

# Quantum Perceptron Models

Nathan Wiebe, Ashish Kapoor, Krysta M. Svore  
*Microsoft Research, One Microsoft Way, Redmond WA 98052*

We demonstrate how quantum computation can provide non-trivial improvements in the computational and statistical complexity of the perceptron model. We develop two quantum algorithms for perceptron learning. The first algorithm exploits quantum information processing to determine a separating hyperplane using a number of steps sublinear in the number of data points  $N$ , namely  $O(\sqrt{N})$ . The second algorithm illustrates how the classical mistake bound of  $O(\frac{1}{\gamma^2})$  can be further improved to  $O(\frac{1}{\sqrt{\gamma}})$  through quantum means, where  $\gamma$  denotes the margin. Such improvements are achieved through the application of quantum amplitude amplification to the version space interpretation of the perceptron model.

## I. INTRODUCTION

Quantum computation is an emerging technology that utilizes quantum effects to achieve significant, and in some cases exponential, speed-ups of algorithms over their classical counterparts. The growing importance of machine learning has in recent years led to a host of studies that investigate the promise of quantum computers for machine learning [1, 2, 12, 13, 17, 21–23].

While a number of important quantum speedups have been found, the majority of these speedups are due to replacing a classical subroutine with an equivalent albeit faster quantum algorithm. The true potential of quantum algorithms may therefore remain underexploited since quantum algorithms have been constrained to follow the same methodology behind traditional machine learning methods [2, 7, 22]. Here we consider an alternate approach: we devise a new machine learning algorithm that is tailored to the speedups that quantum computers can provide.

We illustrate our approach by focusing on perceptron training [18]. The perceptron is a fundamental building block for various machine learning models including neural networks and support vector machines [20]. Unlike many other machine learning algorithms, tight bounds are known for the computational and statistical complexity of traditional perceptron training. Consequently, we are able to rigorously show different performance improvements that stem from either using quantum computers to improve traditional perceptron training or from devising a new form of perceptron training that aligns with the capabilities of quantum computers.

We provide two quantum approaches to perceptron training. The first approach focuses on the computational aspect of the problem and the proposed method quadratically reduces the scaling of the complexity of training with respect to the number of training vectors. The second algorithm focuses on statistical efficiency. In particular, we use the mistake bounds for traditional perceptron training methods and ask if quantum computation lends any advantages. To this end, we propose an algorithm that quadratically improves the scaling of the training algorithm with respect to the margin between the classes in the training data. The latter algorithm combines quantum amplitude estimation in the version space interpretation of the perceptron learning problem. Our approaches showcase the trade-offs that one can consider in developing quantum algorithms, and the ultimate advantages of performing learning tasks on a quantum computer.

The rest of the paper is organized as follows: we first cover the background on perceptrons, version space and Grover search. We then present our two quantum algorithms and provide analysis of their computational and statistical efficiency before concluding.

## II. BACKGROUND

### A. Perceptrons and Version Space

Given a set of  $N$  separable training examples  $\{\phi_1, \dots, \phi_N\} \in \mathbb{R}^D$  with corresponding labels  $\{y_1, \dots, y_N\}$ ,  $y_i \in \{+1, -1\}$ , the goal of perceptron learning is to recover a hyperplane  $w$  that perfectly classifies the training set [18]. Formally, we want  $w$  such that  $y_i \cdot w^T \phi_i > 0$  for all  $i$ . There are various simple online algorithms that start with a random initialization of the hyperplane and make updates as they encounter more and more data [8, 11, 18, 19]; however, the rule that we consider for online perceptron training is, upon misclassifying a vector  $(\phi, y)$ ,  $w \leftarrow w + y\phi$ .

A remarkable feature of the perceptron model is that upper bounds exist for the number of updates that need to be made during this training procedure. In particular, if the training data is composed of unit vectors,  $\phi_i \in \mathbb{R}^D$ , that are separated by a margin of  $\gamma$  then there are perceptron training algorithms that make at most  $O(\frac{1}{\gamma^2})$  mistakes [16],

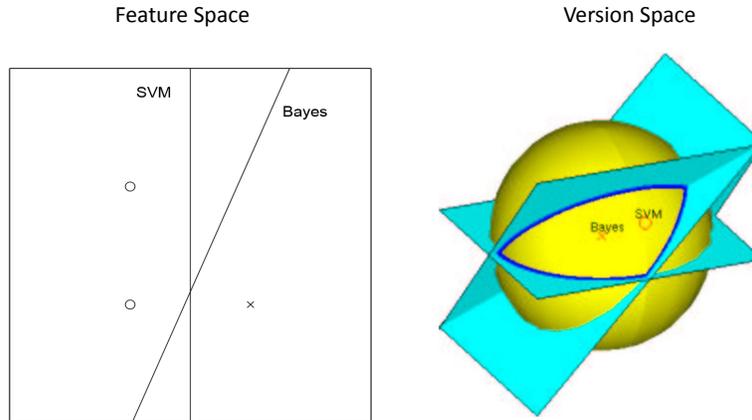


FIG. 1: Version space and feature space views of classification. This figure is from [14].

independent of the dimension of the training vectors. Similar bounds also exist when the data is not separated [6] and also for other generalizations of perceptron training [8, 11, 19]. Note that in the worst case, the algorithm will need to look at all points in the training set at least once, consequently the computation complexity will be  $O(N)$ .

Our goal is to explore if the quantum procedures can provide improvements both in terms of computational complexity (that is better than  $O(N)$ ) and statistical efficiency (improve upon  $O(\frac{1}{\gamma^2})$ ). Instead of solely applying quantum constructs to the feature space, we also consider the version space interpretation of perceptrons which leads to the improved scaling with  $\gamma$ .

Formally, version space is defined as the set of all possible hyperplanes that perfectly separate the data:  $VS := \{w | y_i \cdot w^T \phi_i > 0 \text{ for all } i\}$ . Given a training datum, the traditional representation is to depict data as points in the feature space and use hyperplanes to depict the classifiers. However, there exists a dual representation where the hyperplanes are depicted as points and the data points are represented as hyperplanes that induce constraints on the feasible set of classifiers. Figure 1, which is borrowed from [14], illustrates the version space interpretation of perceptrons. Given three labeled data points in a 2D space, the dual space illustrates the set of normalized hyperplanes as a yellow ball with unit radius. The third dimension corresponds to the weights that multiply the two dimensions of the input data and the bias term. The planes represent the constraints imposed by observing the labeled data as every labeled data renders one-half of the space infeasible. The version space is then the intersection of all the half-spaces that are valid. Naturally, classifiers including SVMs [20] and Bayes point machines [10] lie in the version space.

