Fabio A. González^{1*}, Alejandro Gallego¹, Santiago Toledo-Cortés¹ and Vladimir Vargas-Calderón²

^{1*}MindLab, Depto. de Ing. de Sistemas e Industrial, Universidad Nacional de Colombia, Bogotá, DC, Colombia.
²Grupo de Superconductividad y Nanotecnología, Depto. de Física, Universidad Nacional de Colombia, Bogotá, DC, Colombia.

*Corresponding author(s). E-mail(s): fagonzalezo@unal.edu.co; Contributing authors: jagallegom@unal.edu.co; stoledoc@unal.edu.co; vvargasc@unal.edu.co;

Abstract

A density matrix describes the statistical state of a quantum system. It is a powerful formalism to represent both the quantum and classical uncertainty of quantum systems and to express different statistical operations such as measurement, system combination and expectations as linear algebra operations. This paper explores how density matrices can be used as a building block for machine learning models exploiting their ability to straightforwardly combine linear algebra and probability. One of the main results of the paper is to show that density matrices coupled with random Fourier features could approximate arbitrary probability distributions over \mathbb{R}^n . Based on this finding the paper builds different models for density estimation, classification and regression. These models are differentiable, so it is possible to integrate them with other differentiable components, such as deep learning architectures and to learn their parameters using gradient-based optimization. In addition, the paper presents optimization-less training strategies based on estimation and model averaging. The models are evaluated in benchmark tasks and the results are reported and discussed.

 $\label{eq:Keywords: quantum machine learning, density matrix, density estimation, classification, regression$

1 Introduction

The formalism of density operators and density matrices was developed by von Neumann as a foundation of quantum statistical mechanics (Von Neumann, 1927). From the point of view of machine learning, density matrices have an interesting feature: the fact that they combine linear algebra and probability, two of the pillars of machine learning, in a very particular but powerful way.

The main question addressed by this work is how density matrices can be used in machine learning models. One of the main approaches to machine learning is to address the problem of learning as one of estimating a probability distribution from data: joint probabilities P(x, y) in generative supervised models or conditional probabilities P(y|x) in discriminative models.

The central idea of this work is to use density matrices to represent these probability distributions tackling the important question of how to encode arbitrary probability density functions in \mathbb{R}^n into density matrices.

The quantum probabilistic formalism of von Neumann is based on linear algebra, in contrast with classical probability which is based on set theory. In the quantum formalism the sample space corresponds to a Hilbert space \mathcal{H} and the event space to a set of linear operators in \mathcal{H} , the density operators (Wilce, 2021).

The quantum formalism generalizes classical probability. A density matrix in an *n*-dimensional Hilbert space can be seen as a catalog of categorical distributions on the finite set $\{1 \dots n\}$. A direct application of this fact is not very useful as we want to efficiently model continuous probability distributions in \mathbb{R}^n . One of the main results of this paper is to show that it is possible to model arbitrary probability distributions in \mathbb{R}^n using density matrices of finite dimension in conjunction with random Fourier features (Rahimi and Recht, 2007). In particular the paper presents a method for non-parametric density estimation that combines density matrices and random Fourier features to efficiently learn a probability density function from data and to efficiently predict the density of new samples.

The fact that the probability density function is represented in matrix form and that the density of a sample is calculated by linear algebra operations makes it easy to implement the model in GPU-accelerated machine learning frameworks. This also facilitates using density matrices as a building block for classification and regression models, which can be trained using gradientbased optimization and can be easily integrated with conventional deep neural networks. The paper presents examples of these models and shows how they can be trained using gradient-based optimization as well as optimization-less learning based on estimation.

The paper is organized as follows: Section 2 covers the background on kernel density estimation, random features, and density matrices; Section 5 presents four different methods for density estimation, classification and regression; Section 6 discusses some relevant works; Section 7 presents the experimental evaluation; finally, Section 8 discusses the conclusions of the work.

2 Background and preliminaries

2.1 Kernel density estimation

Kernel Density Estimation (KDE) (Rosenblatt, 1956; Parzen, 1962), also known as Parzen-Rossenblat window method, is a non-parametric density estimation method. This method does not make any particular assumption about the underlying probability density function. Given an iid set of samples $X = \{x_1, \ldots, x_N\}$, the smooth Parzen's window estimate has the form

$$\hat{f}_{\lambda}(x) = \frac{1}{NM_{\lambda}} \sum_{i=1}^{N} k_{\lambda}(x, x_i), \qquad (1)$$

where $k_{\lambda}(\cdot)$ is a kernel function, λ is the smoothing bandwidth parameter of the estimate and M_{λ} is a normalizing constant. A small λ -parameter implies a small grade of smoothing.

Rosenblatt (1956) and Parzen (1962) showed that eq. (1) is an unbiased estimator of the pdf f. If k_{γ} is the Gaussian kernel, eq. (1) takes the form

$$\hat{f}_{\gamma}(x) = \frac{1}{NM_{\gamma}} \sum_{i=1}^{N} e^{-\gamma ||x_i - x||^2},$$
(2)

where $M_{\gamma} = (\pi/\gamma)^{\frac{d}{2}}$.

KDE has several applications: to estimate the underlying probability density function, to estimate confidence intervals and confidence bands (Efron, 1992; Chernozhukov et al, 2014), to find local modes for geometric feature estimation (Chazal et al, 2017; Chen et al, 2016), to estimate ridge of the density function (Genovese et al, 2014), to build cluster trees (Balakrishnan et al, 2013), to estimate the cumulative distribution function (Nadaraya, 1964), to estimate receiver operating characteristic (ROC) curves (McNeil and Hanley, 1984), among others.

One of the main drawbacks of KDE is that it is a memory-based method, i.e. it requires the whole training set to do a prediction, which is linear on the training set size. This drawback is typically alleviated by methods that use data structures that support efficient nearest-neighbor queries. This approach still requires to store the whole training dataset.

2.2 Random features

Random Fourier features (RFF) (Rahimi and Recht, 2007) is a method that builds an embedding $\phi_{\text{rff}} : \mathbb{R}^d \to \mathbb{R}^D$ given a shift-invariant kernel $k : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ such that $\forall x, y \in \mathbb{R}^d$, $k(x, y) \approx \langle \phi_{\text{rff}}(x), \phi_{\text{rff}}(y) \rangle = \phi_{\text{rff}}^T(x)\phi_{\text{rff}}(y)$. One of the main applications of RFF is to speedup kernel methods, being data independence one of its advantages.

The RFF method is based on the Bochner's theorem. In layman's terms, Bochner's theorem shows that a shift invariant positive-definite kernel $k(\cdot)$ is the Fourier transform of a probability measure p(w). Rahimi and Recht

3

(2007) use this result to approximate the kernel function by designing a sample procedure that estimates the integral of the Fourier transform. The first step is to draw D iid samples $\{w_1, \ldots, w_D\}$ from p and D iid samples $\{b_1, \ldots, b_D\}$ from a uniform distribution in $[0, 2\pi]$. Then, define:

$$\phi_{\text{rff}} : \mathbb{R}^d \to \mathbb{R}^D$$

$$x \mapsto \sqrt{\frac{2}{D}} \left(\cos\left(w_1^T x + b_1\right), \dots, \cos\left(w_D^T x + b_D\right) \right).$$
(3)

Rahimi and Recht (2007) showed that the expected value of $\phi_{\text{rff}}^T(x)\phi_{\text{rff}}(y)$ uniformly converges to k(x, y):

Theorem 1 (*Rahimi and Recht, 2007*) Let \mathcal{M} be a compact subset of \mathbb{R}^d with a diameter diam(\mathcal{M}). Then for the mapping ϕ_{rff} defined above, we have

$$\Pr\left[\sup_{x,y\in\mathcal{M}} |\phi_{\mathrm{rff}}^{T}(x)\phi_{\mathrm{rff}}(y) - k(x,y)| \ge \epsilon\right] \le 2^{8} \left(\frac{\sigma_{p}\mathrm{diam}(\mathcal{M})}{\epsilon}\right)^{2} \exp\left(-\frac{D\epsilon^{2}}{4(d+2)}\right), \tag{4}$$

where, σ_p^2 is the second momentum of the Fourier transform of k. In particular, for the Gaussian kernel $\sigma_p^2 = 2d\gamma$, where γ is the spread parameter (see Eq. 2).

Different approaches to compute random features for kernel approximation have been proposed based on different strategies: Monte Carlo sampling (Le et al, 2013; Yu et al, 2016), quasi-Monte-Carlo sampling (Avron et al, 2016; Shen et al, 2017), and quadrature rules (Dao et al, 2017).

