

Global Optimum Search in Quantum Deep Learning

Lanston Hau Man Chu, Tejas Bhojraj, Rui Huang

Project of CS 880: Quantum Computing

Instructor: Dieter van Melkebeek

May 8, 2020

Overview

- 1 Introduction and Theory
 - Gradient Descent and Its Drawbacks
 - Overview of the Two Approaches
- 2 Approach 1: Average Approach
- 3 Approach 2: Partial Swap Test Cut-off Method (PSTC)
- 4 Equivalency Discussion
- 5 Summary
- 6 Extensions and Future Work
 - Additional Restriction
 - Adversarial Example
- 7 Notations Table

Gradient Descent

Gradient Descent

- Gradient descent (GD) is an essential algorithm widely used in machine learning to optimize the objective function

Gradient Descent

- Gradient descent (GD) is an essential algorithm widely used in machine learning to optimize the objective function
- Key idea: Move θ iteratively towards the direction of the steepest gradient to reach the optimum

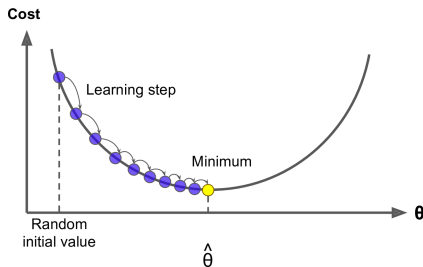
Gradient Descent

- Gradient descent (GD) is an essential algorithm widely used in machine learning to optimize the objective function
- Key idea: Move θ iteratively towards the direction of the steepest gradient to reach the optimum

$$\theta^* = \operatorname{argmin}_{\theta \in \Theta} \mathcal{L}(\theta)$$

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{j=1}^N \ell(\theta, x_j)$$

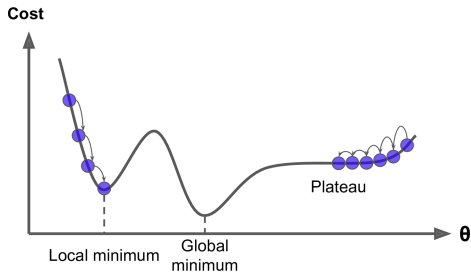
$$\theta_{k+1} = \theta_k - \eta \left. \frac{\partial \mathcal{L}_\theta}{\partial \theta} \right|_{\theta=\theta_k}$$



Drawbacks of Gradient Descent

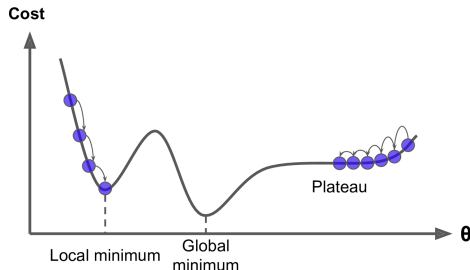
Drawbacks of Gradient Descent

- Depending on the initial point, GD may find a local optimum instead of the global optimum.



Drawbacks of Gradient Descent

- Depending on the initial point, GD may find a local optimum instead of the global optimum.



- The upper bound or expected number of iterations to reach convergence is difficult to determine.

Key Contributions

Key Contributions

- Two quantum approaches (Average Approach and PSTC Approach) to find the global optimum (instead of a local optimum) for optimizing machine learning models.

Key Contributions

- Two quantum approaches (Average Approach and PSTC Approach) to find the global optimum (instead of a local optimum) for optimizing machine learning models.
- Theoretical analyses to show that the expected cost for both approaches are $O(\sqrt{|\Theta|}N)$.

Key Contributions

- Two quantum approaches (Average Approach and PSTC Approach) to find the global optimum (instead of a local optimum) for optimizing machine learning models.
- Theoretical analyses to show that the expected cost for both approaches are $O(\sqrt{|\Theta|}N)$.
- Novel objective function maximizing the number of “cut-off indicators \mathbb{E}_{θ_j} ” to fit the property of quantum computing in optimizing machine learning models

Key Contributions

- Two quantum approaches (Average Approach and PSTC Approach) to find the global optimum (instead of a local optimum) for optimizing machine learning models.
- Theoretical analyses to show that the expected cost for both approaches are $O(\sqrt{|\Theta|}N)$.
- Novel objective function maximizing the number of “cut-off indicators \mathbb{E}_{θ_j} ” to fit the property of quantum computing in optimizing machine learning models
- Potential for PSTC to reduce the cost further to $O(\sqrt{|\Theta|} \cdot \textit{sublinear}(N))$ in future work

Overview of Average Approach

$$|\theta\rangle|j\rangle|\mathbf{0}\rangle \xrightarrow{\text{Quantum parallelism with cost } O(N)} |\theta\rangle|j\rangle \left| \sum_{j=1}^N \ell(\theta, x_j) \right\rangle$$

$$\theta_{avg}^* = \operatorname{argmin}_{\theta \in \Theta} \mathcal{L}_{\theta}^{avg} = \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{j=1}^N \ell(\theta, x_j)$$

Overview of Average Approach

$$|\theta\rangle|j\rangle|\mathbf{0}\rangle \xrightarrow{\text{Quantum parallelism with cost } O(N)} |\theta\rangle|j\rangle\left|\sum_{j=1}^N \ell(\theta, x_j)\right\rangle$$

$$\theta_{avg}^* = \operatorname{argmin}_{\theta \in \Theta} \mathcal{L}_{\theta}^{avg} = \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{j=1}^N \ell(\theta, x_j)$$

- Treat the machine learning model as a quantum blackbox

Overview of Average Approach

$$|\theta\rangle|j\rangle|\mathbf{0}\rangle \xrightarrow{\text{Quantum parallelism with cost } O(N)} |\theta\rangle|j\rangle\left|\sum_{j=1}^N \ell(\theta, x_j)\right\rangle$$

$$\theta_{avg}^* = \operatorname{argmin}_{\theta \in \Theta} \mathcal{L}_{\theta}^{avg} = \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{j=1}^N \ell(\theta, x_j)$$

- Treat the machine learning model as a quantum blackbox
- Utilize Durr & Hoyer (DH) algorithm to find the global minimum

Overview of Average Approach

$$|\theta\rangle|j\rangle|\mathbf{0}\rangle \xrightarrow{\text{Quantum parallelism with cost } O(N)} |\theta\rangle|j\rangle\left|\sum_{j=1}^N \ell(\theta, x_j)\right\rangle$$

$$\theta_{avg}^* = \operatorname{argmin}_{\theta \in \Theta} \mathcal{L}_{\theta}^{avg} = \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{j=1}^N \ell(\theta, x_j)$$

- Treat the machine learning model as a quantum blackbox
- Utilize Durr & Hoyer (DH) algorithm to find the global minimum
- Expected number of steps: $O(N\sqrt{|\Theta|})$

