

Contents

1	Introduction	3
2	Classical Support Vector Machine	5
2.1	General overview	5
2.2	Support vector machine	6
2.2.1	Linear SVM	7
2.2.2	Kernel SVM	10
3	Quantum Support Vector Machine	17
3.1	General overview	17
3.2	Quantum feature map and Quantum Kernels	17
3.3	Quantum encoding circuits	20
3.4	QSVM algorithm and Qiskit implementation	26
3.5	QSVM potential	30
3.6	Issues of QSVM	33
3.7	Quantum kernel alignment	36
3.7.1	Circuit optimization	36
3.7.2	Performance	40
4	Genetic algorithm	43
4.1	General introduction	43
4.2	Genetic algorithm for QSVM	44
4.3	Results of the Experiments	53
4.3.1	Circle dataset	53

4.3.2	Moon dataset	59
5	Support Vector Regressor	62
5.1	Results of the Experiments	63
6	Conclusions and Directions for Future Work	66

1 Introduction

Quantum machine learning (QML) has recently gained attention as an emerging field of research in quantum information science, combining the principles of quantum computing and machine learning. [1] The main goal of QML is to leverage quantum phenomena like superposition, entanglement and interference to perform machine learning tasks. In contrast to traditional machine learning techniques, QML algorithms seem to offer potential benefits, as they can solve complex problems more efficiently and cost-effectively.

One of the most promising and most researched QML method at present is the Quantum Support Vector Machine (QSVM) [8] [18], an extension of the classical Support Vector Machine (SVM), a machine learning algorithm for classification and regression tasks. The classical SVM is one of the most used and popular models, that has been successfully applied to various fields, such as image recognition, text classification and bioinformatics. [4] However, as the complexity and the size of the dataset increases, classical SVMs are no longer able to perform efficiently. QSVM aims to overcome these limitations by leveraging the principles of quantum computing, constructing a quantum circuit that maps classical data into an exponentially large quantum Hilbert space, in this way hoping to achieve an exponential speedup.

However, the design of such a circuit is often done manually, following for example standard practices and patterns, or some sort of rule of thumb. This approach may work for simple problems, but one quickly finds out that using hand-crafted circuits for non-trivial tasks leads to a bad performance of the

model, as finding the correct ansatz by hand is too difficult. It is necessary to automate this process, i.e. we need an algorithm that designs a properly working quantum circuit from scratch.

One way to attack this problem is by using a genetic algorithm. Genetic algorithms are meta-heuristic optimization algorithms inspired by the process of natural selection. In this work we develop a genetic algorithm that automatically designs, from first principles, the circuit to use in a QSVM classification, and we test it on various datasets.

The algorithm was implemented exploiting the Python library Qiskit, for what regards the quantum device simulation, and scikit-learn, for the machine learning part. Qiskit is an open source Software Development Kit for working with quantum computers at the level of extended quantum circuits, operators, and primitives.

This work is organized as follows: in Chapter 2 we introduce how a classical SVM works, in particular we present the fundamental concept of kernelized SVM. In Chapter 3 we discuss the QSVM, we define what it is and see how it performs on some datasets. We will encounter the tedious problem of choosing the quantum embedding circuit. In Chapter 4 we introduce our genetic algorithm, and test its performance on several datasets. In Chapter 5 we draw conclusions and discuss possible future developments.

2 Classical Support Vector Machine

2.1 General overview

Machine learning (ML) is a subset of artificial intelligence, focused on developing algorithms that enable computers to learn from and make decisions based on data. A ML algorithm takes a set of N samples of data, called *dataset*, and automatically predicts properties of unknown data, without the need to give it explicit instructions. The dataset is often split into a *training set*, on which the algorithm is trained on, and a *test set*, utilized to evaluate the performance of the model. Data can be multidimensional, and each dimension is usually called *feature*.

ML algorithms can be classified in different categories:

- **Supervised learning:** each data point has a label that the model aims to predict. The label can be:
 - Discrete, meaning the data can be separated into two or more classes and the goal is to use the training data to predict the labels for the test data. This task is called *classification*. An example is identifying whether an image contains a dog or a cat, or diagnosing a medical condition as positive or negative based on patient records.
 - Continuous, when each data is labelled by a continuous value, meaning the target variable can take on any value within a range and the model attempts to predict this continuous value. This

is called *regression*. For example, predicting the future price of a house based on its features is a regression problem.

- **Unsupervised learning:** the dataset does not have labels, yet we would like to empirically determine something about the data. For example we may want to determine if the data can be represented as belonging to distinct groups (clustering).

In this work we will focus on the Support Vector Machine [4], a supervised classification algorithm.

Datasets vary from small simple dataset with 2 numerical features, up to large and complex datasets with thousands of numerical and categorical features. In our case, due to the limited computational power we possess, we will study only small 2D datasets.

2.2 Support vector machine

The Support Vector Machine (SVM) [4] is a binary classification algorithm, whose goal is to build the maximum margin separator between the two classes, that is, the separator that maximizes the distance of the closest point from each class. The standard SVM algorithm is a linear algorithm, so in particular it will try and build the separating margin as a hyperplane in a d -dimensional space (hence a $(d - 1)$ -dimensional plane), where d is the number of features. The points that touch the margin, or that are on the wrong side of it, are called support vectors. The distance between the decision boundary and the support vectors is called margin. The algorithm will

find the biggest possible margin.

2.2.1 Linear SVM

Let's start assuming that the classes are linearly separable. We are provided with a dataset with N d -dimensional instances $\{\mathbf{x}_i\}_{i=0,\dots,N-1}$. The two classes will be labelled with

$$y \in \{-1, 1\}.$$

The margin will be the set of points

$$\{\mathbf{x} \in \mathbb{R}^d : w_0 + \mathbf{w}^T \cdot \mathbf{x} = 0\}, \quad (1)$$

for appropriate parameters $w_0 \in \mathbb{R}$ and $\mathbf{w} \in \mathbb{R}^d$, which define an hyperplane and must be found by the algorithm. Now we have to find

$$\max_{w_0, \mathbf{w}}(m),$$

with the constraint

$$\frac{1}{\|\mathbf{w}\|} y_i (w_0 + \mathbf{w}^T \cdot \mathbf{x}_i) \geq m, \forall i = 0, \dots, N-1, \quad (2)$$

that can be rewritten as

$$y_i (w_0 + \mathbf{w}^T \cdot \mathbf{x}_i) \geq m \|\mathbf{w}\|, \forall i = 0, \dots, N-1. \quad (3)$$

The constraint prevents data points from falling into the margin. Rescaling \mathbf{w} up to a multiplicative factor does not change the hyperplane it defines, so for convenience we can choose its norm such that

$$\|\mathbf{w}\| = \frac{1}{m}. \quad (4)$$

Therefore, the problem becomes minimizing

$$||\mathbf{w}||,$$

with the constraint

$$y_i(w_0 + \mathbf{w}^T \cdot \mathbf{x}_i) \geq 1, \forall i = 0, \dots, N-1. \quad (5)$$

In the theory of convex optimization one can solve for the Lagrangian dual of this problem. We can introduce the dual variables α_i such that

$$\mathbf{w} = \sum_{i=0}^{N-1} \alpha_i y_i \mathbf{x}_i. \quad (6)$$

One would obtain that the dual problem consists in maximizing, with respect to the weight vector $\alpha \in \mathbb{R}^N$, the expression

$$f(\alpha_0, \dots, \alpha_{N-1}) = \sum_{i=0}^{N-1} \alpha_i - \frac{1}{2} \sum_{i=0}^{N-1} \sum_{j=0}^{N-1} \alpha_i \alpha_j y_i y_j (\mathbf{x}_i, \mathbf{x}_j), \quad (7)$$

with the constraint

$$\alpha_i \geq 0, \quad \forall i = 0, \dots, N-1, \quad (8)$$

$$\sum_{i=0}^{N-1} \alpha_i y_i = 0. \quad (9)$$

The constraint in eq. (9) ensures that the hyperplane is not overly influenced by one of the two classes, like for example if one class has more data points than the other, avoiding overfitting and fairly considering the contribution of both. In eq. (7) we indicated as $(\mathbf{x}_i, \mathbf{x}_j)$ the standard dot product of \mathbb{R}^d , explicitly

$$(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \cdot \mathbf{x}_j = \sum_{k=0}^{d-1} (\mathbf{x}_i)_k (\mathbf{x}_j)_k.$$

Eq. (7) defines a quadratic programming optimization problem, therefore the global maximum of f can be efficiently found in the context of convex analysis. The parameter w_0 can be found by imposing that, for a support vector

$$y_i(\mathbf{w}^T \mathbf{x}_i + w_0) = 1,$$

that yields, considering that $y_i = \pm 1$ and therefore $y_i^{-1} = y_i$,

$$w_0 = y_i - \mathbf{w}^T \mathbf{x}_i. \quad (10)$$

Once the optimal α_i have been found, given a new instance $\tilde{\mathbf{x}}$, according to eq. (6) and eq. (1), we can predict its class calculating

$$\text{sign} \left[\sum_{i=0}^{N-1} \alpha_i y_i(\tilde{\mathbf{x}}, \mathbf{x}_i) + w_0 \right]. \quad (11)$$

What we just described is the so-called hard margin SVM [4], because we did not allow points to fall inside the margin. One could relax this assumption, modifying the constraint in eq. (5) into

$$y_i(w_0 + \mathbf{w}^T \cdot \mathbf{x}_i) \geq 1 - \xi_i, \quad \forall i = 0, \dots, N-1, \quad (12)$$

where we introduced the slack variables ξ_i . We limit the softness of the margin by setting a positive constant C such that

$$\begin{aligned} \xi &\geq 0, \\ \sum_{i=0}^{N-1} \xi_i &\leq C. \end{aligned} \quad (13)$$

This is called soft margin SVM.

Scikit-learn [16] provides a straightforward implementation of the SVM algorithm, which we can use to observe the algorithm in action through an example. We use a mock dataset with 2 features, so we can easily print the data, the decision boundary and the margin.

```
from sklearn.svm import SVC
from sklearn.datasets import make_blobs

X,y = make_blobs(n_samples=200) #create mock dataset
svm = SVC(kernel='linear', C=1) #create svm
svm.fit(X, y) #fit the svm
```

The result of the fit is shown in Figure 1. We can observe how the algorithm built the largest possible margin.

2.2.2 Kernel SVM

We now need to address the issue of dealing with a highly non-linearly separable dataset. Let's consider as an example another mock dataset, shown in Figure 2. It is clear that in this case we cannot use the SVM algorithm in its basic form, not even with a soft margin. We must introduce the idea of kernelization. Let's introduce a function, called feature map, which projects the data in a higher dimensional space. That means a function

$$\phi : \mathbb{R}^d \rightarrow \mathbb{R}^D, \quad (14)$$

with $D > d$. The codomain of the feature map is called feature space. If we choose a suitable feature map we can hope to obtain a linearly separable

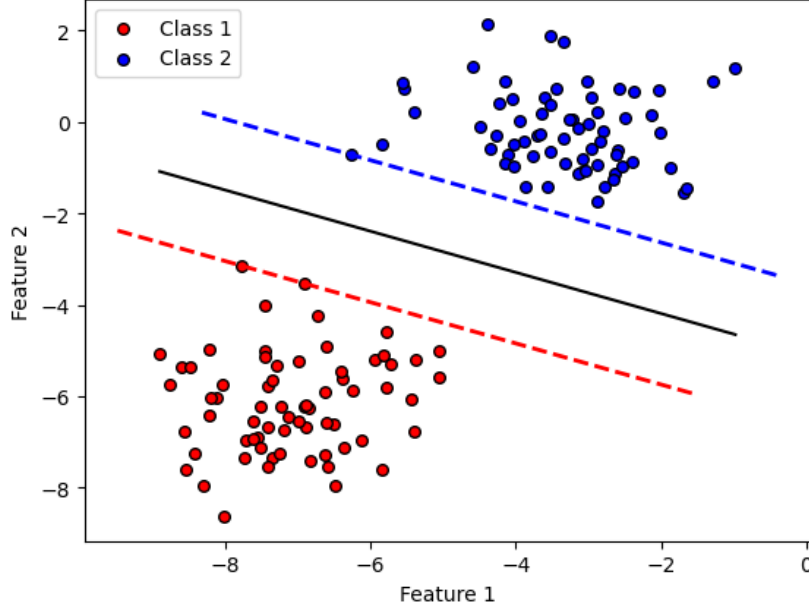


Figure 1: SVM decision boundary and margin border, fitted on a 2 feature mock dataset of 200 instances.

dataset in the feature space. The choice of the feature map is completely arbitrary, as long as it is a bijective function. Therefore, in principle, each time we are given a dataset we must choose an appropriate feature map for this strategy to work. For our example let's consider the feature map

$$\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3,$$

$$\begin{pmatrix} x_0 \\ x_1 \end{pmatrix} \mapsto \begin{pmatrix} x_0^2 \\ x_1^2 \\ \sqrt{2}x_0x_1 \end{pmatrix}. \quad (15)$$

The data of Figure 2 after the application of the feature map ϕ are represented in Figure 3. The dataset is now linearly separable in the feature space, so

