# Quantum Kernel Alignment with Stochastic Gradient Descent

Gian Gentinetta[1], David Sutter[2], Christa Zoufal[2], Bryce Fuller[3], and Stefan Woerner[2]

[1] *Institute of Physics, École Polytechnique Fédérale de Lausanne*
[2] *IBM Quantum, IBM Research Europe – Zurich*
[3] *IBM Quantum, IBM T.J. Watson Research Center, Yorktown Heights*

## Abstract

Quantum support vector machines have the potential to achieve a quantum speedup for solving certain machine learning problems. The key challenge for doing so is finding good quantum kernels for a given data set — a task called kernel alignment. In this paper we study this problem using the Pegasos algorithm, which is an algorithm that uses stochastic gradient descent to solve the support vector machine optimization problem. We extend Pegasos to the quantum case and and demonstrate its effectiveness for kernel alignment. Unlike previous work which performs kernel alignment by training a QSVM within an outer optimization loop, we show that using Pegasos it is possible to simultaneously train the support vector machine and align the kernel. Our experiments show that this approach is capable of aligning quantum feature maps with high accuracy, and outperforms existing quantum kernel alignment techniques. Specifically, we demonstrate that Pegasos is particularly effective for non-stationary data, which is an important challenge in real-world applications.

## 1 Introduction

Finding a practically relevant problem that can be solved better or faster with a quantum computer compared to the best classical implementation is a grand challenge in the field of quantum computing. Quantum machine learning problems are potential candidates for demonstrating a quantum advantage [1–4] and it has been shown for certain artificially constructed data sets that quantum support vector machines (QSVMs) offer an exponential speedup compared to any known classical algorithm [4]. It is an ongoing topic of research to determine whether similar speedups are heuristically available for practical problems.

The performance of QSVMs are dependent upon the degree to which a quantum kernel encodes information about the problem [5], yet little is known about how to engineer quantum kernels which exploit the structure of an arbitrary dataset. One solution to this problem is to choose a parameterized family of quantum kernel functions and attempt to learn a kernel which is well suited to a given dataset. This process of localizing a good kernel for class of given data is referred to as *kernel alignment* [6–8]. A core limitation of this approach is that apriori it is not easy to optimize over a parameterized space of kernels and the overheads introduced can be prohibitive.

In this work, we present a novel method to speedup the kernel alignment problem via the use of the Pegasos algorithm [9, 10]. Pegasos is an iterative gradient-based algorithm that solves the primal SVM optimization problem and allows for a direct integration of kernel alignment. Trainable parameters in the quantum circuit that implements a feature map can be optimized simultaneously with the weights of the SVM problem. This drastically reduces the complexity of quantum kernel alignment (QKA) compared to the dual approach [8] where the full SVM problem has to be solved as a subroutine in the optimization of the kernel parameters. Furthermore, the iterative nature of Pegasos is ideal for online machine learning and non-stationary data as it enables easy continuation of training as well as unlearning the impact of outdated training data.

In an experimental demonstration we show that we can successfully classify the covariant data set proposed in [8] using Pegasos kernel alignment. We further show in a hardware experiment that Pegasos can adapt to non-stationary data in real-time by continuously retraining

the parametrization while keeping the training accuracy high. This may be relevant in practice where non-stationarity is unavoidable.

# 2 Support vector machines

## 2.1 Classical support vector machines

Given an unknown probability distribution $P(\mathbf{x}, y)$ for data vectors $\mathbf{x} \in \mathbb{R}^r$ and class membership labels $y \in \{-1, 1\}$, we draw a set of training data $X = \{\mathbf{x}_1, \ldots, \mathbf{x}_M\}$ with corresponding labels $y = \{y_1, \ldots, y_M\}$. The support vector machine (SVM) [11–13] defines a classification function $c : \mathbb{R}^r \to \{-1, 1\}$ implementing the trade-off between accurately predicting the true labels and maximizing the orthogonal distance between the two classes within the training data.

For a given feature map $\varphi : \mathbb{R}^r \to \mathbb{R}^s$, the decision boundary is

$$c_{\text{SVM}}(\mathbf{x}) = \text{sign}\left[\mathbf{w}^\top \varphi(\mathbf{x})\right],$$

for $\mathbf{w} \in \mathbb{R}^s$. L et $\mathbf{w}^\star$ denote the hyperplane defined as the solution to the primal optimization problem

$$\text{(primal problem)} \quad \min_{\mathbf{w} \in \mathbb{R}^s} \left\{ \frac{\lambda}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^{M} \max\left\{0, 1 - y_i \langle \mathbf{w}^{,} \varphi(\mathbf{x}_i) \rangle \right\} \right\}, \quad (1)$$

where the term containing $\lambda > 0$ provides regularization and the second term is a sum over the hinge losses of the data points from the training set.

