# Polyadic Quantum Classifier

William Cappelletti, Rebecca Erbanni and Joaquín Keller

*Entropica Labs*, Singapore

{william, rebecca, joaquin}@entropicalabs.com

*Abstract*—We introduce here a supervised quantum machine learning algorithm for multi-class classification on NISQ architectures. A parametric quantum circuit is trained to output a specific bit string corresponding to the class of the input datapoint.

We train and test it on an IBMq 5-qubit quantum computer and the algorithm shows good accuracy —compared to a classical machine learning model— for ternary classification of the Iris dataset and an extension of the XOR problem.

Furthermore, we evaluate with simulations how the algorithm fares for a binary and a quaternary classification on resp. a known binary dataset and a synthetic dataset.

*Index Terms*—quantum machine learning, variational quantum algorithm, NISQ architecture, classification algorithm, supervised learning

## I. Introduction

Quantum machine learning (QML) has raised great expectations, it is thought [10] [11] to be one of the first possible applications of quantum computing to be able to run on NISQ[1] computers.

Nonetheless, QML is still in its infancy. We can make a parallel with the dawn of machine learning in the 1950s when the emblematic perceptron [12], [13] was introduced to solve binary classification problems. Today's research in QML has followed the same path and binary classification has been broadly studied.

In 1969 [14] it was shown that the perceptron could not solve the simple XOR problem. In fact it can only classify linearly separable datasets and it wasn't before 1986 that multilayer perceptrons with backpropagation [15] addressed harder problems. For instance, the ternary classification of the Iris flower dataset [16], which is non-linearly separable, could not be solved with the perceptron approach, but is now a typical test case [17, Ch. 2] of machine learning.

In this paper, we introduce a QML algorithm for multi-class classification and challenge it with the Iris flower dataset and a extension of the XOR problem with added Gaussian noise. The Iris flower dataset has three classes, two of which are not linearly separable. Binary QML classifications on this dataset have been addressed, on the linearly separable class against the rest, by [1], [2] and, pairwise, on all classes by [3].

Any n-ary classifier can be implemented with n binary classifiers [9] however, to predict all classes at once, we take a direct approach described in section IV.

There was no guarantee it would be possible to train our classifier on an actual quantum computer. However our simulations of the algorithm on the Iris dataset show that the

quantum computing power needed, for both training and test, can be found in today's hardware. And since models trained on simulators were tested on IBMq by [3] and [2], we were optimist that the inherent noise of the hardware wouldn't be an impediment. Indeed, we run our algorithm —train and test— on IBMq quantum hardware, for the ternary classification of the Iris dataset. Results are in section IX.

Similarly, in section X, we successfully trained a model for the Gaussian XOR problem on the same quantum system.

Other experiments, with a simulator, on binary and quaternary classifications suggest that our approach is flexible enough to be applied to many problems. The outcome of these experiments are in sections XI and XII.

As detailed in sections II and III, our algorithm is based on parametric quantum circuits. Our approach is mostly empirical and the exact design of these circuits relies on experiments and on some guidelines presented in sections VI and VII.

## II. A variational quantum algorithm

Like many quantum machine learning algorithms [1]–[8], our procedure for classification of classical data fits in the general scheme of variational quantum algorithms [18], where a parametric quantum computation $\mathfrak{F}_{\boldsymbol{\theta}}$ is applied to an input vector $\boldsymbol{x}$ to get the result $\hat{y} = \mathfrak{F}_{\boldsymbol{\theta}}(\boldsymbol{x})$.

The variational algorithm per se consists of running $\mathfrak{F}_{\boldsymbol{\theta}}(\boldsymbol{x})$ for different values of $\boldsymbol{\theta}$ and $\boldsymbol{x}$ to find an optimal $\boldsymbol{\theta}^{\star}$ for which the results are satisfactory.

Although some have envisioned hardware architectures with quantum random access memory [19] and other interesting features, our algorithm relies only on the simplest functionalities available in actual quantum hardware.

In most of today's quantum computers or QPUs[2], an elemental computation Q take as input $n$, a number of *shots*, and P, a quantum program, or circuit. A circuit P is a sequence of *gates*, or elemental operations on qubits.

In what is called a shot, the qubits of the quantum computer are all initialized at $|0\rangle$, the sequence P of gates is applied to the qubits, then they are all measured. A *run* or *circuit run* is a sequence of $n$ shots. Runs and shots are quantum computations of different granularity. The result $Q(n, P)$ of a circuit run is the sequence $\hat{R}$ of $n$ bit strings in $\{0, 1\}^{N}$, corresponding to the measurement of the $N$ qubits of the circuit.

We can note that, the computational time complexity of $Q(n, P)$ is $\mathcal{O}(n \times |P|)$ where $|P|$ is the number of gates in P.

---

[1]Noisy Intermediate-Scale Quantum

[2]Quantum Processing Unit

## III. How to Input Data

To input data we resort to a method called variational, or parametric, encoding which was first proposed in [8]. The method consists of using a parametric circuit to encode the input vector $x$ to parameters of the circuit.

A parametric circuit is a circuit where some gates can take continuous angles as parameters and are $2\pi$-periodic regarding them.

We define an encoding function $f$ to map each coordinate or *feature*, of the input vector $x$ to an angle in the interval $]-\pi, \pi[$, which gives us a vector of parameters $\omega = f(x)$ to be used in a parametric circuit $P_\omega$.

