

# Learning with Density Matrices and Random Fourier Features

Fabio González  
Machine Learning 2021-1

# ① Kernel density estimation (KDE)

Density estimation: Given a sample  $\{x_i\}_{i=1..N}$  from an unknown distribution estimate the PDF of the distribution.

KDE (Rosenblatt, 1956) (Parzen, 1962)

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

kernel function

Bandwith

PDF estimator

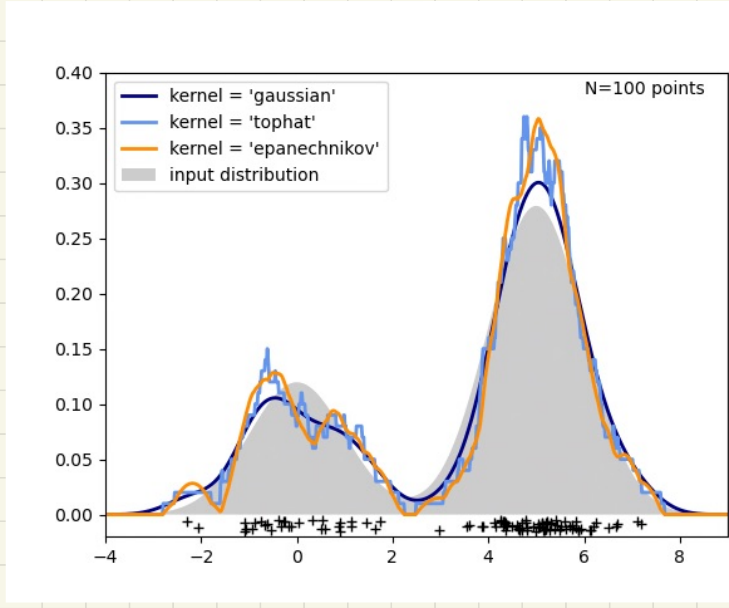
$$\hat{g}_{\gamma, X}(x) = \frac{1}{N(\pi/\gamma)^{\frac{d}{2}}} \sum_{i=1}^N e^{-\gamma \|x_i - x\|^2}$$

Gaussian kernel

$$\gamma = \frac{1}{2\sigma^2}$$

## Drawbacks

- Memory based method: you have to store all the training dataset.
- Prediction time  $O(N)$
- Problems dealing with high-dimensional data



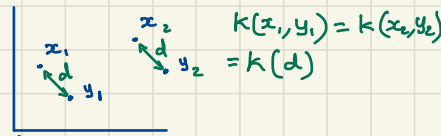
## ② Random Fourier Features (RFF) (Rahini & Recht, 2007)

**Theorem 1** (Bochner [13]). A continuous kernel  $k(\mathbf{x}, \mathbf{y}) = k(\mathbf{x} - \mathbf{y})$  on  $\mathcal{R}^d$  is positive definite if and only if  $k(\delta)$  is the Fourier transform of a non-negative measure.

If a shift-invariant kernel  $k(\delta)$  is properly scaled, Bochner's theorem guarantees that its Fourier transform  $p(\omega)$  is a proper probability distribution. Defining  $\zeta_\omega(\mathbf{x}) = e^{j\omega^T \mathbf{x}}$ , we have

$$k(\mathbf{x} - \mathbf{y}) = \int_{\mathcal{R}^d} p(\omega) e^{j\omega^T (\mathbf{x} - \mathbf{y})} d\omega = E_\omega[\zeta_\omega(\mathbf{x}) \zeta_\omega(\mathbf{y})^*], \quad (2)$$

