

Concept learning of parameterized quantum models from limited measurements



Po-Wei Huang

arXiv:2408.05116 [quant-ph]



Slides adapted from Beng Yee Gan



Beng Yee
Gan



Po-Wei
Huang



Elies
Gil-Fuster




Patrick
Rebentrost

Power of quantum computers

Quantum machine learning

Dequantizing algorithms to understand quantum advantage in machine learning

Ewin Tang 

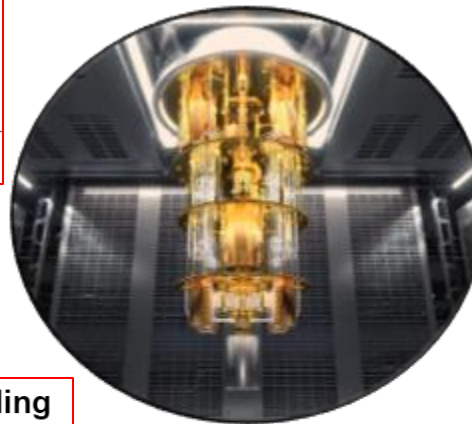
Random circuit sampling

A Polynomial-Time Classical Algorithm for Noisy Random Circuit Sampling

Authors:  Dorit Aharonov,  Xun Gao,  Zeph Landau,  Yunchao Liu,  Umesh Vazirani | [Authors Info & Claims](#)

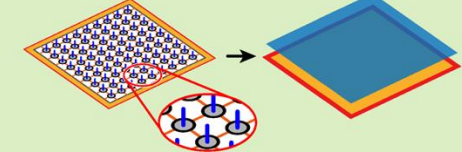
STOC 2023: Proceedings of the 55th Annual ACM Symposium on Theory of Computing • Pages 945 - 957

<https://doi.org/10.1145/3564246.3585234>



Quantum simulations

Tensor network



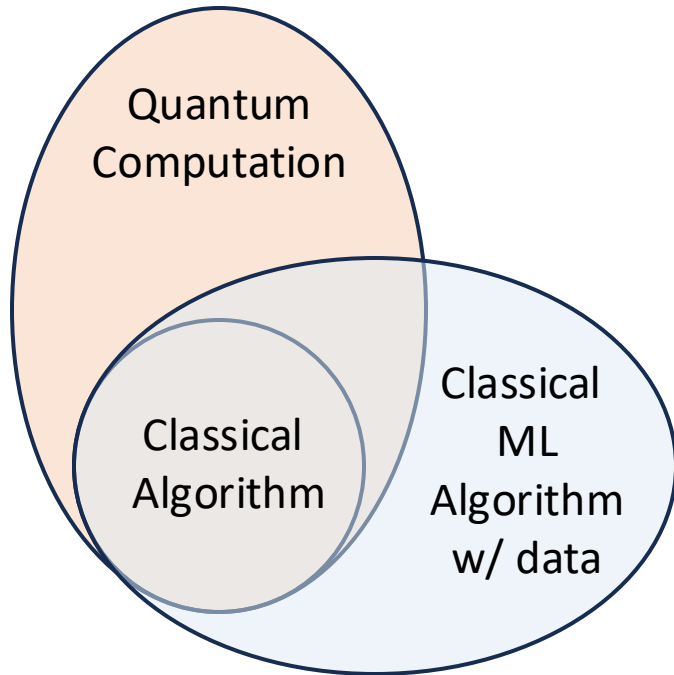
Shor's algorithm

**Quantum Computers
Destroy Internet Security**



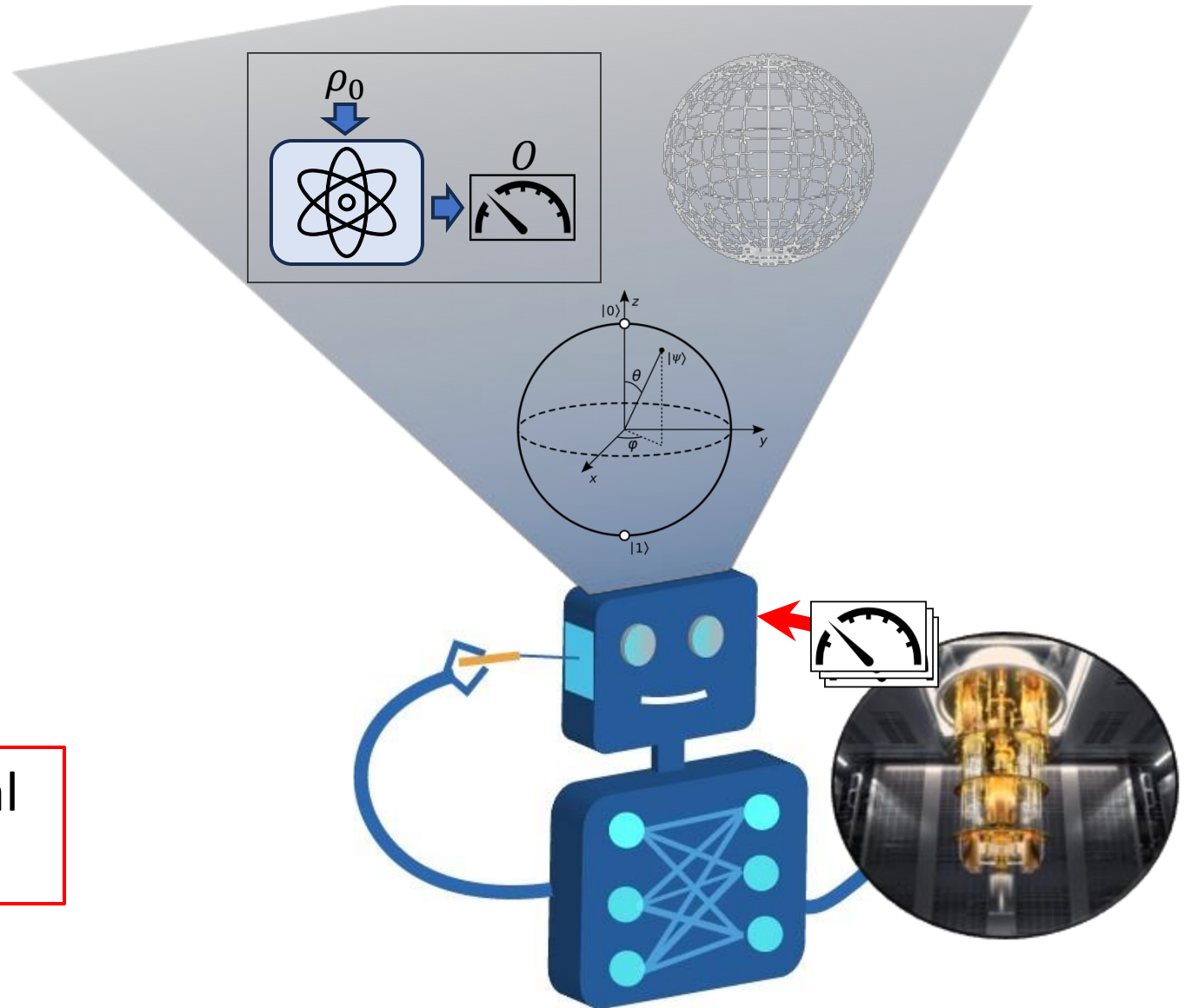
Performance is benchmarked against classical computers.

Replacing (some) quantum with classical

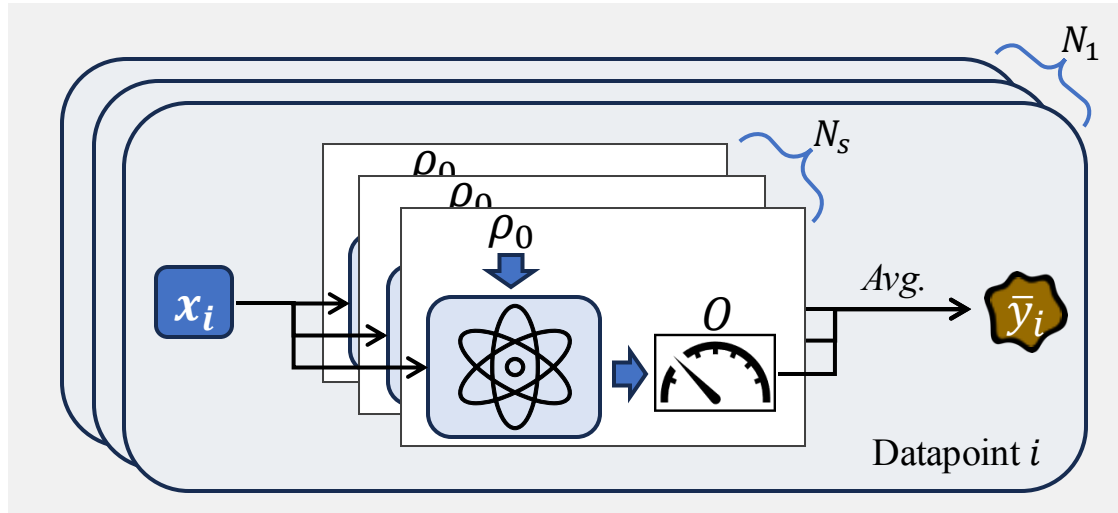


Nat. Commun., 12(1), 2631 (2021).

Access to data makes classical machines more powerful.



Learning from limited measurements

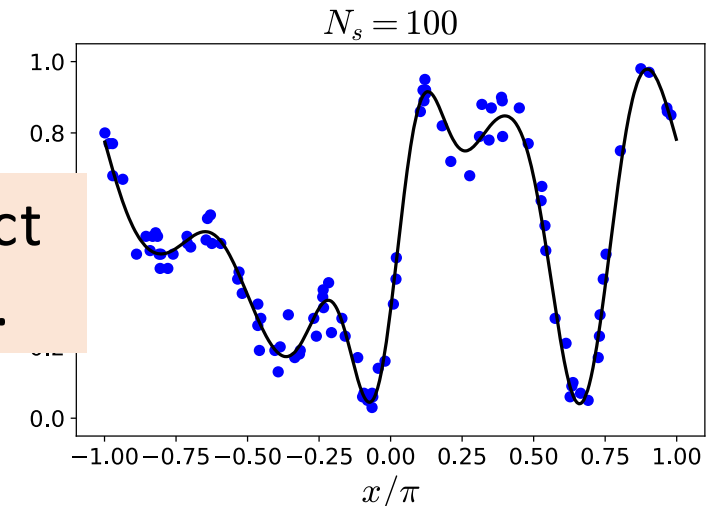
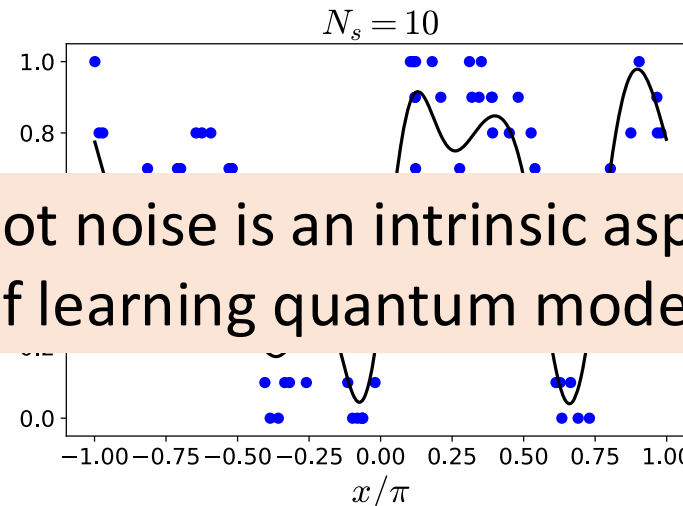
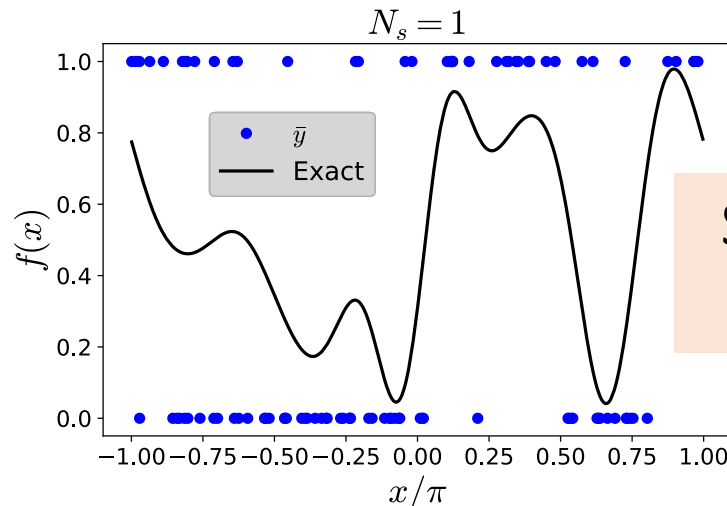


