

Quantumly sampling optimised random features to speed up kernel based methods for Machine Learning

arXiv:2004.10756

Sathyawageeswar Subramanian

Department of Computer Science & DIMAP
University of Warwick, UK

Sathya.Subramanian@warwick.ac.uk

Acknowledgements

Co-authors and Funding

This talk is based on joint work with Dr. Hayata Yamasaki¹, Dr. Sho Sonoda², and Prof. Masato Koashi³, published in NeurIPS 2020:

Learning with Optimized Random Features: Exponential Speedup by Quantum Machine Learning without Sparsity and Low-Rank Assumptions;

H Yamasaki, S Subramanian, S Sonoda, M Koashi; *Advances in Neural Information Processing Systems, 2020.*

Funding: This work was supported by CREST (Japan Science and Technology Agency) JPMJCR1671, Cross-ministerial Strategic Innovation Promotion Program (SIP) (Council for Science, Technology and Innovation (CSTI)), JSPS Overseas Research Fellowships, a Cambridge-India Ramanujan scholarship from the Cambridge Trust and the SERB (Govt. of India), and JSPS KAKENHI 18K18113.

¹ Austrian Academy of Sciences (JST PRESTO), ² RIKEN AIP, ³ The University of Tokyo

1 Introduction

- Supervised Learning
- Random Features
- Motivation for invoking QML

2 Main result

3 Proof Components

- Input model: QRAM
- Implementing matrix functions: QSVT
- Avoiding sparsity and low-rank assumptions: QFT

4 Summary and Outlook

Supervised Learning

labeled examples : $(\vec{x}_1, y_1), \dots, (\vec{x}_N, y_N)$

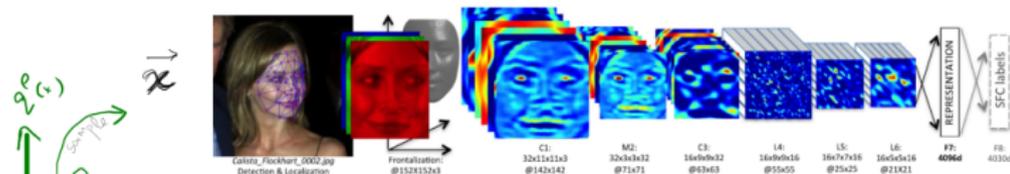
Unknown, true, underlying function

output labels
eg: cat

input data
eg: images
 $\mathcal{X} \ni \vec{x}$

$$f: \mathcal{X} \rightarrow \mathcal{Y}$$

$y \in \mathbb{R}$



→ Probability measure : $dP(x) = q(x)dx$
(DATA DISTRIBUTION)

$\mathcal{X} \subset \mathbb{R}^D$

$\mathcal{Y} \subset \mathbb{R}$

Image from Taigman, Yang, Ranzato, Wolf (2014)

Supervised Learning

Models and Kernels

- Learn (approximation) of f from a few labeled examples; $N =$ sample complexity
- **Model:** Find best $\hat{f} \in \mathcal{F}$ a chosen space of functions
- **Generalisation error:** minimise average error over (unseen) data;
 $\int_{\tilde{\mathcal{X}}} d\rho(x) \left| \hat{f}(x) - f(x) \right| \leq \epsilon$ — learning to desired accuracy ϵ
- **Kernel:** “Similarity function” $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$; way to use known data to predict f on unknown data
- **Translation Invariance:** $k(x, y) = \tilde{k}(|x - y|)$

Random feature expansions

- **Idea:** express the unknown function as an integral / linear combination over “feature functions” determined using input data
- **Kernel** can be written as an average over feature functions, with **some distribution over a ‘feature space’** $\mathcal{V} \subset \mathbb{R}^{D'}$
- **1D Feature (function):** $\varphi : \mathcal{V} \times \mathcal{X} \rightarrow \mathbb{R}$, $\varphi(\mathbf{v}_m, \mathbf{x}) = e^{-2\pi i \mathbf{v}_m \cdot \mathbf{x}}$
- Ultimately, target function can be written as

$$\hat{f}_\alpha(\mathbf{x}) := \sum_{m=0}^{M-1} \alpha_m \varphi(\mathbf{v}_m, \mathbf{x}),$$

with coefficients $\alpha_i \in \mathbb{R}$ to be determined by regression, minimising generalisation error