Overview of Partial Swap Test Cut-off (PSTC) Approach

$$|\theta\rangle|j\rangle|\mathbf{0}\rangle \xrightarrow{\text{Quantum parallelism with cost } O(1)} |\theta\rangle|j\rangle|E_{\theta j}\rangle$$

$$\theta_{\text{PSTC}}^* = \operatorname{argmax}_{\theta \in \Theta} \mathcal{L}_{\theta}^{\text{PSTC}} = \operatorname{argmax}_{\theta \in \Theta} \frac{1}{N} \sum_{j=1}^N \mathbb{1}[\ell(\theta, x_j) \leq \ell_{\text{threshold}}]$$

Overview of Partial Swap Test Cut-off (PSTC) Approach

$$|\theta\rangle|j\rangle|\mathbf{0}\rangle \xrightarrow{\text{Quantum parallelism with cost } O(1)} |\theta\rangle|j\rangle|E_{\theta j}\rangle$$

$$\theta_{\text{PSTC}}^* = \operatorname{argmax}_{\theta \in \Theta} \mathcal{L}_{\theta}^{\text{PSTC}} = \operatorname{argmax}_{\theta \in \Theta} \frac{1}{N} \sum_{j=1}^N \mathbb{1}[\ell(\theta, x_j) \leq \ell_{\text{threshold}}]$$

- The cut-off indicator $E_{\theta j} \triangleq \mathbb{1}[\ell(\theta, x_j) \leq \ell_{\text{threshold}}]$.

Overview of Partial Swap Test Cut-off (PSTC) Approach

$$|\theta\rangle|j\rangle|\mathbf{0}\rangle \xrightarrow{\text{Quantum parallelism with cost } O(1)} |\theta\rangle|j\rangle|E_{\theta j}\rangle$$

$$\theta_{\text{PSTC}}^* = \operatorname{argmax}_{\theta \in \Theta} \mathcal{L}_{\theta}^{\text{PSTC}} = \operatorname{argmax}_{\theta \in \Theta} \frac{1}{N} \sum_{j=1}^N \mathbb{1}[\ell(\theta, x_j) \leq \ell_{\text{threshold}}]$$

- The cut-off indicator $E_{\theta j} \triangleq \mathbb{1}[\ell(\theta, x_j) \leq \ell_{\text{threshold}}]$.
- Employ a “partial” swap test and amplitude amplification to boost the probability of finding a good θ

Overview of Partial Swap Test Cut-off (PSTC) Approach

$$|\theta\rangle|j\rangle|0\rangle \xrightarrow{\text{Quantum parallelism with cost } O(1)} |\theta\rangle|j\rangle|E_{\theta j}\rangle$$

$$\theta_{\text{PSTC}}^* = \operatorname{argmax}_{\theta \in \Theta} \mathcal{L}_{\theta}^{\text{PSTC}} = \operatorname{argmax}_{\theta \in \Theta} \frac{1}{N} \sum_{j=1}^N \mathbb{1}[\ell(\theta, x_j) \leq \ell_{\text{threshold}}]$$

- The cut-off indicator $E_{\theta j} \triangleq \mathbb{1}[\ell(\theta, x_j) \leq \ell_{\text{threshold}}]$.
- Employ a “partial” swap test and amplitude amplification to boost the probability of finding a good θ
- Still $O(N\sqrt{|\Theta|})$, but only at checking process

Overview of Partial Swap Test Cut-off (PSTC) Approach

$$|\theta\rangle|j\rangle|0\rangle \xrightarrow{\text{Quantum parallelism with cost } O(1)} |\theta\rangle|j\rangle|E_{\theta j}\rangle$$

$$\theta_{\text{PSTC}}^* = \operatorname{argmax}_{\theta \in \Theta} \mathcal{L}_{\theta}^{\text{PSTC}} = \operatorname{argmax}_{\theta \in \Theta} \frac{1}{N} \sum_{j=1}^N \mathbb{1}[\ell(\theta, x_j) \leq \ell_{\text{threshold}}]$$

- The cut-off indicator $E_{\theta j} \triangleq \mathbb{1}[\ell(\theta, x_j) \leq \ell_{\text{threshold}}]$.
- Employ a “partial” swap test and amplitude amplification to boost the probability of finding a good θ
- Still $O(N\sqrt{|\Theta|})$, but only at checking process
- Enable future improvement to maybe $O(\text{sublinear}(N) \cdot \sqrt{|\Theta|})$

Overview

- 1 Introduction and Theory
 - Gradient Descent and Its Drawbacks
 - Overview of the Two Approaches
- 2 **Approach 1: Average Approach**
- 3 Approach 2: Partial Swap Test Cut-off Method (PSTC)
- 4 Equivalency Discussion
- 5 Summary
- 6 Extensions and Future Work
 - Additional Restriction
 - Adversarial Example
- 7 Notations Table

Approach 1: Average Approach

- ① Consider the function $f(i) = f(\theta_i) := \sum_j \ell(\theta_i, x_j)$
- ② Use the DH algorithm [2] to find the minimum of the set $\{f(i)\}_i$ with a slight modification: When the DH algorithm queries $f(i)$ for some i (i.e. when the algorithm makes a query on $|i\rangle|0^t\rangle$ and needs $|i\rangle|f(i)\rangle$ as output), do the following subroutine:
 - Classically compute $f(i)$

Approach 1: Average Approach

- ① Consider the function $f(i) = f(\theta_i) := \sum_j \ell(\theta_i, x_j)$
- ② Use the DH algorithm [2] to find the minimum of the set $\{f(i)\}_i$ with a slight modification: When the DH algorithm queries $f(i)$ for some i (i.e. when the algorithm makes a query on $|i\rangle|0^t\rangle$ and needs $|i\rangle|f(i)\rangle$ as output), do the following subroutine:
 - Classically compute $f(i)$
 - Build a unitary operator, U mapping $|i\rangle|0^t\rangle \mapsto |i\rangle|f(i)\rangle$.

Approach 1: Average Approach

- ① Consider the function $f(i) = f(\theta_i) := \sum_j \ell(\theta_i, x_j)$
- ② Use the DH algorithm [2] to find the minimum of the set $\{f(i)\}_i$ with a slight modification: When the DH algorithm queries $f(i)$ for some i (i.e. when the algorithm makes a query on $|i\rangle|0^t\rangle$ and needs $|i\rangle|f(i)\rangle$ as output), do the following subroutine:
 - Classically compute $f(i)$
 - Build a unitary operator, U mapping $|i\rangle|0^t\rangle \mapsto |i\rangle|f(i)\rangle$.
 - Return $U(|i\rangle|0^t\rangle)$ to the DH algorithm.