The optimization problem (1) features a dual formulation as

$$\text{(dual problem)} \quad \begin{cases} \max_{\alpha_i \in \mathbb{R}} & \sum_{i=1}^{M} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{M} \alpha_i \alpha_j y_i y_j \, k(\mathbf{x}_i, \mathbf{x}_j) - \frac{\lambda}{2} \sum_{i=1}^{M} \alpha_i^2 \\ \text{s.t.} & 0 \leq \alpha_i \quad \forall i = 1, \ldots, M, \end{cases} \quad (2)$$

where $k(\mathbf{x}, \mathbf{y}) := \langle \varphi(\mathbf{x}), \varphi(\mathbf{y}) \rangle$ denotes the kernel function. From (2), we see that solving the dual requires the evaluation of the full kernel matrix $K \in \mathbb{R}^{M \times M}$ defined via its entries

$$K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j) \quad \text{for} \quad i, j = 1, \ldots, M.$$

Given $K$, the dual optimization can be restated as a convex quadratic program (see [10]) and therefore be solved in polynomial time regardless of the dimension of the feature vectors [14]. For this reason the dual formulation is often favored in practice.

Interestingly, there is an algorithm to solve the primal problem called Pegasos [15], which employs stochastic sub-gradient descent to minimize the objective. Unlike a direct implementation of (1), Pegasos only requires access to the kernel matrix $K$ and thus avoids explicitly representing high dimensional feature vectors. A description of the original algorithm and its quantum version can be found in [10, 15].

## 2.2 Quantum support vector machines

A quantum support vector machine is simply a support vector machine equipped with a quantum kernel function. Because there are quantum circuits that can be run efficiently on quantum hardware but cannot be simulated efficiently on any classical computer, we can define kernel functions that can only be evaluated efficiently on a quantum computer. Note that this does not yet answer the question if these kernels are useful. Furthermore, while QSVMs show potential for practical quantum advantage [4, 8] they also suffer from disadvantages such as exponential concentration phenomena inducing flat training landscapes (also called barren plateaus) [16–19]. In fact, exponentially vanishing gradients occur for quantum kernels [16] if (i) the quantum feature map forms (an approximate) 2−design; (ii) the chosen quantum feature map leads to product states whose fidelity is on average exponentially small; (iii) the quantum hardware used to evaluate the quantum kernel entries is impacted by too much noise, e.g., Pauli noise acting before and after each gate.

We can, therefore, conclude that quantum kernels have to be designed very carefully to allow for efficient trainability.

More formally, consider a feature map

$$\psi \colon \mathbb{R}^r \to \mathcal{S}(2^q)$$
$$\mathbf{x} \mapsto |\psi(\mathbf{x})\rangle\langle\psi(\mathbf{x})| \, ,$$

where $\mathcal{S}(2^q)$ denotes the space of density matrices on $q$ qubits [2]. The kernel function is then given by

$$k(\mathbf{x}, \mathbf{y}) = \text{tr}\big[\, |\psi(\mathbf{y})\rangle\langle\psi(\mathbf{y})| \, |\psi(\mathbf{x})\rangle\langle\psi(\mathbf{x})| \, \big] = |\langle\psi(\mathbf{x})|\psi(\mathbf{y})\rangle|^2 \,. \tag{3}$$

Figure 1 explains how to define and evaluate a kernel function via a quantum circuit.
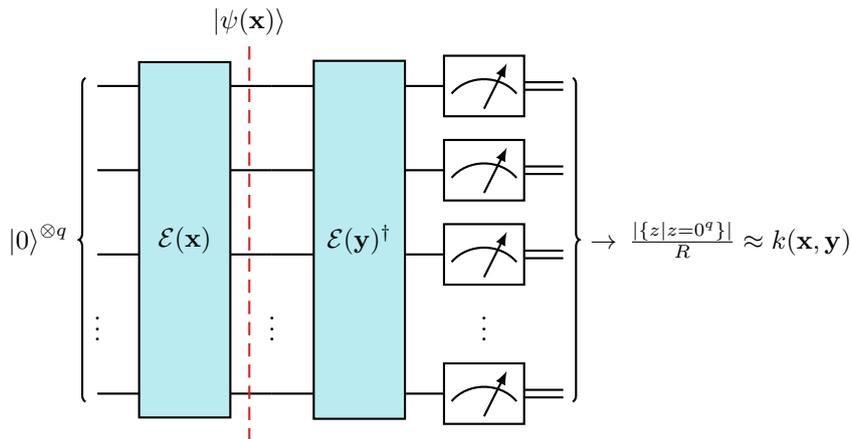