Random feature expansions

Optimised Random Features

- **Aim:** Minimise the number M of random features required and so the length of the linear combination \hat{f} by using feature functions optimised for the input of labeled examples, and the choice of kernel
- **Conventional RF [Rahimi & Recht (2008)]:**
Sample $v_0, \dots, v_{M-1} \sim d\tau$; get features $\varphi(v_i, \cdot)$; Regression gives α_i ;
 $d\tau$ data independent, # features $M = \tilde{O}(1/\epsilon^2)$
- **Optimised RF [Bach (2017)]:**
Sample $v_0, \dots, v_{M-1} \sim q_\epsilon^*(v) d\tau(v)$; ..regression gives α_i ;
 $q_\epsilon^*(v) d\tau(v)$ data-optimised, # features $M = \tilde{O}(\log^2 1/\epsilon)$
- Provably optimal upto logarithmic factors
- After discretisation and other technicalities,

$$Q_\epsilon^*(\tilde{v}) \propto \left\langle \varphi(\tilde{v}, \cdot) \left| \hat{q}^\rho(\hat{\Sigma} + \epsilon \mathbb{1})^{-1} \right| \varphi(\tilde{v}, \cdot) \right\rangle \quad (1)$$

Optimised Random Features

Our Problem setup

- **Integral operator Σ :**

$$(\Sigma f)(x') := \int_{\mathcal{X}} d\rho(x) k(x', x) f(x). \quad (2)$$

- After discretisation and other technicalities, (with standard bracket notation for inner products etc.)

$$Q_{\epsilon}^*(\tilde{v}) \propto \left\langle \varphi(\tilde{v}, \cdot) \left| \hat{q}^{\rho}(\hat{\Sigma} + \epsilon \mathbb{1})^{-1} \right| \varphi(\tilde{v}, \cdot) \right\rangle, \quad (3)$$

where $q^{\rho}(x)$ is the probability density function over the input space

Function / operator on \mathcal{X}	Vector / operator on $\mathcal{H}^{\mathcal{X}}$
$f : \mathcal{X} \rightarrow \mathbb{C}$	$ f\rangle \propto \sum_{\bar{x}} f(\bar{x}) \bar{x}\rangle$
$\varphi(v, \cdot) : \mathcal{X} \rightarrow \mathbb{C}$	$ \varphi(v, \cdot)\rangle \propto \sum_{\bar{x}} \varphi(v, \bar{x}) \bar{x}\rangle$
$\tilde{k} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$	$k := \sum_{\bar{x}', \bar{x}} \tilde{k}(\bar{x}', \bar{x}) \bar{x}'\rangle \langle \bar{x} $
$q^{\rho} : \mathcal{X} \rightarrow \mathbb{R}$	$q^{\rho} := \sum_{\bar{x}} q^{\rho}(\bar{x}) \bar{x}\rangle \langle \bar{x} $
Σ acting on $f : \mathcal{X} \rightarrow \mathbb{C}$	$\Sigma := k q^{\rho}$
$\Sigma f : \mathcal{X} \rightarrow \mathbb{C}$	$\Sigma f\rangle$

Motivation

Why this method, and why Quantum?

- Scope & domain of applicability Bach's method extend beyond supervised ML (eg. ridge regression, clustering can all be 'kernelised'), to numerical computation of high-dimensional integrals in signal processing, applied mathematics, Bayesian inference, function approximation, optimisation

Motivation

Why this method, and why Quantum?

- Scope & domain of applicability Bach's method extend beyond supervised ML (eg. ridge regression, clustering can all be 'kernelised'), to numerical computation of high-dimensional integrals in signal processing, applied mathematics, Bayesian inference, function approximation, optimisation
- Sampling from the optimised distribution Q_ϵ^* over features = **classical bottleneck**; best known algorithm requires inversion of a full-rank, non-sparse, $\mathcal{O}(e^D)$ dimensional matrix, incurring worst case runtime $\mathcal{O}(e^D)$
- Q_ϵ^* involves inverting integral operator. **Inversion \implies HHL09? But dense, full-rank.** We show how to overcome these difficulties.