Quantum models: $f_{\theta}(x) = \text{tr}(\rho_{\theta}(x)O)$

$\mathbb{E}_{\bar{y}}[\bar{y}|x]$

- Dataset: $(x_i, \bar{y}_i)_{i=1}^{N_1}$

Estimated with N_s shots

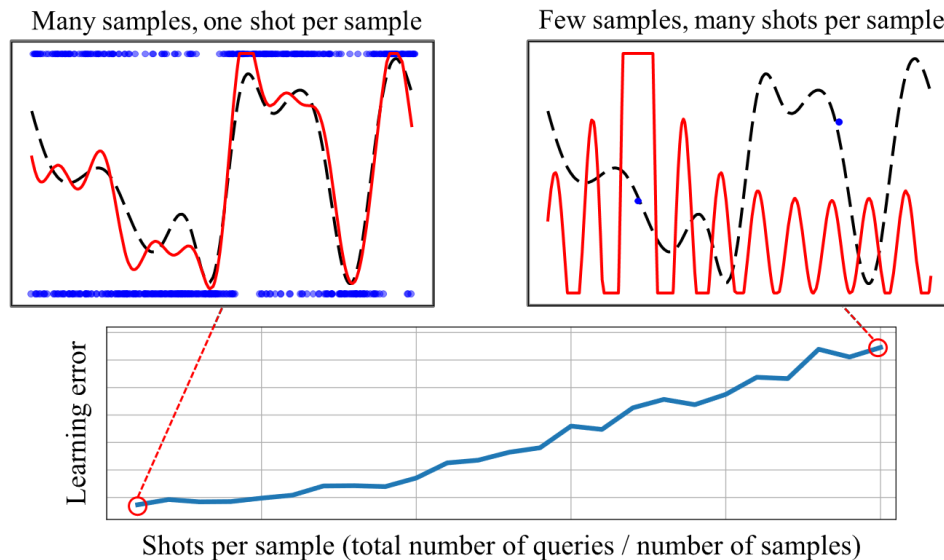


Shot noise is an intrinsic aspect of learning quantum models.

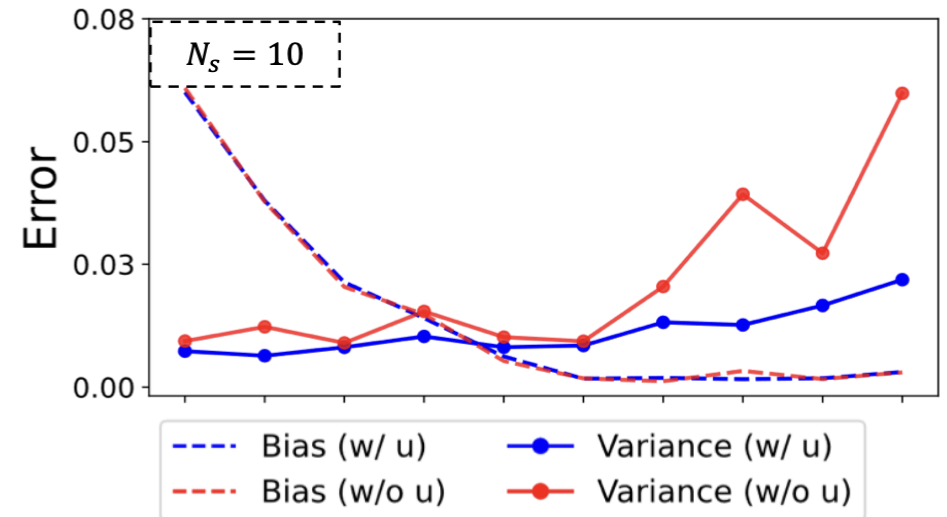
Results overview

Can we obtain provable guarantees of learning that exemplify the relationship between N_1 and N_s ?

(1) Asymmetrical trade-offs between N_1 and N_s



(2) Gradient descent* can be made robust and provide tighter guarantees



Probabilistic concept learning

Probabilistic concept class

$$\mathcal{F} = \{f(\mathbf{x}) = \mathbb{E}_y[y|\mathbf{x}]\}$$

$$\mathcal{D} = p(\mathbf{x})p(y|\mathbf{x})$$

Loss function: $\ell(f_\theta(\mathbf{x}), h_w(\mathbf{x}))$

Hypothesis class

$$\mathcal{H} = \{h_w(\mathbf{x}), \mathbf{w} \in \mathbb{R}^D\}$$

Learning task

➤ Unknown $f \in \mathcal{F}$

➤ Goal:

$$\overbrace{\mathbb{E}_x[\ell(f(\mathbf{x}), h_w(\mathbf{x}))]}^{R(h_w)}$$

➤ Since $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{N_1} \sim \mathcal{D}$

➤ We can only get: $\hat{R}(h_w)$

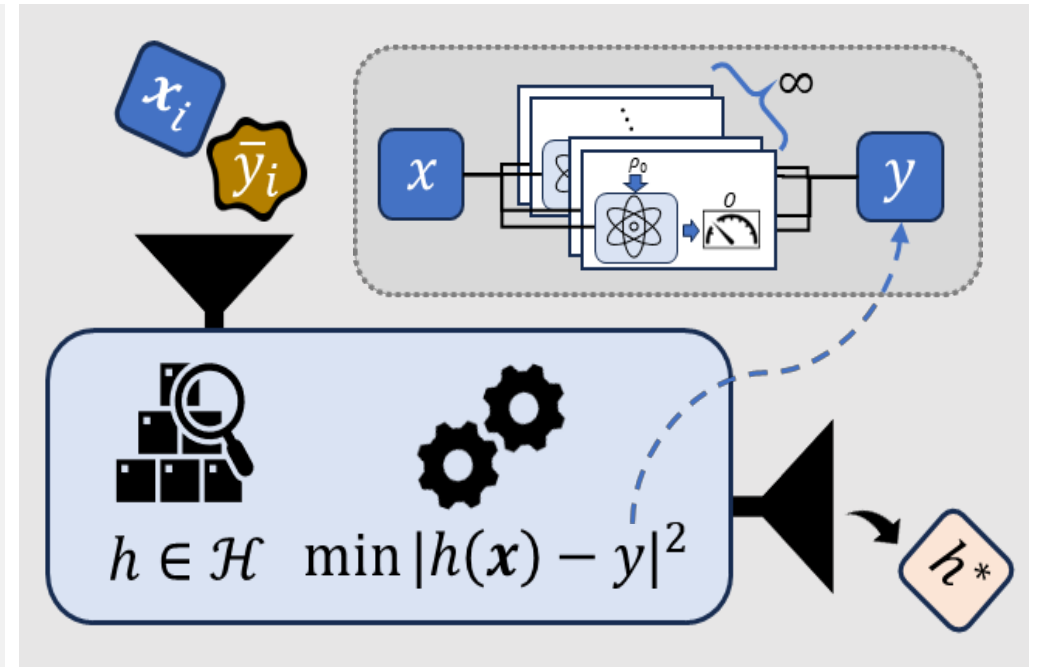
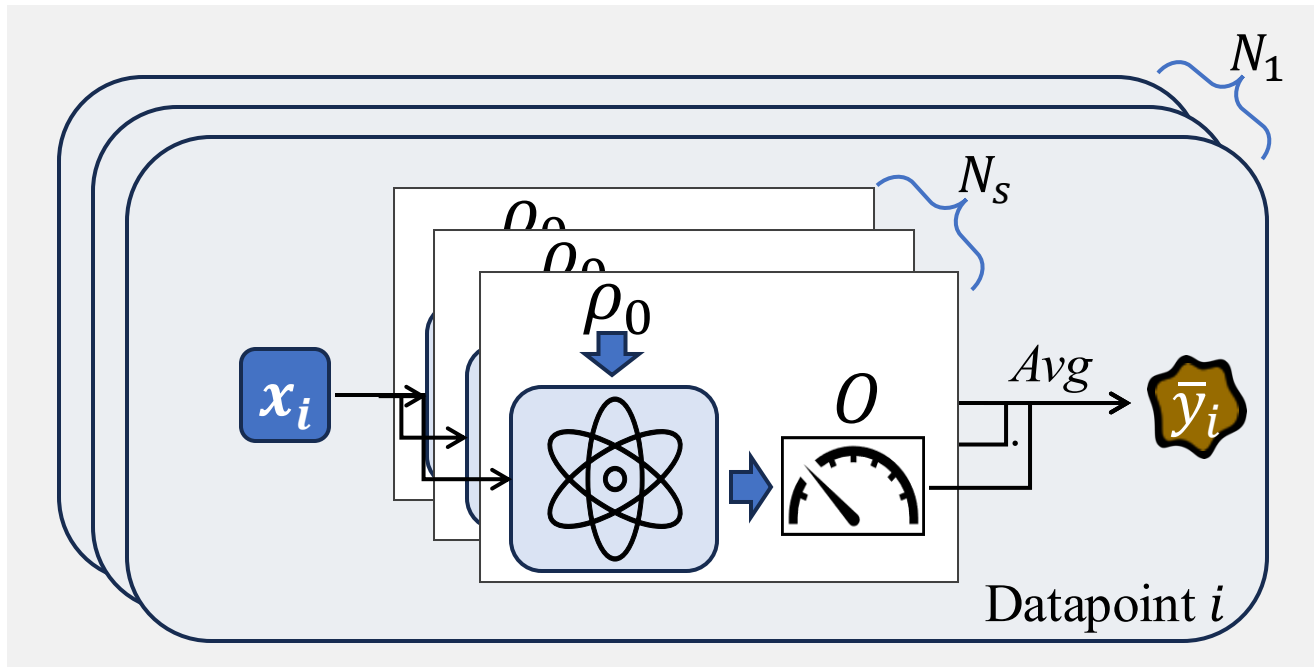
$$\overbrace{\frac{1}{N_1} \sum_{i=1}^{N_1} \ell(\bar{y}, h_w(\mathbf{x}_i))}^{\hat{R}(h_w)}$$

➤ Task: $\hat{R}(h_w) - R(h_{w^*}) \leq \epsilon$

Provable guarantee

Getting data – the black box model

Family of PQC models: $\mathcal{F}_{U,O} = \{f_{\theta}(x) = \underbrace{\langle 0|U^{\dagger}(x,\theta)OU(x,\theta)|0\rangle}_{\text{tr}(\rho_{\theta}(x)O)}\}$



