# Sample Efficient Algorithms for Learning Quantum Channels in PAC Model and the Approximate State Discrimination Problem

## Kai-Min Chung ✉ 🄾
Institute of Information Science, Academia Sinica, Taipei, Taiwan

## Han-Hsuan Lin ✉
Department of Computer Science, National Tsing Hua University, Hsinchu, Taiwan

### ——— Abstract ———

The probably approximately correct (PAC) model [30] is a well studied model in classical learning theory. Here, we generalize the PAC model from concepts of Boolean functions to quantum channels, introducing *PAC model for learning quantum channels*, and give two sample efficient algorithms that are analogous to the classical "Occam's razor" result [12]. The classical Occam's razor algorithm is done trivially by excluding any concepts not compatible with the input-output pairs one gets, but such an approach is not immediately possible with a concept class of quantum channels, because the outputs are unknown quantum states from the quantum channel.

To study the quantum state learning problem associated with PAC learning quantum channels, we focus on the special case where the channels all have constant output. In this special case, learning the channels reduce to a problem of learning quantum states that is similar to the well known quantum state discrimination problem [8], but with the extra twist that we allow $\epsilon$-trace-distance-error in the output. We call this problem *Approximate State Discrimination*, which we believe is a natural problem that is of independent interest.

We give two algorithms for learning quantum channels in PAC model. The first algorithm has sample complexity

$$O\left(\frac{\log|C| + \log(1/\delta)}{\epsilon^2}\right),$$

but only works when the outputs are pure states, where $C$ is the concept class, $\epsilon$ is the error of the output, and $\delta$ is the probability of failure of the algorithm. The second algorithm has sample complexity

$$O\left(\frac{\log^3|C|(\log|C| + \log(1/\delta))}{\epsilon^2}\right),$$

and work for mixed state outputs. Some implications of our results are that we can PAC-learn a polynomial sized quantum circuit in polynomial samples, and approximate state discrimination can be solved in polynomial samples even when the size of the input set is exponential in the number of qubits, exponentially better than a naive state tomography.

## 1    Introduction

In computational learning theory, the Probably Approximately Correct (PAC) model of Valiant [30] gives a complexity-theoretic foundation of what it means for a concept class to be (efficiently) learnable. In the most basic setting of PAC learning model, we want to learn a set of Boolean functions, $C = \{c : \{0,1\}^n \to \{0,1\}\}$, called the concept class. The goal of a learning algorithm $A$ is to guess the identity of an unknown target concept $c^* \in C$ from samples $\{(x_1, c^*(x_1)), (x_2, c^*(x_2)), \dots\}$, where $\{x_1, x_2, \dots\}$ are inputs randomly drawn from a distribution $D$ that is unknown to $A$. Specifically, with error parameters $\epsilon$ and $\delta$, for all concept $c^* \in C$ and probability distribution $D$, $A$ is required to, given access to the samples $\{(x_1, c^*(x_1)), (x_2, c^*(x_2)), \dots\}$, with probability $1 - \delta$, come up with a hypothesis $h \in C$ that is $\epsilon$-close to $c^*$, i.e. $\Pr_{x \leftarrow D}[c(x) \neq h(x)] \leq \epsilon$. Such a learning algorithm is called a proper[1] $(\epsilon, \delta)$-PAC learner for the concept class $C$. Of course, we would like the learner $A$ to be as efficient as possible in terms of both sample complexity (i.e., the number of samples $A$ needs to access) and time complexity, and ideally, polynomial in the input length $n$ and the error parameters $\epsilon^{-1}$ and $\log(1/\delta)$. Since its introduction in the 80's by Valiant, PAC learning theory has been deeply studied to characterize when efficient learning is or is not possible.

Following Valient's PAC learning model on Boolean functions, generalization to different kinds of concept classes has been proposed, including Boolean functions on continuous spaces [13], probabilistic Boolean functions [20, 2], functions with $\{0, \dots, n\}$ outputs [26, 11], and real valued functions [10].

With quantum computers coming closer and closer into reality, it is natural to generalize the PAC learning model to quantum channels, capturing the learnability of quantum circuits or devices that we might build in the near future. Note that quantum states has an inherent "unlearnability", as manifested by the no-cloning theorem and uncertainty principle. Therefore this study of learnability of quantum channels has an interesting interaction between classical learning theory and quantum information theory.

Formally, we define the *PAC learning model for quantum channels* as follows: Let the *concept class $C$* be a finite set of known $d_1$ to $d_2$ dimensional quantum channels. We are trying to learn an unknown quantum channel, the *target concept $c^* \in C$*. In order to do this, we are given *samples* $\{(x_1, c^*(x_1)), (x_2, c^*(x_2)), \dots\}$, where $\{x_1, x_2, \dots\}$ are classical descriptions of the input quantum states to the quantum channel $c^*$ and $\{c^*(x_1), c^*(x_2), \dots\}$ are the corresponding quantum states outputted by $c^*$. The inputs are drawn from a distribution $D$ unknown to the learner. Because of the no-cloning theorem, it is hard to justify holding both the inputs and outputs as unknown quantum states, so we assume that we have full classical description of the input state and keep the outputted states as unknown quantum states, meaning that we hold a copy of the quantum state $c^*(x_i)$ rather than the full classical description of it. A proper $(\epsilon, \delta)$-PAC learner for the concept class $C$ of quantum channels is a quantum algorithm that for all concepts $c^* \in C$ and distribution $D$, takes the description of $C$ and $T$ samples $\{(x_1, c^*(x_1)), (x_2, c^*(x_2)), \dots (x_T, c^*(x_T))\}$ as input[2] and with probability $1 - \delta$, outputs a *hypothesis $h \in C$* that is $\epsilon$-close to the target concept $c^*$, where the distance between two concepts $h, c^*$ depends on the input distribution $D$ and is defined as $\Delta(h, c^*) = \mathbb{E}_{x \in D}[\Delta_{tr}(h(x), c^*(x))]$, i.e. the expected trace distance between the outputs averaged over $D$.

---

[1] Proper means that the hypothesis $h$ must be inside the concept class $C$, whereas an improper learner can output any $h$ as the hypothesis. All learners in this paper are proper, and we sometimes omit the term "proper".

[2] Note that $D$ is not part of the input and is unknown to the learner.

We gave two algorithms for learning quantum channels in PAC model that in a sense generalize the classical Occam's razor algorithm [12]. In particular, our algorithms have poly log sample complexity in the size of the concept class. The first algorithm has sample complexity

$$O\left(\frac{\log |C| + \log(1/\delta)}{\epsilon^2}\right),$$

but requires the outputs to be pure states. The second algorithm has sample complexity

$$O\left(\frac{\log^3 |C|(\log |C| + \log(1/\delta))}{\epsilon^2}\right),$$

while outputs can be mixed.

The Occam's razor algorithm [12] is a classical PAC learner for any finite sized concept class $C$ with sample complexity $O(\log |C|)$. The idea of the algorithm is simple: keep taking samples, check which concepts in the concept class do not agree with the samples and exclude them. One can show that every time a sample is taken, a constant fraction of the concepts that are $\epsilon$-far away from the target concept will be excluded, so an $\epsilon$-close hypothesis can be found in $O(\log |C|)$ samples.

Although the Occam's razor algorithm is simple, generalizing it to our PAC model for quantum channels is troublesome. The main difference is that when learning quantum channels, the outputs from the target concept are copies of unknown (possibly high dimensional) quantum states. By the nature of quantum mechanics, if we just have a few copies of a high dimensional quantum state, we can only learn a tiny fraction of information contained in the quantum state. Since we don't really know what the outputted state is, we cannot simply "exclude all channels that do not output this state." Instead, we need to carefully design the measurement we take on the outputted states, getting the information useful in distinguishing the quantum channels in our concept class. Note that the sample complexities of both of our algorithms do not depend on the dimension of the outputted states.

As a possible application of our result, our algorithms for learning quantum channels in PAC model can be viewed as a sample-efficient way to do quantum process tomography [23] when we know that the target quantum processes comes from a finite set and only care about being correct on average over an input distribution. For example, if we try to PAC-learn a polynomial sized quantum circuit of $n$-qubits, since there are only $2^{\mathrm{poly}(n)}$ possible polynomial sized circuits, our result shows that we can learn it in $\mathrm{poly}(n)$ samples, an exponential improvement over a naive process tomography that has no restriction on concept class size and inputs.

Note that this work studies the sample complexity instead of time complexity of learning. Just like various other cases in theoretical computer science where the oracle-based complexity does not match the time complexity of a problem, sample complexity and time complexity of learning quantum channels in PAC model is unlikely to match. In particular, Arunachalam et al. [5] showed that there is no polynomial time algorithm for learning $\mathrm{TC}^0$ or $\mathrm{AC}^0$ circuit even knowing $D$ is uniform unless LWE can be solved in polynomial time by a quantum computer.

## 1.1 Approximate State Discrimination

As stated previously, the most challenging part of our algorithms is how to extract information from unknown outputted quantum states to distinguish the channels. We isolate and study this problem by focusing on the special case where the channels are "constant," i.e. every

channel in the concept class outputs a fixed quantum state irrespective of the input[3]. Since the input does not matter, we don't need to write it down anymore, so the samples are just copies of the fixed unknown quantum state, and since a concept is fully specified by its unique output state, we might as well describe the concept class as a set of quantum states. In this special case, learning quantum channels in PAC model becomes an interesting hybrid of quantum state discrimination [6, 25, 24, 8, 29] and quantum state tomography [15, 27], and we named it the *approximate state discrimination* problem. The approximate state discrimination problem is formalized as follows: Let $S$ be a known finite set of $d$-dimensional density matrices. We want to learn an unknown target state $\sigma \in S$ using as few identical copies of $\sigma$ as possible. A quantum algorithm is an $(\epsilon, \delta)$-approximate discriminator of $S$ if, for all $\sigma \in S$, it takes the description of $S$ and $T$ copies of $\sigma$ as input and with probability $1 - \delta$ outputs a state $\rho \in S$ with $\Delta_{tr}(\rho, \sigma) \leq \epsilon$. This problem is called approximate state discrimination because it is the same as the state discrimination problem except that $\epsilon$-approximate answers are allowed.