Analysis

- Each time we invoke the subroutine, we need N classical queries to ℓ .

Analysis

- Each time we invoke the subroutine, we need N classical queries to ℓ .
- The DH algorithm makes an expected number of $O(\sqrt{|\Theta|})$ queries.
So, we run the subroutine $O(\sqrt{|\Theta|})$ times on average.

Analysis

- Each time we invoke the subroutine, we need N classical queries to ℓ .
- The DH algorithm makes an expected number of $O(\sqrt{|\Theta|})$ queries.
So, we run the subroutine $O(\sqrt{|\Theta|})$ times on average.
- So, the cost is $O(\sqrt{|\Theta|}N)$.

Analysis

- Each time we invoke the subroutine, we need N classical queries to ℓ .
- The DH algorithm makes an expected number of $O(\sqrt{|\Theta|})$ queries. So, we run the subroutine $O(\sqrt{|\Theta|})$ times on average.
- So, the cost is $O(\sqrt{|\Theta|}N)$.
- Also, we know that the DH algorithm succeeds with probability 0.5 and so by running it many times, we get the parameter minimizing the loss with arbitrarily good accuracy at the cost $O(\sqrt{|\Theta|}N)$.

Overview

- 1 Introduction and Theory
 - Gradient Descent and Its Drawbacks
 - Overview of the Two Approaches
- 2 Approach 1: Average Approach
- 3 Approach 2: Partial Swap Test Cut-off Method (PSTC)
- 4 Equivalency Discussion
- 5 Summary
- 6 Extensions and Future Work
 - Additional Restriction
 - Adversarial Example
- 7 Notations Table

Motivation: Cost reduction

Motivation: Cost reduction

- Average approach: $\sum_{j=1}^N |\theta\rangle|j\rangle|\mathbf{0}\rangle \xrightarrow{\text{Cost } O(N)} |\theta\rangle|j\rangle|\sum_{j=1}^N \ell(\theta, x_j)\rangle$

Motivation: Cost reduction

- Average approach: $\sum_{j=1}^N |\theta\rangle|j\rangle|\mathbf{0}\rangle \xrightarrow{\text{Cost } O(N)} |\theta\rangle|j\rangle|\sum_{j=1}^N \ell(\theta, x_j)\rangle$
- **Want** some function f : $\sum_{j=1}^N |\theta\rangle|j\rangle|\mathbf{0}\rangle \xrightarrow{\text{Cost } O(1)} \sum_{j=1}^N |\theta\rangle|j\rangle|f(\theta, x_j)\rangle$

Motivation: Cost reduction

- Average approach: $\sum_{j=1}^N |\theta\rangle|j\rangle|\mathbf{0}\rangle \xrightarrow{\text{Cost } O(N)} |\theta\rangle|j\rangle \left| \sum_{j=1}^N \ell(\theta, x_j) \right\rangle$
- **Want** some function f : $\sum_{j=1}^N |\theta\rangle|j\rangle|\mathbf{0}\rangle \xrightarrow{\text{Cost } O(1)} \sum_{j=1}^N |\theta\rangle|j\rangle |f(\theta, x_j)\rangle$
- **Cut-off indicator** $E_{\theta j} \triangleq \mathbb{1}[\ell(\theta, x_j) \leq \ell_{threshold}]$ (1-qubit)

$$\boxed{|\theta\rangle|j\rangle|\mathbf{0}\rangle \xrightarrow{\text{Cost } O(1)} |\theta\rangle|j\rangle |E_{\theta j}\rangle}$$

Motivation: Cost reduction

- Average approach: $\sum_{j=1}^N |\theta\rangle|j\rangle|\mathbf{0}\rangle \xrightarrow{\text{Cost } O(N)} |\theta\rangle|j\rangle|\sum_{j=1}^N \ell(\theta, x_j)\rangle$
- **Want** some function f : $\sum_{j=1}^N |\theta\rangle|j\rangle|\mathbf{0}\rangle \xrightarrow{\text{Cost } O(1)} \sum_{j=1}^N |\theta\rangle|j\rangle|f(\theta, x_j)\rangle$
- **Cut-off indicator** $E_{\theta j} \triangleq \mathbb{1}[\ell(\theta, x_j) \leq \ell_{threshold}]$ **(1-qubit)**

$$\boxed{|\theta\rangle|j\rangle|\mathbf{0}\rangle \xrightarrow{\text{Cost } O(1)} |\theta\rangle|j\rangle|E_{\theta j}\rangle}$$

- **Easy** summation (by inner product):

$$\left(\sum_{j=1}^N \langle j | \langle E_{\theta j} | \right) \left(\sum_{k=1}^N |k\rangle |1\rangle \right) = \sum_{j=1}^N \langle j | j \rangle \langle E_{\theta j} | 1 \rangle = \sum_{j=1}^N E_{\theta j}$$

Motivation: Cost reduction

- Average approach: $\sum_{j=1}^N |\theta\rangle|j\rangle|\mathbf{0}\rangle \xrightarrow{\text{Cost } O(N)} |\theta\rangle|j\rangle|\sum_{j=1}^N \ell(\theta, x_j)\rangle$
- **Want** some function f : $\sum_{j=1}^N |\theta\rangle|j\rangle|\mathbf{0}\rangle \xrightarrow{\text{Cost } O(1)} \sum_{j=1}^N |\theta\rangle|j\rangle|f(\theta, x_j)\rangle$
- **Cut-off indicator** $E_{\theta j} \triangleq \mathbb{1}[\ell(\theta, x_j) \leq \ell_{threshold}]$ **(1-qubit)**

$$\boxed{|\theta\rangle|j\rangle|\mathbf{0}\rangle \xrightarrow{\text{Cost } O(1)} |\theta\rangle|j\rangle|E_{\theta j}\rangle}$$

- **Easy** summation (by inner product):

$$\left(\sum_{j=1}^N \langle j | \langle E_{\theta j} | \right) \left(\sum_{k=1}^N |k\rangle |1\rangle \right) = \sum_{j=1}^N \langle j | j \rangle \langle E_{\theta j} | 1 \rangle = \sum_{j=1}^N E_{\theta j}$$

- New objective function (want **larger** $\mathcal{L}_{\theta}^{\text{PSTC}}$):

$$\theta_{\text{PSTC}}^* = \operatorname{argmax}_{\theta \in \Theta} \mathcal{L}_{\theta}^{\text{PSTC}} = \operatorname{argmax}_{\theta \in \Theta} \frac{1}{N} \sum_{j=1}^N \mathbb{1}[\ell(\theta, x_j) \leq \ell_{threshold}]$$

Partial Swap Test ($Q_{1\text{-query}}$)

Partial Swap Test ($Q_{1\text{-query}}$)

