GENERALISATION IN QUANTUM CLASSIFICATION: A QUANTUM-INFO PERSPECTIVE

Leonardo Banchi – University of Florence, Italy

L. Banchi, J. Pereira, S. Pirandola, PRX Quantum 2, 040321 (2021)





«In this world there's two kinds of people, my friend: those with loaded guns and those who dig. You dig.»







«In this world there's two kinds of quantum states, my friend»

- Entangled / separable
- Ground states of critical/not critical models

 $\rho(x)$ c = 1 or c = 2? **c=1** $\rho(x_1)$ $\rho(x_6)$ $\rho(x_5)$ $\rho(x_8)$ $\rho(x_2)$ $\rho(x_7)$ $\rho(x_3)$ $\rho(x_4)$ $\mathbf{C}=\mathbf{2}$



THEORY QUESTION: HOW DIFFICULT IS TO OBTAIN A RELIABLE QUANTUM CLASSIFIER FROM DATA?



MAIN RESULTS

- In **data driven classification** the big question is generalisation: how does the model perform with **new** samples, not present in the training set?
- Main theoretical result (loose): generalisation error / out of sample classification error

$$\approx \sqrt{\frac{2^{I_2(X:Q)}}{T}}$$

"Complexity" quantified by $I_2(X:Q)$, i.e. number of training samples T to get low error with new data

$$\rho(x_1) \qquad \rho(x_6)$$

$$\rho(x_5) \qquad \rho(x_8)$$

 $\rho(x_{2}) \quad \rho(x_{7}) \\ \rho(x_{3}) \quad \rho(x_{4})$



DIFFERENCES WITH QUANTUM STATE DISCRIMINATION

- In quantum state discrimination the states $\rho(x)$ appear with **known** probability distribution P(x) or P(x, c) for classes *c*
 - (Sub)-optimal discrimination strategies are known: Helstrom measurements / pretty-good (square root) measurements
- In machine learning P(x, c) is unknown and must be **learnt** from data
 - In state discrimination we don't have to worry about generalisation. Here we do



EXAMPLE APPLICATIONS



ENTANGLEMENT DETECTION

- Depending on some experimental parameters *x* a device either outputs entangled or separable states *ρ*(*x*)
- The training set is made of very well characterised samples
- What about new samples?

Entangled

Separable



QUANTUM CHANNEL DISCRIMINATION

detect objects xfrom the scattered state of light $\rho(x)$

- Images *x* live in the physical world
- Optimise over the (entangled) input
 POVM

 \mathcal{X} $\rho(x)$ obstacle c = yes/no

• Optimise over the (entangled) input probe state of light and over the detection



QUANTUM BARCODES AND PATTERN RECOGNITION

a)

- Barcode classification must identify each pixel correctly: multiary quantum reading / illumination
- Handwriting classification is easier as errors are tolerated!

error
$$\simeq F(\rho_{\text{black}}, \rho_{\text{white}})^{\text{Hamming}_{4\leftrightarrow 9}}$$

- L. Banchi, Q. Zhuang, S. Pirandola, Phys. Rev. Applied 14, 064026 (2020)
- C Harney, L Banchi, S Pirandola, Phys. Rev. A 103, 052406 (2021)
- JL Pereira, L Banchi, Q Zhuang, S Pirandola, Phys. Rev. A 103, 042614 (2021)









- Classify classical data (e.g. images)
- Embed images x onto a quantum state $x \mapsto \rho(x)$
- Decide the class from a quantum measurement $\{\Pi_{c}\}$
- M Schuld, N Killoran, Phys. rev. lett. 122 (4), 040504, (2019)
- V Havlicek, et al, Nature 567 (7747), 209, (2019)
- S Lloyd, et al, arXiv:2001.03622

QUANTUM EMBEDDINGS

(a) Data distribution





class/label c = "cat"



embedding circuit $x \mapsto \rho(x)$ decision via POVM Π_c

(c) Dilated measurements



optimization of POVM Π_c

optimization of unitary $U_{\mathcal{M}}$



MANY-BODY PHYSICS

Quantum Phase Recognition

• I Cong, S Choi, MD Lukin, Nature Physics 15, 1273 (2019)

- Many-Body Entanglement Measurement from PPT-moments Tr $\left[(\rho_{AB}^{T_B})^n\right]$
 - J Gray, L Banchi, A Bayat, S Bose, Phys. Rev. Lett. 121, 150503, 2018









RESULTS



- We have a training set \mathcal{T} made of T correctly classified states: $\mathcal{T} = \{(\rho(x_t), c_t) \text{ for } t = 1, ..., T\}$
- We can empirically check generalisation using a testing set \mathcal{T}' with T' correctly classified states
- We consider a fixed **quantum embedding** $x \mapsto \rho(x)$ and optimal discrimination via POVM $\{\Pi_c\}$
- What training error / testing error can we expect?

