

QML-Mod5-QML Algorithm Complexity

Riccardo Marega

March 2025

Indice

1 QML Algorithm Complexity -13/06/2025	2
2 QML: practical aspects -14/06/2025	3
3 Noise characterization and mitigation -28/06/2025	50

1 QML Algorithm Complexity -13/06/2025

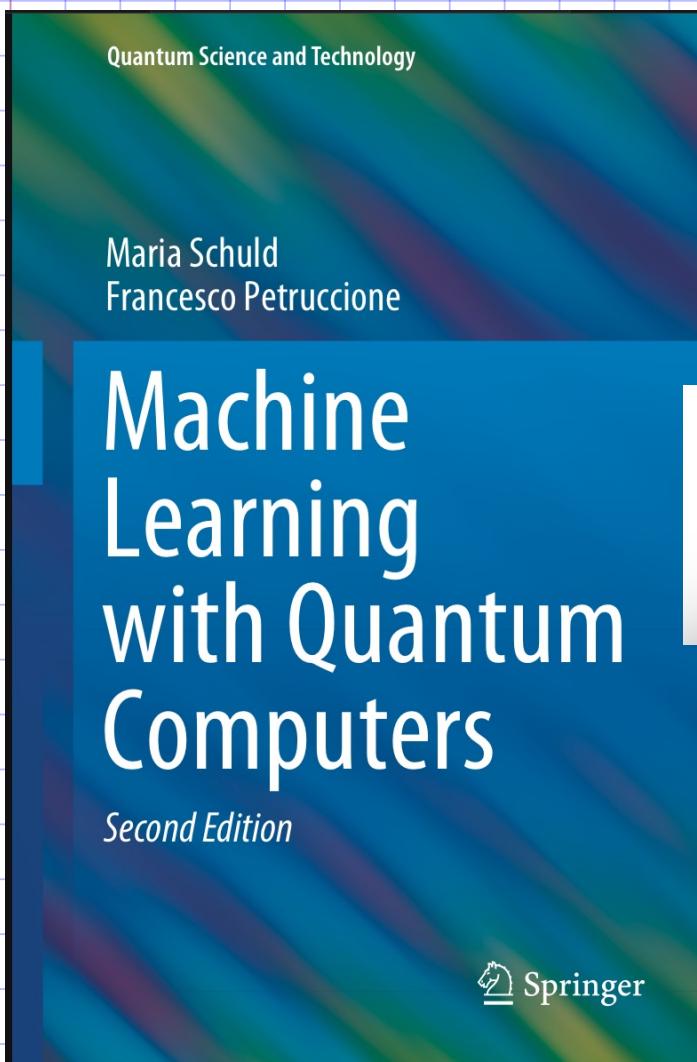
All notes can be found on the ipynb.

2 QML: practical aspects -14/06/2025

QUANTUM MACHINE LEARNING

PRACTICAL ASPECTS

classical
C.S.



PROCEEDINGS A

rspa.royalsocietypublishing.org

Review



Cite this article: Ciliberto C, Herbster M, Lalongo AD, Pontil M, Rocchetto A, Severini S, Wossnig L. 2018 Quantum machine learning: a classical perspective. *Proc. R. Soc. A* **474**: 20170551. <http://dx.doi.org/10.1098/rspa.2017.0551>

Quantum machine learning:
a classical perspective

Carlo Ciliberto¹, Mark Herbster¹, Alessandro Davide Lalongo^{2,3}, Massimiliano Pontil^{1,4}, Andrea Rocchetto^{1,5}, Simone Severini^{1,6} and Leonard Wossnig^{1,7}

¹Department of Computer Science, University College London, London, UK

²Department of Engineering, University of Cambridge, Cambridge, UK

³Max Planck Institute for Intelligent Systems, Tübingen, Germany

<https://doi.org/10.1088/s42254-022-00552-1>

nature reviews physics



Review article



Learning quantum systems

Valentin Gebhart^{1,2}, Raffaele Santagati³, Antonio Andrea Gentile⁴, Erik M. Gauger⁵, David Craig⁶, Natalia Ares⁷, Leonardo Banchi^{8,9}, Florian Marquardt^{10,11}, Luca Pezzè¹² & Cristian Bonato⁶

Check for updates

<https://doi.org/10.1038/s42254-022-00552-1>

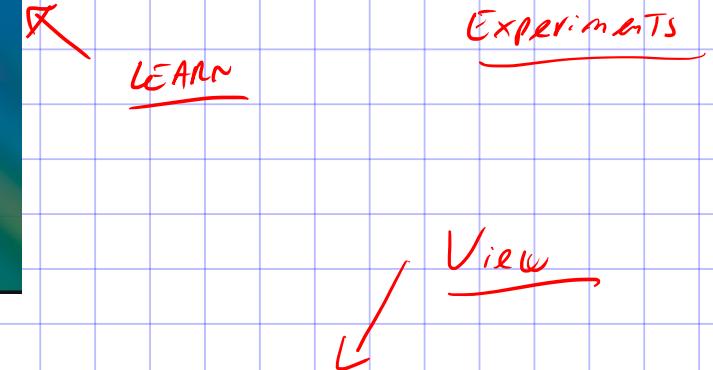
Check for updates

<https://doi.org/10.1038/s42254-022-00552-1>

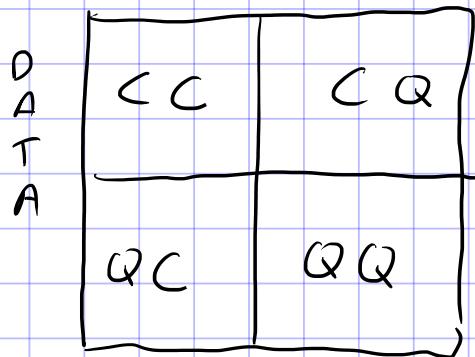
Quantum computing and artificial intelligence: status and perspectives

Giovanni Acampora,¹ Andris Ambainis,² Natalia Ares,³ Leonardo Banchi,^{4,5} Pallavi Bhardwaj,⁶ Daniele Binosi,^{7,8} G. Andrew D. Briggs,⁹ Tommaso Calarco,^{7,10} Vedran Dunjko,¹¹ Jens Eisert,¹² Olivier Ezratty,^{13,14} Paul Erker,^{15,16} Federico Fedele,³ Elies Gil-Fuster,^{12,17} Martin Gärttner,¹⁸ Matt Granath,¹⁹ Markus Heyl,^{20,21} Iordanis Kerenidis,²² Matthias Klusch,²³ Anton Frisk Kockum,²⁴ Richard Kueng,²⁵ Mario Krenn,²⁶ Jörg Lässig,^{27,28} Antonio Macaluso,²³ Sabrina Maniscalco,²⁹ Florian Marquardt,³⁰ Kristel Michelsen,¹⁰ Gorka Muñoz-Gil,³¹ Daniel Müsing,²⁸ Hendrik Poulsen Nautrup,³¹ Evert van Nieuwenburg,^{11,32} Roman Orus,³³ Jörg Schmiedmayer,¹⁵ Markus Schmitt,^{10,34} Philipp Slusallek,^{35,23} Filippo Vicentini,^{36,37,38} Christof Weitenberg,³⁹ and Frank K. Wilhelm^{35,10}

¹University of Naples Federico II, 80126 Naples, Italy



DATA PROCESSING



C : Classical
Q : Quantum

CC : classical Data

Classical Computer
(Classical / Quantum models)

Tensor Network Methods
Dequantized QMC

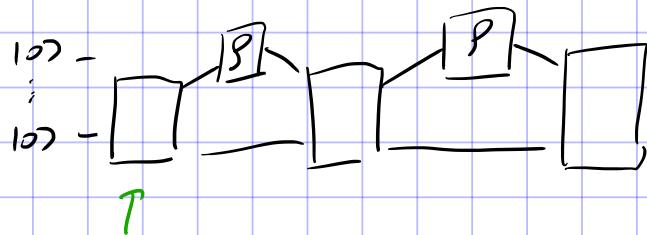
QC : "Quantum Data"

Wave Function collapse

No cloning Theorem

QQ : Quantum Parallelism

✓ STO_n



Quantum operations

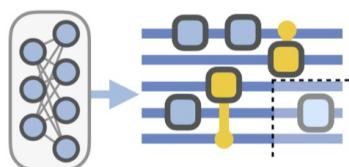
Exponential Advantage

CQ

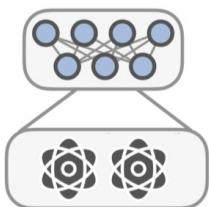
-

Classical Data
Quantum Computer
(images)

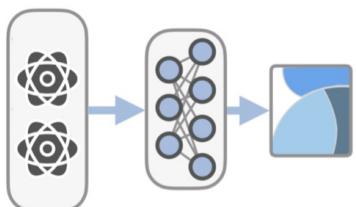
A. Discovery and optimization



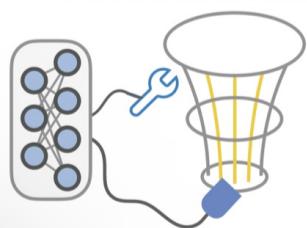
B. Simulation of quantum systems



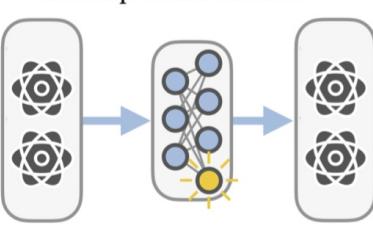
C. Analysis of quantum data



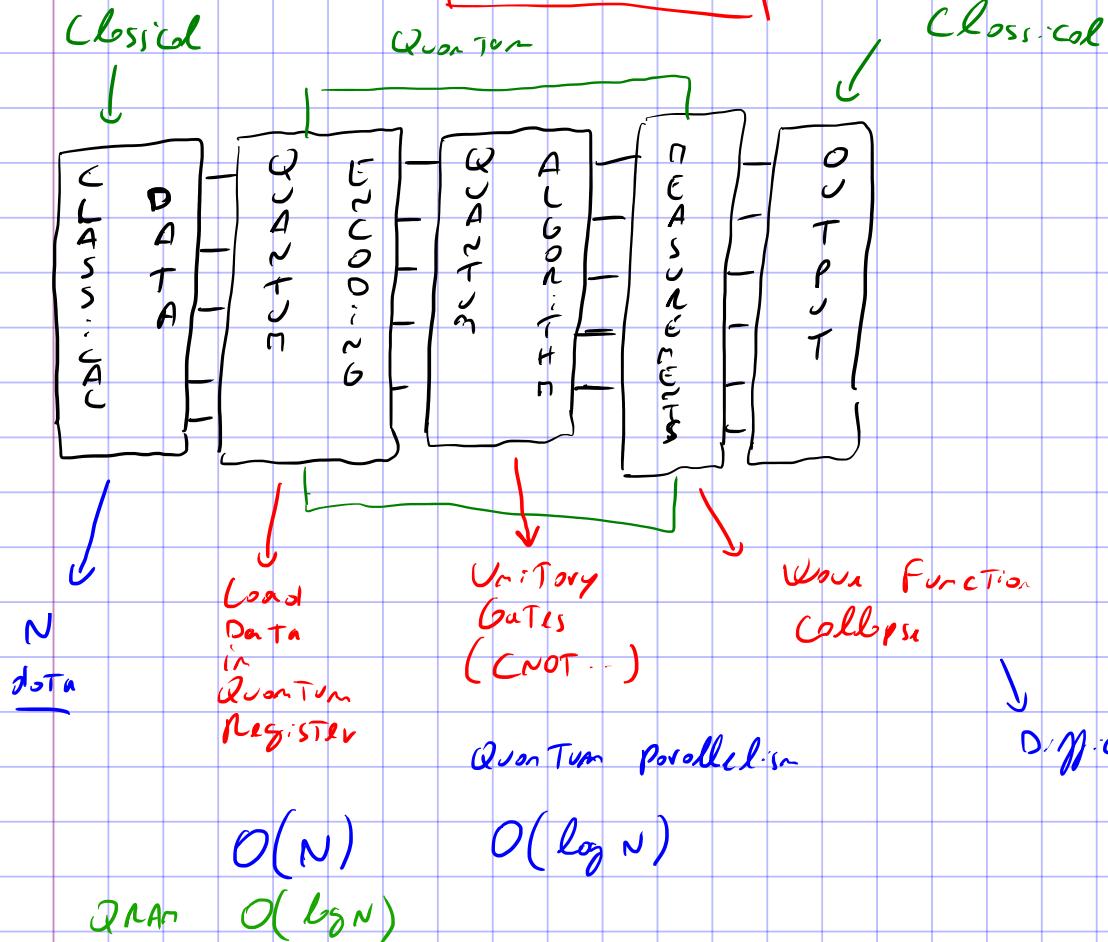
D. Control and certification



E. Interpretable schemes



CQ



Quantum Encoding / Data Loading

Classical data

$$\{x_i\}_{i=1 \dots N}$$

$$x_i \in \mathbb{R}^d$$

$$x \Rightarrow |\psi(x)\rangle$$

Examples

Amplitude Encoding

$$|\psi(x_i)\rangle = \sum_j \underset{\uparrow}{(x_i)_j} |j\rangle$$

$$|j\rangle = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ j \\ 0 \end{pmatrix} \leftarrow j \in \mathbb{Z}$$

$$\prod_{j=1}^d e^{i \hat{G}(x_j) |0\rangle}$$

Hermitian operator / Pauli matrices

Parallel Quantum Strategies

Quantum superposition

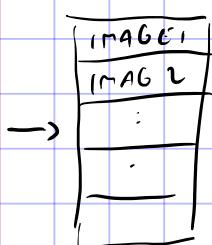
$$\{x_i\}_{i=1 \dots N} \rightarrow \sum_{i=1}^N |\psi(x_i)\rangle$$

Quantum RAM

$$\{x_i\}_{i=1 \dots N} \rightarrow \sum_{i=1}^N |i\rangle |\psi(x_i)\rangle$$

Address

Memory content



ENSEMBLE STRATEGIES

$$\{x_i\} \Rightarrow \frac{1}{N} \sum_{i=1}^N |\psi(x_i)\rangle\langle\psi(x_i)|$$



No Advantages



DENSITY MATRIX

TEST

QUANTUM

IDEAS

Quantum Computing Paradigms

Gate Based

(QUBITS)

Fault Tolerant
Quantum Computer

NISQ

Noisy Intermediate
Semi Quantum Computers

Quantum Annealers



NOT Considered

Continuous Variable / Photonics



Here

FAULT TOLERANT

- Perfect Gates
- Many Qubits
- Algorithms: Everything in Quantum language

Quantum Circuits /
Shor (discrete log)
Grover (database search)

Theory : Asymptotics

$$O(\log N)$$

$$N = 2^m$$

$$O(\log n)$$

Quantum

$$\text{poly}(m)$$

vs

Classical

$$O(2^m) = O(n)$$

QML:

Grover

\downarrow
poly advantage

$$O(n) \rightarrow O(\sqrt{n})$$

HTL (Linear system Eq.)

QSVD \Rightarrow Ex. adv.

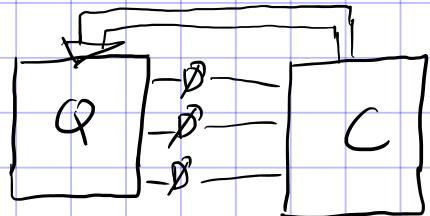
NISQ

Noisy Outcomes

10 - 100 Qubits (Free)

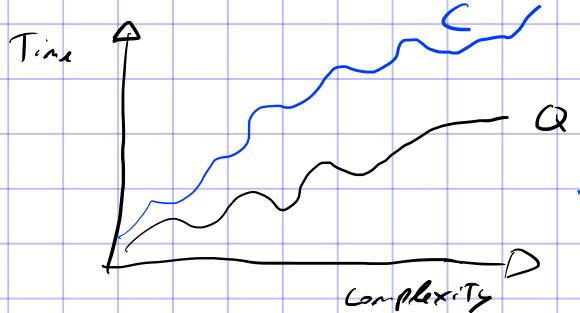
Quantum Resources are Limited

Only use Tel. Quantum Machine for Hard Tasks



Hybrid Quantum / Classical Computation

"Hands On"



QML - Quantum Neural Networks

- Quantum Kernel Methods

"PRACTICAL" QIC



Quantum

Technology

Qiskit

Research

Pricing

Blog

Community

Resources

Sign in to Platform

The new IBM Quantum Platform is here! The current experience will sunset on July 1. [Prepare for the transition →](#)

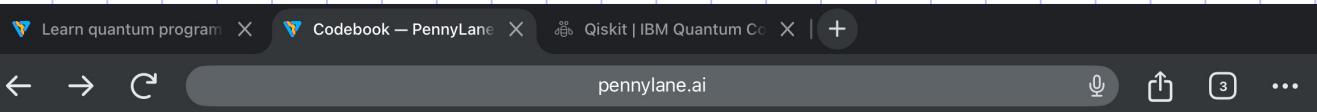
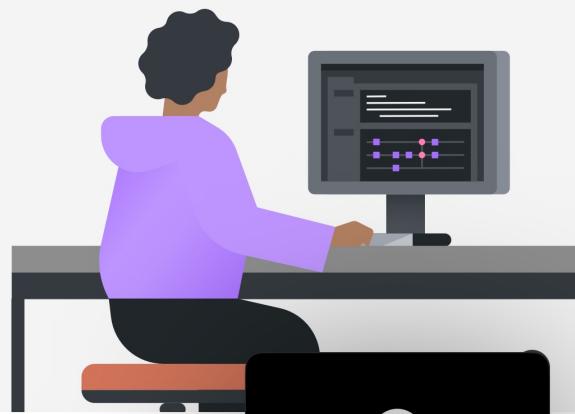
Qiskit SDK v2.0 is here
See what's new →

Qiskit

Qiskit is the world's most popular software stack for quantum computing. Build circuits, leverage Qiskit functions, transpile with AI tools, and execute workloads in an optimized runtime environment.

What is Qiskit?

Get started



Sign in



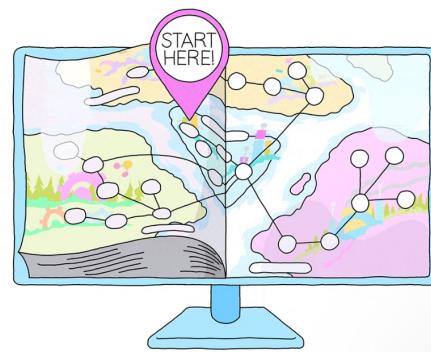
Codebook

Learn quantum computing with PennyLane — the leading tool for programming quantum computers. Explore a specific module or follow a guided path to build your skills step-by-step.