We note $X$ the set of input vectors $x$ from the learning dataset and define the vectors $\overline{X}$ and $\sigma_X$, as its element-wise mean and standard deviation. Similarly, we compute the element-wise standard score $z_X(x) = \frac{x - \overline{X}}{\sigma_X}$,



Fig. 1. The input vector $x$ is encoded as angles $\omega = f(x)$ of a parametric circuit $P_{\omega,\theta}$

In a Gaussian distribution approximation, we define the quantile $q = \Phi^{-1}(1 - \epsilon^{\frac{1}{d}}/2)$, where $d$ is the dimension of $X$ and $\Phi^{-1}$ the quantile function [20]. By definition, the points $x$ such that $\exists i \; |z_X(x)|_i > q$ represent less than an $\epsilon$ fraction of $X$. We fix $\epsilon$ small and ignore such points.

Thus, we define the encoding function as a simple linear rescale and shift of the input:

$$ f(x) = \left(1 - \frac{\alpha}{2}\right) \frac{\pi}{q} z_{X(x)}. $$

By definition, all angles of the encoded vector $\omega = f(x)$ fall within the interval $\left]-\left(1 - \frac{\alpha}{2}\right)\pi, \left(1 - \frac{\alpha}{2}\right)\pi\right[$. This enforces an angular gap $\alpha\pi$ between the extreme values of the encoded dataset, where $\alpha$ is a parameter to be choosen.

Using this encoding function $f$, we can now define our quantum classifier as

$$ \hat{y} = \mathfrak{F}_\theta(x) = g\left(\, \mathbb{Q}\left(n, P_{\omega,\theta}\right)\,\right), $$

where $P_{\omega,\theta}$ is a parametric quantum circuit, $\omega = f(x)$ are the encoded input parameters and $\theta$ the model parameters, i.e. those to optimize. The parameter $n$ is the number of shots, which here is not automatically learned but adjusted "by hand"; and $g$ is the postprocessing function, a classical computation —described in following sections.

We can note that as $f$, the input encoding function, is not parametric, the dataset is encoded only once, prior to the learning process.

## IV. How to Output Data

The output of a basic quantum computation, the run of a circuit $P_{\omega,\theta}$, is a sequence $\mathbb{Q}(n, P_{\omega,\theta}) = \hat{R}$ of $n$ bit strings in $\{0,1\}^N$ where $N$ is the number of quantum bits and $n$ the number of shots.

The outcome of each measurement is a bit string $s$ and from quantum mechanics we know that it follows an underlying probability distribution $P(s)$. We can estimate its probability by $\hat{P}(s) = \hat{C}(s)/n$, where $\hat{C}(s)$ is the number of occurrences of $s$ in $\hat{R}$.

Here, recalling that $\omega$ encodes for the input vector $x$, we use this estimated probability $\hat{P}(s)$ to predict $\hat{y}$ the class of $x$. In order to do so, we associate to each possible class $k$ a bit string $s_k$ in $\{0,1\}^N$ and we note $\hat{P}(s_k)$ as $\hat{P}_k$.

The output of our algorithm, the predicted class is

$$ \hat{y} = \mathfrak{F}_\theta(x) = g\left(\, \mathbb{Q}\left(n, P_{\omega,\theta}\right)\,\right) = \arg\max_k \hat{P}_k. $$

Or, saying it otherwise, the predicted class $\hat{y}$ is the class $k$ for which its associated bit string $s_k$ has the highest number of occurrences.

Note that, even though $\hat{P}_k$ varies from one run to another, if the relative difference between the theoretical $P_k$ associated with each class is high enough, $\arg\max_k$ will return a consistent value even with a low number of shots.

This approach is by essence polyadic, the number of classes being bounded by the total number of possible bit strings. As [7], it does not rely on multiple binary classifications [9].

## V. Learning

The learning process for our classifier consists in finding a set of parameters $\theta^\star$ for which the predictor $g\left(\, \mathbb{Q}\left(n, U_{f(x),\theta^\star}\right)\,\right)$ has high accuracy.

To do so, for each sample $(x, y)$ of the learning dataset $\mathcal{T}$, we define an adequate loss function which correlates with the error of the predictor:

$$ \mathcal{L}_\theta(x, y) = -\log \frac{e^{\hat{P}_y}}{\sum_{k \in K} e^{\hat{P}_k}}, $$

where $\hat{P}_y$ is the probability estimate of the bitstring associated with the real class $y$ of $x$.

The optimization process consists of iteratively exploring the space of parameters $\theta$ with an heuristic. At each step $i$ we compute the loss function $\mathcal{L}_{\theta^i}$ as the average over a random subset $\mathcal{B}_i$ of $\mathcal{T}$:

$$ \mathcal{L}_{\theta^i} = \sum_{(x,y) \in \mathcal{B}_i \subset \mathcal{T}} \frac{\mathcal{L}_{\theta^i}(x, y)}{|\mathcal{B}_i|}. $$

The subset, or *minibatch*, $\mathcal{B}_i$ could be a singleton, the full training set $\mathcal{T}$ or have any cardinality in between.

The result of the training, the vector of parameters $\theta^\star$, is the $\theta^i$ that gives the lowest value for the loss function.

## VI. Optimizing Parametric Quantum Circuits

To specify our circuits we choose an universal set of quantum gates: $s_x^i$ —the $\pi/2$ rotation about the Pauli-X axis, $Z_\phi^i$ —a rotation of an arbitrary angle $\phi$ about the Pauli-Z axis, and $Cz_j^i$ —the controlled-$Z$ gate, where $i$ and $j$ are the qubits to which they apply.

