the expected density operator of the underlying discrete-valued relative to the IID initialization
 349 method in question, which then permits a derivation of Eq. 13.

350 To illustrate this result, we consider the expected prior distributions associated with each of the
 351 features maps considered above. The simplest is the Fourier case, where each complex exponential
 352 $f_k(\mathbf{x}_i) = e^{i2\pi k \mathbf{x}_i}$ will have unit norm, and therefore yield an expected uniform distribution over
 353 the interval $\mathcal{I} = [0, 1]$. We note that even in this simple case though, individual random MPS

354 will generally have single-site marginal distributions that differ from this expected distribution
 355 (Fig. 6d), which only characterizes the average with respect to many different initializations.

356 More interesting is the Legendre case (Fig. 6a), where the initial distribution skews towards
 357 the ends of the interval $\mathcal{I} = [-1, 1]$. In the limit of increasing D , the density of the univariate
 358 PDF $P_{\text{init}}^{(D)}$ diverges at the endpoints of the interval, yet the distribution as a whole converges to an
 359 analytically tractable **arcsin** distribution [50], given by

$$\lim_{D \rightarrow \infty} P_{\text{init}}^{(D)}(x) = \frac{1}{\pi \sqrt{1-x^2}}. \quad (14)$$

360 In practice this bias means the Legendre polynomials lead to significantly worse initialization on
 361 most datasets, and we find better performance with other feature maps.

362 The Laguerre and Hermite cases (Fig. 6b and c) are both associated with a broadening of the
 363 mass of the expected univariate distribution with increasing D , at a rate of $\mathcal{O}(\sqrt{D})$.¹ In this case,
 364 it is sensible to rescale the inputs to these feature maps as the feature dimension is increased,
 365 i.e. using the new feature functions $g_k(x) = f_k(\sqrt{D}x)$. These rescaled feature functions likely
 366 converge to exact analytic forms in the $D \rightarrow \infty$ limit, but we leave this characterization as an
 367 open question.

368 From a practical standpoint, Theorem 1 represents a useful tool for choosing feature maps
 369 based on the marginal distributions associated with each feature of a dataset. Employing a feature
 370 map whose expected prior distribution closely resembles the empirical marginal distribution for
 371 that feature leads to improved performance in training, in that regions of the feature space which
 372 occur more often in the dataset are assigned higher probability at initialization. This could be
 373 compared to importance sampling in Monte Carlo methods, which leaves the same asymptotic
 374 distribution in the high-capacity limit, but accelerates the rate of convergence.

375 5 Universal Approximation with Continuous-valued MPS

376 It is well-known that discrete-valued MPS with sufficiently large bond dimensions can exactly
 377 represent any space of N th order tensors using the truncation-free version of the iterated singular
 378 value decomposition (SVD) protocol of [5, 9]. By extension, any discrete-valued probability
 379 distribution can be exactly represented by an MPS BM whose underlying wavefunction is associ-
 380 ated with the square root of the distribution. The corresponding questions for continuous-valued
 381 MPS and square-integrable functions (or PDFs) of N continuous variables are considerably less
 382 straightforward. It is clear that the exact representation result from the discrete case cannot be
 383 applied here, since the continuous-valued functions of interest live in infinite-dimensional Hilbert
 384 spaces, while the functions describable by a continuous-valued MPS with fixed bond dimension
 385 χ and feature dimension D will necessarily occupy a finite-dimensional manifold [53].

386 We overcome this difficulty by proving *universal approximation theorems*, which bound the
 387 worst-case error in encoding a sufficiently smooth wavefunction (resp. PDF) using a continuous-
 388 valued MPS (resp. MPS BM), as a function of the bond dimension χ and feature dimension D .
 389 These results show in particular that by increasing the values of χ and D , any sufficiently smooth
 390 wavefunction or PDF can be approximated to any desired precision using a continuous-valued
 391 MPS.

392 **Theorem 2.** Consider a family of continuous-valued MPS with polynomial feature functions
 393 $\mathcal{F} = \{f_1, f_2, \dots\}$ forming an orthonormal basis for $[0, 1]$, which is defined on the hypercube

¹Hermite distributions have an asymptotic scaling in amplitude as $\left(1 - \frac{x^2}{2D+1}\right)^{-1/2}$ for large D and $|x| \ll \sqrt{2D+1}$, with an exponentially small weight at $|x| \gg \sqrt{2D+1}$ [51]. A similar scaling holds for Laguerre distributions [52].

394 $\Omega = [0, 1]^N \subseteq \mathbb{R}^N$. Let $k \geq N$ and let $\Phi : \Omega \rightarrow \mathbb{C}$ be any square-integrable function with unit
 395 norm ($\langle \Phi, \Phi \rangle = 1$), whose partial derivatives of order $1, 2, \dots, k$ all exist and are bounded. Then
 396 for every positive $\chi, D \in \mathbb{N}$ there exists a continuous-valued MPS of bond dimension χ and feature
 397 dimension D with unit norm, whose associated function $\Phi_{\text{MPS}}^{(\chi, D)}$ approximates Φ with infidelity

$$1 - |\langle \Phi, \Phi_{\text{MPS}}^{(\chi, D)} \rangle| \leq \gamma_1 \chi^{-k+1} + \gamma_2 D^{-2k}, \quad (15)$$

398 where $\gamma_1, \gamma_2 > 0$ depend on the target function Φ , the assumed degree of smoothness k , and the
 399 feature functions \mathcal{F} .

400 The proof of Theorem 2, along with an overview of the functional analytic concepts used in
 401 the proof and the precise definition of the constants γ_1, γ_2 , are given in Appendix C. The result
 402 makes heavy use of the work of [45], which generalizes the iterated SVD method for computing
 403 discrete MPS representations of tensors to the setting of infinite-dimensional spaces of real-valued
 404 functions.

405 We note that the restriction in Theorem 2, which applies to functions Φ defined on the unit
 406 hypercube $\Omega = [0, 1]^N$ is primarily for ease of presentation, and can be easily relaxed to func-
 407 tions on any product of compact intervals $[a_1, b_1] \times \dots \times [a_N, b_N]$ (i.e. an N -dimensional box).
 408 More generally, although a rigorous proof for the case of functions on non-compact domains (e.g.
 409 $\Phi : \mathbb{R}^N \rightarrow \mathbb{C}$) is not possible with the methods of [45], we give a heuristic argument in Appendix C
 410 for how Theorem 2 can be modified to bound the error involved in approximating functions on
 411 non-compact domains using continuous-valued MPS.

412 The above theorem can be used to prove a similar approximation result for PDFs. In place of
 413 infidelity between wavefunctions, we utilize the Jensen-Shannon (JS) divergence between distri-
 414 butions, which is defined as $\text{JS}(P, Q) = \frac{1}{2} (\text{KL}(P, M) + \text{KL}(Q, M))$ for M the equal-weight mixture
 415 of P and Q taking values $M(\mathbf{x}) = \frac{1}{2} (P(\mathbf{x}) + Q(\mathbf{x}))$. Besides being symmetric in the input PDFs
 416 P and Q , JS divergence takes bounded values (in contrast to KL divergence), and is zero only
 417 when P and Q are identical almost everywhere. The following Theorem therefore guarantees that
 418 any sufficiently smooth PDF can be approximated to arbitrary accuracy using a BM built from a
 419 continuous-valued MPS.

420 **Theorem 3.** Consider a family of continuous-valued MPS with polynomial feature functions
 421 $\mathcal{F} = \{f_1, f_2, \dots\}$ forming an orthonormal basis for $[0, 1]^N$, which is defined on the hypercube
 422 $\Omega = [0, 1]^N \subseteq \mathbb{R}^N$. Let $k \geq N$ and let $P : \Omega \rightarrow \mathbb{R}$ be any Probability Density Function (PDF)
 423 bounded below as $P_{\min} = \min_{\mathbf{x} \in \Omega} P(\mathbf{x}) > 0$, whose partial derivatives of order $1, 2, \dots, k$ all
 424 exist and are bounded. Then for every positive $\chi, D \in \mathbb{N}$ there exists a continuous-valued MPS of
 425 bond dimension χ and feature dimension D with unit norm, whose associated Born machine PDF
 426 $P_{\text{MPS}}^{(\chi, D)}(\mathbf{x}) = |\Phi_{\text{MPS}}^{(\chi, D)}|^2$ approximates P with Jensen-Shannon divergence

$$\text{JS}(P_{\text{MPS}}^{(\chi, D)}, P) \leq \eta_1 \chi^{-\frac{k-1}{2}} + \eta_2 D^{-k}, \quad (16)$$

427 where $\eta_1, \eta_2 > 0$ depend on the target PDF P , the assumed degree of smoothness k , and the
 428 feature functions \mathcal{F} .

429 The proof of Theorem 3 applies Theorem 2 to the approximation of a naive target wavefunction
 430 given by $\Phi_P(\mathbf{x}) = \sqrt{P(\mathbf{x})}$, and then uses standard tools from information theory to translate
 431 bounds in infidelity into bounds in JS divergence. The key technical argument of this proof is
 432 ensuring that the smoothness guarantees assumed of P yield similar smoothness guarantees for
 433 Φ_P , which is complicated by the fact that the derivative of $\sqrt{P(\mathbf{x})}$ becomes infinite in the limit
 434 $P(\mathbf{x}) \rightarrow 0$. To avoid this pathological behavior, we require that $P(\mathbf{x})$ be bounded below by some
 435 P_{\min} , as explained in the proof in Appendix C.

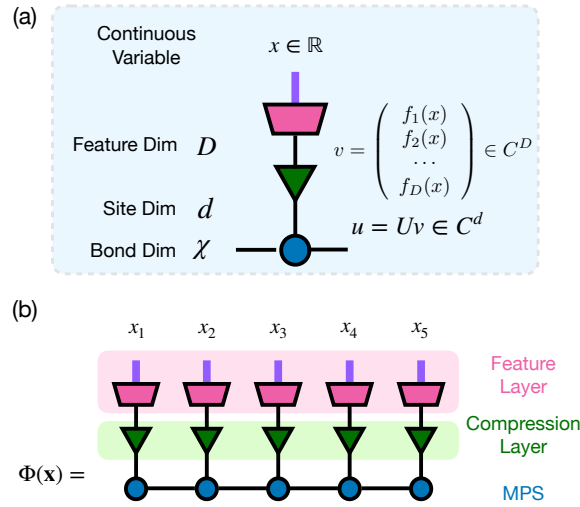


Figure 7: Continuous-valued MPS with compression layer (green). (a) The input x is converted to a vector of dimension D . Then the compression operator (green triangular) is an isometry matrix, which rotate and truncate into a d dimensional site index of the MPS. (b) From top to bottom is the feature mapping layer, compression layer and MPS layer. Both the feature layer and the compression layer is direct product of many local operators.

436 Just as for Theorem 2, the domain in Theorem 3 can be replaced w.l.o.g. with any N -
 437 dimensional box, and by a heuristic argument can be used to bound the error in approximating
 438 PDFs defined on unbounded domains. We note also that different symmetric loss functions can be
 439 used in place of JS divergence in Theorem 3, notably total variation distance.

440 As one final comment on Theorems 2 and 3, the attentive reader might wonder about the
 441 case of smooth target functions, for which the value of k can be made arbitrarily large. While
 442 the bounds in Eqs. 15 and 16 might seem to become arbitrarily small, it is important to note
 443 that the quantities $\gamma_1, \gamma_2, \eta_1, \eta_2$ themselves depend on k , and generally grow very rapidly (e.g.
 444 super-factorially) with increasing k . Consequently, even though smooth PDFs can technically be
 445 approximated with error $\mathcal{O}(\chi^{-\frac{k-1}{2}} + D^{-k})$ for any positive value of k , in practice the large prefactor
 446 in such bounds would make this increasingly favorable scaling only become apparent at values of
 447 χ and D which increase at an astronomical rate.