Search

Open Codebook Map

Browse Modules



Quantum Gates

$$e^{i\theta \hat{\sigma}_z} = R_z(\theta)$$

UNIVERSALITY

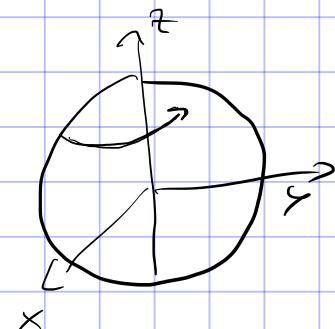


Single Qubit

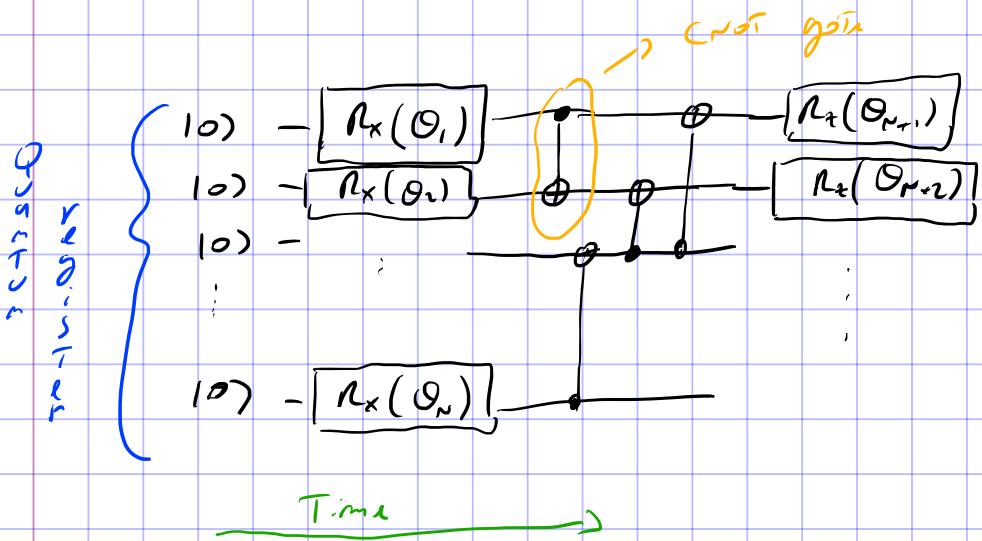
Rotations

R_x, R_y, R_z

Entangling gate ($CNOT$)



Quantum Circuit



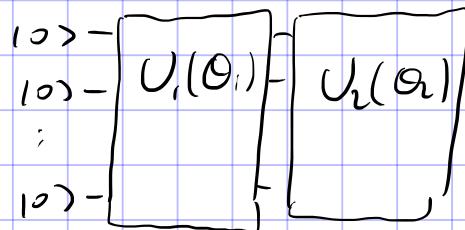
Parametric Quantum Circuit

$$|\Psi(\underline{\theta})\rangle$$

$$\underline{\theta} = \begin{pmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_L \end{pmatrix}$$

parameters

$$|\Psi(\underline{\theta})\rangle$$



$$U_L(\underline{\theta})$$

$$L = \# \text{ Layers}$$

problem dependent "cost function"

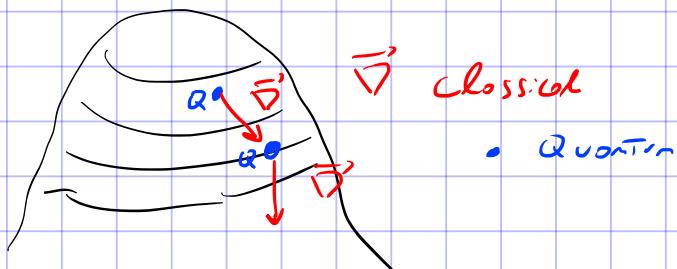
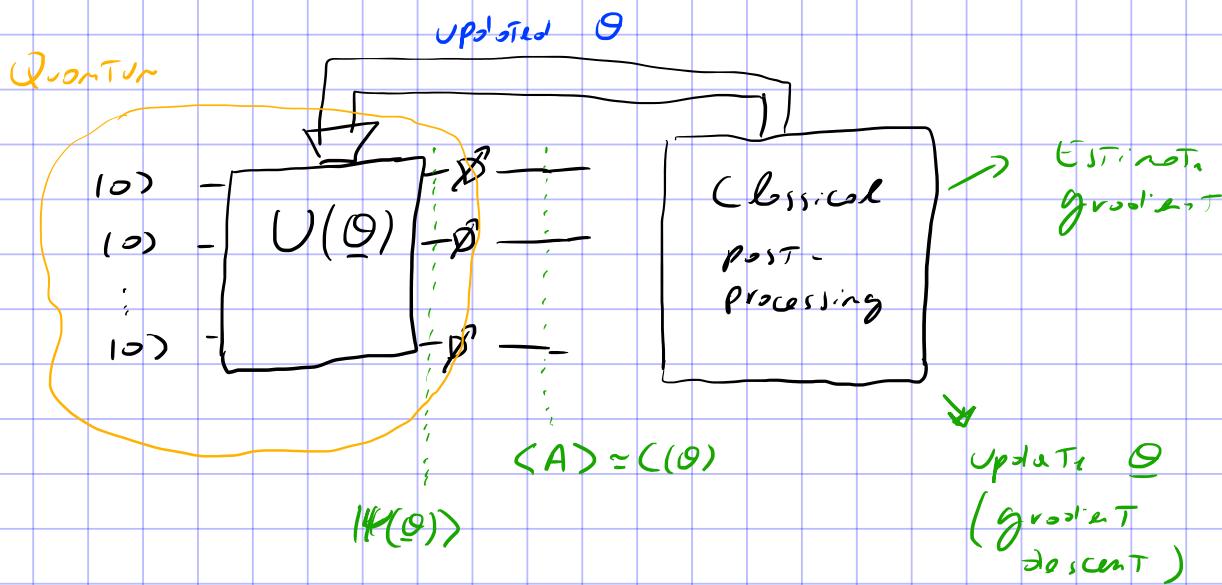
$$C(\underline{\theta}) = \langle \Psi(\underline{\theta}) | \hat{A} | \Psi(\underline{\theta}) \rangle$$

\hat{A}
Observable

Expectation value

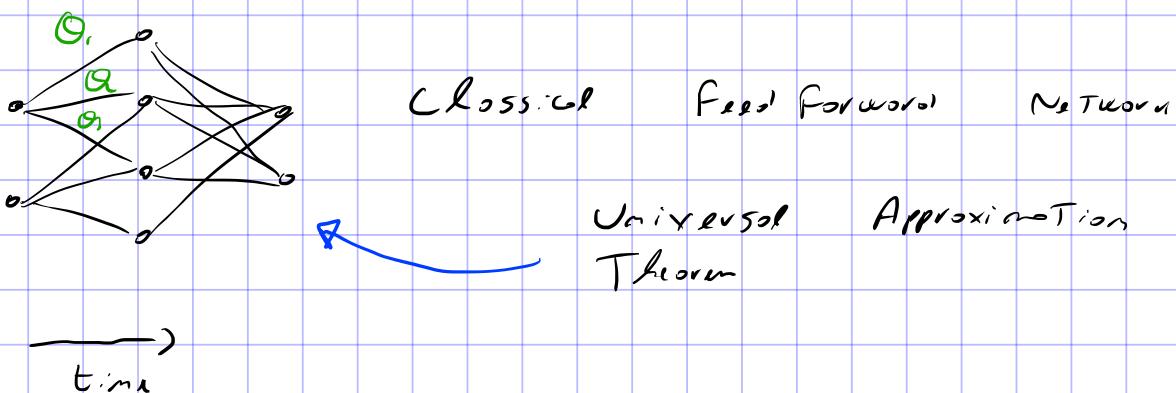
"Training" = minimize $C(\underline{\theta})$

Hybrid Algorithm

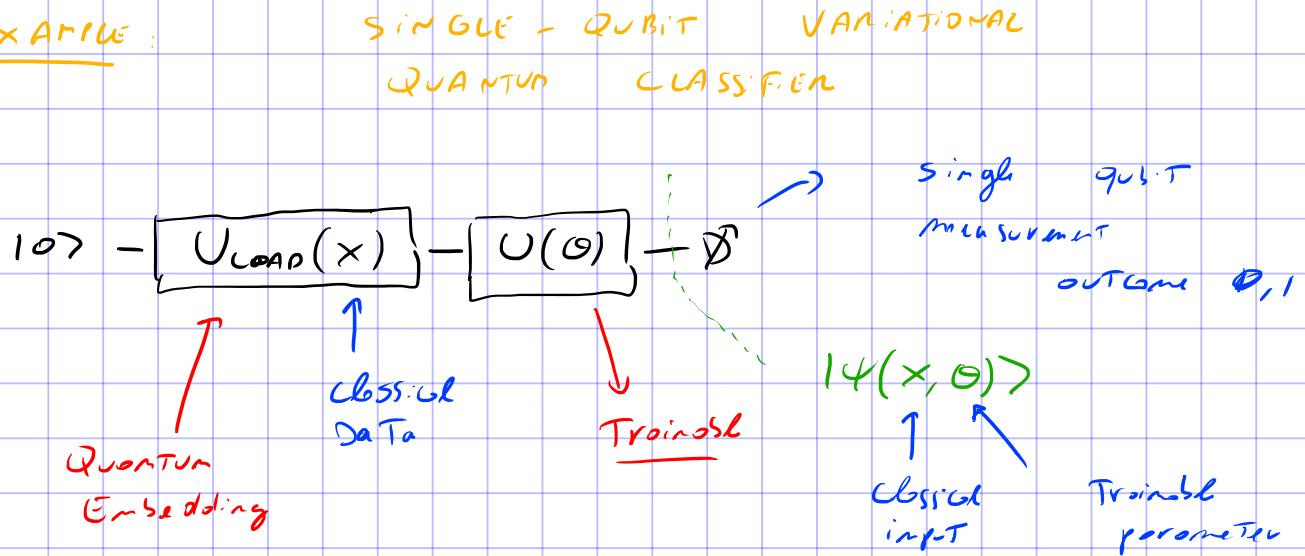


θ

Quantum Neural Networks \equiv Parametric Quantum Circuit



EXAMPLE:



OUTCOMES

$$0 \quad \text{probability} \quad |\langle 0 | \psi(x, \theta) \rangle|^2$$

$$1 \quad \text{probability} \quad |\langle 1 | \psi(x, \theta) \rangle|^2$$

Many layers

$$|0\rangle \xrightarrow{U_L(x)} \xrightarrow{U_1(\theta_1)} \xrightarrow{U_L(x)} \xrightarrow{U_2(\theta_2)} \dots \xrightarrow{\dots} \emptyset$$

L Times

ALTERNATIVE

$$|0\rangle \xrightarrow{U_1(x, \theta_1)} \xrightarrow{U_2(x, \theta_2)} \dots \emptyset$$

Non-linearity

3-vector

3 parameters

$$U_a(x, \theta) = U_a(\underline{\theta}_a + \underline{\omega}^{(a)} \cdot \underline{x})$$

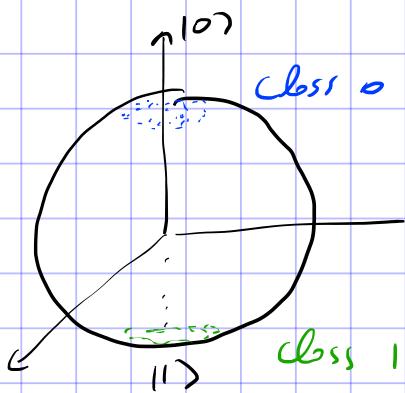
3-vector

d-vector
3x d matrix

Trainable

Ideally

a Trained model



binary classification
problems

To Train a model we use Fidelis

$$F(x, y) = | \langle y | \psi(x, \theta) \rangle |^2 \in [0, 1]$$

↑ ↑
 0,1 generic

$$\text{loss}(x, y) = 1 - F(x, y) \in [0, 1]$$

(infidelity)

Empirical Risk

$$R = \frac{1}{T} \sum_{i=1}^T \text{loss}(x_i, y_i)$$

//

Training Data

Training data

$$\{(x_i, y_i)\}_{i=1, \dots, T}$$

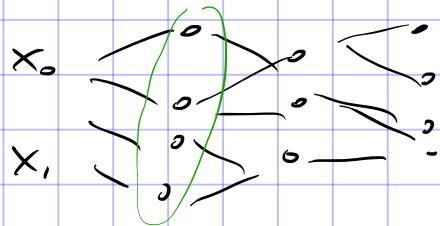
↑ ↑
 input class

Training \equiv Empirical Risk Minimization

$$\underset{\theta}{\operatorname{argmin}} (R)$$

↳ Finding the θ that minimize R

"Data Re-uploading" idea



Non linear activation function

$$10) -|\underline{U_c(x)}| - |\underline{U(\varrho)}| - \delta$$

↑
Non linear

$$10) -\boxed{U} - \boxed{U}$$

Data re-uploading for a universal quantum classifier

Adrián Pérez-Salinas, Alba Cervera-Lierta, Elies Gil-Fuster, José I. Latorre

A single qubit provides sufficient computational capabilities to construct a universal quantum classifier when assisted with a classical subroutine. This fact may be surprising since a single qubit only offers a simple superposition of two states and single-qubit gates only make a rotation in the Bloch sphere. The key ingredient to circumvent these limitations is to allow for multiple data re-uploading. A quantum circuit can then be organized as a series of data re-uploading and single-qubit processing units. Furthermore, both data re-uploading and measurements can accommodate multiple dimensions in the input and several categories in the output, to conform to a universal quantum classifier. The extension of this idea to several qubits enhances the efficiency of the strategy as entanglement expands the superpositions carried along with the classification. Extensive benchmarking on different examples of the single- and multi-qubit quantum classifier validates its ability to describe and classify complex data.

$$U(\phi) = e^{i\phi \cdot \hat{\sigma}}$$

$$\begin{aligned} U(\phi) U(\phi_1) &= e^{i\phi \cdot \hat{\sigma}} e^{i\phi_1 \cdot \hat{\sigma}} \\ &= e^{i\theta(\phi, \phi_1) \cdot \hat{\sigma}} \end{aligned}$$

$$\begin{aligned} &e^{i\mathbf{x} \cdot \hat{\sigma}} e^{i\phi_1 \cdot \hat{\sigma}} e^{i\mathbf{x} \cdot \hat{\sigma}} e^{i\phi_2 \cdot \hat{\sigma}} \\ &= e^{i\mathcal{J}(\mathbf{x}, \phi_1, \phi_2) \cdot \hat{\sigma}} \end{aligned}$$

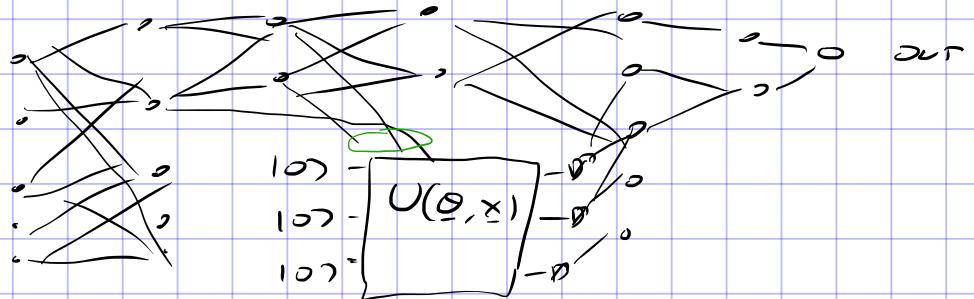
Non linear

TRAINING QUANTUM NEURAL NETWORKS

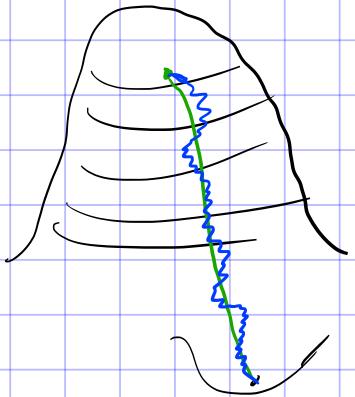
VIA GRADIENT DESCENT

Why gradient descent?

- 1) Deep Learning \rightarrow Gradient descent
- 2) Possible To integrate QNN with NN



- 3) Compatible with stochastic measurement outcomes
 \rightarrow Unbiased estimator of gradient



gradient descent

stochastic gradient
descent

- 4) Works very well in practice
 \rightarrow good generalization

- 5) scaling $O(N)$ $N \#$ Training data

Automatic

Differentiation

F.F. N.N.

$$x^{(l+1)} = f_l(\omega_l \cdot x^{(l)} + b^{(l)})$$

activation function

Trainable weight

$$x^{(0)} = x \quad \text{input}$$

$$x^{(L+1)} = y \quad \text{output}$$

$$C(\theta) = E_{x, y \sim \text{data}} \text{Loss}(x, y, \theta)$$

$$\frac{\partial C}{\partial \omega_{ij}^l} = \frac{\partial C}{\partial y} \frac{\partial y}{\partial \omega_{ij}^l} = \frac{\partial C}{\partial y} \frac{\partial x^{(l+1)}}{\partial x^{(l)}} \dots \frac{\partial x^{(L+1)}}{\partial \omega_{ij}^0}$$

$$\frac{\partial x^{(l+1)}}{\partial x^{(l)}} = f_l' \psi^l$$

Integrate QNN. provided we can compute gradient

Suppose

$$\partial_\theta(x) = \langle \psi(x, \theta) | \hat{A} | \psi(x, \theta) \rangle$$

Trainable parameters

observable / measurement
Trainable Embedding

Question

$$\frac{\partial J}{\partial \theta_i} ?$$

E.g. finite differences?

