Trainability Proof Notes

Zoe Holmes

December 2022

1 Local to global state bound

Proposition 1. Let ρ be an arbitrary quantum state, let U_{θ} be a parameterized unitary matrix and write $\rho_U = U_{\theta}^{\dagger} \rho U_{\theta}$ for short. Consider a global cost of the form

$$C_G(\boldsymbol{\theta}) = 1 - \text{Tr}\left[\rho_U H_G\right] \quad where \quad H_G = |0\rangle\langle 0|^{\otimes n}$$
 (1)

is the all zero projector, and a local cost of the form

$$C_L(\boldsymbol{\theta}) = 1 - \text{Tr}\left[\rho_U H_L\right] \quad where \quad H_L = \frac{1}{n} \sum_{j=1}^{n} |0\rangle\langle 0|_j \otimes \mathbb{I}_{\overline{j}}$$
 (2)

It follows that

$$C_L(\theta) \leqslant C_G(\theta) \leqslant nC_L(\theta)$$
 (3)

Thus $C_L(\boldsymbol{\theta}) = 0$ iff $C_G(\boldsymbol{\theta}) = 0$.

Proof. We can write $H_{\rm L} = \frac{1}{n} \sum_{i=1}^{n} H_{{\rm L},j}$, where

$$H_{L,j} = |0\rangle\langle 0|_j \otimes \mathbb{I}_{\overline{j}} \tag{4}$$

are projectors that mutually commute. Note that $\prod_{j=1}^n H_{L,j} = H_G$. We can associate events E_j with the projectors $H_{L,j}$ such that $\Pr[E_j] = \operatorname{Tr}\left[\rho_U H_{L,j}\right]$. Then, $\operatorname{Tr}\left[\rho\prod_{j=1}^n H_{L,j}\right] = \Pr\bigcap_{j=1}^n E_j$. Recall, from basic probability theory, that for any set of events $\mathcal{A} := \{A_1, A_2, \ldots, A_n\}$, it holds that

$$\Pr[\bigcup_{i=1}^{n} A_i] \geqslant \frac{1}{n} \sum_{i=1}^{n} \Pr[A_i].$$
 (5)

Choosing $A_i = \overline{E_j}$, we see

$$\Pr\left[\bigcup_{j=1}^{n} \overline{E_{j}}\right] \geqslant \frac{1}{n} \sum_{i=j}^{n} \Pr\left[\overline{E_{j}}\right]$$

$$\implies 1 - \Pr\left[\bigcap_{j=1}^{n} E_{j}\right] \geqslant \frac{1}{n} \sum_{j=1}^{n} \left(1 - \Pr[E_{j}]\right)$$

$$\implies 1 - \Pr\left[\bigcap_{j=1}^{n} E_{j}\right] \geqslant 1 - \frac{1}{n} \sum_{j=1}^{n} \operatorname{Tr}[\rho_{U} H_{L,j}]$$

$$\implies 1 - \operatorname{Tr}[\rho_{U} H_{G}] \geqslant 1 - \operatorname{Tr}[\rho_{U} H_{L}].$$
(6)

This is precisely the first desired inequality $C_L \leq C_G$.

To prove the remaining inequality, observe that, via the union bound,

$$\Pr\left[\bigcup_{j=1}^{n} A_j\right] \leqslant \sum_{i=1}^{n} \Pr[A_j] \tag{7}$$

we have

$$\Pr\left[\bigcup_{j=1}^{n} \overline{E_{j}}\right] \leqslant \sum_{j=1}^{n} \Pr\left[\overline{E_{j}}\right]$$

$$\implies 1 - \Pr\left[\bigcap_{j=1}^{n} E_{j}\right] \leqslant \sum_{j=1}^{n} (1 - \Pr[E_{j}])$$

$$\implies 1 - \operatorname{Tr}[\rho_{U} H_{G}] \leqslant n \left(1 - \operatorname{Tr}[\rho_{U} H_{L}]\right)$$
(8)

Thus, $C_G \leq nC_L$ as required.