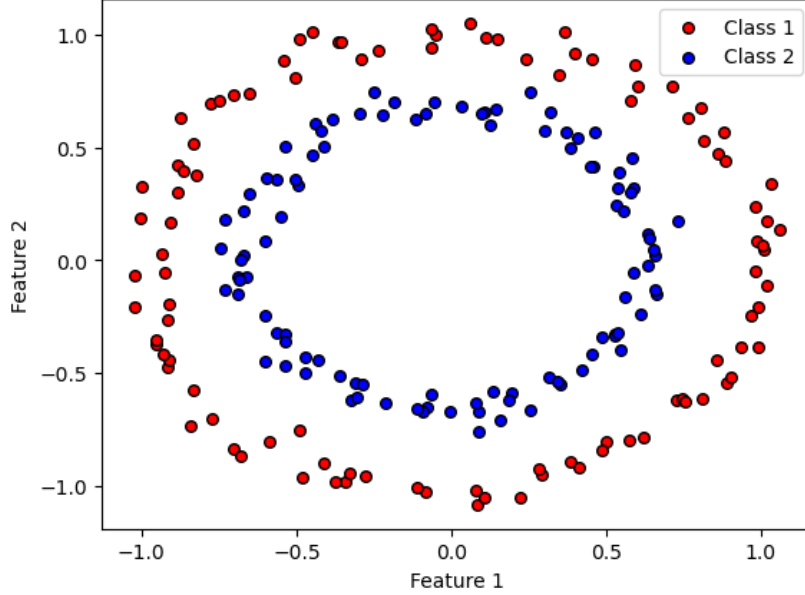


Figure 2: Higly non-linear 2 feature mock dataset with 200 instances.

the strategy worked. We can now apply the SVM algorithm in this space. Consider the two central equations of the algorithm: equation (7), which provides the expression to maximize in order to find the margin, and equation (11), which gives the rule for predicting the class of a new instance. These two equations are now modified into

$$f(\alpha_0, \dots, \alpha_{N-1}) = \sum_{i=0}^{N-1} \alpha_i - \frac{1}{2} \sum_{i=0}^{N-1} \sum_{j=0}^{N-1} \alpha_i \alpha_j y_i y_j (\phi(\mathbf{x}_i), \phi(\mathbf{x}_j)), \quad (16)$$

$$\text{sign} \left[\sum_{i=0}^{N-1} \alpha_i y_i (\phi(\tilde{\mathbf{x}}), \phi(\mathbf{x}_i)) + w_0 \right]. \quad (17)$$

A crucial observation is that in these two expressions only the scalar product of the feature map values appears. Therefore, we can conclude that the specific form of the feature map is not important, but rather the scalar product

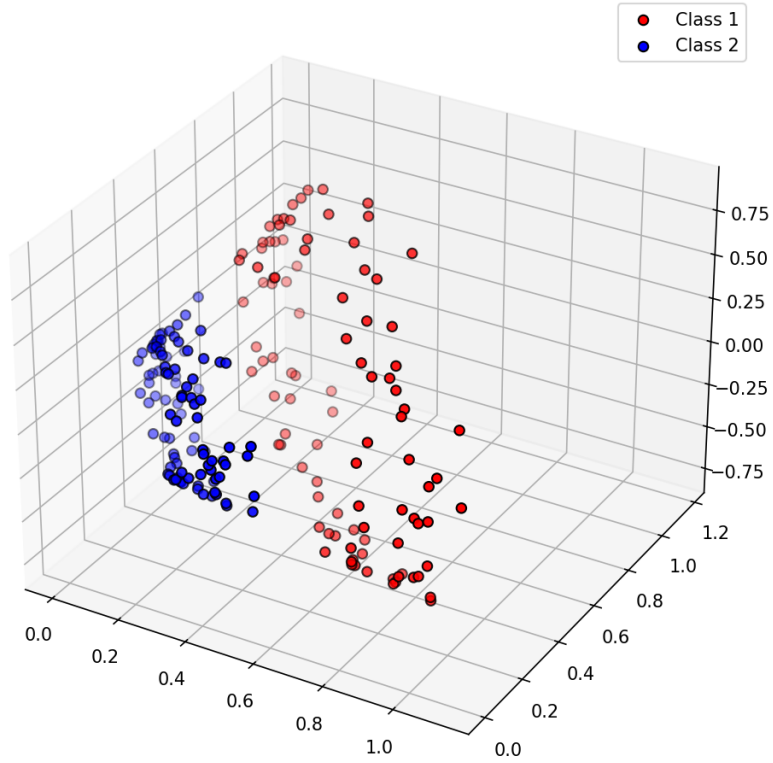


Figure 3: Highly non-linear mock dataset in the feature space after the application of the feature map ϕ . We observe now that the dataset is linearly separable.

it produces. We can define the kernel K as

$$\begin{aligned}
 K : \mathbb{R}^d \times \mathbb{R}^d &\rightarrow \mathbb{R}, \\
 \mathbf{x}, \mathbf{y} &\mapsto (\phi(\mathbf{x}), \phi(\mathbf{y})),
 \end{aligned} \tag{18}$$

where $(\phi(\mathbf{x}), \phi(\mathbf{y})) = \phi(\mathbf{x})^T \cdot \phi(\mathbf{y})$ is the standard scalar product of \mathbb{R}^D . Eq. (16) and eq. (17) now become

$$f(\alpha_0, \dots, \alpha_{N-1}) = \sum_{i=0}^{N-1} \alpha_i - \frac{1}{2} \sum_{i=0}^{N-1} \sum_{j=0}^{N-1} \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j), \quad (19)$$

$$\text{sign} \left[\sum_{i=0}^{N-1} \alpha_i y_i K(\tilde{\mathbf{x}}, \mathbf{x}_i) + w_0 \right]. \quad (20)$$

We see explicitly that the only quantity that matters is the kernel K . In our specific example the value of the kernel is

$$K(\mathbf{x}, \mathbf{y}) = \begin{pmatrix} x_0^2 & x_1^2 & \sqrt{2}x_0x_1 \end{pmatrix} \cdot \begin{pmatrix} y_0^2 \\ y_1^2 \\ \sqrt{2}y_0y_1 \end{pmatrix} = (\mathbf{x}^T \cdot \mathbf{y})^2. \quad (21)$$

Therefore, once we are given a dataset it is sufficient for us to choose an appropriate kernel, and forget about the feature map. Once the kernel has been chosen the SVM can be trained using eq. (19), and we can use it to predict a new class using eq. (20). There are some properties that the kernel must satisfy:

- The kernel must be symmetric, that is

$$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d, K(\mathbf{x}, \mathbf{y}) = K(\mathbf{y}, \mathbf{x}).$$

- The kernel must be positive definite, that is

$$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d, K(\mathbf{x}, \mathbf{y}) \geq 0.$$

Common choices of kernels are

- Linear kernel:

$$K(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \cdot \mathbf{y}.$$

This goes back to the standard SVM we used in Figure 1. It is suitable only for linearly separable (or close to, using soft margin) datasets.

- Polynomial kernel:

$$K(\mathbf{x}, \mathbf{y}) = (\gamma \mathbf{x}^T \cdot \mathbf{y} + c)^\delta.$$

For $c = 0$, $\gamma = 1$ and $\delta = 2$ we obtain the kernel of eq. (21).

- Gaussian kernel:

$$K(\mathbf{x}, \mathbf{y}) = \exp(-\gamma \|\mathbf{x} - \mathbf{y}\|).$$

This is also known as Radial Basis Function (RBF) kernel.

- Sigmoid kernel:

$$K(\mathbf{x}, \mathbf{y}) = \tanh(\mathbf{x}^T \cdot \mathbf{y} + c).$$

Scikit learn offers an easy way to easily implement all these common kernels.

For example, the kernel in eq. (21) can be implemented as

```
svm = SVC(kernel='poly', degree=2, gamma=1, coef0=0)
```

One can also create a custom kernel, passing as an argument a callable function to be used to calculate the kernel. Fitting this SVC function to the non-linear dataset of Figure 2 yields the result shown in Figure 4.

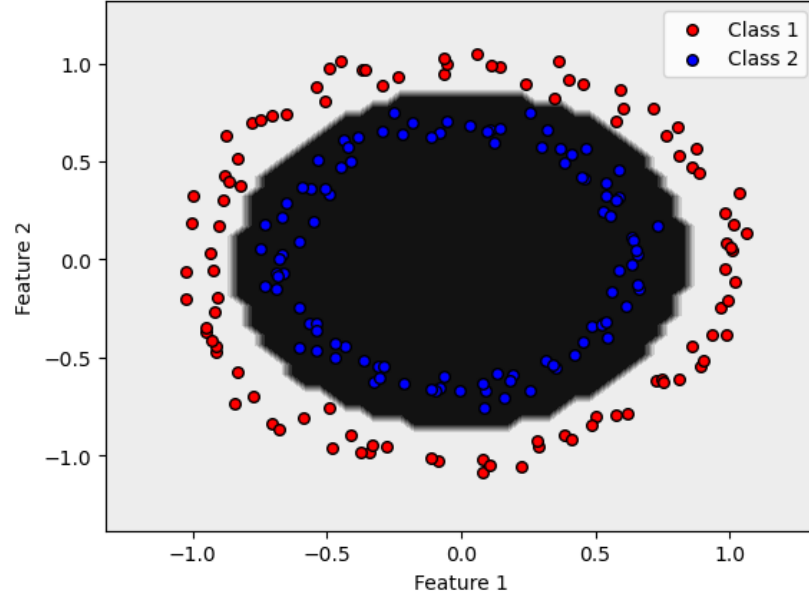


Figure 4: Highly non-linear mock dataset decision boundary, fitted with a SVM using kernel of eq. (21). The white and black areas are the two predicted classes. We observe how, using this kernel, we obtain a non-linear decision boundary.

3 Quantum Support Vector Machine

3.1 General overview

The concept of QSVM as we use today was first introduced in 2019, by Havlicek et al. [8] and, independently, by Schuld [18], that proposed and experimentally implemented the method on a superconducting processor. In 2021 Liu et al. established a rigorous quantum speedup for supervised classification using a general-purpose quantum learning algorithm that only requires classical access to data. Since then several interesting works have been published, but none faces the problem of how to build the quantum circuit to use to encode the classical data in the quantum device. This choice, we will see, is highly problematic, as already noted by Park et al. [15]. In order to use the QSVM for real world applications it is mandatory to have a functioning way of finding the correct circuit to use.

3.2 Quantum feature map and Quantum Kernels

We anticipated how the classical SVM algorithm faces some important limitations when the feature space becomes large, as estimating kernel functions becomes computationally intensive. Quantum computing could enhance the algorithm's performance by providing access to exponentially large Hilbert feature spaces. The idea is to construct a feature map which maps classical data into a quantum state which lives in an exponentially large Hilbert

feature space. Therefore, in this context a feature map is a function

$$\phi : \mathbb{R}^d \rightarrow \mathcal{H}, \quad (22)$$

$$\mathbf{x} \mapsto \phi(\mathbf{x}) \equiv |\phi(\mathbf{x})\rangle.$$

In the framework of quantum computing \mathcal{H} is a n -qubit Hilbert space, that is a space of the form

$$\mathcal{H} = \bigotimes_{i=0}^{n-1} \mathcal{H}_{qubit}, \quad (23)$$

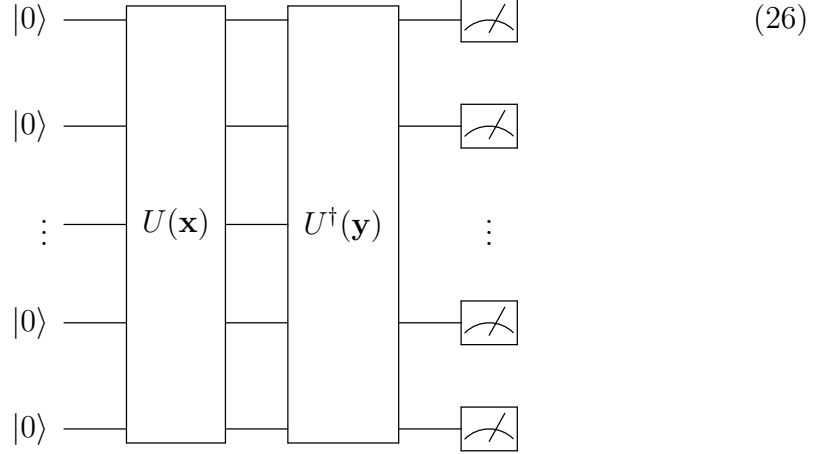
where \mathcal{H}_{qubit} is the Hilbert space of a single qubit. The dimension of \mathcal{H} is thus 2^n . The feature map will be implemented by means of a parametrized quantum circuit. That means that it exists a unitary operator that depends on d classical parameters $U(\mathbf{x}) = U(x_0, \dots, x_{d-1})$ such that

$$|\phi(\mathbf{x})\rangle = U(\mathbf{x})|0\rangle^{\otimes n}. \quad (24)$$

This circuit is called *quantum encoding circuit*, because it encodes classical data into a quantum state. The classical data is passed to the circuit as a parameter. We will later make examples of frequently used encoding circuits. Once we have the feature map, the kernel is constructed as

$$K(\mathbf{x}, \mathbf{y}) = |\langle \phi(\mathbf{x}) | \phi(\mathbf{y}) \rangle|^2. \quad (25)$$

Here $\langle \cdot, \cdot \rangle$ denotes the standard internal scalar product between vectors in \mathcal{H} . This definition clearly yields a kernel that satisfies the two kernel properties. How do we calculate the kernel in practice? Suppose we want to calculate the kernel $K(\mathbf{x}, \mathbf{y})$ and consider the following circuit.