448 6 Compression Layer

449 The feature dimension D plays a crucial role in determining the expressivity of the continuous-
 450 valued model, as it determines the number of basis functions spanning the space of functions on
 451 the continuous variable. This in turn determines the precision of the continuous variable being
 452 modeled, with a dimension D limiting the precision to roughly $\mathcal{O}(D^{-1})$. While a larger feature
 453 dimension enables the MPS to capture finer details of the distribution, it also comes at the cost of
 454 significantly increased computational complexity. As a concrete example, training using two-site
 455 update scheme leads to a memory cost of $\mathcal{O}(\chi^2 D)$ and a computational cost of $\mathcal{O}(\chi^3 D^3)$, making
 456 it impractical to increase the feature dimension beyond a certain limit.

457 To address this issue, we propose the addition of an intermediate compression layer that con-
 458 nects the D -dimensional feature space to a smaller site space of dimension d in the underlying
 459 discrete-valued MPS. It may be the case in practice that the univariate functions needed to de-
 460 scribe each feature of a target distribution or function are easily describable in a low-dimensional

space, but where each basis function is more complex than a predetermined feature function. Our compression layer takes advantage of this possibility by storing a tunable collection of \mathbf{d} basis functions, which are each taken to be a superposition of \mathbf{D} fixed feature functions, where $\mathbf{d} \ll \mathbf{D}$. This allows us to take advantage of the expressive power of high numbers of feature functions while minimizing computational costs, improving the efficiency and performance of the continuous MPS model. While we have so far taken $\mathbf{D} = \mathbf{d}$, when clarity is needed we will refer to \mathbf{D} as the feature dimension of the model and \mathbf{d} as the site dimension.

Adding a compression layer results in a model that is a simple example of a tree tensor network, as shown in Fig. 7. The compression layer consists of many different $\mathbf{D} \times \mathbf{d}$ matrices $\{U_i\}_{i=1}^N$ satisfying the isometric condition $U_i^\dagger U_i = I_{\mathbf{d}}$, which are tunable parameters of the model. In the case of datasets possessing similar kinds of features (e.g. time series data), it may be advantageous to choose all isometries U_i to be equal.

Jointly training the compression layer with the MPS parameters can be done in either the context of gradient-based optimization, or in an alternating manner in the context of DMRG. The former case can be straightforwardly handled by the use of tools for gradient-based optimization on Stiefel manifolds (i.e. families of isometric matrices), so we describe here the latter optimization process. The isometry U_i at a site i is trained to maximize the NLL loss associated to a training dataset \mathcal{D} , where samples from the dataset are associated with continuous features $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$. For a given sample from \mathcal{D} , the $N - 1$ features at all other sites $\mathbf{x}_i = (\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_N)$ are contracted with all cores of the underlying discrete-valued MPS, giving a \mathbf{d} dimensional vector \mathbf{v}_{-i} at site i , while the remaining feature \mathbf{x}_i is embedded as a \mathbf{D} dimensional vector $\zeta(\mathbf{x}_i)$. Given this information, the goal is to find the isometry which minimizes the negative log-likelihood loss over the training dataset, or equivalently:

$$U_i = \operatorname{argmax} \sum_{\mathbf{x} \in \mathcal{D}} \log(|\langle \zeta(\mathbf{x}_i) | U_i | \mathbf{v}_{-i} \rangle|). \quad (17)$$

We can find a good U_i by first linearizing the $\log(|\cdot|)$ term, which turns Eq. 17 into a Procrustes problem [54] of linear alignment under an isometric constraint. Procrustes problems can be easily solved with a singular value decomposition on the effective matrix being contracted with U_i in Eq. 17 (after linearization), where setting all singular values to 1 gives the optimal isometry. Upon reaching a candidate solution U_i , the nonlinearity $\log(|\cdot|)$ is linearized again and the optimization process repeated until convergence, typically within a few iterations. Full pseudocode for this training is presented in Appendix E.

We note that although computing a vector \mathbf{v}_{-i} for each sample $\mathbf{x} \in \mathcal{D}$ may appear expensive, the use of cached environment tensors reduces the incremental cost of this computation to only $\mathcal{O}(\chi^2 \mathbf{d})$ when carried out in the context of the adapted DMRG procedure of [12], making this a very lightweight addition to the basic continuous-valued MPS model.

7 Numerical Results

We test the continuous-valued MPS BM model on five distinct density estimation tasks. The first is a rotated hypercube, a simple linearly transformed multidimensional uniform distribution, which we use this to explore the scaling of accuracy with increasing bond and feature dimensions. The second and third are the synthetic two moons dataset [55] and the non-synthetic Iris dataset [56], both of which contain a mixture of continuous and discrete variables. The fourth is a dataset of samples from the classical 2D XY model at nonzero temperature, a statistical mechanical model whose partition function has previously been shown amenable to TN methods [57, 58]. Finally, we use a specifically designed synthetic dataset to test the dynamic basis compression training algorithm.

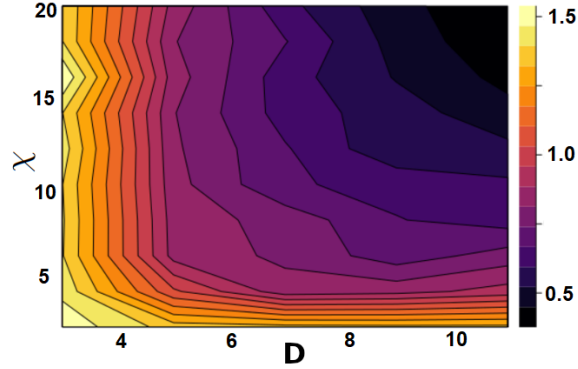


Figure 8: KL divergence of continuous MPS on rotated hypercube dataset, trained with different feature dimensions D and bond dimensions χ .

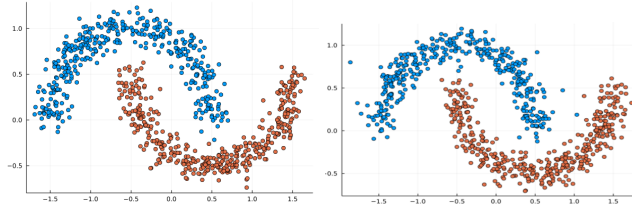


Figure 9: Left plot: 800 samples from the Two Moons distribution, $\sigma = 0.1$. Right plot: 800 samples from our model with $\chi = 8$, $D = 17$.

505 In each setting the continuous-valued MPS BM is trained to minimize the NLL loss on the tar-
 506 get dataset using a two-site DMRG procedure. This is equivalent to minimizing the KL divergence
 507 of the model’s learned PDF relative to the distribution which produced the target dataset, and in
 508 cases where the entropy of the target distribution can be accurately estimated, we will report the
 509 KL divergence of the model. Otherwise we report the raw values of the NLL loss, which can be
 510 negative in the continuous-valued setting. Any experimental details not specified below can be
 511 found in Appendix D.

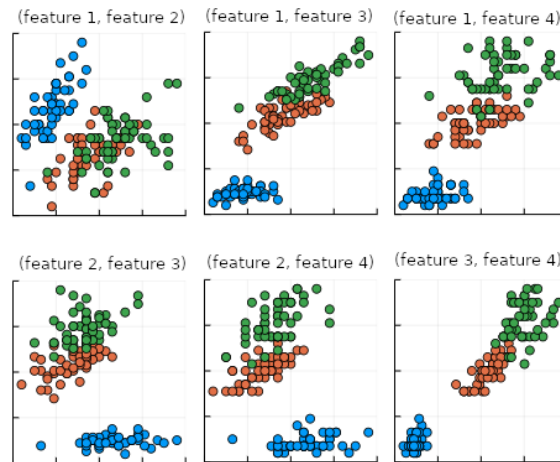
512 7.1 Rotated Hypercube

513 As a simple testbed, we used a distribution drawn uniformly from a rotated hypercube $[-1, 1]^N$
 514 for $N = 5$. This dataset has nontrivial correlations between each pair of variables and sharp jumps
 515 in the overall density, yet still has continuously differentiable marginals.

516 For each feature dimension D and maximum bond dimension χ , the MPS was trained from
 517 an initial random state with 18 DMRG sweeps and a maximum bond dimension that increased
 518 linearly up to χ . The KL divergence of the model for different values of χ and D make use of the
 519 Fourier feature map, which was found to work best in this setting. The KL divergence of the model
 520 for different bond and feature dimensions are plotted in Fig. 8. As expected, the loss decreases as
 521 we improve either dimension, and saturates if one is increased without the other.

522 We note that both real and complex tensor networks can be utilized for continuous-valued
 523 BMs, and during this initial set of experiments, we quickly found that real-valued tensor networks
 524 empirically performed much worse, often failing to converge at all (see Appendix D.1). We at-
 525 tribute this to large jumps in the MPS during the truncation process when using two-site DMRG,
 526 and speculate that better behavior might be observed for real-valued models when training using
 527 gradient descent. Because of this behavior though, all remaining experiments were carried out
 528 using complex-valued tensor networks.

a) True Data



b) Sampled Data

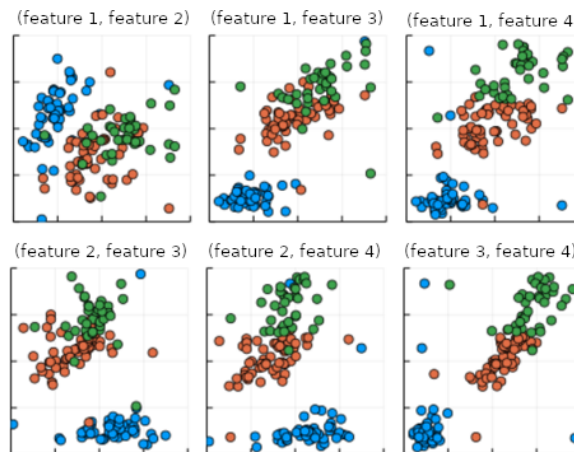


Figure 10: The six different pairwise marginals between each pair of four continuous variables associated to (a) the 150 samples in the Iris dataset, and (b) 150 samples drawn from the continuous MPS model. The three class labels are indicated by color.

529 7.2 Two Moons

530 The **two moons** dataset is a standard synthetic dataset available from scikit-learn [55], containing
 531 two continuous features encoding the position on a 2D plane and one binary feature indicating
 532 which “moon” the sample belongs to. We use a three-site MPS containing two continuous indices
 533 and one discrete index to learn the structure of the dataset in an unsupervised manner, but note that
 534 the efficient conditional sampling permits the trained MPS BM model to be immediately used as
 535 either a supervised classifier or a conditional generative model.

536 Hermite and Fourier feature maps were both tested on the dataset, with an identical training
 537 schedule used as for the rotated cube. We found the Fourier basis to give favorable performance
 538 at all parameters, with a comparison of samples from the trained model with those from the two
 539 moons distribution shown in Fig. 9. More information, including the KL divergence at different
 540 values of D and χ , can be found in Appendix D.2.

541 7.3 Iris Dataset

542 The Iris dataset [59] has four continuous features and a three-class categorical feature. Being a
 543 small dataset of only 150 samples, we must pay attention to overfitting. For each bond dimension
 544 and feature dimension under consideration, we use five-fold cross validation, and report the mean
 545 of the NLL loss on the validation set in each of the five folds. The petal measurements are strictly
 546 positive values, so the Legendre feature map seemed the most natural in this regard. However, we
 547 again found that the Fourier feature map performed the best in practice. Although the Iris dataset
 548 was used here for an unsupervised density modeling task, we note that as in the two moons task,
 549 the MPS BM can immediately be used to either predict the class label given the four continuous
 550 features, or conditionally generate continuous samples given a specific class.