We note that there are quantum constructs such as Grover search and amplitude amplification which provide non-trivial speedups for the search task. This is the main reason why we resort to the version space interpretation. We can use this formalism to simply pose the problem of determining the separating hyperplane as a search problem in the dual space. For example given a set of candidate hyperplanes, our problem reduces to searching amongst the sample set for the classifier that will successfully classify the entire set. Therefore training the perceptron is equivalent to finding *any* feasible point in the version space. We describe these quantum constructs in detail below.

## B. Grover's Search

Both quantum approaches introduced in this work and their corresponding speed-ups stem from a quantum subroutine called Grover's search [4, 9], which is a special case of a more general method referred to as amplitude amplification [5]. Rather than sampling from a probability distribution until a given marked element is found, the Grover search algorithm draws only one sample and then uses quantum operations to modify the distribution from which it sampled. The probability distribution is rotated, or more accurately the quantum state that yields the distribution is rotated, into one whose probability is sharply concentrated on the marked element. Once a sharply peaked distribution is identified, the marked item can be found using just one sample. In general, if the probability of finding such an element is known to be  $a$  then amplitude amplification requires  $O(\sqrt{1/a})$  operations to find the marked item with certainty.

While Grover's search is a *quantum subroutine*, it can in fact be understood using only geometric arguments. The only notions from quantum mechanics used are those of the quantum state vector and that of Born's rule

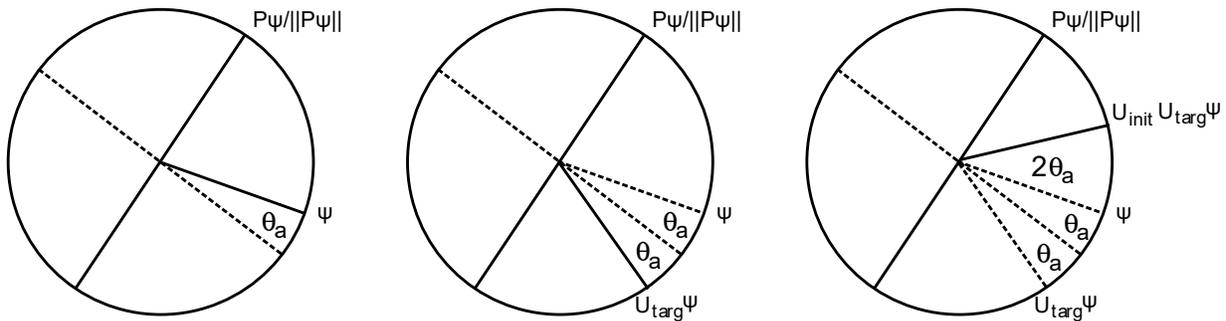


FIG. 2: A geometric description of the action of  $U_{\text{grover}}$  on an initial state vector  $\psi$ .

(measurement). A quantum state vector is a complex unit vector whose components have magnitudes that are equal to the square-roots of the probabilities. In particular, if  $v$  is a quantum state vector and  $p$  is the corresponding probability distribution then

$$p = v^\dagger \circ v, \quad (1)$$

where the unit column vector  $v$  is called the quantum state vector which sits in the vector space  $\mathbb{C}^n$ ,  $\circ$  is the Hadamard (pointwise) product and  $\dagger$  is the complex conjugate transpose. A quantum state can be measured such that if we have a quantum state vector  $v$  and a basis vector  $w$  then the probability of measuring  $v = w$  is  $|\langle v, w \rangle|^2$ , where  $\langle \cdot, \cdot \rangle$  denotes the inner product. One of the main differences between quantum and classical distributions is that the probability distribution resulting from measurement depends strongly on the basis in which the vector is measured. This basis dependence of measurement is the root of many of the differences between quantum and classical probability theory and also gives rise to many celebrated results in the foundations of quantum mechanics such as Bell's theorem [3].

At first glance, introducing the quantum state vector  $v$  may not seem to provide any advantages over working with  $p$  for the purposes of sampling. More careful consideration reveals that the fact that  $v$  is complex valued allows transformations on  $v$  to be performed that cannot be performed on  $p$ . In particular, we can reflect the quantum state vector about any axis, whereas we cannot do the same to  $p$  without violating its positivity. Grover's search, in fact, is a cunning way to perform a series of reflections on  $v$  to bias  $p$  towards the marked state we wish to find. While such reflections may not make sense from a classical perspective, quantum computers can be used to realize them efficiently.

The key feature of a quantum computer is that it permits any unitary transformation to be performed on the unit vector  $v$ , within arbitrarily small approximation error. We define the initial quantum state vector to be  $\psi$  and define  $P$  to be a projection matrix onto a set of configurations that we want to find. In particular, if we define  $\nu_{\text{good}}$  to be the set of all items that we want the quantum algorithm to find then

$$P\psi = \begin{cases} \psi, & \text{if } \psi \in \nu_{\text{good}} \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

Here being able to apply  $P$  does not imply that  $\nu_{\text{good}}$  is known. Instead, it implies that a subroutine that checks to see if  $\psi \in \nu_{\text{good}}$  exists. The fact that  $P$  is implemented by a linear transformation of the state vector also allows it to be simultaneously applied to exponentially many  $v$  via  $P\psi = \frac{1}{\|a\|} \sum_{j=1}^N a_j v = \frac{1}{\|a\|} \sum_{j=1}^N a_j P v$ . These two features allow a single application of  $\mathbb{1} - 2P$  to be efficiently applied, assuming membership in  $\nu_{\text{good}}$  can be efficiently tested, even though  $\psi$  is a sum of exponentially many basis vectors.

In order to perform the search algorithm we need to implement two unitary operations:

$$U_{\text{init}} = 2\psi\psi^\dagger - \mathbb{1}, U_{\text{targ}} = \mathbb{1} - 2P. \quad (3)$$

The operators  $U_{\text{init}}$  and  $U_{\text{targ}}$  can be interpreted geometrically as reflections within a two-dimensional space spanned by the vectors  $\psi$  and  $P\psi$ . If we assume that  $P\psi \neq 0$  and  $P\psi \neq \psi$  then these two reflection operations can be used to rotate  $\psi$  in the space  $\text{span}(\psi, P\psi)$ . Specifically this rotation is  $U_{\text{grover}} = U_{\text{init}}U_{\text{targ}}$ . Its action is illustrated in Figure 2. If the angle between the vector  $\psi$  and  $P\psi/\|P\psi\|$  is  $\pi/2 - \theta_a$ , where  $\theta_a := \sin^{-1}(|\langle \psi, P\psi/\|P\psi\| \rangle|)$ . It then follows from elementary geometry and the rule for computing the probability distribution from a quantum state (known as Born's rule) that after  $j$  iterations of Grover's algorithm the probability of measuring a desirable outcome is

$$p(v \in \nu_{\text{good}} | j) = \sin^2((2j + 1)\theta_a). \quad (4)$$

It is then easy to see that if  $\theta_a \ll 1$  and a probability of success greater than  $1/4$  is desired then  $j \in O(1/\sqrt{\theta_a})$  suffices to find a marked outcome. This is quadratically faster than is possible from statistical sampling which requires  $O(1/\theta_a)$  samples on average.