Suppose we run this circuit R times, and we call A the number of times that we measure the bit string $000 \cdots 0$. We state that

$$\lim_{R \rightarrow +\infty} \frac{A}{R} = K(\mathbf{x}, \mathbf{y}). \quad (27)$$

The proof is straightforward. The left hand side of eq. (27) is the probability of measuring $000 \cdots 0$, which according to quantum mechanics can be calculated as

$$\begin{aligned} |\langle 0 | U(\mathbf{x}) U^\dagger(\mathbf{y}) | 0 \rangle|^2 &= |\langle 0 | U^\dagger(\mathbf{y}) U(\mathbf{x}) | 0 \rangle|^2 = \\ &= |\langle \phi(\mathbf{y}) | \phi(\mathbf{x}) \rangle|^2 = K(\mathbf{x}, \mathbf{y}). \end{aligned}$$

□

Therefore, to evaluate the kernel, it suffices to construct the quantum circuit (26) and measure the frequency with which the string $000 \cdots 00$ occurs. We have to perform this operation for each pair of instances, and construct the

matrix

$$K_{ij} = K(\mathbf{x}_i, \mathbf{x}_j). \quad (28)$$

This kernel matrix is then plugged into eq. (19) to train a classical SVM.

3.3 Quantum encoding circuits

Let's now discuss some possible choices of quantum encoding circuits:

- Basis encoding: this can be applied if the instances live in a finite dimensional space, whose dimension is 2^n , where n is the number of qubits. If the dimension is smaller than 2^n we can still apply this encoding by using padding. The instances can be mapped into bit strings of length n and the feature map corresponds to

$$\phi : \{0, 1\}^n \rightarrow \mathcal{H}, \quad (29)$$

$$x = (x_0 \cdots x_{n-1}) \mapsto |x_0\rangle \otimes \cdots \otimes |x_{n-1}\rangle \equiv |x_0 \cdots x_{n-1}\rangle = |x\rangle.$$

Here $|x\rangle$ is a state of the so called computational basis. For example, the bit string 01001 is mapped to the quantum state

$$|0\rangle \otimes |1\rangle \otimes |0\rangle \otimes |0\rangle \otimes |1\rangle = |01001\rangle.$$

The kernel that originates from this feature map is

$$K(x, y) = |\langle x|y\rangle|^2 = |\delta_{x,y}|^2 = \delta_{x,y} \quad (30)$$

since vectors in the computational basis are orthogonal, and δ is the Kronecker delta.

- Amplitude encoding: this can be applied with instances that are normalized vectors, that is vectors of the form

$$\mathbf{x} = (x_0, \dots, x_{d-1})^T \in \mathbb{R}^d, \quad (31)$$

such that

$$\sum_{i=0}^{d-1} x_i^2 = 1.$$

In this case, the encoding consists in mapping

$$\mathbf{x} = (x_0, \dots, x_{d-1})^T \mapsto |\phi(\mathbf{x})\rangle = \sum_{i=0}^{d-1} x_i |i\rangle, \quad (32)$$

where $|i\rangle$ is the i -th vector of the computational basis. The fact that \mathbf{x} is normalized ensures that the feature vector is normalized as well, and therefore represents a physical state. The kernel generated by this feature map is

$$K(x, y) = \left| \sum_{i,j} x_i y_j \langle i|j \rangle \right|^2 = \left| \sum_{i,j} x_i y_j \delta_{i,j} \right|^2 = (\mathbf{x}^T \cdot \mathbf{y})^2. \quad (33)$$

As we can see we went back to a polynomial kernel.

- By creating more copies of the feature vector in amplitude encoding, we can create higher order polynomial kernels.
- Product encoding: in this case each feature of the input is encoded in the amplitudes of one separate qubit. For example, we can map

$$\mathbf{x} = (x_0, \dots, x_{d-1})^T \mapsto \begin{pmatrix} \cos x_0 \\ \sin x_0 \end{pmatrix} \otimes \dots \otimes \begin{pmatrix} \cos x_{d-1} \\ \sin x_{d-1} \end{pmatrix}, \quad (34)$$

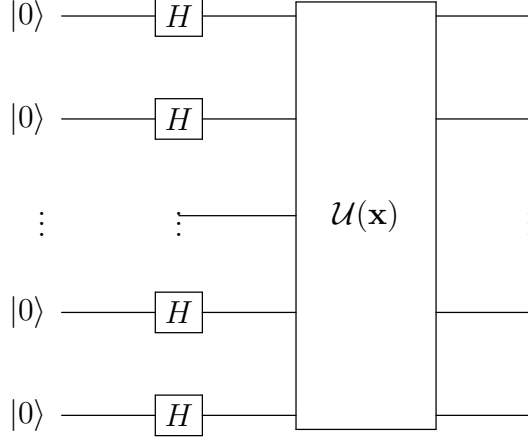
therefore the i -th qubit will be in the state

$$\cos x_i |0\rangle + \sin x_i |1\rangle. \quad (35)$$

This generates the kernel

$$K(x, y) = \prod_{i=0}^{d-1} \cos(x_i - y_i). \quad (36)$$

- **ZZ feature map:** The encoding maps discussed above are quantum maps that ultimately produce a kernel easily computable with classical tools, like a polynomial or cosine kernel. It is important now to discuss quantum encoding circuits which do not produce an easy classical kernel, otherwise there is no advantage in using quantum computing for this task if the kernel could always be easily calculated with classical tools. The hope is that by using a purely quantum kernel we obtain better performance than by using a classical kernel, in terms of classification or regression capability. Since general quantum circuits are not expected to be classically simulable, there are many choices one can make. One example is generated by the following circuit



where H is the standard Hadamard gate and

$$\mathcal{U}(\mathbf{x}) = \exp \left[i \sum_{S \subseteq [n]} \phi_S(\mathbf{x}) \prod_{i \in S} Z_i \right]. \quad (37)$$

Here $[n] = \{0, \dots, n-1\}$, Z_i is the third Pauli matrix acting on the i -th qubit, and $\phi_S(\mathbf{x})$ are arbitrary coefficients, which are the ones that actually encode the classical data. S runs over all possible subsets of $[n]$, and can be thought of as an index that describes connectivities between different qubits or data points. We have 2^n possible choices of the coefficients. It is convenient to choose them in such a way that only the terms with $|S| \leq d$ contribute, in order to obtain a circuit that can be easily implemented on a quantum computer. This is done by asking that

$$\phi_S = 0 \quad \text{if} \quad |S| > d. \quad (38)$$

Let's focus on a simple case where $n = d = 2$. There are two default

choices for the coefficients in this case. The first one is to set

$$\phi_{\{i\}}(\mathbf{x}) = x_i, \quad (39)$$

$$\phi_{\{i,j\}}(\mathbf{x}) = 0. \quad (40)$$

The feature map becomes

$$U(\mathbf{x}) = e^{ix_0 Z_0} e^{ix_1 Z_1} H^{\otimes n}. \quad (41)$$

This is called the *Z feature map*. The corresponding circuit, decomposed in elementary quantum gates, is shown in Figure 5. The other

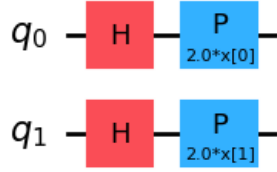


Figure 5: Z feature map, with 2 qubits and 2 features, decomposed in terms of elementary gates. The gate P is defined as $P(\phi) = (1, 0; 0, e^{i\phi})$, where ϕ is the parameter of the gate, and in our case can either be $2x_0$ or $2x_1$. We observe how the circuit is parametrized by the coefficients of eq. (39).

possible choice is

$$\phi_{\{i\}}(\mathbf{x}) = x_i, \quad (42)$$

$$\phi_{\{i,j\}}(\mathbf{x}) = (\pi - x_i)(\pi - x_j). \quad (43)$$

In this case the feature maps becomes

$$U(\mathbf{x}) = e^{i(\pi-x_0)(\pi-x_1)Z_0Z_1} e^{ix_0Z_0} e^{ix_1Z_1} H^{\otimes n}. \quad (44)$$

This is called the *ZZ feature map*, and it is an encoding circuit thought to be hard to simulate classically. It is shown in Figure 6. We will see

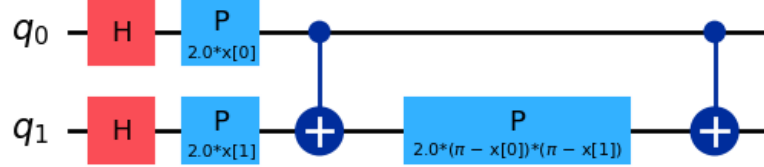


Figure 6: ZZ feature map, with 2 qubits and 2 features, decomposed in terms of elementary gates. The gate P is defined as $P(\phi) = (1, 0; 0, e^{i\phi})$, where ϕ is the parameter of the gate. The values that the parameter take are shown in the Figure, according to eq. (42) and (43).

how it behaves later on. Those two are standard choices, but there are many more circuits that can be built starting from the general structure of eq. (37), especially if we consider a greater number of qubits, and so we allow connectivities of three and more qubits.

- Pauli feature map: Eq. (37) can be generalised to include not only Z gates, but all of the Pauli gates. We can have

$$\mathcal{U}(\mathbf{x}) = \exp \left[i \sum_{S \subseteq [n]} \phi_S(\mathbf{x}) \prod_{i \in S} P_i \right], \quad (45)$$

where $P \in \{1, X, Y, Z\}$. This allows us to generate not only interactions of the ZZ type, but also for example YY type or ZY type. An example with Y and ZY interactions is shown in Figure 7.

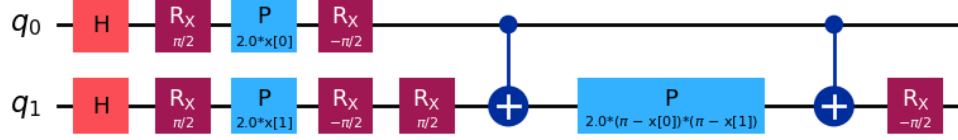


Figure 7: Pauli feature map, with 2 qubits and 2 features, decomposed in term of elementary gates, constructed with Y and YZ interactions.

- In general any parametrized quantum circuit can be used to encode classical data and create a quantum kernel. We will later use generic circuits and study what sort of datasets they separate.

3.4 QSVM algorithm and Qiskit implementation

To sum things up, the step of the Quantum Support Vector Machine (QSVM) algorithm are the following:

- Prepare the classical data into two vectors X and y , where

$$X_i = \mathbf{x}_i \in \mathbb{R}^d, \quad i = 0, \dots, N-1,$$

is the i -th d -dimensional instance, and

$$y_i \in \{-1, 1\},$$

indicates the class that this instance belongs to.

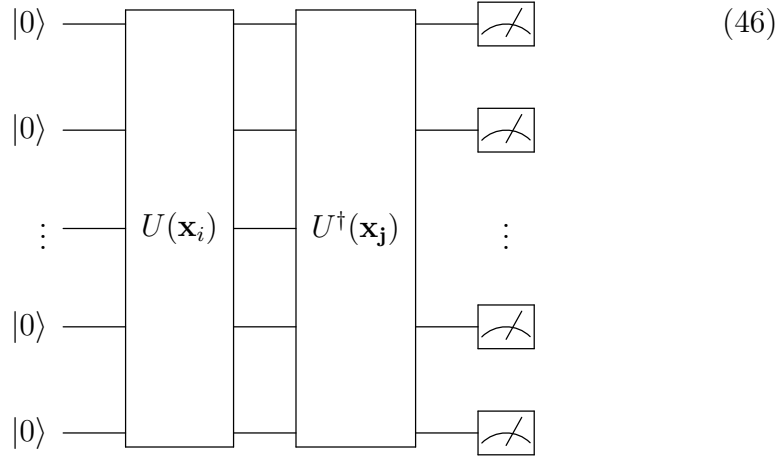
- Choose a parametrized quantum circuit $U(\mathbf{x})$. The circuit accepts d parameters, which are the features of the instance we are considering.

One can choose a standard circuit like the ZZ feature map, or can create a custom circuit. We will extensively revisit the point of choosing the circuit later on.

- For each pair of instances $\mathbf{x}_i, \mathbf{x}_j$ we must evaluate the kernel matrix

$$K_{ij} = K(\mathbf{x}_i, \mathbf{x}_j).$$

This is an $N \times N$ matrix, so for each pair we have to build $U(\mathbf{x}_i)$ and $U^\dagger(\mathbf{x}_j)$, in order to construct the circuit



and evaluate the matrix element

$$K_{ij} = |\langle 0^{\otimes n} | U(\mathbf{x}_i) U^\dagger(\mathbf{x}_j) | 0^{\otimes n} \rangle|^2. \quad (47)$$

As we saw this can be done, up to an error $R^{-1/2}$, by running the circuit (46) R times and counting the number of times we measure the string $000 \dots 0$.