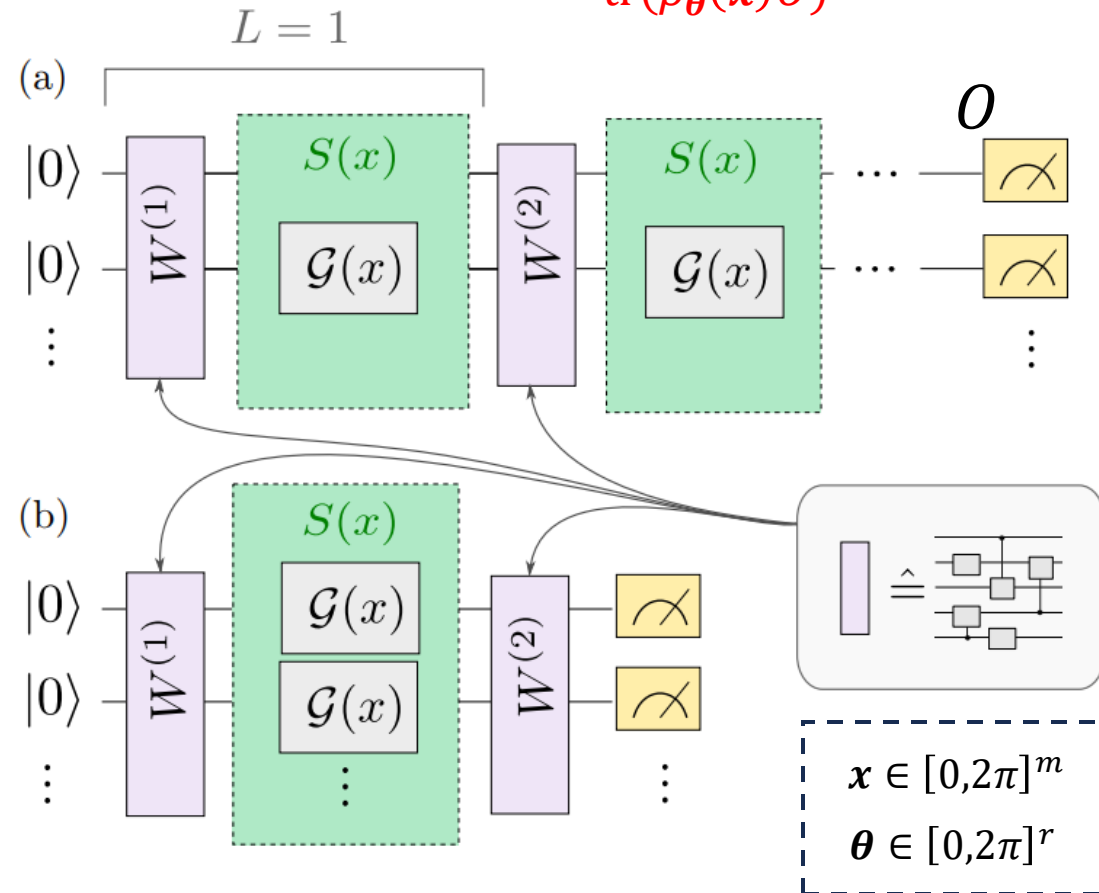
Which PQCs are learnable? How do we get the hypothesis class \mathcal{H} ?

Getting data – the grey box model

Family of PQC models: $\mathcal{F}_{U,O} = \{f_{\theta}(x) = \underbrace{\langle \mathbf{0} | U^\dagger(x, \theta) O U(x, \theta) | \mathbf{0} \rangle}_{\text{tr}(\rho_{\theta}(x) O)}\}$

Which PQCs are learnable?
No clear answer yet.

But we know
some PQCs are learnable,
and we know how to provide
their classical surrogates.



PRA 103, 032430 (2021).

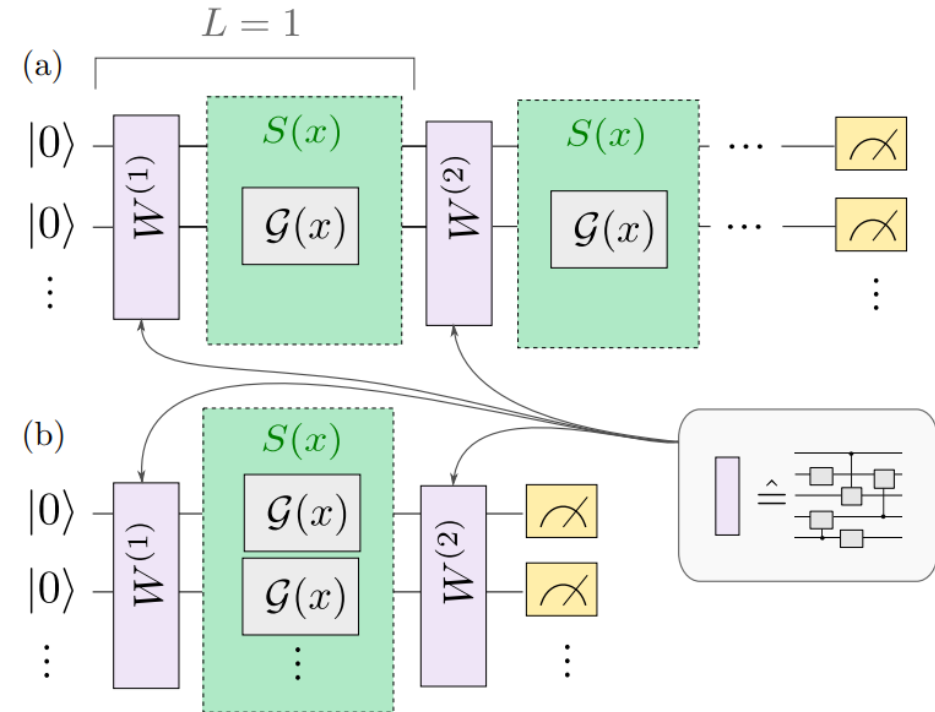
Fourier representation of PQC

$$f_{\theta}(x) = c_{\omega_0}(\theta) + \sum_{i=1}^{|\Omega|-1} a_{\omega_i}(\theta) \cos(\langle \omega_i, x \rangle) + b_{\omega_i}(\theta) \sin(\langle \omega_i, x \rangle)$$

$$= \langle \mathbf{w}_F(\theta), \phi_F(x) \rangle$$

$$\sqrt{|\Omega|} \begin{pmatrix} c_{\omega_0}(\theta) \\ a_{\omega_1}(\theta) \\ b_{\omega_1}(\theta) \\ \vdots \\ a_{\omega_{|\Omega|-1}}(\theta) \\ b_{\omega_{|\Omega|-1}}(\theta) \end{pmatrix} \xrightarrow{\quad} \frac{1}{\sqrt{|\Omega|}} \begin{pmatrix} 1 \\ \cos(\langle \omega_1, x \rangle) \\ \sin(\langle \omega_1, x \rangle) \\ \vdots \\ \cos(\langle \omega_{|\Omega|-1}, x \rangle) \\ \sin(\langle \omega_{|\Omega|-1}, x \rangle) \end{pmatrix}$$

$$\mathcal{F}_{U,O} = \{f_{\theta}(x) = \underbrace{\langle \mathbf{w}_F(\theta), \phi_F(x) \rangle}_{\text{tr}(\rho_{\theta}(x)O)}\}$$



Classical machine learning models

PQC models: $f_{\theta}(\mathbf{x}) \in [0,1]$

$$f_{\theta}(\mathbf{x}) = \langle \mathbf{w}_F, \boldsymbol{\phi}_F(\mathbf{x}) \rangle = \langle \mathbf{w}, \boldsymbol{\phi}(\mathbf{x}) \rangle + \xi(\mathbf{x})$$

$$\begin{aligned}\xi(\mathbf{x}) &= \langle \mathbf{w}_F, \boldsymbol{\phi}_F(\mathbf{x}) \rangle - \langle \mathbf{w}, \boldsymbol{\phi}(\mathbf{x}) \rangle \\ \xi(\mathbf{x}) &\in [-M, M] \\ \mathbb{E}_{\mathbf{x}}[\xi(\mathbf{x})^2] &\leq \epsilon_1\end{aligned}$$

Hypothesis class: $\mathcal{H} = \{h(\mathbf{x}) = \langle \mathbf{w}, \boldsymbol{\phi}(\mathbf{x}) \rangle\}$

$$\|\mathbf{w}\|_2 \leq B, \|\boldsymbol{\phi}(\mathbf{x})\|_2 \leq 1$$

Simplify hypothesis class:

- Directly truncate $\Omega \rightarrow D$
- Random Fourier Features

$$\boldsymbol{\phi}(\mathbf{x}) = \frac{1}{\sqrt{D}} \begin{pmatrix} 1 \\ \cos(\langle \boldsymbol{\omega}_1, \mathbf{x} \rangle) \\ \sin(\langle \boldsymbol{\omega}_1, \mathbf{x} \rangle) \\ \vdots \\ \cos(\langle \boldsymbol{\omega}_D, \mathbf{x} \rangle) \\ \sin(\langle \boldsymbol{\omega}_D, \mathbf{x} \rangle) \end{pmatrix}$$

Empirical risk minimization

PQC models: $f_{\theta}(\mathbf{x}) \in [0,1]$

$$f_{\theta}(\mathbf{x}) = \langle \mathbf{w}, \boldsymbol{\phi}(\mathbf{x}) \rangle + \xi(\mathbf{x})$$

$$\boldsymbol{\phi}(\mathbf{x}) = \frac{1}{\sqrt{D}} \begin{pmatrix} 1 \\ \cos(\langle \boldsymbol{\omega}_1, \mathbf{x} \rangle) \\ \sin(\langle \boldsymbol{\omega}_1, \mathbf{x} \rangle) \\ \vdots \\ \cos(\langle \boldsymbol{\omega}_D, \mathbf{x} \rangle) \\ \sin(\langle \boldsymbol{\omega}_D, \mathbf{x} \rangle) \end{pmatrix}$$

$$\xi(\mathbf{x}) = \langle \mathbf{w}_F, \boldsymbol{\phi}_F(\mathbf{x}) \rangle - \langle \mathbf{w}, \boldsymbol{\phi}(\mathbf{x}) \rangle$$

$$\xi(\mathbf{x}) \in [-M, M]$$

$$\mathbb{E}_{\mathbf{x}}[\xi(\mathbf{x})^2] \leq \epsilon_1$$

Minimize $\hat{R}(h)$ to achieve low $R(h)$

Hypothesis class: $\mathcal{H} = \{h(\mathbf{x}) = \langle \mathbf{w}, \boldsymbol{\phi}(\mathbf{x}) \rangle\}$

$$\mathbf{w}^* = \arg \min_{\|\mathbf{w}\|_2 < B} \frac{1}{N_1} \sum_{i=1}^{N_1} (\langle \mathbf{w}, \boldsymbol{\phi}(\mathbf{x}_i) \rangle - y_i)^2$$

$$R(h_{\mathbf{w}^*}) = \mathbb{E}_{\mathbf{x}}[(h_{\mathbf{w}^*}(\mathbf{x}) - f_{\theta}(\mathbf{x}))^2]$$

- Constrained convex optimization
- Kernel ridge regression with line search

Guarantees for empirical risk minimization

PQC models: $f_{\theta}(\mathbf{x}) \in [0,1]$

$$f_{\theta}(\mathbf{x}) = \langle \mathbf{w}, \boldsymbol{\phi}(\mathbf{x}) \rangle + \xi(\mathbf{x})$$

$$\boldsymbol{\phi}(\mathbf{x}) = \frac{1}{\sqrt{D}} \begin{pmatrix} 1 \\ \cos(\langle \boldsymbol{\omega}_1, \mathbf{x} \rangle) \\ \sin(\langle \boldsymbol{\omega}_1, \mathbf{x} \rangle) \\ \vdots \\ \cos(\langle \boldsymbol{\omega}_D, \mathbf{x} \rangle) \\ \sin(\langle \boldsymbol{\omega}_D, \mathbf{x} \rangle) \end{pmatrix}$$

$$\xi(\mathbf{x}) = \langle \mathbf{w}_F, \boldsymbol{\phi}_F(\mathbf{x}) \rangle - \langle \mathbf{w}, \boldsymbol{\phi}(\mathbf{x}) \rangle$$

$$\xi(\mathbf{x}) \in [-M, M]$$

$$\mathbb{E}_{\mathbf{x}}[\xi(\mathbf{x})^2] \leq \epsilon_1$$

Minimize $\hat{R}(h)$ to achieve low $R(h)$

Hypothesis class: $\mathcal{H} = \{h(\mathbf{x}) = \langle \mathbf{w}, \boldsymbol{\phi}(\mathbf{x}) \rangle\}$

$$\text{Lemma 1: } R(h_{\mathbf{w}^*}) \leq \epsilon_1 + \tilde{\mathcal{O}} \left(D^2 \sqrt{\frac{1}{N_1}} \right)$$

Amount of data to learn: $N_1 \in \mathcal{O}(D^4)$

- Can we do better for the number of data?
- No indication on number of shots N_S
- Does a trade off between N_1 and N_S exist?