Reference

- 1 Francis Bach (2017). On the equivalence between kernel quadrature rules & random feature expansions. Journal of Machine Learning Research, 18, 1–38.

- **Kernels & QML:** Quantum enhanced features [Havlicek et al. (Nature 2019)], Experimental Kernel-based QML [Bartkiewicz et al. (Nature Scientific Reports 2020)], Vedaie et al. (2011.09694), Schuld (2101.11020), Park et al. (2004.03489), Blank et al. (npj QI 2020) ...

- **Kernels & QML:** Quantum enhanced features [Havlicek et al. (Nature 2019)], Experimental Kernel-based QML [Bartkiewicz et al. (Nature Scientific Reports 2020)], Vedaie et al. (2011.09694), Schuld (2101.11020), Park et al. (2004.03489), Blank et al. (npj QI 2020) ...

Category	Method	Title
Quantum version of SVM	Grover algorithm	Quantum optimization for training support vector machines (Anguita et al. 2003)
Quantum version of SVM	HHL algorithm	Quantum support vector machine for big data classification (Rebentrost et al. 2014)
Experimental	NMR 4-qubit quantum processor	Experimental implementation of a quantum support vector machine (Li et al. 2015)
Experimental	IBM quantum experience	Quantum algorithm Implementations for beginners (Patrick et al. 2018)
Quantum version of SVM and ECOC	HHL algorithm	Quantum error-correcting output codes (Windridge et al. 2018)]
Kernel methods	Variational quantum circuit	Quantum machine learning in feature Hilbert spaces (Schuld and Killoran 2019)
Kernel methods	Variational quantum circuit	Supervised learning with quantum-enhanced feature spaces (Havlicek et al. 2019)
Kernel methods	Topological quantum computation	Hamming distance kernelisation via topological quantum computation (Di Pierro et al. 2017)

Table 1 from Mengoni, R., Di Pierro, A. Kernel methods in Quantum Machine Learning. Quantum Machine Intelligence **1**, 65–71 (2019).

Our contributions

- Even with powerful tools like QRAM, attaining genuine quantum speedups has been difficult
- Our work circumvents (1) sparsity and (2) low rank assumptions by exploiting a combination of the QFT and QSVT
- Our asymptotic speedup is exponential over the best known classical algorithm, in many useful parameter regimes
- Technical analysis of Fourier sparse operators and construction of block encodings could be more widely applicable

Main result

Quantumly sampling optimised random features

Theorem

Given D -dimensional data, for any accuracy $\epsilon > 0$ we can sample a (discrete) optimised random feature $\tilde{v} \in \tilde{\mathcal{V}}$ from a weighted distribution $Q_\epsilon^*(\tilde{v})P^T(\tilde{v})$ with

$$\sum_{\tilde{v} \in \tilde{\mathcal{V}}} |(Q(\tilde{v}) - Q_\epsilon^*(\tilde{v}))P^T(\tilde{v})| \leq \delta',$$

in runtime

$$T = \tilde{O}(D \log D) \times \tilde{O}\left(\frac{Q_{\max}^*}{\epsilon} \text{poly log } \frac{1}{\delta'}\right).$$

Main result

Quantumly sampling optimised random features

Theorem

Given D -dimensional data, for any accuracy $\epsilon > 0$ we can sample a (discrete) optimised random feature $\tilde{v} \in \tilde{\mathcal{V}}$ from a weighted distribution $Q_\epsilon^*(\tilde{v})P^T(\tilde{v})$ with

$$\sum_{\tilde{v} \in \tilde{\mathcal{V}}} |(Q(\tilde{v}) - Q_\epsilon^*(\tilde{v}))P^T(\tilde{v})| \leq \delta',$$

in runtime

$$T = \tilde{O}(D \log D) \times \tilde{O}\left(\frac{Q_{\max}^*}{\epsilon} \text{poly log } \frac{1}{\delta'}\right).$$

In particular, T is linear in D , while the best known classical algorithm for estimating $Q_\epsilon^*(\tilde{v})P^T(\tilde{v})$ requires $\mathcal{O}(e^D)$ time.