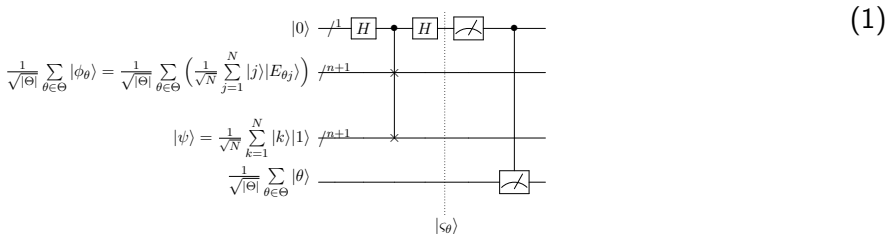
- Want inner product: $\left(\sum_{j=1}^N \langle j | \langle E_{\theta j} | \right) \left(\sum_{k=1}^N |k\rangle |1\rangle \right) = \sum_{j=1}^N E_{\theta j}$

$$\text{Define } \begin{cases} |\phi_{\theta}\rangle \triangleq \frac{1}{\sqrt{N}} \sum_{j=1}^N |j\rangle |E_{\theta j}\rangle \\ |\psi\rangle \triangleq \frac{1}{\sqrt{N}} \sum_{k=1}^N |k\rangle |1\rangle \end{cases}$$

Partial Swap Test ($Q_{1\text{-query}}$)

- Want inner product: $\left(\sum_{j=1}^N \langle j | \langle E_{\theta_j} | \right) \left(\sum_{k=1}^N |k\rangle |1\rangle \right) = \sum_{j=1}^N E_{\theta_j}$

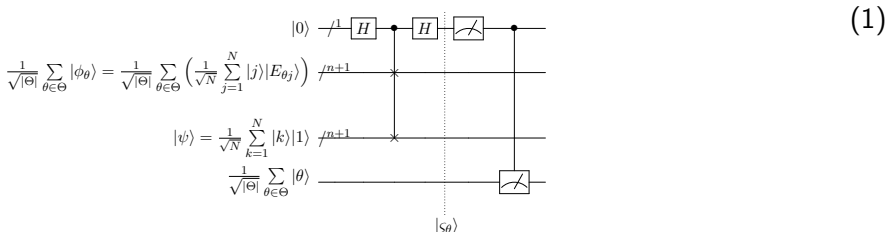
$$\text{Define } \begin{cases} |\phi_{\theta}\rangle \triangleq \frac{1}{\sqrt{N}} \sum_{j=1}^N |j\rangle |E_{\theta_j}\rangle \\ |\psi\rangle \triangleq \frac{1}{\sqrt{N}} \sum_{k=1}^N |k\rangle |1\rangle \end{cases}$$



Partial Swap Test ($Q_{1\text{-query}}$)

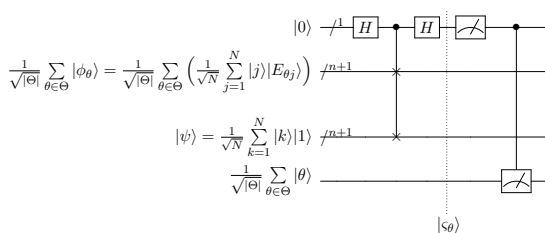
- Want inner product: $\left(\sum_{j=1}^N \langle j | \langle E_{\theta_j} | \right) \left(\sum_{k=1}^N |k\rangle |1\rangle \right) = \sum_{j=1}^N E_{\theta_j}$

Define $\begin{cases} |\phi_{\theta}\rangle \triangleq \frac{1}{\sqrt{N}} \sum_{j=1}^N |j\rangle |E_{\theta_j}\rangle \\ |\psi\rangle \triangleq \frac{1}{\sqrt{N}} \sum_{k=1}^N |k\rangle |1\rangle \end{cases}$



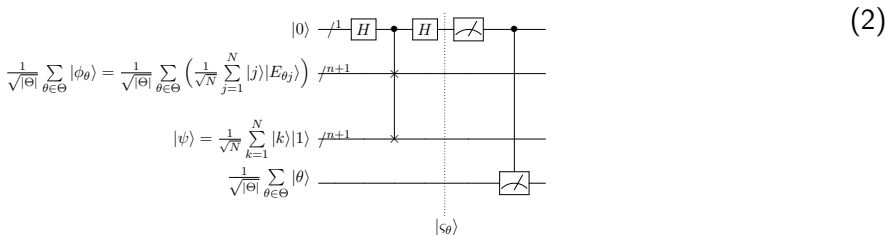
$$\frac{1}{\sqrt{|\Theta|}} \sum_{\theta \in \Theta} |s_{\theta}\rangle |\theta\rangle = \frac{1}{\sqrt{|\Theta|}} \sum_{\theta \in \Theta} \frac{1}{2} |0\rangle (|\phi_{\theta}\rangle |\psi\rangle + |\psi\rangle |\phi_{\theta}\rangle) |\theta\rangle + \frac{1}{2} |1\rangle (|\phi_{\theta}\rangle |\psi\rangle - |\psi\rangle |\phi_{\theta}\rangle) |\theta\rangle$$

Partial Swap Test ($Q_{1\text{-query}}$)



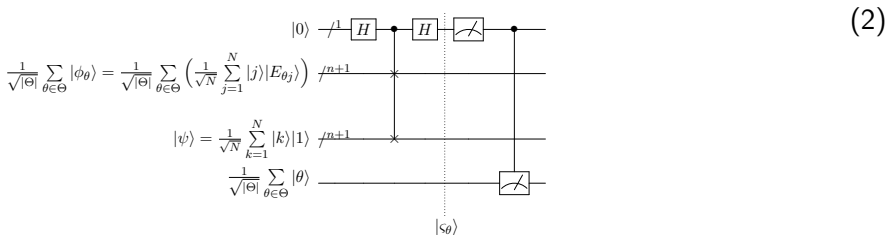
(2)

Partial Swap Test ($Q_{1\text{-query}}$)



$$\Rightarrow \text{Lemma 2.2: } \mathbb{P}[\text{1st qubit} = 0] = \frac{1}{2} + \frac{1}{2|\Theta|} \sum_{\theta \in \Theta} |\langle \phi_\theta, \psi \rangle|^2$$