- Once we have the kernel matrix we use classical optimization to maximize

$$f(\alpha_0, \dots, \alpha_{N-1}) = \sum_{i=0}^{N-1} \alpha_i - \frac{1}{2} \sum_{i=0}^{N-1} \sum_{j=0}^{N-1} \alpha_i \alpha_j y_i y_j K_{ij}. \quad (48)$$

Once we found the optimal α we have constructed the decision boundary.

- Given a new instance $\tilde{\mathbf{x}}$ we can, in the same way, calculate the kernel matrix element

$$\tilde{K}_i = K(\mathbf{x}_i, \tilde{\mathbf{x}}),$$

and predict the class of $\tilde{\mathbf{x}}$ by applying

$$\tilde{y} = \text{sign} \left[\sum_{i=0}^{N-1} \alpha_i y_i \tilde{K}_i + w_0 \right]. \quad (49)$$

Schematically these operations are described in Algorithm 1.

Algorithm 1 Quantum Support Vector Machine (QSVM)

```
1: Input: Data  $\{(\mathbf{x}_i, y_i)\}_{i=0, \dots, N-1}$ , Quantum circuit  $U(\mathbf{x})$ 
2: Parameters: Number of measurement shots  $R$ 
3: for  $i = 0$  to  $N - 1$  do
4:   for  $j = 0$  to  $N - 1$  do
5:     Prepare  $U(\mathbf{x}_i)$  and  $U^\dagger(\mathbf{x}_j)$ 
6:     Run circuit  $U(\mathbf{x}_i)U^\dagger(\mathbf{x}_j)$   $R$  times
7:     Measure the frequency  $f$  of  $00 \dots 0$ 
8:     Set  $f = K_{ij}$ 
9:   end for
10: end for
11: Optimize  $\alpha$  to maximize
```

$$f(\alpha) = \sum_{i=0}^{N-1} \alpha_i - \frac{1}{2} \sum_{i=0}^{N-1} \sum_{j=0}^{N-1} \alpha_i \alpha_j y_i y_j K_{ij}$$

```
12: for each new instance  $\tilde{\mathbf{x}}$  do
13:   Compute  $\tilde{K}_i = K(\mathbf{x}_i, \tilde{\mathbf{x}})$ 
14:   Predict  $\tilde{y} = \text{sign} \left( \sum_{i=0}^{N-1} \alpha_i y_i \tilde{K}_i + w_0 \right)$ 
15: end for
```

This algorithm can be implemented in practice by exploiting Qiskit to perform the circuit simulation and calculate the kernel, and subsequently using scikit-learn to perform the classical optimization once the kernel has been calculated. We can make an example, separating a mock dataset like the one in Figure 1, with a Pauli feature map.

```

#create mock dataset
X,y=make_blobs(n_samples=200)
#since the kernel involves rotation it is better to
    ↪ bring the data between 0 and pi
X=MinMaxScaler(feature_range=(0,np.pi)).
    ↪ fit_transform(X)
#create quantum kernel
kernel=FidelityQuantumKernel(feature_map=ZFeatureMap
    ↪ ())
#pass the kernel as a callable function
svm=SVC(kernel=kernel.evaluate)
#fit the svm
svm.fit(X, y)

```

The result of the fit is shown in Figure 8.

3.5 QSVM potential

In order to study the potential of the QSVM algorithm we could observe the form of the decision boundary of an artificial dataset which is easily separable by a QSVM. If the decision boundary is "classically looking", that is, similar to the one we observed when we studied the linear or polynomial kernel, then it would look like the QSVM algorithm may not provide advantage over the classical SVM. We hope instead to find a complex looking decision boundary, which is not separable by classical kernels.

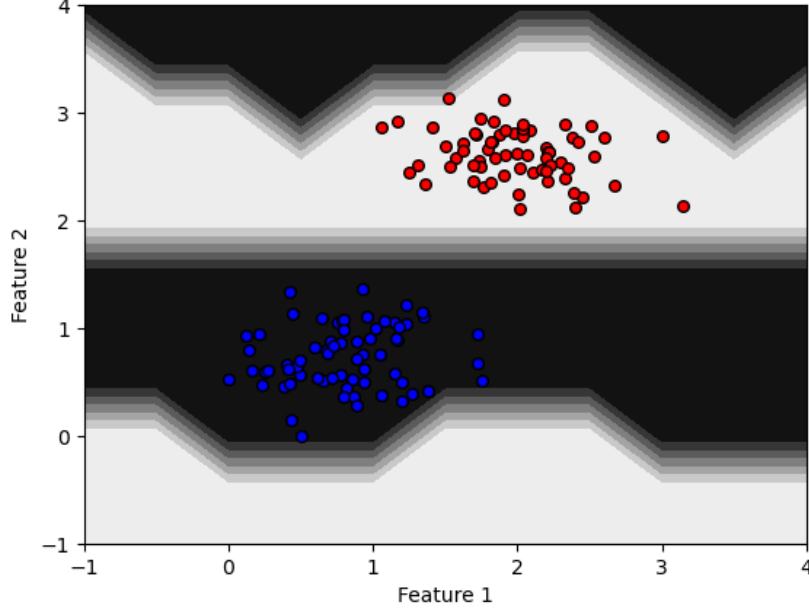


Figure 8: QSVN decision boundary with ZFeatureMap, fitted on a two feature mock dataset of 200 instances.

We use ZZFeatureMap, with $n = d = 2$. Furthermore, we choose $f = Z_1 Z_2$ and $V \in \text{SU}(4)$. We assign $m(\mathbf{x}) = +1$ when

$$\langle \phi(\mathbf{x}) | V^\dagger f V | \phi(\mathbf{x}) \rangle > \Delta, \quad (50)$$

and $m(\mathbf{x}) = -1$ when

$$\langle \phi(\mathbf{x}) | V^\dagger f V | \phi(\mathbf{x}) \rangle < \Delta. \quad (51)$$

Δ controls how big is the gap between the two classes. We show this dataset in Figure 9. As we can see, such dataset shows a very complicated pattern. Classical methods cannot hope to separate these two classes. In fact, if we divide the dataset into a training set and a test set, and we apply the classical

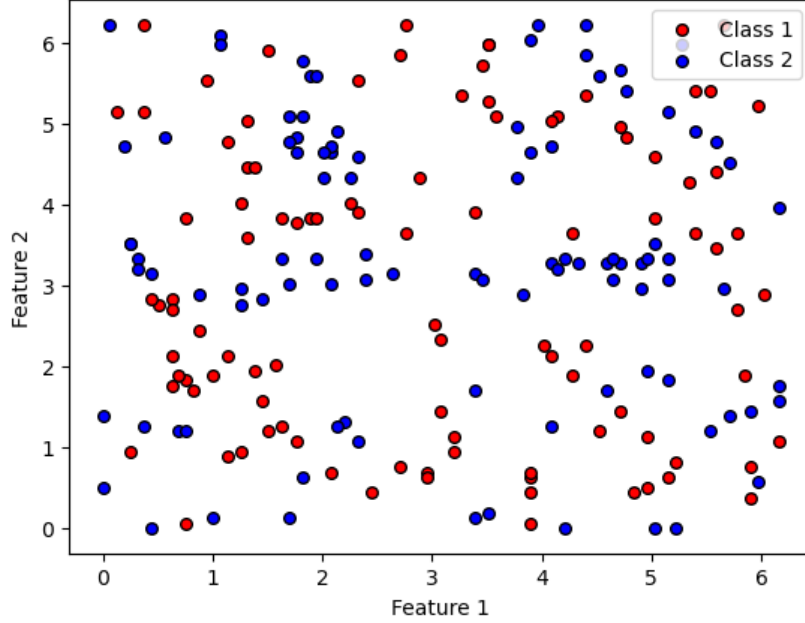


Figure 9: Artificial ad hoc dataset, generated with $\Delta = 0.3$. This dataset, by construction, is easily separable by ZZFeatureMap.

SVM algorithm, for example with the RBF kernel, we obtain an accuracy on the test set of 0.51, which is essentially like assigning the classes with a random coin toss. The result of the classical fit is shown in Figure 10, in which we observe how the correct decision boundary has not been found. In Figure 11 is shown instead the result of the fit with the QSVM with ZZFeatureMap, which yields an accuracy on the test set of 1.0, and a nice looking decision boundary.

Finding a dataset which is separable by quantum methods and not by classical methods leads us into thinking that the QSVM algorithm might outperform the SVM algorithm on datasets with complex patterns, which is

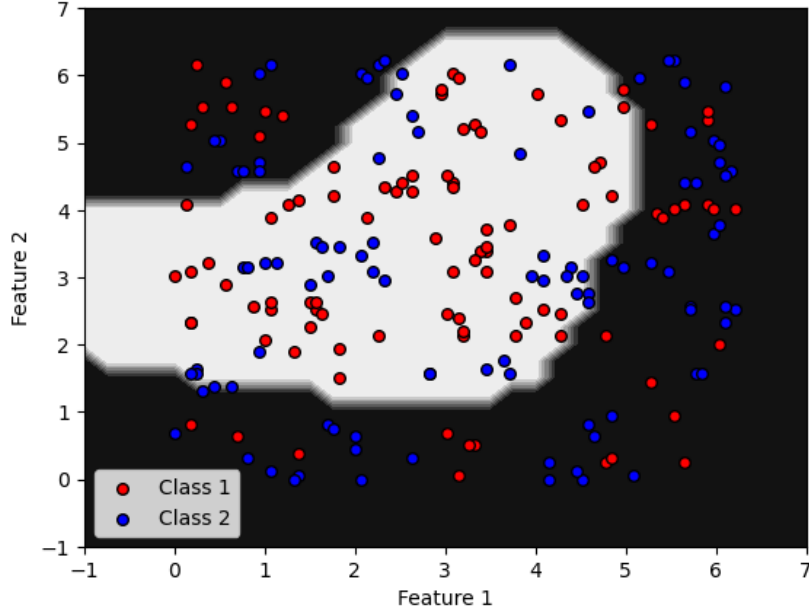


Figure 10: Result of the classical fit with RBF kernel on the artificial dataset. We observe how the SVM could not find the correct separation between the classes.

the result we hoped to find.

3.6 Issues of QSVM

In a QSVM algorithm, the only decision the user needs to make is the selection of the feature map, which corresponds to choosing the quantum encoding circuit. It's clear that this choice is crucial, as it determines whether the algorithm will be successful. In the classical case, we also need to choose a kernel, but this decision is generally less sensitive than in the quantum case. For instance, the RBF kernel performs well across most applications. In the

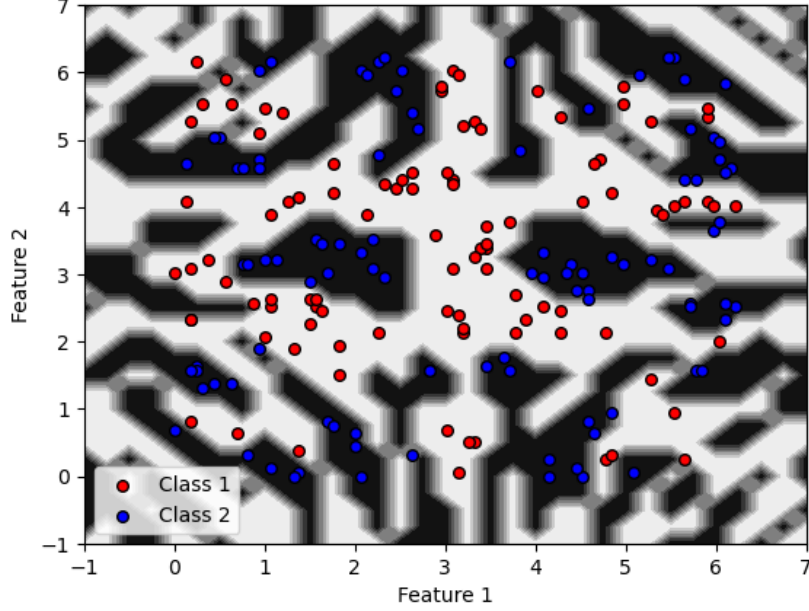


Figure 11: Result of the quantum fit on the artificial dataset, with ZZFeatureMap. We observe how the QSVM found the correct separation between the classes.

quantum case, however, selecting the feature map is particularly challenging: if the wrong feature map is chosen, the QSVM performs very poorly, even on simple datasets, as we will demonstrate with some examples. Let's consider the dataset shown in Figure 2, which is simple enough to be expected to be separated without too many difficulties. We already saw classical kernels separate this dataset very well. Let us now try to apply some of the quantum kernels we discussed in the previous sections.

We find that the results of the fits are not satisfying, as the accuracy is always between 40% and 60%, but most importantly the form of the decision

boundary tells us that the algorithm did not find the correct pattern for the data, as shown in Figure 12.