RFF may be used to formulate a non-memory based version of KDE. For the Gaussian kernel we have:

$$\hat{f}_{\gamma}(x) = \frac{1}{NM_{\gamma}} \sum_{i=1}^{N} k_{\gamma}(x_i, x)$$

$$\approx \frac{1}{NM_{\gamma}} \sum_{i=1}^{N} \langle \phi_{\rm rff}(x_i), \phi_{\rm rff}(x) \rangle$$

$$= \frac{1}{M_{\gamma}} \left\langle \frac{1}{N} \sum_{i=1}^{N} \phi_{\rm rff}(x_i), \phi_{\rm rff}(x) \right\rangle$$

$$= \frac{1}{M_{\gamma}} \langle \Phi_{\rm train}, \phi_{\rm rff}(x) \rangle$$

$$= \frac{1}{M_{\gamma}} \Phi_{\rm train}^{T} \phi_{\rm rff}(x) \qquad (5)$$

 Φ_{train} in eq. (5) can be efficiently calculated during training time, since is just an average of the RFF embeddings of the training samples. The time complexity of prediction, eq. (5), is constant on the size of the training dataset. The price of this efficiency improvement is a loss in precision, since we are using an approximation of the Gaussian kernel.

3 Density estimation with density matrices

The Gaussian kernel satisfy $\forall x, y \in \mathbb{R}^d$, $k_{\gamma}(x, y) > 0$, however the RFF estimation may be negative. To alleviate this we could estimate the square of the kernel and use the fact that $k_{\gamma}(x, y) = k_{\gamma/2}^2(x, y)$. In this case we have:

$$\hat{f}_{\gamma}(x) = \frac{1}{NM_{\gamma}} \sum_{i=1}^{N} k_{\gamma}(x_{i}, x)$$

$$= \frac{1}{NM_{\gamma}} \sum_{i=1}^{N} k_{\gamma/2}^{2}(x_{i}, x)$$

$$\approx \frac{1}{NM_{\gamma}} \sum_{i=1}^{N} \langle \phi_{\text{rff}}(x_{i}), \phi_{\text{rff}}(x) \rangle^{2}$$

$$= \frac{1}{NM_{\gamma}} \sum_{i=1}^{N} \langle \phi_{\text{rff}}(x), \phi_{\text{rff}}(x_{i}) \rangle \langle \phi_{\text{rff}}(x_{i}), \phi_{\text{rff}}(x) \rangle$$

$$= \frac{1}{NM_{\gamma}} \sum_{i=1}^{N} \phi_{\text{rff}}^{T}(x) \phi_{\text{rff}}(x_{i}) \phi_{\text{rff}}^{T}(x_{i}) \phi_{\text{rff}}(x)$$

$$= \frac{1}{M_{\gamma}} \phi_{\text{rff}}^{T}(x) \left(\frac{1}{N} \sum_{i=1}^{N} \phi_{\text{rff}}(x_{i}) \phi_{\text{rff}}^{T}(x_{i}) \right) \phi_{\text{rff}}(x)$$

$$= \frac{1}{M_{\gamma}} \phi_{\text{rff}}^{T}(x) \rho_{\text{train}} \phi_{\text{rff}}(x) =: \hat{f}_{\rho_{\text{train}}}(x)$$
(6)

In eq. (6) it is important to take into account that the parameters of the RFF embedding, ϕ_{rff} , are sampled using a parameter $\gamma/2$ for the Gaussian kernel.

The following proposition shows that $\hat{f}_{\rho_{\text{train}}}$, as defined in eq. (6), uniformly converges to the Gaussian kernel Parzen's estimator \hat{f}_{γ} (eq. (2)).

Proposition 2 Let \mathcal{M} be a compact subset of \mathbb{R}^d with a diameter diam (\mathcal{M}) , let $X = \{x_i\}_{i=1...N} \subset \mathcal{M}$ a set of iid samples, then $\hat{f}_{\rho_{\text{train}}}$ (eq. (6)) and \hat{f}_{γ} satisfy:

$$\Pr\left[\sup_{x\in\mathcal{M}}|\hat{f}_{\rho_{\mathrm{train}}}(x)-\hat{f}_{\gamma}(x)|\geq\epsilon\right]\leq$$

$$2^{8} \left(\frac{\sqrt{2d\gamma}\operatorname{diam}(\mathcal{M})}{3M_{\gamma}\epsilon}\right)^{2} \exp\left(-\frac{D(3M_{\gamma}\epsilon)^{2}}{4(d+2)}\right)$$
(7)

Proof (see Apendix A)

The Parzen's estimator is an unbiased estimator of the true density function from which the samples were generated and Proposition 2 shows that $\hat{f}_{\rho_{\text{train}}}(x)$ can approximate this estimator.

A further improvement to the $\hat{f}_{\rho_{\text{train}}}(x)$ estimator is to normalize the RFF embedding as follows:

$$\left|\bar{\phi}_{\rm rff}(x)\right\rangle = \frac{\phi_{\rm rff}(x)}{\left\|\phi_{\rm rff}(x)\right\|}\tag{8}$$

Π

Here we use the Dirac notation to emphasize the fact that $\bar{\phi}_{\rm rff}$ is a quantum feature map. This has the effect that the estimation $k_{\gamma}(x,x) = \langle \bar{\phi}_{\rm rff}(x) | \bar{\phi}_{\rm rff}(x) \rangle = 1$ will be exact and $\forall x, y \in \mathbb{R}^d, \langle \bar{\phi}_{\rm rff}(x) | \bar{\phi}_{\rm rff}(y) \rangle \leq 1$.

During the training phase ρ_{train} is estimated as the average of the cross product of the normalized RFF embeddings of the training samples:

$$\rho_{\text{train}} = \frac{1}{N} \sum_{i=1}^{N} \left| \bar{\phi}_{\text{rff}}(x_i) \right\rangle \left\langle \bar{\phi}_{\text{rff}}(x_i) \right| \tag{9}$$

The time complexity of calculating ρ_{train} is $O(D^2N)$, i.e. linear on the size of the training dataset. One advantage over conventional KDE is that we do not need to store the whole training dataset, but a more compact representation.

During the prediction phase the density of a new sample is calculated as:

$$\hat{f}_{\rho_{\rm train}}(x) = \frac{1}{M_{\gamma}} \left\langle \bar{\phi}_{\rm rff}(x) \right| \rho_{\rm train} \left| \bar{\phi}_{\rm rff}(x) \right\rangle \tag{10}$$

The $\hat{f}_{\rho_{\text{train}}}$ estimator has an important advantage over the Parzen's estimator, its computational complexity. The time to calculate the Parzen's estimator (eq. (2)) is O(dN) while the time to estimate the density based on the density matrix ρ_{train} (eq. (10)) is $O(D^2)$, which is constant on the size of the training dataset.

The ρ_{train} matrix in eq. (9) is a well known mathematical object in quantum mechanics, a density matrix, and eq. (10) is an instance of the Born rule which calculates the probability that a measurement of a quantum system produces a particular result. This connection and the basic ideas behind density matrices are discussed in the next section.

4 Density matrices

This section introduces some basic mathematical concepts that are part of the mathematical framework that supports quantum mechanics and discusses their connection with the ideas introduced in the previous subsection. The contents of this section are not necessary for understanding the rest of the paper and are included to better explain the connection of the ideas presented in this paper with the quantum mechanics mathematical framework.

The state of a quantum system is represented by a vector $\psi \in \mathcal{H}$, where \mathcal{H} is the Hilbert space of the possible states of the system. Usually¹ $\mathcal{H} = \mathbb{C}^d$.

As an example, consider a system that could be in two possible states, e.g. the spin of an electron that could be up (\uparrow) or down (\downarrow) with respect to some axis z. The state of this system is, in general, represented by a regular column vector $|\psi\rangle = (\alpha, \beta)$, with $|\alpha|^2 + |\beta|^2 = 1$. This state represents a system that is in a superposition of the two basis states $|\psi\rangle = \alpha \uparrow +\beta \downarrow$. The outcome of a measurement of this system, along the z axis, is determined by the Born rule: the spin is up with probability $|\alpha|^2$ and down with probability $|\beta|^2$. Notice that α and β could be negative or complex numbers, but the Born rule guarantees that we get valid probabilities.

The normalized RFF mapping (eq. (8)) can be seen as a function that maps a sample to the state of a quantum system. In quantum machine learning literature, there are different approaches to encode data in quantum states (Schuld, 2018). The use of RFF as a data quantum encoding strategy was first proposed by (González et al, 2020; González et al, 2021).