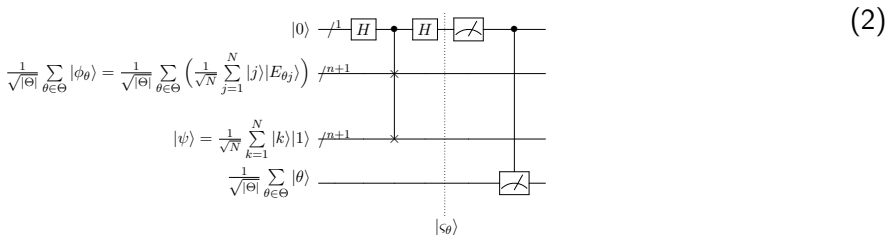
Partial Swap Test ($Q_{1\text{-query}}$)



$$\Rightarrow \text{Lemma 2.2: } \mathbb{P}[\text{1st qubit} = 0] = \frac{1}{2} + \frac{1}{2|\Theta|} \sum_{\theta \in \Theta} |\langle \phi_\theta, \psi \rangle|^2$$

$$\Rightarrow \boxed{\text{Theorem 2.4: } \mathbb{P}[\text{observe } \theta | \text{1st qubit} = 0] \propto \frac{1}{2} + \frac{1}{2} \left| \frac{1}{N} \sum_{j=1}^N E_{\theta j} \right|^2} = \frac{1}{2} + \frac{1}{2} |\mathcal{L}_\theta^{\text{PSTC}}|^2$$

Partial Swap Test ($Q_{1\text{-query}}$)

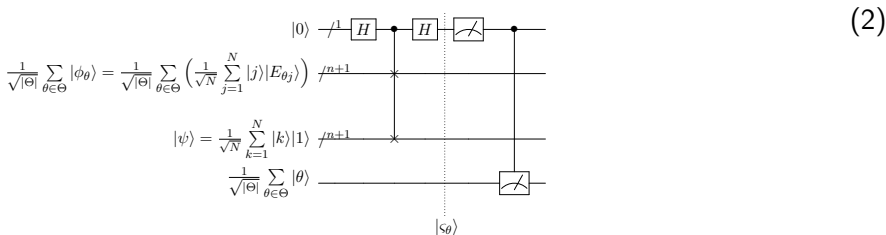


$$\Rightarrow \text{Lemma 2.2: } \mathbb{P}[\text{1st qubit} = 0] = \frac{1}{2} + \frac{1}{2|\Theta|} \sum_{\theta \in \Theta} |\langle \phi_\theta, \psi \rangle|^2$$

$$\Rightarrow \boxed{\text{Theorem 2.4: } \mathbb{P}[\text{observe } \theta | \text{1st qubit} = 0] \propto \frac{1}{2} + \frac{1}{2} \left| \frac{1}{N} \sum_{j=1}^N E_{\theta,j} \right|^2} = \frac{1}{2} + \frac{1}{2} |\mathcal{L}_\theta^{\text{PSTC}}|^2$$

- Want **larger** $\mathcal{L}_\theta^{\text{PSTC}}$

Partial Swap Test ($Q_{1\text{-query}}$)



$$\Rightarrow \text{Lemma 2.2: } \mathbb{P}[\text{1st qubit} = 0] = \frac{1}{2} + \frac{1}{2|\Theta|} \sum_{\theta \in \Theta} |\langle \phi_\theta, \psi \rangle|^2$$

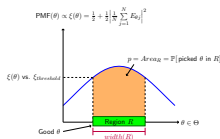
$$\Rightarrow \text{Theorem 2.4: } \mathbb{P}[\text{observe } \theta | \text{1st qubit} = 0] \propto \frac{1}{2} + \frac{1}{2} \left| \frac{1}{N} \sum_{j=1}^N E_{\theta,j} \right|^2 = \frac{1}{2} + \frac{1}{2} |\mathcal{L}_\theta^{\text{PSTC}}|^2$$

- Want **larger** $\mathcal{L}_\theta^{\text{PSTC}}$
- The **better** θ we want, the **higher chance** we can observe θ !

$\mathcal{A}_{\text{Boost}}$: Amplitude Amplification

$\mathcal{A}_{\text{Boost}}$: Amplitude Amplification

Table: Amplitude amplification for the $\text{PMF}(\theta) \propto \xi(\theta)$



Before $\mathcal{A}_{\text{Boost}}$

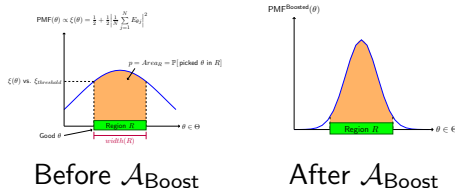


After $\mathcal{A}_{\text{Boost}}$

$$p = \text{Area}_R = \mathbb{P}[\text{picked } \theta \text{ in } R]$$

$\mathcal{A}_{\text{Boost}}$: Amplitude Amplification

Table: Amplitude amplification for the $\text{PMF}(\theta) \propto \xi(\theta)$



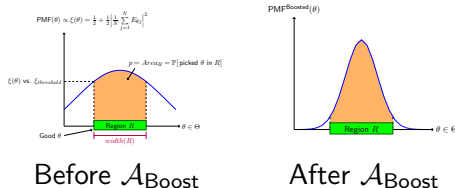
$$p = \text{Area}_R = \mathbb{P}[\text{picked } \theta \text{ in } R]$$

- In the manner \sim Grover's Search [1]

$$\begin{array}{ccc} \#(Q_{1\text{-query}}) = O(\frac{1}{p}) & \xrightarrow{\text{improved}} & \#(Q_{1\text{-query}}) = O(\frac{1}{\sqrt{p}}) \\ \text{via classic algorithm} & & \text{via } \mathcal{A}_{\text{Boost}} \end{array}$$

$\mathcal{A}_{\text{Boost}}$: Amplitude Amplification

Table: Amplitude amplification for the $\text{PMF}(\theta) \propto \xi(\theta)$



$$p = \text{Area}_R = \mathbb{P}[\text{picked } \theta \text{ in } R]$$

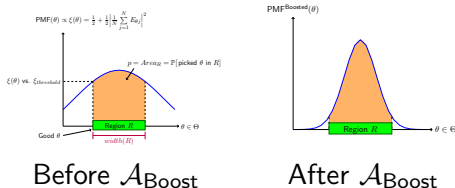
- In the manner \sim Grover's Search [1]

$$\begin{array}{ccc} \#(Q_{1\text{-query}}) = O(\frac{1}{p}) & \xrightarrow{\text{improved}} & \#(Q_{1\text{-query}}) = O(\frac{1}{\sqrt{p}}) \\ \text{via classic algorithm} & & \text{via } \mathcal{A}_{\text{Boost}} \end{array}$$

- Checking process: To determine whether $\theta \in R$

$\mathcal{A}_{\text{Boost}}$: Amplitude Amplification

Table: Amplitude amplification for the $\text{PMF}(\theta) \propto \xi(\theta)$



$$p = \text{Area}_R = \mathbb{P}[\text{picked } \theta \text{ in } R]$$

- In the manner \sim Grover's Search [1]

$$\begin{array}{ccc} \#(Q_{1\text{-query}}) = O\left(\frac{1}{p}\right) & \xrightarrow{\text{improved}} & \#(Q_{1\text{-query}}) = O\left(\frac{1}{\sqrt{p}}\right) \\ \text{via classic algorithm} & & \text{via } \mathcal{A}_{\text{Boost}} \end{array}$$