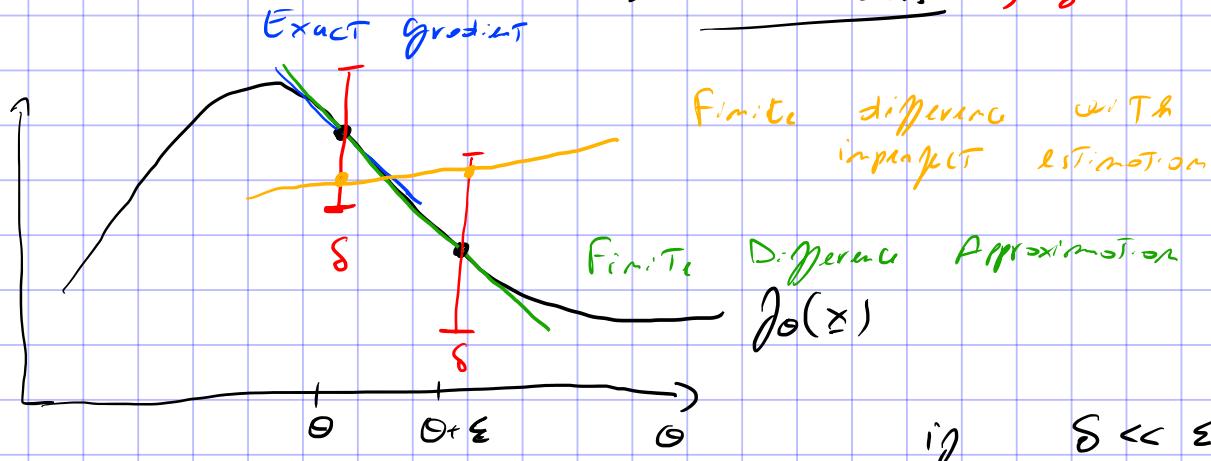
$$\frac{\partial \hat{J}_\theta(x)}{\partial \theta_i} = \frac{\hat{J}_\theta + \varepsilon e_i(x) - \hat{J}_\theta(x)}{\varepsilon} + O(\varepsilon^2)$$

$$e_i = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{pmatrix} \text{ i-th element}$$

Not Good

\hat{J}_θ is estimated from measurements

→ Error bars $\rightarrow \delta$



if $\delta \ll \varepsilon$

Then Finite difference is ok

$$\hat{A} = \sum_e (\alpha_e | \alpha_e \times \alpha_e |)$$

↓ ↓

Eigenvalues Eigen vectors

outcomes α_e with probability $|\langle \alpha_e | \psi(\theta, x) \rangle|^2$

$$\langle A \rangle = \langle \psi(\underline{\Omega}, \underline{x}) | \hat{A} | \psi(\underline{\Omega}, \underline{x}) \rangle$$

(\hookrightarrow loca / exact point)

$$= \frac{1}{M} \sum_{m=1}^M a_m$$

Yellow point

M # measurement sites

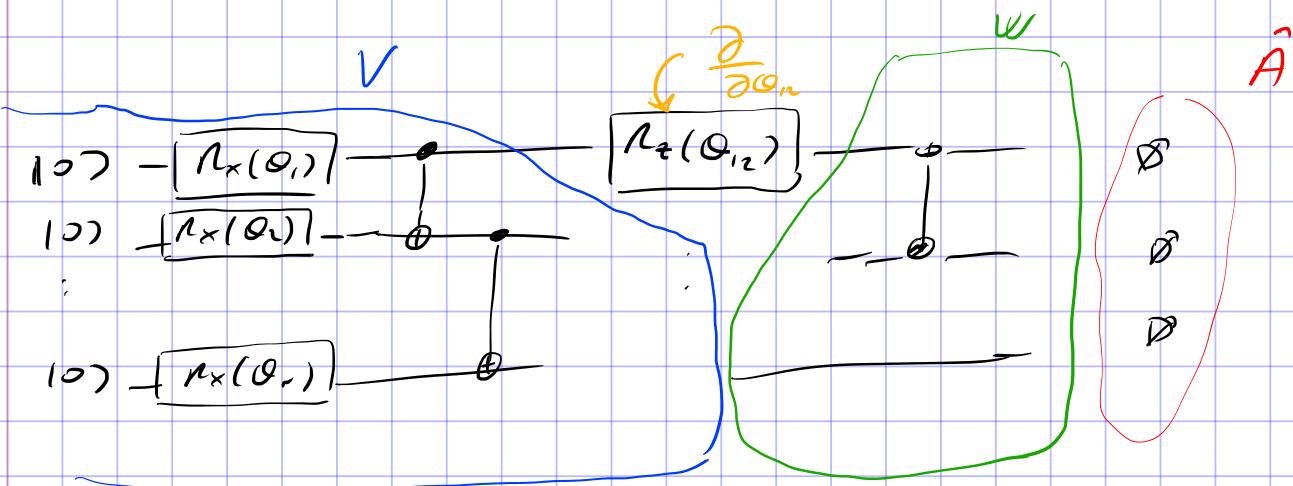
a_m outcome of the m th measurement

Variance $\frac{\langle A^2 \rangle - \langle A \rangle^2}{M} = \delta^2$

$$\delta \sim \frac{1}{\sqrt{M}} \ll \varepsilon$$

$$\Rightarrow M \gg \frac{1}{\varepsilon^2} \quad \varepsilon \sim 10^{-5} \quad M \sim 10^{10}$$

EXACT CALCULATION OF GRADIENT



$$\langle \psi(\underline{\Omega}) \rangle = W e^{i \Omega_n \hat{t}} V |0 \dots 0\rangle$$

$$j_0 = \langle 0 \dots 0 | V^+ e^{-i \Omega_n \hat{t}} W^+ A W e^{i \Omega_n \hat{t}} V |0 \dots 0 \rangle$$

$$\hat{e} = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix} \quad \hat{e}^2 = \mathbb{1}$$

$\hat{\sigma}$ so that $\hat{\sigma}^2 = 1$

$$e^{i\theta \hat{\sigma}} = \sum_m \frac{(i\theta \hat{\sigma})^m}{m!} = \sum_{m \text{ even}} \frac{(i\theta)^m}{m!} \mathbb{1} + \sum_{m \text{ odd}} \frac{(i\theta)^m}{m!} \hat{\sigma}$$

$$= \cos \theta \mathbb{1} + i \sin \theta \hat{\sigma}$$

im general

$$\hat{J}_n(x) = a(x) \cos^2 \theta_n + b(x) \sin^2 \theta_n + c(x) \cos(\theta_n) \sin(\theta_n)$$

$$= \frac{1 + \cos(2\theta_n)}{2} + \frac{1 - \cos(2\theta_n)}{2} + \frac{\sin(2\theta_n)}{2}$$

$$= \frac{(a+b)}{2} + \frac{(a-b)}{2} \cos(2\theta_n) + \frac{1}{2} \sin(2\theta_n)$$

$$\frac{\partial J}{\partial \theta_n} = -\frac{(a-b)}{x} \sin(2\theta_n) \cancel{x} + \frac{1}{x} \cos(2\theta_n) \cancel{x}$$

$\left(\begin{array}{l} \sin(a) = \frac{-\cos(a+\alpha) - \cos(a-\alpha)}{2 \sin \alpha} \\ \cos(a) = \dots \end{array} \right)$

$$\frac{\partial J}{\partial \theta_n} = \frac{J_{\theta+u} - J_{\theta-u}}{\sin(u)} \quad \text{exact}$$

\cup not small

In practice we measure \hat{A}

with different circuits $\theta \rightarrow \theta \pm u$

$a_{\theta_n} \approx$

$$\nabla \hat{\theta} = \mathbb{E} \left(\frac{\alpha_{\theta}^+ - \alpha_{\theta}^-}{\sin \psi} \right)$$

Empirical
Average over
Measurement
Shots

Unbiased estimator of
The gradient

$$\text{Variance} \sim (\sin \psi)^{-2} \quad \text{small}$$

$$\Rightarrow \sin \psi \quad \text{large}$$

$$\psi = \frac{\pi}{2}$$

SUMMARY (PARAMETER SHIFT RULE)

$$\frac{\partial \hat{\theta}_{\underline{\theta}}(\underline{x})}{\partial \theta_j} = \hat{\theta}_{\underline{\theta}} + \frac{\pi}{2} \ell_j - \hat{\theta}_{\underline{\theta}} - \frac{\pi}{2} \ell_j$$

Example

single output binary classif.

$$C(\theta) = \frac{1}{T} \sum_{x, y \in \text{Training}} \text{loss}(x, y, \theta)$$

$$= \frac{1}{T} \sum_{x, y} \left(1 - |\langle \gamma | \varphi(x, \theta) \rangle|^2 \right)$$

$$\langle \varphi(x, \theta) | \gamma \rangle \langle \gamma | \varphi(x, \theta) \rangle$$

$$\gamma = 0 \quad \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}$$

$$\gamma = 1 \quad \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix}$$

$$\hat{z} = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$$

$$\hat{1} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

$$\hat{A} = \frac{\hat{1} + \hat{z}(-1)}{2}$$

$$\begin{array}{c} \nearrow \gamma=0 \\ \searrow \gamma=-1 \end{array} \quad \begin{array}{l} \left(\begin{array}{cc} 1 & 0 \\ 0 & 0 \end{array} \right) \\ \left(\begin{array}{cc} 0 & 0 \\ 0 & 1 \end{array} \right) \end{array}$$

$$\partial_{\Theta}(x) = \langle \varphi(x, \Theta) | \hat{z} | \varphi(x, \Theta) \rangle$$

$$C(\Theta) = \frac{1}{T} \sum_{x \in T} \left(1 - \frac{1 + (-1)^{\partial_{\Theta}(x)}}{2} \right)$$

$$\frac{\partial C}{\partial \Theta_j} = \frac{1}{T} \sum_{x \in T} -(-1)^{\partial_{\Theta + \frac{\pi}{2} \varepsilon_j}(x)} \left(\partial_{\Theta + \frac{\pi}{2} \varepsilon_j}(x) - \partial_{\Theta - \frac{\pi}{2} \varepsilon_j}(x) \right)$$

$$= C(\Theta + \frac{\pi}{2} \varepsilon_j) - C(\Theta - \frac{\pi}{2} \varepsilon_j) \quad \text{Exact}$$

Training

via

gradient descent

descent

iteration

iteration

$$\underline{\Theta}_{n+1} = \underline{\Theta}_n - \gamma \vec{\nabla} C$$

\uparrow

Learning Rate

BARREN

PLATEAU

$$|\psi(x, \theta)\rangle = \hat{U}(x) e^{i\theta_j \hat{H}_j} |\psi(x)\rangle \quad \text{for } \dots \theta_j$$

\downarrow

$\hat{H}_j^2 = \mathbb{1}$

independent
on θ_j

$$\partial \phi(x) = \langle \psi(x, \theta) | \hat{A} | \psi(x, \theta) \rangle$$

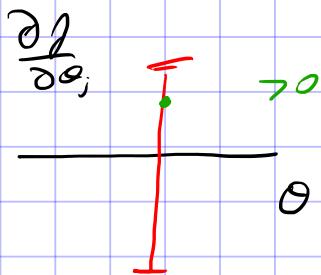
$$\frac{\partial \phi}{\partial \theta_j} = - (a_j - b_j) \sin(2\theta_j) + c_j \cos(2\theta_j)$$

Suppose θ_j initialized randomly

uniform $[-\pi, \pi]$

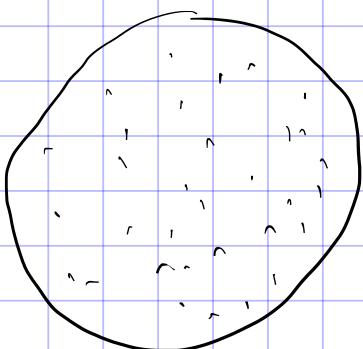
$$\mathbb{E}_{\theta_j} \left(\frac{\partial \phi}{\partial \theta_j} \right) = \frac{1}{2\pi} \int_{-\pi}^{\pi} (- (a_j - b_j) \sin(2\theta_j) + c_j \cos(2\theta_j)) d\theta_j$$

$$= 0$$



with θ_j
Variance $\sim (a-b)^2 + c^2$

IF THE CIRCUIT IS Highly expressive



our circuit can explore
a large part of space
of unitaries

$$\mathbb{E} \left(\frac{\partial J_{\theta}}{\partial \theta_j} \right) = 0$$

$$\text{Var} \left(\frac{\partial J_{\theta}}{\partial \theta_j} \right) = O \left(\text{poly} \left(\frac{1}{\epsilon^m} \right) \right)$$

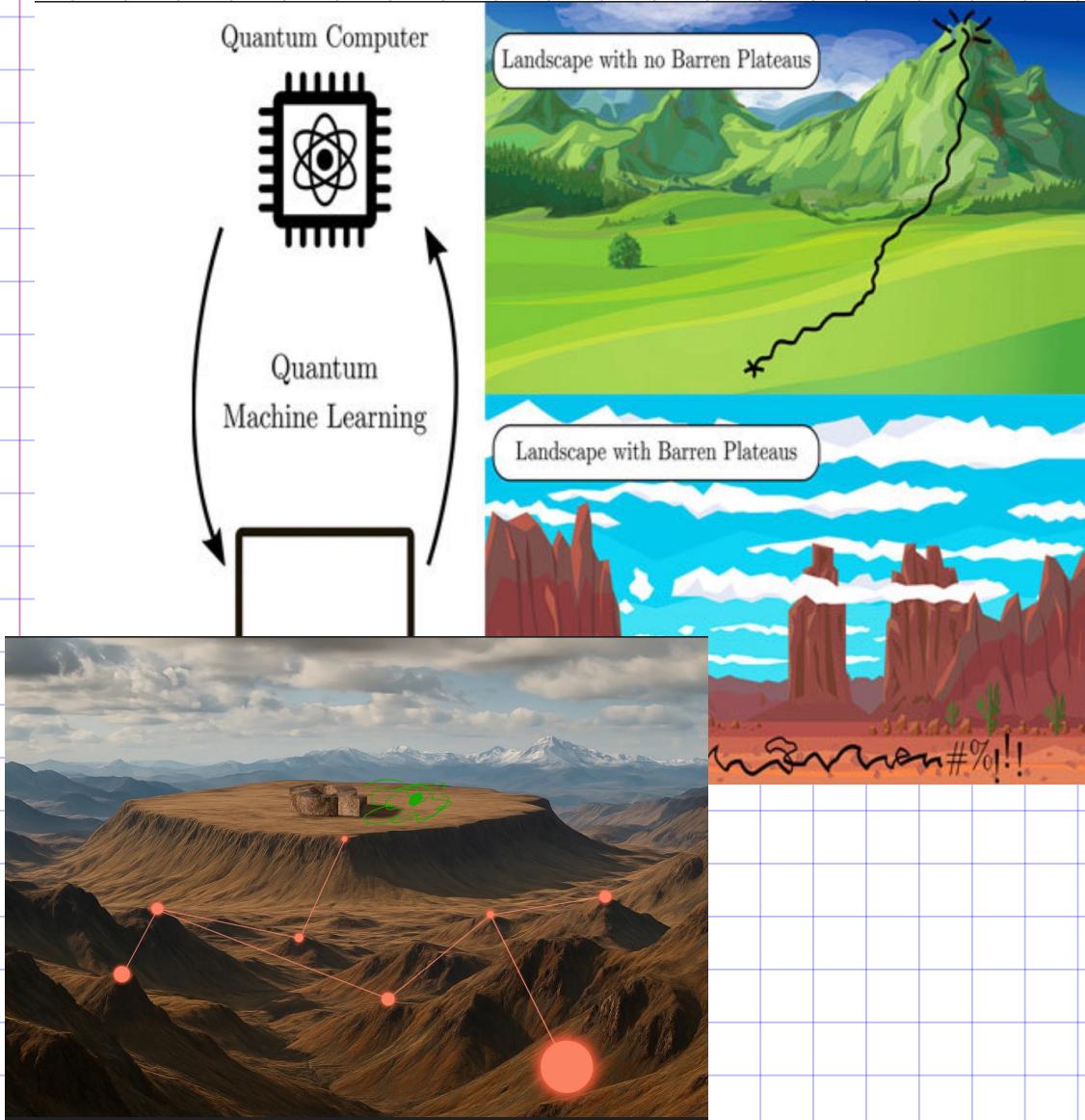
$m \neq \text{const}$

$$\mathbb{E}_{\substack{\text{all other} \\ \text{parameters}}} \left((\alpha - s)^2 + c^2 \right)$$

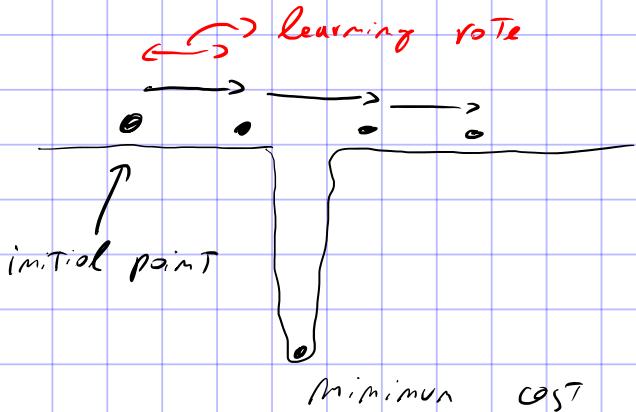
* Concentration of measur.

Each sample with high prob. is close to the mean, but the mean is zero

The interval greatest can be very small



Narrow Gorges



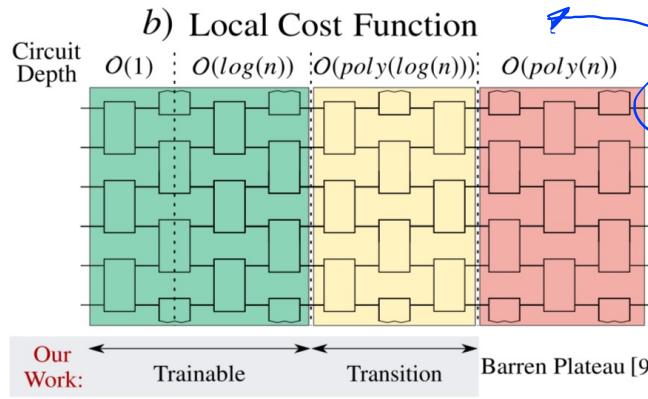
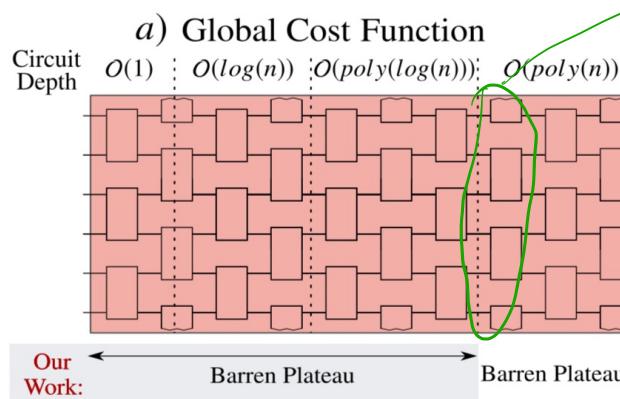
Article | [Open access](#) | Published: 19 March 2021

Cost function dependent barren plateaus in shallow parametrized quantum circuits

M. Cerezo Akira Sone, Tyler Volkoff, Lukasz Cincio & Patrick J. Coles

[Nature Communications](#) 12, Article number: 1791 (2021) | [Cite this article](#)

33k Accesses | 755 Citations | 85 Altmetric | [Metrics](#)



No D.P.
No gorges

Code PLTewrionb

Global GST function

$$\hat{H} = -100 \dots 0 \otimes 000 \dots 01$$

↑
zero energy for
all states except
 $|0 \dots 0\rangle$

Energy -1 for
state $|0 \dots 0\rangle$

Local cost function
WITH THE sum
of min of global
cost function

$$H = -\frac{1}{n} \sum_j \log \otimes \mathbb{1}_j$$

↑
identity
on every $\mathbb{1}_j$
but j th part

An initialization strategy for addressing barren plateaus in parametrized quantum circuits

Edward Grant¹, Leonard Wossnig¹, Mateusz Ostaszewski²,
and Marcello Benedetti³

¹Rahko Limited & Department of Computer Science, University College London

²Institute of Theoretical and Applied Informatics, Polish Academy of Sciences

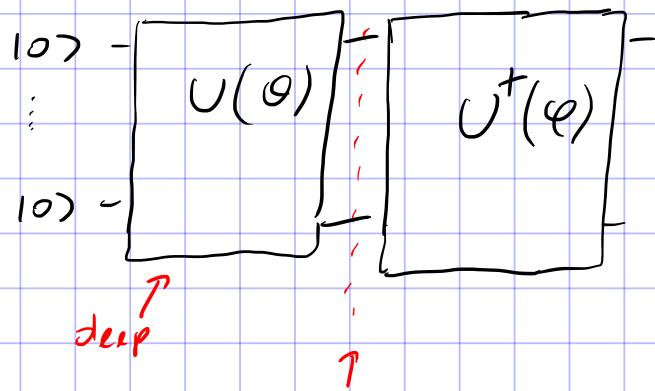
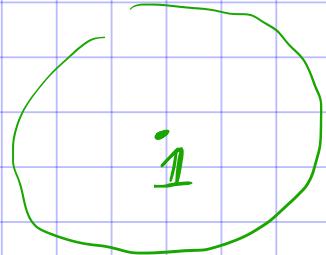
³Cambridge Quantum Computing Limited & Department of Computer Science, University College London

Published: 2019-12-09, volume 3, page 214

Eprint: arXiv:1903.05076v3

Doi: <https://doi.org/10.22331/q-2019-12-09-214>

Citation: Quantum 3, 214 (2019).



initialization

$$\Theta = \varphi$$

$$U(\Theta)U^\dagger(\varphi=\Theta)=\mathbb{I}$$

cancel out
→ No R.P.