The probabilities that arise from the superposition of states in the previous example is a manifestation of the uncertainty that is inherent to the nature of quantum physical systems. We call this kind of uncertainty quantum uncertainty. Other kind of uncertainty comes, for instance, from errors in the measurement or state-preparation processes, we call this uncertainty classical uncertainty. A density matrix is a formalism that allows us to represent both types of uncertainty. To illustrate it, let's go back to our previous example. The density matrix representing the state ψ is:

$$\rho = |\psi\rangle \langle \psi| = \begin{bmatrix} |\alpha|^2 & \alpha\beta^* \\ \beta\alpha^* & |\beta|^2 \end{bmatrix},$$
(11)

As a concrete example, consider $\langle \psi_1 | = \left(\frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}}\right)$ the corresponding density matrix is:

$$\rho_1 = |\psi_1\rangle \langle \psi_1| = \begin{bmatrix} \frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} \end{bmatrix},$$
(12)

which represents a superposition state where we have a $\frac{1}{2}$ probability of measuring any of the two states. Notice that the probabilities for each state are in the diagonal of the density matrix. ρ_1 is a rank-1 density matrix, and this

¹In this paper we use $\mathcal{H} = \mathbb{R}^d$, but most of the methods and results can be extended to the complex case.

means that it represents a pure state. A mixed state, i.e. a state with classical uncertainty, is represented by a density matrix with the form:

$$\rho = \sum_{i=1}^{N} p_i |\psi_i\rangle \langle\psi_i|, \qquad (13)$$

where $p_i > 0 \in \mathbb{R}$, $\sum_{i=1}^{N} p_i = 1$, and $\{\psi_i\}_{i=1...N}$ are the states of a an ensemble of N quantum systems, where each system has an associated probability p_i . The matrix ρ_{train} in eq. (9) is in fact a density matrix that represents the state of an ensemble of quantum systems where each system corresponds to a training data sample. The probability is the same for all the N elements of the ensemble, $\frac{1}{N}$.

As a concrete example of a mixed state consider two pure states $\psi_2 = (1,0)$ and $\psi'_2 = (0,1)$, and consider a system that is prepared in state ψ_2 with probability $\frac{1}{2}$ and in state ψ'_2 with probability $\frac{1}{2}$ as well. The state of this system is represented by the following density matrix:

$$\rho_2 = \frac{1}{2} |\psi_2\rangle \langle\psi_2| + \frac{1}{2} |\psi_2'\rangle \langle\psi_2'| = \begin{bmatrix} \frac{1}{2} & 0\\ 0 & \frac{1}{2} \end{bmatrix}, \tag{14}$$

At first sight, states ρ_1 and ρ_2 may be seen as representing the same quantum system, one where the probability of measuring an up state (or down state) in the z axis is $\frac{1}{2}$. However they are different systems, ρ_1 represents a system with only quantum uncertainty, while ρ_2 corresponds a system with classical uncertainty. To better observe the differences of the two systems we have to perform a measurement along a particular axis. To do so, we use the following version of the Born rule for density matrices:

$$P(\varphi|\rho) = \operatorname{Tr}(\rho |\varphi\rangle \langle \varphi|) = \langle \varphi| \rho |\varphi\rangle$$
(15)

which calculates the probability of measuring the state φ in a system in state ρ . If we set $\varphi = \left(\frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}}\right)$ we get $P(\varphi|\rho_1) = 1$ and $P(\varphi|\rho_2) = \frac{1}{2}$, showing that in fact both systems are different.

5 Methods

5.1 Density matrix kernel density estimation (DMKDE)

In this subsection we present a model for non-parametric density estimation based on the ideas discussed in subsection 3. The model receives an input $x \in \mathbb{R}^d$, represents it using a RFF quantum feature map (eq. (3)) and estimates the density of it using eq. (10). The model can be trained by averaging the density matrices corresponding to the training samples or by using stochastic gradient descent. The second approach requires a re-parametrization of the model that we discuss next. The main parameter of the model is ρ_{train} , which is a Hermitian matrix. To ensure this property, we can represent it using a factorization as follows:

$$\rho_{\text{train}} = V^T \Lambda V, \tag{16}$$

where $V \in \mathbb{R}^{r \times D}$, $\Lambda \in \mathbb{R}^{r \times r}$ is a diagonal matrix and r < D is the reduced rank of the factorization. With this new representation, eq. (10) can be re-expressed as:

$$\hat{f}_{\rho_{\text{train}}}(x) = \frac{1}{M_{\gamma}} \|\Lambda^{\frac{1}{2}} V \bar{\phi}_{\text{rff}}(x)\|^2.$$
(17)

This reduces the time to calculate the density of a new sample to O(Dr).



Fig. 1 Density matrix kernel density estimation (DMKDE).

The model is depicted in Fig. 1 and its function is governed by the following equations:

$$z := \phi_{\rm rff}(x) = \cos(W_{\rm rff}x + b_{\rm rff}), \qquad (18a)$$

$$z' := \frac{z}{\|z\|},\tag{18b}$$

$$\tilde{y} := \frac{1}{M_{\gamma}} \|\Lambda^{\frac{1}{2}} V z'\|^2 \tag{18c}$$

The hyperparameters of the model are the dimension of the RFF representation D, the spread parameter γ of the Gaussian kernel and the rank r of the density matrix factorization. The parameters are the weights and biases of the RFF, $W_{\text{rff}} \in \mathbb{R}^{D \times d}$ and $b_{\text{rff}} \in \mathbb{R}^d$ (corresponding to the w_i and b_i parameters in Eq. 3), and the components of the factorization, $V \in \mathbb{R}^{r \times D}$ and $\lambda \in \mathbb{R}^r$, the vector with the elements in the diagonal of Λ .

The training process of the model is as follows:

- 1. Input. A sample set $X = \{x_i\}_{i=1...N}$ with $x_i \in \mathbb{R}^d$, parameters $\gamma \in \mathbb{R}^+$, $D \in \mathbb{N}$
- 2. Calculate $W_{\text{rff}} = [w_1 \dots w_D]$ and $b_{\text{rff}} = [b_1 \dots b_D]$ using the random Fourier features method described in Section 2.2 for approximating a Gaussian kernel with parameters $\gamma/2$ and D.
- 3. Apply $\overline{\phi}_{\text{rff}}$ (eq. (8)):

$$z_i = \bar{\phi}_{\rm rff}(x_i). \tag{19}$$

4. Estimate ρ_{train} :

$$\rho_{\text{train}} = \frac{1}{N} \sum_{i=1}^{N} z_i z_i^T, \qquad (20)$$

5. Make a spectral decomposition of rank r of ρ_{train} :

$$\rho_{\text{train}} = V^T \Lambda V.$$

Notice that this training procedure does not require any kind of iterative optimization. The training samples are only used once and the time complexity of the algorithm is linear on the number of training samples. The complexity of step 4 is $O(D^2N)$ and of step 5 is $O(D^3)$.

Since the operations defined in eq. (18) are differentiable, it is possible to use gradient-descent to minimize an appropriate loss function. For instance, we can minimize the negative log-likelihood:

$$L = -\sum_{i=1}^{K} \log(\tilde{y}) \tag{21}$$

In contrast with the learning procedure based on density matrix estimation, using SGD does not guarantee that we will approximate the real density function. If we train all the parameters, maximizing the likelihood becomes an ill-posed problem because of singularities (a Gaussian with arbitrary small variance centered in one training point) (Bishop, 2006). Keeping fixed the RFF parameters and optimizing the parameters of the density matrix, V and λ has shown a good experimental performance. The version of the model trained with gradient descent is called DMKDE-SGD.

Something interesting to notice is that the process described by eqs. (19) and (20) generalizes density estimation for variables with a categorical distribution, i.e. $x \in \{1, \ldots, K\}$. To see this, we replace $\bar{\phi}_{\text{rff}}$ in eq. (19) by the well-known one-hot-encoding feature map:

$$\begin{aligned}
\phi_{\text{ohe}} &: D \to \mathbb{R}^K \\
& i \mapsto E_i,
\end{aligned}$$
(22)

where E_i is the unit vector with a 1 in position *i* and 0 in the other positions. It is not difficult to see that in this case

$$\rho_{ii} = \Pr(x = i) = \frac{|\{x_j | j \in \{1, \dots, N\}, x_j = i\}|}{N}.$$
(23)

5.2 Density matrix kernel density classification (DMKDC)

The extension of kernel density estimation to classification is called kernel density classification (Hastie et al, 2009). The posterior probability is calculated as

$$\hat{\Pr}(Y=j|X=x) = \frac{\pi_j f_j(x)}{\sum_{k=1}^K \pi_k \hat{f}_k(x)},$$
(24)

where π_j and \hat{f}_j are respectively the class prior and the density estimator of class j.



Fig. 2 Density matrix kernel density classification (DMKDC).