Figure 1: **Quantum kernel evaluation** [10, 20]**:** Let $\mathcal{E}(\mathbf{x})$ denote a parametrized unitary which defines the feature map $|\psi(\mathbf{x})\rangle = \mathcal{E}(\mathbf{x})|0\rangle^{\otimes q}$. Preparing the state $\mathcal{E}(\mathbf{y})^\dagger \mathcal{E}(\mathbf{x})|0\rangle^{\otimes q}$ and measuring all of the qubits in the computational basis, a bit string $z \in \{0, 1\}^q$ is determined. Repeating this process $R$-times, the frequency of the all zero outcome approximates the kernel value $k(\mathbf{x}, \mathbf{y})$ in (3).

It remains unknown how to reliably find good quantum kernels for practical datasets. Kernel alignment is one method which has shown promise in tailoring quantum kernels to specific datasets. In the remainder of this work, we propose a method for training a QSVM which allows for the quantum kernel to be simultaneously aligned to the target dataset during training.

## 3 Quantum kernel alignment with Pegasos

### 3.1 Quantum kernel alignment

Here we discuss how trainable parameters can be used to optimize a feature map to match the training data in the QSVM problem. Choosing the quantum circuit used as a feature map in the QSVM problem is an important and non-trivial task. The choice of the feature map decides whether the data is (linearly) separable in the feature space and, thus, whether the QSVM succeeds in classifying the data. While there are some tricks, such as exploiting symmetries of the system [21], these usually require some knowledge about the structure of the data. Furthermore, even when it is known that data exhibits particular structure, it is often not obvious how to build a circuit which takes advantage of this knowledge. Instead, Glick *et al.* [8] proposed to automate the fine-tuning of the feature map by optimizing the training loss with respect to additional trainable parameters in the circuit.

In the following, we denote the feature state prepared from our trainable quantum circuit as $|\psi_\theta(\boldsymbol{x})\rangle = |\psi(\theta, \boldsymbol{x})\rangle$. Again, the kernel function is defined as the overlap $k_\theta(\boldsymbol{x}, \boldsymbol{y}) = |\langle\psi_\theta(\boldsymbol{x})|\psi_\theta(\boldsymbol{y})\rangle|^2$. With this notation, we can extend the primal and dual SVM optimization problems given by (1)
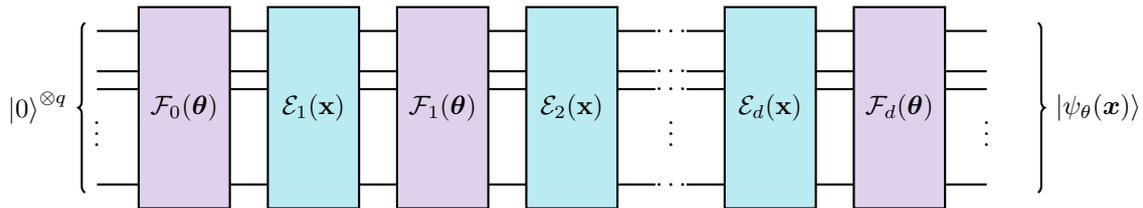
Figure 2: **Trainable feature map**: The feature map composed of unitary gates $\mathcal{E}_i(\mathbf{x})$ used to upload the datum $\mathbf{x}$ is expanded by embedding unitaries $\mathcal{F}_i(\theta)$. The parameters $\boldsymbol{\theta}$ are trained to minimize the training loss.