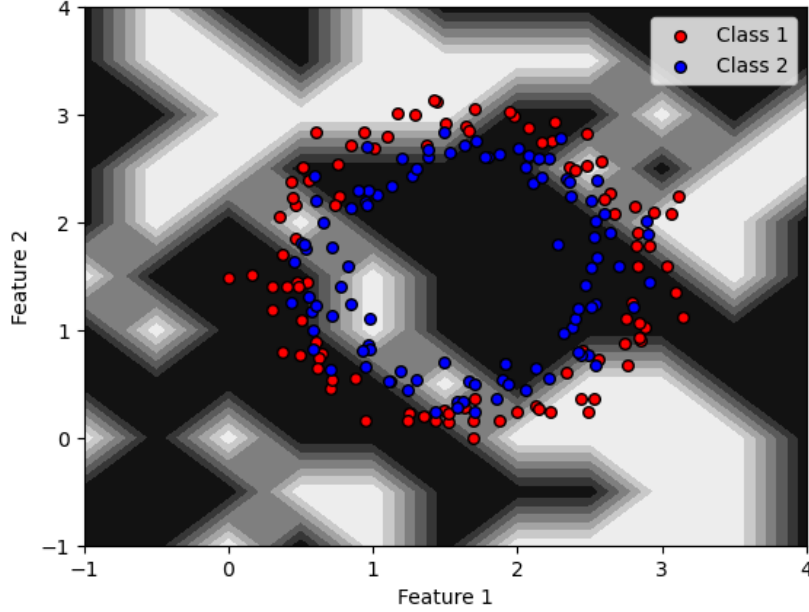


Figure 12: Circle dataset fitted with PauliFeatureMap leads to a very inaccurate decision boundary.

It is clear at this point that the choice of the quantum encoding map is a very delicate part of the algorithm, and cannot be taken lightly. Choosing the feature map by trying the most common ones, as one may do with classical SVM, is not a safe method, as for many datasets, even simple ones, the QSVM may fail. We need a more systematic way to choose the encoding, as the performance of the model is significantly influenced by this choice. We will discuss those methods in the next sections.

3.7 Quantum kernel alignment

3.7.1 Circuit optimization

In the previous sections, we observed how the main challenge with the QSVM algorithm lies in selecting an appropriate feature map. This choice is difficult because there is no clear guidance in constructing the appropriate circuit, and the number of possible circuits for data encoding is potentially unlimited. Without guidance, it is problematic to select a circuit which works, since we saw how the choice of the feature map is very delicate. A first way to address this problem is the so-called Quantum Kernel Alignment (QKA).

QKA is a hybrid quantum-classical approach, and consists in employing a parametrized quantum feature map $U(\mathbf{x}, \lambda)$, which depends upon k parameters $\lambda = (\lambda_0, \dots, \lambda_{k-1})$. The quantum circuit now has $d + k$ parameters, where d are passed through \mathbf{x} and k through λ . However, the latter are fundamentally different from the former. The \mathbf{x} parameters, as we saw, are the way the instance is passed to the circuit. The λ parameters instead must be trained in order to find the best possible circuit. A parametrized feature map gives rise to a parametrized quantum kernel

$$K_{i,j}(\lambda) = K_\lambda(\mathbf{x}_i, \mathbf{x}_j),$$

and consequently to a parametrized objective function of the SVM

$$f(\alpha, \lambda) = \sum_{i=0}^{N-1} \alpha_i - \frac{1}{2} \sum_{i=0}^{N-1} \sum_{j=0}^{N-1} \alpha_i \alpha_j y_i y_j K_{ij}(\lambda). \quad (52)$$

Kernel alignment aims at minimizing the generalization error bound of SVM,

and consists in solving the optimization problem

$$\min_{\lambda} \max_{\alpha} f(\alpha, \lambda). \quad (53)$$

Here again α is the solution to the convex optimization problem of the SVM, while λ is the variational parameter of the quantum circuit which must be trained. The optimization over α is performed as usual for the SVM. The optimization over λ is performed using a Classical Simultaneous Stochastic Approximation (SPSA). [21] The choice of SPSA lies in the fact that it needs less kernel evaluations compared to gradient descent methods, and it is more robust against noise, which is a serious issue of current real world quantum devices.

QKA is a hybrid quantum-classical method, because the quantum circuit is utilized to evaluate the kernel when necessary, while all the other tasks such as optimizations are performed by classical algorithms. The fundamental steps of QKA are:

- Choose a parametrized quantum circuit $U(\mathbf{x}, \lambda)$, where $\lambda \in \mathbb{R}^k$.
- Choose an initial parameter value $\lambda = \lambda_0$, and the maximum number of iterations P of the SPSA algorithm.
- Generate a random vector $\Delta \in \{-1, 1\}^k$.
- Create two new parameters λ_+ and λ_- such that

$$\lambda_{\pm, i} = \lambda_i \pm c_i \Delta_i,$$

where

$$c_i = \frac{c}{(i+1)^\gamma},$$

for chosen constants c and γ .

- Evaluate, using the QSVM standard approach on quantum device, the two kernels

$$K_{\pm,ij} = K_{\lambda_{\pm}}(\mathbf{x}_i, \mathbf{x}_j).$$

- Find the solution α_{\pm} of the convex optimization problem in eq. (52), for the two kernels of the previous point, exploiting the standard SVM optimization.
- Perform two evaluations of the loss function to estimate its gradient with respect to λ_k :

$$g_k(\alpha_k, \lambda_k) = \frac{1}{2c_k\Delta_k} (f(\alpha_{+,k}, \lambda_{+,k}) - f(\alpha_{-,k}, \lambda_{-,k})), \quad (54)$$

where the loss function f is given by eq. (52).

- Update λ depending on the gradient $g_k(\alpha_k, \lambda_k)$ and the learning rate a_k . The update rule is

$$\lambda_{k+1} = \lambda_k - a_k g_k(\alpha_k, \lambda_k). \quad (55)$$

- Repeat the procedure until the cost function has converged or the maximum number of iterations has been reached.
- The final kernel is the optimal quantum kernel $K_{\lambda^*}(\underline{x}_i, \underline{x}_j)$.

There are several possible choices for the embedding of the parameters in the quantum circuit. One common choice is to insert two parametrized rotations before the actual feature map, as shown in Figure 13. Another

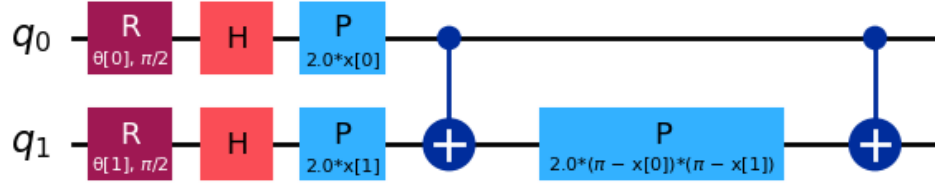


Figure 13: Parametrized quantum circuit, with two rotation inserted before a ZZFeatureMap.

possible choice is to insert the parameters inside a ZZFeatureMap, and then apply the canonical ZZFeatureMap, as shown in Figure 14. The first one is parametrized by λ , while the second one by \mathbf{x} .

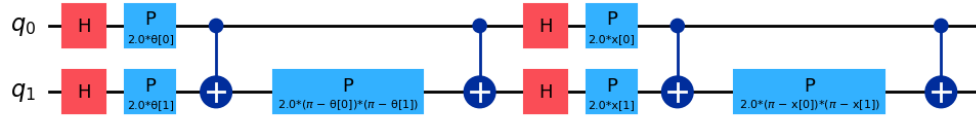


Figure 14: Parametrized quantum circuit, with a ZZFeatureMap inserted before another ZZFeatureMap. The former has the parameters λ for Quantum Kernel Alignment, while the latter has the parameters \mathbf{x} in order to encode the classical instance.

QKA can be implemented in Qiskit using the TrainableFidelityQuan-

tumKernel class and the SPSA optimizer as follows.

```
kernel = TrainableFidelityQuantumKernel()
cb_qkt = QKTCallback()
spsa_opt = SPSA(maxiter=20, callback=cb_qkt.callback
    ↪ , learning_rate=0.05)
loss = SVCLoss() #loss function for the SVM
qkt = QuantumKernelTrainer(quantum_kernel=kernel,
    ↪ loss=loss, optimizer=spsa_opt)
qka_results = qkt.fit(X_train, y_train) #optimize
    ↪ kernel
optimized_kernel = qka_results.quantum_kernel
```

The final kernel is the optimized one, which can be used to train the QSVM and obtain the final result.

3.7.2 Performance

We can evaluate the performance of the QKA-QSVM by observing the behavior on the previously discussed datasets, in which the standard SVM algorithm performed badly. Let us consider for example the circular dataset of Figure 2, and let us use the feature map of Figure 13. The resultinf fit is shown in Figure 15, while the convergence of the SVC loss function and the final kernel are shown in Figure 16.

We observe how even in this case the result is not satisfactory at all. The problem is that, even though now we have the freedom to insert parameters

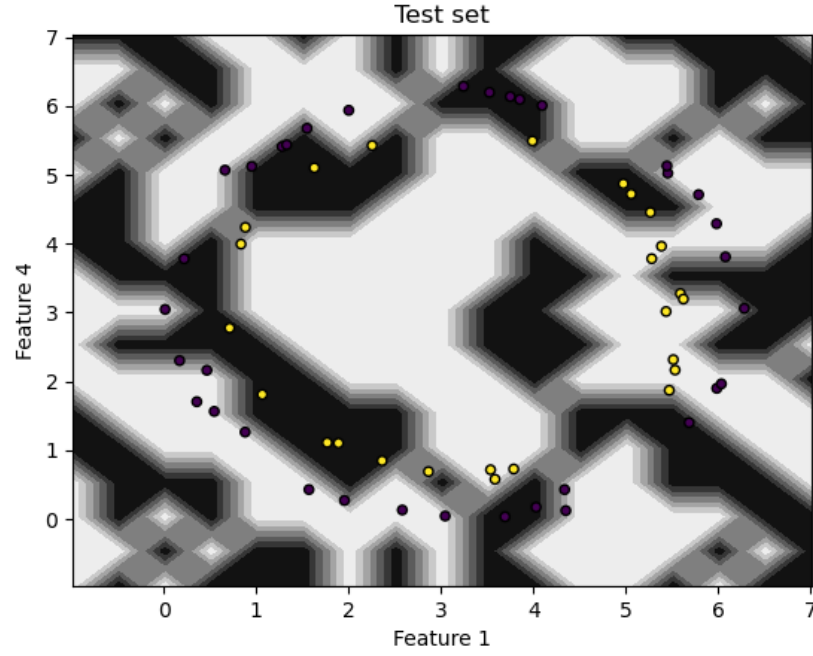


Figure 15: Result of the fit of the circular dataset using QKA-QSVM. An accuracy of 50% is achieved.

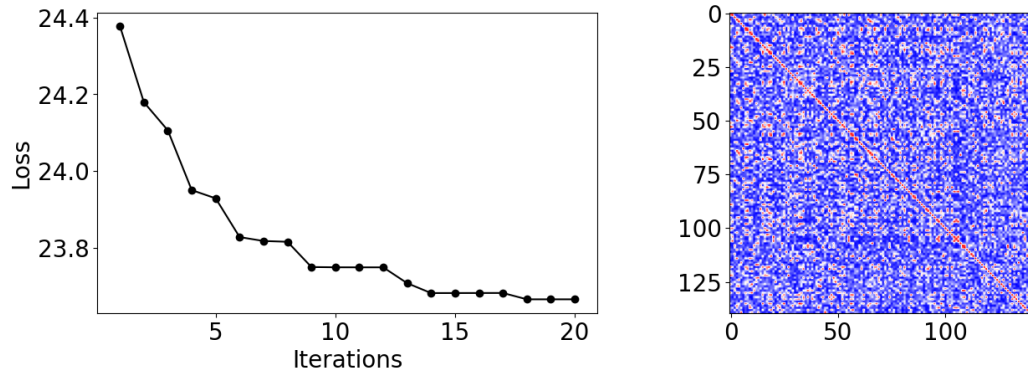


Figure 16: On the left, the convergence of the loss function during the QKA training is shown. On the right, we can see the final trained kernel matrix.

inside the model, it is still not clear how to choose the circuit or the parameters. This problem is not solved, consequently we are not capable of choosing a circuit that works for this particular dataset. We thus need to find a more advanced algorithm to choose the circuit.

4 Genetic algorithm

A genetic algorithm can be used to choose the best feature map. We'll see how these types of algorithms can pick the right feature map based on first principles, without us having to decide the circuit's structure. Genetic algorithms have been previously used in the context of quantum computing, for example in [2][3][22]. We will apply this strategy for the QSVM case.

4.1 General introduction

Genetic algorithms [14][6] are optimization algorithms inspired by the process of natural selection, a key mechanism of biological evolution. The core concept is based on mimicking the biological evolution process, particularly principles like survival of the fittest, reproduction, mutation, and inheritance.

The primary components of the algorithm are the individuals, whose specific characteristics depend on the problem being addressed. Each individual is defined by a genotype and a phenotype. The genotype consists of a collection of genes, typically represented by numbers or character strings. The phenotype is a practical manifestation of the genotype, such as a quantum circuit or a neural network.

Individuals are organized into populations, which are evolved with the passing of generations. The transition from one generation to the next occurs through crossover and mutation, mimicking the process of natural selection. A fitness function is defined to assess how well an individual performs in relation to the specific problem being addressed. The top-performing individuals

in a generation, i.e., those with the highest fitness scores, can be passed directly to the next generation without undergoing crossover or mutation. This process is known as elitism.

