# scientific reports

Check for updates

OPEN

# Robust estimation of the intrinsic dimension of data sets with quantum cognition machine learning

Luca Candelori[1,2✉], Alexander G. Abanov[3,8], Jeffrey Berger[1,8], Cameron J. Hogan[4,8], Vahagn Kirakosyan[1,8], Kharen Musaelian[1,8], Ryan Samson[1,8], James E. T. Smith[1,8], Dario Villani[1,7,8], Martin T. Wells[4,8] & Mengjia Xu[5,6,8]

We propose a new data representation method based on Quantum Cognition Machine Learning and apply it to manifold learning, specifically to the estimation of intrinsic dimension of data sets. The idea is to learn a representation of each data point as a quantum state, encoding both local properties of the point as well as its relation with the entire data. Inspired by ideas from quantum geometry, we then construct from the quantum states a point cloud equipped with a quantum metric. The metric exhibits a spectral gap whose location corresponds to the intrinsic dimension of the data. The proposed estimator is based on the detection of this spectral gap. When tested on synthetic manifold benchmarks, our estimates are shown to be robust with respect to the introduction of point-wise Gaussian noise. This is in contrast to current state-of-the-art estimators, which tend to attribute artificial "shadow dimensions" to noise artifacts, leading to overestimates. This is a significant advantage when dealing with real data sets, which are inevitably affected by unknown levels of noise. We show the applicability and robustness of our method on real data, by testing it on the ISOMAP face database, MNIST, and the Wisconsin Breast Cancer Dataset.

When data is characterized by a large number of features (e.g., zip code, annual income, age, credit card spend, etc. for borrowers; cholesterol, blood pressure, BMI, etc. for patients; or the latent and dependent variables), it tends to lie on a surface that has a smaller dimensionality than the full feature space[1]. Finding this low-dimensional surface is often referred to as *manifold learning*. The lower dimensionality reflects the underlying latent structures in the data, correlations, and a variety of nonlinear relationships[2,3]. Furthermore, data points whose feature vectors are close together should possess similar properties related to the nature of the data. For example, in a regression problem, the output/target variables are expected to depend smoothly on the input variables. These characteristics of real data suggest that any given dataset consisting of $D$ features lies entirely on a smooth manifold $M \subseteq \mathbb{R}^D$ of dimension $d$ (the manifold hypothesis[4]), where $d$ is much smaller than the total number of features $D$, typically $d < 100$. The dimension $d$ of the manifold is referred to as the *intrinsic dimension* of the data[1]. This number represents the minimal number of parameters required to characterize the data. Knowledge of the intrinsic dimension $d$ can be used to effectively choose a target space for dimension-reduction models (such as PCA, Isomap, t-SNE, etc.) or to compress deep neural networks while maintaining the performance[5]. Intrinsic dimension estimation is also widely used in network analysis[6,7], complex materials[8] and health sciences[9].

One of the main challenges for manifold learning is the inevitable presence of noise in real data. A typical "global" projective approach is to impose a functional form (e.g. PCA where the manifold is assumed to be linear) and to assume that the error between the manifold approximation and the actual data is noise, which is then analyzed. Imposing a functional relation immediately gives estimates for the intrinsic dimension, which

[1]Qognitive, Inc., Miami Beach, FL 33139, USA. [2]Department of Mathematics, Wayne State University, Detroit, MI 48202, USA. [3]Department of Physics and Astronomy, Stony Brook University, Stony Brook, NY 11790, USA. [4]Department of Statistics and Data Science, Cornell University, Ithaca, NY 14853, USA. [5]Department of Data Science, New Jersey Institute of Technology, Newark, NJ 07102, USA. [6]Center for Brains, Minds and Machines, Massachusetts Institute of Technology, Cambridge, MA 02139, USA. [7]Department of Mathematics, King's College London, WC2R 2LS London, UK. [8]These authors contributed equally: Alexander G. Abanov, Jeffrey Berger, Cameron J. Hogan, Vahagn Kirakosyan, Kharen Musaelian, Ryan Samson, James E. T. Smith, Dario Villani, Martin T. Wells and Mengjia Xu. ✉email: luca.candelori@qognitive.io; candelori@wayne.edu

tend to be robust to the introduction of additional noise. However, when the data manifold $M$ has a lot of curvature, linear methods will fail. The problem can be somewhat alleviated by nearest neighbors methods that sample locally around each data point, assuming that at a sufficiently small scale all manifolds are close to being linear[10–12]. Indeed, all current state-of-the-art intrinsic dimension estimators (some of which we describe below) are "local", producing estimates that are based on a local sampling around each data point[4]. Such techniques are designed and benchmarked against highly dimensional and highly curved manifolds. While they perform reasonably well in this ideal setup, they often tend to fall apart when noise is re-introduced into the data. Indeed, local methods cannot distinguish *shadow* dimensions that are transversal to the data manifold, and that are only artifacts created by the noise, leading to significant overestimates of intrinsic dimension.

In this paper, we propose a new data representation and manifold learning technique based on Quantum Cognition Machine Learning (QCML)[13] and quantum geometry[14–16]. The main idea is to create a (non-commutative) quantum model for the data manifold itself, from which we can estimate important geometric features, such as intrinsic dimension. Picking a quantum model is similar to what is done in linear methods, in the sense that a functional relation is imposed on the data. But in contrast to linear methods, we *learn* the model from the data, and we make no assumptions about the underlying distribution. Our method gives local estimates of the intrinsic dimension at every data point, but also takes into account the global geometry of the data manifold $M$. To this end, we are able to develop a manifold approximation method that is both robust to noise and flexible enough to capture non-linear geometric features of the data manifold.

Current state-of-the-art intrinsic dimension estimators measure statistics related to the density of nearest neighbors lying within a certain radius $r$ from a data point $x$, and express these statistics as functions of intrinsic dimension (CorrInt[17], MLE[18], DANCo[19], TwoNN[20]). These methods do not make any linearity assumption about the data, but do require the data to be dense in small patches around any given point. As is well-known, this requirement is fundamentally incompatible with the *curse of dimensionality*[21,22], which usually occurs in dimensions when $d$ is greater than the logarithm of the sample size[23], and indeed these methods tend to underestimate the intrinsic dimension when $d$ is large relative to the sample size. The overestimation effect induced by noise combined with the underestimation effect induced by the curse of dimensionality often results in unreliable intrinsic dimension estimates.

Our proposed approach relies on the manifold hypothesis but differs from the current projective and nearest neighbors methods. Our method produces local intrinsic dimension estimates that are not based on neighborhood sampling of the data but include global information from the entire data set. Indeed, our method first learns a model for the entire data manifold $M$, as a semi-classical limit of a quantum matrix configuration (in the sense of quantum geometry[14–16]). In particular, given a data set $X$ containing $D$ features, we train $D$ quantum observables $A = \{A_1, \ldots, A_D\}$ (i.e. a *matrix configuration*) as it is done in QCML[13]. The whole matrix configuration $A$ can be used to map each data point $x \in X$ to a *quasi-coherent* quantum state $\psi_0(x)$, which is then mapped back into the data space, producing a *point cloud* approximation $X_A$ to the actual data manifold $M$. Each element of the point cloud represents the expected position in the feature space of its corresponding data point, and it comes with a "cloud" of uncertainty around its actual position whose shape is determined by the quantum fluctuations of the matrix configuration. The point $x$ is further equipped with a *quantum metric* $g(x)$, which is a $D \times D$ real symmetric positive semi-definite matrix. This metric, already considered by physicists[14,15], encodes much of the local geometry of the data manifold; it can be shown that its rank in particular is approximately equal to the intrinsic dimension of $M$, and that its non-zero eigenvalues are all close to 1. Therefore, intrinsic dimension estimates can be given by detecting the spectral gap of the quantum metric, separating the zero eigenvalues from the non-zero eigenvalues that are close to 1.

