

# Genetic algorithm for Quantum Support Vector Machines

Lorenzo Tasca

October 2024

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Classical machine learning</b>	<b>2</b>
2.1	General overview . . . . .	2
2.2	Support vector machine . . . . .	2
2.2.1	Linear SVM . . . . .	2
2.2.2	Kernel SVM . . . . .	5
<b>3</b>	<b>Quantum machine learning</b>	<b>10</b>
3.1	Quantum Support Vector Machine . . . . .	10
3.1.1	Quantum feature map and Quantum Kernels . . . . .	10
3.1.2	Quantum encoding circuits . . . . .	12
3.1.3	QSVM algorithm and Qiskit implementation . . . . .	17
3.1.4	QSVM potential . . . . .	20
3.1.5	Issues of QSVM (mancano delle immagini) . . . . .	23
3.2	Quantum kernel alignment . . . . .	24

# 1 Introduction

## 2 Classical machine learning

### 2.1 General overview

baa

### 2.2 Support vector machine

The Support Vector Machine (SVM) 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 (so 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

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

for appropriate parameters  $w_0 \in \mathbb{R}$  and  $\mathbf{w} \in \mathbb{R}^d$ , which define the 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

$$\frac{1}{2} \|\mathbf{w}\|,$$

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  $\alpha_i$  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

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

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

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 is

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

Once the optimal  $\alpha_i$  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, because we did not allow points to fall inside the margin. One could relax this assumption, modifying the constraint in eq. (5) into

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

where we introduced the slack variables  $\xi_i$ . 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 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.

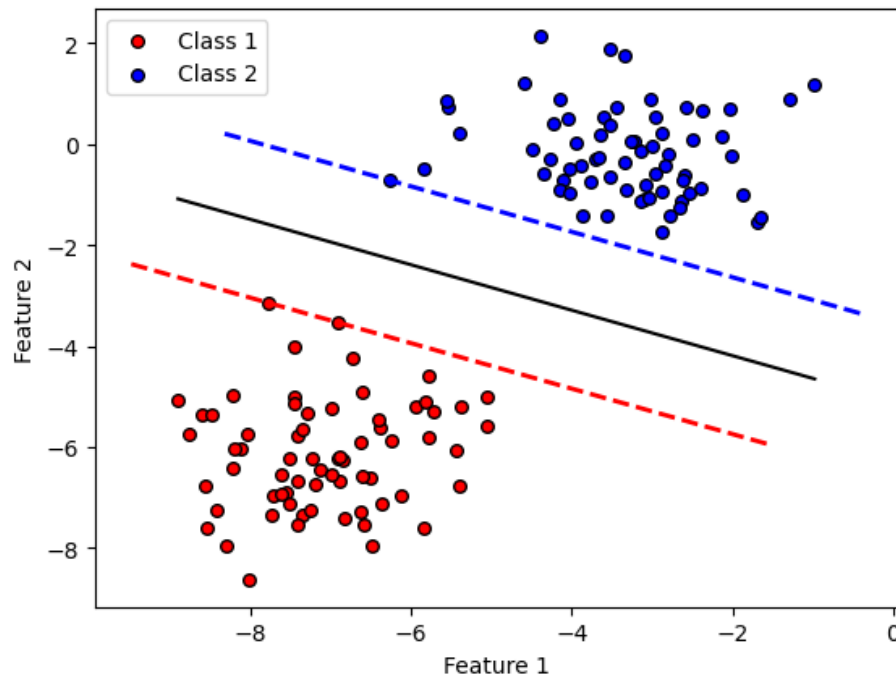


Figure 1: SVM decision boundary and margin border, fitted on a 2 feature mock dataset of 200 instances.

