

Learning DNFs under product distributions via μ -biased quantum Fourier sampling

Varun Kanade*

Andrea Rocchetto[†]Simone Severini[‡]

Abstract

We show that DNF formulae can be quantum PAC-learned in polynomial time under product distributions using a quantum example oracle. The best classical algorithm (without access to membership queries) runs in superpolynomial time. Our result extends the work by Bshouty and Jackson (1998) that proved that DNF formulae are efficiently learnable under the uniform distribution using a quantum example oracle. Our proof is based on a new quantum algorithm that efficiently samples the coefficients of a μ -biased Fourier transform.

1 Introduction

The question whether the class of boolean functions that can be expressed as polynomial size formulae in disjunctive normal form (DNF) is *probably approximately correct* (PAC) learnable in polynomial time is one of the central unresolved problems in the PAC learning framework introduced by Valiant [Val84]. Currently, the best classical algorithm for this algorithm has running time $2^{\tilde{O}(n^{1/3})}$ [KS01]. A number of variants of this problem have been studied, by relaxing the requirements, e.g. learning with respect to restricted classes of distributions, and/or by enhancing the power of the learning algorithms, e.g. providing access to membership query, random walk, or a quantum example oracle [Jac94; AFK13; Bsh+05; BJ98].

Two cases in which it is possible to show DNF learnability under specific assumptions are particularly relevant to our setting. First, when the distribution is uniform, or more generally, a product distribution, a quasi-polynomial $n^{O(\log(n))}$ algorithm is known [Ver90]. Second, in the *membership query* (MQ) model, where the learner can query an oracle for a value of the unknown function at a given point in the domain, Jackson [Jac97] gave a polynomial time algorithm for DNFs that works over both the uniform and product distributions.

The question of PAC learning has been extended to learners with access to quantum resources by Bshouty and Jackson [BJ98]. The main results in quantum learning theory are reviewed by Arunachalam and de Wolf [AW17a]. The two main requirements for a quantum PAC-learner are the ability to query an oracle that can provide examples in quantum superposition and access to a quantum computer to run the learning algorithm. Two measures of interest in the PAC framework are the sample complexity and the time complexity. The *sample complexity* is the worst-case number of examples required to learn a class of functions. The *time complexity* is the worst-case scaling of a learner for that function class. It has been shown that the quantum PAC model gives only a constant factor advantage in terms of sample complexity with respect to the classical version [AW17b]. Certain results suggest that a classical/quantum

*University of Oxford and The Alan Turing Institute. Email: varunk@cs.ox.ac.uk

[†]University of Oxford and University College London. Email: andrea.rocchetto@spc.ox.ac.uk

[‡]University College London and Shanghai Jiao Tong University. Email: s.severini@ucl.ac.uk

separation exists when considering the time complexity of some learning problems. When learning with respect to the uniform distribution, the class of polynomial-size DNF formulae [BJ98] and k -juntas [AS07] under the uniform distribution are known to be efficiently quantum PAC-learnable (note that the learnability of k -juntas is implied by the result on DNFs). In the classical setting, in both these cases current best known algorithms are quasi-polynomial time algorithms; while no formal hardness results are known, it would be highly surprising if a polynomial time algorithm for these algorithms in the classical setting was discovered. In the context of learning in the presence of noise, [CSS15] proved that parity functions under uniform can be efficiently learned using a quantum example oracle. Classically, the problem is widely believed to require subexponential, but superpolynomial, time [BKW03; Lyu05]. The result of Cross, Smith, and Smolin [CSS15] has been extended to linear functions and to more complex error models by [GK17].

2 Overview of the results

We show that DNF formulae under the product distribution can be learned in polynomial time in the quantum PAC model. Our proof builds on the work by Feldman [Fel12] for learning DNFs under product using membership queries. Feldman’s proof is in turn based on a result by Kalai, Samorodnitsky, and Teng [KST09] that shows that DNFs can be approximated by heavy low-degree Fourier coefficients alone. Notably, Feldman’s result also applies to learning settings where the examples are drawn from a product distribution, *i.e.* a distribution that factorises over the elements of the input vector.

The only part of Feldman’s algorithm that makes use of MQ is the subroutine that approximates the Fourier spectrum of f . The approximation is obtained using the *Kushilevitz-Mansour* (KM) algorithm [KM93] (for the case of uniform distributions) and the *Extended Kushilevitz-Mansour* (EKM) algorithm [KST09] (for the case of product distributions). Bshouty and Jackson showed that it is possible to approximate the Fourier coefficients of f using quantum Fourier sampling. This technique, introduced by Bernstein and Vazirani [BV97], allows one to sample efficiently from the distribution defined by $\hat{f}^2(a)$ using the *Quantum Fourier Transform* (QFT).

In order to extend the result by Bshouty and Jackson for learning under product it is necessary to find a quantum technique to sample according to the coefficients of a Fourier transform defined over an inner product where each term is weighted according to the product distribution. Bahadur [Bah61] and Furst, Jackson and Smith [FJS91] showed that the Fourier transform can be extended to product distributions, thus defining the μ -biased Fourier transform.

In this work we introduce the μ -biased quantum Fourier transform. We show the validity of our construction in two steps. At first we explicitly construct a unitary that implements the single qubit transform. Then we argue that this construction can be efficiently implemented on a quantum circuit with logarithmic overhead. By exploiting the factorisation of the product distribution we show that how to build an n -qubit transform as a tensor product of n single qubit transforms.

The main technical contribution of this paper is a quantum algorithm to approximate the heavy μ -biased Fourier spectrum of f without using membership queries. This can be interpreted as an immediate quantum version of the EKM algorithm for approximating the Fourier coefficients of f . We provide rigorous upper bounds on the scaling of the algorithm using the Dvoretzky–Kiefer–Wolfowitz theorem, a concentration inequality that bounds the number of samples required to estimate a probability distribution in infinity norm. The learnability of DNFs under product immediately follows from an application of the quantum EKM algorithm to Corollary 5.1 in [Fel12].

2.1 Related works

The learnability of DNF formulae under the uniform distribution using a quantum example oracle was first studied by Bshouty and Jackson that, in the same paper, also introduced the quantum PAC model [BJ98]. Their approach to learning DNF was built on the *harmonic sieve* algorithm previously developed by Jackson [Jac94]. At the core of Jackson’s algorithm lies a useful property of DNFs which guarantees that, for every s -term DNF and for every probability distribution \mathcal{D} , there exists a parity χ_a such that $|\mathbf{E}_{\mathcal{D}}[f\chi_a]| \geq 1/(2s + 1)$ [Jac94]. This implies that for every f and \mathcal{D} there exists a parity that weakly approximates f . In the harmonic sieve algorithm Freund’s boosting algorithm is then used to turn the weak learner into a strong one [Fre95]. The only part of the harmonic sieve algorithm that requires MQ is the KM algorithm used to find the weakly approximating parity function. Bshouty and Jackson assume that the examples are given by a quantum example oracle and replace the KM algorithm with quantum Fourier sampling [BV97].

Jackson, Tamon, and Yamakami studied the learnability of DNFs in the quantum membership model [JTY02] (where the quantum example oracle is replaced by an oracle that returns $f(x)$ for a given x). By using the quantum Goldreich-Levin algorithm developed by Adcock and Cleve [AC02] they were able to obtain a better bound on the query complexity with respect to the best classical algorithm. We recall that the classical KM algorithm can be derived from the Goldreich-Levin theorem, an important result that reduces the computational problem of inverting a one-way function to the problem of predicting a given hard-predicate associated with that function [GL89]. Adcock and Cleve’s result shows that this reduction can be obtained more efficiently when considering quantum functions and quantum hard-predicates. A different quantum implementation of the Goldreich-Levin algorithm was given by Montanaro and Osborne [MO10].

2.2 Organisation

We introduce notations and important concepts in Section 3. In Section 5 we define the μ -biased quantum Fourier transform and discuss some of its properties. In Section 6 we introduce an efficient quantum algorithm to sample efficiently from the Fourier coefficients of the μ -biased Fourier transform and show how this can be used to prove the PAC-learnability of DNF formulae under product. We conclude in Section 7 where we discuss how to implement the QEX oracle. In Appendix A we bound the error introduced by approximating μ .