As an example, if the initial success probability is  $1/4$  then  $\theta_a = \sin^{-1}(1/2) = \pi/6$ . Therefore if we take  $j = 1$  then  $p(v \in \nu_{\text{good}}|j) = 1$ . As a result a desirable outcome can be found after only 3 quantum operations whereas 4 samples from the initial distribution would be needed on average to find a marked outcome if quantum methods were not used.

If on the other hand, the success probability were  $1/2$  then  $\theta_a = \pi/4$  and  $\sin^2((2j+1)\pi/4) = 1/2$  for all  $j$ . This problem can be easily addressed by doing something that would not make any sense classically: purposefully lowering the success probability to  $1/4$  by requiring a new event  $w = (v, u)$  where we define  $w$  to be a good state if the independent variables  $v$  is good and  $u \sim \text{Bern}(1/2)$  is 0. The independence assumption means that the probability that both conditions are satisfied is  $1/4$  and hence a good  $v$  can be found with certainty by applying amplitude amplification on  $w$ . More generally, if  $\theta_a$  is known then this trick can be applied to make  $\theta_a \mapsto \pi/(2[2j+1])$  (for positive integer  $j$ ) which makes the search procedure deterministic.

On the other hand, if  $\theta_a$  is not known then it isn't clear how  $j$  should be chosen to make the success probability greater than  $1/4$ . Fortunately, methods are known to deal with such issues [4, 5]. The simplest one exploits the fact that the average of  $p$  over a range of  $j = 0, \dots, M-1$  can be easily computed:

$$\begin{aligned} p(v \in \nu_{\text{good}}; M) &= \frac{1}{M} \sum_{j=0}^{M-1} \sin^2((2j+1)\theta_a) \\ &= \frac{1}{2} \left( 1 - \frac{\sin(4M\theta_a)}{2M \sin(2\theta_a)} \right). \end{aligned} \quad (5)$$

If  $M \geq M_0 := \frac{1}{\sin(2\theta_a)}$  then it is straight forward to see that

$$\frac{1}{2} \left( 1 - \frac{\sin(4M\theta_a)}{2M \sin(2\theta_a)} \right) \geq \frac{1}{2} \left( 1 - \frac{1}{2M \sin(2\theta_a)} \right) \geq \frac{1}{4}. \quad (6)$$

The average probability is then guaranteed to be at least  $1/4$  if  $j$  is chosen to be drawn uniformly from  $\{0, \dots, M-1\}$  if  $M \geq M_0$ . If a lower bound on  $\theta_a$  is known a good sample can be drawn, then an appropriate value of  $M$  can be computed.

If no lower bound on  $\theta_a$  is known then a marked element can nonetheless be found with high probability by exponential searching. Exponential searching involves, for step  $i$  taking  $M = c^i$  for some  $c \in (1, 2)$ . After a logarithmic number of applications of amplitude amplification it will attain  $M \geq M_0$  with high probability. After which the average success probability is known to be bounded below by  $1/4$  and the algorithm will succeed with high probability in a constant number of attempts. Thus the quadratic speedup holds even if the success probability is not known a priori.

### III. ONLINE QUANTUM PERCEPTRON

Now that we have discussed Grover's search we turn our attention to applying it to speed up online perceptron training. In order to do so, we first need to define the quantum model that we wish to use as our quantum analogue of perceptron training. While there are many ways of defining such a model but the following approach is perhaps the most direct. Although the traditional feature space perceptron training algorithm is online [16], meaning that the training examples are provided one at a time to it in a streaming fashion, we deviate from this model slightly by instead requiring that the algorithm be fed training examples that are, in effect, sampled uniformly from the training set. This is a slightly weaker model, as it allows for the possibility that some training examples will be drawn multiple times. However, the ability to draw quantum states that are in a uniform superposition over all vectors in the training set enables quantum computing to provide advantages over both classical methods that use either access model.

We assume without loss of generality that the training set consists of  $N$  unit vectors,  $\phi_1, \dots, \phi_N$ . If we then define  $\Phi_1, \dots, \Phi_N$  to be the basis vectors whose indices each coincide with a  $(B+1)$ -bit representation of the corresponding  $(\phi_j, y_j)$  where  $y_j \in \{-1, 1\}$  is the class assigned to  $\phi_j$  and let  $\Phi_0$  be a fixed unit vector that is chosen to represent a blank memory register.

We introduce the vectors  $\Phi_j$  to make it clear that the quantum vectors states used to represent training vectors do not live in the same vector space as the training vectors themselves. We choose the quantum state vectors here to occupy a larger space than the training vectors because the Heisenberg uncertainty principle makes it much more difficult for a quantum computer to compute the class that the perceptron assigns to a training vector in such cases.

For example, the training vector  $(\phi_j, y_j) \equiv ([0, 0, 1, 0]^T, 1)$  can be encoded as an unsigned integer  $00101 \equiv 5$ , which in turn can be represented by the unit vector  $\Phi = [0, 0, 0, 0, 1]^T$ . More generally, if  $\phi_j \in \mathbb{R}^D$  were a vector of

---

**Algorithm 1** Online quantum perceptron training algorithm
 

---

```

for  $k = 1, \dots, \lceil \log_{3/4} \gamma^2 \epsilon \rceil$  do
for  $j = 1 : \lceil \log_c(1/\sin(2\sin^{-1}(1/\sqrt{N}))) \rceil$  do
  Draw  $m$  uniformly from  $\{0, \dots, \lceil c^j \rceil\}$ .
  Prepare quantum state  $\Psi$ .
   $\Psi \leftarrow ((2\Psi\Psi^\dagger - \mathbb{1})F_w)^m \Psi$ .
  Measure  $\Psi$ , assume outcome is  $u_q$ .
   $(\phi, y) \leftarrow U^c(q)$ .
  if  $f_w(\phi, y) = 1$  then
    Return  $w' \leftarrow w + y\phi$ 
  end if
end for
end for
  Return  $w'$ 

```

---

floating point numbers then a similar vector could be constructed by concatenating the binary representations of the  $D$  floating point numbers that comprise it with  $(y_j + 1)/2$  and express the bit string as an unsigned integer,  $Q$ . The integer can then be expressed as a unit vector  $\Phi : [\Phi]_q = \delta_{q,Q}$ . While encoding the training data as an exponentially long vector is inefficient in a classical computer, it is not in a quantum computer because of the quantum computer's innate ability to store and manipulate exponentially large quantum state vectors.

