

# Behavior of the Maximum Likelihood in Quantum State Tomography

Travis L Scholten and Robin Blume-Kohout

Center for Computing Research (CCR), Sandia National Labs and  
Center for Quantum Information and Control (CQuIC), University of New Mexico

(Dated: April 10, 2017)

Quantum state tomography on a  $d$ -dimensional system demands resources that grow rapidly with  $d$ . Model selection can be used to tailor the number of fit parameters to the data, but quantum tomography violates some common assumptions that underly canonical model selection techniques based on ratios of maximum likelihoods (loglikelihood ratio statistics), due to the nature of the state space boundaries. Here, we study the behavior of the maximum likelihood in different Hilbert space dimensions, and derive an expression for a complexity penalty – the expected value of the loglikelihood ratio statistic (roughly, the logarithm of the maximum likelihood) – that can be used to make an appropriate choice for  $d$ .

In quantum information science, an experimentalist may wish to determine the quantum state  $\rho_0$  that is produced by a specific initialization procedure. This can be done using quantum state tomography [1]: many copies of  $\rho_0$  are produced; they are measured in diverse ways; and finally the outcomes of those measurements (data) are collated and analyzed to produce an estimate  $\hat{\rho}$ . This is a straightforward statistical inference process [2, 3], where the data are used to fit the parameters of a statistical model – provided we know what model to use. But this is not always the case. In state tomography, the parameter is  $\rho$ , and the model is the set of all possible density matrices on a Hilbert space  $\mathcal{H}$  (equipped with the Born rule). It is not always *a priori* obvious what  $\mathcal{H}$  or its dimension is; examples include optical modes [4–8] and leakage levels in AMO and superconducting [9, 10] qubits. Choosing an appropriate Hilbert space on the fly is an instance of the general statistical problem known as *model selection*, and while model selection is well-studied in classical statistics [11], applying it to quantum tomography leads to some surprising twists. These problems stem from the positivity constraint ( $\hat{\rho} \geq 0$ ). The techniques we use to resolve them are relevant not just to model selection, but to quantum tomography of rank-deficient states.

Here, we study model selection for choosing the Hilbert space dimension of a quantum system. We define the model  $\mathcal{M}_d$  (for any  $d \geq 1$ ) as

$$\mathcal{M}_d = \{\rho \mid \rho \in \mathcal{B}(\mathcal{H}_d), \text{Tr}(\rho) = 1, \rho \geq 0\}, \quad (1)$$

where  $\mathcal{B}(\mathcal{H}_d)$  is the space of bounded operators on a  $d$ -dimensional Hilbert space  $\mathcal{H}_d$ . Identifying a good choice for  $d$  is particularly relevant for tomography of optical modes, as formally,  $\mathcal{H}$  is infinite-dimensional, necessitating the use of smaller, finite-dimensional models. Even in less extreme cases, discovering  $d$  is larger (or smaller) than expected could help in the detection and diagnosis of leakage, or of couplings between the quantum device and an external environment. Model selection for  $d$  involves evaluating whether  $\mathcal{M}_d$  is a better model than  $\mathcal{M}_{d'}$ . How does one do so?

Many methods for selecting between multiple models involve fitting each model’s parameters using *maximum*

*likelihood estimation* (MLE) [12–14], which reports the parameter values that maximize the likelihood (the probability of the observed data). Classical estimation problems usually satisfy *local asymptotic normality* [15, 16], meaning that: (1) as  $N_{\text{samples}} \rightarrow \infty$ ,  $\hat{\rho}_{\text{MLE}}$  is normally distributed around  $\rho_0$  with covariance matrix  $\mathcal{I}^{-1}$ , and (2) the likelihood function in a neighborhood of  $\hat{\rho}_{\text{MLE}}$  is locally Gaussian with Hessian  $\mathcal{I}$ , where  $\mathcal{I}$  is the *Fisher information*.

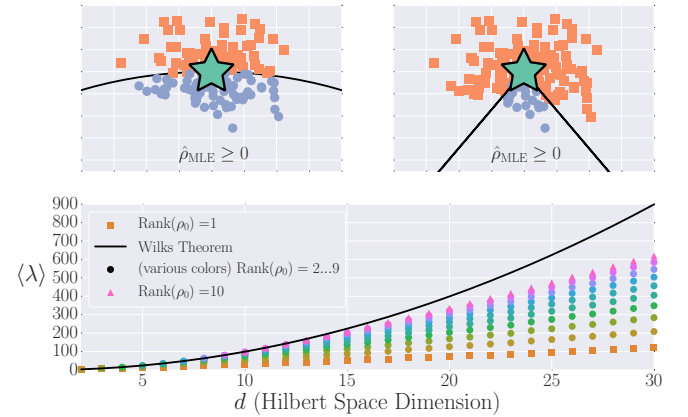


FIG. 1. Impact of the boundary on maximum likelihood tomography. **Top:** Two views through the qutrit state space. Without the positivity constraint, some estimates (orange squares) are not positive semidefinite, and do not represent valid estimates of a quantum state. The distribution of  $\hat{\rho}_{\text{MLE}}$  (blue circles) is generally non-normal, and depends on the true state  $\rho_0$  (star). **Bottom:** Comparison of the classical theory (Wilks Theorem) prediction for the loglikelihood ratio  $\langle \lambda \rangle$  to numerical data for states  $\rho_0$  with ranks  $r = 1, \dots, 10$ . The Wilks Theorem fails badly for low-rank states; our main result (Equation 16) fixes this problem (see Figure 5).

In state tomography, these conditions are violated by the constraint  $\hat{\rho}_{\text{MLE}} \geq 0$ . This constraint distorts the distribution of estimates, even in the asymptotic limit, precisely when  $\rho_0$  is rank-deficient within  $\mathcal{H}_d$ . (That is, when  $\rho_0$  is on the boundary of  $\mathcal{H}_d$ . See Figure 1.) While the behavior of MLEs is well-studied classically, its behavior near boundaries is not. How boundaries affect

$\hat{\rho}_{\text{MLE}}$  – and its derived properties – is of broad interest in state and process tomography [17–20], and is also critical for model selection [21–26].

We will focus on the special and simple case where  $\mathcal{I}$  at  $\rho_0$  is proportional to the Hilbert-Schmidt metric, so the likelihood function (and the distribution of the *unconstrained* estimates  $\hat{\rho}$ ) is given by

$$\mathcal{L}(\rho) = \Pr(\hat{\rho}|\rho) \propto e^{-\|\rho - \hat{\rho}\|_2^2 / 2\epsilon^2} \quad (2)$$

for some  $\epsilon$  that scales as  $1/\sqrt{N_{\text{samples}}}$ . In practice,  $\mathcal{I}$  depends on  $\rho_0$  and the particular tomographic measurement performed, but the interaction of an arbitrary Fisher information with boundaries [27] is complex and intractable. The isotropic assumption greatly simplifies our study of the problem, and permits the derivation of analytic results which capture realistic tomographic scenarios surprisingly well [28].

In this case,  $\hat{\rho}$  is often not positive semidefinite (see Figure 1). For each  $\hat{\rho}$ ,  $\hat{\rho}_{\text{MLE}}$  is the solution to the following optimization problem:

$$\hat{\rho}_{\text{MLE}} = \underset{\substack{\rho \in \mathcal{B}(\mathcal{H}) \\ \text{Tr}(\rho)=1, \rho \geq 0}}{\text{argmin}} \text{Tr}[(\hat{\rho} - \rho)^2]. \quad (3)$$

Note that if  $\hat{\rho} \geq 0$ ,  $\hat{\rho}_{\text{MLE}} = \hat{\rho}$ . Generally, the distribution of  $\hat{\rho}_{\text{MLE}}$  will not be normal; we do not attempt to derive  $\Pr(\hat{\rho}_{\text{MLE}})$  explicitly. Instead, we demonstrate a method for computing the behavior of useful statistics which depend on it.