In a QPU, these elemental gates are implemented using pulses —physical phenomena that modify the state of the qubits. A controlled-$Z$ gate is translated in a 2-qubit pulse, a $s_x$ gate needs a 1-qubit pulse and $Z$-rotations need no pulse
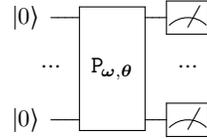
[21] as they are implemented by modifying the subsequent pulses.

Each QPU might have a slightly different set of elemental gates, by choice of the hardware manufacturer. Nonetheless, a circuit written using our set of elemental gates have a translation that preserves the number of 1-qubit pulses and 2-qubit pulses for most of the existing QPUs (see Appendix A).

A quantum program, or, circuit P specifies [22] a unitary matrix, that we note here $\breve{\text{P}}$. The unitary $\breve{\text{P}}$ transforms a quantum state into another one.

Qubits in a given quantum state $\psi$ yield, when measured, a bit string $s$ with probability $P(s) = |\psi_s|^2$, where $\psi_s$ is the amplitude associated with $s$. We note $\mathcal{M}$ the measurement operator $\mathcal{M} : \psi \mapsto s$. From the definition of $\mathcal{M}$ it follows that

$$\exists \alpha \in \mathbb{R} \quad U = e^{i\alpha}U' \quad \implies \quad \mathcal{M}U = \mathcal{M}U'.$$

If this holds we say that $U$ and $U'$ are *equal up to a global phase* and note $U \equiv U'$.

The purpose of a quantum computer is to implement the unitary $\breve{\text{P}}$ and the final measure operation. However, due to engineering limitations, the hardware is imperfect and the actual quantum operations differs from those specified in P. This difference or *noise* grows with the number of pulses and the time needed to run the circuit.

Hence, to minimize the noise, if two circuits P and P' specify for equivalent unitaries $\breve{\text{P}} \equiv \breve{\text{P}}'$ or yield the same result $\mathcal{M}\breve{\text{P}} = \mathcal{M}\breve{\text{P}}'$ we will choose the circuit minimizing the number of pulses.

Fig. 2. A controlled-$Z$ after initialization has no effect, see (3).

We can now list some simple properties on the unitaries of our elemental gates that will help in optimizing our circuits. Since the unitary $\widetilde{AB}$ of the concatenation $AB$ of quantum programs $A$ and $B$ is the composition of their unitaries $\breve{A}\breve{B}$ we will omit the $\breve{\ }$ symbol when no confusion is possible.

(1) As $Z_\phi^i Z_\lambda^i \equiv Z_{\phi+\lambda}^i$, in a circuit, two consecutive Z-rotations on the same qubit can be replaced by just one.
(2) As $Z_\phi^i|0\rangle^{\otimes N} \equiv |0\rangle^{\otimes N}$, a Z-rotation immediately after qubit initialization can removed.
(3) For any 1-qubit unitary $U^i$ applied to qubit $i$, as $Cz_j^i U^i|0\rangle^{\otimes N} \equiv U^i|0\rangle^{\otimes N}$, we can remove a $Cz$ applied to a qubit $j$ immediately after its initialization. See Fig. 2.
(4) As $\mathcal{M}Z_\phi U = \mathcal{M}U$, we can remove a Z-rotation applied to a qubit just before its measurement.
(5) Similarly, as $\mathcal{M}Cz_j^i U = \mathcal{M}U$, we can remove any controlled-Z followed by no gate before measurement.
(6) As $Cz_j^i Cz_j^i$ is equal to the identity, two consecutive controlled-Z on the same qubits can be removed.
(7) Any unitary $U^i$ on one qubit $i$ can be decomposed as $U^i \equiv Z_\phi^i s_x^i Z_\alpha^i s_x^i Z_\lambda^i$.
(8) Properties (7), (2), (6) and (1) imply that any qubit should have at most $s_x Z_\phi s_x$ between initialization and its first $Cz$ gate.
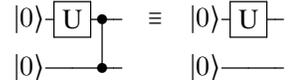
(9) Properties (7), (6) and (1) imply that any qubit should have at most a sequence $s_x Z_\phi s_x Z_\alpha$ between two $Cz$ gates
(10) Properties (7), (4) and (1) imply that any qubit should have at most $s_x Z_\phi s_x Z_\alpha$ after its last $Cz$ gate and before measurement.
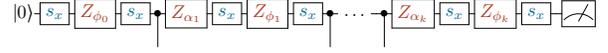
Fig. 3. Maximum succession of atomic gates on a single quantum bit.

Properties (8), (9) and (10) give a maximal sequence Any circuit with more gates than that can be optimized further.

## VII. DESIGNING CIRCUITS

The constraints for optimality leave a high level of freedom on how to design parametric circuits: less elemental gates are still possible; some angles can be fixed as constants; some parameters will encode input features, while others will be parameters to optimize, and a choice has to be made on how to entangle qubits.

Following intuitions and empirical evidences (some circuits perform better than others) we have identified a set of rules to design what we think are good parametric circuits for our algorithm. However as intuitions are often misguided and empirical evidences can be overthrown by new experiments, we hope and expect these rules to be challenged and/or expanded by further research work.

Two main criteria enter in the design of parametric circuits intended to be used as classifiers. First, we want to minimize the number of
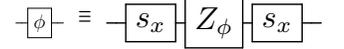
Fig. 4. Compact notation for $s_x Z_\phi s_x$.

gates. Second, we need enough parametric gates to encode the input vector and provide adequate learning capacity.