Some of the existing estimators (so-called 'projective methods', such as PCA) also provide explicit embeddings of the data into $d$-dimensional space, where $d$ is the intrinsic dimension. In this sense, our method can be considered 'locally' projective: at each point $x \in X_A$ of the point cloud, the eigenvectors of the quantum metric $g(x)$ with eigenvalues close to one correspond to the directions that are tangent to the data manifold, therefore providing a set of $d$ local coordinates for the manifold.

We test our intrinsic dimension estimator on both synthetic and real data sets, following the benchmarking framework proposed in ref.[4] and implemented in the `scikit-dimension` Python package[23]. In addition to this standard framework, we stress-test our estimates by introducing increasing levels of Gaussian noise into the data, and compare the results with other state-of-the-art techniques. In all of our testing, higher levels of noise increasingly degrade the quality of the point cloud approximation $X_A$, and the spectral gap detection in the quantum metric becomes increasingly difficult. However, they do not qualitatively alter the intrinsic dimension estimation. This stands in marked contrast to other intrinsic dimension estimators that we tested, whose estimates are highly sensitive to even small amounts of noise.

## Results
### Quantum geometry in data analysis
Consider a $t \times D$ data set $X$ containing $t$ data points $x_1, \ldots, x_t$, where each data point $x_i$ consists of a $D$-dimensional real-valued vector of data features $x_i = (a_i^1, \ldots, a_i^D)$. We do not assume any particular ordering of the features, so that, for example, a digital image would be encoded as a flattened one-dimensional array of pixel values. We allow categorical data features, provided an appropriate embedding (e.g. one-hot encoding or target encoding) has been chosen, although in this article we will only consider numerical features. We assume that $X$ lies entirely on a smooth manifold $M$, called the *data manifold*, of intrinsic dimension $d < D$. We further assume that the $D$ features of the data extend to smooth functions $a^k \in C^\infty(M)$, for $k = 1, \ldots, D$, giving the coordinates of an embedding $(a^1, \ldots, a^D) : M \hookrightarrow \mathbb{R}^D$ of the data manifold into $D$-dimensional Euclidean space. In quantum geometry, the commutative algebra $C^\infty(M)$ of smooth functions on a manifold is replaced by the non-commutative algebra of Hermitian operators on a $N$-dimensional Hilbert space[24,25]. The

choice of dimension $N$ is arbitrary, and independent of $D$. Typically, smaller choices of $N$ will result in quantum geometries exhibiting strong quantum effects, while the limit as $N$ goes to infinity can be considered 'classical' (i.e. not quantum)[14].

For the purposes of this work, any set $A = \{A_1, \ldots, A_D\}$ consisting of $D$ Hermitian matrices on a $N$-dimensional Hilbert space is called a *matrix configuration*, and can be viewed as a non-commutative avatar of the set of $D$ coordinate functions $a^k$ on a manifold $M \hookrightarrow \mathbb{R}^D$. Typically in physics, the matrix configuration $A$ is given by a quantum theory and the goal is to construct a symplectic manifold $M \hookrightarrow \mathbb{R}^D$, so that $A$ represents a quantization of the coordinate functions $x^k$ giving the embedding; that is, a compatibility between the Poisson bracket on $M$ and the commutator bracket on $A$ is required, among other conditions.

In the context of data analysis, the situation is reversed: $M$ is given by the data manifold, and we instead propose to learn a suitable matrix configuration $A$, reflecting as much of the geometry of $M$ as possible. We do so through the formalism of *quasi-coherent states*[14,15]. Recall that in quantum mechanics a *state* is a vector of unit norm in a Hilbert space, and is represented in bra-ket notation by a ket $|\psi\rangle$. The inner product of two states $|\psi_1\rangle, |\psi_2\rangle$ is represented by a bra-ket $\langle\psi_1|\psi_2\rangle$. The *expectation value* of a Hermitian operator $A$ on a state $|\psi\rangle$ is denoted by $\langle\psi|A|\psi\rangle = \langle A\psi|\psi\rangle = \langle\psi|A\psi\rangle$, representing the expected outcome of the measurement corresponding to $A$ on the state $|\psi\rangle$. For any state $|\psi\rangle$ in $N$-dimensional Hilbert space and an $N \times N$ matrix configuration $A = \{A_1, \ldots, A_D\}$, define the state's *position* vector by

$$A(\psi) = (\langle\psi|A_1|\psi\rangle, \ldots, \langle\psi|A_D|\psi\rangle) \in \mathbb{R}^D$$

and the state's *variance* (or *quantum fluctuation*) $\sigma^2(\psi)$ by

$$\sigma_k^2(\psi) = \langle\psi|A_k^2|\psi\rangle - \langle\psi|A_k|\psi\rangle^2, \quad \sigma^2(\psi) = \sum_{k=1}^{D} \sigma_k^2(\psi) \in \mathbb{R}.$$

Intuitively, the matrix configuration $A$ assigns to each quantum state $|\psi\rangle$ a point $A(\psi)$ in Euclidean space $\mathbb{R}^D$, together with a "cloud" around it representing the uncertainty of the measurement of the point's position in space. In this context, $A(\psi)$ represents the center of the cloud, while $\sigma(\psi)$ is a statistical measure of the cloud's dispersion.

Now for any data point $x = (a_k) \in \mathbb{R}^D$, we want to construct a quantum state $\psi_0(x)$ reflecting not only the absolute position of $x$ within feature space, but also its relation to all the other points in the data set $X$. To do so, consider the *error Hamiltonian*

$$H(x) = \frac{1}{2}\sum_{k=1}^{D}(A_k - a_k \cdot I_N)^2, \tag{1}$$

where $I_N$ denotes the $N \times N$ identity matrix. Note that the error Hamiltonian is a positive semi-definite Hermitian operator. We will assume throughout the article non-degeneracy of the Hamiltonian, and we list the eigenvalues of the Hamiltonian $0 \leq E_0(x) < \cdots < E_{N-1}(x)$ in increasing order. For the present purposes, degeneracies of $H(x)$ do not play a role. We also let $|\psi_n(x)\rangle$, $n = 0, \ldots, N-1$, be corresponding choices of normalized eigenvectors, or *eigenstates*. By the non-degeneracy assumption, all the eigenstates are uniquely defined up to multiplication by a phase factor $e^{i\theta}, \theta \in \mathbb{R}$. For each $x$, an eigenstate $|\psi_0(x)\rangle$ associated to the lowest eigenvalue of $H(x)$ is called a *quasi-coherent state* of $x$. A simple calculation shows that

$$E_0(x) = \frac{1}{2}\|A(\psi_0(x)) - x\|^2 + \frac{1}{2}\sigma^2(\psi_0(x)), \tag{2}$$

so that the lowest eigenvalue (i.e. the *ground state energy*) of the error Hamiltonian can be broken down into two contributions: the squared distance between $x$ and the position of its corresponding quasi-coherent state, and the quantum fluctuation of the quasi-coherent state itself. This is analogous to the bias-variance breakdown of the mean-squared error loss function. We can now train a matrix configuration $A$ so as to minimize the combined loss function (2) for all data points $x \in X$. In this way, the matrix configuration captures global features of the data, which are then reflected into the ground state $\psi_0(x)$, for each $x \in X$.

From the trained matrix configuration $A$, we may then calculate the *point cloud*

$$X_A = \{A(\psi_0(x)) : x \in X\} \subseteq \mathbb{R}^D, \tag{3}$$

which can be viewed as an approximate sampling of the data manifold $M$. The original data points $x \in X$ may contain noise, missing features, or otherwise deviate substantially from the idealized underlying data manifold $M$. By choosing an appropriate matrix configuration $A$, capturing enough global information about the data, the set $X_A$ turns out to be much closer to $M$ than the original data set $X$. Key geometric features of the data manifold, such as the intrinsic dimension $d$, can be recovered from $X_A$ in a way that is robust to noise and other artifacts.