On such statistic is the the *loglikelihood ratio* [29],

$$\lambda(\mathcal{M}_1, \mathcal{M}_2) = -2 \log \left( \frac{\max_{\rho \in \mathcal{M}_1} \mathcal{L}(\rho)}{\max_{\rho \in \mathcal{M}_2} \mathcal{L}(\rho)} \right), \quad (4)$$

which is commonly used to compare the goodness of fit between two models [14, 25, 29]. Intuitively, the model with the higher likelihood is more plausible – except that models with more adjustable parameters will almost always fit the data better! This is very clear in the case of *nested* models ( $\mathcal{M}_1 \subset \mathcal{M}_2$ ) [30]. If two models are equally valid – i.e. they both contain  $\rho_0$  – the larger one will usually fit the data better because its extra parameters allow it to fit more of the noise in the data. For the same reason, the larger model’s fit will be less accurate. This makes it imperative to correct for overfitting, by handicapping larger models.

For this reason, any model selection method that relies (explicitly or implicitly) on a statistic to quantify “how well model  $\mathcal{M}$  fits the data” also relies on a *null theory* to predict how that statistic will behave *if*  $\rho_0 \in \mathcal{M}$ . A model selection criterion based on an invalid null theory (or a criterion used in a context where its null theory does not apply) will tend to choose the wrong model. As we show below, state tomography violates the null theory for the behavior of the loglikelihood ratio statistic.

Consider using model selection to evaluate a particular  $d$ -dimensional Hilbert space  $\mathcal{H}_d$ , by comparing  $\mathcal{M}_d$  to

$\mathcal{M}_{d+1}$  using  $\lambda$ . Before we can use the observed value of  $\lambda$  to decide whether  $\mathcal{H}_{d+1}$  is significantly better, we need a null theory for its behavior when it *isn’t* better (i.e., when  $\rho_0 \in \mathcal{M}_d, \mathcal{M}_{d+1}$ ). In what follows, it’s useful to simplify the problem of computing  $\lambda(\mathcal{M}_d, \mathcal{M}_{d+1})$  to that of computing  $\lambda(\rho_0, \mathcal{M}_d)$  and  $\lambda(\rho_0, \mathcal{M}_{d+1})$  using the identity

$$\lambda(\mathcal{M}_d, \mathcal{M}_{d+1}) = \lambda(\rho_0, \mathcal{M}_{d+1}) - \lambda(\rho_0, \mathcal{M}_d), \quad (5)$$

where

$$\begin{aligned} \lambda(\rho_0, \mathcal{M}_d) &= -2 \log \left( \frac{\mathcal{L}(\rho_0)}{\max_{\rho \in \mathcal{M}_d} \mathcal{L}(\rho)} \right) \\ &= (\text{Tr}[(\rho_0 - \hat{\rho})^2] - \text{Tr}[(\hat{\rho}_{\text{MLE}} - \hat{\rho})^2]) / \epsilon^2. \end{aligned} \quad (6)$$

A further simplification is possible by observing that, as  $\epsilon \rightarrow 0$ , the distribution of  $\hat{\rho}$  will become more and more tightly concentrated around  $\rho_0$ . In turn, the local geometry of the Hilbert space will become more and more flat. For such state spaces [31], the loglikelihood ratio statistic reduces to the mean-squared error between  $\hat{\rho}_{\text{MLE}}$  and  $\rho_0$ :

$$\epsilon \rightarrow 0 \implies \lambda(\rho_0, \mathcal{M}_d) = \text{Tr}[(\hat{\rho}_{\text{MLE}} - \rho_0)^2] / \epsilon^2 \quad (7)$$

To summarize, through these identities and approximations, we have reduced the problem of computing  $\lambda(\mathcal{M}_d, \mathcal{M}_{d+1})$  to that of computing the (scaled) mean-squared error  $\text{Tr}[(\hat{\rho}_{\text{MLE}} - \rho_0)^2] / \epsilon^2$  for each  $\hat{\rho}_{\text{MLE}}$  in both  $\mathcal{M}_d$  and  $\mathcal{M}_{d+1}$ .

When  $\rho_0$  is full-rank and has eigenvalues which are substantially larger than  $\epsilon$ , local asymptotic normality holds, and the null theory for  $\lambda$  is given by the *Wilks Theorem* [32]. This theorem says that if  $\rho_0 \in \mathcal{M}_1 \subset \mathcal{M}_2$ , where  $\mathcal{M}_1$  has  $k$  free parameters and  $\mathcal{M}_2$  has  $K + k$  free parameters, then  $\lambda$  is a  $\chi_K^2$  random variable. Thus, when local asymptotic normality holds in state tomography,  $\lambda(\rho_0, \mathcal{M}_d)$  is a  $\chi_{d^2-1}^2$  random variable [33].

However if  $\rho_0$  is rank-deficient, then the boundary looms, and  $\Pr(\hat{\rho}_{\text{MLE}})$  is no longer Gaussian, meaning local asymptotic normality does not hold and the Wilks Theorem does not apply! Crucially, *even if  $\rho_0$  is full-rank in  $\mathcal{M}_d$ , it will be rank-deficient in  $\mathcal{M}_{d+1}$* . Thus, if we want to use model selection to choose  $d$  in state tomography, we must understand the null behavior of  $\lambda(\rho_0, \mathcal{M}_d)$  – i.e., derive a replacement Wilks Theorem – for rank-deficient  $\rho_0$ .

In our derivation, we assume that  $\rho_0, \hat{\rho}_{\text{MLE}} \in \mathcal{M}_d$ , that  $r \equiv \text{Rank}(\rho_0) < d$ , and that the Fisher information at  $\rho_0$  is  $\mathcal{I} = \epsilon^2 \mathbb{I}$ . The loglikelihood ratio that we’re trying to predict is given in Equation (7), where  $\hat{\rho}_{\text{MLE}}$  is defined in Equation (3), and the distribution of the *unconstrained* estimates around  $\rho_0$  is as given in Equation (2). Because  $\Pr(\hat{\rho}_{\text{MLE}})$  is complicated, depending strongly on the local geometry of the state space around  $\rho_0$ , we will not attempt to compute  $\Pr(\lambda)$  directly. Instead, we compute  $\langle \lambda \rangle$ , which may be used as a threshold for choosing  $d$ .

To start to do so, we need a procedure to compute  $\hat{\rho}_{\text{MLE}}$  given  $\hat{\rho}$  – i.e., to solve the optimization problem in Eq.

(3). Fortunately, an algorithm for doing so was presented in Ref. [28]:

1. Subtract  $q\mathbb{1}$  from the unconstrained  $\hat{\rho}$ , for a particular real scalar  $q$ ,
2. “Truncate”  $\hat{\rho} - q\mathbb{1}$ , by replacing each of its negative eigenvalues with zero.

Here,  $q$  is defined implicitly such that  $\text{Tr}[\text{Trunc}(\hat{\rho} - q\mathbb{1})] = 1$ .