Classical machine learning models

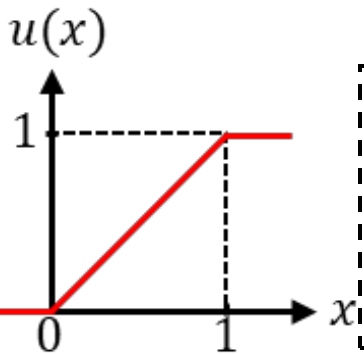
PQC models: $f_{\theta}(x) \in [0,1]$

$$\begin{aligned} f_{\theta}(x) &= \langle \mathbf{w}_F, \boldsymbol{\phi}_F(x) \rangle = \langle \mathbf{w}, \boldsymbol{\phi}(x) \rangle + \xi(x) \\ &= u(\langle \mathbf{w}, \boldsymbol{\phi}(x) \rangle + \xi(x)) \end{aligned}$$

$$\begin{aligned} \xi(x) &= \langle \mathbf{w}_F, \boldsymbol{\phi}_F(x) \rangle - \langle \mathbf{w}, \boldsymbol{\phi}(x) \rangle \\ \xi(x) &\in [-M, M] \\ \mathbb{E}_x[\xi(x)^2] &\leq \epsilon_1 \end{aligned}$$

Hypothesis class: $\mathcal{H} = \{h(x) = u(\langle \mathbf{w}, \boldsymbol{\phi}(x) \rangle)\}$

$$\|\mathbf{w}\|_2 \leq B, \|\boldsymbol{\phi}(x)\|_2 \leq 1$$



$$u(x) = \begin{cases} 0, & \text{if } x < 0, \\ x, & \text{if } 0 \leq x \leq 1, \\ 1, & \text{if } x > 1, \end{cases}$$

$$\boldsymbol{\phi}(x) = \frac{1}{\sqrt{D}} \begin{pmatrix} 1 \\ \cos(\langle \boldsymbol{\omega}_1, x \rangle) \\ \sin(\langle \boldsymbol{\omega}_1, x \rangle) \\ \vdots \\ \cos(\langle \boldsymbol{\omega}_D, x \rangle) \\ \sin(\langle \boldsymbol{\omega}_D, x \rangle) \end{pmatrix}$$

The learning algorithm

Algorithm 1: The learning algorithm

Input: Training data size N_1 , validation data size N_2 , number of measurement shots N_s , parameter setting of quantum model θ , distribution of input $p(\mathbf{x})$, non-decreasing L -Lipschitz function $u : \mathbb{R} \rightarrow \mathcal{Y}$, kernel function k corresponding to feature map ϕ , learning rate $\lambda > 0$, number of iterations T

1 Sample N_1 training data inputs $\mathbf{x}_1, \dots, \mathbf{x}_{N_1} \sim p(\mathbf{x})$.

Just think of it as kernelized gradient descent with a validation dataset.

2 Sample training dataset $(\mathbf{x}_i, y_i)_{i=1}^{N_1}$.
3 Repeat steps 1 and 2 to collect labelled validation data of size N_2 , $(\mathbf{p}_j, \bar{q}_j)_{j=1}^{N_2}$.

4 Initialize $\alpha^i := 0 \in \mathbb{R}^{N_1}$.

5 **for** $t = 1, \dots, T$ **do**

6 $h^t(\mathbf{x}) := u\left(\sum_{i=1}^{N_1} \alpha_i^t k(\mathbf{x}, \mathbf{x}_i)\right)$
7 **for** $i = 1, 2, \dots, N_1$ **do**
8 $\alpha_i^{t+1} := \alpha_i^t + \frac{\lambda}{N_1} (\bar{y}_i - h^t(\mathbf{x}_i))$

Output: h^r where

$$r = \arg \min_{t \in \{1, \dots, T\}} \frac{1}{N_2} \sum_{j=1}^{N_2} (\bar{q}_j - h^t(\mathbf{p}_j))^2$$

Provable guarantee on concept learning

Apply gradient descent* on data to achieve low $R(h)$.

$$\text{Theorem 1: } R(h) \leq \tilde{O} \left(\sqrt{\epsilon_1} + M^4 \sqrt{\frac{1}{N_1}} + D \sqrt{\frac{1}{N_1}} + D \sqrt{\frac{1}{N_1 N_s}} \right)$$

$$\text{Lemma 1: } R(h_{w^*}) \leq \epsilon_1 + \tilde{O} \left(D^2 \sqrt{\frac{1}{N_1}} \right)$$

Provable guarantee on concept learning

Apply gradient descent* on data to achieve low $R(h)$.

$$\text{Theorem 1: } R(h) \leq \tilde{O} \left(\sqrt{\epsilon_1} + M^4 \sqrt{\frac{1}{N_1}} + D \sqrt{\frac{1}{N_1}} + D \sqrt{\frac{1}{N_1 N_s}} \right)$$

$$\mathbb{E}_x[\xi(x)^2] \leq \epsilon_1 \quad \xi(x) \in [-M, M]$$

$$\xi(x) = \langle \mathbf{w}_F, \phi_F(x) \rangle - \langle \mathbf{w}, \phi(x) \rangle$$

$$f_\theta(x) = \langle \mathbf{w}_F, \phi_F(x) \rangle \quad h(x) = u(\langle \mathbf{w}, \phi(x) \rangle)$$

Quantum models cannot be learned without a (fairly) efficient and good classical approximation.

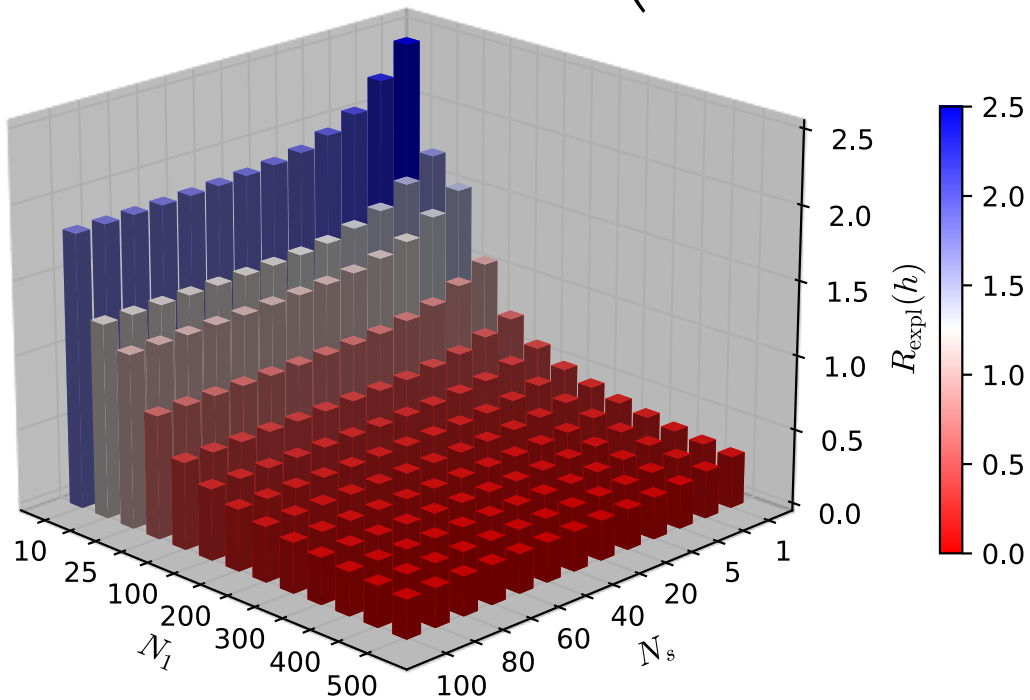
Bias from the model representation.

Asymmetrical effects of N_1 and N_s

Corollary 1: $R(h) \leq \tilde{O} \left(D \sqrt{\frac{1}{N_1}} + D \sqrt{\frac{1}{N_1 N_s}} \right)$

$N_1 \rightarrow \infty$
 $R(h) \leq 0$

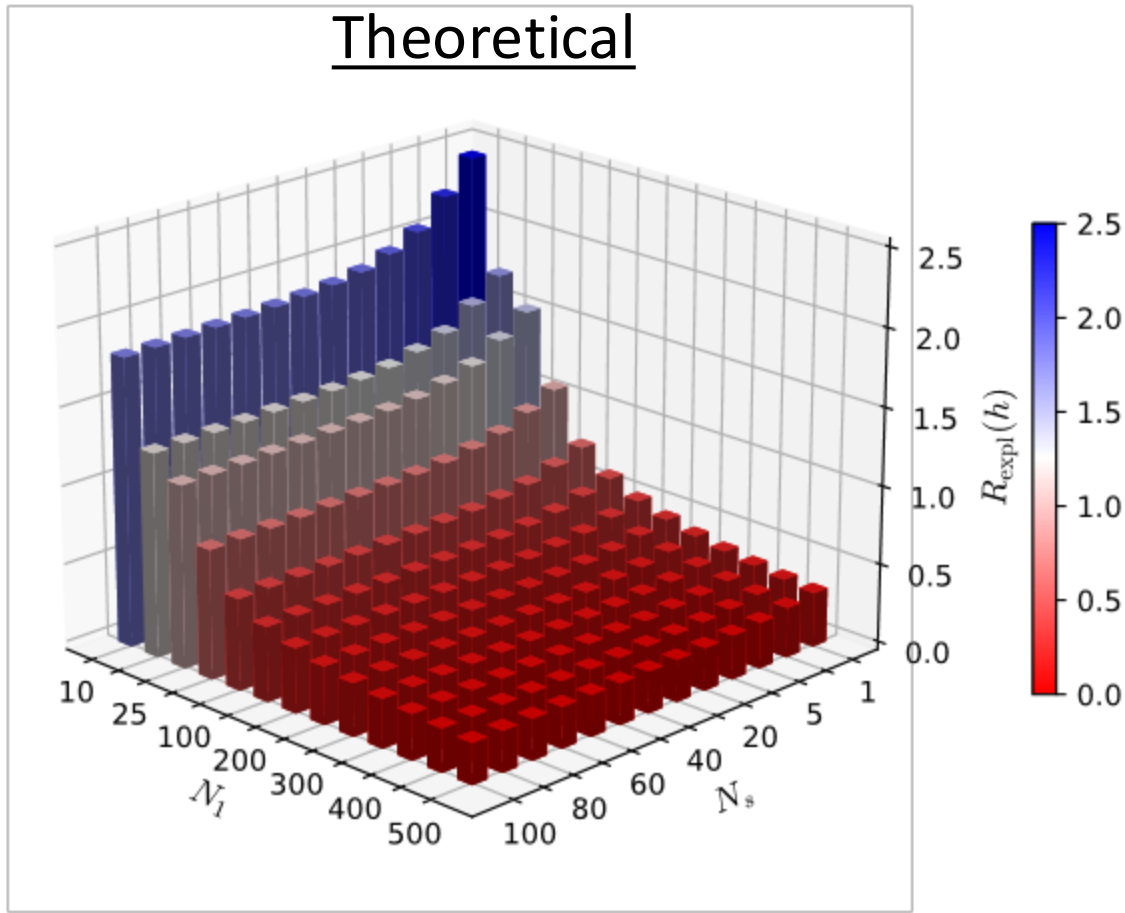
$N_s \rightarrow \infty$
 $R(h) \leq \tilde{O} \left(D \sqrt{\frac{1}{N_1}} \right)$