551 We found that overfitting did occur at higher dimensions, with higher losses being seen on
 552 the held-out data fold (see Appendix D.3 for the NLL loss as a function of χ and D). Optimal
 553 performance was observed at $\chi = 9$ and $D = 7$, with a validation loss of -1.40 ± 0.01 . The
 554 samples in the Iris dataset are compared to a similar number of samples from the trained MPS BM
 555 in Fig. 10, where the four continuous features are displayed as six pairwise marginals. The trained
 556 model shows good agreement with the original Iris dataset, although some outliers are visible.

557 7.4 XY Model

558 The classical XY model [60, 61] is a physical system of 2D unit vectors \vec{v}_i , with an interaction
 559 energy $E_{i,j} = -\vec{v}_i \cdot \vec{v}_j$ between adjacent sites. For an N site system, representing each vector
 560 $\vec{v}_i = (\sin x_i, \cos x_i)$ by its angle $x_i \in [0, 2\pi]$ gives N continuous features for each sample drawn
 561 from the thermal ensemble associated to the interaction Hamiltonian. This feature space has a
 562 natural periodic structure, allowing a further test of the Fourier feature map. We chose $N = 16$,
 563 with the associated sites arranged in a 4×4 grid. To ensure a challenging long-range correla-
 564 tion structure, we trained on a dataset of samples drawn from the model's thermal distribution at
 565 temperature $T = 0.8$ which was close to the model's critical temperature of $T_c \approx 0.882$.

566 The MPS BM model for $\chi = 12$ and $D = 13$ was able to reach a KL divergence of ap-
 567 proximately 0.52 relative to the true XY distribution, which was lower than the KL divergence
 568 of 0.6 found by a variational autoencoder (VAE) benchmark with hidden dimension of 512 and
 569 10-dimensional latent space. The VAE benchmark additionally required careful hyperparameter
 570 tuning and several attempts to reach this value, whereas the continuous-valued MPS was able to
 571 reach a lower KL divergence without any modification. Other derived metrics were used to further
 572 verify the performance of the MPS model, as reported in Appendix D.4.

573 7.5 Compression Test

574 To verify that the performance of the compression layer, we created a synthetic dataset containing
 575 several tightly-grouped variables (see Appendix D.5 for details). The dataset possesses four con-
 576 tinuous features with very different single-site marginals, which are shown in Fig. 11. To assess
 577 the impact of the compression layer, we compared three continuous-valued MPS models: (a) a
 578 larger MPS model with $D = 16$, (b) a smaller MPS model with $D = 3$, and (c) a compression-
 579 enhanced MPS model with distinct feature dimension $D = 16$ and site dimension $d = 3$. Although
 580 model (c) employs the same number of feature functions as the larger model (a), its reduced site
 581 dimension makes its computational cost closer to the smaller model (b).

582 The single-site marginal distributions of the three trained models are shown in Fig. 11(a-c),
 583 where it is evident that models (a) and (c) give a more faithful reconstruction of the dataset struc-
 584 ture than model (b). This is supported by the final NLL loss of the trained models, with model
 585 (a) attaining the best NLL loss of -2.17, followed by model (c) with a comparable loss of -2.05,
 586 and finally model (b) with a much higher loss of 2.04. We therefore see that by using compres-
 587 sion layers, continuous-valued MPS with small site dimensions can deliver performance that is

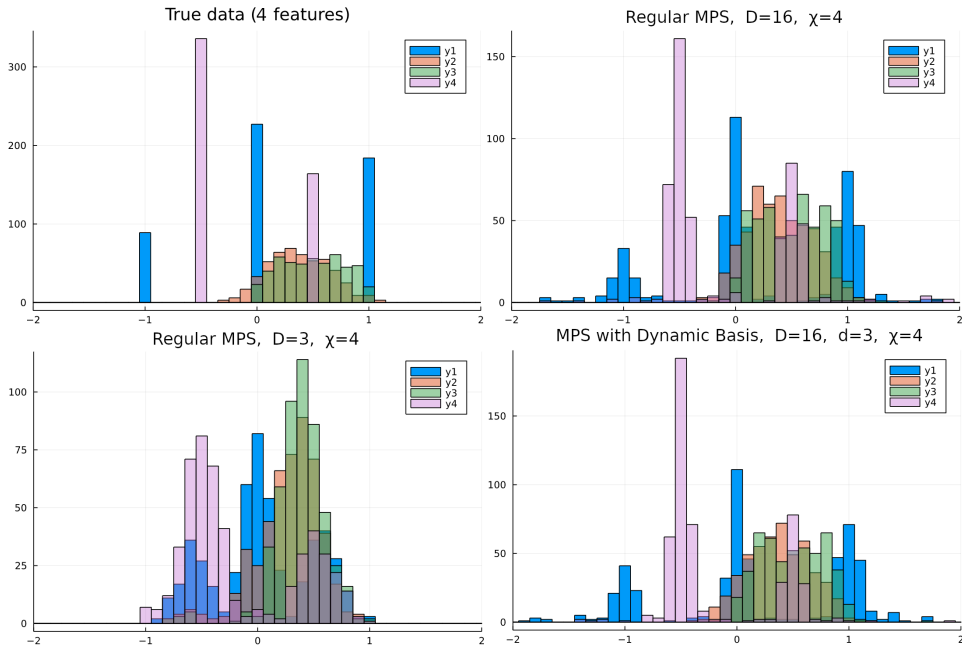


Figure 11: Comparison of the single-feature marginal distributions for three representative MPS models on the synthetic compression dataset, with bond dimension $\chi = 4$. TL: True one-site marginals. TR: an MPS with $D = 16$ obtained the best recovery of the target dataset. BL: The smaller MPS with $D = 3$ gave considerably worse behavior. BT: Using a compression layer with site dimension $d = 3$ and feature $D = 16$ gave comparable performance to the larger model, while maintaining a reduced computational cost.

588 nearly identical to much larger MPS, but without a significant increase in computational cost or
 589 parameter count.

590 8 Conclusions

591 We have introduced a family of continuous-valued TN generative models, which share the per-
 592 fect sampling and conditional generation properties of standard discrete-valued TN BMs, while
 593 also permitting the use of arbitrary combinations of continuous and discrete data. The gener-
 594 ality of these models is proven by a pair of universal approximation theorems, which ensure that
 595 any sufficiently smooth PDF or continuous-valued wavefunction can be efficiently represented to
 596 arbitrary precision using continuous-valued MPS. Benchmarking this model on a broad range of
 597 synthetic and real-world datasets with discrete and continuous variables, we find it able to accu-
 598 rately learn the structure of each dataset, with a programmable compression layer giving enhanced
 599 performance in the presence of limited computational resources.

600 A key ingredient in our continuous generalization is the notion of feature maps to embed con-
 601 tinuous data as finite-dimensional vectors. While feature maps have been used for supervised TN
 602 models since at least [11, 18], a major contribution of our work is the discovery of much richer
 603 structure in feature maps within the context of generative modeling. We prove a general character-
 604 ization of the influence of feature maps on the marginal distributions of continuous-valued MPS
 605 BMs at initialization, and investigate several concrete feature maps in detail from a theoretical and
 606 empirical perspective. Our focus on isometric feature maps, which we prove entails no loss of
 607 generality, lets us derive a canonical form for continuous-valued MPS that preserves the conve-
 608 nient properties of discrete-valued MPS and permits the use of powerful methods like DMRG for

609 optimization.

610 While we have restricted to the use of MPS for convenience, in principle any discrete-valued
611 TN can be extended by our methods into a corresponding continuous-valued model, and bench-
612 marking the performance of more sophisticated TNs (e.g. tree TNs, MERA, and PEPS) in prob-
613 lems with continuous data is an obvious next step. A more open-ended direction is to develop
614 methods for boosting the expressivity of feature maps, or choosing them based on the structure
615 of particular datasets. Our compression layer represents an important contribution along these
616 lines, but using neural networks or other ML models may boost expressivity yet further. Develop-
617 ing heuristics for better choosing the feature dimension D in a given problem, analogous to how
618 entanglement-based area laws guide the choice of bond dimension χ , is another problem deserving
619 future attention. Along similar lines, we anticipate generalizations of two-site update scheme that
620 permit the dynamic variation of both D and χ to be a useful aid for optimizing continuous-valued
621 TN models.

622 9 Acknowledgments

623 The authors would like to thank Vladimir Vargas-Calderón for contributing the VAE benchmark
624 result. G.R.'s research was supported by the Canadian Institute for Advanced Research (CIFAR
625 AI chair program).

626 References

- 627 [1] I. Affleck, T. Kennedy, E. Lieb and H. Tasaki, *Rigorous results on valence-bond ground*
628 *states in antiferromagnets*, Physical Review Letters **59**(7), 799 (1987).
- 629 [2] M. Fannes, B. Nachtergaele and R. F. Werner, *Finitely correlated states on quantum spin*
630 *chains*, Communications in mathematical physics **144**(3), 443 (1992).
- 631 [3] M. M. Wolf, F. Verstraete, M. B. Hastings and J. I. Cirac, *Area laws in quantum systems:*
632 *mutual information and correlations*, Physical review letters **100**(7), 070502 (2008).
- 633 [4] R. Orús, *A practical introduction to tensor networks: Matrix product states and projected*
634 *entangled pair states*, Annals of physics **349**, 117 (2014).
- 635 [5] G. Vidal, *Efficient classical simulation of slightly entangled quantum computations*, Physical
636 review letters **91**(14), 147902 (2003).
- 637 [6] C. Schön, E. Solano, F. Verstraete, J. I. Cirac and M. M. Wolf, *Sequential generation of*
638 *entangled multiqubit states*, Physical review letters **95**(11), 110503 (2005).
- 639 [7] Y.-Y. Shi, L.-M. Duan and G. Vidal, *Classical simulation of quantum many-body systems*
640 *with a tree tensor network*, Physical review A **74**(2), 022320 (2006).
- 641 [8] T. G. Kolda and B. W. Bader, *Tensor decompositions and applications*, SIAM review **51**(3),
642 455 (2009).
- 643 [9] I. V. Oseledets, *Tensor-train decomposition*, SIAM Journal on Scientific Computing **33**(5),
644 2295 (2011).
- 645 [10] A. Cichocki, N. Lee, I. Oseledets, A.-H. Phan, Q. Zhao, D. P. Mandic *et al.*, *Tensor networks*
646 *for dimensionality reduction and large-scale optimization: Part 1 low-rank tensor decompo-*
647 *sitions*, Foundations and Trends® in Machine Learning **9**(4-5), 249 (2016).