Although this was intended as a (very fast) numerical algorithm, we will manipulate it (by a series of approximations) to derive a closed-form expression for the average  $\langle \lambda \rangle$ . We begin by observing that  $\lambda(\rho_0, \mathcal{M}_d)$  can be written as a sum over matrix elements,

$$\begin{aligned} \lambda &= \epsilon^{-2} \text{Tr}[(\hat{\rho}_{\text{MLE}} - \rho_0)^2] = \epsilon^{-2} \sum_{jk} |(\hat{\rho}_{\text{MLE}} - \rho_0)_{jk}|^2 \\ &= \sum_{jk} \lambda_{jk} \quad \text{where} \quad \lambda_{jk} = \epsilon^{-2} |(\hat{\rho}_{\text{MLE}} - \rho_0)_{jk}|^2, \end{aligned} \quad (8)$$

and therefore  $\langle \lambda \rangle = \sum_{jk} \langle \lambda_{jk} \rangle$ . Each term  $\langle \lambda_{jk} \rangle$  quantifies the average mean-squared error of a single matrix element of  $\hat{\rho}_{\text{MLE}}$ , and while the Wilks Theorem predicts  $\langle \lambda_{jk} \rangle = 1$  for all  $j, k$ , numerical simulations (see Figure 2) show that this only holds true for *some* matrix elements. A few contribute more than 1 unit (on average) while many others contribute much less, meaning that the Wilks Theorem predicts too high a value for the total  $\langle \lambda \rangle$ . (See bottom of Figure 1.) Thus motivated, we divide the parameters of  $\hat{\rho}$  into two parts (see Figure 2),

1. The “kite” comprises all diagonal elements *and* all elements on the kernel (null space) of  $\rho_0$ ,
2. The “L” comprises all off-diagonal elements on the support of  $\rho_0$  *and* between the support and the kernel,

and observe that  $\langle \lambda \rangle = \langle \lambda_{\text{L}} \rangle + \langle \lambda_{\text{kite}} \rangle$ . The rationale for this division is simple: small fluctuations on the “L” do not change the zero eigenvalues of  $\hat{\rho}$  to 1<sup>st</sup> order, whereas those on the “kite” do. In what follows, we study how imposing the positivity constraint  $\hat{\rho}_{\text{MLE}} \geq 0$  affects the behavior of the matrix elements in each part.

In doing so, it is helpful to think about the error of the unconstrained estimate  $\delta \equiv \hat{\rho} - \rho_0$ , a normally-distributed *traceless* matrix. To simplify the analysis, we explicitly drop the  $\text{Tr}(\rho) = 1$  constraint and let  $\delta$  be  $\mathcal{N}(0, \epsilon^2 \mathbb{1})$  distributed over the  $d^2$ -dimensional space of Hermitian matrices (a good approximation when  $d \gg 2$ ) [34], which makes  $\delta$  proportional to an element of the Gaussian Unitary Ensemble (GUE) [35].

To first order in  $\epsilon$ , elements of  $\delta$  in the “L” do not affect positivity, so they are unconstrained by the boundary, and behave exactly as expected from classical theory. The  $\delta_{jk}$  in the “L” may be seen as errors which arise due to small unitary perturbations of  $\rho_0$ . Writing  $\hat{\rho} = U^\dagger \rho_0 U$ , where  $U = e^{i\epsilon H}$ , we have

$$\hat{\rho} \approx \rho_0 + i\epsilon[\rho_0, H] + \mathcal{O}(\epsilon^2).$$

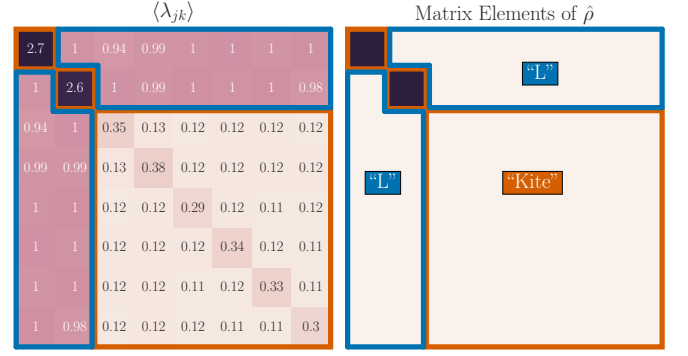


FIG. 2. When a rank-2 state is reconstructed in  $d = 8$  dimensions, the total loglikelihood ratio  $\lambda(\rho_0, \mathcal{M}_8)$  is the sum of terms  $\lambda_{jk}$  from errors in each matrix element  $(\hat{\rho}_{\text{MLE}})_{jk}$ . **Left:** Numerics show a clear division; some matrix elements have  $\langle \lambda_{jk} \rangle \sim 1$  as predicted by the Wilks Theorem, while others are either more or less. **Right:** The numerical results motivate dividing the elements of  $\hat{\rho}$  into two parts: the “kite” and the “L”.

Then,  $\delta \approx i\epsilon[\rho_0, H]$ . If  $j = k$ , then  $\delta_{jj} = 0$ . Thus, small unitaries cannot create errors in the diagonal matrix elements, at  $\mathcal{O}(\epsilon)$ . If  $j \neq k$ , then  $\delta_{jk} \neq 0$ , in general. (Small unitaries *can* introduce errors on off-diagonal elements.)

However, if either  $j$  or  $k$  (or both) lie within the *kernel* of  $\rho_0$  (i.e.,  $\langle k|\rho_0|k \rangle$  or  $\langle j|\rho_0|j \rangle$  is 0), then the corresponding  $\delta_{jk}$  are zero. The only off-diagonal elements where small unitaries can introduce errors are those which are coherent between the kernel of  $\rho_0$  and its support. These off-diagonal elements are precisely the “L”, and are the set  $\{\delta_{jk} \mid \langle j|\rho_0|j \rangle \neq 0, j \neq k, 0 \leq j, k \leq d-1\}$ . Each  $\delta_{jk}$  in the “L” is a  $\mathcal{N}(0, \epsilon^2)$  random variable, and crucially, is identical to the error  $(\hat{\rho}_{\text{MLE}} - \rho_0)_{jk}$ . This is because these  $\delta_{jk}$  are *unaffected by the boundary*, so when we impose the positivity constraint (i.e., compute  $\hat{\rho}_{\text{MLE}}$ ), their values remain the same. Therefore,  $\langle \lambda_{jk} \rangle = \langle \delta_{jk}^2 \rangle / \epsilon^2 = 1$ . As there are  $2rd - r(r+1)$  of them,  $\langle \lambda_{\text{L}} \rangle = 2rd - r(r+1)$ .

Computing  $\langle \lambda_{\text{kite}} \rangle$  is a bit harder, because the boundary *does* constrain its elements. Here, we turn to the truncation algorithm given above for finding  $\hat{\rho}_{\text{MLE}}$ , which is most naturally performed in the eigenbasis of  $\hat{\rho}$ . Exact diagonalization of  $\hat{\rho}$  is not feasible analytically, but only the *small* eigenvalues of  $\hat{\rho}$  are critical in truncation. As long as all the nonzero eigenvalues of  $\rho_0$  are much larger than  $\epsilon$ , the eigenbasis of  $\hat{\rho}$  is accurately approximated by: (1) the eigenvectors of  $\rho_0$  on its support; and (2) the eigenvectors of  $\delta_{\text{ker}} = \Pi_{\text{ker}} \delta \Pi_{\text{ker}}$ , where  $\Pi_{\text{ker}}$  is the projector onto the kernel of  $\rho_0$ .