3 Preliminaries

3.1 Notation

We denote vectors with lower-case letters. For a vector $x \in \mathbb{R}^n$, let x_i denotes the i -th element of x . A vector is sparse if most of its entries are 0. If x is sparse we can describe it using only its non-zero coefficients. We call this the *succint representation* of x . For an integer k , let $[k]$ denotes the set $\{1, \dots, k\}$. We use the following standard norms. The ℓ_0 norm $\|x\|_0 = |\{i \in [k] | x_i \neq 0\}|$, the ℓ_2 norm $\|x\|_2 = \sqrt{\sum_{i \in [k]} x_i^2}$, and the ℓ_∞ norm $\|x\|_\infty = \max_{i \in [k]} \{|x_i|\}$.

Let $f : \mathbb{R} \rightarrow \mathbb{R}^+$ and $g : \mathbb{R} \rightarrow \mathbb{R}^+$. We use $f(n) = O(g(n))$ to indicate that the asymptotic scaling of $|f|$ is upper-bounded, up to a constant factor, by $g(n)$. Similarly, $f(n) = \Omega(g(n))$ indicates that the asymptotic scaling of $|f|$ is lower-bounded, up to a constant factor, by g . The notation $f(n) = \Theta(g(n))$ indicates that f is bounded both above and below by g asymptotically. The notations $\tilde{O}(g(n))$ and $\tilde{\Omega}(g(n))$ hide logarithmic factors.

For a set V and $x \in V$ we denote by $\mathcal{D}(x)$ a probability distribution over V . The notation $x \sim \mathcal{D}$ indicates that x is sampled according to \mathcal{D} . The expected value of a random variable f is denoted as $\mathbf{E}_{x \sim \mathcal{D}}[f(x)]$. We often use $\mathbf{E}_{\mathcal{D}}[\cdot]$ to indicate $\mathbf{E}_{x \sim \mathcal{D}}[\cdot]$. If \mathcal{D} has a subscript, as in \mathcal{D}_μ , we write $\mathbf{E}_\mu[\cdot]$ to indicate $\mathbf{E}_{\mathcal{D}_\mu}[\cdot]$. When \mathcal{D} is the uniform distribution we omit the distribution in the subscript and use $\mathbf{E}[\cdot]$. The probability that an event A occurs is denoted by $\Pr[A]$.

3.2 Fourier analysis over the Boolean cube

Let $x \in \{-1, 1\}^n$ and let f and g be real functions defined over the Boolean hypercube $f, g : \{-1, 1\}^n \rightarrow [-1, 1]$. The space of real functions over the Boolean hypercube is a vector space with inner product $\langle f, g \rangle = \frac{1}{2^n} \sum_{x \in \{0, 1\}^n} f(x)g(x) = \mathbf{E}[f \cdot g]$ where the expectation is taken uniformly over all $x \in \{0, 1\}^n$. A *parity function* $\chi_a : \{-1, 1\}^n \rightarrow \{-1, 1\}$ labels a $x \in \{-1, 1\}^n$ according to a characteristic vector $a \in \{0, 1\}^n$ and is defined as $\chi_a(x) = (-1)^{a \cdot x}$. The set of parity functions $\{\chi_a(x)\}_{a \in \{0, 1\}^n}$ forms an orthonormal base for the space of real functions over the Boolean hypercube. This fact implies that we can uniquely represent every function f as a linear combination of parities, the *Fourier transform* of f . The linear coefficients, known as the *Fourier coefficients*, are given by the projections of the function into the parity base and are denoted with $\hat{f}(a) = \langle f, \chi_a \rangle = \mathbf{E}[f(x)\chi_a(x)]$. We say that a Fourier coefficient is “heavy” if has large magnitude $|\hat{f}(a)|$. The set of Fourier coefficients is called the *Fourier spectrum* of f and is denoted by \hat{f} , which can also be seen as a 2^n dimensional vector in \mathbb{R}^{2^n} . For a set $S \subseteq \{0, 1\}^n$, $\hat{f}(S)$ denotes the vector of all Fourier coefficients with indices in S . The *degree* of a Fourier coefficient $\hat{f}(a)$ is $\|a\|_0$. Let $B_d = \{a \in \{0, 1\}^n \mid \|a\|_0 \leq d\}$. We denote by $\hat{f}(B_d)$ vector of all degree $\leq d$ coefficients of f . The squared Fourier coefficients are related by Parseval’s equality $\mathbf{E}[f^2] = \sum_a \hat{f}(a)^2 = \|\hat{f}\|_2^2$. This implies that for any $f : \{-1, 1\}^n \rightarrow [-1, 1]$, $\sum_a \hat{f}(a)^2 \leq 1$ (the equality holds if f is Boolean-valued).

The Fourier spectrum of a function f can be approximated using the KM algorithm. The KM algorithm, based upon a celebrated result by Goldreich and Levin [GL89], requires oracular access to f (i.e. it requires an oracle that for every $x \in \{-1, 1\}^n$ returns $f(x)$).

Theorem 1 (KM algorithm). *Let $f : \{-1, 1\}^n \rightarrow [-1, 1]$ be a real-valued function and let $\epsilon > 0$, $\delta > 0$, $\mu \in (-1, 1)^n$. Then, with probability at least $1 - \delta$, there exists an algorithm with oracle access to f that returns a succinctly represented vector \tilde{f} such that $\|\hat{f} - \tilde{f}\|_\infty \leq \epsilon$ and $\|\tilde{f}\|_0 \leq 4/\epsilon^2$. The algorithm runs in $\tilde{O}(n^2 \log(1/\delta)/\epsilon^6)$ time and makes $\tilde{O}(n \log(1/\delta)/\epsilon^6)$ queries to f .*

3.3 μ -biased Fourier analysis

A *product distribution* \mathcal{D}_μ over $\{-1, 1\}^n$ is characterised by a real vector $\mu \in (-1, 1)^n$. Such a distribution \mathcal{D} assigns values to each variable independently, so for $x \in \{-1, 1\}^n$ we have $\mathcal{D}_\mu(x) = \prod_{i: x_i=1} (1 + \mu_i)/2 \prod_{i: x_i=-1} (1 - \mu_i)/2$ and $\mathbf{E}_\mu[x_i] = \mu_i$. We say that the distribution \mathcal{D}_μ is c -bounded if $\mu \in [-1 + c, 1 - c]^n$, where $c \in (0, 1]$.

Bahadur [Bah61] and Furst, Jackson and Smith [FJS91] showed that the Fourier transform can be extended to product distributions, thus defining the μ -biased Fourier transform. The book by O’Donnell gives a brief introduction to μ -biased Fourier analysis and its applications [ODo14]. For an inner product $\langle f, g \rangle_\mu = \mathbf{E}_\mu[f(x)g(x)]$, the set of functions $\{\phi_{\mu,a} \mid a \in \{0, 1\}^n\}$, where $\phi_{\mu,a}(x) = \prod_{i: a_i=1} (x_i - \mu_i) / \sqrt{1 - \mu_i^2}$ forms an orthonormal basis for the vector space of real-valued functions on $\{-1, 1\}^n$. In this way every function $f : \{-1, 1\}^n \rightarrow [-1, 1]$ can be represented as $f(x) = \sum_{a \in \{0, 1\}^n} \hat{f}_\mu(a) \phi_{\mu,a}(x)$, where $\hat{f}_\mu(a) = \mathbf{E}_\mu[f(x)\phi_{\mu,a}(x)]$. For vectors of μ -biased Fourier coefficients we extend the same notation introduced for standard Fourier

coefficients. Parseval's equality extends to product distributions $\mathbf{E}_\mu[f^2] = \sum_a \hat{f}_\mu(a)^2 = \|\hat{f}_\mu\|_2^2$. This implies that for any $f : \{-1, 1\}^n \rightarrow [-1, 1]$, $\sum_a \hat{f}_\mu(a)^2 \leq 1$.

The KM algorithm has been extended to product distributions by Bellare [Bel91], Jackson [Jac94] and Kalai *et al.* [KST09]. We follow the presentation of Feldman and give the version of Kalai *et al.* also known as the Extended Kushilevitz Mansour (EKM) algorithm.