and (2), by taking into consideration the optimization of $\theta$. This results in

$$\text{(QKA: primal problem)} \quad \min_{\theta \in \mathbb{R}^d} \min_{\mathbf{w} \in \mathbb{R}^s} \left\{ \frac{\lambda}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^{M} \max\left\{0, 1 - y_i \langle \mathbf{w}, \psi_\theta(\mathbf{x}_i) \rangle \right\} \right\} \tag{4}$$

and

$$\text{(QKA: dual problem)} \left\{ \begin{array}{ll} \min\limits_{\theta \in \mathbb{R}^d} \max\limits_{\alpha_i \in \mathbb{R}} & \sum_{i=1}^{M} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{M} \alpha_i \alpha_j y_i y_j \, k_\theta(\mathbf{x}_i, \mathbf{x}_j) - \frac{\lambda}{2} \sum_{i=1}^{M} \alpha_i^2 \\ \text{s.t.} & 0 \leq \alpha_i \quad \forall i = 1, \ldots, M \,. \end{array} \right. \tag{5}$$

In [8] the dual problem (5) is solved through a nested optimization. In an outer loop the trainable parameters $\theta$ are optimized, where in every optimization step the dual QSVM problem is solved to optimize the hyperplane parameters $\boldsymbol{\alpha}$. While this is a viable option for smaller data sets, solving the dual optimization problem becomes prohibitively expensive if the size of the data set $M$ is large. In fact, to reach an accuracy of $\varepsilon$ with respect to the ideal hyperplane, a total of $\mathcal{O}\left(M^{4.67}/\varepsilon^2\right)$ circuit evaluations are required [10]. Thus, solving the dual problem as a subroutine in the optimization of the trainable parameters makes this algorithm unfeasible for large $M$.

Instead of solving the min-max problem in (5), we utilize the min-min property of the primal formulation of the QKA problem (4). The min-min problem has the advantage that the two minimizations can be done simultaneously where the min-max problem has to be solved sequentially which is considerably more expensive. To solve the primal problem, the Pegasos algorithm is employed and adapted to also optimize the trainable parameters.

## 3.2 Pegasos quantum kernel alignment

In the following, we derive how we can adapt the stochastic gradient descent (SGD) based Pegasos algorithm to solve the primal formulation of the quantum kernel alignment problem (4). Pegasos is a classical algorithm that finds the optimal weights $\mathbf{w}$ in the primal SVM problem through SGD [9]. Compared to the more commonly used dual solvers, Pegasos has been shown to scale favourably for large training data sets in the presence of shot noise [10].

The main idea of the Pegasos algorithm is to write the the weights as a linear combination of the feature vectors, i.e.,

$$\mathbf{w}_\tau = \frac{1}{\lambda \tau} \sum_{t=1}^{\tau} \alpha_t y_{i_t} \psi(\mathbf{x}_{i_t}) \,.$$

Here we sum over the iterations, where $(\mathbf{x}_{i_t}, y_{i_t})$ is the datum (with features and label) sampled in iteration $t$ and $\alpha_t \in \{0, 1\}$ indicates whether the datum has been chosen as a support vector. Using this ansatz for $\mathbf{w}$ allows us to write the inner product $\langle \mathbf{w}, \psi(\mathbf{x}_j) \rangle$ in (4) in terms of the kernel function such that direct access to the feature vectors is not required.

We now extend this ansatz by adding trainable parameters $\theta_t$ to the feature map via

$$\mathbf{w}_\tau = \frac{1}{\lambda \tau} \sum_{t=1}^{\tau} \alpha_t y_{i_t} \psi_{\theta_t}(\mathbf{x}_{i_t}) \,. \tag{6}$$

4

In addition to updating $\mathbf{w}$, we also perform an update step on $\theta$. This is possible thanks to the min-min structure of (4).

Assuming we know the values of $\alpha_t$ for $t < \tau$ and $\theta_t$ for $t \leq \tau$, we next derive the values for $\alpha_\tau$ and $\theta_{\tau+1}$ using SGD. The optimization step starts by uniformly sampling $(\mathbf{x}_{i_\tau}, y_{i_\tau})$ from the training data set. The loss function for this data point is then given as

$$f^\tau(\theta, \mathbf{w}) = \frac{\lambda}{2} \|\mathbf{w}\|^2 + \max\left[0, 1 - y_{i_\tau} \langle \mathbf{w}, \psi_\theta(\mathbf{x}_{i_\tau}) \rangle\right]. \tag{7}$$

We first calculate the gradient with respect to $\mathbf{w}$ as

$$\frac{\partial f^\tau}{\partial \mathbf{w}}(\mathbf{w}, \theta_\tau) = \begin{cases} \lambda \mathbf{w}, & \text{if } y_{i_{\tau+1}} \langle \mathbf{w}, \psi_{\theta_\tau}(\mathbf{x}_{i_\tau}) \rangle > 1 \\ \lambda \mathbf{w} - y_{i_\tau} \psi_{\theta_\tau}(\mathbf{x}_{i_\tau}), & \text{otherwise} . \end{cases}$$