Changing to this basis diagonalizes the “kite” portion of  $\delta$ , and leaves all elements of the “L” unchanged (at  $\mathcal{O}(\epsilon)$ ). The diagonal elements of  $\hat{\rho}$  now fall into two categories:

1.  $r$  elements corresponding to the eigenvalues of  $\rho_0$ , which are given by  $p_j = \rho_{jj} + \delta_{jj}$  where  $\rho_{jj}$  is the  $j^{\text{th}}$  eigenvalue of  $\rho_0$ , and  $\delta_{jj} \sim \mathcal{N}(0, \epsilon^2)$ .
2.  $N \equiv d - r$  elements that are eigenvalues of  $\delta_{\text{ker}}$ ,

which we denote by  $\kappa = \{\kappa_j : j = 1 \dots N\}$ ,

and  $\lambda_{\text{kite}}$  is

$$\epsilon^2 \lambda_{\text{kite}} = \sum_{j=1}^r [\rho_{jj} - (p_j - q)^+]^2 + \sum_{j=1}^N [(\kappa_j - q)^+]^2, \quad (9)$$

where  $(x)^+ = \max(x, 0)$ .  $q$  is implicitly defined such that  $f(q) \equiv \text{Tr}[\text{Trunc}(\hat{\rho} - q\mathbb{I})]$  satisfies  $f(q) = 1$ . In terms of the eigenvalues of  $\hat{\rho}$ , this means  $q$  is the solution to

$$\sum_{j=1}^r (p_j - q)^+ + \sum_{j=1}^N (\kappa_j - q)^+ = 1 \quad (10)$$

To compute  $q$ , we first observe that while the  $\kappa_j$  are random variables, they are not normally distributed. Instead, because  $\delta_{\text{ker}}$  is proportional to a  $\text{GUE}(N)$  matrix, for  $N \gg 1$ , the distribution of any eigenvalue  $\kappa_j$  converges to a Wigner semicircle distribution [36] given by  $\text{Pr}(\kappa) = \frac{2}{\pi R^2} \sqrt{R^2 - \kappa^2}$  for  $|\kappa| \leq R$ , with  $R = 2\epsilon\sqrt{N}$ . The eigenvalues are not independent; they tend to avoid collisions (“level avoidance” [37]), and typically form a surprisingly regular array over the support of the Wigner semicircle. Since our goal is to compute  $\langle \lambda_{\text{kite}} \rangle$ , we can capitalize on this behavior by replacing each random sample of  $\kappa$  with a *typical sample*  $\bar{\kappa}$  given by its order statistics. These are the average values of the *sorted*  $\kappa$ , so  $\bar{\kappa}_j$  is the average value of the  $j^{\text{th}}$  largest value of  $\kappa$ . Large random samples are usually well approximated (for many purposes) by their order statistics even when the elements of the sample are independent, and level avoidance makes the approximation even better.

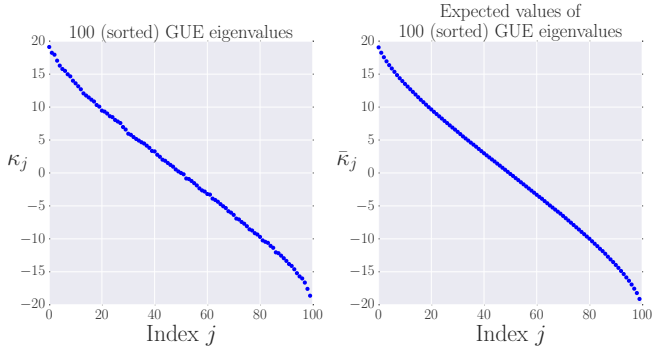


FIG. 3. Typical samples of  $\text{GUE}(N)$  eigenvalues are accurately approximated by order statistics of the distribution (average values of a sorted sample). **Left:** The sorted eigenvalues (i.e., order statistics  $\kappa_j$ ) of one randomly chosen  $\text{GUE}(100)$  matrix. **Right:** Approximate expected values of the order statistics,  $\bar{\kappa}_j$ , of the  $\text{GUE}(100)$  distribution, computed as the average of the sorted eigenvalues of 100 randomly chosen  $\text{GUE}(100)$  matrices.

Suppose that  $\kappa$  are the eigenvalues of a  $\text{GUE}(N)$  matrix, sorted from highest to lowest. Figure 3 illustrates such a sample for  $N = 100$ . It also shows the *average* values of 100 such samples (all sorted). These are

the *order statistics*  $\bar{\kappa}$  of the distribution (more precisely, what is shown is a good *estimate* of the order statistics; the actual order statistics would be given by the average over infinitely many samples). The point of the figure is to show that, while the order statistics *are* slightly more smoothly and predictably distributed than a single (sorted) sample, the two are remarkably similar. A single sample  $\kappa$  will fluctuate around the order statistics, but these fluctuations are relatively small, partly because the sample is large, and partly because the  $\text{GUE}$  eigenvalues experience level repulsion. Thus, the “typical” behavior of a sample – by which we mean the mean value of a statistic of the sample – is well captured by the order statistics (which have no fluctuations at all).

We now turn to the problem of modeling  $\kappa$  quantitatively. We note up front that we are only going to be interested in certain properties of  $\kappa$ : specifically, partial sums of all  $\kappa_j$  greater or less than the threshold  $q$ , or partial sums of functions of the  $\kappa_j$  (e.g.  $(\kappa_j - q)^2$ ). We require only that an ansatz be accurate for such quantities. We do not use this fact explicitly, but it motivates our approach – and we do not claim that our ansatz is accurate for *all* conceivable functions.

In general, if a sample  $\kappa$  of size  $N$  is drawn so that each  $\kappa$  has the same probability density function  $\text{Pr}(\kappa)$ , then a good approximation for the  $j^{\text{th}}$  order statistic is given by the inverse *cumulative* distribution function (CDF):

$$\bar{\kappa}_j \approx \text{CDF}^{-1} \left( \frac{j - 1/2}{N} \right). \quad (11)$$

This is closely related to the observation that the histogram of a sample tends to look similar to the underlying probability density function. More precisely, it is equivalent to the observation that the empirical distribution function (the CDF of the histogram) tends to be (even more) similar to the underlying CDF. (For i.i.d. samples, this is the content of the Glivenko-Cantelli theorem [38]). Figure 4 compares the order statistics of  $\text{GUE}(100)$  and  $\text{GUE}(10)$  eigenvalues (computed as numerical averages over 100 random samples) to the inverse CDF for the Wigner semicircle distribution. Even though the Wigner semicircle model of  $\text{GUE}$  eigenvalues is only exact as  $N \rightarrow \infty$ , it provides a nearly-perfect model for  $\bar{\kappa}$  even at  $N = 10$  (and remains surprisingly good all the way down to  $N = 2$ ).

We make one further approximation, by assuming that  $N \gg 1$ , so the distribution of the  $\bar{\kappa}_j$  is effectively continuous and identical to  $\text{Pr}(\kappa)$ . For the quantities that we compute, this is equivalent to replacing the empirical distribution function (which is a step function) by the CDF of the Wigner semicircle distribution. So, whereas for any given sample the partial sum of all  $\kappa_j > q$  jumps discontinuously when  $q = \kappa_j$  for any  $j$ , in this approximation it changes smoothly. This accurately models the *average* behavior of partial sums.

These approximations provide the ansatz that we use below, for the eigenvalues of  $\hat{\rho}$ , as  $\{p_j\} \cup \{\bar{\kappa}_j\}$ , where the  $p_j$  are  $\mathcal{N}(\rho_{jj}, \epsilon^2)$  random variables, and the  $\bar{\kappa}_j$  are the

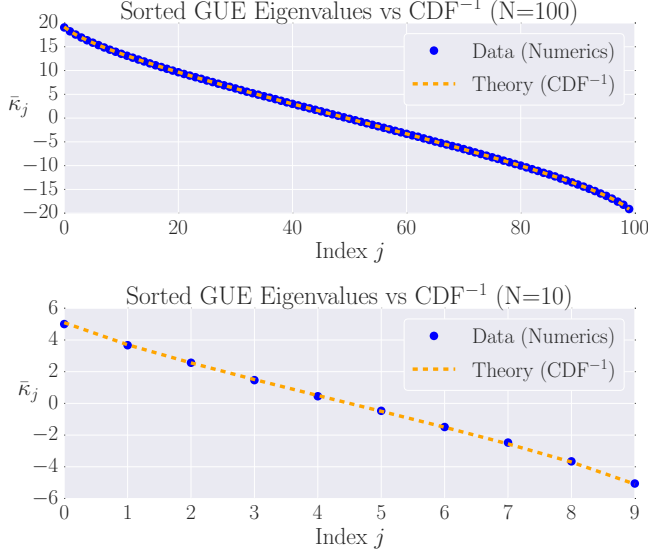


FIG. 4. Order statistics of the GUE( $N$ ) eigenvalue distribution are very well approximated by the inverse CDF of the Wigner semicircle distribution. In both figures, we compare the order statistics of a GUE( $N$ ) distribution to the inverse CDF of the Wigner semicircle distribution. **Top:**  $N = 100$ . **Bottom:**  $N = 10$ . Agreement in both cases is essentially perfect.