- 648 [11] E. Stoudenmire and D. J. Schwab, *Supervised learning with tensor networks*, Advances in
649 Neural Information Processing Systems **29** (2016).
- 650 [12] Z.-Y. Han, J. Wang, H. Fan, L. Wang and P. Zhang, *Unsupervised generative modeling using*
651 *matrix product states*, Phys. Rev. X **8**(3), 031012 (2018).
- 652 [13] W. Huggins, P. Patil, B. Mitchell, K. B. Whaley and E. M. Stoudenmire, *Towards quantum*
653 *machine learning with tensor networks*, Quantum Science and Technology **4**(2), 024001
654 (2019), doi:[10.1088/2058-9565/aea94](https://doi.org/10.1088/2058-9565/aea94).
- 655 [14] R. Orús, *Tensor networks for complex quantum systems*, Nature Reviews Physics **1**(9), 538
656 (2019).
- 657 [15] S. R. White, *Density matrix formulation for quantum renormalization groups*, Physical
658 review letters **69**(19), 2863 (1992).
- 659 [16] N. Cohen, O. Sharir and A. Shashua, *On the expressive power of deep learning: A tensor*
660 *analysis*, In *Conference on learning theory*, pp. 698–728. PMLR (2016).
- 661 [17] V. Khurikov, A. Novikov and I. Oseledets, *Expressive power of recurrent neural networks*,
662 arXiv preprint arXiv:1711.00811 (2017).
- 663 [18] A. Novikov, M. Trofimov and I. Oseledets, *Exponential machines*, arXiv:1605.03795 (2016).
- 664 [19] C. Yin, B. Acun, X. Liu and C.-J. Wu, *Tt-rec: Tensor train compression for deep learning*
665 *recommendation models* (2021), [arXiv:2101.11714](https://arxiv.org/abs/2101.11714).
- 666 [20] Y. Panagakis, J. Kossaiifi, G. G. Chrysos, J. Oldfield, M. A. Nicolaou, A. Anandkumar and
667 S. Zafeiriou, *Tensor methods in computer vision and deep learning*, Proceedings of the IEEE
668 **109**(5), 863 (2021).
- 669 [21] H. Zhou, L. Li and H. Zhu, *Tensor regression with applications in neuroimaging data anal-*
670 *ysis*, Journal of the American Statistical Association **108**(502), 540 (2013).
- 671 [22] Q. Xie, Q. Zhao, D. Meng, Z. Xu, S. Gu, W. Zuo and L. Zhang, *Multispectral images*
672 *denoising by intrinsic tensor sparsity regularization*, In *Proceedings of the IEEE conference*
673 *on computer vision and pattern recognition*, pp. 1692–1700 (2016).
- 674 [23] J. Liu, S. Li, J. Zhang and P. Zhang, *Tensor networks for unsupervised machine learning*,
675 Physical Review E **107**(1), L012103 (2023).
- 676 [24] J. Miller, G. Rabusseau and J. Terilla, *Tensor networks for probabilistic sequence modeling*,
677 In *International Conference on Artificial Intelligence and Statistics*, pp. 3079–3087, (PMLR)
678 (2021).
- 679 [25] S. Cheng, J. Chen and L. Wang, *Information perspective to probabilistic modeling: Boltz-*
680 *mann machines versus born machines*, Entropy **20**(8), 583 (2018), doi:[10.3390/e20080583](https://doi.org/10.3390/e20080583).
- 681 [26] I. Glasser, N. Pancotti and J. I. Cirac, *From probabilistic graphical models to gen-*
682 *eralized tensor networks for supervised learning*, IEEE Access **8**, 68169 (2020),
683 doi:[10.1109/ACCESS.2020.2986279](https://doi.org/10.1109/ACCESS.2020.2986279).
- 684 [27] S. Cheng, L. Wang, T. Xiang and P. Zhang, *Tree tensor networks for generative modeling*,
685 Phys. Rev. B **99**, 155131 (2019).
- 686 [28] Y. Levine, O. Sharir, N. Cohen and A. Shashua, *Quantum entanglement in deep learning*
687 *architectures*, Physical review letters **122**(6), 065301 (2019).

- 688 [29] S. Lu, M. Kanász-Nagy, I. Kukuljan and J. I. Cirac, *Tensor networks and efficient descrip-*
689 *tions of classical data*, arXiv preprint arXiv:2103.06872 (2021).
- 690 [30] I. Glasser, R. Sweke, N. Pancotti, J. Eisert and I. Cirac, *Expressive power of tensor-network*
691 *factorizations for probabilistic modeling*, Advances in neural information processing sys-
692 *tems* **32** (2019).
- 693 [31] S. Adhikary, S. Srinivasan, J. Miller, G. Rabusseau and B. Boots, *Quantum tensor networks,*
694 *stochastic processes, and weighted automata*, In *International Conference on Artificial In-*
695 *telligence and Statistics*, pp. 2080–2088. PMLR (2021).
- 696 [32] A. J. Ferris and G. Vidal, *Perfect sampling with unitary tensor networks*, Physical Review B
697 **85**(16), 165146 (2012).
- 698 [33] J. Martyn, G. Vidal, C. Roberts and S. Leichenauer, *Entanglement and tensor networks for*
699 *supervised image classification*, arXiv preprint arXiv:2007.06082 (2020).
- 700 [34] J. Reyes and M. Stoudenmire, *A multi-scale tensor network architecture for classification*
701 *and regression* (2020).
- 702 [35] I. Convy and K. B. Whaley, *Interaction decompositions for tensor network regression*, Mach.
703 Learn. Sci. Technol. **3**(4), 045027 (2022).
- 704 [36] T. Hao, X. Huang, C. Jia and C. Peng, *A quantum-inspired tensor network method for*
705 *constrained combinatorial optimization problems*, arXiv preprint arXiv:2203.15246 (2022).
- 706 [37] J.-G. Liu, X. Gao, M. Cain, M. D. Lukin and S.-T. Wang, *Computing solution space prop-*
707 *erties of combinatorial optimization problems via generic tensor networks*, arXiv preprint
708 arXiv:2205.03718 (2022).
- 709 [38] J. Eisert, M. Cramer and M. B. Plenio, *Colloquium: Area laws for the entanglement entropy*,
710 Reviews of modern physics **82**(1), 277 (2010).
- 711 [39] M. Schuld and N. Killoran, *Quantum machine learning in feature hilbert spaces*, Physical
712 Review Letters **122** (2019), doi:[10.1103/physrevlett.122.040504](https://doi.org/10.1103/physrevlett.122.040504).
- 713 [40] S. Wahls, V. Koivunen, H. V. Poor and M. Verhaegen, *Learning multidimensional fourier*
714 *series with tensor trains*, In *2014 IEEE Global Conference on Signal and Information Pro-*
715 *cessing (GlobalSIP)*, pp. 394–398. IEEE (2014).
- 716 [41] S.-H. Lin, O. Kuijpers, S. Peterhansl and F. Pollmann, *Distributive pre-training of generative*
717 *modeling using matrix-product states*, arXiv preprint arXiv:2306.14787 (2023).
- 718 [42] M. Ali and A. Nouy, *Approximation with tensor networks. part iii: Multivariate approxima-*
719 *tion*, arXiv preprint arXiv:2101.11932 (2021).
- 720 [43] M. Ali and A. Nouy, *Approximation theory of tree tensor networks: Tensorized univariate*
721 *functions*, Constructive Approximation pp. 1–82 (2023).
- 722 [44] F. Wesel and K. Batselier, *Quantized fourier and polynomial features for more expressive*
723 *tensor network models*, In *International Conference on Artificial Intelligence and Statistics*,
724 pp. 1261–1269. PMLR (2024).
- 725 [45] D. Bigoni, A. P. Engsig-Karup and Y. M. Marzouk, *Spectral tensor-train decomposition*,
726 SIAM Journal on Scientific Computing **38**(4), A2405 (2016).

- 727 [46] A. Gorodetsky, S. Karaman and Y. Marzouk, *A continuous analogue of the tensor-train*
728 *decomposition*, Computer methods in applied mechanics and engineering **347**, 59 (2019).
- 729 [47] Y. Hur, J. G. Hoskins, M. Lindsey, E. M. Stoudenmire and Y. Khoo, *Generative modeling via*
730 *tensor train sketching*, Applied and Computational Harmonic Analysis **67**, 101575 (2023).
- 731 [48] G. S. Novikov, M. E. Panov and I. V. Oseledets, *Tensor-train density estimation*, In *Uncer-*
732 *tainty in artificial intelligence*, pp. 1321–1331. PMLR (2021).
- 733 [49] F. Verstraete and J. I. Cirac, *Continuous matrix product states for quantum fields*, Phys. Rev.
734 Lett. **104**, 190405 (2010), doi:[10.1103/PhysRevLett.104.190405](https://doi.org/10.1103/PhysRevLett.104.190405).
- 735 [50] B. Ellefsen, *Math StackExchange: Averaged value of product of Legendre Polynomials*
736 (2023).
- 737 [51] M. Abramowitz and I. A. Stegun, *Handbook of mathematical functions with formulas,*
738 *graphs, and mathematical tables*, vol. 55, US Government printing office (1968).
- 739 [52] D. Borwein, J. M. Borwein and R. E. Crandall, *Effective laguerre asymptotics*, SIAM Journal
740 on Numerical Analysis **46**(6), 3285 (2008).
- 741 [53] S. Holtz, T. Rohwedder and R. Schneider, *On manifolds of tensors of fixed tt-rank*, Nu-
742 merische Mathematik **120**(4), 701 (2012).
- 743 [54] J. Gower and G. Dijkstra, *Procrustes problems.*, Oxford University Press (2004).
- 744 [55] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel,
745 P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos *et al.*, *Scikit-learn: Machine*
746 *learning in Python*, Journal of Machine Learning Research **12**, 2825 (2011).
- 747 [56] R. A. Fisher, *The use of multiple measurements in taxonomic problems*, Annals of eugenics
748 **7**(2), 179 (1936).
- 749 [57] J. Yu, Z. Xie, Y. Meurice, Y. Liu, A. Denbleyker, H. Zou, M. Qin, J. Chen and T. Xiang,
750 *Tensor renormalization group study of classical x y model on the square lattice*, Physical
751 Review E **89**(1), 013308 (2014).
- 752 [58] L. Vanderstraeten, B. Vanhecke, A. M. Läuchli and F. Verstraete, *Approaching the kosterlitz-*
753 *thouless transition for the classical x y model with tensor networks*, Physical Review E
754 **100**(6), 062136 (2019).
- 755 [59] M. Lichman, *UCI machine learning repository* (2013).
- 756 [60] N. D. Mermin and H. Wagner, *Absence of ferromagnetism or antiferromagnetism in*
757 *one- or two-dimensional isotropic heisenberg models*, Phys. Rev. Lett. **17**, 1133 (1966),
758 doi:[10.1103/PhysRevLett.17.1133](https://doi.org/10.1103/PhysRevLett.17.1133).
- 759 [61] Y. Nambu, *A Note on the Eigenvalue Problem in Crystal Statistics*, Progress of Theoretical
760 Physics **5**(1), 1 (1950), doi:[10.1143/ptp/5.1.1](https://doi.org/10.1143/ptp/5.1.1), <https://academic.oup.com/ptp/article-pdf/5/1/1/5253488/5-1-1.pdf>.
- 761
- 762 [62] M. A. Nielsen and I. Chuang, *Quantum computation and quantum information* (2002).
- 763 [63] N. Metropolis, A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller and E. Teller, *Equation of*
764 *state calculations by fast computing machines*, The Journal of Chemical Physics **21**(6), 1087
765 (1953), doi:[10.1063/1.1699114](https://doi.org/10.1063/1.1699114), <https://doi.org/10.1063/1.1699114>.

766 [64] F. Perez-Cruz, *Kullback-leibler divergence estimation of continuous distributions*,
 767 In *2008 IEEE International Symposium on Information Theory*, pp. 1666–1670,
 768 doi:10.1109/ISIT.2008.4595271 (2008).

769 A Generality of Isometric Feature Map Condition

770 Given an arbitrary feature map $\zeta : \mathcal{I} \rightarrow \mathbb{K}^D$ represented by a basis of feature functions given by
 771 $\mathcal{F} = \{f_1, f_2, \dots, f_D\}$, we present a general procedure to create new feature functions satisfying
 772 Eq. 6, thereby giving an isometric map. We further show how this procedure can be applied to
 773 any continuous-valued MPS without imposing any changes in the associated continuous-valued
 774 function $\Phi : \mathcal{I}^N \rightarrow \mathbb{K}$. We assume first that the functions are linearly independent, as otherwise
 775 we can remove any linearly dependent basis functions without any impact on the feature map’s
 776 expressive power.