The inner product in the if-condition can be evaluated using the kernel trick as

$$\langle \mathbf{w}, \psi_{\theta_\tau}(\mathbf{x}_{i_\tau}) \rangle = \frac{1}{\tau - 1} \sum_{t=1}^{\tau-1} \alpha_t y_{i_t} k_{\theta_t, \theta_\tau}(\mathbf{x}_{i_t}, \mathbf{x}_{i_\tau}) ,$$

where we introduced the pseudo-kernel $k_{\theta, \phi}(\mathbf{x}, \mathbf{y}) = \langle \psi_\theta(\mathbf{x}), \psi_\phi(\mathbf{y}) \rangle = |\langle \psi_\theta(\boldsymbol{x}) | \psi_\phi(\boldsymbol{y}) \rangle|^2$. Next, the weights are updated according to a learning rate of $\mu_\tau = 1/\lambda\tau$, which leads to

$$\begin{aligned} \mathbf{w}_\tau &= \mathbf{w}_{\tau-1} - \frac{1}{\lambda\tau} \frac{\partial f^\tau}{\partial \mathbf{w}}(\theta_\tau, \mathbf{w}_{\tau-1}) \\ &= \begin{cases} (1 - \frac{1}{\tau})\mathbf{w}_{\tau-1}, & \text{if } y_{i_\tau} \langle \mathbf{w}, \psi_{\theta_\tau}(\mathbf{x}_{i_\tau}) \rangle > 1 \\ (1 - \frac{1}{\tau})\mathbf{w}_{\tau-1} + \frac{1}{\lambda\tau} y_{i_\tau} \psi_{\theta_\tau}(\mathbf{x}_{i_\tau}), & \text{otherwise.} \end{cases} \\ &= \frac{1}{\lambda\tau} \sum_{t=1}^{\tau} \alpha_t y_{i_t} \psi_{\theta_t}(\mathbf{x}_{i_t}) , \end{aligned}$$

where $\alpha_\tau = \mathbb{1}\left[y_{i_\tau} \langle \mathbf{w}, \psi_{\theta_\tau}(\mathbf{x}_{i_\tau}) \rangle \leq 1\right]$ for $\mathbb{1}[\cdot]$ being the indicator function. Having found the value of $\alpha_\tau$, we determine the trainable parameters for the next step $\theta_{\tau+1}$. Note that if we are in the regime where the hinge-loss in (7) is 0, the gradient with respect to $\theta$ will vanish. In the non-zero hing-loss regime, we can approximate the gradient in the $\theta$-direction numerically, e.g. with simultaneous perturbation stochastic approximation (SPSA). A detailed instruction of the algorithm is provided in Algorithm 1.

## 3.3 Kernel alignment for non-stationary data

The formulation in (6) provides a clear chronological structure to how the weights $\mathbf{w}$ are defined. This is useful in online machine learning, where the structure of the data is time dependent and training data becomes outdated after a certain period. In standard machine learning models, it is usually easy to add new training data and continue training an existing model. However, it is often unclear how to remove or unlearn the impact of outdated training data. For that reason, models are usually completely retrained after a some time in order to adjust to the change in the data structure. This is not the case for classification with the Pegasos algorithm: If we realise that the structure in our training data has changed significantly, we can simply discard past iterations (and the corresponding data points) and only sum over the relevant time period by setting $\alpha_t = 0$ for $t$ lying far in the past.

# 4 Experimental demonstration

In this section we demonstrate that the proposed algorithm succeeds in solving the kernel alignment task proposed in [8]. For a detailed derivation and motivation of the specific classification problem, we thus refer to [8]. Crucially, given a device with $n$ qubits and connectivity graph $G(V, E)$, a data set $\mathcal{D}(n, G)$ with $n$ features $\mathbf{x} \in \mathbb{R}^n$ and binary labels $y \in \{-1, 1\}$ is created. We classify