(fixed, smoothed) order statistics of a Wigner semicircle distribution. In turn, our expression for  $\langle \lambda_{\text{kite}} \rangle$  (Equation (9)) and the defining equation for  $q$  (Equation (10)) are both well approximated as

$$\langle \lambda_{\text{kite}} \rangle \approx \frac{1}{\epsilon^2} \left\langle \sum_{j=1}^r [\rho_{jj} - (p_j - q)^+]^2 + \sum_{j=1}^N [(\bar{\kappa}_j - q)^+]^2 \right\rangle$$

$$1 \approx \sum_{j=1}^r (p_j - q)^+ + \sum_{j=1}^N (\bar{\kappa}_j - q)^+$$

To proceed with truncation, we observe that the  $\bar{\kappa}_j$  are symmetrically distributed around  $\kappa = 0$ , so half of them are negative. Therefore, with high probability,  $\text{Tr}[\text{Trunc}(\hat{\rho})] > 1$ , and so we will need to subtract  $q$  from  $\hat{\rho}$  before truncating. (This is in distinction to the case where we have to *add*  $q$ .)

We make another assumption; namely, that the eigenvalues of  $\rho_0$  are large compared to the perturbations  $\delta_{jj}$  and  $q$ . This implies  $(p_j - q)^+ = p_j - q$ . Under this

assumption,  $q$  is the solution to

$$1 \approx \sum_{j=1}^r (p_j - q)^+ + \sum_{j=1}^N (\bar{\kappa}_j - q)^+$$

$$\approx 1 - rq + \Delta + N \int_{\kappa=q}^{2\epsilon\sqrt{N}} (\kappa - q) \text{Pr}(\kappa) d\kappa$$

$$\Rightarrow 0 = -rq + \Delta + \frac{\epsilon}{12\pi} \left[ \frac{(q^2 + 8N)\sqrt{-q^2 + 4N}}{-12qN \left( \frac{\pi}{2} - \sin^{-1} \left( \frac{q}{2\sqrt{N}} \right) \right)} \right], \quad (12)$$

where  $\Delta = \sum_{j=1}^r \delta_{jj}$  is a  $\mathcal{N}(0, r\epsilon^2)$  random variable. We choose to replace a discrete sum (line 1) with an integral (line 2). This approximation is valid when  $N \gg 1$ , as we can accurately approximate a discrete collection of closely spaced real numbers by a smooth density or distribution over the real numbers that has approximately the same CDF. It is also remarkably accurate in practice.

In yet another approximation, we replace  $\Delta$  with its average value, which is zero. We could obtain an even more accurate expression by treating the fluctuations in  $\Delta$  more carefully, but this crude approximation turns out to be quite accurate already.

To solve Equation (12), it is necessary to further simplify the complicated expression resulting from the integral (line 3). To do so, we assume  $\rho_0$  is relatively low-rank, so  $r \ll N$ . In this case, the sum of the positive  $\bar{\kappa}_j$  is large compared with  $r$ , almost all of them need to be subtracted away, and therefore  $q$  is close to  $2\epsilon\sqrt{N}$ . [39] We therefore replace the complicated expression with its leading order Taylor expansion around  $q = 2\epsilon\sqrt{N}$ , substitute into Equation (12), and obtain the equation

$$\frac{rq}{\epsilon} = \frac{4}{15\pi} N^{1/4} \left( 2\sqrt{N} - \frac{q}{\epsilon} \right)^{5/2}. \quad (13)$$

This equation is a quintic polynomial, so it has no closed-form solution. However, its roots have a well-defined asymptotic ( $N \rightarrow \infty$ ) expansion that becomes accurate quite rapidly (e.g., for  $N > 4$ ):

$$z \equiv q/\epsilon \approx 2\sqrt{N} - \frac{(240r\pi)^{2/5}}{4} N^{1/10} + \frac{(240r\pi)^{4/5}}{80} N^{-3/10}. \quad (14)$$

Now that we know how much to subtract off in the trun-



cation process, we can compute  $\langle \lambda_{\text{kite}} \rangle$ :

$$\begin{aligned}
\langle \lambda_{\text{kite}} \rangle &\approx \frac{1}{\epsilon^2} \left\langle \sum_{j=1}^r [\rho_{jj} - (p_j - q)^+]^2 + \sum_{j=1}^N [(\bar{\kappa}_j - q)^+]^2 \right\rangle \\
&\approx \frac{1}{\epsilon^2} \left\langle \sum_{j=1}^r [-\delta_{jj} + q]^2 + \sum_{j=1}^N [(\bar{\kappa}_j - q)^+]^2 \right\rangle \\
&\approx r + rz^2 + \frac{N}{\epsilon^2} \int_{\kappa=q}^{2\epsilon\sqrt{N}} \text{Pr}(\kappa)(\kappa - q)^2 d\kappa \\
&= r + rz^2 + \frac{N(N+z^2)}{\pi} \left( \frac{\pi}{2} - \sin^{-1} \left( \frac{z}{2\sqrt{N}} \right) \right) \\
&\quad - \frac{z(z^2 + 26N)}{24\pi} \sqrt{4N - z^2} \quad (15)
\end{aligned}$$

Thus the total expected value,  $\langle \lambda \rangle = \langle \lambda_L \rangle + \langle \lambda_{\text{kite}} \rangle$ , is

$$\begin{aligned}
\langle \lambda(\rho_0, \mathcal{M}_d) \rangle &\approx 2rd - r^2 + rz^2 \\
&\quad + \frac{N(N+z^2)}{\pi} \left( \frac{\pi}{2} - \sin^{-1} \left( \frac{z}{2\sqrt{N}} \right) \right) \\
&\quad - \frac{z(z^2 + 26N)}{24\pi} \sqrt{4N - z^2} \quad (16)
\end{aligned}$$

where  $z$  is given in Equation (14),  $N = d - r$ , and  $r = \text{Rank}(\rho_0)$ .

Equation (16) is our main result. To test its validity, we compare it to numerical simulations for  $d = 2, \dots, 30$  and  $r = 1, \dots, 10$ , in Figure 5. The prediction of the Wilks Theorem is wildly incorrect for  $r \ll d$ . In contrast, Equation (16) is almost perfectly accurate when  $r \ll d$ , but it does begin to break down (albeit fairly gracefully) as  $r$  becomes comparable to  $d$ . We conclude that our analysis (and Equation (16)) correctly models tomography *if* the Fisher information is isotropic ( $\mathcal{I} \propto \mathbb{1}$ ).