After initialization
 $\varphi \neq \Theta$ during

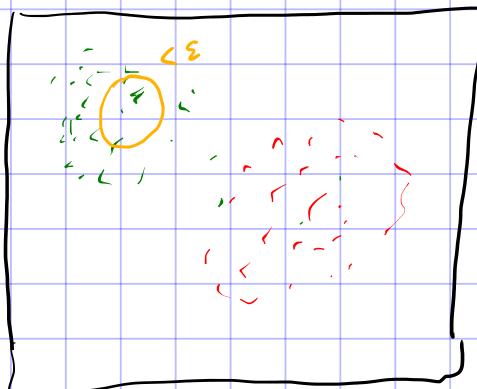
$$\varphi = \Theta$$

Training

Train independently

UNSUPERVISED LEARNING

- ## - Clustering



in T or due to
distance between
points

- ## Generative Models

(Example celebrity dataset)

data

$$x \sim P(x) \stackrel{\text{Sampled From}}{\curvearrowleft} \text{Unknown}$$

Tasq

Final a parametric approximation
of $P(x)$

$$P_{\Theta}(x)$$

Goal

be able To
generate new
samples

$$x \sim p_0(x)$$

Four samples should "look real"

Example "Celebrity Potatos"

Training

minimize θ distance $(P_\theta(x), P(x))$

Quantum Physics

[Submitted on 11 Apr 2018]

Differentiable Learning of Quantum Circuit Born Machine

Jin-Guo Liu, Lei Wang

Quantum circuit Born machines are generative models which represent the probability distribution of classical dataset as quantum pure states. Computational complexity considerations of the quantum sampling problem suggest that the quantum circuits exhibit stronger expressibility compared to classical neural networks. One can efficiently draw samples from the quantum circuits via projective measurements on qubits. However, similar to the leading implicit generative models in deep learning, such as the generative adversarial networks, the quantum circuits cannot provide the likelihood of the generated samples, which poses a challenge to the training. We devise an efficient gradient-based learning algorithm for the quantum circuit Born machine by minimizing the kernelized maximum mean discrepancy loss. We simulated generative modeling of the Bars-and-Stripes dataset and Gaussian mixture distributions using deep quantum circuits. Our experiments show the importance of circuit depth and gradient-based optimization algorithm. The proposed learning algorithm is runnable on near-term quantum device and can exhibit quantum advantages for generative modeling.

Comments: 9 pages, 7 figures, Github page for code [this URL](#)

Subjects: Quantum Physics (quant-ph); Machine Learning (cs.LG); Machine Learning (stat.ML)

Cite as: arXiv:1804.04168 [quant-ph]

(or arXiv:1804.04168v1 [quant-ph] for this version)

<https://doi.org/10.48550/arXiv.1804.04168> ⓘ

Journal reference: Phys. Rev. A 98, 062324 (2018)

Related DOI: <https://doi.org/10.1103/PhysRevA.98.062324> ⓘ

Quantum Physics

[Submitted on 23 Apr 2018 (v1), last revised 30 Apr 2018 (this version, v2)]

Quantum generative adversarial networks

Pierre-Luc Dallaire-Demers, Nathan Killoran

Quantum machine learning is expected to be one of the first potential general-purpose applications of near-term quantum devices. A major recent breakthrough in classical machine learning is the notion of generative adversarial training, where the gradients of a discriminator model are used to train a separate generative model. In this work and a companion paper, we extend adversarial training to the quantum domain and show how to construct generative adversarial networks using quantum circuits. Furthermore, we also show how to compute gradients -- a key element in generative adversarial network training -- using another quantum circuit. We give an example of a simple practical circuit ansatz to parametrize quantum machine learning models and perform a simple numerical experiment to demonstrate that quantum generative adversarial networks can be trained successfully.

Comments: 10 pages, 8 figures

Subjects: Quantum Physics (quant-ph); Machine Learning (cs.LG)

Cite as: arXiv:1804.08641 [quant-ph]

(or arXiv:1804.08641v2 [quant-ph] for this version)

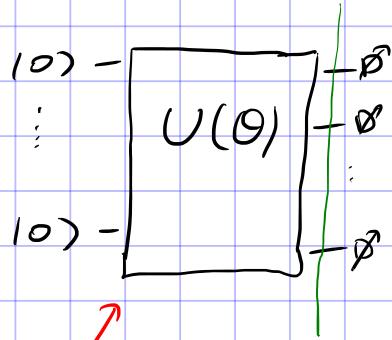
<https://doi.org/10.48550/arXiv.1804.08641> ⓘ

Journal reference: Phys. Rev. A 98, 012324 (2018)

Related DOI: <https://doi.org/10.1103/PhysRevA.98.012324> ⓘ

Main idea

sample data from measurement
of a quantum state

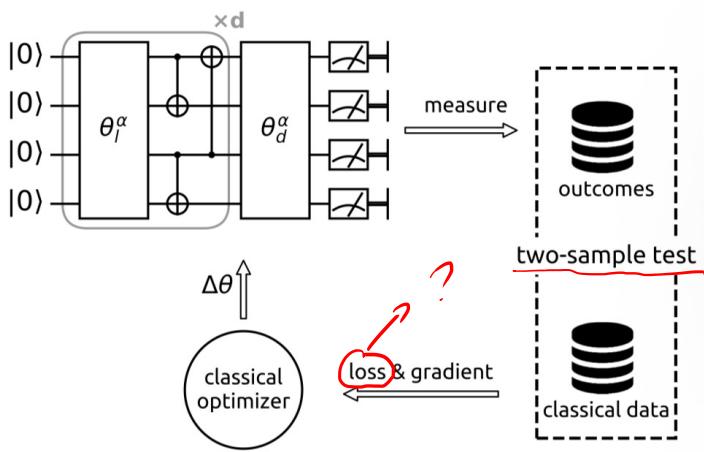


→ measured outcome → X

probability

$$P_\theta(x) = |\langle x | \psi(\theta) \rangle|^2$$

Quantum Circuit Born Machine



$\overbrace{\text{Maximum Mean Discrepancy}}^{\text{MMD}}$ $\overbrace{\text{Embedding Function}}^{\text{E}}$

$$\text{Loss}(\theta) = C(\theta) = \left\| \mathbb{E}_{x \sim P_\theta} [\phi(x)] - \mathbb{E}_{x \sim P} [\phi(x)] \right\|_2^2$$

\uparrow
 parametric distribution
 \uparrow
 true distribution

$$= \mathbb{E}_{x, y \sim P_\theta} [\phi(x) \cdot \phi(y)] + \mathbb{E}_{x, y \sim P} [\phi(x) \cdot \phi(y)] - 2 \mathbb{E}_{\substack{x \sim P_\theta \\ y \sim P}} [\phi(x) \cdot \phi(y)]$$

\uparrow
 inner product

Kernel Function

$$K(x, y) = \phi(x) \cdot \phi(y)$$

Choosing $\phi \Leftrightarrow$ choosing K

good choice Gaussian Kernel

$$K(x, y) = \exp\left(-\frac{\|x-y\|^2}{2\sigma^2}\right)$$

Thm

MMO + Gaussian Kernel

$$C = 0 \quad \text{iff} \quad P(x) = P_\Theta(x)$$

Training

$$P_\Theta(x) = |\langle x | \Phi(\Theta) \rangle|^2$$

$$= \langle \Phi(\Theta) | \underbrace{\langle x | x \rangle}_{\hat{A}} | \Phi(\Theta) \rangle$$

parameter shift rule

$$\frac{\partial P_\Theta(x)}{\partial \Theta_j} = P_\Theta + \frac{\pi}{2} \xi_j(x) - P_\Theta - \frac{\pi}{2} \xi_j(x)$$

↙ orbitrary function ↘

$$\frac{\partial}{\partial \Theta_j} \mathbb{E}_{x \sim P_\Theta} [\delta(x)] = \frac{\partial}{\partial \Theta_j} \sum_x P_\Theta(x) \delta(x)$$

$$= \sum_x \delta(x) (P_{\Theta + \frac{\pi}{2} \xi_j}(x) - P_{\Theta - \frac{\pi}{2} \xi_j}(x))$$

$$= \mathbb{E}_{x \sim P_{\Theta + \frac{\pi}{2} \xi_j}} [\delta(x)] - \mathbb{E}_{x \sim P_{\Theta - \frac{\pi}{2} \xi_j}} [\delta(x)]$$

$$\underline{\Theta}_j^\pm = \underline{\Theta} \pm \frac{\pi}{2} \xi_j$$

$$\frac{\partial C}{\partial \Theta_j} = 2 \left[\mathbb{E}_{\substack{x \sim P_{\Theta_j^+} \\ y \sim P_\Theta}} [u(x, y)] - \mathbb{E}_{\substack{x \sim P_{\Theta_j^-} \\ y \sim P_\Theta}} [u(x, y)] \right]$$

Four

$$+ \mathbb{E}_{\substack{x \sim P_{\Theta_j^+} \\ y \sim P}} [u(x, y)] - \mathbb{E}_{\substack{x \sim P_{\Theta_j^-} \\ y \sim P}} [u(x, y)] \right]$$

True distribution

→ Unknown

→ in practice we use
Training data

Q.SU.R.T

M_j

$e^{i\frac{\theta}{2}\hat{\sigma}}$

$e^{i\theta \hat{\sigma}}$

Algorithm

- Sample S elements "x" from P_0, P_1
using S "shots"
Coll outcomes x_s^0, x_s^1
 $s = 1 \dots S$

- Take a batch of Training Data

x_s^b $s = 1 \dots S$

Training Data has to be bigger than S

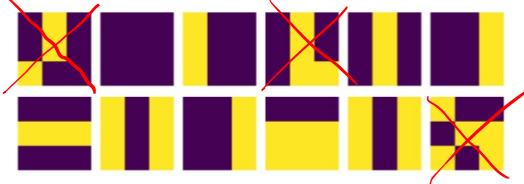
Empirical average

$$\frac{\partial C}{\partial \theta_j} \approx \frac{2}{S^2} \sum_{t,s=1}^S \left[U(x_s^t, x_t^0) - U(x_s^t, x_t^1) + U(x_s^t, \underline{x}_t^b) - U(x_s^t, \underline{x}_t^s) \right]$$

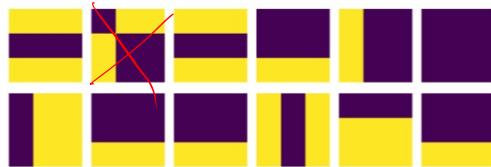
Training Data

everything else is for obta generated
by the quantum computer

$N = 2000$
 $\chi = 88.6\%$



$N = 20000$
 $\chi = 92.4\%$



$N = \infty$
 $\chi = 95.4\%$



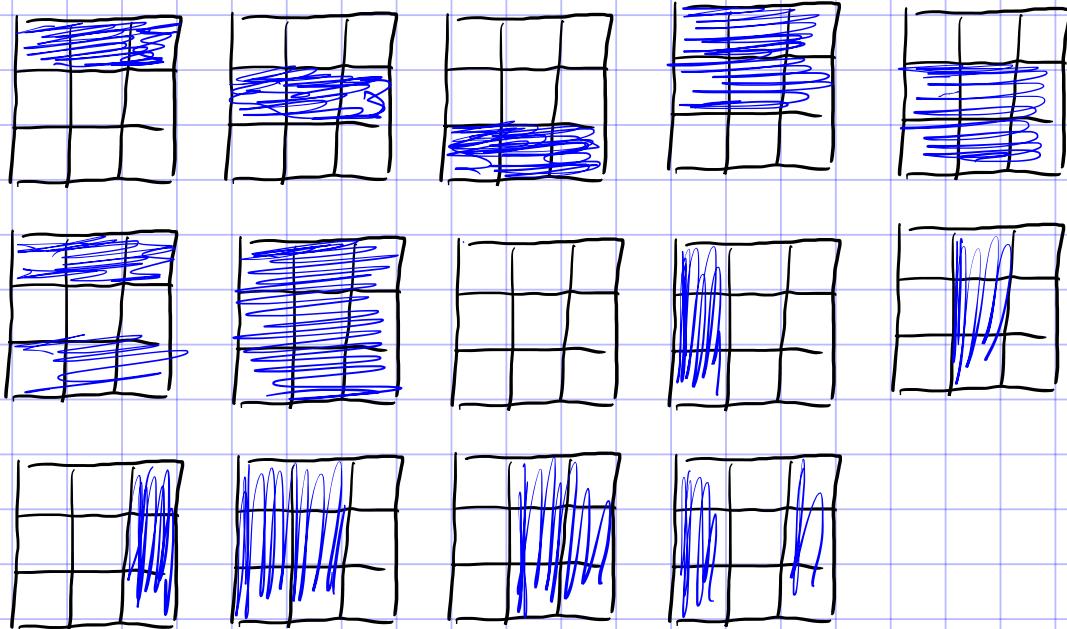
binary 7×3 images

X_{11}	X_{12}	X_{13}
X_{21}	X_{22}	X_{23}
X_{31}	X_{32}	X_{33}

$X_i \in \{0, 1\}$

rule

True images
DANS on stripes
or C,TInv



15

"True images"

2⁹

possible
configurations
~ 500

Quantum Generative Adversarial Networks

Classical

GANs

2 players

Generator G

Discriminator D

G wants to generate your data to fool D

D judges if data is fake or real

Turn based game

G, D iteratively update their strategy to win

Theorem (Nash Equilibrium)

- G produces good data
- D cannot tell if data is real vs fake

Quantum generative adversarial networks with Cirq + TensorFlow



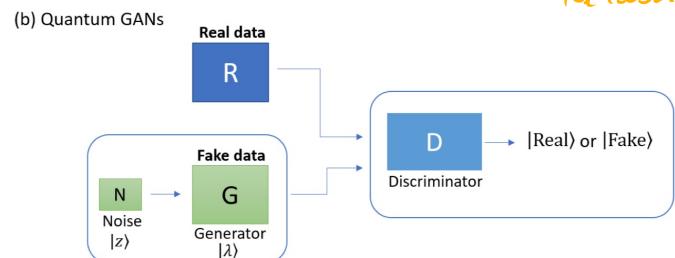
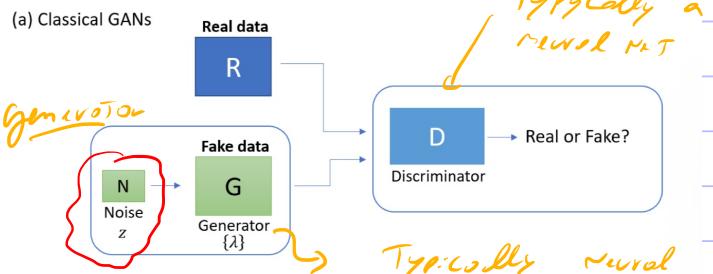
Classical GANs

$$t \sim N(0, 1)$$

Generator

$$x = g(t)$$

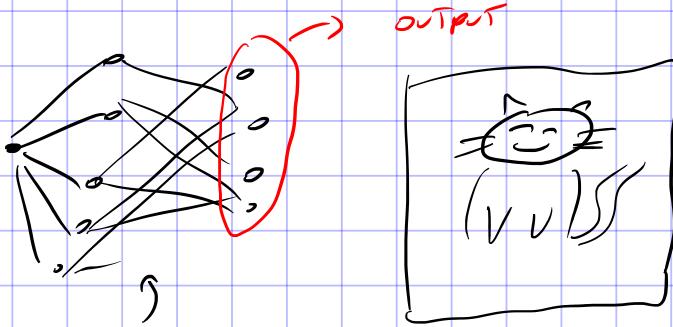
↳ generated fake image



Classical generator

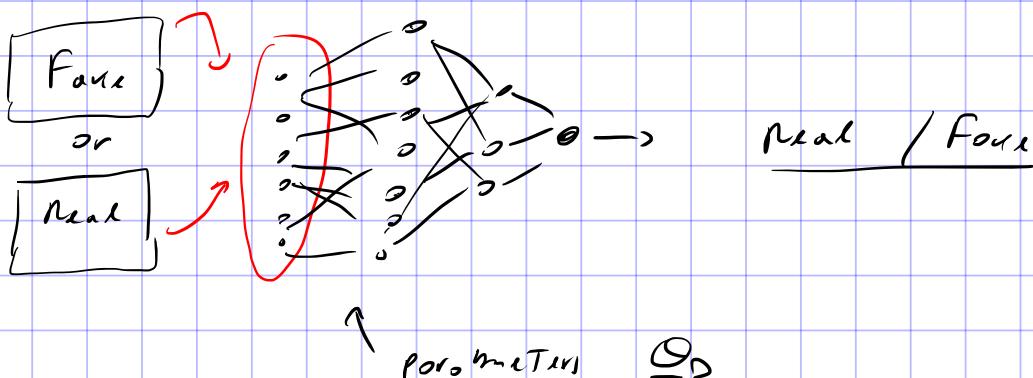
$$z \sim N(0,1)$$

parameters $\underline{\Theta}_G$



Classical Discriminator

input



$$V(\underline{\Theta}_G, \underline{\Theta}_D)$$

Score Function as
misclassification probability

goal of G

$$\text{minimize}_{\underline{\Theta}_G} V(\underline{\Theta}_G, \underline{\Theta}_D)$$

goal of D

$$\text{maximize}_{\underline{\Theta}_D} V(\underline{\Theta}_G, \underline{\Theta}_D)$$

$$\min_{\underline{\Theta}_G} \max_{\underline{\Theta}_D} V(\underline{\Theta}_G, \underline{\Theta}_D)$$