### 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

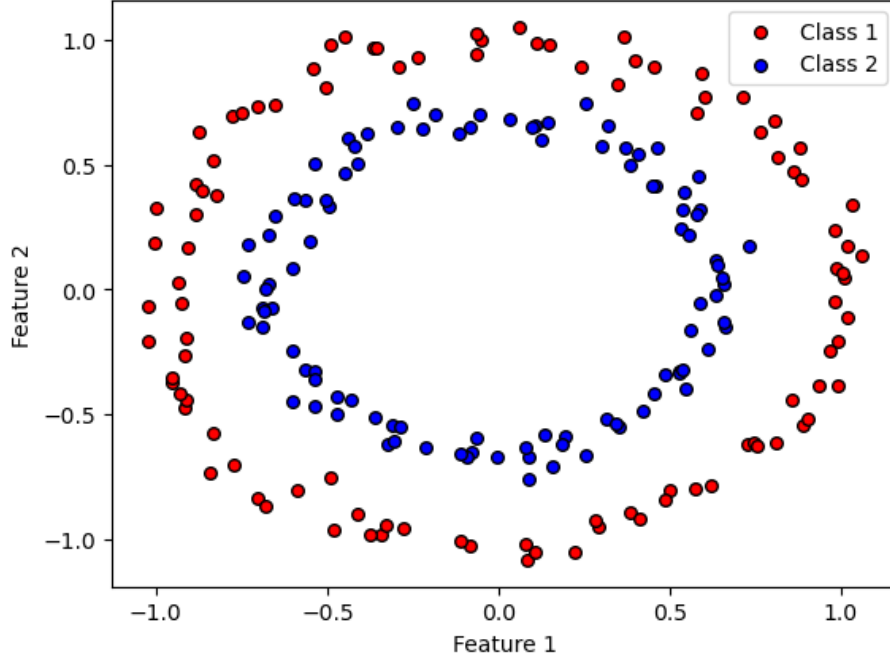


Figure 2: Highly non-linear 2 feature mock dataset with 200 instances.

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 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

$$\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3, \quad \begin{pmatrix} x_0 \\ x_1 \end{pmatrix} \mapsto \begin{pmatrix} x_0^2 \\ x_1^2 \\ \sqrt{2}x_0x_1 \end{pmatrix}. \quad (15)$$

The data of Figure (2) after the application of the feature map  $\phi$  are represented in Figure (3). The dataset is now linearly separable in the feature space, so the

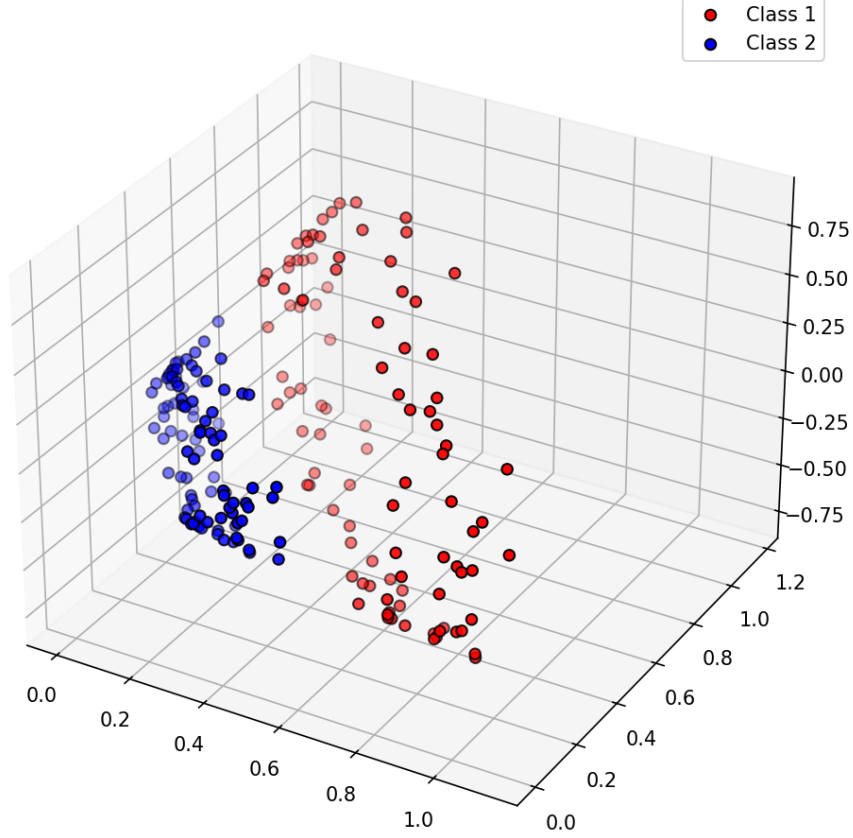


Figure 3: Highly non-linear mock dataset in the feature space after the application of the feature map  $\phi$ . We observe now that the dataset is linearly separable.