2 Barren Plateaus

A parameterized unitary circuit can be represented as

$$U(\boldsymbol{\theta}) = \prod_{l=1}^{L} \exp(-i\theta_l V_l) W_l = \underbrace{\prod_{l=k+1}^{L} \exp(-i\theta_l V_l) W_l}_{U_+} \underbrace{\prod_{l=1}^{k} \exp(-i\theta_l V_l) W_l}_{U_-}$$
(9)

where $\{V_l\}_{l=1}^L$ are hermitian, involutionary $(V_l^2 = \mathbb{I})$ generators of rotations $\exp(-i\theta_l V_l)$ around independent, trainable angles $\boldsymbol{\theta} \in [0, 2\pi)^L$ and $\{W_l\}_{l=1}^L$ are fixed gates.

Taking the partial derivative of the circuit gives

$$\partial_k U(\boldsymbol{\theta}) = -iU_+ V_k U_- \tag{10}$$

$$\partial_k U^{\dagger} = iU_{-}^{\dagger} V_k U_{+}^{\dagger} \tag{11}$$

We will take our cost function to be the standard VQE cost, i.e.

$$E(\boldsymbol{\theta}) = \langle 0|U(\boldsymbol{\theta})^{\dagger}HU(\boldsymbol{\theta})|0\rangle = \text{Tr}\left(U(\boldsymbol{\theta})\rho U(\boldsymbol{\theta})^{\dagger}H\right). \tag{12}$$

This is the expectation value of a Hermitian operator H for a state $U(\boldsymbol{\theta})\rho U(\boldsymbol{\theta})^{\dagger}$ where $\rho = |0\rangle\langle 0|^{\otimes n}$. The derivative of the cost function with respect to θ_k is given by

$$\partial_{\theta_k} E(\boldsymbol{\theta}) = i \operatorname{Tr} \left(\rho (U_-^{\dagger} V_k U_+^{\dagger} H U - U^{\dagger} H U_+ V_k U_-) \right)$$
(13)

$$= i \operatorname{Tr} \left(\rho U_{-}^{\dagger} \left(V_{k} \underbrace{U_{+}^{\dagger} H U_{+}}_{H_{+}} - \underbrace{U_{+}^{\dagger} H U_{+}}_{H_{+}} V_{k} \right) U_{-} \right)$$

$$\tag{14}$$

$$= i \operatorname{Tr} \left(\underbrace{U_{-}\rho U_{-}^{\dagger}}_{\rho_{-}} [V_{k}, H_{+}] \right) \tag{15}$$

$$= i \operatorname{Tr} \left([\rho_{-}, V_{k}] U_{+}^{\dagger} H U_{+} \right) \tag{16}$$

In the final line, we use the cyclicity of the trace operation, i.e. Tr[A[B,C]] = Tr[ABC - ACB] = Tr[ABC - BAC] = Tr[[A,B]C].

2.1 Average of the cost gradients

The average value of the gradient of the cost function, over a random initialisation of the parameters θ , is given by

$$\langle \partial_{\theta_k} E(\boldsymbol{\theta}) \rangle = \int dU p(U) \partial_{\theta_k} E(\boldsymbol{\theta}) = \int dU_+ p(U_+) dU_- p(U_-) \partial_{\theta_k} E(\boldsymbol{\theta})$$
(17)

Thus we have

$$\langle \partial_{\theta_k} E(\boldsymbol{\theta}) \rangle = i \int \int dU_- p(U_-) \operatorname{Tr} \left(\rho_- \left[V_k, \int dU_+ p(U_+) H_+ \right] \right)$$
 (18)

(19)

Assuming U_{-} is at least a 1-design using Eq.(39) (from Section 2.3) we have that

$$\langle \partial_{\theta_k} E(\boldsymbol{\theta}) \rangle = i \int d\mu(U_-) \operatorname{Tr} \left(\rho_- \left[V_k, \int dU_+ p(U_+) H_+ \right] \right)$$
 (20)

$$= i \frac{\operatorname{Tr}(\rho)}{d} \operatorname{Tr}\left(\left[V_k, \int dU_+ p(U_+) H_+\right]\right)$$
(21)

$$=0 (22)$$

as the trace of a commutator is zero.