Since approximate state discrimination is a special case of PAC learning quantum channels[4], it can also be solved with

$$O\left(\frac{\log |S| + \log(1/\delta)}{\epsilon^2}\right)$$

samples if $S$ consists of pure states and

$$O\left(\frac{\log^3 |S|(\log |S| + \log(1/\delta))}{\epsilon^2}\right)$$

samples if $S$ consists of mixed states.

## 1.2    Related Works and Independent Work

There are several works in the literature that study the sample complexity of PAC learning with different ways of generalization to quantum information. Cheng, Hsieh, and Yeh [14] studies the sample complexity of PAC learning arbitrary two outcome measurements, where the inputs are quantum states, and the learner has complete classical description of them. They show an upper of sample complexity linear in the dimension of the Hilbert space. Note that one can trivially get a lower bound of similar order by noticing that Boolean functions is a subset of two outcome measurements. Arunachalam and de Wolf [4] studies the sample complexity of PAC learning classical functions with quantum samples and shows that there is no quantum speed up. See [3] for a survey of quantum learning theory.

### 1.2.0.1   Independent Work

Independent to our work, in [7], Bădescu and O'Donnell formulate the problem of *quantum hypothesis selection*. Quantum hypothesis selection can be viewed as a generalization of our approximate state discrimination problem where the unknown state $\sigma$ might not be in the hypothesis set $|S|$, and the learner what to find the state in $|S|$ that is closest to the

---

[3] The outputs of different concepts are still different.
[4] We choose not to write up stand-alone algorithms for the approximate state discrimination problem as it will be very similar to that of PAC learning quantum channels. However, the reader can read the analysis of our algorithms with constant output assumptions to easily get the intuition behind them.

unknown state $\sigma$ (see Theorem 1.5 of [7] for the formal definition). This is similar to the agnostic learning model [18, 21]. Let $\eta$ be the minimum distance from the unknown state to something in $|S|$, Bădescu and O'Donnell give an algorithm that finds some $\rho \in S$ such that $\Delta_{tr}(\rho, \sigma) \leq 3.01\eta + \epsilon$ using $O\left(\frac{\log^3 |S| + \log(1/\delta)}{\epsilon^2}\right)$ samples. Since quantum hypothesis selection is a generalization of approximate state discrimination, Bădescu and O'Donnell's algorithm supersedes our algorithm for approximate state discrimination for the mixed state.

However, it is important to note that Bădescu and O'Donnell's algorithm requires many identical copies of the unknown state and thus does not generalize to our main result of PAC learning of quantum channels because every channel output might be a different state. On the other hand, as will shown in the following technical overview, our approach for approximate state discrimination involves a binary search through gap amplification and pretty good measurement and generalizes naturally to the PAC learning of quantum channels.

In [1], Aharonov, Cotler, and Qi introduced the notion of *quantum algorithmic measurement*, which broadly captures the query and computational complexity of quantum experiments, including those that generate unknown identical quantum states. In [19], Huang, Kueng, and Preskill compared the complexity of classically or quantumly training a machine learning model for predicting outcomes of physical experiments.

## 1.3 Technical Overview

The intuition behind both of our learning algorithms start with looking at the tensor product of all outputted states. The fidelity between such tensor produces decays exponentially in the number of samples drawn, so with enough samples , the tensor products from $\epsilon$-far concepts will become almost orthogonal (see Lemma 4), so intuitively, we should be able to distinguish between them.
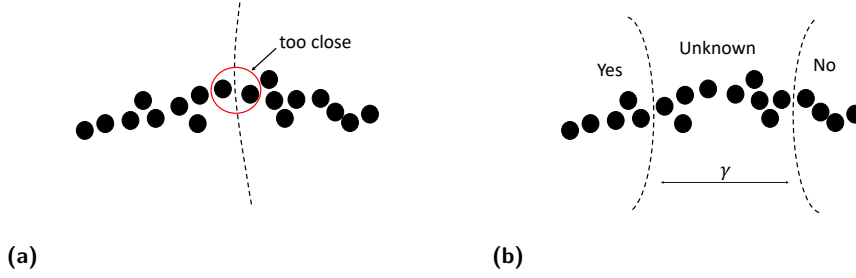
### 1.3.0.1 Pure State algorithm

In the case where the channels always output pure states, we have a rather simple algorithm. The key part is a theorem by Sen [28] on high dimensional random orthonormal measurements, which states that if we do a measurement of random orthonormal basis on two pure states, with high probability[5], the trace distance between the distribution of measurement outcome is lower bounded by a constant times the trace distance between those two states (see Theorem 11). This result might seem counter-intuitive, but remember that a random orthonormal measurement in $d$ dimension has $d$ possible outputs instead of 2. With this theorem in hand, the algorithm is rather easy: take enough samples to amplify the distance between outputted states and do a random orthonormal measurement on each sample. Choose the hypothesis as the channel that most likely to give the measurement result.

### 1.3.0.2 Mixed State algorithm

Our thought process on designing a learner for the channels that output mixed states is the following. In this case, Theorem 11 does not give us a useful result, so we need to find something else. Noticing the connection to the quantum state discrimination problem, we turned to pretty good measurement (PGM, Definition 3), a well studied tool for solving the quantum state discrimination problem. However, the lack of minimum distance between our outputted states is pretty pathological to PGM, so it was pretty easy to self-reject all our

---

[5] The probability goes to 1 as the dimension goes to infinity.

**(a)**                                    **(b)**

■ **Figure 1** (a) Pathological case when trying to cut the concept class into two sets. (b) Cutting the concept class into three sets.

attempts. Following that, we sought guidance from the analysis of classical Occam's razor algorithm, where a constant fraction of concepts are ruled out by each sample. We tried to divide the concept class into two sets, then do a PGM to distinguish those two, so we can recurse this into a binary search. Cutting the concepts into two sets does not work either because there can be concepts really close to any cut, which again is pathological to PGM. At this point, we realized that we need to have some kind of minimum distance for our PGM, so we cut the concept class into *three* sets, $S_{yes}$, $S_{no}$, and $S_{unknown}$. We set a minimum distance $\gamma$ between elements of $S_{yes}$ and $S_{no}$, so those two sets can be distinguished. This is the idea that works out. See Figure 1 for a graphical representation.

To follow our intuition in the previous paragraph, we give a definition about the distance between two sets of quantum states. Actually fidelity is more useful than trace distance, so we give the following definition of fidelity between sets of quantum states, which is the maximum fidelity among all pairs:

$$F\left(S_{yes}, S_{no}\right) = \max\left\{F(\sigma, \rho) | \sigma \in S_{yes}, \rho \in S_{no}\right\}$$

### 1.3.0.3   Bichromatic State Discrimination Problem (BSD)

The key component our mixed state algorithm is solving what we called $(\eta, N)$-*Bichromatic State Discrimination Problem* (BSD). The $(\eta, N)$-Bichromatic State Discrimination Problem is defined as follows: given complete information of two sets of quantum states, $S_{yes}$ and $S_{no}$, with fidelity $F(S_{yes}, S_{no}) \leq \eta$ and size $S_{yes} \leq N$, $S_{no} \leq N$, and one copy of an unknown quantum state $\sigma$, the goal is to decide whether $\sigma \in S_{yes}$ or $\sigma \in S_{no}$. A quantum algorithm solves $(\eta, N)$-BSD with error $\delta$ if for all $S_{yes}$ and $S_{no}$ such that $F(S_{yes}, S_{no}) \leq \eta$, $S_{yes} \leq N$, and $S_{no} \leq N$, given complete information about $S_{yes}$ and $S_{no}$ and one copy of an unknown quantum state $\sigma$ as input to the algorithm, the algorithm output a *yes/no* answer satisfies the following two conditions[6]:

**1.** If $\sigma \in S_{yes}$, the learner outputs *yes* with probability $(1 - \delta)$.

**2.** If $\sigma \in S_{no}$, the learner outputs *no* with probability $(1 - \delta)$.

See Figure 2 for some graphical intuition of BSD.

Note that BSD only requires maximum fidelity *between* the two sets; two states from the same set can be arbitrarily close. This does not violate quantum state discrimination lower bounds because the solver only needs to discriminate between the two sets.

We are able to show that BSD can be solved with good enough parameter:

---

[6] The learner can output anything if $\sigma$ does not come from either of the two sets.

**Figure 2** Bichromatic State Discrimination Problem.

▶ **Theorem 1.** *There exist an algorithm that solves $(\eta, N)$-BSD with error $\delta = N^2\eta$.*

The proof idea of Theorem 1 is trying to apply PGM on $S_{yes} \cup S_{no}$. We start with the observation that the result of [6] and [9], which gives an upper bound on PGM's error probability of mistaking one state as other states, can be generalized to an upper bound on PGM's error probability of mistaking one subset of states to its complement subset (See Appendix C). This almost gives us the required error bound for BSD, except that the PGM result is for the average case, where $\sigma$ is drawn from some probability distribution, so we turned it into a worst case result with the minimax argument of [17].

### 1.3.0.4 Back to Learning Quantum Channels

With BSD solved, we can get an algorithm that recursively exclude a constant fraction of the concept class. In each recursion, the algorithm partition the remain concepts into three sets, $S_{yes}, S_{unknown}$, and $S_{no}$. Ideally, $S_{yes}$ and $S_{no}$ both occupy a constant fraction of the remaining concepts and have minimum distance $\gamma = \Omega(1/\operatorname{poly}\log|C|)$. Noticing that the fidelity between tensor products of outputs decays exponentially with number of samples by lemma 4, the BSD between $O(\log|C|/\gamma)$ samples of $S_{yes}$ or $S_{no}$ can be solved with high probability. If the target concept is in $S_{yes}$, the BSD solver will return *yes* with high probability, and if the target concept is in $S_{no}$, the BSD solver will return *no*. If the target concept is from $S_{unknown}$, the BSD solver might return anything, but what we can be sure is that, if the BSD solver returned *yes*, the target concept is not from $S_{no}$, and if the BSD solver return *no*, the target concept is not from $S_{yes}$. Therefore, we can always exclude either $S_{yes}$ or $S_{no}$ as possible target concept.