Theorem 2 (EKM algorithm). *Let $f : \{-1, 1\}^n \rightarrow [-1, 1]$ be a real-valued function and let $\epsilon > 0$, $\delta > 0$, $\mu \in (-1, 1)^n$. Then, with probability at least $1 - \delta$, there exists an algorithm with oracle access to f that returns a succinctly represented vector \tilde{f}_μ such that $\|\hat{f}_\mu - \tilde{f}_\mu\|_\infty \leq \epsilon$ and $\|\tilde{f}_\mu\|_0 \leq 4/\epsilon^2$. The algorithm runs in time polynomial in $n, 1/\epsilon$ and $\log(1/\delta)$.*

3.4 Quantum computation and quantum Fourier transform

A generic n -qubit state is a complex vector, also known as the *state vector*, acting on a Hilbert space of dimension 2^n equipped with an Hermitian scalar product $\langle \cdot | \cdot \rangle$. We use the Dirac notation to denote quantum states and write $|\psi\rangle$ to denote the quantum state ψ . Given a basis $\{|b_i\rangle\}_{i \in [2^n]}$ the elements of $|\psi\rangle$ correspond to its projections over the basis elements. Each element of $|\psi\rangle$ corresponds to a different measurable outcome. The probability of measurement outcome i is $p(i) = |\psi_i|^2$, where $\psi_i = \langle \psi | b_i \rangle \in \mathbb{C}$ is the projection of $|\psi\rangle$ onto $|b_i\rangle$. Let $|\psi\rangle$ and $|\phi\rangle$ be two interacting quantum states, their joint description $|\tau\rangle$ is given by the tensor product of the respective state vectors $|\tau\rangle = |\psi\rangle \otimes |\phi\rangle$.

The evolution of quantum states is governed by *quantum operators*. Quantum operators acting on a n -qubit states are 2^n dimensional unitary matrices and are denoted with capital letters. Let $x \in \{0, 1\}^n$, the QFT over \mathbb{Z}_2^n is defined as $H^{\otimes n} |x\rangle = 2^{-n/2} \sum_{a \in \{0, 1\}^n} (-1)^{x \cdot a} |a\rangle$, where H is the Hadamard transform $H = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$.

We often work in the *computational basis* $\{|i\rangle\}$ where, for an n -qubit system, each basis element corresponds to an n -bit string. A single qubit system can take two values $|0\rangle$ and $|1\rangle$. When working on the Boolean hypercube $\{-1, 1\}^n$ we take $0 \equiv -1$ and $1 \equiv 1$. A *quantum register* is a collection of qubits. Given a Boolean-valued function $f : \{-1, 1\}^n \rightarrow \{-1, 1\}$ that can be efficiently computed by a classical circuit, a *quantum membership oracle* O_f is a unitary map that applied on $n + 1$ qubits acts as follow: $O_f : |i\rangle |0\rangle \rightarrow |i\rangle |f(i)\rangle$. By combining a membership oracle with the Hadamard transform it is possible to produce a *quantum phase oracle* $OP_f : |i\rangle |-\rangle \rightarrow (-1)^{f(x)} |i\rangle |-\rangle$, where $|-\rangle = \frac{1}{\sqrt{2}}(|0\rangle - |1\rangle)$. This operation is also known as *phase kickback*. For ease of notation, in the following we will not write explicitly the ancilla register $|-\rangle$.

Given a probability distribution whose density is efficiently integrable there exists an efficient technique developed by Grover and Rudolph to generate a quantum superposition which approximate the distribution [GR02].

Lemma 1. *Let $\mathcal{D}(x_i)$ over $\{0, 1\}^n$ be a probability distribution over $\{0, 1\}^n$ whose density is efficiently integrable. Then, there exists an efficient quantum algorithm that returns the quantum state*

$$|\psi\rangle = \sum_i \sqrt{\mathcal{D}(x_i)} |i\rangle.$$

3.5 PAC learning and quantum PAC learning

In the PAC model developed by Valiant [Val84] a learner tries to approximate with high probability an unknown concept f from a training set of m random labelled examples $\{(x_i, f(x_i))\}_{i \in [m]}$. The examples are given by an *example oracle* $\text{EX}(f, \mathcal{D})$ that returns an example $(x, f(x))$, where x is randomly sampled from a probability distribution \mathcal{D} over $\{-1, 1\}^n$. A *concept class* C is a

set of concepts. A learning algorithm \mathcal{A} gets as input the training set and outputs a hypothesis h that is a good approximation of f with probability ϵ . We say that a concept class C is *PAC-learnable* if, for every \mathcal{D} , f , h , δ , when running a learning algorithm L on $m \geq m_C$ examples generated by \mathcal{D} , we have that, with probability at least $1 - \delta$, $\Pr_{x \sim \mathcal{D}}[h(x) \neq f(x)] \leq \epsilon$. PAC theory introduces two parameters to classify the efficiency of a learner. The first one, m_C , is information-theoretic and determines the minimum number of examples required to PAC-learn the class C . We refer to m_C as the *sample complexity* of the concept class C . The second parameter, the *time complexity*, is computational and corresponds to the runtime of the best learner for the class C . We say that a concept class is *efficiently* PAC-learnable if the running time of L is polynomial in n , $1/\epsilon$ and $1/\delta$.

Two extensions of the PAC model are relevant for our purposes. In the MQ model the learner has access, in addition to the example oracle $\text{EX}(f, \mathcal{D})$, to a membership oracle $\text{MQ}(f)$ that for every x returns the value $f(x)$. In the quantum PAC model, the examples are given by a *quantum example oracle* $\text{QEX}(f, \mathcal{D})$ that returns the superposition $\sum_x \sqrt{\mathcal{D}(x)} |x, f(x)\rangle$. It has been proven [BJ98] that membership queries are strictly more powerful than a quantum example oracle (*i.e.* a quantum example oracle can be simulated by a membership oracle but the converse is not true). When $\mathcal{D}(x)$ is the product distribution we use $\text{EX}(f, \mu)$ and $\text{QEX}(f, \mu)$.

A *DNF formula* $f : \{-1, 1\}^n \rightarrow \{-1, 1\}$ is a disjunction of terms where each term is a conjunction of Boolean literals and a literal is either a variable or its negation (*e.g.*, $f(x) = (x_2 \wedge x_3 \wedge \neg x_1) \vee (\neg x_4 \wedge x_1)$). The *size of a DNF* s is the number its conjunctions (also known as *terms*). In our example $s = 2$.

4 Overview of Feldman's algorithm

Our proof of the learnability of DNFs under product builds on an algorithm by Feldman that greatly simplified the learnability of DNFs [Fel12]. At the core of Feldman's algorithm lies a result by Kalai, Samorodnitsky, and Teng [KST09] that shows that DNFs can be approximated by heavy low-degree Fourier coefficients alone. More formally, they proved that, for any s -term DNF f , it is possible to find a function g that is ϵ -close to f provided that $\|\hat{f} - \hat{g}\|_\infty \leq \epsilon/(2s+1)$. This fact gives a direct learnability condition and avoids an involved boosting procedure to turn a weak learner into a strong one (as in Jackson's harmonic sieve [Jac94]). Feldman further refined this fact about DNFs

Theorem 3 (Theorem 3.8 in [Fel12]). *Let $c \in (0, 1]$ be a constant, μ be a c -bounded distribution and $\epsilon > 0$. For an integer $s > 0$ let f be an s -term DNF. For $d = \lceil \log(s/\epsilon) / \log(2/(2-c)) \rceil$ and every bounded function $g : \{-1, 1\}^n$,*

$$\mathbf{E}_\mu[|f(x) - g(x)|] \leq (2 \cdot (2-c)^{d/2} \cdot s + 1) \cdot \left\| \hat{f}_\mu(B_d) - \hat{g}_\mu(B_d) \right\|_\infty + 4\epsilon.$$

By this theorem the learnability of DNF reduces to constructing a g that approximates the heavy low-degree Fourier spectrum of f . This is exactly the approach followed by Feldman that we now proceeded to sketch.

The first step of the procedure is to run the EKM algorithm to estimate the heavy Fourier spectrum of f . The EKM algorithm returns a succinct representation of the spectrum and the learner selects only the coefficients that have degree $\geq d$. This is the only step of the algorithm that requires membership queries and is the subroutine that will be replaced by the quantum EKM algorithm that will be derived in Section 6.