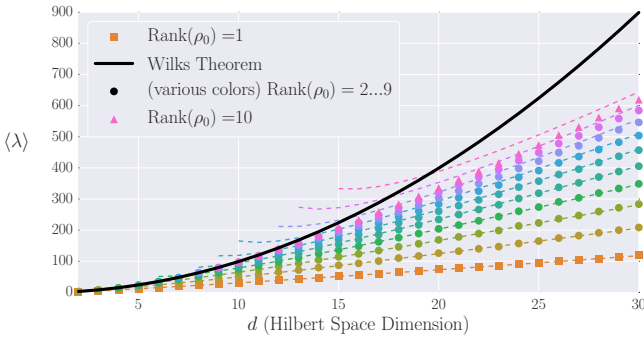


FIG. 5. Numerical results for  $\langle \lambda \rangle$  compared to the prediction of the Wilks Theorem (solid line) and our replacement theory as given in Equation (16), (dashed lines). Our formula depends on the rank  $r$  of  $\rho_0$  (unlike the Wilks prediction), and is nearly perfect for  $r \ll d$ . It becomes less accurate as  $r$  approaches  $d/2$ , and is invalid when  $r \approx d$ .

In practice, the Fisher information is rarely isotropic. So we tested our idealized result by applying it to a realistic, challenging, and experimentally relevant problem:

quantum heterodyne (equivalent to double homodyne) state tomography [5, 6, 8, 40] of a single optical mode. (See Figure 6 for a plot of the *condition number* – the ratio of the largest to smallest eigenvalue – of the estimated Fisher information. It is clear that, for such a tomographic setup,  $\mathcal{I} \not\propto \mathbb{1}$ .) States of this continuous-variable system are described by density operators on the infinite-dimensional Hilbert space  $L^2(\mathbb{R})$ . Fitting these infinitely many parameters to finitely much data demands simpler models.

We consider a family of nested models motivated by a low-energy (few-photon) ansatz, and choose the Hilbert space  $\mathcal{H}_d$  to be that spanned by the photon number states  $\{|0\rangle, \dots, |d-1\rangle\}$ . Heterodyne tomography reconstructs  $\rho_0$  using data from repeated measurements of the coherent-state POVM,  $\{|\alpha\rangle\langle\alpha|/\pi, \alpha = x+ip \in \mathbb{C}\}$ , which corresponds to sampling directly from the state's Husimi  $Q$ -function [41].

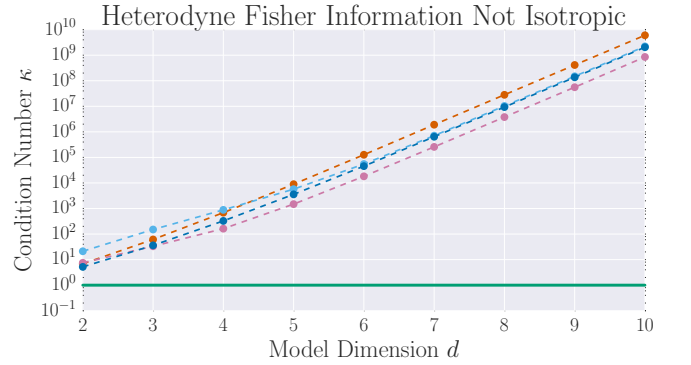


FIG. 6. The condition number  $\kappa$  – the ratio of the largest to smallest eigenvalue – of the estimated heterodyne Fisher information grows with model dimension, indicating an increase in anisotropy. (Estimates are the average over 100 Hessians of the loglikelihood function.) The dashed lines indicate different states  $\rho_0$ , and the solid line is  $\kappa = 1$  (i.e.,  $\mathcal{I} \propto \mathbb{1}$ ).

We examined the behavior of  $\lambda$  for 13 distinct states  $\rho_0$ , both pure and mixed, supported on  $\mathcal{H}_2, \mathcal{H}_3, \mathcal{H}_4$ , and  $\mathcal{H}_5$ . We used rejection sampling to simulate 100 heterodyne datasets with up to  $N_{\text{samples}} = 10^5$ , and found MLEs over each of the 9 models  $\mathcal{M}_2, \dots, \mathcal{M}_{10}$  using numerical optimization [42]. For each  $\rho_0$  and each  $d$ , we averaged  $\lambda(\rho_0, \mathcal{M}_d)$  over all 100 datasets to obtain an empirical average loglikelihood ratio  $\bar{\lambda}$  for each  $(\rho_0, d)$  pair.

Results of this test are shown in Figure 7, where we plot the predictions for  $\langle \lambda \rangle$  given by the Wilks Theorem and Equation (16), against the empirical average  $\bar{\lambda}$ , for a variety of  $\rho_0$  and  $d$ . Our formula correlates very well with the empirical average, while the Wilks Theorem (unsurprisingly) overestimates  $\lambda$  dramatically for low-rank states. Whereas a model selection procedure based on Wilks Theorem would tend to falsely reject larger Hilbert spaces (by setting the threshold for acceptance too high), our formula provides a reliable null theory.

Interestingly, as  $d$  grows, Equation (16) also begins

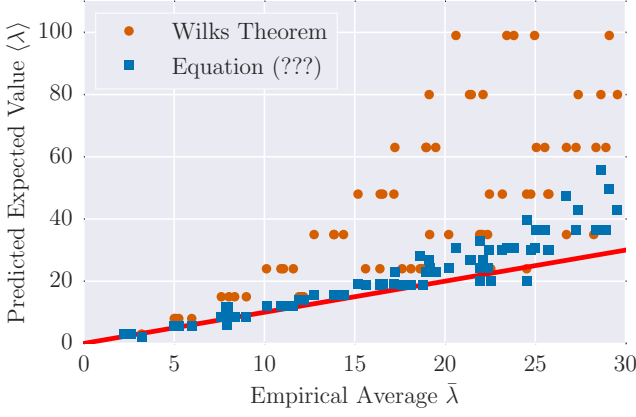


FIG. 7. The Wilks Theorem (orange dots) dramatically overestimates  $\langle \lambda(\rho_0, \mathcal{M}_d) \rangle$  in optical heterodyne tomography. Our formula, Equation 16 (blue squares), is far more accurate. Residual discrepancies occur in large part because  $N_{\text{samples}}$  is not yet “asymptotically large”. The solid red line corresponds to perfect correlation between theory ( $\langle \lambda \rangle$ ) and practice ( $\bar{\lambda}$ ).

to overpredict. As Figure 8 indicates, a more accurate description is that the numerical experiments are *under-achieving*, because  $\bar{\lambda}$  is still growing with  $N_{\text{samples}}$ . Both the Wilks Theorem and our analysis are derived in an asymptotic limit  $N \rightarrow \infty$ ; for finite but large  $N$ , both may be invalid. Figure 8 shows that, even at  $N \sim 10^5$ , the behavior of  $\bar{\lambda}$  has failed to become asymptotic. This is surprising, and suggests heterodyne tomography is a particularly exceptional and challenging case to model statistically.

However, our model *does* get some of the qualitative features correct. In Figure 9, we look at  $\langle \lambda_{jk} \rangle$ , where we assume an isotropic Fisher information, and when we simulate heterodyne tomography. While the numbers given for  $\langle \lambda_{jk} \rangle$  do not agree exactly, they still break down into two groups, the “L” and the “kite”. (See Figure 10 for an analysis of the exact differences in the values.)