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

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

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

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 it produces. We can define the kernel  $K$  as

$$K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R},$$

$$\mathbf{x}, \mathbf{y} \mapsto (\phi(\mathbf{x}), \phi(\mathbf{y})), \quad (18)$$

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).

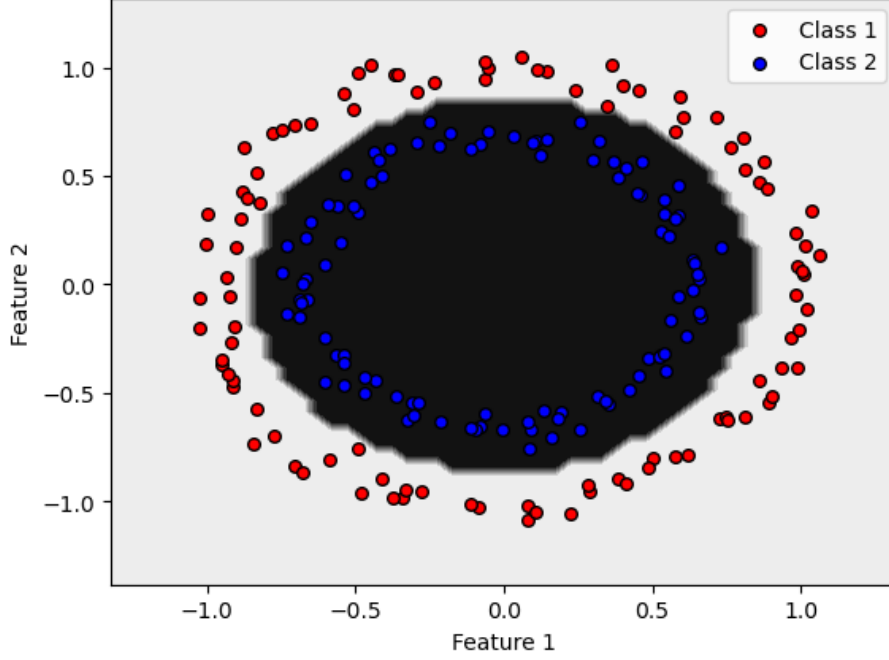


Figure 4: Highly non-linear mock dataset decision boundary, fitted with a SVM using kernel of eq. (21). The white and black areas are the two predicted classes. We observe how, using the kernel, we obtained a non-linear decision boundary.

### 3 Quantum machine learning

Introduction bla bla

#### 3.1 Quantum Support Vector Machine

##### 3.1.1 Quantum feature map and Quantum Kernels

The 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

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

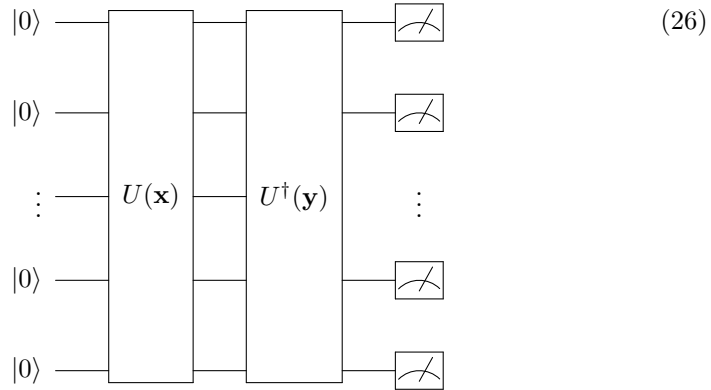
where  $\mathcal{H}_{qubit}$  is the Hilbert space of a single qubit. The dimension of  $\mathcal{H}$  is  $2^n$ . The feature map will be implemented by the means of a parametrized quantum circuit. That means that it exists a unitary operator that depends on  $d$  classical parameters  $U(\mathbf{x}) = U(x_0, \dots, x_{d-1})$  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 LHS of eq. (27) is the probability of measuring  $000 \cdots 0$ , which according to quantum mechanics can be calculated as

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

□

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.