- Checking process: To determine whether $\theta \in R \leftarrow \text{Cost } O(N)$

$\mathcal{A}_{\text{PSTM}}$

Lemma 2.5: $O(\frac{1}{\sqrt{p}}) \leq O(\sqrt{|\Theta|})$

$\mathcal{A}_{\text{PSTM}}$

Lemma 2.5: $O(\frac{1}{\sqrt{p}}) \leq O(\sqrt{|\Theta|})$

- Cost of $\mathcal{A}_{\text{Boost}}$: $O(\sqrt{|\Theta|}N)$

$\mathcal{A}_{\text{PSTM}}$

Lemma 2.5: $O(\frac{1}{\sqrt{p}}) \leq O(\sqrt{|\Theta|})$

- Cost of $\mathcal{A}_{\text{Boost}}$: $O(\sqrt{|\Theta|}N)$
- $\mathcal{A}_{\text{PSTC}}$: Run $\mathcal{A}_{\text{Boost}}$ for $O(\log |\Theta|)$ times to reduce $\text{width}(R)$

$\mathcal{A}_{\text{PSTM}}$

Lemma 2.5: $O\left(\frac{1}{\sqrt{p}}\right) \leq O(\sqrt{|\Theta|})$

- Cost of $\mathcal{A}_{\text{Boost}}$: $O(\sqrt{|\Theta|}N)$
- $\mathcal{A}_{\text{PSTC}}$: Run $\mathcal{A}_{\text{Boost}}$ for $O(\log |\Theta|)$ times to reduce $\text{width}(R)$
- Cost of $\mathcal{A}_{\text{PSTC}}$: Looks like $O(\sqrt{|\Theta|}N \log |\Theta|)$, but in fact $O(\sqrt{|\Theta|}N)$

$\mathcal{A}_{\text{PSTM}}$

Lemma 2.5: $O(\frac{1}{\sqrt{p}}) \leq O(\sqrt{|\Theta|})$

- Cost of $\mathcal{A}_{\text{Boost}}$: $O(\sqrt{|\Theta|}N)$
- $\mathcal{A}_{\text{PSTC}}$: Run $\mathcal{A}_{\text{Boost}}$ for $O(\log |\Theta|)$ times to reduce *width*(R)
- Cost of $\mathcal{A}_{\text{PSTC}}$: Looks like $O(\sqrt{|\Theta|}N \log |\Theta|)$, but in fact $O(\sqrt{|\Theta|}N)$
- \therefore Region R in later iteration is narrower than the previous iterations

$$\text{i.e. } \text{Cost}_{\text{Final iteration}} + \dots + \text{Cost}_{\text{1st iteration}} = O(\sqrt{|\Theta|} + \frac{\sqrt{|\Theta|}}{2} + \frac{\sqrt{|\Theta|}}{4} + \dots) = O(\sqrt{|\Theta|})$$

Overview

- 1 Introduction and Theory
 - Gradient Descent and Its Drawbacks
 - Overview of the Two Approaches
- 2 Approach 1: Average Approach
- 3 Approach 2: Partial Swap Test Cut-off Method (PSTC)
- 4 Equivalency Discussion**
- 5 Summary
- 6 Extensions and Future Work
 - Additional Restriction
 - Adversarial Example
- 7 Notations Table

Equivalency Discussion

- Are the 2 objective functions θ_{avg}^* and θ_{PSTC}^* equivalent?

Equivalency Discussion

- Are the 2 objective functions θ_{avg}^* and θ_{PSTC}^* equivalent?
- I.e., are there constants C, c such that on any sample set and any choice of ℓ ,

$$c\theta_{PSTC}^* \leq \theta_{avg}^* \leq C\theta_{PSTC}^*.$$

Equivalency Discussion

- Suppose such c, C existed. Consider a sample set of size 3 as below:
Let the parameters be $1, 2, \dots, C + 1$. Let $\ell_{threshold} = 2$ and suppose that for all $i \in 2, \dots, C$, we have the losses:

Equivalency Discussion

- Suppose such c, C existed. Consider a sample set of size 3 as below:
Let the parameters be $1, 2, \dots, C + 1$. Let $\ell_{threshold} = 2$ and suppose that for all $i \in 2, \dots, C$, we have the losses:



$$\ell_{i,1} = 2.5, \ell_{i,2} = \ell_{i,3} = 2$$

and suppose that,

$$\ell_{1,1} = \ell_{1,2} = \ell_{1,3} = 1.9$$

$$\ell_{C+1,1} = 2.1, \ell_{C+1,2} = 1, \ell_{C+1,3} = 0$$

Equivalency Discussion

- Suppose such c, C existed. Consider a sample set of size 3 as below:
Let the parameters be $1, 2, \dots, C + 1$. Let $\ell_{threshold} = 2$ and suppose that for all $i \in 2, \dots, C$, we have the losses:

•

$$\ell_{i,1} = 2.5, \ell_{i,2} = \ell_{i,3} = 2$$

and suppose that,

$$\ell_{1,1} = \ell_{1,2} = \ell_{1,3} = 1.9$$

$$\ell_{C+1,1} = 2.1, \ell_{C+1,2} = 1, \ell_{C+1,3} = 0$$

- Here, $\theta_{avg}^* = \underset{\theta}{\operatorname{argmin}} \frac{1}{3} \sum_{j=1}^3 \ell_{\theta j} = C + 1$, but

$$\theta_{PSTC}^* = \underset{\theta}{\operatorname{argmax}} \frac{1}{3} \sum_{j=1}^3 \mathbb{1}[\ell_{\theta j} \leq \ell_{threshold}] = 1 .$$

Equivalency Discussion

- Suppose such c, C existed. Consider a sample set of size 3 as below:
Let the parameters be $1, 2, \dots, C + 1$. Let $\ell_{threshold} = 2$ and suppose that for all $i \in 2, \dots, C$, we have the losses:

•

$$\ell_{i,1} = 2.5, \ell_{i,2} = \ell_{i,3} = 2$$

and suppose that,

$$\ell_{1,1} = \ell_{1,2} = \ell_{1,3} = 1.9$$

$$\ell_{C+1,1} = 2.1, \ell_{C+1,2} = 1, \ell_{C+1,3} = 0$$

- Here, $\theta_{avg}^* = \operatorname{argmin}_{\theta} \frac{1}{3} \sum_{j=1}^3 \ell_{\theta j} = C + 1$, but

$$\theta_{PSTC}^* = \operatorname{argmax}_{\theta} \frac{1}{3} \sum_{j=1}^3 \mathbb{1}[\ell_{\theta j} \leq \ell_{threshold}] = 1 .$$

- So, $C + 1 \leq C$ gives a contradiction.