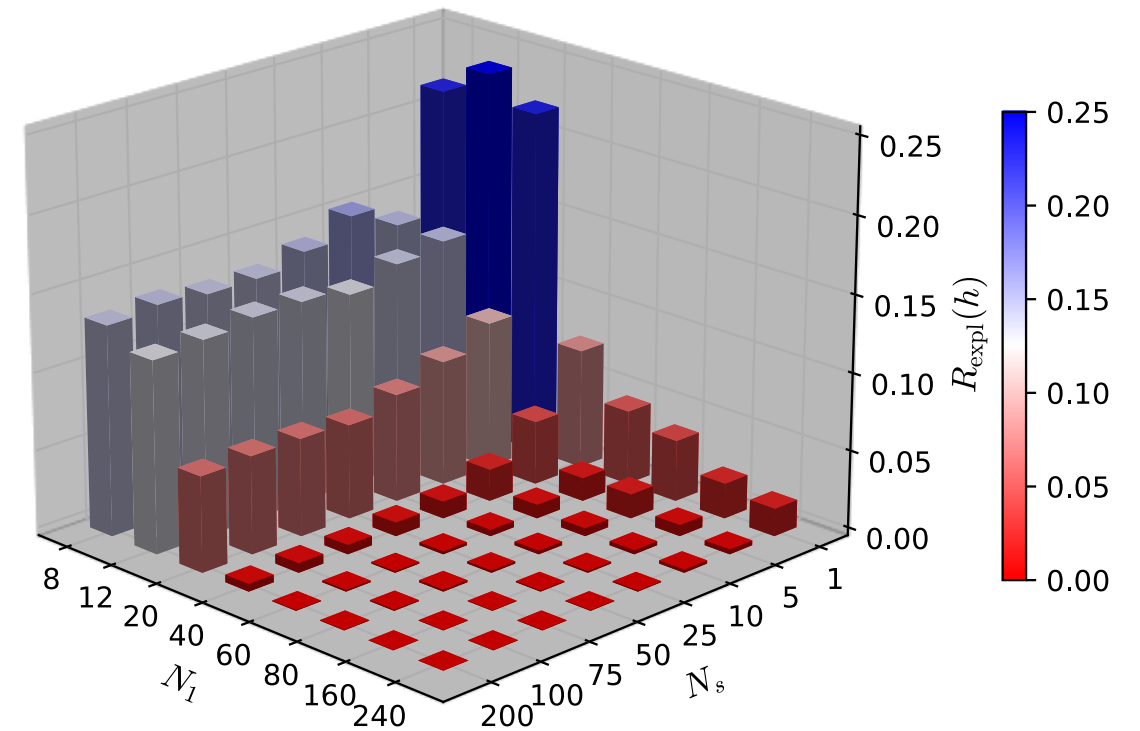
N_1 and N_s have asymmetrical effects on the model's learning performance.

Asymmetrical effects of N_1 and N_s

Theoretical

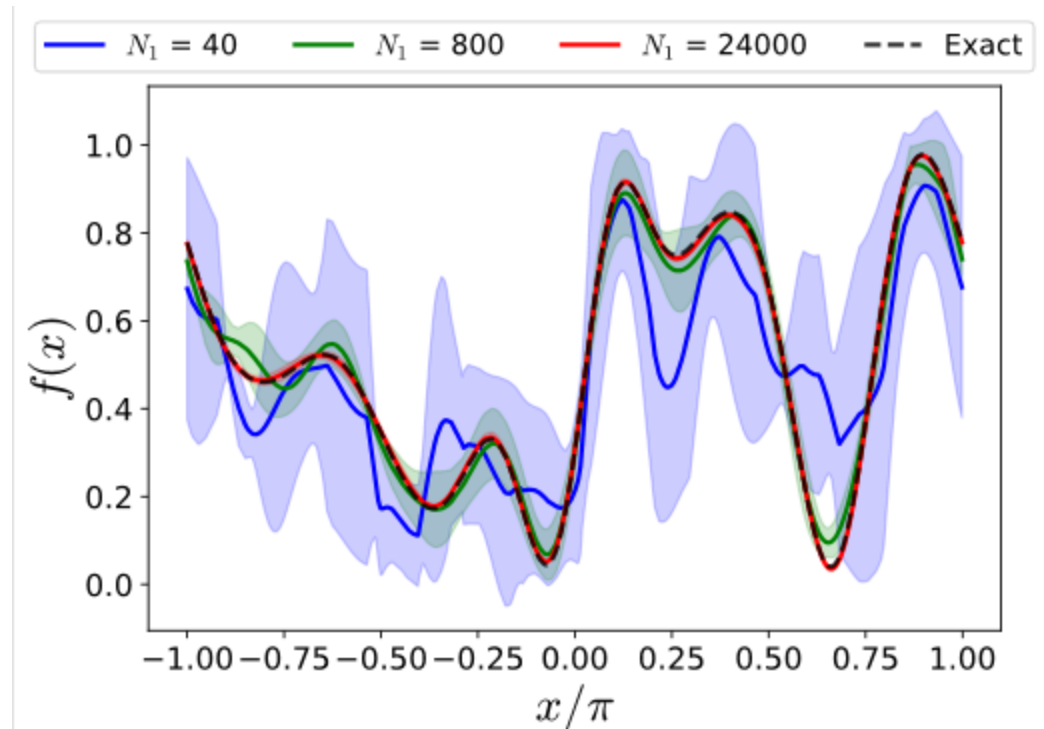
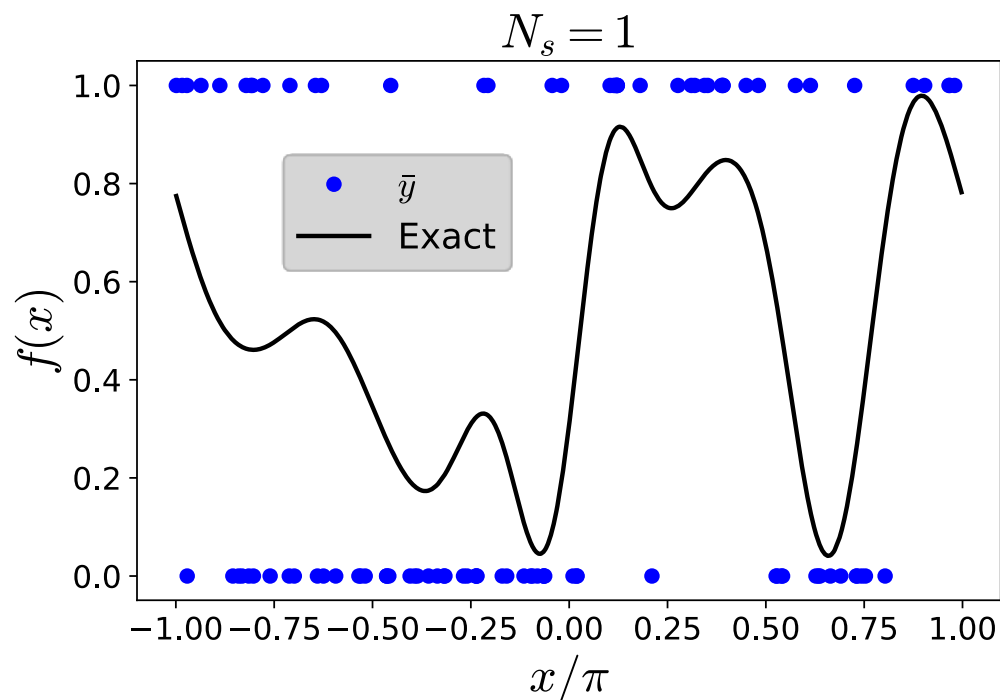


Empirical

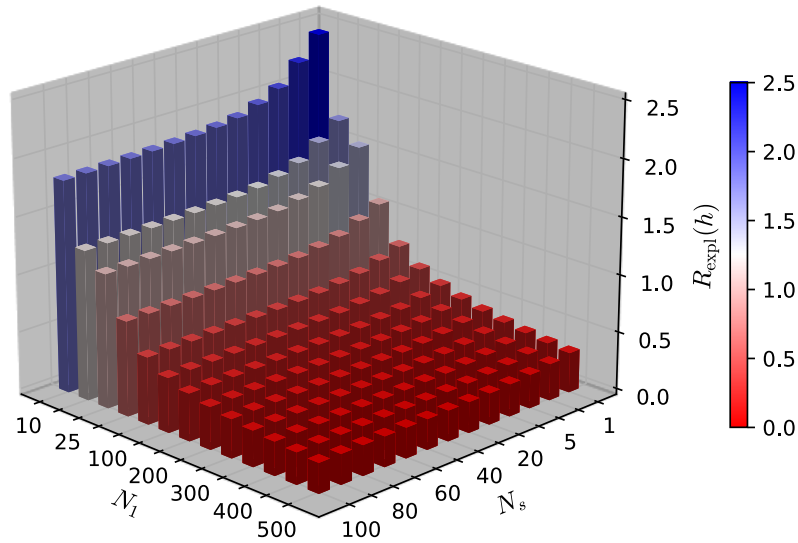


Learning PQC models with $N_s = 1$

$$\text{Corollary 1: } R(h) \leq \tilde{O} \left(D \sqrt{\frac{1}{N_1}} + D \sqrt{\frac{1}{N_1 N_s}} \right) \xrightarrow{N_1 \rightarrow \infty} R(h) \leq 0$$



Trade-off between N_1 and N_s



Assumed we have unlimited queries to quantum systems.

Access to quantum computers is expensive.

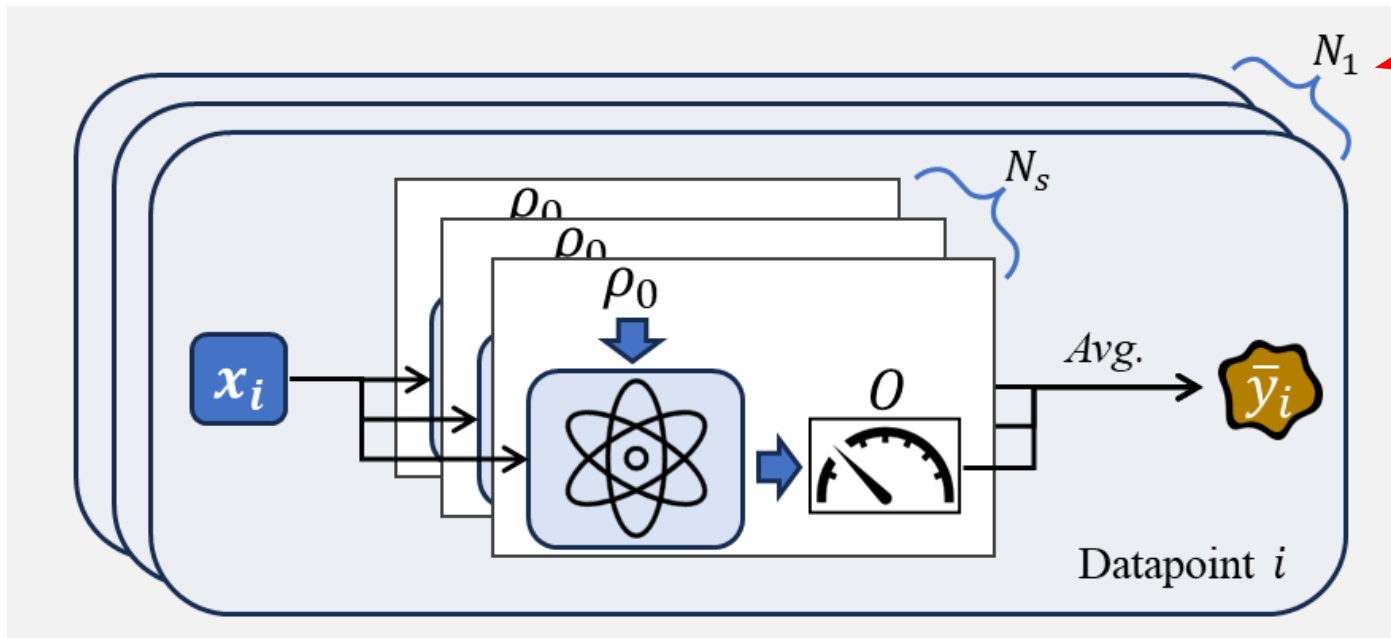
Hardware Provider	QPU family	Per-task price	Per-shot price
IonQ	Harmony	\$0.30000	\$0.01000
IonQ	Aria	\$0.30000	\$0.03000
IQM	Garnet	\$0.30000	\$0.00145
QuEra	Aquila	\$0.30000	\$0.01000
Rigetti	Aspen-M	\$0.30000	\$0.00035

Some hardware is more expensive than others.