→ Isotropic kernel



if  $k$  is the Gaussian kernel,  $k(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle_F$ , the dimension of  $F$  is infinite  
RFF method: Finds an embedding  $\phi_{\text{RFF}}: \mathbb{R}^d \rightarrow \mathbb{R}^D$  such that

$$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d \quad \langle \phi_{\text{RFF}}(\mathbf{x}), \phi_{\text{RFF}}(\mathbf{y}) \rangle \approx k(\mathbf{x}, \mathbf{y})$$

$$\phi_{\text{RFF}}: \mathbb{R}^d \rightarrow \mathbb{R}^D$$

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

(1)

$$w_1 \dots w_D \sim P(\omega)$$

$$b_1 \dots b_D \sim U(0, 2\pi)$$

**Advantage:** kernel methods complexity is typically  $\mathcal{O}(N^2)$   
with RFF you can reduce this complexity

### ③ RFF and KDE

$$\begin{aligned}\hat{g}_{r,x}(x) &= \frac{1}{N} \sum_{i=1}^N k_{\sigma}(x_i, x) \\ &\approx \frac{1}{N} \sum_{i=1}^N \langle \phi_{\text{rff}}(x_i), \phi_{\text{rff}}(x) \rangle \\ &\approx \left\langle \frac{1}{N} \sum_{i=1}^N \phi_{\text{rff}}(x_i), \phi_{\text{rff}}(x) \right\rangle\end{aligned}$$

$$\begin{aligned}g_{r,x}(x) &= \frac{1}{N} \sum_{i=1}^N k_{\sigma/2}^2(x_i, x) \\ &= \frac{1}{N} \sum_{i=1}^N \left( \phi_{\text{rff}}(x_i)^T \phi_{\text{rff}}(x) \right)^2 \\ &= \frac{1}{N} \sum_{i=1}^N \phi_{\text{rff}}(x)^T \underbrace{\phi_{\text{rff}}(x_i) \phi_{\text{rff}}(x_i)^T}_{\rho} \phi_{\text{rff}}(x) \\ &= \phi_{\text{rff}}(x)^T \left[ \frac{1}{N} \sum_{i=1}^N \phi_{\text{rff}}(x_i) \phi_{\text{rff}}(x_i)^T \right] \phi_{\text{rff}}(x)\end{aligned}$$

↓  
 $\rho$

$$\begin{aligned}k_{\sigma/2}^2(x_i, x) &= \left( e^{-\sigma/2 \|x_i - x\|} \right)^2 \\ &= e^{-\sigma \|x_i - x\|} \\ &= k_{\sigma}(x_i, x)\end{aligned}$$

## ④ Density Matrices

The state of a quantum system is represented by a vector  $\psi \in H$  ( $H$  is a Hilbert space, typically  $\mathbb{C}^n$ )

E.g. the spin of an electron  $\{\uparrow, \downarrow\}$

$$\psi = (\alpha, \beta) \quad |\alpha|^2 + |\beta|^2 = 1$$

Superposition: In general the a quantum state is a combination of basis states

$$\uparrow : (1, 0) \quad \downarrow : (0, 1) \quad \psi = \alpha \uparrow + \beta \downarrow$$

$|\alpha|^2$ : Probability of obtaining  $\uparrow$       $|\beta|^2$ : Probability of obtaining  $\downarrow$

Density Matrix: Representation of the state of a quantum system that can represent quantum uncertainty (superposition) and classical uncertainty.

$$\rho = \psi \psi^* = \begin{bmatrix} |\alpha|^2 & \alpha \beta^* \\ \beta^* \alpha & |\beta|^2 \end{bmatrix}$$

Pure

$$\rho = \sum_{i=1}^N p_i \psi_i \psi_i^*$$

Mixed

$$\sum_{i=1}^N p_i = 1$$

Two systems:  $\psi_1 = (\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}})$   $\rho_1 = \psi_1 \psi_1^* = \begin{bmatrix} \frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} \end{bmatrix}$

$\psi_2 = (1, 0)$   $\psi_2' = (0, 1)$   $\rho_2 = \frac{1}{2} \psi_2 \psi_2^* + \frac{1}{2} \psi_2' \psi_2'^* = \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{2} \end{bmatrix}$

Measurement probability (Born Rule):

$P(\varphi | \rho)$ : Given a system in state  $\rho$  the probability of measuring  $\varphi$

