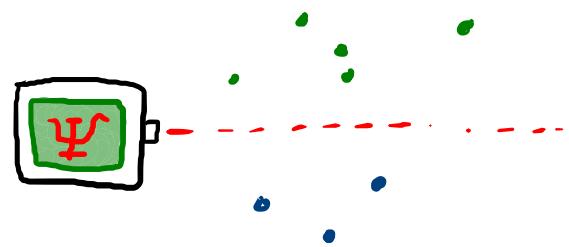


QUANTUM MACHINE LEARNING & KERNEL METHODS



Based on Schuld, 2101.11020

David Wakeham
UBC Strings Talk

Sep 10, 2021

WHAT IS QML?

- Quantum machine learning (QML) is ill-defined.
- It could mean learning quantum data, building models using quantum algorithms, or both!

ALGORITHM / MODEL

		classical	quantum
		classical	CQ
DATA	classical	CC	
	quantum	QC	QQ

- The top right is the hype corner. I'll focus on concrete example of supervised learning.

I

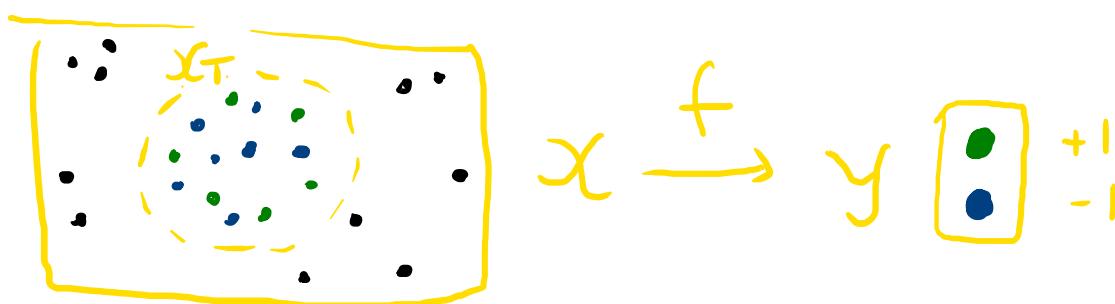
SOME CLASSICAL BACKGROUND

SUPERVISED LEARNING

- Basic idea: data points $x \in \mathcal{X}$ we are trying to classify, e.g. pictures of donuts vs. bagels.



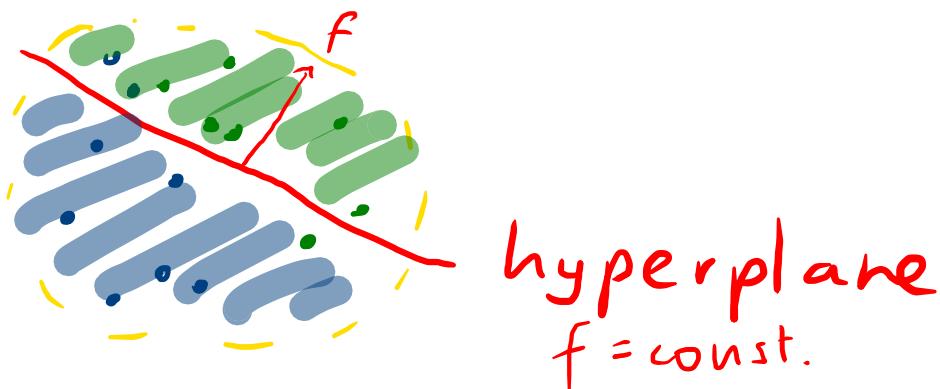
- In supervised learning, we are given access to a training set $x^m \in \mathcal{X}_T$ of M labelled data points.