The Wilks Theorem is not generally reliable in quantum state tomography, but our Equation (16) provides a much more broadly applicable replacement that can be used in model selection methods. This includes protocols like the AIC and BIC [11, 43–45] that do not explicitly use the Wilks Theorem, but rely on the same assumptions (asymptotic normality, etc). Null theories of loglikelihood ratios have many other applications, including hypothesis testing [14, 25] and confidence regions [46], and our result is directly applicable to them. Refs. [25, 46] both point out explicitly that their methods are unreliable near boundaries and therefore cannot be applied to rank-deficient states; our result fixes this outstanding problem. However, our numerical experiments with heterodyne tomography show unexpected behavior, indicating that quantum tomography can still surprise, and may violate *all* asymptotic statistics results. In such cases, bootstrapping [47, 48] may be the only reliable way to construct null theories for  $\lambda$ . Finally, the *meth-*

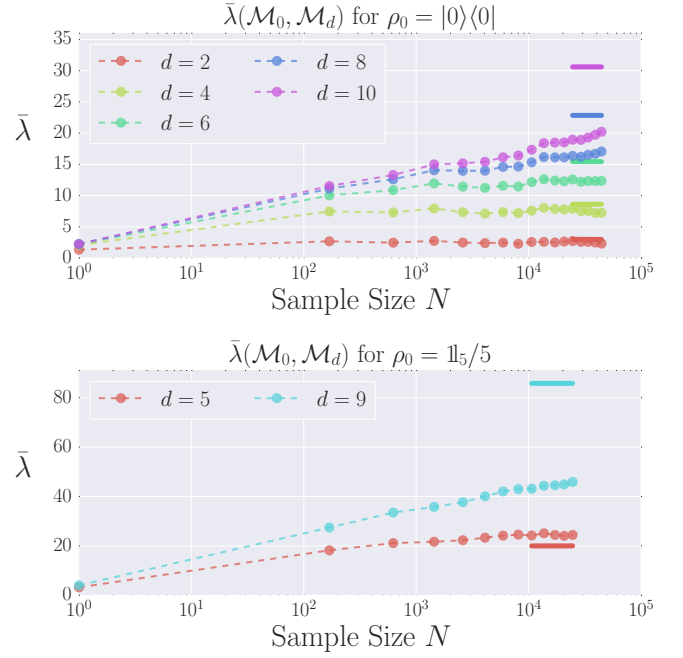


FIG. 8. The empirical average  $\bar{\lambda}$  may have achieved its asymptotic value, or is still growing, depending on the true state  $\rho_0$  and the model dimension  $d$ . Solid lines indicate the value of our formula for the asymptotic expected value, given in Equation (16).

Isotropic Model (10000 Trials)										Heterodyne Tomography (100 Trials)									
0	5.3	0.99	0.97	0.98	1	1	0.99	0.99		0	1.3	1	0.92	0.94	0.85	0.95	0.85	0.67	
1	0.99	0.17	0.06	0.06	0.06	0.06	0.06	0.06		1	1	0.11	0.12	0.05	0.02	0.01	0.01	0.01	
2	0.97	0.06	0.16	0.06	0.06	0.06	0.06	0.06		2	0.92	0.12	0.06	0.06	0.04	0.03	0.02	0.01	
3	0.98	0.06	0.06	0.16	0.06	0.06	0.06	0.06		3	0.94	0.05	0.06	0.05	0.04	0.02	0.02	0.01	
4	1	0.06	0.06	0.06	0.16	0.06	0.06	0.06		4	0.85	0.02	0.04	0.04	0.04	0.02	0.02	0.02	
5	1	0.06	0.06	0.06	0.06	0.16	0.06	0.06		5	0.95	0.01	0.03	0.02	0.02	0.02	0.01	0.02	
6	0.99	0.06	0.06	0.06	0.06	0.06	0.16	0.06		6	0.85	0.01	0.02	0.02	0.02	0.01	0.02	0.02	
7	0.99	0.06	0.06	0.06	0.06	0.06	0.06	0.16		7	0.67	0.01	0.01	0.01	0.02	0.02	0.02	0.02	
	0	1	2	3	4	5	6	7			0	1	2	3	4	5	6	7	
0	2.7	0.99	0.97	0.98	1	1	0.99	0.99		0	0.79	1.1	0.94	0.77	0.51	0.38	0.28	0.35	
1	0.99	2.6	1	0.99	1	1	1	0.99		1	1.1	1.8	1.1	0.79	0.89	0.79	0.7	0.57	
2	0.97	1	0.33	0.12	0.12	0.12	0.12	0.12		2	0.94	1.1	1.14	0.08	0.04	0.03	0.01	0.01	
3	0.98	0.99	0.12	0.34	0.12	0.12	0.12	0.12		3	0.77	0.79	0.08	0.11	0.04	0.03	0.02	0.02	
4	1	1	0.12	0.12	0.33	0.12	0.12	0.12		4	0.51	0.89	0.04	0.04	0.08	0.04	0.03	0.02	
5	1	1	0.12	0.12	0.12	0.34	0.12	0.12		5	0.38	0.79	0.03	0.03	0.04	0.08	0.04	0.03	
6	0.99	1	0.12	0.12	0.12	0.12	0.33	0.12		6	0.28	0.7	0.01	0.02	0.03	0.04	0.05	0.03	
7	0.99	0.99	0.12	0.12	0.12	0.12	0.12	0.34		7	0.35	0.57	0.01	0.02	0.02	0.03	0.03	0.05	
	0	1	2	3	4	5	6	7			0	1	2	3	4	5	6	7	

FIG. 9. The values of  $\langle \lambda_{jk} \rangle$  assuming an isotropic Fisher information (left), and for heterodyne tomography (right). **Top:**  $\rho_0 = |0\rangle\langle 0|$ . **Bottom:**  $\rho_0 = \mathcal{I}_2/2$ . **Discussion:** Qualitatively, the behavior is the same, though there are quantitative differences, particularly within the kite.

*ods* presented here have application beyond the analysis of loglikelihoods. They shed light on the behavior of  $\hat{\rho}_{\text{MLE}}$  for rank-deficient states, and can be used to predict other derived properties such as the average rank of the estimate, which is independently interesting for (e.g.) quantum compressed sensing [18, 49–51].

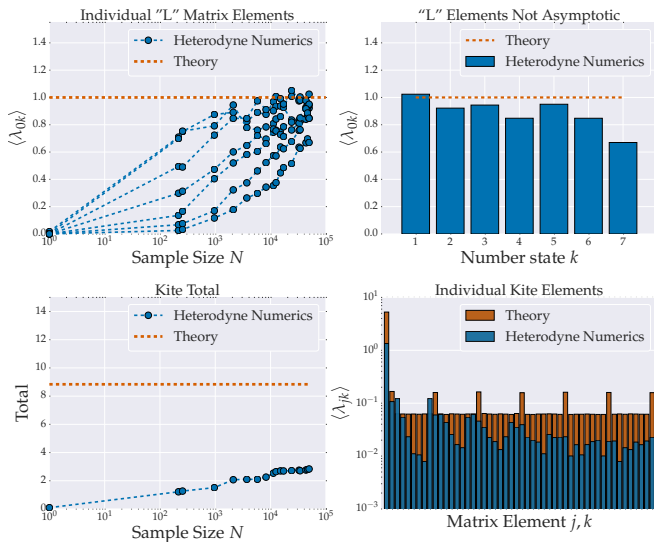


FIG. 10. Examining how our predicted values for  $\langle \lambda_{jk} \rangle$  disagree with simulated heterodyne experiments. We take  $\rho_0 = |0\rangle\langle 0|$  and  $d = 8$ . **Top Left:** The values of  $\langle \lambda_{0k} \rangle$  in the “L” as a function of sample size  $N$ . **Top Right:** Even at the largest  $N$  studied,  $\langle \lambda_{0k} \rangle$  is nontrivially less than 1, especially for the higher number states. **Bottom Left:** The total from the “kite” versus  $N$ . It is clear the total is still growing. **Bottom Right:** The individual “kite” elements  $\langle \lambda_{jk} \rangle$  at the largest  $N$  studied; most are small compared to values they would have in the isotropic case.

**Acknowledgements:** The authors are grateful for those who provide support for the following software packages: iPython/Jupyter [52], matplotlib [53], mpi4py [54], NumPy [55], pandas [56], Python 2.7 [57], seaborn [58], and SciPy [59]. TLS thanks Jonathan A Gross for helpful discussions on software design, coding, and statistics, John King Gamble for useful insights on parallelized computation and feedback on earlier drafts of this paper, and Daniel Seuss for proofreading edits.