We follow this approach to define a classification model that uses the density estimation strategy based on RFF and density matrices described in the previous section. The input to the model is a vector $x \in \mathbb{R}^d$. The model is depicted in Fig. 2 and defined by the following equations:

$$z := \phi_{\rm rff}(x) = \cos(W_{\rm rff}x + b_{\rm rff}), \tag{25a}$$

Springer Nature 2021 LATEX template

12 Learning with Density Matrices and Random Features

$$z' := \frac{z}{\|z\|},\tag{25b}$$

$$\tilde{y}_i := \|\Lambda_i^{\frac{1}{2}} V_i z'\|^2 \quad \forall i = 1 \dots K,$$
(25c)

$$\tilde{y}'_i := \frac{\pi_i y_i}{\sum_{j=i}^K \tilde{y}_j} \quad \forall i = 1 \dots K,$$
(25d)

The hyperparameters of the model are the dimension of the RFF representation D, the spread parameter γ of the Gaussian kernel, the class priors π_i and the rank r of the density matrix factorization. The parameters are the weights and biases of the RFF, $W_{\text{rff}} \in \mathbb{R}^{D \times d}$ and $b_{\text{rff}} \in \mathbb{R}^d$, and the components of the factorization, $V_i \in \mathbb{R}^{r \times D}$ and $\lambda_i \in \mathbb{R}$ for $i = 1 \dots K$.

The model can be trained using two different strategies: one, using DMKDE to estimate the density matrices of each class; two, use stochastic gradient descent over the parameters to minimize an appropriate loss function.

The training process based on density matrix estimation is as follows:

- 1. Use the RFF method to calculate $W_{\rm rff}$ and $b_{\rm rff}$.
- 2. For each class *i*:
 - (a) Estimate π_i as the relative frequency of the class *i* in the dataset.
 - (b) Estimate ρ_i using eq. (20) and the training samples from class *i*.
 - (c) Find a factorization of rank r of ρ_i :

$$\rho_i = V_i^T \Lambda V_i.$$

Notice that this training procedure does not require any kind of iterative optimization. The training samples are only used once and the time complexity of the algorithm is linear on the number of training samples. The complexity of step 2(b) is $O(D^2N)$ and of 2(c) is $O(D^3)$.

Since the operations defined in eqs. (25a) to (25d) are differentiable, it is possible to use gradient-descent to minimize an appropriate loss function. For instance, we can use categorical cross entropy:

$$L = \sum_{i=1}^{K} y_i \log(\tilde{y}'_i) \tag{26}$$

where $y = (y_1, \ldots, y_K)$ corresponds to the one-hot-encoding of the real label of the sample x. The version of the model trained with gradient descent is called DMKDC-SGD.

An advantage of this approach is that the model can be jointly trained with other differentiable architecture such as a deep learning feature extractor.

5.3 Quantum measurement classification (QMC)

In DMKDC we assume a categorical distribution for the output variable. If we want a more general probability distribution we need to define a more general classification model. The idea is to model the joint probability of inputs and outputs using a density matrix. This density matrix represents the state of a bipartite system whose representation space is $\mathcal{H}_{\mathcal{X}} \otimes \mathcal{H}_{\mathcal{Y}}$ where $\mathcal{H}_{\mathcal{X}}$ is the representation space of the inputs, $\mathcal{H}_{\mathcal{Y}}$ is the representation space of the outputs and \otimes is the tensor product. A prediction is made by performing a measurement with an operator specifically prepared from a new input sample.



Fig. 3 Quantum measurement classification (QMC).

This model is based on the one described by González et al (2020) and is depicted in Figure 3 and works as follows:

• Input encoding. The input $x \in \mathbb{R}^d$ is encoded using a feature map $\phi_{\mathcal{X}}$

$$z := \phi_{\mathcal{X}}(x). \tag{27}$$

• Measurement operator. The effect of this measurement operator is to collapse, using a projector to z, the part $\mathcal{H}_{\mathcal{X}}$ of the bipartite system while keeping the $\mathcal{H}_{\mathcal{Y}}$ part unmodified. This is done by defining the following operator:

$$\pi := z z^T \otimes \mathrm{Id}_{\mathcal{H}_{\mathcal{V}}},\tag{28}$$

where $\mathrm{Id}_{\mathcal{H}_{\mathcal{V}}}$ is the identity operator in $\mathcal{H}_{\mathcal{Y}}$.

• Apply the measurement operator to the training density matrix:

$$\rho := \frac{\pi \rho_{\text{train}} \pi}{\text{Tr}[\pi \rho_{\text{train}} \pi]},\tag{29}$$

• Calculate the partial trace of ρ with respect to \mathcal{X} to obtain a density matrix that encodes the prediction:

$$\rho_{\mathcal{Y}} := \operatorname{Tr}_{\mathcal{X}}[\rho]. \tag{30}$$

The parameter of the model, without taking into account the parameters of the feature maps, is the $\rho_{\text{train}} \in \mathbb{R}^{D_{\mathcal{X}} D_{\mathcal{Y}} \times D_{\mathcal{X}} D_{\mathcal{Y}}}$ density matrix, where $D_{\mathcal{X}}$ and $D_{\mathcal{Y}}$ are the dimensions of $\mathcal{H}_{\mathcal{X}}$ and $\mathcal{H}_{\mathcal{Y}}$ respectively. As discussed in Section 5.1, the density matrix ρ_{train} can be factorized as:

$$\rho_{\rm train} = V^T \Lambda V \tag{31}$$

where $V \in \mathbb{R}^{r \times D_{\mathcal{X}} D_{\mathcal{Y}}}$, $\Lambda \in \mathbb{R}^{r \times r}$ is a diagonal matrix and $r < D_{\mathcal{X}} D_{\mathcal{Y}}$ is the reduced rank of the factorization. This factorization not only helps to reduce the space necessary to store the parameters, but learning V and Λ , instead of ρ_{train} , helps to guarantee that ρ_{train} is a valid density matrix.

As in Subsection 5.2, we described two different approaches to train the system: one based on estimation of the ρ_{train} and one based on learning ρ_{train} using gradient descent. QMC can be also trained using these two strategies.

In the estimation strategy, given a training data set $\{(x_i, y_i)\}_{i=1...N}$ the training density matrix is calculated by:

$$\rho_{\text{train}} = \frac{1}{N} \sum_{i=1}^{N} \left(\phi_{\mathcal{X}}(x_i) \otimes \phi_{\mathcal{Y}}(y_i) \right) \left(\phi_{\mathcal{X}}(x_i) \otimes \phi_{\mathcal{Y}}(y_i) \right)^T.$$
(32)

The computational cost is $O(ND_{\mathcal{X}}^2 D_{\mathcal{Y}}^2)$.

For the gradient-descent-based strategy (QMC-SGD) we can minimize the following loss function:

$$L = \sum_{i=1}^{D_{\mathcal{Y}}} y_i \log(\rho_{\mathcal{Y}ii}), \tag{33}$$

where $\rho_{\mathcal{Y}ii}$ is the *i*-th diagonal element of $\rho_{\mathcal{Y}}$.

As in DMKDC-SGD, this model can be combined with a deep learning architecture and the parameters can be jointly learned using gradient descent.

QMC can be used with different feature maps for inputs and outputs. For instance, if $\phi_{\mathcal{X}} = \phi_{\text{rff}}$ (eq. (3)) and $\phi_{\mathcal{Y}} = \phi_{\text{ohe}}$ (eq. (22)), QMC corresponds to DMKDC. However, in this case DMKDC is preferred because of its reduced computational cost.

5.4 Quantum measurement regression (QMR)

In this section we show how to use QMC to perform regression. For this we will use a feature map that allows us to encode continuous values. The feature map is defined with the help of D equally distributed landmarks in the [0, 1] interval²:

 $^{^2\}mathrm{Without}$ loss of generality the continuous variable to be encoded is restricted to the [0,1] interval.

$$\alpha_i = \frac{i-1}{D-1} \text{ for } i = 1 \dots D.$$
(34)

The following function (which is equivalent to a softmax) defines a set of unimodal probability density functions centered at each landmark:

$$p_i(x) = \left(\frac{\exp(-\beta \|x - \alpha_i\|^2)}{\sum_{j=1}^m \exp(-\beta \|x - \alpha_j\|^2)}\right)_{i=1\dots D},$$
(35)

where β controls the shape of the density functions.