777 First, we calculate a Hermitian “overlap matrix” $M \in \mathbb{K}^{D \times D}$ whose elements M_{ij} give the
 778 overlap between feature functions f_i and f_j , namely

$$M_{ij} = \langle f_i, f_j \rangle = \int_{x \in \mathcal{I}} f_i^*(x) f_j(x) dx. \quad (18)$$

779 By the assumption of all functions in \mathcal{F} being linearly independent, M is full-rank, which allows
 780 its matrix inverse square root $X = M^{-\frac{1}{2}}$ to be computed from the spectral decomposition of M :

$$M = Q \Lambda Q^\dagger \quad (19)$$

$$X = Q \Lambda^{-\frac{1}{2}} Q^\dagger \quad (20)$$

781 where Q is a $D \times D$ unitary and $\Lambda^{-\frac{1}{2}}$ denotes the elementwise inverse square root of the diagonal
 782 matrix Λ containing strictly positive diagonal entries. X is itself an invertible Hermitian matrix,
 783 and Eq. 19 and Eq. 20 can be used to verify that $XMX = I$

784 Using X we can generate a new isometric feature map, whose basis of feature functions
 785 $\{g_1, g_2, \dots, g_D\}$ is given by

$$g_k(x) = \sum_{j=1}^D f_j(x) X_{jk}. \quad (21)$$

786 The isometric nature of the new feature map can be verified by the orthonormality of the feature
 787 functions, concretely:

$$\begin{aligned} & \int_{x \in \mathcal{I}} g_i^*(x) g_j(x) dx \\ &= \int_{x \in \mathcal{I}} \left(\sum_a f_a^*(x) X_{ai}^* \right) \left(\sum_b f_b(x) X_{bj} \right) dx \\ &= \sum_{a,b} X_{ai}^* X_{bj} \left(\int f_a^*(x) f_b(x) dx \right) \\ &= \sum_{a,b} X_{ai}^* X_{bj} M_{ab} = (X^\dagger M X)_{ij} = \delta_{ij}. \end{aligned}$$

788 Finally, we note that this transformation can be applied to an existing continuous-valued MPS
 789 model without any change in the associated function Φ . The new feature map defined by Eq. 21
 790 is equivalent to the composite function sending $x \mapsto \zeta(x)X$, and applying the inverse matrix

811 **Lemma 1.** Consider a tensor $A^{(i)} \in \mathbb{K}^{\chi_{i-1} \times \chi_i \times D_i}$ whose elements have mean $\mathbb{E}[A_{\alpha,\beta,k}^{(i)}] = \mathbf{0}$ and
 812 variance $\mathbb{E}[|A_{\alpha,\beta,k}^{(i)}|^2] = t_i$. The sixth-order variance tensor $B^{(i)}$ given by taking two copies of
 813 $A^{(i)}$ and averaging over all IID initializations is given by

$$\begin{aligned} B_{\alpha,\alpha',\beta,\beta',k,k'}^{(i)} &= \mathbb{E}\left[(A^{(i)})_{\alpha,\beta,k}^* A_{\alpha',\beta',k'}^{(i)}\right] \\ &= t \delta_{\alpha,\alpha'} \delta_{\beta,\beta'} \delta_{k,k'}. \end{aligned} \quad (22)$$

814 The proof of Lemma 1 is an immediate consequence of the IID nature of the different elements
 815 of $A^{(i)}$. The elements of $B^{(i)}$ are covariances between pairs of elements in $A^{(i)}$, which by assump-
 816 tion are 0 for different elements, and t_i for identical elements. The result has a convenient graph-
 817 ical form, shown in Fig. 12a, which facilitates many calculations involving randomly-initialized
 818 MPS.

819 As a concrete example, we can compute the expected squared norm of a continuous-valued
 820 MPS Φ whose underlying discrete-valued MPS ψ has been initialized using core tensors with
 821 IID random elements with mean zero. The isometric nature of the model's feature map leads the
 822 squared norm of a continuous-valued MPS to be identical to that of its underlying discrete-valued
 823 MPS, which with the use of Lemma 1 can be verified to equal the product of all feature and bond
 824 dimensions in the model, namely

$$\mathbb{E}[\|\psi\|^2] = \prod_{i=1}^N t_i D_i \chi_i, \quad (23)$$

825 where we take χ_N to be $\chi_N = \mathbf{1}$. In order to ensure proper normalization when the MPS is used as a
 826 probabilistic BM, it is necessary to have the per-core variances t_i to satisfy $\prod_{i=1}^N t_i = \prod_{i=1}^N D_i \chi_i$,
 827 which can be ensured by taking $t_i = (D_i \chi_i)^{-1}$.

828 Given this, the proof of Theorem 1 reduces to taking the definition of the expected univariate
 829 marginal distribution $P_{\text{init}}(\mathbf{x}_i)$, wherein all other variables are traced out, and applying TN iden-
 830 tities to simplify the resultant expression to the form given in Eq. 13 (see Fig. 12b). Applying
 831 the isometric condition of Eq. 6 allows all pairs of traced-over feature maps to be removed (see
 832 Fig. 3), with Lemma 1 permitting a comparable removal of all pairs of matched tensor cores $A^{(i)}$.
 833 The result is the simple diagram on the right side of Fig. 12b, with a proportionality factor equal
 834 to the product of all t_i with a term $\mathbf{a}_i = \prod_{i=1}^{N-1} \chi_i \prod_{j \neq i} D_j$ coming from tracing over all bond
 835 dimensions, as well as all feature dimensions except for D_i . The result is the scalar factor D_i^{-1} ,
 836 which under the typical assumption of constant feature dimension $D_i = D$, gives the proportional-
 837 ity factor appearing in Eq. 13. This completes the proof of Theorem 1.

838 C Proof of Universal Approximation Results

839 We prove the universality approximation results of Theorems 2 and 3 in the following, which are
 840 restated below for ease of reference. In order to prove these Theorems, we must first introduce
 841 some concepts from functional analysis, which are used to introduce and prove Theorem 4, which
 842 generalizes Theorem 2 to characterize a wider range of functions.

843 C.1 Functional Analysis Preliminaries

844 Our results concern the setting of spaces of scalar-valued functions $f : \Omega \rightarrow \mathbb{K}$ defined on the N -
 845 dimensional hypercube $\Omega = [0, 1]^N \subseteq \mathbb{R}^N$ equipped with L_2 -norm $\|f\|_{L_\mu^2} := \left(\int_{\mathbf{x} \in \Omega} |f(\mathbf{x})|^2 d\mu\right)^{1/2}$

846 associated with a positive-valued finite measure μ (i.e. $\mu(\Omega) < \infty$). We use \mathbb{K} to indicate one
 847 of either \mathbb{R} or \mathbb{C} , in the typical case where the choice of field doesn't change the validity of the
 848 definitions or results.

849 If $\mathbf{i} = (i_1, i_2, \dots, i_k)$ for some $k \geq 0$, then we use $\partial^{(\mathbf{i})} f = \frac{\partial^k f}{\partial x_{i_1} \partial x_{i_2} \dots \partial x_{i_k}}$ to indicate a standard
 850 k 'th-order partial derivative of the function f , with the usual caveat that such derivatives only
 851 exist for sufficiently smooth functions. More generally, we use $\mathbf{D}^{(\mathbf{i})} f$ to indicate a k 'th-order *weak*
 852 derivative of f , which is defined as a function satisfying the formula

$$\int_{\mathbf{x} \in \Omega} (\mathbf{D}^{(\mathbf{i})} f(\mathbf{x})) \varphi(\mathbf{x}) d\mu = (-1)^k \int_{\mathbf{x} \in \Omega} f(\mathbf{x}) (\partial^{(\mathbf{i})} \varphi(\mathbf{x})) d\mu \quad (24)$$

853 for all infinite-differentiable functions $\varphi : \Omega \rightarrow \mathbb{K}$ which vanish on the boundary of Ω . The usual
 854 integration by parts formula ensures that each partial derivative $\partial^{(\mathbf{i})} f$ is itself a weak derivative
 855 of f , but the latter can also be defined for functions f whose k 'th-order partial derivatives don't
 856 exist for all $\mathbf{x} \in \Omega$. A function can possess multiple different weak derivatives, but these will agree
 857 almost everywhere in Ω (i.e. everywhere but a measure zero subset of Ω).

858 We are interested in functions that are "sufficiently nice" for proving universal approximation
 859 results, which leads to the concept of *Sobolev spaces*. The k 'th order Sobolev space $\mathcal{H}_{\mathbb{K}}^k$ on Ω
 860 is defined as the collection of all functions $f : \Omega \rightarrow \mathbb{K}$ possessing all weak derivatives $\mathbf{D}^{(\mathbf{i})} f$ of order
 861 $|\mathbf{i}| \leq k$ (i.e. $\mathbf{i} = (i_1, i_2, \dots, i_\ell)$ for $\ell \leq k$), which each have finite L_2 norm. This is equivalent to
 862 the condition

$$\|f\|_{\mathcal{H}_{\mathbb{K}}^k} := \sum_{|\mathbf{i}| \leq k} \|\mathbf{D}^{(\mathbf{i})} f\|_{L_\mu^2} < \infty, \quad (25)$$

863 where the quantity $\|f\|_{\mathcal{H}_{\mathbb{K}}^k}$ is referred to as the k 'th-order *Sobolev norm* of f . For $k = 0$, the
 864 Sobolev norm reduces to the usual L_2 norm on Ω , and more generally $\|f\|_{L_\mu^2} \leq \|f\|_{\mathcal{H}_{\mathbb{K}}^k}$, so that
 865 every $f \in \mathcal{H}_{\mathbb{K}}^k$ necessarily has bounded L_2 norm. We will also employ the so-called *Sobolev*
 866 *seminorm* $|f|_{\mathcal{H}_{\mathbb{K}}^k}$ of a function $f \in \mathcal{H}_{\mathbb{K}}^k$, which is defined as $|f|_{\mathcal{H}_{\mathbb{K}}^k} := \sum_{|\mathbf{i}|=k} \|\mathbf{D}^{(\mathbf{i})} f\|_{L_\mu^2}$.

867 A final concept needed in the following is the notion of α -Hölder continuity, where a function
 868 $f : \Omega \rightarrow \mathbb{K}$ is bounded in variation as

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq C \|\mathbf{x} - \mathbf{y}\|^\alpha, \quad (26)$$

869 for $C \in \mathbb{R}$ a constant holding for any pair of points $\mathbf{x}, \mathbf{y} \in \Omega$. We can without loss of generality
 870 take α to be in the range $\alpha \in (0, 1]$, and note that when $\alpha = 1$ the notion of α -Hölder continuity
 871 reduces to the more familiar definition of Lipschitz continuity. Any function $f : \Omega \rightarrow \mathbb{K}$ whose
 872 first derivatives exist at all points in Ω and are bounded as $\left| \frac{\partial f}{\partial x_i} \right| < \infty$ (for $i = 1, 2, \dots, N$) will
 873 always be Lipschitz continuous, and therefore α -Hölder continuous for any $0 < \alpha \leq 1$.

874 C.2 Proof of Theorem 2

875 **Theorem 2.** Consider a family of continuous-valued MPS with polynomial feature functions
 876 $\mathcal{F} = \{f_1, f_2, \dots\}$ forming an orthonormal basis for $[0, 1]$, which is defined on the hypercube
 877 $\Omega = [0, 1]^N \subseteq \mathbb{R}^N$. Let $k \geq N$ and let $\Phi : \Omega \rightarrow \mathbb{C}$ be any square-integrable function with unit
 878 norm ($\langle \Phi, \Phi \rangle = 1$), whose partial derivatives of order $1, 2, \dots, k$ all exist and are bounded. Then
 879 for every positive $\chi, D \in \mathbb{N}$ there exists a continuous-valued MPS of bond dimension χ and feature
 880 dimension D with unit norm, whose associated function $\Phi_{\text{MPS}}^{(\chi, D)}$ approximates Φ with infidelity

$$1 - |\langle \Phi, \Phi_{\text{MPS}}^{(\chi, D)} \rangle| \leq \gamma_1 \chi^{-k+1} + \gamma_2 D^{-2k}, \quad (15)$$

881 where $\gamma_1, \gamma_2 > 0$ depend on the target function Φ , the assumed degree of smoothness k , and the
 882 feature functions \mathcal{F} .

883 Rather than proving Theorem 2 directly, we instead prove a more general Theorem 4, given
 884 below. The fact that Theorem 4 implies Theorem 2 is immediate from the definitions and facts
 885 above concerning Sobolev spaces and Hölder continuity.