There is another complication in that the distance between the concepts depends on the unknown distribution $D$ and thus cannot be calculated. In stead, we use the *empirical distance* between concepts, $\Delta_{emp}(c_1, c_2) = \frac{1}{T}\sum_{i=1}^{T}[\Delta_{tr}(c_1(x_i), c_2(x_i))]$, where $\{x_i\}$ are the inputs points we drawn in each recursion. Our calculation shows that that the error incurred from this change of distance measure is negligible.

### 1.3.0.5 Partition Sub-algorithm

It is not always possible to have an ideal partition where $S_{yes}$ and $S_{no}$ are both constant-fraction sized[7] and separated by the gap $\gamma$. Therefore, we designed a classical partition sub-algorithm (Algorithm 1) to handle these exceptions.

---

[7] By "constant-fraction sized" we mean "occupies a constant fraction of the remaining concepts".

An example where the ideal partition is not possible is the extreme case where every concept in the concept class is literally identical to each other. Note that in this extreme case can be trivially solved by output anything in the concept class as the hypothesis because everything is $\epsilon$-close to $c^*$.

Our partition sub algorithm builds on the intuition of what happened in the above extreme case. More specifically, our partition algorithm will not reserve a constant-fraction sized $S_{no}$ if a significant fraction of $C$ is clustered around a concept. In such case, we choose the cluster as $S_{yes}$ with a $\gamma$-thick "shell" of $S_{unknown}$ around it. If we measured *no*, we can rule out $S_{yes}$, which is a constant fraction of $|C|$. If we measured yes, we can output the center of the cluster as the hypothesis, and we tune $\gamma$ so that everything in either $S_{yes}$ or $S_{unknown}$ is $\epsilon$-close to the center. This completes our algorithm for mixed state outputs.

## 1.4 Lower Bounds and Agnostic Model

We complement our positive results on the sample complexity of PAC learning quantum channels with two simple lower bounds. First, by adapting a lower bound argument in [15], we prove that $\tilde{\Omega}((\log|C|)/\epsilon^2)$ samples are necessary to PAC learn quantum channels when the outputs are pure states, showing that our positive result is tight in the dependency on $|C|$ and $\epsilon$. In particular, for the dependency on $\epsilon$, this is in contrast with the classical results on the sample complexity for PAC learning concepts with Boolean outputs, where a tight $\Theta((\log|C|)/\epsilon)$ sample complexity is known [22, 16].[8]

Agnostic model is a learning model closely related to the PAC model, and the two models have similar sample complexity [18, 21]. In the agnostic model, the samples comes from a concept $c_s$ that is not necessarily inside the concept class $C$. Accordingly, the goal of the learner is to find, with $\epsilon$-distance error, the target concept $c^* \in C$ that is closest to $c_s$. We introduce the agnostic model for learning quantum channels, see section B.2 for details. Interestingly, in stark contrast to our algorithms that have dimension-independent sample complexity for learning quantum channels in PAC model, we found an $\Omega(\sqrt{d})$ lower bound on the sample complexity for learning quantum channels in agnostic model with output dimension $d$. Thus, in the agnostic model, learning quantum channels requires number of samples polynomial in the dimension, so it is not possible to efficiently learn quantum channels with large output dimension. Also, our negative example is in fact classical in nature, consisting of two concepts that output classical distributions, so learning classical distributions efficiently in agnostic model in large dimension is also impossible. However, since quantum pure states are not generalizations of classical distributions, the possibility of sample efficiently learn quantum channels with *pure state output* in agnostic is still open.

## 2 Preliminary

Throughout this paper, log is base 2 and ln is base $e$.

We use $\|\cdot\|_1$ to denote the trace norm $\|A\|_1 = \mathrm{tr}\sqrt{A^\dagger A}$. We use $\|\cdot\|_2$ or $\|\cdot\|_F$ to denote the Frobenius norm $\|A\|_2 = \sqrt{\mathrm{tr}(A^\dagger A)}$.

Denote the trace distance and fidelity between two distribution $D_1, D_2$ as $\Delta_{tr}(D_1, D_2)$ and $F(D_1, D_2)$, where the trace distance is equal to the total variation distance. Denote the trace distance and fidelity between two quantum states $\rho_1, \rho_2$ as $\Delta_{tr}(\rho_1, \rho_2) = \frac{1}{2}\|\rho_1 - \rho_2\|_1$

---

[8] The classical results show that the sample complexity is characterized by the VC dimension of the concept class $C$. In the case that $C$ is finite, $\log|C|$ is a trivial upper bound on the VC dimension.

and $F(\rho_1, \rho_2) = \left\| \sqrt{\rho_1}\sqrt{\rho_2} \right\|_1$. For a quantum state $\sigma$ and a quantum measurement $M$, denote $M(\sigma)$ as the output probability distribution when applying $M$ on $\sigma$.

Note that fidelity and trace distance are related by

$$1 - F \leq \Delta_{tr} \leq \sqrt{1 - F^2}.$$

For two quantum channel concepts $c_1$, $c_2$, define the distance between them with respect to $D$ as

$$\Delta(c_1, c_2) = \mathbb{E}_{x \in D} \left[ \Delta_{tr}(c_1(x), c_2(x)) \right].$$

We say that $c_1$, $c_2$ are $\epsilon$-close if $\Delta(c_1, c_2) \leq \epsilon$ and $\epsilon$-far if $\Delta(c_1, c_2) \geq \epsilon$. For two sets of concepts $S_1$ and $S_2$, define the distance between them as $\Delta(S_1, S_2) = \min \{ \Delta(c_1, c_2) | c_1 \in S_1, c_2 \in S_2 \}$.

## 2.1 Chernoff Bound

We use the following standard multiplicative version of Chernoff bound.

▶ **Theorem 2.** *Let* $X_1, \ldots, X_T \in [0, 1]$ *be independent random variables with* $\mathbb{E}[X_i] = \mu_i$. *Let* $X = (1/T) \sum_i X_i$, $\mu = (1/T) \sum_i \mu_i$ *and* $\alpha \in (0, 1)$. *We have*

$$\Pr[|X - \mu| \geq \alpha\mu] \leq 2^{-\Omega(\alpha^2 T \mu)}.$$

## 2.2 Pretty Good Measurement

The pretty good measurement (PGM) is defined as follows:

▶ **Definition 3** (pretty good measurement). *Let* $\{\sigma_i\}$ *be a set of density matrices and* $\{p_i\}$ *a probability distribution over* $\{\sigma_i\}$. *Define*

$$A_i = p_i \sigma_i, \ A = \sum_i A_i. \tag{1}$$

*The PGM associated with* $\{\sigma_i\}, \{p_i\}$ *is the measurement* $\{E_i\}$ *with*

$$E_i = A^{-1/2} A_i A^{-1/2}. \tag{2}$$

## 3 Problem Definitions

In this section we describe the PAC model of learning quantum channel and approximate state discrimination.

## 3.1 Classical PAC Learning Model

We start with a review of the classical PAC learning model.

In the classical probably approximately correct (PAC) learning model, a learner tries to learn a *target concept* $c^* \in C$ from a known *concept class* $C$, which is a set of Boolean functions $c : \{0, 1\}^n \to \{0, 1\}$, with respect to an *unknown* distribution $D$ over the input domain $\{0, 1\}^n$. Specifically, the learner is given access to a sample oracle $\mathcal{O}_{c^*, D}$, which generates i.i.d. samples $(x_i, c^*(x_i))$, where each $x_i \leftarrow D$ is drawn according to the distribution

$D$, and outputs a *hypothesis* $h \in C$.[9] The distance between two concepts $c$ and $h$ under the distribution $D$ is defined as $\Delta_D(c, h) = \mathbb{E}_{x \sim D} |c(x) - h(x)|$. The goal of the learner is to find a hypothesis $h$ with sufficiently small distance $\Delta_D(c^*, h)$ to $c^*$.

A learning algorithm $A$ is a *proper $(\epsilon, \delta)$-PAC learner* for a concept class $C$ if the following holds: For every $c^* \in C$ and distribution $D$, given oracle access to $\mathcal{O}_{c^*, D}$, $A^{\mathcal{O}_{c^*, D}}$ outputs an $h \in C$ such that $\Delta_D(c^*, h) \leq \epsilon$ with probability at least $1 - \delta$. The sample complexity of $A$ is the maximum number of samples $T$ that $A$ needs to query $\mathcal{O}_{c^*, D}$ to output $h$. The *proper $(\epsilon, \delta)$-PAC sample complexity* of a concept class $C$ is the minimum sample complexity over all learners. A $\tilde{\Theta}((\log |C|)/\epsilon)$ sample complexity is known [22, 16].[10]

## 3.2   Learning Quantum Channels in PAC model

We now generalize classical PAC learning to the context of learning quantum channels. As above, we consider a learner trying to learn a target concept $c^* \in C$ from a known concept class $C$ with respect to an unknown distribution $D$. Here, we consider the concept class $C$ as a finite set of known $d_1$ to $d_2$ dimensional quantum channels, and $D$ as a distribution over the Hilbert space of dimension $d_1$. Precisely, the learner is given access to a sample oracle $\mathcal{O}_{c^*, D}$ and outputs a hypothesis $h \in C$. The oracle $\mathcal{O}_{c^*, D}$ generates i.i.d. samples $(x_i, c^*(x_i))$, where each $x_i \leftarrow D$ is the classical description of a state drawn according to the distribution $D$, and $c^*(x_i)$ is the (potentially mixed) quantum state outputted by $c^*$ on input $x_i$.