## Quantum cognition machine learning

Training a matrix configuration $A$ on a data set $X$ is the optimization problem at the core of Quantum Cognition Machine Learning (QCML)[13]. QCML has been developed independently of quantum geometry, and this is the first work pointing out the relation between these two areas of study. In the original formulation of QCML, a matrix configuration $A$ is trained so as to minimize the aggregate energy loss function (2) across all data. In the present context, minimizing energy sometimes has the undesired effect of training $A$ so that the aggregate quantum fluctuation $\sum_{x \in X} \sigma^2(\psi(x))$ goes to zero, forcing all the matrices $A_1, \ldots, A_D$ in the matrix configuration to commute. A commutative matrix configuration is highly undesirable. It produces a point cloud approximation $X_A$ consisting of $N$ points, corresponding to the positions of the $N$ common eigenstates of the matrix configuration, with no point cloud around them. Indeed, it can be shown that $X_A$ in this case consists of a $N$-means clustering of the data set $X$, and is therefore entirely classical[26].

Instead, in this work we train the matrix configuration $A = \{A_1, \ldots, A_D\}$ on the data set $X$ by minimizing the mean squared distance between the data set $X$ and the point cloud $X_A$, i.e. by finding

$$A = \mathrm{argmin}_{B=\{B_1,\ldots,B_D\}} \left( \sum_{x \in X} \|B(\psi_0(x)) - x\|^2 \right), \tag{4}$$

where the minimum is taken over the space of all $D$-tuples of $N \times N$ Hermitian matrices. The optimization (4) can be tackled efficiently using gradient descent methods, similar to those employed in state-of-the-art machine learning models. In our study, we find $A$ by implementing the optimization problem as a custom layer in PyTorch [27]. To ensure convergence to a meaningful matrix configuration, it is often beneficial to pre-process the data so that it has a homogeneous scale, for example by standardizing all the features so that they have mean $\mu = 0$ and standard deviation $\sigma = 1$.

Note that the choice of loss function in (4) corresponds to the "squared-bias" term in the bias-variance decomposition of the energy functional $E_0(x)$ in (2). We do not minimize the quantum fluctuation, or "variance" term. Indeed, while the bias term is in general unbounded, the quantum fluctuation $\sigma^2(x)$ has a simple bound in terms of the matrix configuration $A$ only (i.e. independent of $x$), given by

$$\sigma^2(x) \leq \sum_{k=1}^{D} (\mu_k - m_k)^2 \leq \frac{D}{4}(\mu - m)^2,$$

where $\mu_k$ (resp. $m_k$) is the highest (resp. lowest) eigenvalue of $A_k$ and $\mu = \max_k \mu_k$ (resp. $m = \min_k m_k$). This bound has an elementary proof similar to Popoviciu's inequality[28] on variances. Note that the eigenvalues of $A_k$ correspond to possible measurement outcomes of the $k$-th coordinate of the position of a point $x$. Therefore, if we train $A$ so that the positions $X_A$ are close to a compact data set $X$, we expect the quantum fluctuation to be commensurate with the average noise level in the data $X$. This is indeed what we observe in practice.

It is also possible to modify the loss function in (4) by adding back the quantum fluctuation term with a weight $w \in \mathbb{R}_{\geq 0}$, a tunable hyperparameter,

$$A = \mathrm{argmin}_{B=\{B_1,\ldots,B_D\}} \left( \sum_{x \in X} \|B(\psi_0(x)) - x\|^2 + w \cdot \sigma^2(x) \right). \tag{5}$$

In this way, the choice $w = 0$ recovers the bias-only loss function (4) while $w = 1$ corresponds to the original energy loss (2). In applications, small non-zero values of $w$ may lead to more robust point cloud approximations $X_A$, especially in the presence of significant amounts of noise.

It is also possible in principle to replace the error Hamiltonian (1) with the Dirac operator defined in ref.[15]. The advantage of using the Dirac operator is that the energy loss is allowed to reach zero without the matrix configuration $A$ being necessarily commutative. Equivalently, the quasi-coherent states in this case are zero modes. However, the Hilbert space dimension required by the Dirac operator scales exponentially in the number of features $D$, and this is not practical when dealing with data sets containing a large number of features.

## Intrinsic dimension estimation

Suppose now that a matrix configuration $A$ has been trained from a data set $X$ as in (4), so that the data manifold $M$, by construction, lies within a region of $\mathbb{R}^D$ where the energy functional $E_0(x)$ is near-minimal and it has minimal variation (assuming that the quantum fluctuation term in (2) is not too large). We may then apply the technique described in ref.[15] to calculate the intrinsic dimension of $M$. In particular, from formula (2), we see that as $x$ moves away from the manifold $M$ then the energy $E_0(x)$ increases like the squared distance from $x$ to $M$, while in the directions tangent to $M$ the energy is approximately constant. This means that the Hessian matrix

of the energy functional at $x$ should exhibit a clear spectral gap between the lowest $d = \dim M$ eigenvalues, corresponding to the directions tangent to $M$ and near zero, and the highest $D - d$ eigenvalues, of order one and corresponding to the directions that point away from $M$. Detecting the exact location of the spectral gap is, therefore, equivalent to estimating the intrinsic dimension of $M$.

This observation can be turned into an algorithm to estimate the intrinsic dimension. First, the Hessian matrix of the energy functional can be computed in terms of the matrix configuration $A$, using perturbation theory. Its entries are given by the formula

$$\frac{\partial^2 E_0}{\partial x_\mu \partial x_\nu} = \delta_{\mu\nu} - 2 \sum_{n=1}^{N-1} \mathrm{Re} \frac{\langle \psi_0(x)|A_\mu|\psi_n(x)\rangle \langle \psi_n(x)|A_\nu|\psi_0(x)\rangle}{E_n(x) - E_0(x)}, \quad \mu, \nu = 1, \ldots, D \tag{6}$$

where, as before, we write $\psi_n(x)$ and $E_n(x)$ for the eigenstates and energies of the error Hamiltonian $H(x)$ given by (1). Notice that (6) is exact, despite being derived using perturbation theory. In detecting the spectral gap, it is more convenient to consider the second term of (6) only, a real symmetric $D \times D$ matrix $g(x)$ whose entries are given by

$$g_{\mu\nu}(x) = 2 \sum_{n=1}^{N-1} \mathrm{Re} \frac{\langle \psi_0(x)|A_\mu|\psi_n(x)\rangle \langle \psi_n(x)|A_\nu|\psi_0(x)\rangle}{E_n(x) - E_0(x)}, \quad \mu, \nu = 1, \ldots, D. \tag{7}$$

It can be easily shown that the matrix $g(x)$ is positive semi-definite, and in the context of matrix geometry it is called the *quantum metric*[15,16,25]. Indeed, in our context it can be viewed as an approximate Riemannian metric on the data manifold $M$. For a point $x$ belonging to the point cloud $X_A$, the eigenvalues of $g(x)$ tend to be either close to one or close to zero, with a spectral gap occurring between the highest $d$ and the lowest $D - d$ eigenvalues. The eigenvectors corresponding to the the highest $d$ eigenvalues will point in the directions tangential to the data manifold $M$, with the remaining eigenvectors being transversal to the data manifold. In this way, an examination of the spectral gap at $g(x)$ provides an estimate of the intrinsic dimension $d = \dim_x M$. We could, in principle, apply this procedure to estimate the intrinsic dimension at points $x \in X$ directly, bypassing the point cloud. However, as noted in ref.[15], much clearer spectral gaps emerge in practice when calculating the quantum metric on the point cloud $X_A$. This is because $X_A$, as noted earlier, is much more robust to noise and to small perturbations of the data manifold.

The estimation of intrinsic dimension from the point cloud $X_A$ is based on the assumption that the matrix configuration has been trained well enough so that the point cloud forms a good approximate model for the data manifold $M$, in particular so that the intrinsic dimensions of the data set $X$ and $X_A$ are equal. Since the matrix configuration $A$ is trained in such a way as to minimize the squared distance between $X$ and $X_A$, it is reasonable to assume that this is the case. However, the quality of this approximation will depend on many factors, mainly the choice of quantum fluctuation control $w$ in the loss function (5) and the choice of the Hilbert space dimension $N$. The effect of these choices on $X_A$ is discussed in more detail in the supplementary material for this article (see Figures 1 and 2 of the Supplementary Material).