Any machine learning algorithm, be it quantum or classical, needs to have a mechanism to access the training data. We assume that the data is accessed via an oracle that not only accesses the training data but also determines whether the data is misclassified. To clarify, let  $\{u_j : j = 1 : N\}$  be an orthonormal basis of quantum state vectors that serve as addresses for the training vectors in the database. Given an input address for the training datum, the unitary operations  $U$  and  $U^\dagger$  allow the quantum computer to access the corresponding vector. Specifically, for all  $j$

$$\begin{aligned} U[u_j \otimes \Phi_0] &= u_j \otimes \Phi_j \\ U^\dagger[u_j \otimes \Phi_j] &= u_j \otimes \Phi_0. \end{aligned} \quad (7)$$

Given an input address vector  $u_j$ , the former corresponds to a database access and the latter inverts the database access.

Note that because  $U$  and  $U^\dagger$  are linear operators we have that  $U \sum_{j=1}^N u_j \otimes \Phi_0 = \sum_j u_j \otimes \Phi_j$ . A quantum computer can therefore access each training vector simultaneously using a single operation, while only requiring enough memory to store one of the  $\Phi_j$ . The resultant vector is often called in the physics literature a *quantum superposition of states* and this feature of linear transformations is referred to as quantum parallelism within quantum computing.

The next ingredient that we need is a method to test if the perceptron correctly assigns a training vector addressed by a particular  $u_j$ . This process can be pictured as being performed by a unitary transformation that flips the sign of any basis-vector that is misclassified. By linearity, a single application of this process flips the sign of any component of the quantum state vector that coincides with a misclassified training vector. It therefore is no more expensive than testing if a given training vector is misclassified in a classical setting. We denote the operator, which depends on the perceptron weights  $w$ ,  $F_w$  and require that

$$F_w[u_j \otimes \Phi_0] = (-1)^{f_w(\phi_j, y_j)} [u_j \otimes \Phi_0], \quad (8)$$

where  $f_w(\phi_j)$  is a Boolean function that is 1 if and only if the perceptron with weights  $w$  misclassifies training vector  $\phi_j$ . Since the classification step involves computing the dot-products of finite size vectors, this process is efficient given that the  $\Phi_j$  are efficiently computable.

We apply  $F_w$  in the following way. Let  $\mathcal{F}_w$  be a unitary operation such that

$$\mathcal{F}_w \Phi_j = (-1)^{f_w(\phi_j, y_j)} \Phi_j. \quad (9)$$

$\mathcal{F}_w$  is easy to implement in the quantum computer using a multiply controlled phase gate and a quantum implementation of the perceptron classification algorithm,  $f_w$ . We can then write

$$F_w = U^\dagger(\mathbb{1} \otimes \mathcal{F}_w)U. \quad (10)$$

Classifying the data based on the phases (the minus signs) output by  $F_w$  naturally leads to a very memory efficient training algorithm because only one training vector is ever stored in memory during the implementation of  $F_w$  given in (10). We can then use  $F_w$  to perform Grover's search algorithm, by taking  $U_{\text{targ}} = F_w$  and  $U_{\text{init}} = 2\psi\psi^\dagger - \mathbb{1}$

with  $\psi = \Psi := \frac{1}{\sqrt{N}} \sum_{j=1}^N u_j$ , to seek out training vectors that the current perceptron model misclassifies. This leads to a quadratic reduction in the number of times that the training vectors need to be accessed by  $F_w$  or its classical analogue.

In the classical setting, the natural object to query is slightly different. The oracle that is usually assumed in online algorithms takes the form  $U^c : \mathbb{Z} \mapsto \mathbb{C}^D$  where

$$U^c(j) = \phi_j. \quad (11)$$

We will assume that a similar function exists in both the classical and the quantum settings for simplicity. In both cases, we will consider the cost of a query to  $U^c$  to be proportional to the cost of a query to  $F_w$ .

We use these operations in Algorithm 1 to implement a quantum search for training vectors that the perceptron misclassifies. This leads to a quadratic speedup relative to classical methods as shown in the following theorem.

**Theorem 1.** *Given a training set that consists of unit vectors  $\Phi_1, \dots, \Phi_N$  that are separated by a margin of  $\gamma$  in feature space, the number of applications of  $F_w$  needed to infer a perceptron model,  $w$ , such that  $P(\exists j : f_w(\phi_j) = 1) \leq \epsilon$  using a quantum computer is  $N_{\text{quant}}$  where*

$$\Omega(\sqrt{N}) \ni N_{\text{quant}} \in O\left(\frac{\sqrt{N}}{\gamma^2} \log\left[\frac{1}{\epsilon\gamma^2}\right]\right),$$

whereas the number of queries to  $f_w$  needed in the classical setting,  $N_{\text{class}}$ , where the training vectors are found by sampling uniformly from the training data is bounded by

$$\Omega(N) \ni N_{\text{class}} \in O\left(\frac{N}{\gamma^2} \log\left[\frac{1}{\epsilon\gamma^2}\right]\right).$$

We assume in Theorem 1 that the training data in the classical case is accessed in a manner that is analogous to the sampling procedure used in the quantum setting. If instead the training data is supplied by a stream (as in the standard online model) then the upper bound changes to  $N_{\text{class}} \in O(N/\gamma^2)$  because all  $N$  training vectors can be deterministically checked to see if they are correctly classified by the perceptron. A quantum advantage is therefore obtained if  $N \gg \log^2(1/\epsilon\gamma^2)$ .

In order to prove Theorem 1 we need to have two technical lemmas (proven in the appendix). The first bounds the complexity of the classical analogue to our training method:

**Lemma 1.** *Given only the ability to sample uniformly from the training vectors, the number of queries to  $f_w$  needed to find a training vector that the current perceptron model fails to classify correctly, or conclude that no such example exists, with probability  $1 - \epsilon\gamma^2$  is at most  $O(N \log(1/\epsilon\gamma^2))$ .*

The second proves the correctness of Algorithm 1 and bounds the complexity of the algorithm:

**Lemma 2.** *Assuming that the training vectors  $\{\phi_1, \dots, \phi_N\}$  are unit vectors and that they are drawn from two classes separated by a margin of  $\gamma$  in feature space, Algorithm 2 will either update the perceptron weights, or conclude that the current model provides a separating hyperplane between the two classes, using a number of queries to  $F_w$  that is bounded above by  $O(\sqrt{N} \log(1/\epsilon\gamma^2))$  with probability of failure at most  $\epsilon\gamma^2$ .*

After stating these results, we can now provide the proof of Theorem 1.