Once the learner estimated the Fourier spectrum of f it proceeds with the construction of

g . The procedure is simple and based on an iterative process. Note that by Parseval

$$\mathbf{E}_\mu[(f - g)^2] = \sum_b (\hat{f}_\mu(b) - \hat{g}_\mu(b))^2 = \|\hat{f}_\mu - \hat{g}_\mu\|_2^2. \quad (1)$$

Starting with a g such that $|\hat{f}_\mu(a) - \hat{g}_\mu(a)| \geq \gamma$ it is possible to construct a g' such that g' is closer than g to f in l_2 norm with the following rule:

$$g' = g + (\hat{f}_\mu(a) - \hat{g}_\mu(a))\phi_{\mu,a}.$$

Then by Eq.1 we have that

$$\begin{aligned} \mathbf{E}_\mu[(f - g')^2] &= \sum_{b \neq a} (\hat{f}_\mu(b) - \hat{g}_\mu(b))^2 \\ &= \mathbf{E}_\mu[(f - g)^2] - (\hat{f}_\mu(a) - \hat{g}_\mu(a))^2 \\ &\leq \mathbf{E}_\mu[(f - g)^2] - \gamma^2. \end{aligned}$$

The problem with this procedure is that the function g' might have value outside $[-1, 1]$ but Feldman showed that the function can be adjusted to the right range by performing a single projection after all the updates.

Once has been reached a precision such that an application of Theorem 3 gives $\mathbf{E}_\mu[|f(x) - g(x)|] \leq \epsilon$, the algorithm outputs $\text{sign}(g)$ as hypothesis. From this, it follows the learnability of f

$$\Pr_\mu[f \neq \text{sign}(g)] \leq \mathbf{E}_\mu[|f - g|] \leq \epsilon.$$

The runtime of all the above operations is polynomial in n and inverse polynomial in the error parameters resulting in the following corollary

Corollary 4 (Corollary 5.1 in [Fel12]). *Let f compute an s -term DNF. Let $c \in (0, 1]$ be a constant and let $\mathcal{D}_\mu(x)$ be a c -bounded probability distribution. Let $\text{EX}(f, \mu)$ be an example oracle and $\text{MQ}(f)$ a membership oracle. Then, there exists an algorithm with $\text{EX}(f, \mu)$ and $\text{MQ}(f)$ access that efficiently PAC learns f over \mathcal{D}_μ .*

Finally, we note that the requirement of c -bounded distributions is imposed in order to control the magnitude of modulus of the μ -biased Fourier basis $\{|\phi_{\mu,a}(x)|\}$ that, otherwise, would diverge for $\mu \rightarrow \pm 1$.

5 Quantum μ -biased Fourier transform

In this section we introduce the μ -biased quantum Fourier transform and show how this can be used to derive a quantum algorithm for sampling from the probability distribution defined by the Fourier coefficients of the μ -biased transform. We recall that the μ -biased Fourier transform is defined as

$$f(x) = \sum_{a \in \{0,1\}^n} \hat{f}_\mu(a) \phi_{\mu,a}(x), \quad (2)$$

where $\phi_{\mu,a}(x) = \prod_{i:a_i=1} (x_i - \mu_i) / \sqrt{1 - \mu_i^2}$, $\hat{f}_\mu(a) = \mathbf{E}_\mu[f(x) \phi_{\mu,a}(x)]$, and $\mathcal{D}_\mu(x) = \prod_{i:x_i=1} (1 + \mu_i)/2 \prod_{i:x_i=-1} (1 - \mu_i)/2$. Our construction of the n -qubit μ -biased QFT exploits a fundamental property of product distributions, namely that the orthonormal basis $\{\phi_{\mu,a}(x)\}$ it defines can be factorised on the individual bits. This fact allows us to give an explicit form of the n -qubit transform as a tensor product of n single qubit transforms. We begin by constructing the single

qubit transform. Later we will show how to construct efficiently an n -qubit transform out of n single qubit ones. In the following we assume that the function $f : \{-1, 1\}^n \rightarrow \{-1, 1\}$ is a Boolean-valued function defined over the Boolean hypercube (as every DNF). Our results can be extended to real-valued functions over the Boolean hypercube using a discretisation procedure. As shown in [BJ98] the error induced by the approximation can be controlled.

The action of the single qubit μ -biased QFT can be explicitly constructed

$$H_\mu |x\rangle = \sum_{a \in \{0,1\}} \sqrt{\mathcal{D}_\mu(x)} \phi_{\mu,a}(x) |a\rangle. \quad (3)$$

Here we defined H_μ as the single qubit μ -biased QFT operator whose description in the computational basis is readily given by:

$$H_\mu = \begin{bmatrix} \sqrt{\mathcal{D}_\mu(x=-1)} \phi_0(x=-1) & \sqrt{\mathcal{D}_\mu(x=1)} \phi_0(x=1) \\ \sqrt{\mathcal{D}_\mu(x=-1)} \phi_1(x=-1) & \sqrt{\mathcal{D}_\mu(x=1)} \phi_1(x=1) \end{bmatrix}.$$

By taking the functional forms of $\mathcal{D}_\mu(x)$ and $\phi(x)$ we can write

$$H_\mu = \begin{bmatrix} \sqrt{\frac{1-\mu}{2}} & \sqrt{\frac{1+\mu}{2}} \\ -\frac{(1+\mu)\sqrt{1-\mu}}{\sqrt{2-2\mu^2}} & -\frac{(-1+\mu)\sqrt{1+\mu}}{\sqrt{2-2\mu^2}} \end{bmatrix}.$$

It is easy to verify that this matrix is unitary and positive semidefinite. We also note that, as consequence of the Solovay-Kitaev theorem [Kit97], it is possible to approximate H_μ from a fixed finite set of universal quantum gates with logarithmic overhead.

We can construct the extension of the μ -biased QFT to the case of n qubits by taking the tensor product of n single qubit operators. Let $x \in \{-1, 1\}^n$ and $a \in \{0, 1\}^n$, if we denote as a_i the i -th digit of a , $\mathcal{D}_{\mu_i}(x)$ as the probability associated to the i -th bit, and $\phi_{\mu,a_i}(x)$ its respective basis element, we can write:

$$H_\mu \otimes \cdots \otimes H_\mu |x\rangle = \sum_{a_1} \cdots \sum_{a_n} \prod_{i=1}^n \sqrt{\mathcal{D}_{\mu_i}(x)} \phi_{\mu_i,a_i}(x) |a_1\rangle \cdots |a_n\rangle.$$

By exploiting the product structure of $\mathcal{D}_\mu(x)$ and $\{\phi_{\mu,a}(x)\}$ that is, $\mathcal{D}_\mu(x) = \prod_i \mathcal{D}_{\mu_i}(x)$ and $\{\phi_{\mu,a}(x) = \prod_i \phi_{\mu_i,a_i}(x)\}$ we can write the n qubit μ -biased QFT as:

$$H_\mu^n |x\rangle = \sum_{a \in \{0,1\}^n} \sqrt{\mathcal{D}_\mu(x)} \phi_{\mu,a}(x) |a\rangle. \quad (4)$$

We remark that it is possible to construct the n qubit transform only because the product distribution and the $\{\phi_{\mu,a}(x)\}$ basis factorises. Without this factorisation we could still write Eq. 4 but we would not know how to implement this transformation efficiently on a quantum computer (the Solovay-Kitaev theorem guarantees that only single qubit unitary can be efficiently approximated by a universal set of gates).

Finally, we note that the construction of the μ -biased transform assumes knowledge of the vector μ . It is possible to estimate μ_i for each i using random samples from \mathcal{D}_μ . In Appendix A we prove that the error introduced by this approximation can be controlled if $\mathcal{D}_\mu(x)$ is c -bounded.

As a simple application of the μ -biased QFT we show how to sample from the probability distribution defined by the coefficients of the single bit μ -biased Fourier transform (recall that Parseval's equality holds in the μ -biased setting).