Overview

- 1 Introduction and Theory
 - Gradient Descent and Its Drawbacks
 - Overview of the Two Approaches
- 2 Approach 1: Average Approach
- 3 Approach 2: Partial Swap Test Cut-off Method (PSTC)
- 4 Equivalency Discussion
- 5 **Summary**
- 6 Extensions and Future Work
 - Additional Restriction
 - Adversarial Example
- 7 Notations Table

Summary

Approach	Quantum Parallelism	Checking Process	Overall
Average Cut-off (PSTC)	$O(N)$ $O(1)$	N / A $O(N)$	$O(\sqrt{ \Theta }N)$ $O(\sqrt{ \Theta }N)$

Summary

Approach	Quantum Parallelism	Checking Process	Overall
Average	$O(N)$	N / A	$O(\sqrt{ \Theta }N)$
Cut-off (PSTC)	$O(1)$	$O(N)$	$O(\sqrt{ \Theta }N)$

- The overall cost of average approach is same as PSTC's

Summary

Approach	Quantum Parallelism	Checking Process	Overall
Average	$O(N)$	N / A	$O(\sqrt{ \Theta }N)$
Cut-off (PSTC)	$O(1)$	$O(N)$	$O(\sqrt{ \Theta }N)$

- The overall cost of average approach is same as PSTC's
- The cost of PSTC at quantum parallelism is much cheaper

Summary

Approach	Quantum Parallelism	Checking Process	Overall
Average	$O(N)$	N / A	$O(\sqrt{ \Theta }N)$
Cut-off (PSTC)	$O(1)$	$O(N)$	$O(\sqrt{ \Theta }N)$

- The overall cost of average approach is same as PSTC's
- The cost of PSTC at quantum parallelism is much **cheaper**
- Checking process cost (i.e. check whether $\theta \in R$) of PSTC is still **expensive**

Summary

Approach	Quantum Parallelism	Checking Process	Overall
Average	$O(N)$	N / A	$O(\sqrt{ \Theta }N)$
Cut-off (PSTC)	$O(1)$	$O(N)$	$O(\sqrt{ \Theta }N)$

- The overall cost of average approach is same as PSTC's
- The cost of PSTC at quantum parallelism is much **cheaper**
- Checking process cost (i.e. check whether $\theta \in R$) of PSTC is still **expensive**
- The “fundamental cost” of PSTC should be $O(\sqrt{|\Theta|} \cdot \text{sublinear } N)$

Summary

Approach	Quantum Parallelism	Checking Process	Overall
Average	$O(N)$	N / A	$O(\sqrt{ \Theta }N)$
Cut-off (PSTC)	$O(1)$	$O(N)$	$O(\sqrt{ \Theta }N)$

- The overall cost of average approach is same as PSTC's
- The cost of PSTC at quantum parallelism is much **cheaper**
- Checking process cost (i.e. check whether $\theta \in R$) of PSTC is still **expensive**
- The “fundamental cost” of PSTC should be $O(\sqrt{|\Theta|} \cdot \text{sublinear } N)$
- Future work: to lower the cost of PSTC to $O(\sqrt{|\Theta|} \cdot \text{sublinear } N)$ by enhancing the checking process

Overview

- 1 Introduction and Theory
 - Gradient Descent and Its Drawbacks
 - Overview of the Two Approaches
- 2 Approach 1: Average Approach
- 3 Approach 2: Partial Swap Test Cut-off Method (PSTC)
- 4 Equivalency Discussion
- 5 Summary
- 6 Extensions and Future Work**
 - Additional Restriction
 - Adversarial Example
- 7 Notations Table

Additional Restriction

Additional Restriction

- In optimization problem, we would normally search θ over Θ as below:

$$\operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{j=1}^N \ell(\theta, x_j)$$

Additional Restriction

- In optimization problem, we would normally search θ over Θ as below:

$$\operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{j=1}^N \ell(\theta, x_j)$$

- It is also common for us to have additional condition on θ (e.g. fairness, smoothness, or regularization):

$$\operatorname{argmin}_{\theta \in \Theta \cap \mathcal{A}} \frac{1}{N} \sum_{j=1}^N \ell(\theta, x_j)$$

Additional Restriction

- In optimization problem, we would normally search θ over Θ as below:

$$\operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{j=1}^N \ell(\theta, x_j)$$

- It is also common for us to have additional condition on θ (e.g. fairness, smoothness, or regularization):

$$\operatorname{argmin}_{\theta \in \Theta \cap \mathcal{A}} \frac{1}{N} \sum_{j=1}^N \ell(\theta, x_j)$$

- In view of this, we can generalize the cut-off indicator $E_{\theta j}$ to $E^{\mathcal{A}}_{\theta j}$:

$$\mathcal{L}_{\theta}^{\mathcal{A}, \text{PSTC}} \triangleq \frac{1}{N} \sum_{j=1}^N E^{\mathcal{A}}_{\theta j}, \text{ where } E^{\mathcal{A}}_{\theta j} = \mathbb{1}[\ell(\theta, x_j) \leq \tilde{\ell} \text{ and } \theta \in \mathcal{A}]$$

Adversarial Example

- The goal is to seek a perturbed input $x^{adv} = x + \delta^{adv}$
- The change δ^{adv} is not perceivable by human
- x^{adv} will be mis-classified by our target model h_θ

Adversarial Example

- The goal is to seek a perturbed input $x^{adv} = x + \delta^{adv}$
- The change δ^{adv} is not perceivable by human
- x^{adv} will be mis-classified by our target model h_θ

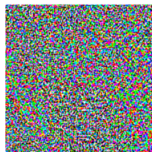


“panda”

57.7% confidence

x

+



noise

δ^{adv}

=



“gibbon”

99.3% confidence

x^{adv}

Adversarial Example

- One version of the adversarial problem can be formulated in this way:

$$x^{adv} = x + \delta^{adv} = x + \operatorname{argmax}_{\|\delta\| \leq \epsilon} \ell(h_{\theta}(x + \delta), y)$$

Adversarial Example

- One version of the adversarial problem can be formulated in this way:

$$x^{adv} = x + \delta^{adv} = x + \operatorname{argmax}_{\|\delta\| \leq \epsilon} \ell(h_{\theta}(x + \delta), y)$$

- In this problem, the role of x and θ are switched: θ is now fixed, and x is the parameters to be optimized.