- First obtain consistent and asymptotically exact discretisation scheme, $\mathcal{X} \mapsto \tilde{\mathcal{X}}, \mathcal{V} \mapsto \tilde{\mathcal{V}}$ etc.
- Design quantum state with $Q_\epsilon^*(\tilde{\mathcal{V}})$ as amplitudes
- Use properties of translation invariant kernel to obtain decomposition of integral operator into simple components: a full rank diagonal operator and the QFT: **technical tools** are perfect reconstruction of the kernel via translation invariance, regularisation
- Implement using QRAM, QSVT, QFT

- **Want to sample from:**

$$Q_\epsilon^*(\tilde{v}) \propto \left\langle \varphi(\tilde{v}, \cdot) \left| \hat{q}^\rho (\hat{\Sigma} + \epsilon \mathbb{1})^{-1} \right| \varphi(\tilde{v}, \cdot) \right\rangle \quad (4)$$

- **Quantum state**² that does the job:

$$|\Psi\rangle^{XX'} \propto \sum_{\tilde{v} \in \tilde{\mathcal{X}}} \hat{\Sigma}_\epsilon^{-\frac{1}{2}} |\tilde{v}\rangle^X \otimes \sqrt{Q^\top F_D^\dagger} \sqrt{\hat{q}^\rho} |\tilde{v}\rangle^{X'} \quad (5)$$

- **Decomposition of $\hat{\Sigma}$:**

$$\hat{\Sigma}_\epsilon \propto \sqrt{\hat{q}^\rho} \cdot F_D^\dagger Q^\top F_D \cdot \sqrt{\hat{q}^\rho} + \epsilon \mathbb{1}, \quad (6)$$

where F_D is the quantum fourier transform (QFT) on \mathbb{C}^D .

²'Aggregated' weights $Q^\top(\tilde{v}) := \sum_{\nu \in \mathbb{Z}^D} q^\top(\tilde{v} + \nu)$

Input model

QRAM for quantum access to Big Data

- **Classically**: memory contents at address k retrievable in $\mathcal{O}(1)$ time

$$\text{RAM} : k \mapsto f(k)$$

- **Quantumly**: unitary version of RAM running in quantum superposition [e.g. Jiang et al. (2019), Hann et al. (2019)]

$$\text{QRAM} : \sum_k \alpha_k |k\rangle |0\rangle \mapsto \sum_k \alpha_k |k\rangle |f(k)\rangle$$

- Linear algebra algorithms using QRAM typically have complexity scaling with the **Frobenius norm** of the input matrix

Input model: \hat{q}^ρ

QRAM for quantum access to Big Data

- probability measure $d\rho(x) = q^\rho(x)dx$; empirical pmf $\hat{q}^\rho(x)$
- With $\mathcal{O}(N)$ classical preprocessing to count the N input data points, construct $\log N$ depth binary tree for addressing data [e.g. Kerenidis & Prakash (2017)]
- We use $\hat{q}^\rho(x)$ embedded into a diagonal operator $\hat{\mathbf{q}}$

$$|0\rangle \mapsto \sum_x \sqrt{\hat{\mathbf{q}}^\rho(x)} |x\rangle$$

- Frobenius norm of $\sqrt{\hat{\mathbf{q}}}$ is unity since q is a probability distribution

Matrix functions and block-encodings

- First introduced to study Hamiltonian simulation, has found a variety of applications in the last few years.
- A block encoding U_A of a Hermitian A is a unitary that encodes a (sub-)normalised version of A in its top left block

$$U_A = \begin{pmatrix} A/\alpha & \cdot \\ \cdot & \cdot \end{pmatrix},$$

where $\alpha \geq \|A\|$

- 1 A Gilyen et al. Quantum singular value transformation and beyond: exponential improvements for quantum matrix arithmetics, 2018 (arXiv:1806.01838).
- 2 S Subramanian et al. Implementing smooth functions of a Hermitian matrix on a quantum computer, 2018 (arXiv:1806.06885).
- 3 A Childs et al. Quantum Algorithm for Systems of Linear Equations with Exponentially Improved Dependence on Precision, 2015 (arXiv:1511.02306).