*Proof of Theorem 1.* The upper bounds follow as direct consequences of Lemma 2 and Lemma 1. Novikoff's theorem [6, 16] states that the algorithms described in both lemmas must be applied at most  $1/\gamma^2$  times before finding the result. However, either the classical or the quantum algorithm may fail to find a misclassified vector at each of the  $O(1/\gamma^2)$  steps. The union bound states that the probability that this happens is at most the sum of the respective probabilities in each step. These probabilities are constrained to be  $\gamma^2\epsilon$ , which means that the total probability of failing to correctly find a mistake is at most  $\epsilon$  if both algorithms are repeated  $1/\gamma^2$  times (which is the worst case number of times that they need to be repeated).

The lower bound on the quantum query complexity follows from contradiction. Assume that there exists an algorithm that can train an arbitrary perceptron using  $o(\sqrt{N})$  query operations. Now we want to show that unstructured search with one marked element can be expressed as a perceptron training algorithm. Let  $w$  be a known set of perceptron weights and assume that the perceptron only misclassifies one vector  $\phi_1$ . Thus if perceptron training succeeds then  $w$  the value of  $\phi_1$  can be extracted from the updated weights. This training problem is therefore equivalent to searching for a misclassified vector. Now let  $\phi_j = [1 \oplus F(j), F(j)]^T \otimes \chi_j$  where  $\chi_j$  is a unit vector that represents the bit string  $j$  and  $F(j)$  is a Boolean function. Assume that  $F(0) = 1$  and  $F(j) = 0$  if  $j \neq 0$ , which is without

loss of generality equivalent to Grover's problem [4, 9]. Now assume that  $\phi_j$  is assigned to class  $2F(j) - 1$  and take  $w = [1/\sqrt{2}, 1/\sqrt{2}]^T \otimes \frac{1}{\sqrt{N}} \sum_j \chi_j$ . This perceptron therefore misclassifies  $\phi_0$  and no other vector in the training set. Thus updating the weights yields  $\phi_j$ , which in turn yields the value of  $j$  such that  $F(j) = 1$ , and therefore Grover's search reduces to perceptron training.

Since Grover's search reduces to perceptron training in the case of one marked item the lower bound of  $\Omega(\sqrt{N})$  queries for Grover's search [4] applies to perceptron training. Since we assumed that perceptron training requires  $o(\sqrt{N})$  queries this is a contradiction. Thus the true lower bound must be  $\Omega(\sqrt{N})$ .

We have assumed that in the classical setting that the user only has access to the training vectors through an oracle that is promised to draw a uniform sample from  $\{(\phi_1, y_1), \dots, (\phi_N, y_N)\}$ . Since we are counting the number of queries to  $f_w$  it is clear that in the worst possible case that the training vector that the perceptron makes a mistake on can be the last unique value sampled from this list. Thus if the query complexity were  $o(N)$  there would be a contradiction, hence the query complexity is  $\Omega(N)$  classically. □

#### IV. QUANTUM VERSION SPACE PERCEPTRON

The strategy for our quantum version space training algorithm is to pose the problem of determining a separating hyperplane as search. Specifically, the idea is to first generate  $K$  sample hyperplanes  $w_1, \dots, w_K$  from a spherical Gaussian distribution  $\mathcal{N}(0, \mathbf{1})$ . Given a large enough  $K$ , we are guaranteed to have at least one hyperplane amongst the samples that would lie in the version space and perfectly separate the data. As discussed earlier Grover's algorithm can provide quadratic speedup over the classical search consequently the efficiency of the algorithm is determined by  $K$ . Theorem 2 provides an insight on how to determine this number of hyperplanes to be sampled.

**Theorem 2.** *Given a training set that consists of  $d$ -dimensional unit vectors  $\Phi_1, \dots, \Phi_N$  with labels  $y_1, \dots, y_N$  that are separated by a margin of  $\gamma$  in feature space, then a  $D$ -dimensional vector  $w$  sampled from  $\mathcal{N}(0, \mathbf{1})$  perfectly separates the data with probability  $\Theta(\gamma)$ .*

The proof of this theorem is provided in the supplementary material. The consequence of Theorem 2 stated below is that the expected number of samples  $K$ , required such that a separating hyperplane exists in the set, only needs to scale as  $O(\frac{1}{\gamma})$ . Thus if amplitude amplification is used to boost the probability of finding a vector in the version space then the resulting quantum algorithm will need only  $O(\frac{1}{\sqrt{\gamma}})$  quantum steps on average.

Next we show how to use Grover's algorithm to search for a hyperplane that lies in the version space. Let us take  $K = 2^m$ , for positive integer  $m$ . Then given  $w_1, \dots, w_K$  be the sampled hyperplanes, we represent  $W_1, \dots, W_K$  to be vectors that encode a binary representation of these random perceptron vectors. In analogy to  $\Phi_0$ , we also define  $W_0$  to be a vector that represents an empty data register. We define the unitary operator  $V$  to generate these weights given an address vector  $u_j$  using the following

$$V[u_j \otimes W_0] = [u_j \otimes W_j]. \quad (12)$$

In this context we can also think of the address vector,  $u_j$ , as representing a seed for a pseudo-random number generator that yields perceptron weights  $W_j$ .

Also let us define the classical analogue of  $V$  to be  $V^c$  which obeys  $V^c(j) = w_j$ . Now using  $V$  (and applying the Hadamard transform [15]) we can prepare the following quantum state

$$\Psi := \frac{1}{\sqrt{K}} \sum_{k=1}^K u_k \otimes W_k, \quad (13)$$

which corresponds to a uniform distribution over the randomly chosen  $w$ .

Now that we have defined the initial state,  $\Psi$ , for Grover's search we need to define an oracle that marks the vectors inside the version space. Let us define the operator  $\hat{\mathcal{F}}_{\phi, y}$  via

$$\hat{\mathcal{F}}_{\phi, y}[u_j \otimes W_0] = (-1)^{1+f_{w_j}(\phi, y)}[u_j \otimes W_0]. \quad (14)$$

This unitary operation looks at an address vector,  $u_j$ , computes the corresponding perceptron model  $W_j$ , flips the sign of any component of the quantum state vector that is in the half space in version space specified by  $\phi$  and then uncomputes  $W_j$ . This process can be realized using a quantum subroutine that computes  $f_w$ , an application of  $V$  and  $V^\dagger$  and also the application of a conditional phase gate (which is a fundamental quantum operation that is usually denoted  $Z$ ) [15].