Lemma 2 (μ -biased quantum Fourier sampling). *Let $f : \{-1, 1\} \rightarrow \{-1, 1\}$ be a Boolean-valued function. Then, there exists a quantum algorithm with quantum membership oracle O_f access that returns $a \in \{0, 1\}$ with probability $\hat{f}^2(a)/2$. The algorithm requires exactly 1 O_f query and $O(1)$ gates.*

Proof. Starting with the $|0\rangle$ state, apply an Hadamard transform to get $\frac{1}{\sqrt{2}} \sum_x |x\rangle$. Then apply Lemma 1 to get $\frac{1}{\sqrt{2}} \sum_x \sqrt{\mathcal{D}_\mu(x)} |x\rangle$. By querying the quantum membership oracle O_f , one obtains $\frac{1}{\sqrt{2}} \sum_x \sqrt{\mathcal{D}_\mu(x)} f(x) |x\rangle$. Finally, applying the μ -biased QFT results in

$$\begin{aligned} \frac{1}{\sqrt{2}} \sum_x \sqrt{\mathcal{D}_\mu(x)} f(x) \left(\sum_a \sqrt{\mathcal{D}_\mu(x)} \phi_{\mu,a}(x) |a\rangle \right) &= \frac{1}{\sqrt{2}} \sum_{x,a} \mathcal{D}_\mu(x) f(x) \phi_{\mu,a}(x) |a\rangle \\ &= \frac{1}{\sqrt{2}} \sum_a \hat{f}_\mu(a) |a\rangle. \end{aligned}$$

Measuring the state, one obtains a with probability $\hat{f}_\mu^2(a)/2$. \square

In order to use this result in the context of quantum PAC learning we need to replace the membership oracle O_f with a quantum example oracle. The following lemma, that extends to the μ -biased case Lemma 1 in [BJ98], serves this purpose. Differently from Lemma 2 we present directly the n -dimensional case.

Lemma 3. *Let $f : \{-1, 1\}^n \rightarrow \{-1, 1\}$ be a Boolean-valued function. Then, there exists a quantum algorithm with quantum example oracle $\text{QEX}(f, \mu)$ access that returns $a \in \{0, 1\}^n$ with probability $\hat{f}^2(a)/2$. The algorithm requires exactly 1 QEX query and $O(n)$ gates.*

Proof. Let $f'(x) = (1 - f(x))/2$ be the truth table representation of $f(x)$ with $(-1)^{f'(x)} = f(x)$. Given access to $\text{QEX}(f, \mu)$ it is always possible to construct an oracle for $\text{QEX}(f', \mu)$ (this is equivalent to a relabelling of the qubits). Apply $\text{QEX}(f', \mu)$ on a $|0, \dots, 0\rangle$ to get $\sum_x \sqrt{\mathcal{D}_\mu(x)} |x, f'(x)\rangle$. Then apply H_μ^n on the first register:

$$\sum_{x \in \{-1, 1\}^n} \sum_{a \in \{0, 1\}^n} \sqrt{\mathcal{D}_\mu(x)} \sqrt{\mathcal{D}_\mu(x)} \phi_{\mu,a}(x) |a, f'(x)\rangle.$$

An application of the standard QFT on the second register gives:

$$\begin{aligned} \sum_{x,a,z} \frac{1}{\sqrt{2}} (-1)^{f'(x)z} \mathcal{D}_\mu(x) \phi_{\mu,a}(x) |a, z\rangle &= \frac{1}{\sqrt{2}} \left(\sum_a \hat{f}_\mu(a) |a, 1\rangle + \sum_a \mathbf{E}_\mu[\phi_{\mu,a}(x)] |a, 0\rangle \right) \\ &= \frac{1}{\sqrt{2}} \left(\sum_a \hat{f}_\mu(a) |a, 1\rangle + \sum_a \mathbf{E}_\mu[\phi_{\mu,a}(x) \phi_{\mu,0}(x)] |a, 0\rangle \right) \\ &= \frac{1}{\sqrt{2}} \left(\sum_a \hat{f}_\mu(a) |a, 1\rangle + |0 \dots 0, 0\rangle \right), \end{aligned}$$

where we used the orthonormality of the $\{\phi_{\mu,a}(x)\}$ basis and $\phi_{\mu,0}(x) = 1$. Measuring the first register we obtain $|a, 1\rangle$ with probability $\hat{f}_\mu^2(a)/2$ \square

6 Quantum computation of μ -biased Fourier spectrum

In this section we give a quantum algorithm to approximate the μ -biased Fourier spectrum of a function. This can be interpreted as a quantum version of the EKM algorithm. As a simple

application of the quantum EKM algorithm we obtain the learnability of DNFs under product distributions in the quantum PAC model.

We will make repeated use of the *Dvoretzky-Kiefer-Wolfowitz* (DKW) theorem, a concentration inequality that bounds the number of samples required to estimate a cumulative distribution in ℓ_∞ norm. The DKW Theorem was first proposed by Dvoretzky-Kiefer-Wolfowitz in 1956 with an almost tight bound [DKW56]. In 1958 Birnbaum and McCarty conjectured that the inequality was tight [BM58]. This conjecture was proved by Massart in 1990 [Mas+90]. The DKW theorem is usually given for continuous probability distribution but its validity extends also to discrete distributions (a detailed discussion can be found in [Kos07]).

Let X_1, \dots, X_m be a sequence of i.i.d. random variables drawn from a distribution $f(x)$ on $\{-1, 1\}^n$ with *Cumulative Distribution Function* (CDF) $F(x) = \sum_{X_i \leq x} f(X_i)$, and let x_1, \dots, x_n be their realizations. Given a set A the indicator function $\mathbf{1}_A : A \rightarrow \{0, 1\}$ takes values $f(x) = 0$ if $x \notin A$ and $f(x) = 1$ if $x \in A$. We denote the empirical probability distribution associated to $f(x)$ as $f_m(x) = \sum_{i=1}^m \mathbf{1}_{\{X_i=x\}}/m$ and its empirical cumulative distribution as $F_m(x) = \frac{1}{m} \sum_{i=1}^m \mathbf{1}_{\{X_i \leq x\}}$. The DKW theorem states that:

Theorem 5 (Dvoretzky-Kiefer-Wolfowitz). *For any i.i.d. sample X_1, \dots, X_m with cumulative distribution $F(x)$ and empirical cumulative distribution $F_m(x) = \frac{1}{m} \sum_{i=1}^m \mathbf{1}_{\{X_i \leq x\}}$*

$$\Pr \left(\max_{x \in \{-1, 1\}^n} |F(x) - F_m(x)| \geq \epsilon \right) \leq 2e^{-2m\epsilon^2}.$$

To make notation consistent, in the following we write $\|F(x) - F_m(x)\|_\infty$ instead of $\max_{x \in \{-1, 1\}^n} |F(x) - F_m(x)|$. By using the DKW theorem we can prove a useful lemma that bounds the number of samples needed to estimate a probability distribution in ℓ_∞ norm.

Lemma 4. *Let $f(x)$ be a probability distribution over $\{-1, 1\}^n$ and let $\tau > 0$, $\delta > 0$. Then, there exists an algorithm that with probability $1 - \delta$ outputs $f_m(x)$ such that $\|f(x) - f_m(x)\|_\infty \leq \tau$ using $m = O(\log(1/\delta)/\tau^2)$ samples.*

Proof. Let $\{e_1, \dots, e_{2^n}\}$ be an ordering of elements of the Boolean hypercube $\{1, 1\}^n$. We have that

$$\|f(x) - f_m(x)\|_\infty = \max_{\{e_1, \dots, e_{2^n}\}} |f(e_i) - f_m(e_i)| = \|F(e_{i+1}) - F_m(e_{i+1}) - (F(e_i) - F_m(e_i))\|_\infty.$$

An application of the triangle inequality gives

$$\begin{aligned} \|F(e_{i+1}) - F_m(e_{i+1}) - (F(e_i) - F_m(e_i))\|_\infty &\leq \|F(e_{i+1}) - F_m(e_{i+1})\|_\infty + \|(F(e_i) - F_m(e_i))\|_\infty \\ &\leq 2\|F(x) - F_m(x)\|. \end{aligned}$$

By Theorem 5 we have that, with probability $1 - \delta$, $\Pr(\|F(x) - F_m(x)\| \geq \gamma) \leq 2e^{-2m\gamma^2}$. Let $\gamma = \tau/2$, then

$$\Pr(\|f(x) - f_m(x)\|_\infty \leq \tau) \leq 1 - 2e^{-m\tau^2/2},$$

from which it is easy to that $m = O(\log(1/\delta)/\tau^2)$. \square

The combined application of Lemma 3 and Lemma 4 allows us to prove the following result:

Theorem 6 (Quantum EKM algorithm). *Let $f : \{-1, 1\} \rightarrow \{-1, 1\}$ be a Boolean-valued function and let $\epsilon > 0$, $\delta > 0$, $\mu \in (-1, 1)^n$. Then, there exists a quantum algorithm with $\text{QEX}(f, \mu)$ access that, with probability at least $1 - \delta$, returns a succinctly represented vector \hat{f}_μ , such that $\|\hat{f}_\mu - \tilde{f}_\mu\|_\infty \leq \epsilon$ and $\|\tilde{f}_\mu\|_0 \leq 4/\epsilon^2$. The algorithm requires $O(\log^2(1/\delta)/\epsilon^8)$ QEX queries and $O(n\log^2(1/\delta)/\epsilon^8)$ gates.*

Proof. We begin by estimating the a 's related to the heavy Fourier coefficients of f . Let $p(a) = |\hat{f}_\mu(a)|^2$ be the probability distribution defined by the μ -biased Fourier coefficients of f . Lemma 3 gives a procedure that, with 1 QEX query and $O(n)$ gates, samples $|a, 0\rangle$ with probability $q(a, 0) = |\hat{f}_\mu(a)|^2/2$ and $|0 \dots 0, 0\rangle$ with probability $q(0, 1) = 1/2$. Applying Lemma 4 on the distribution q with $\tau = \epsilon^2/8$ we obtain that $O(\log(1/\delta)\epsilon^4)$ samples are required to have an estimate $\|q(a, i) - \tilde{q}(a, i)\|_\infty \leq \epsilon^2/8$ in high probability. This implies that $\|\hat{f}_\mu^2(a) - \tilde{f}_\mu^2(a)\|_\infty \leq \epsilon^2/4$. By selecting the characteristic vectors that correspond to coefficients such that $|\tilde{f}_\mu(a)|^2 > \epsilon^2/2$ (and discarding the element $|a, 0\rangle$) we can output a list of a 's such that, with probability $\geq 1 - \delta$, all the corresponding Fourier coefficients have $|\hat{f}_\mu(a)| > \epsilon$ and there are no coefficients such that $|\hat{f}_\mu(a)| \leq \epsilon/2$. By Parseval's equality this implies that the list may contain at most $4/\epsilon^2$ elements.

The final step requires the estimation of the Fourier coefficients. For a given a , the Fourier coefficient $\hat{f}_\mu(a) = \mathbf{E}_\mu[f(x)\chi_a(x)]$ can be obtained by sampling using the QEX oracle to simulate EX (to get an example $(x, f(x))$ it would suffice to measure a state prepared with QEX) in time $O(\log(1/\delta)/\epsilon^2)$ (the number of examples required for the estimate is a standard application of the Chernoff bound).

Combining the bounds to estimate the a 's to the one to estimate the Fourier coefficients we obtain that the algorithm requires $O(\log^2(1/\delta)/\epsilon^8)$ QEX queries and $O(n\log^2(1/\delta)/\epsilon^8)$ gates (each estimate of the $\hat{f}(a)$ must be repeated $O(4/\epsilon^2)$ times). \square

Theorem 6 can be straightforwardly used in the method developed by Feldman (Corollary 5.1 in [Fel12]) to obtain the learnability of DNF under product.

Corollary 7. *Let f compute an s -term DNF. Let $c \in (0, 1]$ be a constant, let $\mathcal{D}_\mu(x)$ be a c -bounded probability distribution and let $\text{QEX}(f, \mu)$ be a quantum example oracle. Then, there exists a quantum algorithm with $\text{QEX}(f, \mu)$ access that efficiently PAC learns f over \mathcal{D}_μ .*

We recall that the collection of the heavy Fourier coefficients of the DNF f is the only step of Feldman's algorithm that requires MQ. The remaining of the algorithm makes use of the coefficients to construct a function g that approximates f .

7 Construction of quantum example oracles

A large class of quantum algorithms for learning problems involving classical data or functions require quantum data oracles that can efficiently update the training set in superposition. Here, we define efficient as logarithmic in the dimension of the training set or in the dimension of the support of the function. The QEX oracle is one of such oracles and can be implemented using the *Quantum Random Access Memory* (QRAM) [GLM08]. The QRAM is a quantum procedure that allows one to encode in superposition N data points into $\log(N)$ qubits in time $O(\log(N))$. More specifically, let $\{m_j\}$ be the content of a memory structure with N elements. The action of the QRAM on a state $\sum_i \alpha_i |i, 0\rangle$ is

$$\text{QRAM} \sum_i \alpha_i |i, 0\rangle \rightarrow \sum_i \alpha_i |i, m_j\rangle.$$

Let f be a Boolean-valued function $f : \{-1, 1\}^n \rightarrow \{-1, 1\}$. The action of a QEX oracle for a target concept f and probability distribution $\mathcal{D}(x)$ is

$$\text{QEX}(f, \mathcal{D}) |0 \dots 0, 0\rangle \rightarrow \sum_x \sqrt{\mathcal{D}(x)} |x, f(x)\rangle.$$

In order to use the QRAM to update classical data we need a classical memory structure that stores the truth table of f . Provided this classical memory, we can construct the QEX oracle in two steps. Starting from the state $|0 \dots 0, 0\rangle$, use Lemma 1 on the first register to obtain $\sum_x \sqrt{\mathcal{D}(x)} |x, 0\rangle$. Because $\mathcal{D}(x)$ is a discrete probability distribution, it is always efficiently integrable and thus it is always possible to apply Lemma 1. Finally, an application of the QRAM returns the state $\sum_x \sqrt{\mathcal{D}(x)} |x, f(x)\rangle$. Note that because the QRAM can load data in logarithmic time, it is possible to build a superposition that encodes a Boolean function supported on $\{-1, 1\}^n$ in polynomial time. In this sense, the total runtime of the algorithm for learning DNFs remains polynomial.

A possible drawback of a QRAM implementation is that it is not clear whether a QRAM can actually be built. A discussion of the challenges related to the construction of a QRAM can be found in the review article by Ciliberto *et al.* [Cil+18].

The QRAM can be substituted with a standard quantum membership oracle in the procedure above. Although the use of the membership oracle would be limited to the construction of the QEX oracle (which Bshouty and Jackson proved to be unable to simulate efficiently a classical membership oracle) it is still unclear whether a QEX oracle can be built without using a membership oracle.

Acknowledgements

We thank Carlo Ciliberto for elucidative comments on the DKW inequality and Leonard Wossnig for helpful conversations on the implementation of the QEX oracle. AR is supported by an EPSRC DTP Scholarship and by QinetiQ Ltd. SS is supported by the Royal Society, EPSRC, the National Natural Science Foundation of China, and the grant ARO-MURI W911NF-17-1-0304 (US DOD, UK MOD and UK EPSRC under the Multidisciplinary University Research Initiative).

References

- [AC02] M. Adcock and R. Cleve. “A quantum Goldreich-Levin theorem with cryptographic applications”. In: *Annual Symposium on Theoretical Aspects of Computer Science*. Springer. 2002, pp. 323–334.
- [AW17a] S. Arunachalam and R. de Wolf. “Guest Column: A Survey of Quantum Learning Theory”. In: *SIGACT News* 48.2 (2017), pp. 41–67.
- [AW17b] S. Arunachalam and R. de Wolf. “Optimal Quantum Sample Complexity of Learning Algorithms”. In: *32nd Computational Complexity Conference, CCC 2017, July 6-9, 2017, Riga, Latvia*. 2017.
- [AS07] A. Atılcı and R. A. Servedio. “Quantum algorithms for learning and testing juntas”. In: *Quantum Information Processing* 6.5 (2007), pp. 323–348.
- [AFK13] P. Awasthi, V. Feldman, and V. Kanade. “Learning using local membership queries”. In: *Conference on Learning Theory*. 2013, pp. 398–431.
- [Bah61] R. R. Bahadur. “A representation of the joint distribution of responses to n dichotomous items”. In: *Studies in item analysis and prediction* 6 (1961), pp. 158–168.
- [Bel91] M. Bellare. *The Spectral Norm of Finite Functions*. Tech. rep. Cambridge, MA, USA, 1991.
- [BV97] E. Bernstein and U. Vazirani. “Quantum complexity theory”. In: *SIAM Journal on Computing* 26.5 (1997), pp. 1411–1473.