The feature map is defined as:

$$\phi_{\rm sm}: [0,1] \to \mathbb{R}^D$$

$$x \mapsto (\sqrt{p_1(x)}, \dots, \sqrt{p_D(x)}).$$
(36)

This feature map is used in QMC as the feature map of the output variable $(\phi_{\mathcal{Y}})$. To calculate the prediction for a new sample x we apply the process described in Subsection 5.3 to obtain $\rho_{\mathcal{Y}}$. Then the prediction is given by:

$$\hat{y} = E_{\rho_{\mathcal{Y}}}[\alpha_i] = \sum_{i=1}^{D} \rho_{\mathcal{Y}ii}\alpha_i.$$
(37)

Note that this framework also allows to easily compute confidence intervals for the prediction. The model can be trained using the strategies discussed in Subsection 5.3. For gradient-based optimization we use a mean squared error loss function:

$$L = \sum_{i=1}^{D} (y - \hat{y})^2 + \alpha \sum_{i=1}^{D} \rho_{\mathcal{Y}ii} (\hat{y} - \alpha_i)^2$$
(38)

where the second term correspond to the variance of the prediction and α controls the trade-off between error and variance.

6 Related Work

The ability of density matrices to represent probability distributions has been used in previous works. The early work by Wolf (2006) uses the density matrix formalism to perform spectral clustering, and shows that this formalism not only is able to predict cluster labels for the objects being classified, but also provides the probability that the object belongs to each of the clusters. Similarly, Tiwari and Melucci (2019) proposed a quantum-inspired binary classifier using density matrices, where samples are encoded into pure quantum states. In a similar fashion, Sergioli et al (2018) proposed a quantum nearest mean

classifier based on the trace distance between the quantum state of a sample, and a quantum centroid that is a mixed state of the pure quantum states of all samples belonging to a single class. Another class of proposals directly combine these quantum ideas with customary machine learning techniques, such as frameworks for multi-modal learning for sentiment analysis (Li et al, 2021; Li et al, 2020; Zhang et al, 2018).

Since its inception, random features have been used to improve the performance of several kernel methods: kernel ridge regression (Avron et al, 2017), support vector machines (SVM) (Sun et al, 2018), and nonlinear component analysis (Xie et al, 2015). Besides, random features have been used in conjunction with deep learning architectures in different works (Arora et al, 2019; Ji and Telgarsky, 2019; Li et al, 2019).

The combination of RFF and density matrices was initially proposed by González et al (2020). In that work, RFF are used as a quantum feature map, among others, and the QMC method (Subsection 5.3) was presented. In González et al (2020) the coherent state kernel showed better performance than the Gaussian kernel. It is important to notice that the coherent state kernel was calculated exactly while the Gaussian kernel was approximated using RFF. It is possible to apply RFF to approximate the coherent state kernel and use it as the quantum feature map in the models presented in this paper. The emphasis of González et al (2020) is to show that quantum measurement can be used to do supervised learning. In contrast, the present paper addresses a wider problem with several new contributions: a new method for density estimation based on density matrices and RFF, the proof of the connection between this method and kernel density estimation, and new differentiable models for density estimation, classification and regression.

The present work can be seen as a type of quantum machine learning (QML), which is generally referred as the field in the intersection of quantum computing and machine learning (Schuld et al, 2015; Schuld, 2018). In particular, the methods in this paper are in the subcategory of QML called quantum inspired classical machine learning, where theory and methods from quantum physics are borrowed and adapted to machine learning methods intended to run in classical computers. Works in this category include: quantum-inspired recommendation systems (Tang, 2019a), quantum-inspired kernel-based classification methods (Tiwari et al, 2020; González et al, 2020), conversational sentiment analysis based on density matrix-like convolutional neural networks (Zhang et al, 2019), dequantised principal component analysis (Tang, 2019b), among others.

Being a memory-based strategy, KDE suffers from large-scale, high dimensional data. Due to this issue, fast approximate evaluation of non-parametric density estimation is an active research topic. Different approaches are proposed in the literature: higher-order divide-and-conquer method (Gray and Moore, 2003), separation of near and far-field (pruning) (March et al, 2015), and hashing based estimators (HBE) (Charikar and Siminelakis, 2017). Even though the purpose of the present work was not to design methods for fast approximation of KDE, the use of RFF to speed KDE seems to be a promising research direction. Comparing DMKDE against fast KDE approximation methods is part of our future work.

7 Experimental Evaluation

In this section we perform some experiments to evaluate the performance of the proposed methods in different benchmark tasks. The experiments are organized in three subsections: density estimation evaluation, classification evaluation and ordinal regression evaluation. The source code of the methods and the scripts of the experiments are available at https://drive.google.com/drive/folders/16pHMLjIvr6v1zY6cMvo11EqMAMqjn3Xa as Jupyter notebooks.

7.1 Density estimation evaluation

The goal of these experiments is to evaluate the efficacy and efficiency of DMKDE to approximate a pdf. We compare it against conventional Gaussian KDE.

7.1.1 Data sets and experimental setup

We used three datasets:

- 1-D synthetic. The first synthetic dataset corresponds to a mixture of univariate Gaussians as shown in Figure 4. The mixture weights are 0.3 and 0.7 respectively and the parameters are $(\mu_1 = 0, \sigma = 1)$ and $(\mu_1 = 5, \sigma = 1)$. We generated 10,000 samples for training and use as test dataset 1,000 samples equally spaced in the interval [-5, 10].
- 2-D synthetic. This dataset corresponds to three spirals as depicted in Figure 6. The training and test datasets have 10,0000 and 1,000 points respectively, all of them generated with the same stochastic procedure.
- MNIST dataset. We used PCA to reduce the original 784 dimension to 40. The resulting vectors were scaled to [0, 1]. We used stratified sampling to choose 10,000 and 1,000 samples for training and testing respectively.

We performed two types of experiments over the three datasets. In the first, we wanted to evaluate the accuracy of DMKDE. In the second, we evaluated the time to predict the density on the test set.

In the first experiment, DMKDE was run with different number of RFF to see how the dimension of the RFF representation affected the accuracy of the estimation. For the 1-D dataset, both the DMKDE prediction and the KDE prediction were compared against the true pdf using root mean squared error (RMSE). For the 2-D dataset the RMSE between the DMKDE prediction and the KDE prediction was evaluated. In the case of MNIST, and because of the small values for the density, we calculated the RMSE between the log density predicted by DMKDE and KDE. The number of eigencomponents (r) was chosen by sorting the eigenvalues in descending order and plotting them to look



Fig. 4 1-D synthetic dataset. The gray zone is the area of the true density. The estimated pdf by DMKDE ($\gamma = 2$) and KDE ($\gamma = 4$) is shown.

for the curve elbow. For the 1-D and 2-D datasets, the γ value was chosen to get a good approximation of the data density, this was visually verified. For the MNIST dataset, the γ value was chosen by looking at a histogram of pairwise distances of the data. The value of the parameters were: ($\gamma = 16, r = 30$) for the 1-D dataset, ($\gamma = 256, r = 100$) for the 2-D dataset, ($\gamma = 1, r = 150$) for the MNIST dataset.

For the second experiment, we measured the time taken to predict 1,000 test samples for both KDE and DMKDE using different number of train samples. KDE was implemented in Python using liner algebra operations accelerated by numpy. At least for the experiments reported, our implementation was faster than other KDE implementations available such as the one provided by scikit learn (https://scikit-learn.org/stable/modules/density.html), which is probably optimized for other use cases. DMKDE was implemented in Python using Tensorflow. The main reason for using Tensorflow was its ability to automatically calculate the gradient of computational graphs. KDE could not benefit from this feature, on the contrary, its performance could be hurt by Tensorflow's larger memory footprint. Another advantage of Tensorflow is its ability to generate code optimized for a GPU, so both methods were run on a 2.20 GHz dual-core Intel(R) Xeon(R) CPU without a GPU to avoid any unfair advantage.

7.1.2 Results and discussion

Figure 5 shows how the accuracy of DMKDE increases with an increasing number of RFF. For each configuration 30 experiments were run and the blue solid line represents the mean RMSE of the experiments and the blue region represents the 95% confidence interval. In all the three datasets, 2^{10} RFF





Fig. 5 Accuracy of the density estimation of DMKDE for different number of RFF for the 1-D dataset (top left), 2-D dataset (top right) and MNIST dataset (bottom). For the 1-D dataset both KDE and DMKDE are compared against the true density. For the two other datasets the difference between KDE and DMKDE is calculated. In all the cases the RMSE is calculated. The blue shaded zone represents the 95% confidence interval.

Figure 6 shows the 2-D spirals dataset (left) and the density estimation of both KDE (center) and DMKDE (right). The density calculated by DMKDE is very close to the one calculated with KDE.

Figure 7 shows a comparison of the log density predicted by KDE and DMKDE. Both models were applied to test samples and samples generated randomly from a uniform distribution. As expected points are clustered around the diagonal. The DMKDE log density of test samples (left) seems to be more accurately predicted than the one of random samples. The reason is that the density of random samples is smaller than the density of test samples and the difference is amplified by the logarithm.