PROBLEM DEFINITION

class 1 $\rho(x_1)$ $\rho(x_6)$ $\rho(x_5)$ $\rho(x_8)$ class 2 $\rho(x_2)$ $\rho(x_7)$ $\rho(x_3)$ $\rho(x_4)$





• Empirical loss / training error

Abstract classification error

$$R(\Pi, \rho) = \mathbb{E}_{(c,x) \sim P(c,x)} \left| \sum_{\substack{c \neq \tilde{c}}} T \right|_{c \neq \tilde{c}}$$

- Real optimal $\Pi^* = \operatorname{argmin}_{\Pi} R(\Pi, \rho)$
- Testing error





$\operatorname{Tr}\left[\Pi_{\tilde{c}}\rho(x)\right] = 1 - \mathbb{E}_{(c,x)\sim P(c,x)}\operatorname{Tr}\left[\Pi_{c}\rho(x)\right]$

• Optimal measurement for empirical risk minimisation $\Pi^{\mathcal{T}} = \operatorname{argmin}_{\Pi} R_{\mathcal{T}}(\Pi, \rho)$





"COMPLEXITY" OF QUANTUM EMBEDDINGS

Classification error \approx Training error

Error

Embedding complexity





VS. DEEP LEARNING









• Bound for the generalisation error

$$G_{\mathcal{T}} \leq 2\sqrt{\frac{\mathscr{B}}{T}} + \sqrt{\frac{2\log(1/\delta)}{T}}$$

• Bound for the approximation error (binary hyp.)

$$\mathscr{A} = \sum_{x} \frac{|P(x|0) - P(x|1)|}{2} - \frac{||\rho_0 - \rho_1|}{2}$$

where
$$\rho_c = \sum P(x \mid c)\rho(x)$$

MAIN RESULT

 $_{1}\|_{1}$



Generalisation bound by combining Rademacher complexity from statistical learning theory, operator inequalities and the following (new) lemma

variable with probability distribution p_i . Then

 $\mathbb{E}_{i \sim p} \left(\|A_i\|_1 \right) \le 7$

where $\mathbb{E}_{i \sim p} f(i) := \sum_{i} p_i f(i)$.

This is an operator version of the generalised mean inequality



Lemma 1. Let A_i be a set of operators and i a random

$$\operatorname{Tr} \sqrt{\mathbb{E}}_{i \sim p} \left(A_i A_i^{\dagger} \right) , \qquad (A12)$$



Approximation bound by combining known results from QIP: Helstrom measurements (binary case), Min-Entropy, Entropy inequalities

Konig, Renner, Schaffner, IEEE Trans. Inf. Th. 55, 4337 (2009). Berta, Seshadreesan, Wilde, J. Math. Phys. 56, 022205 (2015).



Good embeddings should maximise I(C:Q) and minimise $I_2(X:Q)$

Spoiler: Information Bottleneck!

Two extreme cases:

Basis encoding: $x \mapsto \rho(x) = |x\rangle \langle x|$ minimum $\mathscr{A}(\rho) = 0$, maximum $\mathscr{G}_{\mathscr{T}}(\rho)$

Constant embedding : $x \mapsto \rho$ C maximum $\mathscr{A}(\rho)$, minimum $\mathscr{G}_{\mathcal{T}}(\rho) = 0$