---
**Algorithm 1** Kernel Alignment with Pegasos
---
1: **Inputs:**
2: training data $T = \{\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_M\}$
3: labels $L = \{y_1, y_2, ..., y_M\}$
4: regularization parameter $\lambda \in \mathbb{R}^+$
5: pseudo-Kernel of the form $k_{\theta_1, \theta_2}(\mathbf{x}, \mathbf{y}) = |\langle \psi_{\theta_2}(\mathbf{x}) | \psi_{\theta_1}(\mathbf{y}) \rangle|^2$
6: initial kernel parameters $\theta_1 \in \mathbb{R}^d$
7: number of steps $\tau \in \mathbb{N}$
8: number of initialization steps $\tau_{in} < \tau$
9:
10: **for** $t = 1, 2, ..., \tau$ **do**
11:     Choose $i_t \in \{1, ..., M\}$ uniformly at random.
12:     **if** $t = 1$ **then**
13:         $\alpha_t \leftarrow 1$
14:         $\theta_{t+1} \leftarrow \theta_t$
15:     **else**
16:         **if** $y_{i_t} \frac{1}{\lambda t} \sum_{s=1}^{t-1} \alpha_s y_{i_s} k_{\theta_s, \theta_t}(\mathbf{x}_{i_s}, \mathbf{x}_{i_t}) < 1$ **then**
17:             $\alpha_t \leftarrow 1$
18:             **if** $t > \tau_{in}$ **then**
19:                 $f(\theta) \leftarrow -y_{i_t} \frac{1}{\lambda t} \sum_{s=1}^{t-1} \alpha_s y_{i_s} k_{\theta_s, \theta}(\mathbf{x}_{i_s}, \mathbf{x}_{i_t})$
20:                 $\theta_{t+1} \leftarrow$ Peform minimization step on $f(\theta)$ around $\theta_t$, e.g. using SPSA.
21:             **else**
22:                 $\theta_{t+1} \leftarrow \theta_t$
23:             **end if**
24:         **else**
25:             $\alpha_t \leftarrow 0$
26:             $\theta_{t+1} \leftarrow \theta_t$
27:         **end if**
28:     **end if**
29: **end for**
---

$\mathcal{D}(n, G)$ using the feature map $|\psi_\theta(\mathbf{x})\rangle = U(\mathbf{x})V_\theta |0\rangle$, where $U(\mathbf{x}) = \bigotimes_{k=1}^{n} R_X(x_{2k-1})R_Z(x_{2k})$ and $V_\theta = \left(\prod_{(kl) \in E} \text{CNOT}(k,l)\right)\left(\bigotimes_{k=1}^{n} R_Y(\theta)\right)$. For the ideal parameter $\theta = \theta_{\text{opt}} = \pi/2$ the data set becomes linearly separable and the QSVM can thus reach 100% accuracy, however if $\theta$ is chosen badly, the data cannot be separated, resulting in a lower training accuracy.

## 4.1 Stationary data

In this fist experiment, we train Pegasos on $\mathcal{D}(n, G)$. We initialize $\theta = 0$ and perform the simultaneous optimization of $\theta$ and $\boldsymbol{\alpha}$ as described in Section 3. Figure 3 shows that we are able to reach 100% accuracy.



Figure 3: **Quantum kernel alignment with Pegasos** using statevector simulation on 10 qubits using Qiskit. We train Pegasos on the data provided in [8] to classify the training data correctly while aligning the kernel. The top plot shows how the trainable parameter is trained using SPSA with a learning rate of $\mu = 0.1$. The bottom plot shows if the sampled data points are classified correctly for every iteration (single shot accuracy) and averaged over the last 50 iterations, as well as the test accuracy on 25 unseen data points.

## 4.2 Non-stationary data

In a second experiment, we test the algorithm on non-stationary data. For this, we slightly change the setup from the above experiment by changing the structure of the data set as a function of time. Instead of fixing the parameter $\theta_{\text{opt}} = \pi/2$ used to generate the data, we allow $\theta_{\text{opt}}$ to vary every time a new tuple $(\mathbf{x_i}, y_i)$ is sampled like

$$\theta_{\text{opt}}(t) = \sin(2\pi t/T),$$

where $t$ indicates the current iteration count and $T = 1000$. As described in Section 3.3, we only keep a window of $\tau$ iteration steps in memory in order to forget the impact of outdated training data. The goal is to correctly classify new samples while keeping the structure parameter $\theta$ close to $\theta_{\text{opt}}(t)$. Figure 4 shows a hardware experiment demonstrating that Pegasos is able to track