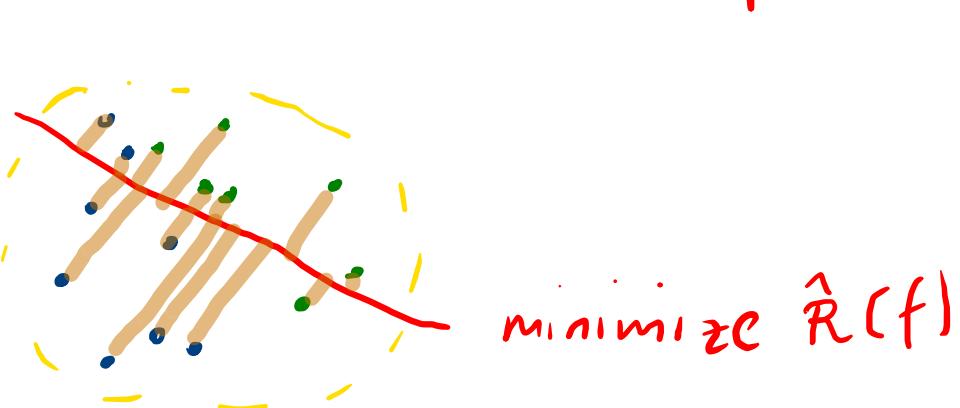
- We use this training data to create a model $f(x)$ which guesses the label $y \in \mathcal{Y}$ for all of \mathcal{X} .

LINEAR MODELS

- Lots of ways to build this predictor $f: \mathcal{X} \rightarrow \mathcal{Y}$, but simplest class is set of linear models.
- Typically, data space \mathcal{X} is a vector space, and I can separate training data using a hyperplane:

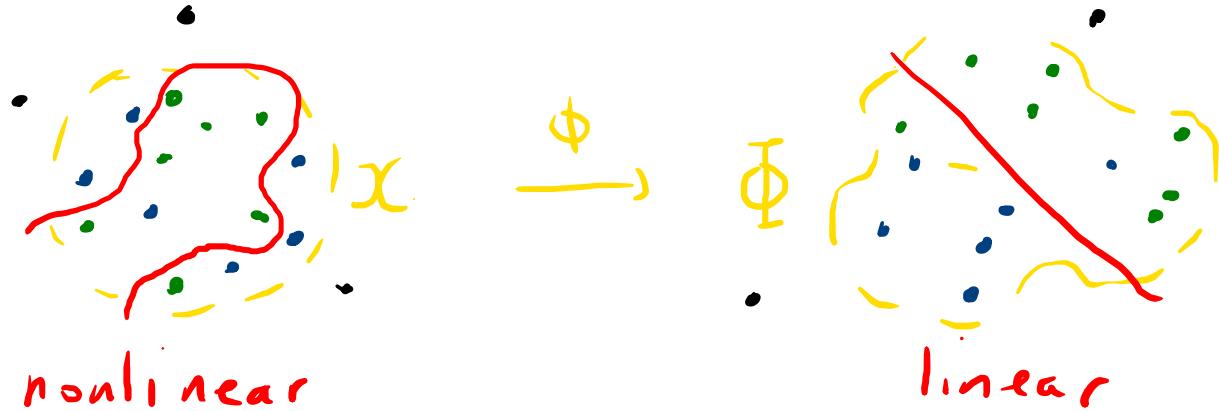


- The hyperplane minimizes an empirical risk function $\hat{R}(f)$:



NONLINEAR MODELS

- Donuts and bagels aren't well-linearly separated in image space.