The distance between two concepts $c$ and $h$ under the distribution $D$ is the expected trace distance $\Delta(c, h) = \mathbb{E}_{x \in D} [\Delta_{tr}(c(x), h(x))]$. The goal of the learner is to find a hypothesis $h \in C$ with sufficiently small $\Delta(c^*, h)$.

A quantum learning algorithm $A$ is a *proper $(\epsilon, \delta)$-PAC learner* for $C$ if the following holds: For every $c^* \in C$ and distribution $D$, given oracle access to $\mathcal{O}_{c^*, D}$, $A^{\mathcal{O}_{c^*, D}}$ outputs an $h \in C$ such that $\Delta_D(c^*, h) \leq \epsilon$ with probability at least $1 - \delta$. The sample complexity of $A$ is the maximum number of samples $T$ that $A$ needs to query $\mathcal{O}_{c^*, D}$ to output $h$. The *proper $(\epsilon, \delta)$-PAC sample complexity* of a concept class $C$ is the minimum sample complexity over all learners.

## 3.3   Approximate State Discrimination

Let $S$ be a finite set of $d$-dimensional density matrices. We want to learn a target state $\sigma \in S$ using as few identical copies of $\sigma$ as possible. A quantum algorithm is an $(\epsilon, \delta)$-approximate discriminator of $S$ if it takes the description of $S$ and $T$ copies of $\sigma$ as input and with probability $1 - \delta$ outputs a state $\rho \in S$ with $\Delta_{tr}(\rho, \sigma) \leq \epsilon$, for any $\sigma \in S$.

Note that approximate state discrimination can be viewed as a special case of PAC learning quantum channels with constant output, so the algorithms for PAC learning quantum channels in Section 4 and Section 5 trivially works for approximate state discrimination.

## 4   PAC Learning Quantum Channels with Pure State Output

See Appendix A

---

[9] The requirement that the hypothesis $h$ is in the concept class $C$ is referred to as proper learning. We focus on proper learning since our algorithms satisfy this property.

[10] We use $\tilde{\Theta}$ to denote $\Theta$ with log factors. The classical results show that the sample complexity is $\Theta((d + \log 1/\delta)/\epsilon)$, where $d$ is the VC dimension of the concept class. In the case where $|C|$ is finite, $\log |C|$ is a trivial upper bound on $d$, and there are concept classes whose VC dimension $d$ matches $\log |C|$.

## 5 PAC Learning Quantum Channels with Mixed State Output

The random orthonormal measurement approach in Section 4 does not work since two high dimensional mixed states with constant trace distance between them can have negligible Frobenius distance between them. Instead, We follow the intuitions detailed in Section 1.3. We define the bichromatic state discrimination problem (BSD), solve BSD with PGM techniques , and build our learner algorithm with the BSD solver and a partition sub-algorithm.

Before we show the algorithms for bicromatic state discrimination, let us first show that we can efficiently amplify the distance between concepts by taking samples.

▶ **Lemma 4** (concept distance amplification). *Let $c$ be a quantum channel concept $\epsilon$-far from the target concept $c^*$. Let $\{x_1, x_2, \ldots, x_T\}$ be $T$ inputs drawn from the distribution $D$. With probability $1 - 2^{-\Omega(T\epsilon)}$ over $\{x_i\}$ drawn, we have*

$$F\left(\bigotimes_{i \in [T]} c(x_i), \bigotimes_{i \in [T]} c^*(x_i)\right) \leq 2^{-\Omega(T\epsilon^2)} \tag{3}$$

*and*

$$\Delta_{tr}\left(\bigotimes_{i \in [T]} c(x_i), \bigotimes_{i \in [T]} c^*(x_i)\right) \geq 1 - 2^{-\Omega(T\epsilon^2)}. \tag{4}$$

**Proof.** By Chernoff bound, with probability $1 - 2^{-\Omega(T\epsilon)}$,

$$\sum_i \Delta_{tr}\left(c(x_i), c^*(x_i)\right) \geq \frac{1}{2}T\epsilon. \tag{5}$$

Then by Cauchy-Schwarz Inequality,

$$\sum_i (\Delta_{tr}\left(c(x_i), c^*(x_i)\right))^2 \geq \frac{1}{4}T\epsilon^2. \tag{6}$$

Then the amplified fidelity is bounded by

$$F\left(\bigotimes_i c(x_i), \bigotimes_i c^*(x_i)\right) = \Pi_i F\left(c(x_i), c^*(x_i)\right)$$

$$\leq \Pi_i \sqrt{1 - (\Delta_{tr}\left(c(x_i), c^*(x_i)\right))^2}$$

$$\leq \exp\left[-\frac{1}{2}\sum_i (\Delta_{tr}\left(c(x_i), c^*(x_i)\right))^2\right] = 2^{-\Omega(T\epsilon^2)}, \tag{7}$$

where the last inequality is true because $1 - x \leq e^{-x}$. And the amplified trace distance is

$$\Delta_{tr}\left(\bigotimes_{i \in [T]} c(x_i), \bigotimes_{i \in [T]} c^*(x_i)\right) \geq 1 - F\left(\bigotimes_i c(x_i), \bigotimes_i c^*(x_i)\right) = 1 - 2^{-\Omega(T\epsilon^2)}. \tag{8}$$

◀

Lemma 4 means that we can amplify the distance between tensor products of samples from quantum channels as efficiently as we do on samples of fixed quantum states. This means that PAC learning quantum channels is really similar to approximate state discrimination even in the mixed state case.

Now back to BSD. The bichromatic state discrimination problem (BSD) is defined as follows:

▶ **Definition 5** (Bichromatic State Discrimination Problem (BSD)). *Given complete information of two sets of quantum states, $S_{yes}$ and $S_{no}$, with fidelity $F(S_{yes}, S_{no}) \leq \eta$ and size $S_{yes} \leq N$, $S_{no} \leq N$, and one copy of an unknown quantum state $\sigma$, the goal is to decide whether $\sigma \in S_{yes}$ or $\sigma \in S_{no}$. We say a quantum algorithm solves $(\eta, N)$-BSD with error $\delta$ if for all $S_{yes}$ and $S_{no}$ such that $F(S_{yes}, S_{no}) \leq \eta$, $S_{yes} \leq N$, and $S_{no} \leq N$, given complete information about $S_{yes}$ and $S_{no}$ and one copy of an unknown quantum state $\sigma$ as input to the algorithm, the algorithm output and yes/no answer satisfies the following two conditions:*
1. *If $\sigma \in S_{yes}$, the learner outputs yes with probability $(1 - \delta)$.*
2. *If $\sigma \in S_{no}$, the learner outputs no with probability $(1 - \delta)$.*
*The learner can output anything if $\sigma$ does not come from either of the two sets.*

We show the existence of a BSD solver by first showing that PGM over $S_{yes} \cup S_{no}$ solves the "average case" BSD and then turn it into a "worst case" result by the minimax theorem.

First by slightly modifying a result of [9] and [6], We show that PGM can solve the "average case" BSD:

▶ **Lemma 6** (PGM for "average BSD"). *Let $S_{yes}, S_{no}$ be two sets of density matrices and $\{p_i\}$ be a probability distribution over $S_{yes} \cup S_{no}$.* [11] *The PGM on $S_{yes} \cup S_{no}, \{p_i\}$ satisfies*

$$\sum_{i \in S_{yes}} \sum_{j \in S_{no}} [p_i \Pr(PGM(\sigma_i) = j) + p_j \Pr(PGM(\sigma_j) = i)] \leq \sum_{i \in S_{yes}} \sum_{j \in S_{no}} F(\sigma_i, \sigma_j). \quad (9)$$

**Proof.** See appendix C. ◀

We can group together the outputs of the PGM in Lemma 6 and define a binary measurement $\{E_{yes}, E_{no}\}$, where $E_{yes} = \sum_{i \in S_{yes}} E_i$, $E_{no} = \sum_{i \in S_{no}} E_i$, and $\{E_i\}$ is the PGM. By Lemma 6, the binary measurement solves "average BSD" with error probability at most $\sum_{i \in S_{yes}} \sum_{j \in S_{no}} F(\sigma_i, \sigma_j)$.[12]

Since the upper bound on error is independent of the distribution $\{p_i\}$, minimax theorem guarantees the existence of a measurement that distinguishes between $S_{yes}$ and $S_{no}$ for any distribution $\{p_i\}$ with error probability less than $\sum_{i \in S_{yes}} \sum_{j \in S_{no}} F(\sigma_i, \sigma_j)$[13]. In particular, if $p_i = 1$ for some $\sigma_i \in S_{yes}$, the probability of the minimax measurement mistaking $\sigma_i$ as something in $S_{no}$ is upper bounded by $\sum_{i \in S_{yes}} \sum_{j \in S_{no}} F(\sigma_i, \sigma_j)$, and vice versa. We formalize this discussion as the following Theorem.

▶ **Theorem 7** (solver for BSD, Theorem 1 restated). *There exist an algorithm that solves $(\eta, N)$-BSD with error $\delta = N^2 \eta$*

**Proof.** Consider the zero sum game between two players where player1 choose a probability distribution $\{p_i\}$ over $S_{yes} \cup S_{no}$ and player2 choose a binary measurement strategy $M$. The score of player1 is given by the following error probability[14]:

$$P_{bi-error} = \sum_{i \in S_{yes}} [p_i \Pr(M(\sigma_i) = no)] + \sum_{j \in S_{no}} [p_j \Pr(M(\sigma_j) = yes)] \quad (10)$$

---

[11] We will slightly abuse the notation and write $i \in S_{yes}$ or $j \in S_{no}$ instead of $\sigma_i \in S_{yes}$ or $\sigma_j \in S_{no}$.