This double constraint –minimize the number of gates and maximize the number of parametric gates– should have led to keep all the $Z$-rotations as parametric gates. Yet, without fully understanding why, we noticed that a parametric $Z$-rotation right after a controlled-$Z$ seem not to add much capacity and tend to impair the learning phase. Hence, we opt for a rule enforcing exactly one parametric gate between entanglements, or more precisely, as shown in Fig. 5, a sequence $s_x Z_\phi s_x$ between any two consecutive controlled-$Z$ gates.
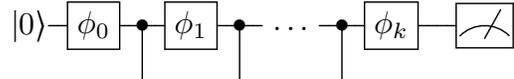
Fig. 5. Succession of operations on a quantum bit, only one parametric gate between entanglements.

If we extend this rule for individual qubits to the whole set of qubits, we have a step where each qubit is rotated by some angle followed by a step where qubits are entangled by pairs. Things are slightly, but not fundamentally, different depending whether the number of qubits is odd or even.

Although we don't have clear rules on how or why choose a given entangling pattern at each step, it is well understood [23]

that to unleash the power on quantum computing qubits need to be highly entangled.

At each step a quantum computer only allows 2-qubit gates between a limited set of qubit pairs, which is the coupling map or qubit connectivity graph.

Assuming that the qubit connectivity graph is connected, by choosing alternating entangling patterns (see Fig. 6) it is possible in few steps to ensure that all qubits are entangled pairwise. In designing our circuits, we will choose the succession of entangling patterns in such a way it minimizes the number of steps to reach a state where all qubits are entangled. Similar rules to design circuits are proposed in [24].
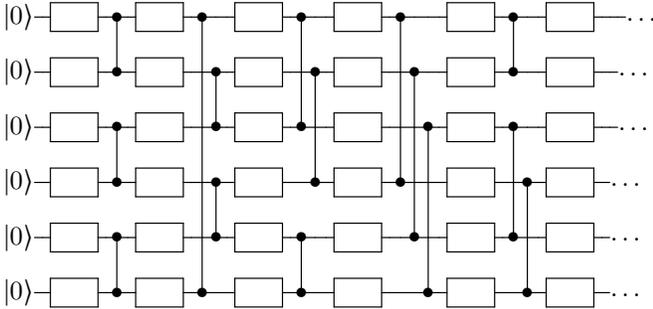


Fig. 6.    Alternating different entangling patterns on 6 fully connected qubits

Once the structure of the parametric circuit is fixed, the next decision we face is which parameters are going to be input parameters $\boldsymbol{\omega}$ (see section III) and which are going to be model parameters $\boldsymbol{\theta}$.

Although the model parameters $\boldsymbol{\theta}$ have a similar role and are indistinguishable, each of the input parameters $\boldsymbol{\omega}$ corresponds to features that have a meaning, so their position in the circuit matters. This choice is today made experimentally on learning performance.

Also we implement the idea of *data re-uploading*, suggested first by [7] and used in [6], which consists in inputting $\boldsymbol{\omega}$ more than once. Data re-uploading proves to be useful to add learning capacity without adding more qubits.

## VIII. EXPERIMENTAL PROTOCOL

To assess the capabilities of our algorithm, we carry out a series of experiments on various classification problems: the Iris flower dataset [16], the XOR problem with Gaussian noise, the skin segmentation dataset [25] and a four-classes synthetic dataset, generated using scikit-learn [26].

In each case, we follow the same experimental procedure; first we randomly split the dataset in a training and a testing subset, preserving the original ratio of each class. Second, we use the training dataset to find the parameters $\boldsymbol{\theta}^\star$ that minimize the loss function $\mathcal{L}_{\boldsymbol{\theta}}$. Finally, we use the test set to evaluate the prediction capabilities on independent data of $\mathfrak{F}_{\boldsymbol{\theta}^\star}$ the trained model (see [17, Ch. 7]).

In the learning phase, we randomly initialize the parameters $\boldsymbol{\theta}_0$ of the circuit and optimize them iteratively as explained in section V. The loss function takes the full training data set as minibatch.

We use the optimization package of SciPy [27], with two different heuristics, namely the algorithms BFGS [28] and COBYLA [29].

The BFGS algorithm uses the gradient of the loss function $\mathcal{L}_{\boldsymbol{\theta}}$ and moves the parameters accordingly. However, as $\mathcal{L}_{\boldsymbol{\theta}}$ depends on the random variable $\hat{P}_k$, it is random and its exact gradient is inaccessible. Nonetheless, with a high number of shots, the variance of $\hat{P}_k$ is low and we can resort to a finite differences estimate.

Moreover, for a small number of qubits, it is possible to simulate a quantum computer and get the theoretical $P_k$. We can then compute the *exact* loss function using $P_k$ instead of its estimate.

On the other hand, as COBYLA algorithm is gradient-free, it can work with a low number of shots. Then, COBYLA is the preferred method using a real QPU and BFGS when using exact probabilities from a simulated QPU.

Once our quantum model is trained, we compare its results on the test set to those obtained using a classical model. For this purpose, we consider the gradient boosting method XGBoost [30] —a state of the art model in many applications.