Fig. 6 2-D spirals dataset (top left) and the density estimation of both KDE (top right) and DMKDE (bottom).

Figure 8 shows the time of both methods for different sizes of the training dataset. The prediction time of KDE depends on the size of the training dataset, while the time of DMKDE does not depend on it. The advantage of DMKDE in terms of computation time is clear for training datasets above 10^4 data samples.

7.2 Classification evaluation

In this set of experiments, we evaluated DMKDC over different well known benchmark classification datasets.

7.2.1 Data sets and experimental setup

Six benchmark data sets were used. The details of these datasets are shown in Table 1. In the case of Gisette and Cifar, we applied a principal component analysis algorithm using 400 principal components in order to reduce the dimension. DMKDC was trained using the estimation strategy (DMKDC) and an ADAM stochastic gradient descent strategy (DMKDC-SGD). As baseline we compared against a linear support vector machine (SVM) trained using the



Fig. 7 Scatter-plots comparing the log density predicted by KDE and DMKDE: test samples (top left), uniformly random generated samples (top right), both test and random samples (bottom).

Table 1 Data sets used for classification evaluation.

Data set	Attributes	CLASSES	TRAIN-TEST
Letters	16	26	14000-6000
USPS	256	10	7291-2007
Forest	54	3	70-30
MNIST	784	10	60000-10000
GISETTE	5000	2	4200-1800
CIFAR	3072	10	60000-10000

same RFF as DMKDC. The SVM was trained using the LinearSVC model from scikit-learn, which is based in an efficient C implementation tailored to linear SVMs. In the case of MNIST and Cifar, we additionally built a union of a LeNet architecture (LeCun et al, 1989), as a feature extraction block, and DMKDC-SGD as the classifier layer. The LeNet block is composed of two





Fig. 8 Evaluation of the prediction time of DMKDE and KDE: 1-D dataset (top left), 2-D dataset (top right) and MNIST dataset (bottom).

convolutional layers, one flatten layer and one dense layer. The first convolutional layer has 20 filters, kernel size of 5, same as padding, and ReLu as the activation function. The second convolutional layer has 50 filters, kernel size of 5, same as padding, and ReLu as the activation function. The dense layer has 84 units and ReLU as the activation function. The dense layer is finally connected to DMKDC. We report results for the combined model (LeNet DMKDC-SGD) and the LeNet model with a softmax output layer (LeNet). To make the comparison with baseline models fair, in all the cases the RFF layer of DMKDC-SGD is frozen, so its weights are not modified by the stochastic gradient descent learning process.

For each data set, we made a hyper parameter search using a five-fold crossvalidation with 25 randomly generated configurations. The number of RFF was set to 1000 for all the methods. For each dataset we calculated the inter-sample median distance μ and defined an interval around $\gamma = \frac{1}{2\sigma^2}$. The *C* parameter of the SVM was explored in an exponential scale from 2^{-5} to 2^{10} . For the ADAM optimizer in DMKDC-SGD with and without LeNet we choose the

Data set	ATTRIBUTES	TRAIN	Test	
DIABETES	2	30	13	
Pyrimidines	27	50	24	
TRIAZINES	60	100	86	
WISCONSIN	32	130	64	
Machine CPU	6	150	59	
Auto MPG	7	200	192	
Boston Housing	13	300	206	
STOCK DOMAIN	9	600	350	
Abalone	8	1000	3177	

 Table 2
 Specifications of the data sets used for ordinal regression evaluation. Train and

 Test indicate the number of samples, which is the same for all the twenty partitions.

learning rate in the interval (0, 0.001]. The number of eigen-components of the factorization was chosen from $\{0.1, 0.2, 0.5, 1\}$ where each number represents a percentage of the RFF. After finding the best hyper-parameter configuration using cross validation, 10 different experiments were performed with different random initialization. The mean and the standard deviation of the accuracy is reported.

7.2.2 Results and discussion

Table 3 shows the results of the classification experiments. DMKDC is a shallow method that uses RFF, so a SVM using the same RFF is fair and strong baseline. In all the cases, except one, DMKDC-SGD outperforms the SVM, which shows that it is a very competitive shallow classification method. DMKDC trained using estimation shows less competitive results, but they are still remarkable taking into account that this is an optimization-less training strategy that only passes once over the training dataset. For MNIST and Cifar the use of a deep learning feature extractor is mandatory to obtain competitive results. The results show that DMKDC-SGD can be integrated with deep neural network architectures to obtain competitive results.

The improvement on classification performance of DMKC-SGD comes at the cost of increased training time. The training of DMKDC is very efficient since it corresponds to do an average of the training density matrices. Linear SVM training is also very efficient. In contrast, DMKDC-SGD requires an iterative training process that has to be tuned to get it to converge to a good local optimum, as is the case for current deep learning models.

7.3 Ordinal regression evaluation

Many multi-class classification problems can be seen as ordinal regression problems. That is, problems where labels not only indicate class membership, but also an order. Ordinal regression problems are halfway between a classification

SVM-RFF DMKDC 0.924±0.002 0.918±0.00 0.940±0.001 0.806±0.00 0.941±0.003 0.836±0.00 0.944±0.003 0.836±0.00 0.950±0.006 0.811±0.00 0.453±0.003 0.271±0.00
SVM-RFF 0.924±0.002 (0.940±0.001 (0.697±0.046 (0.944±0.003 (0.950±0.006 (0.453±0.003 (

Table 3 Accuracy test results for DMKDC and DMKDC-SGD compared against a linear support vector machine over RFF (SVM-RFF). Two deep learning models are also evaluated on the two image datasets: a convolutional neural network (LeNet) and its combination with DMKDC-SGD (LeNet DMKDC).

Learning with Density Matrices and Random Features

problem and a regression problem, and given the discrete probability distribution representation used in QMR, ordinal regression seems to be a suitable problem to test it.

7.3.1 Data sets and experimental setup

Nine standard benchmark data sets for ordinal regression were used. The details of each data set are reported in Table 2. These data sets are originally used in metric regression tasks. To convert the task into an ordinal regression one, the target values were discretized by taking five intervals of equal length over the target range. For each set, 20 different train and test partitions are made. These partitions are the same used by Chu and Ghahramani (2005) and several posterior works, and are publicly available at http://www.gatsby.ucl.ac.uk/~chuwei/ordinalregression.html. The models were evaluated using the mean absolute error (MAE), which is a popular and widely used measure in ordinal regression (Gutiérrez et al, 2016; Garg and Manwani, 2020).

QMR was trained using the estimation strategy (QMR) and an ADAM stochastic gradient descent strategy (QMR-SGD). For each data set, and for each one of the 20 partitions, we made a hyper parameter search using a fivefold cross-validation procedure. The search was done generating 25 different random configuration. The range for γ was computed in the same way as for the classification experiments, $\beta \in (0, 25)$, the number of RFF randomly chosen between the number of attributes and 1024, and the number of eigencomponents of the factorization was chosen from $\{0.1, 0.2, 0.5, 1\}$ where each number represents a percentage of the RFF. For the ADAM optimizer in QMR-SGD we choose the learning rate in the interval (0, 0.001] and $\alpha \in (0, 1)$. The RFF layer was always set to trainable, and the criteria for selecting the best parameter configuration was the MAE performance.

7.3.2 Results and discussion

For each data set, the means and standard deviations of the test MAE for the 20 partitions are reported in Table 4, together with the results of previous state-of-the-art works on ordinal regression: Gaussian Processes (GP) and support vector machines (SVM) (Chu and Ghahramani, 2005), Neural Network Rank (NNRank) (Cheng et al, 2008), Ordinal Extreme Learning Machines (ORELM) (Deng et al, 2010) and Ordinal Regression Neural Network (ORNN) (Fernandez-Navarro et al, 2014).

QMR-SGD shows a very competitive performance. It outperforms the baseline methods in six out of the nine data sets. The training strategy based on estimation, QMR, did not performed as well. This evidences that for this problem a fine tuning of the representation is required and it is successfully accomplished by the gradient descent optimization.