The algorithm for estimating intrinsic dimension can be summarized as follows.

---

**Data:** Data set $X \subseteq \mathbb{R}^D$ lying on a data manifold $M \subseteq \mathbb{R}^D$
**Result:** A list `dlist` of local intrinsic dimension estimates $d_x \approx \dim_x M$

1   Train a matrix configuration $A = \{A_1, \ldots, A_D\}$ on $X$ as in (4) or (5);
2   `dlist` $\leftarrow \emptyset$ ;
3   **for** $x \in X$ **do**
4      calculate the ground state $|\psi_0(x)\rangle$ of the error Hamiltonian $H(x)$ ;
5      calculate the position $y = A(\psi_0(x)) \in X_A$ ;
6      calculate the spectrum $e_0 \leq \ldots \leq e_{D-1}$ of the quantum metric $g(y)$ ;
7      Calculate $\gamma = \mathrm{argmax}_{i=1,\ldots,D}\, e_i/e_{i-1}$, the largest spectral gap ;
8      Append $d = D - \gamma$ to `dlist`
9   **end**
10   Return `dlist`

---

**Algorithm 1.** Quantum Cognition Machine Learning intrinsic dimension estimator.

The Algorithm 1 returns a list of intrinsic dimension estimates for every point $x \in X$. To extract a global estimate, a variety of techniques can be employed, such as taking the mode, median, or geometric mean to more refined $k$-nearest neighbor techniques. A global estimate can also be easily adapted to the case where multiple connected components of $M$ are detected, each of possibly different dimensions. Note that in steps 7-8 of Algorithm 1, we calculate the largest spectral gap by comparing successive ratios of eigenvalues. With this approach, the results $d = 0, D$ cannot be detected. We are indeed assuming through the article that the data

manifold does not have zero-dimensional or codimesion zero connected components. In practice, if the data set is zero-dimensional (i.e. a sparse set of points) then the eigenvalues of the quantum metric will all be zero at every point, which can be easily checked. Similarly, if the data set is $D$-dimensional then the eigenvalues will all be close to 1 at every point, which is also easy to check.

It is possible to replace these crude spectral gap estimates with more advanced methods. For example, if $D$ is large and $d \ll D$, as is typical in real data sets, methods based on random matrix theory[29] are likely to give more robust estimates. Phase transitions in random matrix theory (RMT) refer to abrupt changes in the behavior of eigenvalues of large random matrices as certain parameters are varied. These transitions are particularly interesting because they often separate different regimes of matrix behavior. The eigenvalues of large random matrices follow well-defined distributions (like the Marchenko-Pastur distribution[30]) and as the matrix size grows, eigenvalue behavior exhibits certain regularities, with interesting gaps between signal and noise eigenvalues. There is often a critical threshold phase transition at which the behavior of the eigenvalues changes sharply. The presence of spectral gaps between eigenvalues can signal the existence of a significant phase transition and in high-dimensional problems, RMT can predict the existence of these gaps. Furthermore, the eigenvectors associated with eigenvalues that exhibit an eigen-gap will be informative and uninformative when the eigen-gap vanishes[31].

One approach to recover the true signal matrix is to threshold the singular values of the quantum metric $g(y)$ and keep the singular values that are likely to correspond to the signal and discard those that are likely to be noise[29]. This leads to a singular value thresholding rule, where a threshold $\tau$ is applied to the singular values of the observable matrix, and only the singular values larger than $\tau$ are retained. It was shown that in the asymptotic limit as $t, D \to \infty$ with $t/D \to \gamma$, the optimal threshold is $\tau_{opt} = \frac{4}{\sqrt{3}} \cdot \sigma$ where $\sigma$ is the standard deviation of an underlying Gaussian noise matrix[29]. The noise parameter $\sigma$ can be estimated by $\hat{\sigma}$ using the Marchenko-Pastur bulk singular values. This estimate can then be used to adaptively set the threshold for singular value thresholding. Specifically, the rule $\hat{\tau}_{opt} = \frac{4}{\sqrt{3}} \cdot \hat{\sigma}$ can be applied to the singular values of the quantum metric $g(y)$ for hard thresholding to find the spectral gap. In the following, we will refer to this thresholding method as the "RMT-based" estimate.

The choice of dimension $N$ of the Hilbert space underlying the matrix configuration $A$ is a hyperparameter of the algorithm. As shown in[15,25], we have the rank bound

$$\text{rank } g \leq 2(N-1), \tag{8}$$

so that $N$ should be chosen large enough to ensure $2(N-1) > d$. Since a priori we only know that $d < D$, a sensible choice would be to set $N \geq D/2 + 1$. However, for datasets with a large number of features (i.e. $D$ large), this choice might be impractical, since the number of parameters of a QCML estimator scales quadratically in $N$. Instead, a simple strategy for choosing $N$ that we employ in large real data sets is to first pick $N$ small and gradually increase it until a clear spectral gap emerges and is consistent across different choices of dimensions. In general, larger Hilbert space dimension $N$ will result in point clouds $X_A$ that are closer to the original data $X$ (low bias) but may also model noise artifacts (high variance). A smaller $N$ will result in approximations that may have higher loss/higher energy (high bias) but that may be more robust with respect to noise (low variance).

## Benchmarks

### The fuzzy sphere

We first evaluate Algorithm 1 in the case when the data $X$ is a sample of $T = 2500$ uniformly distributed points on the unit sphere $M = S^2$, embedded in $D = 3$ dimensions. We allow the data to be "noisy", that is, any given data point $x \in X$ might not necessarily lie on $M$ but it could be drawn from a Gaussian distribution whose mean is on $M$ and whose standard deviation is a `noise` parameter. By the rank bound in (8) on the quantum metric, the minimum possible choice of Hilbert space dimension is $N = 3$. Plots of the point cloud $X_A$ and the spectra of the quantum metric $g(x)$ at different points $x \in X_A$ are shown in Figure 1. With zero noise (Figure 1 a-b) the point cloud approximation $X_A$ is very close to the original unit sphere and a clear spectral gap emerges at every point between the top two eigenvalues of the quantum metric and the lowest eigenvalue. The intrinsic dimension estimate is thus $d = 2$ at all points. As the noise level increases, up to `noise` $= 0.2$ (Figure 1 c-d) the point cloud starts picking up some noise artifacts and the variance of the metric spectrum increases. However, even for `noise` $= 0.2$, the intrinsic dimension estimate is $d = 2$ for all 2500 points.