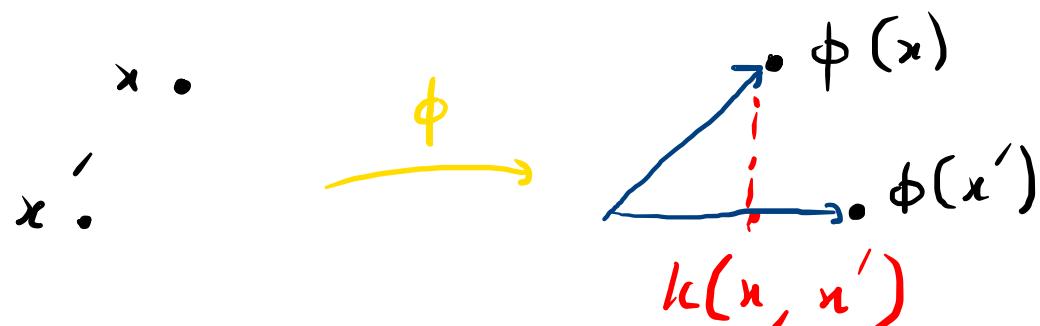


- Workaround: map data to "feature space", $\phi(x) \in \Phi$, where a linear decision boundary is better.
- Unfortunately, useful feature spaces are high-dimensional! Explicitly mapping/manipulating feature vectors is costly.

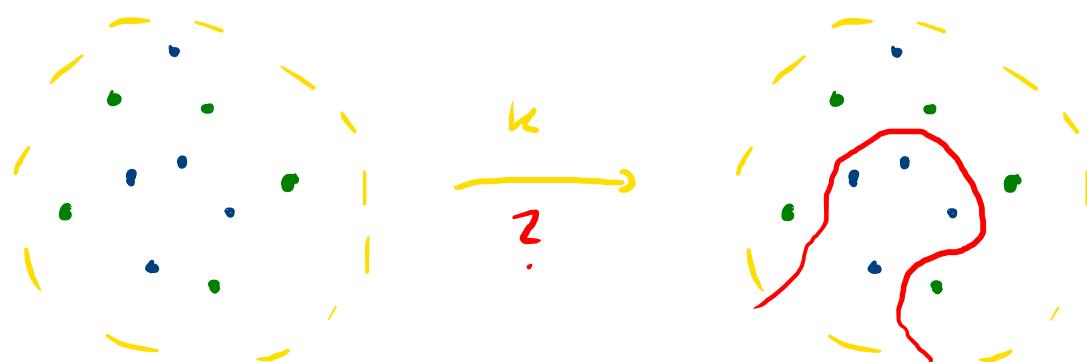


THE KERNEL TRICK

- But there's a **hack** for avoiding costly explicit manipulation.
We work with **inner products** instead of feature vectors:



- The function taking a pair of data point and giving the inner product of feature vectors is the **kernel**:
$$k(x, x') = \langle \phi(x), \phi(x') \rangle_{\Phi}.$$
- These are just numbers! But what can we do with them?



REPRODUCING KERNELS

- It's useful to talk about hyperplanes themselves. Mathematically, we use the reproducing kernel Hilbert space (RKHS).
- Use the kernel to map each data point x to a function λ_x :

$$x \mapsto \lambda_x(\cdot) = k(x, \cdot).$$

$$\begin{array}{c} x \\ \downarrow \lambda \\ \text{hom}(x, \mathbb{R}) \end{array}$$

- Next, define an inner product which reproduces the kernel:

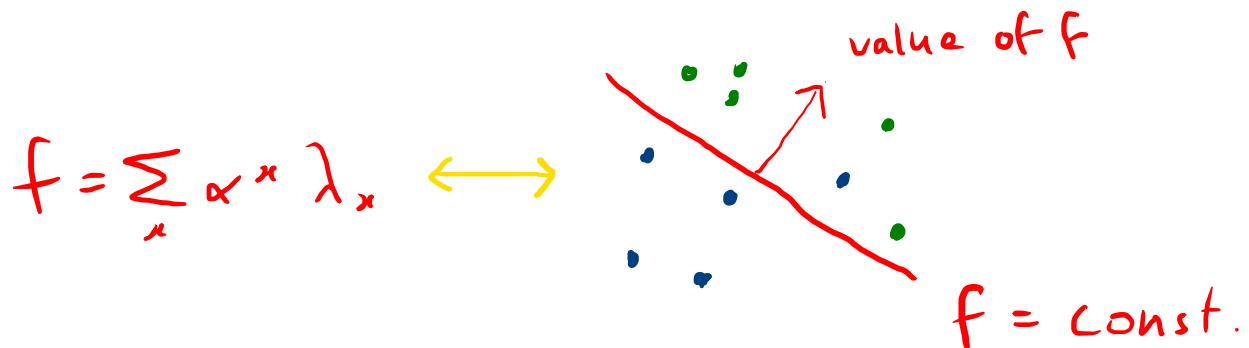
$$\langle \lambda_u, \lambda_{u'} \rangle_F = k(u, u').$$