Crossover emulates the crossing-over mechanism in meiosis. Two parent individuals from the current generation are selected with a probability proportional to their fitness, and two offspring are produced by combining the genotypes of the parents in a specific manner. Additionally, the offspring may undergo random gene mutations. The fitness of both offspring is then evaluated, with only the fittest being retained for the next generation. This way the generations all have the same number of individuals, to avoid an exponential growth of the dimension.

This process is repeated multiple times to generate successive generations, and continues until either the maximum number of generations is reached or an individual achieves a predetermined target fitness level. The best individual of the last generation is the result of the genetic algorithm.

4.2 Genetic algorithm for QSVM

Let's now discuss how the above-mentioned concepts apply to the case of Quantum Support Vector Machine. The goal is to find the best possible feature map, given a dataset.

An individual is meant to represent a possible feature map choice. We define the genes of an individual in the form

$$[G, q_t, q_c, f], \tag{56}$$

where

- G is a gate chosen from a universal set of gates. For example, we can choose the set

$$\{X, \sqrt{X}, \text{CNOT}, R_z\}. \quad (57)$$

Here universal means that every possible n -input/ n -output gate can be generated using some combination of gates (that is, a circuit) from this set. X, \sqrt{X}, CNOT are fixed 1 and 2 qubit gates, whose matrix form is

$$X = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix},$$

$$\sqrt{X} = \frac{1}{2} \begin{pmatrix} 1+i & 1-i \\ 1-i & 1+i \end{pmatrix},$$

$$\text{CNOT} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{pmatrix}.$$

On the other hand, R_z is a parametrized gate of the form

$$R_z(\theta) = \begin{pmatrix} e^{-i\theta/2} & 0 \\ 0 & e^{i\theta/2} \end{pmatrix}.$$

Since the feature map is constituted by a parametrized quantum circuit $U(\mathbf{x})$, the parameters \mathbf{x} will be inserted inside the R_z type gates.

- q_t is the target qubit to which the gate is applied, therefore it will be an integer number between 0 and $n - 1$, where n is the number of qubits.
- q_c is the control qubit, if it exists for that particular gate. If we use the gates set presented above only the CNOT gate has a control qubit, otherwise it is set to None.
- f represents which feature to use as the gate parameter, if the gate accepts a parameter. In our example only the R_z gate accepts one. In the other cases it is set to None. Since we have d features (x_0, \dots, x_{d-1}) , f will be an integer between 0 and $d - 1$.

The genotype will be a collection of genes of this form. The number of genes can vary, in particular it will represent the number of gates of the circuit. The phenotype associated to a genotype will be the corresponding parametrized quantum circuit. Let's make an example of this and consider $n = d = 2$ and the genotype

$$\{[R_z, 0, \text{None}, 0], [R_z, 1, \text{None}, 0], [\text{CNOT}, 1, 0, \text{None}], [X, 1, \text{None}, \text{None}]\}.$$

There are four genes, so the circuit must have four gate. If we look at the first gene we can conclude that the first gate is a R_z with the first qubit as target. The control qubit is set to None since R_z is a single qubit gate. The 0-th feature x_0 is used as the parameter of the circuit. Similar consid-

erations hold for the others three gates. The circuit, i.e. the phenotype, is represented in Figure 17.

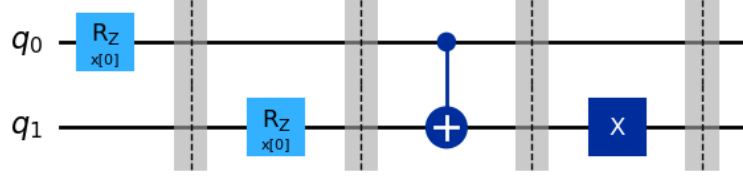


Figure 17: Example of the phenotype of an individual. The phenotype is a quantum circuit.

We saw how parent individuals are passed to the child generation using crossover and mutations. Let's discuss how these operations are performed in the QSVM case.

Crossover consists in cutting the circuit in half and pasting the two halves in two different children circuits. In formulas, assuming the number of genes k to be even, starting from the parents genotypes

$$\{g_0^1, \dots, g_{\frac{k}{2}-1}^1, g_{\frac{k}{2}}^1, \dots, g_{k-1}^1\},$$

$$\{g_0^2, \dots, g_{\frac{k}{2}-1}^2, g_{\frac{k}{2}}^2, \dots, g_{k-1}^2\},$$

we build the two children genotype

$$\{g_0^1, \dots, g_{\frac{k}{2}-1}^1, g_{\frac{k}{2}}^2, \dots, g_{k-1}^2\},$$

$$\{g_0^2, \dots, g_{\frac{k}{2}-1}^2, g_{\frac{k}{2}}^1, \dots, g_{k-1}^1\}.$$

In Figures 18-21 it is shown an example to illustrate what happens on the phenotype.

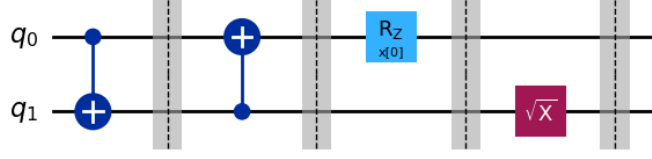


Figure 18: First parent phenotype.

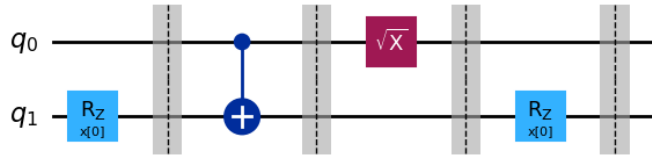


Figure 19: Second parent phenotype.

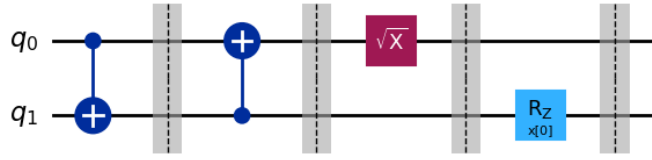


Figure 20: First child phenotype.

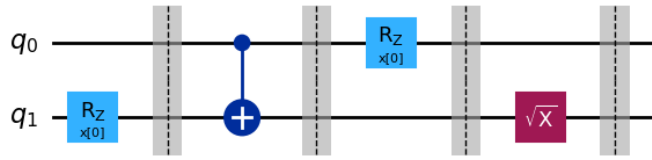


Figure 21: Second child phenotype.

Mutation, instead, consists in a random mutation of the genes. That is, every gene is mutated, with a probability p , into another randomly chosen gene. An alternative way could be to mutate only one randomly chosen

gene. An example is shown in Figures 22-23. This is a possible result of the

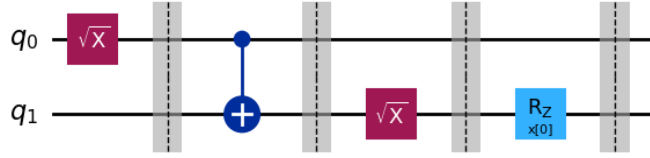


Figure 22: Starting non mutated individual.

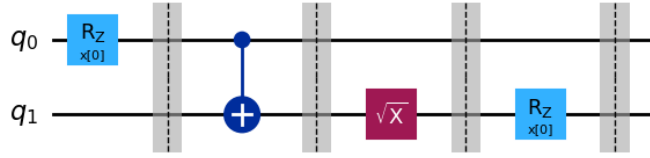


Figure 23: Mutated individual. The first gene is mutated from a \sqrt{X} gate acting on the first qubit into an R_z gate acting again on the first qubit. The other genes in this case were not mutated.

mutation; of course, being an intrinsically random process, another mutation would produce a different result.

The last concept we need to define in the QSVM case is the one of the fitness function. That is, we need a quantitative way to assess how well an individual performs on the particular dataset we are considering. In machine learning there are several ways to evaluate the performance on a classification task. The simplest one is the accuracy on the test set, that is, the ratio of correctly classified points and the total number of points:

$$\text{accuracy} = \frac{\text{number of correctly classified points}}{\text{total number of points}}. \quad (58)$$

For instance, we could use test set accuracy as the fitness of an individual. However, this approach would not sufficiently penalize poorly performing individuals, who are selected to generate offspring with a probability proportional to their fitness. If we use accuracy, an individual with a fitness of 0.5 has only half the probability of being chosen compared to one with a fitness close to 1. This is problematic because the former performs very poorly (its performance is comparable to randomly guessing the labels) whereas the latter correctly classifies almost every label. Moreover, individuals with a fitness below 0.5, which are not uncommon given the potential issues with the QSVM algorithm when an unsuitable feature map is selected, still have a non-negligible chance of being chosen. We would prefer these individuals to be heavily penalized, as their performance is even worse than random guessing. A possible way to achieve this is to use

$$\text{fitness} = \text{accuracy}^6. \quad (59)$$

This function is represented in Figure 24. This choice strongly penalizes individuals with performance scores below 0.5, drastically reducing their chances of being selected for reproduction. Conversely, individuals with higher accuracy have significantly greater fitness, resulting in a much higher probability of being chosen to generate offspring.

Let us now provide a detailed explanation of how the algorithm works. The initial inputs to the algorithm are as follows:

- A dataset, for which the objective is to discover a feature map that accurately classifies the data. This dataset will be split as usual into

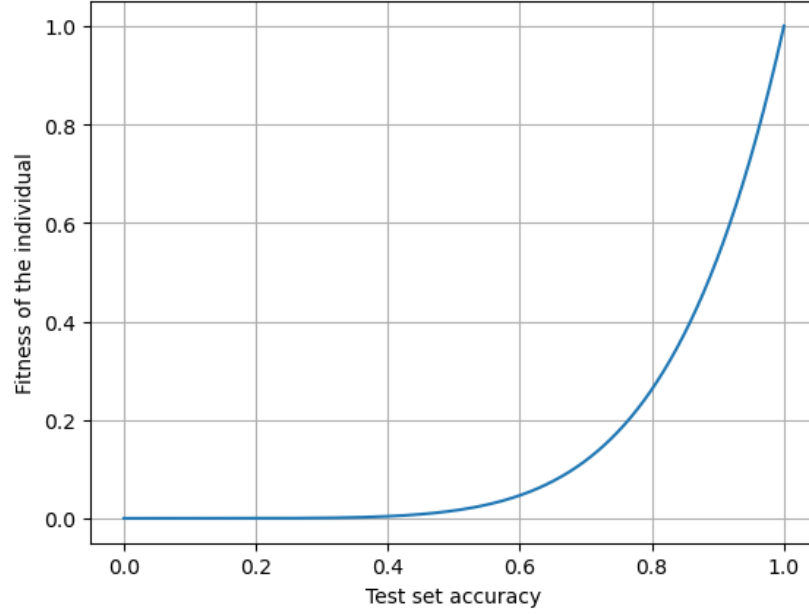


Figure 24: Fitness of an individual as a function of its accuracy on the test set. This choice is made to properly penalize bad performing individuals.

a training set and a test set. We will use the circular mock dataset of Figure 2, as we saw the difficulties of its classification without a genetic algorithm. The dataset automatically defines the number of features d .

- The number k of genes, that is the number of gates of the feature map.
- The number n of qubits that we want the circuit to have.
- The number of individuals m that compose a generation.
- The percentage of individuals that automatically pass to the next generation through elitism. A typical value is 10%.

- The gate set to use to create the feature map. We will use

$$\{X, \sqrt{X}, \text{CNOT}, R_z\},$$

but other choices can be made, depending also on the type of hardware architecture we plan to use for running the algorithm. For example, Google Sycamore processor works by using the gates

$$\{\sqrt{X}, \sqrt{Y}, \sqrt{X+Y}, \text{iSWAP}, \text{CZ}\},$$

so if we want to run the algorithm on that processor we could use this complete set of gates.

- The probability p of a mutation occurring on an individual. A standard value is 10%.
- The target fitness we want to achieve. The algorithm will terminate if an individual reaches that target fitness.
- The maximum number of generations we allow. If no individual achieves the target fitness, the algorithm will terminate when the maximum number of generations is reached.

The first step is to produce the first generation of individuals. This is done by randomly choosing m circuits. Then the genetic algorithm begins, and works as follows:

- The fitness of all the individuals is calculated.

- A fixed percentage of individuals, the ones with the best fitness values, are automatically passed to the next generation.
- The remaining spots in the next generation are filled by children generated with crossover and mutations. For each spot, two parents are chosen by randomly picking in the current generation, with a probability proportional to the fitness. Then two children are generated by crossover, and a random mutation is applied to each of them. The fitness of both children is evaluated and only the best one is kept and placed in the next generation. This procedure is repeated until the next generation is completely filled. Recall that all the generations have the same number of individuals.
- This process is repeated until an individual achieves the target fitness, or the total number of generations is reached.
- The best individual of the last generation is the result of the algorithm. In particular, its phenotype is the best circuit that can be used in a QSVM to separate this particular dataset.

4.3 Results of the Experiments

4.3.1 Circle dataset

We now execute the algorithm on the circular dataset shown in Figure 2. The genetic algorithm was implemented in Python, exploiting Qiskit to run the QSVM when necessary.