GANs

Quantum GANs

S_n

real

Quantum

states

S_F

Fake

Quantum

states

if data is classical

$$\{x_i\}_{i=1}^n \rightarrow |\psi\rangle$$

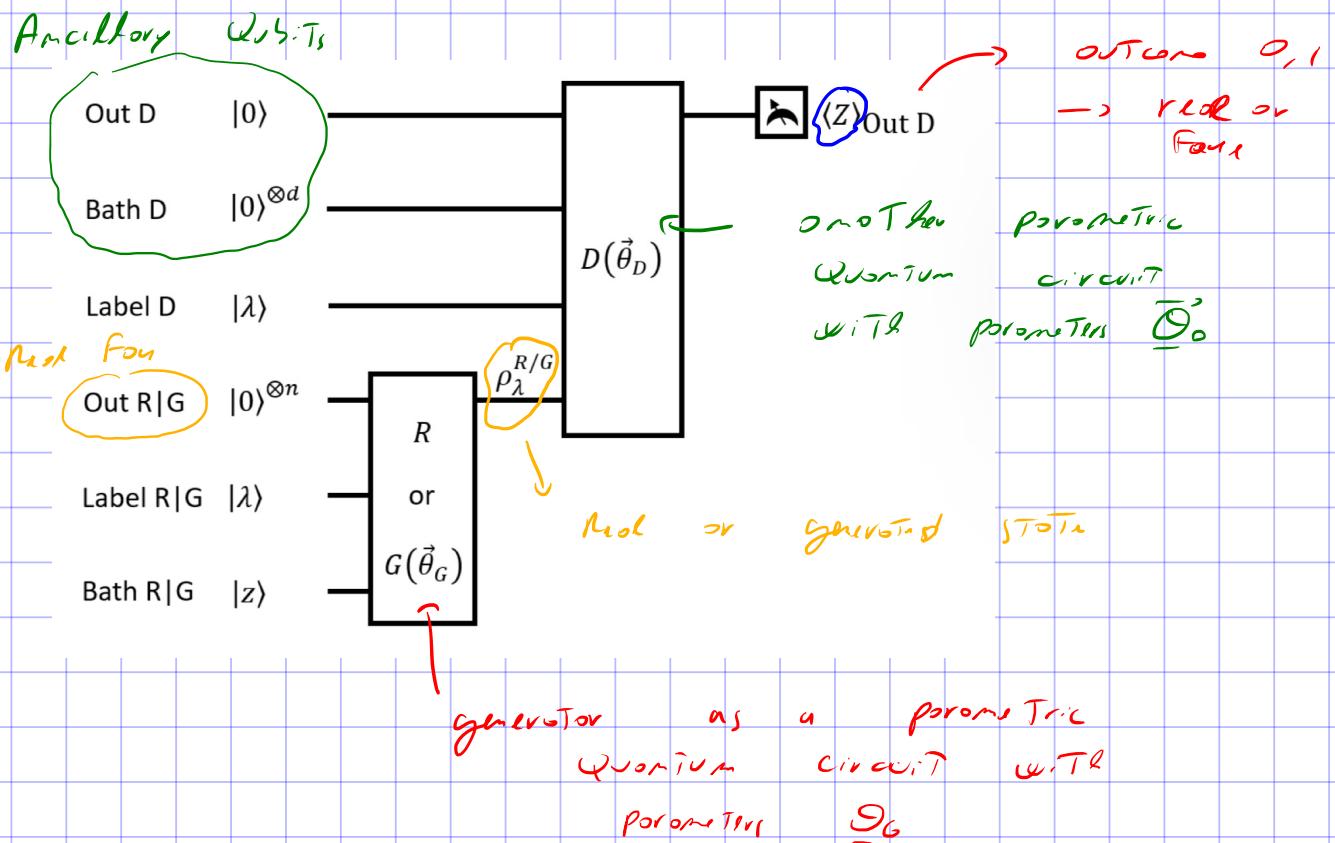
$$f = \frac{1}{n} \sum_{i=1}^n |\psi(x_i)\rangle \langle \psi(x_i)|$$

6 outputs from quantum states

$$G(|\lambda\rangle, z) = |\psi_z(\lambda)\rangle \langle \psi_z(\lambda)|$$

D gets either real or fake state

$$R \text{ real data } n(|\lambda\rangle) = |\psi(\lambda)\rangle$$



$$\hat{t} = 10 \times 01 - 11 \times 11$$

$0 \rightarrow \text{real}$
 $1 \rightarrow \text{fake}$

Why oscillatory quibits?

States discrimination

p_0, p_1 Two possible states

Most general measurement strategy

POVM (Position Operator valued Measurement)

Π_0, Π_1 Hermitian

$$\Pi_0, \Pi_1 \geq 0$$

$$\hat{\Pi}_0 + \hat{\Pi}_1 = \mathbb{1}$$

Theorem (No. mona)

$\text{Tr}(\Pi_i | \rho) \rightarrow$ probabilities of getting outcome i

$\text{D} \rightarrow \boxed{U} \rightarrow \emptyset$ ↗
 $\rho \rightarrow \boxed{U}$
 ↑
 circuit (decompose)

projective measurement with action σ_i

$$V(\underline{\Omega}_G, \underline{\Omega}_D) = \underset{\text{data}}{\mathbb{E}}$$

P (real is real and fake is fake)

P
 discriminator prediction

Measurement of two quantum circuit

→ gradients

Param. Tern
Sigmoid rule

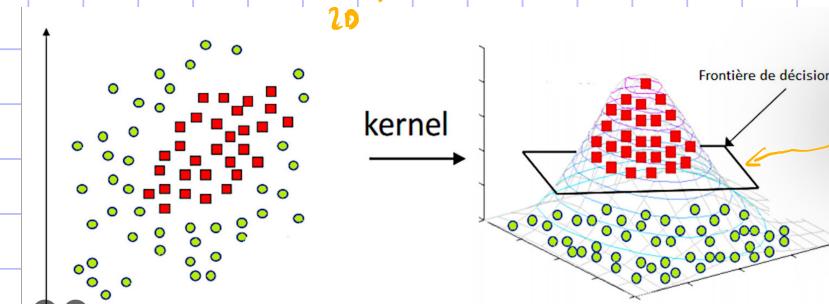
QUANTUM KERNEL METHODS

Linear Loss functions on a bigger space

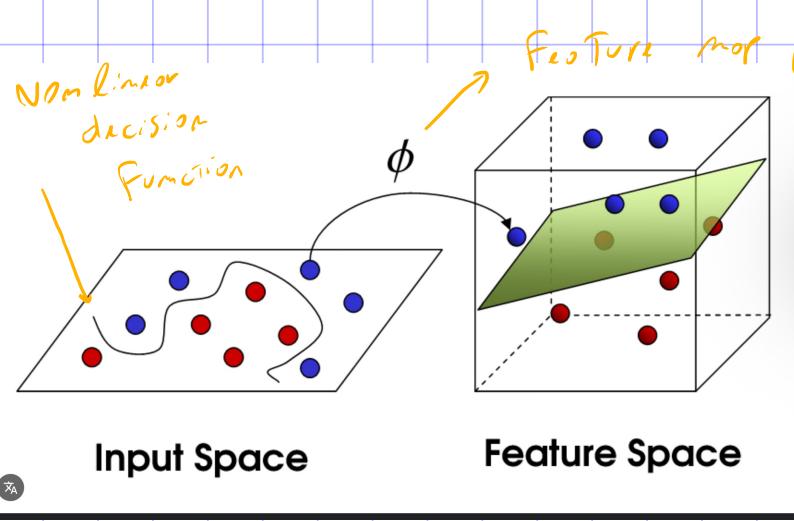
→ intuition → data are easier

To loss in
a larger space

No linear decision line



Linear decision exists in 3D plane



Linear classification
Hyperplane

blue points
above

red points
below

Non linear decision problems as linear
decision problems on "feature space"

feature map between input to "feature space"

$$\phi : x \rightarrow \phi(x)$$

x $\phi(x)$
input feature

Feature space
de J in a Hilbert space
(sum of Quantum states)

=, Quantum Feature map

$$\phi : x \rightarrow |\psi(x)\rangle$$

Quantum Physics

[Submitted on 26 Jan 2021 (v1), last revised 17 Apr 2021 (this version, v2)]

Supervised quantum machine learning models are kernel methods

Maria Schuld

With near-term quantum devices available and the race for fault-tolerant quantum computers in full swing, researchers became interested in the question of what happens if we replace a supervised machine learning model with a quantum circuit. While such "quantum models" are sometimes called "quantum neural networks", it has been repeatedly noted that their mathematical structure is actually much more closely related to kernel methods: they analyse data in high-dimensional Hilbert spaces to which we only have access through inner products revealed by measurements. This technical manuscript summarises and extends the idea of systematically rephrasing supervised quantum models as a kernel method. With this, a lot of near-term and fault-tolerant quantum models can be replaced by a general support vector machine whose kernel computes distances between data-encoding quantum states. Kernel-based training is then guaranteed to find better or equally good quantum models than variational circuit training. Overall, the kernel perspective of quantum machine learning tells us that the way that data is encoded into quantum states is the main ingredient that can potentially set quantum models apart from classical machine learning models.

Comments: 26 pages, 9 figures - Version 2 emphasises focus on supervised learning, adds more references to existing literature, deletes section on state discrimination due to a technical error, and updates the comparison between kernel-based and variational training

Subjects: Quantum Physics (quant-ph); Machine Learning (stat.ML)

Cite as: arXiv:2101.11020 [quant-ph]

(or arXiv:2101.11020v2 [quant-ph] for this version)

<https://doi.org/10.48550/arXiv.2101.11020>

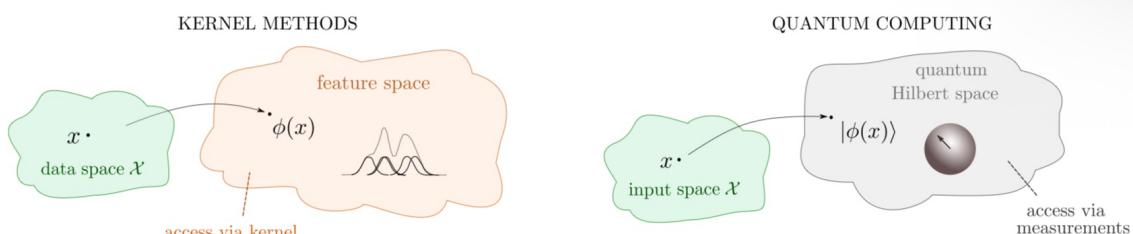


FIG. 1. **Quantum computing and kernel methods are based on a similar principle.** Both have mathematical frameworks in which information is mapped into and then processed in high-dimensional spaces to which we have only limited access. In kernel methods, the access to the feature space is facilitated through *kernels* or inner products of feature vectors. In quantum computing, access to the Hilbert space of quantum states is given by measurements, which can also be expressed by inner products of quantum states.

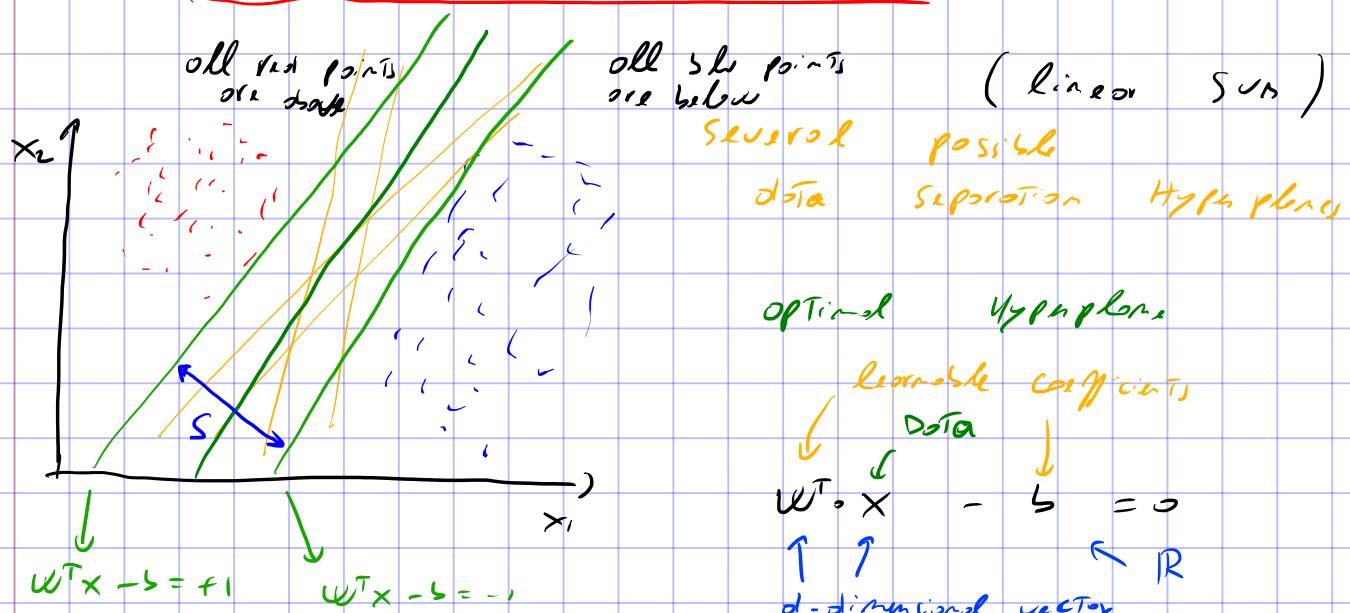
Classically \rightarrow Kernel Trick

$$k(x, y) = \phi(x) \cdot \phi(y)$$

inner product
between two
feature maps

Choosing Feature Map \leftrightarrow Choosing Kernel

SUPPORT VECTOR MACHINES



parallel hyperplanes

separation

$$S = \frac{2}{\|w\|}$$

2D $w_1 x_1 + w_2 x_2 - b = 0$

line ↗

Select the hyperplanes with maximum "margin"

$$\frac{2}{\|w\|}$$

We need to solve

$$w^T x_i - b \geq 1$$

red points

$$w^T x_i - b \leq -1$$

blue points

$$y_i = 1$$

red points

(binary classification)

$$y_i = -1$$

blue points

$$y_i (w^T x_i - b) \geq 1$$

compact form

$$1 - y_i (w^T x_i - b) \leq 0$$

(hard constraint)

Soft Constraint

Loss whenever TH
inequality is violated

Hinge loss

$$\ell(x_i, y_i) = \max \{0, 1 - y_i(\omega^T x_i - b)\}$$

$$\ell(x_i, y_i) = 0 \quad \text{if } \quad y_i(\omega^T x_i - b) \geq 1$$

$$\text{otherwise.} \quad \ell(x_i, y_i) \geq 0$$

Cost function

$$C(\omega, b) = \lambda \|\omega\| + \sum_{i=1}^n \ell(x_i, y_i)$$

↑

Hyper parameter

↑↑

Training pairs

Training / Learning

$$\min_{\omega, b} C(\omega, b)$$

↓

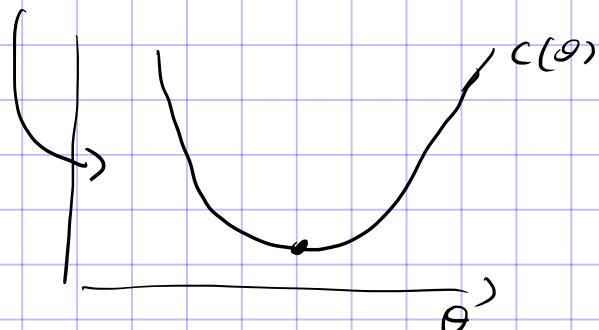
$$= \max \text{ margin } \frac{2}{\|\omega\|}$$

$$\hookrightarrow \min \|\omega\|$$

= minimizes The loss

λ Hyperparameter To model TH Tradeoff

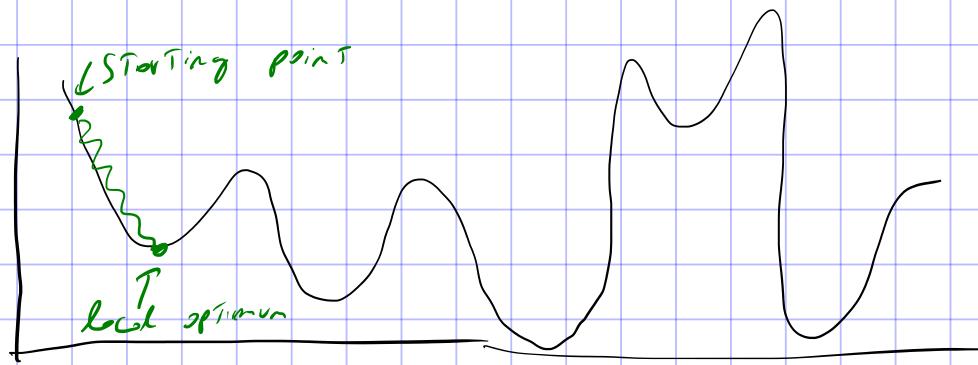
Convex Optimization problem



$$\Theta = (\omega, b)$$

Algorithm converges
To global optimum

in Neural Networks



Non Linear sum \equiv Linear sum in
feature space

$$x \rightarrow \phi(x) \quad \text{Feature map}$$

$$\gamma_i(\omega^\top x_i - b) \geq 1 \rightarrow \gamma_i(\omega^\top \phi(x_i) - b) \geq 1$$

$$\text{loss}(x_i, \gamma_i) = \max\{0, 1 - \gamma_i(\omega^\top \phi(x_i) - b)\}$$

<u>Problem</u>	$\phi(x)$	can live in a very large space (infinite dimensional)
----------------	-----------	---

$$\text{dimension}(\omega) = \text{dimension}(\phi(x))$$

Dual formulation

$$\max_{\alpha \geq 0} \sum_{i=1}^n \alpha_i - \sum_{i,j=1}^n \alpha_i \alpha_j K(x_i, x_j) - \frac{1}{2} \sum_i \alpha_i^2$$

much simpler

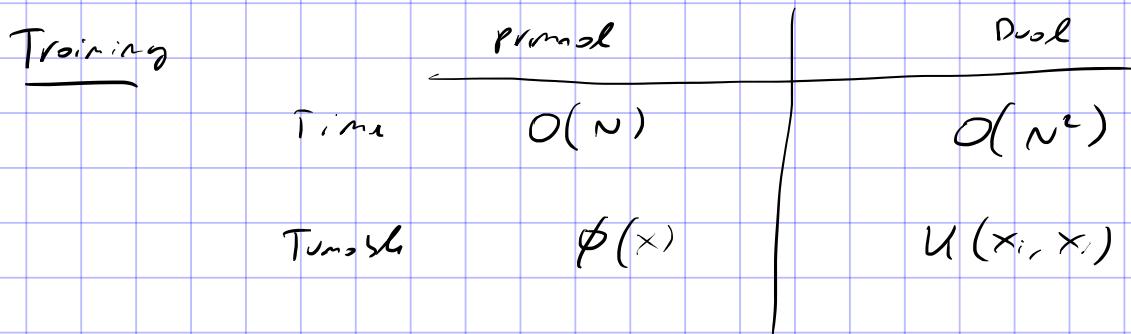
Primal formulation

$$\min_{\omega, b} \left[\lambda \|\omega\| + \sum_{i=1}^n l(x_i, \gamma_i) \right]$$

$d_i \quad i = 1 \dots n \quad \# \text{ Training Data}$

$w_i \quad i = 1 \dots d_f \quad d_f \quad \text{dimension of feature space}$
(can be huge)

Dual formulation is independent on the dimension of feature space



$$U(x, x') = \phi(x) \cdot \phi(x')$$

Predictions

Primal : check the side of the plane

decision Hyperplane $w^T \phi(x) - b = 0$

$$y_{predicted} = \text{sign}(w^T \phi(x) - b) = \pm 1$$

solutions to primal problem

Dual formulation

$$w_{optimal} = \sum_i \alpha_i y_i \phi(x_i)$$

$$y_{\text{predicted}} = \text{sign} \left(\sum_{i=1}^n \alpha_i y_i K(x_i, x) \right)$$