$$\begin{array}{ccc} X^2 & & \\ (\lambda, \lambda) \downarrow & \searrow k & \\ \text{hom}^2 & \xrightarrow{\langle \cdot, \cdot \rangle_F} & \mathbb{R} \end{array}$$

- To create a Hilbert space F , we take the linear span of λ_x , make $\langle \cdot, \cdot \rangle_F$ bilinear, and complete limits as usual.

THE REPRESENTER THEOREM

- The RKHS is (in essence) the space of **linear functionals** on the span of $\phi(x)$, i.e. the space of linear models:



- We might expect that the models involving only training data $x^m \in \mathcal{X}_T$ are **optimal** in some sense. This is true!
- More precisely, the **Representer Theorem** states that any empirical risk-minimizing $f^* \in \mathcal{F}$ takes the form

$$f_m^*(x) = \sum_m \alpha^m \lambda_{x^m}(x) = \sum_m \alpha^m k(x^m | x).$$

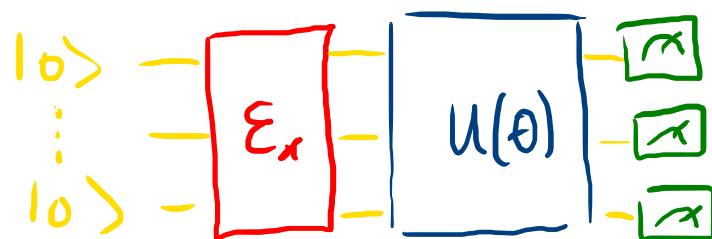
They are expansions in training data!

Π

QUANTUM MODELS

QUANTUM CIRCUITS

- Let's forget about kernels and linear models, and go quantum. Our goal is to train a circuit to do our classification.
- Scheme: apply an encoding unitary E_x , then some unitary $U(\theta)$ we have trained with \mathcal{X}_T , and finally measure.



- The Hilbert space \mathcal{H} of the circuit acts like a feature space:
 $x \mapsto |\phi(x)\rangle = E_x|0\rangle$.



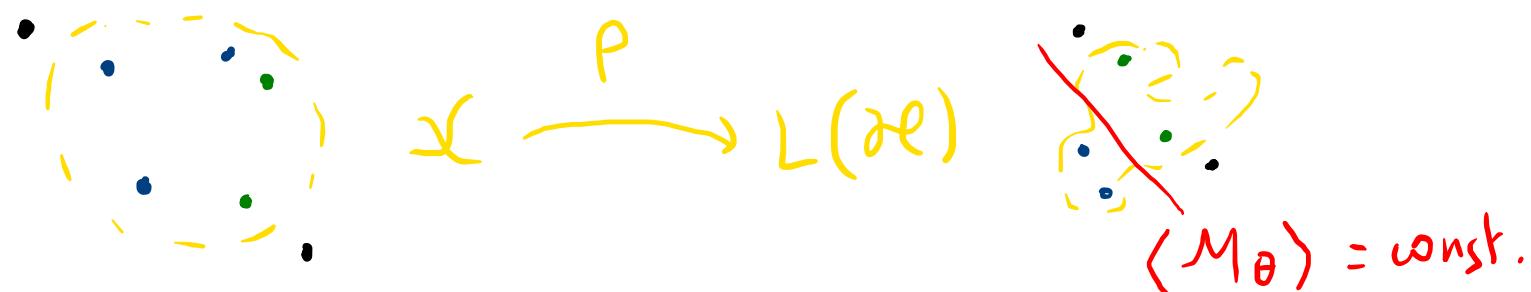
MODELS AND MEASUREMENTS

- What do **models** look like in the "feature space" \mathcal{H} ? Let's massage it until it looks **linear**:

$$\begin{aligned} f(x) &= \langle \phi(x) | u(\theta)^+ u(\theta) | \phi(x) \rangle \\ &= \text{tr} [| \phi(x) \times \phi(x) | u(\theta)^+ u(\theta)] \\ &= \text{tr} [\rho_x M_\theta] \\ &= \langle \rho_x, M_\theta \rangle_{\text{HS}}. \end{aligned}$$

- This is a **linear model**, but not in \mathcal{H} ! Rather, on the **space of linear operators** $L(\mathcal{H})$. This suggests a feature map

$$x \mapsto \rho_x = |\phi(x) \times \phi(x)|.$$



THE QUANTUM KERNEL

- Kernels come from feature maps. A **data-dependent circuit** gives a feature map. What is the "quantum" kernel?

$$\begin{aligned} K(x, x') &= \langle \rho_x, \rho_{x'} \rangle_{HS} \\ &= \text{tr}[\rho_x \rho_{x'}] \\ &= |\langle \phi(x) | \phi(x') \rangle|^2, \end{aligned}$$

i.e. the overlap squared for pure state encoding.

- More generally, the **RKHS** F_{Φ} is the (completed) span of

$$\begin{aligned} \lambda_n &= \kappa(x, \cdot) \\ &= \text{tr}[\rho_x \cdot] \\ &= \langle \cdot \rangle_x. \end{aligned}$$

In other words, F_{Φ} consists of expectations in "data states".

ENCODING AND KERNELS

- Since $K(x, x') = \text{Tr}[\rho_x \rho_{x'}]$, we need to specify the encoding $x \mapsto \rho_x$ to get an actual kernel. The choice depends on \mathcal{X} .
- Some examples:

data	encoding	kernel	cost
$x \in \{0, 1\}^n$	$ x\rangle\langle x $	$\delta_{x, x'}$	$\mathcal{O}(n)$
$x = \sum_i x_i i\rangle \in \mathbb{C}^{2^n}$	$ x\rangle\langle x $	$ x+x' ^2$	$\mathcal{O}(n)$
$x = (x_i) \in [0, 2\pi]^n$	$\bigotimes_i e^{i \frac{x_i}{2} Y} 0\rangle$	$\prod_i \cos(x_i - x'_i) ^2$	$\mathcal{O}(n)$
$x = (x_i) \in \mathbb{R}^n$	$\bigotimes_i e^{x_i (a^\dagger - a)} 0\rangle$	$e^{- x-x' ^2}$	photronics

- A very general class of encodings can be nicely captured in Fourier space, but I won't discuss that now.

QUANTUM REPRESENTERS

- Now we can apply the Representer Theorem to learn what optimal quantum models look like:

$$f = \sum_m \alpha^m \lambda_{x^m} = \sum_m \alpha^m \langle \cdot \rangle_m$$

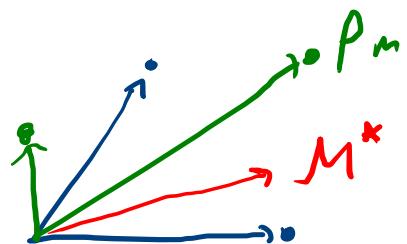
for $x^m \in \mathcal{X}_T$. This is a little abstract.