As a last note, we fix the values introduced in section III. In this paper, we use an angular gap $\alpha\pi$ of $\frac{\pi}{10}$ and fix the quantile $q$ to 3, so we ignore at most $\epsilon = 1\%$ of the points of every dataset[3].
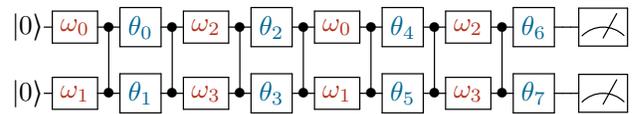
## IX. IRIS DATASET ON QPU



Fig. 7.    The quantum circuit for the classification of Iris dataset. Vector $\boldsymbol{\omega}$ encodes for (sepal length, sepal width, petal length, petal width). Note that the features are *re-uploaded*.

The Iris flower dataset [16], consists in 50 samples from each of three species of Iris (Iris setosa, Iris virginica and Iris versicolor), for a total of 150 data points, which we split in a train and test set of sizes 90 and 60 respectively. The input vector of each sample consists of four features, the length and the width of sepals and petals.

It is a well know test case for ternary classification, and, thanks to its relatively small sample size, we were able to run training and test entirely on an actual QPU, the IBMq-valencia 5-qubit quantum computer.

After trying several different circuits on simulated QPU, we choose to implement the four features in only two qubits. Fig. 7 shows our best performing circuit. The three classes, setosa, virginica and versicolor, are read respectively in the bit strings 00, 01 and 10, where the leftmost bit is the measurement of the upmost qubit.

Since IBMq-valencia has five qubits, to halve the number of QPU calls, we parallelize computations by running two

---

[3]Iris: $\epsilon = 1\%$; Skin: $\epsilon = 0.8\%$; Artificial: $\epsilon = 0.5\%$.

circuits at once on independent pairs of qubits, and we read the output accordingly.

At each learning iteration, to compute the loss function, we run the circuit for all 90 elements of the training set. For the first 20 iterations the number of shots per run is 250, then we increase it to 500 and from iteration 50 we raise it again to 750 shots.
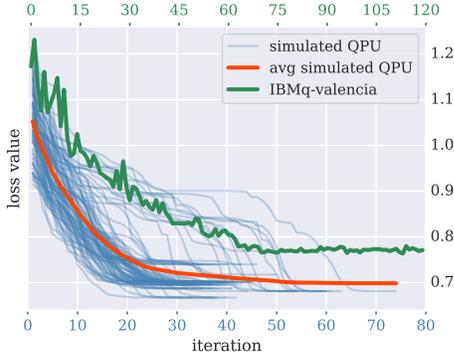


Fig. 8. Evolution of the loss function value during training of the quantum Iris classifier. The green line represents the training we did on IBMq-valencia with gradient-free optimizer COBYLA. Blue lines represent different trainings on simulated QPU using the exact loss function with gradient-based optimizer BFGS, and red line shows their average. Note that the number of iterations for the IBMq-valencia training have their own scale (0-120), shown on top.

We performed 120 optimization steps with COBYLA. Fig. 8 shows the evolution of the loss function, and we see that around iteration 80 it already stopped improving. The experiment took 5400 dual circuit runs, for a total of 3.26 millions shots, and spent approximately one hour and 20 minutes running on the QPU, as reported by IBM Quantum Experience platform.

We assess the resulting trained model over six different IBMq machines[4], with three runs each and 300 shots per run. Fig. 9 shows the confusion matrices and accuracies on the 60 points of the test set for both the quantum model and the classical XGBoost model, trained with the same data.

| QUANTUM MODEL TRAINED ON IBMQ-VALENCIA | | | |
|---|---|---|---|
| | **0** | **1** | **2** |
| **0** | 19.73 | 0.27 | 0.00 |
| **1** | 0.00 | 19.93 | 0.07 |
| **2** | 0.00 | 2.40 | 17.60 |

ACCURACY : $95.44 \pm 3.54\%$

| CLASSICAL MODEL | | | |
|---|---|---|---|
| | **0** | **1** | **2** |
| **0** | 20 | 0 | 0 |
| **1** | 0 | 19 | 1 |
| **2** | 0 | 2 | 18 |

ACCURACY : $95.00\%$

Fig. 9. Compared test scores between our quantum model trained on IBMq valencia and XGBoost a state-of-the-art classical model. We test the quantum model on six different IBMq machines with three runs each and the reported scores show the averages. In the confusion matrices, each entry $C_{i,j}$ is the number of observations actually in class $i$ (row), but predicted to be in class $j$ (column).

These results on quantum hardware are in line with the preliminary experiments with exact loss function on simulated

---

[4]essex, burlington, vigo, yorktown, london, ourense

QPU. We trained the quantum model a hundred times starting with different initial parameters and most of the time the training converges to close minimal values of the loss function (see Fig. 8). The learning process on IBMq-valencia however stagnated at higher values. Nonetheless, the best model on simulated QPU, the one with the lowest training loss value, did not get a higher accuracy when tested on the six IBMq machines.

## X. GAUSSIAN XOR ON QPU

The XOR problem is a decision problem on binary input consisting in learning the exclusive-or function. Given an input in $\{0,1\}^2$, the XOR function outputs 1 for $(0,1)$, $(1,0)$, and 0 for $(1,1)$, $(0,0)$. Minsky and Papert [14] showed that this problem is not linearly separable and cannot be learned by a perceptron [12].