### 3.1.2 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  $2^n$ , 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.

- Amplitude encoding: this can be applied with instances that are normalized vector, 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 sending

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

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

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

We went back to a polynomial kernel.

- By creating more copies of the feature vector of 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 send

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

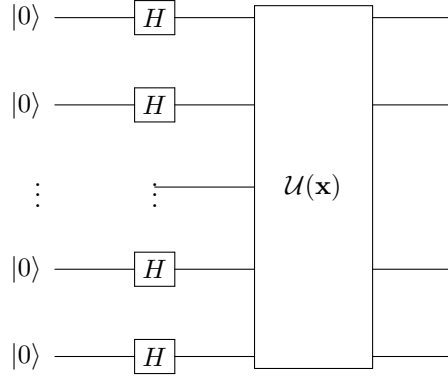
therefore the  $i$ -th qubit will be in the state

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

This generates the kernel

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

- **ZZ feature map:** The previously discussed encoding maps where quantum maps that in the end generated a kernel which is easy to calculate with classical tools, for example a polynomial kernel or a cosine kernel. It is important now to discuss quantum encoding circuits that creates a kernel which is not possible to write as an easy classical kernel, otherwise there is no advantage in using quantum computing for this task if the kernel could have been 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. 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 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 a way such 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 Z feature map. The circuit decomposed in elementary quantum gates is shown in Figure 5. The other possible choice is

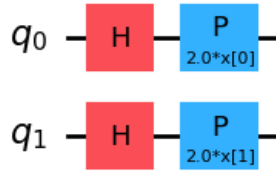


Figure 5: Z feature map, with 2 qubits and 2 features, decomposed in terms of elementary gates. The gate  $P$  is defined as  $P(\phi) = (1, 0; 0, e^{i\phi})$ . We observe how the circuit is parametrized by the coefficients of eq. (39).

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

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

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 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 how

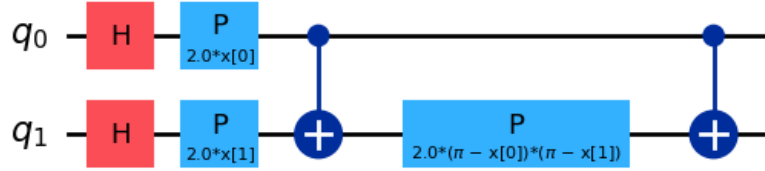


Figure 6: ZZ feature map, with 2 qubits and 2 features, decomposed in terms of elementary gates. The gate  $P$  is defined as  $P(\phi) = (1, 0; 0, e^{i\phi})$ . We observe how the circuit is parametrized by the coefficients of eq. (42) and (43).

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 interaction of the ZZ type, but also for example YY type or ZY type. An example with Y and ZY interactions is shown in Figure 45.



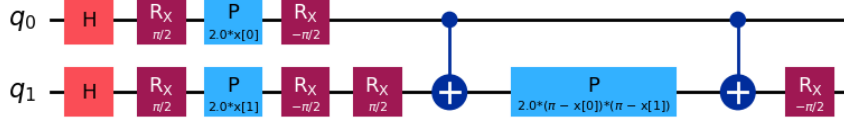


Figure 7: Pauli feature map, with 2 qubits and 2 features, decomposed in term of elementary gates, constructed with Y and YZ interactions.

- 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.