---

**Algorithm 2** Quantum version space perceptron training algorithm

---

```

for  $k = 1, \dots, \lceil \log_{3/4} \epsilon \rceil$  do
for  $j = 1 : \lceil \log_c(1/\sin(2\sin^{-1}(1/\sqrt{K}))) \rceil$  do
  Draw  $m$  uniformly from  $\{0, \dots, \lceil c^j \rceil\}$ .
  Prepare quantum state  $\Psi = \frac{1}{\sqrt{K}} \sum_{p=1}^K u_p \otimes W_0$ .
   $\Psi \leftarrow ((2\Psi\Psi^\dagger - \mathbb{1})G)^m \Psi$ .
  Measure  $\Psi$ , assume outcome is  $u_q$ .
   $w \leftarrow V^c(q)$ .
if  $f_w(\phi_\ell, y_\ell) = 0$  for all  $\ell \in \{1, \dots, N\}$  then
  Return  $w$ 
end if
end for
end for
Return  $w = 0$ 

```

---

The oracle  $\hat{\mathcal{F}}_{\phi, y}$  does not allow us to directly use Grover's search to rotate a quantum state vector that is outside the version space towards the version space boundary because it effectively only checks one of the half-space inequalities that define the version space. It can, however, be used to build an operation,  $\hat{G}$ , that reflects about the version space:

$$\hat{G}[u_j \otimes W_0] = (-1)^{1+(f_{w_j}(\phi_1, y_1) \vee \dots \vee f_{w_j}(\phi_N, y_N))} [u_j \otimes W_0]. \quad (15)$$

The operation  $\hat{G}$  can be implemented using  $2N$  applications of  $\hat{\mathcal{F}}_\phi$  as well as a sequence of  $O(N)$  elementary quantum gates, hence we cost a query to  $\hat{G}$  as  $O(N)$  queries to  $\hat{\mathcal{F}}_{\phi, y}$ .

We use these components in Algorithm 2 to, in effect, amplify the margin between the two classes from  $\gamma$  to  $\sqrt{\gamma}$ . We give the asymptotic scaling of this algorithm in the following theorem (see appendix for proof).

**Theorem 3.** *Given a training set that consists of unit vectors  $\Phi_1, \dots, \Phi_N$  that are separated by a margin of  $\gamma$  in feature space, the number of queries to  $\hat{\mathcal{F}}_{\phi, y}$  needed to infer a perceptron model with probability at least  $1 - \epsilon$ ,  $w$ , such that  $w$  is in the version space using a quantum computer is  $N_{\text{quant}}$  where*

$$N_{\text{quant}} \in O\left(\frac{N}{\sqrt{\gamma}} \log^{3/2}\left[\frac{1}{\epsilon}\right]\right).$$

*Proof.* The proof of the theorem follows directly from bounds on  $K$  and the validity of Algorithm 2. It is clear from previous discussions that Algorithm 2 carries out Grover's search, but instead of searching for a  $\phi$  that is misclassified it instead searches for a  $w$  in version space. Its validity therefore follows by following the exact same steps followed in the proof of Lemma 2 but with  $N = K$ . However, since the algorithm need is not repeated  $1/\gamma^2$  times in this context we can replace  $\gamma$  with 1 in the proof. Thus if we wish to have a probability of failure of at most  $\epsilon'$  then the number of queries made to  $\hat{G}$  is in

$$O(\sqrt{K} \log(1/\epsilon')).$$

This also guarantees that if any of the  $K$  vectors are in the version space then the probability of failing to find that vector is at most  $\epsilon'$ .

Next since one query to  $\hat{G}$  is costed at  $N$  queries to  $\hat{\mathcal{F}}_{\phi, y}$  the query complexity (in units of queries to  $\hat{\mathcal{F}}_{\phi, y}$ ) becomes  $O(N\sqrt{K} \log(1/\epsilon'))$ . The only thing that then remains is to bound the value of  $K$  needed.

The probability of finding a vector in the version space is  $\Theta(\gamma)$  from Theorem 2. This means that there exists  $\alpha > 0$  such that the probability of failing to find a vector in the version space  $K$  times is at most

$$(1 - \alpha\gamma)^K \leq e^{-\alpha\gamma K}. \quad (16)$$

Thus this probability is at most  $\delta$  for

$$K \in \Omega\left(\frac{1}{\gamma} \log(1/\delta)\right). \quad (17)$$

It then suffices to pick  $K \in \Theta\left(\frac{1}{\gamma} \log(1/\delta)\right)$  for the algorithm.

The union bound implies that the probability that either none of the vectors lie in the version space or that Grover's search failing to find such an element is at most  $\epsilon' + \delta \leq \epsilon$ . Thus it suffices to pick  $\epsilon' \in \Theta(\epsilon)$  and  $\delta \in \Theta(\epsilon)$  to ensure that the total probability is at most  $\epsilon$ . Therefore the total number of queries made to  $\hat{\mathcal{F}}_{\phi, y}$  is in  $O\left(\frac{N}{\sqrt{\gamma}} \log^{3/2}(1/\epsilon)\right)$  as claimed.  $\square$

The classical algorithm discussed previously has complexity  $O(N \log(1/\epsilon)/\gamma)$ , which follows from the fact (proven in the appendix) that  $K \in \Theta(\log(1/\epsilon)/\gamma)$  suffices to make the probability of not drawing an element of the version space at most  $\epsilon$ . This demonstrates a quantum advantage if  $\frac{1}{\gamma} \gg \log(1/\epsilon)$ , and illustrates that quantum computing can be used to boost the effective margins of the training data. Quantum models of perceptrons therefore not only provide advantages in terms of the number of vectors that need to be queried in the training process, they also can make the perceptron much more perceptive by making training less sensitive to small margins.

These performance improvements can also be viewed as mistake bounds for the version space perceptron. The inner loop in Algorithm 2 attempts to sample from the version space and then once it draws a sample it tests it against the training vectors to see if it errs on any example. Since the inner loop is repeated  $O(\sqrt{K} \log(1/\epsilon))$  times, the maximum number of misclassified vectors that arises from this training process is from Theorem 2  $O(\frac{1}{\sqrt{\gamma}} \log^{3/2}(1/\epsilon))$  which, for constant  $\epsilon$ , constitutes a quartic improvement over the standard mistake bound of  $1/\gamma^2$  [16].