Sandia National Laboratories is a multi-mission laboratory managed and operated by Sandia Corporation, a wholly owned subsidiary of Lockheed Martin Corporation, for the U.S. Department of Energy’s National Nuclear Security Administration under contract DE-AC04-94AL85000.

- [1] M. G. A. Paris and J. Rehacek, eds., *Quantum State Estimation* (Springer, Berlin-Heidelberg, 2004).
- [2] N. Reid and D. R. Cox, *International Statistical Review* **83**, 293 (2015).
- [3] L. Wasserman, *All of Statistics: A Concise Course in Statistical Inference* (Springer New York, 2004).
- [4] J. B. Altepeter, E. R. Jeffrey, and P. G. Kwiat, in *Advances in Atomic, Molecular and Optical Physics*, Vol. 52 (Elsevier, 2005) pp. 105–159.
- [5] J. I. Bertrand and P. Bertrand, *Foundations of Physics* **17**, 397 (1987).
- [6] A. I. Lvovsky and M. G. Raymer, *Reviews of Modern Physics* **81**, 299 (2009).
- [7] G. Breitenbach, S. Schiller, and J. Mlynek, *Nature* **387**, 471 (1997).
- [8] U. Leonhardt and H. Paul, *Progress in Quantum Electronics* **19**, 89 (1995).
- [9] F. Motzoi, J. M. Gambetta, P. Rebentrost, and F. K. Wilhelm, *Physical Review Letters* **103**, 110501 (2009).
- [10] R. Fazio, G. Palma, and J. Siewert, *Physical Review Letters* **83**, 5385 (1999).
- [11] K. P. Burnham, *Sociological Methods & Research* **33**, 261 (2004).
- [12] Z. Hradil, *Physical Review A* **55**, R1561 (1997).
- [13] D. F. V. James, P. G. Kwiat, W. J. Munro, and A. G. White, *Physical Review A* **64**, 052312 (2001).
- [14] R. Blume-Kohout, J. O. S. Yin, and S. J. van Enk, *Physical Review Letters* **105**, 170501 (2010).
- [15] L. Le Cam, *Annals of Mathematical Statistics* **41**, 802 (1970).
- [16] L. Le Cam, in *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability*, edited by J. Neyman (1956) pp. 129–156.
- [17] E. J. Candes and T. Tao, *IEEE Transactions on Information Theory* **52**, 5406 (2006).
- [18] S. T. Flammia, D. Gross, Y.-K. Liu, and J. Eisert, *New Journal of Physics* **14** (2012), 10.1088/1367-2630/14/9/095022.
- [19] D. Suess, L. Rudnicki, and D. Gross, arXiv:1608.00374 (2016).
- [20] A. Carpentier, J. Eisert, D. Gross, and R. Nickl, arXiv:1504.03234v2 (2015).
- [21] L. Schwarz and S. J. van Enk, *Physical Review A* **88**, 032318 (2013).
- [22] M. Guta, T. Kypraios, and I. Dryden, *New Journal of Physics* **14** (2012), 10.1088/1367-2630/14/10/105002.
- [23] S. J. van Enk and R. Blume-Kohout, *New Journal of Physics* **15**, 025024 (2013).
- [24] J. O. S. Yin and S. J. van Enk, *Physical Review A* **83**, 062110 (2011).
- [25] T. Moroder, M. Kleinmann, P. Schindler, T. Monz, O. Guhne, and R. Blatt, *Physical Review Letters* **110**, 180401 (2013).
- [26] L. Knips, C. Schwemmer, N. Klein, J. Reuter, G. Tóth, and H. Weinfurter, arXiv:1512.06866 (2015).
- [27] By which we mean how the distribution of  $\hat{\rho}_{\text{MLE}}$  differs



- from that of  $\hat{\rho}$ .
- [28] J. A. Smolin, J. M. Gambetta, and G. Smith, Physical Review Letters **108**, 070502 (2012).
  - [29] J. Neyman and E. S. Pearson, Philosophical Transactions of the Royal Society of London **231**, 289 (1933).
  - [30] If  $\mathcal{M}_1 \subset \mathcal{M}_2$ , then the maximum likelihood of  $\mathcal{M}_2$  is at least as high as that of  $\mathcal{M}_1$ .
  - [31] For the *classical* probability simplex, the relationship between  $\lambda$  and mean-squared error is *exact* as  $\epsilon \rightarrow 0$ . Quantum state spaces have non-zero curvature, even as  $\epsilon \rightarrow 0$ , so the relationship is approximate.
  - [32] S. S. Wilks, The Annals of Mathematical Statistics **9**, 60 (1938).
  - [33] Because  $\rho_0$  is *fixed*, the number of free parameters is 0.
  - [34] That is, we let  $\text{Tr}(\delta)$  fluctuate as well.
  - [35] Y. V. Fyodorov, Arxiv:math-ph/0412017 (2004).
  - [36] E. P. Wigner, Annals of Mathematics **67**, 325 (1958).
  - [37] T. Tao and V. Vu, Random Matrices: Theory and Applications **2** (2013), 10.1142/S201032631350007X.
  - [38] A. W. van der Vaart, *Asymptotic Statistics* (Cambridge University Press, 2000).
  - [39] This justifies the assumption that  $\rho_{jj} + \delta_{jj} - q > 0$ .
  - [40] A. I. Lvovsky, H. Hansen, T. Aichele, O. Benson, J. Mlynek, and S. Schiller, Physical Review Letters **87**, 050402 (2001).
  - [41] K. Husimi, in *Proceedings of the Physico-Mathematical Society of Japan*, Vol. 22 (1940) pp. 264 – 314.
  - [42] The model  $\mathcal{M}_1$  is trivial, as  $\mathcal{M}_1 = \{|0\rangle\langle 0|\}$ . This model will almost always be wrong, in general.
  - [43] H. Akaike, IEEE Transactions on Automatic Control **19**, 716 (1974).
  - [44] G. Schwarz, The Annals of Statistics **6**, 461 (1978).
  - [45] R. E. Kass and A. E. Raftery, Journal of the American Statistical Association **90**, 773 (1995).
  - [46] S. Glancy, E. Knill, and M. Girard, New Journal of Physics **14** (2012), 10.1088/1367-2630/14/9/095017.
  - [47] B. Efron, The Annals of Statistics **7**, 1 (1979).
  - [48] J. Higgins, *An Introduction to Modern Nonparametric Statistics* (Brooks/Cole, 2004).
  - [49] A. Steffens, C. Riofrio, W. McCutcheon, I. Roth, B. A. Bell, A. McMillan, M. S. Tame, J. G. Rarity, and J. Eisert, arXiv:1611.01189 (2016).
  - [50] A. Kalev and C. H. Baldwin, arXiv:1511.01433v1 (2015).
  - [51] A. Kalev, R. L. Kosut, and I. H. Deutsch, npj Quantum Information **1**, 15018 (2015).
  - [52] F. Pérez and B. E. Granger, Computing in Science and Engineering **9**, 21 (2007).
  - [53] J. D. Hunter, Computing in Science and Engineering **9**, 90 (2007).
  - [54] L. Dalcin, R. Paz, and M. Storti, Journal of Parallel and Distributed Computing **65**, 1108 (2005).
  - [55] S. Van Der Walt, S. C. Colbert, and G. Varoquaux, Computing in Science and Engineering **13**, 22 (2011).
  - [56] W. McKinney, in *Proceedings of the 9th Python in Science Conference*, edited by S. van der Walt and J. Millman (2010) pp. 51–56.
  - [57] G. van Rossum, “Python Language Reference,” (1995).
  - [58] M. Waskom, “seaborn,” (2016).
  - [59] T. E. Oliphant, Computing in Science and Engineering **9**, 10 (2007).