### 3.1.3 QSVM algorithm and Qiskit implementation

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

- 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

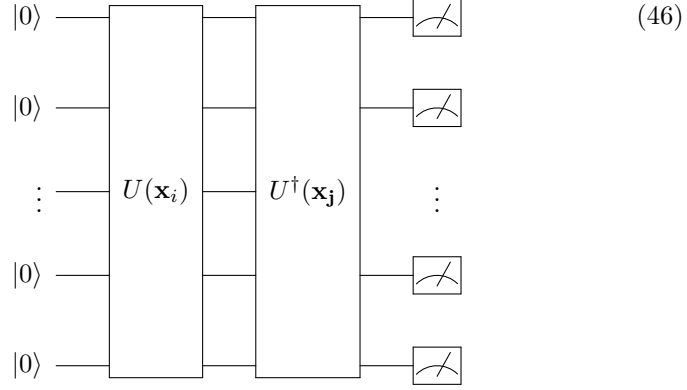
$$y_i \in \{-1, 1\},$$

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 build  $U(\mathbf{x}_i)$  and  $U^\dagger(\mathbf{x}_j)$ , in order to construct the circuit



and evaluate the matrix element

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

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 that 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  $\alpha$  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 we have

---

**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:** 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 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 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.

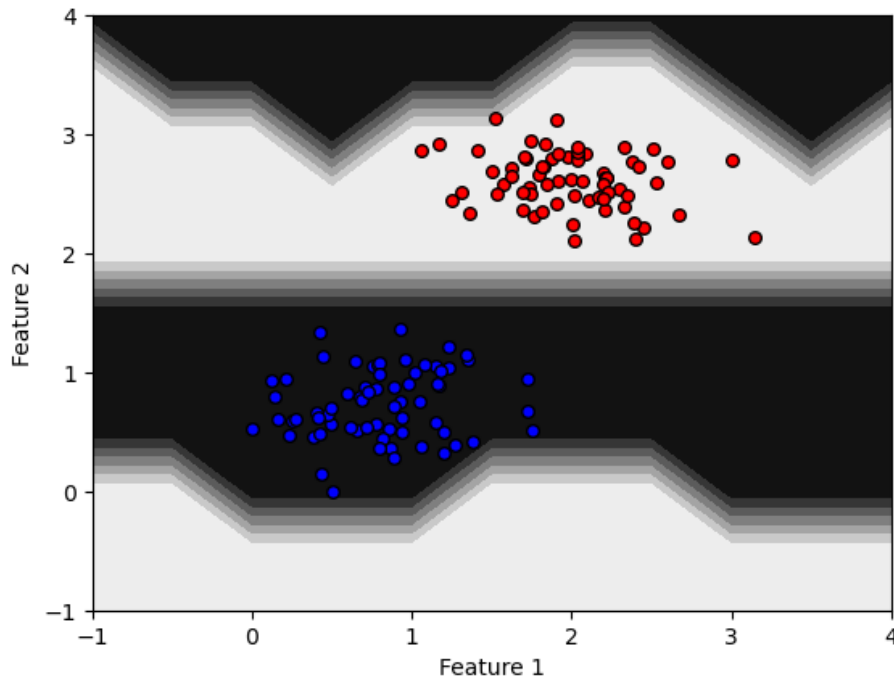


Figure 8: QSVM decision boundary with ZFeatureMap, fitted on a 2 feature mock dataset of 200 instances.

### 3.1.4 QSVM potential

In order to study the potential of the QSVM algorithm we could try and 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.

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

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

$\Delta$  controls how big is the gap between of the two classes. We show this dataset in Figure 9. We observe how the dataset of Figure 9 shows a very complicated