- [BM58] Z. Birnbaum and R. McCarty. “A Distribution-Free Upper Confidence Bound for $\Pr\{Y < X\}$, Based on Independent Samples of X and Y ”. In: *The Annals of Mathematical Statistics* (1958), pp. 558–562.
- [BKW03] A. Blum, A. Kalai, and H. Wasserman. “Noise-tolerant learning, the parity problem, and the statistical query model”. In: *Journal of the ACM (JACM)* 50.4 (2003), pp. 506–519.
- [BJ98] N. H. Bshouty and J. C. Jackson. “Learning DNF over the uniform distribution using a quantum example oracle”. In: *SIAM Journal on Computing* 28.3 (1998), pp. 1136–1153.
- [Bsh+05] N. H. Bshouty et al. “Learning DNF from random walks”. In: *Journal of Computer and System Sciences* 71.3 (2005), pp. 250–265.
- [Cil+18] C. Ciliberto et al. “Quantum machine learning: a classical perspective”. In: *Proc. R. Soc. A*. Vol. 474. 2209. The Royal Society. 2018, p. 20170551.
- [CSS15] A. W. Cross, G. Smith, and J. A. Smolin. “Quantum learning robust against noise”. In: *Physical Review A* 92.1 (2015), p. 012327.
- [DKW56] A. Dvoretzky, J. Kiefer, and J. Wolfowitz. “Asymptotic minimax character of the sample distribution function and of the classical multinomial estimator”. In: *The Annals of Mathematical Statistics* (1956), pp. 642–669.
- [Fel12] V. Feldman. “Learning DNF Expressions from Fourier Spectrum.” In: *COLT*. Vol. 8. 8.4. 2012, pp. 8–4.
- [Fre95] Y. Freund. “Boosting a weak learning algorithm by majority”. In: *Information and computation* 121.2 (1995), pp. 256–285.
- [FJS91] M. L. Furst, J. C. Jackson, and S. W. Smith. “Improved learning of AC 0 functions”. In: *COLT*. Vol. 91. 1991, pp. 317–325.
- [GLM08] V. Giovannetti, S. Lloyd, and L. Maccone. “Quantum random access memory”. In: *Physical review letters* 100.16 (2008), p. 160501.
- [GL89] O. Goldreich and L. A. Levin. “A hard-core predicate for all one-way functions”. In: *Proceedings of the twenty-first annual ACM symposium on Theory of computing*. ACM. 1989, pp. 25–32.
- [GK17] A. B. Grilo and I. Kerenidis. “Learning with Errors is easy with quantum samples”. In: *arXiv preprint arXiv:1702.08255* (2017).
- [GR02] L. Grover and T. Rudolph. “Creating superpositions that correspond to efficiently integrable probability distributions”. In: *arXiv preprint quant-ph/0208112* (2002).
- [Jac94] J. Jackson. “An efficient membership-query algorithm for learning DNF with respect to the uniform distribution”. In: *Foundations of Computer Science, 1994 Proceedings., 35th Annual Symposium on*. IEEE. 1994, pp. 42–53.
- [Jac97] J. C. Jackson. “An efficient membership-query algorithm for learning DNF with respect to the uniform distribution”. In: *Journal of Computer and System Sciences* 55.3 (1997), pp. 414–440.
- [JTY02] J. C. Jackson, C. Tamon, and T. Yamakami. “Quantum DNF learnability revisited”. In: *Lecture notes in computer science* (2002), pp. 595–604.
- [KST09] A. T. Kalai, A. Samorodnitsky, and S.-H. Teng. “Learning and smoothed analysis”. In: *Foundations of Computer Science, 2009. FOCS’09. 50th Annual IEEE Symposium on*. IEEE. 2009, pp. 395–404.

- [Kit97] A. Y. Kitaev. “Quantum computations: algorithms and error correction”. In: *Russian Mathematical Surveys* 52.6 (1997), pp. 1191–1249.
- [KS01] A. R. Klivans and R. Servedio. “Learning DNF in time $2^{\tilde{O}(n^{1/3})}$ ”. In: *Proceedings of the thirty-third annual ACM symposium on Theory of computing*. ACM, 2001, pp. 258–265.
- [Kos07] M. R. Kosorok. *Introduction to empirical processes and semiparametric inference*. Springer Science & Business Media, 2007.
- [KM93] E. Kushilevitz and Y. Mansour. “Learning decision trees using the Fourier spectrum”. In: *SIAM Journal on Computing* 22.6 (1993), pp. 1331–1348.
- [Lyu05] V. Lyubashevsky. “The parity problem in the presence of noise, decoding random linear codes, and the subset sum problem”. In: *Approximation, Randomization and Combinatorial Optimization. Algorithms and Techniques*. Springer, 2005, pp. 378–389.
- [Mas+90] P. Massart et al. “The tight constant in the Dvoretzky-Kiefer-Wolfowitz inequality”. In: *The annals of Probability* 18.3 (1990), pp. 1269–1283.
- [MO10] A. Montanaro and T. J. Osborne. “Quantum Boolean Functions”. In: *Chicago Journal of Theoretical Computer Science* 2010.1 (Jan. 2010).
- [ODo14] R. O’Donnell. *Analysis of boolean functions*. Cambridge University Press, 2014.
- [Val84] L. G. Valiant. “A theory of the learnable”. In: *Communications of the ACM* 27.11 (1984), pp. 1134–1142.
- [Ver90] K. A. Verbeurgt. “Learning DNF Under the Uniform Distribution in Quasi-Polynomial Time.” In: *COLT*. 1990, pp. 314–326.

A Error analysis

In the main section we assumed that the vector μ parametrising the product distribution was given to the learner. Here we prove that, if $\mathcal{D}_\mu(x)$ is c -bounded, it is possible to estimate μ introducing an error that can be made small at a cost that scales polynomially in n . We recall that $\mu \in [-1 + c, 1 - c]^n$, $c \in (0, 1]$, and $\mu_i = \mathbf{E}_\mu[x_i]$. A simple application of the Chernoff bound gives that, with probability $1 - \delta$, $O(\log(1/\delta)/\epsilon^2)$ samples give an estimate $\tilde{\mu}$ such that $|\mu - \tilde{\mu}| \leq \epsilon$.

We want to estimate the error introduced by approximating H_μ^n with $H_{\tilde{\mu}}^n$ (note that the μ -biased QFT is now parametrised by $\tilde{\mu}$) in terms of the operator norm. Let A be an operator, the operators norm $\|A\|$ is defined as:

$$\|A\| = \sup_{|\psi\rangle \neq 0} \frac{\|A|\psi\rangle\|}{\| |\psi\rangle \|}.$$

The error analysis then requires to bound the quantity:

$$\|H_\mu^n - H_{\tilde{\mu}}^n\| \leq \gamma.$$

In order to prove the bound we introduce a useful lemma:

Lemma 5. *Let $A = A_n \cdots A_1$ be a product of unitary operators A_j . Assume that for every A_j there exists an approximation \tilde{A}_j such that $\|A - \tilde{A}_j\| \leq \epsilon_j$. The follow inequality holds*

$$\|A_n \cdots A_1 - \tilde{A}_n \cdots \tilde{A}_1\| \leq \sum_j \epsilon_j.$$

Proof. We prove by induction. The base step follows from the assumptions. For the inductive step let $X_k = A_k \cdots A_1$ and $\tilde{X}_k = \tilde{A}_k \cdots \tilde{A}_1$. Because the inductive hypothesis holds we have

$$\|X_k - \tilde{X}_k\| \leq \sum_{j=1}^k \epsilon_j.$$

By making use of the triangular inequality, the induction hypothesis, and noting that the product of unitaries is unitary we have

$$\begin{aligned} \|A_{k+1}X_k - \tilde{A}_{k+1}\tilde{X}_k\| &= \|A_{k+1}(X_k - \tilde{X}_k) + (A_{k+1} - \tilde{A}_{k+1})\tilde{X}_k\| \\ &\leq \|A_{k+1}(X_k - \tilde{X}_k)\| + \|(A_{k+1} - \tilde{A}_{k+1})\tilde{X}_k\| \\ &= \|A_{k+1}\| \|X_k - \tilde{X}_k\| + \|A_{k+1} - \tilde{A}_{k+1}\| \|\tilde{X}_k\| \\ &\leq \sum_{j=1}^k \epsilon_j + \epsilon_{k+1} \\ &= \sum_{j=1}^{k+1} \epsilon_j. \end{aligned}$$