Pairs in The Training set

Learned during Training

New data point

"Kernel Expansion of The Data"

Both For Training / Testing No calculation
 in Feature space
 aside from
 computing The
 Kernel

General Result : Representer Theorem

$$\begin{aligned} f(x) &= \underset{\gamma}{\operatorname{argmin}} \left[\sum_i \text{loss}(x_i, \gamma_i) + \text{regularization} \right] \\ &= \sum_{i=1}^n y_i \alpha_i K(x_i, x) \end{aligned}$$

$\uparrow \lambda \|\gamma\|^2$

in class. col sum, no calculations in
 Feature space

→ choose Kernel

Example

$$K(x, x') = \exp \left(- \frac{\|x - x'\|^2}{2\sigma^2} \right)$$

QUANTUM SUM

Quantum computers can do offbeat calculations in large Hilbert spaces

Directly work with feature map

choice of feature map

$$1) \quad x \rightarrow |\psi(x)\rangle$$

$$2) \quad x \rightarrow |\psi(x)\rangle\langle\psi(x)|$$



more natural

$\phi(x)$

use ②

$$u(x, x') = \phi(x) \cdot \phi(x')$$

inner product between operators

Hilbert-Schmidt product

$$A \cdot B = \text{Tr}(A^+ B)$$

$$u(x, x') = \text{Tr} \left(|\psi(x)\rangle\langle\psi(x)| \mid\psi(x')\rangle\langle\psi(x')| \right)$$

$$= |\langle\psi(x)|\psi(x')\rangle|^2 \in \mathbb{R}$$

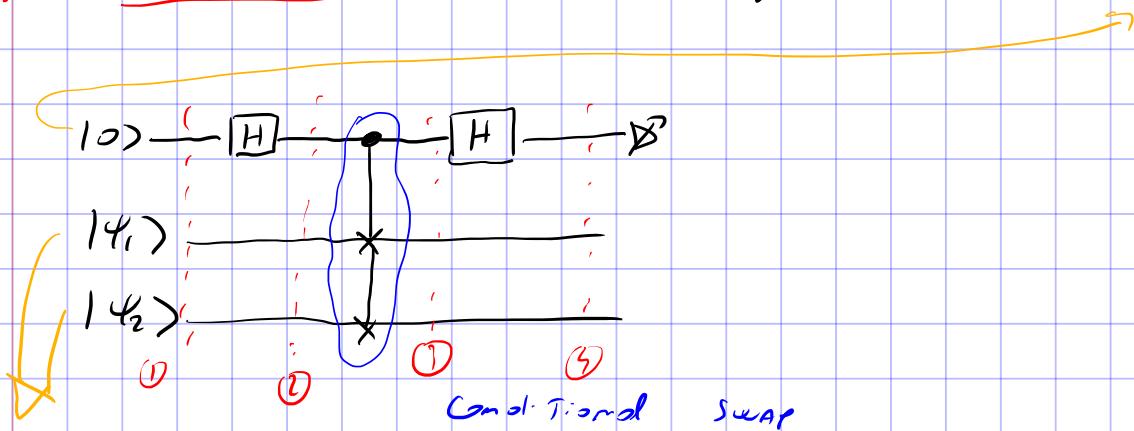
Representer Theorem

$$\begin{aligned}
 f_{\text{opt}}(x) &= \sum_i \alpha_i y_i \mathcal{U}(x_i, x) \\
 &= \sum_i \alpha_i y_i \text{Tr} [| \psi(x_i) \rangle \langle \psi(x_i) | \psi(x) \rangle \langle \psi(x) |] \\
 &= \langle \psi(x) | M | \psi(x) \rangle \\
 M &= \sum_i \alpha_i y_i | \psi(x_i) \rangle \langle \psi(x_i) | \quad (\text{measurable observable}) \\
 &\quad \uparrow \\
 &\quad \text{Optimized during Training}
 \end{aligned}$$

Quantum Circuits for Computing The Kernel

1) SWAP TEST

2 registers + 2 ancilla



Two registers

(possibly with many qubits)

$$① |0\rangle |1\rangle |1\rangle |1\rangle$$

$$② (|0\rangle |1\rangle |1\rangle |1\rangle + |1\rangle |1\rangle |1\rangle) / \sqrt{2}$$

$$③ (|0\rangle |1\rangle |1\rangle + |1\rangle |1\rangle |1\rangle) / \sqrt{2}$$

$$⑤ ((|0\rangle + |1\rangle) |1\rangle |1\rangle + (|0\rangle - |1\rangle) |1\rangle |1\rangle) / 2$$

$$= |0\rangle \left(\frac{|\psi_1\psi_2\rangle + |\psi_2\psi_1\rangle}{2} \right) + |1\rangle \left(\frac{|\psi_1\psi_2\rangle - |\psi_2\psi_1\rangle}{2} \right)$$

$$P_0 = \left\| \frac{|\psi_1\psi_2\rangle + |\psi_2\psi_1\rangle}{2} \right\|^2 = \frac{1}{2} (1 + |\langle\psi_1|\psi_2\rangle|^2 + |\langle\psi_2|\psi_1\rangle|^2)$$

$$P_1 = \left\| \frac{|\psi_1\psi_2\rangle - |\psi_2\psi_1\rangle}{2} \right\|^2 = \frac{1}{2} (1 - |\langle\psi_1|\psi_2\rangle|^2)$$

$$P_0 + P_1 = 1$$

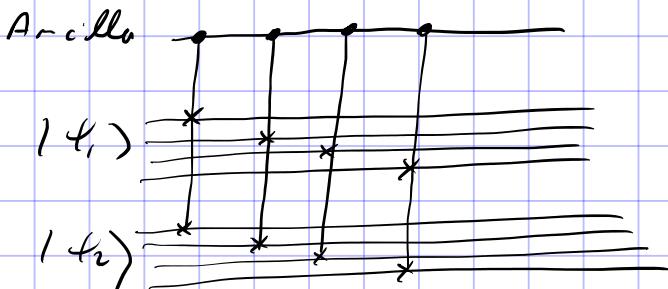
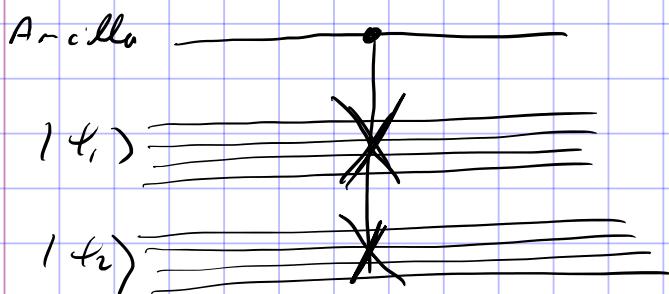
$$|\langle\psi_1|\psi_2\rangle| = P_0 - P_1 = 2P_0 - 1$$

For many Qubits: Tj

$$\text{C-Swap} = \prod_{i=1}^n \text{C-Swap}_{i, i+1}$$

↓

3 qubit FREQ DOMAIN A?



C-Swap_{i, i+1}

S

CNOTS

Single QUBIT

2)

SWAP MEASUREMENT

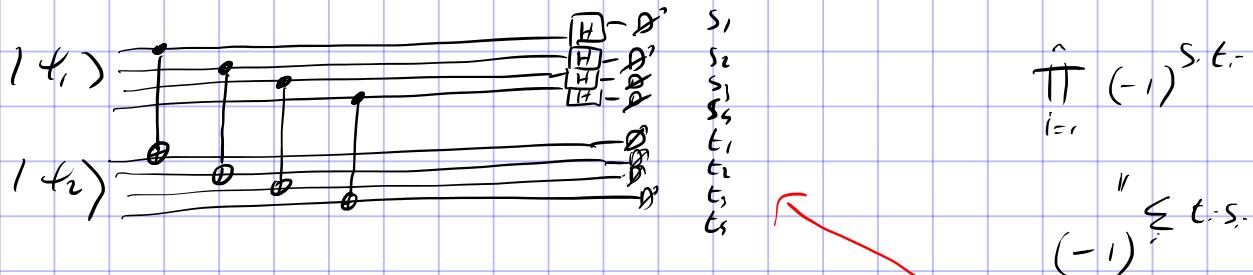
$$|\langle \psi_1 | \psi_2 \rangle|^2 = \langle \psi_1 | \psi_2 \rangle \langle \psi_2 | \psi_1 \rangle = \langle \psi_1 | \text{SWAP} | \psi_2 \rangle$$

$$\text{SWAP} = \prod_{i=1}^n \text{SWAP}_{i,i+1}$$

$$\text{SWAP}_{i,i+1} = U \times U^\dagger$$

$$U = (-1)^{\sum_{i < j} S_i \cdot S_j}$$

$$U = \text{CNOT} \cdot H_i$$



$$|\langle \psi_1 | \psi_2 \rangle|^2 = \mathbb{E}((-1)^{\sum_{i < j} S_i \cdot S_j})$$

↓
average over measurement shots

③

DIRECT FIDELITY ESTIMATION

$$|\psi_1\rangle = U(x_1) |0 \dots 0\rangle$$

$$|\psi_2\rangle = U(x_2) |0 \dots 0\rangle$$

$$U(x_1, x_2) = |\langle \psi_1 | \psi_2 \rangle|^2 = |\langle 0 \dots 0 | U(x_1)^\dagger U(x_2) |0 \dots 0 \rangle|^2$$

$$= \langle 0 \dots 0 | U(x_1)^\dagger U(x_2) |0 \dots 0 \rangle \times \langle 0 \dots 0 | U(x_2)^\dagger U(x_1) |0 \dots 0 \rangle$$

$$10) - \left[\begin{array}{c} U(x_1) \\ \vdots \\ U(x_n) \end{array} \right] \left[\begin{array}{c} U(x_1)^+ \\ \vdots \\ U(x_n)^+ \end{array} \right] = 0$$

Extract
probabil.

$P_0 \dots 0$

\downarrow
probability to get
all outcomes zero

Advantage: single register

Disadvantages: - Higher depth

- Prob. may be small

SUMMARY:

Quantum

Kernel Methods

Hybrid

Algorithms



Classical Kernel



Training

Kernels estimator

From Quantum
Computer

Article | Published: 12 July 2021

A rigorous and robust quantum speed-up in supervised machine learning

[Yunchao Liu](#), [Srinivasan Arunachalam](#) & [Kristan Temme](#)

[Nature Physics](#) 17, 1013–1017 (2021) | [Cite this article](#)

17k Accesses | 409 Citations | 174 Altmetric | [Metrics](#)

Quantum Physics

[Submitted on 5 Oct 2020 ([v1](#)), last revised 30 Nov 2020 (this version, v2)]

A rigorous and robust quantum speed-up in supervised machine learning

[Yunchao Liu](#), [Srinivasan Arunachalam](#), [Kristan Temme](#)

Over the past few years several quantum machine learning algorithms were proposed that promise quantum speed-ups over their classical counterparts. Most of these learning algorithms either assume quantum access to data -- making it unclear if quantum speed-ups still exist without making these strong assumptions, or are heuristic in nature with no provable advantage over classical algorithms. In this paper, we establish a rigorous quantum speed-up for supervised classification using a general-purpose quantum learning algorithm that only requires classical access to data. Our quantum classifier is a conventional support vector machine that uses a fault-tolerant quantum computer to estimate a kernel function. Data samples are mapped to a quantum feature space and the kernel entries can be estimated as the transition amplitude of a quantum circuit. We construct a family of datasets and show that no classical learner can classify the data inverse-polynomially better than random guessing, assuming the widely-believed hardness of the discrete logarithm problem. Meanwhile, the quantum classifier achieves high accuracy and is robust against additive errors in the kernel entries that arise from finite sampling statistics.

Comments: 27 pages, 2 figures

Subjects: Quantum Physics (quant-ph); Machine Learning (cs.LG)

Cite as: arXiv:2010.02174 [quant-ph]

(or arXiv:2010.02174v2 [quant-ph] for this version)

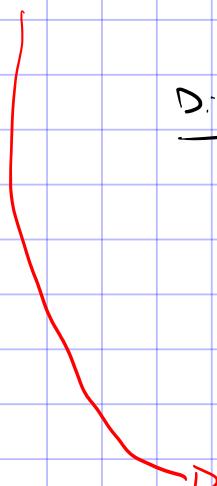
<https://doi.org/10.48550/arXiv.2010.02174>

Journal reference: Nature Physics, 2021

Related DOI: <https://doi.org/10.1038/s41567-021-01287-z>

There Exist
a Dataset
where QuanTum
Kernel Method,
can offer
Exponentional
Advantage

Advantages : Strong Theory, global optimum



Disadvantages

Compute kernel

$$O(N^2)$$



neural networks

$$O(N)$$

Quantum setting No back propagations

QNN

$$O(N^r)$$

r # parameters

QUM

$$O(N^s)$$

\sim # training
data

3 Noise characterization and mitigation -28/06/2025

NOISE CERTIFICATION & CHARACTERIZATION

noise mitigation

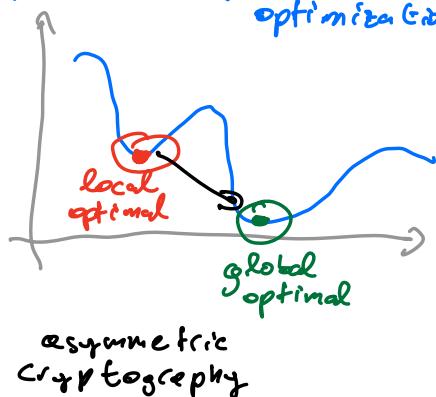
master executive quantum machine learning 2025

Motivation

Quantum computing promises **breakthrough** in a variety of domains = **exastic computational advantage**

material simulation
(Deutsch)

optimization problems (finance: portfolio optimization)



Shor's algorithm (1994)
(integer factorization)

$$N = P_1 \cdot P_2 \quad (\text{prime number})$$

public key private keys

Problem: These algorithms are in principle meant to be carried on ideal, noiseless quantum computers

"noise" is present in every quantum computation

✗
reality

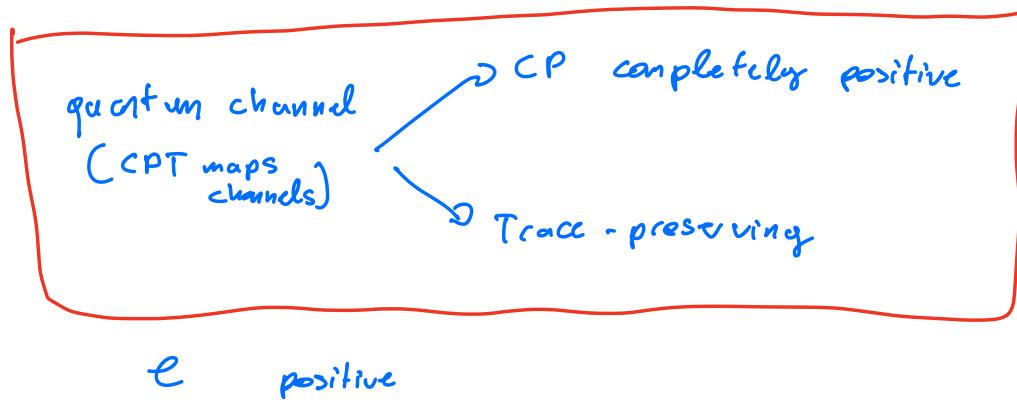
Sources: "channel" environment + noise from device itself

Def (noise) (i) unwanted effects changing the quantum state (where computation/information is encoded) during the process

(ii) difference between ideal vs experimentally implemented quantum gate/operation

(c) Typical noise examples

requisity: noise is modeled as a proper physical quantum channel



⚠ there can be positive channels

$$e \otimes I_m (\tilde{\rho}) \geq 0$$

or

④ Completely positive channels $e (\rho \geq 0)$

$$e \otimes I_m (\rho) \geq 0 \quad \text{for every } m$$

≤ 1

② Trace preservation of $e \quad \text{Tr}[e(\rho)] = \text{Tr}[\rho] = 1$

$$\rho = \begin{pmatrix} p_1 & & \\ & p_2 & \dots \\ & & p_n \end{pmatrix} \quad p_j \geq 0 \quad |\psi\rangle$$

$$e \otimes I (\rho) = \sigma = \begin{pmatrix} q_1 & & \\ & q_2 & \dots \\ & & q_m \end{pmatrix} \quad q_1, q_2, \dots, q_m \geq 0 \quad m \geq n \quad |\phi_j\rangle$$

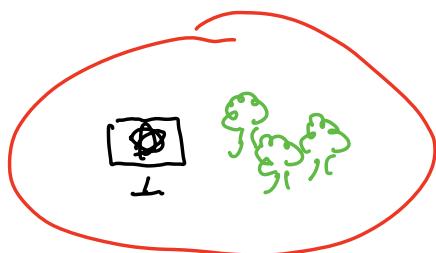
$$\text{Tr}[e(\rho)] = \sum_{j=1}^m q_j = \sum_{j=1}^n p_j$$

Define quantum channels (Kraus operators $\{E_j\}_j$)

$$e(\rho) = \sum_{j=1}^R E_j \rho E_j^\dagger \quad \text{such that } \sum_j E_j^\dagger E_j = 1$$

R = "Kraus Rank"

$$U \otimes U^\dagger \quad E_1 = U \quad U^\dagger U = U U^\dagger = 1$$



Examples

④ Depolarizing noise

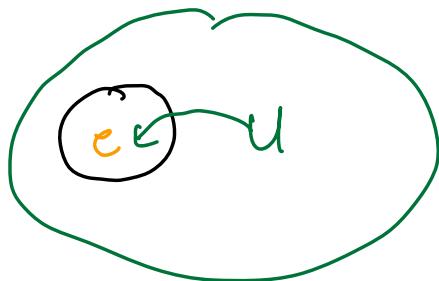
$$e_{\text{dep}}(\rho) = (1-p)\rho + p \frac{11}{d}$$

depolarization strength $0 \leq p \leq 1$

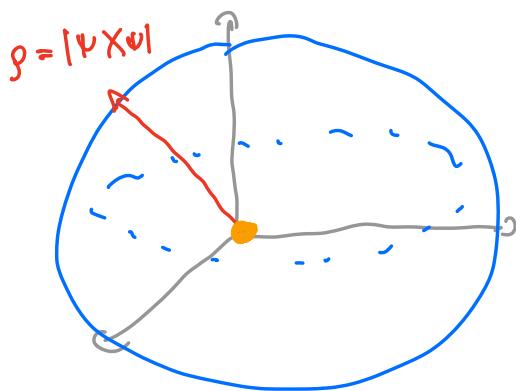
$$\frac{11}{d} = \begin{pmatrix} 1/d & & & \\ & 1/d & \dots & \\ & & \ddots & \\ & & & 1/d \end{pmatrix}$$