RISK/LOSS/ERRORS

• "Single-shot" linear loss, $(c_k, x_k) \in \mathcal{T}$

$\sum_{\substack{c \neq c_k}} \operatorname{Tr} \left[\Pi_c \rho(x_k) \right]$

• Many shots

 $\sum \operatorname{Tr} \left[\Pi_c \rho(x_k)^{\otimes N} \right]$ $c \neq c_k$

in parallel or with different experiments

For large N we find $\mathscr{A}(\rho) \leq KF_{\max}^N$ where $F_{\max} = \max_{\substack{x \neq y}} F(\rho(x), \rho(y))$ **Conjecture:** $\mathscr{B} \simeq \begin{cases} \mathscr{O}(\operatorname{poly}(N)) \\ \mathscr{O}(f^N) & f > 1 \end{cases}$ depending on the data distribution. Related result with classical information $I(X:C_1, ..., C_N) = \begin{cases} O(\log N) \\ O(N) \end{cases}$

Haussler, Opper, Ann. Stat. 1997



NUMERICAL EXPERIMENTS



Bigger Hilbert spaces have lower approximation error, but larger generalisation error

In the numerical example we consider $x \mapsto |\psi(x)\rangle \langle \psi(x)|^{\otimes N}$ $|\psi(x)\rangle = \cos(x)|0\rangle + \sin(x)|1\rangle$





We proved that low dimensional embeddings / Low entropy datasets generalise well!

If $\rho(x)$ "fully scrambles" x in a d-dimensional subspace of the full Hilbert space, then $\mathscr{B} \approx \mathcal{O}(d)$

Geometric characterisation:

$$\mathcal{B}_{c} \leq 1 + \sqrt{(r_{c}^{2} - r_{c})\left(1 - \operatorname{Tr}\left[\rho_{c}^{2}\right]\right)}$$
$$\operatorname{Tr}[\rho_{c}^{2}] = \sum_{x,y} P(x \mid c) P(y \mid c) F(\rho(x), \rho(y))^{2}$$
$$R \leq F(\rho_{c}, \rho_{1})/2.$$



Conjecture from: S Lloyd, et al, arXiv:2001.03622





QUANTUM KERNELS

For pure state embeddings $\rho(x) = |\psi(x)\rangle\langle\psi(x)|$ we find

calculation \mathscr{B} easier for large-dimensional embeddings.

Quantum kernels are used in

- Quantum support vector machines
- Quantum enhanced-feature space

Take home message: avoid $K \propto$ identity (bad generalisation)

- $\mathscr{B} = \left[\mathrm{Tr}\sqrt{K} \right]^2$
- where $K_{xy} = \sqrt{p(x)p(y)} |\langle \psi(x) | \psi(y) \rangle|$ is a (normalised) **kernel** matrix. This makes the





To favour generalisation the **final** Hilbert space must be small, but the initial one can be big!

We may iteratively discard information via pooling layers (e.g. QCNN)

Pooling favours generalisation but harms the accuracy (via data processing)

Take home message: if low training error is achievable with pooling layers, then generalisation can only be better!



INFORMATION BOTTLENECK FOR QUANTUM CLASSIFIERS

 $\rho(x)$ as "bottleneck" that squeezes the relevant information that *x* provides about *c*

IB principle (loss independent): minimise

$$\mathscr{L}_{IB} = I(X:Q) - \beta I(C:Q)$$

Self-consistent solutions (similar for ρ)

 $\tilde{\lambda}_{z} | \psi(z) \rangle = e^{(1-\beta)\log\bar{\rho} + \beta \sum_{c} P(c|z)\log\rho_{c}} | \psi(z) \rangle$

See also: Salek et al. IEEE Trans. Inf. Th. (2018), Datta et al., IEEE-ISIT, (2019)







APPLICATIONS



QUANTUM PHASE RECOGNITION

Task: recognize the phases of matter of a quantum many-body system by taking measurements on the quantum system itself

$$H = -\sum_{i=1}^{L} (\sigma_i^x \sigma_{i+1}^x + h \sigma_i^z), \qquad \Longrightarrow$$

Ordered (|h| < 1) / disordered (|h| > 1) phases

T : number of training samples per class S : number of measurement shots



VARIATIONAL QUANTUM INFORMATION BOTTLENECK



Single-qubit "Data Re-uploading" Classifier

 $|\psi_w(x)\rangle = \prod [R^z (w^{z\ell} \cdot x + w_0^{z\ell}) R^y (w^{y\ell} \cdot x + w_0^{y\ell})]|0\rangle,$

For two-qubit states the "re-upoading" embedding can be trained with an efficient variational minimisation of the IB Lagrangian



SUMMARY



Main open questions: deep learning regime / asymptotics with multiple shots