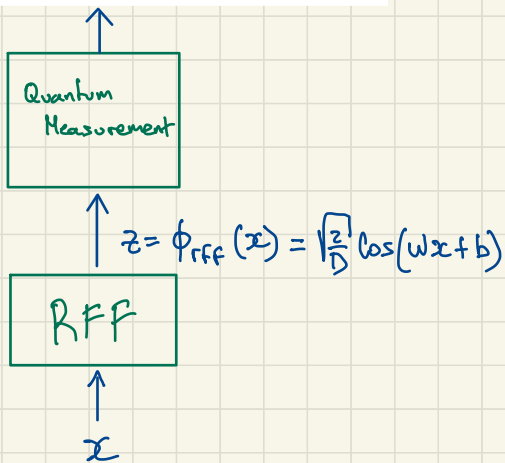
$$P(\varphi | \rho) = \text{Tr}(\rho \varphi \varphi^*) = \varphi^* \rho \varphi$$

$$\varphi = (\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}}) \quad P(\varphi | \rho_1) = 1 \quad P(\varphi | \rho_2) = \frac{1}{2}$$

# ⑤ Density Matrix Kernel Density Estimation (DMKDE)

## Prediction

$$\hat{f}_\rho(x) = \frac{\text{Tr}(\rho \phi_{\text{rff}}(x) \phi_{\text{rff}}(x)^*)}{Z} = \frac{\phi_{\text{rff}}(x)^* \rho \phi_{\text{rff}}(x)}{Z}, \quad (12)$$



## Time Complexity

Parzen's Estimator (KDE):  $O(dN)$

DMKDE :  $O(D^2)$

## Training

- Input. A sample set  $X = \{x_i\}_{i=1 \dots N}$  with  $x_i \in \mathbb{R}^d$ , parameters  $\gamma \in \mathbb{R}^+$ ,  $D \in \mathbb{N}$
- 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 2.1 for approximating a Gaussian kernel with parameters  $\gamma$  and  $D$ .
- Apply  $\phi_{\text{rff}}$  (eq. (1)):

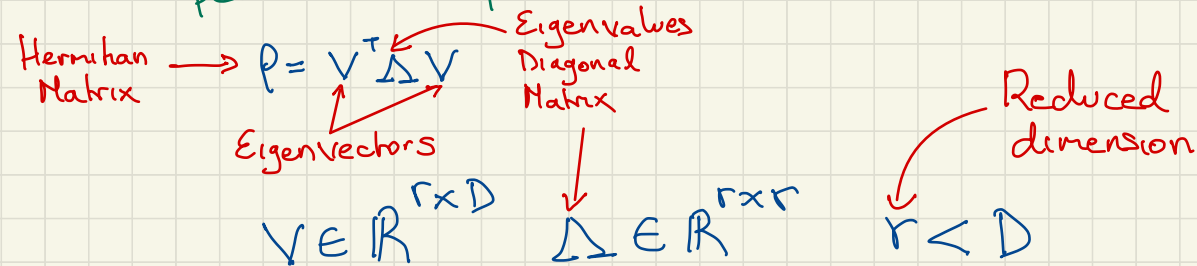
$$z_i = \phi_{\text{rff}}(x_i). \quad (10)$$

- Density matrix estimation:

$$\rho = \frac{1}{N} \sum_{i=1}^N z_i z_i^*, \quad (11)$$

## ⑥ Factorized DMKDE

Spectral Decomposition:



$$\hat{f}_\rho(x) = \frac{1}{Z} \|\Lambda^{\frac{1}{2}} V \phi_{\text{rff}}(x)\|^2$$

Time Complexity

$$\mathcal{O}(Dr)$$



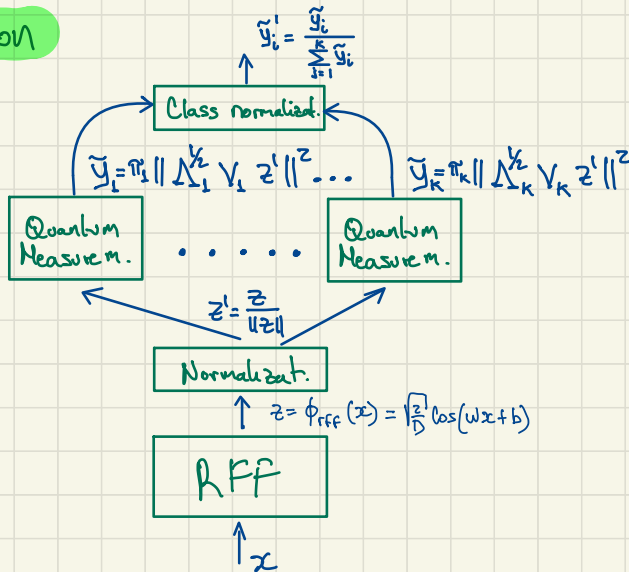
# ⑧ Density Matrix Kernel Density Classification (DMKDC)