For comparison, we selected some of the best-performing state-of-the-art algorithms for intrinsic dimension estimation (DANCo, MLE, CorrInt, TwoNN, as implemented in ref.[23] ) and tested them at different levels of noise, and for three different data set sizes $T = 250, 2500, 25000$ (Figure 2). In Figure 2, the slope of the intrinsic dimension estimate for the QCML model is zero, so that the estimate is unaffected by noise in the range `noise` $\in [0, 0.2]$. The dimension estimate is also consistent across different number of samples $T = 250, 2500, 25000$, indicating additional robustness with respect to data distribution and density. By comparison, the estimates of all other baseline algorithms quickly converge to $d = 3$, creating a "shadow" dimension out of the noise. Increasing the size of the sample does not seem to aid the state-of-the-art algorithms in detecting the correct intrinsic dimension. In fact, the slopes of the "shadow dimension" graphs in Figure 2 get noticeably steeper for $T = 25000$ samples, indicating an even faster degradation of the intrinsic dimension estimate as the data density increases.
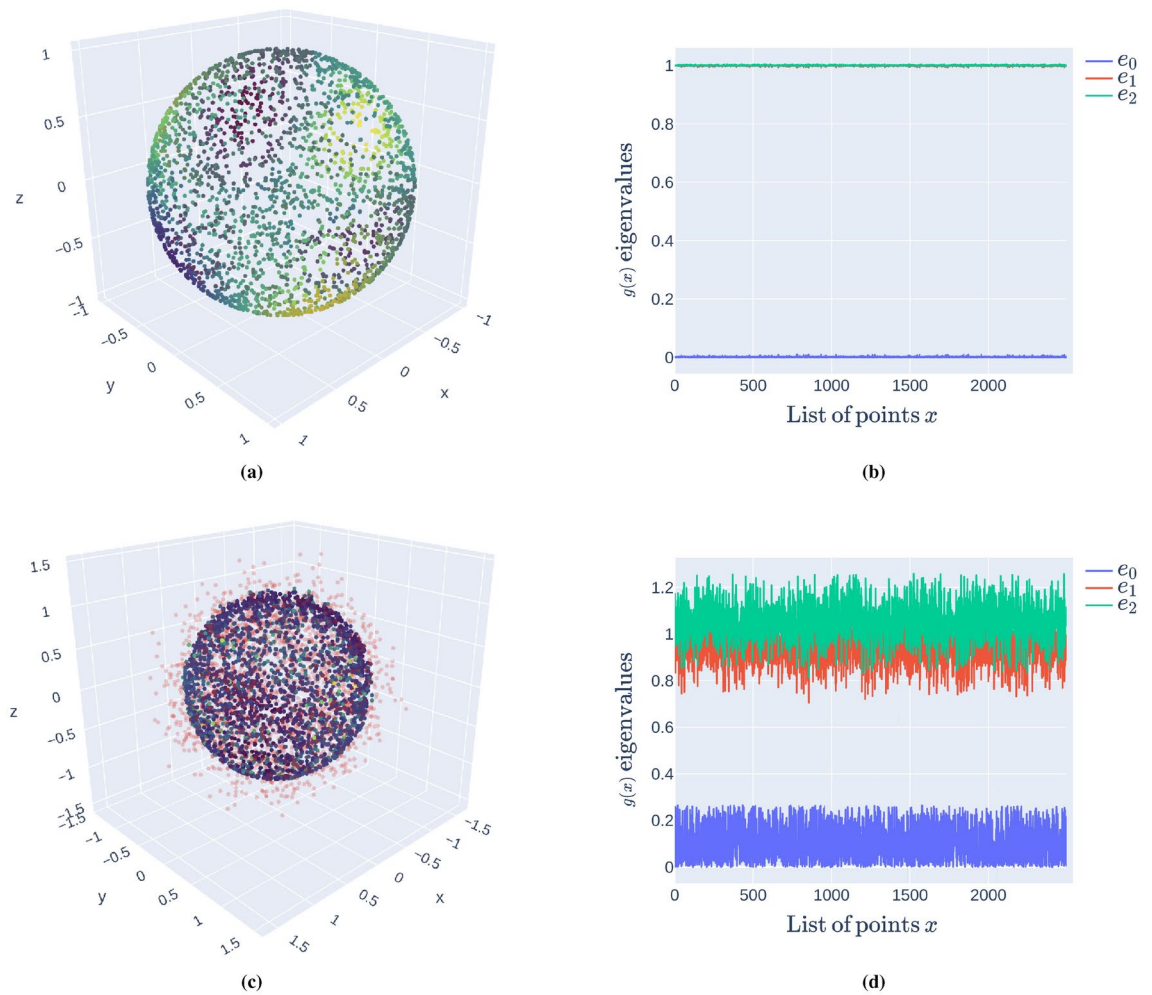
**Fig. 1**. Two configurations are shown for a data set $X$ consisting of $T = 2500$ points uniformly distributed on the unit sphere with different levels of noise. (**a,c**) Scatter plot of the point cloud $X_A$ for (**a**) `noise = 0`, and (**c**) `noise = 0.2`, for two corresponding matrix configurations $A$ trained with Hilbert space dimension $N = 3$. The original dataset is overlayed in red. Darker points correspond to lower relative error energy $E_0$. (**b,d**) Spectral gaps for (**b**) `noise = 0` and (**d**) `noise = 0.2`. The $x$-axis corresponds to points $x \in X_A$ and on the $y$-axis the eigenvalues of the quantum metric $g(x)$ are plotted.



**Fig. 2**. Intrinsic dimension estimates for the unit sphere $S^2$ as a function of `noise` level. Varying data set sizes of (**a**) $T = 250$, (**b**) $T = 2500$, (**c**) $T = 25000$ points are tested. For the QCML estimator, the average estimate across all $T$ points is plotted. The QCML estimate is robust to noise when compared to other methods.

It is worth noting that the optimal matrix configurations $A = \{A_1, A_2, A_3\}$ obtained by the QCML estimator in this case are well-known to physicists as "fuzzy spheres"[14–16]. Up to a change of basis and a re-scaling factor, the elements of $A$ are given by the angular momentum operators in quantum mechanics.

*Higher-dimensional synthetic manifolds*
Next, we test the QCML intrinsic dimension estimator on three higher-dimensional manifolds included in the benchmarking framework of ref.[4]. The first is the 17-dimensional standard hypercube embedded into $D = 18$ dimensions (Figure 3 (a), (d) ), and labeled $M_{10b}$ in the scikit-dimension Python package. The second is the 10-dimensional manifold $M_\beta$ (Figure 3 (b), (e) ), embedded in $D = 40$ dimensions, and the third is the 18-dimensional manifold $MN_1$ (Figure 3 (c), (f) ) embedded into $D = 72$ dimensions. These benchmarks are considered among the most difficult for intrinsic dimension estimation, due to both the non-uniform density of the data (for $M_{10b}$ and $M_\beta$) and the non-linearity of the embedding (for $M_\beta$ and $MN_1$). In our testing, we trained the QCML estimator with Hilbert space dimension $N = 16$ on each of these manifolds, and plotted the distribution of the eigenvalues of the quantum metric $g(x)$ across all data points $x \in X$ (Figure 3 (a-c) ). In all cases, a clear spectral gap emerges between the top $d$ eigenvalues that are near 1, and the remaining $D - d$ bottom eigenvalues that are near 0, where $d$ is the correct intrinsic dimension.

These higher-dimensional manifolds can also be used as a testing ground for the random matrix theory (RMT) estimate of the spectral gap. Recall that this technique can be applied whenever the quantum metric is of low rank and of high dimension, and is therefore not suitable for the $S^2$ or $M_{10b}$ examples. For $M_\beta$, the RMT estimate returns the correct dimension $d = 10$. For $MN_1$, the artificial rank bound of 30 imposed by our choice of $N = 16$ implies that the metric is not actually of low rank, and therefore the RMT estimate cannot be applied with this choice of $N$. We re-tested this example with a higher value $N = 37$, the smallest dimension for which the rank bound is equal to the embedding dimension $D = 72$, and obtained an estimate of $d = 15$.