## V. CONCLUSION

We have provided two distinct ways to look at quantum perceptron training that each afford different speedups relative to the other. The first provides a quadratic speedup with respect to the size of the training data. We further show that this algorithm is asymptotically optimal in that if a super-quadratic speedup were possible then it would violate known lower bounds for quantum searching. The second provides a quadratic reduction in the scaling of the training time (as measured by the number of interactions with the training data) with the margin between the two classes. This latter result is especially interesting because it constitutes a quartic speedup relative to the typical perceptron training bounds that are usually seen in the literature.

Perhaps the most significant feature of our work is that it demonstrates that quantum computing can provide provable speedups for perceptron training, which is a foundational machine learning method. While our work gives two possible ways of viewing the perceptron model through the lens of quantum computing, other quantum variants of the perceptron model may exist. Seeking new models for perceptron learning that deviate from these classical approaches may not only provide a deeper understanding of what form learning takes within quantum systems, but also may lead to richer classes of quantum models that have no classical analogue and are not efficiently simulatable on classical hardware. Such models may not only revolutionize quantum learning but also lead to a deeper understanding of the challenges and opportunities that the laws of physics place on our ability to learn.

## Appendix A: Proofs

Here we provide proofs of several of our results stated in the main body. In particular, we give the proofs of Theorem 2, Lemma 1 and Lemma 2 here.

*Proof of Theorem 2.* Given that the margin of the training set is  $\gamma$  there exist a hyperplane  $u$  such that  $y_i \cdot u^T \Phi_i > \gamma$  for all  $i$ . If  $w$  be a sample from  $\mathcal{N}(0, \mathbb{1})$ , then let's first compute what is the probability that perturbing the maximum margin classifier  $u$  by amount  $w$  would still lead to a perfect separation. If we consider a data point  $\Phi^*$  that lies on the margin, i.e.  $y_i \cdot u^T \Phi^* = \gamma$ , we are interested in the probability that  $y_i \cdot (u+w)^T \Phi^* > 0$  and  $y_i \cdot (u+w)^T \Phi^* < 2\gamma$ . The first inequality corresponds to preventing misclassification of  $\Phi^*$ , while the second one corresponds to preventing misclassification of the point belonging to the other class and on the margin. This is same as asking what is the probability that:

$$-\gamma < y_i \cdot w^T \Phi^* < \gamma \quad (\text{A1})$$

Let us define  $z_i := y_i \cdot w^T \Phi^*$ . Since  $w \sim \mathcal{N}(0, \mathbb{1})$  and  $\|\Phi^*\| = 1$  we can show that  $z_i \sim \mathcal{N}(0, 1)$ . Thus, we can write the probability that  $-\gamma < z_i < \gamma$  as:

$$P(-\gamma < z_i < \gamma) = \text{erf}\left(\frac{\gamma}{\sqrt{2}}\right) \quad (\text{A2})$$

Here  $\text{erf}(z) = \frac{1}{\sqrt{\pi}} \int_{-z}^z e^{-x^2/2} dx$  is the error function for the standard normal distribution. Since  $\Phi^*$  is on the margin, the probability that the sample  $w$  will lie in the version space can be simply characterized as the above probability  $P(-\gamma < z_i < \gamma)$ . It is straightforward to show using Maclaurin series expansion that:

$$P(w \in VS) = \frac{2}{\sqrt{\pi}} \left( \frac{\gamma}{\sqrt{2}} - \frac{\gamma^3}{2^3/2^3} + \frac{\gamma^5}{2^5/2^10} - \frac{\gamma^7}{2^7/2^42} \dots \right) \quad (\text{A3})$$

Note, that in our case  $\Phi_i$  are unit normalized for all  $i$ , thus  $\gamma < 1$ . Which in turn implies that most of the higher order terms will be close to zero in the limit of small  $\gamma$  and:

$$P(w \in VS) = \frac{\gamma}{\sqrt{2\pi}} + O(\gamma^3), \quad (\text{A4})$$

which proves our theorem for  $\gamma < 1$ .  $\square$

*Proof of Lemma 1.* There exists a simple algorithm for achieving this upper bound. Draw  $N \lceil \log(1/\epsilon\gamma^2) \rceil$  samples from the set of training vectors. If any are misclassified perform the update, otherwise report that the model classifies all the data.

The proof of validity of this algorithm is trivial and the success probability claim is also quite simple. Given that we draw  $k$  samples from the distribution the probability that any of them fail to detect a mistake, given such a mistake exists, is at most

$$(1 - 1/N)^k \leq \exp(-k/N). \quad (\text{A5})$$

If we want this error to be at most  $\delta$  then it suffices to take

$$k = \lceil N \log(1/\delta) \rceil. \quad (\text{A6})$$

One query to  $f_w$  is required per  $k$ , which means that  $k$  is also equal to the query complexity. Thus if at least one mistake occurs then the algorithm will find it with the aforementioned probability if  $\delta = \epsilon\gamma^2$ . If such an example does not exist, then the algorithm will correctly conclude that a separating hyperplane has already been found. Therefore in either case the success probability is at least  $1 - \epsilon\gamma^2$  as required.  $\square$

*Proof of Lemma 2.* In order to see this, let us first examine the inner loop of Algorithm 2, which involves performing the update  $\Psi \leftarrow ((2\Psi\Psi^\dagger - \mathbf{1})F_w)^m\Psi$ . We know from our discussion of Grover's algorithm in the main body that if we define the initial probability of successfully find a mistake to be  $\sin^2(\theta_a)$  then the probability of finding a  $j$  such that  $\mathcal{F}_w\Phi_j = -1$  after  $m$  updates is  $\sin^2((2m+1)\theta_a)$ . Since this corresponds to finding a vector that the perceptron fails to classify properly, these steps amplify the probability of finding a perceptron error. The query to  $U^c$  that follows identifying the index of this training vector then converts this result into a classical bitstring that can then be used to perform a perceptron update. Therefore the inner loop performs a perceptron update with probability  $\sin^2((2m+1)\theta_a)$  using  $m$  queries to  $F_w$ .