[12] A careful reader might notice that since we only want a binary answer, we are essentially distinguishing the states $A_{yes} = \sum_{i \in S_{yes}} p_i \sigma_i$ and $A_{no} = \sum_{j \in S_{no}} p_j \sigma_j$, and thus the optimal error probability is characterized by trace distance between $A_{yes}$ and $A_{no}$. However, to our knowledge there is no inequality in the literature giving a *lower bound* on trace distance between on linear combinations of density matrices, so actually, the other direction of the trace-distance characterization is the relevant one: Lemma 6 gives a new lower bound on $\Delta_{tr}(A_{yes}, A_{no})$.

[13] This argument was used in [17]

[14] We will slightly abuse the notation and write $i \in S_{yes}$ or $j \in S_{no}$ instead of $\sigma_i \in S_{yes}$ or $\sigma_j \in S_{no}$.

It is easy to check that that strategies of both sides are linear, so we can apply the minimax theorem to get

$$\min_{M} \max_{\{p_i\}} P_{bi-error} = \max_{\{p_i\}} \min_{M} P_{bi-error} \le \sum_{i \in S_{yes}} \sum_{j \in S_{no}} F(\sigma_i, \sigma_j) \le N^2 \eta, \tag{11}$$

where the second inequality is from the promises of $(\eta, N)$ BSD, and the first inequality is shown by considering the binary measurement $\{E_{yes}, E_{no}\}$, where $E_{yes} = \sum_{i \in S_{yes}} E_i$, $E_{no} = \sum_{i \in S_{no}} E_i$, and $\{E_i\}$ is the PGM of Lemma 6. This means that there is a measurement $M$ whose error probability is less than $\sum_{i \in S_{yes}} \sum_{j \in S_{no}} F(\sigma_i, \sigma_j)$ for all probability distribution $\{p_i\}$. In particular, the error probability is at most $N^2 \eta$ when player1 uses the deterministic strategy of always choosing some specific state $\sigma_i \in S_{yes} \cup S_{no}$. Therefore, algorithm of applying the measurement $M$ solves $(\eta, N)$-BSD with error $N^2\delta$.

◄

Theorem 7 implies that if we amplify the maximum fidelity between $S_{yes}$ and $S_{no}$ by Lemma 4 to less than $O(1/|C|^2)$, we have a constant error probability in distinguishing whether a state is from $S_{yes}$ or $S_{no}$. By lemma 4 this requires $\Theta\left(\log |C|/\gamma^2\right)$ samples if the distance between $S_{yes}$ and $S_{no}$ is $\gamma$.

Now we present the partition sub-algorithm. Let $C_r$ be the set of remaining concepts that have not been cut off by the main algorithm. The sub-algorithm partitions the remaining concepts into three disjoint subsets: $(S_{yes}, S_{unknown}, S_{no})$, such that $|S_{yes}| \ge \frac{1}{9}|C_r|$[15], and $\Delta(S_{yes}, S_{no}) \ge \gamma = \Theta(\epsilon/\log|C_r|)$. The sub-algorithm might or might not found an extreme case. If no extreme case is found, $|S_{no}| \ge \frac{1}{9}|C_r|$. If an extreme case is found, more than $\frac{1}{3}|C_r|$ concepts are $\epsilon$-close to some concept. The sub-algorithm initialized with every concept in $S_{no}$. It then repeatedly picks a concept $c_c$ from $S_{no}$ and adds concepts within the ball around $c_c$ to $S_{yes}$ and concepts in a $\gamma$-shell around the ball to $S_{unknown}$. The $\gamma$-shell of $S_{no}$ ensures that $\Delta(S_{yes}, S_{no}) \ge \gamma$ and we choose the radius of the ball so that the number of concepts added to $S_{yes}$ is greater than half the number of concepts added to $S_{unknown}$ to ensure that $|S_{yes}| > \frac{1}{2}|S_{unknown}|$ in the end. The sub-algorithm keeps adding concepts to $S_{yes}$ and $S_{unknown}$ until $|S_{yes}| + |S_{unknown}| > \frac{1}{3}|C_r|$ or the loop is breaked by an extreme case. The sub-algorithm reports an extreme case if the number of concepts to be added to $S_{yes}$ and $S_{unknown}$ in the current iteration is greater than $\frac{1}{3}|C_r|$. In this case we know that more than $\frac{1}{3}|C_r|$ concepts are around $c_c$. If no extreme case is found, since the loop stops when $|S_{yes}| + |S_{unknown}| > \frac{1}{3}|C_r|$ and the last iteration cannot add more than $\frac{1}{3}|C_r|$ concepts to $S_{yes}$ or $S_{unknown}$, there are at least $(1 - \frac{1}{3} - \frac{1}{3})|C_r| > \frac{1}{9}|C_r|$ concepts left in $S_{no}$, and $|S_{yes}| > \frac{1}{3}(|S_{yes}| + |S_{unknown}|) > \frac{1}{9}|C_r|$.

There is another complication in that the distance between the concepts depends on the unknown distribution $D$ and thus cannot be calculated. In stead, we calculate the *empirical distance* between concepts, $\Delta_{emp}(c_1, c_2) = \frac{1}{T} \sum_{i=1}^{T} [\Delta_{tr}(c_1(x_i), c_2(x_i))]$, which depends on the input points drawn from $D$. We also tune $\epsilon$ into $\epsilon/2$ to accommodate for the extra error incurred.

The sub-algorithm is detailed in Algorithm 1.

▶ **Lemma 8.** *The output of Algorithm 1 satisfies the following conditions:*

*$(S_{yes}, S_{unknown}, S_{no})$ is a partition of $C_r$. $\Delta_{emp}(S_{yes}, S_{no}) \ge \gamma = \epsilon/4 \log|C_r|$. $|S_{yes}| \ge \frac{1}{9}|C_r|$. If flag_extreme = false, $|S_{no}| \ge \frac{1}{9}|C_r|$. If flag_extreme = true, $\Delta_{emp}(c, c_c) \le \epsilon/2$, $\forall c \in (S_{yes} \cup S_{unknown})$.*

---

[15] $\frac{1}{9}$ is an arbitrary constant and can be further optimized

◼ **Algorithm 1** partition sub-algorithm.

---

    **Data:**  concepts class $C_r$, real number $\epsilon$.
    **Result:**  Set of concepts $S_{yes}, S_{unknown}, S_{no}$, boolean variable *flag_extreme*,
             concept $c_c$

**1**  $S_{no} \leftarrow C_r$, $S_{yes} \leftarrow \emptyset$, $S_{unknown} \leftarrow \emptyset$, *flag_extreme* $\leftarrow false$, $\gamma \leftarrow \epsilon/(4\log|C_r|)$.
**2**  **while**  $|S_{yes}| + |S_{unknown}| < \frac{1}{3}|C_r|^{16}$ **do**
**3**     |  $c_c \leftarrow$ a random concept in $S_{no}$;
**4**     |  Count the number of concept in $S_{no}$ whose distance to $c_c$ is in the interval
         |  $[(m-1)\gamma, m\gamma)$ for all $m \in [1/\gamma]$ and record the number as $b_m$. I.e.
         |  $b_m \leftarrow |\{c|\Delta(c, c_c) \in [(m-1)\gamma, m\gamma), c \in S_{no}\}|$;
**5**     |  Find the smallest $i^* \geq 2$ such that $b_{i^*} < 2\sum_{i \in [i^*-1]} b_i$;
**6**     |  **if**  $\sum_{i \in [i^*-1]} b_i + b_{i^*} > \frac{1}{3}|C_r|$ **then**
**7**     |   |  *flag_extreme* $\leftarrow true$;
**8**     |   |  move everything in $S_{yes}$ and $S_{unknown}$ back to $S_{no}$;
**9**     |   |  run line 12 once;
**10**    |   |  Terminate;
**11**    |  **end**
**12**    |  For the concepts in $S_{no}$, move the concepts within distance $(i^*-1)\gamma$ of $c_c$ to $S_{yes}$,
        |  and move the concepts whose distance to $c_c$ is in $[(i^*-1)\gamma, i^*\gamma)$ to $S_{unknown}$.
        |  I.e. move $\{c|\Delta(c, c_c) \in [0, (i^*-1)\gamma), c \in S_{no}\}$ to $S_{yes}$ and move
        |  $\{c|\Delta(c, c_c) \in [(i^*-1)\gamma, i^*\gamma), c \in S_{no}\}$ to $S_{unknown}$;
**13** **end**

---

**Proof.** First note that in line 5, $\gamma = \epsilon/(4\log|C_r|)$ ensures that $i^*$ exists and $i^* \leq \epsilon/(2\gamma)$. This can be proved by contradiction: if $b_i^* \geq 2\sum_{i \in [i^*-1]} b_i, \forall i^* \leq \epsilon/(2\gamma)$, then $b_i^* > 2b_{i^*-1}, \forall i^* \leq \epsilon/(2\gamma)$. Together with $b_1 \geq 1$ because $\Delta(c_c, c_c) = 0$, we have $b_{\lfloor \epsilon/2\gamma \rfloor} \geq 2 \cdot 2^{\log|C_r|} b_1 \geq |C_r|$, a contradiction.

$(S_{yes}, S_{unknown}, S_{no})$ is a partition because it is initialized as a partition and we only moves elements between them. Note that whenever we move something to $S_{yes}$, we move a $\gamma$-thick shell around it to $S_{unknown}$. By triangle inequality of empirical distances between concepts, $\Delta_{emp}(S_{yes}, S_{no}) \geq \gamma = \epsilon/4\log|C_r|$ at the end of every step.