□

Let $H_j = I \otimes \cdots I \otimes H_{\mu_j} \otimes I \otimes \cdots I$ and $\tilde{H}_j = I \otimes \cdots I \otimes H_{\tilde{\mu}_j} \otimes I \otimes \cdots I$. By Lemma 5 we have that

$$\|H_\mu^n - H_{\tilde{\mu}}^n\| \leq \sum_{i=1}^n \|H_i - \tilde{H}_i\|. \quad (5)$$

The bound on $\|H_j - \tilde{H}_j\|$ can be simplified using the following property of the operator norm $\|A \otimes B\| = \|A\| \|B\|$,

$$\begin{aligned} \|H_j - \tilde{H}_j\| &= \|I \otimes \cdots I \otimes (H_{\mu_j} - H_{\tilde{\mu}_j}) \otimes I \otimes \cdots I\| \\ &= \|I\| \cdots \|I\| \|H_{\mu_j} - H_{\tilde{\mu}_j}\| \|I\| \cdots \|I\| \\ &= \|H_{\mu_j} - H_{\tilde{\mu}_j}\|. \end{aligned}$$

The problem of bounding Eq. 5 is then equivalent to bounding $\|H_{\mu_i} - H_{\tilde{\mu}_i}\|$. Let $|\psi\rangle = \sum_{x \in \{-1,1\}} \alpha_x |x\rangle$, we have that

$$\|(H_{\mu_i} - H_{\tilde{\mu}_i})|\psi\rangle\| = \left\| \sum_{x \in \{-1,1\}} \sum_{a \in \{0,1\}} \left(\sqrt{\mathcal{D}_{\mu_i}(x)} \phi_{\mu_i,a}(x) - \sqrt{\mathcal{D}_{\tilde{\mu}_i}(x)} \phi_{\tilde{\mu}_i,a}(x) \right) \alpha_x |a\rangle \right\|$$

where $\phi_{\mu,a}(x) = \prod_{i:a_i=1} (x_i - \mu_i) / \sqrt{1 - \mu_i^2}$ and $\mathcal{D}_\mu(x) = \prod_{i:x_i=1} (1 + \mu_i) / 2 \prod_{i:x_i=-1} (1 - \mu_i) / 2$. We have to estimate the following quantity for a generic a, x

$$\begin{aligned} S &= \left| \sqrt{\mathcal{D}_{\mu_i}(x)} \phi_{\mu_i,a}(x) - \sqrt{\mathcal{D}_{\tilde{\mu}_i}(x)} \phi_{\tilde{\mu}_i,a}(x) \right| \\ &= \left| \frac{(x_i - \mu_i) \sqrt{1 - \tilde{\mu}_i^2} \sqrt{\mathcal{D}_{\mu_i}(x)} - (x_i - \tilde{\mu}_i) \sqrt{1 - \mu_i^2} \sqrt{\mathcal{D}_{\tilde{\mu}_i}(x)}}{\sqrt{1 - \mu_i^2} \sqrt{1 - \tilde{\mu}_i^2}} \right|. \end{aligned}$$

Recall that for every i it holds $1 - \mu_i^2 \geq c^2$, $1 - \tilde{\mu}_i^2 \geq c^2$, $|\mu_i - \tilde{\mu}_i| \leq \epsilon$, $x_i \in \{-1, 1\}$. By the triangle inequality we have that

$$\begin{aligned}
S &\leq \frac{1}{c^2} \left| (x_i - \mu_i) \sqrt{1 - \tilde{\mu}_i^2} \sqrt{\mathcal{D}_{\mu_i}(x)} - (x_i - \tilde{\mu}_i) \sqrt{1 - \mu_i^2} \sqrt{\mathcal{D}_{\tilde{\mu}_i}(x)} \right| \\
&= \frac{1}{c^2} \left| (x_i - \mu_i) \left(\sqrt{1 - \tilde{\mu}_i^2} \sqrt{\mathcal{D}_{\mu_i}(x)} - \sqrt{1 - \mu_i^2} \sqrt{\mathcal{D}_{\tilde{\mu}_i}(x)} \right) + (\tilde{\mu}_i - \mu_i) \sqrt{1 - \mu_i^2} \sqrt{\mathcal{D}_{\tilde{\mu}_i}(x)} \right| \\
&\leq \frac{1}{c^2} \left(\left| (x_i - \mu_i) \left(\sqrt{1 - \tilde{\mu}_i^2} \sqrt{\mathcal{D}_{\mu_i}(x)} - \sqrt{1 - \mu_i^2} \sqrt{\mathcal{D}_{\tilde{\mu}_i}(x)} \right) \right| + \left| (\tilde{\mu}_i - \mu_i) \sqrt{1 - \mu_i^2} \sqrt{\mathcal{D}_{\tilde{\mu}_i}(x)} \right| \right) \\
&\leq \frac{1}{c^2} \left((2 - c) \left| \left(\sqrt{1 - \tilde{\mu}_i^2} \sqrt{\mathcal{D}_{\mu_i}(x)} - \sqrt{1 - \mu_i^2} \sqrt{\mathcal{D}_{\tilde{\mu}_i}(x)} \right) \right| + \epsilon \right) \\
&\leq \frac{1}{c^2} \left((2 - c) \left(\left| \sqrt{\mathcal{D}_{\mu_i}(x)} - \sqrt{\mathcal{D}_{\tilde{\mu}_i}(x)} \right| + \left| \sqrt{1 - \mu_i^2} - \sqrt{1 - \tilde{\mu}_i^2} \right| \right) + \epsilon \right).
\end{aligned}$$

If we note that

$$\left| \sqrt{\mathcal{D}_{\mu_i}(x)} - \sqrt{\mathcal{D}_{\tilde{\mu}_i}(x)} \right| = \left| \frac{\tilde{\mu} - \mu}{2(\sqrt{\mathcal{D}_{\mu_i}(x)} + \sqrt{\mathcal{D}_{\tilde{\mu}_i}(x)})} \right| \leq \frac{\epsilon}{\sqrt{8c}}$$

and

$$\left| \sqrt{1 - \mu_i^2} - \sqrt{1 - \tilde{\mu}_i^2} \right| = \left| \frac{\tilde{\mu}^2 - \mu^2}{\sqrt{1 - \mu_i^2} + \sqrt{1 - \tilde{\mu}_i^2}} \right| \leq \frac{\epsilon}{2c}$$

we have

$$S \leq t\epsilon \tag{6}$$

where $t = ((2 - c)(\frac{1}{\sqrt{8c}} + \frac{1}{2c}) + 1)/c^2$. By making use of Eq. 6 and noting that $\sum_x |\alpha_x| \leq 1$ we have

$$\begin{aligned}
\|H_\mu - H_{\tilde{\mu}}\| &\leq \sum_{x,a} \left\| \left(\sqrt{\mathcal{D}_{\mu_i}(x)} \phi_{\mu_i,a}(x) - \sqrt{\mathcal{D}_{\tilde{\mu}_i}(x)} \phi_{\tilde{\mu}_i,a}(x) \right) \alpha_x |a\rangle \right\| \\
&\leq t\epsilon \left(2 \sum_x |\alpha_x| \right) \\
&\leq 2t\epsilon.
\end{aligned}$$

From which it follows that

$$\|H_\mu^n - H_{\tilde{\mu}}^n\| \leq 2nt\epsilon. \tag{7}$$

Eq. 7 guarantees that if one can obtain an approximation μ with linear precision it is possible to control the approximation error. An approximation of ϵ in linear precision can be obtained using a polynomial, in n , number of examples $m = O(\log^2(1/\delta)/n^2 t^2)$.