Next, we plotted the intrinsic dimension estimate returned by Algorithm 1 as a function of Gaussian noise (Figure 3 (d-f) ) and compared the estimate to other standard intrinsic dimension estimators. A common theme among the standard estimators is to first under-estimate the intrinsic dimension, in the presence of zero or low noise. As explained in the introduction, this is a well-known effect due to the "curse of dimensionality", whereby neighboring points in high dimension tend to be very far apart. As the noise is increased, however, the "shadow dimension" effect overcomes the underestimating effects due to sparsity and the standard algorithms begin to
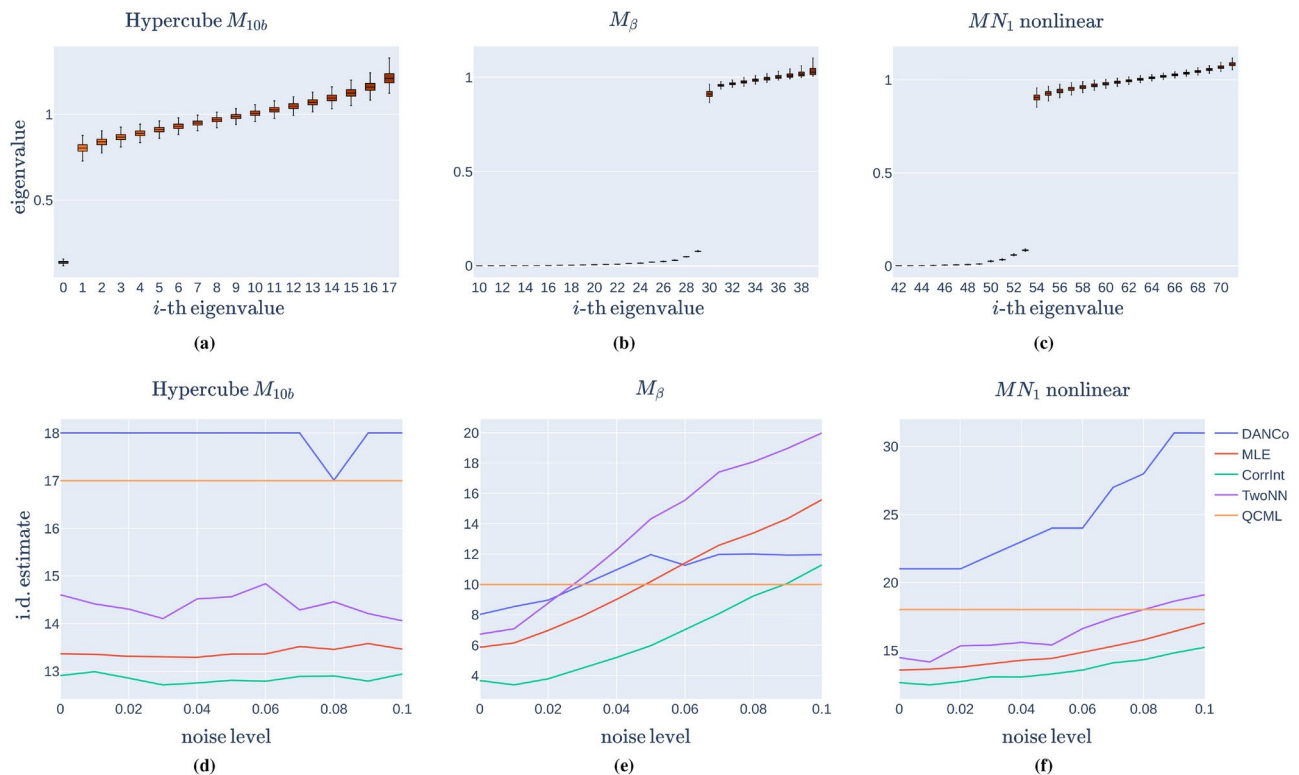


**Fig. 3**. Intrinsic dimension estimates for $T = 2500$ points on three higher-dimensional benchmark manifolds[4]: the 17-dimensional hypercube $M_{10b}$, the 10-dimensional $M_\beta$ manifold embedded into $D = 40$ dimensions, and the 18-dimensional manifold $MN_1$ embedded non-linearly into $D = 72$ dimensions. In the boxplots (**a-c**) the $i$-th box plot represents the distribution of the eigenvalue $e_i$ across all $T = 2500$ points. The outliers have been omitted from the plot for clarity. The plots (**d-f**) show the intrinsic dimension estimates for each manifold as functions of the `noise` parameter. In these examples a global estimate of dimension for the QCML estimator was obtained by taking the median of the local dimension estimates.

overestimate intrinsic dimension. This is particularly evident in the plots for $M_\beta$ and $MN_1$. In contrast, the spectral gap estimate of the QCML estimator is robust with respect to both these effects.

*Image recognition data sets*

We next test the QCML estimator on two of the real data sets suggested in the benchmarking framework of ref.[4], the ISOMAP face database and the MNIST handwritten digits database. The ISOMAP face database consists of 698 grayscale images of size $64 \times 64$ representing the face of a sculpture (Figure 4 (a)). Each image is represented as a vector in $D = 64^2 = 4096$ dimensions and it corresponds to a different rotation with respect to two axes and a different lighting direction, so that the intrinsic dimension of the data manifold in this case is expected to be $d = 3$[4,20]. In Figure 4 (b) a well-defined spectral gap indeed emerges between the top 3 eigenvalues of the quantum metric and the remaining 4093. This result was obtained by training with Hilbert space dimension $N = 32$. The value of $N = 32$ was chosen after experimenting with different Hilbert space dimensions until a clear spectral gap emerged. The RMT-based intrinsic dimension estimate for ISOMAP faces is $d = 3$.

The MNIST database consists of 70000 pictures of handwritten digits, each stored as a $28 \times 28$ grayscale picture. This data set is usually divided into a 60000 training images and 10000 testing images. The overall intrinsic dimension is unknown, but it is expected that each digit has its own intrinsic dimension. For example, in ref.[32] estimates for the dimension of each digit are in the range $d \in [8, 14]$. For our testing, we selected all the 1135 testing samples of the digit "1" (Figure 4 (c) ) and trained with Hilbert space dimension $N = 32$. In Figure 4 (d), a spectral gap can be identified either at $d = 9$ or $d = 10$. The QCML estimator in this case returns $d = 10$ at all points. This estimate however is not very stable with respect to changing $N$. For example, for $N = 24$ the estimate is $d = 12$ and for $N = 16$ we obtained $d = 9$. Similarly, the RMT-based intrinsic dimension estimates for varying $N$ ranges from $d = 9$ to $d = 13$. This instability is to be expected on real data, where a range of intrinsic dimension estimates for the data manifold is perhaps more appropriate.

*Wisconsin Breast Cancer data set*

We also test our intrinsic dimension estimator on the Diagnostic Wisconsin Breast Cancer Database[33]. This database consists of 569 data points representing images of fine needle aspirates (FNA) of a breast tumor. For each image, 30 features are extracted that describe characteristics of the cell nuclei present in the image. Therefore
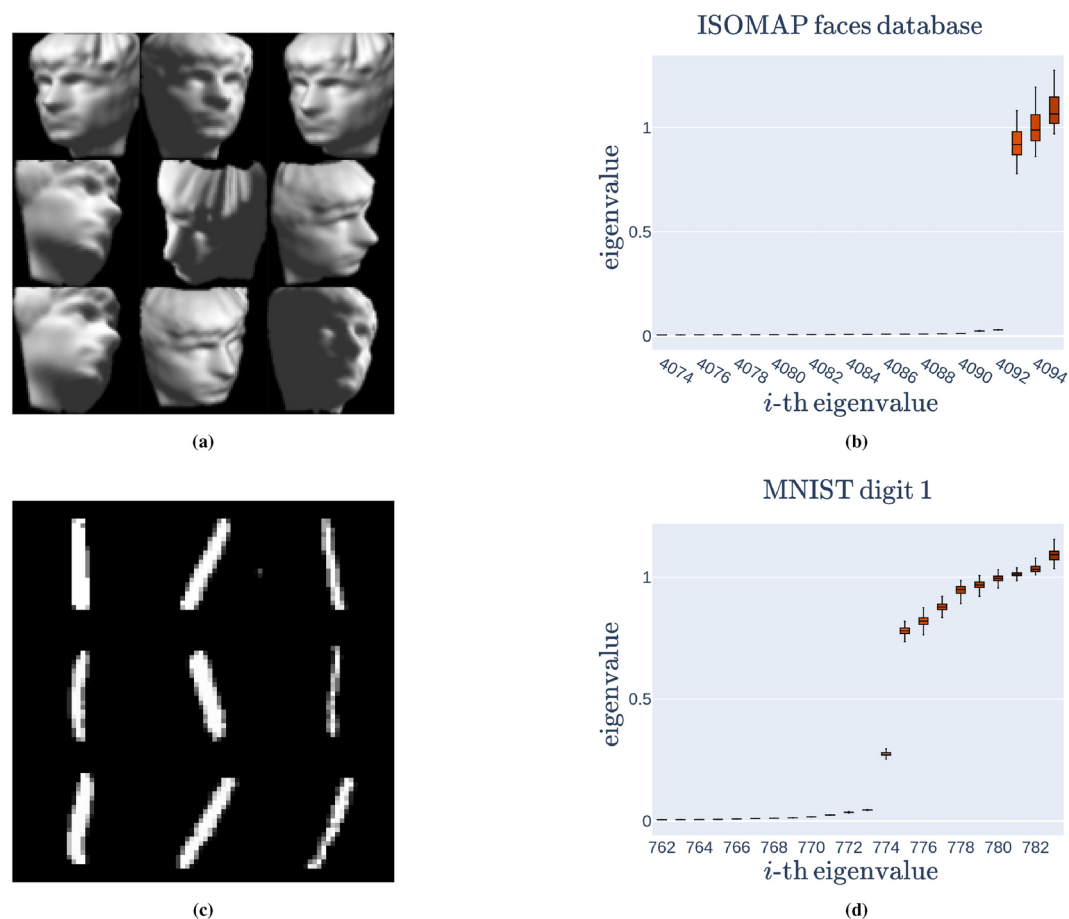