To make the problem more interesting and challenging, we extend the input space to continuous values and generate a dataset to train a model. We take four symmetric points on
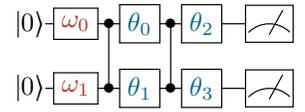


Fig. 10. Circuit to solve the Gaussian XOR problem, vector $\boldsymbol{\omega}$ is the input.

the Cartesian axes, and assign label 0 to the points on $x$-axis and label 1 to those on $y$-axis —this is an affine transformation of the original input points. Then, we use these points as centers for four 2-dimensional Gaussian distributions, from each of which we sample 20 datapoints, assigning them the corresponding label.

Figure 11 shows the scatter plot of this generated dataset, along with the distribution centers. We refer to this mixture of Gaussians as the *Gaussian XOR* problem.

Our quantum classifier uses the two-qubit circuit with four parameters shown in Fig. 10. The datapoints coordinates are encoded in $\boldsymbol{\omega}$ without preprocessing. We associate label 0 with bit string 00 and label 1 with bit string 10.

We train the model on IBMq-valencia with these 80 datapoints, proceeding as for the Iris dataset in section IX. The optimization process converged after 69 steps, using COBYLA, for a total of 2760 dual circuit runs and 1.37 million shots.

Since we know the generating distribution $\mathcal{D}$ of the dataset, we can compare our solution to the optimal-possible decision boundary. This, is given by the *Bayes classifier* [17, Sec. 2.4], which predicts the class directly from the posteriori distribution $\hat{y} = \arg\max_y P(y|\boldsymbol{\omega})$, where $P(y|\boldsymbol{\omega})$ is the conditional probability of label $y$ given observation $\boldsymbol{\omega}$.

The generalization error of the Bayes classifier for this distribution is 3.33%, so the best possible accuracy is 96.67%. Using one million newly generated datapoints we test our quantum model on a simulated QPU with 300 shots. The accuracy of this particular model —trained with 80 points on actual hardware— is 96.31%, not far from the theoretical maximum.

The decision boundaries of the Bayes classifier are the two perpendicular bisectors of the centers of the Gaussian
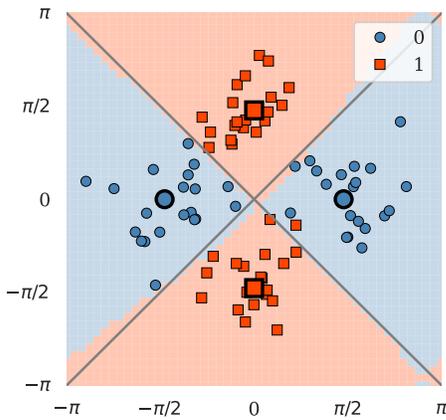
Fig. 11. Scatter plot of the training dataset generated from the Gaussian XOR problem. The points with bold borders are the centers of the Gaussian distributions. The background shows, color coded, the predicted class w.r.t. $\boldsymbol{\omega}$, while the grey lines represent the Bayes-optimal decision boundaries.

distributions. Figure 11 shows the decision boundaries of our quantum model along with the Bayes-optimal ones.

## XI. SKIN SEGMENTATION DATASET

The skin dataset consists in a large number of RGB-color values sampled from face images; each sample is labeled as skin or non-skin. We use a random subset of size 1000, with equally represented classes, which we further separate in a test and train set containing 400 and 600 observations respectively.

The dataset is not linearly separable (see Fig. 13) and we want to know if our algorithm could do the binary classification with a minimalist circuit. After several attempts we manage to get good results with the 3-qubit circuit in Fig. 12, with no data re-uploading and only



Fig. 12. The circuit for skin segmentation quantum model. Vector $\boldsymbol{\omega}$ encodes for (B,G,R).

3 entanglements and 6 parameters. We associate the classes with the permutation-invariant bit strings 000 and 111. We train and test the model with the simulated QPU using exact probabilities. The test set consists of 400 samples, with 200 examples per class. Our simple quantum model on 3-qubit scored a 94% accuracy, close to the 98% of the classical XGBoost model.
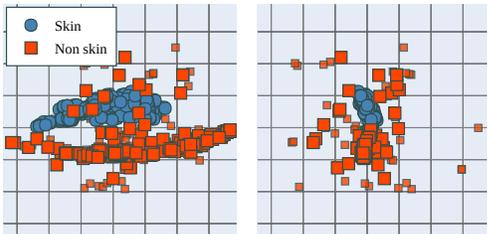


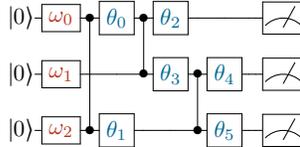Fig. 13. Skin segmentation dataset seen from two orthogonal angles.

## XII. SYNTHETIC DATASET WITH 4 CLASSES

In previous sections, we treated two real-life datasets, which are well known examples of binary and ternary classification. Nonetheless, our algorithm can, in principle, discriminate as many classes as bit strings the circuit can output. So a 2-qubit circuit should be able to solve a quaternary classification problem.

Then, we generate a 4-class bidimensional synthetic dataset of 5000 samples using the `make_classification` function from scikit-learn [26]. We randomly split the dataset in train and test sets representing respectively 60% and 40% of the samples, preserving each class ratio.

Out of many 2-qubit circuits, we choose a circuit with twelve parameters and four data uploadings as defined in Fig. 14. We assign each of the possible bit strings –00, 01, 10, 11– to a class.
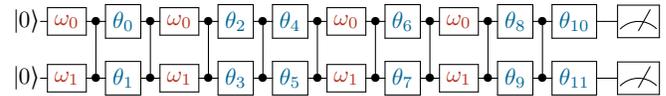


Fig. 14. The quantum circuit for quaternary classification of the two-dimensional synthetic dataset.