If no extreme case is found, at each iteration of the loop at line 12, $(\sum_{i \in [i^*-1]} b_i)$ concepts are moved to $S_{yes}$ from $S_{no}$, and $b_{i^*}$ concepts are moved to $S_{unknown}$ from $S_{no}$. Before the last iteration of the loop $|S_{yes}| + |S_{unknown}| \leq \frac{1}{3}|C_r|$, and the number of concepts moved to $S_{yes}$ and $S_{unknown}$ in the last iteration is $\sum_{i \in [i^*-1]} b_i + b_i^* \leq \frac{1}{3}|C_r|$, so $|S_{no}| \geq (1 - \frac{1}{3} - \frac{1}{3})|C_r| > \frac{1}{9}|C_r|$. Because of the requirement $b_{i^*} < 2\sum_{i \in [i^*-1]} b_i$ in line 5, $\sum_{i \in [i^*-1]} b_i > \frac{1}{3}(\sum_{i \in [i^*-1]} b_i + b_{i^*})$ and thus $|S_{yes}| > \frac{1}{3}(|S_{yes}| + |S_{unknown}|)$ at the end of every loop. Combined with the loop-termination condition $|S_{yes}| + |S_{unknown}| > \frac{1}{3}|C_r|$, we have $|S_{yes}| > \frac{1}{9}|C_r|$.

If an extreme case is found at line 7, because we moved everything back to $S_{no}$, all concepts in $S_{yes}$ or $S_{unknown}$ are added in that one call of line 12, and thus they are all $(i^*\gamma)$-close to $c_c$ . Recall that $i^*\gamma \leq \epsilon/2$, so everything in $S_{yes}$ or $S_{unknown}$ is $\epsilon/2$-close to $c_c$. The analysis on $|S_{yes}|$ is a bit subtle. Similar to the previous paragraph, we have $\sum_{i \in [i^*-1]} b_i > \frac{1}{3}(\sum_{i \in [i^*-1]} b_i + b_{i^*})$. Combined with $\sum_{i \in [i^*-1]} b_i + b_{i^*} > \frac{1}{3}|C_r|$ to trigger line 7, we have $\sum_{i \in [i^*-1]} b_i > \frac{1}{9}|C_r|$. Since the wiping of $S_{yes}$ and $S_{unknown}$ at the beginning of line 7 only opens more possible concepts to be added to $S_{yes}$, we have $|S_{yes}| \geq \sum_{i \in [i^*-1]} b_i > \frac{1}{9}|C_r|$. ◀

With the partition sub-algorithm described, we detail the main algorithm for mixed state case in Algorithm 2.

▌ **Algorithm 2** algorithm for mixed state case.

---

**Data:** Concept class $C$, Sampling Oracle $\mathcal{O}_{c^*,D}$
**Result:** hypothesis $h$

**1** $C_r \leftarrow C$ ;

**2** $T \leftarrow \Theta\left(\frac{\log^2 |C|(\log|C|+\log(1/\delta))}{\epsilon^2}\right)$ ;

**3 while do**

**4**    Call $\mathcal{O}_{c^*,D}$ $T$ times, getting $T$ samples
     $\{(x_1, c^*(x_1)), (x_2, c^*(x_2)), \ldots (x_T, c^*(x_T))\}$;

**5**    $(S_{yes}, S_{unknown}, S_{no}, flag\_extreme, c_c) \leftarrow$ (Algorithm 1)$(C_r, \epsilon)$;

**6**    Construct the measurement $M$ in Theorem 7 between $S_{yes}$ and $S_{no}$ with the
     state $\sigma_i$ corresponding to concept $c_i$ being $\sigma_i = \bigotimes_{j\in[T]} c_i(x_j)$;

**7**    $Measure\_result \leftarrow M(\bigotimes_{j\in[T]} c^*(x_j))$.;

**8**    **if** $Measure\_result = no$ **then**

**9**     | remove $S_{yes}$ from $C_r$;

**10**    **end**

**11**    **if** $Measure\_result = yes$ and $flag\_extreme = false$ **then**

**12**     | remove $S_{no}$ from $C_r$.;

**13**    **end**

**14**    **if** $Measure\_result = yes$ and $flag\_extreme = true$ **then**

**15**     | $h \leftarrow c_c$;

**16**     | Terminate;

**17**    **end**

**18 end**

---

Now we state and prove our result for mixed state case:

▶ **Theorem 9.** *Algorithm 2 is a proper $(\epsilon, \delta)$-PAC learner for any quantum circuit concept class $C$, using*

$$O\left(\frac{\log^3 |C|(\log|C| + \log(1/\delta))}{\epsilon^2}\right)$$

*samples.*

**Proof.** By Lemma 8, Algorithm 2 removes at least $\frac{1}{9}|C_r|$ concepts from $C_r$ in each loop unless it terminates, so it terminates in $O(\log|C|)$ loops at line 16. Combined with the fact that Algorithm 2 takes $O\left(\frac{\log^2 |C|(\log|C|+\log(1/\delta))}{\epsilon^2}\right)$ samples each loop, its sample complexity is

$$O\left(\frac{\log^3 |C|(\log|C| + \log(1/\delta))}{\epsilon^2}\right)$$

.

As for the correctness of the algorithm, first note that by Lemma 8 the empirical distance between any pair of concepts in $S_{yes}$ and $S_{no}$ is at least $\gamma_0 = \epsilon/(4\log|C|)$.

Consider any pair of concepts $c_i \in S_{yes}$ and $C_j \in S_{no}$, with the corresponding states $\sigma_i$ and $\sigma_j$. By definition of empirical distance,

$$\sum_{k=1}^{T} \Delta_{tr}(c_i(x_k), c_j(x_k)) \geq T\gamma_0 \tag{12}$$

Then by Cauchy-Schwarz Inequality,

$$\sum_{k=1}^{T} \Delta_{tr}\left(c_i(x_k), c_j(x_k)\right)^2 \geq T\gamma_0^2. \tag{13}$$

Then the fidelity between $\sigma_i$ and $\sigma_j$ is bounded by

$$
\begin{aligned}
F(\sigma_i, \sigma_j) &= F\left(\bigotimes_k c_i(x_k), \bigotimes_k c_j(x_k)\right) \\
&= \Pi_k F\left(c_i(x_k), c_j(x_k)\right) \\
&\leq \Pi_k \sqrt{1 - \left(\Delta_{tr}\left(c_i(x_k), c_j(x_k)\right)\right)^2} \\
&\leq \exp\left[-\frac{1}{2}\sum_k \left(\Delta_{tr}\left(c_i(x_k), c_j(x_k)\right)\right)^2\right] \\
&= 2^{-\Omega\left(T\gamma_0^2\right)}
\end{aligned}
\tag{14}
$$

where the last inequality is true because $1 - x \leq e^{-x}$.

There are only two possible ways for Algorithm 2 to make an error: first is to remove $c^*$ from $C_r$ in line 9 or line 12, and second is to output a far-away concept at line 15 because of the mismatch between empirical distance and true distance.

For the first error, note that $c^*$ always has empirical distance zero to it self, no matter what $\{x_1, x_2, \ldots, x_T\}$ are sampled. By Theorem 7 and Equation 14 the error probability in each loop is bounded by

$$P_{error,1} \leq |C_r|^2 \cdot 2^{-\Omega(T\gamma_0^2)}. \tag{15}$$

Apply union bound over $O(\log|C|)$ loop we can bound the total error probability by

$$P_{total\,error,1} \leq \log|C||C|^2 \cdot 2^{-\Omega(T\gamma^2)} \leq O\left(\frac{\delta|C|^2\log|C|}{\text{poly}(|C|)}\right) \leq O(\delta) \tag{16}$$

For the second error, consider a pair of concepts that has distance bigger than $\epsilon$. By Chernoff bound, the probability that their empirical distance is less than $\frac{1}{2}\epsilon$ is less than $2^{-\Omega(T\epsilon^2)}$. Union bound over all $O(|C|^2)$ pairs of concepts, we have

$$P_{total\,error,2} \leq |C|^2 \cdot 2^{-\Omega(T\epsilon^2)} \ll P_{total\,error,1}. \tag{17}$$

◀

## References

1   Dorit Aharonov, Jordan Cotler, and Xiao-Liang Qi. Quantum algorithmic measurement. *arXiv preprint arXiv:2101.04634*, 2021.

2   Noga Alon, Shai Ben-David, Nicolo Cesa-Bianchi, and David Haussler. Scale-sensitive dimensions, uniform convergence, and learnability. *Journal of the ACM (JACM)*, 44(4):615–631, 1997.

3   Srinivasan Arunachalam and Ronald de Wolf. Guest column: A survey of quantum learning theory. *SIGACT News*, 48(2):41–67, 2017. `doi:10.1145/3106700.3106710`.

4   Srinivasan Arunachalam and Ronald de Wolf. Optimal quantum sample complexity of learning algorithms. In *Computational Complexity Conference*, volume 79 of *LIPIcs*, pages 25:1–25:31. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, 2017.

**5**    Srinivasan Arunachalam, Alex B Grilo, and Aarthi Sundaram. Quantum hardness of learning shallow classical circuits. *arXiv preprint arXiv:1903.02840*, 2019.

**6**    Koenraad MR Audenaert and Milán Mosonyi. Upper bounds on the error probabilities and asymptotic error exponents in quantum multiple state discrimination. *Journal of Mathematical Physics*, 55(10):102201, 2014.

**7**    Costin Bădescu and Ryan O'Donnell. Improved quantum data analysis. *arXiv preprint arXiv:2011.10908*, 2020.

**8**    Joonwoo Bae and Leong-Chuan Kwek. Quantum state discrimination and its applications. J. Phys. A: Math. Theor. 48 083001 (2015), 2017. `doi:10.1088/1751-8113/48/8/083001`.

**9**    Howard Barnum and Emanuel Knill. Reversing quantum dynamics with near-optimal quantum and classical fidelity. *Journal of Mathematical Physics*, 43(5):2097–2106, 2002.

**10**   Peter L Bartlett, Philip M Long, and Robert C Williamson. Fat-shattering and the learnability of real-valued functions. *journal of computer and system sciences*, 52(3):434–452, 1996.

**11**   Shai Bendavid, Nicolo Cesabianchi, David Haussler, and Philip M Long. Characterizations of learnability for classes of $\{0, ..., n\}$-valued functions. *Journal of Computer and System Sciences*, 50(1):74–86, 1995.