**Fig. 4**. (**a**) Examples of images from the ISOMAP face database, (**b**) Spectral gap for ISOMAP, showing an intrinsic dimension estimate of $d = 3$ at all points. (**c**) Examples of digit "1" in the MNIST data set, (d) spectral gap for MNIST digit "1" with $N = 32$. The intrinsic dimension estimate is $d = 10$ at all points.

in this case we have $T = 569$ points from a manifold sitting inside $D = 30$ dimensional Euclidean space. The intrinsic dimension $d$ of this data set has been estimated to be in the range $d \in [3.5, 6]$[34].

The thirty data features all have significantly different scales, with standard deviations in the range $2.65 \times 10^{-3}$ to $5.69 \times 10^2$. For this reason we chose to normalize the data by a standard scaling so that every feature has mean zero and standard deviation one. Since the data consists of $D = 30$ features, we choose $N = 16$ for the Hilbert space dimension, according to the rank bound (8). For the loss function, this time we chose to introduce a quantum fluctuation term with weight factor $w = 0.75$, as in (5). During testing, this choice led to sharper and more consistent spectral gaps (Figure 5), indicating a gap corresponding to $d = 2$. In general, the effects of the quantum fluctuation term on the loss function are analyzed more thoroughly in the supplementary material for this article (see Figure 1 of supplementary material).

To test the robustness of our estimate, we add synthetic Gaussian `noise` with increasing standard deviation. The goal is to provide an intrinsic dimension estimate that is constant across different levels of noise, just like we did for the synthetic manifold examples. The results are shown in Figure 5 (b), where we tested on 21 equally spaced noise levels from noise $= 0$ to noise $= 0.5$ in increments of 0.05.

The QCML estimator consistently returns an intrinsic dimension estimate of $d = 2$ across all levels of noise tested. We also plot in Figure 5 (b) the results for other estimators. The estimates of these other models tend to slope upwards as the noise level increases, precisely as in the synthetic manifold examples. If we assume that the dataset carries a natural level of noise, then Figure 5 suggests that the estimates of all the other methods should be revised downwards, and thus be closer to $d = 2$.

## Discussion

In this article we introduce a new data representation paradigm based on Quantum Cognition Machine Learning, with an application towards intrinsic dimension estimation. The idea is to learn a non-commutative quantum model[14–16] for the data manifold itself. This quantum model has the ability to abstract out the fundamental features of the geometry of the data manifold. In particular, we demonstrate how the intrinsic dimension of the data can be estimated from the point cloud produced by the quantum model. Because the point cloud reflects global properties of the data, our method is fundamentally robust to noise, as demonstrated on synthetic benchmarks. This is in contrast to other state-of-the-art techniques, which tend to overestimate intrinsic dimension by including "shadow" dimensionality from noise artifacts. In light of our results, we suggest a new paradigm for testing intrinsic dimension estimators: instead of focusing on noise-free synthetic benchmarks of increasing non-linearity and dimensionality, it is perhaps more relevant to focus on the development of estimators that are robust to noise. For practical applications, no real data is immune to noise, and not much meaning can be attached to an intrinsic dimension estimate that is highly dependent on noise levels.

While not a quantum algorithm in itself, it is possible in principle to implement part of the QCML intrinsic dimension estimator on a quantum computer, which could prove advantageous for very large Hilbert space dimension $N$. Developing robust algorithms that extract reliable quantum estimates is an important topic in quantum machine learning[35–37], and it would be interesting to apply QCML further in this direction.

## Methods

All the results and figures for this article have been obtained on a 32-core 13th Gen Intel Core i9-13950HX CPU with 64GB of memory, supplemented by a NVIDIA RTX 5000 Ada Generation Laptop GPU. Training the QCML models involves iterative updates to the quasi-coherent states $|\psi_0(x)\rangle$ and the matrix configuration $A$ to



**Fig. 5**. Intrinsic dimension estimates for the Wisconsin Breast Cancer Dataset using a QCML estimator of dimension $N = 16$ and quantum fluctuation weight $w = 0.75$ in the loss function. (**a**) Spectral gap with zero `noise`. Outliers omitted for clarity. (**b**) Intrinsic dimension estimates for different estimators as function of `noise`. For the QCML estimator, a global estimate of dimension is obtained by taking the mode of the local estimates.

| Dataset | $D$= no. features | $T$= no.samples | $N$= Hilbert space dimension | Running time |
|---|---|---|---|---|
| Sphere | 3 | 2500 | 3 | 2.9s |
| $M_{10b}$ | 18 | 2500 | 16 | 3.5s |
| $M_\beta$ | 40 | 2500 | 16 | 5.4s |
| $MN_1$ | 72 | 2500 | 16 | 7.2s |
| ISOMAP faces | 4096 | 698 | 32 | 142s |
| MNIST - digit 1 | 784 | 1135 | 32 | 28.2s |
| Breast Cancer | 30 | 569 | 16 | 1.9s |

**Table 1**. QCML estimator running times for the data sets analyzed.

lower the loss function until desired convergence is obtained. The specifics of each optimization step depend on the particular loss function used and the choice of initialization of the matrix configuration $A$. A typical training loop would consist, for each epoch, of:

(1) Calculate the quasi-coherent states $|\psi_0(x)\rangle$ for all data points $x \in X$ (or batch of data).
(2) Compute the loss function (4) or (5) and its gradients with respect to $A$.
(3) Update the matrix configuration $A$ with gradient descent.

The above training loop was implemented in PyTorch[27] to obtain all the matrix configurations shown in this article. All other intrinsic dimension estimators (DANCo, MLE, CorrInt, TwoNN) were tested through their implementation in the `scikit-dimension` Python package[23]. A summary of running times for the examples tested in this article are given in Table 1.

The running times represent the total of both times required for 1) training the matrix configuration $A$ and 2) calculate the eigenvalues of the quantum metric $g(x)$, at each point $x \in X_A$. The main bottleneck for the first step is the calculation of the ground states at each train iteration. In the current implementation this is calculated in PyTorch by diagonalizing the error Hamiltonian, which will be slow for large Hilbert space dimensions $N \gg 100$. Note however that in practice it is desirable to keep $N$ as small as possible to control variance, as discussed in the supplementary material (see Figure 2 of the supplementary material). For the second step, the running time is mostly driven by the number of features $D$, since finding the eigenvalues of the quantum metric involves diagonalizing a $D \times D$ matrix. For the ISOMAP dataset, it was necessary for example to batch the calculation of the eigenvalues because of GPU memory constraints. In both steps, the timing with respect to the number of samples $T$ can be greatly accelerated by using a GPU, and by doing so the algorithm is suitable for scaling to large data sets.

## Data availability
The datasets analysed in this article are all publicly available and are listed in the References section.