- Really, this just means optimal measurements are built from density matrices ρ_m of training data:

$$\begin{aligned} f(x) &= \sum_m \alpha^m \langle \rho_x \rangle_m \\ &= \sum_m \alpha^m \text{tr}[\rho_m \rho_x] \\ &= \text{tr} \left[\sum_m \alpha^m \rho_m \rho_x \right] \\ &= \text{tr} \left[M^* \rho_x \right]. \end{aligned}$$

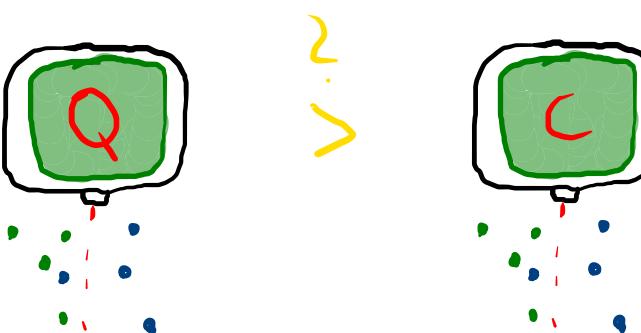
SO WHAT?

- The take-home message is that **classical kernel theory** gives us **optimal quantum models**. This is cool!



$$M^* = \sum_m \alpha^m p_m$$

- This is interesting from a theory perspective, but also **practical**:
 - it tells us what to optimize, i.e. the α^m , and
 - what controls the scaling, i.e. M , the size of \mathcal{X}_+ .
- Question:** does this translate into a quantum speedup?



LOSS AND OPTIMIZATION

- The regularized empirical risk \hat{R} takes the form

$$\hat{R}(f) = \frac{1}{M} \sum_m L[x^m, y^m, f(x^m)] + g(\|f\|).$$

\uparrow loss function \uparrow regularizer

- For $g = \lambda \|f\|^2$ and a candidate optimal quantum model f^* ,

$$\begin{aligned}\|f^*\|^2 &= \sum_{mm'} \alpha_m \alpha_{m'} \operatorname{tr} [\rho_m \rho_{m'}] \\ &= \sum_{mm'} \alpha_m K(x^m, x^{m'}) \alpha_{m'} \\ &= \vec{\alpha}^\top \underbrace{K}_{\text{Gram matrix}} \vec{\alpha}.\end{aligned}$$

- The optimal measurement is then

$$f_{\text{opt}} = \inf_{f^*} \left[\frac{1}{M} \sum_m L + \vec{\alpha}^\top K \vec{\alpha} \right]$$

which is a convex optimization problem for α^m if L is convex, e.g. hinge loss $L = \max(0, 1 - f(x)y)$ for "quantum" SVM.

SPEEDUPS WITH QRAM

- There is a speedup for training qSVMs provided we can efficiently implement quantum Random Access Memory (qRAM):

$$U_{\text{LOAD}}: |j\rangle|0\rangle \mapsto |j\rangle|\text{mem}(j)\rangle$$

address register memory register

- In particular, if the address register refers to training data, and 'loading' takes $\mathcal{O}(n)$ steps, building the Gram matrix $K_{m,m'} = k(x^m, x^{m'})$ takes $\mathcal{O}(M)$ steps. [For $x \in \mathbb{R}^N$, its $\mathcal{O}(MN)$. This is the bottleneck.]
- In contrast, classical SVMs take $\mathcal{O}(M^2)$ steps to train. [A little more carefully, it is $\mathcal{O}(n^2 N + M^3)$.]
- Thus, efficient qRAM, e.g. "Bucket Brigade" model, gives a quadratic speedup for training SVMs!

THANKS!
QUESTIONS?



1. "Supervised QML models are kernel methods" (2021), Schuld.
2. "QML in feature Hilbert spaces" (2018), Schuld & Killoran.
3. "Supervised learning with QC" (2018), Schuld & Petruccione.
4. "Quantum embeddings for ML" (2020), Lloyd et al.
5. "qSUMs for big data classification" (2014), Rebentrost, Mohseni, Lloyd.
6. "qRAM" (2007), Giovannetti, Lloyd & Maccone.