**12**   Anselm Blumer, Andrzej Ehrenfeucht, David Haussler, and Manfred K Warmuth. Occam's razor. *Information processing letters*, 24(6):377–380, 1987.

**13**   Anselm Blumer, Andrzej Ehrenfeucht, David Haussler, and Manfred K Warmuth. Learnability and the vapnik-chervonenkis dimension. *Journal of the ACM (JACM)*, 36(4):929–965, 1989.

**14**   Hao-Chung Cheng, Min-Hsiu Hsieh, and Ping-Cheng Yeh. The learnability of unknown quantum measurements. QIC, Vol. 16, No. 7-8, 0615-0656 (2016), 2015. `arXiv:arXiv:1501.00559`.

**15**   Jeongwan Haah, Aram W. Harrow, Zhengfeng Ji, Xiaodi Wu, and Nengkun Yu. Sample-optimal tomography of quantum states. In *Proceedings of the Forty-eighth Annual ACM Symposium on Theory of Computing*, STOC '16, pages 913–925, New York, NY, USA, 2016. ACM. `doi:10.1145/2897518.2897585`.

**16**   Steve Hanneke. The optimal sample complexity of pac learning. *The Journal of Machine Learning Research*, 17(1):1319–1333, 2016.

**17**   Aram W Harrow and Andreas Winter. How many copies are needed for state discrimination? *IEEE Transactions on Information Theory*, 58(1):1–2, 2012.

**18**   David Haussler. Decision theoretic generalizations of the PAC model for neural net and other learning applications. *Inf. Comput.*, 100(1):78–150, 1992. `doi:10.1016/0890-5401(92)90010-D`.

**19**   Hsin-Yuan Huang, Richard Kueng, and John Preskill. Information-theoretic bounds on quantum advantage in machine learning. *arXiv preprint arXiv:2101.02464*, 2021.

**20**   Michael J Kearns and Robert E Schapire. Efficient distribution-free learning of probabilistic concepts. *Journal of Computer and System Sciences*, 48(3):464–497, 1994.

**21**   Michael J. Kearns, Robert E. Schapire, and Linda Sellie. Toward efficient agnostic learning. *Machine Learning*, 17(2-3):115–141, 1994. `doi:10.1007/BF00993468`.

**22**   Michael J. Kearns and Umesh V. Vazirani. *An Introduction to Computational Learning Theory*. MIT Press, Cambridge, MA, USA, 1994.

**23**   Masoud Mohseni, AT Rezakhani, and DA Lidar. Quantum-process tomography: Resource analysis of different strategies. *Physical Review A*, 77(3):032322, 2008.

**24**   Ashley Montanaro. On the distinguishability of random quantum states. Comm. Math. Phys. 273(3), pp. 619-636, 2007, 2006. `doi:10.1007/s00220-007-0221-7`.

**25**   Ashley Montanaro. A lower bound on the probability of error in quantum state discrimination. In *Information Theory Workshop, 2008. ITW'08. IEEE*, pages 378–380. IEEE, 2008.

**26**   Balas K Natarajan. On learning sets and functions. *Machine Learning*, 4(1):67–97, 1989.

**27**   Ryan O'Donnell and John Wright. Efficient quantum tomography. In *Proceedings of the Forty-eighth Annual ACM Symposium on Theory of Computing*, STOC '16, pages 899–912, New York, NY, USA, 2016. ACM. `doi:10.1145/2897518.2897544`.

**28**   Pranab Sen. Random measurement bases, quantum state distinction and applications to the hidden subgroup problem. In *Computational Complexity, 2006. CCC 2006. Twenty-First Annual IEEE Conference on*, pages 14–pp. IEEE, 2005.

**29**   J. Prabhu Tej, Syed Raunaq Ahmed, A. R. Usha Devi, and A. K. Rajagopal. Quantum hypothesis testing and state discrimination. arXiv:1803.04944, 2018.

**30**   Leslie G. Valiant. A theory of the learnable. *Commun. ACM*, 27(11):1134–1142, 1984. `doi:10.1145/1968.1972`.

## A    PAC Learning Quantum Channels with Pure State Output

The algorithm follows ideas by Sen [28], who shows that random orthonormal measurement preserves trace distance between pure states. One can then apply random orthonormal measurements on each sampled output and take enough samples to amplify the distance between $\epsilon$-far concepts to $1 - O(1/|C|)$ and show that the probability for the maximum likelihood estimate to select a $\epsilon$-far concept over the target concept is less than $O(1/|C|)$. Take a union bound and we have a bounded error probability.

▶ **Theorem 10.** *Algorithm 3 is a proper $(\epsilon, \delta)$-PAC learner for any concept class $C$ of quantum channels with pure state outputs, using*

$$O\left(\frac{(\log|C|) + \log(1/\delta)}{\epsilon^2}\right)$$

*samples.*

**Algorithm 3**  algorithm for pure state output.

---

**1**  Take $T = \Theta((\log|C| + \log(1/\delta))/\epsilon^2)$ samples $(x_1, \sigma_1), (x_2, \sigma_2), \ldots, (x_T, \sigma_T)$ ;

**2**  Do a random orthonormal measurement[a] $M_i$ on each output state $\sigma_i$. Let the measured outputs be $\{z_i\}$ ;

**3**  Output the concept $h \in C$ that is most likely to give the measured result of line 2:

$$h = \underset{c \in C}{\arg\max} \, \Pi_{i \in [T]} \Pr[M_i(c(x_i)) = z_i]$$

---

[a] The measurement has $d_2$ outcomes, where $d_2$ is the dimension of output quantum state.

We need the following theorem to prove the correctness of Algorithm 3. First we state the result 1 of [28] (lemma 4 of arxiv version):

▶ **Theorem 11** (random orthonormal measurement [28]). *Let $\sigma_1$, $\sigma_2$ be two density matrices in $\mathbb{C}^d$. Define $r := \mathrm{rank}(\sigma_1 - \sigma_2)$. There exists a universal constant $k > 0$ such that if $r < k\sqrt{d}$ then with probability at least $1 - \exp(-kd/r)$ over the choice of a random orthonormal measurement basis $M$ in $\mathbb{C}^d$, $\|M(\sigma_1) - M(\sigma_2)\|_1 > k \|\sigma_1 - \sigma_2\|_F$.* [17]

Note that if $\sigma_1$, $\sigma_2$ are pure states, $r < 2 < k\sqrt{n}$ for large enough $n$ and $\|\sigma_1 - \sigma_2\|_1 \leq \sqrt{2} \|\sigma_1 - \sigma_2\|_F$ so that $\Delta_{tr}(M(\sigma_1), M(\sigma_2)) > k/\sqrt{2}\Delta_{tr}(\sigma_1, \sigma_2)$.

The following lemma shows how trace distance of the measured result grows when we take multiple samples.

---

[17] Recall that $M(\sigma)$ is the output distribution of the measurement $M$ on state $\sigma$.

▶ **Lemma 12** (trace distance amplification). *Let $X_1, X_2, \ldots, X_T$ be $T$ independent distributions and so are $Y_1, Y_2, \ldots, Y_T$. Denote the joint distribution $(X_1, X_2, \ldots, X_T)$ as $X$ and $(Y_1, Y_2, \ldots, Y_T)$ as $Y$. Suppose that*

$$\sum_i \Delta_{tr}(X_i, Y_i) = T\epsilon, \tag{18}$$

*then*

$$\Delta_{tr}(X, Y) \geq 1 - 2^{-\Omega(T\epsilon^2)} \tag{19}$$

**Proof.** By Cauchy-Schwarz inequality,

$$\sum_i (\Delta_{tr}(X_i, Y_i))^2 \geq T\epsilon^2, \tag{20}$$

Then the joint fidelity is bounded by

$$
\begin{aligned}
F(X, Y) &= \Pi_i F(X_i, Y_i) \\
&\leq \Pi_i \sqrt{1 - (\Delta_{tr}(X_i, Y_i))^2} \\
&\leq \exp\left[-\frac{1}{2} \sum_i (\Delta_{tr}(X_i, Y_i))^2\right] = 2^{-\Omega(T\epsilon^2)},
\end{aligned}
\tag{21}
$$

where the last inequality is true because $1 - x \leq e^{-x}$. And the joint trace distance is

$$\Delta_{tr}(X, Y) \geq 1 - F(X, Y) = 1 - 2^{-\Omega(T\epsilon^2)}. \tag{22}$$

◀

The following lemma analyzes the effectiveness of maximum likelihood estimate.

▶ **Lemma 13.** *For any two distributions $D, D^*$ have total variation distance $\alpha$, $\Pr_{i \sim D^*}(D(i) \leq D^*(i)) \geq \alpha$*

**Proof.**

$$
\begin{aligned}
0 &\leq \sum_{i:D(i) \leq D^*(i)} D(i) \\
&= \sum_{i:D(i) \leq D^*(i)} D(i) - D^*(i) + \sum_{i:D(i) \leq D^*(i)} D^*(i) \\
&= \frac{1}{2}\left[\sum_{i:D(i) \leq D^*(i)} (D(i) - D^*(i)) + \sum_{i:D^*(i) \leq D(i)} (D^*(i) - D(i))\right] + \sum_{i:D(i) \leq D^*(i)} D^*(i) \\
&= -\alpha + \Pr_{i \sim D^*}(D(i) \leq D^*(i)) \\
\Rightarrow \Pr_{i \sim D^*}&(D(i) \leq D^*(i)) \geq \alpha
\end{aligned}
\tag{23}
$$

The third line is true because $\sum_{i:D(i) \leq D^*(i)} (D(i) - D^*(i)) = \sum_{i:D^*(i) \leq D(i)} (D^*(i) - D(i))$. ◀

We think $D^*$ as the correct distribution and $D$ is a distribution far away, with the total variation distance between them being $\alpha = 1 - \epsilon$. When we use maximum likelihood estimation to distinguish $D^*$ from $D$, Lemma 13 says that the probability of error is less than $\epsilon$. Now we are ready to prove theorem 10.