We train and test this quantum model using exact probabilities on simulated QPU. The average accuracy over a hundred experiments with different initial parameters is $85\pm2.8\%$, close to the $87.98 \pm 0\%$ of the classical XGBoost model.

Fig. 15 shows the prediction for the whole feature space of one the trained models. We note non trivial decision boundaries.

## XIII. FUTURE WORK AND CONCLUSION

We obtain satisfying quantum models for the challenging Iris flower dataset and the other polyadic problems, with good test scores compared with a classical model.

Furthermore, we were able to train the model for the Iris dataset and the Gaussian XOR problem on actual hardware. This was possible because of the low number of shots needed and the noise resilience of the model. Moreover, even though different machines have different noise patterns, the model trained on one machine performed well on all others (see Fig. 9).

It is important to remind that the optimization heuristics used in this work have severe limitations. The loss function is quantum computed and thus random; hence, we can only approximate its gradient to some extent, at the expense of more shots. Usual gradient based methods as BFGS perform poorly with this sort of randomness. In fact, we use them successfully on simulations, but cannot do the same on actual hardware. For QPUs we resort to the gradient-free COBYLA. However, it is well known [31] that COBYLA's performance drastically degrades with the dimension of the parameter space.

To improve the algorithm, we will need to explore optimization techniques adapted to quantum loss functions. One salient idea is to use circuit modifications to approximate the gradient, as proposed by [32] [33].
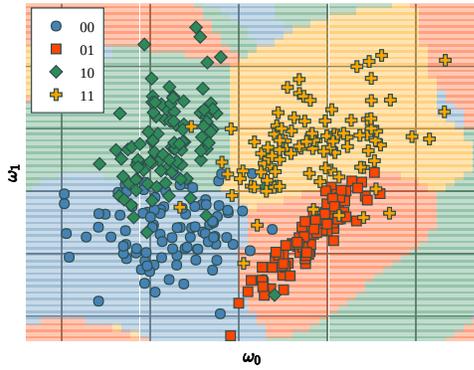
Fig. 15. Scatter plot of the synthetic dataset. The background shows, color coded, the predicted class w.r.t. $\omega_0$ and $\omega_1$.

Another avenue worth exploring is to work in a different parameter space. Here, we optimize directly the angles $\boldsymbol{\theta}$ of the parametric quantum circuit; instead, we could train the model using an intermediate parameter vector $\boldsymbol{\rho}$ mapped to $\boldsymbol{\theta}$. For instance, a non-linear mapping would reshape the optimization landscape. Also, this would allow constraints on $\boldsymbol{\theta}$, in particular parameter sharing, as in convolutional neural networks [34].

However, when addressing higher dimensional data and increasing the number of qubits, the most challenging issue becomes circuit design. We cannot rely anymore on sheer guess-and-check and need to understand better why some circuits outperform others.

As the computational capability of quantum hardware keeps increasing, we expect the empirical study of more complex quantum machine learning models to give new insights and lead to increasingly better algorithms.

## APPENDIX A
### TRANSLATING TO HARDWARE SPECIFIC NATIVE GATES

Quantum hardware manufacturers implement different sets of native gates to specify the circuits to run on their machines. In section VI we define our own set of gates to specify circuits: the 0-pulse $Z$-rotation and the 1-pulse $s_x$ one-qubit gates and the controlled-$Z$ two-qubit gate.

The one-qubit gates are native in IBM [35], Rigetti [36], Honeywell [37] and Google [38] quantum architectures and implemented with the same number of qubits.

The controlled-$Z$ is native for Rigetti's QPU and for most of Google's —with the notable exception of the Sycamore QPU. The tunable two-qubit gate [39] of the Sycamore processor has not been studied in this paper.

Honeywell uses the two-qubit $Zz$ gate and $Cz_j^i$ translates in $Z_{-\pi/2}^i Z_{-\pi/2}^j Zz_j^i$ which preserves the number of pulses.

IBM uses the controlled-not as two-qubit gate, noted $C_j^i$. Controlled-$Z$ can be translated in $H_j C_j^i H_j$, where $H$ is the Hadamard gate. With the following translations

$$
\begin{aligned}
H s_x Z_\phi s_x H &= s_x Z_{\pi - \phi} s_x \\
H s_x Z_\phi s_x &= Z_\pi s_x Z_{\phi - \frac{\pi}{2}} s_x \\
s_x Z_\phi s_x H &= s_x Z_{\phi - \frac{\pi}{2}} s_x Z_\pi
\end{aligned}
$$

the Hadamard gates disappear, thus preserving the number of pulses.

## REFERENCES

[1] A. Daskin, "A simple quantum neural net with a periodic activation function," 2018, arXiv:1804.07633.

[2] E. Grant, M. Benedetti et al., "Hierarchical quantum classifiers," npj Quantum Information, vol. 4, Dec. 2018, arXiv:1804.03680.

[3] M. Schuld, M. Fingerhuth, and F. Petruccione, "Implementing a distance-based classifier with a quantum interference circuit," Europhysics Letters, vol. 119, Sep. 2017, arXiv:1703.10793.

[4] E. Farhi and H. Neven, "Classification with quantum neural networks on near term processors," Feb. 2018, arXiv:1802.06002.

[5] M. Schuld, A. Bocharov et al., "Circuit-centric quantum classifiers," Physical Review A, vol. 101, Mar. 2020, arXiv:1804.00633.

[6] S. Lloyd, M. Schuld et al., "Quantum embeddings for machine learning," Feb. 2020, arXiv:2001.03622.