886 **Theorem 4.** Consider a family of continuous-valued MPS with polynomial feature functions
 887 $\mathcal{F} = \{f_1, f_2, \dots\}$ forming an orthonormal basis for $[0, 1]$, which is defined on the hypercube
 888 $\Omega = [0, 1]^N \subseteq \mathbb{R}^N$. Let $k \geq N$ and let $\Phi : \Omega \rightarrow \mathbb{C}$ be any square-integrable function in the
 889 Sobolev space \mathcal{H}_C^k with unit L_2 norm which is α -Hölder continuous for $\alpha > \frac{1}{2}$. Then for every
 890 positive $\chi, D \in \mathbb{N}$ there exists a continuous-valued MPS of bond dimension χ and feature dimen-
 891 sion D of unit norm, whose associated function $\Phi_{\text{MPS}}^{(\chi, D)}$ approximates Φ with infidelity

$$1 - |\langle \Phi, \Phi_{\text{MPS}}^{(\chi, D)} \rangle| \leq \gamma_1 \chi^{-k+1} + \gamma_2 D^{-2k},$$

892 where $\gamma_1, \gamma_2 > 0$ depend on the target function Φ , the assumed degree of smoothness k , and the
 893 feature functions \mathcal{F} .

894 This more general formulation allows us to make use of an invaluable result from [45], which
 895 applies to functional tensor train (FTT) decompositions that are almost identical to the continuous-
 896 valued MPS considered here. The result in question comes from the fundamental FTT approxi-
 897 mation characterization given in their Theorem 13 with a polynomial interpolation method, as
 898 expressed in their Eqs. 66, 70, and 73². Rephrased in our terminology and notation, this result
 899 takes the form of:

900 **Lemma 2** ([45]). Let $\Phi : \Omega \rightarrow \mathbb{R}$ be a \mathcal{H}_R^k function on $\Omega = [0, 1]^N \subseteq \mathbb{R}^N$ which is α -Hölder
 901 continuous for $\alpha > \frac{1}{2}$, and where $k \geq N$. Then for any collection $\mathcal{F} = \{f_1, f_2, \dots\}$ of polynomial
 902 feature functions which form an orthonormal basis for $[0, 1]$, for every positive $\chi, D \in \mathbb{N}$ there
 903 exists a continuous-valued MPS with bond dimension χ and feature dimension D which computes
 904 a function $\Phi_{\text{MPS}}^{(\chi, D)} : \Omega \rightarrow \mathbb{R}$ satisfying

$$\begin{aligned} \|\Phi - \Phi_{\text{MPS}}^{(\chi, D)}\|_{L_\mu^2} &\leq \sqrt{\frac{N-1}{k-1}} \|\Phi\|_{\mathcal{H}_R^k} (\chi + 1)^{-\frac{k-1}{2}} \\ &\quad + C(k) |\Phi_{\text{MPS}}^{(\chi, D)}|_{\mathcal{H}_R^k} D^{-k}, \end{aligned} \quad (27)$$

905 with $C(k)$ depending on k and (implicitly) on the choice of \mathcal{F} .

906 The RHS of Lemma 2 contains two polynomials of χ and D whose dependence on k is of
 907 the same order as the two polynomials on the RHS of the bound of Theorem 4. In order to use
 908 the former result to prove the latter though, we must do several things: (a) Replace $\chi + 1$ by χ
 909 in the RHS of Eq. 27; (b) Generalize the setting of Lemma 2 from real-valued to complex-valued
 910 functions; (c) Ensure that $\Phi_{\text{MPS}}^{(\chi, D)}$ can be chosen to have unit L_2 norm whenever Φ does; (d) Bound
 911 the quantity $|\Phi_{\text{MPS}}^{(\chi, D)}|_{\mathcal{H}_C^k}$ by a function of k and \mathcal{F} alone; and (e) Convert the L_2 bound derived from
 912 Lemma 2 into the infidelity bound appearing in Theorem 4. We will proceed to do each of these
 913 in the following.

914 **Replace $\chi + 1$ by χ** This is straightforward, as for any positive values $K > 0$ and $m \geq 1$,
 915 the inequality $K(\chi + 1)^{-m} \leq 2K\chi^{-m}$ holds for all bond dimensions $\chi \geq 1$. This replacment
 916 therefore adds a factor of 2 to the first term on the RHS of Eq. 27.

²All equation and theorem references are relative to the published version of [45]

917 **Complex-valued generalization** Although Lemma 2 is phrased in terms of real-valued func-
 918 tions, generalizing this result to complex-valued functions is straightforward. The target function
 919 $\Phi^c : \Omega \rightarrow \mathbb{C}$ can be represented as a weighted sum of two real-valued functions $\Phi^r, \Phi^i : \Omega \rightarrow \mathbb{R}$
 920 via $\Phi^c(\mathbf{x}) = \Phi^r(\mathbf{x}) + i\Phi^i(\mathbf{x})$, and each function approximated separately by continuous MPS
 921 $\Phi_{MPS}^r, \Phi_{MPS}^i$ of bond dimension $\frac{\chi}{2}$ (assuming wlog that χ is even). The two approximating MPS
 922 can be summed together as a single continuous MPS Φ_{MPS}^c of bond dimension χ , giving

$$\begin{aligned} \|\Phi^c - \Phi_{MPS}^c\|_{L_\mu^2} &\leq \|\Phi^r - \Phi_{MPS}^r\|_{L_\mu^2} + \|\Phi^i - \Phi_{MPS}^i\|_{L_\mu^2} \\ &\leq 2\sqrt{\frac{N-1}{k-1}} \left(\|\Phi^r\|_{\mathcal{H}_R^k} + \|\Phi^i\|_{\mathcal{H}_R^k} \right) \left(\frac{\chi}{2} \right)^{-\frac{k-1}{2}} \\ &\quad + C(k) \left(|\Phi_{MPS}^r|_{\mathcal{H}_R^k} + |\Phi_{MPS}^i|_{\mathcal{H}_R^k} \right) D^{-k} \end{aligned} \quad (28)$$

$$\begin{aligned} &\leq 2^{\frac{k}{2}+1} \sqrt{\frac{N-1}{k-1}} \|\Phi^c\|_{\mathcal{H}_C^k} \chi^{-\frac{k-1}{2}} \\ &\quad + \sqrt{2}C(k) |\Phi_{MPS}^c|_{\mathcal{H}_C^k} D^{-k}, \end{aligned} \quad (29)$$

923 where we have used the identity $\|\Phi^r\|_{\mathcal{H}_R^k} + \|\Phi^i\|_{\mathcal{H}_R^k} \leq \sqrt{2}\|\Phi^c\|_{\mathcal{H}_C^k}$ (a basic consequence of complex
 924 versus real L_2 norms) for the Sobolev norm and seminorm.

925 **Ensure $\Phi_{MPS}^{(\chi,D)}$ has unit norm** Theorem 4 not only assumes a target function Φ with unit norm,
 926 but also ensures a continuous MPS approximation with unit norm. This guarantee is not provided
 927 by Lemma 2, whose approximating function $\Phi_{MPS}^{(\chi,D)}$ is not guaranteed to have the same norm as
 928 the target Φ . While we can always rescale $\Phi_{MPS}^{(\chi,D)}$ to have unit norm, we must understand how this
 929 impacts the approximation error, something which can be done through inequalities which hold
 930 in any normed vector space. Given a target vector \mathbf{u} satisfying $\|\mathbf{u}\| = 1$, suppose there exists a
 931 vector \mathbf{v} which approximates \mathbf{u} to within distance $\|\mathbf{u} - \mathbf{v}\|$. The unit vector $\hat{\mathbf{v}} = \mathbf{v}/\|\mathbf{v}\|$ will then
 932 necessarily approximate \mathbf{u} to within distance

$$\begin{aligned} \|\mathbf{u} - \hat{\mathbf{v}}\| &= \|(\mathbf{u} - \mathbf{v}) + (\mathbf{v} - \hat{\mathbf{v}})\| \leq \|\mathbf{u} - \mathbf{v}\| + \|\mathbf{v} - \hat{\mathbf{v}}\| \\ &= \|\mathbf{u} - \mathbf{v}\| + |1 - \|\mathbf{v}\|| = \|\mathbf{u} - \mathbf{v}\| + |\|\mathbf{u}\| - \|\mathbf{v}\|| \\ &\leq 2\|\mathbf{u} - \mathbf{v}\|. \end{aligned}$$

933 **Bound $|\Phi_{MPS}^{(\chi,D)}|_{\mathcal{H}_C^k}$ as a function of k and \mathcal{F}** We utilize the fact that the spatial dependence of
 934 $\Phi_{MPS}^{(\chi,D)}$ is entirely mediated by the first D polynomial embedding functions from \mathcal{F} , which form an
 935 orthonormal basis over the finite-dimensional vector space of polynomials with degree less than
 936 D . This arrangement means that $\Phi_{MPS}^{(\chi,D)}$ has all partial derivatives of arbitrary order, which can be
 937 used to directly compute its Sobolev seminorm $|\Phi_{MPS}^{(\chi,D)}|_{\mathcal{H}_C^k}$ without invoking the notion of weak
 938 derivative. Given the simple rules for taking partial derivatives of multivariate polynomials, we
 939 can see that any single spatial derivative $\frac{\partial}{\partial x_i}$ will preserve the space of polynomials spanned by
 940 the D first embedding functions in \mathcal{F} , and consequently be equivalent to a $D \times D$ matrix acting on
 941 the i 'th mode of the discrete MPS $\psi_{MPS}^{(\chi,D)}$ underlying the continuous MPS $\Phi_{MPS}^{(\chi,D)}$. More generally,
 942 every partial derivative $\partial^{(i)}$ will be equivalently to a bounded linear operator $M^{(i)}$ acting on the
 943 vector space of N 'th order tensors $\mathbb{C}^{D \times \dots \times D} \simeq \mathbb{C}^{D^N}$ where $\psi_{MPS}^{(\chi,D)}$ lives.

944 With these details in place, the seminorm can be explicitly bounded as

$$\begin{aligned}
 |\Phi_{\text{MPS}}^{(\chi, D)}|_{\mathcal{H}_C^k} &= \sum_{|i|=k} \|\partial^{(i)} \Phi_{\text{MPS}}^{(\chi, D)}\|_{L_\mu^2} = \sum_{|i|=k} \|M^{(i)} \psi_{\text{MPS}}^{(\chi, D)}\| \\
 &\leq \sum_{|i|=k} |M^{(i)}| \cdot \|\psi_{\text{MPS}}^{(\chi, D)}\| = \sum_{|i|=k} |M^{(i)}|,
 \end{aligned} \tag{30}$$

945 where in the second equality we have invoked the orthonormality of the feature functions and
 946 the above-remarked equivalence between the action of $\partial^{(i)}$ on continuous MPS and a finite-
 947 dimensional matrix $M^{(i)}$ acting on the underlying discrete MPS $\psi_{\text{MPS}}^{(\chi, D)}$. The notation $|M^{(i)}|$ refers
 948 to the spectral norm (i.e. the largest singular value) of $M^{(i)}$, and the final equality uses the unit
 949 norm assumption $\|\Phi_{\text{MPS}}^{(\chi, D)}\|_{L_\mu^2} = \|\psi_{\text{MPS}}^{(\chi, D)}\| = 1$. Although the value of the spectral norms $|M^{(i)}|$
 950 will depend on the choice of basis functions \mathcal{F} , it is clear that the RHS of Eq. 30 is finite and
 951 depends on nothing else besides k , giving us the desired bound on $|\Phi_{\text{MPS}}^{(\chi, D)}|_{\mathcal{H}_C^k}$.