Alternatively to Kraus representation,
Stirling dilation Th

$$e = \text{Tr}_R [U g \otimes \sigma U^+]$$



2 axiom of QM: every evolution in a
closed system is unitary

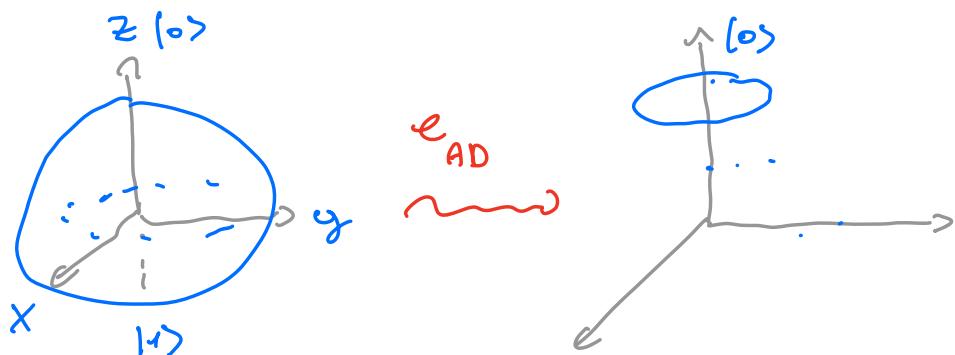


(2) Amplitude damping

$$e_{AD} = E_0 g E_0^+ + E_1 g E_1^+ \quad \text{with} \quad E_0 = \begin{pmatrix} 1 & 0 \\ 0 & \sqrt{\mu-\gamma t} \end{pmatrix}$$

$$E_1 = \begin{pmatrix} 0 & \sqrt{\gamma t} \\ 0 & 0 \end{pmatrix}$$

effect: leave $|0\rangle$ unchanged, but
reduces the amplitude of $|1\rangle$ to a factor
of γ
as the quantum system is skewed towards $|0\rangle$



(3) bit flip /: X gate
phase flip

$$e_{\text{bf}}(g) = (1-p)g + pXgX$$

$$X = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

$$e_{\text{pif}}(g) = (1-p)g + ZgZ$$

$$Z = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$$

$$X \begin{pmatrix} 1 & 0 \\ 0 & \beta \end{pmatrix} X = \begin{pmatrix} p & 0 \\ 0 & \alpha \end{pmatrix}$$

$$X |\psi\rangle = X \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix} = \begin{pmatrix} \alpha_2 \\ \alpha_1 \end{pmatrix}$$

(c) How to access / certify / understand noise

≈ most direct method "quantum process tomography"

Notation: when referring to noise channel, we use Λ (unknown)

reconstructing the noise channel Λ entirely

Step 1: choose a set of pure quantum states

$$|\Psi_1\rangle, |\Psi_2\rangle, \dots, |\Psi_{d^2}\rangle$$

$d =$ dimension
of our Hilbert
space

such that this is a basis
of the set of matrices

$$d = 2^n$$

"tomographically complete set of
states" $n = \#$ qubits

Step 2: apply Λ on each density operator

$$|\Psi_j\rangle\langle\Psi_j| \text{ separately}$$

$$\Lambda(|\Psi_j\rangle\langle\Psi_j|) \quad \text{for all } j=1, \dots, d^2$$

Step 3: perform quantum state tomography on
each, resulting output
"measuring in complete basis"

→ reconstruct a matrix

$$\lambda_{jn} = \text{Tr} [\langle \tilde{\Psi}_j | \tilde{\Psi}_j | \wedge (\langle \tilde{\Psi}_n | \tilde{\Psi}_n |)] =$$
$$\langle \tilde{\Psi}_j | \wedge (\langle \tilde{\Psi}_n | \tilde{\Psi}_n |) | \tilde{\Psi}_j \rangle \quad j=1, \dots, d^2$$
$$d^2 \times d^2 \text{ matrix} \quad n=1, \dots, d^2$$

$$\Lambda(g) = \bigotimes_j X_j \quad g = \perp$$
$$d^2 \times d^2 \quad \sim X = d \times d$$
$$j, n = 1, \dots, d^2$$

→ quantum process tomography delivers a **complete** information about the noise channel Λ

② SPAW errors

(state preparation and measurement errors)

④ Scalability $d^2 \times d^2 = d^4$ operations

$$d = 2^n \quad n = \text{qubits}$$

$$\sim d^4 = 2^{4n} = 16^n$$

$$3 \text{ qubits} \quad 16^3 = 4096$$

Addressing scalability : Randomized Benchmarking

References:

- (1) Emerson, Al�chi, Zyczkowski
"scalable noise estimation with random unitary gates"
(2005)
- (2) Magesan, Gambetta, Emerson
"Characterizing quantum gates via randomized benchmarking"
(2012)
- (3) Helsen, Roth, Ourouti, Werner, Eisert
"A general framework for randomized benchmarking"
(2022)

Setting of RB with Clifford gates

Def (Pauli group)

$$P_1 = \langle \{ -1, 1, i, -i \} \times \{ I, X, Y, Z \} \rangle \quad \begin{matrix} 1\text{-qubit Pauli} \\ \text{group} \end{matrix}$$

$$P_n = \{P_1\}^{\otimes n}$$

A

$$\underbrace{x \otimes y \otimes z \otimes \dots \otimes z}_{n \text{ tensor product}}$$

Def (Clifford group)

$$\mathcal{C}_n := \{U \text{ unitary: } U \Theta U^\dagger \in P_n \text{ for all } \Theta \in P_n\}$$

"the Clifford group is the normalizer of the n-qubit Pauli group"

$$U_1, U_2 \in \mathcal{C}_n \Rightarrow U_1 U_2 \in \mathcal{C}_n$$

$$U_1 \underbrace{U_2 \Theta U_2^\dagger}_{\Theta^1 \in P_n} U_1 = U_1 \Theta^1 U_1 \in P_n$$

dimension of the Pauli group

$n =$	1	2	3	4
	2 ⁴	14'520	~ 82 millions	$> 12 \cdot 10^{12}$

Why do we consider the Clifford group?

because of the twirling property:

$$\frac{1}{|\mathcal{C}|} \sum_{g \in \mathcal{C}} g \circ \lambda \circ g^+ (g) = (1-p) \rho + p \frac{\mathbb{I}}{d}$$
$$= \lambda^{\text{dep}}(g)$$

p is the "average gate fidelity of λ "

$$F(\lambda) = \int_{\text{pure states } \psi} \text{Tr} [|\psi\rangle\langle\psi| \lambda (|\psi\rangle\langle\psi|)] d\mu(\psi)$$



Ex - cursus

universal quantum computation = reproduce every possible unitary gate / evolution

Why the Clifford group is important?

Clifford group + one single non-Clifford gate = universal set
(whatever)

Assumption: we implement ideal Clifford gates $g \in \mathcal{C}$

$$g^2 = 1 \circ g \quad \text{for all } g$$

"gate independent noise assumption

$$\tilde{g} = 1 \quad \text{og}$$

Random procedure for RB

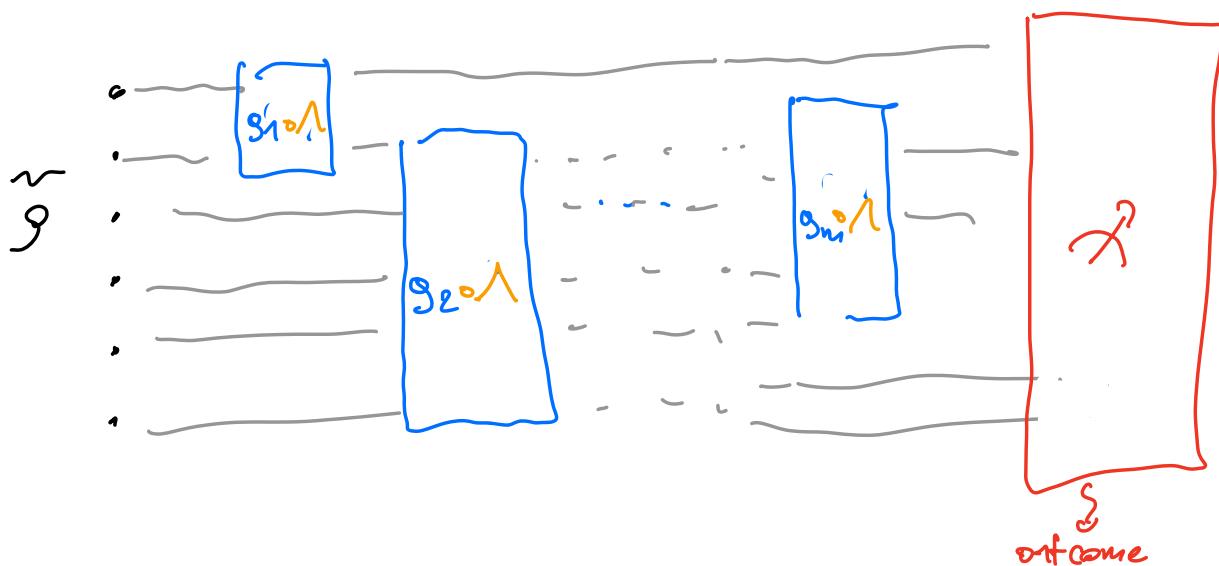
"build random quantum circuits, measure the outcome to obtain the average gate fidelity of the Clifford group"

(1) set a length for the circuit, m

(2) Generate a random sequence of Clifford gates

$$g_1, g_2, \dots, g_m$$

(3) Implement the random sequence as quantum circuit
and measure the final state



(4) repeat steps 2-3 for multiple random sequences
for fixed m

~ compute the "measurement average"

$$Q_{\text{avg}}(m) = \frac{1}{\text{samples}} \sum_j \text{measurements outcomes of random circuit } j$$

(5) increase the sequence length m and repeat the above

~ $Q_{\text{avg}}(m), Q_{\text{avg}}(m+1), \dots, Q_{\text{avg}}(m+k)$

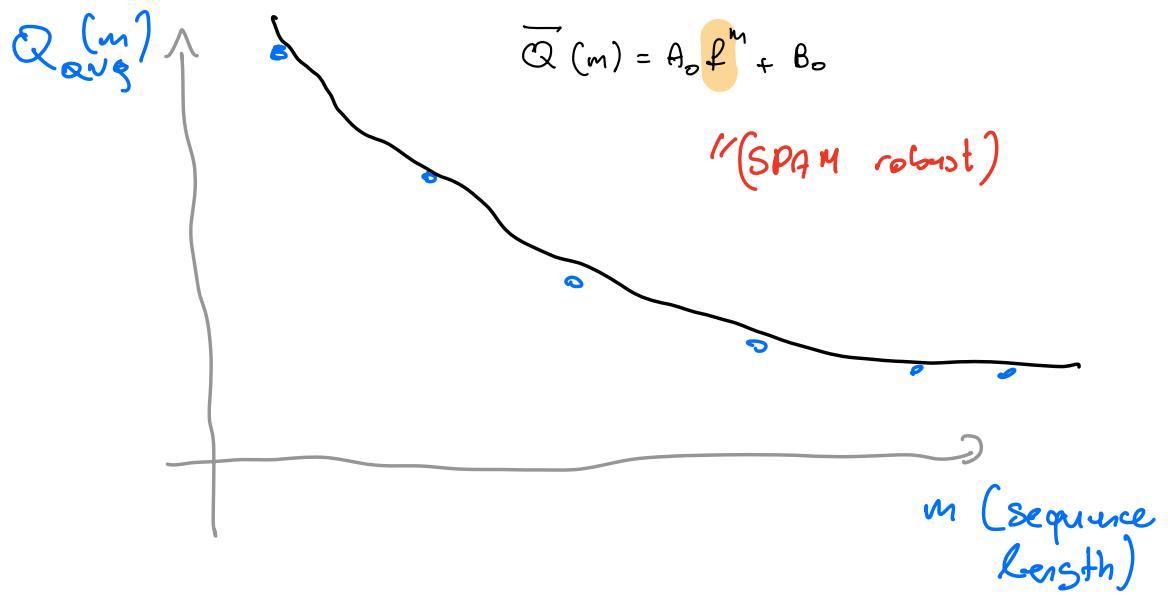
(6) fit the data $\{Q_{\text{avg}}(m)\}_m$ into a fitting model

$$\overrightarrow{Q}(m) = A_0 f^m + B_0$$

to extract the parameter f which is connected to the average gate fidelity $F(\lambda)$ of λ

$$F(\lambda) = f + \frac{(1-f)}{\lambda}$$

this is our final target quantity telling us how close is λ to identity (or $\sigma \approx \hat{\sigma}$)



Theoretical guarantee of RB

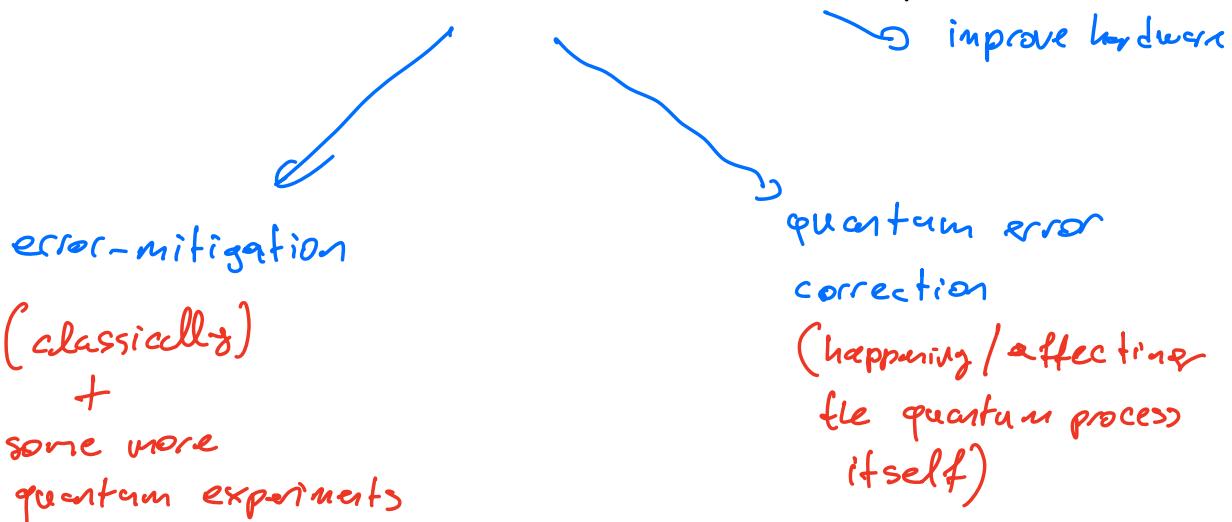
assume we take N samples with \dots

$$N = \text{poly}(n)$$

$$Q_{\text{average}}(m, n) \approx \bar{Q}_{\text{avg}}(m)$$

↑ theoretical mean over
all possible
realizations

How can we use information we retrieve about noise to improve/correct quantum computation?



Simple example: bit flip error channel

$$|\Psi\rangle = \alpha|0\rangle + \beta|1\rangle$$

trick redundancy

$$\text{encode } |\Psi\rangle \text{ into } |\phi\rangle = \alpha|000\rangle + \beta|111\rangle$$

$$\text{and assume } C_{bf} = p X |\phi\rangle \langle \phi| + (1-p)I$$

assume $p < 1$, so that only one qubit at most is flipped per operation

strategy: use projections or leverage "majority vote"

$$P_0 = |000X000| + |111X111| \quad P^2 = P$$

$$P_1 = |000X100| + |011X011|$$

$$P_2 = |010X010| + |101X101|$$

$$P_3 = |001X001| + |110X110|$$

$$\epsilon_{corr} = P_0 g P_0 + \sum_{j=1}^3 X_j P_j g P_j X_j$$

Pauli X on qubit j=1,2,3

$$\Psi = |0\rangle$$

$$\rightarrow |\phi\rangle = |000\rangle$$

↓
bit-flip error
channel

$$|\psi\rangle = |010\rangle$$

$$\epsilon_{corr} (|\psi X \psi|)$$

$$= \epsilon_{corr} (|010X010|) =$$

$$= P_0 |010X010| P_0 = 0$$

$$X_1 P_1 |010X010| P_1 X_1 = 0$$

$$+ X_2 P_2 |010X010| P_2 X_2 = 0$$

$$+ X_3 P_3 (|010X010|) P_3 X_3 = 0$$

$$= 11 \otimes X \otimes 11 (|010X010| + |101X101|) |010X010| \times$$

$$(|010X010| + |101X101|) 11 \otimes X_2 \otimes 11$$

$$= 11 \otimes X \otimes 11 (|010X010|) 11 \otimes X \otimes 11$$

$$X|0\rangle = |1\rangle \quad X|1\rangle = |0\rangle$$

~~$|100\rangle \otimes |000\rangle$~~

Exercise

$$\epsilon_{corr}(|100\rangle \otimes |000\rangle) = |000\rangle \otimes |000\rangle$$

$$\epsilon(|001\rangle \otimes |001\rangle) = |000\rangle \otimes |000\rangle$$

$$\epsilon(|000\rangle \otimes |000\rangle) = |000\rangle \otimes |000\rangle$$

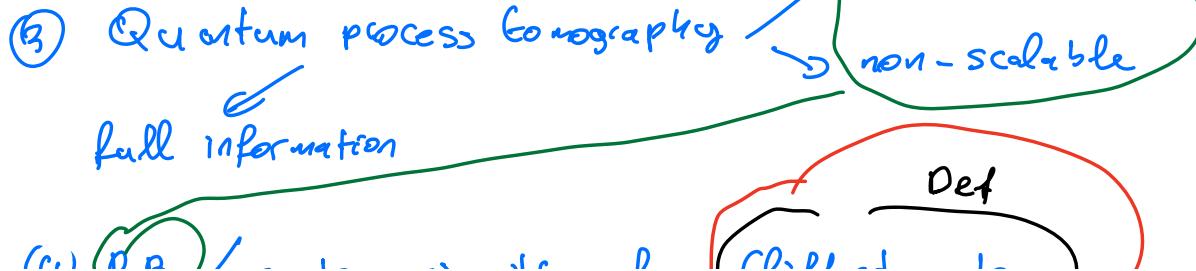
$$|110\rangle \rightarrow |000\rangle$$

Observation: overhead
 ↓
 physical vs logical qubits
 3

↑
 a single logical error-corrected qubit

Morning recap

- ① Short overview about quantum platforms
 - ② noise concept + examples (depolarising, AD, BF)
- quantum channels "CPT maps"



(4) ICD / random circuits for Clifford gate
↓
we learn the "average gate fidelity"
of the gate set + one single non Clifford = universality

sample average vs theoretical average
($\text{poly}(n)$) (unfeasible)

(5) noise mitigation vs error correction
this afternoon correcting single bit-flip error thanks to redundancy

logical error-corrected qubits vs physical qubits

End goal of QC : Fault-tolerant quantum computer
"holy grail of QC"

= "QC that can carry out arbitrarily long and precise quantum computation"

"error is corrected faster than it is generated"