[7] A. Prez-Salinas, A. Cervera-Lierta et al., "Data re-uploading for a universal quantum classifier," Quantum, vol. 4, Feb. 2020, arXiv:1907.02085.

[8] J. Romero and A. Aspuru-Guzik, "Variational quantum generators: Generative adversarial quantum machine learning for continuous distributions," Jan. 2019, arXiv:1901.00848.

[9] M. Aly, "Survey on multiclass classification methods," Neural Netw, vol. 19, 2005, doi:10.1.1.175.107.

[10] J. Preskill, "Quantum computing in the nisq era and beyond," Quantum, vol. 2, Aug. 2018, arXiv:1801.00862.

[11] M. Broughton, G. Verdon et al., "Tensorflow quantum: A software framework for quantum machine learning," Mar. 2020, arXiv:2003.02989.

[12] F. Rosenblatt, "The perceptron: a probabilistic model for information storage and organization in the brain." Psychological review, vol. 65, 1958, doi:10.1037/h0042519.

[13] D. Cardon, J.-P. Cointet, and A. Mazieres, "Neurons spike back: The invention of inductive machines and the artificial intelligence controversy," Reseaux, vol. 36, 2018, doi:10.3917/res.211.0173.

[14] M. Minsky and S. A. Papert, Perceptrons: An introduction to computational geometry. MIT press, 1969.

[15] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning representations by back-propagating errors," nature, vol. 323, 1986.

[16] R. A. Fisher, "The use of multiple measurements in taxonomic problems," Annals of eugenics, vol. 7, 1936, hdl:2440/15227.

[17] J. Friedman, T. Hastie, and R. Tibshirani, The elements of statistical learning. Springer series in statistics New York, 2001, vol. 1, no. 10.

[18] A. Peruzzo, J. McClean et al., "A variational eigenvalue solver on a photonic quantum processor," Nature Communications, vol. 5, no. 1, Jul. 2014, arXiv1304.3061.

[19] V. Giovannetti, S. Lloyd, and L. Maccone, "Quantum random access memory," Phys. Rev. Lett., vol. 100, Apr. 2008, arXiv:0708.1879.

[20] L. Wasserman, All of statistics: a concise course in statistical inference. Springer Science & Business Media, 2013.

[21] D. C. McKay, C. J. Wood et al., "Efficient z-gates for quantum computing," 2016, arXiv:1612.00858.

[22] M. A. Nielsen and I. L. Chuang, Quantum Computation and Quantum Information. Cambridge University Press, 2000, online.

[23] R. Jozsa and N. Linden, "On the role of entanglement in quantum-computational speed-up," Proceedings of the Royal Society of London. Series A: Mathematical, Physical and Engineering Sciences, vol. 459, no. 2036, Aug 2003, arXiv:quant-ph/0201143. [Online]. Available: http://dx.doi.org/10.1098/rspa.2002.1097

[24] M. Cerezo, A. Sone et al., "Cost-function-dependent barren plateaus in shallow quantum neural networks," Jan. 2020, arXiv:2001.00550.

[25] R. Bhatt and A. Dhall, "Skin segmentation dataset," UCI Machine Learning Repository, 2009, archive.ics.uci.edu/ml/datasets/skin+segmentation.

[26] F. Pedregosa, G. Varoquaux *et al.*, "Scikit-learn: Machine learning in Python," *Journal of Machine Learning Research*, vol. 12, 2011, scikit-learn.org.

[27] P. Virtanen, R. Gommers *et al.*, "SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python," *Nature Methods*, vol. 17, 2020, scipy.org.

[28] R. Byrd, P. Lu *et al.*, "A limited memory algorithm for bound constrained optimization," *SIAM Journal on Scientific Computing*, vol. 16, 1995.

[29] M. J. Powell, "A direct search optimization method that models the objective and constraint functions by linear interpolation," 1994.

[30] T. Chen and C. Guestrin, "XGBoost: A scalable tree boosting system," in *Proceedings of the 22nd ACM SIGKDD*, 2016, doi:10.1145/2939672.2939785.

[31] M. J. Powell, "A view of algorithms for optimization without derivatives," *Mathematics Today-Bulletin of the Institute of Mathematics and its Applications*, vol. 43, no. 5, 2007, doi:10.1.1.591.6481.

[32] G. E. Crooks, "Gradients of parameterized quantum gates using the parameter-shift rule and gate decomposition," May 2019, arXiv:1905.13311.

[33] M. Schuld, V. Bergholm *et al.*, "Evaluating analytic gradients on quantum hardware," Nov. 2018, arXiv:1811.11184.

[34] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016, deeplearningbook.org.

[35] H. Abraham, A. Offei *et al.*, "Qiskit: An open-source framework for quantum computing," 2019, qiskit.org.

[36] R. S. Smith, M. J. Curtis, and W. J. Zeng, "A practical quantum instruction set architecture," 2016, arXiv:1608.03355.

[37] J. M. Pino, J. M. Dreiling *et al.*, "Demonstration of the qccd trapped-ion quantum computer architecture," 2020, arXiv:2003.01293.

[38] "Cirq, a python framework for creating, editing, and invoking noisy intermediate scale quantum (nisq) circuits," github.com/quantumlib/Cirq.

[39] F. Arute, K. Arya *et al.*, "Quantum supremacy using a programmable superconducting processor," *Nature*, vol. 574, 2019, doi:10.1038/s41586-019-1666-5.