Adversarial Example

- One version of the adversarial problem can be formulated in this way:

$$x^{adv} = x + \delta^{adv} = x + \underset{\|\delta\| \leq \epsilon}{\operatorname{argmax}} \ell(h_{\theta}(x + \delta), y)$$

- In this problem, the role of x and θ are switched: θ is now fixed, and x is the parameters to be optimized.
- Generalize our cut-off-qubit from $E_{\theta j}$ to $E_{\theta j}^A = E_{\theta j}^{\|\delta\| \leq \epsilon}$, where $E_{\theta j}^{\|\delta\| \leq \epsilon} = \mathbb{1}[\ell(\theta, x_j) \geq \tilde{\ell} \text{ and } \|\delta\| \leq \epsilon]$

Adversarial Example

- One version of the adversarial problem can be formulated in this way:

$$x^{adv} = x + \delta^{adv} = x + \operatorname{argmax}_{\|\delta\| \leq \epsilon} \ell(h_{\theta}(x + \delta), y)$$

- In this problem, the role of x and θ are switched: θ is now fixed, and x is the parameters to be optimized.
- Generalize our cut-off-qubit from $E_{\theta j}$ to $E_{\theta j}^A = E_{\theta j}^{\|\delta\| \leq \epsilon}$, where $E_{\theta j}^{\|\delta\| \leq \epsilon} = \mathbb{1}[\ell(\theta, x_j) \geq \tilde{\ell} \text{ and } \|\delta\| \leq \epsilon]$
- Apply PSTC by searching for “good” $\delta \in X$ and we will be able to obtain an adversarial example in the cost of $O(\sqrt{|X|})$.

Thank you!

Overview

- 1 Introduction and Theory
 - Gradient Descent and Its Drawbacks
 - Overview of the Two Approaches
- 2 Approach 1: Average Approach
- 3 Approach 2: Partial Swap Test Cut-off Method (PSTC)
- 4 Equivalency Discussion
- 5 Summary
- 6 Extensions and Future Work
 - Additional Restriction
 - Adversarial Example
- 7 Notations Table

Notations Table: Basics

Notation	Space	Meaning
θ	Θ	Parameter of model $h_\theta(\cdot)$
$h_\theta(\cdot)$	$X \rightarrow Y$	Model with parameter θ
x	X	An input
y	Y	A label
S	$S \subseteq X$	$\{x_1, \dots, x_N\}$
N	\mathbb{Z}^+	#(samples in S)
n	\mathbb{Z}^+	$N = 2^n$
m	\mathbb{Z}^+	$ \Theta = 2^m$
$\ell(\cdot, \cdot)$, or $loss(\cdot, \cdot)$	$\Theta \times X \rightarrow \mathbb{R}^+$ (or $Y \times Y \rightarrow \mathbb{R}^+$)	Loss function (depending on context)
ℓ_{ij} , or $\ell_{\theta_i j}$, or $\ell_{\theta j}$	\mathbb{R}^+	$\ell_{\theta_i j} = \ell(\theta_i, x_j) = \ell(h_{\theta_i}(x_j), y_j)$
\mathcal{L}_θ^{avg} or $\mathcal{L}^{avg}(\theta)$	\mathbb{R}^+	Average loss $\frac{1}{N} \sum_{j=1}^N \ell_{\theta j}$; The smaller the better
$\mathcal{L}_\theta^{PSTC}$ or $\mathcal{L}^{PSTC}(\theta)$	\mathbb{R}^+	Cut-off loss $\frac{1}{N} \sum_{j=1}^N \mathbb{1}[\ell_{\theta j} \leq \tilde{\ell}]$; The larger the better
$\tilde{\ell}$, or $\ell_{threshold}$	\mathbb{R}^+	Threshold value for the cut-off approach
$E_{\theta j}$	$\{0, 1\}$	Cut-off indicator: Value $\mathbb{1}[\ell_{\theta j} \leq \tilde{\ell}]$ to be stored in 1 qubit
$ \phi\rangle, \psi\rangle$	\mathbb{C}^{n+1}	Pure states used in PSTC

Notations Table: PSTC

$Q_{1\text{-query}}$ $ \varsigma_\theta\rangle$ $ \varsigma_\theta^{(0)}\rangle, \varsigma_\theta^{(1)}\rangle$ $\text{PMF}(\cdot)$ $\mathcal{A}_{1\text{-query}}, \mathcal{A}_\xi, \mathcal{A}_{\text{Boost}}, \mathcal{A}_{\text{PSTC}}$ $\theta_{1\text{-query}}, \theta_{\text{Boost}}, \theta_{\text{Best}}$ $\text{flag}(\cdot)$ O_{flag} \mathcal{R} R $\text{ref}(B\rangle)$ $\text{ref}(U\rangle)$ p $\xi(\cdot)$ $\tilde{\xi}$, or $\xi_{\text{threshold}}$ $ 0\rangle$	Q-Circuits \mathbb{C}^{2n+3} \mathbb{C}^{2n+2} $\Theta \rightarrow [0, 1]$ Algorithms Θ $\Theta \rightarrow \{0, 1\}$ $\mathbb{C}^{m \times m}$ $\mathbb{C}^{m \times m}$ $R \subseteq \Theta$ $\mathbb{C}^{m \times m}$ $\mathbb{C}^{m \times m}$ $[0, 1] \in \mathbb{R}$ $\Theta \rightarrow \mathbb{R}^+$ \mathbb{R}^+ $\mathbb{C}_{\text{not care}}$	The Q-circuit for partial swap test The pure state after the 2nd Hadamard gate in $Q_{1\text{-query}}$ given θ The $ 0\rangle$ and $ 1\rangle$ parts of $ \varsigma_\theta\rangle$ Probability mass function (PMF) of θ Names of some cut-off approach algorithms θ being used in the respective algorithms Flag to indicate whether θ is good or bad Gate of oracles on $\text{flag}(\theta)$ Gate of “keeping zero, flipping else” Region R that “good” θ s concentrate Reflection on the “Bad” vector in $\mathcal{A}_{\text{Boost}}$ Reflection on the “Uniform” vector in $\mathcal{A}_{\text{Boost}}$ $\mathbb{P}[\text{picked } \theta \text{ in } R]$ A function that proportional to pdf of $\theta : P[\text{observe } \theta \text{1st qubit} = 0]$ Threshold value of ξ for region R High dimensional $ 0\rangle$ with dimension not mentioned
x^{adv} δ^{adv} $ \chi^{\text{Type I}}\rangle, \chi^{\text{Type II}}\rangle$	\mathbb{X} \mathbb{X} $\mathbb{C}_{\text{depends}}$	Adversarial example Perturbation Pure states of x referring to Type I and Type II uniformity.

References

- [1] R. de Wolf. Quantum computing: Lecture notes, 2019. URL <https://arxiv.org/abs/1907.09415>.
- [2] C. Durr and P. Hoyer. A Quantum Algorithm for finding the Minimum. 1999. URL <https://arxiv.org/abs/quant-ph/9607014>.