Kernel density classification

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

Posterior probability →  $\hat{\Pr}(Y = j | X = x)$   
Prior →  $\pi_j$   
Density estimation →  $f_j(x)$  and  $\hat{f}_k(x)$

Prediction



Training

Density Matrix Estimation

1. Use the RFF method to calculate  $W_{\text{rff}}$  and  $b_{\text{rff}}$ .
2. For each class  $i$ :
  - (a) Estimate  $\pi_i$  as the relative frequency of the class  $i$  in the dataset.
  - (b) Estimate  $\rho_i$  using eq. (11) and the training samples from class  $i$ .
  - (c) Find a factorization of rank  $r$  of  $\rho_i$ :

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

Stochastic Gradient Descent

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

# 9 Quantum measurement classification (QMC)

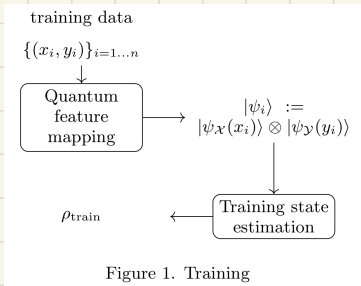
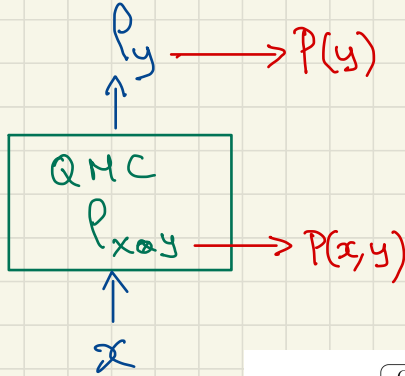


Figure 1. Training

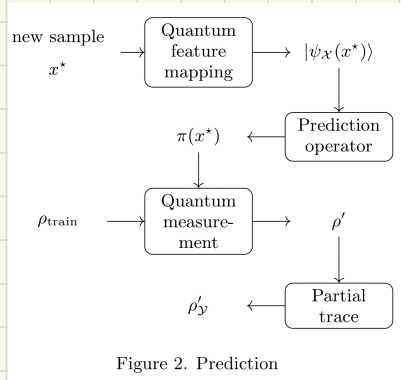


Figure 2. Prediction

Generalizes Bayesian Inference

**Proposition 1.** Let  $T = \{(x_i, y_i)\}_{i=1, \dots, n}$  be a set of training samples,  $x^*$  a sample to classify, with  $x_i, x^* \in \{1, \dots, m\}$  and  $y_i \in \{1, 2\}$ . Let  $\rho_{\text{train}}$  be the state calculated using the mixed state, eq. (8) or equivalently the classic mixture eq. (9), and a one-hot encoding feature map for both  $x_i$  and  $y_i$ . Then the diagonal elements of the density matrix  $\rho'_y$ , calculated using eq. (12) correspond to an estimation, using Bayesian inference, of the conditional probabilities  $P(y = i|x^*)$ :

$$\rho'_{y_i, i} = P(y = i|x^*) = \frac{P(x^*|y = i)P(y = i)}{P(x^*)}, \quad (13)$$

where  $P(x^*|y = i)$ ,  $P(y = i)$  and  $P(x^*)$  are estimated from  $T$ .