Trade-off between N_1 and N_s

Limit the total number of queries to quantum systems.

Total measurement budget: $N_{\text{tot}} = N_1 \cdot (N_s + \gamma)$



Extra $\gamma \in \mathbb{R}^+$ for changing parameter settings

- $\gamma = 0: N_{\text{tot}} = N_1 \cdot N_s$
- $\gamma_{\text{Trapped Ions}} > \gamma_{\text{Superconducting}}$

Trade-off between N_1 and N_s

Corollary 1:

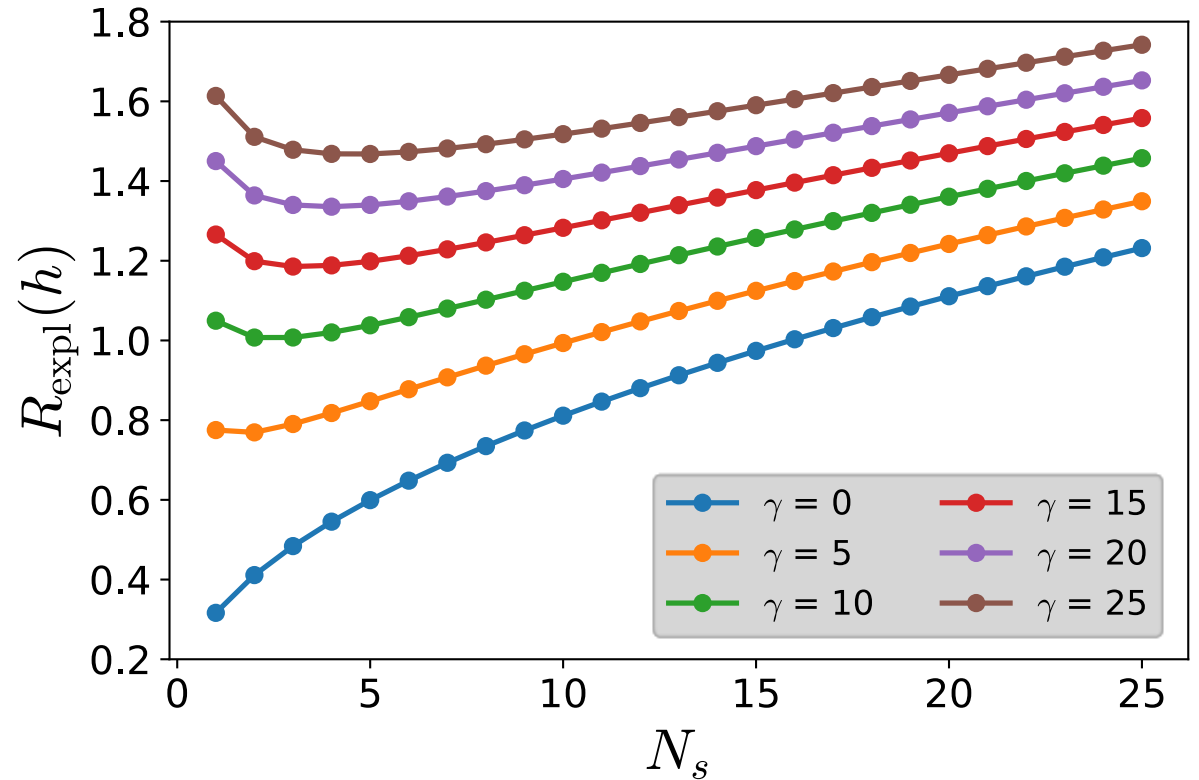
$$R(h) \leq \tilde{\mathcal{O}} \left(D \sqrt{\frac{1}{N_1}} + D \sqrt{\frac{1}{N_1 N_s}} \right)$$

$N_{\text{tot}} = N_1 \cdot (N_s + \gamma)$

$$\frac{1}{N_1} = \frac{N_s + \gamma}{N_{\text{tot}}}$$

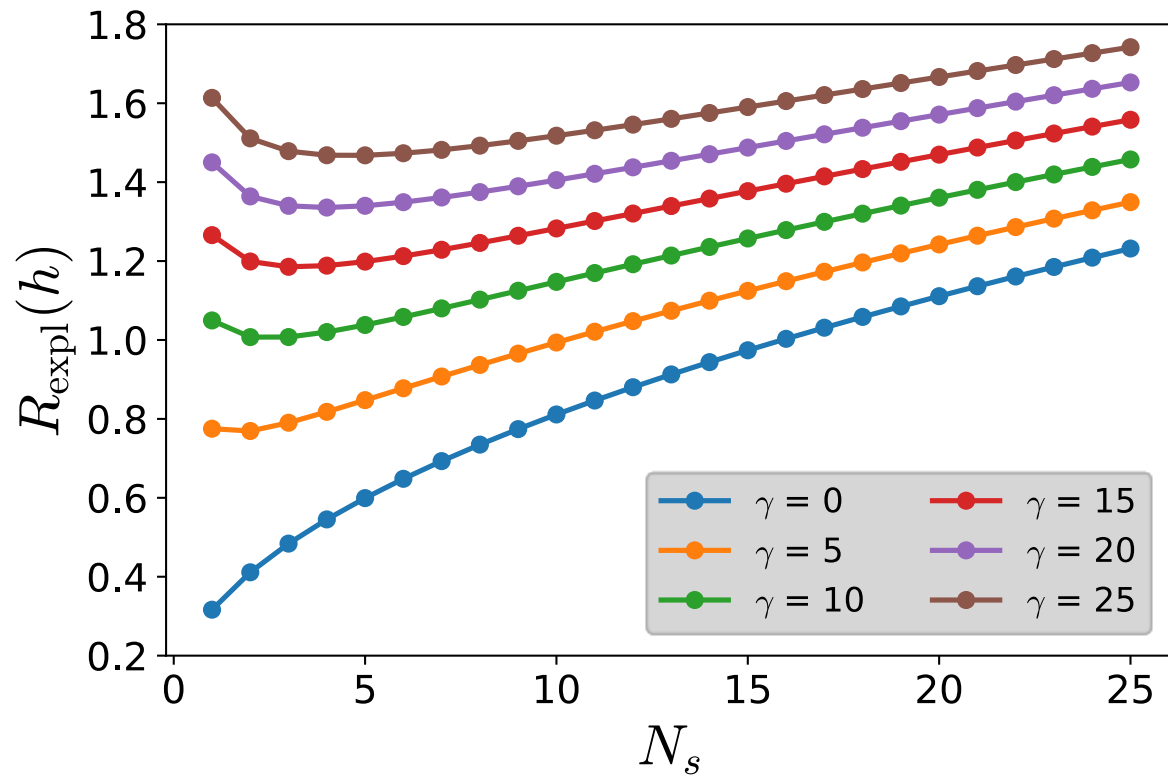
Corollary 2:

$$R(h) \leq \tilde{\mathcal{O}} \left(D \sqrt{\frac{N_s + \gamma}{N_{\text{tot}}}} + D \sqrt{\frac{N_s + \gamma}{N_{\text{tot}} N_s}} \right)$$