An important goal is to construct the smallest possible quantum circuit, minimizing both the number of gates and qubits. This is motivated by the technical limitations of current quantum devices, which struggle with larger, more complex circuits. Consequently we begin with a simple circuit, increasing its depth or number of qubits when necessary.

For our experiment, we set a target fitness of 0.95, corresponding to an accuracy of 0.99, using eq. (59). We begin with a circuit with 2 qubits and 10 genes (gates), thus keeping the circuit compact.

The population size directly influences the efficiency of the algorithm. While larger populations provide more diversity and improve convergence, they significantly slow down the algorithm, particularly when using quantum simulators, like in our case. Since running a single QSVM instance already requires substantial time on a simulator, we opted for a reduced population size of 10 individuals. Although a standard starting point would be 100 individuals, we can begin with a small population and scale up if the results are unsatisfactory. We set a maximum number of 25 generations.

The initial values of parameters are shown in Table (1).

After 25 generations the algorithm reaches a fitness of 86%. The final individual phenotype is represented in Figure (25). This is the best-performing circuit identified by the genetic algorithm. Before exploring its potential, let's first assess whether the algorithm performed as expected. In Figure 26 we plot how the mean fitness of each generation evolves with the passing of generations. We observe an increase in the mean fitness, starting from a very low value in the first generation, which consists of randomly selected circuits,

Parameter	Value
Number of qubits	2
Number of features	2
Number of genes	10
Target accuracy	99%
Maximum number of generation	25
Individuals in a generation	10
Individuals passed through elitism	1
Probability of a mutation	10%

Table 1: Parameters of the genetic algorithm, and their initial value.

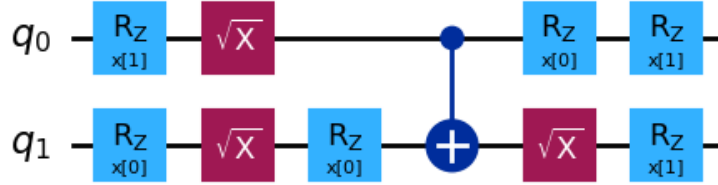


Figure 25: Final circuit produced by the genetic algorithm. Being the phenotype of the best individual of the last generation, this is the circuit with the best performance on this dataset that the genetic algorithm could find.

to a significantly higher value in the final generation. This confirms that the genetic algorithm is effectively improving the fitness of individuals through its internal crossover mechanics. If the mean fitness of the first generation

was similar to that of the last one, it would mean that the circuit was found by random chance, making the genetic algorithm useless. Instead, we clearly see how the algorithm consistently increases the fitness with each generation. The maximum fitness of each generation is also shown in Figure 27.

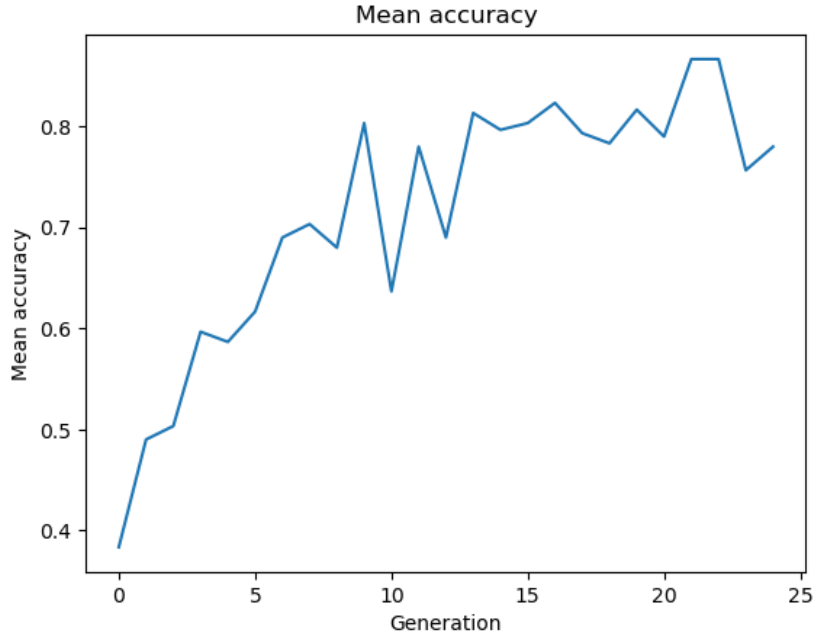


Figure 26: Mean accuracy of each generation. We observe a consistent increase of the accuracy, ensuring that the genetic algorithm is improving the performance of the feature maps from generation to generation.

To verify that the circuit in Figure 25 can truly separate circle-like datasets, we tested its performance on a new artificial dataset. This dataset was generated using the same parameters as the one the circuit was trained on, but with new data points. It is important to ensure that the circuit does not perform well only on the specific dataset it was trained with, but also on

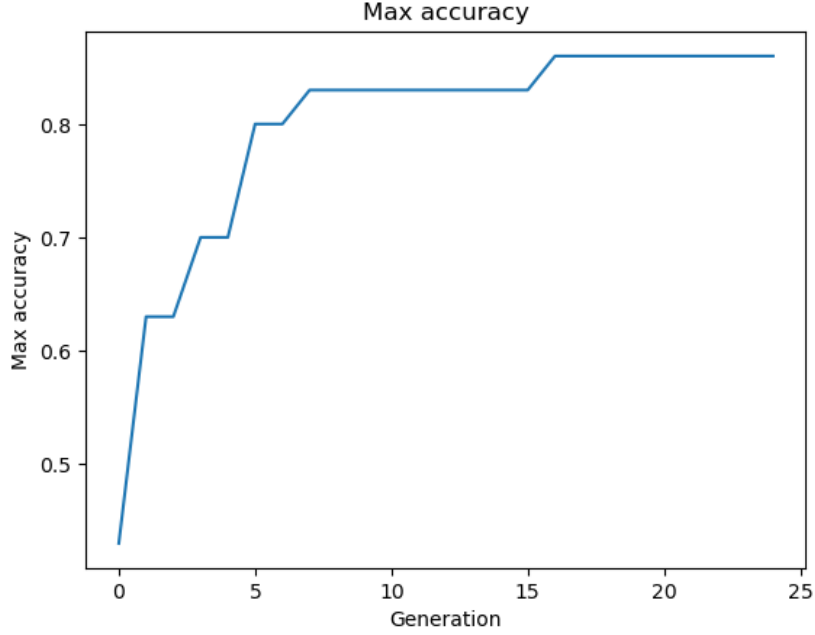


Figure 27: Maximum accuracy of each generation.

slightly varied datasets. The result of the fit is shown in Figure 28, where we achieved 94% accuracy. This indicates that the circuit is highly effective at separating circle-like datasets.

To further demonstrate how the genetic algorithm improved the performance of the individuals, let us classify the dataset with an individual from the first generation. The result, shown in Figure 29, reveals poor performance, similar to the outcomes we observed when using standard features map . So by just randomly picking circuits we are not able to construct well performing individuals. However, by evolving these individuals through crossover, the genetic algorithm is able to generate circuits with high accuracy for this type of datasets.

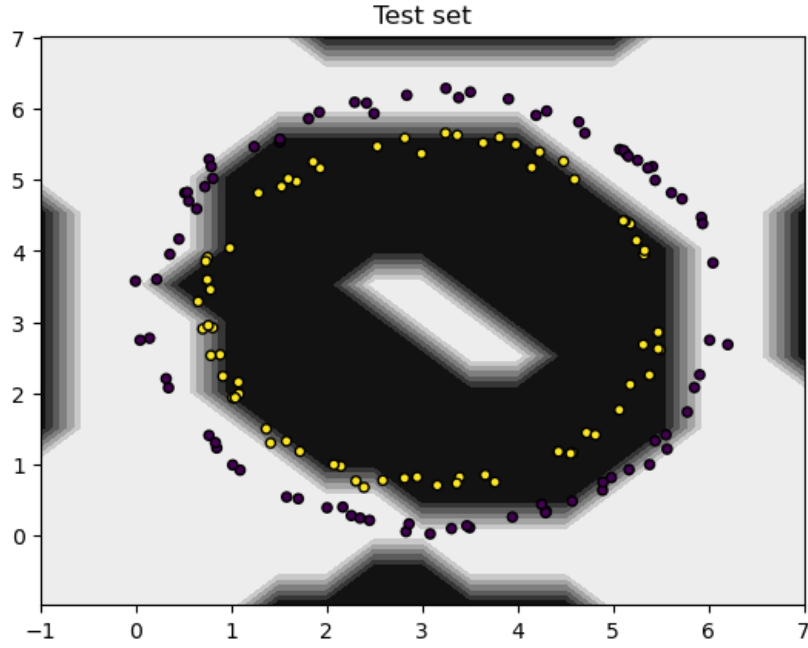


Figure 28: A further test of the resulting circuit of the genetic algorithm on a new dataset, with respect to the one it was trained on. The fact that the circuit is well performing also on this dataset indicates that we actually found a circuit well performing on this sort of datasets.

The fact that it performs well on this dataset is a strong indication of the potential of the final circuit we identified. However, it's natural to question how far this circuit can be pushed, as the dataset, while new, was still generated using the same parameters the genetic algorithm was trained on. What happens if we alter the shape by, for instance, adding noise or changing the radius of the two circles? Will the circuit still be able to classify correctly? Is the circuit able to classify all circle-like datasets? In Figure 30 we used it to classify a different circle dataset, with a gaussian noise added to the points.

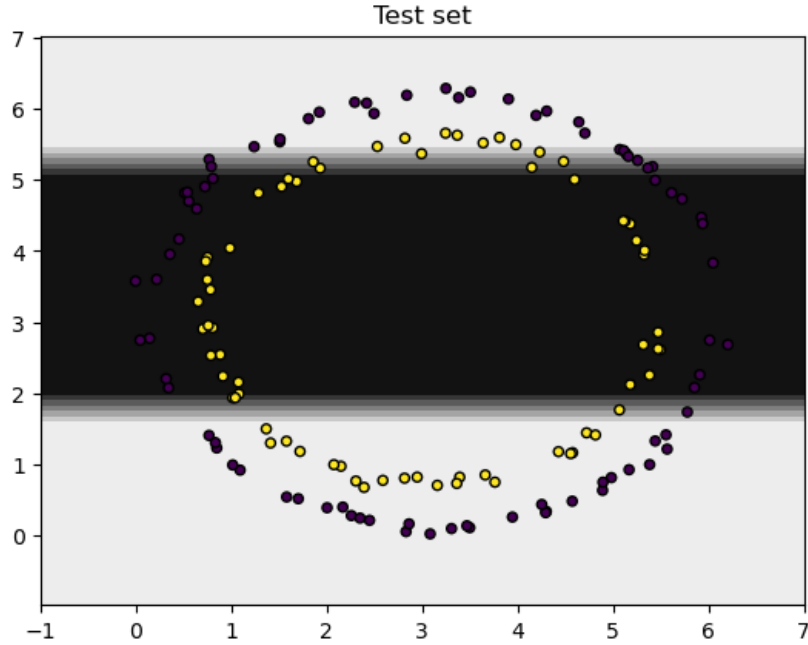


Figure 29: Test of the performance of the best fitness individual from the first generation on the same dataset. We observe how its performance is very low, indicating that the first generation of individuals performs poorly as we expected.

The classification reached a 91% accuracy, indicating a good performance even on this dataset. We can answer affirmatively to the questions we asked.

4.3.2 Moon dataset

Let's now test the performance of the genetic algorithm on another artificial dataset, the moon dataset.[16] The algorithm conceptually works in the same way of the circle dataset, we utilized the same initial parameters as well,

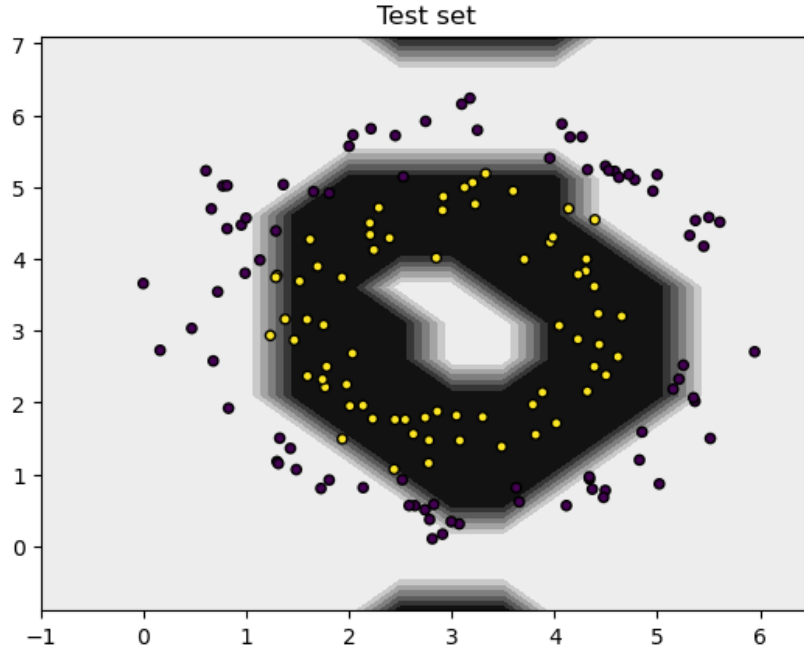


Figure 30: Test of the performance of the best individual on a new dataset with noise. The circuit is revealed to be robust against noise, correctly classifying 91% of the points.

shown in Table (1).