952 **Convert L_2 bound to infidelity bound** Summarizing our results up to this point, we have proved
 953 that for any unit-norm target function $\Phi \in \mathcal{H}_C^k$ and an orthonormal basis of polynomial feature
 954 functions \mathcal{F} , there exist quantities $\gamma'_1, \gamma'_2 > 0$ depending only on Φ , k , and \mathcal{F} for which there
 955 exist unit-norm continuous MPS $\Phi_{\text{MPS}}^{(\chi, D)}$ of arbitrary bond dimension χ and feature dimension D
 956 approximating Φ with L_2 error $\|\Phi - \Phi_{\text{MPS}}^{(\chi, D)}\|_{L_\mu^2} \leq \gamma'_1 \chi^{-\frac{k-1}{2}} + \gamma'_2 D^{-k}$. However Theorem 4 requires
 957 a bound on the infidelity $1 - |\langle \Phi, \Phi_{\text{MPS}}^{(\chi, D)} \rangle|$. This can be achieved by the straightforward inequality
 958 $1 - |\langle u, v \rangle| \leq \frac{1}{2} \|u - v\|^2$, which holds for any u and v in a normed vector space. Combining this
 959 with our L_2 bound gives

$$\begin{aligned}
 1 - |\langle \Phi, \Phi_{\text{MPS}}^{(\chi, D)} \rangle| &\leq \frac{1}{2} \|\Phi - \Phi_{\text{MPS}}^{(\chi, D)}\|_{L_\mu^2}^2 \\
 &\leq \frac{1}{2} \left(\gamma'_1 \chi^{-\frac{k-1}{2}} + \gamma'_2 D^{-k} \right)^2 \\
 &\leq \gamma_1 \chi^{-k+1} + \gamma_2 D^{-2k},
 \end{aligned} \tag{31}$$

960 where we use the identity $\gamma'_1 \gamma'_2 \chi^{-\frac{k-1}{2}} D^{-k} \leq \gamma'_1 \gamma'_2 (\chi^{-k+1} + D^{-2k})$ to simplify the cross-terms
 961 arising from the expansion of the square above to arrive at the constants $\gamma_1 := \frac{1}{2} \gamma_1'^2 + \gamma_1' \gamma_2'$ and
 962 $\gamma_2 := \frac{1}{2} \gamma_2'^2 + \gamma_1' \gamma_2'$. This gives us our desired infidelity bound, completing our proof of Theorem 4,
 963 and by extension Theorem 2.

964 As a final note, we consider the case where the domain of the target function Φ is unbounded
 965 (e.g. all of \mathbb{R}^N). Although the methods of [45] don't apply in this setting (for reasons related
 966 to certain functional analytic lemmas used in the proof of Lemma 2), we can instead consider
 967 a sequence of approximations of Φ by functions Φ_ϵ supported on boxes Ω_ϵ of increasing size,
 968 which each approximate Φ to within a distance of ϵ . By approximating this sequence of functions
 969 of bounded domain using Theorem 2, we can approximate our target Φ to arbitrary precision, albeit
 970 at the cost of introducing another ϵ -dependent term into the error bound of Eq. 15. Although this
 971 argument leaves some technical details to be worked out, it is clear that in practice this method
 972 offers a concrete means of using continuous-valued MPS as universal function approximators for
 973 functions on unbounded domains.

974 C.3 Proof of Theorem 3

975 **Theorem 3.** Consider a family of continuous-valued MPS with polynomial feature functions
 976 $\mathcal{F} = \{f_1, f_2, \dots\}$ forming an orthonormal basis for $[0, 1]$, which is defined on the hypercube

977 $\Omega = [0, 1]^N \subseteq \mathbb{R}^N$. Let $k \geq N$ and let $P : \Omega \rightarrow \mathbb{R}$ be any Probability Density Function (PDF)
 978 bounded below as $P_{\min} = \min_{\mathbf{x} \in \Omega} P(\mathbf{x}) > 0$, whose partial derivatives of order $1, 2, \dots, k$ all
 979 exist and are bounded. Then for every positive $\chi, D \in \mathbb{N}$ there exists a continuous-valued MPS of
 980 bond dimension χ and feature dimension D with unit norm, whose associated Born machine PDF
 981 $P_{\text{MPS}}^{(\chi, D)}(\mathbf{x}) = |\Phi_{\text{MPS}}^{(\chi, D)}|^2$ approximates P with Jensen-Shannon divergence

$$\text{JS}\left(P_{\text{MPS}}^{(\chi, D)}, P\right) \leq \eta_1 \chi^{-\frac{k-1}{2}} + \eta_2 D^{-k}, \quad (16)$$

982 where $\eta_1, \eta_2 > 0$ depend on the target PDF P , the assumed degree of smoothness k , and the
 983 feature functions \mathcal{F} .

984 The proof of Theorem 3 applies Theorem 2 to the artificial wavefunction $\Phi_P(\mathbf{x}) = \sqrt{P(\mathbf{x})}$,
 985 which requires first proving that (a) The partial derivatives of Φ_P of orders $1, 2, \dots, k$ (where
 986 $k \geq N$) all exist and are bounded. With this established, Theorem 2 gives us an approximating
 987 wavefunction $\Phi_{\text{MPS}}^{(\chi, D)}$ with a bounded infidelity relative to Φ_P , and we must (b) Convert the infi-
 988 delity bound between Φ_P and $\Phi_{\text{MPS}}^{(\chi, D)}$ into a bound on the Jensen-Shannon (JS) divergence between
 989 P and $P_{\text{MPS}}^{(\chi, D)}$. We tackle these issues in turn.

990 **Prove the partial derivatives of Φ_P of orders $1, 2, \dots, k$ (where $k \geq N$) all exist and are**
 991 **bounded.** We employ the multivariate version of Faà di Bruno's formula, which is a generaliza-
 992 tion of the standard chain rule to higher-order partial derivatives, stated here as

993 **Lemma 3** (Faà di Bruno). Consider a multivariate function $g : \Omega \rightarrow \mathbb{K}$ for $\Omega \subseteq \mathbb{R}^N$ whose partial
 994 derivatives $\partial^{(i)}g$ up to order k exist and are bounded, as well as a univariate function $h : \mathbb{K} \rightarrow \mathbb{K}$
 995 which is k -times differentiable within the range of g (i.e. $g(\Omega) \subseteq \mathbb{K}$). Then the partial derivative
 996 of the composite function $h \circ g : \mathbf{x} \mapsto h(g(\mathbf{x}))$ wrt the ℓ variables $\mathbf{i} = (x_{i_1}, x_{i_2}, \dots, x_{i_\ell})$ (with
 997 $\ell \leq k$) at a point $\mathbf{x} \in \Omega$ is

$$\partial^{(i)}h(g(\mathbf{x})) = \sum_{\pi \in \Pi} \frac{D^{|\pi|}h}{dy^{|\pi|}}(g(\mathbf{x})) \cdot \prod_{B \in \pi} \partial^{(B)}g(\mathbf{x}), \quad (32)$$

998 where (i) π runs through the set Π of all partitions of the set $\{i_1, i_2, \dots, i_\ell\}$; (ii) $B \in \pi$ denotes an
 999 iteration over all "blocks" of the partition π ; (iii) $\partial^{(B)}$ denotes the partial derivative with respect
 1000 to all of the variables x_i with $i \in B$; (iv) and $|\pi|$ indicates the number of blocks in the partition π .

1001 The details of Eq. 32 are of little interest to us, as we only use it to bound the partial derivatives
 1002 $\partial^{(i)}\Phi_P$. To this end, we first use the assumption $P(\mathbf{x}) \geq P_{\min}$ and chain rule for the square root
 1003 function $h(y) : y \mapsto \sqrt{y}$ to show that

$$\begin{aligned} \left| \frac{D^n h}{dy^n}(g(\mathbf{x})) \right| &= \left| \frac{(2n-3)(2n-5) \cdots (-1)}{2^n} g(\mathbf{x})^{-n+\frac{1}{2}} \right| \\ &\leq 2^n \left(\frac{1}{P_{\min}} \right)^{n-\frac{1}{2}} =: S_{\max}(n). \end{aligned} \quad (33)$$

1004 The fact that $S_{\max}(n)$ is an increasing function of n tells us that $S_{\max} := S_{\max}(k)$ is an upper
 1005 bound for every derivative of h up to order k . Denoting the largest partial derivative of P by
 1006 $G_P := \max_{|j| \leq k} \max_{\mathbf{x} \in \Omega} |\partial^{(j)}P(\mathbf{x})|$, which is finite by assumption (see Theorem 3), we can use
 1007 Eq. 32 to give the bound

$$\begin{aligned} \partial^{(i)}\Phi_P(\mathbf{x}) &= \partial^{(i)}h(P(\mathbf{x})) \leq \sum_{\pi \in \Pi} S_{\max} \prod_{B \in \pi} \partial^{(B)}P(\mathbf{x}) \\ &\leq \sum_{\pi \in \Pi} S_{\max} \prod_{B \in \pi} G_P \leq \sum_{\pi \in \Pi} S_{\max} (G_P)^k \\ &\leq S_{\max} (k G_P)^k, \end{aligned} \quad (34)$$

1008 which suffices to prove the existence and boundedness of the partial derivatives of Φ_P .

1009 **Convert the infidelity bound between Φ_P and $\Phi_{\text{MPS}}^{(\chi, D)}$ into a bound on the Jensen-Shannon**
 1010 **(JS) divergence between P and $P_{\text{MPS}}^{(\chi, D)}$.** We proceed in three steps, using the quantum trace dis-
 1011 tance and the classical total variation (TV) distance as intermediate quantities. The trace distance
 1012 $T(\Phi, \Phi')$ between pure quantum states Φ and Φ' takes the form of $T(\Phi, \Phi') := \sqrt{1 - |\langle \Phi, \Phi' \rangle|^2}$,
 1013 which can be expressed in terms of the infidelity $\mathcal{I}(\Phi, \Phi') = 1 - |\langle \Phi, \Phi' \rangle|^2$ as $T(\Phi, \Phi') = \sqrt{2\mathcal{I} + \mathcal{I}^2}$.
 1014 Thus, the infidelity bound of Eq. 31 gives us a bound on our quantum trace distance of interest.

1015 A well-known interpretation of the quantum trace distance between states Φ, Φ' is a bound
 1016 on the classical TV distance $TV(P, P_{\text{MPS}}^{(\chi, D)}) = \sup_{A \subseteq \Omega} |P(A) - P_{\text{MPS}}^{(\chi, D)}(A)|$ between any classical
 1017 distributions $P_\Phi, P_{\Phi'}$ which arise from von Neumann measurements of the corresponding quan-
 1018 tum states [62]. Given that the Born machine distributions are precisely those arising from von-
 1019 Neumann measurements of the underlying wavefunctions, we have $TV(P, P_{\text{MPS}}^{(\chi, D)}) \leq T(\Phi_P, \Phi_{\text{MPS}}^{(\chi, D)})$,
 1020 and thereby a bound on the TV distance,

$$\begin{aligned} TV(P, P_{\text{MPS}}^{(\chi, D)}) &= \sup_{A \subseteq \Omega} |P(A) - P_{\text{MPS}}^{(\chi, D)}(A)| \leq T(\Phi_P, \Phi_{\text{MPS}}^{(\chi, D)}) \\ &\leq \sqrt{\frac{3}{2}} \left(\gamma'_1 \chi^{-\frac{k-1}{2}} + \gamma'_2 D^{-k} \right). \end{aligned} \quad (35)$$

1021 Finally, the Jensen-Shannon divergence is known to be bounded by the TV distance, written
 1022 as $JS(P, Q) \leq \frac{\ln(2)}{2} TV(P, Q)$ which, combined with the above results, give

$$\begin{aligned} JS(P, P_{\text{MPS}}^{(\chi, D)}) &\leq \frac{\ln(2)}{2} TV(P, P_{\text{MPS}}^{(\chi, D)}) \leq \frac{\ln(2)}{2} T(\Phi_P, \Phi_{\text{MPS}}^{(\chi, D)}) \\ &= \frac{\ln(2)}{2} \sqrt{2\mathcal{I} + \mathcal{I}^2} \leq \frac{\sqrt{3} \ln(2)}{2} \sqrt{\mathcal{I}} \\ &\leq \sqrt{\frac{3}{8}} \ln(2) \left(\gamma'_1 \chi^{-\frac{k-1}{2}} + \gamma'_2 D^{-k} \right). \end{aligned} \quad (36)$$

1023 Taking $\eta_1 := \sqrt{\frac{3}{8}} \ln(2) \gamma'_1$ and $\eta_2 := \sqrt{\frac{3}{8}} \ln(2) \gamma'_2$ completes the proof of Theorem 3.