Matrix functions and block encodings: $\hat{\Sigma}^{-1/2}$

Quantum Singular Value Transformations

- Quantum Singular Value Transformations can use block encodings to probabilistically implement non-unitary operators

$$U_A = \begin{pmatrix} A/\alpha & \cdot \\ \cdot & \cdot \end{pmatrix} \mapsto \begin{pmatrix} \tilde{f}(A) & \cdot \\ \cdot & \cdot \end{pmatrix} \approx U_{f(A)},$$

$$\left\| \tilde{f}(A) - f(A) \right\| < \delta.$$

- Using the block encoding of U_A roughly $\tilde{O}(\kappa \cdot \text{poly log } \frac{D}{\delta})$ times, where κ lower bounds the condition number of A , one can obtain a block-encoding $A^{-1/2}$ to precision δ via the method of QSVT using polynomial approximations of the target function $x^{-1/2}$

Avoiding sparsity and low-rank assumptions

Quantum Fourier Transform

- The structure in this sampling problem and its deep relation to Fourier analysis arising from the use of translation invariant kernels allows us to avoid assuming sparse or low rank input, in contrast to seminal quantum algorithms such as [Harrow et al. \(matrix inversion, 2009\)](#), [Lloyd et al. \(PCA, 2013\)](#), or [Kerenidis & Prakash \(recommendation systems, 2016\)](#)
- Classical FFT on dimension N requires $\mathcal{O}(N \log N)$ time, whereas QFT requires only $\tilde{\mathcal{O}}(\log^2 N)$ time
- Essentially, translation-invariance indicates circulant matrix structure, can be diagonalised by the FFT and then easily inverted in the Fourier basis. [Source of classical bottleneck could be the \$\mathcal{O}\(e^D\)\$ dimensionality of \$\hat{\Sigma}\$](#)

Avoiding sparsity and low-rank assumptions

Resistance to dequantisation by low-rank methods

- Low rank + QRAM recently found to be dequantisable (Tang et al. (2018), Le Gall et al. (2019), Chia et al. (2019, [morning session](#))) via low-rank + classical ℓ_2 sampling
- Because we work with a full rank operator that might also be dense and have large spectral norm, these existing dequantisation methods are not directly applicable

Odds & Ends: after sampling the \tilde{v}_i and deciding $\varphi(\tilde{v}_i, \cdot)$, we do doubly stochastic gradient descent for regression to obtain the coefficients, classically in linear $\mathcal{O}(D)$ time, without losing our speedup

Summary

- **Novelty:** (1) A niche where big data has small Frobenius norm, suitable for QRAM
(2) Non-sparse and full rank operator inverted by taking advantage of Fourier sparsity + QFT + Quantum Singular Value Transformations
- We also show that careful application of (doubly) stochastic gradient descent allows regression to learn the coefficients α_m in $\mathcal{O}(D)$ time, without canceling out the quantum speedup
- Hence widely applicable, promising candidate for 'killer applications' of Quantum Computing / Quantum Machine Learning
- **Drawbacks:** Practicality - not NISQ friendly, we focus on asymptotic complexity, resulting circuits are huge and require thousands of qubits, long coherence times and fault tolerance for QSVT

Outlook

Interesting open directions

- **Specific applications of our framework:** (binary) classification, SVMs, regression
- Other applications where an operator is sparse in a Fourier transformed representation
- **Non-QRAM input models:** classical access [e.g. Arunachalam et al. (arXiv:2010.02174)]
- Investigating the potential for **NISQ** application in a QRAM-free model, optimising circuit constructions

- **NeurIPS 2020:** H Yamasaki, S Subramanian, S Sonoda, M Koashi;
[Learning with Optimized Random Features: Exponential Speedup by Quantum Machine Learning without Sparsity and Low-Rank Assumptions](#)
- **Full(er) paper:** [arXiv:2004.10756](#)
- **More verbose explanations and details:** S Subramanian (2020). [Quantum Algorithms for Matrix Problems and Machine Learning \(Doctoral thesis\)](#).

Thank you!