Can be seen as a kernel method

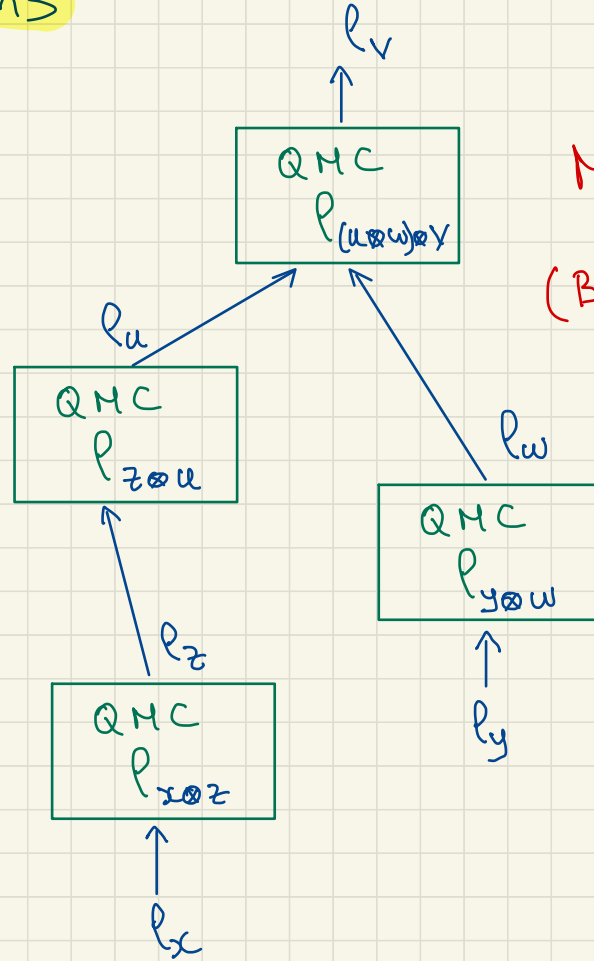
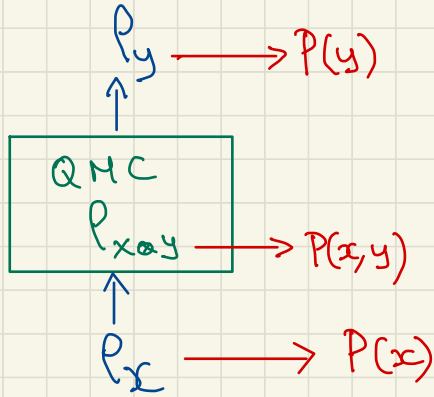
**Proposition 2.** Let  $T = \{(x_i, y_i)\}$  be a set of training samples,  $x^*$  a sample to classify, with  $x_i, x^* \in \mathcal{X}$  and  $y_i \in \mathcal{Y}$ . Let  $\rho_{\text{train}}$  be the state calculated using a mixed state (eq. (8)) and quantum feature maps  $\psi_{\mathcal{X}}$  and  $\psi_{\mathcal{Y}}$ . Then the density matrix  $\rho'_y$ , calculated with eq. (12), can be expressed as:

$$\rho'_y = \mathcal{M} \sum_{i=1}^N |k(x^*, x_i)|^2 |\psi_{\mathcal{Y}}(y_i)\rangle \langle \psi_{\mathcal{Y}}(y_i)|, \quad (14)$$

where  $k(x^*, x_i) = \langle \psi_{\mathcal{X}}(x^*) | \psi_{\mathcal{X}}(x_i) \rangle$  and  $\mathcal{M}^{-1} = \text{Tr}[\pi(x^*) \rho_{\text{train}} \pi(x^*)]$ .

- Generalizes DMKDC
- Produces a density matrix as output

# 10 Further Generalizations



Multilayer Model  
(Bayesian network)

# References

González, F. A., Vargas-Calderón, V., & Vinck-Posada, H. (2021). Classification with Quantum Measurements. *Journal of the Physical Society of Japan*, 90(4), 044002.  
<https://arxiv.org/pdf/2004.01227.pdf>

*González, F. A., Gallego, A., Toledo-Cortés, S., & Vargas-Calderón, V. (2021). Learning with Density Matrices and Random Features. arXiv preprint arXiv:2102.04394*  
<https://arxiv.org/abs/2102.04394>

<https://github.com/fagonzalezo/qmc>