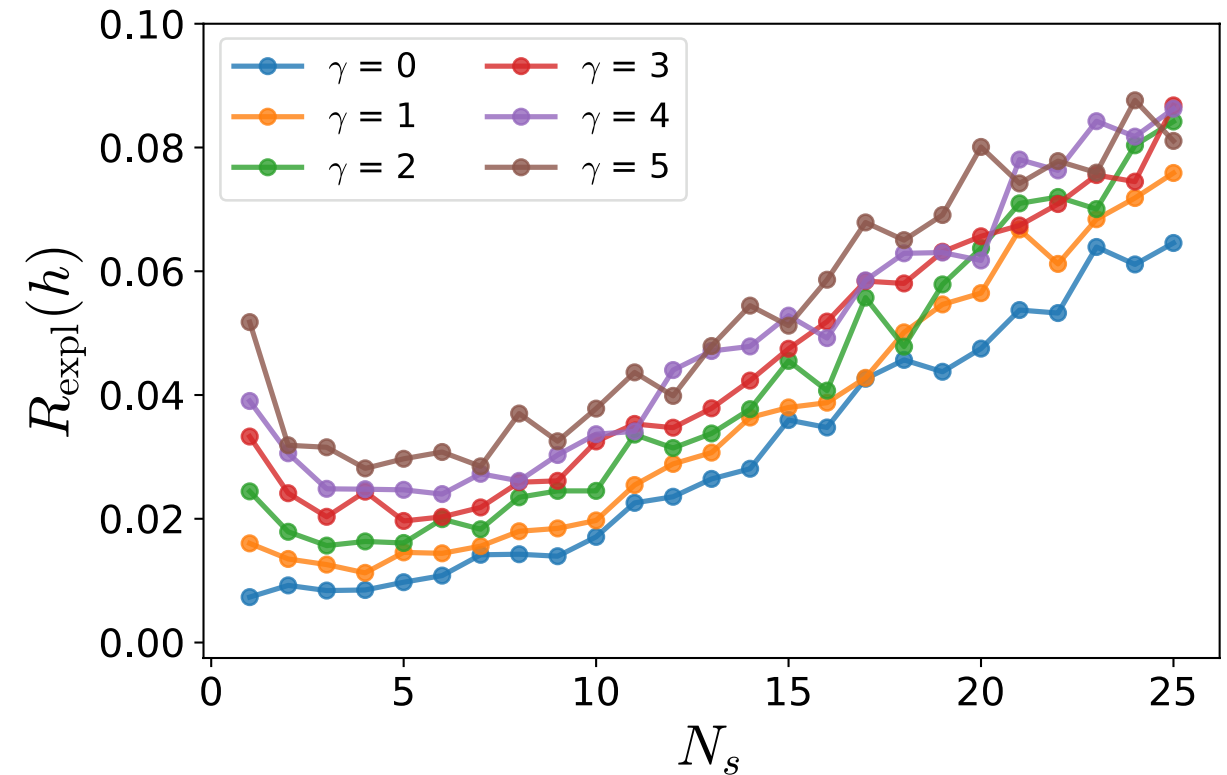


Trade-off between N_1 and N_s

Theoretical



Empirical



Role of the link function u

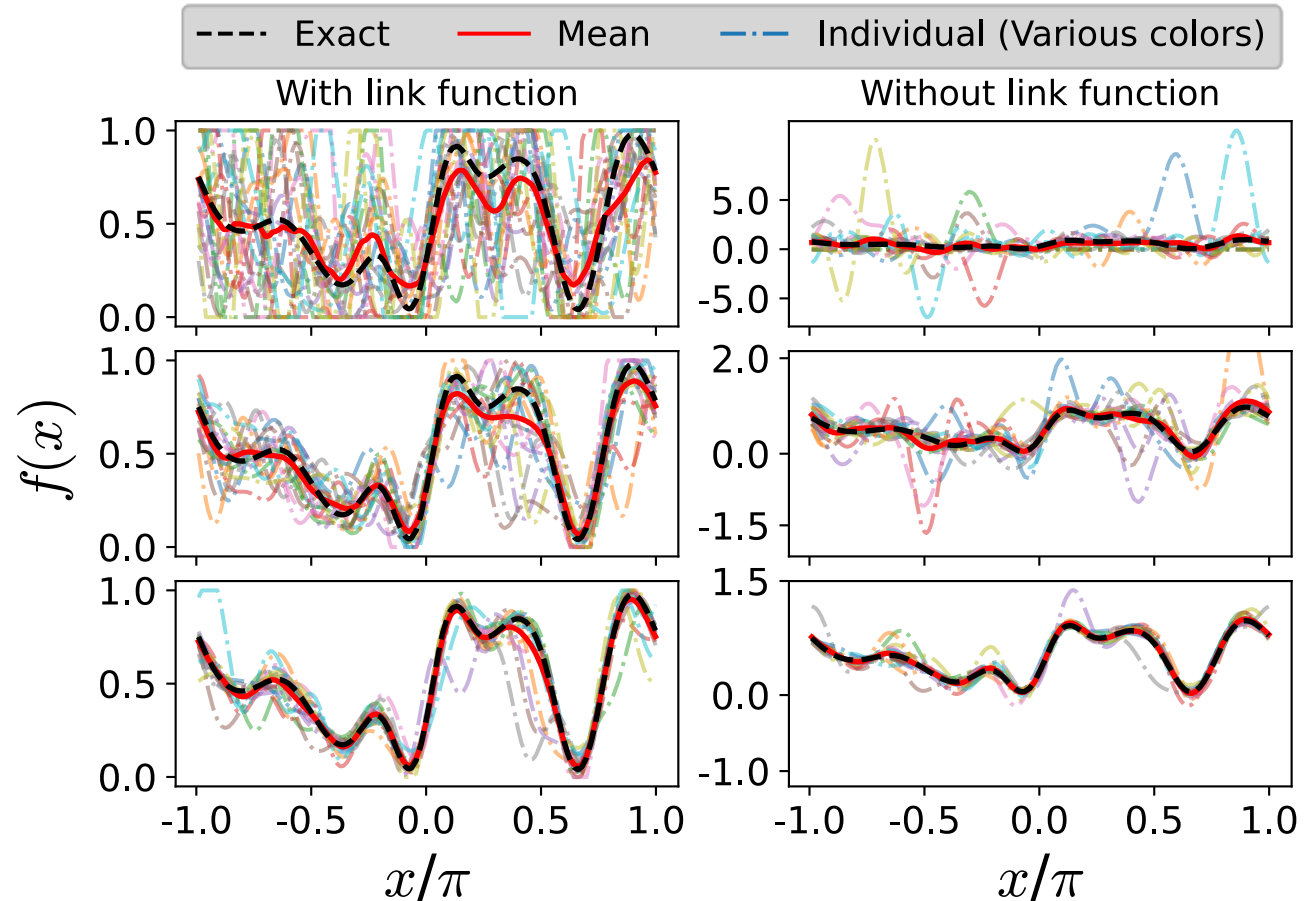
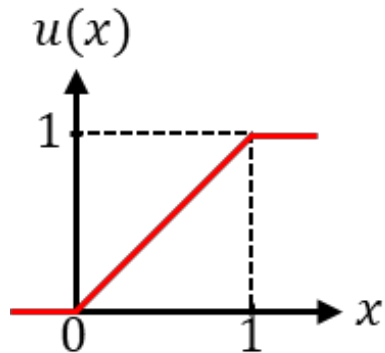
With link function

$$h(x) = \underbrace{u}_{\text{link function}}(\langle \mathbf{w}, \phi(x) \rangle)$$

Without link function

$$g(x) = \langle \mathbf{w}, \phi(x) \rangle$$

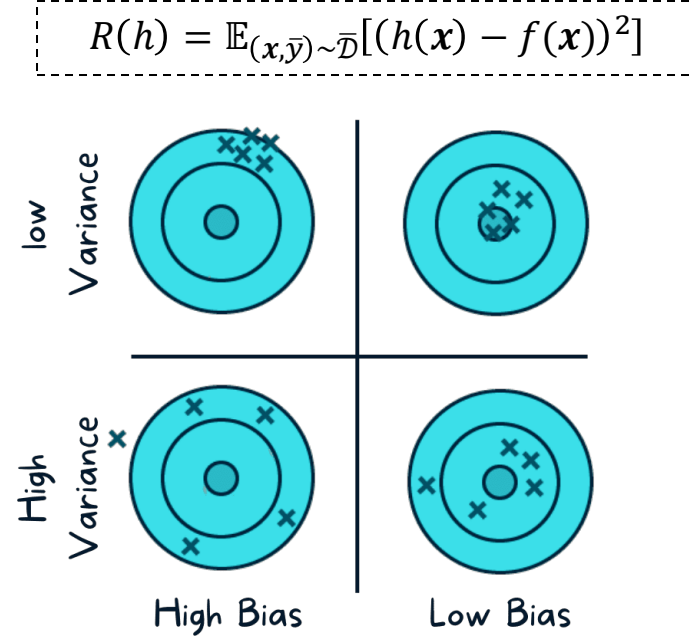
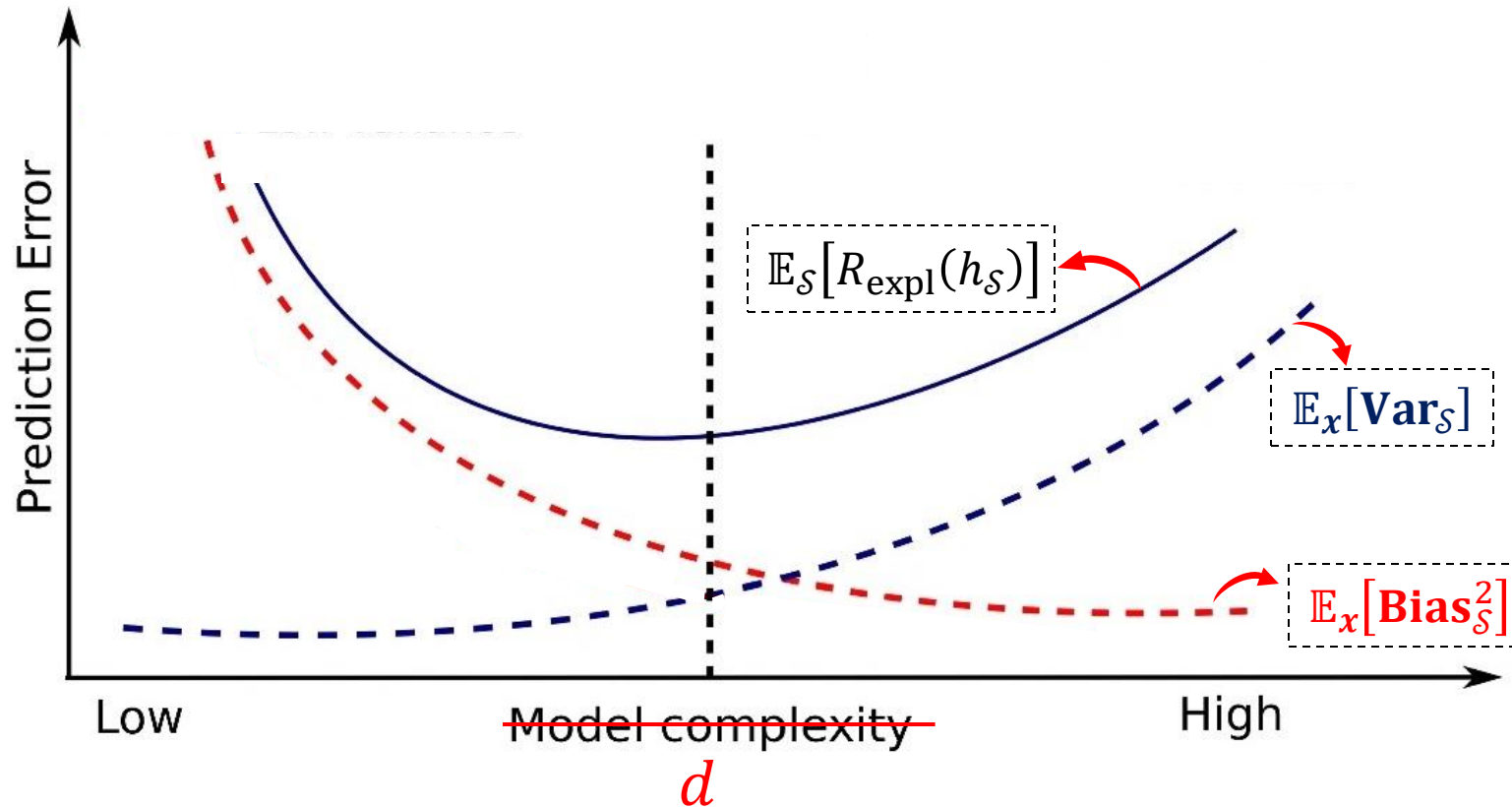
$$u(x) = \begin{cases} 0, & \text{if } x < 0, \\ x, & \text{if } 0 \leq x \leq 1, \\ 1, & \text{if } x > 1, \end{cases}$$



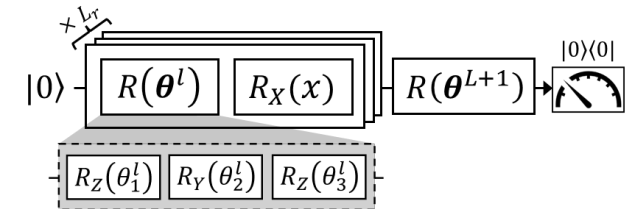
Bias-variance trade-off

Analyze the average $R(h)$ over all possible training datasets

$$\mathbb{E}_{\mathcal{S}}[R(h_{\mathcal{S}})] = \mathbb{E}_x[\mathbf{Bias}_{\mathcal{S}}^2] + \mathbb{E}_x[\mathbf{Var}_{\mathcal{S}}]$$



Data re-uploading models:



Concept class:

$$\mathcal{F}_{L_r} = \left\{ f_{\theta}(x) = c_0(\theta) + \sum_{\omega=1}^{L_r} a_{\omega}(\theta) \cos(\omega x) + b_{\omega}(\theta) \sin(\omega x) \right\}$$

Hypothesis class:

$$\mathcal{H}_d = \left\{ h_d(x) = u \left(v_0 + \sum_{\omega=1}^d \alpha_{\omega} \cos(\omega x) + \beta_{\omega} \sin(\omega x) \right) \right\}$$

Provable guarantee on concept learning

Apply gradient descent* on data to achieve low $R(h)$.

$$\text{Theorem 1: } R(h) \leq \tilde{O} \left(\underbrace{\sqrt{\epsilon_1} + M^4 \sqrt{\frac{1}{N_1}}}_{\text{Bias}} + \underbrace{D \sqrt{\frac{1}{N_1}} + D \sqrt{\frac{1}{N_1 N_s}}}_{\text{Variance}} \right)$$

$$\text{Lemma 1: } R(h_{w^*}) \leq \epsilon_1 + \tilde{O} \left(D^2 \sqrt{\frac{1}{N_1}} \right)$$