1024 D Detailed Methods

1025 D.1 Rotated Cube

1026 The cube was rotated by a random orthogonal transformation, and then scaled per-axis to the range
 1027 $[-1, 1]$ to standardize the range. This resulted in a linear transformation

$$M = \begin{bmatrix} 1.33 & 0.155 & 0.074 & 0.411 & 0.029 \\ -0.072 & 1.181 & 0.029 & 0.375 & -0.342 \\ 0.306 & 0.303 & 0.862 & -0.226 & 0.302 \\ -0.363 & 0.217 & -0.297 & 0.998 & 0.125 \\ 0.024 & 0.229 & 0.358 & 0.514 & 0.875 \end{bmatrix}$$

1028 which acted on the set $[-\frac{1}{2}, \frac{1}{2}]^5$. The simple form allowed use to compute the exact entropy as
 1029 $\log(\det(M)) = -0.4246$.

1030 The training set was 80k sampled points. No minibatching was used. Eighteen sweeps of
 1031 DMRG were performed. At each site, 4 steps of gradient descent were performed, each with a

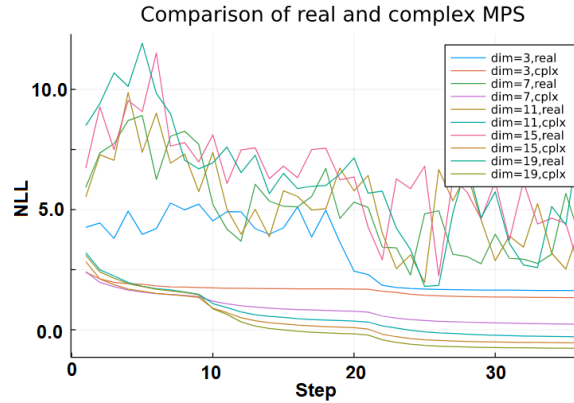


Figure 13: A comparison of real-entry and complex-entry matrix product state performance on the rotated cube dataset. Bond dimension χ and feature dimension D were equal, and tried at 5 different values from 3 up to 19, with both real and complex entries. The bond dimension was initially lower and increased at training epochs 9 and 22, leading to kinks in the loss curve. Complex-entry MPS trained smoothly, while real-entry MPS did not, due to the sharp truncation of the SVD.

1032 learning rate of **0.05**. The maximum bond dimension in the first sweep was $\max(\chi_{\max}/2, 5)$, and
 1033 increased in subsequent sweeps linearly up to χ_{\max} where it stayed for the last five sweeps.

1034 In Fig. 13 we present a comparison of the training performance for MPSs with real and com-
 1035 plex entries on this specific data set.

1036 D.2 Two Moons

1037 For a given noise parameter $\sigma \ll 1$, the entropy of the two moons dataset can be approximated as

$$S \approx \frac{3 \ln(2\pi) + 1}{2} + \ln(\sigma) + \frac{1.81}{\pi} \sigma. \quad (37)$$

1038 This approximation can be understood as $\ln(2)$ for choosing a curve to lie on, $\log(\pi)$ for a uniform
 1039 distribution on a curve of length π , and $\log(\sigma \sqrt{2\pi e})$ for a radial uncertainty σ . The final $\frac{1.81}{\pi} \sigma$
 1040 accounts for extending the curve of length π at the tips by a blur of σ , where

$$1.81 \approx \int_{-\infty}^{\infty} -\sqrt{2}(1 + \operatorname{erf}(x)) \log\left(\frac{1 + \operatorname{erf}(x)}{2}\right) dx. \quad (38)$$

1041 For our experimental results, we used a value of $\sigma = 0.1$, for which $S \approx 0.96$.

1042 To use the Fourier basis, we first scaled the x and y values to the range $[-0.9, 0.9]$. This
 1043 rescaling adds a small constant factor to the NLL, but this was corrected for when comparing
 1044 to the true entropy of the distribution. We used a training set of 10k sampled points. The KL
 1045 divergence as a function of χ and D is presented in Fig. 14, where a minimum value of 0.022 was
 1046 reached. It is apparent that for this dataset, the bond dimension quickly saturated its usefulness
 1047 past $\chi = 4 \sim 5$, with the largest improvement coming from increasing D .

1048 D.3 Iris

1049 We used the Iris dataset available in the UCI Machine Learning Repository [59], consisting of 150
 1050 points with four continuous features describing petal shapes of different Iris flowers, supplemented
 1051 with a categorical feature describing which of three varieties the flower belongs to. The Iris dataset
 1052 was normalized by rescaling each feature to lie in the range $[-1, 1]$, before applying a feature map
 1053 to each. The NLL loss for different bond and feature dimensions is shown in Fig. 15.

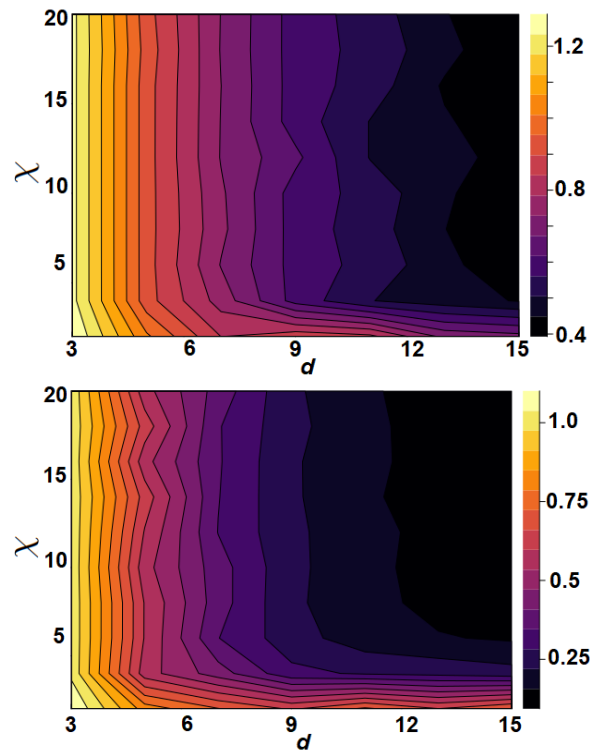


Figure 14: Excess loss (NLL minus distribution entropy) on Two Moons with $\sigma = 0.1$, trained with different embedding dimensions and bond dimensions. Upper plot shows a Hermite embedding. Lower plot shows a Fourier embedding.

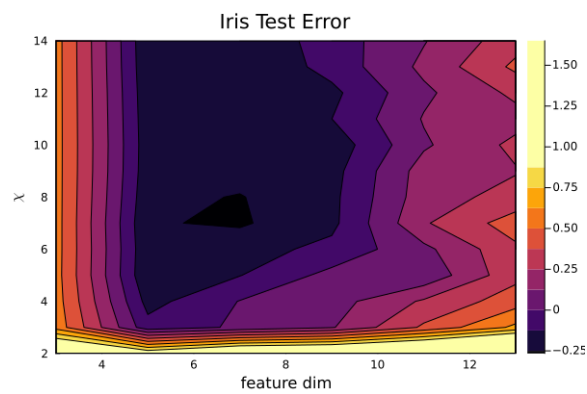


Figure 15: NLL loss on Iris dataset at different bond dimensions and feature dimensions. Each point is the mean of 5 values from a 5-fold cross validation.

1054 D.4 XY Model

1055 The continuous-valued MPS model was trained to minimize the NLL loss on a dataset of samples
 1056 drawn from the finite temperature XY model generated using Markov Chain Monte Carlo [63].
 1057 The conversion of this score to a KL divergence was made using the test of [64] with 100k samples
 1058 and a $k = 10$ neighborhood. To further examine the model's behavior, we also measured the KL
 1059 divergence between the model and true marginal distributions of the angle-invariant quantities
 1060 $C_{\text{neigh}} = \cos(\mathbf{x}_{1,1} - \mathbf{x}_{1,2})$ and $C_{\text{corn}} = \cos(\mathbf{x}_{1,1} - \mathbf{x}_{4,4})$, which measure the correlations between
 1061 neighbors and opposite corners, respectively. These pairwise correlations were both learned very
 1062 well, with a KL divergence of 0.0027 for corner-to-corner correlations (which are the hardest for
 1063 the linear MPS to learn), and even lower values for closer pairs of sites.

1064 D.5 Compressible Data

1065 The deliberately compressible data for Sec. 7.5 was a 4-feature dataset generated by the following
 1066 formulas from four samples \mathbf{x}_i from the uniform distribution on $[0, 1]$:

$$\begin{aligned} y_1 &= -1 + [0.6 + 2.2x_1] \\ y_3 &= x_3 \\ y_4 &= -\frac{1}{2} + [1.4x_4] \\ y_2 &= \frac{y_1 + 2x_2 + y_3 + y_4}{4}, \end{aligned}$$

1067 where $[x]$ denotes the largest integer k such that $k \leq x$. This produces a dataset where each
 1068 feature has a very different marginal (implying that each feature would make best use of a different
 1069 compression map), the features y_1 and y_4 are discrete with only three or two values respectively,
 1070 and y_2 is correlated with the other three (so that the MPS correlation structure is not trivial). The
 1071 single-site marginal distribution is shown in Fig. 11.

1072 E Dynamic Basis Training

1073 The following pseudocode details in more detail the process for optimizing the $D \times d$ isometric
 1074 compression matrices $\{U_i\}_{i=1}^N$ using a dataset of samples $\mathcal{D} = \{\mathbf{x}^{(j)}\}_{j=1}^T$, where each sample has
 1075 features $\mathbf{x}^{(j)} = (x_1^{(j)}, x_2^{(j)}, \dots, x_N^{(j)})$. In practice the steps in this process will be interspersed with
 1076 DMRG updates, in order to benefit from caching of intermediate hidden states.

Algorithm 1 Dynamic Basis Adjustment

```

 $\epsilon \leftarrow 0.5$ 
for  $i \leftarrow 1 \dots N$  do
  for  $j \leftarrow 1 \dots T$  do
     $\mathbf{u}_{i,j} \leftarrow \zeta(\mathbf{x}_i^{(j)})$ 
     $\mathbf{v}_{i,j} \leftarrow \text{MPSCContract}(\mathbf{x}_{(j)}, \neg i)$ 
     $\mathbf{c}_j \leftarrow \mathbf{u}_{i,j}^\dagger U_i \mathbf{v}_{i,j}$ 
     $p_j \leftarrow |\mathbf{c}_j|$ 
     $\phi_j \leftarrow \mathbf{c}_j / p_j$ 
  end for
   $B \leftarrow \sum_{j=1}^T (p_j^\epsilon \phi_j)^{-1} \mathbf{u}_{i,j} \mathbf{v}_{i,j}^\dagger$ 
   $B_U, B_S, B_V = \text{SVD}(B)$ 
   $U_i \leftarrow B_U B_V^T$ 
end for

```

- ▶ Controls stability
- ▶ Loop over each site
- ▶ Loop over each sample in batch
- ▶ D -dim embedding
- ▶ d -dim embedding
- ▶ Loss probability
- ▶ Current phase
- ▶ $D \times d$ matrix
- ▶ Update isometric matrix U_i

Increase ϵ towards 1. If loop is unstable, decrease ϵ towards 0.