## References
1. Bishop, C. *Neural Networks for Pattern Recognition*, Oxford University Press, USA (1995).
2. Johnson, W. B. & Lindenstrauss, J. Extensions of Lipshitz mapping into Hilbert space. *Contemporary Mathematics* **26**, 189–206 (1984).
3. Donoho, D. L. High-dimensional data analysis: The curses and blessings of dimensionality. *AMS Math Challenges Lecture* **1**, 32 (2000).
4. Campadelli, P., Casiraghi, E., Ceruti, C. & Rozza, A. Intrinsic dimension estimation: Relevant techniques and a benchmark framework. *Mathematical Problems in Engineering* **1–21**, 2015. https://doi.org/10.1155/2015/759567 (2015).
5. Li, C., Farkhoor, H., Liu, R. & Yosinski, J. Measuring the intrinsic dimension of objective landscapes. ArXiv preprint arXiv:1804.08838 (2018).
6. Macocco, I., Mira, A. & Laio, A. Intrinsic dimension as a multi-scale summary statistics in network modeling. *Scientific Reports* **14**, 17756 (2024).
7. Xu, M. et al. A new graph gaussian embedding method for analyzing the effects of cognitive training. *PLoS Computational Biology* **16**, e1008186 (2020).
8. Zhou, S., Tordesillas, A., Pouragha, M., Bailey, J. & Bondell, H. On local intrinsic dimensionality of deformation in complex materials. *Scientific Reports* **11**, 10216 (2021).
9. Varghese, A., Santos-Fernandez, E., Denti, F., Mira, A. & Mengersen, K. A global perspective on the intrinsic dimensionality of covid-19 data. *Scientific Reports* **13**, 9761 (2023).
10. Fukunaga, K. & Olsen, D. An algorithm for finding intrinsic dimensionality of data. *IEEE Transactions on Computers* **C–20**, 176–183. https://doi.org/10.1109/T-C.1971.223208 (1971).
11. Cangelosi, R. & Goriely, A. Component retention in principal component analysis with application to cdna microarray data. *Biology Direct* **2**, 2. https://doi.org/10.1186/1745-6150-2-2 (2007).
12. Little, A. V., Lee, J., Jung, Y.-M. & Maggioni, M. Estimation of intrinsic dimensionality of samples from noisy low-dimensional manifolds in high dimensions with multiscale SVD. In *2009 IEEE/SP 15th Workshop on Statistical Signal Processing*, 85–88, https://doi.org/10.1109/SSP.2009.5278634 (2009).
13. Musaelian, K. et al. Quantum cognition machine learning: AI needs quantum. Tech. Rep., Qognitive, Inc, Miami Beach, Florida (2024). Available at https://www.qognitive.io/QCML-Qognitive,Inc.pdf.
14. Ishiki, G. Matrix geometry and coherent states. *Phys. Rev. D* **92**, 046009. https://doi.org/10.1103/PhysRevD.92.046009 (2015).

15. Schneiderbauer, L. & Steinacker, H. C. Measuring finite quantum geometries via quasi-coherent states. *Journal of Physics A: Mathematical and Theoretical* **49**, 285301. https://doi.org/10.1088/1751-8113/49/28/285301 (2016).
16. Steinacker, H. C. Quantum (matrix) geometry and quasi-coherent states. *Journal of Physics A: Mathematical and Theoretical* **54**, 055401. https://doi.org/10.1088/1751-8121/abd735 (2021).
17. Grassberger, P. & Procaccia, I. Measuring the strangeness of strange attractors. *Physica D: Nonlinear Phenomena* **9**, 189–208. https://doi.org/10.1016/0167-2789(83)90298-1 (1983).
18. Levina, E. & Bickel, P. Maximum likelihood estimation of intrinsic dimension. In Saul, L., Weiss, Y. & Bottou, L. (eds.) *Advances in Neural Information Processing Systems*, vol. 17 (MIT Press, 2004).
19. Ceruti, C. et al. Danco: An intrinsic dimensionality estimator exploiting angle and norm concentration. *Pattern Recognition* **47**, 2569–2581. https://doi.org/10.1016/j.patcog.2014.02.013 (2014).
20. Facco, E., d'Errico, M., Rodriguez, A. & Laio, A. Estimating the intrinsic dimension of datasets by a minimal neighborhood information. *Scientific Reports* **7** (2017).
21. Eckmann, J.-P & Ruelle, D. Fundamental limitations for estimating dimensions and lyapunov exponents in dynamical systems. *Physica D: Nonlinear Phenomena* **56**, 185–187. https://doi.org/10.1016/0167-2789(92)90023-G (1992).
22. Jiang, H., Kim, B., Guan, M. & Gupta, M. To trust or not to trust a classifier. *Advances in Neural Information Processing Systems* **31** (2018).
23. Bac, J., Mirkes, E. M., Gorban, A. N., Tyukin, I. & Zinovyev, A. Scikit-dimension: A python package for intrinsic dimension estimation. *Entropy* **23**, https://doi.org/10.3390/e23101368 (2021).
24. Steinacker, H. C. *Quantum Geometry, Matrix Theory, and Gravity* (Cambridge University Press, 2024).
25. Felder, L. O. & Steinacker, H. C. Oxidation, reduction and semi-classical limit for quantum matrix geometries. *Journal of Geometry and Physics* **199**, 105163. https://doi.org/10.1016/j.geomphys.2024.105163 (2024).
26. Canas, G., Poggio, T. & Rosasco, L. Learning manifolds with k-means and k-flats. *Advances in Neural Information Processing Systems* **25** (2012).
27. Paszke, A. et al. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, 8024–8035 (Curran Associates, Inc., 2019).
28. Popoviciu, T. Sur les équations algébriques ayant toutes leurs racines réelles. *Mathematica* **9**, 20 (1935).
29. Gavish, M. & Donoho, D. L. The optimal hard threshold for singular values is $4/\sqrt{3}$. *IEEE Transactions on Information Theory* **60**, 5040–5053. https://doi.org/10.1109/TIT.2014.2323359 (2014).
30. Marchenko, V. A. & Pastur, L. A. Distribution of eigenvalues for some sets of random matrices. *Matematicheskii Sbornik* **114**, 507–536 (1967).
31. Nadakuditi, R. R. When are the most informative components for inference also the principal components? arXiv preprint arXiv:1302.1232 (2013).
32. Hein, M. & Audibert, J.-Y. Intrinsic dimensionality estimation of submanifolds in rd. In *Proceedings of the 22nd International Conference on Machine Learning*, ICML '05, 289-296, https://doi.org/10.1145/1102351.1102388 (Association for Computing Machinery, New York, NY, USA, 2005).
33. Wolberg, W., Mangasarian, O., Street, N. & Street, W. Breast Cancer Wisconsin (Diagnostic). UCI Machine Learning Repository (1995). https://doi.org/10.24432/C5DW2B.
34. Mirkes, E. M., Allohibi, J. & Gorban, A. Fractional norms and quasinorms do not help to overcome the curse of dimensionality. *Entropy* **22**, https://doi.org/10.3390/e22101105 (2020).
35. Ren, W. et al. Experimental quantum adversarial learning with programmable superconducting qubits. *Nature Computational Science* **2**, 711–717 (2022).
36. Gong, W., Yuan, D., Li, W. & Deng, D.-L. Enhancing quantum adversarial robustness by randomized encodings. *Phys. Rev. Res.* **6**, 023020. https://doi.org/10.1103/PhysRevResearch.6.023020 (2024).
37. Li, W. & Deng, D.-L. Extracting reliable quantum outputs for noisy devices. *Nature Computational Science* **4**, 811–812 (2024).

## Acknowledgements

## Author contributions

L.C., A.A., K.M, D.V., C.H. and M.W. designed and performed the research. L.C., A.A. and M.W. prepared the manuscript. L.C. prepared the figures. All authors reviewed the manuscript.

## Declarations

## Competing Interests

The authors declare that they have no competing interests.

## Additional information

**Supplementary Information** The online version contains supplementary material available at https://doi.org/10.1038/s41598-025-91676-8.

**Correspondence** and requests for materials should be addressed to L.C.

**Reprints and permissions information** is available at www.nature.com/reprints.

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.