SVM	$.14 0.746\pm0.14$.07 0.450±0.11	0.02 0.698 ± 0.03	$09 1.003\pm0.07$	0.192 ± 0.04	$0.02 0.260\pm0.02$	$.02 0.267 \pm 0.02$	$0.02 0.108 \pm 0.02$	$00 0.229 \pm 0.00$
GP	$0.662{\pm}0$	0.392 ± 0	0.687 ± 0	1.010 ± 0	0.185 ± 0	0.241 ± 0	0.259 ± 0	0.120 ± 0	0.232 ± 0
NNRANK	0.546 ± 0.15	$0.450{\pm}0.10$	$0.730{\pm}0.07$	I	$0.186{\pm}0.04$	$0.281{\pm}0.02$	$0.295{\pm}0.04$	$0.127{\pm}0.02$	$0.226{\pm}0.01$
ORNN	I	$0.677{\pm}0.20$	Ι	Ι	$0.451{\pm}0.03$	Ι	Ι	$0.127{\pm}0.01$	$0.635{\pm}0.01$
QMR	0.611 ± 0.02	$0.946{\pm}0.06$	0.695 ± 0.01	1.114 ± 0.04	0.995 ± 0.07	$0.710{\pm}0.02$	0.679 ± 0.02	$0.971{\pm}0.00$	$0.307{\pm}0.00$
QMR-SGD	$0.511 {\pm} 0.07$	$0.408 {\pm} 0.07$	$0.674{\pm}0.02$	$0.985{\pm}0.04$	$0.171{\pm}0.03$	$0.230{\pm}0.02$	$0.270{\pm}0.02$	$0.103{\pm}0.01$	$0.233 {\pm} 0.01$
DATA SET	DIABETES	Pyrimidines	TRIAZINES	WISCONSIN	MACHINE	Auto	BOSTON	Stocks	Abalone

Table 4 MAE test results for QMR, QMR-SGD and different baseline methods: support vector machines (SVM), Gaussian Processes (GP), Neural Network Rank (NNRank), Ordinal Extreme Learning Machines (ORELM) and Ordinal Regression Neural Network (ORNN). The results are the mean and standard deviation of the MAE for the twenty partitions of each dataset. The best result for each data set is in bold face.

Learning with Density Matrices and Random Features

8 Conclusions

The mathematical framework underlying quantum mechanics is a powerful formalism that harmoniously combine linear algebra and probability in the form of density matrices. This paper has shown how to use these density matrices as a building block for designing different machine learning models. The main contribution of this work is to show a novel perspective to learning that combines two very different and seemingly unrelated tools, random features and density matrices. The, somehow surprising, connection of this combination with kernel density estimation provides a new way of representing and learning probability density functions from data. The experimental results showed some evidence that this building block can be used to build competitive models for some particular tasks. However, the full potential of this new perspective is still to be explored. Examples of directions of future inquire include using complex valued density matrices, exploring the role of entanglement and exploiting the battery of practical and theoretical tools provided by quantum information theory.

Statements and Declarations

The authors declare that they have no known competing interests.

References

- Arora S, Du SS, Hu W, et al (2019) On exact computation with an infinitely wide neural net. arXiv preprint arXiv:190411955
- Avron H, Sindhwani V, Yang J, et al (2016) Quasi-monte carlo feature maps for shift-invariant kernels. Journal of Machine Learning Research 17(120):1–38. URL http://jmlr.org/papers/v17/14-538.html
- Avron H, Kapralov M, Musco C, et al (2017) Random fourier features for kernel ridge regression: Approximation bounds and statistical guarantees.
 In: International Conference on Machine Learning, PMLR, pp 253–262
- Balakrishnan S, Narayanan S, Rinaldo A, et al (2013) Cluster trees on manifolds. arXiv preprint arXiv:13076515
- Bishop CM (2006) Pattern recognition and machine learning (information science and statistics)
- Charikar M, Siminelakis P (2017) Hashing-based-estimators for kernel density in high dimensions. In: 2017 IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS), IEEE, pp 1032–1043
- Chazal F, Fasy B, Lecci F, et al (2017) Robust topological inference: Distance to a measure and kernel distance. The Journal of Machine Learning Research

18(1):5845-5884

- Chen YC, Genovese CR, Wasserman L, et al (2016) A comprehensive approach to mode clustering. Electronic Journal of Statistics 10(1):210–241
- Cheng J, Wang Z, Pollastri G (2008) A neural network approach to ordinal regression. Proceedings of the International Joint Conference on Neural Networks (May 2014):1279–1284. https://doi.org/10.1109/IJCNN.2008. 4633963, https://arxiv.org/abs/arXiv:0704.1028
- Chernozhukov V, Chetverikov D, Kato K, et al (2014) Gaussian approximation of suprema of empirical processes. Annals of Statistics 42(4):1564–1597
- Chu W, Ghahramani Z (2005) Gaussian Processes for Ordinal Regression. Journal of Machine Learning Research 6:1019–1041. URL http://www.jmlr. org/papers/volume6/chu05a/chu05a.pdf
- Dao T, De Sa C, Ré C (2017) Gaussian quadrature for kernel features. Advances in neural information processing systems 30:6109
- Deng WY, Zheng QH, Lian S, et al (2010) Ordinal extreme learning machine. Neurocomputing 74(1-3):447–456. https://doi.org/10.1016/j.neucom.2010. 08.022, URL http://dx.doi.org/10.1016/j.neucom.2010.08.022
- Efron B (1992) Bootstrap methods: another look at the jackknife. In: Breakthroughs in statistics. Springer, p 569–593
- Fernandez-Navarro F, Riccardi A, Carloni S (2014) Ordinal neural networks without iterative tuning. IEEE Transactions on Neural Networks and Learning Systems 25(11):2075–2085. https://doi.org/10.1109/TNNLS.2014. 2304976
- Garg B, Manwani N (2020) Robust deep ordinal regression under label noise. In: Pan SJ, Sugiyama M (eds) Proceedings of The 12th Asian Conference on Machine Learning, Proceedings of Machine Learning Research, vol 129. PMLR, Bangkok, Thailand, pp 782–796, URL http://proceedings.mlr.press/ v129/garg20a.html
- Genovese CR, Perone-Pacifico M, Verdinelli I, et al (2014) Nonparametric ridge estimation. Annals of Statistics 42(4):1511–1545
- González FA, Vargas-Calderón V, Vinck-Posada H (2020) Supervised Learning with Quantum Measurements. arXiv preprint arXiv:200401227 URL http: //arxiv.org/abs/2004.01227, https://arxiv.org/abs/arXiv:2004.01227
- González FA, Vargas-Calderón V, Vinck-Posada H (2021) Classification with quantum measurements. Journal of the Physical Society of Japan

90(4):044,002

- Gray AG, Moore AW (2003) Nonparametric density estimation: Toward computational tractability. In: Proceedings of the 2003 SIAM International Conference on Data Mining, SIAM, pp 203–211
- Gutiérrez PA, Pérez-Ortiz M, Sánchez-Monedero J, et al (2016) Ordinal Regression Methods: Survey and Experimental Study. IEEE Transactions on Knowledge and Data Engineering 28(1):127–146. https://doi.org/10.1109/ TKDE.2015.2457911
- Hastie T, Tibshirani R, Friedman J (2009) The elements of statistical learning: data mining, inference, and prediction. Springer Science & Business Media
- Ji Z, Telgarsky M (2019) Polylogarithmic width suffices for gradient descent to achieve arbitrarily small test error with shallow relu networks. arXiv preprint arXiv:190912292
- Le Q, Sarlos T, Smola A (2013) Fastfood approximating kernel expansions in loglinear time. In: 30th International Conference on Machine Learning (ICML), URL http://jmlr.org/proceedings/papers/v28/le13.html
- LeCun Y, Boser B, Denker JS, et al (1989) Backpropagation applied to handwritten zip code recognition. Neural computation 1(4):541–551
- Li CL, Chang WC, Mroueh Y, et al (2019) Implicit kernel learning. In: The 22nd International Conference on Artificial Intelligence and Statistics, PMLR, pp 2007–2016
- Li Q, Stefani A, Toto G, et al (2020) Towards multimodal sentiment analysis inspired by the quantum theoretical framework. In: 2020 IEEE Conference on Multimedia Information Processing and Retrieval (MIPR), pp 177–180, https://doi.org/10.1109/MIPR49039.2020.00044
- Li Q, Gkoumas D, Lioma C, et al (2021) Quantum-inspired multimodal fusion for video sentiment analysis. Information Fusion 65:58 – 71. https: //doi.org/https://doi.org/10.1016/j.inffus.2020.08.006, URL http://www. sciencedirect.com/science/article/pii/S1566253520303365
- March WB, Xiao B, Biros G (2015) Askit: Approximate skeletonization kernel-independent treecode in high dimensions. SIAM Journal on Scientific Computing 37(2):A1089–A1110
- McNeil BJ, Hanley JA (1984) Statistical approaches to the analysis of receiver operating characteristic (roc) curves. Medical decision making 4(2):137–150