Under the assumption that  $c \in (1, 2)$  the next loop repeats this sampling process until  $m \geq M_0$  in order to ensure that the probability of finding a misclassified element is at least  $1/4$  [5]. This can be seen using the following argument. First we need to show that the exponential search heuristic requires  $O(M_0)$  queries. Each iteration of the middle loop requires requires a number of queries that is at most proportional to  $\lceil c^i \rceil$ . Therefore the total number of queries is at most proportional to

$$\begin{aligned} \sum_{i=0}^{\lceil \log_c M_0 \rceil} \lceil c^i \rceil &\leq \frac{c^{\lceil \log_c M_0 \rceil + 1}}{c-1} + \lceil \log_c M_0 \rceil \\ &\leq \frac{c}{c-1} M_0 + \lceil \log_c M_0 \rceil. \end{aligned} \quad (\text{A7})$$

Given  $c$  is a constant we have that  $(c-1) \in \Theta(1)$  and thus  $\sum_{i=0}^{\lceil \log_c M_0 \rceil} c^i \in O(1/\sin(\theta_a))$  from (4) in the main body. If there exists an element that the algorithm makes a mistake on then  $\theta_a \geq \sin^{-1}(\sqrt{1/N}) \in \Omega(1/\sqrt{N})$  because the lowest probability of success corresponds to the case where there is only one training vector that is misclassified out of  $N$ . From this we see that, if a misclassified vector exists, then the middle loop is repeated at least  $\log_c(M_0)$  times which means that the final iteration taken corresponds to  $m \geq M_0$  for the purposes of (6) in the main body. Therefore, under these assumptions, the probability that the middle loop updates the perceptron weights is at least  $1/4$  from [5]. Given that a mistake exists to be found the middle loop outputs such an element with a probability of failure that is at most  $3/4$  from (6) in the main body. Furthermore,  $O(M_0) = O(1/\sqrt{N})$  queries to  $F_w$  are required by the inner loop.

The outer loop serves to amplify the success probability to at least  $1 - \epsilon$  from the (average) success probability for  $m \geq M_0$ , which is at least  $1/4$ , given that the perceptron makes a mistake on at least one training vector [5]. Let us assume that we repeat the middle loop of Algorithm 1  $k$  times past this point and terminate searching for a marked state if the probability of failing to detect the element is at most  $\delta$ . Since the probability of the middle loop failing to find such an element, given that it exists, is at most  $3/4$  the probability of failing to find a marked state all  $k$  times is at most  $(3/4)^k$  which implies that it suffices to choose

$$k = \lceil \log_{3/4}(\delta) \rceil. \quad (\text{A8})$$

Given this error bound, the number of Grover iterations needed for the algorithm to find the marked element is

$$O(\sqrt{N} \log_{3/4} \delta) \in O(\sqrt{N} \log(1/\delta)).$$

The result then follows by taking  $\delta = \epsilon\gamma^2$ . Therefore the lemma holds if  $\theta_a > 0$ .

If  $\theta_a = 0$  then the algorithm will never find a quantum state vector that the perceptron misclassifies and will successfully conclude that there is not a marked state after  $O(\sqrt{N} \log(1/\epsilon\gamma^2))$  queries. Therefore the lemma also holds in the trivial case.  $\square$

- 
- [1] Aïmeur, Esma, Brassard, Gilles, and Gambs, Sébastien. Machine learning in a quantum world. In *Advances in artificial intelligence*, pp. 431–442. Springer, 2006.
  - [2] Amin, Mohammad H, Andriyash, Evgeny, Rolfe, Jason, Kulchytskyy, Bohdan, and Melko, Roger. Quantum boltzmann machine. *arXiv preprint arXiv:1601.02036*, 2016.
  - [3] Bell, John S. On the einstein podolsky rosen paradox, 1964.
  - [4] Boyer, Michel, Brassard, Gilles, Høyer, Peter, and Tapp, Alain. Tight bounds on quantum searching. *arXiv preprint quant-ph/9605034*, 1996.
  - [5] Brassard, Gilles, Hoyer, Peter, Mosca, Michele, and Tapp, Alain. Quantum amplitude amplification and estimation. *Contemporary Mathematics*, 305:53–74, 2002.
  - [6] Freund, Yoav and Schapire, Robert E. Large margin classification using the perceptron algorithm. *Machine learning*, 37(3):277–296, 1999.
  - [7] Garnerone, Silvano, Zanardi, Paolo, and Lidar, Daniel A. Adiabatic quantum algorithm for search engine ranking. *Physical review letters*, 108(23):230506, 2012.
  - [8] Gentile, Claudio. A new approximate maximal margin classification algorithm. *The Journal of Machine Learning Research*, 2:213–242, 2002.
  - [9] Grover, Lov K. A fast quantum mechanical algorithm for database search. In *Proceedings of the twenty-eighth annual ACM symposium on Theory of computing*, pp. 212–219. ACM, 1996.
  - [10] Herbrich, Ralf, Graepel, Thore, and Campbell, Colin. Bayes point machines: Estimating the bayes point in kernel space. In *IJCAI Workshop SVMs*, pp. 23–27, 1999.
  - [11] Li, Yaoyong, Zaragoza, Hugo, Herbrich, Ralf, Shawe-Taylor, John, and Kandola, Jaz. The perceptron algorithm with uneven margins. In *ICML*, volume 2, pp. 379–386, 2002.
  - [12] Lloyd, Seth, Mohseni, Masoud, and Rebentrost, Patrick. Quantum algorithms for supervised and unsupervised machine learning. *arXiv preprint arXiv:1307.0411*, 2013.
  - [13] Lloyd, Seth, Mohseni, Masoud, and Rebentrost, Patrick. Quantum principal component analysis. *Nature Physics*, 10(9):631–633, 2014.
  - [14] Minka, Thomas P. *A family of algorithms for approximate Bayesian inference*. PhD thesis, Massachusetts Institute of Technology, 2001.
  - [15] Nielsen, Michael A and Chuang, Isaac L. *Quantum computation and quantum information*. Cambridge university press, 2010.
  - [16] Novikoff, Albert BJ. On convergence proofs for perceptrons. Technical report, DTIC Document, 1963.
  - [17] Rebentrost, Patrick, Mohseni, Masoud, and Lloyd, Seth. Quantum support vector machine for big data classification. *Physical review letters*, 113(13):130503, 2014.
  - [18] Rosenblatt, Frank. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological review*, 65(6):386, 1958.
  - [19] Shalev-Shwartz, Shai and Singer, Yoram. A new perspective on an old perceptron algorithm. In *Learning Theory*, pp. 264–278. Springer, 2005.
  - [20] Suykens, Johan AK and Vandewalle, Joos. Least squares support vector machine classifiers. *Neural processing letters*, 9(3):293–300, 1999.
  - [21] Wiebe, Nathan and Granade, Christopher. Can small quantum systems learn? *arXiv preprint arXiv:1512.03145*, 2015.
  - [22] Wiebe, Nathan, Kapoor, Ashish, and Svore, Krysta M. Quantum deep learning. *arXiv preprint arXiv:1412.3489*, 2014.
  - [23] Wiebe, Nathan, Kapoor, Ashish, and Svore, Krysta. Quantum nearest-neighbor algorithms for machine learning. *Quantum Information and Computation*, 15:318–358, 2015.