Similarly, assuming U_{+} is at least a 1-design, results in:

$$\langle \partial_{\theta_k} E(\boldsymbol{\theta}) \rangle = i \int dU_- p(U_-) \operatorname{Tr} \left(\rho_- \left[V_k, \int d\mu(U_+) H_+ \right] \right)$$
 (23)

$$= i \int dU_{-}p(U_{-}) \operatorname{Tr}\left(\rho_{-} \left[V_{k}, \frac{\operatorname{Tr}(H)}{d}\mathbb{I}\right]\right)$$
(24)

$$=0 (25)$$

Hence, the gradients are not biased in any single direction.

2.2 Variance of the cost gradients

The probability that the cost function gradient deviates from its average of zero can be bounded using Chebyshev's inequality,

$$P(|\partial_{\theta_k} E| \geqslant \epsilon) \leqslant \frac{\operatorname{Var}(\partial_{\theta_k} E)}{\epsilon^2},$$
 (26)

where the variance of the gradient is given as:

$$\operatorname{Var}(\partial_{\theta_k} E) = \left\langle (\partial_{\theta_k} E)^2 \right\rangle - \underbrace{\left\langle \partial_{\theta_k} E \right\rangle^2}_{0} = \int dU p(U) (\partial_{\theta_k} E)^2 \tag{27}$$

$$= \int dU_{-}p(U_{-}) \int dU_{+}p(U_{+})i^{2} \operatorname{Tr}\left(\rho_{-}^{\bigotimes^{2}}[V_{k}, H_{+}]^{\bigotimes^{2}}\right)$$
(28)

Assuming U_{-} is a 2-design, using Eq.(40), results in:

$$\operatorname{Var}(\partial_{\theta_{k}}E) = -\int dU_{+}p(U_{+}) \int d\mu(U_{-}) \operatorname{Tr}\left(U_{-}^{\otimes^{2}}\rho^{\otimes^{2}}U_{-}^{\dagger\otimes^{2}}[V_{k}, H_{+}]^{\otimes^{2}}\right)$$

$$= -\int dU_{+}p(U_{+}) \left(\frac{\operatorname{Tr}\left([V_{k}, H_{+}]^{\otimes^{2}}\right) \operatorname{Tr}\left(\rho^{\otimes^{2}}\right) + \operatorname{Tr}\left([V_{k}, H_{+}]^{2}\right) \operatorname{Tr}\left(\rho^{2}\right)}{d^{2} - 1}$$

$$- \frac{\operatorname{Tr}\left([V_{k}, H_{+}]^{2}\right) \operatorname{Tr}\left(\rho^{\otimes^{2}}\right) + \operatorname{Tr}\left([V_{k}, H_{+}]^{\otimes^{2}}\right) \operatorname{Tr}\left(\rho^{2}\right)}{d(d^{2} - 1)}$$

$$= -\left(1 - \frac{1}{d}\right) \frac{\left\langle \operatorname{Tr}\left([V_{k}, H_{u}]^{2}\right)\right\rangle_{U_{+}}}{d^{2} - 1}$$

$$(29)$$

As $\operatorname{Tr}(\rho) = 1$ and $\operatorname{Tr}([V_k, H_+]^{\otimes^2}) = \operatorname{Tr}([V_k, H_+])^2 = 0$.