- 30 Learning with Density Matrices and Random Features
- Nadaraya EA (1964) Some new estimates for distribution functions. Theory of Probability & Its Applications 9(3):497–500
- Parzen E (1962) On estimation of a probability density function and mode. The annals of mathematical statistics 33(3):1065-1076
- Rahimi A, Recht B (2007) Random features for large-scale kernel machines. In: Proceedings of the 20th International Conference on Neural Information Processing Systems. Curran Associates Inc., Red Hook, NY, USA, NIPS'07, p 1177–1184
- Rosenblatt M (1956) Remarks on some nonparametric estimates of a density function. Ann Math Statist 27(3):832–837. https://doi.org/10.1214/aoms/ 1177728190, URL https://doi.org/10.1214/aoms/1177728190
- Schuld M (2018) Supervised learning with quantum computers. Springer
- Schuld M, Sinayskiy I, Petruccione F (2015) An introduction to quantum machine learning. Contemporary Physics 56(2):172–185
- Sergioli G, Santucci E, Didaci L, et al (2018) A quantum-inspired version of the nearest mean classifier. Soft Computing 22(3):691–705. https://doi.org/10. 1007/s00500-016-2478-2, URL https://doi.org/10.1007/s00500-016-2478-2
- Shen W, Yang Z, Wang J (2017) Random features for shift-invariant kernels with moment matching. In: Proceedings of the AAAI Conference on Artificial Intelligence
- Sun Y, Gilbert A, Tewari A (2018) But how does it work in theory? linear svm with random features. arXiv preprint arXiv:180904481
- Tang E (2019a) A quantum-inspired classical algorithm for recommendation systems. In: Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing, pp 217–228
- Tang E (2019b) Quantum-inspired classical algorithms for principal component analysis and supervised clustering. 1811.00414
- Tiwari P, Melucci M (2019) Towards a quantum-inspired binary classifier. IEEE Access 7:42,354–42,372. https://doi.org/10.1109/ACCESS.2019. 2904624
- Tiwari P, Dehdashti S, Obeid AK, et al (2020) Kernel method based on nonlinear coherent state. arXiv preprint arXiv:200707887
- Von Neumann J (1927) Wahrscheinlichkeitstheoretischer aufbau der quantenmechanik. Nachrichten von der Gesellschaft der Wissenschaften zu Göttingen, Mathematisch-Physikalische Klasse 1927:245–272

- Wilce A (2021) Quantum Logic and Probability Theory. In: Zalta EN (ed) The Stanford Encyclopedia of Philosophy, Fall 2021 edn. Metaphysics Research Lab, Stanford University
- Wolf L (2006) Learning using the born rule. Tech. rep., Massachusetts Institute of Technology Computer Science and Artificial Intelligence Laboratory
- Xie B, Liang Y, Song L (2015) Scale up nonlinear component analysis with doubly stochastic gradients. arXiv preprint arXiv:150403655
- Yu FXX, Suresh AT, Choromanski KM, et al (2016) Orthogonal random features. In: Lee D, Sugiyama M, Luxburg U, et al (eds) Advances in Neural Information Processing Systems, vol 29. Curran Associates, Inc., pp 1975–1983, URL https://proceedings.neurips.cc/paper/2016/file/ 53adaf494dc89ef7196d73636eb2451b-Paper.pdf
- Zhang Y, Song D, Zhang P, et al (2018) A quantum-inspired multimodal sentiment analysis framework. Theoretical Computer Science 752:21 – 40. https://doi.org/https://doi.org/10.1016/j.tcs.2018.04.029, URL http: //www.sciencedirect.com/science/article/pii/S0304397518302639, quantum structures in computer science: language, semantics, retrieval
- Zhang Y, Li Q, Song D, et al (2019) Quantum-inspired interactive networks for conversational sentiment analysis. In: 28th International Joint Conference on Artificial Intelligence (IJCAI2019), URL http://oro.open.ac.uk/61755/

Appendix A Proofs

Proposition 3 Let \mathcal{M} be a compact subset of \mathbb{R}^d with a diameter diam (\mathcal{M}) , let $X = \{x_i\}_{i=1...N} \subset \mathcal{M}$ a set of iid samples, then $\hat{f}_{\rho_{\text{train}}}$ (eq. (6)) and \hat{f}_{γ} satisfy:

$$\Pr\left[\sup_{x \in \mathcal{M}} |\hat{f}_{\rho_{\text{train}}}(x) - \hat{f}_{\gamma}(x)| \ge \epsilon\right] \le 2^{8} \left(\frac{\sqrt{2d\gamma} \operatorname{diam}(\mathcal{M})}{3M_{\gamma}\epsilon}\right)^{2} \exp\left(-\frac{D(3M_{\gamma}\epsilon)^{2}}{4(d+2)}\right)$$
(A1)

Proof

e

$$\hat{f}_{\rho_{\text{train}}}(x) = \frac{1}{M_{\gamma}} \phi_{\text{rff}}^{T}(x) \rho_{\text{train}} \phi_{\text{rff}}(x)$$

$$= \frac{1}{M_{\gamma}} \phi_{\text{rff}}^{T}(x) \left(\frac{1}{N} \sum_{i=1}^{N} \phi_{\text{rff}}(x_i) \phi_{\text{rff}}^{T}(x_i)\right) \phi_{\text{rff}}(x)$$

$$= \frac{1}{M_{\gamma}N} \sum_{i=1}^{N} \phi_{\text{rff}}^{T}(x) \phi_{\text{rff}}(x_i) \phi_{\text{rff}}^{T}(x_i) \phi_{\text{rff}}(x)$$

$$= \frac{1}{M_{\gamma}N} \sum_{i=1}^{N} (\phi_{\text{rff}}^{T}(x) \phi_{\text{rff}}(x_i))^2$$
(A2)

Remembering that in eq. (6) we used a spread parameter of $\frac{\gamma}{2}$ to calculate the parameters of $\phi_{\rm rff}$ and because of Theorem 1 we know that

$$\Pr\left[\sup_{x,y\in\mathcal{M}} |\phi_{\mathrm{rff}}^{T}(x)\phi_{\mathrm{rff}}(y) - e^{-\frac{\gamma}{2}\|x-y\|^{2}}| \ge \epsilon\right] \le 2^{8} \left(\frac{\sqrt{d\gamma}\mathrm{diam}(\mathcal{M})}{\epsilon}\right)^{2} \exp\left(-\frac{D\epsilon^{2}}{4(d+2)}\right) = B$$
By construction $|\phi_{\mathrm{rff}}^{T}(x)\phi_{\mathrm{rff}}(y) + e^{-\frac{\gamma}{2}\|x-y\|^{2}}| \le 3$, then $|(\phi_{\mathrm{rff}}^{T}(x)\phi_{\mathrm{rff}}(y))^{2} - e^{-\gamma\|x-y\|^{2}}| = |(\phi_{\mathrm{rff}}^{T}(x)\phi_{\mathrm{rff}}(y) - e^{-\frac{\gamma}{2}\|x-y\|^{2}})(\phi_{\mathrm{rff}}^{T}(x)\phi_{\mathrm{rff}}(y) + e^{-\frac{\gamma}{2}\|x-y\|^{2}})| \le 3|(\phi_{\mathrm{rff}}^{T}(x)\phi_{\mathrm{rff}}(y) - e^{-\frac{\gamma}{2}\|x-y\|^{2}})|$. Then
$$\Pr\left[\sup_{x,y\in\mathcal{M}} |(\phi_{\mathrm{rff}}^{T}(x)\phi_{\mathrm{rff}}(y))^{2} - e^{-\gamma\|x-y\|^{2}}| \ge 3\epsilon\right] \le B$$
(A3)

$$\Pr\left[\sup_{x,y\in\mathcal{M}} |(\phi_{\mathrm{rff}}^T(x)\phi_{\mathrm{rff}}(y))^2 - e^{-\gamma||x-y||^2}| \ge 3\epsilon\right] \le B$$
(A3)

Combining Equations eq. (A2) and eq. (A3) we get:

$$\Pr\left[\sup_{x\in\mathcal{M}}|\hat{f}_{\rho}(x)-\hat{f}_{\gamma}(x)|\geq 3M_{\gamma}\epsilon\right]\leq B$$

Making a variable change we get:

$$\Pr\left[\sup_{x \in \mathcal{M}} |\hat{f}_{\rho_{\text{train}}}(x) - \hat{f}_{\gamma}(x)| \ge \epsilon\right] \le 2^{8} \left(\frac{\sqrt{2d\gamma}\operatorname{diam}(\mathcal{M})}{3M_{\gamma}\epsilon}\right)^{2} \exp\left(-\frac{D(3M_{\gamma}\epsilon)^{2}}{4(d+2)}\right)$$
(A4)