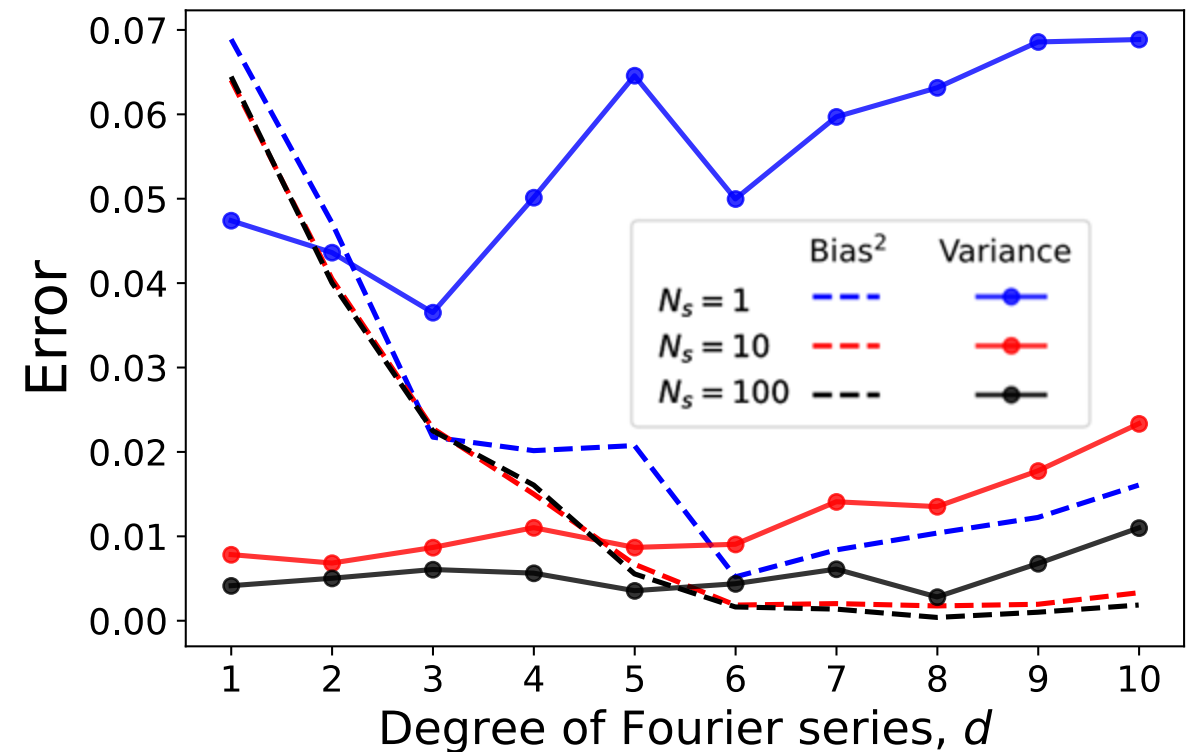
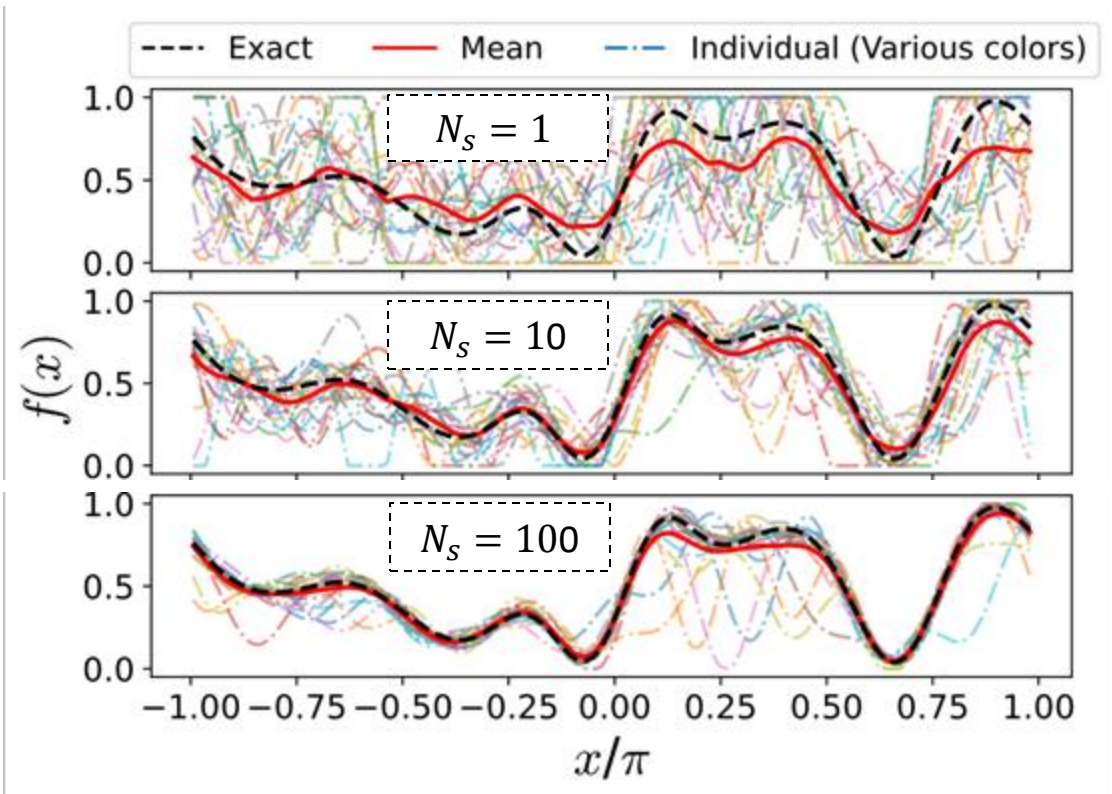
The link function limits the model complexity and forcibly compresses the variance upper bound.

$$R_{\text{expl}}(h) = \mathbb{E}_{(x,y) \sim \bar{\mathcal{D}}}[(h(x) - f(x))^2]$$

Effects of shot noise on bias and variance

Analyze the average $R_{\text{expl}}(h)$ over all possible training datasets

$$\mathbb{E}_{\mathcal{S}}[R_{\text{expl}}(h_{\mathcal{S}})] = \mathbb{E}_x[\mathbf{Bias}_{\mathcal{S}}^2] + \mathbb{E}_x[\mathbf{Var}_{\mathcal{S}}]$$



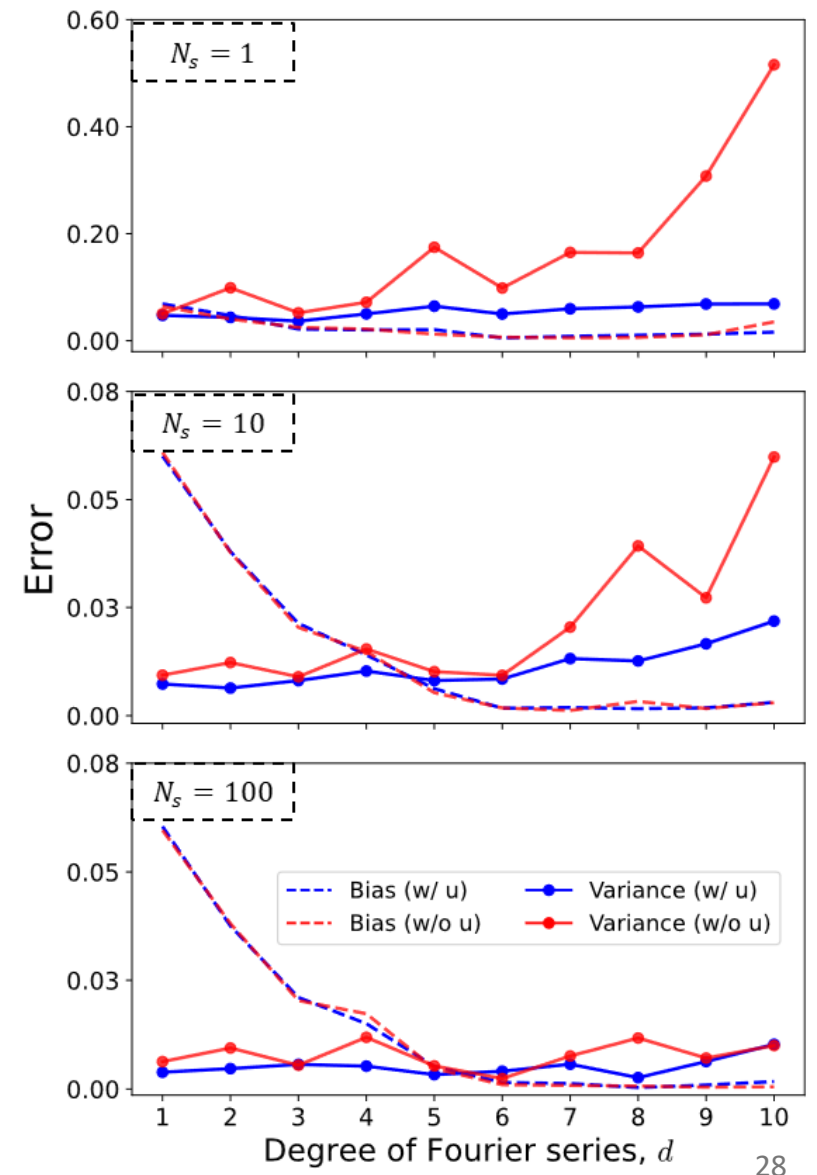
Learning with and without the link function

Hypothesis class (with link function):

$$\mathcal{H}_d = \left\{ h_d(x) = u \left(v_0 + \sum_{\omega=1}^d \alpha_{\omega} \cos(\omega x) + \beta_{\omega} \sin(\omega x) \right) \right\}$$

Hypothesis class (without link function):

$$G_d = \left\{ g_d(x) = v_0 + \sum_{\omega=1}^d \alpha_{\omega} \cos(\omega x) + \beta_{\omega} \sin(\omega x) \right\}$$



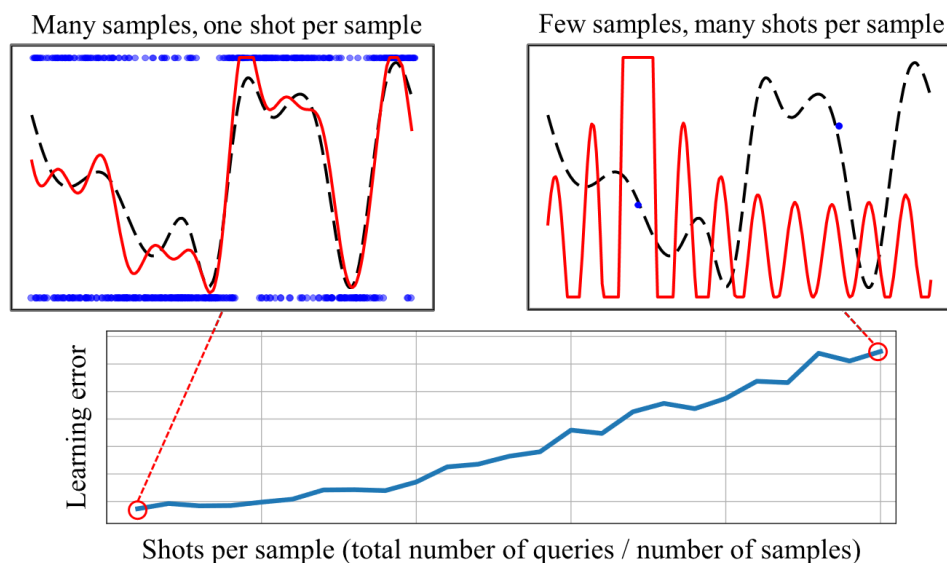
Conclusion

Can we obtain provable guarantees of learning that exemplify the relationship between N_1 and N_s ?



arXiv:2408.05116 [quant-ph]

(1) Asymmetrical trade-offs between N_1 and N_s



(2) Gradient descent* can be made robust and provide tighter guarantees