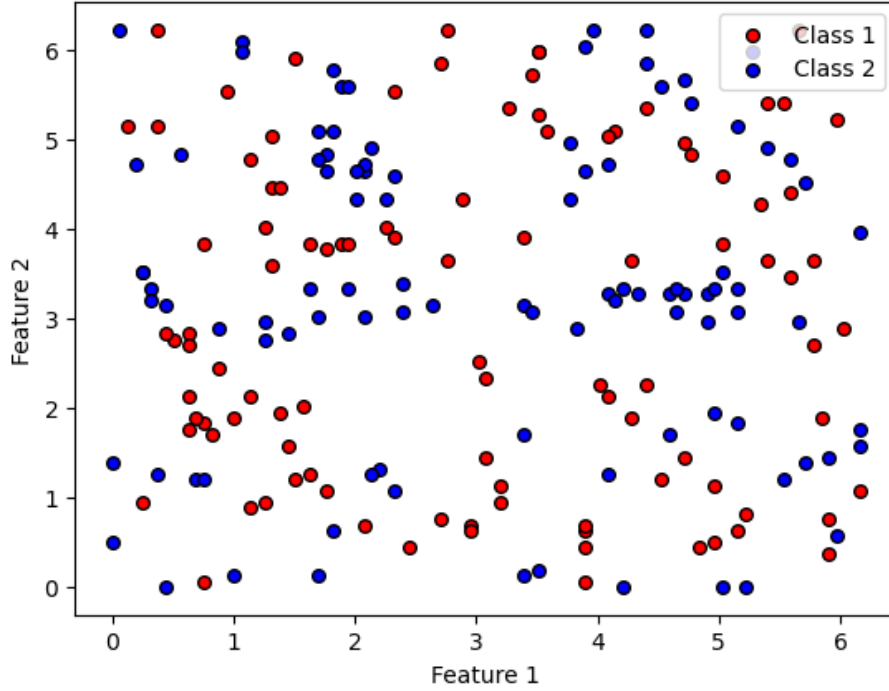


Figure 9: Artificial ad hoc dataset, generated with  $\Delta = 0.3$ . This dataset, by construction, is easily separable by ZZFeatureMap.

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 apply the classical 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

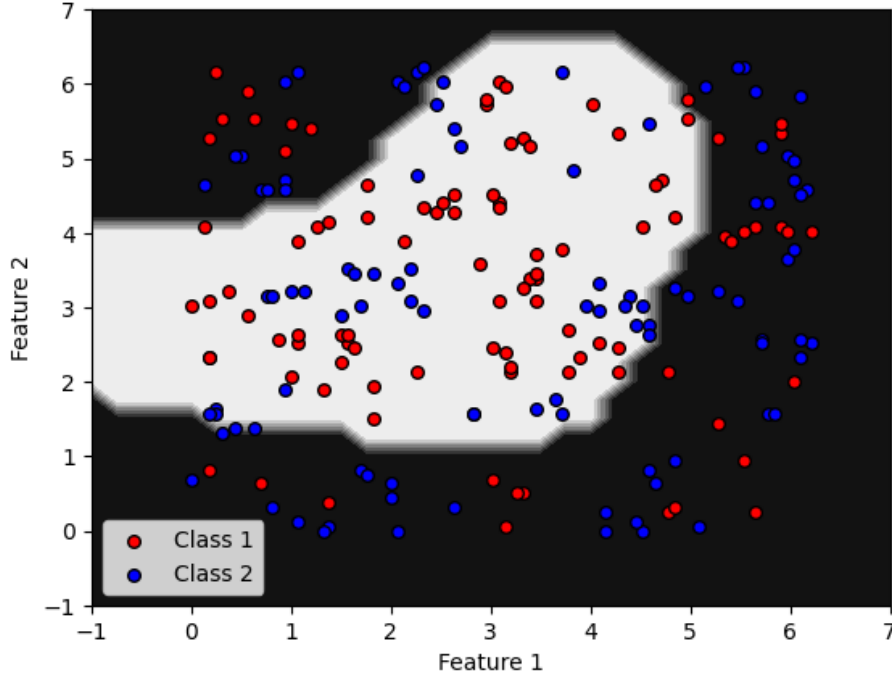


Figure 10: Result of the classical fit with RBF kernel on the artificial dataset. We observe how the SVM could not find the correct separation between the classes.

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 the result we hoped to find.