The resulting final circuit is shown in Figure 31, and its performance on the dataset in Figure 32, achieving an accuracy of 90%.

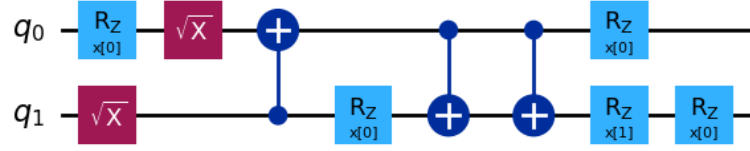


Figure 31: Final circuit produced by the genetic algorithm used on the moon dataset.

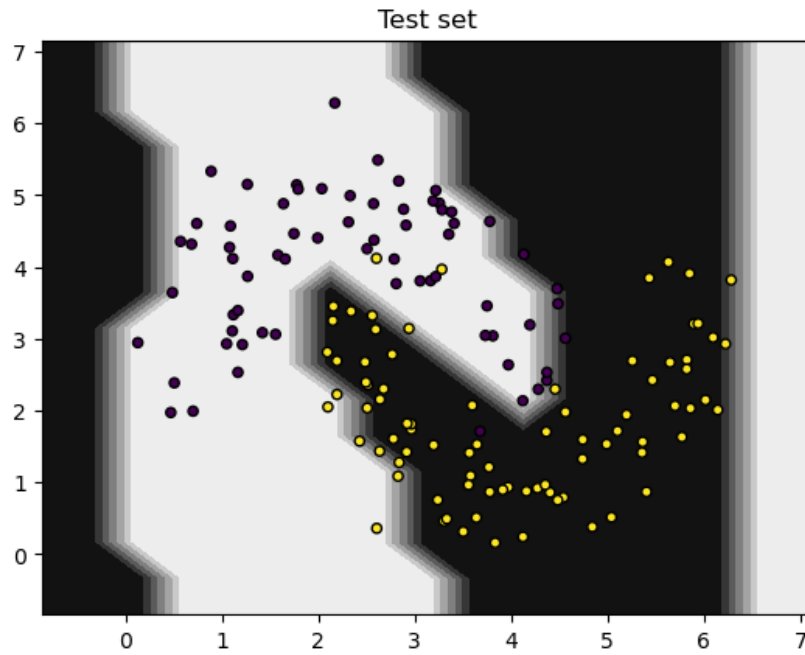


Figure 32: Performance of the final circuit on the moon dataset. Even in this case a good accuracy (90%) is achieved.

5 Support Vector Regressor

The Support Vector Machine can be also used to perform regression tasks: in this case we call it *Support Vector Regressor* (SVR) [5]. We are provided with a dataset with N d -dimensional instances $\{\mathbf{x}_i\}_{i=0,\dots,N-1}$, but this time the labels y_i are continuous values, that is, $y_i \in \mathbb{R}$. As usual the goal is, given a new instance $\tilde{\mathbf{x}}$, to predict its label \tilde{y} .

The model produced by Support Vector Classification depends only on a subset of the training data, because the cost function for building the model does not consider the training points that lie beyond the margin. Analogously, the model produced by SVR depends only on a subset of the training data, because the cost function for building the model ignores any training data close to the model prediction.

As for the Support Vector Classifier (SVC), also the SVR algorithm can be ultimately reduced to a dual minimization problem. The dual variables α_i are just fictitious variables of the dual problem. In this case the function to minimize is [5]:

$$\min_{\alpha, \alpha^*} \frac{1}{2} \sum_{i=0}^{N-1} \sum_{j=0}^{N-1} (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) K_{ij} - \sum_{i=0}^{N-1} (\alpha_i - \alpha_i^*) + \epsilon \sum_{i=0}^{N-1} (\alpha_i + \alpha_i^*), \quad (60)$$

where K is the kernel matrix, ϵ is a hyperparameter of the model, and the α_i are subject to the constraints

$$\alpha_i, \alpha_i^* \geq 0, \quad \forall i = 0, \dots, N-1, \quad (61)$$

$$\sum_{i=0}^{N-1} (\alpha_i - \alpha_i^*) = 0. \quad (62)$$

The prediction rule given a new instance $\tilde{\mathbf{x}}$ is

$$\tilde{y} = \sum_{i=0}^{N-1} (\alpha_i - \alpha_i^*) K(\mathbf{x}_i, \tilde{\mathbf{x}}). \quad (63)$$

Exactly as in the QSVC case, one can build a QSVR by constructing the quantum kernel K with a quantum feature map, and the feature map can be chosen with the genetic algorithm.

5.1 Results of the Experiments

In this example, we utilize the Levy function N.13 in two dimensions [24] to test the QSVR model. The Levy function N.13 is a widely used function to asses the proficiency of machine learning models. Its analytical form is

$$f(x_0, x_1) = \sin^2(3\pi x_0) + (x_0 - 1)^2(1 + \sin^2(3\pi x_1)) + (x_1 - 1)^2(1 + \sin^2(2\pi x_1)). \quad (64)$$

The dataset is artificially generated by applying eq. (64) to a grid of randomly chosen points in the square $[-10, 10]$.

The genetic algorithm is carried out exactly as in the QSVC case. The only difference is that now we need a new way to assess the fitness of an individual, as the accuracy does not make sense in the regression case. One of the most widely used metrics to assess the performance of a regression model is the Root Mean Squared Error (RMSE), defined as

$$\text{RMSE} = \sqrt{\sum_{i=0}^{N-1} \frac{(y_i - f(x_i))^2}{N}}. \quad (65)$$

A higher RMSE means a worse performance, and vice versa. We can thus set

$$\text{fitness} = \frac{1}{\text{RMSE}}. \quad (66)$$

The genetic algorithm is run with the parameters shown in Table 2. The

Parameter	Value
Number of qubits	2
Number of features	2
Number of genes	10
Target fitness	170
Maximum number of generation	25
Individuals in a generation	10
Individuals passed through elitism	1
Probability of a mutation	10%

Table 2: Parameters of the genetic algorithm, and their initial value. The target fitness has been chosen by looking at the performance of the classical SVR with RBF kernel. The other parameters have been chosen with the same philosophy of the QSCV case, that is keeping the circuit compact, and having a maximum number of generations compatible with the computational power available.

mean accuracy of the generations is shown in Figure 33. The final individual fitness is 160, comparable to the one obtained with the RBF classical kernel of 167. Therefore, in the regression task the genetic algorithm was able to

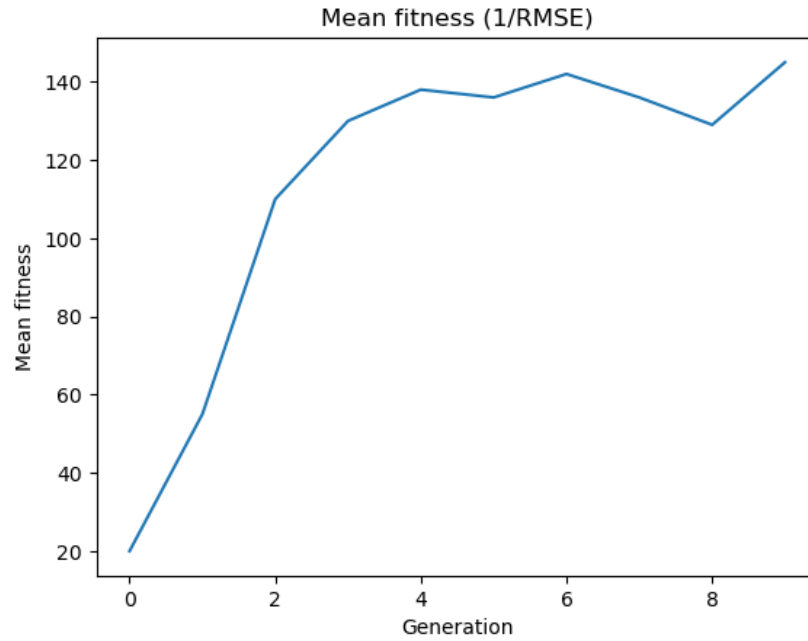


Figure 33: Mean fitness ($1/\text{RMSE}$) of each generation in the QSVR genetic algorithm.

equalize the performance of the classical case, but not outperform it. It is likely that a greater number of individuals and generations is needed, and thus a greater computational power. The genetic QSVR therefore seems to be less powerful than the genetic QSVC.

6 Conclusions and Directions for Future Work

The Quantum Support Vector Machine is an algorithm with great potential to outperform classical supervised classification models. Once fully implemented on a large scalable quantum device, we think it will successfully classify large datasets with many features and complicated patterns, that are intractable with classical methods.

However, we found it to be a very peaky algorithm in the choice of the quantum embedding circuit, so it is essential to use a proper algorithm to construct such a circuit, as using standard rule of thumb is often insufficient. The possibility we explored in this work is the one of using a genetic algorithm.

The algorithm showed great performance on simple 2D datasets, finding the optimal circuit that we were not able to find with other methods, opening the road to a broader application of genetic algorithms for QSVM, and in general in the field of Quantum Machine Learning.

The goal now would be to further study the performance of the genetic algorithm we developed on larger multidimensional datasets, where the QSVM is expected to outperform classical algorithms. Conducting such tests on a quantum simulator, like Qiskit, would require significant computational resources. The ultimate goal therefore is to test the algorithm on a real scalable quantum device, its natural field of application. We are confident that it would show great performance, demonstrating its viability for real-world applications.

References

- [1] Biamonte, Jacob, et al. *Quantum machine learning*. Nature 549.7671 (2017): 195-202.
- [2] Creevey, Floyd M., Charles D. Hill, and Lloyd CL Hollenberg. *GASP: a genetic algorithm for state preparation on quantum computers*. Scientific reports 13.1 (2023): 11956.
- [3] Creevey, Floyd M., et al. *Kernel Alignment for Quantum Support Vector Machines Using Genetic Algorithms*. arXiv preprint arXiv:2312.01562 (2023).
- [4] Cristianini, Nello. *An introduction to support vector machines and other kernel-based learning methods*. Cambridge University Press, 2000.
- [5] Drucker, Harris, et al. *Support vector regression machines*. Advances in neural information processing systems 9, 779-784 (1996).
- [6] Eiben, A. & Smith, Jim. (2003). *Introduction To Evolutionary Computing*, Springer.
- [7] Glick, Jennifer R., et al. *Covariant quantum kernels for data with group structure*. Nature Physics 20.3 (2024): 479-483.
- [8] Havlíček, Vojtěch, et al. *Supervised learning with quantum-enhanced feature spaces*. Nature 567.7747 (2019): 209-212.

- [9] Hubregtsen, Thomas, et al. *Training quantum embedding kernels on near-term quantum computers*. Physical Review A 106.4 (2022): 042431.
- [10] Innan, Nouhaila, et al. *Enhancing quantum support vector machines through variational kernel training*. Quantum Information Processing 22.10 (2023): 374.
- [11] Ivezić, Željko, et al. *Statistics, Data Mining, and Machine Learning in Astronomy*. Princeton University Press. (2014): 134.
- [12] Lloyd, Seth, et al. *Quantum embeddings for machine learning*. arXiv preprint arXiv:2001.03622 (2020).
- [13] Liu, Yunchao, Srinivasan Arunachalam, and Kristan Temme. *A rigorous and robust quantum speed-up in supervised machine learning*. Nature Physics 17.9 (2021): 1013-1017.
- [14] Mitchell, Melanie. *An introduction to genetic algorithms*. MIT press, 1998.
- [15] Park, Jae-Eun, et al. *Practical application improvement to Quantum SVM: theory to practice*. arXiv preprint arXiv:2012.07725 (2020).
- [16] Pedregosa, Fabian, et al. *Scikit-learn: Machine learning in Python*. the Journal of machine Learning research 12 (2011): 2825-2830.
- [17] Rebentrost, Patrick, Masoud Mohseni, and Seth Lloyd. *Quantum support vector machine for big data classification*. Physical review letters 113.13 (2014): 130503.

- [18] Schuld, Maria, and Nathan Killoran. *Quantum machine learning in feature Hilbert spaces*. Physical review letters 122.4 (2019): 040504.
- [19] Schuld, Maria. *Supervised quantum machine learning models are kernel methods*. arXiv preprint arXiv:2101.11020 (2021).
- [20] Schuld, Maria, and Francesco Petruccione. *Machine learning with quantum computers*. Vol. 676. Berlin: Springer, 2021.
- [21] Spall, James C. *Multivariate stochastic approximation using a simultaneous perturbation gradient approximation*. IEEE transactions on automatic control 37.3 (1992): 332-341.
- [22] Sünkel, Leo, et al. *GA4QCO: genetic algorithm for quantum circuit optimization*. arXiv preprint arXiv:2302.01303 (2023).
- [23] Yang, Jiaying, Ahsan Javed Awan, and Gemma Vall-Lloera. *Support vector machines on noisy intermediate scale quantum computers*. arXiv preprint arXiv:1909.11988 (2019).
- [24] Zhou, Xiaojian, et al. *Quantum kernel estimation-based quantum support vector regression*. Quantum Information Processing 23.1 (2024): 29.