Aharanov, Bein-or

"Fault-tolerant quantum computation with constant error rate" '99

Threshold theorem for fault-tolerant QC

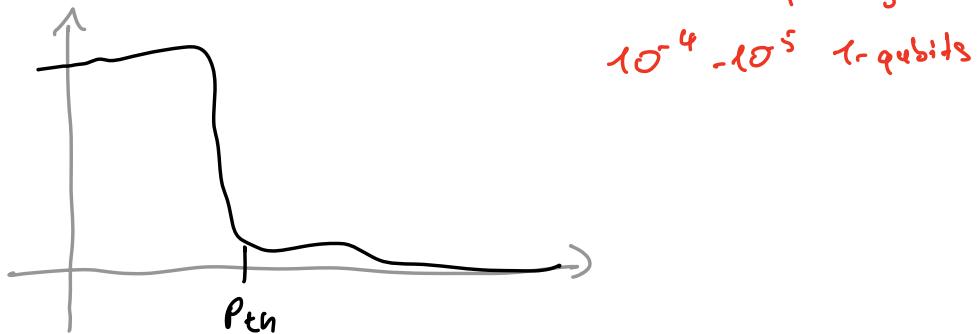
A quantum circuit containing $p(n)$ gates may be simulated with probability error at most ϵ using

$$\mathcal{O}\left(\text{poly}\left(\log \frac{p(n)}{\epsilon}\right) p(n)\right)$$

gates on hardware whose components fails with probability at most $P_{\text{th}} < 10^{-5} \sim 10^{-6}$

10^{-3} 2-qubit gates

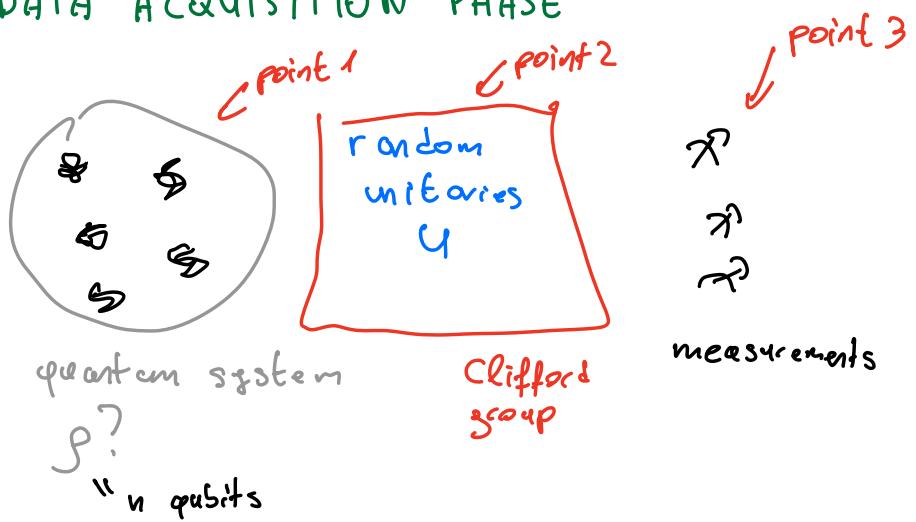
$10^{-4} - 10^{-5}$ 1-qubits



SHADOW ESTIMATION

Reference: "Predicting many properties of a quantum system from very few measurements"
 Huang, Richard Küng, John Preskill
Nature 2020

DATA ACQUISITION PHASE



CLASSICAL POST-PROCESSING PHASE



$$\langle \psi | \theta | \psi \rangle$$

$$\langle \psi \rangle_\theta$$

we are able
to prepare it
many times

Goal: efficiently estimate $\text{Tr}[\theta_i \rho]$ for unknown ρ
 for many $\{\theta_i\}_i$ at once

$$\text{Tr}[\sigma_x |\psi\rangle\langle\psi|] = \langle \psi | \sigma_x | \psi \rangle$$

Protocol: (1) prepare ρ (n -qubit states)

(2) randomly select a Clifford gate U and rotate ρ

$$U \rho U^\dagger$$

(3) Measure in computational basis $\{b\}$

$$\begin{aligned} |0\dots 0\rangle \\ |010\dots\rangle \\ |0111\rangle \end{aligned}$$

→ obtain an outcome b' $P_{b'}[b' \approx b] = \langle b' | U \rho U^\dagger | b' \rangle$

Quantum

(4) Classically reconstruct $U^\dagger |b'\rangle \langle b'| U$

(5) store efficiently the classical shadow

$$\tilde{\rho} = U^{-1} (U^\dagger |b'\rangle \langle b'| U)$$

frame operator (important!)

consisting in the weighted average of the process

$$\mu(\rho) = \mathbb{E}_{\substack{U \sim \text{Clifford} \\ b \in \{0,1\}^n}} \langle b | U \rho U^\dagger | b \rangle U^\dagger |b'\rangle \langle b'| U$$

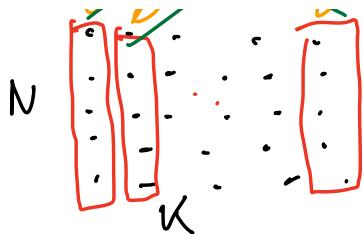
* later: we will see how to reconstruct the inverse easily

PREDICTION PHASE

(A) repeat the data acquisition phase $R = N \cdot K$ times

median

$$\sim \{ \hat{s}_{j,u} \}_{\substack{j=1, \dots, N \\ u=1, \dots, K}}$$

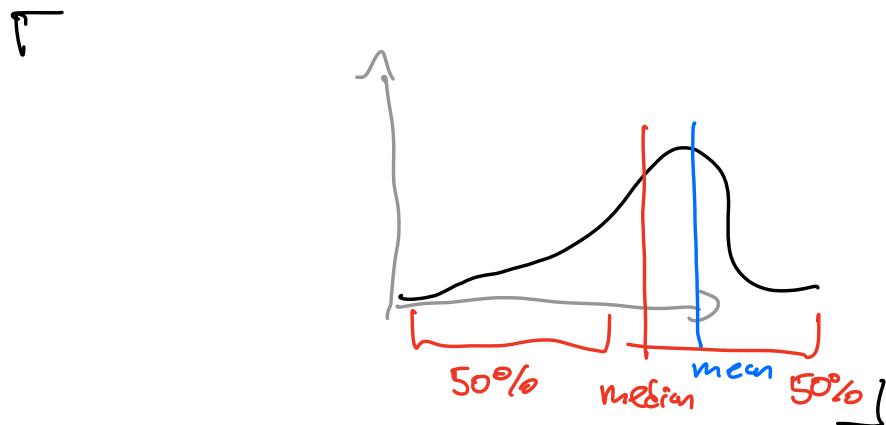


(B) Prepare the mean for K different sets of N data points each

$$\hat{\sigma}_i(K) = \frac{1}{N} \sum_{j=1}^N \text{Tr} [\sigma_i \hat{s}_{j,K}] \quad i=1, \dots, K$$

(C) Estimate the observable σ_i by a "median of means"

$$\hat{\sigma}_i(N, K) = \text{median} \{ \sigma_i(1), \sigma_i(2), \dots, \sigma_i(K) \}$$



(D) repeat it for all wished $\{ \sigma_i \}_{i=1}^u$ observables

Theorem (Assessment prediction guarantee for shadow estimation)

Set accuracy parameters ϵ and $s \in [0,1]$

$$\text{Set } K = 2 \log\left(\frac{2M}{\delta}\right) \quad \text{and} \quad N = \frac{34}{\epsilon^2} \max_i \left\| \frac{\sigma_i - \text{Tr}[\sigma_i]}{2^n} \right\|_{\text{shadow}}$$

shadow norm

Then a collection of $N \cdot K$ independent classical shadows accurately estimates

$$\left| \tilde{\sigma}_i(N, k) - \underbrace{\text{Tr}[\sigma_i]}_{\text{target}} \right| \leq \epsilon \quad \text{for all } i=1, \dots, M$$

with probability at least $1-\delta$.

$$\|\sigma\|_{\text{shadow}} = \max_{\sigma} \left(\mathbb{E} \sum_b \langle b | U \circ U^\dagger | b \rangle \langle b | U \sigma U^{-1}(s) U^\dagger | b \rangle \right)^{1/2}$$

clifford group

$$\left\| \sigma - \frac{\text{Tr}[\sigma]}{2^n} \right\|_{\text{shadow}} \leq 3 \text{Tr}[\sigma^2]$$

for the local 1-qubit Clifford group

$$\left\| \sigma - \frac{\text{Tr}[\sigma]}{2^n} \right\|_{\text{shadow}} \leq 4^{\text{locality}(\sigma)} \|\sigma\|_2$$

Interpretation: + learn exponentially many observables
 of a n -qubit quantum state with
 $R = \Theta(1) \log(M) \frac{1}{\epsilon^2}$ (independent of system size)

- we can hear only local observables

most scalable
procedure in the
system size

Theoretical guarantee comes from

Weingarten calculus
(unitary 3-designs)

Clifford group! \sim scaling in
 n disappear

unitary design

" α distribution over the
unitary group which approximates
the uniform (Haar) measure"

Collins, Matsumoto, NovaU

"The Weingarten calculus"

2022

median of means
no guarantees on
the variance

Berryman, Valion, Vazirani
"Random generation of
combinatorial structures
from a uniform distribution"
1986

Nevirovský, Yudin
"Problem complexity and
method efficiency in
optimization" 1983

Computing the inverse of frame operator M
and address noise

Frame operator

$$M(\rho) = \mathbb{E} \sum_{\text{Un-Clifford } b \in \{0,1\}^{n^2}} \langle b | U \rho U^\dagger | b \rangle U^\dagger | b \rangle \langle b | U$$

taking on average \mathbb{E} over a group (Clifford group)

By group theory, M is symmetric w.r.t action of the group itself

\Rightarrow representation theory / Schur's Lemma

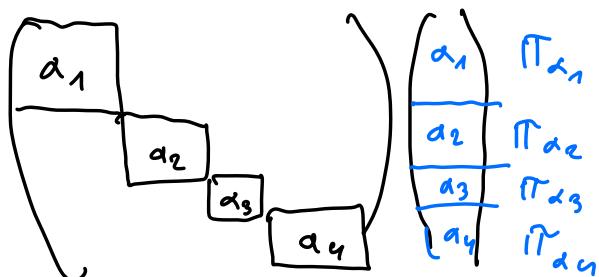
$$M = \sum_{\substack{\alpha \\ \text{irreps}}} \frac{\text{Tr}[B \Pi_\alpha]}{\text{Tr}[\Pi_\alpha]} \quad \Pi_\alpha := \sum_{\alpha} f_\alpha \circ \Pi_\alpha$$

in the "gate representation"

irreducible representation = building blocks

representations = matrices

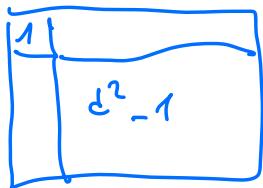
projection onto the invariant subspace of the irrep α



$$B = \sum_b |b, b\rangle \langle b, b|$$

Fact: the defining representation of a 3-design (like the Clifford group) is divided in 2 irreps only

- trivial irrep
- adjoint irrep



$$|\omega \times \omega| \quad \omega = \sum_j |j, \omega\rangle$$

$$M = \Pi_{\text{trivial}} + \frac{\Pi_{\text{adj}}}{d+1}$$

$$\Pi_{\text{adj}} = \mathbb{I}_{d^2} - \Pi_{\text{trivial}}$$

$$M^{-1} = \Pi_{\text{trivial}} + (d+1) \Pi_{\text{adj}}$$

$$\begin{aligned} \Pi^2 &= \Pi \\ \Pi_j \Pi_N &= 0 \end{aligned}$$

$$\begin{aligned} M M^{-1} &= \left(\Pi_{\text{trivial}} + \frac{1}{d+1} \Pi_{\text{adj}} \right) \left(\Pi_{\text{trivial}} + (d+1) \Pi_{\text{adj}} \right) \\ &= \Pi_{\text{trivial}}^2 + (d+1) \cancel{\Pi_{\text{trivial}} \Pi_{\text{adj}}} + \frac{1}{d+1} \cancel{\Pi_{\text{adj}} \Pi_{\text{trivial}}} + \\ &\quad + \frac{d+1}{d+1} \Pi_{\text{adj}}^2 = \\ &= \Pi_{\text{trivial}} + \Pi_{\text{adj}} = \Pi_{\text{trivial}} + \mathbb{I} - \Pi_{\text{trivial}} = \mathbb{I} \end{aligned}$$

Frame operator

$$M(\rho) = \mathbb{E}_{U \sim \text{Clifford}} \sum_{b \in \{0, 1\}^{2n}} \langle b | \tilde{U} \rho \tilde{U}^\dagger | b \rangle U^\dagger |b\rangle \langle b| U$$

Can we calibrate the frame operator to cancel the effect of noise, like if the procedure was carried out with no noise at all?

There are different methods:

(1) "Robust shadow estimation" Chen et. al (2021)

(2) "Noise-mitigated randomized measurements and self-calibrating shadow estimation" Oshita et al (2024)

assume the noise to be gate independent

$$\hat{g} = \lambda \circ g \quad \text{for all } g \quad \lambda(g) = \lambda$$

$$\hat{\mu} = \frac{\text{Tr}[\Pi_{\text{triv}} B \lambda]}{\text{Tr}[\Pi_{\text{triv}}]} \Pi_{\text{triv}} + \frac{\text{Tr}[\Pi_{\text{adj}} B \lambda]}{\text{Tr}[\Pi_{\text{adj}}]} \Pi_{\text{adj}}$$

$$= \Pi_{\text{trivial}} + \frac{\chi_z(\lambda)}{d+1} \Pi_{\text{adj}}$$

where we defined the $\boxed{\chi_z(\lambda) := \frac{\text{Tr}[\Pi_z \lambda]}{d-1}}$

and Π_z is the projection onto Pauli words made of identities and Pauli-Z operators only

$$\Pi_z = B \Pi_{\text{adj}} = \sum_{P \in \{1, Z\}^{\otimes n} / \{I_n\}} \frac{|P\rangle \langle P|}{2^n}$$

Lemma (error bias)

$$\begin{aligned} \|\hat{\boldsymbol{\mu}}_{\text{ideal}}^{-1} \hat{\boldsymbol{\mu}} - \mathbf{I}\| &= \left[\|\boldsymbol{\Pi}_{\text{trivial}} + (\mathbf{I} + \lambda \boldsymbol{\Pi}_{\text{adj}})\| \right] \left[\|\boldsymbol{\Pi}_{\text{trivial}} + \frac{\lambda_2(\lambda)}{\lambda+1} \boldsymbol{\Pi}_{\text{adj}}\| \right] \\ &= \|\boldsymbol{\Pi}_{\text{trivial}} + \lambda_2(\lambda) \boldsymbol{\Pi}_{\text{adj}} - \mathbf{I}\| = (\lambda_2(\lambda) - 1) \|\boldsymbol{\Pi}_{\text{adj}}\| \end{aligned}$$

Idea to noise-mitigate shadow estimation is to substitute

$$\hat{\boldsymbol{\mu}}_{\text{ideal}} \text{ with } \hat{\boldsymbol{\mu}} = \frac{\lambda_2(\lambda)}{\lambda+1} \boldsymbol{\Pi}_{\text{adj}} + \boldsymbol{\Pi}_{\text{trivial}}$$

$$\begin{aligned} \|\hat{\boldsymbol{\mu}}_{\text{ideal}}^{-1} \hat{\boldsymbol{\mu}} - \mathbf{I}\| &= \left(\frac{\lambda+1}{\lambda_2(\lambda)} \|\boldsymbol{\Pi}_{\text{adj}} + \boldsymbol{\Pi}_{\text{trivial}}\| \right) \left(\|\boldsymbol{\Pi}_{\text{trivial}} + \frac{\lambda_2(\lambda)}{\lambda+1} \boldsymbol{\Pi}_{\text{adj}}\| \right) \\ &\approx \|\boldsymbol{\Pi}_{\text{trivial}} + \frac{\lambda_2(\lambda)}{\lambda_2(\lambda)} \boldsymbol{\Pi}_{\text{adj}} - \mathbf{I}\| \end{aligned}$$

assume perfect estimation $\hat{\lambda}(\lambda) = \lambda(\lambda)$

$$\begin{aligned} &= \|\boldsymbol{\Pi}_{\text{trivial}} + \|\boldsymbol{\Pi}_{\text{adj}}\| - \mathbf{I}\| = \\ &= \|\boldsymbol{\Pi}_{\text{trivial}} + \mathbf{I} - \|\boldsymbol{\Pi}_{\text{trivial}}\| - \mathbf{I}\| = \\ &= \|\mathbf{I} - \mathbf{I}\| = 0 // \end{aligned}$$

Recap: ① we assumed we were able to access
the parameter $\lambda_2(\lambda)$ or estimate $\hat{\lambda}_2(\lambda)$

② classically apply $\hat{\boldsymbol{\mu}}^{-1}$ instead of

the noise-unaware M_{ideal}^{-1}

- ③ Assuming our estimation is very close to
the real $\lambda_2(\lambda)$, then
 $\hat{\mu}^{-1} \mu^k - 1 \approx 0$

Goal: how to actually retrieve $\lambda_2(\lambda)$?

④ AE-hoc calibration

(i) prepare $|0\rangle = |0\rangle^{\otimes n}$, sample at random a
Clifford gate

(ii) measure in computational basis \rightarrow get b

(iii) Compute $\hat{\lambda}_2^{(r)} = \frac{4Kb/U|0\rangle|^2}{d-1}$
classically

(iv) repeat steps 1-3 for $R=N \cdot K$ rounds

(v) compute $\hat{\lambda}_2(\lambda)$ as a median-of-means (R, N, K)
(like shadow estimation)

② Use RB

\nearrow scalable fashion
 \searrow SPAM robust

Recall: that the parameters of RB

$$\text{are } \mu_\alpha = \frac{\text{Tr}[\Pi_\alpha \lambda]}{\text{Tr}[\Pi_\alpha]}$$

compared
to the
parameter
of shadow
estimation

$$f_\alpha = \frac{\text{Tr}[\Pi_\alpha B \lambda]}{\text{Tr}[\Pi_\alpha]}$$

(c) for depolarizing noise:

$$\mu_\alpha = f_\alpha \Rightarrow \text{perfect inversion}$$

(d) for the bit-flip noise

$$\mu_{\text{adj}} = \frac{d^2(1-p)^n - 1}{d^2 - 1} \quad \text{vs} \quad f_\alpha = \frac{d(1-p)^n - 1}{d - 1}$$

no error bias

$$\widehat{\mu}_{\text{RB}}^{-1} M^\lambda - 1 = \begin{pmatrix} \frac{d^2 - 1}{d^2(1-p)^n - 1} & \frac{d(1-p)^n}{d - 1} - 1 \\ & \end{pmatrix} \Pi_{\text{adj}}$$

$$= \left\{ (d+1) \frac{d(1-p)^n - 1}{d^2(1-p)^n - 1} - 1 \right\} \Pi_{\text{adj}}$$

$$\xrightarrow{d \gg 1} \left\{ \frac{d+1}{d+1} - 1 \right\} \Pi_{\text{adj}} = 0$$

(e) RB with CNOT-dihedral (substituting Clifford group)

WRAP- UP

Discussed shadow estimation

highly scalable method to learn many properties
of unknown quantum states

randomization + statistical
mechanics
techniques

- (•) It is affected by noise, introducing an error bias



Separate RB to correct the error bias
calibration in classical post-processing

- (•) We discussed how to combine shadow estimation as classical post-processing of RB to learn more granular information about a noise channel λ .

- (•) Broader take-home message

interplay between quantum experiments
and classical post-processing to push forward
quantum learning/noise characterization