To bound $\operatorname{Tr}([V_k, H_+]^2)$ we use the triangle inequality and then Cauchy-Schwarz to show that

$$\left| \text{Tr}([V_k, H_+]^2) \right| = \left| 2 \, \text{Tr}((V_k H_+)^2) - 2 \, \text{Tr}\left(\underbrace{V_k^2}_{\text{\tiny T}} H_+^2\right) \right|$$
 (30)

$$\leq 2\left|\operatorname{Tr}\left((V_k H_+)^2\right)\right| + 2\underbrace{\left|\operatorname{Tr}\left(H_+^2\right)\right|}_{\|H\|_2^2} \tag{31}$$

$$\leq 2||V_k||_2||V_kH_+V_k||_2 + 2||H||_2^2$$
 (32)

$$= 2 \operatorname{Tr} \left[V_k H_+ V_k V_k^{\dagger} H_+ V_k^{\dagger} \right] + 2 \|H\|_2^2$$
 (33)

$$= 2\operatorname{Tr}[H_{+}^{2}] + 2\|H\|_{2}^{2} \tag{34}$$

$$=4\|H\|_2^2\tag{35}$$

Assuming that $||H||_2^2 \in \mathcal{O}(d)$, it follows that $\operatorname{Var}(\partial_{\theta_k} E) \in \mathcal{O}(\frac{1}{d})$.

Assuming U_+ is a 2-design, results in:

$$\operatorname{Var}(\partial_{\theta_k} E) = -\int dU_- p(U_-) \int d\mu(U_+) \operatorname{Tr}\left(U_+^{\bigotimes^2} [\rho_-, V_k]^{\bigotimes^2} U_+^{\dagger \bigotimes^2} H^{\bigotimes^2}\right)$$
(36)

$$= -\frac{1}{d^2 - 1} \left(\operatorname{Tr}(H^2) - \frac{\operatorname{Tr}(H)^2}{d} \right) \left\langle \operatorname{Tr}([\rho_u, V_k]^2) \right\rangle_{U_-}$$
 (37)

For simplicity let's assume that H is expressed as a sum of Pauli operators and so Tr(H) = 0. Then using Cauchy-Schwarz (41) we get:

$$\left| \operatorname{Tr} \left(\left[\rho_{-}, V_{k} \right]^{2} \right) \right| \leqslant 4 \operatorname{Tr} \left(\rho_{-}^{2} \right) = 4 \tag{38}$$

Therefore, $\operatorname{Var}(\partial_{\theta_k} E) \in \mathcal{O}(\frac{1}{d})$.

Therefore the probability that the cost gradient deviates from 0 is suppressed exponentially in the number of qubits.

2.3 Useful formulas

Let $\mathcal{U}(d)$ denote the unitary group of degree $d=2^n$. Let $d\mu(U)$ be the volume element of the Haar measure, where $U \in \mathcal{U}(d)$. The following identities hold.

$$\int d\mu(U) \operatorname{Tr}(UAU^{\dagger}B) = \frac{\operatorname{Tr}(A) \operatorname{Tr}(B)}{d}$$
(39)

$$\int d\mu(U) \operatorname{Tr}\left(AU^{\bigotimes^2}BU^{\dagger\bigotimes^2}\right) = \frac{\operatorname{Tr}(A)\operatorname{Tr}(B) + \operatorname{Tr}(AW)\operatorname{Tr}(BW)}{d^2 - 1} - \frac{\operatorname{Tr}(AW)\operatorname{Tr}(B) + \operatorname{Tr}(A)\operatorname{Tr}(BW)}{d(d^2 - 1)}$$
(40)

where W is the swap operator $W|i\rangle |j\rangle = |j\rangle |i\rangle$.

Note that $\operatorname{Tr}(A \bigotimes C)W) = \operatorname{Tr}(AC)$ (prove this to your self!) and therefore $\operatorname{Tr}\left(A^{\bigotimes^2}W\right) = \operatorname{Tr}(A^2)$.

The Cauchy-Schwarz inequality in trace form is given by:

$$\left|\operatorname{Tr}(A^{\dagger}B)\right| \leqslant ||A||_2||B||_2 = \sqrt{\operatorname{Tr}(A^{\dagger}A)}\sqrt{\operatorname{Tr}(B^{\dagger}B)}$$
 (41)