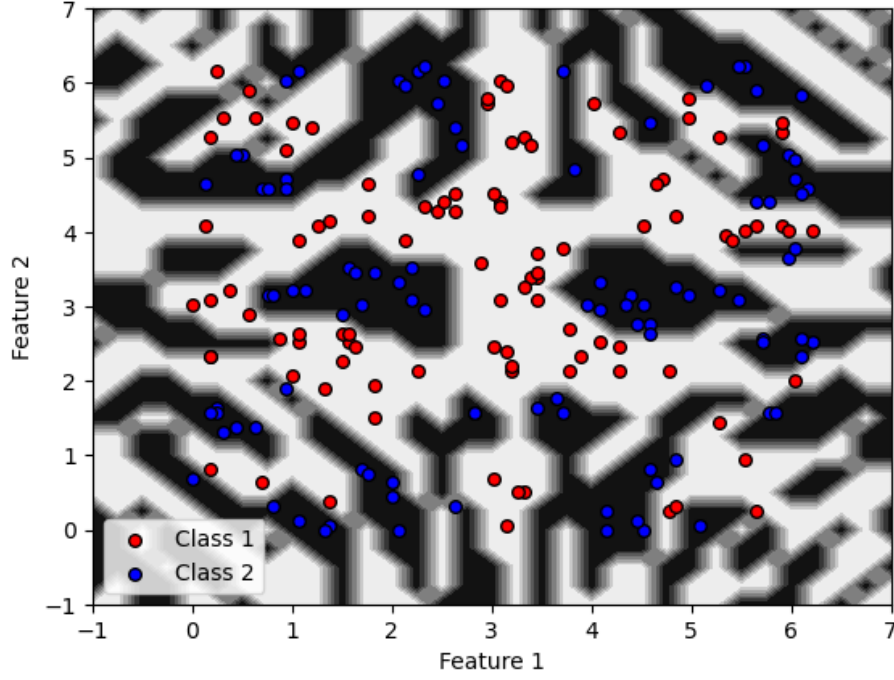


Figure 11: Result of the quantum fit on the artificial dataset, with ZZFeatureMap. We observe how the QSVM found the correct separation between the classes.

### 3.1.5 Issues of QSVM (mancano delle immagini)

In a QSVM algorithm the only choice the user has to make is the choice of the feature map, that is the choice of the quantum encoding circuit. It is easy to convince ourselves that this choice is crucial, and determines if the algorithm will be successful or not. In the classical case we also have to choose a kernel, but we found that this choice is much less delicate than in the quantum case. In fact, for example, the RBF kernel works very well for most applications. In the quantum case instead, we observed that the choice of the feature map is very tricky. If the wrong feature map is chosen, the QSVM performs very poorly even on very simple datasets, as we will now illustrate with some examples. Let's consider the dataset shown in Figure , which is simple enough to be expected to be separated without too many difficulties. The RBF kernel separates this

dataset very well. Let's try to apply some of the quantum kernels we discussed in the previous sections.

We see how the results of the fits are not satisfying. The accuracy may not be too low in some cases, but the form of the decision boundary tells us that the algorithm did not find the correct pattern for the data. With more complex dataset the algorithm may fail completely, like shown in Figure .

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 a lot of 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.

### 3.2 Quantum kernel alignment

In the previous sections, we discussed that 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  $U(\mathbf{x}, \lambda)$ , where we introduced  $k$  parameters  $\lambda = (\lambda_0, \dots, \lambda_{k-1})$  inside the feature map. The quantum circuit now has  $d + k$  parameters,  $d$  are passed through  $\mathbf{x}$  and  $k$  through  $\lambda$ . 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  $\lambda$  parameters instead are just parameters which 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  $\alpha$  is the solution to the convex optimization problem of the SVM, while  $\lambda$  is the variational parameter of the quantum circuit which must be trained. The optimization over  $\alpha$  is performed as usual for the SVM. The optimization over  $\lambda$  is performed using a Classical Simultaneous Stochastic Approximation (SPSA). 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 are:

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

- 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  $\alpha_{\pm}$  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  $\lambda_k$ :

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

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

- Update  $\lambda$  depending on the gradient  $g_k(\alpha_k, \lambda_k)$  and the learning rate  $a$ . 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  $K_{\lambda^*(\underline{x}_i, \underline{x}_j)}$ .