the change in the data structure by continuously tracking $\theta_{\text{opt}}$ and classifying the newly sampled points correctly with an accuracy over 90%.
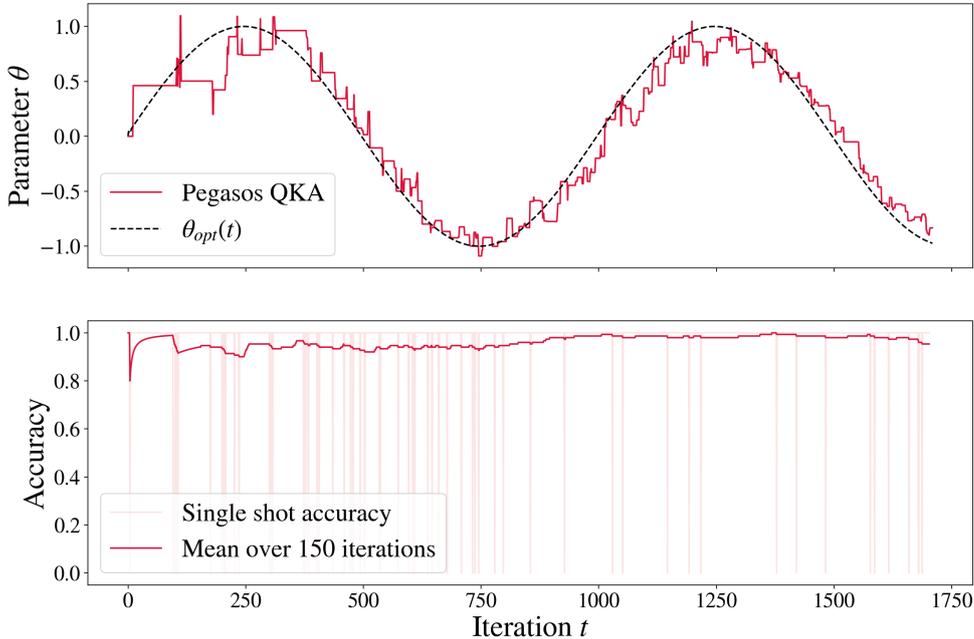


Figure 4: **Quantum kernel alignment with Pegasos for non-stationary data**. In every iteration we modify $\theta_{\text{opt}}$ defining the structure in the data and draw a random data point from the corresponding training set as defined in [8]. We let $\theta_{\text{opt}}$ change continuously according to a sine wave and see how PegasosQKA is able to track this parameter, while keeping up the accuracy of the classification. Experiment on 7-qubit device `ibm_nairobi` with SPSA learning rate $\mu = 0.1$ and training window $\tau = 100$.

## 5  Conclusion

Quantum kernel alignment is a promising technique for fine-tuning quantum kernels and the Pegasos algorithm allows for an elegant implementation which simultaneously performs kernel alignment as the primal QSVM problem is being solved. The speedup obtained over embedding one optimization loop within another and extends the applicability of kernel alignment to larger QSVM problems. In addition, Pegasos naturally supports online machine learning and non-stationary data as new data can easily be added to the training set and old data can be easily unlearned at any point. Our experiments demonstrate that Pegasos with simultaneous kernel alignment works well in practice and we expect this technique to be of value as QSVM models are scaled to larger and more realistic problems.

**Code availability**   The code for our experiments presented in Section 4 has been written using Qiskit [22] and is available at [23].

# References

[1] J. Biamonte, P. Wittek, N. Pancotti, P. Rebentrost, N. Wiebe, and S. Lloyd. Quantum machine learning. *Nature*, 549(7671):195–202, 2017. `DOI: 10.1038/nature23474`.

[2] V. Havlíček, A. D. Córcoles, K. Temme, A. W. Harrow, A. Kandala, J. M. Chow, and J. M. Gambetta. Supervised learning with quantum-enhanced feature spaces. *Nature*, 567(7747):209–212, 2019. `DOI: 10.1038/s41586-019-0980-2`.

[3] A. Abbas, D. Sutter, C. Zoufal, A. Lucchi, A. Figalli, and S. Woerner. The power of quantum neural networks. *Nature Computational Science*, 1(June), 2020. `DOI: 10.1038/s43588-021-00084-1`.

[4] Y. Liu, S. Arunachalam, and K. Temme. A rigorous and robust quantum speed-up in supervised machine learning. *Nature Physics*, 17(9):1013–1017, 2021. `DOI: 10.1038/s41567-021-01287-z`.