**Proof.** Let $c^*$ be the target concept, and $c$ a concept such that $\Delta(c^*, c) > \epsilon$. Recall that we took

$$T = \Theta\left(\frac{\log|C| + \log(1/\delta)}{\epsilon^2}\right)$$

samples. For all $i \in [T]$, apply Theorem 11 to the pair of states $(c^*(x_i), c(x_i))$, we get that with probability $1 - \exp(-kd_2/2)$ over random orthonormal measurements $M_i$,

$$\Delta_{tr}(M_i(c^*(x_i)), M_i(c(x_i))) > k/\sqrt{2}\Delta_{tr}(c^*(x_i), c(x_i)), \tag{24}$$

where $k$ is a universal constant. Since you can pad some ancilla states to increase $d_2$ without changing trace distances if $\exp(-kd_2/2)$ is not small enough, we ignore this term. By Chernoff bound, with probability at least $1 - 2^{-\Omega(T\epsilon)}$ over $\{x_i\}$ sampled from $D$,

$$(1/T) \cdot \sum_i \Delta_{tr}(M_i(c^*(x_i)), M_i(c(x_i))) > (1/T) \cdot \sum_i k/\sqrt{2}\Delta_{tr}(c^*(x_i), c(x_i)) \geq \frac{k}{2\sqrt{2}}\epsilon. \tag{25}$$

So we can apply Lemma 12 to get that with probability at least $1 - 2^{-\Omega(T\epsilon)}$,

$$[\Delta_{tr}(\{M_i(c^*(x_i))\}, \{M_i(h(x_i))\})] \geq 1 - 2^{-\Omega(T\epsilon^2)}. \tag{26}$$

Now, note that by Lemma 13, the probability that the maximal likelihood estimation (incorrectly) selects $c$ is at most $(2^{-\Omega(T\epsilon^2)} + 2^{-\Omega(T\epsilon)})$. By taking a union bound over all such $c$, we get

$$\Pr[\Delta(c^*, h) > \epsilon] \leq (2^{-\Omega(T\epsilon^2)} + 2^{-\Omega(T\epsilon)}) \cdot |C| \leq \delta. \tag{27}$$

◀

## B    Lower Bounds

In this section we describe two simple lower bounds. One is an $\Omega((1-\delta)\ln|C|/\epsilon^2)/\ln(\ln|C|/\epsilon)$ lower bound on the sample complexity of approximate state discrimination for pure states, which in turn gives lower bounds on the sample complexity of PAC learning quantum channels. The other is an $\Omega(\sqrt{d})$ lower bound on the sample complexity of learning large dimensional *classical distribution* in the *agnostic* model, which in turn lower bounds approximate state discrimination and PAC learning quantum state in the agnostic model [18, 21].

### B.1    Lower Bound for Pure State Case

▶ **Theorem 14.** *The sample complexity of $(\epsilon, \delta)$-approximate state discrimination on a set $C$ of pure states is $\Omega((1-\delta)\ln|C|/\epsilon^2)/\ln(\ln|C|/\epsilon)$.*

**Proof.** This lower bound uses the $\epsilon$-packing-net construction of [15]. In Lemma 5 of the arxiv version of [15], the authors showed the existence of a set $C$ of $d$-dimensional pure states with the following three properties: the distance between each state is at least $\epsilon$, the Holevo information $\chi_0$ for states uniformly drawn from the set is $O(\epsilon^2 \ln(d/\epsilon))$, and $\ln|C| = \Omega(d)$. With a simple reduction to communication protocol and Holevo theorem, [15] showed that to distinguish states in $C$ with probability $\delta$, $\frac{(1-\delta)\ln|C| - \ln 2}{\chi_0} = \Omega((1-\delta)\ln|C|/\epsilon^2)/\ln(\ln|C|/\epsilon)$ samples are required. Since every state in $C$ is $\epsilon$-far from each other, an $(\epsilon, \delta)$-approximate

discriminator should be able to distinguish each state in $C$ with probability $\delta$, therefore the discriminator must take $\Omega((1-\delta)\ln|C|/\epsilon^2)/\ln(\ln|C|/\epsilon)$ samples. This matches the sample complexity of our pure state algorithm in terms of $\epsilon$ and $|C|$ with some logarithmic factors. ◄

▶ **Remark 15.** Unfortunately, running the same argument with the mixed state $\epsilon$-packing nets of [15] does not give us tighter lower bound, so we don't have a matching lower bound for the mixed state case.

▶ **Corollary 16.** *The proper $(\epsilon, \delta)$-PAC sample complexity of a concept class $C$ of pure states is $\Omega((1-\delta)\ln|C|/\epsilon^2)/\ln(\ln|C|/\epsilon)$.*

## B.2 Agnostic Model

Agnostic model [18, 21] is a learning model related to the PAC model. In agnostic model, the target concept does not need to come from the concept class. We formally define the agnostic model for learning quantum channels as follows:

We consider a learner trying to learn a target concept $c^*$ with respect to an unknown distribution $D$. The learner is also given a concept class $C$. Since the target concept might not be in the concept class $C$, the learner tries the output the concept $c_{opt}$ that minimize the distance to the target concept $c^*$. Here, we consider the concept class $C$ as a finite set of known $d_1$ to $d_2$ dimensional quantum channels, and $D$ as a distribution over the Hilbert space of dimension $d_1$. Precisely, the learner is given access to a sample oracle $\mathcal{O}_{c^*, D}$ and outputs a hypothesis $h \in C$. The oracle $\mathcal{O}_{c^*, D}$ generates i.i.d. samples $(x_i, c^*(x_i))$, where each $x_i \leftarrow D$ is the classical description of a state drawn according to the distribution $D$, and $c^*(x_i)$ is the (potentially mixed) quantum state outputted by $c^*$ on input $x_i$.

The distance between two concepts $c$ and $h$ under the distribution $D$ is the expected trace distance to the target concept $\Delta(c, h) = \mathbb{E}_{x \in D}[\Delta_{tr}(c(x), h(x))]$. Let $c_{opt}$ be the optimal output, $c_{opt} = \arg\min[\Delta(c, c^*)|c \in C]$. The goal of the learner is to find a hypothesis $h \in C$ with $\Delta(c^*, h) \leq \Delta(c^*, c_{opt}) + \epsilon$.

A quantum learning algorithm $A$ is a $(\epsilon, \delta)$-*agnostic learner* for $C$ if the following holds: For every $c^*$ and distribution $D$, given oracle access to $\mathcal{O}_{c^*, D}$, $A^{\mathcal{O}_{c^*, D}}$ outputs an $h \in C$ such that $\Delta(c^*, h) \leq \Delta(c^*, c_{opt}) + \epsilon$ with probability at least $1 - \delta$. The sample complexity of $A$ is the maximum number of samples $T$ that $A$ needs to query $\mathcal{O}_{c^*, D}$ to output $h$. The $(\epsilon, \delta)$-*agnostic sample complexity* of a concept class $C$ is the minimum sample complexity over all learners.

We show that there is no efficient quantum agnostic learner in the following theorem.

▶ **Theorem 17.** *For all $\epsilon < \frac{1}{10}$ and positive integer $d$, there exist a concept class $C$ of dimension 0 to $d$ whose $(\epsilon, \delta)$-agnostic sample complexity is $\Omega(\sqrt{d})$.*

**Proof.** We can get the $\Omega(\sqrt{d})$ lower bound with a simple concept class of that only has two concepts. Both of the concepts are constant channels that output classical distributions. Consider distributions on $d+1$ dimensions $e_0, e_1, \ldots, e_d$. The first concept $C_1$ has all weight on $e_0$. The second concept $C_2$ has weight uniformly distributed over $e_1, \ldots, e_d$. Now consider the following two set of distribution to be learned. $D_1 = \{D_{1,i}\}$ has weight $1/3$ on $e_0$ and weight $1/d$ on $2/3$ of dimensions $e_1 \ldots, e_d$. Anything in $D_1$ has distance $2/3$ to $C_1$ and distance $1/3$ to $C_2$, so it should be learned as $C_2$. $D_2 = \{D_{2,i}\}$ has weight $1/3$ on $e_0$ and weight $100/d$ on $2/300$ of dimensions $e_1 \ldots, e_d$. Anything in $D_2$ has distance $2/3$ to $C_1$ and distance $(1/3 + 99 * 2/300 + 1 * (1 - 2/300))/2 \sim 0.993$ to $C_2$, so it should be learned as $C_1$. However, $D_1$ and $D_2$ both looks pretty much like a uniform distribution on $e_1, \ldots, e_d$. To

distinguish them we need to see a collision on $e_1, \ldots, e_d$. By a standard birthday bound, we need at least $\Omega(\sqrt{d/100})$ samples to see a collision. Therefore we need $\Omega(\sqrt{d})$ samples to learn classical distributions in agnostic model with constant error. In the regime of $|C| = \text{poly}(d)$, the lower bound means that it's impossible to find an efficient algorithm of sample complexity $O(\text{polylog}\,|C|)$. ◄

▶ Remark 18. Note that the construction of Theorem 17 is based on a classical distribution, so it means that agnostic learning of a classical distribution of many outputs efficiently is also impossible. To the knowledge of the authors, agnostic learning of classical distribution of many output has not been studied in the literature. Also note that classical distribution is not a subclass of pure quantum states, so Theorem 17 does not rule out that quantum channel with *pure* state outcomes can be efficiently agnostically learned.

▶ Remark 19. As mentioned in section 1.2.0.1, [7] studied the problem of quantum hypothesis selection, which can be viewed as a a relax version of agnostic learning, outputting a state $h$ such that $\Delta(c^*, h) \leq 3.01\Delta(c^*, c_{opt}) + \epsilon$.

## C     PGM for "average BSD"

Please refer to the full version at `https://arxiv.org/abs/1810.10938`.