[5] J. Kübler, S. Buchholz, and B. Schölkopf. The inductive bias of quantum kernels. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 12661–12673. Curran Associates, Inc., 2021. Available online: `https://proceedings.neurips.cc/paper_files/paper/2021/file/69adc1e107f7f7d035d7baf04342e1ca-Paper.pdf`.

[6] N. Cristianini, J. Shawe-Taylor, A. Elisseeff, and J. Kandola. On kernel-target alignment. In T. Dietterich, S. Becker, and Z. Ghahramani, editors, *Advances in Neural Information Processing Systems*, volume 14. MIT Press, 2001. Available online: `https://proceedings.neurips.cc/paper_files/paper/2001/file/1f71e393b3809197ed66df836fe833e5-Paper.pdf`.

[7] C. Cortes, M. Mohri, and A. Rostamizadeh. Algorithms for learning kernels based on centered alignment. *J. Mach. Learn. Res.*, 13(1):795–828, 2012. `DOI: 10.5555/2503308.2188413`.

[8] J. R. Glick, T. P. Gujarati, A. D. Corcoles, Y. Kim, A. Kandala, J. M. Gambetta, and K. Temme. Covariant quantum kernels for data with group structure. pages 1–9, 2021. `arXiv: 2105.03406`.

[9] S. Shalev-Shwartz and N. Srebro. SVM optimization: Inverse dependence on training set size. *Proceedings of the 25th International Conference on Machine Learning*, pages 928–935, 2008.

[10] G. Gentinetta, A. Thomsen, D. Sutter, and S. Woerner. The complexity of quantum support vector machines, 2022. `DOI: 10.48550/ARXIV.2203.00031`.

[11] B. E. Boser, I. M. Guyon, and V. N. Vapnik. A Training Algorithm for Optimal Margin Classifiers. In *Proceedings of the Fifth Annual Workshop on Computational Learning Theory*, COLT '92, pages 144–152, New York, NY, USA, 1992. Association for Computing Machinery. `DOI: 10.1145/130385.130401`.

[12] V. N. Vapnik. *The Nature of Statistical Learning Theory*. Springer Science+Business Media, LLC, 2000.

[13] C. Cortes and V. Vapnik. Support-vector networks. In *Machine Learning*, pages 273–297, 1995.

[14] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.

[15] S. Shalev-Shwartz, Y. Singer, N. Srebro, and A. Cotter. Pegasos: Primal estimated sub-gradient solver for SVM. *Mathematical Programming*, 127(1):3–30, 2011. `DOI: 10.1007/s10107-010-0420-4`.

[16] S. Thanasilp, S. Wang, M. Cerezo, and Z. Holmes. Exponential concentration and untrainability in quantum kernel methods, 2022. `DOI: 10.48550/ARXIV.2208.11060`.

[17] M. Cerezo, A. Sone, T. Volkoff, L. Cincio, and P. J. Coles. Cost function dependent barren plateaus in shallow parametrized quantum circuits. *Nature Communications*, 12(1), 2021. `DOI: 10.1038/s41467-021-21728-w`.

[18] S. Wang, E. Fontana, M. Cerezo, K. Sharma, A. Sone, L. Cincio, and P. J. Coles. Noise-induced barren plateaus in variational quantum algorithms. *Nature Communications*, 12(1):6961, 2021. `DOI: 10.1038/s41467-021-27045-6`.

[19] C. Ortiz Marrero, M. Kieferová, and N. Wiebe. Entanglement-induced Barren plateaus. *PRX Quantum*, 2:040316, 2021. `DOI: 10.1103/PRXQuantum.2.040316`.

[20] Y. Liu, S. Arunachalam, and K. Temme. A rigorous and robust quantum speed-up in supervised machine learning. *Nature Physics*, 2021. `DOI: 10.1038/s41567-021-01287-z`.

[21] J. J. Meyer, M. Mularski, E. Gil-Fuster, A. A. Mele, F. Arzani, A. Wilms, and J. Eisert. Exploiting symmetry in variational quantum machine learning, 2022. `DOI: 10.48550/ARXIV.2205.06217`.

[22] M. D. S. Anis et al. Qiskit: An Open-source Framework for Quantum Computing, 2023. `DOI: 10.5281/zenodo.2573505`.

[23] G. Gentinetta, D. Sutter, C. Zoufal, and S. Woerner. Code for manuscript "Quantum Kernel Alignment with Stochastic Gradient Descent", 2023. `DOI: 10.5281/zenodo.7804477`.